

## UvA-DARE (Digital Academic Repository)

# Operator preconditioning and space-time methods for parabolic evolution equations

van Venetië, R.

Publication date 2021 Document Version Final published version

#### Link to publication

#### Citation for published version (APA):

van Venetië, R. (2021). *Operator preconditioning and space-time methods for parabolic evolution equations*. [Thesis, fully internal, Universiteit van Amsterdam].

#### General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

#### **Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: https://uba.uva.nl/en/contact, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

Ореі	ator	oreco	nditic	oning	
and	Space	-time	meth	nods	
for p	arabo	olic ev	oluti	on	
equa	tions				
			Ray V	/mond van enetië	

### **Operator preconditioning**

### AND

## Space-time methods for parabolic evolution equations

Raymond van Venetië

About the cover: the displayed mesh is found by running the adaptive method from Chapter 9 on the 'boundary' of this book. The solution is strongly singular for t = 0, which explains the refinements towards the bottom of the cover.

This research was funded by the Netherlands Organization for Scientific Research (NWO) under contract. no. 613.001.652.

Printed by GVO drukkers & vormgevers ISBN: 978-94-6332-788-6 © 2021 Raymond van Venetië

## Operator preconditioning and Space-time methods for parabolic evolution equations

#### ACADEMISCH PROEFSCHRIFT

ter verkrijging van de graad van doctor aan de Universiteit van Amsterdam op gezag van de Rector Magnificus prof. dr. ir. K.I.J. Maex ten overstaan van een door het College voor Promoties ingestelde commissie, in het openbaar te verdedigen in de Agnietenkapel op dinsdag 28 september 2021, te 15.00 uur door

> Raymond van Venetië geboren te Zoetermeer

### Promotiecommissie

Promotor:	prof. dr. R.P. Stevenson	Universiteit van Amsterdam
Copromotor:	dr. J.H. Brandts	Universiteit van Amsterdam
Overige leden:	prof. dr. H. Peters prof. dr. A.J. Homburg dr. C.C. Stolk dr. G. Gantner prof. dr. U. Langer dr. C.A. Urzúa-Torres	Universiteit van Amsterdam Universiteit van Amsterdam Universiteit van Amsterdam Universiteit van Amsterdam Johannes Kepler University Linz Technische Universiteit Delft

Faculteit der Natuurwetenschappen, Wiskunde en Informatica

## Contents

1	Introduction	1				
	1.1 Numerical methods for operator equations	. 1				
	1.2 About Part I: Operator preconditioning	. 5				
	1.3 About Part II: Parabolic evolution equations	. 9				
I	Operator preconditioning	13				
2	Problems of negative order	15				
	2.1 Introduction	. 15				
	2.2 Operator preconditioning	. 18				
	2.3 Piecewise constant discretization space	. 21				
	2.4 Continuous piecewise linear discretization space	. 34				
	2.5 Higher order case	. 38				
	2.6 Numerical experiments	. 45				
	2.7 Conclusion	. 48				
3	Problems of negative order: preconditioning at linear cost	49				
	3.1 Introduction	. 49				
	3.2 Operator preconditioning	. 51				
	3.3 An operator $B_{\mathcal{T}}^{\mathscr{S}}$ of multi-level type	. 54				
	3.4 Manifold case	. 58				
	3.5 Numerical experiments	. 60				
4	Problems of positive order					
	4.1 Introduction	. 63				
	4.2 Operator preconditioning	. 66				
	4.3 Continuous piecewise linear discretization space	. 68				
	4.4 Construction of $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$	. 71				
	4.5 Extensions	. 79				
	4.6 Numerical experiments	. 82				
	4.7 Conclusion	. 85				
5	The simplest case					
	5.1 Introduction	. 87				
	5.2 Construction of $D_{\mathcal{T}}$ in the domain case	. 89				
	5.3 Manifold case	. 94				
	5.4 Numerical experiments	. 94				

	5.5	Conclusion	97			
Π	Spa	ce-time methods for parabolic evolution equations	99			
6	An adaptive method					
	6.1	Introduction	101			
	6.2	Space-time formulations of a parabolic evolution problem	106			
	6.3	Discretizations	108			
	6.4	Convergent adaptive solution method	111			
	6.5	Wavelets-in-time tensorized with finite-elements-in-space	122			
	6.6	A concrete realization	131			
	6.7	Numerical experiments	140			
	6.8	Conclusion	147			
7	Ada	ptivity: an efficient implementation	149			
	7.1	Introduction	149			
	7.2	Space-time adaptivity for a parabolic model problem	151			
	7.3	The application of linear operators in linear complexity	157			
	7.4	The heat equation and practical realization	165			
	7.5	Implementation	171			
	7.6	Numerical experiments	176			
	7.7		181			
	7.A	Proofs of Theorems in $\frac{5}{.3}$	181			
8	A pa	arallel algorithm	185			
	8.1	Introduction	185			
	8.2	Quasi-optimal approximations to the parabolic problem	187			
	8.3	Solving efficiently on tensor-product discretizations	189			
	8.4	A concrete setting: the reaction-diffusion equation	192			
	8.5	Numerical experiments	196			
	8.6	Conclusion	198			
9	Ada	ptive BEM for the heat equation	201			
	9.1	Introduction	201			
	9.2	Preliminaries	203			
	9.3	A posteriori error estimation	207			
	9.4	Numerical experiments	212			
Su	mma	ıry	219			
Sa	Samenvatting					
Da	Dankwoord					
Bi	Bibliography					

Partial differential equations are used for the modeling of a wide variety of (natural) phenomena appearing in biology, physics, chemistry, engineering, finance and many more fields. Typically, closed-form solutions to these differential equations are unknown or do not exist, despite their importance in practice. A remedy is provided by numerical methods that construct *approximations* of the solutions to these differential equations. Ideally, these numerical methods provide good approximations at a small computational cost. This is the main topic of investigation in this thesis, where we will focus on *linear* partial differential equations of elliptic or parabolic type.

The contents of this thesis can be roughly divided into two parts. In the first part we study *preconditioning*, a technique that is used to accelerate solvers for systems of linear equations arising in our approximation schemes. In the second part, we focus on (adaptive) numerical methods for parabolic evolution equations in a simultaneous *space-time* approach.

#### 1.1 Numerical methods for operator equations

We start with a brief description of the general setting that we study in this thesis. For a thorough introduction, we refer the reader to the literature, e.g. [Bra01, EG04, Ste08a].

#### 1.1.1 Operator equations

We consider *linear* operator equations of the following type. For some Hilbert space  $\mathscr{V}$ , a linear map  $A: \mathscr{V} \to \mathscr{V}'$  and data  $f \in \mathscr{V}'$ , we seek  $u \in \mathscr{V}$  that solves

(1.1) Au = f or equivalently (Au)(v) = f(v)  $(v \in \mathscr{V}).$ 

We require the problem (1.1) to be *well-posed*, meaning that there is a unique solution that depends continuously on the given data. Typically,  $\mathcal{V}$  is a Sobolev space of functions on some bounded domain, and the operator A corresponds to the variational formulation of some partial differential equation, but we also encounter situations where A is a boundary integral operator.

Let us restrict ourselves to A being a *bounded* and *coercive* operator, meaning that for all  $v \in \mathcal{V}$  we have  $||Av||_{\mathcal{V}'} \leq C ||v||_{\mathcal{V}}$  and  $(Av)(v) \geq \alpha ||v||_{\mathcal{V}}^2$ , for some  $C < \infty$  and  $\alpha > 0$ . In this setting, the Lax–Milgram theorem asserts that the problem (1.1) is well-posed with  $||A^{-1}||_{\mathcal{L}(\mathcal{V}',\mathcal{V})} \leq 1/\alpha$ .

*Example* 1.1.1. As a model problem for (elliptic) partial differential equations, one may look at Poisson's equation with homogeneous Dirichlet boundary conditions. For some bounded Lipschitz domain  $\Omega \subset \mathbb{R}^d$ , this is the problem of finding *u* that solves

$$-\Delta u = f \quad \text{on } \Omega, \qquad u = 0 \quad \text{on } \partial \Omega,$$

for given forcing data f.

In variational formulation this problem reads as finding u from the Sobolev space  $H_0^1(\Omega) := \{v \in H^1(\Omega) : v | \partial \Omega = 0\}$  that satisfies

(1.2) 
$$(Au)(v) := \langle \nabla u, \nabla v \rangle_{L_2(\Omega)} = \langle f, v \rangle_{L_2(\Omega)} \quad (v \in H^1_0(\Omega)).$$

The operator A is a bounded linear map  $H_0^1(\Omega) \to H_0^1(\Omega)'$ . Moreover, the Poincaré inequality shows that A is coercive. We conclude that (1.2) is a well-posed operator equation for forcing data  $f \in H_0^1(\Omega)'$ , and falls in the abstract setting of (1.1).

*Example* 1.1.2. Another model problem of interest is Laplace's equation with inhomogeneous Dirichlet boundary conditions. For some bounded Lipschitz domain  $\Omega \subset \mathbb{R}^d$ , this is the problem of finding *u* that solves

(1.3) 
$$\Delta u = 0 \quad \text{on } \Omega, \qquad u = g \quad \text{on } \partial \Omega,$$

for given Dirichlet boundary data *g*. This problem admits an (alternative) weak formulation in terms of boundary integral operators.

This requires some notation. Set  $\Gamma := \partial \Omega$ , and consider the fractional Sobolev space  $H^{\frac{1}{2}}(\Gamma)$ , which can be seen as the trace space of  $H^{1}(\Omega)$ , with its dual that we denote by  $H^{-\frac{1}{2}}(\Gamma)$ . Let *G* be the fundamental solution to the Laplace equation (1.3), e.g., for d = 2 we have  $G(\cdot) := -\frac{1}{2\pi} \log |\cdot|$ . We can now introduce the *Single Layer* integral operator  $A: H^{-\frac{1}{2}}(\Gamma) \to H^{\frac{1}{2}}(\Gamma)$  by

$$(A\tilde{u})(\tilde{v}) := \int_{\Gamma} \tilde{v}(\boldsymbol{x}) \int_{\Gamma} G(\boldsymbol{x} - \boldsymbol{y}) \tilde{u}(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y} \, \mathrm{d}\boldsymbol{x} \quad \left(\tilde{u}, \tilde{v} \in H^{-\frac{1}{2}}(\Gamma)\right).$$

The Single Layer operator A is bounded and coercive, so for  $g \in H^{\frac{1}{2}}(\Gamma)$  we may consider the well-posed problem of finding  $\tilde{u} \in H^{-\frac{1}{2}}(\Gamma)$  that solves

$$A\tilde{u} = g.$$

It turns out that this problem is equivalent to the original problem (1.3), in that we can recover the weak solution  $u \in H^1(\Omega)$  via the representation formula  $u(\boldsymbol{x}) = \int_{\Gamma} G(\boldsymbol{x} - \boldsymbol{y}) \tilde{u}(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y}.$ 

Finally, we note that solving the Laplace problem with Neumann boundary data—so replacing the boundary condition by  $\frac{\partial u}{\partial n} = g$  in (1.3)—allows for a similar weak formulation in terms of boundary integral operators. Here one solves Bu = g, with the *Hypersingular* integral operator  $B: H^{\frac{1}{2}}(\Gamma) \to H^{-\frac{1}{2}}(\Gamma)$ . This operator maps in the 'opposite direction' of *A*, which will be important later in this introduction.

#### 1.1.2 Numerical approximation

For the numerical approximation of the operator equation (1.1) we can use the Ritz–Galerkin method. For some finite-dimensional subspace  $\mathscr{V}_{\mathcal{T}} \subset \mathscr{V}$ , also called the *trial space*, we consider the discretized operator  $A_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{V}_{\mathcal{T}}'$  given by  $(A_{\mathcal{T}}u)(v) := (Au)(v)$   $(u, v \in V_{\mathcal{T}})$  and the naturally embedded right hand side  $f \in \mathscr{V}_{\mathcal{T}}'$ . The Galerkin approximation of (1.1) is then the function  $u_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}$  that solves

(1.4) 
$$A_{\mathcal{T}}u_{\mathcal{T}} = f.$$

Since  $\mathscr{V}_{\mathcal{T}}$  is a closed subspace of  $\mathscr{V}$ , the operator  $A_{\mathcal{T}}$  is again coercive, and hence the Lax–Milgram theorem asserts that this discretized problem is well-posed.

The advantage of the Galerkin method is that it provides us with a concrete and tight error bound. Indeed, Céa's lemma shows that the error satisfies

$$\|u - u_{\mathcal{T}}\|_{\mathscr{V}} \le \|A\|_{\mathcal{L}(\mathscr{V},\mathscr{V}')} \|A^{-1}\|_{\mathcal{L}(\mathscr{V}',\mathscr{V})} \inf_{v \in \mathscr{V}_{\mathcal{T}}} \|u - v\|_{\mathscr{V}},$$

which tells us that  $u_{\mathcal{T}}$  is a *quasi-best* approximation to u from  $\mathscr{V}_{\mathcal{T}}$ .

To design a convergent numerical method, we actually have to construct a family of trial spaces  $(\mathscr{V}_{\mathcal{T}})_{\mathcal{T}\in\mathbb{T}}$ . Here we wish to achieve two things simultaneously. On the one hand, we want our method to have a good approximation rate, that is, we want the approximations to satisfy  $||u - u_{\mathcal{T}}||_{\mathscr{V}} \leq C(\dim \mathscr{V}_{\mathcal{T}})^{-\gamma} ||u||_{\mathscr{V}}$  for all  $\mathcal{T} \in \mathbb{T}$  with some large  $\gamma > 0$  (and some constant C > 0). On the other hand, we want the number of operations required to calculate  $u_{\mathcal{T}}$  from (1.4) to be small, so ideally of order  $\mathcal{O}(\dim \mathscr{V}_{\mathcal{T}})$ . Both these motives play a central role in this thesis.

In the *finite element method*, the trial spaces  $\mathscr{V}_{\mathcal{T}}$  are constructed as spaces of (dis)continuous piecewise polynomials of fixed degree with respect to meshes  $\mathcal{T}$  of the underlying domain. If the solution u is *smooth*, i.e., its derivatives of sufficiently high order are bounded in a suitable norm, then optimal approximation rates are often obtained by simply considering a family of trial spaces with respect to *quasi-uniform* meshes  $\mathcal{T}$ , i.e., meshes having a uniform mesh width. The situation changes drastically when the (unknown) solution u contains *singularities*, which may be induced by the geometry or the data. In such cases, the approximation rate offered by trial spaces with respect to quasi-uniform meshes can drop significantly, making the numerical scheme

converge slowly. Luckily, in many situations the approximation rate can be improved significantly by considering trial spaces with respect to meshes that are *locally refined* at these singularities.

Throughout this thesis we will focus on this latter situation, where the solution u contains singularities. Clearly, in terms of accuracy, using trial spaces adapted to the singularities is preferable. This comes at a price however, as the mathematics and implementation of such adaptive numerical methods is often more difficult than for the quasi-uniform case.

#### 1.1.3 Solving the discretized system

An imported question is how to *efficiently* solve (1.4). To analyze this, we first reformulate the problem in coordinates. With  $\Phi_{\mathcal{T}} := \{\phi_1, \ldots, \phi_n\}$  being a basis for  $\mathscr{V}_{\mathcal{T}}$ , one infers that solving (1.4) for  $u_{\mathcal{T}} = u^{\top} \Phi_{\mathcal{T}}$ , is equivalent to finding  $u \in \mathbb{R}^n$  that solves

$$(1.5) A_{\mathcal{T}} u = f$$

with the *system matrix*  $A_{\mathcal{T}} \in \mathbb{R}^{n \times n}$  and right hand side  $f \in \mathbb{R}^n$  given by

$$\boldsymbol{A}_{\mathcal{T}} := (A_{\mathcal{T}} \Phi_{\mathcal{T}})(\Phi_{\mathcal{T}}) = [(A_{\mathcal{T}} \phi_j)(\phi_i)]_{ij}, \qquad \boldsymbol{f} := f(\Phi_{\mathcal{T}}) = [f(\phi_i)]_i.$$

The computational complexity of solving the above matrix-vector system using a direct method, e.g. using *LU*-factorization, is  $O(n^3)$ . This is prohibitively expensive, as we are interested in constructing methods that run in optimal (linear) time. Another problem is that direct methods require the matrix  $A_{\mathcal{T}}$  to be available, which is not the case for *matrix-free* methods, where one has access to only the application of  $A_{\mathcal{T}}$ . An example of such a matrix-free method will be given in this thesis, where we devise an algorithm to apply a system matrix  $A_{\mathcal{T}}$  in linear complexity, even though the matrix itself is not sparsely populated.

A first efficiency gain can be obtained by solving (1.5) approximately yielding some  $\mathbb{R}^n \ni \hat{u} \approx u$ . For  $\hat{u}_T := \hat{u}^\top \Phi_T$  we have

$$\|u - \hat{u}_{\mathcal{T}}\|_{\mathscr{V}} \leq \underbrace{\|u - u_{\mathcal{T}}\|_{\mathscr{V}}}_{\text{discretization error}} + \underbrace{\|u_{\mathcal{T}} - \hat{u}_{\mathcal{T}}\|_{\mathscr{V}}}_{\text{algebraic error}},$$

so as long as the algebraic error is dominated by the discretization error, our approximation  $\hat{u}_{\tau}$  is still quasi-optimal. Such approximate solutions to (1.5) are given by *iterative solvers*, being methods that produce a sequence of vectors  $(u_1, u_2, ...)$  converging to the solution u.

If the matrix  $A_{\mathcal{T}}$  is symmetric and positive definite, as when the operator A is self-adjoint (A = A') and coercive, then of particular interest is the *Conjugate Gradient* (CG) method. This iterative solver has favourable convergence properties and has minimal computation costs, i.e., the cost of a single iteration is dominated by the cost of a single application of  $A_{\mathcal{T}}$ . Write  $\rho(\cdot)$  for

the spectral radius of an operator, and denote the spectral condition number by  $\kappa_S(\mathbf{A}_{\mathcal{T}}) = \rho(\mathbf{A}_{\mathcal{T}})\rho(\mathbf{A}_{\mathcal{T}}^{-1})$ . Starting from an initial guess, the number of iterations required by CG to reduce the initial algebraic error by a factor  $\epsilon$  is bounded by  $\sqrt{\kappa_S(\mathbf{A})} \log(1/\epsilon)$ .

Unfortunately, for standard finite (or boundary) element bases, the system matrices  $A_{\mathcal{T}}$  are (generally) ill-conditioned, in the sense that the condition number  $\kappa_S(A_{\mathcal{T}})$  increases for decreasing mesh width. This implies that the number of iterations required by CG to reach a certain accuracy increases upon mesh refinement.

One way of fixing this conditioning issue, is to make use of *preconditioning*: instead of solving (1.5), one considers  $G_{\mathcal{T}}A_{\mathcal{T}}u = G_{\mathcal{T}}f$  where  $G_{\mathcal{T}}$  is some preconditioner, i.e., an approximation of  $A_{\mathcal{T}}^{-1}$  that can be applied efficiently. The number of iterations required by CG applied to this preconditioned system to reduce the initial algebraic error by a factor  $\epsilon$ , is bounded by  $\sqrt{\kappa_S(G_{\mathcal{T}}A_{\mathcal{T}})}\log(1/\epsilon)$ .

#### 1.2 About Part I: Operator preconditioning

In the first part of this thesis we focus on constructing *uniformly optimal* preconditioners  $G_{\mathcal{T}}$  for the family  $\mathcal{T} \in \mathbb{T}$ . This means that we want the condition number of the preconditioned matrix  $\kappa_S(G_{\mathcal{T}}A_{\mathcal{T}})$ , and therefore also the number of iterations required by CG, to be bounded independently of the trial space  $\mathscr{V}_{\mathcal{T}}$ .

One technique for constructing such preconditioners is so-called *operator preconditioning* ([Hip06]). This approach hinges on the availability of an *opposite order* operator *B*, being a bounded and coercive linear map from  $\mathcal{V}'$  to  $\mathcal{V}$ . On a continuous level we have that *BA* is a boundedly invertible map  $\mathcal{V} \to \mathcal{V}$ , suggesting that *B* may be used to construct a preconditioner for  $A_{\mathcal{T}} : \mathcal{V}_{\mathcal{T}} \to \mathcal{V}'_{\mathcal{T}}$ .

Let  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W} := \mathscr{V}'$  be some finite dimensional subspace, and consider the discretized operator  $B_{\mathcal{T}} : \mathscr{W}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$  given by  $(B_{\mathcal{T}}u)(v) := (Bu)(v) (u, v \in \mathscr{W}_{\mathcal{T}})$ . If one additionally has a boundedly invertible operator  $D_{\mathcal{T}} : \mathscr{V}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$ , then the composition  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D'_{\mathcal{T}})^{-1} : \mathscr{V}_{\mathcal{T}}' \to \mathscr{V}_{\mathcal{T}}$  is a bounded and coercive mapping, and therefore serves as a preconditioner for  $A_{\mathcal{T}}$ . See also the diagram:

(1.6) 
$$\begin{array}{ccc} \mathscr{V}_{\mathcal{T}} & \xrightarrow{A_{\mathcal{T}}} & \mathscr{V}_{\mathcal{T}}' \\ D_{\mathcal{T}}^{-1} \uparrow & & \downarrow^{(D_{\mathcal{T}}')^{-1}} \\ \mathscr{W}_{\mathcal{T}}' & \xleftarrow{B_{\mathcal{T}}} & \mathscr{W}_{\mathcal{T}} \end{array}$$

The typical example to keep in mind is where  $A \colon H^{-\frac{1}{2}}(\Gamma) \to H^{\frac{1}{2}}(\Gamma)$  is the Single Layer operator and  $B \colon H^{\frac{1}{2}}(\Gamma) \to H^{-\frac{1}{2}}(\Gamma)$  the Hypersingular operator from Example 1.1.2.

Now, for the construction of a suitable  $D_{\mathcal{T}}$ , we assume  $\mathscr{H}$  to be some Hilbert space (we take  $L_2(\Gamma)$  in the above example) for which we have a Gelfand triple

$$\mathscr{W} \hookrightarrow \mathscr{H} \simeq \mathscr{H}' \hookrightarrow \mathscr{W}'.$$

If the trial spaces satisfy  $\mathscr{V}_{\mathcal{T}} \subset \mathscr{H}$ , we can consider the operator  $D_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$ defined by  $(D_{\mathcal{T}}v)(w) \coloneqq \langle v, w \rangle_{\mathscr{H}} \quad (v \in \mathscr{V}_{\mathcal{T}}, w \in \mathscr{W}_{\mathcal{T}})$ . This is a uniformly boundedly invertible operator if the subspaces  $\mathscr{W}_{\mathcal{T}}$  satisfy

(1.7) 
$$\dim \mathscr{W}_{\mathcal{T}} = \dim \mathscr{V}_{\mathcal{T}}$$

and

(1.8) 
$$\inf_{\mathcal{T}\in\mathbb{T}}\inf_{0\neq v\in\mathscr{V}_{\mathcal{T}}}\sup_{0\neq w\in\mathscr{W}_{\mathcal{T}}}\frac{\langle v,w\rangle_{\mathscr{H}}}{\|v\|_{\mathscr{V}}\|w\|_{\mathscr{W}}}>0.$$

Assume these constraints to be satisfied. By equipping  $\mathcal{V}_{\mathcal{T}}$  and  $\mathcal{W}_{\mathcal{T}}$  with bases  $\Phi_{\mathcal{T}}$  and  $\Psi_{\mathcal{T}}$ , respectively, the matrix representation of the preconditioned system reads as

$$(1.9) D_{\mathcal{T}}^{-1} B_{\mathcal{T}} D_{\mathcal{T}}^{-\top} A_{\mathcal{T}},$$

with system matrices  $A_{\mathcal{T}} := (A_{\mathcal{T}} \Phi_{\mathcal{T}})(\Phi_{\mathcal{T}})$  and  $B_{\mathcal{T}} := (B_{\mathcal{T}} \Psi_{\mathcal{T}})(\Psi_{\mathcal{T}})$ , and 'generalized mass matrix'  $D_{\mathcal{T}} := \langle \Phi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{\mathscr{H}}$ . The spectral condition number of this preconditioned matrix system equals that of  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D_{\mathcal{T}}')^{-1}A_{\mathcal{T}}$ , and is thus uniformly bounded.

The real challenge is constructing suitable subspaces  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W}$  that satisfy (1.7) and (1.8). Moreover, we require a basis  $\Psi_{\mathcal{T}}$  for  $\mathscr{W}_{\mathcal{T}}$  that allows the preconditioner (1.9) to be applied efficiently.

Special attention has to be paid to the inverse matrix  $D_{\mathcal{T}}^{-1}$  appearing in the preconditioner. If the matrix  $D_{\mathcal{T}}$  is not diagonal, its inverse has to be approximated, and it can generally be expected that, in order to obtain a uniform preconditioner, the accuracy with which  $D_{\mathcal{T}}^{-1}$  has to be approximated increases with a decreasing minimal mesh-size. As a result, an application of  $D_{\mathcal{T}}^{-1}$  cannot be expected to execute in linear time.

In this thesis, we will propose suitable spaces  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W}$ , and circumvent this latter issue by constructing bases  $\Psi_{\mathcal{T}}$  that are  $\mathscr{H}$ -orthogonal to  $\Phi_{\mathcal{T}}$ , making  $D_{\mathcal{T}}$  a diagonal matrix whose inverse can be exactly evaluated.

#### 1.2.1 Contributions and Outline of Part I

The contents of Chapters 2–5 is essentially that of the following papers:

[SvV20a] R.P. Stevenson and R. van Venetië. Uniform preconditioners for problems of negative order. *Mathematics of Computation*, 89(322):645–674, 2020.

- [SvV21a] R.P. Stevenson and R. van Venetië. Uniform preconditioners of linear complexity for problems of negative order. *Computational Methods in Applied Mathematics*, 21(2):469–478, 2021.
- [SvV20b] R.P. Stevenson and R. van Venetië. Uniform preconditioners for problems of positive order. *Computers & Mathematics with Applications*, 79(12):3516–3530, 2020.
- [SvV21b] R.P. Stevenson and R. van Venetië. Operator preconditioning: the simplest case. Submitted to *Applied Numerical Mathematics*, 2021.

On average, the authors contributed equally to these works.

#### Chapter 2 ([SvV20a])

For some domain (or manifold)  $\Omega$  and  $s \in [0, 1]$ , we consider the fractional Sobolev space  $H^s(\Omega)$  (possibly with homogeneous Dirichlet boundary conditions incorporated) and its dual that we denote here by  $H^{-s}(\Omega)$ .

In this chapter we consider operator preconditioning for a bounded and coercive operator  $A: H^{-s}(\Omega) \to H^s(\Omega)$ , so of *negative* order -2s, discretized by a family of trial spaces  $\mathscr{V}_{\mathcal{T}} \subset H^{-s}(\Omega)$  being discontinuous piecewise constants w.r.t.  $\mathcal{T}$ . We consider a family  $\mathbb{T}$  of uniformly shape regular, possibly locally refined, meshes of  $\Omega$ . In order to apply the aforementioned framework, we assume availability of a suitable opposite order operator  $B: H^s(\Omega) \to H^{-s}(\Omega)$ .

We propose a family of subspaces  $\mathscr{W}_{\mathcal{T}} = \operatorname{span} \Psi_{\mathcal{T}} \subset H^s(\Omega)$  satisfying both (1.7) and (1.8). We achieve this by constructing  $\Psi_{\mathcal{T}}$  as a collection that is  $L_2(\Omega)$ -biorthogonal to  $\Phi_{\mathcal{T}}$ , the (canonical) piecewise constant basis for  $\mathscr{V}_{\mathcal{T}}$ . As a consequence, the matrix  $D_{\mathcal{T}} := \langle \Phi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{L_2(\Omega)}$  is *diagonal*, and thus its inverse, which appears in the preconditioner (1.9), can be evaluated exactly.

The functions  $\Psi_{\mathcal{T}}$  are constructed in  $\mathscr{S}_{\mathcal{T}} \oplus \mathscr{B}_{\mathcal{T}}$ , where  $\mathscr{S}_{\mathcal{T}}$  is the space of continuous piecewise linears w.r.t.  $\mathcal{T}$ , and  $\mathscr{B}_{\mathcal{T}}$  is a space containing bubble functions. This allows us to construct a bounded and coercive operator  $B_{\mathcal{T}} \colon \mathscr{W}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$  as the sum of  $B_{\mathcal{T}}^{\mathscr{S}} \colon \mathscr{S}_{\mathcal{T}} \to \mathscr{S}_{\mathcal{T}}'$ , being the opposite order operator B discretized on  $\mathscr{S}_{\mathcal{T}} \subset H^s(\Omega)$ , and an invertible diagonal scaling operator on the bubble space. Besides the cost of the discretized operator  $B_{\mathcal{T}}^{\mathscr{S}}$ , the cost of the resulting preconditioner scales linearly in dim  $\mathscr{V}_{\mathcal{T}}$ .

Our approach has a few advantages over earlier proposals: it does not require the inverse of a non-diagonal matrix; it applies without any mildly grading assumption on the mesh; and it does not require a barycentric refinement of the mesh underlying the trial space. Furthermore, we will show that our approach extends to the general case where  $\mathscr{V}_{\mathcal{T}}$  is chosen as a space of (dis)continuous piecewise polynomials of any order.

#### Chapter 3 ([SvV21a])

We continue with the setting from the previous chapter. In this chapter we construct a multi-level type operator that both fulfills the role of the opposite order operator *B* and can be applied in optimal *linear complexity*. For this construction, we require  $\mathbb{T}$  to be a family of conforming partitions created by newest vertex bisection. Together with the results from the previous chapter, this provides uniformly optimal preconditioners for negative order operators *A* discretized on a family of possibly locally refined meshes, that can be applied in linear complexity.

#### Chapter 4 ([SvV20b])

In this chapter consider the operator preconditioning framework for *positive* order operators, that is, we switch the roles of *A* and *B*.

More precisely, we construct preconditioners for a bounded and coercive operator  $A: H^s(\Omega) \to H^{-s}(\Omega)$ , being of positive order 2*s*, discretized by  $\mathscr{V}_{\mathcal{T}} \subset H^s(\Omega)$  being continuous piecewise linears w.r.t  $\mathcal{T}$ . Again, the mesh family  $\mathbb{T}$  is supposed to be uniformly shape regular, which allows for locally refined meshes. Assuming the availability of an opposite order operator  $B: H^{-s}(\Omega) \to H^s(\Omega)$ , we explore the operator preconditioning framework, and aim to get results similar to the setting studied in Chapter 2.

We introduce a family of subspaces  $\mathscr{W}_{\mathcal{T}} \subset H^{-s}(\Omega)$  satisfying both (1.7) and (1.8). Similarly as before, we do this by constructing  $\mathscr{W}_{\mathcal{T}}$  as the span of a collection  $\Psi_{\mathcal{T}}$  that is  $L_2(\Omega)$ -biorthogonal to  $\Phi_{\mathcal{T}}$ , the Lagrange basis for  $\mathscr{V}_{\mathcal{T}}$ . Besides an easy proof of the inf-sup condition (1.8), this biorthogonality has the advantage that the matrix  $D_{\mathcal{T}} := \langle \Phi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{L_2(\Omega)}$  is *diagonal*, and thus its inverse, which appears in the preconditioner (1.9), can be evaluated exactly.

Since  $\mathscr{W}_{\mathcal{T}}$  is a non-standard discretization space, we wish to simplify the implementation of an operator  $B_{\mathcal{T}} \colon \mathscr{W}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$ , similarly as the construction given in Chapter 2. We achieve this by constructing  $\mathscr{W}_{\mathcal{T}}$  as a subspace of  $\mathscr{U}_{\mathcal{T}} \oplus \mathscr{B}_{\mathcal{T}}$ , where  $\mathscr{U}_{\mathcal{T}} \subset H^{-s}(\Omega)$  is the space of discontinuous piecewise constants w.r.t.  $\mathcal{T}$ , and  $\mathscr{B}_{\mathcal{T}}$  is some bubble space for which the  $H^{-s}(\Omega)$ -norm is equivalent to a weighted  $L_2(\Omega)$ -norm. This allows us to construct a suitable  $B_{\mathcal{T}} \colon \mathscr{W}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$  as the sum of  $B_{\mathcal{T}}^{\mathscr{U}} \colon \mathscr{U}_{\mathcal{T}} \to \mathscr{U}_{\mathcal{T}}'$ , being our opposite order operator discretized on  $\mathscr{U}_{\mathcal{T}}$ , and an invertible diagonal scaling operator on the bubble space. Besides the cost of the discretized operator  $B_{\mathcal{T}}^{\mathscr{U}}$ , the cost of the resulting preconditioner scales linearly in dim  $\mathscr{V}_{\mathcal{T}}$ .

This construction has the same advantages as that of Chapter 2: it avoids the inverse of a non-diagonal matrix, it circumvents the need for a barycentric refinement of the mesh, it applies without any mildly grading assumption on the mesh, and it can be extended to higher order trial spaces.

#### Chapter 5 ([SvV21b])

By restricting the setting from the previous chapters, we are able to construct uniform preconditioners with an even simpler implementation. Consider some closed manifold (or domain)  $\Omega$  and trial spaces  $\mathscr{V}_{\mathcal{T}}$  that are *continuous* piecewise polynomials of some fixed degree w.r.t.  $\mathcal{T}$ . Let some bounded and coercive operators  $A: H^{-s}(\Omega) \to H^s(\Omega)$  and  $B: H^s(\Omega) \to H^{-s}(\Omega)$  be given, and consider their corresponding discretizations  $A_{\mathcal{T}}: \mathscr{V}_{\mathcal{T}} \to \mathscr{V}'_{\mathcal{T}}$  and  $B_{\mathcal{T}}: \mathscr{V}_{\mathcal{T}} \to \mathscr{V}'_{\mathcal{T}}.$ 

In this chapter we introduce a uniformly boundedly invertible operator  $D_{\mathcal{T}}: \mathscr{V}_{\mathcal{T}} \to \mathscr{V}'_{\mathcal{T}}$ , allowing us to take  $\mathscr{W}_{\mathcal{T}}$  equal to  $\mathscr{V}_{\mathcal{T}}$  in the operator preconditioning framework (1.6). The resulting preconditioned system  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D'_{\mathcal{T}})^{-1}A_{\mathcal{T}}$  and  $(D'_{\mathcal{T}})^{-1}A_{\mathcal{T}}D_{\mathcal{T}}^{-1}B_{\mathcal{T}}$  are uniformly boundedly invertible. Moreover, the matrix representation of  $D_{\mathcal{T}}$  with respect to the Lagrange basis of  $\mathscr{V}_{\mathcal{T}}$  is diagonal, making the implementation of this preconditioner surprisingly simple.

#### **1.3** About Part II: Parabolic evolution equations

The second topic that will be discussed in this thesis is the (adaptive) numerical solution of parabolic evolution equations written in a simultaneous space-time variational formulation.

As an illustrative example, let us introduce the model problem of the *heat equation* with homogeneous Dirichlet boundary conditions. For some time interval I := (0, T) and some spatial domain  $\Omega \subset \mathbb{R}^d$ , the heat equation reads as finding  $u : I \times \Omega \to \mathbb{R}$  that solves

(1.10) 
$$\begin{aligned} \partial_t u - \Delta_{\boldsymbol{x}} u &= g \quad \text{on } I \times \Omega, \\ u &= 0 \quad \text{on } I \times \partial \Omega, \\ u &= u_0 \quad \text{on } \{0\} \times \Omega, \end{aligned}$$

for some forcing function  $g: I \times \Omega \to \mathbb{R}$  and initial data  $u_0: \Omega \to \mathbb{R}$ .

In order to apply our approximation scheme we first need to derive a weak formulation of the above differential equation. We multiply the first equation by a test function v that vanishes on  $I \times \partial \Omega$ , integrate over  $I \times \Omega$ , and apply integration by parts in space to find

$$(Bu)(v) := \int_{I \times \Omega} (\partial_t u) v + \nabla_{\boldsymbol{x}} u \cdot \nabla_{\boldsymbol{x}} v \, \mathrm{d} \boldsymbol{x} \, \mathrm{d} t = \int_{I \times \Omega} gv \, \mathrm{d} \boldsymbol{x} \, \mathrm{d} t.$$

To enforce the initial condition, we introduce the trace map  $\gamma_0: u \mapsto u(0, \cdot)$ , and test it against some additional test function w to find that

$$\int_{\Omega} (\gamma_0 u) w \, \mathrm{d} \boldsymbol{x} = \int_{\Omega} u_0 w \, \mathrm{d} \boldsymbol{x}$$

9

Clearly, we must yet find suitable function spaces for u, v and w. Define  $X := L_2(I; H_0^1(\Omega)) \cap H^1(I; H^{-1}(\Omega)), Y := L_2(I; H_0^1(\Omega))$ , and  $H := L_2(\Omega)$ , assume that  $g \in Y'$  and  $u_0 \in H$ , and denote  $\mathcal{B} := \begin{bmatrix} B & \gamma_0 \end{bmatrix}^\top$ . Finding  $u \in X$  that solves

(1.11) 
$$\mathcal{B}u = \begin{bmatrix} g \\ u_0 \end{bmatrix}$$

is then a well-posed variational formulation of (1.10), meaning that the operator  $\mathcal{B}$  is a boundedly invertible map  $X \to Y' \times H$  ([SS09]).

A difficulty of the operator equation (1.11) appears when we consider discretizations. As the function space on the trial side does not coincide with that of the test side, we cannot simply apply the Galerkin method like we did in (1.4). Suppose that we have a family of trial spaces  $(X^{\delta})_{\delta \in \Delta} \subset X$ , the question is how to construct a family of test spaces  $(Z^{\delta})_{\delta \in \Delta} \subset Y \times H$  for which the discretized operator  $\mathcal{B}^{\delta} \colon X^{\delta} \to Z^{\delta'}$ , given by  $(\mathcal{B}^{\delta}u)(v) := (\mathcal{B}u)(v)$   $((u,v) \in X^{\delta} \times Z^{\delta})$ , is uniformly boundedly invertible. For the latter to hold, it turns out that the pairs  $(X^{\delta}, Z^{\delta})$  must satisfy dim  $X^{\delta} = \dim Z^{\delta}$  and  $\inf_{\delta \in \Delta} \inf_{0 \neq u \in Z^{\delta}} \frac{(\mathcal{B}u)(v)}{||u||_X ||v||_{Y \times H}} > 0$ ; cf. (1.7)–(1.8). The construction of such a test space  $Z^{\delta}$  is hard, and we do not proceed this way.

Instead, we note that an equivalent problem to (1.11) is to compute

$$u = \underset{w \in X}{\operatorname{argmin}} \|Bw - g\|_{Y'}^2 + \|\gamma_0 w - u_0\|_H^2$$

suggesting an approximation approach by restricting the minimization to some finite dimensional trial space  $X^{\delta} \subset X$ . Unfortunately this is not feasible in practice, because the *Y*'-norm cannot be computed. To resolve this, we consider some finite dimensional test space  $Y^{\delta} \subset Y$  and replace the *Y*'-norm by the (computable)  $Y^{\delta'}$ -norm, yielding the approximation

(1.12) 
$$u^{\delta} = \operatorname*{argmin}_{w \in X^{\delta}} \|Bw - g\|_{Y^{\delta'}}^2 + \|\gamma_0 w - u_0\|_H^2.$$

In [And13] it is shown that, if the family of pairs  $(X^{\delta}, Y^{\delta})_{\delta \in \Delta}$  satisfy (1.13)

$$X^{\delta} \subseteq Y^{\delta} \quad (\delta \in \Delta) \quad \text{and} \quad \gamma_{\Delta} := \inf_{\delta \in \Delta} \inf_{0 \neq v \in X^{\delta}} \sup_{0 \neq v \in Y^{\delta}} \frac{(\partial_t w)(v)}{\|\partial_t w\|_{Y'} \|v\|_Y} > 0,$$

then the approximation  $u^{\delta}$  satisfies  $||u-u^{\delta}||_X \leq \gamma_{\Delta}^{-1} \inf_{w \in X^{\delta}} ||u-w||_X$ , making it is a quasi-best approximation to u from  $X^{\delta}$ . The advantage of this approach is that  $Y^{\delta}$  can be chosen larger in relation to  $X^{\delta}$ , making it easier to satisfy the inf-sup condition.

#### 1.3.1 Contributions and Outline of Part II

The contents of Chapters 6–9 is essentially that of the following papers:

- [SvVW21] R.P. Stevenson, R. van Venetië, and J. Westerdiep. A wavelet-intime, finite element-in-space adaptive method for parabolic evolution equations. Submitted to *Advances in Computational Mathematics*, 2021.
- [vVW21b] R. van Venetië and J. Westerdiep. Efficient space-time adaptivity for parabolic evolution equations using wavelets in time and finite elements in space. Submitted to *Numerical Linear Algebra with Applications*, 2021.
- [vVW21a] R. van Venetië and J. Westerdiep. A parallel algorithm for solving linear parabolic evolution equations. Accepted in 9th Parallel-in-Time Workshop, 2021.
- [GvV21] G. Gantner and R. van Venetië. Adaptive space-time BEM for the heat equation. Submitted to *Computers & Mathematics with Applications*, 2021.

On average, the authors contributed equally to these works.

#### Chapter 6 ([SvVW21])

For trial spaces  $X^{\delta}$  that are *full* (or *sparse*) tensor products of finite element spaces in time and space, in [And13] it was shown how to construct corresponding test spaces  $Y^{\delta} \subset Y$  such that (1.13) holds. Unfortunately, neither family allows for adaptive refinements both locally in time and space.

In this chapter we solve this issue by equipping X with a tensor product basis of a wavelet basis in time and a hierarchical finite element basis in space. We then construct  $X^{\delta}$  as the span of a (finite) subset of this tensor product basis, and construct  $Y^{\delta}$  of a similar type such that (1.13) holds, with the dimension of  $Y^{\delta}$  being proportional to that of  $X^{\delta}$ .

Using properties of the wavelets in time and applying multigrid preconditioners in space, we construct optimal preconditioners  $K_X^{\delta} \colon X^{\delta'} \to X^{\delta}$  and  $K_Y^{\delta} \colon Y^{\delta'} \to Y^{\delta}$ , allowing to solve the discrete problem (1.12) efficiently.

We propose an adaptive algorithm using a standard solve-estimate-markrefine loop. Let  $X^{\delta}$  be the current trial space. In the solve step we find its corresponding approximation  $u^{\delta}$  from (1.12). In the estimate step, we introduce a neighborhood  $X^{\delta} \supset X^{\delta}$  and evaluate the residual on a basis of  $X^{\delta} \setminus X^{\delta}$ . In the mark step, we select the basis functions for which the residual is large. Finally, in the refine step, we build a new trial space  $X^{\delta} \supset X^{\delta}$  containing all the marked basis functions. Under a saturation assumption, we prove that the adaptive loop produces an *r*-linearly converging sequence to the solution.

#### Chapter 7 ([vVW21b])

In this chapter we discuss an implementation of the aforementioned adaptive method in which every step is of linear complexity.

The downside of having bases for  $X^{\delta}$  and  $Y^{\delta}$  of wavelet-type is that the system matrices appearing in the implementation of (1.12) are not sparsely populated. By imposing a *double-tree* constraint on the index sets of  $X^{\delta}$  and  $Y^{\delta}$ , we are able to derive a matrix-free algorithm that can apply the system matrices in linear complexity.

In order to build an actual linear complexity implementation, we based our implementation on tree- and double-tree traversals. We conclude this chapter with extensive results that demonstrate the linear runtime of our code in practice.

#### Chapter 8 ([vVW21a])

One of the advantages of a simultaneous space-time solver over classical timestepping approaches is that they are much better suited for a massively parallel implementation. In this chapter we investigate such a parallel algorithm.

We consider trial spaces  $X^{\delta}$  being the tensor product of finite element spaces in time and space, and show how this tensor product assumption simplifies the implementation of (1.12). After introducing suitable preconditioners, we investigate the parallel complexity of the resulting algorithm. We illustrate our theoretical findings with massively time-parallel computations done in practice.

#### Chapter 9 ([GvV21])

Deviating from the previous setting, here we construct an adaptive space-time *boundary element method* for the heat equation. We consider the heat equation with homogeneous forcing data and prescribed Dirichlet data: for given initial condition  $u_0: \Omega \to \mathbb{R}$  and Dirichlet data  $u_D: I \times \partial\Omega \to \mathbb{R}$ , we seek u that solves

$$\begin{array}{ll} \partial_t u - \Delta_{\boldsymbol{x}} u = 0 & \text{ on } I \times \Omega, \\ u = u_D & \text{ on } I \times \partial \Omega, \\ u = u_0 & \text{ on } \{0\} \times \Omega. \end{array}$$

Since the fundamental solution for the heat equation is known, we can proceed similarly as in Example 1.1.2 to find an equivalent formulation of the problem by solving an integral equation on the lateral space-time boundary  $I \times \partial \Omega$ , see e.g. [AN87, Cos90]. In contrast to (1.11), this formulation is *coercive*.

In this chapter, we propose an a posteriori error estimator for the Galerkin approximation to this integral equation. We show that the estimator is a lower bound for the approximation error, and up to weighted  $L_2$ -terms, also an upper bound. In the numerical results, we let this error estimator drive an adaptive loop that allows for anisotropic refinement. We observe that this loop is able to effectively resolve singularities in time and space, and that it recovers the optimal error decay rate in all of our examples.

## Part I

# **Operator preconditioning**

#### 2.1 Introduction

This chapter is about the construction of preconditioners for discretized boundedly invertible linear operators of negative order using the concept of 'operator preconditioning' ([Hip06]). The idea is to precondition the discretized operator by a discretized operator of opposite order. This is an appealing idea, but it turns out that in order to get a uniformly well-conditioned system, as well as a preconditioner that can be implemented efficiently, the second discretization has to be carefully chosen dependent on the first one.

For a Hilbert space  $\mathscr{H}$ , and a densely embedded reflexive Banach space  $\mathscr{W} \hookrightarrow \mathscr{H}$ , consider the Gelfand triple

$$\mathscr{W} \hookrightarrow \mathscr{H} \simeq \mathscr{H}' \hookrightarrow \mathscr{W}'.$$

For *A* being a boundedly invertible coercive linear operator  $\mathscr{W}' \to \mathscr{W}$ , and  $\mathscr{V}_{\mathcal{T}} \subset \mathscr{H}$  being a finite dimensional subspace of  $\mathscr{W}'$ , let  $(A_{\mathcal{T}}v)(\tilde{v}) := (Av)(\tilde{v})$  $(v, \tilde{v} \in \mathscr{V}_{\mathcal{T}})$ . For *B* being a boundedly invertible coercive linear operator  $\mathscr{W} \to \mathscr{W}'$ , and  $\mathscr{W}_{\mathcal{T}}$  being a finite dimensional subspace of  $\mathscr{W}$ , let  $(B_{\mathcal{T}}w)(\tilde{w}) := (Bw)(\tilde{w})$   $(w, \tilde{w} \in \mathscr{W}_{\mathcal{T}})$ .

A typical example is given by the case that for the boundary  $\Gamma$  of some domain,  $\mathscr{H} = L_2(\Gamma)$ ,  $\mathscr{W} = H^{\frac{1}{2}}(\Gamma)$ , A is the single layer integral operator arising from the Laplacian, B is the corresponding hypersingular integral operator,  $\mathcal{T}$  is a partition from an infinite collection of partitions  $\mathbb{T}$ ,  $\mathscr{V}_{\mathcal{T}}$  is a trial space of discontinuous piecewise polynomials w.r.t.  $\mathcal{T}$ , and  $\mathscr{W}_{\mathcal{T}}$  is a suitable subspace of  $\mathscr{W}$ , which thus cannot be equal to  $\mathscr{V}_{\mathcal{T}}$ . Besides as boundary integral equations, coercive linear operators of order -1 also appear in various domain decomposition type methods in the equations for normal fluxes on interfaces.

Although less frequently, coercive linear operators of order -2 also appear in the literature (e.g. see [FH19]).

In order to precondition  $A_{\mathcal{T}} : \mathscr{V}_{\mathcal{T}} \to \mathscr{V}'_{\mathcal{T}}$  with  $B_{\mathcal{T}} : \mathscr{W}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$  we need to be able to 'identify'  $\mathscr{V}'_{\mathcal{T}}$  with  $\mathscr{W}_{\mathcal{T}}$ , similar to the identification of  $\mathscr{W}''$  with  $\mathscr{W}$ .

Let  $\dim \mathscr{W}_{\mathcal{T}} = \dim \mathscr{V}_{\mathcal{T}}$  and

(2.1) 
$$\inf_{\mathcal{T}\in\mathbb{T}}\inf_{0\neq v\in\mathscr{V}_{\mathcal{T}}}\sup_{0\neq w\in\mathscr{W}_{\mathcal{T}}}\frac{\langle v,w\rangle_{\mathscr{H}}}{\|v\|_{\mathscr{W}'}\|w\|_{\mathscr{W}}}>0.$$

Then  $D_{\mathcal{T}}$  defined by  $(D_{\mathcal{T}}v)(w) := \langle v, w \rangle_{\mathscr{H}}$   $(v \in \mathscr{V}_{\mathcal{T}}, w \in \mathscr{W}_{\mathcal{T}})$  is a uniformly boundedly invertible linear map  $\mathscr{V}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$ , and so its adjoint  $D'_{\mathcal{T}}$  is such a map  $\mathscr{V}_{\mathcal{T}}' \to \mathscr{W}_{\mathcal{T}}$ . We conclude that the preconditioned system  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D'_{\mathcal{T}})^{-1}A_{\mathcal{T}}$  is uniformly boundedly invertible  $\mathscr{V}_{\mathcal{T}} \to \mathscr{V}_{\mathcal{T}}$ .

Equipping  $\mathscr{V}_{\mathcal{T}}$  and  $\mathscr{W}_{\mathcal{T}}$  with bases  $\Xi_{\mathcal{T}}$  and  $\Psi_{\mathcal{T}}$ , respectively, the matrix representation of the preconditioned system reads as  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}D_{\mathcal{T}}^{-T}A_{\mathcal{T}}$ , with 'stiffness matrices'  $A_{\mathcal{T}} := (A_{\mathcal{T}}\Xi_{\mathcal{T}})(\Xi_{\mathcal{T}})$  and  $B_{\mathcal{T}} := (B_{\mathcal{T}}\Psi_{\mathcal{T}})(\Psi_{\mathcal{T}})$ , and 'generalized mass matrix'  $D_{\mathcal{T}} := \langle \Xi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{\mathscr{H}}$ . Regardless of the choice of the bases, the spectral condition number of this matrix is equal to that of  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D'_{\mathcal{T}})^{-1}A_{\mathcal{T}}$ , and thus uniformly bounded.

After an earlier proposal from [Ste02], the currently commonly followed construction of a suitable pair  $(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}})$  is the one from [BC07]. Here  $\mathscr{V}_{\mathcal{T}}$  is the space of piecewise constants w.r.t. a partition  $\mathcal{T}$  of a two-dimensional domain or manifold equipped with the usual basis  $\Xi_{\mathcal{T}}$ , and  $\mathscr{W}_{\mathcal{T}}$ , defined as the span of a collection  $\Psi_{\mathcal{T}}$ , is a subspace of the space of continuous piecewise linears w.r.t. a barycentric refinement of  $\mathcal{T}$  constructed by subdividing each triangle into 6 subtriangles by connecting its vertices and midpoints with its barycenter. In [HUT16] the inf-inf-sup condition (2.1) was demonstrated for families of partitions including locally refined ones that satisfy a certain mildly-grading condition from [Ste03a].

A problem with the constructions from both [Ste02, BC07] is that the matrix  $D_{\mathcal{T}}$  is *not* diagonal, so that its inverse has to be approximated. Knowing that  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}D_{\mathcal{T}}^{-T}$  is *not* well-conditioned, because  $A_{\mathcal{T}}$  is not whereas their product is uniformly well-conditioned, the accuracy with which  $D_{\mathcal{T}}^{-1}$  has to be approximated such that it gives rise to a uniform preconditioner increases with an increasing (minimal) mesh-size.

#### 2.1.1 Contributions

For the aforementioned  $\mathscr{V}_{\mathcal{T}}$  and  $\Xi_{\mathcal{T}}$ , in this chapter a space  $\mathscr{W}_{\mathcal{T}}$ , given as the span of a collection  $\Psi_{\mathcal{T}}$ , will be constructed such that (2.1) is valid, and  $D_{\mathcal{T}} = \langle \Xi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{\mathscr{H}}$  is *diagonal*. i.e.,  $\Xi_{\mathcal{T}}$  and  $\Psi_{\mathcal{T}}$  are *biorthogonal*. Thanks to both  $\Xi_{\mathcal{T}}$  and  $\Psi_{\mathcal{T}}$  being 'locally supported', the corresponding biorthogonal projector onto  $\mathscr{W}_{\mathcal{T}}$  is *local*, which allows to demonstrate the inf-inf-sup stability *without any mildly grading assumption* on the partitions.

Each function in  $\Psi_{\mathcal{T}}$  equals a function from the space  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  of continuous piecewise linears<sup>1</sup> w.r.t.  $\mathcal{T}$ , plus a linear combination of 'bubble func-

<sup>&</sup>lt;sup>1</sup>The subscript 0 in the notation  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  refers to possible boundary conditions that are incorporated.

tions' from a space denoted as  $\mathscr{B}_{\mathcal{T}}$ . Since the decomposition of  $\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}}$ into  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  and  $\mathscr{B}_{\mathcal{T}}$  is stable w.r.t. the  $\mathscr{W}$ -norm, instead of simply defining  $(B_{\mathcal{T}}w)(\tilde{w}) := (Bw)(\tilde{w}) \ (w, \tilde{w} \in \mathscr{W}_{\mathcal{T}})$ , a suitable boundedly invertible linear operator  $B_{\mathcal{T}} : \mathscr{W}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$  will be constructed from an invertible diagonal scaling on the bubble space and a boundedly invertible linear operator  $B_{\mathcal{T}}^{\mathscr{G}} : \mathscr{S}_{\mathcal{T},0}^{0,1} \to (\mathscr{S}_{\mathcal{T},0}^{0,1})'$ , e.g.  $(B_{\mathcal{T}}^{\mathscr{G}}w)(\tilde{w}) := (Bw)(\tilde{w})$  with B the hypersingular operator. Other than in [Ste02, BC07], by this use of the stable decomposition there is no need to discretize the hypersingular operator on a refinement of  $\mathcal{T}$ . The whole approach relies on *existence* of bubble functions with certain properties, whereas these functions themselves do not enter the implementation.

The total cost of the resulting preconditioner is the sum of the cost of the application of  $B_{\mathcal{T}}^{\mathscr{S}}$  plus cost that scales linearly in  $\#\mathcal{T}$ . For  $\mathcal{T}$  being a uniform refinement of some initial coarse partition, a  $B_{\mathcal{T}}^{\mathscr{S}}$  of multi-level type can be found whose cost scales linearly in  $\#\mathcal{T}$  ([BPV00]). Such  $B_{\mathcal{T}}^{\mathscr{S}}$  that also apply on locally refined partitions will be discussed in Chapter 3.

The construction of the biorthogonal collection  $\Psi_{\mathcal{T}}$ , and with that of the preconditioner, is based on a general principle. It applies in any space dimension, and, as we will see, it applies equally well when  $\mathscr{V}_{\mathcal{T}}$  is the space of *continuous piecewise linears*. Higher order discretizations will be covered as well.

The construction applies equally well on manifolds. The coefficients of the functions from  $\Psi_{\mathcal{T}}$  in terms of functions from  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  and the bubble functions are given as inner products between functions of  $\mathscr{V}_{\mathcal{T}}$  and  $\mathscr{S}_{\mathcal{T},0}^{0,1}$ . Since in the manifold case, however, generally these inner products cannot be evaluated exactly, we present an alternative slightly modified construction in which the true  $L_2$ -inner product is replaced by a mesh-dependent one by an element-wise freezing of the Jacobian. It still yields a uniform preconditioner on general, possibly locally refined partitions, while the same explicit formula for the expansion coefficients of the functions of  $\Psi_{\mathcal{T}}$  that was derived in the domain case, now also applies in the manifold case.

#### 2.1.2 Notations

In this work, by  $\lambda \leq \mu$  we will mean that  $\lambda$  can be bounded by a multiple of  $\mu$ , independently of parameters which  $\lambda$  and  $\mu$  may depend on, with the sole exception of the space dimension d, or in the manifold case, on the parametrization of the manifold that is used to define the finite element spaces on it. Obviously,  $\lambda \geq \mu$  is defined as  $\mu \leq \lambda$ , and  $\lambda \equiv \mu$  as  $\lambda \leq \mu$  and  $\lambda \geq \mu$ .

For normed linear spaces  $\mathscr{Y}$  and  $\mathscr{Z}$ , in this work for convenience over  $\mathbb{R}$ ,  $\mathcal{L}(\mathscr{Y}, \mathscr{Z})$  will denote the space of bounded linear mappings  $\mathscr{Y} \to \mathscr{Z}$  endowed with the operator norm  $\|\cdot\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Z})}$ . The subset of invertible operators in  $\mathcal{L}(\mathscr{Y}, \mathscr{Z})$  with inverses in  $\mathcal{L}(\mathscr{Z}, \mathscr{Y})$  will be denoted as  $\mathcal{L}is(\mathscr{Y}, \mathscr{Z})$ . The *condition number* of a  $C \in \mathcal{L}is(\mathscr{Y}, \mathscr{Z})$  is defined as  $\kappa_{\mathscr{Y}, \mathscr{Z}}(C) := \|C\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Z})} \|C^{-1}\|_{\mathcal{L}(\mathscr{Z}, \mathscr{Y})}$ . For  $\mathscr{Y}$  a reflexive Banach space and  $C \in \mathcal{L}(\mathscr{Y}, \mathscr{Y}')$  being *coercive*, i.e.,

$$\inf_{0 \neq y \in \mathscr{Y}} \frac{(Cy)(y)}{\|y\|_{\mathscr{Y}}^2} > 0,$$

both *C* and  $\Re(C) := \frac{1}{2}(C + C')$  are in  $\operatorname{Lis}(\mathscr{Y}, \mathscr{Y}')$  with

$$\begin{aligned} \|\Re(C)\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')} &\leq \|C\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')}, \\ \|C^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} &\leq \|\Re(C)^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} = \Big(\inf_{0\neq y\in\mathscr{Y}} \frac{(Cy)(y)}{\|y\|_{\mathscr{Y}}^2}\Big)^{-1}. \end{aligned}$$

The subset of coercive operators in  $\mathcal{L}is(\mathscr{Y}, \mathscr{Y}')$  is denoted as  $\mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$ . If  $C \in \mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$ , then  $C^{-1} \in \mathcal{L}is_c(\mathscr{Y}', \mathscr{Y})$  and  $\|\Re(C^{-1})^{-1}\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Y}')} \leq \|C\|^2_{\mathcal{L}(\mathscr{Y}, \mathscr{Y})} \|\Re(C)^{-1}\|_{\mathcal{L}(\mathscr{Y}', \mathscr{Y})}$ .

Two countable collections  $\Upsilon = (v_i)_i$  and  $\tilde{\Upsilon} = (\tilde{v}_i)_i$  in a Hilbert space will be called *biorthogonal* when  $\langle \Upsilon, \tilde{\Upsilon} \rangle = [\langle v_j, \tilde{v}_i \rangle]_{ij}$  is an *invertible diagonal* matrix, and *biorthonormal* when it is the *identity* matrix.

#### 2.1.3 Organization

In Sect. 2.2 the general principles of operator preconditioning are recalled. In Sect. 2.3, it is applied to operators of negative order discretized with *discontinuous* piecewise constants, first in the domain- and then in the manifold-case. In Sect. 2.4, the same program is followed for trial spaces of *continuous* piecewise linears. In Sect. 2.5 the results from Sect. 2.3-2.4 will be extended to higher order finite element spaces. This will be done by both applying the operator preconditioning framework directly to the higher order spaces, and by using the preconditioner found for the lowest order case in a subspace correction approach. Finally, in Sect. 2.6 we report on some numerical results obtained with the new preconditioners, and compare them with those obtained with the preconditioner from [BC07, HUT16].

#### 2.2 Operator preconditioning

The exposition in this section largely follows [Hip06, Sect. 2] closely. Let  $\mathscr{V}, \mathscr{W}$  be reflexive Banach spaces. We will search a 'preconditioner' *G* for an  $A \in \mathcal{L}is(\mathscr{V}, \mathscr{V}')$ , i.e. an operator  $G \in \mathcal{L}is(\mathscr{V}', \mathscr{V})$  (whose application is 'easy' compared to that of  $A^{-1}$ ). It is often useful, e.g. for the application of Conjugate Gradients, when the preconditioner is coercive, i.e., being an operator in  $\mathcal{L}is_c(\mathscr{V}', \mathscr{V})$ . The following result is easily verified.

**Proposition 2.2.1.** If  $B \in \mathcal{L}$ is $(\mathcal{W}, \mathcal{W}')$  and  $D \in \mathcal{L}$ is $(\mathcal{V}, \mathcal{W}')$ , then

$$G := D^{-1}B(D')^{-1} \in \mathcal{L}is(\mathscr{V}', \mathscr{V}),$$

and

$$\|G\|_{\mathcal{L}(\mathcal{V}',\mathcal{V})} \le \|D^{-1}\|_{\mathcal{L}(\mathcal{W}',\mathcal{V})}^{2}\|B\|_{\mathcal{L}(\mathcal{W},\mathcal{W}')},$$
  
$$\|G^{-1}\|_{\mathcal{L}(\mathcal{V},\mathcal{V}')} \le \|D\|_{\mathcal{L}(\mathcal{V},\mathcal{W}')}^{2}\|B^{-1}\|_{\mathcal{L}(\mathcal{W}',\mathcal{W})}.$$

If additionally  $B \in \mathcal{L}is_c(\mathcal{W}, \mathcal{W}')$ , then  $G \in \mathcal{L}is_c(\mathcal{V}', \mathcal{V})$ , and

$$\|\Re(G)^{-1}\|_{\mathcal{L}(\mathscr{V},\mathscr{V}')} \le \|D\|_{\mathcal{L}(\mathscr{V},\mathscr{W}')}^2 \|\Re(B)^{-1}\|_{\mathcal{L}(\mathscr{W}',\mathscr{W})}.$$

*Remark* 2.2.2. We recall that by an application of the *closed range theorem*,  $D \in \mathcal{L}(\mathcal{V}, \mathcal{W}')$  is in  $\mathcal{L}is(\mathcal{V}, \mathcal{W}')$  if and only if for all  $w \in \mathcal{W}$  there exists a  $v \in \mathcal{V}$  with  $(Dv)(w) \neq 0$ , and

$$0 < \inf_{0 \neq v \in \mathscr{V}} \sup_{0 \neq w \in \mathscr{W}} \frac{(Dv)(w)}{\|v\|_{\mathscr{V}} \|w\|_{\mathscr{W}}} \quad \left( = \|D^{-1}\|_{\mathcal{L}(\mathscr{W}',\mathscr{V})}^{-1} \right)$$

In particular we are interested in finding a preconditioner for an operator  $A_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}}')$  of the form  $G_{\mathcal{T}} = D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D'_{\mathcal{T}})^{-1}$ , where  $\mathscr{V}_{\mathcal{T}}$  is some *finite dimensional* (finite- or boundary element) space. In view of Proposition 2.2.1, for that goal we search some finite dimensional space  $\mathscr{W}_{\mathcal{T}}$  with

(2.2) 
$$\dim \mathscr{W}_{\mathcal{T}} = \dim \mathscr{V}_{\mathcal{T}},$$

and operators  $B_{\mathcal{T}} \in \mathcal{L}is(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  and  $D_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$ .

A typical setting is that, for some reflexive Banach spaces  $\mathscr{V}$  and  $\mathscr{W}$ , and operators  $A \in \operatorname{Lis}_c(\mathscr{V}, \mathscr{V}')$  and  $B \in \operatorname{Lis}_c(\mathscr{W}, \mathscr{W}')$ , we have  $\mathscr{V}_{\mathcal{T}} \subset \mathscr{V}$  (thus equipped with  $\| \|_{\mathscr{V}}$ ),  $(A_{\mathcal{T}}u)(v) := (Au)(v)$   $(u, v \in \mathscr{V}_{\mathcal{T}})$  and, for a suitable  $\mathscr{W}_{\mathcal{T}} \subset$  $\mathscr{W}$  (thus equipped with  $\| \|_{\mathscr{W}}$ ), take  $(B_{\mathcal{T}}w)(z) := (Bw)(z)$   $(w, z \in \mathscr{W}_{\mathcal{T}})$ . In this case  $A_{\mathcal{T}} \in \operatorname{Lis}_c(\mathscr{V}_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}}')$  and  $B_{\mathcal{T}} \in \operatorname{Lis}_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  with

$$\|A_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}}')} \le \|A\|_{\mathcal{L}(\mathscr{V},\mathscr{V}')}, \|\Re(A_{\mathcal{T}})^{-1}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}',\mathscr{V}_{\mathcal{T}})} \le \|\Re(A)^{-1}\|_{\mathcal{L}(\mathscr{V}',\mathscr{V})}, \\ \|B_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}}')} \le \|B\|_{\mathcal{L}(\mathscr{W},\mathscr{W}')}, \|\Re(B_{\mathcal{T}})^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}',\mathscr{W}_{\mathcal{T}})} \le \|\Re(B)^{-1}\|_{\mathcal{L}(\mathscr{W}',\mathscr{W})}.$$

An obvious construction of a suitable  $D_{\mathcal{T}}$  is discussed in the next proposition.

**Proposition 2.2.3** (Fortin projector ([For77])). For some  $D \in \mathcal{L}is(\mathcal{V}, \mathcal{W}')$ , let  $D_{\mathcal{T}} \in \mathcal{L}(\mathcal{V}_{\mathcal{T}}, \mathcal{W}'_{\mathcal{T}})$  be defined by  $(D_{\mathcal{T}}v)(w) := (Dv)(w)$ . Then

$$\|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}}')} \le \|D\|_{\mathcal{L}(\mathscr{V},\mathscr{W}')}.$$

Assuming (2.2), additionally one has  $D_{\mathcal{T}} \in \mathcal{L}is(\mathcal{V}_{\mathcal{T}}, \mathcal{W}_{\mathcal{T}}')$  if, and for  $\mathcal{W}$  being a Hilbert space, only if there exists a projector  $P_{\mathcal{T}} \in \mathcal{L}(\mathcal{W}, \mathcal{W})$  onto  $\mathcal{W}_{\mathcal{T}}$  with  $(D\mathcal{V}_{\mathcal{T}})((\mathrm{Id} - P_{\mathcal{T}})\mathcal{W}) = 0$ , in which case

(2.3) 
$$\|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}',\mathscr{V}_{\mathcal{T}})} \leq \|P_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W},\mathscr{W})}\|D^{-1}\|_{\mathcal{L}(\mathscr{W}',\mathscr{V})}.$$

19

*Proof.* The first statement is obvious. Now let us assume existence of a (Fortin) projector  $P_{\mathcal{T}}$ . Then for  $v_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}$ ,

$$\begin{split} \|D^{-1}\|_{\mathcal{L}(\mathcal{W}',\mathcal{V})}^{-1}\|v_{\mathcal{T}}\|_{\mathcal{V}} &\leq \sup_{0 \neq w \in \mathscr{W}} \frac{(Dv_{\mathcal{T}})(w)}{\|w\|_{\mathscr{W}}} = \sup_{0 \neq w \in \mathscr{W}} \frac{(Dv_{\mathcal{T}})(P_{\mathcal{T}}w)}{\|w\|_{\mathscr{W}}} \\ &\leq \|P_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W},\mathscr{W})} \sup_{0 \neq w_{\mathcal{T}} \in \mathscr{W}_{\mathcal{T}}} \frac{(Dv_{\mathcal{T}})(w_{\mathcal{T}})}{\|w_{\mathcal{T}}\|_{\mathscr{W}}}, \end{split}$$

which together with Remark 2.2.2 and (2.2) shows that  $D_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$ , in particular (2.3).

Conversely (cf. [Bra01, Remark 4.9]), assume  $D_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$ , and let  $\mathscr{W}$  be a Hilbert space. Then given  $w \in \mathscr{W}$ , let  $w_{\mathcal{T}}$  be the first component of the solution  $(w_{\mathcal{T}}, v_{\mathcal{T}}) \in \mathscr{W}_{\mathcal{T}} \times \mathscr{V}_{\mathcal{T}}$  of the well-posed saddle point problem

Then  $P_{\mathcal{T}} := w \mapsto w_{\mathcal{T}}$  is a valid Fortin projector.

In applications, one usually has a *family* of spaces  $\mathscr{V}_{\mathcal{T}}$  and aims at a *uni-form* preconditioner  $G_{\mathcal{T}}$ . In the setting of Proposition 2.2.3 it means that one searches a Fortin projector  $P_{\mathcal{T}}$  such that  $\|P_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W},\mathscr{W})}$  is *uniformly* bounded.

#### 2.2.1 Implementation

Given a finite collection  $\Upsilon = \{v\}$  in a linear space, we set the *synthesis operator* 

$$\mathcal{F}_{\Upsilon}: \mathbb{R}^{\#\Upsilon} \to \operatorname{span} \Upsilon: \mathbf{c} \mapsto \mathbf{c}^{\top} \Upsilon := \sum_{\upsilon \in \Upsilon} c_{\upsilon} \upsilon.$$

Equipping  $\mathbb{R}^{\#\Upsilon}$  with the Euclidean scalar product  $\langle , \rangle$ , and identifying  $(\mathbb{R}^{\#\Upsilon})'$  with  $\mathbb{R}^{\#\Upsilon}$  using the corresponding Riesz map, we infer that the adjoint of  $\mathcal{F}_{\Upsilon}$ , known as the *analysis operator*, satisfies

$$\mathcal{F}'_{\Upsilon} : (\operatorname{span} \Upsilon)' \to \mathbb{R}^{\#\Upsilon} : f \mapsto f(\Upsilon) := [f(v)]_{v \in \Upsilon}.$$

A collection  $\Upsilon$  is a *basis* for its span when  $\mathcal{F}_{\Upsilon} \in \mathcal{L}is(\mathbb{R}^{\#\Upsilon}, \operatorname{span} \Upsilon)$  (and so  $\mathcal{F}'_{\Upsilon} \in \mathcal{L}is((\operatorname{span} \Upsilon)', \mathbb{R}^{\#\Upsilon})$ .)

Now let  $\Xi_{\mathcal{T}} = \{\xi\}$  and  $\Psi_{\mathcal{T}} = \{\psi\}$  be bases for  $\mathscr{V}_{\mathcal{T}}$  and  $\mathscr{W}_{\mathcal{T}}$ , respectively. Then in coordinates the preconditioned system reads as

$$\mathcal{F}_{\Xi_{\mathcal{T}}}^{-1}G_{\mathcal{T}}A_{\mathcal{T}}\mathcal{F}_{\Xi_{\mathcal{T}}} = G_{\mathcal{T}}A_{\mathcal{T}} := D_{\mathcal{T}}^{-1}B_{\mathcal{T}}D_{\mathcal{T}}^{-\top}A_{\mathcal{T}},$$

where

$$\boldsymbol{A}_{\mathcal{T}} := \mathcal{F}'_{\Xi_{\mathcal{T}}} A_{\mathcal{T}} \mathcal{F}_{\Xi_{\mathcal{T}}}, \quad \boldsymbol{B}_{\mathcal{T}} := \mathcal{F}'_{\Psi_{\mathcal{T}}} B_{\mathcal{T}} \mathcal{F}_{\Psi_{\mathcal{T}}}, \quad \boldsymbol{D}_{\mathcal{T}} := \mathcal{F}'_{\Psi_{\mathcal{T}}} D_{\mathcal{T}} \mathcal{F}_{\Xi_{\mathcal{T}}}.$$

20

By identifying a map in  $\mathcal{L}(\mathbb{R}^{\#\Xi_{\mathcal{T}}}, \mathbb{R}^{\#\Xi_{\mathcal{T}}})$  with a  $\#\Xi_{\mathcal{T}} \times \#\Xi_{\mathcal{T}}$  matrix by equipping  $\mathbb{R}^{\#\Xi_{\mathcal{T}}}$  with the canonical basis  $\{e_i\}$ , and by enumerating the elements of  $\Xi_{\mathcal{T}}$  one has

$$(\boldsymbol{A}_{\mathcal{T}})_{ij} = \langle \mathcal{F}'_{\Xi_{\mathcal{T}}} A_{\mathcal{T}} \mathcal{F}_{\Xi_{\mathcal{T}}} \boldsymbol{e}_j, \boldsymbol{e}_i \rangle = (A_{\mathcal{T}} \mathcal{F}_{\Xi_{\mathcal{T}}} \boldsymbol{e}_j) (\mathcal{F}_{\Xi_{\mathcal{T}}} \boldsymbol{e}_i) = (A_{\mathcal{T}} \xi_j) (\xi_i),$$

and similarly,

$$(\boldsymbol{B}_{\mathcal{T}})_{ij} = (B_{\mathcal{T}}\psi_j)(\psi_i), \quad (\boldsymbol{D}_{\mathcal{T}})_{ij} = (D_{\mathcal{T}}\xi_j)(\psi_i),$$

Preferably  $D_{\mathcal{T}}$  is such that its inverse can be applied in linear complexity, as is the case when  $D_{\mathcal{T}}$  is *diagonal*.

*Remark* 2.2.4. Using  $\sigma()$  and  $\rho()$  to denote the spectrum and spectral radius of an operator, clearly  $\sigma(\mathbf{G}_{\mathcal{T}}\mathbf{A}_{\mathcal{T}}) = \sigma(G_{\mathcal{T}}A_{\mathcal{T}})$ . So for the spectral condition number we have

$$\kappa_S(\boldsymbol{G}_{\mathcal{T}}\boldsymbol{A}_{\mathcal{T}}) := \rho(\boldsymbol{G}_{\mathcal{T}}\boldsymbol{A}_{\mathcal{T}})\rho((\boldsymbol{G}_{\mathcal{T}}\boldsymbol{A}_{\mathcal{T}})^{-1}) \le \kappa_{\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}}}(G_{\mathcal{T}}A_{\mathcal{T}}),$$

which thus holds true *independently* of the choice of the basis  $\Xi_{\mathcal{T}}$  for  $\mathscr{V}_{\mathcal{T}}$ . Furthermore, in view of an application of Conjugate Gradients, if  $A_{\mathcal{T}}$  and  $B_{\mathcal{T}}$  are coercive and *self-adjoint*, then  $A_{\mathcal{T}}$  and  $G_{\mathcal{T}}$  are positive definite and symmetric. Equipping  $\mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}}$  with  $\|\!|\cdot|\!| := \|(G_{\mathcal{T}})^{-\frac{1}{2}} \cdot \|$  or  $\|\!|\cdot|\!| := \|(A_{\mathcal{T}})^{\frac{1}{2}} \cdot \|$ , in that case we have

$$\kappa_{(\mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}}, \|\cdot\|), (\mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}}, \|\cdot\|)}(G_{\mathcal{T}}A_{\mathcal{T}}) = \kappa_{S}(G_{\mathcal{T}}A_{\mathcal{T}}).$$

#### 2.3 Piecewise constant discretization space

For a bounded polytopal domain  $\Omega \subset \mathbb{R}^d$ , a measurable, closed, possibly empty  $\gamma \subset \partial \Omega$ , and an  $s \in [0, 1]$ , we take

$$\mathscr{W} := [L_2(\Omega), H^1_{0,\gamma}(\Omega)]_{s,2}, \quad \mathscr{V} := \mathscr{W}',$$

where  $H^1_{0,\gamma}(\Omega)$  is the closure in  $H^1(\Omega)$  of the  $C^{\infty}(\Omega) \cap H^1(\Omega)$  functions that vanish at  $\gamma$ .<sup>2</sup> The role of  $D \in \mathcal{L}is(\mathcal{V}, \mathcal{W}')$  in Proposition 2.2.3 is going to be played by the unique extension to  $\mathcal{V} \times \mathcal{W}$  of the duality pairing

$$(Dv)(w) := \langle v, w \rangle_{L_2(\Omega)},$$

which satisfies  $||D||_{\mathcal{L}(\mathscr{V},\mathscr{W}')} = ||D^{-1}||_{\mathcal{L}(\mathscr{W}',\mathscr{V})} = 1.$ 

Let  $(\mathcal{T})_{\mathcal{T}\in\mathbb{T}}$  be a family of *conforming* partitions of  $\Omega$  into (open) *uniformly* shape regular *d*-simplices, where we assume that  $\gamma$  is the (possibly empty) union of (d-1)-faces of  $T \in \mathcal{T}$ . Thanks to the conformity and the uniform shape

<sup>&</sup>lt;sup>2</sup>In the domain case, it is easy to generalize the results to Sobolev spaces with smoothness index  $s \in [0, \frac{3}{2})$ , or even to  $s \in (-\frac{1}{2}, \frac{3}{2})$ .

regularity, for d > 1 we know that neighbouring  $T, T' \in \mathcal{T}$ , i.e.  $\overline{T} \cap \overline{T'} \neq \emptyset$ , have uniformly comparable sizes. For d = 1, we impose this uniform '*K*-mesh property' explicitly.<sup>3</sup>

For  $\mathcal{T} \in \mathbb{T}$ , we define  $N_{\mathcal{T}}^0$  as the set of vertices of  $\mathcal{T}$  that are not on  $\gamma$ , and for  $\nu \in N_{\mathcal{T}}^0$  we set its *valence* 

$$d_{\mathcal{T},\nu} := \#\{T \in \mathcal{T} \colon \nu \in \overline{T}\}.$$

For  $T \in \mathcal{T}$ , and with  $N_T$  denoting the set of its vertices, we set  $N^0_{\mathcal{T},T} := N^0_{\mathcal{T}} \cap N_T$ , and define  $h_T := |T|^{1/d}$ .

We take

$$\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,0} := \{ u \in L_2(\Omega) \colon u|_T \in \mathcal{P}_0 \ (T \in \mathcal{T}) \} \subset \mathscr{V},$$

and, as a first ingredient in the forthcoming construction of a suitable  $\mathscr{W}_{\mathcal{T}}$ , define the space of continuous piecewise linears, zero on  $\gamma$ , by

$$\mathscr{S}_{\mathcal{T},0}^{0,1} := \{ u \in H^1_{0,\gamma}(\Omega) \colon u|_T \in \mathcal{P}_1 \ (T \in \mathcal{T}) \},\$$

equipped with the usual bases

(2.4) 
$$\Xi_{\mathcal{T}} = \{\xi_T \colon T \in \mathcal{T}\}, \quad \Phi_{\mathcal{T}} = \{\phi_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}^0\},$$

respectively, defined by

(2.5) 
$$\xi_T := \begin{cases} 1 & \text{on } T, \\ 0 & \text{on } \Omega \setminus T, \end{cases} \quad \phi_{\mathcal{T},\nu}(\nu') = \delta_{\nu,\nu'} \quad (\nu,\nu' \in N^0_{\mathcal{T}}).$$

#### **2.3.1** Construction of $\mathscr{W}_{\mathcal{T}}$ and $D_{\mathcal{T}}$ .

Aiming at the construction of a (uniform) preconditioner  $G_{\mathcal{T}} \in \mathcal{L}is_c(\mathcal{V}'_{\mathcal{T}}, \mathcal{V}_{\mathcal{T}})$ using the framework of operator preconditioning, we are going to construct a collection  $\Psi_{\mathcal{T}} \subset H^1_{0,\gamma}(\Omega)$  that is biorthogonal to  $\Xi_{\mathcal{T}}$ , whose elements are 'locally supported', and for which

$$\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}} \subset \mathscr{W}$$

has an 'approximation property'. These three properties of  $\Psi_{\mathcal{T}}$  will allow us to construct a suitable Fortin projector, and they will give rise to a matrix  $D_{\mathcal{T}}$  that is *diagonal*.

The construction of  $\Psi_{\mathcal{T}}$  builds on two collections  $\Theta_{\mathcal{T}}$  and  $\Sigma_{\mathcal{T}}$  of 'locally supported' functions in  $H^1_{0,\gamma}(\Omega)$  whose cardinalities are equal to that of  $\Xi_{\mathcal{T}}$ ,

<sup>&</sup>lt;sup>3</sup>For our convenience, throughout this chapter we consider trial spaces w.r.t. conforming partitions into uniformly shape regular *d*-simplices. It will however become clear that families of non-conforming partitions into uniformly shape regular *d*-simplices or hyperrectangles that satisfy a uniform *K*-mesh property can be dealt with as well.

the first being biorthogonal to  $\Xi_{\mathcal{T}}$ , and the second whose span has an 'approximation property' and is inside  $\mathscr{S}^{0,1}_{\mathcal{T},0}$ .

Let  $\Theta_{\mathcal{T}} = \{\theta_T \colon T \in \mathcal{T}\} \subset H^1_{0,\gamma}(\Omega)$  be such that  $\theta_T \ge 0$ ,  $\operatorname{supp} \theta_T \subset \overline{T}$ ,

(2.6) 
$$\langle \theta_T, \xi_{T'} \rangle_{L_2(\Omega)} \approx \delta_{TT'} \| \theta_T \|_{L_2(\Omega)} \| \xi_{T'} \|_{L_2(\Omega)}, \quad (T, T' \in \mathcal{T}),$$

and, for convenience only, that is scaled such that

(2.7) 
$$\langle \theta_T, \xi_T \rangle_{L_2(\Omega)} = |T|.$$

One obvious possible construction of such  $\Theta_{\mathcal{T}}$  is to take  $\theta_T$  to be the 'bubble function' defined by  $\theta_T(x) = \frac{(2d+1)!}{d} \prod_{i=1}^{d+1} \lambda_i(x)$  for  $x \in T$ , and zero elsewhere, where  $(\lambda_1(x), \dots, \lambda_{d+1}(x))$  are the barycentric coordinates of x w.r.t. T (see e.g. [VS18] for (2.7)). Two forthcoming (harmless) conditions (2.26) and (2.27) on  $\Theta_{\mathcal{T}}$  will be satisfied as well by the above specification of  $\theta_T$ .

Another, equally suited construction is, after making a uniform barycentric refinement of T, to take  $\theta_T$  as a the continuous piecewise linear hat function associated to the barycenter of T, multiplied by a factor d + 1.

We emphasize that the resulting preconditioner will *not* depend on the actual construction of  $\Theta_{T}$ , but that only *existence* of a collection with the aforementioned properties is relevant.

Defining  $\Sigma_{\mathcal{T}} = \{\sigma_{\mathcal{T},T} \colon T \in \mathcal{T}\} \subset \mathscr{S}^{0,1}_{\mathcal{T},0}$  by

$$\sigma_{\mathcal{T},T} := \sum_{\nu \in N^0_{\mathcal{T},T}} d_{\mathcal{T},\nu}^{-1} \phi_{\mathcal{T},\nu},$$

we have

$$\sum_{T\in\mathcal{T}}\sigma_{\mathcal{T},T}=\sum_{\nu\in N^0_{\mathcal{T}}}\phi_{\mathcal{T},\nu},$$

being equal to the constant function 1 on  $\Omega \setminus \bigcup_{\{T \in \mathcal{T}: \overline{T} \cap \gamma \neq \emptyset\}} \overline{T}$ , which yields the aforementioned 'approximation property' (cf. footnote 4).

We now define

$$\Psi_{\mathcal{T}} := \{\psi_{\mathcal{T},T} \colon T \in \mathcal{T}\} \subset \mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \operatorname{span} \Theta_{\mathcal{T}},$$

by

$$(2.8) \quad \psi_{\mathcal{T},T} \coloneqq \sigma_{\mathcal{T},T} + \frac{\langle \mathbb{1} - \sigma_{\mathcal{T},T}, \xi_T \rangle_{L_2(\Omega)}}{\langle \theta_T, \xi_T \rangle_{L_2(\Omega)}} \theta_T - \sum_{T' \in \mathcal{T} \setminus \{T\}} \frac{\langle \sigma_{\mathcal{T},T}, \xi_{T'} \rangle_{L_2(\Omega)}}{\langle \theta_{T'}, \xi_{T'} \rangle_{L_2(\Omega)}} \theta_{T'},$$

The third term at the right-hand side corrects  $\sigma_{\mathcal{T},T}$  such that it becomes orthogonal to  $\xi_{T'}$  for  $T' \neq T$ , whereas the second term ensures that the  $\psi_{\mathcal{T},T}$  sum up to 1, possibly except on a strip along the Dirichlet boundary:

Lemma 2.3.1. It holds that

(2.9) 
$$\sum_{T\in\mathcal{T}}\psi_{\mathcal{T},T} = \sum_{T\in\mathcal{T}}\sigma_{\mathcal{T},T} + \sum_{T\in\mathcal{T}}\frac{\langle\mathbb{1}-\sum_{T'\in\mathcal{T}}\sigma_{\mathcal{T},T'},\xi_T\rangle_{L_2(\Omega)}}{\langle\theta_T,\xi_T\rangle_{L_2(\Omega)}}\theta_T,$$

and

(2.10) 
$$\langle \Xi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{L_2(\Omega)} = \operatorname{diag}\{\langle \mathbb{1}, \xi_T \rangle_{L_2(\Omega)} \colon T \in \mathcal{T}\}.$$

*Proof.* Writing  $1 - \sigma_{\mathcal{T},T} = \sum_{T' \in \mathcal{T} \setminus \{T\}} \sigma_{\mathcal{T},T'} + (1 - \sum_{T' \in \mathcal{T}} \sigma_{\mathcal{T},T'})$ , (2.9) follows from (2.8) by using that

$$\sum_{T \in \mathcal{T}} \sum_{T' \in \mathcal{T} \setminus \{T\}} \frac{\langle \sigma_{\mathcal{T}, T'}, \xi_T \rangle_{L_2(\Omega)}}{\langle \theta_T, \xi_T \rangle_{L_2(\Omega)}} \theta_T - \sum_{T \in \mathcal{T}} \sum_{T' \in \mathcal{T} \setminus \{T\}} \frac{\langle \sigma_{\mathcal{T}, T}, \xi_{T'} \rangle_{L_2(\Omega)}}{\langle \theta_{T'}, \xi_{T'} \rangle_{L_2(\Omega)}} \theta_{T'} = 0.$$

The biorthonormality of  $\Xi_{\mathcal{T}}$  and  $\{\theta_T / \langle \theta_T, \xi_T \rangle_{L_2(\Omega)} \colon T \in \mathcal{T}\}$  shows (2.10).  $\Box$ 

By expanding  $\sigma_{\mathcal{T},T}$  in terms of the nodal basis, and by using  $\int_T \phi_{\mathcal{T},\nu} dx = \frac{|T|}{d+1}$  and the normalization (2.7), we arrive at the explicit expression (2.11)

$$\psi_{\mathcal{T},T} := \sum_{\nu \in N^0_{\mathcal{T},T}} d_{\mathcal{T},\nu}^{-1} \phi_{\mathcal{T},\nu} + \left(1 - \frac{1}{d+1} \sum_{\nu \in N^0_{\mathcal{T},T}} d_{\mathcal{T},\nu}^{-1}\right) \theta_T - \sum_{T' \in \mathcal{T} \setminus \{T\}} \left(\frac{1}{d+1} \sum_{\nu \in N^0_{\mathcal{T},T} \cap N^0_{\mathcal{T},T'}} d_{\mathcal{T},\nu}^{-1}\right) \theta_{T'},$$

see Figure 2.1 for an illustration.



FIGURE 2.1.  $\psi_{T,T}$  in one dimension (with bubbles constructed using a barycentric refinement).

As a consequence of  $\langle \Xi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{L_2(\Omega)}$  being invertible, the biorthogonal 'Fortin' projector  $P_{\mathcal{T}} : L_2(\Omega) \to H^1_{0,\gamma}(\Omega)$  with ran  $P_{\mathcal{T}} = \mathscr{W}_{\mathcal{T}}$  and ran(Id  $- P_{\mathcal{T}}) = \mathscr{V}_{\mathcal{T}}^{\perp_{L_2(\Omega)}}$  exists, and is, thanks to (2.10), given by

$$P_{\mathcal{T}}u = \sum_{T \in \mathcal{T}} \frac{\langle u, \xi_T \rangle_{L_2(\Omega)}}{\langle \mathbb{1}, \xi_T \rangle_{L_2(\Omega)}} \psi_{\mathcal{T},T}.$$

To prepare for the proof of (uniform) boundedness of  $P_{\mathcal{T}}$ , we list a few properties of the collections  $\Xi_{\mathcal{T}}$ ,  $\Theta_{\mathcal{T}}$  and  $\Sigma_{\mathcal{T}}$ . For  $T \in \mathcal{T}$ , we set  $\omega_{\mathcal{T}}^{(0)}(T) := \overline{T}$ , and for  $i = 0, 1, \ldots$ , define the 'rings'

$$R_{\mathcal{T}}^{(i+1)}(T) := \{ T' \in \mathcal{T} : \overline{T'} \cap \omega_{\mathcal{T}}^{(i)}(T) \neq \emptyset \}, \quad \omega_{\mathcal{T}}^{(i+1)}(T) := \cup_{T' \in R_{\mathcal{T}}^{(i+1)}(T)} \overline{T'}.$$

It holds that

(2.12) (a) 
$$\operatorname{supp} \xi_T \subset \omega_{\mathcal{T}}^{(0)}(T)$$
, (b)  $\operatorname{supp} \sigma_{\mathcal{T},T} \subset \omega_{\mathcal{T}}^{(1)}(T)$ , (c)  $\operatorname{supp} \theta_T \subset \omega_{\mathcal{T}}^{(0)}(T)$ ,

- (2.13)  $\|\sigma_T\|_{H^k(\Omega)} \lesssim h_T^{d/2-k} \quad (k \in \{0,1\}),$
- (2.14)  $\|\theta_T\|_{H^1(\Omega)} \lesssim h_T^{-1} \|\theta_T\|_{L_2(\Omega)},$

whilst moreover  $\Xi_T$  is such that

(2.15) 
$$\langle \mathbb{1}, \xi_T \rangle_{L_2(\Omega)} = h_T^{d/2} \|\xi_T\|_{L_2(\Omega)}$$

From (2.12)(b,c), we obtain that

(2.16) 
$$\operatorname{supp} \psi_{\mathcal{T},T} \subset \omega_{\mathcal{T}}^{(1)}(T).$$

By using (2.6), (2.14), (2.13) and (2.12)(a) we infer that for  $k \in \{0, 1\}$ 

$$\left\|\frac{\langle \sigma_{\mathcal{T},T} - \delta_{TT'}\mathbb{1}, \xi_{T'}\rangle_{L_2(\Omega)}}{\langle \theta_{T'}, \xi_{T'}\rangle_{L_2(\Omega)}}\theta_{T'}\right\|_{H^k(\Omega)} \lesssim h_T^{-k} \|\sigma_{\mathcal{T},T} - \delta_{TT'}\mathbb{1}\|_{L_2(\operatorname{supp}\xi_{T'})} \lesssim h_T^{d/2-k},$$

which, by again using (2.13) and (2.14), shows that

(2.17) 
$$\|\psi_{\mathcal{T},T}\|_{H^k(\Omega)} \lesssim h_T^{d/2-k} \quad (k \in \{0,1\}).$$

**Theorem 2.3.2.** It holds that  $\sup_{\mathcal{T}\in\mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}(\mathcal{W},\mathcal{W})} < \infty$ .

*Proof.* From (2.16), (2.17), and (2.15), we have

(2.18) 
$$\begin{aligned} \|P_{\mathcal{T}}u\|_{H^{k}(T)} &\leq \sum_{T' \in R_{\mathcal{T}}^{(1)}(T)} \|\psi_{\mathcal{T},T'}\|_{H^{k}(\Omega)} \frac{\|u\|_{L_{2}(T')}\|\xi_{T'}\|_{L_{2}(\Omega)}}{|\langle \mathbb{1}, \xi_{T'} \rangle_{L_{2}(\Omega)}|} \\ &\lesssim h_{T}^{-k} \|u\|_{L_{2}(\omega_{\mathcal{T}}^{(1)}(T))} \quad (k \in \{0,1\}), \end{aligned}$$

which in particular shows that

(2.19) 
$$\sup_{\mathcal{T}\in\mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}(L_2(\Omega), L_2(\Omega))} < \infty.$$

To continue, we revisit the construction of  $\mathscr{W}_{\mathcal{T}}$  and its basis  $\Psi_{\mathcal{T}}$  by temporarily including in  $N_{\mathcal{T}}^0$  also vertices of  $\mathcal{T}$  that lie on the Dirichlet boundary

 $\gamma$ . Denoting the extended set of vertices by  $N_{\mathcal{T}}$ , consequently for the 'new'  $\psi_{\mathcal{T},T}$ , (2.9) shows that

(2.20) 
$$\sum_{T \in \mathcal{T}} \psi_{\mathcal{T},T} = \sum_{\nu \in N_{\mathcal{T}}} \phi_{\mathcal{T},\nu} = \mathbb{1} \text{ on } \Omega.$$

For any  $\nu \in N_{\mathcal{T}}$ , we select a (d-1)-face e of a  $T \in \mathcal{T}$  with  $\nu \in e$  and  $e \subset \gamma$  if  $\nu \in \gamma$ , and define the functional

$$g_{\mathcal{T},\nu}(u) := \oint_e u \, ds.$$

By the trace theorem and homogeneity arguments (see e.g [SZ90, (3.6)]), one infers that

$$|g_{\mathcal{T},\nu}(u)| \le |e|^{-1} ||u||_{L_1(e)} \lesssim h_T^{-\frac{d}{2}} ||u||_{L_2(T)} + h_T^{-\frac{d}{2}+1} |u|_{H^1(T)}$$

For  $T \in \mathcal{T}$ , we select a  $\nu \in N_T$  with  $\nu \in \gamma$  if  $\overline{T} \cap \gamma \neq \emptyset$ , and define

$$g_{\mathcal{T},T} := g_{\mathcal{T},\nu},$$

and a Scott-Zhang ([SZ90]) type quasi-interpolator  $\Pi_{\mathcal{T}}: H^1(\Omega) \to \mathscr{W}_{\mathcal{T}}^4$  by

$$\Pi_{\mathcal{T}} u = \sum_{T \in \mathcal{T}} g_{\mathcal{T},T}(u) \psi_{\mathcal{T},T}$$

It satisfies

$$\|\Pi_{\mathcal{T}} u\|_{H^{k}(T)} \lesssim h_{T}^{-k} \|u\|_{L_{2}(\omega_{\mathcal{T}}^{(2)}(T))} + h_{T}^{1-k} |u|_{H^{1}(\omega_{\mathcal{T}}^{(2)}(T))} \quad (k \in \{0, 1\}).$$

Invoking (2.20) and using that  $g_{\mathcal{T},T}(1) = 1$ , we infer that for  $k \in \{0,1\}$ 

$$\begin{aligned} \|(\mathrm{Id} - \Pi_{\mathcal{T}})u\|_{H^{k}(T)} &= \inf_{p \in \mathcal{P}_{0}} \|(\mathrm{Id} - \Pi_{\mathcal{T}})(u - p)\|_{H^{k}(T)} \\ &\leq \inf_{p \in \mathcal{P}_{0}} \|u - p\|_{H^{k}(T)} + h_{T}^{-k} \|u - p\|_{L_{2}(\omega_{\mathcal{T}}^{(2)}(T))} + h_{T}^{1-k} |u|_{H^{1}(\omega_{\mathcal{T}}^{(2)}(T))} \\ &\approx \inf_{p \in \mathcal{P}_{0}} h_{T}^{-k} \|u - p\|_{L_{2}(\omega_{\mathcal{T}}^{(2)}(T))} + h_{T}^{1-k} |u|_{H^{1}(\omega_{\mathcal{T}}^{(2)}(T))} \\ &\approx h_{T}^{1-k} |u|_{H^{1}(\omega_{\mathcal{T}}^{(2)}(T))} \end{aligned}$$

by an application of the Bramble-Hilbert lemma (cf. [SZ90, (4.2)]).

Noting that the 'new'  $\psi_{\mathcal{T},T}$  differs only from the 'old', original one when  $\overline{T} \cap \gamma \neq \emptyset$ , and that for those T and  $u \in H^1_{0,\gamma}(\Omega)$  it holds that  $g_{\mathcal{T},T}(u) = 0$ , we conclude that  $\operatorname{ran} \Pi_{\mathcal{T}}|_{H^1_{0,\gamma}(\Omega)}$  is included in the original space  $\mathscr{W}_{\mathcal{T}}$ , which we

<sup>&</sup>lt;sup>4</sup>The existence of such a  $\Pi_{\mathcal{T}}$  which satisfies an estimate of type (2.21) for k = 0 can be used as a definition of a (lowest order) approximation property of  $\mathcal{W}_{\mathcal{T}}$ .

consider again from here on. Using that  $P_{\mathcal{T}}$  is a projector onto this  $\mathscr{W}_{\mathcal{T}}$ , for  $u \in H^1_{0,\gamma}(\Omega)$  writing  $P_{\mathcal{T}}u = \Pi_{\mathcal{T}}u + P_{\mathcal{T}}(\mathrm{Id} - \Pi_{\mathcal{T}})u$ , using (2.18) and (2.21) for  $k \in \{0, 1\}$  we arrive at

$$\begin{split} \|P_{\mathcal{T}}u\|_{H^{1}(T)} &\lesssim \|\Pi_{\mathcal{T}}u\|_{H^{1}(T)} + h_{T}^{-1}\|(\mathrm{Id} - \Pi_{\mathcal{T}})u\|_{L_{2}(\omega_{\mathcal{T}}^{(1)}(T))} \\ &\lesssim \|u\|_{H^{1}(\omega_{\mathcal{T}}^{(2)}(T))} + h_{T}^{-1}\|(\mathrm{Id} - \Pi_{\mathcal{T}})u\|_{L_{2}(\omega_{\mathcal{T}}^{(1)}(T))} \\ &\lesssim \|u\|_{H^{1}(\omega_{\mathcal{T}}^{(3)}(T))}, \end{split}$$

and consequently,

$$\sup_{\mathcal{T}\in\mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}(H^1_{0,\gamma}(\Omega),H^1_{0,\gamma}(\Omega)))} < \infty.$$

In combination with (2.19), the proof is completed by an application of the Riesz-Thorin interpolation theorem.  $\hfill \Box$ 

From Proposition 2.2.3 and Theorem 2.3.2 we conclude the following:

**Corollary 2.3.3.** For  $D_{\mathcal{T}}: \mathscr{V}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$  defined by  $(D_{\mathcal{T}}v)(w) := (Dv)(w) = \langle v, w \rangle_{L_{2}(\Omega)}$ , it holds that  $D_{\mathcal{T}} \in \mathcal{L}$ is $(\mathscr{V}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$  with  $\|D_{\mathcal{T}}\|_{\mathcal{L}}(\mathscr{V}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}}) \leq 1$  and  $\sup_{\mathcal{T} \in \mathbb{T}} \|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}}(\mathscr{W}'_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}}) \leq \sup_{\mathcal{T} \in \mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}}(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}) < \infty.$ 

This result is thus valid *without any additional assumptions on the mesh grading*. The latter is a consequence of the fact that we were able to equip  $\mathscr{V}_{\mathcal{T}}$  and  $\mathscr{W}_{\mathcal{T}}$  with local biorthogonal bases. (Compare [Ste03a, eq. (2.30)] for conditions on the mesh grading without having local biorthogonal bases). Additionally, the biorthogonality has the important advantage of the matrix

$$\boldsymbol{D}_{\mathcal{T}} = \langle \Xi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{L_2(\Omega)} = \operatorname{diag}\{|T| \colon T \in \mathcal{T}\}$$

being diagonal.

Before we discuss in §2.3.3 the construction of  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$ , being the last ingredient of our preconditioner, in the following subsection §2.3.2 we revisit the construction of  $\mathscr{W}_{\mathcal{T}}$  and  $D_{\mathcal{T}}$  in the manifold case.

#### **2.3.2** Construction of $\mathscr{W}_{\mathcal{T}}$ and $D_{\mathcal{T}}$ in the manifold case

Let  $\Gamma$  be a compact *d*-dimensional Lipschitz, piecewise smooth manifold in  $\mathbb{R}^{d'}$  for some  $d' \geq d$  with or without boundary  $\partial \Gamma$ . For some closed measurable  $\gamma \subset \partial \Gamma$  and  $s \in [0, 1]$ , let

$$\mathscr{W} := [L_2(\Gamma), H^1_{0,\gamma}(\Gamma)]_{s,2}, \quad \mathscr{V} := \mathscr{W}'.$$

We assume that  $\Gamma$  is given as the essentially disjoint union of  $\bigcup_{i=1}^{p} \overline{\chi_i(\Omega_i)}$ , with, for  $1 \leq i \leq p, \chi_i \colon \mathbb{R}^d \to \mathbb{R}^{d'}$  being some smooth regular parametrization, and
$\Omega_i \subset \mathbb{R}^d$  an open polytope. W.l.o.g. assuming that for  $i \neq j$ ,  $\overline{\Omega}_i \cap \overline{\Omega}_j = \emptyset$ , we define

$$\chi: \Omega:=\cup_{i=1}^p \Omega_i \to \cup_{i=1}^p \chi_i(\Omega_i)$$
 by  $\chi|_{\Omega_i}=\chi_i$ .

Let  $\mathbb{T}$  be a family of conforming partitions  $\mathcal{T}$  of  $\Gamma$  into 'panels' such that, for  $1 \leq i \leq p$ ,  $\chi^{-1}(\mathcal{T}) \cap \Omega_i$  is a uniformly shape regular conforming partition of  $\Omega_i$  into *d*-simplices (that for d = 1 satisfies a uniform *K*-mesh property). We assume that  $\gamma$  is a (possibly empty) union of 'faces' of  $T \in \mathcal{T}$  (i.e., sets of type  $\chi_i(e)$ , where e is a (d-1)-dimensional face of  $\chi_i^{-1}(T)$ ).

As in Sect. 2.3, for  $\mathcal{T} \in \mathbb{T}$ , we define  $N^0_{\mathcal{T}}$  as the set of vertices of  $\mathcal{T}$  that are not on  $\gamma$ , set  $d_{\mathcal{T},\nu} := \#\{T \in \mathcal{T} : \nu \in \overline{T}\}$ , and for  $T \in \mathcal{T}$ , define  $h_T := |T|^{1/d}$ and  $N_{T,T}^0 := N_T^0 \cap N_T$ , with  $N_T$  being the set of the vertices of T.

We set

$$\begin{aligned} \mathscr{V}_{\mathcal{T}} &= \mathscr{S}_{\mathcal{T}}^{-1,0} \coloneqq \{ u \in L_2(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_0 \ (T \in \mathcal{T}) \} \subset \mathscr{V}, \\ &\qquad \mathscr{S}_{\mathcal{T},0}^{0,1} \coloneqq \{ u \in H^1_{0,\gamma}(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_1 \ (T \in \mathcal{T}) \}, \end{aligned}$$

equipped with  $\Xi_{\mathcal{T}} = \{\xi_T \colon T \in \mathcal{T}\}$  and  $\Phi_{\mathcal{T}} = \{\phi_{\mathcal{T},\nu} \colon \nu \in N^0_{\mathcal{T}}\}$ , respectively, defined by  $\xi_T := 1$  on T,  $\xi_T := 0$  elsewhere, and  $\phi_{\mathcal{T},\nu}(\nu') = \delta_{\nu,\nu'}(\nu,\nu' \in N^0_{\mathcal{T}})$ . Furthermore, we define  $\Sigma_{\mathcal{T}} = \{\sigma_{\mathcal{T},T} : T \in \mathcal{T}\} \subset \mathscr{S}^{0,1}_{\mathcal{T},0}$  and  $\Theta_{\mathcal{T}} = \{\theta_T : T \in \mathcal{T}\} \subset H^1_{0,\gamma}(\Gamma)$  by  $\sigma_{\mathcal{T},T} := \sum_{\nu \in N^0_{\mathcal{T},T}} d^{-1}_{\mathcal{T},\nu} \phi_{\mathcal{T},\nu'}, \theta_T := \theta_{\chi^{-1}(T)} \circ \chi^{-1}$  on T and  $\theta_T := 0$  elsewhere. Thanks to our assumption of  $\theta_{\chi^{-1}(T)} \ge 0$ , it holds that  $\langle \theta_T, \xi_T \rangle_{L_2(\Gamma)} \approx \langle \theta_{\chi^{-1}(T)}, \xi_{\chi^{-1}(T)} \rangle_{L_2(\chi^{-1}(T))} \approx \|\theta_T\|_{L_2(\Gamma)} \|\xi_T\|_{L_2(\Gamma)} \text{ (cf. (2.6)).}$ Now defining  $\Psi_{\mathcal{T}} := \{\psi_{\mathcal{T},T} \colon T \in \mathcal{T}\}$  and  $\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}} \subset \mathscr{W}$  by

$$(2.22) \quad \psi_{\mathcal{T},T} := \sigma_{\mathcal{T},T} + \frac{\langle \mathbb{1} - \sigma_{\mathcal{T},T}, \xi_T \rangle_{L_2(\Gamma)}}{\langle \theta_T, \xi_T \rangle_{L_2(\Gamma)}} \theta_T - \sum_{T' \in \mathcal{T} \setminus \{T\}} \frac{\langle \sigma_{\mathcal{T},T}, \xi_{T'} \rangle_{L_2(\Gamma)}}{\langle \theta_{T'}, \xi_{T'} \rangle_{L_2(\Gamma)}} \theta_{T'},$$

and  $D_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$  by  $(D_{\mathcal{T}}v)(w) \coloneqq (Dv)(w) = \langle v, w \rangle_{L_2(\Gamma)}$ , the analysis from Sect. 2.3 applies verbatim by only changing  $\langle , \rangle_{L_2(\Omega)}$  into  $\langle , \rangle_{L_2(\Gamma)}$ . It yields that  $\|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}}')} \leq 1$ ,  $\sup_{\mathcal{T}\in\mathbb{T}}\|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}',\mathscr{V}_{\mathcal{T}})} < \infty$ , and  $D_{\mathcal{T}} =$ diag{ $\langle \mathbb{1}, \xi_T \rangle_{L_2(\Gamma)} \colon T \in \mathcal{T}$ }.

A hidden problem, however, is that the computation of  $D_{\mathcal{T}}$ , and that of the scalar products in (2.22) involve integrals over  $\Gamma$  that generally have to be approximated using numerical quadrature. Recalling that, for s > 0, the preconditioner  $G_{\mathcal{T}} = D_{\mathcal{T}}^{-1} B_{\mathcal{T}} D_{\mathcal{T}}^{-\top}$  is *not* a uniformly well-conditioned matrix, it is a priorily not clear which quadrature errors are allowable, in particular when T is far from being quasi-uniform. For this reason, in the following §2.3.2 we propose a slightly modified construction of  $\mathscr{W}_{\mathcal{T}}$  and  $D_{\mathcal{T}}$  that does not require the evaluation of integrals over  $\Gamma$ .

As a preparation, in the following lemma we present a non-standard inverse inequality on the family  $(\mathscr{V}_{\mathcal{T}})_{\mathcal{T}\in\mathbb{T}}$ . Proofs of this inequality for  $d \leq 3$  can be found in [DFG<sup>+</sup>04, GHS05]. It turns out that our construction of a 'local' collection  $\Psi_{\mathcal{T}} \subset H^1_{0,\gamma}(\Omega)$  that is biorthogonal to  $\Xi_{\mathcal{T}}$  allows for a very simple proof.

**Lemma 2.3.4** (inverse inequality). With  $h_T|_T := h_T$ , it holds that

$$\|h_{\mathcal{T}} v_{\mathcal{T}}\|_{L_2(\Gamma)} \lesssim \|v_{\mathcal{T}}\|_{H^1_{0,\gamma}(\Gamma)'} \quad (v_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}).$$

*Proof.* For  $P_{\mathcal{T}}: L_2(\Gamma) \to H^1_{0,\gamma}(\Gamma)$  defined by

$$P_{\mathcal{T}}u = \sum_{T \in \mathcal{T}} \frac{\langle u, \xi_T \rangle_{L_2(\Gamma)}}{\langle \mathbb{1}, \xi_T \rangle_{L_2(\Gamma)}} \psi_{\mathcal{T},T}.$$

we have  $\operatorname{ran}(\operatorname{Id} - P_{\mathcal{T}}) = \mathscr{V}_{\mathcal{T}}^{\perp_{L_2(\Gamma)}}$ , and as follows from (2.18),

$$||P_{\mathcal{T}}u||_{H^1(\Gamma)} \lesssim ||h_{\mathcal{T}}^{-1}u||_{L_2(\Gamma)} \quad (u \in L_2(\Gamma)).$$

The proof is completed by

$$\|v_{\mathcal{T}}\|_{H^{1}_{0,\gamma}(\Gamma)'} = \sup_{0 \neq w \in H^{1}_{0,\gamma}(\Gamma)} \frac{\langle v_{\mathcal{T}}, w \rangle_{L_{2}(\Gamma)}}{\|w\|_{H^{1}(\Gamma)}} \geq \frac{\langle v_{\mathcal{T}}, P_{\mathcal{T}}h^{2}_{\mathcal{T}}v_{\mathcal{T}} \rangle_{L_{2}(\Gamma)}}{\|P_{\mathcal{T}}h^{2}_{\mathcal{T}}v_{\mathcal{T}}\|_{H^{1}(\Gamma)}} \gtrsim \frac{\langle h_{\mathcal{T}}v_{\mathcal{T}}, h_{\mathcal{T}}v_{\mathcal{T}} \rangle_{L_{2}(\Gamma)}}{\|h_{\mathcal{T}}v_{\mathcal{T}}\|_{L_{2}(\Gamma)}}$$

### Modified construction for manifolds

Given  $\mathcal{T} \in \mathbb{T}$ , on  $L_2(\Gamma)$  we define an additional, 'mesh-dependent' scalar product

$$\langle u, v \rangle_{\mathcal{T}} \coloneqq \sum_{T \in \mathcal{T}} \frac{|T|}{|\chi^{-1}(T)|} \int_{\chi^{-1}(T)} u(\chi(x)) v(\chi(x)) dx.$$

It is constructed from

<

$$\langle u, v \rangle_{L_2(\Gamma)} = \int_{\Omega} u(\chi(x))v(\chi(x))|\partial\chi(x)|dx$$

by replacing on each  $\chi^{-1}(T)$ , the Jacobian  $|\partial \chi|$  by its average  $\frac{|T|}{|\chi^{-1}(T)|}$  over  $\chi^{-1}(T)$ .<sup>5</sup>

We now redefine  $\Psi_{\mathcal{T}} := \{\psi_{\mathcal{T},T} \colon T \in \mathcal{T}\}, \mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}} \subset \mathscr{W}$  by

$$\psi_{\mathcal{T},T} := \sigma_{\mathcal{T},T} + \frac{\langle \mathbb{1} - \sigma_{\mathcal{T},T}, \xi_T \rangle_{\mathcal{T}}}{\langle \theta_T, \xi_T \rangle_{\mathcal{T}}} \theta_T - \sum_{T' \in \mathcal{T} \setminus \{T\}} \frac{\langle \sigma_{\mathcal{T},T}, \xi_{T'} \rangle_{\mathcal{T}}}{\langle \theta_{T'}, \xi_{T'} \rangle_{\mathcal{T}}} \theta_{T'},$$

and  $D_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$  by  $(D_{\mathcal{T}}v_{\mathcal{T}})(w_{\mathcal{T}}) := \langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{\mathcal{T}}$ . Then, as in the domain case, we get the explicit formulas

$$D_{\mathcal{T}} = \langle \Xi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{\mathcal{T}} = \operatorname{diag}\{ \langle \mathbb{1}, \xi_T \rangle_{\mathcal{T}} \colon T \in \mathcal{T} \} = \operatorname{diag}\{ |T| \colon T \in \mathcal{T} \},$$

<sup>&</sup>lt;sup>5</sup>It will be clear from the following that  $\frac{|T|}{|\chi^{-1}(T)|}$  can be read as *any* constant approximation to  $|\partial\chi|$  on  $L_{\infty}(\chi^{-1}(T))$ -distance  $\leq h_{\chi^{-1}(T)}$ , for example  $|\partial\chi(z)|$  for some  $z \in \chi^{-1}(T)$ . Then in the following, the volumes |T| in the expression for  $D_{\mathcal{T}}$  should be read as  $|\chi^{-1}(T)||\partial\chi(z)|$ , with which also the computation of |T| is avoided.

and  
(2.23)  
$$\psi_{\mathcal{T},T} = \sum_{\nu \in N^0_{\mathcal{T},T}} d_{\mathcal{T},\nu}^{-1} \phi_{\mathcal{T},\nu} + \left(1 - \frac{1}{d+1} \sum_{\nu \in N^0_{\mathcal{T},T}} d_{\mathcal{T},\nu}^{-1}\right) \theta_T - \sum_{T' \in \mathcal{T} \setminus \{T\}} \left(\frac{1}{d+1} \sum_{\nu \in N^0_{\mathcal{T},T} \cap N^0_{\mathcal{T},T'}} d_{\mathcal{T},\nu}^{-1}\right) \theta_{T'}.$$

thus with coefficients that are *independent* of  $\chi$ .

What remains is to prove the uniform boundedness of  $||D_{\mathcal{T}}||_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}}')}$ , and that of  $||D_{\mathcal{T}}^{-1}||_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}',\mathscr{V}_{\mathcal{T}})}$ . Because of the definition of  $D_{\mathcal{T}}$  in terms of the mesh-dependent scalar product, for doing so we cannot simply rely on the 'Fortin criterion' from Proposition 2.2.3.

**Lemma 2.3.5.** It holds that  $\sup_{\mathcal{T}\in\mathbb{T}} \|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}}')} < \infty$ .

*Proof.* If s = 0, i.e., when  $\mathscr{W} = L_2(\Gamma) \simeq L_2(\Gamma)' = \mathscr{V}$ , then the uniform boundedness of  $\|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{W}'_{\mathcal{T}})}$  follows directly from  $\langle \cdot, \cdot \rangle_{\mathcal{T}} \approx \|\cdot\|^2_{L_2(\Gamma)}$ .

By an interpolation argument, in the following it suffices to consider the case s = 1, i.e.,  $\mathcal{W} = H^1_{0,\gamma}(\Gamma)$  and  $\mathcal{V} = H^1_{0,\gamma}(\Gamma)'$ . By definition of  $\langle , \rangle_{\mathcal{T}}$ , it holds that

$$(2.24) \qquad |\langle v, u \rangle_{\mathcal{T}} - \langle v, u \rangle_{L_2(\Gamma)}| \lesssim ||h_{\mathcal{T}}v||_{L_2(\Gamma)} ||u||_{L_2(\Gamma)} \quad (v, u \in L_2(\Gamma)).$$

By writing  $(D_{\mathcal{T}}v_{\mathcal{T}})(w_{\mathcal{T}}) = \langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{L_2(\Gamma)} + \langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{\mathcal{T}} - \langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{L_2(\Gamma)}$ , the uniform boundedness of  $\|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')}$  (for s = 1) now follows by combining (2.24) and Lemma 2.3.4.

The  $\langle , \rangle_{\mathcal{T}}$ -biorthogonal projector  $\check{P}_{\mathcal{T}} : L_2(\Omega) \to H^1_{0,\gamma}(\Omega)$  with ran  $\check{P}_{\mathcal{T}} = \mathscr{W}_{\mathcal{T}}$ and ran $(\mathrm{Id}-\check{P}_{\mathcal{T}}) = \mathscr{V}_{\mathcal{T}}^{\perp_{\langle , \rangle_{\mathcal{T}}}}$  exists and is given by  $\check{P}_{\mathcal{T}}u = \sum_{T \in \mathcal{T}} |T|^{-1} \langle u, \xi_T \rangle_{\mathcal{T}} \psi_{\mathcal{T},T}$ . Since  $\langle , \rangle_{\mathcal{T}}$  gives rise to a norm that is uniformly equivalent to  $\| \|_{L_2(\Gamma)}$ , the proof of Theorem 2.3.2 again applies, and shows that

$$\sup_{\mathcal{T}\in\mathbb{T}}\|\check{P}_{\mathcal{T}}\|_{\mathcal{L}(L_{2}(\Gamma),L_{2}(\Gamma))}<\infty,\quad \sup_{\mathcal{T}\in\mathbb{T}}\|\check{P}_{\mathcal{T}}\|_{\mathcal{L}(H^{1}_{0,\gamma}(\Gamma),H^{1}_{0,\gamma}(\Gamma))}<\infty.$$

as well as

(2.25) 
$$\|\check{P}_{\mathcal{T}}u\|_{H^1(\Gamma)} \lesssim \|h_{\mathcal{T}}^{-1}u\|_{L_2(\Gamma)} \quad (u \in L_2(\Gamma)).$$

These properties of  $\check{P}_{\mathcal{T}}$  will be the key to prove the uniform boundedness of  $\|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}})}$ . Indeed, for s = 0 uniform boundedness of  $\|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}})}$  follows from

$$(D_{\mathcal{T}}v_{\mathcal{T}})(\check{P}_{\mathcal{T}}v_{\mathcal{T}}) = \langle v_{\mathcal{T}}, v_{\mathcal{T}} \rangle_{\mathcal{T}} = \|v_{\mathcal{T}}\|_{L_{2}(\Gamma)}^{2} \gtrsim \|v_{\mathcal{T}}\|_{L_{2}(\Gamma)} \|\check{P}_{\mathcal{T}}v_{\mathcal{T}}\|_{L_{2}(\Gamma)}$$

To conclude, by an interpolation argument, uniform boundedness of  $||D_{\mathcal{T}}^{-1}||_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}}')}$  for any  $s \in [0, 1]$ , it is sufficient to verify the case s = 1, which will be done using the following modified inverse inequality.

Lemma 2.3.6. It holds that

$$\|h_{\mathcal{T}}v_{\mathcal{T}}\|_{L_{2}(\Gamma)} \lesssim \sup_{0 \neq w \in H^{1}_{0,\gamma}(\Gamma)} \frac{\langle v_{\mathcal{T}}, w \rangle_{\mathcal{T}}}{\|w\|_{H^{1}(\Gamma)}} \quad (v_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}).$$

*Proof.* Similar to proof of Lemma 2.3.4, using (2.25) for  $v_T \in \mathscr{V}_T$  we estimate

 $\sup_{0 \neq w \in H^1_{0,\gamma}(\Gamma)} \frac{\langle v_{\mathcal{T}}, w \rangle_{\mathcal{T}}}{\|w\|_{H^1(\Gamma)}} \geq \frac{\langle v_{\mathcal{T}}, \check{P}_{\mathcal{T}} h^2_{\mathcal{T}} v_{\mathcal{T}} \rangle_{\mathcal{T}}}{\|\check{P}_{\mathcal{T}} h^2_{\mathcal{T}} v_{\mathcal{T}}\|_{H^1(\Gamma)}} \gtrsim \frac{\langle h_{\mathcal{T}} v_{\mathcal{T}}, h_{\mathcal{T}} v_{\mathcal{T}} \rangle_{\mathcal{T}}}{\|h_{\mathcal{T}} v_{\mathcal{T}}\|_{L_2(\Gamma)}} \approx \|h_{\mathcal{T}} v_{\mathcal{T}}\|_{L_2(\Gamma)}. \Box$ 

Corollary 2.3.7. It holds that

$$\|v_{\mathcal{T}}\|_{H^{1}_{0,\gamma}(\Gamma)'} \approx \sup_{0 \neq w_{\mathcal{T}} \in \mathscr{W}_{\mathcal{T}}} \frac{\langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{\mathcal{T}}}{\|w_{\mathcal{T}}\|_{H^{1}(\Gamma)}} \quad (v_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}),$$

(with ' $\leq$ ' being the statement  $\sup_{\mathcal{T}\in\mathbb{T}} \|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathcal{W}_{\mathcal{T}}',\mathcal{V}_{\mathcal{T}})} < \infty$  for s = 1).

*Proof.* The inequality ' $\gtrsim$ ' is the statement of Lemma 2.3.5 for s = 1. To prove the other direction, for  $v \in L_2(\Gamma)$ , (2.24) shows that

$$\left| \|v\|_{H^{1}_{0,\gamma}(\Gamma)'} - \sup_{0 \neq w \in H^{1}_{0,\gamma}(\Gamma)} \frac{\langle v, w \rangle_{\mathcal{T}}}{\|w\|_{H^{1}(\Gamma)}} \right| \lesssim \|h_{\mathcal{T}}v\|_{L_{2}(\Gamma)},$$

Taking  $v = v_T \in \mathscr{V}_T$ , from Lemma 2.3.6 we conclude that

$$\|v_{\mathcal{T}}\|_{H^{1}_{0,\gamma}(\Gamma)'} \lesssim \sup_{0 \neq w \in H^{1}_{0,\gamma}(\Gamma)} \frac{\langle v_{\mathcal{T}}, w \rangle_{\mathcal{T}}}{\|w\|_{H^{1}(\Gamma)}} = \sup_{0 \neq w \in H^{1}_{0,\gamma}(\Gamma)} \frac{\langle v_{\mathcal{T}}, \check{P}_{\mathcal{T}}w \rangle_{\mathcal{T}}}{\|w\|_{H^{1}(\Gamma)}}$$
$$\leq \|\check{P}_{\mathcal{T}}\|_{\mathcal{L}(H^{1}_{0,\gamma}(\Gamma), H^{1}_{0,\gamma}(\Gamma))} \sup_{0 \neq w_{\mathcal{T}} \in \mathscr{W}_{\mathcal{T}}} \frac{\langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{\mathcal{T}}}{\|w_{\mathcal{T}}\|_{H^{1}(\Gamma)}} \lesssim \sup_{0 \neq w_{\mathcal{T}} \in \mathscr{W}_{\mathcal{T}}} \frac{\langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{\mathcal{T}}}{\|w_{\mathcal{T}}\|_{H^{1}(\Gamma)}}$$

by  $\sup_{\mathcal{T}\in\mathbb{T}} \|\check{P}_{\mathcal{T}}\|_{\mathcal{L}(H^1_{0,\gamma}(\Gamma),H^1_{0,\gamma}(\Gamma))} < \infty.$ 

# **2.3.3** Construction of $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$ .

Having established  $\sup_{T \in \mathbb{T}} \max \left( \|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')}, \|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}', \mathscr{V}_{\mathcal{T}})} \right) < \infty$ , in both domain and manifold case, for the construction of uniform preconditioners it remains to find  $B_{\mathcal{T}} \in \mathcal{L}$ is<sub>c</sub> $(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  with  $\sup_{T \in \mathbb{T}} \|B_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')} < \infty$  and  $\sup_{T \in \mathbb{T}} \|\Re(B_{\mathcal{T}})^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}', \mathscr{W}_{\mathcal{T}})} < \infty$ .

We add the following two assumptions on the collection  $\Theta_{\mathcal{T}}$  of 'bubbles' and their span  $\mathscr{B}_{\mathcal{T}} := \operatorname{span} \Theta_{\mathcal{T}}$ . For  $k \in \{0, 1\}$  it holds that

(2.26) 
$$\|\sum_{T\in\mathcal{T}} c_T \theta_T\|_{H^k(\Omega)}^2 \approx \sum_{T\in\mathcal{T}} h_T^{-2k} |c_T|^2 \|\theta_T\|_{L_2(\Omega)}^2, \quad ((c_T)_{T\in\mathcal{T}}\subset\mathbb{R}),$$

(2.27) 
$$\|u+v\|_{H^{k}(\Omega)}^{2} \gtrsim \|u\|_{H^{k}(\Omega)}^{2} + \|v\|_{H^{k}(\Omega)}^{2} \quad (u \in \mathscr{S}_{\mathcal{T},0}^{0,1}, v \in \mathscr{B}_{\mathcal{T}})$$

(Here and in the following,  $\Omega$  should be read as  $\Gamma$  in the manifold case). Both properties are easily verified by a standard homogeneity argument for both our earlier specifications of possible  $\Theta_{\mathcal{T}}$ . From (2.27) it follows that  $\mathscr{S}_{\mathcal{T},0}^{0,1} \cap \mathscr{B}_{\mathcal{T}} = \{0\}$ . Let  $I_{\mathcal{T}}^{\mathscr{G}}$  be the linear projector defined on  $\mathscr{S}^{0,1} \oplus \mathscr{B}_{\mathcal{T}}$  by ran  $I_{\mathcal{T}}^{\mathscr{G}} = \mathscr{S}_{\mathcal{T},0}^{0,1}$  and ran  $I_{\mathcal{T}}^{\mathscr{B}} = \mathscr{B}_{\mathcal{T}}$ , where  $I_{\mathcal{T}}^{\mathscr{B}} := \mathrm{Id} - I_{\mathcal{T}}^{\mathscr{G}}$ . Below we give a construction of suitable  $B_{\mathcal{T}}$  that is *independent* of the

Below we give a construction of suitable  $B_{\mathcal{T}}$  that is *independent* of the particular bubbles  $\Theta_{\mathcal{T}}$  being chosen. Like  $\mathscr{W}_{\mathcal{T}}$ , we equip  $\mathscr{S}_{\mathcal{T},0}^{0,1}$ ,  $\mathscr{B}_{\mathcal{T}}$ , and  $\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}}$  with  $\| \|_{\mathscr{W}}$ .

**Proposition 2.3.8.** Given  $B_{\mathcal{T}}^{\mathscr{S}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  and  $B_{\mathcal{T}}^{\mathscr{B}} \in \mathcal{L}is_c(\mathscr{B}_{\mathcal{T}}, \mathscr{B}_{\mathcal{T}}')$ , let  $B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}} \colon \mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}} \to (\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}})'$  be defined by

$$(B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}}w)(\tilde{w}) \coloneqq (B_{\mathcal{T}}^{\mathscr{S}}I_{\mathcal{T}}^{\mathscr{S}}w)(I_{\mathcal{T}}^{\mathscr{S}}\tilde{w}) + (B_{\mathcal{T}}^{\mathscr{B}}I_{\mathcal{T}}^{\mathscr{B}}w)(I_{\mathcal{T}}^{\mathscr{B}}\tilde{w}).$$

Then thanks to (2.27), one has  $B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}}, (\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}})')$ , and

$$\begin{split} \|\Re(B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}})^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},0}^{0,1}\oplus\mathscr{B}_{\mathcal{T}})',\mathscr{S}_{\mathcal{T},0}^{0,1}\oplus\mathscr{B}_{\mathcal{T}})} \\ &\leq 2\max(\|\Re(B_{\mathcal{T}}^{\mathscr{S}})^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},0}^{0,1})',\mathscr{S}_{\mathcal{T},0}^{0,1})}, \|\Re(B_{\mathcal{T}}^{\mathscr{B}})^{-1}\|_{\mathcal{L}(\mathscr{B}_{\mathcal{T}}',\mathscr{B}_{\mathcal{T}})}), \\ \|B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}}\|_{\mathcal{L}(\mathscr{S}_{\mathcal{T},0}^{0,1}\oplus\mathscr{B}_{\mathcal{T}}, (\mathscr{S}_{\mathcal{T},0}^{0,1}\oplus\mathscr{B}_{\mathcal{T}})')} \lesssim \max(\|B_{\mathcal{T}}^{\mathscr{S}}\|_{\mathcal{L}(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')}, \|B_{\mathcal{T}}^{\mathscr{B}}\|_{\mathcal{L}(\mathscr{B}_{\mathcal{T}}, \mathscr{B}_{\mathcal{T}}')}). \end{split}$$

Proof. One has

$$\begin{split} |(B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}}w)(w)| \geq \min(\|\Re(B_{\mathcal{T}}^{\mathscr{S}})^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},0}^{0,1})',\mathscr{S}_{\mathcal{T},0}^{0,1})}^{-1}, \|\Re(B_{\mathcal{T}}^{\mathscr{B}})^{-1}\|_{\mathcal{L}(\mathscr{B}_{\mathcal{T}}',\mathscr{B}_{\mathcal{T}})}^{-1}) \\ & \times (\|I_{\mathcal{T}}^{\mathscr{S}}w\|_{\mathscr{W}}^{2} + \|I_{\mathcal{T}}^{\mathscr{B}}w\|_{\mathscr{W}}^{2}), \end{split}$$

and

$$\begin{split} |(B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}}w)(\tilde{w})| &\leq \max(\|B_{\mathcal{T}}^{\mathscr{S}}\|_{\mathcal{L}(\mathscr{S}_{\mathcal{T},0}^{0,1},(\mathscr{S}_{\mathcal{T},0}^{0,1})')}, \|B_{\mathcal{T}}^{\mathscr{B}}\|_{\mathcal{L}(\mathscr{B}_{\mathcal{T}},\mathscr{B}_{\mathcal{T}}')}) \\ &\times \sqrt{\|I_{\mathcal{T}}^{\mathscr{S}}w\|_{\mathscr{W}}^{2}} + \|I_{\mathcal{T}}^{\mathscr{B}}w\|_{\mathscr{W}}^{2}} \sqrt{\|I_{\mathcal{T}}^{\mathscr{S}}\tilde{w}\|_{\mathscr{W}}^{2}} + \|I_{\mathcal{T}}^{\mathscr{B}}\tilde{w}\|_{\mathscr{W}}^{2}} \end{split}$$

From the triangle inequality and (2.27), one has  $\frac{1}{2} \|w\|_{\mathscr{W}}^2 \leq \|I_{\mathcal{T}}^{\mathscr{S}}w\|_{\mathscr{W}}^2 + \|I_{\mathcal{T}}^{\mathscr{B}}w\|_{\mathscr{W}}^2 \lesssim \|w\|_{\mathscr{W}}^2$ , which completes the proof.

By equipping  $\mathscr{W}_{\mathcal{T}}$ ,  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  and  $\mathscr{B}_{\mathcal{T}}$  by  $\Psi_{\mathcal{T}}$ ,  $\Phi_{\mathcal{T}}$ , and  $\Theta_{\mathcal{T}}$ , respectively, the applications of  $I_{\mathcal{T}}^{\mathscr{S}}|_{\mathscr{W}_{\mathcal{T}}}$  and  $asily determined in linear complexity. Therefore a suitable definition of <math>B_{\mathcal{T}}: \mathscr{W}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$  is given by  $(B_{\mathcal{T}}w)(\tilde{w}):=(B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}}w)(\tilde{w})$ . Clearly,

$$\|B_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}}')} \leq \|B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}}\|_{\mathcal{L}(\mathscr{S}_{\mathcal{T},0}^{0,1}\oplus\mathscr{B}_{\mathcal{T}},(\mathscr{S}_{\mathcal{T},0}^{0,1}\oplus\mathscr{B}_{\mathcal{T}})')},$$
$$\|\Re(B_{\mathcal{T}})^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}',\mathscr{W}_{\mathcal{T}})} \leq \|\Re(B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}})^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},0}^{0,1}\oplus\mathscr{B}_{\mathcal{T}})',\mathscr{S}_{\mathcal{T},0}^{0,1}\oplus\mathscr{B}_{\mathcal{T}})}$$

An obvious choice for  $B_{\mathcal{T}}^{\mathscr{B}} \in \mathcal{L}is_c(\mathscr{B}_{\mathcal{T}}, (\mathscr{B}_{\mathcal{T}})')$  such that

$$\max\left(\sup_{\mathcal{T}\in\mathbb{T}}\|B_{\mathcal{T}}^{\mathscr{B}}\|_{\mathcal{L}(\mathscr{B}_{\mathcal{T}},\mathscr{B}_{\mathcal{T}}')},\|\Re(B_{\mathcal{T}}^{\mathscr{B}})^{-1}\|_{\mathcal{L}(\mathscr{B}_{\mathcal{T}}',\mathscr{B}_{\mathcal{T}})}\right)<\infty,$$

is, in view of (2.26) and  $\|\theta_T\|_{L_2(\Omega)} \stackrel{(2.6)}{\approx} \frac{\langle \theta_T, \xi_T \rangle_{L_2(\Omega)}}{\|\xi_T\|_{L_2(\Omega)}} \stackrel{(2.7)}{=} |T| \|\xi_T\|_{L_2(\Omega)}^{-1} = h^{d/2}$ , given by

(2.28) 
$$\left(B_{\mathcal{T}}^{\mathscr{B}}\sum_{T\in\mathcal{T}}c_{T}\theta_{T}\right)\left(\sum_{T\in\mathcal{T}}d_{T}\theta_{T}\right) := \beta \sum_{T\in\mathcal{T}}h_{T}^{d-2s}c_{T}d_{T}.$$

for some constant  $\beta$ .

Possible choices for  $B_{\mathcal{T}}^{\mathscr{S}} \in \mathcal{L}$ is<sub>c</sub> $(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  with

$$\sup_{\mathcal{T}\in\mathbb{T}} \left( \|B_{\mathcal{T}}^{\mathscr{S}}\|_{\mathcal{L}(\mathscr{S}_{\mathcal{T},0}^{0,1},(\mathscr{S}_{\mathcal{T},0}^{0,1})')}, \|\Re(B_{\mathcal{T}}^{\mathscr{S}})^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},0}^{0,1})',\mathscr{S}_{\mathcal{T},0}^{0,1})}\right) < \infty$$

include  $(B_{\mathcal{T}}^{\mathscr{S}}u)(v) := (Bu)(v) \ (u, v \in \mathscr{S}_{\mathcal{T},0}^{0,1})$  for some  $B \in \mathcal{L}is_c(\mathscr{W}, \mathscr{W}')$ .

For  $d \in \{2,3\}$  and  $\mathscr{W} = H_{00}^{\frac{1}{2}}(\Gamma) := [L_2(\Gamma), H_0^1(\Gamma)]_{\frac{1}{2},2}$ , for this *B* one may take the hypersingular integral operator, whereas for  $\partial\Gamma \neq \emptyset$ , and  $\mathscr{W} = H^{\frac{1}{2}}(\Gamma) = [L_2(\Gamma), H^1(\Gamma)]_{\frac{1}{2},2}$  the recently introduced modified hypersingular integral operator can be applied (see [HJHUT18]). (Note that  $H_0^1(\Gamma) = H^1(\Gamma)$  when  $\partial\Gamma = \emptyset$ .)

For a family of quasi-uniform partitions generated by a repeated application of *uniform refinements* starting from some given initial partition, a computationally attractive alternative choice for  $B_{T}^{\mathscr{S}}$  is provided by the multi-level operator from [BPV00], whose application can be performed in *linear complexity*. In Chapter 3, such operators will be discussed that also apply on locally refined meshes.

For  $\mathscr{W} = H^1_{0,\gamma}(\Omega)$ , i.e., when A is an operator of order -2 (cf. [FH19]), one obviously takes  $(B^{\mathscr{S}}_{\mathcal{T}}u)(v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx$ , or  $(B^{\mathscr{S}}_{\mathcal{T}}u)(v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx + \int_{\Omega} uv \, dx$  when meas $(\gamma) = 0$ , whose application can be performed in linear complexity.

# 2.3.4 Implementation

For both the domain case and the construction in the manifold case in §2.3.2, the matrix representation  $G_{\mathcal{T}} = \mathcal{F}_{\Xi_{\mathcal{T}}}^{-1} G_{\mathcal{T}} (\mathcal{F}_{\Xi_{\mathcal{T}}}')^{-1}$  of our preconditioner  $G_{\mathcal{T}}$ reads as  $G_{\mathcal{T}} = D_{\mathcal{T}}^{-1} B_{\mathcal{T}} D_{\mathcal{T}}^{-\top}$  with

$$\boldsymbol{D}_{\mathcal{T}} = \operatorname{diag}\{|T| \colon T \in \mathcal{T}\},\$$

and

$$\begin{aligned} \boldsymbol{B}_{\mathcal{T}} &\coloneqq \mathcal{F}'_{\Psi_{\mathcal{T}}} B_{\mathcal{T}} \mathcal{F}_{\Psi_{\mathcal{T}}} \\ &= \mathcal{F}'_{\Psi_{\mathcal{T}}} \left( (I_{\mathcal{T}}^{\mathscr{S}} |_{\mathscr{W}_{\mathcal{T}}})' B_{\mathcal{T}}^{\mathscr{S}} I_{\mathcal{T}}^{\mathscr{S}} |_{\mathscr{W}_{\mathcal{T}}} + (I_{\mathcal{T}}^{\mathscr{B}} |_{\mathscr{W}_{\mathcal{T}}})' B_{\mathcal{T}}^{\mathscr{B}} I_{\mathcal{T}}^{\mathscr{B}} |_{\mathscr{W}_{\mathcal{T}}} \right) \mathcal{F}_{\Psi_{\mathcal{T}}} \\ &= \boldsymbol{p}_{\mathcal{T}}^{\top} \boldsymbol{B}_{\mathcal{T}}^{\mathscr{S}} \boldsymbol{p}_{\mathcal{T}} + \boldsymbol{q}_{\mathcal{T}}^{\top} \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} \boldsymbol{q}_{\mathcal{T}}, \end{aligned}$$

where

$$\begin{split} & \boldsymbol{B}_{\mathcal{T}}^{\mathscr{S}} \coloneqq \mathcal{F}_{\Phi_{\mathcal{T}}} B_{\mathcal{T}}^{\mathscr{S}} \mathcal{F}_{\Phi_{\mathcal{T}}}, \quad \boldsymbol{p}_{\mathcal{T}} \coloneqq \mathcal{F}_{\Phi_{\mathcal{T}}}^{-1} I_{\mathcal{T}}^{\mathscr{S}} |_{\mathscr{W}_{\mathcal{T}}} \mathcal{F}_{\Psi_{\mathcal{T}}}, \\ & \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} \coloneqq \mathcal{F}_{\Theta_{\mathcal{T}}} B_{\mathcal{T}}^{\mathscr{B}} \mathcal{F}_{\Theta_{\mathcal{T}}}, \quad \boldsymbol{q}_{\mathcal{T}} \coloneqq \mathcal{F}_{\Theta_{\mathcal{T}}}^{-1} I_{\mathcal{T}}^{\mathscr{B}} |_{\mathscr{W}_{\mathcal{T}}} \mathcal{F}_{\Psi_{\mathcal{T}}}. \end{split}$$

By substituting the definition of  $B_{\mathcal{T}}^{\mathscr{B}}$  from (2.28), the definition of the basis  $\Psi_{\mathcal{T}} = \{\psi_{\mathcal{T},T}\}_{T \in \mathcal{T}}$  for  $\mathscr{W}_{\mathcal{T}}$  from (2.11) and (2.23), and that of the bases  $\Phi_{\mathcal{T}} = \{\phi_{\mathcal{T},\nu}\}_{\nu \in N_{\mathcal{T}}^{0}}$  and  $\Theta_{\mathcal{T}} = \{\theta_{T}\}_{T \in \mathcal{T}}$  for  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  and  $\mathscr{B}_{\mathcal{T}}$ , respectively, we find that

$$\boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} = \beta \boldsymbol{D}_{\mathcal{T}}^{1-\frac{2s}{d}}, \quad (\boldsymbol{p}_{\mathcal{T}})_{\nu T} = \begin{cases} d_{\mathcal{T},\nu}^{-1} & \text{if } \nu \in N_{\mathcal{T},T}^{0}, \\ 0 & \text{if } \nu \notin N_{\mathcal{T},T}^{0}, \end{cases}$$
$$(\boldsymbol{q}_{\mathcal{T}})_{T'T} = \delta_{T'T} - \frac{1}{d+1} \sum_{\nu \in N_{\mathcal{T},T}^{0} \cap N_{\mathcal{T},T'}^{0}} d_{\mathcal{T},\nu}^{-1}, \end{cases}$$

whereas  $B_{\mathcal{T}}^{\mathscr{S}}$  depends on  $B_{\mathcal{T}}^{\mathscr{S}} \in \mathcal{L}$ is<sub>*c*</sub> $(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  being chosen. The cost of the application of  $G_{\mathcal{T}}$  is the cost of the application of  $B_{\mathcal{T}}^{\mathscr{S}}$  plus cost that scales linearly in  $\#\mathcal{T}$ .

# 2.4 Continuous piecewise linear discretization space

Let a bounded polytopal domain  $\Omega \subset \mathbb{R}^d$ ,  $\gamma \subset \partial \Omega$ ,  $s \in [0, 1]$ ,  $\mathscr{W} := [L_2(\Omega), H_{0,\gamma}^1(\Omega)]_{s,2}$ ,  $\mathscr{V} := \mathscr{W}', D \in \mathcal{L}$ is $(\mathscr{V}, \mathscr{W}'), (\mathcal{T})_{\mathcal{T} \in \mathbb{T}}, N^0_{\mathcal{T}}, d_{\mathcal{T},\nu}, N_T$ , and  $N^0_{\mathcal{T},T}$  be all as in Sect. 2.3. In addition, for  $\mathcal{T} \in \mathbb{T}$  let  $N_{\mathcal{T}}$  be the set of *all* vertices of  $\mathcal{T}$ , so including those on a possibly non-empty  $\gamma$ , and for  $\nu \in N_{\mathcal{T}}$  let  $\omega_{\mathcal{T}}(\nu) := \cup_{\{T \in \mathcal{T} : \nu \in N_T\}} \overline{T}$ .

We take

$$\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1} := \{ u \in H^1(\Omega) \colon u |_T \in \mathcal{P}_1 \ (T \in \mathcal{T}) \} \subset \mathscr{V},$$

and, as in Sect. 2.3,

$$\mathscr{S}_{\mathcal{T},0}^{0,1} := \{ u \in H^1_{0,\gamma}(\Omega) \colon u|_T \in \mathcal{P}_1 \ (T \in \mathcal{T}) \},\$$

equipped with nodal bases  $\Xi_{\mathcal{T}} = \{\xi_{\mathcal{T},\nu} : \nu \in N_{\mathcal{T}}\}$  and  $\Phi_{\mathcal{T}} = \{\phi_{\mathcal{T},\nu} : \nu \in N_{\mathcal{T}}^0\}$ , respectively, defined by

$$\xi_{\mathcal{T},\nu}(\nu') = \delta_{\nu,\nu'} \quad (\nu,\nu' \in N_{\mathcal{T}}),$$

and  $\phi_{\mathcal{T},\nu} = \xi_{\mathcal{T},\nu}$  for  $\nu \in N^0_{\mathcal{T}}$ .

Analogously to the case of *discontinuous* piecewise constant trial spaces in  $\mathscr{V}$  studied in Sect. 2.3, using the framework of operator preconditioning outlined in Sect. 2.2 we are going to construct a family of preconditioners  $G_{\mathcal{T}} \in$  $\mathcal{L}$ is<sub>c</sub> $(\mathscr{V}'_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}})$  of type  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D'_{\mathcal{T}})^{-1}$  with uniformly bounded  $||G_{\mathcal{T}}||_{\mathcal{L}}(\mathscr{V}'_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}})$ and  $||\Re(G_{\mathcal{T}})^{-1}||_{\mathcal{L}}(\mathscr{V}_{\mathcal{T}}, \mathscr{V}'_{\mathcal{T}})$ .

The roles played in Sect. 2.3 by |T| (= | supp  $\xi_T$ |) and  $h_T = |T|^{1/d}$ , are in this section going to be played by  $|\omega_T(\nu)|$  (= | supp  $\xi_{T,\nu}|$ ) and  $h_{T,\nu} := |\omega_T(\nu)|^{1/d}$ .

#### **2.4.1** Construction of $\mathscr{W}_{\mathcal{T}}$ and $D_{\mathcal{T}}$

To construct a collection  $\Psi_{\mathcal{T}} = \{\psi_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}\} \subset H^1_{0,\gamma}(\Omega)$  that is biorthogonal to  $\Xi_{\mathcal{T}}$ , consists of locally supported functions, and for which

$$\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}} \subset \mathscr{W}$$

has an 'approximation property', as in Sect. 2.3 we need two collections  $\Sigma_{\mathcal{T}} \subset \mathscr{S}_{\mathcal{T},0}^{0,1}$  and  $\Theta_{\mathcal{T}} \subset H^1_{0,\gamma}(\Omega)$  of locally supported functions with  $\#\Sigma_{\mathcal{T}} = \#\Theta_{\mathcal{T}} = \#\Xi_{\mathcal{T}}$ , where  $\Theta_{\mathcal{T}}$  is biorthogonal to  $\Xi_{\mathcal{T}}$ , and  $\Sigma_{\mathcal{T}}$  has an 'approximation property'.

We define  $\Sigma_{\mathcal{T}} = \{\sigma_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}\}$  by  $\sigma_{\mathcal{T},\nu} \coloneqq \phi_{\mathcal{T},\nu}$  when  $\nu \in N_{\mathcal{T}}^0$ , and  $\sigma_{\mathcal{T},\nu} \coloneqq 0$  when  $\nu \in N_{\mathcal{T}} \setminus N_{\mathcal{T}}^0$ . Then, obviously,  $\sum_{\nu \in N_{\mathcal{T}}} \sigma_{\mathcal{T},\nu}$  equals 1 on  $\Omega \setminus \bigcup_{\{T \in \mathcal{T} : \overline{T} \cap \gamma \neq \emptyset\}} \overline{T}$ .

For constructing  $\Theta_{\mathcal{T}}$ , on a reference *d*-simplex  $\hat{T}$  for  $\varepsilon > 0$  we consider a smooth  $\eta_{\varepsilon} \in [0, 1]$ , symmetric in the barycentric coordinates, with  $\eta_{\varepsilon}(x) = 0$  when  $d(x, \partial \hat{T}) < \varepsilon$ , and  $\eta_{\varepsilon}(x) = 1$  when  $d(x, \partial \hat{T}) > 2\varepsilon$ . Then for some fixed  $\varepsilon > 0$  small enough, it holds that

$$\inf_{0\neq p\in\mathcal{P}_1(\hat{T})}\sup_{0\neq q\in\mathcal{P}_1(\hat{T})}\frac{\langle p,\eta_{\varepsilon}q\rangle_{L_2(\hat{T})}}{\|p\|_{L_2(\hat{T})}\|\eta_{\varepsilon}q\|_{L_2(\hat{T})}}>0,$$

meaning that the biorthogonal projector  $P_{\varepsilon} \in \mathcal{L}(L_2(\hat{T}), L_2(\hat{T}))$  with ran  $P_{\varepsilon} = \eta_{\varepsilon} \mathcal{P}_1(\hat{T})$  and ran $(\mathrm{Id} - P_{\varepsilon}) = \mathcal{P}_1(\hat{T})^{\perp_{L_2(\hat{T})}}$  exists. Consequently, with  $\Phi_{\hat{T}} = \{\phi_{\hat{T},\nu}: \nu \in N_{\hat{T}}\}$  being the nodal basis for  $\mathcal{P}_1(\hat{T})$ , we have that

$$\{\tilde{\phi}_{\hat{T},\varepsilon,\nu}\colon\nu\in N_{\hat{T}}\}\coloneqq\langle\Phi_{\hat{T}},\Phi_{\hat{T}}\rangle_{L_{2}(\hat{T})}^{-1}P_{\varepsilon}\Phi_{\hat{T}}\subset H^{1}_{0}(\hat{T})$$

is  $L_2(\hat{T})$ -biorthonormal to  $\{\phi_{\hat{T},\nu} \colon \nu \in N_{\hat{T}}\}$ .

Now for  $T \in \mathcal{T}$ , let  $F_{\hat{T},T} : T \to \hat{T}$  be an affine bijection. Then  $\{\tilde{\phi}_{T,\varepsilon,\nu} : \nu \in N_T\}$  defined by

(2.29) 
$$\tilde{\phi}_{T,\varepsilon,\nu} := \frac{|\hat{T}|}{|T|} \tilde{\phi}_{\hat{T},\varepsilon,F_{\hat{T},T}(\nu)}$$

is  $L_2(T)$ -biorthonormal to the nodal basis for  $P_1(T)$ .

By selecting for  $\nu \in N_{\mathcal{T}}$ , a  $T(\nu) \in \mathcal{T}$  with  $\nu \in N_T$ , and defining  $\Theta_{\mathcal{T}} = \{\theta_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}\} \subset H^1_{0,\gamma}(\Omega)$  by

$$\theta_{\mathcal{T},\nu} := |\omega_{\mathcal{T}}(\nu)| \hat{\phi}_{T(\nu),\varepsilon,\nu},$$

where the specific scaling is chosen for convenience, we have for  $\nu, \nu' \in N_T$ ,

(2.30) 
$$\begin{aligned} \delta_{\nu\nu'} |\omega_{\mathcal{T}}(\nu)| &= \langle \theta_{\mathcal{T},\nu}, \xi_{\mathcal{T},\nu'} \rangle_{L_2(\Omega)} = \delta_{\nu\nu'} \|\theta_{\mathcal{T},\nu}\|_{L_2(\Omega)} \|\xi_{\mathcal{T},\nu'}\|_{L_2(\Omega)}, \\ \sup \theta_{\mathcal{T},\nu} \subset \overline{T(\nu)}, \ |\theta_{\mathcal{T},\nu}|_{H^1(\Omega)} \lesssim h_{\mathcal{T},\nu}^{-1} \|\theta_{\mathcal{T},\nu}\|_{L_2(\Omega)}, \end{aligned}$$

35

i.e., properties analogous to (2.6), (2.12)(c), and (2.14).

Since furthermore  $\Sigma_T$  and  $\Xi_T$  satisfy properties analogous to (2.12)(a,b), (2.13) and (2.15), defining similarly to (2.8)

$$\begin{aligned} (2.31) \\ \psi_{\mathcal{T},\nu} &\coloneqq \sigma_{\mathcal{T},\nu} + \frac{\langle \mathbb{1} - \sigma_{\mathcal{T},\nu}, \xi_{\mathcal{T},\nu} \rangle_{L_2(\Omega)}}{\langle \theta_{\mathcal{T},\nu}, \xi_{\mathcal{T},\nu} \rangle_{L_2(\Omega)}} \theta_{\mathcal{T},\nu} - \sum_{\nu' \in N_{\mathcal{T}} \setminus \{\nu\}} \frac{\langle \sigma_{\mathcal{T},\nu}, \xi_{\mathcal{T},\nu'} \rangle_{L_2(\Omega)}}{\langle \theta_{\mathcal{T},\nu'}, \xi_{\mathcal{T},\nu'} \rangle_{L_2(\Omega)}} \theta_{\mathcal{T},\nu'} \\ &= \begin{cases} \frac{\theta_{\mathcal{T},\nu}}{d+1} & \nu \in N_{\mathcal{T}} \setminus N_{\mathcal{T}}^0, \\ \phi_{\mathcal{T},\nu} + \frac{d}{(d+2)(d+1)} \theta_{\mathcal{T},\nu} - \sum_{\nu' \in N_{\mathcal{T}} \setminus \{\nu\}} \frac{|\omega_{\mathcal{T}}(\nu) \cap \omega_{\mathcal{T}}(\nu')|}{(d+2)(d+1)|\omega_{\mathcal{T}}(\nu')|} \theta_{\mathcal{T},\nu'} & \nu \in N_{\mathcal{T}}^0, \end{cases} \end{aligned}$$

we infer that  $\sum_{\nu \in N_T} \psi_{T,\nu}$  equals 1 possibly except on a strip along the Dirichlet boundary, and similarly to Theorem 2.3.2, that the biorthogonal projector

(2.32) 
$$P_{\mathcal{T}} \colon u \mapsto \sum_{\nu \in N_{\mathcal{T}}} \frac{\langle u, \xi_{\mathcal{T},\nu} \rangle_{L_2(\Omega)}}{\langle \mathbb{1}, \xi_{\mathcal{T},\nu} \rangle_{L_2(\Omega)}} \psi_{\mathcal{T},\nu},$$

satisfies  $\sup_{\mathcal{T}\in\mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W},\mathscr{W})} < \infty$ . With  $(D_{\mathcal{T}}v)(w) := (Dv)(w)$   $((v,w) \in \mathscr{V}_{\mathcal{T}} \times \mathscr{W}_{\mathcal{T}})$ , we have  $\|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}})} \leq 1$  and  $\sup_{\mathcal{T}\in\mathbb{T}} \|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}',\mathscr{V}_{\mathcal{T}})} < \infty$ , and

$$\boldsymbol{D}_{\mathcal{T}} = \mathcal{F}'_{\Psi_{\mathcal{T}}} D_{\mathcal{T}} \mathcal{F}_{\Xi_{\mathcal{T}}} = \operatorname{diag}\{\langle \mathbb{1}, \xi_{\mathcal{T}, \nu} \rangle_{L_2(\Omega)} \colon \nu \in N_{\mathcal{T}}\} = \operatorname{diag}\left\{\frac{1}{d+1} |\omega_{\mathcal{T}}(\nu)| \colon \nu \in N_{\mathcal{T}}\right\}$$

# **2.4.2** Construction of $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$ .

Since  $\Theta_{\mathcal{T}}$  additionally satisfies, for  $k \in \{0, 1\}$ ,

$$\Big\|\sum_{\nu\in N_{\mathcal{T}}}c_{\nu}\theta_{\mathcal{T},\nu}\Big\|_{H^{k}(\Omega)}^{2} \approx \sum_{\nu\in N_{\mathcal{T}}}h_{\mathcal{T},\nu}^{-k}\|\theta_{\mathcal{T},\nu}\|_{L_{2}(\Omega)}^{2}|c_{\nu}|^{2},$$

where  $\|\theta_{\mathcal{T},\nu}\|_{L_2(\Omega)} \stackrel{(2.30)}{=} |\omega_{\mathcal{T}}(\nu)| \|\xi_{\mathcal{T},\nu}\|_{L_2(\Omega)}^{-1} \approx |\omega_{\mathcal{T}}(\nu)|^{\frac{1}{2}}$ , and

$$||u+v||^2_{H^k(\Omega)} \gtrsim ||u||^2_{H^k(\Omega)} + ||v||^2_{H^k(\Omega)} \quad (u \in \mathscr{S}^{0,1}_{\mathcal{T},0}, v \in \mathscr{B}_{\mathcal{T}} := \operatorname{span} \Theta_{\mathcal{T}}).$$

(cf. (2.26)-(2.27)), we construct  $B_{\mathcal{T}}$  analogously as in §2.3.3: Assuming that we have a  $B_{\mathcal{T}}^{\mathscr{S}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  available with  $\sup_{\mathcal{T}\in\mathbb{T}} \|B_{\mathcal{T}}^{\mathscr{S}}\|_{\mathcal{L}(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')} < \infty$  and  $\sup_{\mathcal{T}\in\mathbb{T}} \|\Re(B_{\mathcal{T}}^{\mathscr{S}})^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},0}^{0,1})', \mathscr{S}_{\mathcal{T},0}^{0,1})} < \infty$ , for some constant  $\beta > 0$  we take

$$\left(B_{\mathcal{T}}^{\mathscr{B}}\sum_{\nu\in N_{\mathcal{T}}}c_{\nu}\theta_{\mathcal{T},\nu}\right)\left(\sum_{\nu\in N_{\mathcal{T}}}d_{\nu}\theta_{\mathcal{T},\nu}\right) := \beta\sum_{\nu\in N_{\mathcal{T}}}|\omega_{\mathcal{T}}(\nu)|^{1-\frac{2s}{d}}c_{\nu}d_{\nu},$$

and

$$B_{\mathcal{T}} := (I_{\mathcal{T}}^{\mathscr{S}}|_{\mathscr{W}_{\mathcal{T}}})' B_{\mathcal{T}}^{\mathscr{S}} I_{\mathcal{T}}^{\mathscr{S}}|_{\mathscr{W}_{\mathcal{T}}} + (I_{\mathcal{T}}^{\mathscr{B}}|_{\mathscr{W}_{\mathcal{T}}})' B_{\mathcal{T}}^{\mathscr{B}} I_{\mathcal{T}}^{\mathscr{B}}|_{\mathscr{W}_{\mathcal{T}}},$$

36

where  $I_{\mathcal{T}}^{\mathscr{S}}$  is the linear projector defined on  $\mathscr{S}_{\mathcal{T}_0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}}$  by ran  $I_{\mathcal{T}}^{\mathscr{S}} = \mathscr{S}_{\mathcal{T}_0}^{0,1}$  and  $\operatorname{ran} I_{\mathcal{T}}^{\mathscr{B}} = \mathscr{B}_{\mathcal{T}}$ , where  $I_{\mathcal{T}}^{\mathscr{B}} := \operatorname{Id} - I_{\mathcal{T}}^{\mathscr{G}}$ . Then one has  $\sup_{\mathcal{T} \in \mathbb{T}} \|B_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}})} < \infty$ and  $\sup_{\mathcal{T}\in\mathbb{T}} \|\Re(B_{\mathcal{T}})^{-1}\|_{\mathcal{L}(\mathscr{W}'_{\mathcal{T}},\mathscr{W}_{\mathcal{T}})} < \infty.$ 

Substituting the definition of  $\psi_{\mathcal{T},\nu}$  one infers that  $G_{\mathcal{T}} = D_{\mathcal{T}}^{-1} B_{\mathcal{T}} D_{\mathcal{T}}^{-\top}$ , where

$$oldsymbol{B}_{\mathcal{T}} = oldsymbol{p}_{\mathcal{T}}^{ op}oldsymbol{B}_{\mathcal{T}}^{ op}oldsymbol{p}_{\mathcal{T}} + oldsymbol{q}_{\mathcal{T}}^{ op}oldsymbol{B}_{\mathcal{T}}^{\mathcal{B}}oldsymbol{q}_{\mathcal{T}},$$

and

$$(\boldsymbol{q}_{\mathcal{T}})_{\nu'\nu} \coloneqq \begin{cases} \frac{\delta_{\nu'\nu}}{d+1} & \nu \in N_{\mathcal{T}} \setminus N_{\mathcal{T}}^{0}, \\ \frac{d}{(d+2)(d+1)} & \nu \in N_{\mathcal{T}}^{0}, \nu' = \nu, \\ -\frac{|\omega_{\mathcal{T}}(\nu) \cap \omega_{\mathcal{T}}(\nu')|}{(d+2)(d+1)|\omega_{\mathcal{T}}(\nu')|} & \nu \in N_{\mathcal{T}}^{0}, \nu' \neq \nu, \end{cases}$$
$$(\boldsymbol{p}_{\mathcal{T}})_{\nu'\nu} \coloneqq \delta_{\nu'\nu} \left(\nu' \in N_{\mathcal{T}}^{0}, \nu \in N_{\mathcal{T}}\right), \qquad \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} \coloneqq \operatorname{diag}\{\beta | \omega_{\mathcal{T}}(\nu)|^{1-\frac{2s}{d}} \colon \nu \in N_{\mathcal{T}}\}$$

#### 2.4.3 Manifold case.

From Sect. 2.3.2 recall the definitions of  $\Gamma$ ,  $\gamma$ ,  $\mathscr{W}$ ,  $\mathscr{V}$ ,  $\chi \colon \Omega \to \bigcup_{i=1}^{p} \chi_i(\Omega_i)$ , and that of the family of conforming partitions  $\mathbb{T}$  of  $\Gamma$ .

As in the domain case discussed in Sect. 2.4.1, for  $\mathcal{T} \in \mathbb{T}$  let  $N_{\mathcal{T}}$  be the set of vertices of  $\mathcal{T}$ , and  $N^0_{\mathcal{T}}$  its subset of vertices not on  $\gamma$ , for  $T \in \mathcal{T}$  let  $N_T$  be the vertices of *T*,  $N_{\mathcal{T},T}^0 := N_{\mathcal{T}}^0 \cap N_T$ , and for  $\nu \in N_{\mathcal{T}}$  let  $\omega_{\mathcal{T}}(\nu) := \bigcup_{\{T \in \mathcal{T}: \nu \in N_T\}} \overline{T}$ . We take

$$\begin{aligned} \mathscr{V}_{\mathcal{T}} &= \mathscr{S}_{\mathcal{T}}^{0,1} := \{ u \in H^1(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_1 \left( T \in \mathcal{T} \right) \} \subset \mathscr{V}, \\ \mathscr{S}_{\mathcal{T},0}^{0,1} &:= \{ u \in H^1_{0,\gamma}(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_1 \left( T \in \mathcal{T} \right) \}, \end{aligned}$$

equipped with nodal bases  $\Xi_{\mathcal{T}} = \{\xi_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}\}$  and  $\Phi_{\mathcal{T}} = \{\phi_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}^0\},\$ respectively, defined by

$$\xi_{\mathcal{T},\nu}(\nu') = \delta_{\nu,\nu'} \quad (\nu,\nu' \in N_{\mathcal{T}}),$$

and  $\phi_{\mathcal{T},\nu} = \xi_{\mathcal{T},\nu}$  for  $\nu \in N^0_{\mathcal{T}}$ .

Actually exclusively for the deriving an inverse inequality analogous to Lemma 2.3.4, first we construct a collection  $\Psi_{\mathcal{T}} = \{\psi_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}\} \subset H^1_{0,\gamma}(\Gamma)$ that has an 'approximation property' and that is biorthogonal to  $\Xi_T$  w.r.t. the *true*  $L_2(\Gamma)$ -scalar product: We define  $\Sigma_{\mathcal{T}} = \{\sigma_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}\}$  by  $\sigma_{\mathcal{T},\nu} \coloneqq \phi_{\mathcal{T},\nu}$ when  $\nu \in N^0_{\mathcal{T}}$ , and  $\sigma_{\mathcal{T},\nu} := 0$  when  $\nu \in N_{\mathcal{T}} \setminus N^0_{\mathcal{T}}$ . Then, obviously,  $\sum_{\nu \in N_{\mathcal{T}}} \sigma_{\mathcal{T},\nu}$ equals 1 on  $\Gamma \setminus \bigcup_{\{T \in \mathcal{T} : \overline{T} \cap \gamma \neq \emptyset\}} \overline{T}$ .

Given a *d*-simplex  $T \subset \mathbb{R}^d$ , by means of an affine bijection we transport the function  $\eta_{\varepsilon}$ , defined in Sect. 2.3.2 on a reference *d*-simplex  $\hat{T}$ , to a function on *T* and denote it by  $\eta_{T,\varepsilon}$ . Then for any panel  $T \in \mathcal{T} \in \mathbb{T}$ , for some  $\varepsilon > 0$  small enough it holds that

$$\inf_{0 \neq p \in \mathcal{P}_1(\chi^{-1}(T))} \sup_{0 \neq q \in \mathcal{P}_1(\chi^{-1}(T))} \frac{\langle p \circ \chi^{-1}, (\eta_{T,\varepsilon}q) \circ \chi^{-1} \rangle_{L_2(T)}}{\| p \circ \chi^{-1} \|_{L_2(T)} \| (\eta_{T,\varepsilon}q) \circ \chi^{-1} \|_{L_2(T)}} > 0$$

Moreover, since the panels T get increasingly flat when diam  $T \to 0$ , there exists an  $\varepsilon > 0$  such that above inf-sup condition is satisfied *uniformly* over all  $T \in \mathcal{T} \in \mathbb{T}$ .

By selecting for each  $\nu \in N_{\mathcal{T}}$  a  $T(\nu) \in \mathcal{T}$  with  $\nu \in N_T$ , as in Sect. 2.4.1 we obtain a collection  $\Theta_{\mathcal{T}} = \{\theta_{\mathcal{T},\nu} : \nu \in N_{\mathcal{T}}\}$  with  $\theta_{\mathcal{T},\nu} \subset H_0^1(T(\nu))$  that is biorthogonal to  $\Xi_{\mathcal{T}}$ , in particular that satisfies (2.30), after which we define the  $\psi_{\mathcal{T},\nu}$  by means of formula (2.31). Having constructed the biorthogonal collections  $\Xi_{\mathcal{T}}$  and  $\Psi_{\mathcal{T}}$ , we set the biorthogonal projector  $P_{\mathcal{T}} : L_2(\Gamma) \to H_{0,\gamma}^1(\Gamma) : u \mapsto$  $\sum_{\nu \in N_{\mathcal{T}}} \frac{\langle u, \xi_{\mathcal{T},\nu} \rangle_{L_2(\Gamma)}}{\langle 1, \xi_{\mathcal{T},\nu} \rangle_{L_2(\Gamma)}} \psi_{\mathcal{T},\nu}$  which satisfies  $\|P_{\mathcal{T}}u\|_{H^1(\Gamma)} \lesssim \|h_{\mathcal{T}}^{-1}u\|_{L_2(\Gamma)}$ . With the aid of this projector, as in Lemma 2.3.4 one infers that

$$(2.33) \|h_{\mathcal{T}}v_{\mathcal{T}}\|_{L_2(\Gamma)} \lesssim \|v_{\mathcal{T}}\|_{(H^1_{0,\infty})'} \quad (v_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}).$$

Having established this inverse inequality, to arrive at a construction of  $\Psi_{\mathcal{T}}$  that does not require the evaluation of integrals over  $\Gamma$ , as in Sect. 2.3.2 we replace  $\langle , \rangle_{L_2(\Gamma)}$  by  $\langle , \rangle_{\mathcal{T}}$ . We redefine  $\Theta_{\mathcal{T}} = \{\theta_{\mathcal{T},\nu} : \nu \in N_{\mathcal{T}}\}$  by

$$\theta_{\mathcal{T},\nu} := |\omega_{\mathcal{T}}(\nu)| \frac{|\chi^{-1}(T)|}{|T|} \tilde{\phi}_{\chi^{-1}(T),\varepsilon,\chi^{-1}(\nu)} \circ \chi^{-1}$$

with the  $\tilde{\phi}$ 's defined in (2.29), and following (2.31) set  $\Psi_{\mathcal{T}} = \{\psi_{\mathcal{T},\nu} : \nu \in N_{\mathcal{T}}\}$ and  $\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}}$  by

$$\psi_{\mathcal{T},\nu} = \begin{cases} \frac{\theta_{\mathcal{T},\nu}}{d+1} & \nu \in N_{\mathcal{T}} \setminus N_{\mathcal{T}}^{0}, \\ \phi_{\mathcal{T},\nu} + \frac{d}{(d+2)(d+1)} \theta_{\mathcal{T},\nu} - \sum_{\nu' \in N_{\mathcal{T}} \setminus \{\nu\}} \frac{|\omega_{\mathcal{T}}(\nu) \cap \omega_{\mathcal{T}}(\nu')|}{(d+2)(d+1)|\omega_{\mathcal{T}}(\nu')|} \theta_{\mathcal{T},\nu'} & \nu \in N_{\mathcal{T}}^{0}, \end{cases}$$

As in Sect. 2.3.2, we set  $(D_{\mathcal{T}}v_{\mathcal{T}})(w_{\mathcal{T}}) := \langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{\mathcal{T}} (v_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}, w_{\mathcal{T}} \in \mathscr{W}_{\mathcal{T}})$ , and as in Sect. 2.3.2, using (2.33) one shows that  $\sup_{\mathcal{T} \in \mathbb{T}} \|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')} < \infty$ . Similarly as in Lemma 2.3.6, one proves that

$$\|h_{\mathcal{T}} v_{\mathcal{T}}\|_{L_2(\Gamma)} \lesssim \sup_{0 \neq w \in H_{0,\gamma}^1(\Gamma)} \frac{\langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{\mathcal{T}}}{\|w\|_{H^1(\Gamma)}}$$

and with that  $\sup_{\mathcal{T}\in\mathbb{T}} \|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathcal{W}_{\mathcal{T}}',\mathcal{V}_{\mathcal{T}})} < \infty$ .

Constructing  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}})$  as in Sect. 2.4.2, one arrives at the same expressions for  $D_{\mathcal{T}}, G_{\mathcal{T}}, B_{\mathcal{T}}, q_{\mathcal{T}}, B_{\mathcal{T}}^{\mathscr{G}}, p_{\mathcal{T}}$ , and  $B_{\mathcal{T}}^{\mathscr{B}}$  as in Sect. 2.4.1-2.4.2 in the domain case.

# 2.5 Higher order case

In this section, we discuss the construction of an uniform preconditioner for  $\mathscr{V}_{\mathcal{T}}$  being either the space  $\mathscr{S}_{\mathcal{T}}^{-1,\ell}$  of *discontinuous* piecewise polynomials of degree  $\ell > 0$  w.r.t.  $\mathcal{T}$ , or the space  $\mathscr{S}_{\mathcal{T}}^{0,\ell}$  of *continuous* piecewise polynomials of degree  $\ell > 1$  w.r.t.  $\mathcal{T}$ .

We write the spaces  $\mathscr{V}_{\mathcal{T}}$ ,  $\mathscr{B}_{\mathcal{T}}$ ,  $\mathscr{W}_{\mathcal{T}}$ , and their bases  $\Xi_{\mathcal{T}}$ ,  $\Theta_{\mathcal{T}}$ ,  $\Psi_{\mathcal{T}}$  from Sect. 2.3 or 2.4 as  $\mathscr{V}_{\mathcal{T}}^0$ ,  $\mathscr{B}_{\mathcal{T}}^0$ ,  $\mathscr{W}_{\mathcal{T}}^0$ , and  $\Xi_{\mathcal{T}}^0$ ,  $\Theta_{\mathcal{T}}^0$ ,  $\Psi_{\mathcal{T}}^0$ , respectively. The biorthogonal projector formerly denoted as  $P_{\mathcal{T}}$  will now be denoted as  $P_{\mathcal{T}}^0$ , and the matrices  $\mathbf{B}_{\mathcal{T}}$  and  $\mathbf{D}_{\mathcal{T}}$  as  $\mathbf{B}_{\mathcal{T}}^0$  and  $\mathbf{D}_{\mathcal{T}}^0$ .

Although we consider the domain case, the results extend to the manifold case following the approach outlined in §2.3.2 or §2.4.3.

In order to construct an uniform preconditioner, obvious possibilities are to apply the framework of operator preconditioning directly to the higher order polynomial space  $\mathcal{V}_{\mathcal{T}}$ , or to use the preconditioner developed for the lowest order case within a subspace correction framework. We investigate both possibilities.

# 2.5.1 Application of the operator preconditioning framework

#### Discontinuous piecewise polynomials

Given  $\ell > 0$ , for  $\mathcal{T} \in \mathbb{T}$ , let  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,\ell}$ . With  $m = m(\ell) := \binom{d+\ell}{\ell} - 1$ , we equip  $\mathscr{V}_{\mathcal{T}}$  with  $\Xi_{\mathcal{T}} = \{\xi_{T,i} : T \in \mathcal{T}, 0 \le i \le m\}$ , where for each  $T \in \mathcal{T}$ ,  $\{\xi_{T,i} : 0 \le i \le m\}$  is constructed by the common affine lifting approach from a basis for the polynomials of degree  $\ell$  on a reference *d*-simplex, such that  $\{\xi_{T,0} : T \in \mathcal{T}\} = \Xi_{\mathcal{T}}^0$ ,  $\operatorname{supp} \xi_{T,i} \subset \overline{T}$ , and  $\|\xi_{T,i}\|_{L_2(\Omega)} = |T|^{\frac{1}{2}.6}$ 

A straightforward generalization of the construction in the first paragraphs of §2.4.1 of a collection in  $H_0^1(T)$  that is biorthogonal to the nodal basis of  $P_1(T)$ shows the following: There exists a set of 'bubbles'  $\Theta_T = \{\theta_{T,i} : T \in T, 0 \le i \le m\} \subset H_{0,\gamma}^1(\Omega)$  such that for  $T, T' \in T, 0 \le i, i' \le m, k \in \{0, 1\}$ ,

- $(2.34) \quad \delta_{(T,i),(T',i')}|T| = \langle \theta_{T,i}, \xi_{T',i'} \rangle_{L_2(\Omega)} = \delta_{(T,i),(T',i')} \|\theta_{T,i}\|_{L_2(\Omega)} \|\xi_{T',i'}\|_{L_2(\Omega)},$
- (2.35)  $\|\theta_{T,i}\|_{H^1(\Omega)} \lesssim h_T^{-1} \|\theta_{T,i}\|_{L_2(\Omega)},$
- (2.36)  $\sup \theta_{T,i} \subset \overline{T},$
- (2.37)  $\left\|\sum_{\{T\in\mathcal{T},0\leq i\leq m\}} c_{T,i}\theta_{T,i}\right\|_{H^{k}(\Omega)}^{2} \approx \sum_{\{T\in\mathcal{T},0\leq i\leq m\}} h_{T}^{-2k} |c_{T,i}|^{2} \|\theta_{T,i}\|_{L_{2}(\Omega)}^{2},$

(2.38)  $\|u+v\|_{H^{k}(\Omega)}^{2} \gtrsim \|u\|_{H^{k}(\Omega)}^{2} + \|v\|_{H^{k}(\Omega)}^{2} \quad (u \in \mathscr{S}_{\mathcal{T},0}^{0,1}, v \in \mathscr{B}_{\mathcal{T}} := \operatorname{span} \Theta_{\mathcal{T}}),$ 

$$(2.39) \quad \{\theta_{T,0} \colon T \in \mathcal{T}\} = \Theta^0_{\mathcal{T}}.$$

Writing  $\Psi^0_{\mathcal{T}} = \{\psi^0_{\mathcal{T},T} : T \in \mathcal{T}\}$ , we define  $\Psi_{\mathcal{T}} := \{\psi_{\mathcal{T},T,i} : T \in \mathcal{T}, 0 \leq i \leq m\}$ ,  $\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}}$  by

$$\psi_{\mathcal{T},T,i} := \begin{cases} \psi_{\mathcal{T},T}^0 - \sum_{\{T' \in \mathcal{T}, 1 \le i' \le m\}} \frac{\langle \psi_{\mathcal{T},T}^0, \xi_{T',i'} \rangle_{L_2(\Omega)}}{\langle \theta_{T',i'}, \xi_{T',i'} \rangle_{L_2(\Omega)}} \theta_{T',i'} & i = 0, \\ \theta_{T,i} & 1 \le i \le m. \end{cases}$$

<sup>6</sup>For i > 1, it is allowed that  $\xi_{T,i}$  is (nearly) orthogonal to  $\mathbb{1}|_T$ , i.e., (2.15) is not required for these  $\xi_{T,i}$ .

Knowing that  $\Psi^0_{\mathcal{T}}$  and  $\Xi^0_{\mathcal{T}}$ , and  $\Theta_{\mathcal{T}}$  and  $\Xi_{\mathcal{T}}$  are biorthogonal, the correction made to the  $\psi^0_{\mathcal{T},T}$  ensures that  $\Psi_{\mathcal{T}}$  and  $\Xi_{\mathcal{T}}$  are biorthogonal, in particular that

$$\langle \psi_{\mathcal{T},T,i}, \xi_{T',i'} \rangle = \delta_{(T,i),(T',i')} |T|$$

For use later, notice that  $\mathscr{W}_{\mathcal{T}}^0 \subset \mathscr{W}_{\mathcal{T}}$ , and that by definition of  $\psi_{\mathcal{T},T}^0$  and  $\sigma_{\mathcal{T},T}$ , for i' > 0,

(2.40) 
$$(\mathbf{R}_{\mathcal{T}})_{(T',i'),T} := -\frac{\langle \psi_{\mathcal{T},T}^{0}, \xi_{T',i'} \rangle_{L_{2}(\Omega)}}{\langle \theta_{T',i'}, \xi_{T',i'} \rangle_{L_{2}(\Omega)}} = -|T'|^{-1} \langle \sigma_{\mathcal{T},T}, \xi_{T',i'} \rangle_{L_{2}(\Omega)} \\ = -|T'|^{-1} \sum_{\nu \in N_{\mathcal{T},T}^{0} \cap N_{\mathcal{T},T'}^{0}} d_{\mathcal{T},\nu}^{-1} \langle \phi_{\mathcal{T},\nu}, \xi_{T',i'} \rangle_{L_{2}(\Omega)}.$$

The biorthogonal projector  $P_{\mathcal{T}} \colon L_2(\Omega) \to H^1_{0,\gamma}(\Omega)$  with ran  $P_{\mathcal{T}} = \mathscr{W}_{\mathcal{T}}$  and ran(Id  $-P_{\mathcal{T}}) = \mathscr{V}_{\mathcal{T}}^{\perp_{L_2(\Omega)}}$  is given by

$$P_{\mathcal{T}}u = \sum_{\{T \in \mathcal{T}, 0 \le i \le m\}} \frac{\langle u, \xi_{T,i} \rangle_{L_2(\Omega)}}{\langle \psi_{\mathcal{T},T,i}, \xi_{T,i} \rangle_{L_2(\Omega)}} \psi_{\mathcal{T},T,i}.$$

**Theorem 2.5.1.** It holds that  $\sup_{\mathcal{T}\in\mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W},\mathscr{W})} < \infty$ .

*Proof.* For  $T'' \in \mathcal{T}$ ,  $k \in \{0, 1\}$ , from  $\operatorname{supp} \xi_{T,i} \subset \overline{T}$  and  $\operatorname{supp} \psi_{\mathcal{T},T'',i} \subset \omega_{\mathcal{T}}^{(1)}(T'')$  we have

$$(2.41) ||P_{\mathcal{T}}u||_{H^{k}(T'')} \leq \sum_{\{T \in R^{(1)}_{\mathcal{T}}(T''), 0 \leq i \leq m\}} ||u||_{L_{2}(T)} \frac{||\psi_{\mathcal{T},T,i}||_{H^{k}(\Omega)} ||\xi_{T,i}||_{L_{2}(\Omega)}}{|\langle\psi_{\mathcal{T},T,i},\xi_{T,i}\rangle_{L_{2}(\Omega)}|}$$

To bound the right-hand side we distinguish between terms with i = 0 and those with i > 0.

From  $\|\psi_{\mathcal{T},T}^0\|_{H^k(\Omega)} \lesssim h_T^{d/2-k}$  ((2.17)), and for  $T' \in \mathcal{T}, 1 \le i' \le m$ ,

$$\left\| \frac{\langle \psi_{\mathcal{T},T}^{0},\xi_{T',i'}\rangle_{L_{2}(\Omega)}}{\langle \theta_{T',i'},\xi_{T',i'}\rangle_{L_{2}(\Omega)}} \theta_{T',i'} \right\|_{H^{k}(\Omega)} \lesssim \frac{\|\psi_{\mathcal{T},T}^{0}\|_{L_{2}(\Omega)} \|\xi_{T',i'}\|_{L_{2}(\Omega)}}{\|\theta_{T',i'}\|_{L_{2}(\Omega)} \|\xi_{T',i'}\|_{L_{2}(\Omega)}} h_{T'}^{-k} \|\theta_{T',i'}\|_{L_{2}(\Omega)} \lesssim h_{T}^{d/2-k},$$

we infer that

$$\|\psi_{\mathcal{T},T,0}\|_{H^k(\Omega)} \lesssim h_T^{d/2-k},$$

while, thanks to (2.15),

$$\frac{\|\xi_{T,0}\|_{L_2(\Omega)}}{|\langle\psi_{\mathcal{T},T,0},\xi_{T,0}\rangle_{L_2(\Omega)}|} = \frac{\|\xi_{T,0}\|_{L_2(\Omega)}}{|\langle\psi_{\mathcal{T},T}^0,\xi_{T,0}\rangle_{L_2(\Omega)}|} = \frac{\|\xi_{T,0}\|_{L_2(\Omega)}}{|\langle\mathbb{1},\xi_{T,0}\rangle_{L_2(\Omega)}|} \approx h_T^{-d/2}.$$

For  $1 \leq i \leq m$ ,

$$\frac{\|\psi_{\mathcal{T},T,i}\|_{H^{k}(\Omega)}\|\xi_{T,i}\|_{L_{2}(\Omega)}}{|\langle\psi_{\mathcal{T},T,i},\xi_{T,i}\rangle_{L_{2}(\Omega)}|} = \frac{\|\theta_{T,i}\|_{H^{k}(\Omega)}\|\xi_{T,i}\|_{L_{2}(\Omega)}}{|\langle\theta_{T,i},\xi_{T,i}\rangle_{L_{2}(\Omega)}|} \approx \frac{\|\theta_{T,i}\|_{H^{k}(\Omega)}}{\|\theta_{T,i}\|_{L_{2}(\Omega)}} \lesssim h_{T}^{-k}.$$

Combining the above inequality with (2.41) shows that

(2.42) 
$$\|P_{\mathcal{T}}u\|_{H^{k}(T'')} \lesssim h_{T}^{-k} \|u\|_{L_{2}(\omega_{\mathcal{T}}^{(1)}(T''))},$$

which is the analogue of estimate (2.18) that was proven for  $P_{\mathcal{T}}^0$ . Since  $\mathscr{W}_{\mathcal{T}}^0 \subset \mathscr{W}_{\mathcal{T}}$ , by making use of the *same* Scott-Zhang type quasi-interpolator  $\Pi_{\mathcal{T}}$  as has been used in the proof of Theorem 2.3.2, copying the remainder of that proof yields the claim.

Thanks this theorem,  $D_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$  defined by  $(D_{\mathcal{T}}v)(w) = \langle v, w \rangle_{L_2(\Omega)}$ satisfies  $\sup_{\mathcal{T} \in \mathbb{T}} \max \left( \|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})}, \|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathscr{W}'_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}})} \right) < \infty$ . By enumerating all  $\xi_{T,i}$  and  $\psi_{\mathcal{T},T,i}$  with index i = 0 before those with index i > 0, its matrix representation reads as  $\mathbf{D}_{\mathcal{T}} = \begin{bmatrix} \mathbf{D}_{\mathcal{T}}^0 & 0\\ 0 & \mathbf{D}_{\mathcal{T}}^1 \end{bmatrix}$ , where  $\mathbf{D}_{\mathcal{T}}^1 := \operatorname{diag}\{|T| \operatorname{Id}_{m \times m} \colon T \in \mathcal{T}\}$ .

Thanks to (2.37)-(2.38), a suitable  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathcal{W}_{\mathcal{T}}, \mathcal{W}_{\mathcal{T}}')$  can be defined similarly as in §2.3.3: With  $I_{\mathcal{T}}^{\mathscr{S}}$  being the linear projector defined on  $\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}}$ by ran  $I_{\mathcal{T}}^{\mathscr{S}} = \mathscr{S}_{\mathcal{T},0}^{0,1}$  and ran  $I_{\mathcal{T}}^{\mathscr{B}} = \mathscr{B}_{\mathcal{T}}$ , where  $I_{\mathcal{T}}^{\mathscr{B}} := \mathrm{Id} - I_{\mathcal{T}}^{\mathscr{S}}$ , we define  $B_{\mathcal{T}}^{\mathscr{S} \oplus \mathscr{B}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}}, (\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}})')$  by

$$(B_{\mathcal{T}}^{\mathscr{S}\oplus\mathscr{B}}w)(\tilde{w}) = (B_{\mathcal{T}}^{\mathscr{S}}I_{\mathcal{T}}^{\mathscr{S}}w)(I_{\mathcal{T}}^{\mathscr{S}}\tilde{w}) + (B_{\mathcal{T}}^{\mathscr{B}}I_{\mathcal{T}}^{\mathscr{B}}w)(I_{\mathcal{T}}^{\mathscr{B}}\tilde{w}),$$

where

$$(B_{\mathcal{T}}^{\mathscr{B}}\sum_{\{T\in\mathcal{T},\,0\leq i\leq m\}}c_{T,i}\theta_{\mathcal{T},i})(\sum_{\{T\in\mathcal{T},\,0\leq i\leq m\}}d_{T,i}\theta_{\mathcal{T},i}):=\beta\sum_{\{T\in\mathcal{T},\,0\leq i\leq m\}}h_{T}^{d-2s}c_{T,i}d_{T,i},$$

and  $B_{\mathcal{T}}^{\mathscr{S}}$  is as in §2.3.3, and define  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  by  $(B_{\mathcal{T}}w)(\tilde{w}) = (B_{\mathcal{T}}^{\mathscr{S} \oplus \mathscr{B}}w)(\tilde{w}).$ 

Using that with  $\mathbf{R}_{\mathcal{T}}$  as defined in (2.40),  $\begin{bmatrix} \mathrm{Id} & 0 \\ \mathbf{R}_{\mathcal{T}} & \mathrm{Id} \end{bmatrix}$  is the basis transformation from  $\Psi_{\mathcal{T}}$  to  $\Psi_{\mathcal{T}}^0 \cup \{\theta_{T,i} \colon T \in \mathcal{T}, 1 \leq i \leq m\}$ , one infers that the representation of the resulting uniform preconditioner reads as (2.43)

$$\mathbf{G}_{\mathcal{T}} = \begin{bmatrix} \mathbf{D}_{\mathcal{T}}^{0} & 0\\ 0 & \mathbf{D}_{\mathcal{T}}^{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathrm{Id} & \mathbf{R}_{\mathcal{T}}^{\top}\\ 0 & \mathrm{Id} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{\mathcal{T}}^{0} & 0\\ 0 & \beta(\mathbf{D}_{\mathcal{T}}^{1})^{1-\frac{2s}{d}} \end{bmatrix} \begin{bmatrix} \mathrm{Id} & 0\\ \mathbf{R}_{\mathcal{T}} & \mathrm{Id} \end{bmatrix} \begin{bmatrix} \mathbf{D}_{\mathcal{T}}^{0} & 0\\ 0 & \mathbf{D}_{\mathcal{T}}^{1} \end{bmatrix}^{-1}$$

which is thus independent of the particular bubbles  $\Theta_{\mathcal{T}}$  being chosen.

#### Continuous piecewise polynomials

Given  $\ell > 1$ , for  $\mathcal{T} \in \mathbb{T}$ , let  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}^{0,\ell}_{\mathcal{T}}$ . We equip  $\mathscr{V}_{\mathcal{T}}$  with a basis  $\Xi_{\mathcal{T}} = \Xi^0_{\mathcal{T}} \cup (\Xi_{\mathcal{T}} \setminus \Xi^0_{\mathcal{T}})$ , where for each  $T \in \mathcal{T}$ , the set of restrictions to T of those basis functions that do not identically vanish on T is a lifted version of a fixed basis for

the polynomials of degree  $\ell$  on a reference *d*-simplex under an affine bijection; the support of each basis function  $\xi \in \Xi_{\mathcal{T}}$  is connected and extends to a uniformly bounded number of  $T \in \mathcal{T}$ ; and finally,  $\|\xi\|_{L_2(\Omega)} \approx |\operatorname{supp} \xi|^{\frac{1}{2}} \approx |T|^{\frac{1}{2}}$ for some  $T \in \mathcal{T}$  with  $\operatorname{supp} \xi \cap T \neq \emptyset$ .

Similar to the previous §2.5.1, a biorthogonal collection  $\Theta_{\mathcal{T}} = \{\theta(\xi) : \xi \in \Xi_{\mathcal{T}}\}$  of 'bubbles' exists that has properties analogous to (2.34)-(2.39), reading |T| as  $|\operatorname{supp} \xi|$ , and  $h_T$  as  $|\operatorname{supp} \xi|^{1/d}$ . Writing the collection  $\Psi_{\mathcal{T}}^0$  as found in §2.4 as  $\{\psi^0(\xi) : \xi \in \Xi_{\mathcal{T}}\}$ , we define  $\Psi_{\mathcal{T}} = \{\psi(\xi) : \xi \in \Xi_{\mathcal{T}}\}$ , biorthogonal to  $\Xi_{\mathcal{T}}$ , and  $\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}}$  by

$$\psi(\xi) := \begin{cases} \psi^0(\xi) - \sum_{\xi' \in \Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0} \frac{\langle \psi^0(\xi), \xi' \rangle_{L_2(\Omega)}}{\langle \theta(\xi'), \xi' \rangle_{L_2(\Omega)}} \theta(\xi') & \xi \in \Xi_{\mathcal{T}}, \\ \theta(\xi) & \xi \in \Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0. \end{cases}$$

Theorem 2.5.1 extends to the current setting and shows that the corresponding biorthogonal projector is uniformly bounded. The representation of resulting uniform preconditioner reads as (2.43), obviously now with  $\mathbf{D}_{\mathcal{T}}^0$ and  $\mathbf{B}_{\mathcal{T}}^0$  as found in the continuous piecewise linear case, and the matrix  $\mathbf{D}_{\mathcal{T}}^1 := \text{diag}\{|\sup p\xi|: \xi \in \Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0\}$ , and  $(\mathbf{R}_{\mathcal{T}}^1)_{\xi,\nu} := -|\sup p\xi|^{-1}\langle \phi_{\mathcal{T},\nu}, \xi \rangle_{L_2(\Omega)}$ for  $(\xi, \nu) \in (\Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0) \times N_{\mathcal{T}}^0$ , and  $(\mathbf{R}_{\mathcal{T}}^1)_{\xi,\nu} = 0$  for  $\nu \in N_{\mathcal{T}} \setminus N_{\mathcal{T}}^0$  (recall  $\#\Xi_{\mathcal{T}}^0 = \#N_{\mathcal{T}}$ ).

# 2.5.2 Application of a subspace correction framework

For  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,\ell}$ , in §2.5.1 we have demonstrated existence of a biorthogonal projector that satisfies (2.42). In §2.5.1 we have shown that a similar result holds true for  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,\ell}$ . From the proof of Lemma 2.3.4 we learn that this implies that for either choice of  $\mathscr{V}_{\mathcal{T}}$ ,

(2.44) 
$$\|h_{\mathcal{T}} v_{\mathcal{T}}\|_{L_2(\Omega)} \lesssim \|v_{\mathcal{T}}\|_{(H^1_{\Omega,\infty})'} \quad (v_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}).$$

Using this inverse inequality, we are going to decompose  $\mathscr{V}_{\mathcal{T}}$  in a uniformly stable way into  $\mathscr{V}_{\mathcal{T}}^0$  and a complement space on which the  $\| \|_{\mathscr{V}}$ -norm is equivalent to a scaled  $L_2(\Omega)$ -norm.

**Proposition 2.5.2.** Let  $Q_{\mathcal{T}}^0 \in \mathcal{L}(L_2(\Omega), L_2(\Omega))$  be a projector with  $\operatorname{ran} Q_{\mathcal{T}}^0 = \mathscr{V}_{\mathcal{T}}^0$ ,  $\sup_{\mathcal{T} \in \mathbb{T}} \|Q_{\mathcal{T}}^0\|_{\mathcal{L}(L_2(\Omega), L_2(\Omega))} < \infty$ , and

$$\sup_{\mathcal{T}\in\mathbb{T}} \|h_{\mathcal{T}}^{-1}(\mathrm{Id}-(Q_{\mathcal{T}}^{0})^{*})\|_{\mathcal{L}(H^{1}_{0,\gamma}(\Omega),L_{2}(\Omega))} < \infty.$$

Then  $\sup_{\mathcal{T}\in\mathbb{T}} \|Q^0_{\mathcal{T}}\|_{\mathscr{V}_{\mathcal{T}}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}})} < \infty$ , and

(2.45) 
$$\|h_{\mathcal{T}}^s \cdot \|_{L_2(\Omega)} \approx \| \cdot \|_{\mathscr{V}} \quad on \ \mathscr{V}_{\mathcal{T}}^1 \coloneqq \operatorname{ran}(\operatorname{Id} - Q_{\mathcal{T}}^0)|_{\mathscr{V}_{\mathcal{T}}}.$$

*Proof.* For  $u \in \mathscr{V}_{\mathcal{T}}$ , thanks to (2.44) we have

$$\begin{aligned} \| (\mathrm{Id} - Q_{\mathcal{T}}^{0}) u \|_{H^{1}_{0,\gamma}(\Omega)'} &= \sup_{v \in H^{1}_{0,\gamma}(\Omega)} \frac{\langle u, (\mathrm{Id} - (Q_{\mathcal{T}}^{0})^{*}) v \rangle_{L_{2}(\Omega)}}{\| v \|_{H^{1}(\Omega)}} \lesssim \| h_{\mathcal{T}} u \|_{L_{2}(T)} \\ &\lesssim \| u \|_{H^{1}_{0,\gamma}(\Omega)'}, \end{aligned}$$

so that the first statement follows by interpolation.

Again by interpolation, the second statement needs only to be proven for s = 1. For that case it follows from the above equation and (2.44).

Next we use the decomposition  $\mathscr{V}_{\mathcal{T}} = \mathscr{V}_{\mathcal{T}}^0 \oplus \mathscr{V}_{\mathcal{T}}^1$  from Proposition 2.5.2 to build a preconditioner on  $\mathscr{V}_{\mathcal{T}}$  from preconditioners on the subspaces.

**Proposition 2.5.3.** In the situation of Proposition 2.5.2, for i = 0, 1, let  $I_{\mathcal{T}}^i$  denote the embedding of  $\mathcal{V}_{\mathcal{T}}^i$  into  $\mathcal{V}_{\mathcal{T}}$ , and let  $G_{\mathcal{T}}^i \in \mathcal{L}is_c((\mathcal{V}_{\mathcal{T}}^i)', \mathcal{V}_{\mathcal{T}}^i)$ . Then  $G_{\mathcal{T}} := \sum_{i=0}^{1} I_{\mathcal{T}}^i G_{\mathcal{T}}^i(I_{\mathcal{T}}^i)' \in \mathcal{L}is_c((\mathcal{V}_{\mathcal{T}})', \mathcal{V}_{\mathcal{T}})$  with

$$\begin{aligned} \|G_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}',\mathscr{V}_{\mathcal{T}})} &\leq 2 \max_{i} \|G_{\mathcal{T}}^{i}\|_{\mathcal{L}((\mathscr{V}_{\mathcal{T}}^{i})',\mathscr{V}_{\mathcal{T}}^{i})}, \\ \|\Re(G_{\mathcal{T}})^{-1}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}}')} &\leq 2 \|Q_{\mathcal{T}}^{0}|_{\mathscr{V}_{\mathcal{T}}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}})} \max_{i} \|\Re(G_{\mathcal{T}}^{i})^{-1}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}^{i},(\mathscr{V}_{\mathcal{T}}^{i})')} \end{aligned}$$

*Proof.* The result follows as an easy case from the general theory of (additive) subspace correction methods (e.g. [Osw94] + references cited there), together with the inequalities

$$\frac{1}{2} \| \cdot \|_{\mathscr{V}}^2 \le \|Q_{\mathcal{T}}^0 \cdot \|_{\mathscr{V}}^2 + \|(\mathrm{Id} - Q_{\mathcal{T}}^0) \cdot \|_{\mathscr{V}}^2 \le 2 \|Q_{\mathcal{T}}^0\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}})} \| \cdot \|_{\mathscr{V}}^2 \quad \text{on } \mathscr{V}_{\mathcal{T}},$$

where we used that  $\|(\text{Id} - Q_{\mathcal{T}}^0)|_{\mathscr{V}_{\mathcal{T}}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}})} = \|Q_{\mathcal{T}}^0|_{\mathscr{V}_{\mathcal{T}}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}})}$ , because  $Q_{\mathcal{T}}^0|_{\mathscr{V}_{\mathcal{T}}}$  is a projector unequal to both the zero map and the identity ([Kat60, XZ03])

On  $\mathscr{V}_{\mathcal{T}}^0$  we already have our uniform preconditioner  $G_{\mathcal{T}}^0$  available, so it remains to construct such a preconditioner on the complement space  $\mathscr{V}_{\mathcal{T}}^1$ . In the situation of Proposition 2.5.2, let  $\Xi_{\mathcal{T}} = \Xi_{\mathcal{T}}^0 \cup (\Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0)$  be a basis for  $\mathscr{V}_{\mathcal{T}}$ such that  $\Xi_{\mathcal{T}}^1 := \Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0$  is a basis for  $\mathscr{V}_{\mathcal{T}}^1$  for which

(2.46) 
$$\|h_{\mathcal{T}}^{s} \sum_{\xi \in \Xi_{\mathcal{T}}^{1}} c_{\xi} \xi \|_{L_{2}(\Omega)}^{2} \approx \sum_{\xi \in \Xi_{\mathcal{T}}^{1}} |c_{\xi}|^{2} \|h_{\mathcal{T}}^{s} \xi \|_{L_{2}(\Omega)}^{2}.$$

Then

$$(H^1_{\mathcal{T}}\sum_{\xi\in \Xi^1_{\mathcal{T}}}c_{\xi}\xi)(\sum_{\xi\in \Xi^1_{\mathcal{T}}}d_{\xi}\xi):=\sum_{\xi\in \Xi^1_{\mathcal{T}}}c_{\xi}d_{\xi}\|h^s_{\mathcal{T}}\xi\|^2_{L_2(\Omega)}$$

43

satisfies  $\sup_{\mathcal{T}\in\mathbb{T}} \max(\|H_{\mathcal{T}}^{1}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}^{1},(\mathscr{V}_{\mathcal{T}}^{1})')}, \|\Re(H_{\mathcal{T}}^{1})^{-1}\|_{\mathcal{L}((\mathscr{V}_{\mathcal{T}}^{1})',\mathscr{V}_{\mathcal{T}})}) < \infty$  thanks to (2.45), and so  $G_{\mathcal{T}}^{1} := (H_{\mathcal{T}}^{1})^{-1}$  is a suitable choice. The implementation of the resulting uniform preconditioner  $G_{\mathcal{T}}$  reads as

(2.47) 
$$\mathbf{G}_{\mathcal{T}} = \begin{bmatrix} \mathbf{G}_{\mathcal{T}}^{0} & 0\\ 0 & \operatorname{diag}\{\|h_{\mathcal{T}}^{s}\xi\|_{L_{2}(\Omega)}^{-2} \colon \xi \in \Xi_{\mathcal{T}}^{1}\} \end{bmatrix}.$$

What remains is, for both options for  $\mathscr{V}_{\mathcal{T}}$ , to specify a  $Q^0_{\mathcal{T}}$  that satisfies the conditions from Proposition 2.5.2, and to equip  $\mathscr{V}_{\mathcal{T}}$  with a basis  $\Xi_{\mathcal{T}}$  that is the union of the basis  $\Xi^0_{\mathcal{T}}$  for  $\mathscr{V}^0_{\mathcal{T}}$ , and a basis  $\Xi^1_{\mathcal{T}}$  for  $\mathscr{V}^1_{\mathcal{T}} = \operatorname{ran}(\operatorname{Id} - Q^0_{\mathcal{T}})|_{\mathscr{V}_{\mathcal{T}}}$ , the latter being *uniformly stable* w.r.t.  $\|h^s_{\mathcal{T}} \cdot \|_{L_2(\Omega)}$ , i.e., one that satisfies (2.46).

#### Discontinuous piecewise polynomials

For  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,\ell}$ , the conditions of Proposition 2.5.2 are fulfilled by taking  $Q_{\mathcal{T}}^0$  to be the  $L_2(\Omega)$ -orthogonal projector onto  $\mathscr{V}_{\mathcal{T}}^0 = \mathscr{S}_{\mathcal{T}}^{-1,0}$ .<sup>7</sup>

By taking  $\Xi_{\mathcal{T}} = \{\xi_{T,i}\}_{T \in \mathcal{T}, 0 \le i \le m}$  to be an  $L_2(\Omega)$ -orthogonal basis for  $\mathscr{V}_{\mathcal{T}}$  such that  $\xi_{\mathcal{T},0} = \xi_T$  and  $\operatorname{supp} \xi_{\mathcal{T},i} \subset \overline{T}$ , (2.46) is valid.

Remarkably, with these specifications and by scaling  $\|\xi_{\mathcal{T},i}\|_{L_2(\Omega)} = |T|^{\frac{1}{2}}$ , the resulting  $\mathbf{G}_{\mathcal{T}}$  is given by (2.43) with  $\mathbf{R}_{\mathcal{T}}$  reading as the zero map ( $\mathbf{R}_{\mathcal{T}}$  from (2.40) is non-zero).

#### Continuous piecewise polynomials

Let  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,\ell}$  for  $\ell > 1$ . It is no option to take  $Q_{\mathcal{T}}^0$  to be the  $L_2(\Omega)$ -orthogonal projector onto  $\mathscr{V}_{\mathcal{T}}^0 = \mathscr{S}_{\mathcal{T}}^{0,1}$  because in that case we will not be able to equip  $\operatorname{ran}(\operatorname{Id} - Q_{\mathcal{T}}^0)|_{\mathscr{V}_{\mathcal{T}}}$  with a locally supported basis.

From (2.32) recall the biorthogonal projector  $P_{\mathcal{T}}^0$  onto  $\mathscr{W}_{\mathcal{T}}^0$  with  $\operatorname{ran}(\operatorname{Id} - P_{\mathcal{T}}^0) = \mathscr{V}_{\mathcal{T}}^0$ . Writing  $(\operatorname{Id} - P_{\mathcal{T}}^0) = (\operatorname{Id} - P_{\mathcal{T}}^0)(\operatorname{Id} - \Pi_{\mathcal{T}}^0)$  with  $\Pi_{\mathcal{T}}^0$  being a Scott-Zhang type interpolator, one shows that  $\sup_{\mathcal{T} \in \mathbb{T}} \|h_{\mathcal{T}}^{-1}(\operatorname{Id} - P_{\mathcal{T}}^0)\|_{\mathcal{L}(H_{0,\gamma}^1(\Omega), L_2(\Omega))} < \infty$ . Since furthermore  $\sup_{\mathcal{T} \in \mathbb{T}} \|P_{\mathcal{T}}^0\|_{\mathcal{L}(L_2(\Omega), L_2(\Omega))} < \infty$ , the conditions of Proposition 2.5.2 are satisfied by taking  $Q_{\mathcal{T}}^0 := (P_{\mathcal{T}}^0)^* = u \mapsto \sum_{\nu \in N_{\mathcal{T}}} \frac{(d+1)\langle u, \psi_{\mathcal{T}, \nu}^0 \rangle_{L_2(\Omega)}}{|\omega_{\mathcal{T}}(\nu)|} \xi_{\mathcal{T}, \nu}^0$ .

Let us denote the *weighted*  $L_2(\Omega)$ -norm  $||h_{\mathcal{T}}^s \cdot ||_{L_2(\Omega)}$  by  $||| \cdot |||$ . We need to equip  $\mathscr{V}_{\mathcal{T}}^1 = \operatorname{ran}(\operatorname{Id} - Q_{\mathcal{T}}^0)|_{\mathscr{V}_{\mathcal{T}}}$  with a basis that is uniformly stable w.r.t.  $||| \cdot |||$ . Since the supports of the basis functions will extend to multiple  $T \in \mathcal{T}$ , this task is more complex than for the discontinuous piecewise polynomial case. Let  $\Xi_{\mathcal{T}} = \Xi_{\mathcal{T}}^0 \cup (\Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0)$  be a basis for  $\mathscr{V}_{\mathcal{T}}$  such that for each  $T \in \mathcal{T}$ ,  $\{\xi|_T \colon \xi \in \Xi, \xi|_T \neq 0\}$  is a uniformly  $L_2(T)$ -stable basis for its span. Common affine equivalent constructions yield such a basis. Then  $\Xi_{\mathcal{T}}$  is stable w.r.t.  $||| \cdot |||$ , uniformly in  $\mathcal{T} \in \mathbb{T}$ , and so in particular, with  $\widetilde{\mathscr{V}_{\mathcal{T}}^1} := \operatorname{span} \Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0$ ,  $\mathscr{V}_{\mathcal{T}} = \mathscr{V}_{\mathcal{T}}^0 \oplus \widetilde{\mathscr{V}_{\mathcal{T}}^1}$  is a uniformly stable decomposition w.r.t.  $||| \cdot |||$ .

<sup>&</sup>lt;sup>7</sup>Notice that for  $s \geq \frac{1}{2}$ ,  $Q_T^0 \notin \mathcal{L}(\mathcal{V}, \mathcal{V})$ , so taking its restriction to  $\mathscr{V}_T$  is essential in Proposition 2.5.2.

Corresponding to this decomposition, for  $v \in \mathscr{V}_{\mathcal{T}}$  we write  $v = v^0 + \bar{v}^1$ . Taking  $v \in \mathscr{V}_{\mathcal{T}}^1$ , it holds that  $0 = Q_{\mathcal{T}}^0 v = Q_{\mathcal{T}}^0 v^0 + Q_{\mathcal{T}}^0 \bar{v}^1 = v^0 + Q_{\mathcal{T}}^0 \bar{v}^1$ , or  $v^0 = -Q_{\mathcal{T}}^0 \bar{v}^1$  i.e.  $v = (\mathrm{Id} - Q_{\mathcal{T}}^0) \bar{v}^1$ , showing that  $\mathrm{Id} - Q_{\mathcal{T}}^0 \colon \mathscr{V}_{\mathcal{T}}^1 \to \mathscr{V}_{\mathcal{T}}^1$  is surjective. Injectivity of this map follows from  $||v||| \approx ||Q_{\mathcal{T}}^0 \bar{v}^1|| + ||\bar{v}^1||| \ge ||\bar{v}^1|||$ , and bounded invertibility, uniformly in  $\mathcal{T}$ , will follow from  $Q_{\mathcal{T}}^0$  being uniformly bounded w.r.t.  $||| \cdot |||$ . The latter holds true because of  $||\psi_{\mathcal{T},\nu}^0||_{L_2(\Omega)}||\xi_{\mathcal{T},\nu}^0||_{L_2(\Omega)} \approx \frac{|\omega_{\mathcal{T}}(\nu)|}{d+1}$ , the local supports of the  $\psi_{\mathcal{T},\nu}^0$  and  $\xi_{\mathcal{T},\nu}^0$ , and the uniform K-mesh property of  $\mathcal{T}$ . We conclude that  $(\mathrm{Id} - Q_{\mathcal{T}}^0)(\Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0)$  is a basis for  $\mathscr{V}_{\mathcal{T}}^1$  that is uniformly stable w.r.t.  $||| \cdot |||$ .

Since  $\Xi_{\mathcal{T}}^{0} \cup (\operatorname{Id} - Q_{\mathcal{T}}^{0})(\Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^{0})$  is not the basis of choice to set up the stiffness matrix, we give the implementation  $\mathbf{G}_{\mathcal{T}}$  of the uniform preconditioner  $G_{\mathcal{T}}$  for  $\mathscr{V}_{\mathcal{T}}$  being equipped with  $\Xi_{\mathcal{T}}$ , partitioned into  $\Xi_{\mathcal{T}}^{0} = \{\xi_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}\}$  and  $\Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^{0}$ . It reads as

$$\mathbf{G}_{\mathcal{T}} = \begin{bmatrix} \mathrm{Id} & \mathbf{S}_{\mathcal{T}} \\ 0 & \mathrm{Id} \end{bmatrix} \begin{bmatrix} \mathbf{G}_{\mathcal{T}}^{0} & 0 \\ 0 & \mathrm{diag}\{\|h_{\mathcal{T}}^{s}\xi\|_{L_{2}(\Omega)}^{-2} \colon \xi \in \Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^{0}\} \end{bmatrix} \begin{bmatrix} \mathrm{Id} & 0 \\ \mathbf{S}_{\mathcal{T}}^{\top} & \mathrm{Id} \end{bmatrix},$$

where for  $(\nu,\xi) \in N_{\mathcal{T}} \times (\Xi_{\mathcal{T}} \setminus \Xi_{\mathcal{T}}^0)$ ,  $(\mathbf{S}_{\mathcal{T}})_{\nu\xi} := -\frac{(d+1)\langle \xi, \psi_{\mathcal{T},\nu}^0 \rangle_{L_2(\Omega)}}{|\omega_{\mathcal{T}}(\nu)|}$ , and where in (2.47) we have replaced  $\|h_{\mathcal{T}}^s(\mathrm{Id} - Q_{\mathcal{T}}^0)\xi\|_{L_2(\Omega)}^{-2}$  by the equivalent  $\|h_{\mathcal{T}}^s\xi\|_{L_2(\Omega)}^{-2}$ .

# 2.6 Numerical experiments

Let  $\Gamma = \partial [0,1]^3 \subset \mathbb{R}^3$  be the two-dimensional manifold without boundary given as the boundary of the unit cube,  $\mathscr{W} := H^{1/2}(\Gamma)$ ,  $\mathscr{V} := H^{-1/2}(\Gamma)$ , and  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,\ell} \subset \mathscr{V}$ .

The role of the opposite order operator  $B_{\mathcal{T}}^{\mathscr{S}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  from Section 2.3.3 will be fulfilled by  $(B_{\mathcal{T}}^{\mathscr{S}}u)(v) := (Bu)(v)$  for an adapted hypersingular operator  $B \in \mathcal{L}is_c(\mathcal{W}, \mathcal{W}')$ . The hypersingular operator  $\tilde{B} \in \mathcal{L}(\mathcal{W}, \mathcal{W}')$ itself is only semi-coercive, but there are various options to change it into a coercive operator ([SW98]). We consider  $B \in \mathcal{L}is_c(\mathcal{W}, \mathcal{W}')$  given by (Bu)(v) = $(\tilde{B}u)(v) + \alpha \langle u, 1 \rangle_{L_2(\Gamma)} \langle v, 1 \rangle_{L_2(\Gamma)}$  for some  $\alpha > 0$ . By comparing different values numerically, we find  $\alpha = 0.05$  to give good results in our examples.

With  $m := \binom{2+\ell}{\ell} - 1$ , as in §2.5.2 we equip  $\mathscr{V}_{\mathcal{T}}$  with a usual  $L_2(\Gamma)$ -orthogonal basis  $\{\mathbb{1}|_T : T \in \mathcal{T}\} \cup \{\xi_{T,i} : T \in \mathcal{T}, 1 \leq i \leq m\}$  where  $\operatorname{supp} \xi_{T,i} \subset \overline{T}$ ,  $\|\xi_{T,i}\|_{L_2(\Omega)} = |T|^{\frac{1}{2}}$ , and denote the resulting stiffness matrix by  $\mathbf{A}_{\mathcal{T}}$ . The lowest order case  $\ell = 0$  corresponds to m = 0.

Equipping  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  with the nodal basis  $\Phi_{\mathcal{T}}$  defined in (2.4)-(2.5), for  $\ell = 0$  the matrix representation of the preconditioner reads as

$$\boldsymbol{G}_{\mathcal{T}} = \boldsymbol{D}_{\mathcal{T}}^{-1} \big( \boldsymbol{p}_{\mathcal{T}}^{\top} \boldsymbol{B}_{\mathcal{T}}^{\mathscr{S}} \boldsymbol{p}_{\mathcal{T}} + \beta \boldsymbol{q}_{\mathcal{T}}^{\top} \boldsymbol{D}_{\mathcal{T}}^{1/2} \boldsymbol{q}_{\mathcal{T}} \big) \boldsymbol{D}_{\mathcal{T}}^{-1}$$

with  $D_{\mathcal{T}} = \text{diag}\{|T|: T \in \mathcal{T}\}$  and uniformly sparse  $p_{\mathcal{T}}$  and  $q_{\mathcal{T}}$  as given in Sect. 2.3.4.

Denoting above  $G_{\mathcal{T}}$  by  $G_{\mathcal{T}}^0$ , by applying for  $\ell > 0$  the subspace correction method from §2.5.2, the matrix representation of the resulting uniform preconditioner is given by

$$\boldsymbol{G}_{\mathcal{T}} = \begin{bmatrix} \boldsymbol{G}_{\mathcal{T}}^{0} & 0\\ 0 & \beta \operatorname{diag}\{|T|^{-3/2} \operatorname{Id}_{m \times m} \colon T \in \mathcal{T}\} \end{bmatrix}$$

The (full) matrix representations of the discretized singular integral operators  $A_{\mathcal{T}}$  and  $B_{\mathcal{T}}^{\mathscr{S}}$  are calculated using the BETL2 software package [HK12] (alternatively, one may apply low rank approximations in a hierarchical format). Condition numbers are determined using Lanczos iteration with respect to  $\|\|\cdot\| := \|A_{\mathcal{T}}^{\frac{1}{2}} \cdot \|$ . The constant  $\beta$  is approximately optimized by comparing different choices numerically.

We will compare our preconditioner to the diagonal preconditioner diag $(\mathbf{A}_{\mathcal{T}})^{-1}$ , and in the piecewise constant case, also to the related preconditioner  $\hat{\mathbf{G}}_{\mathcal{T}}$  from [HUT16], where  $\hat{\mathbf{G}}_{\mathcal{T}} = \hat{\mathbf{D}}_{\mathcal{T}}^{-1} \mathbf{E}_{\mathcal{T}}^{\top} \mathbf{B}_{\hat{\mathcal{T}}}^{\mathscr{G}} \mathbf{E}_{\mathcal{T}} \hat{\mathbf{D}}_{\mathcal{T}}^{-\top}$  is defined as follows. With  $\hat{\mathcal{T}}$ being the barycentric refinement of  $\mathcal{T}$ , a collection  $\hat{\Psi}_{\mathcal{T}} \subset \mathscr{S}_{\hat{\mathcal{T}},0}^{0,1}$  is constructed in [BC07] such that the Fortin projector  $\hat{P}_{\mathcal{T}}$  with ran  $\hat{P}_{\mathcal{T}} = \hat{\mathcal{W}}_{\mathcal{T}} := \operatorname{span} \hat{\Psi}_{\mathcal{T}}$ and ran(Id  $-\hat{P}_{\mathcal{T}}) = \mathscr{V}_{\mathcal{T}}^{\perp_{L_2(\Gamma)}}$  exists, and, under an additional sufficiently mildly-grading condition on the partition, has a uniformly bounded norm  $\|\hat{P}_{\mathcal{T}}\|_{\mathcal{L}(\mathcal{W},\mathcal{W})}$  (cf. Theorem 2.3.2);  $\hat{\mathbf{D}}_{\mathcal{T}} := \langle \Xi_{\mathcal{T}}, \hat{\Psi}_{\mathcal{T}} \rangle_{L_2(\Gamma)}$ ;  $\mathbf{E}_{\mathcal{T}}$  is the representation of the embedding  $\widehat{\mathscr{W}}_{\mathcal{T}} \hookrightarrow \mathscr{S}_{\hat{\mathcal{T}},0}^{0,1}$  equipped with  $\hat{\Psi}_{\mathcal{T}}$  and the nodal basis of  $\mathscr{S}_{\hat{\mathcal{T}},0}^{0,1}$ , respectively; and  $B_{\hat{\mathcal{T}}}^{\mathscr{G}} \in \mathcal{L}\mathrm{is}_c(\mathscr{S}_{\hat{\mathcal{T}}}^{0,1}, (\mathscr{S}_{\hat{\mathcal{T}}}^{0,1})')$  is an opposite order operator that we take as  $(B_{\hat{\mathcal{T}}}^{\mathscr{G}}u)(v) := (Bu)(v)$ , with B the adapted hypersingular operator.

Compared to our  $G_{\mathcal{T}} = G_{\mathcal{T}}^0$ , the preconditioner  $\hat{G}_{\mathcal{T}}$  has the disadvantages that, besides the aforementioned mildly grading condition, the matrix  $\hat{D}_{\mathcal{T}}$ , although uniformly sparse, is not diagonal, so that the (sufficiently accurate) application of its inverse cannot be performed in linear complexity; furthermore that it requires evaluating the adapted hypersingular operator on the larger space  $\mathscr{S}_{\mathcal{T},0}^{0,1} \supset \mathscr{S}_{\mathcal{T},0}^{0,1}$  ( $\#\hat{\mathcal{T}} = 6\#\mathcal{T}$ ); and finally that the non-standard barycentric refinement  $\hat{\mathcal{T}}$  has to be generated.

## 2.6.1 Uniform refinements

Consider a conforming triangulation  $\mathcal{T}_1$  of  $\Gamma$  consisting of 2 triangles per side, so 12 triangles in total. We let  $\mathbb{T}$  be the sequence  $\{\mathcal{T}_k\}_{k\geq 1}$  of uniform redrefinements, where  $\mathcal{T}_k \succ \mathcal{T}_{k-1}$  is found by subdividing each triangle from  $\mathcal{T}_{k-1}$  into 4 congruent subtriangles.

For  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,\ell}$ , Tables 2.1 and 2.2 show the condition numbers of the preconditioned system for  $\ell = 0$  and  $\ell = 2$ , respectively. Aside from being uni-

dofs	$\kappa_S(\operatorname{diag}(\boldsymbol{A}_{\mathcal{T}})^{-1}\boldsymbol{A}_{\mathcal{T}})$	$\kappa_S({oldsymbol G}_{\mathcal T}{oldsymbol A}_{\mathcal T})$	$\kappa_S(\hat{oldsymbol{G}}_{\mathcal{T}}oldsymbol{A}_{\mathcal{T}})$
12	14.56	2.51	1.29
48	29.30	2.52	1.58
192	58.25	2.66	1.77
768	116.3	2.71	1.89
3072	230.0	2.74	1.94
12288	444.8	2.79	

TABLE 2.1. Spectral condition numbers of the preconditioned single layer system, using uniform refinements, discretized by piecewise constants  $\mathscr{S}_{\mathcal{T}}^{-1,0}$ . Both matrices  $G_{\mathcal{T}}$  and  $\hat{G}_{\mathcal{T}}$  are constructed using the adapted hypersingular operator with  $\alpha = 0.05$ ; and  $\beta = 1.25$  in  $G_{\mathcal{T}}$ .

TABLE 2.2. Spectral condition numbers of the preconditioned single layer system, using uniform refinements, discretized by discontinuous piecewise quadratics  $\mathscr{S}_{\mathcal{T}}^{-1,2}$ . The matrix  $G_{\mathcal{T}}$  is constructed using the adapted hypersingular operator, with  $\alpha = 0.05$ , and  $\beta = 1.25$ .

dofs	$\kappa_S(\operatorname{diag}(\boldsymbol{A}_{\mathcal{T}})^{-1}\boldsymbol{A}_{\mathcal{T}})$	$\kappa_S(oldsymbol{G}_\mathcal{T}oldsymbol{A}_\mathcal{T})$
72	167.16	9.58
288	309.12	10.4
1152	616.03	11.1
4608	1211.3	11.3
18432	2337.2	11.4

formly bounded, the condition numbers of our preconditioner  $G_{\mathcal{T}}$  are of modest size. In the constant case,  $\ell = 0$ , Table 2.1 reveals that the preconditioner  $\hat{G}_{\mathcal{T}}$  from [BC07, HUT16] gives better condition numbers. As described above, this quantitative gain comes at a price. In the result of dim  $\mathscr{S}_{\mathcal{T}}^{-1,0} = 3072$ , using full matrices for the discretized adapted hypersingular operator, we found a setup and application time of 1816s and 0.0971s for  $\hat{G}_{\mathcal{T}}$ , compared to 385s and 0.00284s for  $G_{\mathcal{T}}$ . These differences are due to numerical inversion of  $\hat{D}_{\mathcal{T}}$ by LU factorization with partial pivoting, and the enlargement  $\mathscr{S}_{\mathcal{T},0}^{0,1} \supset \mathscr{S}_{\mathcal{T},0}^{0,1}$ , also causing our test machine to go out of memory in calculating  $\hat{G}_{\mathcal{T}}$  for the last refinement. Although we expect them to be in any case significant, these differences can be made smaller when the exact inversion of  $\hat{D}_{\mathcal{T}}$  is avoided, and  $\mathcal{B}_{\mathcal{T}}^{\mathscr{S}}$  and  $\mathcal{B}_{\mathcal{T}}^{\mathscr{S}}$  are replaced by suitable low rank approximations.

TABLE 2.3. Spectral condition numbers of the preconditioned single layer system discretized by piecewise constants  $\mathscr{S}_{\mathcal{T}}^{-1,0}$  using local refinements at each of the eight cube corners. Both matrices  $G_{\mathcal{T}}$  and  $\hat{G}_{\mathcal{T}}$  are constructed using the adapted hypersingular operator with  $\alpha = 0.05$ ; and  $\beta = 1.2$  in  $G_{\mathcal{T}}$ . The second column is defined by  $h_{\mathcal{T},min} := \min_{T \in \mathcal{T}} h_T$ .

dofs	$h_{\mathcal{T},min}$	$\kappa_S(\operatorname{diag}(\boldsymbol{A}_{\mathcal{T}})^{-1}\boldsymbol{A}_{\mathcal{T}})$	$\kappa_S(oldsymbol{G}_\mathcal{T}oldsymbol{A}_\mathcal{T})$	$\kappa_S(\hat{m{G}}_{\mathcal{T}}m{A}_{\mathcal{T}})$
12	$7.0 \cdot 10^{-1}$	14.56	2.61	1.29
432	$2.2\cdot10^{-2}$	68.66	2.64	2.91
912	$6.9\cdot 10^{-4}$	73.15	2.64	3.14
1872	$6.7\cdot 10^{-7}$	73.70	2.64	3.25
2352	$2.1\cdot 10^{-8}$	73.80	2.64	3.26
2976	$2.3\cdot 10^{-10}$	73.66	2.64	

#### 2.6.2 Local refinements

Here we take  $\mathbb{T}$  to be the sequence  $\{\mathcal{T}_k\}_{k\geq 1}$  of locally refined triangulations, where  $\mathcal{T}_k \succ \mathcal{T}_{k-1}$  is constructed using conforming newest vertex bisection to refine all triangles in  $\mathcal{T}_{k-1}$  that touch a corner of the cube.

As noted before, the preconditioner  $\hat{G}_{\mathcal{T}}$  provides uniformly bounded condition numbers if the family  $\mathbb{T}$  satisfies some sufficiently mildly-grading condition on the partition [Ste03a, HUT16]. It is not directly clear whether  $\mathbb{T}$  satisfies this condition, but we included the results nonetheless.

Table 2.3 gives the results for the preconditioned single layer operator discretized by piecewise constants  $\mathscr{S}_{\mathcal{T}}^{-1,0}$ . The condition numbers  $\kappa_S(G_{\mathcal{T}}A_{\mathcal{T}})$  are nicely bounded under local refinements. In this case our preconditioner gives condition numbers slightly smaller than the ones found with  $\hat{G}_{\mathcal{T}}$ . The calculation of the LU decomposition with partial pivoting of  $\hat{D}_{\mathcal{T}}$  turns out to break down in the last result (dim  $\mathscr{S}_{\mathcal{T}}^{-1,0} = 2976$ ).

# 2.7 Conclusion

In this chapter, we have seen how a uniformly boundedly invertible operator  $B_{\mathcal{T}}^{\mathscr{G}}$  from the space of continuous piecewise linears w.r.t. any conforming shape regular partition  $\mathcal{T}$ , equipped with the norm of  $H^s(\Omega)$  (or  $H^s(\Gamma)$ ) for some  $s \in [0, 1]$ , to its dual can be used to uniformly precondition a boundedly invertible operator of opposite order discretized by discontinuous or continuous polynomials of any fixed degree w.r.t.  $\mathcal{T}$ . The cost of the resulting preconditioner is the sum of a cost that scales linearly in  $\#\mathcal{T}$  and the cost of the application of  $B_{\mathcal{T}}^{\mathscr{G}}$ . For  $\mathcal{T}$  being member of a nested sequence of quasiuniform partitions,  $B_{\mathcal{T}}^{\mathscr{G}}$  has been constructed so that it requires linear cost. In the following chapter, we will realize this also for locally refined partitions.

# 3 Problems of negative order: preconditioning at linear cost

# 3.1 Introduction

In this chapter, we construct a multi-level type preconditioner for operators of negative orders  $-2s \in [-2, 0]$  that can be applied in linear time and yields uniformly bounded condition numbers. The preconditioner will be constructing using the framework of 'operator preconditioning' discussed in Chapter 2. The role of the 'opposite order operator' will be fulfilled by a multi-level type operator, based on the work of Wu and Zheng in [WZ17].

For some *d*-dimensional domain (or manifold)  $\Omega$ , a measurable, closed, possibly empty  $\gamma \subset \partial \Omega$ , and an  $s \in [0, 1]$ , we consider the Sobolev spaces

$$\mathscr{W} := [L_2(\Omega), H^1_{0,\gamma}(\Omega)]_{s,2}, \quad \mathscr{V} := \mathscr{W}'.$$

with  $H^1_{0,\gamma}(\Omega)$  being the closure in  $H^1(\Omega)$  of the smooth functions on  $\Omega$  that vanish at  $\gamma$ . Let  $(\mathscr{V}_{\mathcal{T}})_{\mathcal{T}\in\mathbb{T}} \subset \mathscr{V}$  be a family of *piecewise* or *continuous piecewise* polynomials of some fixed degree w.r.t. uniformly shape regular, possibly locally refined partitions. With, for  $\mathcal{T} \in \mathbb{T}$ ,  $A_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{V}'_{\mathcal{T}}$  being some boundedly invertible linear operator, we are interested in constructing a *preconditioner*  $G_{\mathcal{T}} \colon \mathscr{V}'_{\mathcal{T}} \to \mathscr{V}_{\mathcal{T}}$  such that the preconditioned operator  $G_{\mathcal{T}}A_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{V}_{\mathcal{T}}$  is uniformly boundedly invertible, and an application of  $G_{\mathcal{T}}$  can be evaluated in  $\mathcal{O}(\dim \mathscr{V}_{\mathcal{T}})$  arithmetic operations.

In order to create such a preconditioner, we will use the framework described in Chapter 2. Given  $\mathscr{V}_{\mathcal{T}}$ , we constructed an auxiliary space  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W}$ with dim  $\mathscr{W}_{\mathcal{T}} = \dim \mathscr{V}_{\mathcal{T}}$ , such that for  $D_{\mathcal{T}}$  defined by  $(D_{\mathcal{T}}v)(w) := \langle v, w \rangle_{L_2(\Omega)}$  $(v \in \mathscr{V}_{\mathcal{T}}, w \in \mathscr{W}_{\mathcal{T}})$  and some suitable 'opposite order' operator  $B_{\mathcal{T}}^{\mathscr{W}} : \mathscr{W}_{\mathcal{T}} \to \mathscr{W}_{\mathcal{T}}'$ , a preconditioner  $G_{\mathcal{T}}$  of the form  $G_{\mathcal{T}} := D_{\mathcal{T}}^{-1}B_{\mathcal{T}}^{\mathscr{W}}(D_{\mathcal{T}}')^{-1}$  is found. The space  $\mathscr{W}_{\mathcal{T}}$  is equipped with a basis that, modulo a scaling, is biorthogonal to the canonical basis for  $\mathscr{V}_{\mathcal{T}}$ , so that the representation of  $D_{\mathcal{T}}$  is an invertible diagonal matrix.

With  $\mathscr{S}_{\mathcal{T},0}^{0,1} \subset \mathscr{W}$  being the space of continuous piecewise linears w.r.t.  $\mathcal{T}$ , zero on  $\gamma$ , the above preconditioning approach hinges on the availability of

a uniformly boundedly invertible operator  $B_{\mathcal{T}}^{\mathscr{G}}: \mathscr{S}_{\mathcal{T},0}^{0,1} \to (\mathscr{S}_{\mathcal{T},0}^{0,1})'$ , which is generally the most demanding requirement. For example, if  $s = \frac{1}{2}$  and  $\gamma = \emptyset$ , a viable option is to take  $B_{\mathcal{T}}^{\mathscr{G}}$  as the discretized hypersingular operator. While this induces a uniform preconditioner, the application of  $B_{\mathcal{T}}^{\mathscr{G}}$  cannot be evaluated in linear complexity.

In this chapter we construct a suitable multi-level type operator  $B_{\mathcal{T}}^{\mathscr{S}}$  that *can* be applied in linear complexity. For this construction we require  $\mathbb{T}$  to be a family of conforming partitions created by Newest Vertex Bisection ([Mau95, Tra97]). In the aforementioned setting of having an arbitrary  $s \in [0, 1]$ , this multi-level operator  $B_{\mathcal{T}}^{\mathscr{S}}$  induces a uniform preconditioner  $G_{\mathcal{T}}$ , i.e.,  $G_{\mathcal{T}}A_{\mathcal{T}}$  is uniformly well-conditioned, where the cost of applying  $G_{\mathcal{T}}$  scales linearly in dim  $\mathcal{V}_{\mathcal{T}}$ . We also show that the preconditioner extends to the more general *manifold* case, where  $\Omega$  is a *d*-dimensional (piecewise) smooth Lipschitz manifold, and the trial space  $\mathcal{V}_{\mathcal{T}}$  is the parametric lift of a space of piecewise or continuous piecewise polynomials.

Finally, we remark that common multi-level preconditioners based on overlapping subspace decompositions are known not to work well for operators of negative order. A solution is provided by resorting to direct sum multi-level subspace decompositions. Examples are given by wavelet preconditioners, or closely related, the preconditioners from [BPV00], for the latter assuming quasi-uniform partitions.

For  $-s = -\frac{1}{2}$ , an optimal multi-level preconditioner based on a *non*overlapping subspace decomposition for operators defined on the boundary of a 2- or 3-dimensional Lipschitz polyhedron was recently introduced in [FHPS19].

# 3.1.1 Outline

In Sect. 3.2 we summarize the (operator) preconditioning framework from Chapter 2. In Sect. 3.3 we provide the multi-level type operator that can be used as the 'opposite order' operator inside the preconditioner framework. In Sect. 3.4 we comment on how to generalize the results to the case of piecewise smooth manifolds. In Sect. 3.5 we conclude with numerical results.

# 3.1.2 Notation

In this work, by  $\lambda \leq \mu$  we mean that  $\lambda$  can be bounded by a multiple of  $\mu$ , independently of parameters which  $\lambda$  and  $\mu$  may depend on, with the sole exception of the space dimension d, or in the manifold case, on the parametrization of the manifold that is used to define the finite element spaces on it. Obviously,  $\lambda \geq \mu$  is defined as  $\mu \leq \lambda$ , and  $\lambda = \mu$  as  $\lambda \leq \mu$  and  $\lambda \geq \mu$ .

For normed linear spaces  $\mathscr{Y}$  and  $\mathscr{Z}$ , in this work for convenience over  $\mathbb{R}$ ,  $\mathcal{L}(\mathscr{Y}, \mathscr{Z})$  will denote the space of bounded linear mappings  $\mathscr{Y} \to \mathscr{Z}$  endowed with the operator norm  $\|\cdot\|_{\mathcal{L}(\mathscr{Y},\mathscr{Z})}$ . The subset of invertible operators in

 $\mathcal{L}(\mathscr{Y}, \mathscr{Z})$  with inverses in  $\mathcal{L}(\mathscr{Z}, \mathscr{Y})$  will be denoted as  $\mathcal{L}is(\mathscr{Y}, \mathscr{Z})$ . The *condition* number of a  $C \in \mathcal{L}is(\mathscr{Y}, \mathscr{Z})$  is defined as  $\kappa_{\mathscr{Y}, \mathscr{Z}}(C) := \|C\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Z})} \|C^{-1}\|_{\mathcal{L}(\mathscr{Z}, \mathscr{Y})}$ . For  $\mathscr{Y}$  a reflexive Banach space and  $C \in \mathcal{L}(\mathscr{Y}, \mathscr{Y})$  being *coercive*, i.e.,

$$\inf_{0 \neq y \in \mathscr{Y}} \frac{(Cy)(y)}{\|y\|_{\mathscr{Y}}^2} > 0,$$

both C and  $\Re(C) := \frac{1}{2}(C + C')$  are in  $\mathcal{L}is(\mathscr{Y}, \mathscr{Y}')$  with

$$\|\Re(C)\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')} \leq \|C\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')},$$
  
$$\|C^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} \leq \|\Re(C)^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} = \left(\inf_{0 \neq y \in \mathscr{Y}} \frac{(Cy)(y)}{\|y\|_{\mathscr{Y}}^2}\right)^{-1}.$$

The subset of coercive operators in  $\mathcal{L}is(\mathscr{Y}, \mathscr{Y}')$  is denoted as  $\mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$ . If  $C \in \mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$ , then  $C^{-1} \in \mathcal{L}is_c(\mathscr{Y}', \mathscr{Y})$  and  $\|\Re(C^{-1})^{-1}\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Y}')} \leq \|C\|^2_{\mathcal{L}(\mathscr{Y}, \mathscr{Y})} \|\Re(C)^{-1}\|_{\mathcal{L}(\mathscr{Y}', \mathscr{Y})}$ .

Given a family of operators  $C_i \in \mathcal{L}is(\mathscr{Y}_i, \mathscr{Z}_i)$  ( $\mathcal{L}is_c(\mathscr{Y}_i, \mathscr{Z}_i)$ ), we will write  $C_i \in \mathcal{L}is(\mathscr{Y}_i, \mathscr{Z}_i)$  ( $\mathcal{L}is_c(\mathscr{Y}_i, \mathscr{Z}_i)$ ) uniformly in *i*, or simply 'uniform', when

$$\sup_{i} \max(\|C_i\|_{\mathcal{L}(\mathscr{Y}_i,\mathscr{Z}_i)}, \|C_i^{-1}\|_{\mathcal{L}(\mathscr{Z}_i,\mathscr{Y}_i)}) < \infty,$$

or

$$\sup \max(\|C_i\|_{\mathcal{L}(\mathscr{Y}_i,\mathscr{Z}_i)}, \|\Re(C_i)^{-1}\|_{\mathcal{L}(\mathscr{Z}_i,\mathscr{Y}_i)}) < \infty.$$

# 3.2 Operator preconditioning

Let  $(\mathcal{T})_{\mathcal{T}\in\mathbb{T}}$  be a family of *conforming* partitions of a domain  $\Omega \subset \mathbb{R}^d$  into (open) *uniformly shape regular d*-simplices, where we assume that  $\gamma$  is the (possibly empty) union of (d-1)-faces of  $T \in \mathcal{T}$ . For  $d \geq 2$ , such partitions automatically satisfy a uniform *K*-mesh property, and for d = 1 we impose this as an additional condition. The discussion of the manifold case is postponed to Sect. 3.4.

Recalling that  $\mathscr{V}_{\mathcal{T}} \subset \mathscr{V}$  is a family of piecewise or continuous piecewise polynomials of some fixed degree w.r.t.  $\mathcal{T}$ , let  $A_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}}')$  uniformly in  $\mathcal{T} \in \mathbb{T}$ . A common setting is that  $(A_{\mathcal{T}}v)(\tilde{v}) := (Av)(\tilde{v}) (v, \tilde{v} \in \mathscr{V}_{\mathcal{T}})$  for some  $A \in \mathcal{L}is_c(\mathscr{V}, \mathscr{V}')$ . We are interested in finding optimal *preconditioners*  $G_{\mathcal{T}}$  for  $A_{\mathcal{T}}$ , i.e.,  $G_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}_{\mathcal{T}}', \mathscr{V}_{\mathcal{T}})$  uniformly in  $\mathcal{T} \in \mathbb{T}$ , whose application moreover requires  $\mathcal{O}(\dim \mathscr{V}_{\mathcal{T}})$  arithmetic operations.

Recall the space

$$\mathscr{S}^{0,1}_{\mathcal{T},0} := \{ u \in H^1_{0,\gamma}(\Omega) \colon u|_T \in \mathcal{P}_1 \ (T \in \mathcal{T}) \} \subset \mathscr{W}$$

(thus equipped with  $\|\cdot\|_{\mathscr{W}}$ ). In Chapter 2, using operator preconditioning, we reduced the issue of constructing such preconditioners  $G_{\mathcal{T}}$  to the issue of constructing  $B_{\mathcal{T}}^{\mathscr{S}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  uniformly. In the next section we summarize this reduction.

# 3.2.1 Construction of optimal preconditioners

For the moment, consider the lowest order case of  $\mathscr{V}_{\mathcal{T}}$  being either the space of piecewise constants or continuous piecewise linears. In Chapter 2 a space  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W}$  was constructed with dim  $\mathscr{W}_{\mathcal{T}} = \dim \mathscr{V}_{\mathcal{T}}$  and

(3.1) 
$$\inf_{\mathcal{T}\in\mathbb{T}}\inf_{0\neq v\in\mathscr{V}_{\mathcal{T}}}\sup_{0\neq w\in\mathscr{W}_{\mathcal{T}}}\frac{\langle v,w\rangle_{L_{2}(\Omega)}}{\|v\|_{\mathscr{V}}\|w\|_{\mathscr{W}}}>0.$$

Moreover,  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W}$  was equipped with a locally supported basis  $\Psi_{\mathcal{T}}$  that, modulo a scaling, is  $L_2(\Omega)$ -biorthogonal to the canonical basis  $\Xi_{\mathcal{T}}$  of  $\mathscr{V}_{\mathcal{T}}$ .

As a consequence of (3.1),  $D_{\mathcal{T}}$  defined by  $(D_{\mathcal{T}}v)(w) := \langle v, w \rangle_{L_2(\Omega)}$   $(v \in \mathcal{V}_{\mathcal{T}}, w \in \mathcal{W}_{\mathcal{T}})$  is in  $\mathcal{L}is(\mathcal{V}_{\mathcal{T}}, \mathcal{W}_{\mathcal{T}}')$  uniformly. We infer that once we have constructed  $B_{\mathcal{T}}^{\mathcal{W}} \in \mathcal{L}is(\mathcal{W}_{\mathcal{T}}, \mathcal{W}_{\mathcal{T}}')$  uniformly, then by taking

(3.2) 
$$G_{\mathcal{T}} := D_{\mathcal{T}}^{-1} B_{\mathcal{T}}^{\mathscr{W}} (D_{\mathcal{T}}')^{-1},$$

we have  $G_{\mathcal{T}} \in \mathcal{L}$ is $(\mathscr{V}_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}})$  uniformly. Biorthogonality, modulo a scaling, of the bases  $\Psi_{\mathcal{T}}$  and  $\Xi_{\mathcal{T}}$  implies that the matrix representation of  $D_{\mathcal{T}}$  is diagonal, so that  $D_{\mathcal{T}}^{-1}$  and its adjoint can be applied in linear complexity.

The aforementioned space  $\mathscr{W}_{\mathcal{T}}$  is a subspace of  $\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}} \subset \mathscr{W}$ , where  $\mathscr{B}_{\mathcal{T}}$  is a 'bubble space' with dim  $\mathscr{B}_{\mathcal{T}} = \mathcal{O}(\#\mathcal{T})$ , such that the projector  $I_{\mathcal{T}}$  on  $\mathscr{S}_{\mathcal{T},0}^{0,1} \oplus \mathscr{B}_{\mathcal{T}}$ , defined by ran  $I_{\mathcal{T}} = \mathscr{S}_{\mathcal{T},0}^{0,1}$  and ran $(\mathrm{Id} - I_{\mathcal{T}}) = \mathscr{B}_{\mathcal{T}}$ , is 'local' and uniformly bounded, and the canonical basis  $\Theta_{\mathcal{T}}$  of 'bubbles' for  $\mathscr{B}_{\mathcal{T}}$  is, when normalized in  $\|\cdot\|_{\mathscr{W}}$ , a uniformly Riesz basis for  $\mathscr{B}_{\mathcal{T}}$ . Because of the latter,  $B_{\mathcal{T}}^{\mathscr{B}}$  defined by

$$(B_{\mathcal{T}}^{\mathscr{B}}\mathbf{c}^{\top}\Theta_{\mathcal{T}})(\mathbf{d}^{\top}\Theta_{\mathcal{T}}) \coloneqq \beta(\mathbf{\Delta}_{\mathcal{T}}\mathbf{c})^{\top}\mathbf{d}$$

for some diagonal  $\Delta_{\mathcal{T}} = \operatorname{diag}(\langle \Theta_{\mathcal{T}}, \Theta_{\mathcal{T}} \rangle_{\mathscr{W}})$  and constant  $\beta > 0$  is in  $\mathcal{L}$ is<sub>c</sub>( $\mathscr{B}_{\mathcal{T}}, \mathscr{B}'_{\mathcal{T}}$ ) uniformly.

Given some 'opposite order' operator  $B_{\mathcal{T}}^{\mathscr{S}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$ , by taking

$$(3.3) \qquad \qquad B_{\mathcal{T}}^{\mathscr{W}} := I_{\mathcal{T}}' B_{\mathcal{T}}^{\mathscr{S}} I_{\mathcal{T}} + (\mathrm{Id} - I_{\mathcal{T}})' B_{\mathcal{T}}^{\mathscr{B}} (\mathrm{Id} - I_{\mathcal{T}}),$$

it holds that  $B_{\mathcal{T}}^{\mathscr{W}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  uniformly ([SvV20b, Prop. 5.1]), which makes  $G_{\mathcal{T}}$  a uniform preconditioner.

# **3.2.2** Implementation of $G_T$

Recalling the aforementioned bases  $\Xi_{\mathcal{T}}$ ,  $\Psi_{\mathcal{T}}$ , and  $\Theta_{\mathcal{T}}$  for  $\mathscr{V}_{\mathcal{T}}$ ,  $\mathscr{W}_{\mathcal{T}}$  and  $\mathscr{B}_{\mathcal{T}}$ , respectively, equipping  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  with the nodal basis  $\Phi_{\mathcal{T}}$ , and equipping  $\mathscr{V}_{\mathcal{T}}', \mathscr{W}_{\mathcal{T}}',$  $\mathscr{B}_{\mathcal{T}}'$ , and  $(\mathscr{S}_{\mathcal{T},0}^{0,1})'$  with the dual bases  $\Xi_{\mathcal{T}}', \Psi_{\mathcal{T}}', \Theta_{\mathcal{T}}'$ , and  $\Phi_{\mathcal{T}}'$ , respectively, the representation of  $A_{\mathcal{T}} \in \mathcal{L}(\mathscr{V}_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}}')$  is the stiffness matrix  $A_{\mathcal{T}} := (A_{\mathcal{T}}\Xi_{\mathcal{T}})(\Xi_{\mathcal{T}}) := [(A_{\mathcal{T}}\eta)(\xi)]_{(\xi,\eta)\in\Xi_{\mathcal{T}}'}$ , and the representation of  $G_{\mathcal{T}} \in \mathcal{L}(\mathscr{V}_{\mathcal{T}}', \mathscr{V}_{\mathcal{T}})$  is the matrix  $G_{\mathcal{T}} := (G\Xi_{\mathcal{T}}')(\Xi_{\mathcal{T}}')$ . It is given by

(3.4) 
$$\boldsymbol{G}_{\mathcal{T}} = \boldsymbol{D}_{\mathcal{T}}^{-1} \big( \boldsymbol{p}_{\mathcal{T}}^{\top} \boldsymbol{B}_{\mathcal{T}}^{\mathscr{S}} \boldsymbol{p}_{\mathcal{T}} + \boldsymbol{q}_{\mathcal{T}}^{\top} \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} \boldsymbol{q}_{\mathcal{T}} \big) \boldsymbol{D}_{\mathcal{T}}^{-\top},$$

where both

$$\boldsymbol{D}_{\mathcal{T}} := (D_{\mathcal{T}} \Xi_{\mathcal{T}})(\Psi_{\mathcal{T}}), \quad \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} := (B_{\mathcal{T}}^{\mathscr{B}} \Theta_{\mathcal{T}})(\Theta_{\mathcal{T}})$$

are diagonal, both

$$\boldsymbol{p}_{\mathcal{T}} := (I_{\mathcal{T}} \Psi_{\mathcal{T}})(\Phi_{\mathcal{T}}'), \qquad \boldsymbol{q}_{\mathcal{T}} := ((\mathrm{Id} - I_{\mathcal{T}})\Psi_{\mathcal{T}})(\Theta_{\mathcal{T}}')$$

are *uniformly sparse*, and

$$\mathbf{B}_{\mathcal{T}}^{\mathscr{S}} := (B_{\mathcal{T}}^{\mathscr{S}} \Phi_{\mathcal{T}})(\Phi_{\mathcal{T}}).$$

Note that the cost of the application of  $G_{\mathcal{T}}$  scales *linearly* in  $\#\mathcal{T}$  as soon as this holds true for the application of  $B_{\mathcal{T}}^{\mathscr{S}}$ .

The above preconditioning approach is summarized in the following theorem.

**Theorem 3.2.1** (Sect. 2.3). Given a family  $B_{\mathcal{T}}^{\mathscr{S}} \in \operatorname{Lis}_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  uniformly in  $\mathcal{T} \in \mathbb{T}$ . Then for  $B_{\mathcal{T}}^{\mathscr{W}}$  as described in (3.3), the operator  $G_{\mathcal{T}}$  from (3.2) is a uniform preconditioner. Furthermore, if the matrix representation  $B_{\mathcal{T}}^{\mathscr{T}}$ , cf. (3.5), can be applied in  $\mathcal{O}(\#\mathcal{T})$  operations, then the matrix representation of the preconditioner  $G_{\mathcal{T}}$ , cf. (3.4), can be applied in  $\mathcal{O}(\#\mathcal{T})$  operations.

Because  $B_{\mathcal{T}}^{\mathscr{W}}$  in (3.3) is given as the sum of two operators that 'act' on different subspaces of  $\mathscr{W}_{\mathcal{T}}$ , the condition number of the preconditioned system depends on the relative scaling of both these operators which can be steered by selecting the parameter  $\beta$ . A suitable  $\beta$  will be selected experimentally.

Alternatively, Proposition 4.4.1 shows that a value of  $\beta$  is reasonable if it is chosen such that the interval bounded by the coercivity and boundedness constants of  $B_{\mathcal{T}}^{\mathscr{S}}$  is included in that interval corresponding to  $B_{\mathcal{T}}^{\mathscr{B}}$  or vice versa. Also these coercivity and boundedness constants can be approximated experimentally or by making some theoretical estimates.

Constructions of  $\Psi_{\mathcal{T}}$ ,  $\Theta_{\mathcal{T}}$ , and  $\Delta_{\mathcal{T}}$ , and resulting explicit formulas for matrices  $D_{\mathcal{T}}$ ,  $B_{\mathcal{T}}^{\mathscr{B}}$ ,  $p_{\mathcal{T}}$ ,  $q_{\mathcal{T}}$  are derived in Chapter 2. For ease of reading we recall these formulas below for the case that  $\mathscr{V}_{\mathcal{T}}$  is the space of *piecewise constants*. For the continuous piecewise linear case we refer to Sect. 2.4.2.

#### Piecewise constant trial space $\mathscr{V}_{\mathcal{T}}$

For  $\mathcal{T} \in \mathbb{T}$ , we define  $N_{\mathcal{T}}$  as the set of vertices of  $\mathcal{T}$ , and  $N_{\mathcal{T}}^0$  as the set of vertices of  $\mathcal{T}$  that are not on  $\gamma$ . For  $\nu \in N_{\mathcal{T}}$  we set its *valence* 

$$d_{\mathcal{T},\nu} := \#\{T \in \mathcal{T} \colon \nu \in \overline{T}\}.$$

For  $T \in \mathcal{T}$ , and with  $N_T$  denoting the set of its vertices, we set  $N_{\mathcal{T},T}^0 := N_{\mathcal{T}}^0 \cap N_T$ .

If one considers  $\mathscr{V}_{\mathcal{T}}$  as the space of discontinuous piecewise constants, i.e.

$$\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,0} := \{ u \in L_2(\Omega) \colon u |_T \in P_0(T \in \mathcal{T}) \} \subset \mathscr{V},$$

equipped with the canonical basis  $\Xi_{\mathcal{T}} := \{\mathbb{1}_T : T \in \mathcal{T}\}$ , then we find, for arbitrary constant  $\beta > 0$ ,

$$\begin{aligned} \boldsymbol{D}_{\mathcal{T}} &= \operatorname{diag}\{|T| \colon T \in \mathcal{T}\}, \qquad (\boldsymbol{p}_{\mathcal{T}})_{\nu T} = \begin{cases} d_{\mathcal{T},\nu}^{-1} & \text{if } \nu \in N_{\mathcal{T},T}^{0}, \\ 0 & \text{if } \nu \notin N_{\mathcal{T},T}^{0}, \end{cases} \\ \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} &= \beta \boldsymbol{D}_{\mathcal{T}}^{1-\frac{2s}{d}}, \qquad (\boldsymbol{q}_{\mathcal{T}})_{T'T} = \delta_{T'T} - \frac{1}{d+1} \sum_{\nu \in N_{\mathcal{T},T}^{0} \cap N_{\mathcal{T},T'}^{0}} d_{\mathcal{T},\nu}^{-1}. \end{aligned}$$

# 3.2.3 Higher order case

For higher order discontinuous or continuous finite element spaces  $\mathscr{V}_{\mathcal{T}}$ , suitable preconditioners  $G_{\mathcal{T}}$  can be built either from the current preconditioner  $G_{\mathcal{T}}$  for the lowest order case by application of a subspace correction method (most conveniently in the discontinuous case where on each element the space of polynomials of some fixed degree is split into the space of constants and its orthogonal complement), or by expanding  $\mathscr{W}_{\mathcal{T}}$  by enlarging the bubble space  $\mathscr{B}_{\mathcal{T}}$ . While referring to Chapter 2 for details, we recall that with either option the construction of an optimal preconditioner  $G_{\mathcal{T}}$  that can be applied in linear complexity hinges on the availability of an operator  $B_{\mathcal{T}}^{\mathscr{L}} \in \operatorname{Lis}_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  uniformly in  $\mathcal{T} \in \mathbb{T}$ , that can be applied in linear complexity.

# **3.3** An operator $B_T^{\mathscr{S}}$ of multi-level type

In this section we will introduce an operator  $B_{\mathcal{T}}^{\mathscr{S}} \in \operatorname{Lis}_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  of multi-level type. The operator  $B_{\mathcal{T}}^{\mathscr{S}}$  is based on a stable multi-level decomposition of  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  given by Wu and Zheng [WZ17]. Usually such a stable multi-level decomposition is used as a theoretical tool for proving optimality of an additive (or multiplicative) Schwarz type preconditioner for an operator in  $\operatorname{Lis}_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$ . In this work, however, we are going to use their results for the construction of the operator  $B_{\mathcal{T}}^{\mathscr{S}} \in \operatorname{Lis}_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  for which it is crucial that its application can be implemented in linear complexity.

# **3.3.1** Definition and analysis of $B_T^{\mathscr{S}}$

For  $d \ge 2$ , let  $\mathbb{T}$  be the family of all conforming partitions of  $\Omega$  into *d*-simplices that can be created by Newest Vertex Bisection starting from some given conforming initial partition  $\mathcal{T}_{\perp}$  that satisfies a *matching condition* ([Ste08b]).

With  $\mathfrak{T} := \bigcup_{\mathcal{T} \in \mathbb{T}} \{T : T \in \mathcal{T}\}$  and  $\mathfrak{N} := \bigcup_{\mathcal{T} \in \mathbb{T}} N_{\mathcal{T}}$ , for  $T \in \mathfrak{T}$  let gen(T) be the number of bisections needed to create T from its ancestor  $T' \in \mathcal{T}_{\perp}$ , and for

 $\nu \in \mathfrak{N}$  let  $gen(\nu) := min\{gen(T) : T \in \mathfrak{T}, \nu \in N_T\}$ . Notice that  $|T| = 2^{-gen(T)}$ . For  $T \in \mathfrak{T}$ , let  $Q_T$  denote the  $L_2(T)$ -orthogonal projector onto  $\mathcal{P}_1(T)$ .

The case d = 1 can be included by letting  $\mathbb{T}$  be the family of a partitions of  $\Omega$  that can be constructed by bisections from  $\mathcal{T}_{\perp} = \{\Omega\}$  such that the generations of any two neighbouring subintervals in any  $\mathcal{T} \in \mathbb{T}$  differ by not more than one.

For  $\mathcal{T} \in \mathbb{T}$ , set  $L = L(\mathcal{T}) := \max_{T \in \mathcal{T}} \operatorname{gen}(T)$  and define

$$\mathcal{T}_{\perp} = \mathcal{T}_0 \prec \mathcal{T}_1 \prec \cdots \prec \mathcal{T}_L = \mathcal{T} \subset \mathbb{T}$$

by constructing  $\mathcal{T}_{j-1}$  from  $\mathcal{T}_j$  by removing all  $\nu \in N_j := N_{\mathcal{T}_j}$  from the latter for which gen $(\nu) = j$ . For  $\nu \in N_j^0 := N_{\mathcal{T}_j}^0$ , we define  $\omega_j(\nu) = \bigcup \{T \in \mathcal{T}_j : \nu \in N_T\}$ .

With this hierarchy of partitions, we define an averaging quasi-interpolator  $\Pi_j \in \mathcal{L}(\mathscr{S}^{0,1}_{\mathcal{T},0}, \mathscr{S}^{0,1}_{\mathcal{T}_j,0})$  by

(3.6) 
$$(\Pi_{j}u)(\nu) := \frac{\sum_{\{T \in \mathcal{T}_{j} : \nu \in N_{T}\}} |T|(Q_{T}u)(\nu)}{\sum_{\{T \in \mathcal{T}_{j} : \nu \in N_{T}\}} |T|} \quad (u \in \mathscr{S}_{\mathcal{T},0}^{0,1}, \nu \in N_{j}^{0}).$$

Since  $\mathscr{S}_{\mathcal{T}_{j},0}^{0,1}$  is a space of continuous piecewise linears, it indeed suffices to define  $\Pi_{j}u$  at the vertices  $N_{j}^{0}$ . Recall that  $\mathscr{S}_{\mathcal{T},0}^{0,1} \subset \mathscr{W} := [L_{2}(\Omega), H_{0,\gamma}^{1}(\Omega)]_{s,2}$  for some  $s \in [0,1]$ . The next theorem shows that  $\Pi_{j}$  induces a stable multi-level decomposition of  $\mathscr{S}_{\mathcal{T},0}^{0,1}$ .

**Theorem 3.3.1** ([WZ17, Lemma 3.7]). For the averaging quasi-interpolator  $\Pi_j$  from (3.6), and  $\Pi_{-1} := 0$ , it holds that

$$\|u\|_{\mathscr{W}}^2 \approx \sum_{j=0}^{L} 4^{js/d} \|(\Pi_j - \Pi_{j-1})u\|_{L_2(\Omega)}^2 \quad (u \in \mathscr{S}_{\mathcal{T},0}^{0,1}).$$

*Proof.* In [WZ17], the inequality ' $\gtrsim$ ' was proven for the case  $s = 1, d \in \{2, 3\}$ , and  $\gamma = \partial \Omega$ . The arguments, however, immediately extend to  $s \in [0, 1], d \ge 1$ , and  $\gamma \subsetneq \partial \Omega$ .

The proof of the other inequality ' $\lesssim$ ' follows from well-known arguments: For some  $t \in (1, \frac{3}{2})$ , let  $\mathscr{H}^{rt} := [L_2(\Omega), H^1_{0,\gamma}(\Omega) \cap H^t(\Omega)]_{r,2}$  for  $r \in [0, 1]$ . Then  $\mathscr{H}^s \simeq \mathscr{W}$  by the *reiteration theorem*, and for  $r \in [0, 1]$ ,  $\|\cdot\|_{\mathscr{H}^{rt}} \lesssim 2^{jrt/d} \|\cdot\|_{L_2(\Omega)}$ on  $\mathscr{F}^{0,1}_{\mathcal{T}_j,0}$ .

Let  $u \in \mathscr{S}_{\mathcal{T},0}^{0,1}$  be written as  $\sum_{j=0}^{L} u_j$  with  $u_j \in \mathscr{S}_{\mathcal{T}_j,0}^{0,1}$ . Then for  $\varepsilon \in (0,s)$ ,  $\varepsilon \leq t-s$ , we have

$$\begin{aligned} \|u\|_{\mathscr{H}^{s}(\Omega)}^{2} &\lesssim \sum_{j=0}^{L} \sum_{i=j}^{L} \|u_{j}\|_{\mathscr{H}^{s+\varepsilon}(\Omega)} \|u_{i}\|_{\mathscr{H}^{s-\varepsilon}(\Omega)} \\ &\lesssim \sum_{j=0}^{L} \sum_{i=j}^{L} 2^{j(s+\varepsilon)/d} 2^{i(s-\varepsilon)/d} \|u_{j}\|_{L_{2}(\Omega)} \|u_{i}\|_{L_{2}(\Omega)} \lesssim \sum_{j=0}^{L} 4^{js/d} \|u_{j}\|_{L_{2}(\Omega)}^{2}. \ \end{aligned}$$

55



FIGURE 3.1. For d = 3, a tetrahedron  $T \in \mathcal{T}_{j-1}$  and its bisection. The dots indicate all vertices in  $N_i^0 \setminus M_j^0$ .

The relevance of the multi-level decomposition from Theorem 3.3.1 by Wu and Zheng lies in the fact that  $(\Pi_j u)(\nu)$  can only differ from  $(\Pi_{j-1}u)(\nu)$  in any  $\nu \in N_j^0 \setminus N_{j-1}^0$  as well as in only  $two^1$  of its neighbours in  $N_{j-1}^0$  (the endpoints of the edge on which  $\nu$  was inserted):

**Proposition 3.3.2** ([WZ17, Lemma 3.1]). *With for*  $j \ge 1$ ,

$$M_j^0 := \{ \nu \in N_{j-1}^0 \colon \omega_j(\nu) = \omega_{j-1}(\nu) \},\$$

it holds that for  $\nu \in M_{j}^{0}$ ,  $((\Pi_{j} - \Pi_{j-1})u)(\nu) = 0$ , see Figure 3.1.

*Remark* 3.3.3. The proof from [WZ17] given for  $d \in \{2, 3\}$  generalizes to  $d \ge 1$ . Indeed the arguments that are used are based on the fact that the basis for  $S_1(T)$  that is dual to the nodal basis takes equal values in all but one nodal point. This is a consequence of the fact that the mass matrix of the nodal basis for  $S_1(T)$ , and so its inverse, is invariant under permutations of the barycentric coordinates, which holds true in any dimension.

As a consequence of Proposition 3.3.2, we have

$$\|(\Pi_j - \Pi_{j-1})u\|_{L_2(\Omega)}^2 \approx 2^{-j} \sum_{\nu \in N_j^0 \setminus M_j^0} |((\Pi_j - \Pi_{j-1})u)(\nu)|^2.$$

From Theorem 3.3.1, we conclude that  $B_{\mathcal{T}}^{\mathscr{S}} = (B_{\mathcal{T}}^{\mathscr{S}})' \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  defined by

(3.7) 
$$(B_{\mathcal{T}}^{\mathscr{S}}u)(v) := \sum_{j=0}^{L} 2^{j(\frac{2s}{d}-1)} \sum_{\nu \in N_{j}^{0} \setminus M_{j}^{0}} ((\Pi_{j} - \Pi_{j-1})u)(\nu)((\Pi_{j} - \Pi_{j-1})v)(\nu)$$

is uniform, i.e.

$$\sup_{\mathcal{T}\in\mathbb{T}} \max\left( \|B_{\mathcal{T}}^{\mathscr{S}}\|_{\mathcal{L}(\mathscr{S}_{\mathcal{T},0}^{0,1},(\mathscr{S}_{\mathcal{T},0}^{0,1})')}, \|(B_{\mathcal{T}}^{\mathscr{S}})^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},0}^{0,1})',\mathscr{S}_{\mathcal{T},0}^{0,1})}\right) < \infty.$$

<sup>&</sup>lt;sup>1</sup>As pointed out in [WZ17], for  $d \ge 3$  the number of such neighbours will be larger when employing the Scott-Zhang quasi-interpolator. Moreover, this interpolator is not suited for  $s \le \frac{1}{2}$ .

# **3.3.2** Implementation of $B_T^{\mathscr{S}}$

Since the operator  $\Pi_j$  is a weighted local  $L_2(\Omega)$  projection, it allows for a natural implementation by considering  $\mathscr{S}_{\mathcal{T}}^{-1,1}$ , the space of discontinuous piecewise linears w.r.t.  $\mathcal{T}$ . Recall the nodal basis  $\Phi_{\mathcal{T}}$  for  $\mathscr{S}_{\mathcal{T},0}^{0,1}$ , and equip  $\mathscr{S}_{\mathcal{T}}^{-1,1}$  with the element-wise nodal basis.

Denote  $E_{\mathcal{T}}$  for the representation of the embedding  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  into  $\mathscr{S}_{\mathcal{T}}^{-1,1}$ .

For  $0 \leq j \leq L$ , let  $\mathbf{R}_j$  be the representation of the  $L_2(\Omega)$ -orthogonal projector of  $\mathscr{S}_{\mathcal{T}}^{-1,1}$  onto  $\mathscr{S}_{\mathcal{T}_j}^{-1,1}$ , and let  $\mathbf{R}_{-1} := 0$ .

For  $0 \leq j \leq L$ , let  $H_j^{'_j}$  be the representation of the averaging operator  $H_j: \mathscr{S}_{\mathcal{T}_j}^{-1,1} \to \mathscr{S}_{\mathcal{T}_j,0}^{0,1}$  defined by

(3.8) 
$$(H_j u)(\nu) = \frac{\sum_{\{T \in \mathcal{T}_j : \nu \in N_T\}} |T| \, u|_T(\nu)}{\sum_{\{T \in \mathcal{T}_j : \nu \in N_T\}} |T|}, \quad (\nu \in N_{\mathcal{T}}^0).$$

and let  $H_{-1} \coloneqq 0$ .

For  $1 \leq j \leq L$ , let  $P_j$  be the representation of the embedding  $\mathscr{S}_{\mathcal{T}_{j-1},0}^{0,1} \rightarrow \mathscr{S}_{\mathcal{T}_{j,0}}^{0,1}$  (often called *prolongation*), and let  $P_0 := 0$ .

Then the representation  $B_T^{\mathscr{S}}$  of  $B_T^{\mathscr{S}}$  from (3.7) is given by

$$\boldsymbol{B}_{\mathcal{T}}^{\mathscr{S}} = \boldsymbol{E}_{\mathcal{T}}^{\top} \Big( \sum_{j=0}^{L} (\boldsymbol{H}_{j} \boldsymbol{R}_{j} - \boldsymbol{P}_{j} \boldsymbol{H}_{j-1} \boldsymbol{R}_{j-1})^{\top} 2^{j(\frac{2s}{d}-1)} (\boldsymbol{H}_{j} \boldsymbol{R}_{j} - \boldsymbol{P}_{j} \boldsymbol{H}_{j-1} \boldsymbol{R}_{j-1}) \Big) \boldsymbol{E}_{\mathcal{T}}.$$

Applying  $E_{\mathcal{T}}$  amounts to duplicating values at any internal node with a number equal to the valence of that node.

By representing  $\mathcal{T}$  as the leaves of a binary tree with roots being the simplices of  $\mathcal{T}_{\perp}$ , computing for  $\mathbf{x} \in \operatorname{ran} \mathbf{E}_{\mathcal{T}}$  the sequence  $(\mathbf{R}_j \mathbf{x})_{0 \leq j \leq L}$  amounts to computing, while traversing from the leaves to the root, for any parent and both its children the orthogonal projection of a piecewise linear function on the children to the space of linears on the parent. For d = 2, the matrix representation of the latter projection is given in Figure 3.2.

# **Proposition 3.3.4.** The application of $B_{\mathcal{T}}^{\mathscr{S}}$ can be computed in $\mathcal{O}(\#\mathcal{T})$ operations.

*Proof.* Because the number of nodes in a binary tree is less than 2 times the number of its leaves, for  $\mathbf{x} \in \mathbb{R}^{\dim \mathscr{S}_{T,0}^{0,1}}$  the computation of the sequence  $(\mathbf{R}_j \mathbf{E}_T \mathbf{x})_{0 \leq j \leq L}$  takes  $\mathcal{O}(\#\mathcal{T})$  operations. From Proposition 3.3.2 recall that any vector in ran  $\mathbf{H}_j \mathbf{R}_j - \mathbf{P}_j \mathbf{H}_{j-1} \mathbf{R}_{j-1}$  vanishes at  $M_j^0$ , so that the number of its non-zero entries is bounded by  $\#(N_j^0 \setminus M_j^0) \leq 3\#(N_j^0 \setminus N_{j-1}^0)$ . Knowing already  $\mathbf{R}_j \mathbf{E}_T \mathbf{x}$  and  $\mathbf{R}_{j-1} \mathbf{E}_T \mathbf{x}$ , computing any non-zero entry of  $(\mathbf{H}_j \mathbf{R}_j - \mathbf{P}_j \mathbf{H}_{j-1} \mathbf{R}_{j-1}) \mathbf{E}_T \mathbf{x}$  requires  $\mathcal{O}(1)$  operations.

We conclude that the operator  $B_{\mathcal{T}}^{\mathscr{S}}$ , with above matrix representation  $B_{\mathcal{T}}^{\mathscr{S}}$ , satisfies the requirements of Theorem 3.2.1.



FIGURE 3.2. Numbering of the vertices of the parent and that of both children for d = 2, and the resulting matrix representation of the orthogonal projection of the space of piecewise linears on the children to the space of linears on the parent.

# 3.4 Manifold case

Let  $\Gamma$  be a compact *d*-dimensional Lipschitz, piecewise smooth manifold in  $\mathbb{R}^{d'}$  for some  $d' \geq d$  with or without boundary  $\partial \Gamma$ . For some closed measurable  $\gamma \subset \partial \Gamma$  and  $s \in [0, 1]$ , let

$$\mathscr{W} := [L_2(\Gamma), H^1_{0,\gamma}(\Gamma)]_{s,2}, \quad \mathscr{V} := \mathscr{W}'.$$

We assume that  $\Gamma$  is given as the essentially disjoint union of  $\bigcup_{i=1}^{p} \overline{\chi_i(\Omega_i)}$ , with, for  $1 \leq i \leq p, \chi_i \colon \mathbb{R}^d \to \mathbb{R}^{d'}$  being some smooth regular parametrization, and  $\Omega_i \subset \mathbb{R}^d$  an open polytope. W.l.o.g. assuming that for  $i \neq j, \overline{\Omega}_i \cap \overline{\Omega}_j = \emptyset$ , we define

$$\chi \colon \Omega := \bigcup_{i=1}^p \Omega_i \to \bigcup_{i=1}^p \chi_i(\Omega_i) \text{ by } \chi|_{\Omega_i} = \chi_i.$$

Let  $\mathbb{T}$  be a family of conforming partitions  $\mathcal{T}$  of  $\Gamma$  into 'panels' such that, for  $1 \leq i \leq p, \chi^{-1}(\mathcal{T}) \cap \Omega_i$  is a uniformly shape regular conforming partition of  $\Omega_i$  into *d*-simplices (that for d = 1 satisfies a uniform *K*-mesh property). We assume that  $\gamma$  is a (possibly empty) union of 'faces' of  $T \in \mathcal{T}$  (i.e., sets of type  $\chi_i(e)$ , where *e* is a (d-1)-dimensional face of  $\chi_i^{-1}(T)$ ).

We set

$$\mathscr{V}_{\mathcal{T}} := \{ u \in L_2(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_0 \ (T \in \mathcal{T}) \} \subset \mathscr{V},$$

or

$$\mathscr{V}_{\mathcal{T}} := \{ u \in C(\Gamma) \colon u \circ \chi |_{\chi^{-1}(T)} \in \mathcal{P}_1 \ (T \in \mathcal{T}) \} \subset \mathscr{V},$$

equipped with canonical basis  $\Xi_T$ , and, for the construction of a preconditioner,

$$\mathscr{S}_{\mathcal{T},0}^{0,1} := \{ u \in H^1_{0,\gamma}(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_1 \ (T \in \mathcal{T}) \} \subset \mathscr{W},$$

equipped with canonical basis  $\Phi_{\mathcal{T}}$ .

As in the domain case, a space  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W}$  can be constructed with  $\dim \mathscr{W}_{\mathcal{T}} = \dim V_{\mathcal{T}}$  and  $\inf_{\mathcal{T} \in \mathbb{T}} \inf_{0 \neq v \in \mathscr{V}_{\mathcal{T}}} \sup_{0 \neq w \in \mathscr{W}_{\mathcal{T}}} \frac{\langle v, w \rangle_{L_2(\Gamma)}}{\|v\|_{\mathscr{V}} \|w\|_{\mathscr{W}}} > 0$ , which can be equipped with a locally supported basis  $\Psi_{\mathcal{T}}$  that, modulo a scaling, is  $L_2(\Gamma)$ -biorthogonal to  $\Xi_{\mathcal{T}}$ . Now assuming that a family of  $B_{\mathcal{T}}^{\mathscr{S}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  uniformly is available, the construction of an optimal preconditioner  $G_{\mathcal{T}}$  follows exactly the same lines as outlined in Sect. 3.2 for the domain case.

For the case that  $\Gamma$  is not piecewise polytopal, a hidden problem is, however, that above construction of  $\Psi_{\mathcal{T}}$  requires exact integration of lifted polynomials over the manifold. To circumvent this problem, in Sect. 2.3.2 we have relaxed the condition of  $L_2(\Gamma)$ -biorthogonality of  $\Xi_{\mathcal{T}}$  and  $\Psi_{\mathcal{T}}$  to biorthogonality w.r.t. to a mesh-dependent scalar product obtained from the  $L_2(\Gamma)$ -scalar product by replacing the Jacobian on the pull back of each panel by its mean. It was shown that the resulting preconditioner is still optimal, and that the expression for its matrix representation (for the moment without the representation of  $B_{\mathcal{T}}^{\mathscr{S}}$ ), that was recalled in Sect. 3.2.2 for the piecewise constant case, applies verbatim by only reading |T| as the volume of the panel.<sup>2</sup>

It remains to discuss the construction of an operator  $B_{\mathcal{T}}^{\mathscr{P}}$  of multi-level type, where it is now assumed that  $\mathbb{T}$  is a family corresponding to newest vertex bisection. An exact copy of the construction of  $B_{\mathcal{T}}^{\mathscr{P}}$  given in the domain case would require the application of the panel-wise  $L_2(T)$ -orthogonal projector  $Q_T$ , cf. (3.6), which generally poses a quadrature problem. Reconsidering the domain case, the proof of [WZ17, Lemma 3.7] (which provides the proof of the inequality ' $\gtrsim$ ' in our Theorem 3.3.1) builds on the fact that for  $\mathcal{T}_0 \prec \mathcal{T}_1 \prec \cdots$ being a sequence of uniformly refined partitions, the decomposition  $\mathscr{P}_{\mathcal{T}_L,0}^{0,1} =$  $\sum_{j=0}^{L} \mathscr{P}_{\mathcal{T}_j,0}^{0,1} \cap (\mathscr{P}_{\mathcal{T}_{j-1},0}^{0,1})^{\perp_{L_2(\Omega)}}$ , where  $\mathscr{P}_{\mathcal{T}_{-1},0}^{0,1} := \{0\}$ , is stable, uniformly in L, w.r.t. the norm on  $\mathscr{W}$ . This stability holds also true when the orthogonal complements are taken w.r.t. a weighted  $L_2(\Omega)$ -scalar product, for any weight w with w,  $1/w \in L_{\infty}(\Omega)$ .

This has the consequence that for the construction of the multi-level operator  $B_{\tau}^{\mathscr{S}}$  in the manifold case, we may equip  $L_2(\Gamma)$  with scalar product

$$\sum_{i=1}^p \int_{\Omega_i} u(\chi_i(x)) v(\chi_i(x)) \, dx,$$

which is constructed from the canonical  $L_2(\Gamma)$ -scalar product by simply omitting the Jacobians  $|\partial \chi_i(x)|$ . With this modified scalar product, the panel-wise orthogonal projector  $Q_T$  is the same as in the domain case. We conclude that the resulting  $B_{\mathcal{T}}^{\mathscr{S}}$  as in (3.7) is in  $\mathcal{L}$ is<sub> $c</sub>(<math>\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})'$ ) uniformly, and that its application can be performed in linear complexity. Indeed, its implementation is equal as in the domain case as described in Sect. 3.3.2 when |T| in (3.8) is read as  $|\chi^{-1}(T)|$ .</sub>

<sup>&</sup>lt;sup>2</sup>In order to avoid the exact computation of this volume, actually it may read as  $|\chi^{-1}(T)||\partial\chi(z)|$  for arbitrary  $z \in \chi^{-1}(T)$ .

#### 3.5 Numerical experiments

Let  $\Gamma = \partial [0,1]^3 \subset \mathbb{R}^3$  be the two-dimensional manifold without boundary given as the boundary of the unit cube,  $\mathscr{W} := H^{1/2}(\Gamma)$ ,  $\mathscr{V} := H^{-1/2}(\Gamma)$ . We consider the trial space  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,0} \subset \mathscr{V}$  of discontinuous piecewise constants. We will evaluate preconditioning of the discretized single layer operator  $A_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{V}_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}}')$ .

The role of the opposite order operator in  $\mathcal{L}is_c(\mathscr{S}_{\mathcal{T},0}^{0,1}, (\mathscr{S}_{\mathcal{T},0}^{0,1})')$  from Sect. 3.2.1 will be fulfilled by the multi-level operator  $B_{\mathcal{T}}^{\mathscr{S}}$  from (3.7). Equipping  $\mathscr{S}_{\mathcal{T},0}^{0,1}$  with the nodal basis  $\Phi_{\mathcal{T}}$ , the matrix representation of the preconditioner  $G_{\mathcal{T}}$  from Sect. 3.2.1 reads as

$$\boldsymbol{G}_{\mathcal{T}} = \boldsymbol{D}_{\mathcal{T}}^{-1} \big( \boldsymbol{p}_{\mathcal{T}}^{\top} \boldsymbol{B}_{\mathcal{T}}^{\mathscr{S}} \boldsymbol{p}_{\mathcal{T}} + \beta \boldsymbol{q}_{\mathcal{T}}^{\top} \boldsymbol{D}_{\mathcal{T}}^{1/2} \boldsymbol{q}_{\mathcal{T}} \big) \boldsymbol{D}_{\mathcal{T}}^{-1},$$

for  $D_{\mathcal{T}} = \text{diag}\{|T|: T \in \mathcal{T}\}$ , uniformly sparse  $p_{\mathcal{T}}$  and  $q_{\mathcal{T}}$  as given in Sect. 3.2.1, and with the representation of the multi-level operator  $B_{\mathcal{T}}^{\mathscr{S}}$  given by

$$\boldsymbol{B}_{\mathcal{T}}^{\mathscr{S}} = \boldsymbol{E}_{\mathcal{T}}^{\top} \Big( \sum_{j=0}^{L} (\boldsymbol{H}_{j} \boldsymbol{R}_{j} - \boldsymbol{P}_{j} \boldsymbol{H}_{j-1} \boldsymbol{R}_{j-1})^{\top} 2^{-j/2} (\boldsymbol{H}_{j} \boldsymbol{R}_{j} - \boldsymbol{P}_{j} \boldsymbol{H}_{j-1} \boldsymbol{R}_{j-1}) \Big) \boldsymbol{E}_{\mathcal{T}},$$

for the representations  $E_{\mathcal{T}}$ ,  $H_j$ ,  $R_j$  and  $P_j$  as provided in Sect. 3.3.2 (the minor adaptations in the manifold case described in Sect. 3.4 to the matrix representations from Sections 3.2.1 and 3.3.2 vanish in the current simple case).

The BEM++ software package [SBA<sup>+</sup>15] is used to approximate the matrix representation of the discretized single layer operator  $A_{T}$  by hierarchical matrices based on adaptive cross approximation [Hac99, Beb00].

Equipping  $\mathscr{V}_{\mathcal{T}}$  and  $\mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}}$  with 'energy-norms'  $\sqrt{(A_{\mathcal{T}} \cdot)(\cdot)}$  or  $\|A_{\mathcal{T}}^{\frac{1}{2}} \cdot\|$ , respectively, we calculated the (spectral) condition numbers  $\kappa_{\mathcal{L}}(\mathscr{V}_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}})(G_{\mathcal{T}}A_{\mathcal{T}}) = \kappa_{\mathcal{L}}(\mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}}, \mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}})(G_{\mathcal{T}}A_{\mathcal{T}}) = \rho(G_{\mathcal{T}}A_{\mathcal{T}})\rho((G_{\mathcal{T}}A_{\mathcal{T}})^{-1})$ , where  $\rho(\cdot)$  is the spectral radius, using the Lanczos method.

As initial partition  $\mathcal{T}_{\perp} = \mathcal{T}_1$  of  $\Gamma$  we take a conforming partition consisting of 2 triangles per side, so 12 triangles in total, with an assignment of the newest vertices that satisfies the matching condition. We fixed  $\beta = 5.3$ , being the value for which, for a relative small uniform refinement  $\mathcal{T}$  of  $\mathcal{T}_{\perp}$ , we found  $\rho(D_{\mathcal{T}}^{-1}p_{\mathcal{T}}^{\top}B_{\mathcal{T}}^{\mathscr{T}}p_{\mathcal{T}}D_{\mathcal{T}}^{-1}A_{\mathcal{T}}) = \rho(D_{\mathcal{T}}^{-1}\beta q_{\mathcal{T}}^{\top}D_{\mathcal{T}}^{-1}T^{2}q_{\mathcal{T}}D_{\mathcal{T}}^{-1}A_{\mathcal{T}}).$ 

#### 3.5.1 Uniform refinements

Here we let  $\mathbb{T}$  be the sequence  $\{\mathcal{T}_k\}_{k\geq 1}$  of (conforming) uniform refinements, that is,  $\mathcal{T}_k \succ \mathcal{T}_{k-1}$  is found by bisecting each triangle from  $\mathcal{T}_{k-1}$  into 2 subtriangles using Newest Vertex Bisection.

Table 3.1 shows the condition numbers of the preconditioned system in this situation. The condition numbers are relatively small, and the timing results show that the implementation of the preconditioner is indeed linear.

TABLE 3.1. Spectral condition numbers of the preconditioned single layer system discretized by piecewise constants  $\mathscr{S}_{\mathcal{T}}^{-1,0}$ , using uniform refinements. Preconditioner  $G_{\mathcal{T}}$  is constructed using the multi-level operator with  $\beta = 5.3$ . The last column indicates the number of seconds per degree of freedom per application of  $G_{\mathcal{T}}$ .

dofs	$\kappa_S(oldsymbol{A}_\mathcal{T})$	$\kappa_S(oldsymbol{G}_\mathcal{T}oldsymbol{A}_\mathcal{T})$	sec / dof
12	14.5	2.6	$2.6\cdot 10^{-5}$
48	31.0	2.7	$1.4 \cdot 10^{-5}$
192	59.9	2.8	$4.9\cdot 10^{-6}$
768	118.7	3.3	$1.4\cdot10^{-6}$
3072	234.6	3.8	$6.3\cdot10^{-7}$
12288	450.4	4.1	$3.3\cdot10^{-6}$
49152	852.5	4.3	$6.5 \cdot 10^{-7}$
196608	1566.4	4.5	$7.3 \cdot 10^{-7}$
786432	2730.5	4.6	$7.8 \cdot 10^{-7}$

# 3.5.2 Local refinements

Here we take  $\mathbb{T}$  as a sequence  $\{\mathcal{T}_k\}_{k\geq 1}$  of (conforming) locally refined partitions, where  $\mathcal{T}_k \succ \mathcal{T}_{k-1}$  is constructed by applying Newest Vertex Bisection to all triangles in  $\mathcal{T}_{k-1}$  that touch a corner of the cube.

Table 3.2 contains results for the preconditioned single layer operator discretized by piecewise constants  $\mathscr{S}_{\mathcal{T}}^{-1,0}$ . The preconditioned condition numbers are nicely bounded, and the timing results confirm that our implementation of the preconditioner is of linear complexity, also in the case of locally refined partitions.

TABLE 3.2. Spectral condition numbers of the preconditioned single layer system discretized by piecewise constants  $\mathscr{S}_{\mathcal{T}}^{-1,0}$ , using local refinements at each of the eight cube corners. Operator  $G_{\mathcal{T}}$  is applied using the multi-level operator with  $\beta = 5.3$ . The second column is defined by  $h_{\mathcal{T},min} := \min_{T \in \mathcal{T}} \sqrt{|T|}$ . The last column indicates the number of seconds per degree of freedom per application of  $G_{\mathcal{T}}$ .

dofs	$h_{\mathcal{T},min}$	$\kappa_S({oldsymbol G}_{\mathcal T}{oldsymbol A}_{\mathcal T})$	sec / dof
12	$1.4\cdot 10^0$	2.63	$2.5 \cdot 10^{-5}$
336	$8.8\cdot10^{-2}$	2.73	$2.4\cdot10^{-6}$
720	$5.5\cdot10^{-3}$	2.91	$1.8\cdot10^{-6}$
1104	$3.4\cdot10^{-4}$	2.96	$1.8\cdot10^{-6}$
1488	$2.1 \cdot 10^{-5}$	2.99	$2.2\cdot 10^{-6}$
1872	$1.3 \cdot 10^{-6}$	2.98	$2.0\cdot 10^{-6}$
2256	$8.4 \cdot 10^{-8}$	3.00	$2.3\cdot 10^{-6}$
2640	$5.2 \cdot 10^{-9}$	3.00	$2.0\cdot 10^{-6}$
3024	$3.2\cdot10^{-10}$	3.01	$2.3\cdot 10^{-6}$
3408	$2.0 \cdot 10^{-11}$	3.01	$2.5\cdot10^{-6}$
3696	$2.5\cdot10^{-12}$	3.01	$2.6 \cdot 10^{-6}$

# 4.1 Introduction

This chapter deals with the construction of *uniform* preconditioners for operators of *positive* order, using the framework of 'operator preconditioning' as described in [Hip06], see e.g. [SW98, CN00] for earlier work. It will build on our experiences with this approach for problems of negative order developed in Chapter 2.

For some *d*-dimensional domain (or manifold)  $\Omega$ , a measurable, closed, possibly empty  $\gamma \subset \partial \Omega$ , and an  $s \in [0, 1]$ , we consider the Sobolev space

$$\mathscr{V} := \left[ L_2(\Omega), H_{0,\gamma}^1(\Omega) \right]_{s,2},$$

with  $H^1_{0,\gamma}(\Omega)$  being the closure in  $H^1(\Omega)$  of the smooth functions on  $\Omega$  that vanish at  $\gamma$ . For  $\mathscr{V}_{\mathcal{T}} \subset \mathscr{V}$  a closed, e.g. finite dimensional subspace, and  $A_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{V}'_{\mathcal{T}}$  some boundedly invertible linear operator, we are interested in constructing a *preconditioner*  $G_{\mathcal{T}} \colon \mathscr{V}'_{\mathcal{T}} \to \mathscr{V}_{\mathcal{T}}$ . More specifically, thinking of a *family* of spaces  $\mathscr{V}_{\mathcal{T}}$  and operators  $A_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{V}'_{\mathcal{T}}$ , our aim is to construct preconditioners  $G_{\mathcal{T}}$  such that  $G_{\mathcal{T}}A_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}} \to \mathscr{V}'_{\mathcal{T}}$  is uniformly boundedly invertible.

It is well-known that such preconditioners of multi-level type are available. The advantage of operator preconditioning is, however, that it does not require a hierarchy of trial spaces.

In order to apply the operator preconditioning framework, one needs to construct families of closed subspaces  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W} := \mathscr{V}'$ , uniformly boundedly invertible  $B_{\mathcal{T}} : \mathscr{W}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$ , and uniformly boundedly invertible  $D_{\mathcal{T}} : \mathscr{V}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$ . Then the resulting preconditioners  $G_{\mathcal{T}}$  are of the form

$$G_{\mathcal{T}} := D_{\mathcal{T}}^{-1} B_{\mathcal{T}} (D_{\mathcal{T}}')^{-1}.$$

The canonical setting is that for  $A: \mathscr{V} \to \mathscr{V}'$ , i.e., an operator of order 2*s*, and an *opposite order* operator  $B: \mathscr{W} \to \mathscr{W}'$ , both boundedly invertible and coercive, it holds that  $(A_{\mathcal{T}}u)(v) := (Au)(v) (u, v \in \mathscr{V}_{\mathcal{T}}), (B_{\mathcal{T}}u)(v) := (Bu)(v)$  $(u, v \in \mathscr{W}_{\mathcal{T}})$ , and  $(D_{\mathcal{T}}u)(v) := \langle u, v \rangle_{L_2(\Omega)} (u \in \mathscr{V}_{\mathcal{T}}, v \in \mathscr{W}_{\mathcal{T}})$ . A typical example
for s = 1/2 is that A is the Hypersingular Integral operator, and B is the Weakly Singular Integral operator, see [SW98].

A careful selection of  $\mathscr{W}_{\mathcal{T}}$  has to be made to ensure that  $D_{\mathcal{T}}: \mathscr{V}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$ is uniformly boundedly invertible. A suitable family of  $(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}})$  pairs has been introduced in [Ste02, BC07]. Here  $\mathcal{T}$  is a triangular partition of a *twodimensional* domain or manifold,  $\mathscr{V}_{\mathcal{T}}$  is the space of continuous piecewise linears w.r.t.  $\mathcal{T}$ , and  $\mathscr{W}_{\mathcal{T}}$  is a subspace of the space of piecewise constants w.r.t. a barycentric refinement of  $\mathcal{T}$ , constructed by subdividing each triangle into 6 subtriangles by connecting its vertices and midpoints with its barycenter. It has been shown in [Ste02, HUT16] that the preconditioner arising from these pairs  $(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}})$  is a uniform preconditioner for families of partitions that satisfy a certain mildly-grading condition.

A problem with the constructions from [Ste02, BC07] appears when one considers the matrix representation  $G_{\mathcal{T}}$  in the standard bases, i.e.  $G_{\mathcal{T}} = D_{\mathcal{T}}^{-1}B_{\mathcal{T}}D_{\mathcal{T}}^{-\top}$ . Indeed, this matrix  $D_{\mathcal{T}}$  is *not* diagonal, and its inverse is densely populated so that it has to be approximated. Moreover, in order to get a uniform preconditioner, since  $G_{\mathcal{T}}$ , being spectrally equivalent with  $A_{\mathcal{T}}^{-1}$ , gets increasingly ill-conditioned with a decreasing minimal mesh-size, the accuracy with which  $D_{\mathcal{T}}^{-1}$  has to be approximated increases with a decreasing minimal mesh-size. As a result, an application of  $D_{\mathcal{T}}^{-1}$  cannot be expected to execute in linear time.

Another (practical) issue with these constructions is the need for the construction of the non-standard *barycentrical* refinement of  $\mathcal{T}$ . This refinement increases the number of elements by a factor 6, and therefore also increases the cost of evaluating  $B_{\mathcal{T}} : \mathscr{W}_{\mathcal{T}} \to \mathscr{W}'_{\mathcal{T}}$ .

#### 4.1.1 Contributions

With  $\mathscr{V}_{\mathcal{T}}$  being the space of continuous piecewise linears, the construction of  $\mathscr{W}_{\mathcal{T}}$  presented in this chapter improves on the existing approach from [Ste02, BC07] concerning the following aspects:

- The matrix representation D<sub>T</sub> of D<sub>T</sub> will be diagonal, allowing one to (exactly) evaluate D<sub>T</sub><sup>-1</sup> in linear time;
- The operator G<sub>T</sub> will be a uniformly well-conditioned preconditioner for families of uniformly shape regular partitions, without requiring a mildly-grading assumption on the partitions;
- By using a stable decomposition of an enclosing space of *W*<sub>T</sub> into a standard finite element space *W*<sub>T</sub> w.r.t. *T* (either being the space of piecewise constants or *W*<sub>T</sub> = *V*<sub>T</sub>) and some bubble space, our *B*<sub>T</sub> will be the *sum* of the corresponding Galerkin discretization operator of the opposite order operator *B*, and an operator whose representation is a diagonal, with which the undesired barycentrical refinement is avoided;

 The construction of *W<sub>T</sub>* applies in any space dimension, and extends to non piecewise planar manifolds.

We will extend the preconditioners to higher order finite element spaces by applying a subspace correction framework.

Due to the interchanged roles of primal and dual spaces, compared to our work in Chapter 2 on preconditioning operators of negative order, here the stable construction of  $W_T$  is simpler, but, on the other hand, the stable decomposition of an enclosing space of  $W_T$  is more delicate.

# 4.1.2 Outline

Sect. 4.1.3 recalls some notation that will be used throughout the article. In Sect. 4.2 the general theory of operator preconditioning is summarized. In Sect. 4.3, the framework is specialized to operators of positive order discretized with *continuous* piecewise linears. Sect. 4.4 give two constructions of  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathcal{W}_{\mathcal{T}}, \mathcal{W}_{\mathcal{T}}')$  that avoid any refinement of the partition  $\mathcal{T}$  that underlies the trial space  $\mathcal{V}_{\mathcal{T}}$ . In Sect. 4.5 the preconditioners are generalized to higher order finite element spaces, and to spaces defined on manifolds. Finally, in Sect. 4.6 we report some numerical results obtained with the new preconditioners.

#### 4.1.3 Notations

By  $\lambda \leq \mu$  we will mean that  $\lambda$  can be bounded by a multiple of  $\mu$ , independently of parameters which  $\lambda$  and  $\mu$  may depend on, with the sole exception of the space dimension d, or in the manifold case, on the parametrization of the manifold that is used to define the finite element spaces on it. Obviously,  $\lambda \geq \mu$  is defined as  $\mu \leq \lambda$ , and  $\lambda = \mu$  as  $\lambda \leq \mu$  and  $\lambda \geq \mu$ .

For normed linear spaces  $\mathscr{Y}$  and  $\mathscr{Z}$ , in this work for convenience over  $\mathbb{R}$ ,  $\mathcal{L}(\mathscr{Y}, \mathscr{Z})$  will denote the space of bounded linear mappings  $\mathscr{Y} \to \mathscr{Z}$  endowed with the operator norm  $\|\cdot\|_{\mathcal{L}(\mathscr{Y},\mathscr{Z})}$ . The subset of invertible operators in  $\mathcal{L}(\mathscr{Y}, \mathscr{Z})$  with inverses in  $\mathcal{L}(\mathscr{Z}, \mathscr{Y})$  will be denoted as  $\mathcal{L}is(\mathscr{Y}, \mathscr{Z})$ . The *condition number* of a  $C \in \mathcal{L}is(\mathscr{Y}, \mathscr{Z})$  is defined as  $\kappa_{\mathscr{Y},\mathscr{Z}}(C) := \|C\|_{\mathcal{L}(\mathscr{Y},\mathscr{Z})} \|C^{-1}\|_{\mathcal{L}(\mathscr{Z},\mathscr{Y})}$ .

For  $\mathscr{Y}$  a reflexive Banach space and  $C \in \mathcal{L}(\mathscr{Y}, \mathscr{Y}')$  being *coercive*, i.e.,

$$\inf_{0 \neq y \in \mathscr{Y}} \frac{(Cy)(y)}{\|y\|_{\mathscr{Y}}^2} > 0,$$

both C and  $\Re(C) := \frac{1}{2}(C + C')$  are in  $\mathcal{L}is(\mathscr{Y}, \mathscr{Y}')$  with

$$\begin{aligned} \|\Re(C)\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')} &\leq \|C\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')}, \\ \|C^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} &\leq \|\Re(C)^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} = \Big(\inf_{0\neq y\in\mathscr{Y}} \frac{(Cy)(y)}{\|y\|_{\mathscr{Y}}^2}\Big)^{-1}. \end{aligned}$$

The subset of coercive operators in  $\mathcal{L}is(\mathscr{Y}, \mathscr{Y}')$  is denoted as  $\mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$ . If  $C \in \mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$ , then  $C^{-1} \in \mathcal{L}is_c(\mathscr{Y}', \mathscr{Y})$  and  $\|\Re(C^{-1})^{-1}\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Y}')} \leq \|C\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Y})}^2 \|\Re(C)^{-1}\|_{\mathcal{L}(\mathscr{Y}', \mathscr{Y})}$ .

Given a family of operators  $C_i \in \mathcal{L}is(\mathscr{Y}_i, \mathscr{Z}_i)$  ( $\mathcal{L}is_c(\mathscr{Y}_i, \mathscr{Z}_i)$ ), we will write  $C_i \in \mathcal{L}is(\mathscr{Y}_i, \mathscr{Z}_i)$  ( $\mathcal{L}is_c(\mathscr{Y}_i, \mathscr{Z}_i)$ ) uniformly in *i*, or simply 'uniform', when

$$\sup_{i} \max(\|C_i\|_{\mathcal{L}(\mathscr{Y}_i,\mathscr{Z}_i)}, \|C_i^{-1}\|_{\mathcal{L}(\mathscr{Z}_i,\mathscr{Y}_i)}) < \infty,$$

or

$$\sup_{i} \max(\|C_i\|_{\mathcal{L}(\mathscr{Y}_i,\mathscr{Z}_i)}, \|\Re(C_i)^{-1}\|_{\mathcal{L}(\mathscr{Z}_i,\mathscr{Y}_i)}) < \infty.$$

Given a finite collection  $\Upsilon = \{v\}$  in a linear space, we set the *synthesis operator* 

$$\mathcal{F}_{\Upsilon}: \mathbb{R}^{\#\Upsilon} \to \operatorname{span} \Upsilon: \mathbf{c} \mapsto \mathbf{c}^{\top} \Upsilon := \sum_{\upsilon \in \Upsilon} c_{\upsilon} \upsilon.$$

Equipping  $\mathbb{R}^{\#\Upsilon}$  with the Euclidean scalar product  $\langle \cdot, \cdot \rangle$ , and identifying  $(\mathbb{R}^{\#\Upsilon})'$  with  $\mathbb{R}^{\#\Upsilon}$  using the corresponding Riesz map, we infer that the adjoint of  $\mathcal{F}_{\Upsilon}$ , known as the *analysis operator*, satisfies

$$\mathcal{F}'_{\Upsilon} : (\operatorname{span} \Upsilon)' \to \mathbb{R}^{\#\Upsilon} : f \mapsto f(\Upsilon) := [f(\upsilon)]_{\upsilon \in \Upsilon}.$$

A collection  $\Upsilon$  is a *basis* for its span when  $\mathcal{F}_{\Upsilon} \in \mathcal{L}is(\mathbb{R}^{\#\Upsilon}, \operatorname{span} \Upsilon)$  (and so  $\mathcal{F}'_{\Upsilon} \in \mathcal{L}is((\operatorname{span} \Upsilon)', \mathbb{R}^{\#\Upsilon})$ .)

Two countable collections  $\Upsilon = (v_i)_i$  and  $\tilde{\Upsilon} = (\tilde{v}_i)_i$  in a Hilbert space will be called *biorthogonal* when  $\langle \Upsilon, \tilde{\Upsilon} \rangle = [\langle v_j, \tilde{v}_i \rangle]_{ij}$  is an *invertible diagonal* matrix, and *biorthonormal* when it is the *identity* matrix.

# 4.2 Operator preconditioning

We shortly recap the idea of opposite order preconditioning, which is based on the following result, see [Hip06, Sect. 2].

**Proposition 4.2.1.** Let  $\mathscr{V}, \mathscr{W}$  be reflexive Banach spaces. If  $B \in \mathcal{L}is(\mathscr{W}, \mathscr{W}')$  and  $D \in \mathcal{L}is(\mathscr{V}, \mathscr{W}')$ , then

$$G := D^{-1}B(D')^{-1} \in \mathcal{L}is(\mathscr{V}', \mathscr{V}),$$

and

$$\|G\|_{\mathcal{L}(\mathscr{V}',\mathscr{V})} \leq \|D^{-1}\|_{\mathcal{L}(\mathscr{W}',\mathscr{V})}^{2}\|B\|_{\mathcal{L}(\mathscr{W},\mathscr{W}')},$$
  
$$\|G^{-1}\|_{\mathcal{L}(\mathscr{V},\mathscr{V}')} \leq \|D\|_{\mathcal{L}(\mathscr{V},\mathscr{W}')}^{2}\|B^{-1}\|_{\mathcal{L}(\mathscr{W}',\mathscr{W})}.$$

If additionally  $B \in \mathcal{L}is_c(\mathcal{W}, \mathcal{W}')$ , then  $G \in \mathcal{L}is_c(\mathcal{V}', \mathcal{V})$ , and

$$\|\Re(G)^{-1}\|_{\mathcal{L}(\mathscr{V},\mathscr{V}')} \le \|D\|^2_{\mathcal{L}(\mathscr{V},\mathscr{W}')}\|\Re(B)^{-1}\|_{\mathcal{L}(\mathscr{W}',\mathscr{W})}.$$

Let be given families of finite dimensional spaces  $\mathscr{V}_{\mathcal{T}}$  for  $\mathcal{T} \in \mathbb{T}$ , and operators  $A_{\mathcal{T}} \in \mathcal{L}$ is $(\mathscr{V}_{\mathcal{T}}, \mathscr{V}'_{\mathcal{T}})$  uniformly in  $\mathcal{T} \in \mathbb{T}$ . Then in light of Proposition 4.2.1 we will seek preconditioners for  $A_{\mathcal{T}}$  of the form

$$G_{\mathcal{T}} = D_{\mathcal{T}}^{-1} B_{\mathcal{T}} (D_{\mathcal{T}}')^{-1},$$

where  $B_{\mathcal{T}} \in \mathcal{L}is(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  and  $D_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  (both uniformly in  $\mathcal{T} \in \mathbb{T}$ ), and

(4.1) 
$$\dim \mathscr{W}_{\mathcal{T}} = \dim \mathscr{V}_{\mathcal{T}}.$$

A typical situation is that for some reflexive Banach space  $\mathscr{V}$  and  $A \in \mathcal{L}is_c(\mathscr{V}, \mathscr{V}')$ , it holds that  $\mathscr{V}_{\mathcal{T}} \subset \mathscr{V}$  (thus equipped with  $\|\cdot\|_{\mathscr{V}}$ ) and  $(A_{\mathcal{T}}u)(v) := (Au)(v)$  ( $u, v \in \mathscr{V}_{\mathcal{T}}$ ), so that indeed  $A_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{V}_{\mathcal{T}}, \mathscr{V}'_{\mathcal{T}})$  uniformly in  $\mathcal{T} \in \mathbb{T}$ . Then for a suitable reflexive Banach space  $\mathscr{W}$ , an operator  $B \in \mathcal{L}is_c(\mathscr{W}, \mathscr{W}')$ , and a subspace  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{W}$  (thus equipped with  $\|\cdot\|_{\mathscr{W}}$ ), one can take  $(B_{\mathcal{T}}w)(z) := (Bw)(z)$  ( $w, z \in \mathscr{W}_{\mathcal{T}}$ ), giving  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$  uniformly. A possible construction of  $D_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$  uniformly is discussed in the next proposition.

**Proposition 4.2.2** (Fortin projector ([For77])). For some  $D \in \mathcal{L}is(\mathscr{V}, \mathscr{W}')$ , let  $D_{\mathcal{T}} \in \mathcal{L}(\mathscr{V}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$  be defined by  $(D_{\mathcal{T}}v)(w) := (Dv)(w)$ . Then

$$\|D_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{W}_{\mathcal{T}}')} \le \|D\|_{\mathcal{L}(\mathscr{V},\mathscr{W}')}.$$

Assuming (4.1), additionally one has  $D_{\mathcal{T}} \in \mathcal{L}is(\mathcal{V}_{\mathcal{T}}, \mathcal{W}_{\mathcal{T}}')$  if, and for  $\mathcal{W}$  being a Hilbert space, only if there exists a projector  $P_{\mathcal{T}} \in \mathcal{L}(\mathcal{W}, \mathcal{W})$  onto  $\mathcal{W}_{\mathcal{T}}$  with  $(D\mathcal{V}_{\mathcal{T}})((\mathrm{Id} - P_{\mathcal{T}})\mathcal{W}) = 0$ , in which case

$$\|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}(\mathscr{W}_{\mathcal{T}}',\mathscr{V}_{\mathcal{T}})} \leq \|P_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{W},\mathscr{W})}\|D^{-1}\|_{\mathcal{L}(\mathscr{W}',\mathscr{V})}.$$

In our applications, the choices for  $\mathcal{W}$  and D will be obvious, and the key ingredient for the construction of a uniform preconditioner  $G_{\mathcal{T}}$  will be the selection of  $\mathcal{W}_{\mathcal{T}}$  that allows for a uniformly bounded Fortin projector  $P_{\mathcal{T}}$ .

#### 4.2.1 Implementation

Let  $\Phi_{\mathcal{T}} = (\phi_i)_i$  and  $\Psi_{\mathcal{T}} = (\psi_i)_i$  be bases for  $\mathscr{V}_{\mathcal{T}}$  and  $\mathscr{W}_{\mathcal{T}}$ , respectively. Then in coordinates the preconditioned system reads as

$$\mathcal{F}_{\Phi_{\mathcal{T}}}^{-1}G_{\mathcal{T}}A_{\mathcal{T}}\mathcal{F}_{\Phi_{\mathcal{T}}} = G_{\mathcal{T}}A_{\mathcal{T}} \coloneqq D_{\mathcal{T}}^{-1}B_{\mathcal{T}}D_{\mathcal{T}}^{-\top}A_{\mathcal{T}},$$

where

$$\boldsymbol{A}_{\mathcal{T}} := \mathcal{F}'_{\Phi_{\mathcal{T}}} A_{\mathcal{T}} \mathcal{F}_{\Phi_{\mathcal{T}}}, \quad \boldsymbol{B}_{\mathcal{T}} := \mathcal{F}'_{\Psi_{\mathcal{T}}} B_{\mathcal{T}} \mathcal{F}_{\Psi_{\mathcal{T}}}, \quad \boldsymbol{D}_{\mathcal{T}} := \mathcal{F}'_{\Psi_{\mathcal{T}}} D_{\mathcal{T}} \mathcal{F}_{\Phi_{\mathcal{T}}}.$$

By identifying a map in  $\mathcal{L}(\mathbb{R}^{\#\Phi_{\mathcal{T}}}, \mathbb{R}^{\#\Phi_{\mathcal{T}}})$  with a  $\#\Phi_{\mathcal{T}} \times \#\Phi_{\mathcal{T}}$  matrix by equipping  $\mathbb{R}^{\#\Phi_{\mathcal{T}}}$  with the canonical basis  $(e_i)_i$  one has,

$$(\boldsymbol{A}_{\mathcal{T}})_{ij} = \langle \mathcal{F}_{\Phi_{\mathcal{T}}}' A_{\mathcal{T}} \mathcal{F}_{\Phi_{\mathcal{T}}} \boldsymbol{e}_j, \boldsymbol{e}_i \rangle = (A_{\mathcal{T}} \mathcal{F}_{\Phi_{\mathcal{T}}} \boldsymbol{e}_j) (\mathcal{F}_{\Phi_{\mathcal{T}}} \boldsymbol{e}_i) = (A_{\mathcal{T}} \phi_j) (\phi_i)_{ij}$$

and similarly,

$$(\boldsymbol{B}_{\mathcal{T}})_{ij} = (B_{\mathcal{T}}\psi_j)(\psi_i), \quad (\boldsymbol{D}_{\mathcal{T}})_{ij} = (D_{\mathcal{T}}\phi_j)(\psi_i).$$

Preferably  $D_{\mathcal{T}}$  is such that its inverse can be applied in linear complexity, as is the case when  $D_{\mathcal{T}}$  is *diagonal*. A goal of this work is to construct such a diagonal  $D_{\mathcal{T}}$ .

*Remark* 4.2.3. Using  $\sigma(\cdot)$  and  $\rho(\cdot)$  to denote the spectrum and spectral radius of an operator, clearly  $\sigma(\mathbf{G}_{\mathcal{T}}\mathbf{A}_{\mathcal{T}}) = \sigma(G_{\mathcal{T}}A_{\mathcal{T}})$ . So for the spectral condition number we have

$$\kappa_S(\boldsymbol{G}_{\mathcal{T}}\boldsymbol{A}_{\mathcal{T}}) \coloneqq \rho(\boldsymbol{G}_{\mathcal{T}}\boldsymbol{A}_{\mathcal{T}})\rho((\boldsymbol{G}_{\mathcal{T}}\boldsymbol{A}_{\mathcal{T}})^{-1}) \le \kappa_{\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}}}(\boldsymbol{G}_{\mathcal{T}}\boldsymbol{A}_{\mathcal{T}}),$$

which thus holds true *independently* of the choice of the basis  $\Phi_{\mathcal{T}}$  for  $\mathscr{V}_{\mathcal{T}}$ . Furthermore, in view of an application of Conjugate Gradients, if  $A_{\mathcal{T}}$  and  $B_{\mathcal{T}}$  are coercive and *self-adjoint*, then  $A_{\mathcal{T}}$  and  $G_{\mathcal{T}}$  are positive definite and symmetric. Equipping  $\mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}}$  with  $\|\cdot\| := \|(G_{\mathcal{T}})^{-\frac{1}{2}} \cdot \|$  or  $\|\cdot\| := \|(A_{\mathcal{T}})^{\frac{1}{2}} \cdot \|$ , in that case we have

$$\kappa_{(\mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}}, \|\cdot\|), (\mathbb{R}^{\dim \mathscr{V}_{\mathcal{T}}}, \|\cdot\|)}(G_{\mathcal{T}}A_{\mathcal{T}}) = \kappa_{S}(G_{\mathcal{T}}A_{\mathcal{T}}).$$

# 4.3 Continuous piecewise linear discretization space

For a bounded polytopal domain  $\Omega \subset \mathbb{R}^d$ , a measurable, closed, possibly empty  $\gamma \subset \partial \Omega$ , and an  $s \in [0, 1]$ , we take

$$\mathscr{V} := [L_2(\Omega), H^1_{0,\gamma}(\Omega)]_{s,2}, \quad \mathscr{W} := \mathscr{V}',$$

which forms the Gelfand triple  $\mathscr{V} \hookrightarrow L_2(\Omega) \simeq L_2(\Omega)' \hookrightarrow \mathscr{W}$ . We define the operator  $D \in \mathcal{L}is(\mathscr{V}, \mathscr{W}')$  as the unique extension to  $\mathscr{V} \times \mathscr{W}$  of the duality pairing

$$(Dv)(w) := \langle v, w \rangle_{L_2(\Omega)}$$

which satisfies  $||D||_{\mathcal{L}(\mathscr{V},\mathscr{W}')} = ||D^{-1}||_{\mathcal{L}(\mathscr{W}',\mathscr{V})} = 1.$ 

Let  $(\mathcal{T})_{\mathcal{T}\in\mathbb{T}}$  be a family of conforming partitions of  $\Omega$  into closed uniformly shape regular *d*-simplices. Thanks to the conformity and the uniform shape regularity, for d > 1 we know that neighbouring  $T, T' \in \mathcal{T}$ , i.e.  $T \cap T' \neq \emptyset$ , have uniformly comparable sizes. For d = 1, we impose this uniform '*K*-mesh property' explicitly.

For some  $\mathcal{T} \in \mathbb{T}$ , denote  $N_{\mathcal{T}}^0$  as the subset of vertices that are not on  $\gamma$ , where we assume that  $\gamma$  is the (possibly empty) union of (d-1)-faces of  $T \in \mathcal{T}$ . For  $T \in \mathcal{T}$ , write  $N_T$  for the set of its vertices, set  $N_T^0 := N_{\mathcal{T}}^0 \cap N_T$ ,  $h_T := |T|^{1/d}$ , and the piecewise constant function  $h_{\mathcal{T}}$  by  $h_{\mathcal{T}}|_T = h_T$  ( $T \in \mathcal{T}$ ). For any vertex  $\nu \in N_{\mathcal{T}}^0$ , define the patch  $\omega_{\mathcal{T},\nu} := \bigcup_{\{T \in \mathcal{T}: \nu \in T\}} T$  and the local mesh size  $h_{\mathcal{T},\nu} := |\omega_{\mathcal{T},\nu}|^{1/d}$ . We omit notational dependence on  $\mathcal{T}$  if it is clear from the context, and simply write  $\omega_{\nu}$  and  $h_{\nu}$ .

Let the discretization space  $\mathscr{V}_{\mathcal{T}}$  be the space of continuous piecewise linears, zero on  $\gamma$ ,

$$\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1} := \{ u \in H^1_{0,\gamma}(\Omega) : u |_T \in \mathcal{P}_1 \ (T \in \mathcal{T}) \} \subset \mathscr{V}_T$$

equipped with the nodal bases

$$\Phi_{\mathcal{T}} = \{\phi_{\nu} \colon \nu \in N_{\mathcal{T}}^0\}$$

defined by  $\phi_{\nu}(\nu') := \delta_{\nu\nu'}$  ( $\nu, \nu' \in N^0_{\mathcal{T}}$ ). For future reference, define the space of discontinuous piecewise constants by

$$\mathscr{S}_{\mathcal{T}}^{-1,0} := \{ u \in L_2(\Omega) \colon u |_T \in \mathcal{P}_0 \, (T \in \mathcal{T}) \} \subset \mathscr{W},$$

equipped with the basis

$$\Sigma_{\mathcal{T}} := \{\mathbb{1}_T \colon T \in \mathcal{T}\}$$

where  $\mathbb{1}_K$  is defined by, for any  $K \subseteq \Omega$ ,  $\mathbb{1}_K := 1$  on K, and  $\mathbb{1}_K := 0$  elsewhere.

#### **4.3.1** The subspace $\mathscr{W}_{\mathcal{T}}$

We will construct the preconditioning space  $\mathscr{W}_{\mathcal{T}}$  as

$$\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}} \subset \mathscr{W}, \text{ with } \dim \mathscr{W}_{\mathcal{T}} = \dim \mathscr{V}_{\mathcal{T}}$$

for a collection  $\Psi_{\mathcal{T}} \subset L_2(\Omega)$  that is biorthogonal to  $\Phi_{\mathcal{T}}$ , and for which the biorthogonal projector  $P_{\mathcal{T}} \in \mathcal{L}(\mathcal{W}, \mathcal{W})$  onto  $\mathcal{W}_{\mathcal{T}}$  is uniformly bounded. We require the collection  $\Psi_{\mathcal{T}} := \{\psi_{\nu} \in \mathcal{W} : \nu \in N_{\mathcal{T}}^0\}$  to satisfy

(4.3) 
$$\begin{aligned} \left| \langle \phi_{\nu}, \psi_{\nu'} \rangle_{L_{2}(\Omega)} \right| &\approx \delta_{\nu\nu'} \| \phi_{\nu} \|_{L_{2}(\Omega)} \| \psi_{\nu'} \|_{L_{2}(\Omega)} \quad (\nu, \nu' \in N_{\mathcal{T}}^{0}), \\ & \operatorname{supp} \psi_{\nu} \subseteq \omega_{\nu} \quad (\nu \in N_{\mathcal{T}}^{0}). \end{aligned} \end{aligned}$$

Existence of such collections will be shown later in Sect. 4.4.

#### 4.3.2 Bounded Fortin projector

From (4.3) it follows that the biorthogonal Fortin projector  $P_{\mathcal{T}} \colon H^1_{0,\gamma}(\Omega)' \to L_2(\Omega)$  onto  $\mathscr{W}_{\mathcal{T}}$  with  $\operatorname{ran}(I - P_{\mathcal{T}}) = \mathscr{V}_{\mathcal{T}}^{\perp_{L_2(\Omega)}}$  exists, and is given by

$$P_{\mathcal{T}}u = \sum_{\nu \in N_{\mathcal{T}}^0} \frac{\langle u, \phi_{\nu} \rangle_{L_2(\Omega)}}{\langle \phi_{\nu}, \psi_{\nu} \rangle_{L_2(\Omega)}} \psi_{\nu}.$$

Uniform boundedness of  $||P_{\mathcal{T}}||_{\mathcal{L}(\mathcal{W},\mathcal{W})}$  follows from uniform boundedness of its adjoint  $P'_{\mathcal{T}}$ , which can be shown similarly as in Theorem 2.5.1<sup>1</sup>:

<sup>&</sup>lt;sup>1</sup>Note that the roles of  $\mathscr{V}$  and  $\mathscr{W}$  are interchanged compared to Chapter 2.

**Theorem 4.3.1.** It holds that  $\sup_{\mathcal{T}\in\mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}(\mathcal{W},\mathcal{W})} = \sup_{\mathcal{T}\in\mathbb{T}} \|P'_{\mathcal{T}}\|_{\mathcal{L}(\mathcal{V},\mathcal{V})} < \infty.$ 

*Proof.* Let  $\mathcal{T} \in \mathbb{T}$ . Define  $\omega_T^{(0)} := T$  for  $T \in \mathcal{T}$ , and for i = 1, ..., denote  $\omega_T^{(i)} := \bigcup_{\{T' \in \mathcal{T}: T' \cap \omega_T^{(i-1)} \neq \emptyset\}} T'$ . The adjoint  $P'_{\mathcal{T}}: L_2(\Omega) \to H^1_{0,\gamma}(\Omega)$  onto  $\mathscr{V}_{\mathcal{T}}$  is given by

$$P'_{\mathcal{T}}u = \sum_{\nu \in N^0_{\mathcal{T}}} \frac{\langle u, \psi_{\nu} \rangle_{L_2(\Omega)}}{\langle \phi_{\nu}, \psi_{\nu} \rangle_{L_2(\Omega)}} \phi_{\nu}.$$

Properties of the nodal basis functions,  $\|\phi_{\nu}\|_{L_2(\Omega)}^2 = h_{\nu}^d$  and  $\|\phi_{\nu}\|_{H^1(\Omega)}^2 \leq h_{\nu}^{d-2}$ , in combination with (4.3), can be used to show that, for  $T \in \mathcal{T}$  and  $k \in \{0, 1\}$ ,

(4.4) 
$$\|P_{\mathcal{T}}'u\|_{H^{k}(T)} \leq \sum_{\nu \in N_{T}^{0}} \|\phi_{\nu}\|_{H^{k}(T)} \frac{\|u\|_{L_{2}(\operatorname{supp}\psi_{\nu})}\|\psi_{\nu}\|_{L_{2}(\Omega)}}{|\langle\phi_{\nu},\psi_{\nu}\rangle_{L_{2}(\Omega)}|} \\ \lesssim \sum_{\nu \in N_{T}^{0}} h_{\nu}^{-k} \|u\|_{L_{2}(\operatorname{supp}\psi_{\nu})} \lesssim h_{T}^{-k} \|u\|_{L_{2}(\omega_{T}^{(1)})},$$

from which we may directly conclude that

$$\sup_{\mathcal{T}\in\mathbb{T}} \|P_{\mathcal{T}}'\|_{\mathcal{L}(L_2(\Omega),L_2(\Omega))} < \infty.$$

For proving boundedness in  $H^1_{0,\gamma}(\Omega)$ , we consider the Scott-Zhang ([SZ90]) interpolator  $\Pi_{\mathcal{T}} : H^1_{0,\gamma}(\Omega) \to \mathscr{V}_{\mathcal{T}}$ . From (4.4) and properties of the  $\Pi_{\mathcal{T}}$  [SZ90, (3.8) and (4.3)], we deduce that

$$\begin{aligned} \|P'_{\mathcal{T}}u\|_{H^{1}(T)} &= \|\Pi_{\mathcal{T}}u + P'_{\mathcal{T}}(\mathrm{Id} - \Pi_{\mathcal{T}})u\|_{H^{1}(T)} \\ &\lesssim \|u\|_{H^{1}(\omega_{\mathcal{T}}^{(1)}(T))} + h_{T}^{-1}\|(\mathrm{Id} - \Pi_{\mathcal{T}})u\|_{L_{2}(\omega_{\mathcal{T}}^{(1)}(T))} \\ &\lesssim \|u\|_{H^{1}(\omega_{\mathcal{T}}^{(2)}(T))}, \end{aligned}$$

and consequently

$$\sup_{\mathcal{T}\in\mathbb{T}} \|P_{\mathcal{T}}'\|_{\mathcal{L}(H^1_{0,\gamma}(\Omega),H^1_{0,\gamma}(\Omega))} < \infty.$$

An application of the Riesz-Thorin interpolation theorem yields the result.  $\Box$ 

The basis  $\Psi_T$  has the crucial benefit that the matrix representation of  $D_T$ , i.e.

$$\boldsymbol{D}_{\mathcal{T}} = \langle \Phi_{\mathcal{T}}, \Psi_{\mathcal{T}} \rangle_{L_2(\Omega)},$$

is diagonal, and thus easily invertible, cf. Sect. 4.2.1.

Combining the theorem with Proposition 4.2.2 gives the following corollary (without requiring additional assumptions on the family of partitions  $\mathbb{T}$ ).

**Corollary 4.3.2.** Suppose we have  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathcal{W}_{\mathcal{T}}, \mathcal{W}'_{\mathcal{T}})$  uniformly. With  $D_{\mathcal{T}} \colon \mathcal{V}_{\mathcal{T}} \to \mathcal{W}_{\mathcal{T}}$  defined by  $(D_{\mathcal{T}}v)(w) := \langle v, w \rangle_{L_2(\Omega)}$ , we find that  $G_{\mathcal{T}} = D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D'_{\mathcal{T}})^{-1} \in \mathcal{L}is_c(\mathcal{V}'_{\mathcal{T}}, \mathcal{V}_{\mathcal{T}})$  is a uniform preconditioner of  $A_{\mathcal{T}} \in \mathcal{L}is_c(\mathcal{V}_{\mathcal{T}}, \mathcal{V}'_{\mathcal{T}})$ .

Given some  $B \in \mathcal{L}is_c(\mathcal{W}, \mathcal{W}')$ , a possible choice for  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathcal{W}_{\mathcal{T}}, \mathcal{W}'_{\mathcal{T}})$ uniformly in  $\mathcal{T} \in \mathbb{T}$ , is  $(B_{\mathcal{T}}u)(v) := (Bu)(v)$   $(u, v \in \mathcal{W}_{\mathcal{T}})$ . For  $d \in \{2, 3\}$  and  $\mathcal{W}' = \mathcal{V} = H^{\frac{1}{2}}(\Omega)$ , a suitable *B* is given by the Weakly Singular Integral operator, whereas for  $\mathcal{W}' = \mathcal{V} = H^{\frac{1}{2}}_{00}(\Omega) := [L_2(\Omega), H^1_0(\Omega)]_{\frac{1}{2},2}$ , the recently in [HJHUT18] introduced *Modified* Weakly Singular Integral operator can be applied. Similar comments apply to screens.

# 4.4 Construction of $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$

We expect it to be impossible to construct a basis  $\Psi_{\mathcal{T}}$  in the (standard) spaces  $\mathscr{S}_{\mathcal{T}}^{-1,0}$ or  $\mathscr{S}_{\mathcal{T}}^{0,1}$  that is *local* and *biorthogonal* to  $\Phi_{\mathcal{T}}$  as required in (4.3). One remedy is to construct  $\Psi_{\mathcal{T}}$  in a (finite element) space w.r.t. a refined partition  $\mathcal{T}_* \succ \mathcal{T}$ . However, this implies that some opposite order operator  $B \in \mathcal{L}is_c(\mathcal{W}, \mathcal{W}')$  has to be discretized on a space w.r.t. the *refined* partition  $\mathcal{T}_*$ . This increases the cost of the preconditioner, and moreover, increases implementational complexity as one has to actually construct this refined partition.

To circumvent (explicit) dependence on the refined partition  $\mathcal{T}_*$ , we shall apply the idea described in Sect. 2.3. That is, we will construct an operator  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  by decomposing an enclosing space of space  $\mathscr{W}_{\mathcal{T}}$  into a a standard finite element space  $\mathscr{U}_{\mathcal{T}}$ , either  $\mathscr{S}_{\mathcal{T}}^{-1,0}$  (in Sect. 4.4.2) or  $\mathscr{S}_{\mathcal{T}}^{0,1}$  (in Sect. 4.4.3), and some bubble space  $\mathscr{B}_{\mathcal{T}}$ . On  $\mathscr{U}_{\mathcal{T}}$  we will apply the Galerkin discretization operator of the opposite order operator B, whereas on the bubble space  $\mathscr{B}_{\mathcal{T}}$  a diagonal scaling will suffice.

In the first subsection we present this construction of  $B_{\mathcal{T}}$  for some abstract  $\mathscr{W}_{\mathcal{T}}$ . In the subsequent subsections, we will present two viable options for  $\mathscr{W}_{\mathcal{T}}$ , leading to two different preconditioners.

#### 4.4.1 Stable decomposition

The role of the space ' $\mathscr{Y}$ ' is the next proposition is going to be played by  $\mathscr{W}_{\mathcal{T}}$ .

**Proposition 4.4.1.** Let  $\mathscr{Z}$  be an inner product space,  $Q \in \mathcal{L}(\mathscr{Z}, \mathscr{Z})$  a projector, and with  $\mathscr{U} := \operatorname{ran} Q$ , let  $\mathscr{B} := \operatorname{ran}(\operatorname{Id} - Q)$ ,  $B^{\mathscr{U}} \in \operatorname{Lis}_{c}(\mathscr{U}, \mathscr{U}')$ , and  $B^{\mathscr{B}} \in \operatorname{Lis}_{c}(\mathscr{B}, \mathscr{B}')$ . Then for any subspace  $\mathscr{Y} \subset \mathscr{Z}$ ,

$$(By)(\tilde{y}) := (B^{\mathscr{U}}Qy)(Q\tilde{y}) + (B^{\mathscr{B}}(\mathrm{Id} - Q)y)((\mathrm{Id} - Q)\tilde{y}) \quad (y, \tilde{y} \in \mathscr{Y})$$

is bounded and coercive —  $B \in \mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$  — with

$$\begin{split} \|B\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')} &\leq \\ \left(\|Q\|^2 + \sqrt{\|Q\|^4 - \|Q\|^2}\right) \max(\|B^{\mathscr{U}}\|_{\mathcal{L}(\mathscr{U},\mathscr{U}')}, \|B^{\mathscr{B}}\|_{\mathcal{L}(\mathscr{B},\mathscr{B}')}), \\ \|\Re(B)^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} &\leq \\ \left(1 + \sqrt{1 - \|Q\|^{-2}}\right) \max(\|\Re(B^{\mathscr{U}})^{-1}\|_{\mathcal{L}(\mathscr{U}',\mathscr{U})}, \|\Re(B^{\mathscr{B}})^{-1}\|_{\mathcal{L}(\mathscr{B}',\mathscr{B})}), \end{split}$$

where  $||Q|| := ||Q||_{\mathcal{L}(\mathscr{Z},\mathscr{Z})}$ .

*Proof.* Let  $y, \tilde{y} \in \mathscr{Y}$ . Write u = Qy,  $b = (\mathrm{Id} - Q)y$ , and similarly  $\tilde{u} = Q\tilde{y}$ ,  $\tilde{b} = (\mathrm{Id} - Q)\tilde{y}$ . We have

$$\begin{split} |(B(y))(\tilde{y})| &\leq \max\left( \|B^{\mathscr{U}}\|_{\mathcal{L}(\mathscr{U},\mathscr{U}')}, \|B^{\mathscr{B}}\|_{\mathcal{L}(\mathscr{B},\mathscr{B}')} \right) \cdot \left( \|u\|_{\mathscr{Z}} \|\tilde{u}\|_{\mathscr{Z}} + \|b\|_{\mathscr{Z}} \|\tilde{b}\|_{\mathscr{Z}} \right) \\ &\leq \max(\cdots) \sqrt{\|u\|_{\mathscr{Z}}^2 + \|b\|_{\mathscr{Z}}^2} \cdot \sqrt{\|\tilde{u}\|_{\mathscr{Z}}^2 + \|\tilde{b}\|_{\mathscr{Z}}^2}, \end{split}$$

and

$$|(B(y))(y)| \ge \min\left(||\Re(B^{\mathscr{U}})^{-1}||_{\mathcal{L}(\mathscr{U}',\mathscr{U})}^{-1}, ||\Re(B^{\mathscr{B}})^{-1}||_{\mathcal{L}(\mathscr{B}',\mathscr{B})}^{-1}\right) \cdot (||u||_{\mathscr{Z}}^{2} + ||b||_{\mathscr{Z}}^{2})$$

With  $\gamma := \sup_{0 \neq (u,b) \in \mathscr{U} \times \mathscr{B}} \frac{|\langle u, b \rangle_{\mathscr{U}}|}{\|u\|_{\mathscr{U}} \|b\|_{\mathscr{U}}}$ , for  $0 \neq (u, b) \in \mathscr{U} \times \mathscr{B}$  we have  $\frac{\|u+b\|_{\mathscr{U}}^2}{\|u\|_{\mathscr{U}}^2 + \|b\|_{\mathscr{U}}^2} \in [1 - \gamma, 1 + \gamma]$ . Using that  $\sqrt{\frac{1}{1 - \gamma^2}} = \|Q\|$  (see e.g. [Szy06, (5.5), (5.7), (6.2)]), the proof is easily completed.

*Remark* 4.4.2. For a quantitatively weaker result as Proposition 4.4.1 to hold it is actually sufficient when Q is only defined on  $\mathscr{Y}$ , and neither is it needed that it is a projector. Under these relaxed conditions, obvious estimates show bounds as in Proposition 4.4.1 with the factors  $||Q||^2 + \sqrt{||Q||^4 - ||Q||^2}$  and  $1 + \sqrt{1 - ||Q||^{-2}}$  reading as  $||Q|_{\mathscr{Y}}||^2 + (1 + ||Q|_{\mathscr{Y}}||)^2$  and 2, respectively. Both original factors are equal to 1 when Q is an orthogonal projector.

We are going to apply this abstract proposition with  $\mathscr{Y}' = \mathscr{W}_{\mathcal{T}}$ ,  $\mathscr{U}' = \mathscr{U}_{\mathcal{T}}$ being a standard finite element space,  $\mathscr{B}' = \mathscr{B}_{\mathcal{T}}$  being a suitably constructed 'bubble space', and  $\mathscr{Z}' = \mathscr{Z}_{\mathcal{T}} := \mathscr{U}_{\mathcal{T}} + \mathscr{B}_{\mathcal{T}}$ , all equipped with the norm on  $\mathscr{W}$ . The resulting 'B' will be the  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  we are seeking.

In order to apply above proposition, what is left is the construction of a (uniformly) bounded projector defined on  $\mathscr{Z}_{\mathcal{T}}$ . Furthermore, to allow for a simple preconditioner on  $\mathscr{B}_{\mathcal{T}}$  we would like to find a setting in which on this bubble space the  $\mathscr{W}$ -norm is equivalent to a weighted  $L_2$ -norm. Both issues will be dealt with in the next two lemmas. The operator  $Q_{\mathcal{T}}|_{\mathscr{Z}_{\mathcal{T}}}$  in the first lemma will play the role of 'Q' in Proposition 4.4.1.

**Lemma 4.4.3.** Let  $Q_{\mathcal{T}} \in \mathcal{L}(L_2(\Omega), H^1_{0,\gamma}(\Omega)')$  be a projector,  $\mathscr{U}_{\mathcal{T}} \subseteq \operatorname{ran} Q_{\mathcal{T}}$  and  $\mathscr{B}_{\mathcal{T}} \subseteq \operatorname{ran}(\operatorname{Id} - Q_{\mathcal{T}})$  be subspaces of  $L_2(\Omega)$ , and with  $\mathscr{B}_{\mathcal{T}} := \mathscr{U}_{\mathcal{T}} + \mathscr{B}_{\mathcal{T}}$ , let

(1) 
$$\|h_{\mathcal{T}}^{-1}(\mathrm{Id} - Q'_{\mathcal{T}})\|_{\mathcal{L}(H^{1}_{0,\gamma}(\Omega), L_{2}(\Omega))} \lesssim 1$$
, (approximation property)  
(2)  $\sup_{\mathcal{T} \in \mathbb{T}} \|Q_{\mathcal{T}}\|_{\mathscr{X}_{\mathcal{T}}}\|_{\mathcal{L}((\mathscr{X}_{\mathcal{T}}, \|\cdot\|_{L_{2}(\Omega)}), L_{2}(\Omega))} \lesssim 1$ , (boundedness in  $L_{2}(\Omega)$ )  
(3)  $\|h_{\mathcal{T}} \cdot \|_{L_{2}(\Omega)} \lesssim \|\cdot\|_{H^{1}_{0,\gamma}(\Omega)'}$  on  $\mathscr{X}_{\mathcal{T}}$ . (inverse inequality)

Then  $Q_{\mathcal{T}}|_{\mathscr{Z}_{\mathcal{T}}} \colon \mathscr{Z}_{\mathcal{T}} \to \mathscr{Z}_{\mathcal{T}} \text{ is a projector, } \operatorname{ran} Q_{\mathcal{T}}|_{\mathscr{Z}_{\mathcal{T}}} = \mathscr{U}_{\mathcal{T}}, \operatorname{ran}(\operatorname{Id} - Q_{\mathcal{T}}|_{\mathscr{Z}_{\mathcal{T}}}) = \mathscr{B}_{\mathcal{T}},$ and

- (i)  $\sup_{\mathcal{T}\in\mathbb{T}} \|Q_{\mathcal{T}}\|_{\mathscr{Z}_{\mathcal{T}}}\|_{\mathcal{L}((\mathscr{Z}_{\mathcal{T}},\|\cdot\|_{\mathscr{W}}),\mathscr{W})} < \infty,$
- (ii)  $\|\cdot\|_{\mathscr{W}} \approx \|h^s_{\mathcal{T}} \cdot \|_{L_2(\Omega)}$  on  $\mathscr{B}_{\mathcal{T}}$ .

*Proof.* The first three statements are easily verified. From (1) it follows that for  $u \in H^1_{0,\gamma}(\Omega)'$ :

$$\begin{aligned} \| (\mathrm{Id} - Q_{\mathcal{T}}) u \|_{H^{1}_{0,\gamma}(\Omega)'} &= \sup_{v \in H^{1}_{0,\gamma}(\Omega)} \frac{\langle u, (\mathrm{Id} - Q'_{\mathcal{T}}) v \rangle_{L_{2}(\Omega)}}{\| v \|_{H^{1}_{0,\gamma}(\Omega)}} \\ &\leq \sup_{v \in H^{1}_{0,\gamma}(\Omega)} \frac{\| h_{\mathcal{T}} u \|_{L_{2}(\Omega)} \| h_{\mathcal{T}}^{-1} (\mathrm{Id} - Q'_{\mathcal{T}}) v \|_{L_{2}(\Omega)}}{\| v \|_{H^{1}_{0,\gamma}(\Omega)}} \\ &\lesssim \| h_{\mathcal{T}} u \|_{L_{2}(\Omega)}. \end{aligned}$$

Together with the inverse inequality on  $\mathscr{Z}_{\mathcal{T}}$ , this gives boundedness of  $\|(\mathrm{Id} - Q_{\mathcal{T}})|_{\mathscr{Z}_{\mathcal{T}}}\|_{\mathcal{L}((\mathscr{Z}_{\mathcal{T}}, \|\cdot\|_{H^{1}_{0,\gamma}(\Omega)'}), H^{1}_{0,\gamma}(\Omega)')}$  and thus of  $\|Q_{\mathcal{T}}|_{\mathscr{Z}_{\mathcal{T}}}\|_{\mathcal{L}((\mathscr{Z}_{\mathcal{T}}, \|\cdot\|_{H^{1}_{0,\gamma}(\Omega)'}), H^{1}_{0,\gamma}(\Omega)')}$ . The first result then follows from (2) and an interpolation argument.

By the inverse inequality on  $\mathscr{B}_{\mathcal{T}}$  and the previously derived inequality, we have for  $b_{\mathcal{T}} \in \mathscr{B}_{\mathcal{T}} \subseteq \operatorname{ran}(\operatorname{Id} - Q_{\mathcal{T}})$  that

$$\|b_{\mathcal{T}}\|_{H^{1}_{0,\gamma}(\Omega)'} = \|(\mathrm{Id} - Q_{\mathcal{T}})b_{\mathcal{T}}\|_{H^{1}_{0,\gamma}(\Omega)'} \lesssim \|h_{\mathcal{T}}b_{\mathcal{T}}\|_{L_{2}(\Omega)} \lesssim \|b_{\mathcal{T}}\|_{H^{1}_{0,\gamma}(\Omega)'}.$$

Another interpolation argument yields the second result.

**Lemma 4.4.4.** Suppose that  $\|\cdot\|_{\mathscr{W}} \approx \|h_{\mathcal{T}}^s \cdot \|_{L_2(\Omega)}$  holds on  $\mathscr{B}_{\mathcal{T}}$ , and that  $\Theta_{\mathcal{T}}$  is a uniformly  $\|h_{\mathcal{T}}^s \cdot \|_{L_2(\Omega)}$ -stable basis for  $\mathscr{B}_{\mathcal{T}}$ , *i.e.* 

$$\mathscr{B}_{\mathcal{T}} = \operatorname{span} \Theta_{\mathcal{T}} \quad and \quad \left\| h^s_{\mathcal{T}} \sum_{\theta \in \Theta_{\mathcal{T}}} c_{\theta} \theta \right\|_{L_2(\Omega)}^2 \approx \sum_{\theta \in \Theta_{\mathcal{T}}} |c_{\theta}|^2 \|h^s_{\mathcal{T}} \theta\|_{L_2(\Omega)}^2,$$

then, for any  $\beta_1 > 0$ , an operator  $B_{\mathcal{T}}^{\mathscr{B}} \in \mathcal{L}is_c(\mathscr{B}_{\mathcal{T}}, \mathscr{B}'_{\mathcal{T}})$  is given by

(4.5) 
$$\left(B_{\mathcal{T}}^{\mathscr{B}}\sum_{\theta\in\Theta_{\mathcal{T}}}c_{\theta}\theta\right)\left(\sum_{\theta\in\Theta_{\mathcal{T}}}d_{\theta}\theta\right) = \beta_{1}\sum_{\theta\in\Theta_{\mathcal{T}}}c_{\theta}d_{\theta}\|h_{\mathcal{T}}^{s}\theta\|_{L_{2}(\Omega)}^{2}.$$

*Remark* 4.4.5. It is not possible to construct  $B_{\mathcal{T}} \in \mathcal{L}is(\mathscr{W}_{\mathcal{T}}, \mathscr{W}_{\mathcal{T}}')$  directly as a diagonal scaling operator. Indeed, this would require  $||w_{\mathcal{T}}||_{\mathscr{W}} \lesssim ||h_{\mathcal{T}}^s w_{\mathcal{T}}||_{L_2(\Omega)}$  for  $w_{\mathcal{T}} \in \mathscr{W}_{\mathcal{T}}$ . Suppose this to be true, then by  $L_2(\Omega)$ -boundedness of the biorthogonal projector  $P_{\mathcal{T}}$ , we would find for  $v_{\mathcal{T}} \in \mathscr{V}_{\mathcal{T}}$  that

$$\begin{split} \|h_{\mathcal{T}}^{-s}v_{\mathcal{T}}\|_{L_{2}(\Omega)} &= \sup_{w \in L_{2}(\Omega)} \frac{\langle h_{\mathcal{T}}^{-s}v_{\mathcal{T}}, P_{\mathcal{T}}w \rangle_{L_{2}(\Omega)}}{\|w\|_{L_{2}(\Omega)}} \lesssim \sup_{w \in L_{2}(\Omega)} \frac{\langle h_{\mathcal{T}}^{-s}v_{\mathcal{T}}, P_{\mathcal{T}}w \rangle_{L_{2}(\Omega)}}{\|P_{\mathcal{T}}w\|_{L_{2}(\Omega)}} \\ &= \sup_{w_{\mathcal{T}} \in \mathscr{W}_{\mathcal{T}}} \frac{\langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{L_{2}(\Omega)}}{\|h_{\mathcal{T}}^{s}w_{\mathcal{T}}\|_{L_{2}(\Omega)}} \lesssim \sup_{w_{\mathcal{T}} \in \mathscr{W}_{\mathcal{T}}} \frac{\langle v_{\mathcal{T}}, w_{\mathcal{T}} \rangle_{L_{2}(\Omega)}}{\|w_{\mathcal{T}}\|_{\mathscr{W}}} \le \|v_{\mathcal{T}}\|_{\mathscr{V}}, \end{split}$$

which is known not to be true for smooth functions in  $\mathscr{V}_{\mathcal{T}}$ .

Concluding: If, given a family of subspaces  $\mathscr{W}_{\mathcal{T}} \subset L_2(\Omega)$ , one can find a family of projectors  $Q_{\mathcal{T}} \in \mathcal{L}(L_2(\Omega), H^1_{0,\gamma}(\Omega)')$ , subspaces  $\mathscr{U}_{\mathcal{T}} \subseteq \operatorname{ran} Q_{\mathcal{T}}$  (of finite element type) and  $\mathscr{B}_{\mathcal{T}} \subseteq \operatorname{ran}(\operatorname{Id} - Q_{\mathcal{T}})$  such that

$$(4.6) \mathscr{W}_{\mathcal{T}} \subset \mathscr{Z}_{\mathcal{T}} := \mathscr{U}_{\mathcal{T}} + \mathscr{B}_{\mathcal{T}}$$

(with these spaces equipped with  $\|\cdot\|_{\mathscr{W}}$ -norm) and the conditions of Lemma 4.4.3 are satisfied, then given  $B_{\mathcal{T}}^{\mathscr{U}} \in \mathcal{L}is_{c}(\mathscr{U}_{\mathcal{T}}, \mathscr{U}_{\mathcal{T}}')$  and  $B_{\mathcal{T}}^{\mathscr{B}} \in \mathcal{L}is_{c}(\mathscr{B}_{\mathcal{T}}, \mathscr{B}_{\mathcal{T}}')$ , the operator  $B_{\mathcal{T}}^{\mathscr{W}}$  defined by (4.7)

$$(B_{\mathcal{T}}^{\mathscr{W}})(\tilde{w}) := (B_{\mathcal{T}}^{\mathscr{W}} Q_{\mathcal{T}} w)(Q_{\mathcal{T}} \tilde{w}) + (B_{\mathcal{T}}^{\mathscr{B}} (\mathrm{Id} - Q_{\mathcal{T}}) w)((\mathrm{Id} - Q_{\mathcal{T}}) \tilde{w}) \quad (w, \tilde{w} \in \mathscr{W}_{\mathcal{T}}),$$

is in  $\mathcal{L}is_c(\mathscr{W}_{\mathcal{T}}, \mathscr{W}'_{\mathcal{T}})$ . Moreover, assuming a uniformly  $||h^s_{\mathcal{T}} \cdot ||_{L_2(\Omega)}$ -stable basis for  $\mathscr{B}_{\mathcal{T}}$ , the operator  $B^{\mathscr{B}}_{\mathcal{T}}$  can be of simple diagonal scaling type, where a natural definition for  $B^{\mathscr{U}}_{\mathcal{T}}$  is by  $(B_{\mathcal{T}}u)(\tilde{u}) := (Bu)(\tilde{u}) \ (u, \tilde{u} \in \mathscr{U}_{\mathcal{T}})$  for some opposite order operator  $B \in \mathcal{L}is_c(\mathscr{W}, \mathscr{W}')$ . Finally, since  $Q_{\mathcal{T}}$  enters the implementation, we search this projector to be of local type.

# **4.4.2** A space $\mathscr{W}_{\mathcal{T}}$ enclosed in a space decomposable into the piecewise constants and bubbles

In this subsection, we construct  $\mathscr{W}_{\mathcal{T}} = \operatorname{span} \Psi_{\mathcal{T}}$  such that both  $\Psi_{\mathcal{T}}$  is biorthogonal to  $\Phi_{\mathcal{T}}$  ((4.3)), and  $\mathscr{W}_{\mathcal{T}}$  is enclosed in a space that allows an appropriate decomposition into the space of piecewise constants  $\mathscr{U}_{\mathcal{T}} := \mathscr{S}_{\mathcal{T}}^{-1,0}$  and a bubble space  $\mathscr{B}_{\mathcal{T}}$ .

Fix  $\mathcal{T} \in \mathbb{T}$  and let  $\mathcal{T}_* \succ \mathcal{T}$  be a uniform red-refinement, i.e. every simplex  $T \in \mathcal{T}$  is subdivided into  $2^d$  subsimplices.<sup>2</sup> We define  $\Psi_{\mathcal{T}} = \{\psi_{\mathcal{T},\nu} \colon \nu \in N_{\mathcal{T}}^0\} \subset \mathscr{S}_{\mathcal{T}_*}^{-1,0}$  by taking a weighted difference of 'patch indicator' functions:

(4.8) 
$$\psi_{\mathcal{T},\nu} := 2^{d+1} \mathbb{1}_{\omega_{\mathcal{T}_*,\nu}} - \mathbb{1}_{\omega_{\mathcal{T},\nu}} \quad (\nu \in N^0_{\mathcal{T}}).$$

**Lemma 4.4.6.** The collection  $\Psi_{\mathcal{T}}$  satisfies (4.3) with supp  $\psi_{\mathcal{T},\nu} = \omega_{\mathcal{T},\nu}$  and

(4.9) 
$$\langle \psi_{\mathcal{T},\nu}, \phi_{\mathcal{T},\nu'} \rangle_{L_2(\Omega)} = \delta_{\nu\nu'} |\omega_{\mathcal{T},\nu}| \quad (\nu,\nu' \in N^0_{\mathcal{T}}).$$

*Proof.* Clearly  $\sup \psi_{\mathcal{T},\nu} = \omega_{\mathcal{T},\nu}$ , so we are left to show the biorthogonality condition. Fix some vertex  $\nu \in N^0_{\mathcal{T}}$ . For a simplex  $T_{\nu} \in \mathcal{T}$  with  $\nu \in T_{\nu}$ , we have

$$\langle \mathbb{1}_{T_{\nu}}, \phi_{\mathcal{T},\nu} \rangle_{L_2(\Omega)} = \frac{|T_{\nu}|}{d+1}.$$

<sup>&</sup>lt;sup>2</sup>Red-refinement is not uniquely defined for  $d \ge 3$ , but the refined simplices at the corners of the 'parent simplex' are uniquely determined which suffices for our goal.

Let  $T_{*,\nu} \in \mathcal{T}_*$  be the (unique) simplex with  $\nu \in T_{*,\nu} \subset T_{\nu}$ . From the refinement equation satisfied by the nodal hats, and  $|T_{*,\nu}| = 2^{-d} |T_{\nu}|$ , it follows that

$$\langle \mathbb{1}_{T_{*,\nu}}, \phi_{\mathcal{T},\nu} \rangle_{L_2(\Omega)} = \langle \mathbb{1}_{T_{*,\nu}}, \phi_{\mathcal{T}_{*,\nu}} + \sum_{\nu \neq \tilde{\nu} \in N_{T_{*,\nu}}} 2^{-1} \phi_{\mathcal{T}_{*},\tilde{\nu}} \rangle_{L_2(\Omega)} = \frac{2^{-d} |T_{\nu}|}{d+1} (1+2^{-1}d),$$
  
$$\langle \mathbb{1}_{T_{*,\nu}}, \phi_{\mathcal{T},\nu'} \rangle_{L_2(\Omega)} = \dots = \frac{2^{-d} |T_{\nu}|}{d+1} 2^{-1} \quad (\nu \neq \nu' \in N_{T_{\nu}}^0).$$

From these relations (4.9) follows.

By Lemma 4.4.6 it has been established that the Fortin interpolator is uniformly bounded, and that  $D_{\mathcal{T}}$  is represented by a diagonal matrix. The next proposition verifies the conditions imposed in Sect. 4.4.1 for the construction of  $B_{\mathcal{T}}$ .

**Proposition 4.4.7.** Let  $\mathscr{U}_{\mathcal{T}} := \mathscr{S}_{\mathcal{T}}^{-1,0}$ ,  $\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}}$  as constructed above,  $Q_{\mathcal{T}}$  be the  $L_2(\Omega)$ -orthogonal projector onto  $\mathscr{U}_{\mathcal{T}}, \Theta_{\mathcal{T}} := (\operatorname{Id} - Q_{\mathcal{T}})\Psi_{\mathcal{T}}$ , and  $\mathscr{B}_{\mathcal{T}} := \operatorname{span} \Theta_{\mathcal{T}}$ . Then  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{Z}_{\mathcal{T}} := \mathscr{U}_{\mathcal{T}} + \mathscr{B}_{\mathcal{T}}$  ((4.6)), the conditions of Lemma 4.4.3 are satisfied, in particular  $Q_{\mathcal{T}}\psi_{\nu} = \mathbb{1}_{\omega_{\nu}}$ , and  $\Theta_{\mathcal{T}}$  is a uniformly  $\|h_{\mathcal{T}}^s \cdot \|_{L_2(\Omega)}$ -stable basis for  $\mathscr{B}_{\mathcal{T}}$  as required for Lemma 4.4.4.

*Proof.* The first statement follows from  $\mathscr{W}_{\mathcal{T}} \subset L_2(\Omega)$ . The first two conditions of Lemma 4.4.3 are obviously valid. Concerning the third condition, the inverse inequality  $\|h_{\mathcal{T}} \cdot \|_{L_2(\Omega)} \lesssim \| \cdot \|_{H^1_{0,\gamma}(\Omega)'}$  holds, for general d, on  $\mathscr{S}_{\mathcal{T}_*}^{-1,0}$ , see e.g. Lemma 2.3.4, and thus in particular on  $\mathscr{Z}_{\mathcal{T}}$ . The property  $Q_{\mathcal{T}}\psi_{\nu} = \mathbb{1}_{\omega_{\nu}}$  is easily checked.

We are left to show that the collection of bubbles  $\{\theta_{\nu} := (\mathrm{Id} - Q_{\mathcal{T}})\psi_{\nu} : \nu \in N_{\mathcal{T}}^0\}$  is  $\|h_{\mathcal{T}}^s \cdot \|_{L_2(\Omega)}$ -stable. Pick some  $T \in \mathcal{T}$ , then the normalized 'bubble element matrix' satisfies

(4.10)  
$$\frac{\frac{1}{4}|T|^{-1}\langle\theta_{\nu},\theta_{\nu'}\rangle_{L_{2}(T)} = |T|^{-1}\langle 2^{d}\mathbb{1}_{\omega_{\mathcal{T}_{*},\nu}} - \mathbb{1}_{\omega_{\mathcal{T}_{*},\nu'}} - \mathbb{1}_{\omega_{\mathcal{T}_{*},\nu'}}\rangle_{L_{2}(T)} \\= \begin{cases} 2^{d} - 1 & \nu = \nu' \in N_{T}^{0}, \\ -1 & \nu \neq \nu' \in N_{T}^{0}. \end{cases}$$

For  $d \ge 2$ , this constant (symmetric)  $(d+1) \times (d+1)$  matrix is strictly diagonally dominant, and therefore positive definite. We conclude this proposition by

$$\begin{split} & \left\| \sum_{\nu \in N_{\mathcal{T}}^{0}} h_{\mathcal{T}}^{s} c_{\nu} \theta_{\nu} \right\|_{L_{2}(\Omega)}^{2} = \sum_{T \in \mathcal{T}} h_{T}^{2s} \left\| \sum_{\nu \in N_{T}^{0}} c_{\nu} \theta_{\nu} \right\|_{L_{2}(T)}^{2} \approx \sum_{T \in \mathcal{T}} h_{T}^{2s} \sum_{\nu \in N_{T}^{0}} |c_{\nu}|^{2} \|\theta_{\nu}\|_{L_{2}(T)}^{2} \\ & = \sum_{\nu \in N_{\mathcal{T}}^{0}} |c_{\nu}|^{2} \sum_{T \in \mathcal{T}} \|h_{T}^{s} \theta_{\nu}\|_{L_{2}(T)}^{2} = \sum_{\nu \in N_{\mathcal{T}}^{0}} |c_{\nu}|^{2} \|h_{\mathcal{T}}^{s} \theta_{\nu}\|_{L_{2}(\Omega)}^{2}. \end{split}$$

75

*Remark* 4.4.8. For d = 1, the bubbles arising from  $\Psi_{\mathcal{T}}$  as given in (4.8) do not form a  $\|h^s_{\mathcal{T}} \cdot \|_{L_2(\Omega)}$ -stable collection. Instead, with  $\mathcal{T}_{**} \succ \mathcal{T}$  being the two times uniform red-refinement, one can consider  $\psi_{\mathcal{T},\nu} = \frac{16}{3} \mathbb{1}_{\omega_{\mathcal{T}*,\nu}} - \frac{1}{3} \mathbb{1}_{\omega_{\mathcal{T},\nu}}$  for which the statements of Lemma 4.4.6 and Proposition 4.4.7 are again valid.

#### Implementation

The matrix representation of preconditioner  $\mathcal{F}_{\Phi_{\tau}}^{-1}G_{\tau}(\mathcal{F}_{\Phi_{\tau}}')^{-1}$  is given by

$$G_{\mathcal{T}} = D_{\mathcal{T}}^{-1} B_{\mathcal{T}} D_{\mathcal{T}}^{- op}$$

With  $\Psi_{\mathcal{T}}$  as constructed in (4.8), we find that  $D_{\mathcal{T}} = \mathcal{F}'_{\Psi_{\mathcal{T}}} D_{\mathcal{T}} \mathcal{F}_{\Phi_{\mathcal{T}}}$  is given by

$$\boldsymbol{D}_{\mathcal{T}} = \operatorname{diag}\{|\omega_{\nu}| \colon \nu \in N_{\mathcal{T}}^0\}.$$

Given some  $B_{\mathcal{T}}^{\mathscr{U}} \in \mathcal{L}is_c(\mathscr{U}_{\mathcal{T}}, \mathscr{U}_{\mathcal{T}}')$  (recall that  $\mathscr{U}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,0}$ ), then by taking  $B_{\mathcal{T}}$  as described in (4.7), we have

$$\begin{aligned} \boldsymbol{B}_{\mathcal{T}} &\coloneqq \mathcal{F}'_{\Psi_{\mathcal{T}}} B_{\mathcal{T}} \mathcal{F}_{\Psi_{\mathcal{T}}} \\ &= \mathcal{F}'_{\Psi_{\mathcal{T}}} (Q'_{\mathcal{T}} B_{\mathcal{T}}^{\mathscr{U}} Q_{\mathcal{T}} + (\mathrm{Id} - Q_{\mathcal{T}})' B_{\mathcal{T}}^{\mathscr{B}} (\mathrm{Id} - Q_{\mathcal{T}})) \mathcal{F}_{\Psi_{\mathcal{T}}} \\ &= \boldsymbol{p}_{\mathcal{T}}^{\mathcal{T}} B_{\mathcal{T}}^{\mathscr{U}} \boldsymbol{p}_{\mathcal{T}} + B_{\mathcal{T}}^{\mathscr{B}}, \end{aligned}$$

where, using that  $\mathcal{F}_{\Theta_{\mathcal{T}}}^{-1}(\mathrm{Id} - Q_{\mathcal{T}})\mathcal{F}_{\Psi_{\mathcal{T}}} = \mathrm{Id}$  by  $\Theta_{\mathcal{T}} = (I - Q_{\mathcal{T}})\Psi_{\mathcal{T}}$ ,

$$\boldsymbol{B}_{\mathcal{T}}^{\mathscr{U}} \coloneqq \mathcal{F}_{\Sigma_{\mathcal{T}}} B_{\mathcal{T}}^{\mathscr{U}} \mathcal{F}_{\Sigma_{\mathcal{T}}}, \quad \boldsymbol{p}_{\mathcal{T}} \coloneqq \mathcal{F}_{\Sigma_{\mathcal{T}}}^{-1} Q_{\mathcal{T}} \mathcal{F}_{\Psi_{\mathcal{T}}}, \quad \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} \coloneqq \mathcal{F}_{\Theta_{\mathcal{T}}} B_{\mathcal{T}}^{\mathscr{B}} \mathcal{F}_{\Theta_{\mathcal{T}}},$$

Recall the canonical basis  $\Sigma_{\mathcal{T}}$  for  $\mathscr{U}_{\mathcal{T}}$  from (4.2). Using  $Q_{\mathcal{T}}\psi_{\nu} = \mathbb{1}_{\omega_{\nu}}$  shows that

$$(\boldsymbol{p}_{\mathcal{T}})_{T\nu} = \begin{cases} 1 & \text{if } T \subset \omega_{\nu}, \\ 0 & \text{else.} \end{cases}$$

From (4.10), we infer that  $\|h_{\mathcal{T}}^s \theta_{\nu}\|_{L_2(\Omega)}^2 \approx |\omega_{\nu}|^{1+\frac{2s}{d}}$ . By making a harmless modification to the definition of  $B_{\mathcal{T}}^{\mathscr{B}}$  in (4.5) based on this equivalency, we obtain that

$$\boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} = \beta_1 \boldsymbol{D}_{\mathcal{T}}^{1+\frac{2s}{d}}.$$

The matrix  $B_{\mathcal{T}}^{\mathscr{U}}$  depends on the operator  $B_{\mathcal{T}}^{\mathscr{U}} \in \mathcal{L}is_c(\mathscr{U}_{\mathcal{T}}, \mathscr{U}_{\mathcal{T}}')$  that is selected. The canonical choice is the Galerkin discretization operator on  $\mathscr{U}_{\mathcal{T}}$  of a  $B \in \mathcal{L}is_c(\mathscr{W}, \mathscr{W}')$ . The cost of the application of  $G_{\mathcal{T}}$  is the cost of the application of  $B_{\mathcal{T}}^{\mathscr{U}}$  plus cost that scales linearly in  $\#\mathcal{T}$ .

# **4.4.3** A space $\mathscr{W}_{\mathcal{T}}$ that is enclosed in a space decomposable into the continuous piecewise linears and bubbles

We follow the same program as in the previous subsection Sect. 4.4.2 but now with  $\mathscr{U}_{\mathcal{T}} := \mathscr{S}_{\mathcal{T}}^{0,1}$ , being the space of continuous piecewise linears.

Other than in Sect. 4.4.2, we cannot apply Proposition 4.4.1 for  $Q_T$  being the orthogonal projector onto  $\mathscr{U}_T$ , since with the current choice of this space it will not be a local projector. As an alternative, we take  $Q_T$  to be some biorthogonal projector. The question whether it enjoys an approximation property is answered in the following lemma.

**Lemma 4.4.9.** For  $\nu \in N_{\mathcal{T}}$ , so including vertices on  $\gamma$ , let  $\phi_{\nu} \in L_2(\Omega)$  be such that

(4.11) 
$$\|\widetilde{\phi}_{\nu}\|_{L_{2}(\Omega)} \lesssim h_{\nu}^{d/2}, \quad \sum_{\nu \in N_{\mathcal{T}}} \widetilde{\phi}_{\nu} = \mathbb{1}_{\Omega}, \quad \operatorname{supp} \widetilde{\phi}_{\nu} \subset B(\nu; Rh_{\nu})$$

for some constant R > 0, and

$$\left| \langle \widetilde{\phi}_{\nu}, \phi_{\nu'} \rangle_{L_2(\Omega)} \right| = \delta_{\nu\nu'} |\omega_{\nu}| \quad (\nu, \nu' \in N^0_{\mathcal{T}}).$$

Denote  $\widetilde{\mathscr{U}_{\mathcal{T}}} := \operatorname{span}\{\widetilde{\phi}_{\nu} : \nu \in N_{\mathcal{T}}^{0}\}$ , so without vertices on  $\gamma$ . The biorthogonal projector  $Q_{\mathcal{T}} : u \mapsto \sum_{\nu \in N_{\mathcal{T}}^{0}} \frac{\langle u, \widetilde{\phi}_{\nu} \rangle_{L_{2}(\Omega)}}{\langle \phi_{\nu}, \widetilde{\phi}_{\nu} \rangle_{L_{2}(\Omega)}} \phi_{\nu}$ , for which  $\operatorname{ran} Q_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1}$  and  $\operatorname{ran}(\operatorname{Id} - Q_{\mathcal{T}}) = \widetilde{\mathscr{U}_{\mathcal{T}}^{\perp}}_{\mathcal{T}}^{L_{2}(\Omega)}$ , satisfies the approximation property

$$\|h_{\mathcal{T}}^{-1}(\mathrm{Id}-Q_{\mathcal{T}}')v\|_{\mathcal{L}(H^{1}_{0,\gamma}(\Omega),L_{2}(\Omega))} \lesssim 1,$$

and  $\|Q_{\mathcal{T}}\|_{\mathcal{L}(L_2(\Omega), L_2(\Omega))} \lesssim 1.$ 

*Proof.* We use the same strategy as in Chapter 2. That is, we define a Scott-Zhang type quasi-interpolator  $\Pi_{\mathcal{T}} \colon H^1(\Omega) \to L_2(\Omega)$ , cf. [SZ90]. For every  $\nu \in N_{\mathcal{T}}$ , select a (d-1)-face  $e_{\nu}$  of some  $T \in \mathcal{T}$  with  $\nu \in e_{\nu}$  and  $e_{\nu} \subset \gamma$  if  $\nu \in \gamma$ . We define  $\Pi_{\mathcal{T}}$  by

$$\Pi_{\mathcal{T}} u := \sum_{\nu \in N_{\mathcal{T}}} g_{\mathcal{T},\nu}(u) \widetilde{\phi}_{\mathcal{T},\nu}, \quad g_{\mathcal{T},\nu}(u) := \oint_{e_{\nu}} u \, \mathrm{d}s.$$

Since  $g_{\mathcal{T},\nu}(1) = 1$ , using the properties from (4.11) one can show, cf. proof of Theorem 2.5.1 for details, that

$$\|h_{\mathcal{T}}^{-1}(\mathrm{Id} - \Pi_{\mathcal{T}})(u)\|_{L_2(\Omega)} \lesssim \|u\|_{H^1(\Omega)} \quad (u \in H^1(\Omega)).$$

By construction,  $g_{\mathcal{T},\nu}(u) = 0$  for  $\nu$  on  $\gamma$  and  $u \in H^1_{0,\gamma}(\Omega)$ , and therefore ran  $\Pi_{\mathcal{T}}|_{H^1_{0,\gamma}(\Omega)} \subset \widetilde{\mathscr{U}_{\mathcal{T}}}$ . Finally, combined with  $L_2(\Omega)$ -boundedness and locality of  $Q'_{\mathcal{T}}$ , and the fact that  $Q'_{\mathcal{T}}$  reproduces  $\widetilde{\mathscr{U}_{\mathcal{T}}}$ , we find that

$$\begin{split} \|h_{\mathcal{T}}^{-1}(\mathrm{Id} - Q_{\mathcal{T}}')v\|_{L_{2}(\Omega)} &= \inf_{w_{\mathcal{T}} \in \widetilde{\mathscr{U}_{\mathcal{T}}}} \|h_{\mathcal{T}}^{-1}(\mathrm{Id} - Q_{\mathcal{T}}')(v - w_{\mathcal{T}})\|_{L_{2}(\Omega)} \\ &\lesssim \|h_{\mathcal{T}}^{-1}(\mathrm{Id} - \Pi_{\mathcal{T}})(v)\|_{L_{2}(\Omega)} \lesssim \|v\|_{H^{1}_{0,\gamma}(\Omega)} \ (v \in H^{1}_{0,\gamma}(\Omega)). \end{split}$$

The last statement can be proven similarly as in the proof of Theorem 4.3.1.  $\Box$ 

As before, let  $\mathcal{T}_* \succ \mathcal{T}$  denote a uniform red-refinement of  $\mathcal{T}$ , and for any  $T \in \mathcal{T}$  and  $\nu \in N_T$ , let  $T_{*,\nu} \in \mathcal{T}_*$  denote the simplex with  $\nu \in T_{*,\nu} \subset T$ . For  $\nu \in N_T$ , so including boundary vertices, define

$$\widetilde{\phi}_{\mathcal{T},\nu} \coloneqq \frac{1}{d+1} \sum_{\substack{T \in \mathcal{T} \\ T \subset \omega_{\nu}}} \left( \mathbb{1}_T + \frac{d2^{1+d}}{d+1} \mathbb{1}_{T_{*,\nu}} - \frac{2^{1+d}}{d+1} \sum_{\substack{\nu' \in N_T \\ \nu' \neq \nu}} \mathbb{1}_{T_{*,\nu'}} \right) \in \mathscr{S}_{\mathcal{T}_*}^{-1,0}$$

These functions satisfy (4.11), and

$$\langle \phi_{\mathcal{T},\nu}, \phi_{\mathcal{T},\nu'} \rangle_{L_2(\Omega)} = \delta_{\nu\nu'} (d+1)^{-1} |\omega_{\mathcal{T},\nu}|,$$

and so determine a valid biorthogonal projector  $Q_T$  via Lemma 4.4.9.

For  $\mathcal{T}_{**} \succ \mathcal{T}_*$  a uniform red-refinement of  $\mathcal{T}_*$ , we define  $\Theta_{\mathcal{T}} := \{\theta_{\mathcal{T},\nu} : \nu \in N_{\mathcal{T}}^0\}$  by

$$\theta_{\mathcal{T},\nu} \coloneqq \frac{2^{d+2}}{d+2} \left( 2^d \mathbb{1}_{\omega_{\mathcal{T}_{**},\nu}} - \mathbb{1}_{\omega_{\mathcal{T}_{*},\nu}} \right)$$

Since red-refinement subdivides each simplex into d subsimplices, one infers that

(4.12) 
$$\mathscr{B}_{\mathcal{T}} \coloneqq \operatorname{span} \Theta_{\mathcal{T}} \perp_{L_2(\Omega)} \mathscr{S}_{\mathcal{T}_*}^{-1,0},$$

so that in particular  $\mathscr{B}_{\mathcal{T}} \subset \ker Q_{\mathcal{T}}$ .

Defining  $\Psi_{\mathcal{T}} := \{ \psi_{\mathcal{T},\nu} \colon \nu \in N^0_{\mathcal{T}} \}$  by

$$\psi_{\mathcal{T},\nu} := \phi_{\mathcal{T},\nu} + \theta_{\mathcal{T},\nu},$$

calculations as in the proof of Lemma 4.4.6 show the following result.

**Lemma 4.4.10.** The collection  $\Psi_{\mathcal{T}}$  satisfies (4.3) with supp  $\psi_{\mathcal{T},\nu} = \omega_{\mathcal{T},\nu}$  and

$$\langle \psi_{\mathcal{T},\nu}, \phi_{\mathcal{T},\nu'} \rangle_{L_2(\Omega)} = \delta_{\nu\nu'} (d+1)^{-1} |\omega_{\mathcal{T},\nu}| \quad (\nu,\nu' \in N^0_{\mathcal{T}}).$$

So the Fortin interpolator is uniformly bounded, and  $D_T$  is represented by a diagonal matrix. Next we verify the conditions imposed in Sect. 4.4.1 for the construction of  $B_T$ .

**Proposition 4.4.11.** Let  $\mathscr{U}_{\mathcal{T}}$ ,  $Q_{\mathcal{T}}$ ,  $\mathscr{B}_{\mathcal{T}}$ , and  $\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}}$  be defined as above. Then  $\mathscr{W}_{\mathcal{T}} \subset \mathscr{Z}_{\mathcal{T}} := \mathscr{U}_{\mathcal{T}} + \mathscr{B}_{\mathcal{T}}$  ((4.6)), the conditions of Lemma 4.4.3 are satisfied, in particular  $\Phi_{\mathcal{T}} = Q_{\mathcal{T}}\Psi_{\mathcal{T}}$  and so  $\Theta_{\mathcal{T}} = (\operatorname{Id} - Q_{\mathcal{T}})\Psi_{\mathcal{T}}$ , and lastly,  $\Theta_{\mathcal{T}}$  is an  $\|h_{\mathcal{T}}^s \cdot \|_{L_2(\Omega)}$ -orthogonal basis for  $\mathscr{B}_{\mathcal{T}}$  as required for Lemma 4.4.4.

*Proof.* The first statement is obviously true. We have already verified the first two conditions of Lemma 4.4.3. The third condition follows from this inverse inequality on  $\mathscr{S}_{\mathcal{T}_{**}}^{-1,1}$  (see e.g. (2.44)), and  $\Phi_{\mathcal{T}} = Q_{\mathcal{T}} \Psi_{\mathcal{T}}$  is a consequence of (4.12). The last statement follows from  $|\operatorname{supp} \theta_{\nu} \cap \operatorname{supp} \theta_{\nu'}| = 0$  when  $\nu \neq \nu'$ .

#### Implementation

Suppose that we have some operator  $B_{\mathcal{T}}^{\mathscr{U}} \in \mathcal{L}is_c(\mathscr{U}_{\mathcal{T}}, \mathscr{U}_{\mathcal{T}}')$  uniformly (here  $\mathscr{U}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1}$ ). The matrix representation of the preconditioner  $G_{\mathcal{T}}$ , with  $B_{\mathcal{T}}$  from (4.7) and the bases from Proposition 4.4.11, becomes

$$\begin{aligned} \boldsymbol{G}_{\mathcal{T}} &= \boldsymbol{D}_{\mathcal{T}}^{-1} \boldsymbol{B}_{\mathcal{T}} \boldsymbol{D}_{\mathcal{T}}^{-\top}, \\ \boldsymbol{B}_{\mathcal{T}} &\coloneqq \mathcal{F}'_{\Psi_{\mathcal{T}}} (\boldsymbol{Q}'_{\mathcal{T}} \boldsymbol{B}_{\mathcal{T}}^{\mathscr{U}} \boldsymbol{Q}_{\mathcal{T}} + (\mathrm{Id} - \boldsymbol{Q}_{\mathcal{T}})' \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} (\mathrm{Id} - \boldsymbol{Q}_{\mathcal{T}})) \mathcal{F}_{\Psi_{\mathcal{T}}} \\ &= \boldsymbol{B}_{\mathcal{T}}^{\mathscr{U}} + \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}}, \end{aligned}$$

with these matrices given by

$$\begin{aligned} \boldsymbol{D}_{\mathcal{T}} &:= \mathcal{F}'_{\Psi_{\mathcal{T}}} D_{\mathcal{T}} \mathcal{F}_{\Phi_{\mathcal{T}}} = \operatorname{diag} \left\{ \frac{|\omega_{\nu}|}{d+1} \colon \nu \in N_{\mathcal{T}}^{0} \right\}, \\ \boldsymbol{B}_{\mathcal{T}}^{\mathscr{U}} &:= \mathcal{F}'_{\Phi_{\mathcal{T}}} B_{\mathcal{T}}^{\mathscr{U}} \mathcal{F}_{\Phi_{\mathcal{T}}}, \quad \boldsymbol{B}_{\mathcal{T}}^{\mathscr{B}} \coloneqq \mathcal{F}'_{\Theta_{\mathcal{T}}} B_{\mathcal{T}}^{\mathscr{B}} \mathcal{F}_{\Theta_{\mathcal{T}}} = \beta_{1} \boldsymbol{D}_{\mathcal{T}}^{1+\frac{2s}{d}}, \end{aligned}$$

where we used that  $\mathcal{F}_{\Phi_{\mathcal{T}}}^{-1}Q_{\mathcal{T}}\mathcal{F}_{\Psi_{\mathcal{T}}} = \mathbf{Id}$  and  $\mathcal{F}_{\Theta_{\mathcal{T}}}^{-1}(\mathrm{Id} - Q_{\mathcal{T}})\mathcal{F}_{\Psi_{\mathcal{T}}} = \mathbf{Id}$ , and where, based on  $\|h_{\mathcal{T}}^s\theta_{\nu}\|_{L_2(\Omega)}^2 \approx |\omega_{\nu}|^{1+\frac{2s}{d}}$ , we made an harmless modification to the operator  $B_{\mathcal{T}}^{\mathscr{B}}$  from Lemma 4.4.4.

# 4.5 Extensions

#### 4.5.1 Higher order

Add the superscript 1 to the spaces defined so far, e.g. write  $\mathscr{V}_{\mathcal{T}}^1$  for  $\mathscr{S}_{\mathcal{T}}^{0,1}$  with its nodal basis  $\Phi_{\mathcal{T}}^1$ , and similarly use  $G_{\mathcal{T}}^1$  for the associated preconditioner from either Sect. 4.4.2 or Sect. 4.4.3.

We will now consider a (family of) higher order continuous piecewise polynomials, i.e. for some  $\ell \in \{2, 3, ...\}$  let

$$\mathscr{V}_{\mathcal{T}} = \mathscr{S}^{0,\ell}_{\mathcal{T}} := \{ u \in H^1_{0,\gamma}(\Omega) \colon u |_T \in \mathcal{P}_{\ell} \ (T \in \mathcal{T}) \} \subset \mathscr{V}.$$

Because we have an inverse inequality on  $\mathscr{V}_{\mathcal{T}}$ , we can construct a uniform preconditioner  $G_{\mathcal{T}} \in \mathcal{L}is(\mathscr{V}'_{\mathcal{T}}, \mathscr{V}_{\mathcal{T}})$  using an additive subspace correction method. That is, we consider the overlapping decomposition  $\mathscr{V}_{\mathcal{T}} = \mathscr{V}_{\mathcal{T}}^1 + \mathscr{V}_{\mathcal{T}}^2$ , where these spaces are given by

$$\mathscr{V}_{\mathcal{T}} = (\mathscr{V}_{\mathcal{T}}, \|\cdot\|_{\mathscr{V}}), \quad \mathscr{V}_{\mathcal{T}}^{1} = (\mathscr{V}_{\mathcal{T}}^{1}, \|\cdot\|_{\mathscr{V}}), \quad \mathscr{V}_{\mathcal{T}}^{2} = (\mathscr{V}_{\mathcal{T}}, \|h_{\mathcal{T}}^{-s} \cdot\|_{L_{2}(\Omega)}).$$

**Proposition 4.5.1.** For  $k \in \{1, 2\}$ , let  $G_{\mathcal{T}}^k \in \mathcal{L}is_c((\mathscr{V}_{\mathcal{T}}^k)', \mathscr{V}_{\mathcal{T}}^k)$ , then for  $I_{\mathcal{T}}^k : \mathscr{V}_{\mathcal{T}}^k \to \mathscr{V}_{\mathcal{T}}$  the trivial embedding, we find that  $G_{\mathcal{T}} := \sum_{k=1}^2 I_{\mathcal{T}}^k G_{\mathcal{T}}^k(I_{\mathcal{T}}^k)' \in \mathcal{L}is_c(\mathscr{V}_{\mathcal{T}}^{\prime}, \mathscr{V}_{\mathcal{T}})$ , with

$$\begin{aligned} \|G_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}',\mathscr{V}_{\mathcal{T}})} &\lesssim \max_{k=1,2} \|G_{\mathcal{T}}^{k}\|_{\mathcal{L}((\mathscr{V}_{\mathcal{T}}^{k})',\mathscr{V}_{\mathcal{T}}^{k})}, \\ \|\Re(G_{\mathcal{T}})^{-1}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}},\mathscr{V}_{\mathcal{T}}')} &\lesssim \max_{k=1,2} \|\Re(G_{\mathcal{T}}^{k})^{-1}\|_{\mathcal{L}(\mathscr{V}_{\mathcal{T}}^{k},(\mathscr{V}_{\mathcal{T}}^{k})')} \end{aligned}$$

*Proof.* We have the (standard) inverse inequality  $||u||_{\mathscr{V}} \lesssim ||h_{\mathcal{T}}^{-s}u||_{L_2(\Omega)}$  for  $u \in \mathscr{V}_{\mathcal{T}}$ . Let  $u \in \mathscr{V}_{\mathcal{T}}$ , then for any  $(u_1, u_2) \in \mathscr{V}_{\mathcal{T}}^1 \times \mathscr{V}_{\mathcal{T}}^2$  with  $u_1 + u_2 = u$  we find

$$||u||_{\mathscr{V}} \leq ||u_1||_{\mathscr{V}} + ||u_2||_{\mathscr{V}} \lesssim ||u_1||_{\mathscr{V}} + ||h_{\mathcal{T}}^{-s}u_2||_{L_2(\Omega)}.$$

Denote  $\Pi^1_{\mathcal{T}} \colon H^1_{0,\gamma}(\Omega) \to \mathscr{V}^1_{\mathcal{T}}$  for the Scott-Zhang interpolator ([SZ90]). For  $u \in \mathscr{V}_{\mathcal{T}}$ , take  $u_1 = \Pi^1_{\mathcal{T}} u \in \mathscr{V}^1_{\mathcal{T}}$  and  $u_2 = u - \Pi^1_{\mathcal{T}} u \in \mathscr{V}^2_{\mathcal{T}}$ , then from approximation properties of the interpolator we infer

$$\begin{aligned} \|u_1\|_{\mathscr{V}} + \|h_{\mathcal{T}}^{-s} u_2\|_{L_2(\Omega)} &\leq \|u\|_{\mathscr{V}} + \|u_2\|_{\mathscr{V}} + \|h_{\mathcal{T}}^{-s} u_2\|_{L_2(\Omega)} \\ &\lesssim \|u\|_{\mathscr{V}} + \|h_{\mathcal{T}}^{-s} u_2\|_{L_2(\Omega)} \lesssim \|u\|_{\mathscr{V}}. \end{aligned}$$

Since apparently for  $u \in \mathscr{V}_{\mathcal{T}}$ ,

$$\|u\|_{\mathscr{V}} = \inf \left\{ \|u_1\|_{\mathscr{V}} + \|h_{\mathcal{T}}^{-s} u_2\|_{L_2(\Omega)} \colon u_1 \in \mathscr{V}_1, \, u_2 \in \mathscr{V}_2, \, u_1 + u_2 = u \right\},$$

the result follows from subspace correction methods theory, e.g. [Osw94].  $\Box$ 

On the space  $\mathscr{V}_{\mathcal{T}}^1$  we can apply the operator  $G_{\mathcal{T}}^1$  constructed earlier, whereas on  $\mathscr{V}_{\mathcal{T}}^2$  a simple scaling operator suffices. Denote  $N_{\mathcal{T}}^{0,\ell}$  for the set of canonical Lagrange evaluation points of  $\mathscr{S}_{\mathcal{T}}^{0,\ell}$ , and let  $\Phi_{\mathcal{T}}^\ell = \{\phi_{\nu}^\ell \colon \nu \in N_{\mathcal{T}}^{0,\ell}\}$  be the corresponding nodal basis. For some constant  $\beta_2 > 0$ , define an operator  $R_{\mathcal{T}} \colon \mathscr{V}_{\mathcal{T}}^2 \to (\mathscr{V}_{\mathcal{T}}^2)'$  by

$$(R_{\mathcal{T}}u)(w) := \beta_2^{-1} \sum_{\nu \in N_{\mathcal{T}}^{0,\ell}} \|h_{\mathcal{T}}^{-s} \phi_{\nu}^{\ell}\|_{L_2(\Omega)}^2 u(\nu) w(\nu).$$

**Proposition 4.5.2.** The operator  $G_{\mathcal{T}}^2 := R_{\mathcal{T}}^{-1}$  satisfies  $G_{\mathcal{T}}^2 \in \mathcal{L}is_c((\mathscr{V}_{\mathcal{T}}^2)', \mathscr{V}_{\mathcal{T}}^2)$  uniformly.

*Proof.* It is not hard to see that the result follows if  $\Phi_{\mathcal{T}}^{\ell}$  is a (uniformly)  $\|h_{\mathcal{T}}^{-s} \cdot \|_{L_2(\Omega)}$ -stable basis. Writing  $N_T^{0,\ell} := N_{\mathcal{T}}^{0,\ell} \cap T$ , this stability can be deduced from

$$\begin{split} \left\| h_{\mathcal{T}}^{-s} \sum_{\nu \in N_{\mathcal{T}}^{0,\ell}} c_{\nu} \phi_{\nu}^{\ell} \right\|_{L_{2}(\Omega)}^{2} &= \sum_{T \in \mathcal{T}} h_{T}^{-2s} \left\| \sum_{\nu \in N_{T}^{0,\ell}} c_{\nu} \phi_{\nu}^{\ell} \right\|_{L_{2}(T)}^{2} \approx \sum_{T \in \mathcal{T}} h_{T}^{-2s} \sum_{\nu \in N_{T}^{0,\ell}} |c_{\nu}|^{2} \| \phi_{\nu}^{\ell} \|_{L_{2}(T)}^{2} \\ &= \sum_{\nu \in N_{\mathcal{T}}^{0,\ell}} |c_{\nu}|^{2} \| h_{\mathcal{T}}^{-s} \phi_{\nu}^{\ell} \|_{L_{2}(\Omega)}^{2}. \end{split}$$

#### Implementation

Equipping  $\mathscr{V}_{\mathcal{T}}$  and  $\mathscr{V}_{\mathcal{T}}^2$  with  $\Phi_{\mathcal{T}}^l$ , and  $\mathscr{V}_{\mathcal{T}}^1$  with  $\Phi_{\mathcal{T}}^1$ , the matrix representation of  $G_{\mathcal{T}} := \sum_{k=1}^2 I_{\mathcal{T}}^k G_{\mathcal{T}}^k (I_{\mathcal{T}}^k)' \in \mathcal{L}is_c(\mathscr{V}_{\mathcal{T}}', \mathscr{V}_{\mathcal{T}})$  is given by

$$\boldsymbol{G}_{\mathcal{T}} = \boldsymbol{q}_{\mathcal{T}} \boldsymbol{G}_{\mathcal{T}}^1 \boldsymbol{q}_{\mathcal{T}}^\top + \boldsymbol{G}_{\mathcal{T}}^2,$$

with  $G_{T}^{1}$  either from Sect. 4.4.2 or Sect. 4.4.3,

$$(q_{\mathcal{T}})_{\nu'\nu} = \phi_{\nu'}^{\ell}(\nu) \quad (\nu' \in N_{\mathcal{T}}^{0,\ell}, \nu \in N_{\mathcal{T}}^{0,1}).$$

and

$$\boldsymbol{G}_{\mathcal{T}}^{2} = \beta_{2} \operatorname{diag}\{\|h_{\mathcal{T}}^{-s}\phi_{\nu}^{\ell}\|_{L_{2}(\Omega)}^{-2} \colon \nu \in N_{\mathcal{T}}^{0,\ell}\}.$$

#### 4.5.2 Manifolds

Let  $\Gamma$  be a compact *d*-dimensional Lipschitz, piecewise smooth manifold in  $\mathbb{R}^{d'}$  for some  $d' \geq d$  with or without boundary  $\partial \Gamma$ . For some closed measurable  $\gamma \subset \partial \Gamma$  and  $s \in [0, 1]$ , let

$$\mathscr{V} := [L_2(\Gamma), H^1_{0,\gamma}(\Gamma)]_{s,2}, \quad \mathscr{W} := \mathscr{V}'.$$

We assume that  $\Gamma$  is given as the closure of the disjoint union of  $\bigcup_{i=1}^{p} \chi_i(\Omega_i)$ , with, for  $1 \leq i \leq p, \chi_i \colon \mathbb{R}^d \to \mathbb{R}^{d'}$  being some smooth regular parametrization, and  $\Omega_i \subset \mathbb{R}^d$  an open polytope. W.l.o.g. assuming that for  $i \neq j, \overline{\Omega}_i \cap \overline{\Omega}_j = \emptyset$ , we define

$$\chi \colon \Omega := \cup_{i=1}^p \Omega_i \to \cup_{i=1}^p \chi_i(\Omega_i) \text{ by } \chi|_{\Omega_i} = \chi_i.$$

Let  $\mathbb{T}$  be a family of conforming partitions  $\mathcal{T}$  of  $\Gamma$  into 'panels' such that, for  $1 \leq i \leq p, \chi^{-1}(\mathcal{T}) \cap \Omega_i$  is a uniformly shape regular conforming partition of  $\Omega_i$  into *d*-simplices (that for d = 1 satisfies a uniform *K*-mesh property). We assume that  $\gamma$  is a (possibly empty) union of 'faces' of  $T \in \mathcal{T}$  (i.e., sets of type  $\chi_i(e)$ , where *e* is a (d-1)-dimensional face of  $\chi_i^{-1}(T)$ ).

The usual lowest order boundary element spaces are defined by

$$\mathscr{S}_{\mathcal{T}}^{-1,0} := \{ u \in L_2(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_0 \ (T \in \mathcal{T}) \}, \\ \mathscr{S}_{\mathcal{T}}^{0,1} := \{ u \in H^1_{0,\gamma}(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_1 \ (T \in \mathcal{T}) \},$$

with their canonical bases denoted as  $\Sigma_{\mathcal{T}} = \{\mathbb{1}_T : T \in \mathcal{T}\}\$  and  $\Phi_{\mathcal{T}} = \{\phi_{\nu} : \nu \in N^0_{\mathcal{T}}\}\$ , respectively, with  $N^0_{\mathcal{T}}$  the vertices of  $\mathcal{T}$  not on  $\gamma$ .

The construction of the preconditioners in the domain case relied on the explicit construction of a collection  $\Psi_{\mathcal{T}}$  biorthogonal to  $\Phi_{\mathcal{T}}$ , and on the explicit computation of a (bi)orthogonal projection of  $\mathscr{W}_{\mathcal{T}} := \operatorname{span} \Psi_{\mathcal{T}}$  onto either  $\mathscr{S}_{\mathcal{T}}^{-1,0}$  or  $\mathscr{S}_{\mathcal{T}}^{0,1}$ , where orthogonality was interpreted w.r.t. the  $L_2(\Omega)$ -scalar product. Both the construction of  $\Psi_{\mathcal{T}}$  and the computation of the (bi)orthogonal projection could be reduced to computations on the individual elements in the partition, which yielded explicit expressions.

When attempting to transfer everything to the manifold case, a problem is the appearance of a generally non-constant weight  $x \mapsto |\partial \chi(x)|$  in the  $L_2(\Gamma)$ scalar product

$$\langle u, v \rangle_{L_2(\Gamma)} = \int_{\Omega} u(\chi(x)) v(\chi(x)) |\partial \chi(x)| \, dx$$

To deal with this, following Sect. 2.3.2, on  $L_2(\Gamma)$  we define an additional 'meshdependent' scalar product

$$\langle u, v \rangle_{\mathcal{T}} := \sum_{T \in \mathcal{T}} \frac{|T|}{|\chi^{-1}(T)|} \int_{\chi^{-1}(T)} u(\chi(x)) v(\chi(x)) dx$$

which is constructed by replacing on each  $\chi^{-1}(T)$ , the Jacobian  $|\partial \chi|$  by its average  $\frac{|T|}{|\chi^{-1}(T)|}$  over  $\chi^{-1}(T)$ , and interpret (bi)orthogonality with respect to this scalar product.

Now all steps in the *construction* of the preconditioners carry over, and yield preconditioners for the manifold case whose implementations are exactly as described in Sect. 4.4.2 and Sect. 4.4.3, where the patch volumes  $|\omega_{\mathcal{T},\nu}|$  now should be read as the volumes of the patches on  $\Gamma$ .

To *prove* that the constructed preconditioners are indeed uniform preconditioners requires additional work due to the use of the mesh-dependent scalar product. We refer to Chapter 2 for details. The key ingredient is that not only the norm associated to  $\langle \cdot, \cdot \rangle_{L_2(\Gamma)}$  is uniformly equivalent to  $\| \cdot \|_{L_2(\Gamma)}$ , but also that  $\langle \cdot, \cdot \rangle_{L_2(\Gamma)}$  and  $\langle \cdot, \cdot \rangle_{T}$  are close in the sense that

$$|\langle v, u \rangle_{\mathcal{T}} - \langle v, u \rangle_{L_2(\Gamma)}| \lesssim \|h_{\mathcal{T}}v\|_{L_2(\Gamma)} \|u\|_{L_2(\Gamma)} \quad (v, u \in L_2(\Gamma)).$$

# 4.6 Numerical experiments

Let  $\Gamma = \partial [0,1]^3 \subset \mathbb{R}^3$  be the boundary of the unit cube,  $\mathscr{V} := H^{1/2}(\Gamma)$ ,  $\mathscr{W} := H^{-1/2}(\Gamma)$ , and  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,\ell} \subset \mathscr{V}$  the trial space of continuous piecewise polynomials of degree  $\ell$  w.r.t. a partition  $\mathcal{T}$ . We shall evaluate preconditioning of essentially a discretized Hypersingular Integral operator.

The Hypersingular Integral operator  $A \in \mathcal{L}(\mathcal{V}, \mathcal{V}')$  is only semi-coercive, since it has a non-trivial kernel equal to span{1}. Solving  $\tilde{A}u = f$  for fwith f(1) = 0 is, however, equivalent to solving Au = f with A given by  $(Au)(v) = (\tilde{A}u)(v) + \alpha \langle u, 1 \rangle_{L_2(\Gamma)} \langle v, 1 \rangle_{L_2(\Gamma)}$  for some  $\alpha > 0$ . This operator A is in  $\mathcal{L}is_c(\mathcal{V}, \mathcal{V}')$ , and we shall consider preconditioning discretizations  $A_{\mathcal{T}} \in \mathcal{L}is_c(\mathcal{V}_{\mathcal{T}}, \mathcal{V}'_{\mathcal{T}})$  of A. By comparing different values numerically, we found  $\alpha = 0.05$  to give good results in our examples.

As opposite order operator *B* we take the Weakly Singular integral operator, which on compact 2-dimensional manifolds is known to be in  $\mathcal{L}is_c(\mathcal{W}, \mathcal{W}')$ . We will compare preconditioners  $G_{\mathcal{T}}$  based on the discretizations  $B_{\mathcal{T}}^{\mathcal{U}} \in \mathcal{L}is_c(\mathcal{U}_{\mathcal{T}}, \mathcal{U}_{\mathcal{T}}')$  of *B*, for  $\mathcal{U}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,0}$  or  $\mathcal{U}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1}$  equipped with the canonical bases  $\Sigma_{\mathcal{T}} = \{\mathbb{1}_T : T \in \mathcal{T}\}$  and  $\Phi_{\mathcal{T}} = \{\phi_{\nu} : \nu \in N_{\mathcal{T}}\}$ , respectively, cf. Sect. 4.4.2 or Sect. 4.4.3.

For  $\ell = 1$  (the lowest order case) and  $\mathscr{V}_{\mathcal{T}}$  being equipped with  $\Phi_{\mathcal{T}}$ , the matrix representation of the preconditioner  $G_{\mathcal{T}}$  reads either as (Sect. 4.4.2)

$$oldsymbol{G}_{\mathcal{T}} = oldsymbol{G}_{\mathcal{T}}^{-1,0} = oldsymbol{D}_{\mathcal{T}}^{-1} ig( oldsymbol{p}_{\mathcal{T}}^{ op} oldsymbol{B}_{\mathcal{T}}^{\mathscr{U}} oldsymbol{p}_{\mathcal{T}} + eta_1 oldsymbol{D}_{\mathcal{T}}^{3/2} ig) oldsymbol{D}_{\mathcal{T}}^{-1}$$

where  $\boldsymbol{B}_{\mathcal{T}}^{\mathscr{U}} = (B\Sigma_{\mathcal{T}})(\Sigma_{\mathcal{T}}), \boldsymbol{D}_{\mathcal{T}} = \operatorname{diag}\{|\omega_{\nu}| : \nu \in N_{\mathcal{T}}\}, (\boldsymbol{p}_{\mathcal{T}})_{T\nu} = \begin{cases} 1 & \text{if } T \subset \omega_{\nu}, \\ 0 & \text{otherwise,} \end{cases}$ and  $\beta_1 > 0$  is some constant, or as (Sect. 4.4.3)

$$oldsymbol{G}_{\mathcal{T}} = oldsymbol{G}_{\mathcal{T}}^{0,1} = oldsymbol{D}_{\mathcal{T}}^{-1} igl( oldsymbol{B}_{\mathcal{T}}^{\mathscr{U}} + eta_1 oldsymbol{D}_{\mathcal{T}}^{3/2} igr) oldsymbol{D}_{\mathcal{T}}^{-1}$$

where  $\boldsymbol{B}_{\mathcal{T}}^{\mathscr{U}} = (B\Phi_{\mathcal{T}})(\Phi_{\mathcal{T}}), \boldsymbol{D}_{\mathcal{T}} = \text{diag}\{|\frac{\omega_{\nu}}{d+1}| : \nu \in N_{\mathcal{T}}\}, \text{ and } \beta_1 > 0 \text{ is some constant.}$ 

For  $\ell > 1$  denote the above  $G_{\mathcal{T}}$  by either  $G_{\mathcal{T}}^{1,-1,0}$  or  $G_{\mathcal{T}}^{1,0,1}$ , then, with  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,\ell}$  being equipped with the standard nodal basis  $\{\phi_{\mathcal{L}}^{\ell} : \nu \in N_{\mathcal{T}}^{\ell}\}$ , the matrix representation of the preconditioner  $G_{\mathcal{T}} \in \mathcal{L}is_{c}((\mathscr{S}_{\mathcal{T}}^{0,\ell})', \mathscr{S}_{\mathcal{T}}^{0,\ell})$  from Sect. 4.5.1 is

$$\boldsymbol{G}_{\mathcal{T}}^{*} = \boldsymbol{q}_{\mathcal{T}}\boldsymbol{G}_{\mathcal{T}}^{1,*}\boldsymbol{q}_{\mathcal{T}}^{\top} + \beta_{2}\operatorname{diag}\{\|\boldsymbol{h}_{\mathcal{T}}^{-\frac{1}{2}}\boldsymbol{\phi}_{\nu}^{\ell}\|_{L_{2}(\Omega)}^{-2} \colon \nu \in N_{\mathcal{T}}^{\ell}\},\$$

where either \* = -1, 0 or \* = 0, 1, and  $(q_T)_{\nu'\nu} = \phi_{\nu'}^{\ell}(\nu) \ (\nu' \in N_T^{\ell}, \nu \in N_T^1)$ .

The (full) matrix representations of the discretized singular integral operators  $A_{\mathcal{T}}$  and  $B_{\mathcal{T}}^{\mathscr{U}}$  are calculated using the BEM++ software package [ŚBA<sup>+</sup>15]. Condition numbers are determined using Lanczos iteration with respect to  $\|\cdot\| := \|A_{\mathcal{T}}^{\frac{1}{2}} \cdot \|$ .

#### 4.6.1 Uniform refinements

Consider a conforming triangulation  $\mathcal{T}_1$  of  $\Gamma$  consisting of 2 triangles per side, so 12 triangles with 8 vertices in total. We let  $\mathbb{T}$  be the sequence  $\{\mathcal{T}_k\}_{k\geq 1}$  of uniform newest vertex bisections, where  $\mathcal{T}_k \succ \mathcal{T}_{k-1}$  is found by bisecting each triangle from  $\mathcal{T}_{k-1}$ .

With  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1}$ , Table 4.1 compares the condition numbers for the preconditioned system given by Sect. 4.4.2 ( $\mathscr{U}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,0}$ ) and by Sect. 4.4.3 ( $\mathscr{U}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1}$ ). We see that the condition numbers remain nicely bounded, and that both choices give similar condition numbers.

Instead of using the 'full matrices', we can consider compressed hierarchical matrices to approximate the stiffness matrices  $\mathbf{A}_{\mathcal{T}}$  and  $\mathbf{B}_{\mathcal{T}}^{\mathscr{U}}$  for finer partitions. Table 4.2 gives the condition numbers, again for uniform refinements, but now using hierarchical matrices based on adaptive cross approximation [Hac99, Beb00]. We see that even for large systems, our preconditioner gives very satisfactory results.

Finally, consider the (higher order) trial space  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,3}$ . Table 4.3 gives condition numbers for the preconditioned system, using the method described in Sect. 4.5.1.

TABLE 4.1. Spectral condition numbers of the preconditioned hypersingular system, using uniform refinements, discretized by continuous piecewise linears  $\mathscr{S}_{\mathcal{T}}^{0,1}$ , with  $\alpha = 0.05$ . The preconditioners  $G_{\mathcal{T}}^{-1,0}$  and  $G_{\mathcal{T}}^{0,1}$  are constructed using the single layer operator discretized on  $\mathscr{U}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{-1,0}$  (Sect. 4.4.2) and  $\mathscr{U}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1}$  (Sect 4.4.3), respectively, where have used  $\beta_1 = 0.65$  in the first case and  $\beta_1 = 0.34$  in the second case.

dofs	$\kappa_S({oldsymbol A}_{\mathcal T})$	$\kappa_S(oldsymbol{G}_\mathcal{T}^{-1,0}oldsymbol{A}_\mathcal{T})$	$\kappa_S(oldsymbol{G}_\mathcal{T}^{0,1}oldsymbol{A}_\mathcal{T})$
14	3.0	2.71	2.64
50	7.1	2.36	2.37
194	14.2	2.25	2.26
770	28.7	2.30	2.27
3074	57.8	2.29	2.27
12290	115.7	2.29	2.27
49154	231.4	2.30	2.27

TABLE 4.2. In the same setting as Table 4.1, but using compressed hierarchical matrices.

dofs	$\kappa_S({oldsymbol A}_{\mathcal T})$	$\kappa_S(oldsymbol{G}_\mathcal{T}^{-1,0}oldsymbol{A}_\mathcal{T})$	$\kappa_S(oldsymbol{G}_\mathcal{T}^{0,1}oldsymbol{A}_\mathcal{T})$
12290	115.6	2.29	2.27
24578	168.7	2.24	2.24
49154	231.3	2.30	2.27
98306	336.9	2.25	2.25
196610	461.7	2.30	2.28
393218	671.9	2.27	2.28
786434	751.6	2.30	2.30

TABLE 4.3. Spectral condition numbers of the preconditioned hypersingular system, using uniform refinements, discretized by continuous piecewise cubics  $\mathscr{S}_{\mathcal{T}}^{0,3}$ , with  $\alpha = 0.05$ . The higher order preconditioners  $G_{\mathcal{T}}^{-1,0}$  and  $G_{\mathcal{T}}^{0,1}$  are constructed as described in Sect. 4.5.1, by using the preconditioners from Table 4.1 with constants  $\beta_1 = 0.65, \beta_2 = 0.065$  in the first case and  $\beta_1 = 0.34, \beta_2 = 0.065$  in the second case.

dofs	$\kappa_S(oldsymbol{A}_\mathcal{T})$	$\kappa_S(oldsymbol{G}_\mathcal{T}^{-1,0}oldsymbol{A}_\mathcal{T})$	$\kappa_S(oldsymbol{G}_\mathcal{T}^{0,1}oldsymbol{A}_\mathcal{T})$
56	19.49	4.75	4.72
218	36.27	5.18	5.17
866	74.78	6.23	6.20
3458	150.73	6.55	6.48
13826	301.97	6.63	6.57
55298	603.86	6.65	6.58

TABLE 4.4. Spectral condition numbers of the preconditioned hypersingular system discretized by  $\mathscr{S}_{\mathcal{T}}^{0,1}$  using local refinements at each of the eight cube corners. Both preconditioners  $G_{\mathcal{T}}^{-1,0}$  and  $G_{\mathcal{T}}^{0,1}$  are constructed with same parameters as in Table 4.1, and are compared against diagonal preconditioning. The second column is defined by  $h_{\mathcal{T},min} := \min_{T \in \mathcal{T}} h_T$ .

dofs	$h_{\mathcal{T},min}$	$\kappa_S(\operatorname{diag}(\boldsymbol{A}_{\mathcal{T}})^{-1}\boldsymbol{A}_{\mathcal{T}})$	$\kappa_S(oldsymbol{G}_\mathcal{T}^{-1,0}oldsymbol{A}_\mathcal{T})$	$\kappa_S(oldsymbol{G}_{\mathcal{T}}^{0,1}oldsymbol{A}_{\mathcal{T}})$
8	$1.4\cdot 10^0$	2.15	2.83	2.68
14	$1.0\cdot 10^0$	2.79	2.71	2.64
314	$1.1\cdot 10^{-2}$	12.11	2.21	2.20
626	$1.2\cdot10^{-4}$	13.18	2.31	2.30
938	$1.3 \cdot 10^{-6}$	13.43	2.36	2.36
1250	$1.4 \cdot 10^{-8}$	13.51	2.39	2.38
1562	$1.6 \cdot 10^{-10}$	13.53	2.41	2.39
1850	$2.5\cdot10^{-12}$	13.55	2.41	2.40

# 4.6.2 Local refinements

Here we take  $\mathbb{T}$  to be the sequence  $\{\mathcal{T}_k\}_{k\geq 1}$  of locally refined triangulations, where  $\mathcal{T}_k \succ \mathcal{T}_{k-1}$  is constructed using conforming newest vertex bisection to refine all triangles in  $\mathcal{T}_{k-1}$  that touch a corner of the cube.

Table 4.4 gives condition numbers of the preconditioned hypersingular system discretized by continuous piecewise linears, i.e.  $\mathscr{V}_{\mathcal{T}} = \mathscr{S}_{\mathcal{T}}^{0,1}$ . The condition numbers remain bounded under local refinements, confirming uniformity of the preconditioner w.r.t.  $\mathbb{T}$ .

# 4.7 Conclusion

Using the framework of operator preconditioning, we have constructed uniform preconditioners for elliptic operators of orders  $2s \in [0, 2]$  discretized by continuous finite (or boundary) elements. The evaluation of the preconditioners requires the application of an opposite order operator plus minor cost of linear complexity. Compared to earlier proposals, both the construction of a so-called dual-mesh and the inversion of a non-diagonal matrix are avoided, and our results are valid without constraints on the mesh-grading. For lowest order finite elements the computed condition numbers of the preconditioned system are below 2.5.

# 5.1 Introduction

This chapter deals with the construction of uniform preconditioners for negative and positive order operators, discretized by continuous piecewise polynomial trial spaces, using the framework of 'operator preconditioning' [Hip06], see also [SW98, Ste02, BC07, HJHUT20].

For some *d*-dimensional closed domain (or manifold)  $\Omega$  and an  $s \in [0, 1]$ , we consider the (fractional) Sobolev space  $H^s(\Omega)$  and its dual that we denote by  $H^{-s}(\Omega)$ . Let  $(\mathscr{S}_{\mathcal{T}})_{\mathcal{T}\in\mathbb{T}}$  be a family of *continuous piecewise* polynomials of some fixed degree  $\ell$  w.r.t. uniformly shape regular, possibly locally refined, partitions.

Given some families of uniformly boundedly invertible operators

$$A_{\mathcal{T}} \colon \left(\mathscr{S}_{\mathcal{T}}, \|\cdot\|_{H^{-s}(\Omega)}\right) \to \left(\mathscr{S}_{\mathcal{T}}, \|\cdot\|_{H^{-s}(\Omega)}\right)', \\ B_{\mathcal{T}} \colon \left(\mathscr{S}_{\mathcal{T}}, \|\cdot\|_{H^{s}(\Omega)}\right) \to \left(\mathscr{S}_{\mathcal{T}}, \|\cdot\|_{H^{s}(\Omega)}\right)',$$

we are interested in constructing a *preconditioner* for  $A_{\mathcal{T}}$  using operator preconditioning with  $B_{\mathcal{T}}$ , and vice versa. To this end, we introduce a uniformly boundedly invertible operator  $D_{\mathcal{T}}: (\mathscr{S}_{\mathcal{T}}, \|\cdot\|_{H^{-s}(\Omega)}) \to (\mathscr{S}_{\mathcal{T}}, \|\cdot\|_{H^{s}(\Omega)})'$ , yielding preconditioned systems  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D'_{\mathcal{T}})^{-1}A_{\mathcal{T}}$  and  $(D'_{\mathcal{T}})^{-1}A_{\mathcal{T}}D_{\mathcal{T}}^{-1}B_{\mathcal{T}}$  that are uniformly boundedly invertible.

In Chapters 2 and 4 we already constructed such preconditioners in a more general setting where *different* ansatz spaces were used to define  $A_{\mathcal{T}}$  and  $B_{\mathcal{T}}$ . The setting studied in the current work, however, allows for preconditioners with a remarkably simple implementation.

A typical setting is that for some  $A: H^{-s}(\Omega) \to H^s(\Omega)$  and  $B: H^s(\Omega) \to H^{-s}(\Omega)$ , both boundedly invertible and coercive, it holds that  $(A_{\mathcal{T}}u)(v) := (Au)(v)$  and  $(B_{\mathcal{T}}u)(v) := (Bu)(v)$  with  $u, v \in \mathscr{S}_{\mathcal{T}}$ . An example for  $s = \frac{1}{2}$  is that A is the Single Layer Integral operator and B is the Hypersingular Integral operator. For this case, continuity of piecewise polynomial trial functions is required for discretizing B, but not for A, for which often discontinuous piecewise polynomials are employed. Nevertheless, when the solution of

the Single Layer Integral equation is expected to be smooth, e.g., when  $\Omega$  is a smooth manifold, then it is advantageous to take an ansatz space of continuous (or even smoother) functions also for *A*.

An obvious choice for  $D_{\mathcal{T}}$  would be to consider  $(D_{\mathcal{T}}u)(v) := \langle u, v \rangle_{L_2(\Omega)}$ . However, a problem becomes apparent when one considers the matrix representation  $D_{\mathcal{T}}$  of  $D_{\mathcal{T}}$  in the standard basis being the mass matrix: the inverse matrix  $D_{\mathcal{T}}^{-1}$ , that appears in the preconditioned system, is densely populated. In view of application cost, this inverse matrix has to be approximated, where it generally can be expected that, in order to obtain a uniform preconditioner, approximation errors have to decrease with a decreasing (minimal) mesh size, which will be confirmed in a numerical experiment. To circumvent this issue, we will introduce a  $D_{\mathcal{T}}$  that has a *diagonal* matrix representation, so that its inverse can be exactly evaluated.

#### 5.1.1 Notation

In this work, by  $\lambda \leq \mu$  we mean that  $\lambda$  can be bounded by a multiple of  $\mu$ , independently of parameters which  $\lambda$  and  $\mu$  may depend on, with the sole exception of the space dimension d, or in the manifold case, on the parametrization of the manifold that is used to define the finite element spaces on it. Obviously,  $\lambda \geq \mu$  is defined as  $\mu \leq \lambda$ , and  $\lambda = \mu$  as  $\lambda \leq \mu$  and  $\lambda \geq \mu$ .

For normed linear spaces  $\mathscr{Y}$  and  $\mathscr{Z}$ , in this chapter for convenience over  $\mathbb{R}$ ,  $\mathcal{L}(\mathscr{Y}, \mathscr{Z})$  will denote the space of bounded linear mappings  $\mathscr{Y} \to \mathscr{Z}$  endowed with the operator norm  $\|\cdot\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Z})}$ . The subset of invertible operators in  $\mathcal{L}(\mathscr{Y}, \mathscr{Z})$  with inverses in  $\mathcal{L}(\mathscr{Z}, \mathscr{Y})$  will be denoted as  $\mathcal{L}is(\mathscr{Y}, \mathscr{Z})$ .

For  $\mathscr{Y}$  a reflexive Banach space and  $C \in \mathcal{L}(\mathscr{Y}, \mathscr{Y}')$  being *coercive*, i.e.,

$$\inf_{0 \neq y \in \mathscr{Y}} \frac{(Cy)(y)}{\|y\|_{\mathscr{U}}^2} > 0,$$

both C and  $\Re(C) := \frac{1}{2}(C + C')$  are in  $\mathcal{L}is(\mathscr{Y}, \mathscr{Y}')$  with

$$\begin{aligned} \|\Re(C)\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')} &\leq \|C\|_{\mathcal{L}(\mathscr{Y},\mathscr{Y}')}, \\ \|C^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} &\leq \|\Re(C)^{-1}\|_{\mathcal{L}(\mathscr{Y}',\mathscr{Y})} = \Big(\inf_{0\neq y\in\mathscr{Y}} \frac{(Cy)(y)}{\|y\|_{\mathscr{Y}}^2}\Big)^{-1}. \end{aligned}$$

The subset of coercive operators in  $\mathcal{L}is(\mathscr{Y}, \mathscr{Y}')$  is denoted as  $\mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$ . If  $C \in \mathcal{L}is_c(\mathscr{Y}, \mathscr{Y}')$ , then  $C^{-1} \in \mathcal{L}is_c(\mathscr{Y}', \mathscr{Y})$  and  $\|\Re(C^{-1})^{-1}\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Y}')} \leq \|C\|_{\mathcal{L}(\mathscr{Y}, \mathscr{Y})}^2 \|\Re(C)^{-1}\|_{\mathcal{L}(\mathscr{Y}', \mathscr{Y})}$ .

Given a family of operators  $C_i \in \mathcal{L}is(\mathscr{Y}_i, \mathscr{Z}_i)$  ( $\mathcal{L}is_c(\mathscr{Y}_i, \mathscr{Z}_i)$ ), we will write  $C_i \in \mathcal{L}is(\mathscr{Y}_i, \mathscr{Z}_i)$  ( $\mathcal{L}is_c(\mathscr{Y}_i, \mathscr{Z}_i)$ ) uniformly in *i*, or simply 'uniform', when

$$\sup\max(\|C_i\|_{\mathcal{L}(\mathscr{Y}_i,\mathscr{Z}_i)},\|C_i^{-1}\|_{\mathcal{L}(\mathscr{Z}_i,\mathscr{Y}_i)})<\infty,$$

or

$$\sup_{i} \max(\|C_i\|_{\mathcal{L}(\mathscr{Y}_i,\mathscr{Z}_i)}, \|\Re(C_i)^{-1}\|_{\mathcal{L}(\mathscr{Z}_i,\mathscr{Y}_i)}) < \infty$$

# **5.2** Construction of $D_T$ in the domain case

For some *d*-dimensional domain  $\Omega$  and an  $s \in [0, 1]$ , we consider the Sobolev spaces

$$H^{s}(\Omega) := [L_{2}(\Omega), H^{1}(\Omega)]_{s,2}, \quad H^{-s}(\Omega) := H^{s}(\Omega)',$$

which form the Gelfand triple  $H^{s}(\Omega) \hookrightarrow L_{2}(\Omega) \simeq L_{2}(\Omega)' \hookrightarrow H^{-s}(\Omega)$ .

*Remark* 5.2.1. In this work, for convenience we restrict ourselves to Sobolev spaces with positive smoothness index which do not incorporate homogeneous Dirichlet boundary conditions and their duals. The proofs given below can however be extended to the setting with boundary conditions, see the arguments found in Chapters 2 and 4.

Let  $(\mathcal{T})_{\mathcal{T}\in\mathbb{T}}$  be a family of *conforming* partitions of  $\Omega$  into (open) *uniformly* shape regular *d*-simplices. Thanks to the conformity and the uniform shape regularity, for d > 1 we know that neighbouring  $T, T' \in \mathcal{T}$ , i.e.  $\overline{T} \cap \overline{T'} \neq \emptyset$ , have uniformly comparable sizes. For d = 1, we impose this uniform '*K*-mesh property' explicitly.

Fix  $\ell > 0$ . For  $\mathcal{T} \in \mathbb{T}$ , let  $\mathscr{S}_{\mathcal{T}}$  denote the space of *continuous piecewise* polynomials of degree  $\ell$  w.r.t.  $\mathcal{T}$ , i.e.,

$$\mathscr{S}_{\mathcal{T}} := \{ u \in H^1(\Omega) : u |_T \in \mathcal{P}_\ell \ (T \in \mathcal{T}) \}.$$

Additionally, for  $r \in [-1, 1]$ , we will write  $\mathscr{S}_{\mathcal{T},r}$  as shorthand notation for the normed linear space  $(\mathscr{S}_{\mathcal{T}}, \|\cdot\|_{H^r(\Omega)})$ .

Denote  $N_{\mathcal{T}}$  for the set of the usual Lagrange evaluation points of  $\mathscr{S}_{\mathcal{T}}$ , and equip the latter space with  $\Phi_{\mathcal{T}} = \{\phi_{\mathcal{T},\nu} : \nu \in N_{\mathcal{T}}\}$ , being the *canonical nodal basis* defined by  $\phi_{\mathcal{T},\nu}(\nu') := \delta_{\nu\nu'} (\nu,\nu' \in N_{\mathcal{T}})$ . For  $T \in \mathcal{T}$ , set  $h_T := |T|^{1/d}$  and let  $N_T := \overline{T} \cap N_{\mathcal{T}}$  be the set of evaluation points in  $\overline{T}$ . We will omit notational dependence on  $\mathcal{T}$  if it is clear from the context, e.g., we will simply write  $\phi_{\nu}$ .

#### 5.2.1 Operator preconditioning

Given some family of opposite order operators  $A_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},-s},(\mathscr{S}_{\mathcal{T},-s})')$ and  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},s},(\mathscr{S}_{\mathcal{T},s})')$ , both uniformly in  $\mathcal{T} \in \mathbb{T}$ , we are interested in constructing *optimal* preconditioners for both  $A_{\mathcal{T}}$  and  $B_{\mathcal{T}}$ , using the idea of opposite order preconditioning ([Hip06]).

That is, if one has an additional family of operators  $D_{\mathcal{T}} \in \mathcal{L}is(\mathscr{S}_{\mathcal{T},-s},(\mathscr{S}_{\mathcal{T},s})')$ uniformly in  $\mathcal{T} \in \mathbb{T}$ , then uniformly preconditioned systems for  $A_{\mathcal{T}}$  and  $B_{\mathcal{T}}$ are given by

(5.1) 
$$D_{\mathcal{T}}^{-1}B_{\mathcal{T}}(D_{\mathcal{T}}')^{-1}A_{\mathcal{T}} \in \mathcal{L}is(\mathscr{S}_{\mathcal{T},-s},\mathscr{S}_{\mathcal{T},-s}), \\ (D_{\mathcal{T}}')^{-1}A_{\mathcal{T}}D_{\mathcal{T}}^{-1}B_{\mathcal{T}} \in \mathcal{L}is(\mathscr{S}_{\mathcal{T},s},\mathscr{S}_{\mathcal{T},s}),$$

see the following diagram:

$$\begin{array}{ccc} \mathscr{S}_{\mathcal{T},-s} & \xrightarrow{A_{\mathcal{T}}} (\mathscr{S}_{\mathcal{T},-s})' \\ D_{\mathcal{T}}^{-1} & & \downarrow (D_{\mathcal{T}}')^{-1} \\ (\mathscr{S}_{\mathcal{T},s})' & \xleftarrow{B_{\mathcal{T}}} \mathscr{S}_{\mathcal{T},s} \end{array}$$

In the following we shall be concerned with constructing a suitable family  $D_{\mathcal{T}}$ .

#### An obvious but unsatisfactory choice for $D_T$

An option would be to consider  $(D_{\mathcal{T}}u)(v) := \langle u, v \rangle_{L_2(\Omega)} \ (u, v \in \mathscr{S}_{\mathcal{T}})$ , being uniformly in  $\mathcal{L}(\mathscr{S}_{\mathcal{T},-s}, (\mathscr{S}_{\mathcal{T},s})')$ . For showing boundedness of its inverse, let  $Q_{\mathcal{T}}$  be the  $L_2(\Omega)$ -orthogonal projector onto  $\mathscr{S}_{\mathcal{T}}$  then

$$\begin{split} \|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},s})',\mathscr{S}_{\mathcal{T},-s})}^{-1} &= \inf_{0 \neq u \in \mathscr{S}_{\mathcal{T},-s}} \sup_{0 \neq v \in H^{s}(\Omega)} \frac{\langle u, v \rangle_{L_{2}(\Omega)}}{\|u\|_{H^{-s}(\Omega)} \|Q_{\mathcal{T}}v\|_{H^{s}(\Omega)}} \\ &\geq \|Q_{\mathcal{T}}\|_{\mathcal{L}(H^{s}(\Omega),H^{s}(\Omega))}^{-1}, \end{split}$$

As follows from Proposition 2.2.3, the converse is also true, i.e., uniform boundedness of  $\|D_{\mathcal{T}}^{-1}\|_{\mathcal{L}((\mathscr{S}_{\mathcal{T},s})',\mathscr{S}_{\mathcal{T},-s})}$  is actually *equivalent* to uniform boundedness of  $\|Q_{\mathcal{T}}\|_{\mathcal{L}(H^s(\Omega),H^s(\Omega))}$ .

This uniform boundedness of  $||Q_{\mathcal{T}}||_{\mathcal{L}(H^s(\Omega), H^s(\Omega))}$  is well-known for families of *quasi-uniform*, uniformly shape regular conforming partitions of  $\Omega$  into say *d*-simplices. It has also been demonstrated for families of locally refined partitions, for d = 2 including those that are generated by the newest vertex bisection (NVB) algorithm, see [Car02, GHS16, DST20]. On the other hand, in [BY14] a one-dimensional counterexample was presented in which the  $L_2(\Omega)$ -orthogonal projector on a family of sufficiently strongly graded, although uniform *K* meshes, is not  $H^1(\Omega)$ -stable. Thus, in any case uniform  $H^1(\Omega)$ -stability cannot hold without assuming some sufficiently mild grading of the meshes.

Aside from this latter theoretical shortcoming, more importantly, there is a computational problem with the current choice of  $D_{\mathcal{T}}$ . The matrix representation of  $D_{\mathcal{T}}$  w.r.t.  $\Phi_{\mathcal{T}}$  is the 'mass matrix'  $D_{\mathcal{T}} := \langle \Phi_{\mathcal{T}}, \Phi_{\mathcal{T}} \rangle_{L_2(\Omega)}$ . Its inverse  $D_{\mathcal{T}}^{-1}$ , appearing in the preconditioner, is densely populated, and therefore has to be approximated, where generally the error in such approximations has to decrease with a decreasing (minimal) mesh-size in order to arrive at a uniform preconditioner.

## 5.2.2 Constructing a practical $D_T$

To avoid the preceding problems, we shall construct  $D_{\mathcal{T}} \in \mathcal{L}is(\mathscr{S}_{\mathcal{T},-s},(\mathscr{S}_{\mathcal{T},s})')$  with a diagonal matrix representation. To this end, we require some auxiliary

space  $\widetilde{\mathscr{P}}_{\mathcal{T}} \subset H^1(\Omega)$  equipped with a local basis  $\widetilde{\Phi}_{\mathcal{T}}$  that is  $L_2(\Omega)$ -biorthogonal to  $\Phi_{\mathcal{T}}$  and that has 'approximation properties'. To be precise, let  $\widetilde{\Phi}_{\mathcal{T}} := \{\widetilde{\phi}_{\nu} \in H^1(\Omega) : \nu \in N_{\mathcal{T}}\}$  be some collection that satisfies:

(5.2) 
$$\begin{split} &\langle \widetilde{\phi}_{\nu}, \phi_{\nu'} \rangle_{L_{2}(\Omega)} = \delta_{\nu\nu'} \langle \mathbb{1}, \phi_{\nu} \rangle_{L_{2}(\Omega)}, \quad \sum_{\nu \in N_{\mathcal{T}}} \widetilde{\phi}_{\nu} = \mathbb{1}_{\Omega}, \\ & \| \widetilde{\phi}_{\nu} \|_{H^{k}(\Omega)} \lesssim \| \phi_{\nu} \|_{H^{k}(\Omega)} \quad \left( k \in \{0, 1\} \right), \quad \operatorname{supp} \widetilde{\phi}_{\nu} \subseteq \operatorname{supp} \phi_{\nu}.^{1} \end{split}$$

We will take  $D_{\mathcal{T}} := I'_{\mathcal{T}} \widetilde{D}_{\mathcal{T}}$  with  $\widetilde{D}_{\mathcal{T}}$  and  $I_{\mathcal{T}}$  being defined and analyzed in the next two theorems.

**Theorem 5.2.2.** The operator  $\widetilde{D}_{\mathcal{T}}: \mathscr{S}_{\mathcal{T},-s} \to (\widetilde{\mathscr{S}}_{\mathcal{T},s})'$ , defined by  $(\widetilde{D}_{\mathcal{T}}u)(v) := \langle u, v \rangle_{L_2(\Omega)}$ , satisfies  $\widetilde{D}_{\mathcal{T}} \in \mathcal{L}$ is $(\mathscr{S}_{\mathcal{T},-s}, (\widetilde{\mathscr{S}}_{\mathcal{T},s})')$  uniformly in  $\mathcal{T} \in \mathbb{T}$ .

*Proof.* This proof largely follows Sect. 2.3, but because here we consider a Sobolev space  $H^{s}(\Omega)$  that does not incorporate homogeneous boundary conditions, it allows for an easier proof.

From the assumptions (5.2), it follows that the biorthogonal 'Fortin' projector  $P_{\mathcal{T}}: L_2(\Omega) \to H^1(\Omega)$  onto  $\widetilde{\mathscr{S}_{\mathcal{T}}}$  with  $\operatorname{ran}(\operatorname{Id} - P_{\mathcal{T}}) = \mathscr{S}_{\mathcal{T}}^{\perp_{L_2(\Omega)}}$  exists, and is given by

$$P_{\mathcal{T}}u = \sum_{\nu \in N_{\mathcal{T}}} \frac{\langle u, \phi_{\nu} \rangle_{L_{2}(\Omega)}}{\langle \widetilde{\phi}_{\nu}, \phi_{\nu} \rangle_{L_{2}(\Omega)}} \widetilde{\phi}_{\nu}.$$

Let  $T \in \mathcal{T}$ , by (5.2) and the fact that  $\langle \mathbb{1}, \phi_{\nu} \rangle_{L_2(\Omega)} \approx \|\phi_{\nu}\|_{L_2(\Omega)}^2$ , we find for  $k \in \{0, 1\}$ 

(5.3) 
$$\|P_{\mathcal{T}}u\|_{H^{k}(T)} \lesssim \sum_{\nu \in N_{T}} \frac{\|\widetilde{\phi}_{\nu}\|_{H^{k}(T)}}{\|\phi_{\nu}\|_{L_{2}(\Omega)}} \|u\|_{L_{2}(\operatorname{supp}\phi_{\nu})} \lesssim h_{T}^{-k} \|u\|_{L_{2}(\omega_{\mathcal{T}}(T))}$$

with  $\omega_{\mathcal{T}}(T) := \bigcup_{\{\nu \in N_T\}} \operatorname{supp} \phi_{\nu}$ . This shows  $\sup_{\mathcal{T} \in \mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}(L_2(\Omega), L_2(\Omega))} < \infty$ . From the above inequality, and  $\sum_{\nu \in N_{\mathcal{T}}} \widetilde{\phi}_{\nu} = \mathbb{1}$ , we deduce that

$$\begin{split} \| (\mathrm{Id} - P_{\mathcal{T}}) u \|_{H^{1}(T)} &= \inf_{p \in \mathcal{P}_{0}} \| (\mathrm{Id} - P_{\mathcal{T}}) (u - p) \|_{H^{1}(T)} \\ &\lesssim \inf_{p \in \mathcal{P}_{0}} \| u - p \|_{H^{1}(T)} + h_{T}^{-1} \| u - p \|_{L_{2}(\omega_{\mathcal{T}}(T))} \\ &\lesssim \inf_{p \in \mathcal{P}_{0}} h_{T}^{-1} \| u - p \|_{L_{2}(\omega_{\mathcal{T}}(T))} + |u|_{H^{1}(T)} \\ &\lesssim |u|_{H^{1}(\omega_{\mathcal{T}}(T))}, \end{split}$$

with the last step following from the Bramble-Hilbert lemma. We conclude that  $\sup_{\mathcal{T} \in \mathbb{T}} \|P_{\mathcal{T}}\|_{\mathcal{L}(H^1(\Omega), H^1(\Omega))} < \infty$ , and consequently by the Riesz-Thorin

<sup>&</sup>lt;sup>1</sup>This last condition can be replaced by  $\tilde{\phi}_{\nu}$  having (uniformly) local support.

interpolation theorem, that

$$\sup_{\mathcal{T}\in\mathbb{T}}\|P_{\mathcal{T}}\|_{\mathcal{L}(H^s(\Omega),H^s(\Omega))}<\infty.$$

This latter property guarantees that  $\widetilde{D}_{\mathcal{T}}$  is uniformly boundedly invertible:

. .

$$\begin{split} \|\widetilde{D}_{\mathcal{T}}\|_{\mathcal{L}(\mathscr{S}_{\mathcal{T},-s},(\widetilde{\mathscr{S}}_{\mathcal{T},s})')} &= \sup_{0 \neq u \in \mathscr{S}_{\mathcal{T},-s}} \sup_{0 \neq v \in \widetilde{\mathscr{S}}_{\mathcal{T},s}} \frac{\langle u, v \rangle_{L_{2}(\Omega)}}{\|u\|_{H^{-s}(\Omega)} \|v\|_{H^{s}(\Omega)}} \leq 1, \\ \|\widetilde{D}_{\mathcal{T}}^{-1}\|_{\mathcal{L}((\widetilde{\mathscr{F}}_{\mathcal{T},s})',\mathscr{S}_{\mathcal{T},-s})}^{-1} &= \inf_{0 \neq u \in \mathscr{S}_{\mathcal{T},-s}} \sup_{0 \neq v \in \widetilde{\mathscr{S}}_{\mathcal{T},s}} \frac{\langle u, v \rangle_{L_{2}(\Omega)}}{\|u\|_{H^{-s}(\Omega)} \|v\|_{H^{s}(\Omega)}} \\ &= \inf_{0 \neq u \in \mathscr{S}_{\mathcal{T},-s}} \sup_{0 \neq v \in H^{s}(\Omega)} \frac{\langle u, v \rangle_{L_{2}(\Omega)}}{\|u\|_{H^{-s}(\Omega)} \|P_{\mathcal{T}}v\|_{H^{s}(\Omega)}} \\ &\geq \|P_{\mathcal{T}}\|_{\mathcal{L}(H^{s}(\Omega),H^{s}(\Omega))}^{-1}. \end{split}$$

**Theorem 5.2.3.** For  $I_{\mathcal{T}}: \mathscr{S}_{\mathcal{T},s} \to \widetilde{\mathscr{S}}_{\mathcal{T},s}$  being the bijection given by  $I_{\mathcal{T}}\phi_{\nu} = \widetilde{\phi}_{\nu}$  $(\nu \in N_{\mathcal{T}})$ , it holds that  $I_{\mathcal{T}} \in \mathcal{L}$ is $(\mathscr{S}_{\mathcal{T},s}, \widetilde{\mathscr{S}}_{\mathcal{T},s})$  uniformly in  $\mathcal{T} \in \mathbb{T}$ .

Proof. Note that we may write

$$I_{\mathcal{T}}u = \sum_{\nu \in N_{\mathcal{T}}} \frac{\langle u, \phi_{\nu} \rangle_{L_{2}(\Omega)}}{\langle \phi_{\nu}, \widetilde{\phi}_{\nu} \rangle_{L_{2}(\Omega)}} \widetilde{\phi}_{\nu} \quad \text{and} \quad I_{\mathcal{T}}^{-1}u = \sum_{\nu \in N_{\mathcal{T}}} \frac{\langle u, \phi_{\nu} \rangle_{L_{2}(\Omega)}}{\langle \widetilde{\phi}_{\nu}, \phi_{\nu} \rangle_{L_{2}(\Omega)}} \phi_{\nu}.$$

Equivalently to (5.3), we see for  $k \in \{0, 1\}$  that

$$\|I_{\mathcal{T}}u\|_{H^{k}(T)} \lesssim \sum_{\nu \in N_{T}} \frac{\|\phi_{\nu}\|_{H^{k}(T)} \|\phi_{\nu}\|_{L_{2}(\Omega)}}{\|\phi_{\nu}\|_{L_{2}(\Omega)}^{2}} \|u\|_{L_{2}(\operatorname{supp}\phi_{\nu})} \lesssim h_{T}^{-k} \|u\|_{L_{2}(\omega_{\mathcal{T}}(T))}.$$

Following the same arguments as in the proof of Theorem 5.2.2, using that  $I_T \mathbb{1} = \mathbb{1}$ , then reveals that  $I_T$  is uniformly bounded. Uniformly boundedness of  $I_T^{-1}$  follows similarly.

As announced earlier, we define  $D_{\mathcal{T}} \in \mathcal{L}(\mathscr{S}_{\mathcal{T},-s}, (\mathscr{S}_{\mathcal{T},s})')$  by  $D_{\mathcal{T}} := I'_{\mathcal{T}} \widetilde{D}_{\mathcal{T}}$ , so  $(D_{\mathcal{T}}u)(v) := \langle u, I_{\mathcal{T}}v \rangle_{L_2(\Omega)} (u, v \in \mathscr{S}_{\mathcal{T}})$ . Combining the previous theorems gives the following corollary.

**Corollary 5.2.4.** The operator  $D_{\mathcal{T}}$  is in  $\mathcal{L}$ is $(\mathscr{S}_{\mathcal{T},-s}, (\mathscr{S}_{\mathcal{T},s})')$  uniformly in  $\mathcal{T} \in \mathbb{T}$ .

*Remark* 5.2.5. The matrix representation of  $D_{\mathcal{T}}$  w.r.t.  $\Phi_{\mathcal{T}}$  given by

$$\boldsymbol{D}_{\mathcal{T}} = \langle \Phi_{\mathcal{T}}, I_{\mathcal{T}} \Phi_{\mathcal{T}} \rangle_{L_2(\Omega)} = \operatorname{diag}\{ \langle \mathbb{1}, \phi_{\nu} \rangle_{L_2(\Omega)} : \nu \in N_{\mathcal{T}} \}\}$$

which is *diagonal* and therefore easily invertible. The matrix  $D_T$  is known as the *lumped mass matrix*.

*Remark* 5.2.6. The operator  $D_{\mathcal{T}}$  depends merely on the *existence* of a biorthogonal basis  $\tilde{\Phi}_{\mathcal{T}}$  that satisfies (5.2). Indeed, this basis does not appear in the implementation of  $D_{\mathcal{T}}$ .

A possible construction of  $\widetilde{\Phi}_{\mathcal{T}}$  can be given using techniques from Chapter 2. Consider some collection of local 'bubble' functions  $\Theta_{\mathcal{T}} = \{\theta_{\nu} \in H^1(\Omega) : \nu \in N_{\mathcal{T}}\}$  that satisfy:  $|\langle \theta_{\nu}, \phi_{\nu'} \rangle_{L_2(\Omega)}| \approx \delta_{\nu\nu'} \|\phi_{\nu}\|_{L_2(\Omega)}^2$ ,  $\|\theta_{\nu}\|_{H^k(\Omega)} \lesssim \|\phi_{\nu}\|_{H^k(\Omega)}$  $(k \in \{0, 1\})$ , and  $\operatorname{supp} \theta_{\nu} \subseteq \operatorname{supp} \phi_{\nu}$ . Existence of such a collection can be shown by a construction on a reference *d*-simplex, and then using an affine bijection to transfer it to general elements, see Sect. 2.4. A suitable  $\widetilde{\Phi}_{\mathcal{T}}$  that satisfies (5.2) is then given by

$$\widetilde{\phi}_{\nu} := \phi_{\nu} + \frac{\langle \mathbb{1}, \phi_{\nu} \rangle_{L_2(\Omega)}}{\langle \theta_{\nu}, \phi_{\nu} \rangle_{L_2(\Omega)}} \theta_{\nu} - \sum_{\nu' \in N_{\tau}} \frac{\langle \phi_{\nu}, \phi_{\nu'} \rangle_{L_2(\Omega)}}{\langle \theta_{\nu'}, \phi_{\nu'} \rangle_{L_2(\Omega)}} \theta_{\nu'}.$$

We emphasize that the construction of a uniform preconditioner outlined in this subsection does not assume some sufficiently mild grading of the meshes.

#### Implementation

Taking  $\Phi_{\mathcal{T}}$  as basis for both  $\mathscr{S}_{\mathcal{T},-s}$  and  $\mathscr{S}_{\mathcal{T},s}$ , the *matrix representation* of the preconditioned systems from (5.1) read as

$$oldsymbol{D}_{\mathcal{T}}^{-1}oldsymbol{B}_{\mathcal{T}}oldsymbol{D}_{\mathcal{T}}^{- op}oldsymbol{A}_{\mathcal{T}}$$
 and  $oldsymbol{D}_{\mathcal{T}}^{- op}oldsymbol{A}_{\mathcal{T}}oldsymbol{D}_{\mathcal{T}}^{-1}oldsymbol{B}_{\mathcal{T}},$ 

where

$$\boldsymbol{A}_{\mathcal{T}} := (A_{\mathcal{T}} \Phi_{\mathcal{T}})(\Phi_{\mathcal{T}}), \quad \boldsymbol{B}_{\mathcal{T}} := (B_{\mathcal{T}} \Phi_{\mathcal{T}})(\Phi_{\mathcal{T}}), \\ \boldsymbol{D}_{\mathcal{T}} = \boldsymbol{D}_{\mathcal{T}}^{\top} := (D_{\mathcal{T}} \Phi_{\mathcal{T}})(\Phi_{\mathcal{T}}) = \operatorname{diag}\{\langle \mathbb{1}, \phi_{\nu} \rangle_{L_{2}(\Omega)} : \nu \in N_{\mathcal{T}}\}.$$

Alternatively, we could equip the spaces with the *scaled* nodal basis  $\check{\Phi}_{\mathcal{T}} := \mathbf{D}_{\mathcal{T}}^{-\frac{1}{2}} \Phi_{\mathcal{T}}$ , so that the  $L_2(\Omega)$ -norm of any basis function is proportional to 1, yielding

$$\begin{split} \breve{\boldsymbol{A}}_{\mathcal{T}} &:= (A_{\mathcal{T}} \breve{\Phi}_{\mathcal{T}}) (\breve{\Phi}_{\mathcal{T}}) = (\boldsymbol{D}_{\mathcal{T}}^{-\frac{1}{2}})^{\top} \boldsymbol{A}_{\mathcal{T}} \boldsymbol{D}_{\mathcal{T}}^{-\frac{1}{2}}, \\ \breve{\boldsymbol{B}}_{\mathcal{T}} &:= (B_{\mathcal{T}} \breve{\Phi}_{\mathcal{T}}) (\breve{\Phi}_{\mathcal{T}}) = (\boldsymbol{D}_{\mathcal{T}}^{-\frac{1}{2}})^{\top} \boldsymbol{B}_{\mathcal{T}} \boldsymbol{D}_{\mathcal{T}}^{-\frac{1}{2}}, \\ \breve{\boldsymbol{D}}_{\mathcal{T}} &:= (D_{\mathcal{T}} \breve{\Phi}_{\mathcal{T}}) (\breve{\Phi}_{\mathcal{T}}) = (\boldsymbol{D}_{\mathcal{T}}^{-\frac{1}{2}})^{\top} \boldsymbol{D}_{\mathcal{T}} \boldsymbol{D}_{\mathcal{T}}^{-\frac{1}{2}} = \mathbf{Id} \end{split}$$

showing that  $\breve{B}_{\mathcal{T}}$  is a uniform preconditioner for  $\breve{A}_{\mathcal{T}}$  (and vice versa). To the best of our knowledge, so far this most easy form of operator preconditioning, where the stiffness matrix of some operator w.r.t. some basis is preconditioned by stiffness matrix of an opposite order operator w.r.t. the same basis, has not been shown to be optimal.

# 5.3 Manifold case

Let  $\Gamma$  be a compact *d*-dimensional Lipschitz, piecewise smooth manifold in  $\mathbb{R}^{d'}$  for some  $d' \ge d$  without boundary  $\partial \Gamma$ . For  $s \in [0, 1]$ , we consider the Sobolev spaces

$$H^s(\Gamma) := [L_2(\Gamma), H^1(\Gamma)]_{s,2}, \quad H^{-s}(\Gamma) := H^s(\Gamma)'.$$

We assume that  $\Gamma$  is given as the closure of the disjoint union of  $\bigcup_{i=1}^{p} \chi_i(\Omega_i)$ , with, for  $1 \leq i \leq p$ ,  $\chi_i \colon \mathbb{R}^d \to \mathbb{R}^{d'}$  being some smooth regular parametrization, and  $\Omega_i \subset \mathbb{R}^d$  an open polytope. W.l.o.g. assuming that for  $i \neq j$ ,  $\overline{\Omega}_i \cap \overline{\Omega}_j = \emptyset$ , we define

$$\chi \colon \Omega := \cup_{i=1}^p \Omega_i \to \cup_{i=1}^p \chi_i(\Omega_i) \text{ by } \chi|_{\Omega_i} = \chi_i.$$

Let  $\mathbb{T}$  be a family of conforming partitions  $\mathcal{T}$  of  $\Gamma$  into 'panels' such that, for  $1 \leq i \leq p, \chi^{-1}(\mathcal{T}) \cap \Omega_i$  is a uniformly shape regular conforming partition of  $\Omega_i$  into *d*-simplices (that for d = 1 satisfies a uniform *K*-mesh property).

Fix  $\ell > 0$ , we set

$$\mathscr{S}_{\mathcal{T}} := \{ u \in H^1(\Gamma) \colon u \circ \chi|_{\chi^{-1}(T)} \in \mathcal{P}_{\ell} \ (T \in \mathcal{T}) \},\$$

equipped with the canonical nodal basis  $\Phi_{\mathcal{T}} = \{\phi_{\nu} : \nu \in N_{\mathcal{T}}\}.$ 

For construction of an operator  $D_{\mathcal{T}} \in \mathcal{L}$ is  $(\mathscr{S}_{\mathcal{T},-s}, (\mathscr{S}_{\mathcal{T},s})')$  one can proceed as in the domain case. A suitable collection  $\Phi_{\mathcal{T}}$  that is  $L_2(\Gamma)$ -biorthogonal to  $\Phi_{\mathcal{T}}$  exists. Moreover, the analysis from the domain case applies verbatim by only changing  $\langle \cdot, \cdot \rangle_{L_2(\Omega)}$  into  $\langle \cdot, \cdot, \rangle_{L_2(\Gamma)}$ . A *hidden problem*, however, is that the computation of  $D_{\mathcal{T}} = \text{diag}\{\langle \mathbb{1}, \phi_{\nu} \rangle_{L_2(\Gamma)} : \nu \in N_{\mathcal{T}}\}$  involves integrals over  $\Gamma$ that generally have to be approximated using numerical quadrature.

In Sect. 2.3.2 we solved this issue by defining an additional 'mesh-dependent' scalar product

$$\langle u, v \rangle_{\mathcal{T}} := \sum_{T \in \mathcal{T}} \frac{|T|}{|\chi^{-1}(T)|} \int_{\chi^{-1}(T)} u(\chi(x)) v(\chi(x)) dx.$$

This is constructed by replacing on each  $\chi^{-1}(T)$ , the Jacobian  $|\partial \chi|$  by its average  $\frac{|T|}{|\chi^{-1}(T)|}$  over  $\chi^{-1}(T)$ .

By considering  $\Phi_{\mathcal{T}}$  that is biorthogonal to  $\Phi_{\mathcal{T}}$  with respect to  $\langle \cdot, \cdot \rangle_{\mathcal{T}}$ , and the linear bijection  $I_{\mathcal{T}}$  given by  $I_{\mathcal{T}}\phi_{\nu} = \widetilde{\phi}_{\nu}$ , one is able to show that the operator  $D_{\mathcal{T}}$  defined as  $(D_{\mathcal{T}}u)(v) := \langle u, I_{\mathcal{T}}v \rangle_{\mathcal{T}}$  satisfies the necessary requirements. For details we refer to Chapter 2. The resulting matrix representation of  $D_{\mathcal{T}}$  w.r.t.  $\Phi_{\mathcal{T}}$  is then given by  $D_{\mathcal{T}} = \text{diag}\{\langle \mathbb{1}, \phi_{\nu} \rangle_{\mathcal{T}} : \nu \in N_{\mathcal{T}}\}.$ 

# 5.4 Numerical experiments

Let  $\Gamma = \partial [0,1]^3 \subset \mathbb{R}^3$  be the two-dimensional manifold without boundary given as the boundary of the unit cube,  $s = \frac{1}{2}$ , and  $\mathscr{S}_{\mathcal{T}}$  the space of continuous piecewise polynomials of degree  $\ell$  w.r.t. a partition  $\mathcal{T}$ . We will

evaluate preconditioning of the discretized Single Layer Integral operator  $A_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},-s}, (\mathscr{S}_{\mathcal{T},-s})')$  and an (essentially) discretized Hypersingular Integral operator  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},s}, (\mathscr{S}_{\mathcal{T},s})')$ .

The Hypersingular Integral operator  $\tilde{B} \in \mathcal{L}(H^{\frac{1}{2}}(\Gamma), H^{-\frac{1}{2}}(\Gamma))$ , is only-semi coercive, but solving  $\tilde{B}u = f$  for f with  $f(\mathbb{1}) = 0$  is equivalent to solving Bu = f with B given by  $(Bu)(v) = (\tilde{B}u)(v) + \alpha \langle u, \mathbb{1} \rangle_{L_2(\Gamma)} \langle v, \mathbb{1} \rangle_{L_2(\Gamma)}$ , for some fixed  $\alpha > 0$ . This operator B is in  $\mathcal{L}is_c(H^{\frac{1}{2}}(\Gamma), H^{-\frac{1}{2}}(\Gamma))$ , and we shall consider discretizations  $B_{\mathcal{T}} \in \mathcal{L}is_c(\mathscr{S}_{\mathcal{T},s}, (\mathscr{S}_{\mathcal{T},s})')$  of B. We found  $\alpha = 0.05$  to give good results in our examples.

Equipping both  $\mathscr{S}_{\mathcal{T},s}$  and  $\mathscr{S}_{\mathcal{T},-s}$  with the standard nodal basis  $\Phi_{\mathcal{T}} = \{\phi_{\nu} : \nu \in N_{\mathcal{T}}\}$ , the matrix representations of the preconditioned systems from Sect. 5.2.2 read as

$$D_{\mathcal{T}}^{-1}B_{\mathcal{T}}D_{\mathcal{T}}^{-\top}A_{\mathcal{T}} \text{ and } D_{\mathcal{T}}^{-\top}A_{\mathcal{T}}D_{\mathcal{T}}^{-1}B_{\mathcal{T}},$$

for  $D_{\mathcal{T}} = \operatorname{diag}\{\langle \mathbb{1}, \phi_{\nu} \rangle_{L_{2}(\Gamma)} : \nu \in N_{\mathcal{T}}\}, A_{\mathcal{T}} = (A_{\mathcal{T}} \Phi_{\mathcal{T}})(\Phi_{\mathcal{T}}) \text{ and } B_{\mathcal{T}} := (B_{\mathcal{T}} \Phi_{\mathcal{T}})(\Phi_{\mathcal{T}}).$ 

We calculated (spectral) condition numbers of these preconditioned systems, where this condition number is given by  $\kappa_S(\mathbf{X}) := \rho(\mathbf{X})\rho(\mathbf{X}^{-1})$  with  $\rho(\cdot)$  denoting the spectral radius. Note that the condition numbers of the preconditioned systems coincide, i.e.,

$$\kappa_S(\boldsymbol{D}_{\mathcal{T}}^{-1}\boldsymbol{B}_{\mathcal{T}}\boldsymbol{D}_{\mathcal{T}}^{-\top}\boldsymbol{A}_{\mathcal{T}}) = \kappa_S(\boldsymbol{D}_{\mathcal{T}}^{-\top}\boldsymbol{A}_{\mathcal{T}}\boldsymbol{D}_{\mathcal{T}}^{-1}\boldsymbol{B}_{\mathcal{T}}),$$

so we may restrict ourselves to results for preconditioning of  $A_{\mathcal{T}}$ .

We used the BEM++ software package [SBA<sup>+</sup>15] to approximate the matrix representation of  $A_{\mathcal{T}}$  and  $B_{\mathcal{T}}$  by hierarchical matrices based on adaptive cross approximation [Hac99, Beb00].

As initial partition  $\mathcal{T}_{\perp} = \mathcal{T}_1$  of  $\Gamma$  we take a conforming partition consisting of 2 triangles per side, so 12 triangles in total, with an assignment of the newest vertices that satisfies the so-called matching condition. We let  $\mathbb{T}$  be the sequence  $\{\mathcal{T}_k\}_{k\geq 1}$  where the (conforming) partition  $\mathcal{T}_k$  is found by applying both uniform and local refinements. To be precise,  $\mathcal{T}_k$  is constructed by first applying k uniform bisections to  $\mathcal{T}_{\perp}$ , and then 4k local refinements by repeatedly applying NVB to all triangles that touch a corner of the cube. These partitions share both the difficulties of locally refined partitions (the presence of triangles with strongly different sizes) and that of uniform partitions (the diagonally scaled stiffness matrix has a condition number  $\geq 2^{k|s|}$ ).

# 5.4.1 Comparison preconditioners

Write  $G_{\mathcal{T}}^{D} := D_{\mathcal{T}}^{-1} B_{\mathcal{T}} D_{\mathcal{T}}^{-\top}$  for the preconditioner constructed in Sect. 5.2.2. We will compare this with the preconditioner described in Sect. 5.2.1, for which the matrix representation is given by  $G_{\mathcal{T}}^{M} := M_{\mathcal{T}}^{-1} B_{\mathcal{T}} M_{\mathcal{T}}^{-\top}$  with mass matrix  $M_{\mathcal{T}} = M_{\mathcal{T}}^{\top} = \langle \Phi_{\mathcal{T}}, \Phi_{\mathcal{T}} \rangle_{L_{2}(\Gamma)}$ . Because our partitions of the two-dimensional

TABLE 5.1. Spectral condition numbers,  $\kappa_S(G^{\circ}_{\mathcal{T}}A_{\mathcal{T}})$  for  $\circ \in \{D, M\}$ , of the preconditioned Single Layer system discretized on  $\{\mathcal{T}_k\}_{k\geq 1}$ , by continuous piecewise linears  $(\ell = 1)$  in the middle columns and discretized by continuous piecewise cubics  $(\ell = 3)$  in the right columns. Here  $G^D_{\mathcal{T}}$  is the preconditioner introduced in Sect. 5.2.2, whereas  $G^M_{\mathcal{T}}$  is the preconditioner described in Sect. 5.2.1 whose application requires an application of  $M^{-1}_{\mathcal{T}}$ , which we implemented using an LU-factorization.

Partition $\mathcal{T}$		Linears ( $\ell = 1$ )			Cubics $(\ell = 3)$		
$h_{min}$	$h_{max}$	dofs	$oldsymbol{G}_{\mathcal{T}}^{D}oldsymbol{A}_{\mathcal{T}}$	$oldsymbol{G}_{\mathcal{T}}^{M}oldsymbol{A}_{\mathcal{T}}$	dofs	$oldsymbol{G}_{\mathcal{T}}^{D}oldsymbol{A}_{\mathcal{T}}$	$oldsymbol{G}_{\mathcal{T}}^{M}oldsymbol{A}_{\mathcal{T}}$
$1.4\cdot 10^0$	$1.4\cdot 10^0$	8	16.2	1.20	56	90.5	1.68
$4.4\cdot10^{-2}$	$5.0\cdot10^{-01}$	218	14.9	1.91	1946	87.9	2.08
$1.3\cdot10^{-3}$	$3.5\cdot10^{-01}$	482	14.7	2.04	4322	86.1	2.17
$4.3\cdot10^{-5}$	$1.7\cdot10^{-01}$	962	14.7	2.10	8642	85.0	2.21
$1.3 \cdot 10^{-6}$	$8.8 \cdot 10^{-02}$	2306	15.4	2.14	20738	84.9	2.23
$4.2 \cdot 10^{-8}$	$4.4 \cdot 10^{-02}$	7106	15.6	2.16	63938	84.9	2.24
$1.3 \cdot 10^{-9}$	$2.2 \cdot 10^{-02}$	25730	15.8	2.17	231554	84.8	2.25
$4.1 \cdot 10^{-11}$	$1.1 \cdot 10^{-02}$	99650	15.8	2.17	896834	84.7	2.25

surface are created with NVB, we know that also the latter preconditioner provides uniformly bounded condition numbers. In contrast to  $D_{\mathcal{T}}^{-1}$ , the inverse  $M_{\mathcal{T}}^{-1}$  cannot be evaluated in linear complexity. We implemented the application of  $M_{\mathcal{T}}^{-1}$  by computing an LU-factorization of  $M_{\mathcal{T}}$ .

Table 5.1 compares the spectral condition numbers for the preconditioned Single Layer systems with trial spaces given by continuous piecewise linears and those by continuous piecewise cubics. The condition numbers  $\kappa_S(G_T^D A_T)$ are uniformly bounded, but quantitatively the condition numbers  $\kappa_S(G_T^M A_T)$ are better.

# 5.4.2 Improving the preconditioner quality

As observed in Table 5.1, the preconditioner  $G_{\mathcal{T}}^M$  appears to be of superior quality, but it has unfavourable computational complexity. It does suggest a way for improving  $G_{\mathcal{T}}^D$ : by replacing  $D_{\mathcal{T}}^{-1}$  with a better approximation of  $M_{\mathcal{T}}^{-1}$ , one may hope to improve the quality. To this end, we introduce damped (preconditioned) Richardson. Let  $0 < \lambda_{-} \leq \lambda_{min}(D_{\mathcal{T}}^{-1}M_{\mathcal{T}}), \lambda_{max}(D_{\mathcal{T}}^{-1}M_{\mathcal{T}}) \leq \lambda_{+},$  $R_{\mathcal{T}}^{(0)} := 0$  and for  $k \geq 0$  define

$$\boldsymbol{R}_{\mathcal{T}}^{(k+1)} \coloneqq \boldsymbol{R}_{\mathcal{T}}^{(k)} + \omega \boldsymbol{D}_{\mathcal{T}}^{-1} (\mathrm{Id} - \boldsymbol{M}_{\mathcal{T}} \boldsymbol{R}_{\mathcal{T}}^{(k)}), \quad \omega = \frac{2}{\lambda_{-} + \lambda_{+}},$$

being the result of k Richardson iterations. Correspondingly define

(5.4) 
$$\boldsymbol{G}_{\mathcal{T}}^{(k)} \coloneqq \boldsymbol{R}_{\mathcal{T}}^{(k)} \boldsymbol{B}_{\mathcal{T}} \boldsymbol{R}_{\mathcal{T}}^{(k)}$$

2.52

2.52

2.52

2.52

2.52

co	etized b ntinuou	y continu 1s piecew	ious piec rise cubio	ewise lii cs in the	nears in f right co	the left col lumns.	umns ar	nd discre	etiz
		Linears	$\ell (\ell = 1)$			Cubics	$(\ell = 3)$		
	dofs	k = 2	k = 4	k = 6	dofs	k = 2	k = 4	k = 6	-
	8	2.26	1.29	1.22	56	10.1	3.99	2.65	-
	218	3.05	2.07	1.94	1946	8.96	3.57	2.52	
	482	3.53	2.28	2.08	4322	8.80	3.59	2.52	

8642

20738

63938

231554

896834

8.63

8.54

8.54

8.54

8.54

3.59

3.59

3.59

3.59

3.59

2.19

2.24

2.27

2.28

2.29

TABLE 5.2. Spectral condition numbers  $\kappa_S(\mathbf{G}_{\mathcal{T}}^{(k)}\mathbf{A}_{\mathcal{T}})$  with  $\mathbf{G}_{\mathcal{T}}^{(k)}$  the preconditioner from (5.4) that incorporates k Richardson iterations. The systems are discretized by continuous piecewise linears in the left columns and discretized by continuous piecewise cubics in the right columns.

It follows that  $G_{\mathcal{T}}^{(1)} = G_{\mathcal{T}}^D$  and  $\lim_{k\to\infty} G_{\mathcal{T}}^{(k)} = G_{\mathcal{T}}^M$ . Although we have no proof, we suspect that  $G_{\mathcal{T}}^{(k)}$  provides a uniform preconditioner for  $A_{\mathcal{T}}$  due to the fact that  $R_{\mathcal{T}}^{(k)}$  approximates  $M_{\mathcal{T}}^{-1}$ , while preserving constant functions, being a key ingredient in the proofs of Theorems 5.2.2 and 5.2.3.

Values for  $\lambda_{-}$  and  $\lambda_{+}$  can be found by calculating the extremal eigenvalues of the corresponding preconditioned mass matrix on a reference simplex, see e.g. [Wat87]. For  $\ell = 1$  this gives  $\omega = \frac{2(d+2)}{d+3}$ , whereas for  $\ell = 3$  and d = 2 we computed  $\omega = 0.836$ .

Table 5.2 compares the condition numbers  $\kappa_S(\mathbf{G}_{\mathcal{T}}^{(k)} \mathbf{A}_{\mathcal{T}})$  for  $k \in \{2, 4, 6\}$ . We see that a few Richardson iterations drastically improves our preconditioner, making its quality on par with that of  $\mathbf{G}_{\mathcal{T}}^M$  while having a favourable linear application cost.

Finally, to show that one cannot simply use any (iterative) method for approximating  $M_{\tau}^{-1}$ , we consider the case where one approximates this inverse using a Jacobi preconditioner. The resulting preconditioner is then given by

(5.5) 
$$\boldsymbol{G}_{\mathcal{T}}^{J} \coloneqq (\operatorname{diag} \boldsymbol{M}_{\mathcal{T}})^{-1} \boldsymbol{B}_{\mathcal{T}} (\operatorname{diag} \boldsymbol{M}_{\mathcal{T}})^{-\top}$$

Table 5.3 clearly displays that this is not a uniformly bounded preconditioner, which we assume is due to the fact that  $(\operatorname{diag} M_{\mathcal{T}})^{-1}$  does not preserve constant functions for  $\ell > 1$ .

# 5.5 Conclusion

962

2306

7106

25730

99650

3.79

3.98

4.18

4.35

4.47

2.44

2.52

2.57

2.61

2.65

Considering discretized opposite order operators  $A_T$  and  $B_T$  using the same ansatz space of continuous piecewise polynomial w.r.t. a possibly locally re-

TABLE 5.3. Spectral of	condition numb	pers $\kappa_S(\boldsymbol{G}_{\mathcal{T}}^J \boldsymbol{A}_{\mathcal{T}})$	with $G_{\mathcal{T}}^{J}$	from	(5.5),	and
systems discretized	by continuous p	piecewise cubics	$(\ell = 3).$			

dofs	$oldsymbol{G}_{\mathcal{T}}^{J}oldsymbol{A}_{\mathcal{T}}$
56	62.6
1946	377.1
4322	495.6
8642	1016.9
20738	3067.8
63938	10928.3

fined partition  $\mathcal{T}$ , we consider matrices  $D_{\mathcal{T}}$  such that  $D_{\mathcal{T}}^{-1}B_{\mathcal{T}}D_{\mathcal{T}}^{-\top}$  is a uniform preconditioner for  $A_{\mathcal{T}}$ , and  $D_{\mathcal{T}}^{-\top}A_{\mathcal{T}}D_{\mathcal{T}}^{-1}$  for  $B_{\mathcal{T}}$ . The obvious choice for  $D_{\mathcal{T}}$ would be the mass matrix, however, it yields uniformly bounded condition numbers only under a mildly grading assumption on the mesh, and more importantly, it has the disadvantage that its inverse is dense. We proved that when taking  $D_{\mathcal{T}}$  as the lumped mass matrix the condition numbers are uniformly bounded, remarkably without a sufficiently mild grading assumption on the mesh, while obviously its inverse can be applied in linear cost.

In our experiments with locally refined meshes generated by Newest Vertex Bisection, the condition numbers with  $D_{\mathcal{T}}$  being the mass matrix are quantitatively better than those found with  $D_{\mathcal{T}}$  being the lumped mass matrix though. Constructing  $D_{\mathcal{T}}^{-1}$  as an approximation for the inverse mass matrix by a few preconditioned damped Richardson steps with the lumped mass matrix as a preconditioner, both the resulting matrix can be applied at linear cost and the observed condition numbers are essentially as good as with the inverse mass matrix.

# Part II

# Space-time methods for parabolic evolution equations
# 6.1 Introduction

This chapter is about the adaptive numerical solution of parabolic evolution equations written in a simultaneous space-time variational formulation. In comparison to the usually applied time-marching schemes, simultaneous space-time solvers offer the following potential advantages:

- local, adaptive refinements simultaneous in space and time ([SY18, RS19, GS19]),
- quasi-best approximation from the selected trial space ('Cea's lemma') ([And13, LM17, SW21b]), being a necessary requirement for proving optimal rates for adaptive routines ([CS11, KSU16, RS19]),
- superior parallel performance ([DGVdZ18, NS19, HLNS19, vVW21a]),
- using the product structure of the space-time cylinder, sparse tensor product approximation ([GO07, CS11, KSU16, RS19]) which allows to solve the whole time evolution at a complexity of solving the corresponding stationary problem.

Other relevant publications on space-time solvers include [Ste15, LMN16, SZ20, Dev20, DS20].

In any case without applying sparse tensor product approximation, a disadvantage of the space-time approach is the larger memory consumption because instead of solving a sequence of PDEs on a *d*-dimensional space, one has to solve one PDE on a (d + 1)-dimensional space. This disadvantage, however, disappears when one needs simultaneously the whole time evolution as for example with problems of optimal control ([GK11, BRU20]) or data-assimilation ([DSW21]).

# 6.1.1 Parabolic problem in a simultaneous space-time variational formulation

For some separable Hilbert spaces  $V \hookrightarrow H$  with dense embedding (e.g.  $H_0^1(\Omega)$ and  $L_2(\Omega)$  for the model problem of the heat equation on a spatial domain  $\Omega \subset \mathbb{R}^d$ ), and a boundedly invertible  $A(t) = A(t)' \colon V \to V'$  with  $(A(t) \cdot)(\cdot) = \|\cdot\|_V^2$  (e.g.  $(A(t)\eta)(\zeta) = \int_{\Omega} \nabla \eta \cdot \nabla \zeta \, d\mathbf{x}$ ), we consider

$$\begin{cases} \frac{du}{dt}(t) + A(t)u(t) = g(t) & (t \in (0,T)), \\ u(0) = u_0. \end{cases}$$

An application of a variational formulation of the PDE over space and time leads to an equation

(6.1) 
$$\begin{bmatrix} B\\ \gamma_0 \end{bmatrix} u = \begin{bmatrix} g\\ u_0 \end{bmatrix}$$

where, with  $X := L_2(I; V) \cap H^1(I; V')$  and  $Y := L_2(I; V)$ , the operator at the left hand side is boundedly invertible  $X \to Y' \times H$ .

### 6.1.2 Our previous work

In [CS11, RS19] we equipped *X* and *Y* with Riesz bases being tensor products of wavelet bases in space in time, and *H* with some spatial Riesz basis. Consequently, the equation (6.1) got an equivalent formulation as a bi-infinite well-posed matrix-vector equation  $\begin{bmatrix} \mathbf{B} \\ \gamma_0 \end{bmatrix} \mathbf{u} = \begin{bmatrix} \mathbf{g} \\ \mathbf{u}_0 \end{bmatrix}$  (actually, in [RS19], we considered a formulation of first order, and in [CS11] we used a variational formulation with essentially interchanged roles of *X* and *Y*, which however is irrelevant for the current discussion). To get a coercive bilinear form we formed normal equations to which we applied an adaptive wavelet scheme ([CDD01]). With such a scheme the norm of a sufficiently accurate approximation of the (infinite) residual vector of a current approximation is used as an a posteriori error estimator. The coefficients in modulus of this vector are applied as local error indictors in a bulk chasing (or Dörfler) marking procedure. The resulting adaptive algorithm converges at the best possible rate in linear computational complexity.

The goal of the current work is to investigate to what extent similar optimal theoretical results can be shown for finite element discretizations, whilst realizing a quantitatively superior implementation.

# 6.1.3 Least squares minimization

Without having Riesz bases for X and Y, already the step of first discretizing and then forming normal equations does not apply, and we reverse their order. A problem equivalent to (6.1) is to compute

(6.2) 
$$u = \underset{w \in X}{\operatorname{argmin}} \|Bw - g\|_{Y'}^2 + \|\gamma_0 w - u_0\|_H^2.$$

An obvious approach for the numerical approximation is to consider the minimization over finite dimensional subspaces  $X^{\delta}$  of X, which however is not feasible because of the presence of the dual norm.

For trial spaces  $X^{\delta}$  that are 'full' (or 'sparse') tensor products of finite element spaces in space and time, in [And13] it was shown how to construct corresponding test spaces  $Y^{\delta} \subset Y$  of similar type and dimension, such that  $(X^{\delta}, Y^{\delta})$  is uniformly *inf-sup stable* meaning that when the continuous dual norm  $\|\cdot\|_{Y'}$  is replaced by the discrete dual norm  $\|\cdot\|_{Y^{\delta'}}$ , a minimization over  $X^{\delta}$  yields a quasi-best approximation to *u* from  $X^{\delta}$ . Such a family of trial spaces however does not allow to create a nested sequence of trial spaces by adaptive *local* refinements.

#### 6.1.4 Family of inf-sup stable pairs of trial and test spaces

To construct an alternative, essentially larger family, let  $\Sigma$  be a wavelet Riesz basis for  $L_2(0, T)$  that, after renormalization, is also a Riesz basis for  $H^1(0, T)$ . We equip this basis with a tree structure where every wavelet that is not on the coarsest level has a parent on the next coarser level. In space, we consider the collection of all linear finite element spaces that can be generated by conforming newest vertex bisection starting from an initial conforming partition of a polytopal  $\Omega$  into *d*-simplices. The restriction to linear finite elements is not essential and is made for simplicity only. Now we consider trial spaces  $X^{\delta}$  that are spanned by a number of wavelets each of them tensorized with a finite element space from the aforementioned collection. In order to be able to apply the arising system matrices in linear complexity, see [KS14] or Chapter 7, we impose the condition that if a wavelet tensorized with a finite element space is in the spanning set, then so is its parent wavelet tensorized with a finite element space that includes the former one.

The infinite collection of finite element spaces can be associated to a hierarchical 'basis' that can be equipped with a tree structure. Each hierarchical basis function, except those on the coarsest level, is associated to a node  $\nu$  that was inserted as the midpoint of an edge connecting two nodes on the next coarser level, which nodes we call the parents of  $\nu$ . With this definition there is a one-to-one correspondence between the finite element spaces from our collection and the spans of the sets of hierarchical basis functions that form trees. Consequently, our collection of trial spaces  $X^{\delta}$  consists of the spans of sets of tensor products of wavelets-in-time and hierarchical basis functions-inspace which sets are *downwards closed*, also known as *lower*, in the sense that if a pair of a wavelet and a hierarchical basis function is in the set, then so are all its parents in time and space. Spaces from this collection can be 'locally' expanded by adding the span of a tensor product of a wavelet and hierarchical basis function one-by-one.

For this family of spaces  $X^{\delta}$  we construct a corresponding family of spaces  $Y^{\delta} \subset Y$  of similar type such that each pair  $(X^{\delta}, Y^{\delta})$  is uniformly inf-sup stable,

with the dimension of  $Y^{\delta}$  being proportional to that of  $X^{\delta}$ . Furthermore, using the properties of the wavelets in time and by applying multigrid preconditioners in space we construct optimal preconditioners at X and Y-side which allow a fast solution of the discrete problems  $\operatorname{argmin}_{w \in X^{\delta}} \|Bw - g\|_{Y^{\delta'}}^2 + \|\gamma_0 w - u_0\|_H^2$ .

# 6.1.5 Adaptive algorithm

Having fixed the family of trial spaces, it remains to develop an algorithm that selects a suitable, preferably quasi-optimal, nested sequence of spaces from the family adapted to the solution u of (6.2). The theory about adaptive (Ritz-) Galerkin approximations for such quadratic minimization problems is in a mature state. As noticed before, however, Galerkin approximations for (6.2) are not computable.

Therefore given  $X^{\delta}$ , let  $X^{\underline{\delta}} \supset X^{\delta}$  be such that saturation holds, i.e., for some constant  $\zeta < 1$ , it holds that  $\inf_{w \in X^{\underline{\delta}}} \|u - w\|_X \leq \zeta \inf_{w \in X^{\delta}} \|u - w\|_X$ . We now replace problem (6.2) by

(6.3) 
$$u^{\underline{\delta}\underline{\delta}} = \operatorname*{argmin}_{w \in X^{\underline{\delta}}} \|Bw - g\|^2_{Y^{\underline{\delta}'}} + \|\gamma_0 w - u_0\|^2_H,$$

where in the notation  $u^{\underline{\delta}\underline{\delta}}$  the first instance of  $\underline{\delta}$  refers to the space  $Y^{\underline{\delta}}$  and the second to the space  $X^{\underline{\delta}}$ . Its (computable) Galerkin approximation from  $X^{\overline{\delta}}$  is given by

$$u^{\delta\delta} = \underset{w \in X^{\delta}}{\operatorname{argmin}} \|Bw - g\|_{Y^{\delta'}}^2 + \|\gamma_0 w - u_0\|_H^2.$$

By a standard adaptive procedure, described below, we expand  $X^{\delta}$  to some  $X^{\tilde{\delta}} \subseteq X^{\underline{\delta}}$  such that  $u^{\underline{\delta}\tilde{\delta}}$  is closer to  $u^{\underline{\delta}\underline{\delta}}$  than  $u^{\underline{\delta}\delta}$ . Next, we replace  $Y^{\underline{\delta}}$  by  $Y^{\underline{\tilde{\delta}}}$  (being the test space corresponding to  $X^{\underline{\tilde{\delta}}}$ ) and repeat (i.e. consider (6.3) with  $(\underline{\delta}, \underline{\delta})$  reading as  $(\underline{\tilde{\delta}}, \underline{\tilde{\delta}})$ , and improve its Galerkin approximation  $u^{\underline{\tilde{\delta}}\tilde{\delta}}$  from  $X^{\underline{\tilde{\delta}}}$  by an adaptive enlargement of the latter space).

The adaptive expansion of the trial space  $X^{\delta}$  to  $X^{\overline{\delta}}$  will be by the application of the usual solve-estimate-mark-refine paradigm, where the error indicators are the coefficients of the residual vector w.r.t. (modified) tensor product basis functions that were added to  $X^{\delta}$  to create  $X^{\underline{\delta}}$ . In order for this collection of additional tensor product basis functions to be stable in *X*-norm, for this step we modify the hierarchical basis functions such that they get a vanishing moment, and therefore become closer to 'real' wavelets.

Under the aforementioned saturation assumption, we prove that the overall adaptive procedure produces an *r*-linearly converging sequence to the solution.

# 6.1.6 Numerical results

We tested the adaptive algorithm in several examples with a two-dimensional spatial domain. In all but one case, we observed a convergence rate equal to

1/2, being the best that can be expected in view of the piecewise polynomial degree of the trial functions *and* the tensor product construction, and for non-smooth solutions improving upon usual non-adaptive approximation. Only for the case where  $u_0 = 1$  and homogenous Dirichlet boundary conditions are prescribed, the observed rate was reduced to 0.4. It is unknown whether or not this is the best non-linear approximation rate for our family of trial spaces.

Thanks to the use of optimal preconditioners and that of a carefully designed matrix-vector multiplication routine, which generalizes such a routine for sparse-grids introduced in [BZ96] to adaptive settings, we observe that throughout the whole execution of the adaptive loop the total runtime remains proportional to the current number of unknowns.

Recently in [FK21], see also [GS21], a first order system least squares (FOSLS) of second order parabolic PDEs was proven to be well-posed. This formulation has the very attractive property that the several components of the residual are all measured in  $L_2$ -norms. So other than with (6.2) there is no need to discretize a dual-norm, and so to guarantee an inf-sup condition. Minimization over any conforming trial space yields a quasi-best approximation from that space in the corresponding 'energy-norm'. This norm, however, is stronger than the norm on X. For the aforementioned example of a discontinuity between initial and boundary conditions, and with the application of continuous piecewise linear finite elements w.r.t. tetrahedral meshes of the space-time cylinder it results in a convergence rate of 0.07 for uniform refinements, which is not visibly improved using adaptive refinements.

#### 6.1.7 Organization

This chapter is organized as follows: In Sect. 6.2 the well-posed space-time variational formulation of the parabolic problem is discussed, and in Sect. 6.3 we discuss its inf-sup stable discretisation. The adaptive solution procedure is presented in Sect. 6.4, and its convergence is proven. The construction of the trial and test spaces is detailed in Sect. 6.5, and optimal preconditioners are presented. In Sect. 6.6, the definition of the enlarged space  $X^{\delta}$  is given, and the construction of a stable basis of a stable complement space of  $X^{\delta}$  in  $X^{\delta}$  is outlined. Numerical results are presented in Sect. 6.7, and a conclusion is formulated in Sect. 6.8.

#### 6.1.8 Notations

In this work, by  $C \leq D$  we will mean that C can be bounded by a multiple of D, independently of parameters which C and D may depend on. Obviously,  $C \geq D$  is defined as  $D \leq C$ , and C = D as  $C \leq D$  and  $C \geq D$ .

For normed linear spaces E and F, by  $\mathcal{L}(E, F)$  we will denote the normed linear space of bounded linear mappings  $E \to F$ , and by  $\mathcal{L}is(E, F)$  its subset of boundedly invertible linear mappings  $E \to F$ . We write  $E \hookrightarrow F$  to denote

that *E* is continuously embedded into *F*. For simplicity only, we exclusively consider linear spaces over the scalar field  $\mathbb{R}$ .

# 6.2 Space-time formulations of a parabolic evolution problem

Let V, H be separable Hilbert spaces of functions on some "spatial domain" such that  $V \hookrightarrow H$  with dense embedding. Identifying H with its dual, we obtain the Gelfand triple  $V \hookrightarrow H \simeq H' \hookrightarrow V'$ .

For a.e.

$$t \in I := (0, T),$$

let  $a(t; \cdot, \cdot)$  denote a bilinear form on  $V \times V$  such that for any  $\eta, \zeta \in V, t \mapsto a(t; \eta, \zeta)$  is measurable on *I*, and such that for some  $\varrho \in \mathbb{R}$ , for a.e.  $t \in I$ ,

(6.4)  $|a(t;\eta,\zeta)| \lesssim \|\eta\|_V \|\zeta\|_V \quad (\eta,\zeta \in V) \quad (boundedness),$ 

(6.5) 
$$a(t;\eta,\eta) + \varrho\langle\eta,\eta\rangle_H \gtrsim \|\eta\|_V^2$$
  $(\eta \in V)$  (Gårding inequality).

With  $A(t) \in \mathcal{L}is(V, V')$  being defined by  $(A(t)\eta)(\zeta) = a(t; \eta, \zeta)$ , given a forcing function g and an initial value  $u_0$ , we are interested in solving the *parabolic initial value problem* to finding u such that

(6.6) 
$$\begin{cases} \frac{du}{dt}(t) + A(t)u(t) = g(t) & (t \in I), \\ u(0) = u_0. \end{cases}$$

In a simultaneous space-time variational formulation, the parabolic PDE reads as finding u from a suitable space of functions of time and space such that

$$(Bw)(v) := \int_{I} \langle \frac{dw}{dt}(t), v(t) \rangle + a(t; w(t), v(t)) dt = \int_{I} \langle g(t), v(t) \rangle =: g(v)$$

for all v from another suitable space of functions of time and space. One possibility to enforce the initial condition is by testing it against additional test functions. A proof of the following result can be found in [SS09], cf. [DL92, Ch.XVIII, §3] and [Wlo82, Ch. IV, §26] for slightly different statements.

**Theorem 6.2.1.** With  $X := L_2(I; V) \cap H^1(I; V')$ ,  $Y := L_2(I; V)$ , under conditions (6.4) and (6.5) it holds that

$$\begin{bmatrix} B\\ \gamma_0 \end{bmatrix} \in \mathcal{L}is(X, Y' \times H),$$

where for  $t \in \overline{I}$ ,  $\gamma_t : u \mapsto u(t, \cdot)$  denotes the trace map. That is, assuming  $g \in Y'$  and  $u_0 \in H$ , finding  $u \in X$  such that

(6.7) 
$$\begin{bmatrix} B\\ \gamma_0 \end{bmatrix} u = \begin{bmatrix} g\\ u_0 \end{bmatrix}$$

is a well-posed simultaneous space-time variational formulation of (6.6).

With  $\tilde{u}(t) := u(t)e^{-\varrho t}$ , (6.6) is equivalent to  $\frac{d\tilde{u}}{dt}(t) + (A(t) + \varrho \operatorname{Id})\tilde{u}(t) = g(t)e^{-\varrho t}$   $(t \in I), \tilde{u}(0) = u_0$ . Since  $((A(t) + \varrho \operatorname{Id})\eta)(\eta) \gtrsim ||\eta||_V^2$ , w.l.o.g. we assume that (6.5) is valid for  $\rho = 0$ , i.e.,  $a(t; \cdot, \cdot)$  is *coercive* uniformly for a.e.  $t \in I$ .

For simplicity, cf. discussion in Remark 6.3.5, additionally we assume that  $a(t; \cdot, \cdot)$  is *symmetric*, and define  $A = A' \in \mathcal{L}is(Y, Y')$  by  $(Aw)(v) = \int_{I} (A(t)w(t))v(t) dt$ .

Because  $\begin{bmatrix} A & 0 \\ 0 & \text{Id} \end{bmatrix} \in \mathcal{L}$ is $(Y \times H, Y' \times H)$ , an equivalent formulation of (6.7) as a self-adjoint saddle point equation reads as finding  $(\mu, \sigma, u) \in Y \times H \times X$  (where  $\mu = 0 = \sigma$ ) such that

(6.8) 
$$\begin{bmatrix} A & 0 & B \\ 0 & \mathrm{Id} & \gamma_0 \\ B' & \gamma'_0 & 0 \end{bmatrix} \begin{bmatrix} \mu \\ \sigma \\ u \end{bmatrix} = \begin{bmatrix} g \\ u_0 \\ 0 \end{bmatrix},$$

or equivalently

(6.9) 
$$\begin{bmatrix} A & B \\ B' & -\gamma'_0 \gamma_0 \end{bmatrix} \begin{bmatrix} \mu \\ u \end{bmatrix} = \begin{bmatrix} g \\ -\gamma'_0 u_0 \end{bmatrix}$$

or

(6.10) 
$$\underbrace{(B'A^{-1}B + \gamma'_0\gamma_0)}_{S:=} u = \underbrace{B'A^{-1}g + \gamma'_0u_0}_{f:=}.$$

We equip *Y* and *X* with 'energy'-norms

 $\|\cdot\|_{Y}^{2} := (A \cdot)(\cdot), \quad \|\cdot\|_{X}^{2} := \|\cdot\|_{Y}^{2} + \|\partial_{t} \cdot\|_{Y'}^{2} + \|\gamma_{T} \cdot\|_{H}^{2},$ 

which are equivalent to the canonical norms on Y and X. Notice that (6.8)–(6.10) are the Euler-Langrange equations that result from the minimization problem

$$u = \underset{w \in X}{\operatorname{argmin}} \|Bw - f\|_{Y'}^2 + \|\gamma_0 w - u_0\|_H^2.$$

**Lemma 6.2.2.** We have  $\|\cdot\|_X^2 = (S \cdot)(\cdot)$ .

*Proof.* It holds that

$$\begin{aligned} \|w\|_X^2 &= \sup_{0 \neq v_1 \in Y} \frac{(Bw)(v_1)}{\|v_1\|_Y^2} + \|\gamma_0 w\|_H^2 = \sup_{0 \neq (v_1, v_2) \in Y \times H} \frac{((Bw)(v_1) + \langle \gamma_0 w, v_2 \rangle_H)^2}{\|v_1\|_Y^2 + \|v_2\|_H^2} \\ &= (Sw)(w), \end{aligned}$$

where the first equality can be found in e.g. [ESV17, Thm. 2.1], and, when realising that *S* is the Schur complement of the operator in (6.8), the last one in e.g. [KS08, Lemma 2.2].

## 6.3 Discretizations

#### 6.3.1 Galerkin discretization of the Schur complement equation

Let  $(X^{\delta})_{\delta \in \Delta}$  be a collection of closed, e.g, finite dimensional, subspaces of *X*, so equipped with  $\|\cdot\|_X$ . We will specify such a family in Sect. 6.5-6.6.1. We define a *partial order* on  $\Delta$  by

$$\delta \preceq \tilde{\delta} \Longleftrightarrow X^{\delta} \subseteq X^{\tilde{\delta}}.$$

For  $\delta \in \Delta$ , let  $u_{\delta} \in X^{\delta}$  denote the *Galerkin approximation to the solution* u of (6.10), i.e., the solution of

(6.11) 
$$(Su_{\delta})(v) = f(v) \quad (v \in X^{\delta}).$$

being the best approximation to u from  $X^{\delta}$  w.r.t.  $\|\cdot\|_X$ .

For proving convergence of an adaptive solution routine, as well as for a posteriori error estimation, we shall make the following assumption.

**Assumption 6.3.1** (Saturation). There exists a collection of subspaces  $({}^{\delta}G \times {}^{\delta}U_0)_{\delta \in \Delta} \subseteq$  $Y' \times H$ , a mapping  $:: \Delta \to \Delta: \delta \mapsto \underline{\delta}$  where  $\underline{\delta} \succeq \delta$ , and some fixed constant  $\zeta < 1$  such that for all  $\delta \in \Delta$ , assuming that  $(g, u_0) \in {}^{\delta}G \times {}^{\delta}U_0$ ,

$$(6.12) \|u - u_{\underline{\delta}}\|_X \le \zeta \|u - u_{\delta}\|_X.$$

*Remark* 6.3.2. Notice that above assumption cannot be valid without a restriction on the right-hand side  $f = B'A^{-1}g + \gamma'_0u_0 \in X'$ . Indeed given any  $X^{\delta} \subset X^{\underline{\delta}} \subsetneq X$ , consider a non-zero  $f \in X'$  that vanishes on  $X^{\underline{\delta}}$ . Then  $u_{\delta} = u_{\delta} = 0 \neq u$ , meaning that (6.12) does not hold.

For the time being we will operate under the restrictive assumption that whenever we apply (6.12) (visible by the appearance of the constant  $\zeta$ ) we simply assume that  $(g, u_0) \in {}^{\delta}G \times {}^{\delta}U_0$ . Later, in Sect. 6.4.3, we will remove this assumption.

The discretized problem from (6.11) only serves theoretical purposes. Indeed, since the Schur complement operator *S* contains the inverse of *A*, there is no way to determine  $u_{\delta}$  exactly. The reason to introduce (6.11) is that *S* is an *elliptic* operator, so that for  $\delta \leq \tilde{\delta}$  we can make use of  $||u - u_{\tilde{\delta}}||_X^2 =$  $||u - u_{\delta}||_X^2 - ||u_{\tilde{\delta}} - u_{\delta}||_X^2$ , being a crucial tool for proving convergence of adaptive algorithms.

# 6.3.2 Uniformly stable Galerkin discretization of the saddle-point formulation

Our numerical approximations will be based on Galerkin discretizations of the saddle-point formulation (6.9). Let  $(Y^{\delta})_{\delta \in \Delta}$  be a collection of closed subspaces

of *Y*, so equipped with  $\|\cdot\|_Y$ , such that

(6.13) 
$$X^{\delta} \subseteq Y^{\delta} \quad (\delta \in \Delta),$$

and

(6.14) 
$$1 \ge \gamma_{\Delta} := \inf_{\delta \in \Delta} \inf_{0 \ne w \in X^{\delta}} \sup_{0 \ne v \in Y^{\delta}} \frac{(\partial_t w)(v)}{\|\partial_t w\|_{Y'} \|v\|_Y} > 0.$$

Notice that  $1 - \gamma_{\Delta}$  can be made arbitrarily small by selecting, for each  $\delta \in \Delta$ ,  $Y^{\delta}$  sufficiently large in relation to  $X^{\delta}$ .

For  $\delta, \hat{\delta} \in \Delta$  with  $Y^{\hat{\delta}} \supseteq Y^{\delta}$ , and  $E_Y^{\hat{\delta}}, E_X^{\delta}$  denoting the embeddings  $Y^{\hat{\delta}} \to Y$ ,  $X^{\delta} \to X$ , let  $(\mu^{\hat{\delta}\delta}, u^{\hat{\delta}\delta}) \in Y^{\hat{\delta}} \times X^{\delta}$  be the solution of

(6.15) 
$$\begin{bmatrix} E_Y^{\hat{\delta}'} A E_Y^{\hat{\delta}} & E_Y^{\hat{\delta}'} B E_X^{\delta} \\ E_X^{\delta'} B' E_Y^{\hat{\delta}} & -E_X^{\delta'} \gamma_0' \gamma_0 E_X^{\delta} \end{bmatrix} \begin{bmatrix} \mu^{\hat{\delta}\delta} \\ u^{\hat{\delta}\delta} \end{bmatrix} = \begin{bmatrix} E_Y^{\hat{\delta}'} g \\ -E_X^{\delta'} \gamma_0' u_0 \end{bmatrix},$$

or, equivalently,

(6.16) 
$$\underbrace{E_X^{\delta'}(B'E_Y^{\hat{\delta}}(E_Y^{\hat{\delta}'}AE_Y^{\hat{\delta}})^{-1}E_Y^{\hat{\delta}'}B + \gamma'_0\gamma_0)E_X^{\delta}}_{S^{\hat{\delta}\delta:=}} = \underbrace{E_X^{\delta'}(B'E_Y^{\hat{\delta}}(E_Y^{\hat{\delta}'}AE_Y^{\hat{\delta}})^{-1}E_Y^{\hat{\delta}'}g + \gamma'_0u_0)}_{f^{\hat{\delta}\delta:=}}.$$

Below we will see that (6.15)-(6.16) are uniquely solvable. Formulated in 'operator language', (6.15) is the Galerkin discretization of (6.9) on the closed subspace  $Y^{\hat{\delta}} \times X^{\delta} \subseteq Y \times X$ . Unless  $Y^{\hat{\delta}} = Y$ , it holds that  $S^{\hat{\delta}\delta} \neq E_X^{\delta'} S E_X^{\delta}$  and  $f^{\hat{\delta}\delta} \neq E_X^{\delta'} f$ , and so generally  $u^{\hat{\delta}\delta} \neq u_{\delta}$ .

As we will see, however, for  $Y^{\delta}$ , and thus  $Y^{\hat{\delta}}$ , 'large' in relation to  $X^{\delta}$ ,  $u^{\hat{\delta}\delta}$  will be 'close' to  $u_{\delta}$ . This will allow us to show that (*r*-linear) convergence of a sequence of Galerkin solutions  $u_{\delta}$  of (6.11) implies (*r*-linear) convergence of the corresponding sequence  $u^{\hat{\delta}\delta}$ .

We equip  $X^{\delta}$  with a family of 'energy' norms

$$||w||_{X^{\hat{\delta}\delta}}^2 := ||w||_Y^2 + \sup_{0 \neq v \in Y^{\hat{\delta}}} \frac{(\partial_t w)(v)^2}{||v||_Y^2} + ||\gamma_T w||^2.$$

By definition of  $\gamma_{\Delta}$  it holds that

(6.17) 
$$\gamma_{\Delta} \| \cdot \|_X \le \| \cdot \|_{X^{\hat{\delta}\delta}} \le \| \cdot \|_X \quad \text{on } X^{\delta}.$$

As follows from [SW21b, Lemma 3.3], similar to Lemma 6.2.2 we have the following result.

**Lemma 6.3.3.** Thanks to (6.13) (and  $Y^{\delta} \subseteq Y^{\hat{\delta}}$ ), for  $w \in X^{\delta}$  it holds that

$$\|w\|_{X^{\hat{\delta}\delta}}^2 = (S^{\hat{\delta}\delta}w)(w) = \sup_{0 \neq v \in Y^{\hat{\delta}}} \frac{(Bw)(v)^2}{\|v\|_Y^2} + \|\gamma_0 w\|_H^2$$

By using additionally (6.14) this result shows that  $(S^{\hat{\delta}\delta} \cdot)(\cdot)$  is coercive on  $X^{\delta} \times X^{\delta}$  so that (6.16), and thus (6.15), has a unique solution.

Moreover, we have the following result.

**Theorem 6.3.4** ([SW21b, Thm. 3.7]). *Thanks to* (6.13) (and  $Y^{\delta} \subseteq Y^{\hat{\delta}}$ ) and (6.14), *it holds that* 

(6.18) 
$$\|u - u_{\delta}\|_{X} \le \|u - u^{\delta\delta}\|_{X} \le \gamma_{\Delta}^{-1} \|u - u_{\delta}\|_{X}.$$

*Remark* 6.3.5. Without the assumption of  $a(t; \cdot, \cdot)$  being symmetric, the operator A in (6.8), (6.9), (6.10), (6.15), (6.16), and in the definition of  $\|\cdot\|_Y$  should be replaced by  $A_s := \frac{1}{2}(A+A')$ , whereas  $\partial_t$  in the definitions of  $\|\cdot\|_X$ ,  $\gamma_\Delta$  in (6.18), and  $\|\cdot\|_{X^{\hat{\delta}\delta}}$  should be replaced by  $\partial_t + A_a$ , where  $A_a := \frac{1}{2}(A - A')$ . Then, as shown in [SW21a, Thm. 6.1], it holds that  $\|u - u^{\hat{\delta}\delta}\|_X \le \gamma_\Delta^{-2} \|u - u_\delta\|_X$ .

Still without assuming that  $a(t; \cdot, \cdot)$  is symmetric, it is interesting that under the original, easier to demonstrate inf-sup condition (6.14) in terms of  $\partial_t$ , a quasi-optimality result similar to (6.18) can be shown, where then the upper bound for  $||u-u\hat{\delta}\delta||_X/||u-u_\delta||_X$  depends on  $||A_a||_{\mathcal{L}(Y,Y')}$ , and cannot be driven to 1 by taking  $Y^{\delta}$  sufficiently large in relation to  $X^{\delta}$ . The latter, however, will be essential for the analysis in the current work, being the reason why we consider only symmetric  $a(t; \cdot, \cdot)$ .

#### 6.3.3 Modified discretized saddle-point

In view of obtaining an efficient implementation, in the definition of  $(\mu^{\hat{\delta}\delta}, u^{\hat{\delta}\delta})$ in (6.15), and so in that of  $S^{\hat{\delta}\delta}$  and  $f^{\hat{\delta}\delta}$  in (6.16), we replace  $(E_Y^{\hat{\delta}'}AE_Y^{\hat{\delta}})^{-1}$  by some  $K_Y^{\hat{\delta}} = K_Y^{\hat{\delta}'} \in \mathcal{L}is(Y^{\hat{\delta}'}, Y^{\hat{\delta}})$  for which both, for some constant  $\kappa_{\Delta} \geq 1$ ,

(6.19) 
$$\frac{((K_Y^{\hat{\delta}})^{-1}v)(v)}{(Av)(v)} \in [\kappa_{\Delta}^{-1}, \kappa_{\Delta}] \quad (\delta \in \Delta, \, v \in Y^{\hat{\delta}})$$

(i.e.  $K_Y^{\hat{\delta}}$  is an optimal (self-adjoint and coercive) *preconditioner* for  $E_Y^{\hat{\delta}'}AE_Y^{\hat{\delta}}$ ), and which can be applied at linear cost. The resulting system (6.16) is now amenable to the application of the (preconditioned) conjugate residuals iteration.

Despite this modification, we keep using the old notations for  $\mu^{\hat{\delta}\delta}$ ,  $u^{\hat{\delta}\delta}$ ,  $S^{\hat{\delta}\delta}$ ,  $\|\cdot\|_{X^{\hat{\delta}\delta}} := (S^{\hat{\delta}\delta} \cdot)(\cdot)^{\frac{1}{2}}$ , and  $f^{\hat{\delta}\delta}$ .

As shown in [SW21b, Remark 3.8], instead of (6.18) now it holds that

(6.20) 
$$\|u-u_{\delta}\|_{X} \leq \|u-u^{\hat{\delta}\delta}\|_{X} \leq \frac{\kappa_{\Delta}}{\gamma_{\Delta}}\|u-u_{\delta}\|_{X},$$

whereas one deduces that (6.17) now should be read as

(6.21) 
$$\frac{\gamma_{\Delta}}{\sqrt{\kappa_{\Delta}}} \| \cdot \|_{X} \le \| \cdot \|_{X^{\delta\delta}} \le \sqrt{\kappa_{\Delta}} \| \cdot \|_{X} \quad \text{on } X^{\delta}.$$

For our forthcoming analysis, we will need  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1$  to be sufficiently small.

*Remark* 6.3.6. Later, in the proof of Proposition 6.4.5, temporarily we will consider the system (6.15) with  $\hat{\delta} = \delta$  (i.e.  $Y^{\hat{\delta}} = Y^{\delta}$ ), but with  $X^{\delta}$  replaced by X, and, as we do in the current subsection,  $E_Y^{\delta'}AE_Y^{\delta}$  replaced by  $(K_Y^{\delta})^{-1}$ . The resulting Schur operator  $B'E_Y^{\delta}K_Y^{\delta}E_Y^{\delta'}B + \gamma'_0\gamma_0$  will be denoted as  $S^{\delta\infty}$ .

Notice that the exact solution u solves  $S^{\delta \infty} u = B' E_Y^{\delta} K_Y^{\delta} E_Y^{\delta'} g + \gamma'_0 u_0$ . Observing that  $S^{\delta \delta} = E_X^{\delta'} S^{\delta \infty} E_X^{\delta}$ , we have the Galerkin orthogonality  $(S^{\delta \infty}(u - u^{\delta \delta}))(X^{\delta}) = 0$ . It holds that

$$\| \cdot \|_{X^{\delta\infty}} := (S^{\delta\infty} \cdot)(\cdot)^{\frac{1}{2}} = \sqrt{\sup_{0 \neq v \in Y^{\delta}} \frac{((B \cdot)(v))^{2}}{((K_{Y}^{\delta})^{-1}v)(v)}} + \|\gamma_{0} \cdot \|_{H}^{2}$$
$$= \sqrt{(E_{Y}^{\delta'}B \cdot)(K_{Y}^{\delta}E_{Y}^{\delta'}B \cdot)} + \|\gamma_{0} \cdot \|_{H}^{2}$$

is only a *semi*-norm on *X*, which is equal to  $\|\cdot\|_{X^{\delta\delta}}$  on  $X^{\delta}$ , and

(6.22)  $\|\cdot\|_{X^{\delta\infty}} \leq \sqrt{\kappa_{\Delta}} \|\cdot\|_X$  on X.

# 6.4 Convergent adaptive solution method

## 6.4.1 Preliminaries

For  $\delta \in \Delta$ , we consider the modified discretized saddle point problem (i.e. (6.16) with  $(E_Y^{\hat{\delta}'}AE_Y^{\hat{\delta}})^{-1}$  replaced by  $K_Y^{\hat{\delta}}$ ) taking  $\hat{\delta} := \underline{\delta}$  from Assumption 6.3.1. So for a given 'trial space'  $X^{\delta}$ , we employ  $Y^{\hat{\delta}}$  as 'test space', which is known to be sufficiently large to give stability even when employed with trial space  $X^{\hat{\delta}} \supseteq X^{\delta}$ . We will use this room to (adaptively) expand  $X^{\delta}$  to some  $X^{\hat{\delta}} \subset X^{\hat{\delta}}$  while keeping  $Y^{\hat{\delta}}$  fixed. Then in a second step we adapt the test space to the new trial space, i.e., replace  $Y^{\hat{\delta}}$  by  $Y^{\hat{\delta}}$ . By doing so will construct a sequence  $(\delta_i) \subseteq \Delta$  with  $\delta_i \preceq \delta_{i+1}$  such that  $(u^{\hat{\delta}_i \delta_i})_i$  converges *r*-linearly to *u*.

As a first step, in the next lemma it is shown that if one constructs from  $w \in X^{\delta}$  a  $v \in X^{\delta}$  that is closer to the best approximation  $u_{\delta}$  to u from  $X^{\delta}$ , then, thanks to Assumption 6.3.1, v is also closer to u.

**Lemma 6.4.1.** Let  $w \in X^{\delta}$ ,  $v \in X^{\underline{\delta}}$  be such that for some  $\rho \leq 1$ ,

$$\|u_{\underline{\delta}} - v\|_X \le \rho \|u_{\underline{\delta}} - w\|_X$$

Then

$$||u - v||_X \le \sqrt{\zeta^2 + \rho^2 (1 - \zeta^2)} ||u - w||_X$$

*Proof.* Using  $u - u_{\delta} \perp_X X^{\underline{\delta}}$  twice, we obtain

$$\begin{split} \|u - v\|_X^2 &= \|u - u_{\bar{\delta}}\|_X^2 + \|u_{\bar{\delta}} - v\|_X^2 \\ &\leq \|u - u_{\bar{\delta}}\|_X^2 + \rho^2 \|u_{\bar{\delta}} - w\|_X^2 \\ &= \|u - u_{\bar{\delta}}\|_X^2 + \rho^2 (\|u - w\|_X^2 - \|u - u_{\bar{\delta}}\|_X^2) \\ &= (1 - \rho^2) \|u - u_{\bar{\delta}}\|_X^2 + \rho^2 \|u - w\|_X^2 \\ &\leq (\zeta^2 (1 - \rho^2) + \rho^2) \|u - w\|_X^2, \end{split}$$

where we used Assumption 6.3.1 and  $||u - u_{\delta}||_X \le ||u - w||_X$ .

Notice that  $u^{\delta\delta}$  is the Galerkin approximation from  $X^{\delta}$  to the solution  $u^{\delta\delta} \in X^{\delta}$  of the system  $S^{\delta\delta}u^{\delta\delta} = f^{\delta\delta}$ , i.e., it is its best approximation from  $X^{\delta}$  w.r.t.  $\|\cdot\|_{X^{\delta\delta}}$ . In the next proposition it is shown that an improved Galerkin approximation from an intermediate space  $X^{\delta} \supseteq X^{\delta} \supseteq X^{\delta}$ , i.e., the function  $u^{\delta\delta}$ , is, for  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1$  sufficiently small, also an improved approximation to u, and furthermore that this holds true also for  $u^{\delta\delta}$ . The latter function will be the successor of  $u^{\delta\delta}$  in our converging sequence.

**Proposition 6.4.2.** Let  $\delta \preceq \delta$  be such that

(6.23) 
$$\|u^{\underline{\delta}\underline{\delta}} - u^{\underline{\delta}\overline{\delta}}\|_{X^{\underline{\delta}\underline{\delta}}} \le \rho \|u^{\underline{\delta}\underline{\delta}} - u^{\underline{\delta}\delta}\|_{X^{\underline{\delta}\underline{\delta}}}.$$

Then it holds that

$$\|u - u^{\tilde{\delta}\tilde{\delta}}\|_X \le \underbrace{\frac{\kappa_\Delta}{\gamma_\Delta}\sqrt{\zeta^2 + \hat{\rho}^2(1 - \zeta^2)}}_{\bar{\rho}:=} \|u - u^{\tilde{\delta}\delta}\|_X,$$

where  $\hat{\rho} := (1 + \rho \frac{\sqrt{\kappa_{\Delta}}}{\gamma_{\Delta}}) \sqrt{\frac{\kappa_{\Delta}^2}{\gamma_{\Delta}^2}} - 1 \sqrt{\frac{\zeta^2}{1-\zeta^2}} + \rho \frac{\sqrt{\kappa_{\Delta}}}{\gamma_{\Delta}}$ . Notice that  $\hat{\rho}$  and  $\bar{\rho}$  are < 1 when  $\rho < 1$  and  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1$  is sufficiently small dependent on  $\rho$  with  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1 \downarrow 0$  when  $\rho \uparrow 1$ .

*Proof.* Using that  $u - u_{\underline{\delta}} \perp_X X^{\underline{\delta}}$ , it follows that  $||u - u^{\underline{\delta}\delta}||_X^2 \leq \frac{\kappa_{\underline{\Delta}}^2}{\gamma_{\underline{\Delta}}^2} ||u - u_{\delta}||_X^2$ ((6.20)) is equivalent to  $||u_{\underline{\delta}} - u^{\underline{\delta}\delta}||_X \leq \sqrt{\frac{\kappa_{\underline{\Delta}}^2}{\gamma_{\underline{\Delta}}^2} - 1} ||u - u_{\delta}||_X$ . Similarly, Assumption 6.3.1 is equivalent to  $||u - u_{\delta}||_X \leq \sqrt{\frac{\zeta^2}{1-\zeta^2}} ||u_{\underline{\delta}} - w||_X$  for any  $w \in X^{\delta}$ .

Additionally using (6.21), we infer that

$$\begin{split} \|u_{\underline{\delta}} - u^{\delta\bar{\delta}}\|_{X} &\leq \|u_{\underline{\delta}} - u^{\delta\bar{\delta}}\|_{X} + \|u^{\underline{\delta}\underline{\delta}} - u^{\delta\bar{\delta}}\|_{X} \\ &\leq \|u_{\underline{\delta}} - u^{\delta\underline{\delta}}\|_{X} + \frac{\sqrt{\kappa_{\Delta}}}{\gamma_{\Delta}}\|u^{\underline{\delta}\underline{\delta}} - u^{\delta\bar{\delta}}\|_{X^{\underline{\delta}\underline{\delta}}} \\ &\leq \|u_{\underline{\delta}} - u^{\underline{\delta}\underline{\delta}}\|_{X} + \rho\frac{\sqrt{\kappa_{\Delta}}}{\gamma_{\Delta}}\|u^{\underline{\delta}\underline{\delta}} - u^{\delta\bar{\delta}}\|_{X} \\ &\leq \left(1 + \rho\frac{\sqrt{\kappa_{\Delta}}}{\gamma_{\Delta}}\right)\|u_{\underline{\delta}} - u^{\underline{\delta}\underline{\delta}}\|_{X} + \rho\frac{\sqrt{\kappa_{\Delta}}}{\gamma_{\Delta}}\|u_{\underline{\delta}} - u^{\underline{\delta}\underline{\delta}}\|_{X} \\ &\leq \left[\left(1 + \rho\frac{\sqrt{\kappa_{\Delta}}}{\gamma_{\Delta}}\right)\sqrt{\frac{\kappa_{\Delta}^{2}}{\gamma_{\Delta}^{2}} - 1}\sqrt{\frac{\zeta^{2}}{1-\zeta^{2}}} + \rho\frac{\sqrt{\kappa_{\Delta}}}{\gamma_{\Delta}}\right]\|u_{\underline{\delta}} - u^{\underline{\delta}\underline{\delta}}\|_{X} \end{split}$$

From Lemma 6.4.1 we conclude that  $\|u-u^{\delta \tilde{\delta}}\|_X \leq \sqrt{\zeta^2 + \hat{\rho}^2(1-\zeta^2)} \|u-u^{\delta \delta}\|_X$ . Thanks to (6.20), it holds that

$$\|u-u^{\tilde{\delta}\tilde{\delta}}\|_{X} \leq \frac{\kappa_{\Delta}}{\gamma_{\Delta}} \|u-u_{\tilde{\delta}}\|_{X} \leq \frac{\kappa_{\Delta}}{\gamma_{\Delta}} \|u-u^{\delta\tilde{\delta}}\|_{X}$$

which completes the proof.

#### 6.4.2 Bulk chasing and a posteriori error estimation

To realize (6.23), i.e., to construct from the Galerkin approximation  $u^{\delta\delta}$  to  $u^{\delta\delta}$  an improved Galerkin approximation  $u^{\delta\delta}$ , we apply the concept of bulk chasing, also known as Dörfler marking, on a collection of a posteriori error indicators that constitute an efficient and reliable error estimator. We will apply an estimator of 'hierarchical basis' type ([ZMD<sup>+</sup>11]):

Let  $\Theta_{\delta} = \{\theta_{\lambda} : \lambda \in J_{\delta}\} \subseteq X^{\delta}$  be such that  $X^{\delta} + \operatorname{span} \Theta_{\delta} = X^{\delta}$  and, for some constants  $0 < m \le M$ , for all  $\delta \in \Delta$ ,  $z \in X^{\delta}$  and  $\mathbf{c} := (c_{\lambda})_{\lambda \in J_{\delta}} \subset \mathbb{R}$ .

(6.24) 
$$m^2 \|z + \mathbf{c}^\top \Theta_\delta\|_X^2 \le \|z\|_X^2 + \|\mathbf{c}\|^2 \le M^2 \|z + \mathbf{c}^\top \Theta_\delta\|_X^2.$$

A suitable collection  $\Theta_{\delta}$  will be constructed in Sect. 6.6.1.

**Proposition 6.4.3.** Assume (6.24). Let  $\mathbf{r}_{\delta}^{\underline{\delta}} := (f^{\underline{\delta}\underline{\delta}} - S^{\underline{\delta}\underline{\delta}}u^{\underline{\delta}\delta})(\Theta_{\delta})$ , being the residual vector of  $u^{\underline{\delta}\underline{\delta}}$ . Let  $J \subseteq J_{\delta}$  be such that for some constant  $\vartheta \in (0, 1]$ ,

$$\|\mathbf{r}_{\delta}^{\underline{\delta}}|_{J}\| \geq \vartheta \|\mathbf{r}_{\delta}^{\underline{\delta}}\|,$$

and, for some  $\tilde{\delta} \leq \underline{\delta}$ , let  $X^{\delta} + \operatorname{span} \Theta_{\delta}|_{J} \subseteq X^{\tilde{\delta}}$ . Then with  $\rho := \sqrt{1 - \left(\frac{m}{M} \frac{\gamma_{\Delta}}{\kappa_{\Delta}^{2}} \vartheta\right)^{2}}$ , (6.23) is valid, i.e.,

(6.25) 
$$\|u^{\delta\delta} - u^{\delta\bar{\delta}}\|_{X^{\underline{\delta}\underline{\delta}}} \le \rho \|u^{\underline{\delta}\underline{\delta}} - u^{\delta\delta}\|_{X^{\underline{\delta}\underline{\delta}}};$$

113

and so, when  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1$  is sufficiently small dependent on  $\vartheta$ , with  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1 \downarrow 0$  when  $\vartheta \downarrow 0$ , for some constant  $\rho < \bar{\rho} < 1$ ,

$$\|u - u^{\underline{\tilde{\delta}}} \|_X \le \bar{\rho} \|u - u^{\underline{\delta}} \|_X.$$

*Proof.* As a consequence of (6.24) and (6.21), we have

$$\frac{\gamma_{\Delta}}{\kappa_{\Delta}^2}m^2\|z+\mathbf{c}^{\top}\Theta_{\delta}\|_{X^{\underline{\delta}\underline{\delta}}}^2 \leq \|z\|_{X^{\underline{\delta}\underline{\delta}}}^2 + \|\mathbf{c}\|^2 \leq \frac{\kappa_{\Delta}^2}{\gamma_{\Delta}}M^2\|z+\mathbf{c}^{\top}\Theta_{\delta}\|_{X^{\underline{\delta}\underline{\delta}}}^2.$$

We infer that

$$\begin{split} \|u^{\delta\tilde{\delta}} - u^{\delta\delta}\|_{X^{\delta\tilde{\delta}}} &= \sup_{0 \neq (z, \mathbf{c}) \in X^{\delta} \times \mathbb{R}^{\#J_{\delta}}} \frac{(S^{\delta\delta}(u^{\delta\delta} - u^{\delta\delta}))(z + \mathbf{c}^{\top}\Theta_{\delta})}{\|z + \mathbf{c}^{\top}\Theta_{\delta}\|_{X^{\delta\tilde{\delta}}}} \\ &\geq m \frac{\sqrt{\gamma_{\Delta}}}{\kappa_{\Delta}} \sup_{0 \neq (z, \mathbf{c}) \in X^{\delta} \times \mathbb{R}^{\#J_{\delta}}} \frac{(S^{\delta\delta}(u^{\delta\tilde{\delta}} - u^{\delta\delta}))(\mathbf{c}^{\top}\Theta_{\delta})}{\sqrt{\|z\|_{X^{\delta\tilde{\delta}}}^{2}} + \|\mathbf{c}\|^{2}}} \\ &\geq m \frac{\sqrt{\gamma_{\Delta}}}{\kappa_{\Delta}} \sup_{0 \neq \mathbf{c} \in \mathbb{R}^{\#J_{\delta}}} \frac{\langle \mathbf{c}|_{J}, (f^{\delta\tilde{\delta}} - S^{\delta\delta}u^{\delta\delta})(\Theta_{\delta}|_{J})\rangle}{\|\mathbf{c}|_{J}\|} \\ &= m \frac{\sqrt{\gamma_{\Delta}}}{\kappa_{\Delta}} \|\mathbf{r}_{\delta}^{\delta}|_{J}\| \geq m \frac{\sqrt{\gamma_{\Delta}}}{\kappa_{\Delta}} \vartheta \|\mathbf{r}_{\delta}^{\delta}\| \\ &= m \frac{\sqrt{\gamma_{\Delta}}}{\kappa_{\Delta}} \vartheta_{0 \neq (z, \mathbf{c}) \in X^{\delta} \times \mathbb{R}^{\#J_{\delta}}} \frac{(S^{\delta\tilde{\delta}}(u^{\delta\tilde{\delta}} - u^{\delta\tilde{\delta}}))(\mathbf{c}^{\top}\Theta_{\delta})}{\sqrt{\|z\|_{X^{\delta\tilde{\delta}}}^{2}} + \|\mathbf{c}\|^{2}} \\ &\geq \frac{m}{M} \frac{\gamma_{\Delta}}{\kappa_{\Delta}} \vartheta \|u^{\delta\delta} - u^{\delta\delta}\|_{X^{\delta\delta}}. \end{split}$$

so that

$$\begin{split} \|u^{\underline{\delta}\underline{\delta}} - u^{\underline{\delta}\overline{\delta}}\|_{X^{\underline{\delta}\underline{\delta}}}^2 &= \|u^{\underline{\delta}\underline{\delta}} - u^{\underline{\delta}\delta}\|_{X^{\underline{\delta}\underline{\delta}}}^2 - \|u^{\underline{\delta}\overline{\delta}} - u^{\underline{\delta}\delta}\|_{X^{\underline{\delta}\underline{\delta}}}^2 \\ &\leq \Big(1 - \big(\frac{m}{M}\frac{\gamma\Delta}{\kappa_{\Delta}^2}\vartheta\big)^2\Big)\|u^{\underline{\delta}\underline{\delta}} - u^{\underline{\delta}\delta}\|_{X^{\underline{\delta}\underline{\delta}}}^2, \end{split}$$

which completes the proof of (6.25).

The final statement follows from an application of Proposition 6.4.2.  $\Box$ 

Additionally we have that  $\|\mathbf{r}_{\delta}^{\delta}\|$  provides an efficient and reliable a posteriori estimator for  $\|u - u^{\delta\delta}\|_X$ :

**Proposition 6.4.4.** Assume (6.24). Recalling that  $\zeta < 1$ , let  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} < \frac{1}{\zeta}$ . Then for  $\delta \in \Delta$ ,

$$\frac{m\frac{\sqrt{\gamma_{\Delta}}}{\kappa_{\Delta}^{3/2}}}{1+\zeta\sqrt{\frac{\kappa_{\Delta}^{2}}{\gamma_{\Delta}^{2}}-1}}\|\mathbf{r}_{\delta}^{\delta}\| \leq \|u-u^{\delta\delta}\|_{X} \leq \frac{M\frac{\kappa_{\Delta}^{5/2}}{\gamma_{\Delta}^{3/2}}}{\sqrt{1-\zeta^{2}}-\zeta\sqrt{\frac{\kappa_{\Delta}^{2}}{\gamma_{\Delta}^{2}}-1}}\|\mathbf{r}_{\delta}^{\delta}\|.$$

*Proof.* Assumption 6.3.1 gives  $||u-u_{\delta}||_X \leq \zeta ||u-u^{\delta\delta}||_X$ , which by  $u-u_{\delta} \perp_X X^{\delta}$  yields

(6.27) 
$$||u_{\underline{\delta}} - u^{\underline{\delta}\delta}||_X \le ||u - u^{\underline{\delta}\delta}||_X \le \frac{1}{\sqrt{1-\zeta^2}} ||u_{\underline{\delta}} - u^{\underline{\delta}\delta}||_X$$

As we already have noted in the proof of Proposition 6.4.2, (6.20) is equivalent to  $\|u_{\underline{\delta}} - u^{\underline{\delta}\underline{\delta}}\|_X \le \sqrt{\frac{\kappa_{\underline{\Delta}}^2}{\gamma_{\underline{\Delta}}^2} - 1} \|u - u_{\underline{\delta}}\|_X$ . Together with Assumption 6.3.1, it gives

$$\left| \|u_{\underline{\delta}} - u^{\underline{\delta}\delta}\|_{X} - \|u^{\underline{\delta}\delta} - u^{\underline{\delta}\delta}\|_{X} \right| \le \zeta \sqrt{\frac{\kappa_{\Delta}^{2}}{\gamma_{\Delta}^{2}} - 1} \|u - u^{\underline{\delta}\delta}\|_{X}$$

which in combination with (6.27) and  $\zeta \frac{\kappa_{\Delta}}{\gamma_{\Delta}} < 1$  yields

$$\frac{1}{1+\zeta\sqrt{\frac{\kappa_{\Delta}^2}{\gamma_{\Delta}^2}-1}}\|u^{\underline{\delta}\underline{\delta}}-u^{\underline{\delta}\underline{\delta}}\|_X \le \|u-u^{\underline{\delta}\underline{\delta}}\|_X \le \frac{1}{\sqrt{1-\zeta^2}-\zeta\sqrt{\frac{\kappa_{\Delta}^2}{\gamma_{\Delta}^2}-1}}\|u^{\underline{\delta}\underline{\delta}}-u^{\underline{\delta}\underline{\delta}}\|_X.$$

The proof is completed by (6.21), and  $m \frac{\sqrt{\gamma_{\Delta}}}{\kappa_{\Delta}} \|\mathbf{r}_{\delta}^{\underline{\delta}}\| \le \|u^{\underline{\delta}\underline{\delta}} - u^{\underline{\delta}\delta}\|_{X^{\underline{\delta}\underline{\delta}}} \le M \frac{\kappa_{\Delta}}{\sqrt{\gamma_{\Delta}}} \|\mathbf{r}_{\delta}^{\underline{\delta}}\|$ where the latter inequalities were shown in (6.26) when reading  $(J, \vartheta, \tilde{\delta})$  as  $(J_{\delta}, 1, \underline{\delta})$ .

Next we present an alternative a posteriori error estimator that does not rely on (6.24), that we expect to be more accurate, and that can be a computed at the cost of one additional inner product.

**Proposition 6.4.5.** Let  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} < \frac{1}{\zeta}$ , and for  $v \in X^{\delta}$ , define

$$\mathcal{E}^{\delta}(v) = \mathcal{E}^{\delta}(v; g, u_0) := \sqrt{(E_Y^{\delta'}(g - Bv))(K_Y^{\delta} E_Y^{\delta'}(g - Bv)) + \|u_0 - \gamma_0 v\|_H^2}$$

Then

$$\frac{\gamma_{\Delta} - \zeta \kappa_{\Delta}}{\sqrt{\kappa_{\Delta}}} \|u - v\|_X \le \mathcal{E}^{\delta}(v) \le \sqrt{\kappa_{\Delta}(\zeta^2 + (1 + \zeta \frac{\kappa_{\Delta}}{\gamma_{\Delta}})^2)} \|u - v\|_X$$

*Proof.* From Remark 6.3.6, recall that the semi-norm  $\|\cdot\|_{X^{\delta\infty}}$  on X equals  $\|\cdot\|_{X^{\delta\delta}}$  on  $X^{\delta}$ , and  $(S^{\delta\infty}(u-u^{\delta\delta}))(X^{\delta}) = 0$ , which implies

(6.28) 
$$||u - w||^2_{X^{\delta\infty}} = ||u - u^{\delta\delta}||^2_{X^{\delta\infty}} + ||u^{\delta\delta} - w||^2_{X^{\delta\infty}} \quad (w \in X^{\delta}),$$

and  $\|\cdot\|_{X^{\delta\infty}} \leq \sqrt{\kappa_{\Delta}} \|\cdot\|_X$  ((6.22)).

From (6.20) and Assumption 6.3.1, we have for  $v \in X^{\delta}$ 

(6.29) 
$$\|u - u^{\underline{\delta}\underline{\delta}}\|_X \le \frac{\kappa_{\Delta}}{\gamma_{\Delta}} \|u - u_{\underline{\delta}}\|_X \le \zeta \frac{\kappa_{\Delta}}{\gamma_{\Delta}} \|u - v\|_X$$

From (6.29), the triangle-inequality, (6.21), and (6.28) we obtain for  $v \in X^{\delta}$ 

$$\begin{split} \|u - v\|_X &\leq \frac{1}{1 - \zeta \frac{\kappa_\Delta}{\gamma_\Delta}} \|u^{\underline{\delta}\underline{\delta}} - v\|_X \\ &\leq \frac{\sqrt{\kappa_\Delta}}{\gamma_\Delta - \zeta\kappa_\Delta} \|u^{\underline{\delta}\underline{\delta}} - v\|_{X^{\underline{\delta}\underline{\delta}}} \leq \frac{\sqrt{\kappa_\Delta}}{\gamma_\Delta - \zeta\kappa_\Delta} \|u - v\|_{X^{\underline{\delta}\infty}}. \end{split}$$

Conversely, we have

$$\begin{aligned} \|u-v\|_{X^{\delta\infty}}^2 \stackrel{(6.28)}{=} \|u-u^{\delta\delta}\|_{X^{\delta\infty}}^2 + \|u^{\delta\delta}-v\|_{X^{\delta\infty}}^2 \\ &\stackrel{(6.28),(6.22)}{\leq} \|u-u_{\delta}\|_{X^{\delta\infty}}^2 + \kappa_{\Delta}\|u^{\delta\delta}-v\|_X^2 \\ &\stackrel{(6.22),(6.29)}{\leq} \kappa_{\Delta}\|u-u_{\delta}\|_X^2 + \kappa_{\Delta}(1+\zeta\frac{\kappa_{\Delta}}{\gamma_{\Delta}})^2\|u-v\|_X^2 \\ &\stackrel{\leq}{\leq} \kappa_{\Delta}(\zeta^2 + (1+\zeta\frac{\kappa_{\Delta}}{\gamma_{\Delta}})^2)\|u-v\|_X^2. \end{aligned}$$

by again applying Assumption 6.3.1.

Noting that  $||u-v||^2_{X^{\delta\infty}} := E_Y^{\delta'}(g-Bv)(K_Y^{\delta}E_Y^{\delta'}(g-Bv)) + ||u_0-\gamma_0 v||^2_{H'}$ the proof is completed.

Notice that the estimator of  $||u - v||_X$  from Proposition 6.4.5 is exact when  $\zeta = 0$  and  $\kappa_{\Delta} = 1 = \gamma_{\Delta}$ , where the one from Proposition 6.4.3, for  $v = u^{\delta\delta}$ , is exact only when additionally m = 1 = M.

#### 6.4.3 Data oscillation

In view of the discussion following Assumption 6.3.1, notice that *all* results obtained so far that depend on the 'saturation constant'  $\zeta$ , i.e., Lemma 6.4.1, and Propositions 6.4.2, 6.4.3, and 6.4.4, are only valid under the condition that  $(g, u_0) \in {}^{\delta}G \times {}^{\delta}U_0$ .

Let us now consider the situation that the solutions and residuals in these statements refer to solutions and residuals with the true data  $(g, u_0) \in Y' \times H$  being *replaced* by an approximation  $({}^{\delta}g, {}^{\delta}u_0) \in {}^{\delta}G \times {}^{\delta}U_0$ . In the following we denote such solutions and residuals with an *additional left superscript*  $\delta$ , or more generally  $\tilde{\delta}$  when a right-hand side  $({}^{\tilde{\delta}}g, {}^{\tilde{\delta}}u_0) \in {}^{\tilde{\delta}}G \times {}^{\tilde{\delta}}U_0$  has been used for their computation.

**Proposition 6.4.6.** Assume (6.24), and let  $\vartheta \in (0,1]$  be a constant. Then for  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1$  and a constant  $\hat{\omega} > 0$  both being sufficiently small dependent on  $\vartheta$ , with  $\max(\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1, \hat{\omega}) \downarrow 0$  when  $\vartheta \downarrow 0$ , there exists a constant  $\check{\rho} < 1$  such that for  $J \subseteq J_{\delta}$  with  $\|\delta \mathbf{r}_{\delta}^{\underline{\delta}}\|_{J} \| \ge \vartheta \|\delta \mathbf{r}_{\delta}^{\underline{\delta}}\|$ , and  $X^{\delta} + \operatorname{span} \Theta_{\delta}|_{J} \subseteq X^{\tilde{\delta}}$ , and

$$\max\left(\|g - {}^{\delta}g\|_{Y'} + \|u_0 - {}^{\delta}u_0\|_H, \|g - {}^{\tilde{\delta}}g\|_{Y'} + \|u_0 - {}^{\tilde{\delta}}u_0\|_H\right) \le \hat{\omega}\|^{\delta}\mathbf{r}_{\delta}^{\frac{\delta}{2}}\|,$$

it holds that

$$\|^{\tilde{\delta}}u - {}^{\tilde{\delta}}u^{\underline{\tilde{\delta}}}\delta \|_X \le \check{\rho}\|^{\delta}u - {}^{\delta}u^{\underline{\delta}}\delta \|_X.$$

*Proof.* In the newly introduced notations, the statements of Propositions 6.4.3 and 6.4.4 read as

$$\|{}^{\delta}u - {}^{\delta}u^{\underline{\delta}\delta}\|_X \le \bar{\rho}\|{}^{\delta}u - {}^{\delta}u^{\underline{\delta}\delta}\|_X,$$

and

(6.30)  $\|^{\delta} \mathbf{r}_{\delta}^{\delta}\| \approx \|^{\delta} u - {}^{\delta} u^{\delta} \|_{X}$ 

The proof is easily completed by

$$\| {}^{\delta u} - {}^{\delta u} \|_X \\ \| {}^{\tilde{\delta}} u^{\tilde{\delta} \tilde{\delta}} - {}^{\delta u} {}^{\tilde{\delta} \tilde{\delta}} \|_X \\ \right\} \lesssim \| {}^{\tilde{\delta}} g - {}^{\delta} g \|_{Y'} + \| {}^{\tilde{\delta}} u_0 - {}^{\delta} u_0 \|_H \le 2\hat{\omega} \| {}^{\delta} \mathbf{r}_{\tilde{\delta}}^{\delta} \| \approx \hat{\omega} \| {}^{\delta} u - {}^{\delta} u^{\delta \delta} \|_X. \square$$

In view of the latter proposition, we make the following assumption.

**Assumption 6.4.7.** We assume to have maps of the following types available:

$$\begin{split} &\Delta \to Y' \times H \colon \delta \mapsto ({}^{\delta}g, {}^{\delta}u_0) \in {}^{\delta}G \times {}^{\delta}U_0, \\ &\eta \colon \Delta \to \mathbb{R} \text{ such that } \|g - {}^{\delta}g\|_{Y'} + \|u_0 - {}^{\delta}u_0\|_H \leq \eta(\delta), \text{ and } \eta(\tilde{\delta}) \leq \eta(\delta) \text{ when } \tilde{\delta} \succeq \delta, \\ &\mathbb{R}_{>0} \to \Delta \colon \varepsilon \mapsto \delta(\varepsilon) \text{ such that } \eta(\delta(\varepsilon)) \leq \varepsilon. \end{split}$$

Notice that this in particular means that for any  $\varepsilon > 0$  we are able to find a  $\delta \in \Delta$  and  $({}^{\delta}g, {}^{\delta}u_0) \in {}^{\delta}G \times {}^{\delta}U_0$  with  $||g - {}^{\delta}g||_{Y'} + ||u_0 - {}^{\delta}u_0||_H \le \varepsilon$ . A specification of a suitable family  $({}^{\delta}G, {}^{\delta}U_0)_{\delta \in \Delta}$  will be given in Sect. 6.6.4.

Given a  $\delta \in \Delta$ , and thinking of  $({}^{\delta}g, {}^{\delta}u_0)$  being a quasi-best approximation to  $(g, u_0)$  from  ${}^{\delta}G \times {}^{\delta}U_0$ , the difference  $(g, u_0) - ({}^{\delta}g, {}^{\delta}u_0)$  is often referred to as *data-oscillation*.

#### 6.4.4 A convergent algorithm

In view of the statement from Proposition 6.4.6, in the following we will use the short-hand notations

$$u^{\delta} = {}^{\delta} u^{\underline{\delta}\delta}, \quad \mathbf{r}^{\delta} = {}^{\delta} \mathbf{r}^{\underline{\delta}}_{\overline{\delta}},$$

i.e.,  $u^{\delta}$  is the solution of

(6.31) 
$$\underbrace{E_X^{\delta'}(B'E_Y^{\delta}K_Y^{\delta}E_Y^{\delta'}B + \gamma_0'\gamma_0)E_X^{\delta}}_{S^{\delta} =} u^{\delta} = \underbrace{E_X^{\delta'}(B'E_Y^{\delta}K_Y^{\delta}E_Y^{\delta'}g + \gamma_0'^{\delta}u_0)}_{f^{\delta} = \delta f_{\delta}^{\delta} :=}.$$

(cf. (6.16) and Sect. 6.3.3), and

(6.32) 
$$\mathbf{r}^{\delta} = E_X^{\delta'} \left[ B' E_Y^{\delta} K_Y^{\delta} E_Y^{\delta'} ({}^{\delta}g - Bu^{\delta}) + \gamma_0' ({}^{\delta}u_0 - \gamma_0 u^{\delta}) \right] (\Theta_{\delta})$$

Instead of solving (6.31) exactly, we will allow it to be solved approximately with a sufficiently small relative tolerance by the application of an iterative method. To that end, we assume to have available a  $K_X^{\delta} = K_X^{\delta'} \in \mathcal{L}is(X^{\delta'}, X^{\delta})$  for which both

(6.33) 
$$((K_X^{\delta})^{-1}w)(w) \approx \|w\|_X^2 \quad (w \in X^{\delta})$$

(i.e.  $K_X^{\delta}$  is an optimal (self-adjoint and coercive) *preconditioner* for  $S^{\delta\delta}$ ), and which can be applied at linear cost. Besides for an efficient iterative solving of (6.31), we will use this preconditioner to compute a quantity that is equivalent to the *X*-norm of the (algebraic) error in any approximation from  $X^{\delta}$  to  $u^{\delta}$ .

We denote such an *approximate* solution of (6.31) by  $\tilde{u}^{\delta} \in X^{\delta}$ , with corresponding residual vector  $\tilde{\mathbf{r}}^{\delta}$  defined as in (6.32) by replacing  $u^{\delta}$  by  $\tilde{u}^{\delta}$ .

### Algorithm 6.4.8.

Let  $\omega > 0$ ,  $\vartheta \in (0, 1]$ ,  $0 < \xi < 1$  be constants, and let  $\varepsilon > 0$ .  $\delta := \delta_{\text{init}} \in \Delta$ ,  $t_{\delta} \approx \|g\|_{Y'} + \|u_0\|_H$ . do

do compute  $\tilde{u}^{\delta} \in X^{\delta}$  with  $\tilde{t}_{\delta} := \sqrt{(f^{\delta} - S^{\underline{\delta}\delta}\tilde{u}^{\delta})(K_X^{\delta}(f^{\delta} - S^{\underline{\delta}\delta}\tilde{u}^{\delta}))} \leq \frac{t_{\delta}}{2}; t_{\delta} := \tilde{t}_{\delta}$ if  $e_{\delta} := \|\tilde{\mathbf{r}}^{\delta}\| + \eta(\delta) + t_{\delta} \leq \varepsilon$  then stop end if until  $t_{\delta} \leq \xi e_{\delta}$ if  $\eta(\delta) > \omega \|\tilde{\mathbf{r}}^{\delta}\|$ then select  $\tilde{\delta} \in \Delta$  s.t.  $X^{\tilde{\delta}} \supseteq X^{\delta}$  is (a near-smallest) space with  $\eta(\tilde{\delta}) \leq \eta(\delta)/2$ . else determine  $\delta \preceq \tilde{\delta} \preceq \underline{\delta}$  s.t.  $X^{\tilde{\delta}}$  is (a near-smallest) space that for a  $J \subseteq I_{\delta}^{\tilde{\delta}}$ contains  $X^{\delta} + \operatorname{span} \Theta_{\delta}|_{J}$  where  $\|\tilde{\mathbf{r}}^{\delta}|_{J}\| \geq \vartheta \|\tilde{\mathbf{r}}^{\delta}\|$ . end if  $t_{\tilde{\delta}} := e_{\delta}, \delta := \tilde{\delta}$ enddo

**Theorem 6.4.9.** Assume (6.24), and let the constants  $\gamma_{\Delta}$  and  $\kappa_{\Delta}$  be as defined in (6.14) and (6.19), respectively. For constants  $\vartheta$ ,  $\omega/\vartheta$ ,  $\xi/\vartheta$ ,  $(\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1)/\omega$  that are sufficiently small, with additionally  $\omega$  and  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1$  sufficiently small dependent on  $\vartheta$  with  $\max(\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1, \omega) \downarrow 0$  when  $\vartheta \downarrow 0$ , there exists a constant  $\check{\rho} < 1$  such that between any two successive passings of the until-clause the value of  $\vartheta \| u - \tilde{u}^{\delta} \|_X + \eta(\delta)$  decreases with at least a factor  $\check{\rho}$ . For any  $\varepsilon > 0$ , Algorithm 6.4.8 terminates, and at termination it holds that  $\| u - \tilde{u}^{\delta} \|_X + \eta(\delta) \lesssim \varepsilon$ .

*Remark* 6.4.10. Minor adaptations to the proof show that the statement remains true when one takes  $e_{\delta} := \mathcal{E}^{\delta}(\tilde{u}^{\delta}; {}^{\delta}g, {}^{\delta}u_0) + \eta(\delta)$  in Algorithm 6.4.8. Having to compute  $\tilde{\mathbf{r}}^{\delta}$  anyway, the additional cost of computing this  $e_{\delta}$  is small, and it can be expected to be closer to  $||u - \tilde{u}^{\delta}||_X + \eta(\delta)$ .

*Proof.* By replacing w by  $K_X^{\delta} S^{\underline{\delta}\delta} w$  in (6.33), one infers that  $(S^{\underline{\delta}\delta}v)(K_X^{\delta} S^{\underline{\delta}\delta}v) \approx \|v\|_X^2$  for  $v \in X^{\delta}$ , and so

(6.34) 
$$\|u^{\delta} - \tilde{u}^{\delta}\|_X^2 \approx (f^{\delta} - S^{\underline{\delta}\delta}\tilde{u}^{\delta})(K_X^{\delta}(f^{\delta} - S^{\underline{\delta}\delta}\tilde{u}^{\delta})).$$

From (6.21) and (6.24), we have (6.35)

$$\|\mathbf{r}^{\delta} - \tilde{\mathbf{r}}^{\delta}\| = \sup_{0 \neq \mathbf{c} \in \mathbb{R}^{\#J^{\delta}}} \frac{(S^{\underline{\delta}\underline{\delta}}(\tilde{u}^{\delta} - u^{\delta}))(\mathbf{c}^{\top}\Theta_{\delta})}{\|\mathbf{c}\|} \le \frac{\kappa_{\Delta}^{\frac{1}{2}}}{m} \|u^{\delta} - \tilde{u}^{\delta}\|_{X^{\underline{\delta}\underline{\delta}}} \le \frac{\kappa_{\Delta}}{m} \|u^{\delta} - \tilde{u}^{\delta}\|_{X}.$$

If the algorithm stops, then

$$\begin{aligned} \|u - \tilde{u}^{\delta}\|_{X} &\leq \|u - {}^{\delta}u\|_{X} + \|{}^{\delta}u - u^{\delta}\|_{X} + \|u^{\delta} - \tilde{u}^{\delta}\|_{X} \\ &\lesssim \eta(\delta) + \|{}^{\delta}u - u^{\delta}\|_{X} + \|u^{\delta} - \tilde{u}^{\delta}\|_{X} \quad \text{(by Thm. 6.2.1 \& Ass. 6.4.7)} \\ &\stackrel{(6.30)}{\approx} \eta(\delta) + \|\mathbf{r}^{\delta}\| + \|u^{\delta} - \tilde{u}^{\delta}\|_{X} \leq \eta(\delta) + \|\mathbf{\tilde{r}}^{\delta}\| + \|\mathbf{r}^{\delta} - \mathbf{\tilde{r}}^{\delta}\| + \|u^{\delta} - \tilde{u}^{\delta}\|_{X} \\ &\stackrel{(6.35)}{\lesssim} \eta(\delta) + \|\mathbf{\tilde{r}}^{\delta}\| + \|u^{\delta} - \tilde{u}^{\delta}\|_{X} \lesssim \eta(\delta) + \|\mathbf{\tilde{r}}^{\delta}\| + t_{\delta} \leq \varepsilon. \end{aligned}$$

The inner do-loop always terminates either by passing the until-clause or by the stop-statement. Indeed, inside this loop the value of  $t_{\delta}$  is driven to 0, so that  $\|\tilde{\mathbf{r}}^{\delta}\| + \eta(\delta)$  tends to  $\|\mathbf{r}^{\delta}\| + \eta(\delta)$ . So if  $\|\mathbf{r}^{\delta}\| + \eta(\delta) \neq 0$ , then at some moment  $t_{\delta} \leq \xi(\|\tilde{\mathbf{r}}^{\delta}\| + \eta(\delta) + t_{\delta})$ , whereas if  $\|\mathbf{r}^{\delta}\| + \eta(\delta) = 0$  then at some moment  $e_{\delta} \leq \varepsilon$ .

When passing the until-clause, it holds that  $t_{\delta} \leq \xi(\|\mathbf{\tilde{r}}^{\delta}\| + \eta(\delta) + t_{\delta})$ , and so by using  $\xi < 1$  kicking back  $t_{\delta}$ ,

(6.36) 
$$t_{\delta} \lesssim \xi(\|\tilde{\mathbf{r}}^{\delta}\| + \eta(\delta))$$
$$\leq \xi(\|\mathbf{r}^{\delta}\| + \|\tilde{\mathbf{r}}^{\delta} - \mathbf{r}^{\delta}\| + \eta(\delta))$$
$$\lesssim \xi(\|^{\delta}u - u^{\delta}\|_{X} + \|\tilde{u}^{\delta} - u^{\delta}\|_{X} + \eta(\delta))$$
$$\lesssim \xi(\|u - u^{\delta}\|_{X} + t_{\delta} + \eta(\delta))$$

Taking  $\xi$  small enough and kicking back  $t_{\delta}$ , we obtain  $t_{\delta} \leq \xi(\|u - u^{\delta}\|_X + \eta(\delta))$ , and similarly

(6.37) 
$$t_{\delta} \lesssim \xi(\|u - \tilde{u}^{\delta}\|_X + \eta(\delta)),$$

(6.38) 
$$t_{\delta} \lesssim \xi(\|^{\delta}u - u^{\delta}\|_{X} + \eta(\delta)).$$

When passing the until-clause, furthermore we have

$$\begin{aligned} \|u - \tilde{u}^{\delta}\|_{X} &\lesssim t_{\delta} + \|u - u^{\delta}\|_{X} \lesssim t_{\delta} + \|^{\delta}u - u^{\delta}\|_{X} + \eta(\delta) \approx t_{\delta} + \|\mathbf{r}^{\delta}\| + \eta(\delta) \\ \end{aligned}$$

$$(6.39) \qquad \leq t_{\delta} + \|\mathbf{r}^{\delta} - \mathbf{\tilde{r}}^{\delta}\| + \|\mathbf{\tilde{r}}^{\delta}\| + \eta(\delta) \lesssim t_{\delta} + \|u^{\delta} - \tilde{u}^{\delta}\|_{X} + \|\mathbf{\tilde{r}}^{\delta}\| + \eta(\delta) \\ \lesssim t_{\delta} + \|\mathbf{\tilde{r}}^{\delta}\| + \eta(\delta) \overset{(6.36)}{\lesssim} \|\mathbf{\tilde{r}}^{\delta}\| + \eta(\delta). \end{aligned}$$

Denoting  $\delta$  at the subsequent passing of the until-clause as  $\tilde{\delta}$ , we have

$$\begin{split} \|u - \tilde{u}^{\tilde{\delta}}\|_{X} &\lesssim t_{\tilde{\delta}} + \|^{\tilde{\delta}}u - u^{\tilde{\delta}}\|_{X} + \eta(\tilde{\delta}) \\ &\stackrel{(6.20)}{\leq} t_{\tilde{\delta}} + \frac{\kappa_{\Delta}}{\gamma_{\Delta}} \|^{\tilde{\delta}}u - ^{\tilde{\delta}}u_{\tilde{\delta}}\|_{X} + \eta(\tilde{\delta}) \\ &\stackrel{\delta \leq \tilde{\delta}}{\leq} t_{\tilde{\delta}} + \frac{\kappa_{\Delta}}{\gamma_{\Delta}} \|^{\tilde{\delta}}u - \tilde{u}^{\delta}\|_{X} + \eta(\tilde{\delta}) \\ &\leq t_{\tilde{\delta}} + \frac{\kappa_{\Delta}}{\gamma_{\Delta}} \|u - \tilde{u}^{\delta}\|_{X} + \mathcal{O}(\eta(\tilde{\delta})). \end{split}$$

By using  $t_{\tilde{\delta}} \lesssim \xi(\|u - \tilde{u}^{\tilde{\delta}}\|_X + \eta(\tilde{\delta}))$  ((6.37)) and kicking back  $\|u - \tilde{u}^{\tilde{\delta}}\|_X$ , we infer that for  $\xi$  sufficiently small,

(6.40) 
$$\|u - \tilde{u}^{\tilde{\delta}}\|_X \le \left(\frac{\kappa_{\Delta}}{\gamma_{\Delta}} + \mathcal{O}(\xi)\right)\|u - \tilde{u}^{\delta}\|_X + \mathcal{O}(\eta(\tilde{\delta}))$$

In the case that

(6.41) 
$$\eta(\delta) > \omega \| \tilde{\mathbf{r}}^{\delta} \|_{2}$$

it holds that

$$\eta(\delta) \le \eta(\delta)/2$$

and so, thanks to (6.39) and (6.41),

$$(6.42) ||u - \tilde{u}^{\delta}||_X \lesssim \eta(\delta)/\omega.$$

For any constant  $\rho_1 < 1$ , using (6.40) and (6.42) we have

$$\begin{aligned} \|u - \tilde{u}^{\tilde{\delta}}\|_{X} &\leq \rho_{1} \|u - \tilde{u}^{\delta}\|_{X} + \frac{\kappa_{\Delta}/\gamma_{\Delta} + \mathcal{O}(\xi) - \rho_{1}}{\omega} \|u - \tilde{u}^{\delta}\|_{X} + \mathcal{O}(\eta(\tilde{\delta})) \\ &\leq \rho_{1} \|u - \tilde{u}^{\delta}\|_{X} + \left(\frac{\kappa_{\Delta}/\gamma_{\Delta} + \mathcal{O}(\xi) - \rho_{1}}{\omega} + 1\right) C \eta(\delta), \end{aligned}$$

for some constant C > 0, and so

$$\vartheta \| u - \tilde{u}^{\tilde{\delta}} \|_{X} + \eta(\tilde{\delta}) \le \rho_1 \vartheta \| u - \tilde{u}^{\delta} \|_{X} + \left( \vartheta \left( \frac{\kappa_{\Delta} / \gamma_{\Delta} - \rho_1 + \mathcal{O}(\xi)}{\omega} + 1 \right) C + \frac{1}{2} \right) \eta(\delta)$$

Now let  $\vartheta > 0$  be such that  $2\vartheta C + \frac{1}{2} < 1$ . Given a constant  $\omega$  (which later will be selected such that  $\omega/\vartheta$  is sufficiently small), let  $(\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1)/\omega$ ,  $(1 - \rho_1)/\omega$ ,  $\xi/\omega$  be sufficiently small such that the expression  $\frac{\kappa_{\Delta}/\gamma_{\Delta} - \rho_1 + \mathcal{O}(\xi)}{\omega} \leq 1$ . We conclude that in this case  $\vartheta \| u - \tilde{u}^{\delta} \|_X + \eta(\delta)$  is reduced by at least a factor  $\max(\rho_1, 2\vartheta C + \frac{1}{2}) < 1$ .

Next we consider the other case that

(6.43) 
$$\eta(\delta) \le \omega \| \tilde{\mathbf{r}}^{\delta} \|,$$

so that

(6.44) 
$$\|\tilde{\mathbf{r}}^{\delta}|_{J}\| \geq \vartheta \|\tilde{\mathbf{r}}^{\delta}\|$$

We have

(6.45) 
$$\|\mathbf{r}^{\delta} - \tilde{\mathbf{r}}^{\delta}\| \lesssim \|u^{\delta} - \tilde{u}^{\delta}\|_{X} \lesssim t_{\delta}^{(6.36)} \leq \xi(\|\tilde{\mathbf{r}}^{\delta}\| + \eta(\delta)) \lesssim \xi \|\tilde{\mathbf{r}}^{\delta}\| \\ \leq \xi(\|\mathbf{r}^{\delta}\| + \|\mathbf{r}^{\delta} - \tilde{\mathbf{r}}^{\delta}\|),$$

and so by taking  $\xi$  sufficiently small and kicking back  $\|\mathbf{r}^{\delta} - \tilde{\mathbf{r}}^{\delta}\|$ , also

(6.46) 
$$\|\mathbf{r}^{\delta} - \tilde{\mathbf{r}}^{\delta}\| \lesssim \xi \|\mathbf{r}^{\delta}\|,$$

which together with (6.43) implies that

(6.47) 
$$\eta(\delta) \lesssim \omega \|\mathbf{r}^{\delta}\|$$

From (6.44)-(6.46) we infer that

$$\|\mathbf{r}^{\delta}\| \stackrel{(6.45)}{\lesssim} \|\tilde{\mathbf{r}}^{\delta}\| \stackrel{(6.44)}{\lesssim} \|\tilde{\mathbf{r}}^{\delta}|_{J}\| \leq \|\mathbf{r}^{\delta}|_{J}\| + \|\mathbf{r}^{\delta} - \tilde{\mathbf{r}}^{\delta}\| \stackrel{(6.46)}{\lesssim} \|\mathbf{r}^{\delta}|_{J}\| + \xi \|\mathbf{r}^{\delta}\|.$$

By taking  $\xi$  sufficiently small and kicking back  $\|\mathbf{r}^{\delta}\|$ , we conclude that there exists a constant  $\tilde{\vartheta} > 0$  such that

$$\|\mathbf{r}^{\delta}|_{J}\| \geq \tilde{\vartheta} \|\mathbf{r}^{\delta}\|.$$

Assuming that  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1$  and  $\omega$  are small enough, using (6.47) an application of Proposition 6.4.6 shows that there exists a constant  $\rho_3 < 1$  such that

(6.48) 
$$\|\tilde{\delta}u - u^{\tilde{\delta}}\|_X \le \rho_3 \|^{\delta}u - u^{\delta}\|_X.$$

Furthermore we have

$$\|\tilde{\mathbf{r}}^{\delta}\| \leq \|\mathbf{r}^{\delta}\| + \|\tilde{\mathbf{r}}^{\delta} - \mathbf{r}^{\delta}\| \lesssim \|^{\delta}u - u^{\delta}\|_{X} + t_{\delta} \lesssim \|u - \tilde{u}^{\delta}\|_{X} + \eta(\delta) + t_{\delta}$$

$$\stackrel{(6.37)}{\lesssim} \|u - \tilde{u}^{\delta}\|_{X} + \eta(\delta),$$

and so from (6.43) by kicking back  $\eta(\delta)$  and taking  $\omega$  sufficiently small,

(6.49) 
$$\eta(\delta) \lesssim \omega \|u - \tilde{u}^{\delta}\|_X.$$

#### We conclude that

$$\begin{split} \vartheta \| u - \tilde{u}^{\tilde{\delta}} \|_{X} + \eta(\tilde{\delta}) &\leq \vartheta \|^{\tilde{\delta}} u - u^{\tilde{\delta}} \|_{X} + \mathcal{O}(\eta(\tilde{\delta}) + t_{\tilde{\delta}}) \\ & \stackrel{(6.38)}{\leq} (\vartheta + \mathcal{O}(\xi)) \|^{\tilde{\delta}} u - u^{\tilde{\delta}} \|_{X} + \mathcal{O}(\eta(\tilde{\delta})) \\ & \stackrel{(6.48)}{\leq} (\vartheta + \mathcal{O}(\xi)) \rho_{3} \|^{\delta} u - u^{\delta} \|_{X} + \mathcal{O}(\eta(\tilde{\delta})) \\ &\leq (\vartheta + \mathcal{O}(\xi)) \rho_{3} \big[ \| u - \tilde{u}^{\delta} \|_{X} + \mathcal{O}(\eta(\delta)) + t_{\delta} \big] + \mathcal{O}(\eta(\tilde{\delta})) \\ & \stackrel{(6.37)}{\leq} (\vartheta + \mathcal{O}(\xi)) \rho_{3} \big[ (1 + \mathcal{O}(\xi)) \| u - \tilde{u}^{\delta} \|_{X} + \mathcal{O}(\eta(\delta)) \big] + \mathcal{O}(\eta(\delta)) \\ & \stackrel{(6.49)}{\leq} \big[ (\vartheta + \mathcal{O}(\xi)) \rho_{3} (1 + \mathcal{O}(\xi)) + \mathcal{O}(\omega) \big] \| u - \tilde{u}^{\delta} \|_{X} \\ &= \big[ \rho_{3} + \mathcal{O}(\frac{\omega + \xi}{\vartheta}) \big] \vartheta \| u - \tilde{u}^{\delta} \|_{X}. \end{split}$$

So for  $\omega/\vartheta$  and  $\xi/\vartheta$  sufficiently small, also in this case we established a reduction of  $\vartheta \| u - \tilde{u}^{\tilde{\delta}} \|_X + \eta(\tilde{\delta})$  by at least a constant factor less than 1.

What is left to show is that the algorithm terminates. We have shown that the value of  $\vartheta \| u - \tilde{u}^{\delta} \|_X + \eta(\delta)$  at passing the until-clause is *r*-linearly converging. We consider the corresponding value of  $e_{\delta} = \| \mathbf{\tilde{r}}^{\delta} \| + \eta(\delta) + t_{\delta}$ . Arguments that we have used multiple times show that  $\| \mathbf{\tilde{r}}^{\delta} \| \leq \| u - \tilde{u}^{\delta} \|_X + \eta(\delta) + t_{\delta}$ , and so  $e_{\delta} \leq \| u - \tilde{u}^{\delta} \|_X + \eta(\delta) + t_{\delta}$ . Using that  $t_{\delta} \leq \xi(\| \mathbf{\tilde{r}}^{\delta} \| + \eta(\delta)) \leq \xi e_{\delta}$ , for  $\xi$  sufficiently small kicking back  $e_{\delta}$  shows that

$$e_{\delta} \lesssim \|u - \tilde{u}^{\delta}\|_X + \eta(\delta),$$

and so

$$e_{\delta} \lesssim \vartheta \| u - \tilde{u}^{\delta} \|_X + \eta(\delta).$$

This last statement implies that at some moment  $e_{\delta} \leq \varepsilon$ , meaning that the algorithm stops.

# 6.5 Wavelets-in-time tensorized with finite-elements-in-space

We specify the parabolic problem at hand, as well as the type of families  $(X^{\delta})_{\delta \in \Delta}$  and  $(Y^{\delta})_{\delta \in \Delta}$  of 'trial' and 'test' spaces. A likely harmless minor further restriction to these families that will be needed for the construction of an *X*-stable collection  $\Theta_{\delta}$  that spans an *X*-stable complement space of  $X^{\delta}$  in  $X^{\delta}$ , specifically condition (6.24), will be postponed to Sect. 6.6.1.

#### 6.5.1 Continuous problem

For some bounded domain  $\Omega \subset \mathbb{R}^d$ , we take  $H = L_2(\Omega)$  and, for some closed  $\emptyset \subseteq \Gamma_D \subseteq \partial\Omega$ ,  $V = H^1_{0,\Gamma_D}(\Omega) := \operatorname{clos}_{H^1(\Omega)} \{ u \in C^{\infty}(\Omega) \cap H^1(\Omega) : u|_{\Gamma_D} = 0 \}$ ,

and

$$a(t;\eta,\zeta) := \int_{\Omega} \mathbf{K} \nabla \eta \cdot \nabla \zeta + c \eta \zeta \, d\mathbf{x},$$

where  $\mathbf{K} = \mathbf{K}^{\top} \in L_{\infty}(I \times \Omega)$  with  $\mathbf{K}(\cdot) = \operatorname{Id} \text{ a.e., and } c \in L_{\infty}(I \times \Omega)$ . W.l.o.g. we take T = 1.

#### 8

### 6.5.2 Wavelets in time

We will construct the trial and test spaces as the span of wavelets-in-time tensorized with finite element spaces-in-space. In this subsection we collect some assumptions on the wavelets.

At the 'trial side' we consider a countable collection  $\Sigma = \{\sigma_{\lambda} : \lambda \in \vee_{\Sigma}\}$ of functions  $I \to \mathbb{R}$  known as *wavelets*. To each  $\lambda \in \vee_{\Sigma}$  we associate a value  $|\lambda| \in \mathbb{N}_0$ , called the *level* of  $\lambda$ . We assume that the wavelets are *locally supported* meaning that  $\sup_{n \in \mathbb{N}, \ell \in \mathbb{N}_0} \#\{\lambda \in \vee_{\Sigma} : |\lambda| = \ell, |\sup \sigma_{\lambda} \cap 2^{-\ell}(n + [0, 1])| > 0\} < \infty$  and  $\operatorname{diam } \operatorname{supp} \sigma_{\lambda} \lesssim 2^{-|\lambda|}$ . To each  $\lambda \in \vee_{\Sigma}$  with  $|\lambda| > 0$ , we associate one or more  $\tilde{\lambda} \in \vee_{\Sigma}$  with  $|\tilde{\lambda}| = |\lambda| - 1$  and  $|\operatorname{supp} \sigma_{\lambda} \cap \operatorname{supp} \sigma_{\tilde{\lambda}}| > 0$  which we call the *parent*(*s*) of  $\lambda$ . We denote this relation between a parent  $\tilde{\lambda}$  and a child  $\lambda$  by

 $\tilde{\lambda} \triangleleft_{\Sigma} \lambda.$ 

The definitions of parents and children give rise to obvious notions of ancestors and descendants.

To each  $\lambda \in \vee_{\Sigma}$  we associate some neighbourhood  $S_{\Sigma}(\lambda)$  of  $\operatorname{supp} \sigma_{\lambda}$  with  $\operatorname{diam} S_{\Sigma}(\lambda) \leq 2^{-|\lambda|}$  and

$$\tilde{\lambda} \triangleleft_{\Sigma} \lambda \Longrightarrow \mathcal{S}_{\Sigma}(\tilde{\lambda}) \supseteq \mathcal{S}_{\Sigma}(\lambda).$$

For some wavelets bases, e.g. Alpert wavelets ([Alp93]), it suffices to take  $S_{\Sigma}(\lambda) = \operatorname{supp} \sigma_{\lambda}$ . With  $C_{\Sigma} := \operatorname{sup}_{\lambda \in \vee_{\Sigma}} 2^{|\lambda|} \operatorname{diam} \operatorname{supp} \sigma_{\lambda}$ , a neighbourhood that in any case is sufficiently large is  $\{t \in I : \operatorname{dist}(t, \operatorname{supp} \sigma_{\lambda}) \leq C_{\Sigma} 2^{-|\lambda|}\}$ . Indeed, if with this definition  $t \in S_{\Sigma}(\lambda)$  and  $\tilde{\lambda} \triangleleft_{\Sigma} \lambda$ , then  $\operatorname{dist}(t, \operatorname{supp} \sigma_{\tilde{\lambda}}) \leq \operatorname{dist}(t, \operatorname{supp} \sigma_{\lambda}) + \operatorname{diam}(\operatorname{supp} \sigma_{\lambda}) \leq 2C_{\Sigma} 2^{-|\lambda|}$ , i.e.  $t \in S_{\Sigma}(\tilde{\lambda})$ .

We assume that  $\Sigma$  is a Riesz basis for  $L_2(I)$ , and, when renormalized in  $H^1(I)$ -norm, it is a Riesz basis for  $H^1(I)$ . Although not essential, thinking of wavelets being (essentially) constructed by means of *dilation*, we assume that

$$\|\sigma_{\lambda}\|_{H^{1}(I)} \approx 2^{|\lambda|}.$$

At the 'test side' we consider a similar collection  $\Psi = \{\psi_{\mu} : \mu \in \forall_{\Psi}\}$  of wavelets, with the difference though that this one has to be an even *orthonormal* basis for  $L_2(\Omega)$ , whilst, renormalized in  $H^1(I)$ -norm, it does not need to be a Riesz basis for  $H^1(I)$ .

We will assume that  $\Sigma$  and  $\Psi$  are selected such that for any  $\ell \in N_0$ ,

$$\operatorname{span}\{\sigma_{\lambda} \colon |\lambda| \le \ell\} \cup \operatorname{span}\{\sigma_{\lambda}' \colon |\lambda| \le \ell\} \subseteq \operatorname{span}\{\psi_{\mu} \colon |\mu| \le \ell\},$$

so that in particular

(6.50) 
$$|\mu| > |\lambda| \Longrightarrow \langle \sigma_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)} = 0 = \langle \sigma'_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)}.$$

#### 6.5.3 Uniform stability

In the following proposition, we further specify the type of families of trial and test spaces that we consider, and formulate sufficient conditions for the requirements (6.13)-(6.14), which implied uniform stability of the Galerkin discretizations of our saddle-point problem (6.9).

**Proposition 6.5.1.** *For*  $\delta \in \Delta$ *, let* 

$$X^{\delta} = \sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes W^{\delta}_{\lambda}, \quad Y^{\delta} = \sum_{\mu \in \vee_{\Psi}} \psi_{\mu} \otimes V^{\delta}_{\mu}$$

for subspaces  $W^{\delta}_{\lambda}, V^{\delta}_{\mu} \subseteq V$  of which finitely many are non-zero. Let

(6.51)  $\langle \sigma_{\lambda}, \psi_{\mu} \rangle_{L_2(I)} \neq 0 \Longrightarrow V_{\mu}^{\delta} \supseteq W_{\lambda}^{\delta},$ 

and, for some constant  $\gamma_{\Delta} > 0$ , for any  $\mu \in \vee_{\Psi}$ ,

(6.52) 
$$\inf_{\substack{0\neq w\in \sum : \langle \sigma'_{\lambda},\psi_{\mu}\rangle_{L_{2}(I)}\neq 0\}}} \sup_{\substack{\delta \in v \in V_{\mu}^{\delta}}} \frac{w(v)}{\|w\|_{V'}\|v\|_{V}} \geq \gamma_{\Delta}.$$

*Then*  $X^{\delta} \subseteq Y^{\delta}$  *and* 

$$\inf_{0 \neq w \in X^{\delta}} \sup_{0 \neq v \in Y^{\delta}} \frac{(\partial_t w)(v)}{\|\partial_t w\|_{Y'} \|v\|_Y} \ge \gamma_{\Delta},$$

i.e., the conditions (6.13)-(6.14) for uniform stability are satisfied.

*Proof.* For  $w_{\lambda} \in W_{\lambda}^{\delta}$  and  $w := \sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes w_{\lambda} \in X^{\delta}$ ,  $\sigma_{\lambda} = \sum_{\mu \in \vee_{\Psi}} \langle \sigma_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)} \psi_{\mu}$ shows that  $w = \sum_{\mu \in \vee_{\Psi}} \psi_{\mu} \otimes \sum_{\lambda \in \vee_{\Sigma}} \langle \sigma_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)} w_{\lambda} \in Y^{\delta}$  by the first assumption. Similarly  $\partial_{t}w = \sum_{\mu \in \vee_{\Psi}} \psi_{\mu} \otimes \tilde{v}_{\mu}$  where  $\tilde{v}_{\mu} := \sum_{\lambda \in \vee_{\Sigma}} \langle \sigma'_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)} w_{\lambda}$ . For any  $\varepsilon > 0$ , the second assumption shows that for any  $\mu \in \vee_{\Psi}$ , there exists a  $v_{\mu} \in V_{\mu}^{\delta}$  with  $\tilde{v}_{\mu}(v_{\mu}) \geq (\gamma_{\Delta} - \varepsilon) \|\tilde{v}_{\mu}\|_{V'} \|v_{\mu}\|_{V}$  and  $\|\tilde{v}_{\mu}\|_{V'} = \|v_{\mu}\|_{V}$ . With  $v := \sum_{\mu \in \vee_{\Psi}} \psi_{\mu} \otimes v_{\mu} \in Y^{\delta}$ , we infer that  $(\partial_{t}w)(v) = \sum_{\mu \in \vee_{\Psi}} \tilde{v}_{\mu}(v_{\mu}) \geq (\gamma_{\Delta} - \varepsilon) \|\partial_{t}w\|_{Y'} \|v\|_{Y}$ . In order to be able to apply at linear cost the arising linear operators in (6.31)-(6.32), we will restrict the type of trial spaces  $X^{\delta} = \sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes W_{\lambda}^{\delta}$  by imposing the following *tree condition* 

(6.53) 
$$\tilde{\lambda} \triangleleft_{\Sigma} \lambda \Longrightarrow W^{\delta}_{\tilde{\lambda}} \supseteq W^{\delta}_{\lambda}.$$

For the same reason the analogous condition will be needed for  $Y^{\delta}$ . For  $X^{\delta}$  that satisfies (6.53), below the latter will be verified, and sufficient, more easily verifiable conditions for (6.51)-(6.52) are derived.

**Proposition 6.5.2.** Let  $X^{\delta} = \sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes W^{\delta}_{\lambda}$  satisfy (6.53). For  $\mu \in \vee_{\Psi}$ , set

$$\hat{W}^{\delta}_{\mu} := \sum_{\{\lambda \in \vee_{\Sigma} \colon |\lambda| = |\mu|, \, |\mathcal{S}_{\Psi}(\mu) \cap \mathcal{S}_{\Sigma}(\lambda)| > 0\}} W^{\delta}_{\lambda}.$$

Build  $Y^{\delta} = \sum_{\mu \in \lor_{\Psi}} \psi_{\mu} \otimes V^{\delta}_{\mu}$  by taking  $V^{\delta}_{\mu} = \{0\}$  when  $\hat{W}^{\delta}_{\mu} = \{0\}$ , and otherwise

$$V^{\delta}_{\mu} \supseteq \hat{W}^{\delta}_{\mu}$$

where

(6.54) 
$$\inf_{0 \neq w \in \hat{W}^{\delta}_{\mu}} \sup_{0 \neq v \in V^{\delta}_{\mu}} \frac{w(v)}{\|w\|_{V'}} \ge \gamma_{\Delta}$$

for some constant  $\gamma_{\Delta} > 0$ . Then the conditions (6.51) and (6.52) from Proposition 6.5.1 for uniform stability are satisfied.

When  $\dim V_{\mu}^{\delta} \lesssim \dim \hat{W}_{\mu}^{\delta}$ , then  $\dim Y^{\delta} \lesssim \dim X^{\delta}$ , and under the natural condition that a larger  $\hat{W}_{\mu}^{\delta}$  gives rise to a larger (more precisely, not smaller)  $V_{\mu}^{\delta}$ , the constructed  $Y^{\delta}$  satisfies the tree condition

(6.55) 
$$\tilde{\mu} \triangleleft_{\Psi} \mu \Longrightarrow V_{\tilde{\mu}} \supseteq V_{\mu}$$

*Proof.* Let  $\langle \sigma_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)} \neq 0$  or  $\langle \sigma'_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)} \neq 0$ . Then  $|\mathcal{S}_{\Sigma}(\lambda) \cap \mathcal{S}_{\Psi}(\mu)| > 0$ and  $|\lambda| \geq |\mu|$  by (6.50). When  $|\lambda| > |\mu|$ ,  $\lambda$  has an ancestor  $\tilde{\lambda}$  with  $|\tilde{\lambda}| = |\mu|$ ,  $W^{\delta}_{\tilde{\lambda}} \supseteq W^{\delta}_{\lambda'}$ , and  $\mathcal{S}_{\Sigma}(\tilde{\lambda}) \supseteq \mathcal{S}_{\Sigma}(\lambda)$ , and thus  $|\mathcal{S}_{\Sigma}(\tilde{\lambda}) \cap \mathcal{S}_{\Psi}(\mu)| > 0$ . We conclude that both  $\sum_{\{\lambda \in \vee_{\Sigma} : \langle \sigma_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)} \neq 0\}} W^{\delta}_{\lambda}$  and  $\sum_{\{\lambda \in \vee_{\Sigma} : \langle \sigma'_{\lambda}, \psi_{\mu} \rangle_{L_{2}(I)} \neq 0\}} W^{\delta}_{\lambda}$  are included in  $\hat{W}^{\delta}_{\mu}$ , so that (6.51) and (6.52) are guaranteed by the selection of  $V^{\delta}_{\mu}$ .

The statement dim  $Y^{\delta} \leq \dim X^{\delta}$  when dim  $V_{\mu}^{\delta} \leq \dim \hat{W}_{\mu}^{\delta}$  follows from  $\dim \hat{W}_{\mu}^{\delta} \leq \sum_{\{\lambda \in \vee_{\Sigma} : |\lambda| = |\mu|, |S_{\Psi}(\mu) \cap S_{\Sigma}(\lambda)| > 0\}} \dim W_{\lambda}^{\delta}$ , and the fact that for any  $\lambda \in \vee_{\Sigma}$ , the number of  $\mu \in \vee_{\Psi}$  with  $|\mu| = |\lambda|$  and  $|S_{\Psi}(\mu) \cap S_{\Sigma}(\lambda)| > 0$  is uniformly bounded.

Let  $\tilde{\mu} \triangleleft_{\Psi} \mu$ , and so  $S_{\Psi}(\tilde{\mu}) \supseteq S_{\Psi}(\mu)$ . For each  $\lambda \in \bigvee_{\Sigma}$  with  $|\lambda| = |\mu|$  and  $|S_{\Psi}(\mu) \cap S_{\Sigma}(\lambda)| > 0$ , there exists a  $\tilde{\lambda} \triangleleft_{\Sigma} \lambda$ , thus with  $S_{\Sigma}(\tilde{\lambda}) \supseteq S_{\Sigma}(\lambda)$ , and  $W_{\tilde{\lambda}}^{\delta} \supseteq W_{\lambda}^{\delta}$  by (6.53). We conclude that  $\hat{W}_{\mu}^{\delta} \supseteq \hat{W}_{\mu}^{\delta}$ , which completes the proof of (6.55).

As follows from [DSW21, Thm. 3.10] (taking *B* to be the Riesz map  $H \rightarrow H'$ ) condition (6.54) has the following equivalent formulation.

**Proposition 6.5.3.** Condition (6.54) is equivalent to existence of a projector  $Q \in \mathcal{L}(V, V)$  with ran  $Q \subseteq V_{\mu}^{\delta}$ , ran  $Q^* \supseteq \hat{W}_{\lambda}^{\delta}$ , and  $\|Q\|_{\mathcal{L}(V, V)} \leq \frac{1}{\gamma_{\Delta}}$ .

# 6.5.4 Selection of the spatial approximation spaces as finite element spaces

We will select the spaces  $W_{\lambda}^{\delta}$  from a collection  $\mathcal{O}$  of finite element spaces in V, which collection is closed under taking (finite) sums, and for which

(6.56) 
$$\inf_{W \in \mathcal{O}} \inf_{0 \neq w \in W} \sup_{0 \neq v \in W} \frac{w(v)}{\|w\|_{V'} \|v\|_{V}} > 0.$$

Consequently, the stability conditions (6.13)-(6.14) are satisfied for some  $\gamma_{\Delta} > 0$  by simply taking in Proposition 6.5.2.

$$(6.57) V_{\mu}^{\delta} := \hat{W}_{\mu}^{\delta} \in \mathcal{O}.$$

As follows from Proposition 6.5.3, (6.56) is equivalent to uniform boundedness w.r.t. the norm on *V* of the *H*-orthogonal projector onto  $W \in \mathcal{O}$ . It is well-known that an example of such a collection  $\mathcal{O}$  is given by the set of all finite element spaces  $W_{\lambda}^{\delta}$  w.r.t. quasi-uniform, uniformly shape regular conforming partitions of  $\Omega$  into, say, *d*-simplices.

It is known that the uniform boundedness w.r.t. the *V*-norm of the *H*-orthogonal projector holds also true for finite element spaces w.r.t. *locally re-fined* partitions as long as the grading of the partitions is sufficiently mild. In [GHS16] it has been shown that for d = 2 spatial dimensions, and polynomial orders up to 12, the collection of all conforming partitions that can be generated by *newest vertex bisection* (NVB), starting from a fixed conforming initial partition  $\mathcal{T}_{\perp}$  with an assignment of the newest vertices that satisfies a so-called matching condition, is sufficiently mildly graded in the above sense. Since the overlay of two conforming NVB partitions is a conforming NVB partition, this collection is closed under taking (finite) sums. In other words, with this collection of finite element spaces, which we will employ in our experiments, again the choice (6.57) guarantees uniform stability.

In [Car04] a result similar to that from [GHS16] has been shown for redblue-green refinement and lowest order finite elements again for d = 2. Unfortunately, for d > 2 such results seem not yet to be available.

*Remark* 6.5.4 (Getting  $\gamma_{\Delta}$  close to 1). We discussed uniform boundedness w.r.t. the *V*-norm of the *H*-orthogonal projectors onto a family of finite element spaces, which, by taking  $V_{\mu}^{\delta} := \hat{W}_{\mu}^{\delta}$  in Proposition 6.5.2, yields the uniform inf-sup condition (6.14) *some* value  $\gamma_{\Delta} > 0$ , and so uniform stability of the Galerkin discretizations of the saddle-point (6.9).

For proving convergence of our adaptive routine Algorithm 6.4.8, however, we needed a value of  $\gamma_{\Delta} > 0$  that is sufficiently close to 1. Although in our numerical experiments, reported on in Sect. 6.7, with continuous piecewise linear finite element spaces generated by conforming NVB and  $V_{\mu}^{\delta} = \hat{W}_{\mu}^{\delta}$ , the adaptive routine is *r*-linearly converging, there is no guarantee that  $1 - \gamma_{\Delta}$  is sufficiently small.

Restricting to quasi-uniform partitions, below we show that  $1 - \gamma_{\Delta}$  can be made arbitrarily small by taking the mesh underlying  $V_{\mu}^{\delta}$  to be a sufficiently deep, but fixed *refinement* of the mesh underlying  $\hat{W}_{\mu}^{\delta}$ . One may conjecture that the same result holds true for sufficiently mildly graded locally refined meshes.

Let the diameters of any *d*-simplex in the partitions underlying  $\hat{W}^{\delta}_{\mu}$  and  $V^{\delta}_{\mu}$  be proportional to  $h_c$  and  $h_f$ , respectively. For  $s \in [0, 1]$ , let  $\mathcal{H}^s := [H, V]_{s,2}$ . In any case when  $\Omega$  is a Lipschitz domain, it is known that there exists an  $s \in (0, 1]$  such that the solution  $u \in V$  of  $\langle u, v \rangle_V = f(v)$  ( $v \in V$ ) satisfies  $\|u\|_{H^{1+s}(\Omega)} \lesssim \|f\|_{(\mathcal{H}^{1-s})'}$ , assuming the right-hand side is bounded. From this, the Aubin-Nitsche duality argument shows that the *V*-orthogonal projector  $P_{\mu}$  onto  $V^{\delta}_{\mu}$  satisfies  $\|\mathrm{Id} - P_{\mu}\|_{\mathcal{L}(V,\mathcal{H}^{1-s})} \lesssim h^s_f$ . On the other hand, on  $\hat{W}^{\delta}_{\mu}$  we have the following inverse inequality  $\|\cdot\|_{(\mathcal{H}^{1-s})'} \lesssim h^{-s}_c \|\cdot\|_{V'}$  (e.g. (2.44)).

Given  $w \in \hat{W}_{\mu}^{\delta}$ , for any  $\varepsilon > 0$ , there exists a  $v \in V$  with  $w(v) \ge (1 - \varepsilon) ||w||_{V'} ||v||_V$ . We infer that, for some constant C > 0,

$$w(P_{\mu}v) = w(v) + w((\mathrm{Id} - P_{\mu})v)$$
  

$$\geq (1 - \varepsilon) ||w||_{V'} ||v||_{V} - ||w||_{(\mathcal{H}^{1-s})'} ||(\mathrm{Id} - P_{\mu})v||_{\mathcal{H}^{1-s}}$$
  

$$\geq (1 - (\varepsilon + C(h_{f}/h_{c})^{s}) ||w||_{V'} ||v||_{V}$$
  

$$\geq (1 - (\varepsilon + C(h_{f}/h_{c})^{s}) ||w||_{V'} ||P_{\mu}v||_{V}.$$

Since  $\varepsilon > 0$  was arbitrary, we conclude that  $\gamma_{\Delta} \ge (1 - C(h_f/h_c)^s)$  which proves our assertion.

#### 6.5.5 Best possible rates

Although so far we have not proved it, we expect that the sequence of approximations generated by our adaptive Algorithm 6.4.8 is not only *r*-linearly converging, but, ignoring data oscillation, that it is a sequence of approximations from a sequence of spaces from the family  $(X^{\delta})_{\delta \in \Delta}$  that converges with the best possible rate. In this subsection, we show that with our selection of the  $(X^{\delta})_{\delta \in \Delta}$ , under some (mild) smoothness conditions on the solution *u* this best possible rate *equals* the rate of best approximation to the solution of the corresponding stationary problem from the spatial finite element spaces w.r.t. the *V*-norm.

Consider a family of spaces  $X^{\delta} = \sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes W^{\delta}_{\lambda}$  that satisfies (6.53), with the  $W^{\delta}_{\lambda}$  selected from a collection of finite element spaces  $\mathcal{O}$  that in any

case contains all such spaces that correspond to uniform refinements of some initial partition of  $\Omega$ . Let  $\Sigma$  a collection of wavelets of order  $d_t$ , and assume that the finite element spaces are of order  $d_x$ . When for each  $X^{\delta}$ , the space  $Y^{\delta}$  is selected as in Proposition 6.5.2, then the combination of (6.20) and the analysis from [SS09, Sect. 7.1] shows that if the exact solution u of our parabolic problem satisfies the mixed regularity condition  $u \in H^{d_t}(I) \otimes H^{d_x}(\Omega)$ , then a suitable (non-adaptive) choice of the spaces  $W^{\delta}_{\lambda}$  yields a sequence of solutions  $u^{\hat{\delta}\delta} \in X^{\delta}$  (for arbitrary  $Y^{\hat{\delta}} \supset Y^{\delta}$ ) of the modified discretized saddle-point from Sect 6.3.3, for which

$$\|u - u^{\hat{\delta}\delta}\|_X \lesssim (\dim X^{\delta})^{-\min(d_t - 1, \frac{d_x - 1}{d})}.$$

Note that for  $d_t - 1 \ge \frac{d_x - 1}{d}$ , the rate  $\frac{d_x - 1}{d}$  equals the best rate in the *V*-norm that can be expected when the finite element spaces are employed for solving the corresponding *stationary* problem, which is posed on a *d*-dimensional domain instead over the d + 1-dimensional space-time cylinder.

For an optimal *adaptive* choice of the  $W_{\lambda}^{\delta}$  as finite element spaces w.r.t. a sufficiently 'rich' collection of locally refined partitions, as the collection of all conforming NVB partitions, it can be expected that the rate  $\min(d_t - 1, \frac{d_x - 1}{d})$  is realized under *much* milder regularity conditions on *u*.

When instead of being finite element spaces, the spaces  $W^{\delta}_{\lambda}$  can be selected as the spans of some wavelets from a Riesz basis for V of order  $d_x$ , and additionally the tree condition (6.53) is dropped, a precise characterization of those u that can be approximated at a rate  $s < \min(d_t - 1, \frac{d_x - 1}{d})$  in terms of tensor products of Besov spaces can be deduced from [Nit06, SU09]. The collection of finite element spaces w.r.t. locally refined meshes, as those generated by NVB, is very resemblant to the collection of spans of sets of such wavelets when on these sets a tree condition is imposed similar to the tree constraint (6.53) that we imposed in the temporal direction. In other words, the collection of spaces  $X^{\delta}$  that we consider is similar to the collection of spans of sets of tensor products of temporal and spatial wavelets when these sets satisfy a 'doubletree' constraint. In view of results from [BDDP02], we do not expect that this constraint makes the resulting approximation classes much smaller.

#### 6.5.6 Preconditioners

Our adaptive solution method of the parabolic problem requires optimal preconditioners for  $E_Y^{\delta'}AE_Y^{\delta}$  and  $S^{\delta\delta}$ , i.e., for both Z = Y and Z = X and  $\delta \in \Delta$ , we need operators  $K_Z^{\delta} = K_Z^{\delta'} \in \mathcal{L}(Z^{\delta'}, Z^{\delta})$  with  $h(K_Z^{\delta}h) \approx ||h||_{Z^{\delta'}}^2$   $(h \in Z^{\delta'})$ , moreover which should be applicable at linear cost.

To construct these preconditioners, for  $Z \in \{Y, X\}$  we will select a symmetric, bounded, and coercive bilinear form on  $Z \times Z$ , and after selecting *some* basis for  $Z^{\delta}$ , we will construct a matrix  $\mathbf{K}_{Z}^{\delta} = \mathbf{K}_{Z}^{\delta^{\top}}$  that can be applied in linear complexity, and that is *uniformly* spectrally equivalent to the inverse

of the *stiffness matrix* corresponding to this bilinear form (being the matrix representation of the linear mapping  $Z^{\delta} \to Z^{\delta'}$  defined by the bilinear form w.r.t. the basis for  $Z^{\delta}$  being chosen and the corresponding dual basis for  $Z^{\delta'}$ ). Then  $K_Z^{\delta} \in \mathcal{L}(Z^{\delta'}, Z^{\delta})$ , defined as the operator whose matrix representation is  $\mathbf{K}_Z^{\delta}$  w.r.t. the aforementioned bases of  $Z^{\delta'}$  and  $Z^{\delta}$ , is the preconditioner that satisfies our needs.

Notice that the choice of the basis for  $Z^{\delta}$  is irrelevant. Indeed, denoting the aforementioned stiffness matrix as  $\mathbf{C}_{Z}^{\delta}$  with corresponding operator  $C_{Z}^{\delta} = C_{Z}^{\delta'} \in \mathcal{L}is(Z^{\delta}, Z^{\delta'})$ , one may verify that

$$\begin{split} \|K_{Z}^{\delta}\|_{\mathcal{L}(Z^{\delta'},Z^{\delta})}\|(K_{Z}^{\delta})^{-1}\|_{\mathcal{L}(Z^{\delta},Z^{\delta'})} &\approx \|K_{Z}^{\delta}C_{Z}^{\delta}\|_{\mathcal{L}(Z^{\delta},Z^{\delta})}\|(K_{Z}^{\delta}C_{Z}^{\delta})^{-1}\|_{\mathcal{L}(Z^{\delta},Z^{\delta})} \\ &= \frac{\lambda_{\max}(\mathbf{K}_{Z}^{\delta}\mathbf{C}_{Z}^{\delta})}{\lambda_{\min}(\mathbf{K}_{Z}^{\delta}\mathbf{C}_{Z}^{\delta})}. \end{split}$$

#### Preconditioner at the 'test side'

Let Y = Z. Since  $\Psi$  is an orthonormal basis for  $L_2(I)$ , any  $y \in Y$  is of the form  $\sum_{\mu \in \vee_{\Psi}} \psi_{\mu} \otimes v_{\mu}$  for some  $v_{\mu} \in V$  with  $\sum_{\mu \in \vee_{\Psi}} ||v_{\mu}||_{V}^{2} < \infty$ . Taking as bilinear form on  $Y \times Y$  simply the scalar product on  $Y \times Y$ , we have

$$\langle \sum_{\mu_1 \in \vee_{\Psi}} \psi_{\mu_1} \otimes v_{\mu_1}^{(1)}, \sum_{\mu_2 \in \vee_{\Psi}} \psi_{\mu_2} \otimes v_{\mu_2}^{(2)} \rangle_Y = \sum_{\mu \in \vee_{\Psi}} \langle v_{\mu}^{(1)}, v_{\mu}^{(2)} \rangle_V.$$

Equipping  $Y^{\delta} = \sum_{\mu \in \vee_{\Psi}} \psi_{\mu} \otimes V_{\mu}^{\delta}$  with a basis of type  $\cup_{\mu \in \vee_{\Psi}} \psi_{\mu} \otimes \Phi_{\mu}^{\delta}$ , the resulting stiffness matrix reads as blockdiag  $[\mathbf{A}_{\mu}^{\delta}]_{\mu \in \vee_{\Psi}}$ , where  $\mathbf{A}_{\mu}^{\delta} = \langle \Phi_{\mu}^{\delta}, \Phi_{\mu}^{\delta} \rangle_{V}$  is the stiffness matrix of  $\langle \cdot, \cdot \rangle_{V}$  w.r.t.  $\Phi_{\mu}^{\delta}$ . Selecting  $\mathbf{K}_{\mu}^{\delta} \approx (\mathbf{A}_{\mu}^{\delta})^{-1}$ , the matrix representation of the optimal preconditioner reads as

$$\mathbf{K}_{Y}^{\delta} = \text{blockdiag}[\mathbf{K}_{\mu}^{\delta}]_{\mu \in \vee_{\Psi}}.$$

It is well-known that when  $V_{\mu}^{\delta}$  is a finite element space, possibly w.r.t. a locally refined partition, suitable  $\mathbf{K}_{\mu}^{\delta}$  of *multi-grid* type are available. These  $\mathbf{K}_{\mu}^{\delta}$  can be applied in linear complexity, and so can  $\mathbf{K}_{Y}^{\delta}$ .

To show, in Theorem 6.4.9, that our adaptive Algorithm 6.4.8 is *r*-linearly converging we required  $C_{\Delta} - 1$  to be sufficiently small, which requires that  $\|(E_Y^{\delta'}AE_Y^{\delta})^{-1} - K_Y^{\delta}\|_{\mathcal{L}(Y^{\delta'},Y^{\delta})}$  or, equivalently,  $\|\text{Id} - K_Y^{\delta}E_Y^{\delta'}AE_Y^{\delta}\|_{\mathcal{L}(Y^{\delta},Y^{\delta})}$  is sufficiently small, i.e. the eigenvalues of  $K_Y^{\delta}E_Y^{\delta'}AE_Y^{\delta}$  are sufficiently close to 1. Given an initial optimal, self-adjoint, and coercive preconditioner  $K_Y^{\delta}$ , and some upper and lower bounds on the spectrum of the preconditioned system, one can satisfy the latter condition by polynomial acceleration using Chebychev polynomials of sufficiently high degree. In our numerical experiments, it turned out that it was not needed to apply this 'acceleration'.

#### Preconditioner at the 'trial side'

The preconditioner presented in this section is inspired by constructions of preconditioners in [And16, NS19] for parabolic problems discretized on a tensor product of temporal and spatial spaces.

Thanks to  $\Sigma$  and  $\{2^{-|\lambda|}\sigma_{\lambda}: \lambda \in \vee_{\Sigma}\}$  being Riesz bases for  $L_2(I)$  and  $H^1(I)$ , any  $x \in X$  is of the form  $\sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes w_{\lambda}$  for some  $w_{\lambda} \in V$  with  $\sum_{\lambda \in \vee_{\Sigma}} \|w_{\lambda}\|_{V}^{2} + 4^{|\lambda|} \|w_{\lambda}\|_{V}^{2} < \infty$ , and

$$\langle \sum_{\lambda_1 \in \vee_{\Sigma}} \sigma_{\lambda_1} \otimes w_{\lambda_1}^{(1)}, \sum_{\lambda_2 \in \vee_{\Sigma}} \sigma_{\lambda_2} \otimes w_{\lambda_2}^{(2)} \rangle \coloneqq \sum_{\lambda \in \vee_{\Sigma}} \langle w_{\lambda}^{(1)}, w_{\lambda}^{(2)} \rangle_V + 4^{|\lambda|} \langle w_{\lambda}^{(1)}, w_{\lambda}^{(2)} \rangle_{V'}$$

is a symmetric, bounded, and coercive bilinear form on  $X \times X$ . Equipping  $X^{\delta} = \sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes W_{\lambda}^{\delta}$  with a basis of type  $\cup_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes \Phi_{\lambda}^{\delta}$ , the resulting stiffness matrix reads as

$$\operatorname{blockdiag}[\mathbf{A}^{\delta}_{\lambda} + 4^{|\lambda|} \langle \Phi^{\delta}_{\lambda}, \Phi^{\delta}_{\lambda} \rangle_{V'}]_{\lambda \in \vee_{\Sigma}}$$

where  $\mathbf{A}_{\lambda}^{\delta} = \langle \Phi_{\lambda}^{\delta}, \Phi_{\lambda}^{\delta} \rangle_{V}.$ 

Thanks to our assumption (6.56), for  $u \in W_{\lambda}^{\delta}$  it holds that  $||u||_{V'} \lesssim \sup_{0 \neq w \in W_{\lambda}^{\delta}} \frac{\langle u, w \rangle}{||w||_{V}} \leq ||u||_{V'}$ . With u denoting the representation of u w.r.t.  $\Phi_{\lambda}^{\delta}$ , we have

$$\sup_{0 \neq w \in W_{\lambda}^{\delta}} \frac{\langle u, w \rangle}{\|w\|_{V}} = \|(\mathbf{A}_{\lambda}^{\delta})^{-\frac{1}{2}} \mathbf{M}_{\lambda}^{\delta} \mathbf{u}\|,$$

where  $\mathbf{M}_{\lambda}^{\delta} = \langle \Phi_{\lambda}^{\delta}, \Phi_{\lambda}^{\delta} \rangle$ , so that

C

$$\langle \Phi_{\lambda}^{\delta}, \Phi_{\lambda}^{\delta} \rangle_{V'} \lesssim \mathbf{M}_{\lambda}^{\delta} (\mathbf{A}_{\lambda}^{\delta})^{-1} \mathbf{M}_{\lambda}^{\delta} \leq \langle \Phi_{\lambda}^{\delta}, \Phi_{\lambda}^{\delta} \rangle_{V'}.$$

Since both  $\mathbf{A}_{\lambda}^{\delta}$  and  $\mathbf{M}_{\lambda}^{\delta}$  are symmetric and positive definite, [PW12, Thm. 4] shows that

$$\begin{split} \frac{1}{2} \big( \mathbf{A}_{\lambda}^{\delta} + 4^{|\lambda|} \mathbf{M}_{\lambda}^{\delta} (\mathbf{A}_{\lambda}^{\delta})^{-1} \mathbf{M}_{\lambda}^{\delta} \big) &\leq (\mathbf{A}_{\lambda}^{\delta} + 2^{|\lambda|} \mathbf{M}_{\lambda}^{\delta}) (\mathbf{A}_{\lambda}^{\delta})^{-1} (\mathbf{A}_{\lambda}^{\delta} + 2^{|\lambda|} \mathbf{M}_{\lambda}^{\delta}) \\ &\leq \mathbf{A}_{\lambda}^{\delta} + 4^{|\lambda|} \mathbf{M}_{\lambda}^{\delta} (\mathbf{A}_{\lambda}^{\delta})^{-1} \mathbf{M}_{\lambda}^{\delta}. \end{split}$$

Now assuming that

(6.58) 
$$\mathbf{K}_{\lambda}^{\delta} \approx (\mathbf{A}_{\lambda}^{\delta} + 2^{|\lambda|} \mathbf{M}_{\lambda}^{\delta})^{-1},$$

we infer that

$$\mathbf{K}_{X}^{\delta} = \text{blockdiag} \left[ \mathbf{K}_{\lambda}^{\delta} \mathbf{A}_{\lambda}^{\delta} \mathbf{K}_{\lambda}^{\delta} \right]_{\lambda}$$

is the matrix representation of an optimal preconditioner.

Notice that (6.58) requires an optimal preconditioner of a discretized reactiondiffusion equation that is robust w.r.t. to the size of the (constant) reaction term. In [OR00] it was shown that, under a 'full-regularity' assumption, for quasi-uniform meshes multiplicative multi-grid yields such a preconditioner, moreover whose application can be performed at linear cost. Although we expect that using the theory of subspace correction methods the full regularity assumption can be avoided, and furthermore that the optimality, robustness and linear complexity result extends to locally refined meshes, proofs of such extensions seem not to be available.

## 6.6 A concrete realization

#### **6.6.1** The collection $\mathcal{O}$ of finite element spaces, and the map $\delta \rightarrow \underline{\delta}$

We further specify the collection  $\mathcal{O}$  of finite element spaces, construct a linearly independent set in  $H^1_{0,\Gamma_D}(\Omega)$ , known as the hierarchical basis, and equip it with a tree structure such that there exists a 1-1 correspondence between the finite element spaces in  $\mathcal{O}$ , and the spans of subsets of the hierarchical basis that form trees.

With this specification of  $\mathcal{O}$ , there will be a 1-1 correspondence between the spaces  $X^{\delta} = \sum_{\lambda \in \sigma_{\lambda}} \sigma_{\lambda} \otimes W_{\lambda}^{\delta}$  with  $W_{\lambda}^{\delta} \in \mathcal{O}$  that satisfy (6.53), and the spans of collections of tensor products of wavelets  $\sigma_{\lambda}$  and hierarchical basis functions whose sets of index pairs are *lower*, also known as *downward closed*. Given such a  $X^{\delta}$ , we will define  $X^{\delta}$  by a certain enlargement the lower set.

For  $d \geq 2$ , let  $\mathbb{T}$  be the family of all *conforming* partitions of a polytope  $\Omega \subset \mathbb{R}^d$  into (closed) *d*-simplices that can be created by NVB starting from some given conforming initial partition  $\mathcal{T}_{\perp}$  with an assignment of the newest vertices that satisfies the matching condition, see [Ste08b]. We define a *partial order* on  $\mathbb{T}$  by writing  $\mathcal{T} \preceq \tilde{\mathcal{T}}$  when  $\tilde{\mathcal{T}}$  is a refinement of  $\mathcal{T}$ .

With some small adaptations that we leave to the reader, in the following the case d = 1 can be included by letting  $\mathbb{T}$  to be the family of a partitions of  $\Omega$  into (closed) subintervals that can be constructed by bisections from  $\mathcal{T}_{\perp} = {\Omega}$  such that the generations of any two neighbouring subintervals in any  $\mathcal{T} \in \mathbb{T}$  differ by not more than one.

The collection  $\mathcal{O}$  that we will consider is formed by the spaces  $W = W_{\mathcal{T}}$  of *continuous piecewise linears* w.r.t.  $\mathcal{T} \in \mathbb{T}$ , zero on a possible Dirichlet boundary  $\Gamma_D$  being the union of  $\partial T \cap \partial \Omega$  for some  $T \in \mathcal{T}_{\perp}$ . We expect that generalizations to finite element spaces of higher order do not impose essential difficulties.

For  $T \in \mathfrak{T} := \bigcup_{\mathcal{T} \in \mathbb{T}} \{T : T \in \mathcal{T}\}\)$ , we set gen(T) to be the number of bisections needed to create T from its 'ancestor'  $T' \in \mathcal{T}_{\perp}$ . With  $\mathfrak{N}$  being the set of all vertices (or nodes) of all  $T \in \mathfrak{T}$ , for  $\nu \in \mathfrak{N}$  we set  $gen(\nu) := \min\{gen(T) : T \in \mathfrak{T}, \nu \in T\}$ .

Any  $\nu \in \mathfrak{N}$  with  $gen(\nu) > 0$  is the midpoint of an edge of one or more  $T \in \mathfrak{T}$ with  $gen(T) = gen(\nu) - 1$ . The vertices  $\tilde{\nu}$  of these T with  $gen(\tilde{\nu}) = gen(\nu) - 1$ are defined as the parents of  $\nu$ . We denote this relation between a parent  $\tilde{\nu}$ and a child  $\nu$  by  $\tilde{\nu} \triangleleft_{\mathfrak{N}} \nu$ , see Figure 6.1. Vertices  $\nu \in \mathfrak{N}$  with  $gen(\nu) = 0$  have no parents.



FIGURE 6.1.  $\tilde{\nu}_1, \tilde{\nu}_2 \triangleleft_{\mathfrak{N}} \nu$ , and  $\mathcal{T}$  and its refinement  $\mathcal{T}^{d+}$  (for d = 2).

An (essentially) non-overlapping partition  $\mathcal{T}$  of  $\overline{\Omega}$  into  $T \in \mathfrak{T}$  is in  $\mathbb{T}$  if and only if the set  $N_{\mathcal{T}}$  of vertices of all  $T \in \mathcal{T}$  forms a *tree*, meaning that it contains all  $\nu \in \mathfrak{N}$  with gen $(\nu) = 0$  as well as all parents of any  $\nu \in N_{\mathcal{T}}$  with gen $(\nu) > 0$ , cf. [DKS16] for the d = 2 case.

**Definition 6.6.1.** For any  $\mathcal{T} \in \mathbb{T}$ , we define  $\mathcal{T}^{d+} \in \mathbb{T}$  (denoted as  $\mathcal{T}^{++}$  in [DKS16] for the d = 2 case) by replacing any  $T \in \mathcal{T}$  by its  $2^d$  'descendants' of the *d*th generation, see Figure 6.1.

Since this refinement adds exactly one vertex at the midpoint on any edge of all  $T \in \mathcal{T}$ , one infers that indeed  $\mathcal{T}^{d+} \in \mathbb{T}$ . The corresponding tree  $N_{\mathcal{T}^{d+}}$  is created from  $N_{\mathcal{T}}$  by the addition of all descendants up to generation d of all  $\nu \in N_{\mathcal{T}}$ .<sup>1</sup>

For  $\nu \in \mathfrak{N}$ , we set  $\phi_{\nu}$  as the continuous piecewise linear function w.r.t. the *uniform partition*  $\{T \in \mathfrak{T}: \text{gen}(T) = \text{gen}(\nu)\} \in \mathbb{T}$ , which function is 1 at  $\nu$  and 0 at all other vertices of this partition. Setting  $\mathfrak{N}_0 := \mathfrak{N} \setminus \Gamma_D$  and, for any  $\mathcal{T} \in \mathbb{T}$ ,  $N_{\mathcal{T},0} := N_{\mathcal{T}} \setminus \Gamma_D$ , the collection  $\{\phi_{\nu} : \nu \in \mathfrak{N}_0\}$  is known as the *hierarchical basis*, and for any  $\mathcal{T} \in \mathbb{T}$ , it holds that  $W_{\mathcal{T}} = \text{span}\{\phi_{\nu} : \nu \in N_{\mathcal{T},0}\}$ .

With above specification of the collection  $\mathcal{O}$  of finite element spaces, there exists a 1-1 correspondence between the spaces  $\sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes W_{\lambda}^{\delta}$  with  $W_{\lambda}^{\delta} \in \mathcal{O}$  that satisfy (6.53), and the spaces of the form

(6.59) 
$$X^{\delta} = \operatorname{span}\{\sigma_{\lambda} \otimes \phi_{\nu} \colon (\lambda, \nu) \in I_{\delta,0} \coloneqq I_{\delta} \setminus (\vee_{\Sigma} \times \Gamma_{D})\}$$

for some finite  $I_{\delta} \subset \vee_{\Sigma} \times \mathfrak{N}$  being a *lower set* in the sense

(6.60) 
$$(\lambda,\nu) \in I_{\delta} \text{ and } \begin{cases} \tilde{\lambda} \triangleleft_{\Sigma} \lambda \implies (\tilde{\lambda},\nu) \in I_{\delta}, \\ \tilde{\nu} \triangleleft_{\mathfrak{N}} \nu \text{ or } \operatorname{gen}(\tilde{\nu}) = 0 \implies (\lambda,\tilde{\nu}) \in I_{\delta}. \end{cases}$$

For above specification of  $X^{\delta}$ , from Proposition 6.5.2 with the specification (6.57) one infers that the corresponding space

(6.61) 
$$Y^{\delta} = \operatorname{span}\{\psi_{\mu} \otimes \phi_{\nu} \colon (\mu, \nu) \in I_{\delta, 0}^{Y}\},\$$

<sup>&</sup>lt;sup>1</sup>The addition of only all children of all  $\nu \in N_{\mathcal{T}}$  yields a tree only if  $\mathcal{T}$  is a uniform partition.

where

$$(6.62) \quad I_{\delta,0}^Y := \{(\mu,\nu) \colon \exists (\lambda,\nu) \in I_{\delta,0}, \ \mu \in \lor_{\Psi}, \ |\mu| = |\lambda|, \ |\mathcal{S}_{\Psi}(\mu) \cap \mathcal{S}_{\Sigma}(\lambda)| > 0\}$$

which index set is a lower set.

Remark 6.6.2 (Complexity of matrix-vector multiplications). The fact that the index sets of the bases for  $X^{\delta}$  and  $Y^{\delta}$  are lower sets is the key why it is possible to compute residuals of the system  $S^{\delta} \delta u^{\delta} = f^{\delta}$  ((6.31)) in  $\mathcal{O}(\dim X^{\delta})$  operations. Indeed, when one has a bilinear form that is 'local' and equals the tensor product of bilinear forms in time and space, and two spaces spanned by tensor product multi-level bases corresponding to lower sets, then the resulting generalized system matrix w.r.t. both bases can be applied in a number of operations that is proportional to the sum of the dimensions of both spaces. The algorithm that realizes this complexity makes a clever use of multi- to single-scale transformations alternately in time and space. In a 'uniform' sparse-grid setting, i.e., without 'local refinements', this algorithm was introduced in [BZ96], and it was later extended to general lower sets in [KS14]. The definition of a lower set in [KS14], there called multi-tree, is more restrictive than our current definition that allows more localized refinements. Details about the matrixvector multiplication and a proof of its optimal computational complexity is given in Chapter 7.

**Definition 6.6.3.** Given  $X^{\delta} = \operatorname{span}\{\sigma_{\lambda} \otimes \hat{\phi}_{\nu} : (\lambda, \nu) \in I_{\delta,0}\}$  for some lower set  $I_{\delta} \subset \bigvee_{\Sigma} \times \mathfrak{N}$ , we define the lower set  $I_{\delta}$ , and with that  $X^{\underline{\delta}}$ , by adding, for each  $(\lambda, \nu) \in I_{\delta}$  and any child  $\tilde{\lambda}$  of  $\lambda$  and any descendant  $\tilde{\nu}$  of  $\nu$  up to generation d, all pairs  $(\tilde{\lambda}, \nu)$  and  $(\lambda, \tilde{\nu})$  to  $I_{\delta}$ .

# 6.6.2 The collection $\Theta_{\delta}$ such that $X^{\underline{\delta}} = X^{\delta} \oplus \Theta_{\delta}$

Recall that for the bulk chasing process we need an '*X*-stable' basis  $\Theta_{\delta}$  that spans an '*X*-stable' complement space of  $X^{\delta}$  in  $X^{\delta}$ , i.e., a collection that satisfies (6.24). For that goal we define a *modified hierarchical basis* { $\hat{\phi}_{\nu} : \nu \in \mathfrak{N}_0$ } by  $\hat{\phi}_{\nu} := \phi_{\nu}$  when gen( $\nu$ ) = 0, and

$$\hat{\phi}_{\nu} := \phi_{\nu} - \frac{\sum_{\{\tilde{\nu} \in \mathfrak{N}_0 : \ \tilde{\nu} \triangleleft_{\mathfrak{N}}\nu\}} \frac{\int_{\Omega} \phi_{\nu} dx}{\int_{\Omega} \phi_{\tilde{\nu}} dx} \phi_{\tilde{\nu}}}{\#\{\tilde{\nu} \in \mathfrak{N} : \ \tilde{\nu} \triangleleft_{\mathfrak{N}}\nu\}}$$

otherwise. Notice that for those  $\nu$  with  $gen(\nu) > 0$  that have all their parents not on  $\Gamma_D$  it holds that  $\int_{\Omega} \hat{\phi}_{\nu} dx = 0$ , i.e.,  $\hat{\phi}_{\nu}$  has a *vanishing moment*, and furthermore that for any  $\mathcal{T} \in \mathbb{T}$ ,  $W_{\mathcal{T}} = \operatorname{span}\{\hat{\phi}_{\nu} : \nu \in N_{\mathcal{T},0}\}$ .

For any  $\mathcal{T} \in \mathbb{T}$ , it holds that

$$W_{\mathcal{T}} = \operatorname{span}\{\phi_{\nu} \colon \nu \in N_{\mathcal{T},0}\} = \operatorname{span}\{\hat{\phi}_{\nu} \colon \nu \in N_{\mathcal{T},0}\},\$$

and thus for any lower set  $I_{\delta} \subset \vee_{\Sigma} \times \mathfrak{N}$ ,

$$X^{\delta} = \operatorname{span} \{ \sigma_{\lambda} \otimes \hat{\phi}_{\nu} \colon (\lambda, \nu) \in I_{\delta, 0} \} = \operatorname{span} \{ \sigma_{\lambda} \otimes \phi_{\nu} \colon (\lambda, \nu) \in I_{\delta, 0} \}.$$

Moreover, for any  $\mathcal{T} \in \mathbb{T}$ , the basis transformation from the modified to unmodified hierarchical basis for  $W_{\mathcal{T}}$  can be applied in linear complexity traversing from the leaves to the roots.

Given  $\delta$ , the collection  $\Theta_{\delta}$  will be the set of properly normalized functions  $\sigma_{\lambda} \otimes \hat{\phi}_{\nu}$  for  $(\lambda, \nu) \in I_{\underline{\delta},0} \setminus I_{\delta,0}$ . In order to demonstrate (6.24), we have to impose some gradedness assumption on the lower sets  $I_{\delta}$ .

**Definition 6.6.4.** The gradedness constant of a lower set  $I_{\delta} \subset \vee_{\Sigma} \times \mathfrak{N}$  is the smallest  $L_{\delta} \in \mathbb{N}$  such that for all  $(\lambda, \nu) \in I_{\delta}$  for which  $\nu$  has an ancestor  $\tilde{\nu} \in \mathfrak{N}$  with  $gen(\nu) - gen(\tilde{\nu}) = L_{\delta}$ , it holds that  $(\check{\lambda}, \tilde{\nu}) \in I_{\delta}$  for any child  $\check{\lambda} \in \vee_{\Sigma}$  of  $\lambda$ .

*Remark* 6.6.5 (Uniform boundedness of the gradedness constants). Under the (unproven) assumption that our adaptive method creates a sequence of spaces  $X^{\delta}$  which are quasi-optimal for the approximation of the solution of the the parabolic PDE, one may hope that these spaces have a *uniformly bounded grad-edness constant*, unless (locally) the solution u is extremely more smooth as function of t than as function of the spatial variables.

To see this, consider the non-adaptive sparse grid index sets of the form  $\{(\lambda, \nu) \in \forall_{\Sigma} \times \mathfrak{N} : \tilde{L}|\lambda| + \operatorname{gen}(\nu) \leq N\}$  for some constant  $\tilde{L}$  and  $N \in \mathbb{N}$ , which are appropriate when the behaviour of u as function of t on the one hand and that of the spatial variables on the other is globally similar. Then for  $\tilde{L} \leq L$ , the gradedness constant of this index set is  $\leq L$ , where the smallest spatial resolution in the 'sparse-grid mesh' equals the smallest temporal resolution in this mesh to the power  $\tilde{L}/d$ . So only when a polynomial decay of the spatial resolution as function of the temporal resolution does not suffice for a proper approximation of u, one cannot expect to have a gradedness constant that is uniformly bounded.

**Proposition 6.6.6.** For  $(\lambda, \nu) \in \bigvee_{\Sigma} \times \mathfrak{N}_{0}$ , let  $e_{\lambda\mu} := 1/\sqrt{2^{(\frac{2}{d}-1)\operatorname{gen}(\nu)} + 4^{|\lambda|}2^{(-\frac{2}{d}-1)\operatorname{gen}(\nu)}}$ and  $\theta_{\lambda\nu} := e_{\lambda\mu}\sigma_{\lambda} \otimes \hat{\phi}_{\nu}$ . For any  $\delta \in \Delta$ , let  $\Theta_{\delta} := \{\theta_{\lambda\nu} : (\lambda, \nu) \in J_{\delta} := I_{\underline{\delta},0} \setminus I_{\delta,0}\}$ . Then  $X^{\delta} \oplus \operatorname{span} \Theta_{\delta} = X^{\underline{\delta}}$ , and there exist constants  $0 < m_{\delta} \leq M_{\delta}$ , only dependent on the gradedness constant  $L_{\delta}$ , such that for any  $z \in X^{\delta}$  and  $\mathbf{c} = (c_{\lambda\nu})_{(\lambda,\nu)\in I_{\delta,0}} \setminus I_{\delta,0} \subset \mathbb{R}$ ,

$$m_{\delta}(\|z\|_{X}^{2} + \|\mathbf{c}\|^{2}) \le \|z + \mathbf{c}^{\top}\Theta_{\delta}\|_{X}^{2} \le M_{\delta}(\|z\|_{X}^{2} + \|\mathbf{c}\|^{2})$$

So under the mild assumption that the gradedness constants of the sets  $X^{\delta}$  that we encounter are uniformly bounded, we have shown that the condition (6.24) is satisfied.

*Proof.* Setting  $c_{\lambda\nu} := 0$  when  $(\lambda, \nu) \notin I_{\delta,0} \setminus I_{\delta,0}$ , and writing  $z = \sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes w_{\lambda}$ where  $w_{\lambda} \in \text{span}\{\hat{\phi}_{\nu} : (\lambda, \nu) \in I_{\delta,0}\}$ , from  $\Sigma$  and  $\{2^{-|\lambda|}\sigma_{\lambda} : \lambda \in \vee_{\Sigma}\}$  being Riesz bases for  $L_2(I)$  and  $H^1(I)$ , and  $\mathbf{c}^{\top}\Theta_{\delta} = \sum_{\lambda} \sigma_{\lambda} \otimes \sum_{\nu} e_{\lambda\mu}c_{\lambda\nu}\hat{\phi}_{\nu}$ , an application of Lemma 6.6.7 given below shows that

$$\begin{split} \|z + \mathbf{c}^{\top} \Theta_{\delta}\|_{X}^{2} &\approx \sum_{\lambda} \left\{ \|w_{\lambda} + \sum_{\nu} e_{\lambda\mu} c_{\lambda\nu} \hat{\phi}_{\nu}\|_{H^{1}(\Omega)}^{2} + 4^{|\lambda|} \|w_{\lambda} + \sum_{\nu} e_{\lambda\mu} c_{\lambda\nu} \hat{\phi}_{\nu}\|_{H^{1}_{0,\Gamma_{D}}(\Omega)'}^{2} \right\} \\ &\approx \left\{ \sum_{\lambda} \|w_{\lambda}\|_{H^{1}(\Omega)}^{2} + 4^{|\lambda|} \|w_{\lambda}\|_{H^{1}_{0,\Gamma_{D}}(\Omega)'}^{2} + \sum_{\nu} \left( 2^{(\frac{2}{d}-1)\operatorname{gen}(\nu)} + 4^{|\lambda|} 2^{(-\frac{2}{d}-1)\operatorname{gen}(\nu)} \right) |e_{\lambda\mu} c_{\lambda\nu}|^{2} \right\} \\ &= \sum_{\lambda} \|w_{\lambda}\|_{H^{1}(\Omega)}^{2} + 4^{|\lambda|} \|w_{\lambda}\|_{H^{1}_{0,\Gamma_{D}}(\Omega)'}^{2} + \sum_{\nu} |c_{\lambda\nu}|^{2} \approx \|z\|_{X}^{2} + \|\mathbf{c}\|^{2}, \end{split}$$

with the  $\approx$ -symbol in the second line dependent on the gradedness constant.  $\Box$ 

**Lemma 6.6.7.** For  $\tilde{\mathcal{T}} \in \mathbb{T}$ , and either  $\mathbb{T} \ni \mathcal{T} \preceq \tilde{\mathcal{T}}$  and  $v \in W_{\mathcal{T}}$ , or  $\mathcal{T} = \emptyset$ ,  $N_{\mathcal{T},0} := \emptyset$ , and v = 0, and scalars  $(d_{\nu})_{\nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}}$ , it holds that

(6.63) 
$$\|v + \sum_{\nu} d_{\nu} \hat{\phi}_{\nu}\|_{H^{1}(\Omega)}^{2} \approx \|v\|_{H^{1}(\Omega)}^{2} + \sum_{\nu} 2^{(\frac{2}{d}-1)\operatorname{gen}(\nu)} |d_{\nu}|^{2}$$

(6.64) 
$$\|v + \sum_{\nu} d_{\nu} \hat{\phi}_{\nu}\|_{H^{1}_{0,\Gamma_{D}}(\Omega)'}^{2} \approx \|v\|_{H^{-1}(\Omega)}^{2} + \sum_{\nu} 2^{(-\frac{2}{d}-1)\operatorname{gen}(\nu)} |d_{\nu}|^{2}$$

with the constants hidden in the  $\approx$ -symbols only dependent on  $M_{\tilde{\mathcal{T}}\mathcal{T}} := \max\{\operatorname{gen}(\tilde{T}) - \operatorname{gen}(T) : \tilde{\mathcal{T}} \ni \tilde{T} \subset T \in \mathcal{T}\}$  or  $M_{\tilde{\mathcal{T}}\mathcal{T}} := \max\{\operatorname{gen}(\tilde{T}) : \tilde{T} \in \tilde{\mathcal{T}}\}$  for  $\mathcal{T} = \emptyset$ .

*Proof.* Once the equivalences are shown uniformly in any  $\mathcal{T} \leq \tilde{\mathcal{T}}$  for which  $M_{\tilde{\mathcal{T}}\mathcal{T}} = 1$ , a repeated application of these equivalences shows them for the general case, with constants that are only dependent on  $M_{\tilde{\mathcal{T}}\mathcal{T}}$ . So in the following, it suffices to consider the case that  $M_{\tilde{\mathcal{T}}\mathcal{T}} = 1$ . The case  $\mathcal{T} = \emptyset$  is easy, so we will consider the case that  $\mathcal{T} \in \mathbb{T}$ .

Let  $\Phi_{\tilde{\tau}} = \{\phi_{\tilde{\tau},\nu} \colon \nu \in N_{\tilde{\tau},0}\}$  denote the standard nodal basis for  $W_{\tilde{\tau}}$ . For any weight function  $0 < w_{\tilde{\tau}} \in \prod_{T \in \tilde{\tau}} P_0(T)$ , with  $\|\cdot\|_{L_{2,w_{\tilde{\tau}}}(\Omega)} \coloneqq \|w_{\tilde{\tau}}^{\frac{1}{2}}\cdot\|_{L_2(\Omega)}$ it holds that  $\|\sum_{\nu} c_{\nu}\phi_{\tilde{\tau},\nu}\|_{L_{2,w_{\tilde{\tau}}}(\Omega)}^2 \equiv \sum_{\nu} |c_{\nu}|^2 \|\phi_{\tilde{\tau},\nu}\|_{L_{2,w_{\tilde{\tau}}}(\Omega)}^2$ , only dependent on the spectrum of the element mass matrix on a reference element, i.e., on the space dimension d, so independent of the weight function  $w_{\tilde{\tau}}$ . We refer to this equivalence by saying that  $\Phi_{\tilde{\tau}}$  is (uniformly) stable w.r.t.  $\|\cdot\|_{L_{2,w_{\tilde{\tau}}}(\Omega)}$ .

Notice that for  $\nu \in N_{\tilde{\tau},0} \setminus N_{\tau,0}$ , it holds that  $\phi_{\tilde{\tau},\nu} = \phi_{\nu}$ . W.r.t. the splitting  $N_{\tilde{\tau},0} = N_{\tau,0} + N_{\tilde{\tau},0} \setminus N_{\tau,0}$ , the basis transformation from  $\Phi_{\tau} \cup \{\hat{\phi}_{\nu} : \nu \in N_{\tilde{\tau},0} \setminus N_{\tau,0}\}$  to  $\Phi_{\tau} \cup \{\phi_{\nu} : \nu \in N_{\tilde{\tau},0} \setminus N_{\tau,0}\}$  is of the form  $\begin{bmatrix} \mathrm{Id} & * \\ 0 & \mathrm{Id} \end{bmatrix}$ , and the
basis transformation from the latter basis to  $\Phi_{\tilde{\mathcal{T}}}$  is of the form  $\begin{bmatrix} \mathrm{Id} & 0 \\ * & \mathrm{Id} \end{bmatrix}$ . The entries in both non-zero off-diagonal blocks are uniformly bounded, where non-zeros can only occur for index pairs  $(\nu, \tilde{\nu})$  that are vertices of the same  $\tilde{T} \in \tilde{\mathcal{T}}$ . Consequently, for a family of weight functions  $(w_{\tilde{\mathcal{T}}})_{\tilde{\mathcal{T}} \in \mathbb{T}}$  that have *uniformly bounded jumps* in the sense that

(6.65) 
$$\sup_{\tilde{\mathcal{T}}\in\mathbb{T}}\sup_{\{T,T'\in\tilde{\mathcal{T}}:\ T\cap T'\neq\emptyset\}}\frac{w_{\tilde{\mathcal{T}}}|_T}{w_{\tilde{\mathcal{T}}}|_{T'}}<\infty,$$

all basis transformations between the  $L_{2,w_{\tilde{\tau}}}(\Omega)$ -normalized bases  $\Phi_{\mathcal{T}} \cup \{\hat{\phi}_{\nu} \colon \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$ ,  $\Phi_{\mathcal{T}} \cup \{\phi_{\nu} \colon \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$  and  $\Phi_{\tilde{\mathcal{T}}}$  are uniformly bounded.

Since, as we have seen,  $\Phi_{\tilde{\mathcal{T}}}$  is (uniformly) stable w.r.t.  $\|\cdot\|_{L_{2,w_{\tilde{\mathcal{T}}}}(\Omega)}$ , we conclude that also  $\Phi_{\mathcal{T}} \cup \{\hat{\phi}_{\nu} : \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$  and  $\Phi_{\mathcal{T}} \cup \{\phi_{\nu} : \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$  are (uniformly) stable w.r.t.  $\|\cdot\|_{L_{2,w_{\tilde{\mathcal{T}}}}(\Omega)}$ . Because of the *uniform K-mesh property* of  $\mathcal{T} \in \mathbb{T}$ , examples of families of weights that satisfy (6.65) are given by  $(h_{\tilde{\mathcal{T}}}^s)_{\tilde{\mathcal{T}} \in \mathbb{T}}$  for any  $s \in \mathbb{R}$ , where  $h_{\tilde{\mathcal{T}}}|_T := 2^{-\operatorname{gen}(T)/d} (\eqsim |T|^{1/d})$   $(T \in \tilde{\mathcal{T}})$ .

For showing (6.63), let  $P_{\mathcal{T}}: W_{\tilde{\mathcal{T}}} \to W_{\mathcal{T}}$  be the projector with ran  $P_{\mathcal{T}} = W_{\mathcal{T}}$ and ran(Id  $-P_{\mathcal{T}}$ ) = span{ $\hat{\phi}_{\nu}: \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}$ }. Using the form of the basis transformation from  $\Phi_{\mathcal{T}} \cup \{\phi_{\nu}: \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$  to  $\Phi_{\mathcal{T}} \cup \{\hat{\phi}_{\nu}: \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$ , one infers that

$$\mathrm{Id} - P_{\mathcal{T}} = J_{\mathcal{T}} \circ (\mathrm{Id} - I_{\mathcal{T}}),$$

where  $I_{\mathcal{T}}$  is the nodal interpolator onto  $W_{\mathcal{T}}$ , and  $J_{\mathcal{T}}$  is defined by  $J_{\mathcal{T}}\phi_{\nu} = \hat{\phi}_{\nu}$ . Since both  $\{\hat{\phi}_{\nu} \colon \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$  and  $\{\phi_{\nu} \colon \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$  are uniformly stable w.r.t.  $\|h_{\tilde{\mathcal{T}}}^{-1} \cdot \|_{L_2(\Omega)}$ , and  $\|h_{\tilde{\mathcal{T}}}^{-1}\hat{\phi}_{\nu}\|_{L_2(\Omega)} \approx \|h_{\tilde{\mathcal{T}}}^{-1}\phi_{\nu}\|_{L_2(\Omega)}$ , it follows that  $J_{\mathcal{T}}$  is uniformly bounded w.r.t.  $\|h_{\tilde{\mathcal{T}}}^{-1} \cdot \|_{L_2(\Omega)}$ , i.e.  $\|h_{\tilde{\mathcal{T}}}^{-1}J_{\mathcal{T}}h_{\tilde{\mathcal{T}}}\|_{\mathcal{L}(L_2(\Omega),L_2(\Omega))} \lesssim 1$ , and so

$$\|h_{\tilde{\mathcal{T}}}^{-1}(\mathrm{Id}-P_{\mathcal{T}})v\|_{L_{2}(\Omega)} \lesssim \|h_{\tilde{\mathcal{T}}}^{-1}(\mathrm{Id}-I_{\mathcal{T}})v\|_{L_{2}(\Omega)} \lesssim |v|_{H^{1}(\Omega)} \quad (v \in W_{\tilde{\mathcal{T}}}).$$

Using the common inverse inequality  $\|\cdot\|_{H^1(\Omega)} \lesssim \|h_{\tilde{\tau}}^{-1}\cdot\|_{L_2(\Omega)}$  on  $W_{\tilde{\tau}}$ , we infer that  $(\mathrm{Id} - P_{\mathcal{T}})$  is uniformly bounded in the  $H^1(\Omega)$ -norm, and that  $\|\cdot\|_{H^1(\Omega)} \approx \|h_{\tilde{\tau}}^{-1}\cdot\|_{L_2(\Omega)}$  on  $\mathrm{ran}(\mathrm{Id} - P_{\mathcal{T}})$ . The proof of (6.63) is completed by the uniform stability of  $\{\hat{\phi}_{\nu} \colon \nu \in N_{\tilde{\tau},0} \setminus N_{\mathcal{T},0}\}$  w.r.t.  $\|h_{\tilde{\tau}}^{-1}\cdot\|_{L_2(\Omega)}$ , and the fact that  $\|h_{\tilde{\tau}}^{-1}\hat{\phi}_{\nu}\|_{L_2(\Omega)}^2 \approx 2^{(\frac{2}{d}-1)\operatorname{gen}(\nu)}$ .

Moving to (6.64), either by  $\int_{\Omega} \hat{\phi}_{\nu} dx = 0$ , or otherwise using the proximity of the Dirichlet boundary  $\Gamma_D$  by an application of Poincaré's inequality, it holds that

$$|\langle \hat{\phi}_{\nu}, v \rangle_{L_2(\Omega)}| \lesssim 2^{-\operatorname{gen}(\nu)/d} \| \hat{\phi}_{\nu} \|_{L_2(\Omega)} |v|_{H^1(\operatorname{supp} \hat{\phi}_{\nu})} \quad (\nu \in \mathfrak{N}_0 \setminus N_{\mathcal{T}_{\perp}, 0}).$$

By using that for  $T \in \tilde{\mathcal{T}}$  the number of  $\nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}$  for which  $\operatorname{supp} \hat{\phi}_{\nu}$  has non-empty intersection with T is uniformly bounded, and furthermore that  $\Phi_{\mathcal{T}} \cup \{\hat{\phi}_{\nu} \colon \nu \in N_{\tilde{\mathcal{T}},0} \setminus N_{\mathcal{T},0}\}$  is uniformly stable w.r.t.  $\|h_{\mathcal{T}} \cdot \|_{L_2(\Omega)}$ , we infer that for any  $z = \sum_{\nu \in N_{\mathcal{T},0}} z_{\nu} \phi_{\mathcal{T},\nu} \in W_{\mathcal{T}}$  it holds that

the last inequality by application of a less common inverse inequality which proof can be found in Lemma 2.3.4 for general dimensions *d*. From (6.66) it follows that Id  $-P_{\mathcal{T}}$  is uniformly bounded in the  $H^1_{0,\Gamma_D}(\Omega)'$ -norm, and also that  $\|\sum_{\nu\in N_{\tilde{\mathcal{T}},0}\setminus N_{\mathcal{T},0}} d_{\nu}\hat{\phi}_{\nu}\|^2_{H^1_{0,\Gamma_D}(\Omega)'} \approx \sum_{\nu\in N_{\tilde{\mathcal{T}},0}\setminus N_{\mathcal{T},0}} |d_{\nu}|^2 2^{(-\frac{2}{d}-1)\operatorname{gen}(\nu)}$ , where we used that  $\|h_{\tilde{\mathcal{T}}}\hat{\phi}_{\nu}\|^2_{L_2(\Omega)} \approx 2^{(-\frac{2}{d}-1)\operatorname{gen}(\nu)}$ . The proof of (6.64) is completed.  $\Box$ 

## 6.6.3 The wavelet collections $\Sigma$ and $\Psi$

As wavelet basis  $\Sigma = \{\sigma_{\lambda} : \lambda \in \vee_{\Sigma}\}$  we select the three-point hierarchical basis illustrated in Figure 6.2. This basis is known to be a Riesz basis for  $L_2(I)$ , and, after re-normalization, for  $H^1(I)$  (see [Ste96]). It also satisfies the other assumptions made in Sect. 6.5.2. The wavelets up to level  $\ell$  span the space of continuous piecewise linear functions on I w.r.t. the uniform partition into  $2^{\ell}$  subintervals.

As wavelet basis  $\Psi = \{\psi_{\mu} : \mu \in \lor_{\Psi}\}$  we take the orthonormal (discontinuous) piecewise linear wavelets, see Figure 6.3. The wavelets up to level  $\ell$  span the space of (discontinuous) piecewise linear functions on *I* w.r.t. the uniform partition into  $2^{\ell}$  subintervals.

# 6.6.4 The family $({}^{\delta}G, {}^{\delta}U_0)_{\delta \in \Delta}$

The index set  $\forall_{\Sigma}$  is naturally identified with the set of 'nodal dyadic' points, see Figure 6.4, which is the natural index set for the one-dimensional hierarchical basis that we denote by  $\{\phi_{\lambda} : \lambda \in \forall_{\Sigma}\}$ . Recalling that for  $\delta \in \Delta$ ,  $X^{\delta} =$  $\operatorname{span}\{\sigma_{\lambda} \otimes \phi_{\nu} : (\lambda, \nu) \in I_{\delta,0} = I_{\delta} \setminus (\forall_{\Sigma} \times \Gamma_{D})\}$  for some lower set  $I_{\delta} \subset \forall_{\Sigma} \times \mathfrak{N}$ ,



FIGURE 6.2. Three-point hierarchical basis  $\Sigma$ . On level 0 there are two wavelets, and on level 1 there is one wavelet, whose parents are both wavelets on level 0. On each level  $\ell > 1$  there are  $2^{\ell-1}$  wavelets, among them near each boundary one boundary-adapted wavelet, where each wavelet has one parent being the wavelet on level  $\ell - 1$  whose support includes the support of its child (so  $S_{\Sigma}(\lambda)$  can be taken equal to supp  $\sigma_{\lambda}$ ). All but one wavelets have one (bdr. wav) or two vanishing moments.



FIGURE 6.3.  $L_2(I)$ -orthonormal (discontinuous) piecewise linear wavelet basis  $\Psi$ . On level  $\ell = 0$  there are 2 wavelets. On each level  $\ell \ge 1$  there are  $2^{\ell}$  wavelets of two types, each of them having 2 parents being the wavelets on level  $\ell - 1$  whose supports include the supports of their children (so  $S_{\Psi}(\mu)$  can be taken equal to supp  $\psi_{\mu}$ ). The wavelets on level 0 have either 0 or 1 vanishing moment, all other wavelets have two vanishing moments.



FIGURE 6.4. Index set  $\vee_{\Sigma}$  with parent-child relations, and the one-dimensional hierarchical basis.

we define

$${}^{\delta}G := \operatorname{span}\{\phi_{\lambda} \otimes \phi_{\nu} \colon (\lambda, \nu) \in I_{\delta}\}, \ {}^{\delta}U_0 := \operatorname{span}\{\phi_{\nu} \colon (\lambda, \nu) \in I_{\delta}, \ \phi_{\lambda}(0) \neq 0\}.$$

Since the level of resolution of these spaces is comparable to that of  $X^{\delta}$ , based on our experiences with wavelet and finite element methods we *expect* that with this choice of  $({}^{\delta}G, {}^{\delta}U_0)$  and the definition of  $X^{\underline{\delta}\delta}$ , that *saturation* holds, i.e., that Assumption 6.3.1 assumption is valid.

Given  $g \in Y'$  and  $u_0 \in L_2(\Omega)$ , it remains to define their approximations  $({}^{\delta}g, {}^{\delta}u_0) \in ({}^{\delta}G, {}^{\delta}U_0)$ . In general, the construction of these approximations depends on the data at hand. Below we give a construction that applies to general continuous g and  $u_0$ , and that avoids quadrature issues.

For  $\nu \in \mathfrak{N}$  with  $\operatorname{gen}(\nu) = 0$ , let  $\overline{\phi}_{\nu} := \delta_{\nu}$ . Each  $\nu \in \mathfrak{N}$  with  $\operatorname{gen}(\nu) > 0$ is the midpoint of an edge of a  $T \in \mathcal{T}$  with  $\operatorname{gen}(T) = \operatorname{gen}(\nu) - 1$ . Denoting the endpoints of this edge as  $\nu_1, \nu_2 \in \mathfrak{N}$ , let  $\overline{\phi}_{\nu} := \delta_{\nu} - \frac{1}{2}(\delta_{\nu_1} + \delta_{\nu_2})$ . Then  $\{\overline{\phi}_{\nu} : \nu \in \mathfrak{N}\} \subset C(\overline{\Omega})'$  is biorthogonal to  $\{\phi_{\nu} : \nu \in \mathfrak{N}\}$ . With  $\{\overline{\phi}_{\lambda} : \lambda \in \bigvee_{\Sigma}\} \subset C(\overline{I})'$  defined analogously for the one-dimensional case, for  $g \in C(\overline{I \times \Omega})$  and  $u_0 \in C(\overline{\Omega})$  we define the interpolants

$${}^{\delta}g := \sum_{(\lambda,\nu)\in I_{\delta}} (\tilde{\phi}_{\lambda}\otimes \tilde{\phi}_{\nu})(g)\phi_{\lambda}\otimes \phi_{\nu}, \ {}^{\delta}u_0 := \sum_{\{\nu: \ (\lambda,\nu)\in I_{\delta}, \ \phi_{\lambda}(0)\neq 0\}} \tilde{\phi}_{\nu}(u_0)\phi_{\nu}.$$

Since we expect that for sufficiently smooth g and  $u_0$ , the errors  $||g - \delta g||_{Y'}$  and  $||u_0 - \delta u_0||_{L_2(\Omega)}$  are of higher order than the approximation error  $\inf_{w \in X^{\delta}} ||u - w||_X$ , for our convenience in the adaptive Algorithm 6.4.8 we ignore errors caused by data-oscillation by setting  $\eta(\cdot) \equiv 0$ .

Notice that setting up the matrix vector formulation of the system (6.31) that defines our approximation  $u^{\delta}$  requires computing the vectors

$$\left[\langle {}^{\delta}g,\psi_{\mu}\otimes\phi_{\nu}\rangle_{L_{2}(I\otimes\Omega)}\right]_{(\mu,\nu)\in I_{\delta,0}^{Y}},\quad \left[\langle {}^{\delta}u_{0},\phi_{\nu}\rangle_{L_{2}(\Omega)}\right]_{\{\nu:\ (\lambda,\nu)\in I_{\delta,0},\ \sigma_{\lambda}(0)\neq0\}}$$

which can be performed in  $\mathcal{O}(\dim X^{\delta})$  operations because  $I_{\delta}$  and  $I_{\underline{\delta},0}^{Y}$  are lower sets (and  $\#I_{\delta,0}^{Y} \lesssim \#I_{\delta}$ ).

# 6.7 Numerical experiments

We test our algorithm on the heat equation, i.e., the parabolic problem with  $a(t; \eta, \zeta) = \int_{\Omega} \nabla \eta \cdot \nabla \zeta d\mathbf{x}$ , posed on a two-dimensional polygonal spatial domain  $\Omega$ , and Dirichlet boundary  $\Gamma_D = \partial \Omega$ . Recall from §6.6.3 the three-point continuous piecewise linear temporal wavelet basis  $\Sigma$ , the orthonormal discontinuous piecewise linear temporal wavelet basis  $\Psi$ , and the hierarchical continuous piecewise linear spatial basis  $\Xi := \{\phi_{\nu} : \nu \in \mathfrak{N}_0\}$ .

We consider 'trial' spaces  $X^{\delta}$  which are spanned by finite subsets of  $\Sigma \otimes \Xi$ whose index sets are lower sets (more precisely, satisfy (6.59)-(6.60)), and corresponding 'test' spaces  $Y^{\delta}$  spanned by finite subsets of  $\Psi \otimes \Xi$  as defined in (6.61)-(6.62). We construct the enlarged trial space  $X^{\delta}$  as defined in Def. 6.6.3, with corresponding test space  $Y^{\delta}$ .

For a given level  $N \in \mathbb{N}$ , span{ $\sigma_{\lambda} : |\lambda| \leq N$ } coincides with the span of the continuous piecewise linears on an N-times recursive dyadic refinement of I, and span{ $\phi_{\nu} \in \Xi : \text{gen}(\nu) \leq 2N$ } coincides with that of the continuous piecewise linears, zero at  $\partial\Omega$ , on a 2N-times recursive bisection refinement of an initial partition  $\mathcal{T}_{\perp}$ . Therefore, the span of the *'full'* tensor product { $\sigma_{\lambda} : |\lambda| \leq N$ }  $\otimes \{\phi_{\nu} : \text{gen}(\nu) \leq 2N\}$  equals a space of lowest order continuous finite elements w.r.t. a quasi-uniform shape regular product mesh into prismatic elements.

Taking only those index pairs  $(\lambda, \nu)$  for which  $2|\lambda| + \text{gen}(\nu) \leq 2N$  produces a 'sparse' tensor product on level N. Sparse tensor products allow to overcome the *curse of dimensionality* in the sense that for smooth solutions they achieve a rate in X-norm that is equal to the best rate in the  $H^1(\Omega)$ -norm that can be expected for the corresponding stationary problem on the spatial domain, here the Poisson equation; see also Sect. 6.5.5.

We run our adaptive Algorithm 6.4.8 with  $\theta = 0.5$  and  $\xi = \frac{1}{2}$ , computing  ${}^{\delta}g$  and  ${}^{\delta}u_0$  as in Sect. 6.6.4. Since we envisage that in our experiments dataoscillation errors are not dominant, for our convenience we took  $\omega = \infty$ . We solve the arising linear system of (6.31) using Preconditioned CG, using the previous solution as initial guess. We then perform Dörfler marking on the residual, yielding a minimal set J, and finally choose  $I_{\delta}$  as the smallest lower set containing  $J \cup I_{\delta}$ . Due to this constraint generally we add index pairs outside of the marked set, i.e.  $I_{\delta} \setminus I_{\delta} \supseteq J$ . Still, in our experiments, we *observe*  $\#I_{\delta} - \#I_{\delta} \leq \#J$  with a moderate constant.

*Remark* 6.7.1. Rather we would have applied an algorithm that produces a  $I_{\delta}$  such that  $I_{\delta} \setminus I_{\delta}$  is *guaranteed* to have an, up to a multiplicative factor, smallest cardinality among all lower sets  $I_{\delta} \supset I_{\delta}$  that realize the bulk criterion. Such an algorithm was introduced in [BD04, BFV19] for 'single-tree' approximation, but seems not to be available for the 'double-tree' (i.e. lower set) constraint that we need here.

We compare adaptive refinement with non-adaptive full- and sparse tensor

products, and monitor the error estimator  $\mathcal{E}^{\delta}(\tilde{u}^{\delta})$  from Proposition 6.4.5, the residual error estimator from Proposition 6.4.3, and the  $L_2(\Omega)$  trace error at t = 0.

#### 6.7.1 Condition numbers of preconditioner

For the calibration of our preconditioners, we consider  $\Omega := [0, 1]^2$ , and compare *uniformly refined* space-time meshes with *locally refined* meshes with refinements towards  $\{0\} \times \partial \Omega$ .

The replacement of the nonlocal operator  $(E_Y^{\delta'}AE_Y^{\delta})^{-1}$  in the forward application of  $S^{\delta\delta}$  by the block-diagonal preconditioner  $K_Y^{\delta}$  from Sect. 6.5.6 is only guaranteed to result in a convergent algorithm when the eigenvalues of  $K_Y^{\delta'}E_Y^{\delta''}AE_Y^{\delta'}$  are sufficiently close to one.

In Table 6.1, we investigate the values

$$\kappa_{\delta} := \max\{\lambda_{\max}(\mathbf{K}_{Y}^{\delta}\mathbf{A}_{Y}^{\delta}), 1/\lambda_{\min}(\mathbf{K}_{Y}^{\delta}\mathbf{A}_{Y}^{\delta})\}$$

with  $A_Y^{\delta}$  the matrix representation of  $E_Y^{\delta'}AE_Y^{\delta}$ , and  $\mathbf{K}_Y^{\delta}$  built from spatial multigrid preconditioners  $\mathbf{K}_{\mu}^{\delta}$  corresponding to *n* V-cycles. In each V-cycle we applied one pre- and one post Gauss-Seidel smoother. In case of a locally refined spatial mesh, on each level these Gauss-Seidel updates were restricted to the vertices whose generation is equal to that level as well as both endpoints of the edge on which these vertices were inserted ([WZ17]). We see that for both uniform and locally refined space-time meshes,  $\kappa_{\delta}$  converges to 1 rapidly in *n*, and is essentially independent of dim  $X^{\delta}$ . In our examples,  $\kappa_{\delta}$  is sufficiently close to one already for n = 1.

Fixing n = 1 for the forward application of  $S^{\delta\delta}$ , we want to precondition  $S^{\delta\delta}$  itself as well. Following Sect. 6.5.6, we build a block-diagonal preconditioner taking  $\mathbf{K}^{\delta}_{\lambda}$  to correspond to m V-cycles of the aforementioned multigrid method now applied to  $\mathbf{A}^{\delta}_{\lambda} + 2^{|\lambda|} \mathbf{M}^{\delta}_{\lambda}$  with  $\mathbf{A}^{\delta}_{\lambda}$  and  $\mathbf{M}^{\delta}_{\lambda}$  being stiffness- or mass-matrices.

$\dim X^{\delta}$		n = 1	n=2	n = 3	n = 4	n = 5	n = 6
uniform	729	1.343	1.070	1.017	1.004	1.001	1.000
35937		1.360	1.075	1.019	1.004	1.001	1.000
2146689		1.365	1.077	1.019	1.004	1.001	1.000
local	766	1.306	1.058	1.013	1.003	1.001	1.000
30151		1.307	1.058	1.013	1.003	1.001	1.000
1964797		1.307	1.058	1.013	1.003	1.001	1.000

TABLE 6.1. Values  $\kappa_{\delta} := \max\{\lambda_{\max}(\mathbf{K}_{Y}^{\delta}\mathbf{A}_{Y}^{\delta}), 1/\lambda_{\min}(\mathbf{K}_{Y}^{\delta}\mathbf{A}_{Y}^{\delta})\}$  using spatial multigrid with *n* V-cycles.

Table 6.2 shows the condition numbers of the preconditioned matrix. We again see fast stabilization in m as well as in dim  $X^{\delta}$ . We fix m = 3 in the sequel. Most interestingly, in every of our example problems, the adaptive algorithm only needs one or two PCG iterations to reach the error tolerance  $t_{\delta}$ .

### 6.7.2 Smooth problem

We consider the square domain  $\Omega := [0, 1]^2$  and prescribe

$$u(t, x, y) := (1 + t^2)x(1 - x)y(1 - y)$$

with derived data  $u_0$  and g. For this smooth solution, full and sparse tensor products are expected to yield the best possible error decays proportional to  $(\dim X^{\delta})^{-1/3}$  and  $(\dim X^{\delta})^{-1/2}$ , respectively.

The left side of Figure 6.5 shows the error progressions for the smooth problem. We plot the error estimator  $\mathcal{E}^{\delta}(\tilde{u}^{\delta}) := \mathcal{E}^{\delta}(\tilde{u}^{\delta}; \delta g, \delta u_0) \approx \|\delta u - \tilde{u}^{\delta}\|_X$  from Proposition 6.4.5, the residual error estimator  $\|\mathbf{r}^{\delta}\|$ , and  $\|\gamma_0(\delta u - \tilde{u}^{\delta})\|_{L_2(\Omega)}$ . We see that the error progressions are as expected. For this solution, adaptive refinement yields no advantage over sparse grid refinement. We observe a higher order of convergence for the trace at t = 0 measured in  $L_2(\Omega)$ .

### 6.7.3 Moving peak problem

We consider a square domain  $\Omega := [0, 1]^2$  and select

$$u(t, x, y) := x(1 - x)y(1 - y)\exp(-100[(x - t)^{2} + (y - t)^{2}]).$$

We took this example from [LS20]. The solution is smooth, and almost zero everywhere except on a small strip near the diagonal from (0,0,0) to (1,1,1) of the space-time cylinder. As *u* is smooth, we expect sparse grid refinements

d	$\lim X^{\delta}$	m = 1	m = 2	m = 3	m = 4	m = 5	m = 6
uniform	4913	9.196	6.119	6.048	6.042	6.041	6.041
35937		9.718	6.315	6.263	6.260	6.260	6.260
- 	274625	9.991	6.750	6.749	6.751	6.752	6.752
2146689		10.115	7.080	7.087	7.088	7.088	7.088
local	3520	5.707	5.132	5.110	5.111	5.111	5.111
	30151	6.355	5.734	5.706	5.704	5.704	5.704
244870		7.619	6.879	6.843	6.841	6.841	6.841
1964797		9.353	8.734	8.703	8.701	8.701	8.701

TABLE 6.2. Spectral condition numbers of  $K_X^{\delta}S^{\underline{\delta}\delta}$ , using spatial multigrid with m V-cycles.



FIGURE 6.5. Error progressions for (left) the *smooth problem* and (right) the *moving peak* problem. Shown: estimated *X*-norm error (solid line), residual norm (dashed), and t = 0 trace error (dotted) as a function of dim  $X^{\delta}$  for adaptive (black), sparse grid (red), and full grid refinement (orange).

to asymptotically yield the optimal error decay proportional to  $(\dim X^{\delta})^{-1/2}$ , albeit with a terrible constant. Adaptive refinement should be able to achieve the same rate at quantitatively smaller doubletrees.

From the right of Figure 6.5, we see that the sparse grid rate is not (yet) optimal, while our adaptive routine is able to find the optimal rate from  $\dim X^{\delta} \approx 10^3$  onwards. Figure 6.6 shows the number of basis functions  $\sigma_{\lambda} \otimes \phi_{\nu}$  whose supports intersect given points in the time-space cylinder. We see the adaptation to the moving peak.

#### 6.7.4 Cylinder problem

Selecting the L-shaped domain  $\Omega := [-1,1]^2 \setminus [-1,0]^2$  with data  $u_0 \equiv 0$  and  $g(t,x,y) := t \cdot \mathbb{1}_{\{x^2+y^2 < 1/4\}}$ , the true solution is known to be singular at the re-entrant corner and at the wall of the cylinder  $\{(t,x,y) : x^2 + y^2 = 1/4\}$ . We took this example from [FK21]. The left side of Figure 6.7 shows the error progression for this cylinder problem. We see that the full grid error decay proportional to  $(\dim X^{\delta})^{-1/4}$  is improved to an error decay proportional to  $(\dim X^{\delta})^{-1/4}$  is error decay proportional to  $(\dim X^{\delta})^{-1/2}$ , recovering the rate for a smooth solution.



FIGURE 6.6. *Moving peak* problem, adaptive lower set with dim  $X^{\delta} = 89401$ . Shown:  $\#\{(\lambda, \nu) \in I_{\delta} : (t, x, y) \in \operatorname{supp} \sigma_{\lambda} \otimes \phi_{\nu}\}$  for a selection of times *t*.

# 6.7.5 Singular problem

We again select the L-shaped domain  $\Omega := [-1, 1]^2 \setminus [-1, 0]^2$  with data  $u_0 \equiv 1$ and  $g \equiv 0$ . The solution has a strong singularity along  $\{0\} \times \partial\Omega$  due to the incompatibility of initial- and boundary conditions, in addition to the singularity at the re-entrant corner (0, 0). At the right of Figure 6.7, for uniform refinement, we see the extremely slow error decay proportional to  $(\dim X^{\delta})^{-1/11}$ , already found in [FK21]. Interestingly, sparse grid refinement offers no rate improvement over full grid refinement. The adaptive algorithm yields a much better error decay proportional to  $(\dim X^{\delta})^{-2/5}$ . We observed that increasing the Dörfler marking parameter to  $\theta = 0.7$  decreases the convergence rate to -1/3, whereas a  $\theta$  smaller than 0.5 did not improve the rate beyond -2/5. Looking at Figure 6.8, we see strong adaptivity towards  $\{0\} \times \partial\Omega$  and  $I \times \{(0,0)\}$ , and observe basis functions  $\sigma_{\lambda} \otimes \phi_{\nu}$  that span  $X^{\delta}$  whose barycenter is at  $t = 2^{-14} \approx 10^{-4}$ .

# 6.7.6 Gradedness and error reduction

In Sect. 6.4 we used (6.24) to demonstrate proportionality of  $\|\mathbf{r}^{\delta}\|$  and  $\|u-u^{\delta}\|_X$ , as well as a constant error reduction in each iteration of the adaptive algorithm. In Proposition 6.6.6, we showed that (6.24) holds when the gradedness  $L_{\delta}$  of Definition 6.6.4 is uniformly bounded.



FIGURE 6.7. Error progressions for (left) the *cylinder* problem and (right) the *singular* problem. Shown: estimated *X*-norm error (solid line), residual norm (dashed), and t = 0 trace error (dotted) as a function of dim  $X^{\delta}$  for adaptive (black), sparse grid (red), and full grid refinement (orange).



FIGURE 6.8. Barycenters of supports of basis functions  $\sigma_{\lambda} \otimes \phi_{\nu}$  spanning  $X^{\delta}$  generated by Algorithm 6.4.8 of dimension 81074 for the *singular* problem. Left: a top-down view, with a 10× zoom to the origin; right: centers in spacetime, logarithmic in time.



FIGURE 6.9. Gradedness and estimated *X*-norm error at every iteration of the adaptive loop, for the four different model problems under consideration.

In the left picture of Figure 6.9, we see however a more than expected increase in gradedness, where in particular for the singular problem we observe a logarithmic increase in terms of dim  $X^{\delta}$ . However, this turns out not to be a problem in practice: Figures 6.5 and 6.7 demonstrate that the residual error  $\|\mathbf{r}^{\delta}\|$  and the estimated *X*-norm error  $\mathcal{E}^{\delta}(\tilde{u}^{\delta})$  are very close, and even converge for the singular problem. Moreover, in the right picture of Figure 6.9, we see a constant error reduction of  $\check{\rho} \approx 0.89$  at every step of the adaptive algorithm, and hence, that the conclusion of Theorem 6.4.9 holds in practice.

#### 6.7.7 Total runtime and memory consumption

Figure 6.10 shows the total runtime and peak memory consumption after every iteration of the adaptive algorithm. The top row shows absolute values, and the bottom row values relative to dim  $X^{\delta}$ .

The left of the figure shows that the adaptive algorithm runs in optimal linear time in the dimension of the current trial space.

The right of the figure shows that the peak memory is linear as well, stabilizing to around 15kB per degree of freedom. This is relatively high, mainly because our implementation uses trees rather than hash maps to represent vectors to ensure a linear-time implementation of the matrix-vector products (cf. Rem. 6.6.2).



FIGURE 6.10. Total runtime and peak memory consumption as function of  $\dim X^{\delta}$ , measured after every iteration of the adaptive loop, for the four different model problems.

# 6.8 Conclusion

We have constructed an adaptive solver for a space-time variational formulation of parabolic evolution problems. The collection of trial spaces are given by the spans of sets of tensor products of wavelets-in-time and hierarchical basis functions-in-space. Compared to our previous works [CS11, RS19] where we employed 'true' wavelets also in space, the theoretical results are weaker. We have demonstrated *r*-linear convergence of the adaptive routine, but have not shown optimal rates at linear complexity. On the other hand, the runtimes that we obtained with the current approach are much better.

# 7.1 Introduction

This chapter is about an *efficient* adaptive method for parabolic evolution equations using a simultaneous space-time variational formulation. Compared to the more classical time-stepping schemes, these space-time methods are very flexible. Among other things, they are especially well-suited for massively parallel computation ([NS19, vVW21a]), and some can guarantee quasi-best approximations from the trial space ([And13, FK21, SZ20]).

We are interested in those space-time methods that permit adaptive refinement locally in space *and* time. Within this class, wavelet-based methods (see [SS09, GK11, KSU16]) are attractive, as they can be shown to be *quasioptimal*: they produce a sequence of solutions that converges at the best possible rate, at optimal linear computational cost. Moreover, they can overcome the *curse of dimensionality* using a form of *sparse tensor-product approximation*, solving the whole time evolution at a runtime proportional to that of solving the corresponding *stationary* problem.

In Chapter 6, we constructed an *r*-linearly converging space-time adaptive solver for parabolic evolution equations that exploits the product structure of the space-time cylinder to construct a family of trial spaces given as the spans of wavelets-in-time tensorized with (locally refined) finite element spaces-in-space.

The principal difference between this and other wavelet-based methods is that we use wavelets in time *only*, and standard finite elements in space. This eases implementation, and alleviates the need for a suitable spatial wavelet basis, which is generally difficult for general domains ([RS18]). Unfortunately, there is no free lunch: a proof of optimal convergence is, for our method, not yet available.

In this chapter we discuss an implementation of the adaptive algorithm from Chapter 6, in which the different steps (each iteration of the linear algebraic solver, the error estimation, Dörfler marking, and refinement of trial- and test spaces) are of linear complexity.

Special care has to be taken for matrix-vector products. For a bilinear form

that is 'local' and equals (a sum of) tensor-product(s) of bilinear forms in time and space, and 'trial' and 'test' spaces spanned by tensor-product multi-level bases with *double-tree* index sets, the resulting system matrix w.r.t. both bases can be applied in linear complexity, even though this matrix is not sparse. The algorithm that realizes this complexity makes a clever use of multi- to single-scale transformations alternately in time and space. This *unidirectional principle* was introduced in [BZ96] for 'uniform' sparse grids, so without 'local refinements', and it was later extended to general *downward closed* or *lower sets*, also called *adaptive sparse grids*, in [KS14]. The definition of a lower set in [KS14], there called multi-tree, is more restrictive than our current definition that allows more localized refinements.

To the best of our knowledge, other implementations for the efficient evaluation of tensor-product bilinear forms (see [Pf110, KS14, Pab15, Rek18]) are based on the concept of hash maps. There, a hash function is used to map basis functions to array indices. In an adaptive loop, the final set of basis functions is unknown in advance so it is impossible to construct a hash function that guarantees an upper bound on the number of hash collisions. Aiming at true linear complexity, we implement these operations by traversing *trees* and *double-trees*, so without the use of hash maps.

### Organization

In §7.2, we look at the abstract parabolic problem, its stable discretization, and the adaptive routine. In §7.3, we provide an abstract algorithm for the efficient evaluation of tensor-product bilinear forms w.r.t. multilevel bases indexed on *double-trees*. In §7.4, we take the *heat equation* as a model problem, and provide a concrete family of trial- and test spaces with bases indexed by double-trees that permits local space-time adaptivity. In §7.5, we discuss the practical implementation of the adaptive algorithm. Finally, in §7.6, we provide extensive numerical experiments to demonstrate the linear runtime of the algorithm.

#### Notation

In this work, by  $C \leq D$  we will mean that C can be bounded by a multiple of D, independently of parameters which C and D may depend on. Obviously,  $C \gtrsim D$  is defined as  $D \leq C$ , and  $C \approx D$  as  $C \leq D$  and  $C \gtrsim D$ .

For normed linear spaces E and F, by  $\mathcal{L}(E, F)$  we will denote the normed linear space of bounded linear mappings  $E \to F$ , and by  $\mathcal{L}is(E, F)$  its subset of boundedly invertible linear mappings  $E \to F$ . We write  $E \hookrightarrow F$  to denote that E is continuously embedded into F. For simplicity only, we exclusively consider linear spaces over the scalar field  $\mathbb{R}$ .

# 7.2 Space-time adaptivity for a parabolic model problem

In this section, we summarize the relevant parts of Chapter 6.

Let *V*, *H* be separable Hilbert spaces of functions on some "spatial domain" such that  $V \hookrightarrow H$  with dense and compact embedding. Identifying *H* with its dual, we obtain the Gelfand triple  $V \hookrightarrow H \simeq H' \hookrightarrow V'$ .

For a.e.

$$t \in I := (0, T),$$

let  $a(t; \cdot, \cdot)$  denote a bilinear form on  $V \times V$  so that for any  $\eta, \zeta \in V, t \mapsto a(t; \eta, \zeta)$  is measurable on I, and such that for a.e.  $t \in I$ ,

$$\begin{aligned} |a(t;\eta,\zeta)| &\lesssim \|\eta\|_V \|\zeta\|_V \quad (\eta,\zeta \in V) \quad (boundedness), \\ a(t;\eta,\eta) &\gtrsim \|\eta\|_V^2 \qquad (\eta \in V) \quad (coercivity). \end{aligned}$$

With  $(A(t)\cdot)(\cdot) := a(t; \cdot, \cdot) \in \mathcal{L}is(V, V')$ , given a forcing function g and initial value  $u_0$ , we want to solve the *parabolic initial value problem* of

(7.1) finding 
$$u: I \to V$$
 such that 
$$\begin{cases} \frac{\mathrm{d}u}{\mathrm{d}t}(t) + A(t)u(t) &= g(t) \quad (t \in I), \\ u(0) &= u_0. \end{cases}$$

*Example* 7.2.1. For the model problem of the *heat equation* on some spatial domain  $\Omega \subset \mathbb{R}^d$  we select  $V := H_0^1(\Omega)$ ,  $H := L_2(\Omega)$ , and  $a(t; \eta, \zeta) := \int_{\Omega} \nabla_{\mathbf{x}} \eta \cdot \nabla_{\mathbf{x}} \zeta \, \mathrm{d} \mathbf{x}$ .

In our simultaneous space-time variational formulation, the parabolic problem is to find u s.t.

$$(Bu)(v) := \int_{I} \langle \frac{\mathrm{d}u}{\mathrm{d}t}(t), v(t) \rangle_{H} + a(t; u(t), v(t)) \mathrm{d}t = \int_{I} \langle g(t), v(t) \rangle_{H} =: g(v)$$

for all *v* from some suitable space of functions of time and space. One possibility to enforce the initial condition is by testing against additional test functions.

**Theorem 7.2.2 ([SS09]).** With  $X := L_2(I; V) \cap H^1(I; V')$ ,  $Y := L_2(I; V)$ , we have

$$\begin{bmatrix} B\\ \gamma_0 \end{bmatrix} \in \mathcal{L}\mathrm{is}(X, Y' \times H),$$

where for  $t \in \overline{I}$ ,  $\gamma_t : u \mapsto u(t, \cdot)$  denotes the trace map. In other words,

(7.2) finding  $u \in X$  s.t.  $(Bu, \gamma_0 u) = (g, u_0)$  given  $(g, u_0) \in Y' \times H$ 

is a well-posed simultaneous space-time variational formulation of (7.1).

We define  $A \in \mathcal{L}is(Y, Y')$  and  $\partial_t \in \mathcal{L}is(X, Y')$  as

$$(Au)(v) := \int_I a(t; u(t), v(t)) dt$$
, and  $\partial_t := B - A$ .

Following [SW21b], we assume that *A* is *self-adjoint*. Morever, in view of an efficient implementation, we assume that *A* is a finite sum of tensor-product operators. If *A* does not have this structure, one may alternatively consider (low-rank) tensor-product approximations of *A*, see e.g. [Hac12] for an overview.

We equip *Y* and *X* with 'energy'-norms

$$\|\cdot\|_{Y}^{2} := (A \cdot)(\cdot), \quad \|\cdot\|_{X}^{2} := \|\partial_{t} \cdot\|_{Y'}^{2} + \|\cdot\|_{Y}^{2} + \|\gamma_{T} \cdot\|_{H}^{2},$$

which are equivalent to the canonical norms on Y and X.

The solution u of (7.2) equals the solution of the following minimization problem

(7.3) 
$$u = \underset{w \in X}{\operatorname{argmin}} \|Bw - g\|_{Y'}^2 + \|\gamma_0 w - u_0\|_H^2,$$

which in turn is the second component of the solution of

(7.4) finding 
$$(\mu, u) \in Y \times X$$
 s.t.  $\begin{bmatrix} A & B \\ B' & -\gamma'_0 \gamma_0 \end{bmatrix} \begin{bmatrix} \mu \\ u \end{bmatrix} = \begin{bmatrix} g \\ -u_0 \end{bmatrix}$ .

Indeed, taking the Schur complement of (7.4) w.r.t. the *Y*-block results in the Euler-Lagrange equations of (7.3).

## 7.2.1 Discretizations

Take a family  $(X^{\delta})_{\delta \in \Delta}$  of closed subspaces of *X*, and define

(7.5) 
$$u_{\delta} = \underset{w \in X^{\delta}}{\operatorname{argmin}} \|Bw - g\|_{Y'}^{2} + \|\gamma_{0}w - u_{0}\|_{H^{2}}^{2}$$

being the best approximation to u from  $X^{\delta}$  w.r.t.  $\|\cdot\|_X$ . Solving this problem, however, is not feasible because of the presence of the dual norm. Therefore, take  $(Y^{\delta})_{\delta \in \Delta}$  to be a family of closed subspaces of Y such that (7.6)

$$X^{\delta} \subseteq Y^{\delta} \quad (\delta \in \Delta), \quad \text{and} \quad \gamma_{\Delta} := \inf_{\delta \in \Delta} \inf_{0 \neq v \in X^{\delta}} \sup_{0 \neq v \in Y^{\delta}} \frac{(\partial_t w)(v)}{\|\partial_t w\|_{Y'} \|v\|_Y} > 0.$$

For  $\underline{\delta} \in \Delta$  with  $Y^{\underline{\delta}} \supseteq Y^{\delta}$ , we replace Y' by  $Y^{\underline{\delta}'}$  in (7.5) yielding the approximation

$$u^{\delta\delta} = \underset{w \in X^{\delta}}{\operatorname{argmin}} \|Bw - g\|_{Y^{\delta'}}^{2} + \|\gamma_{0}w - u_{0}\|_{H}^{2}.$$

Notice that  $u^{\underline{\delta}\delta}$  approximates  $u_{\delta}$  in that  $u^{\underline{\delta}\delta} = u_{\delta}$  when  $Y^{\underline{\delta}} = Y$ .

With  $E_Y^{\delta}: Y^{\delta} \to Y$  and  $E_X^{\delta}: X^{\delta} \to X$  denoting the trivial embeddings,  $u^{\delta\delta}$  is the second component of the solution of

$$\begin{bmatrix} E_Y^{\delta'} A E_Y^{\delta} & E_Y^{\delta'} B E_X^{\delta} \\ E_X^{\delta'} B' E_Y^{\delta} & -E_X^{\delta'} \gamma_0' \gamma_0 E_X^{\delta} \end{bmatrix} \begin{bmatrix} \mu^{\underline{\delta}\delta} \\ u^{\underline{\delta}\delta} \end{bmatrix} = \begin{bmatrix} E_Y^{\delta'} g \\ -E_X^{\delta'} \gamma_0' u_0 \end{bmatrix}.$$

Taking the Schur complement w.r.t. the  $Y^{\underline{\delta}}$ -block then leads to the equation

(7.7) 
$$\begin{split} E_X^{\delta}{}'(B'E_Y^{\delta}(E_Y^{\delta}{}'AE_Y^{\delta})^{-1}E_Y^{\delta}{}'B + \gamma_0'\gamma_0)E_X^{\delta}u^{\delta\delta} \\ &= E_X^{\delta}{}'(B'E_Y^{\delta}(E_Y^{\delta}{}'AE_Y^{\delta})^{-1}E_Y^{\delta}{}'g + \gamma_0'u_0), \end{split}$$

which has a unique solution (cf. Lemma 6.3.3) that satisfies  $||u - u^{\delta \delta}||_X \leq \gamma_{\Delta}^{-1} ||u - u_{\delta}||_X$  whenever  $Y^{\delta} \supseteq Y^{\delta}$ ; cf. [SW21b, Thm. 3.7]. For now, we assume the right-hand side of (7.7) to be evaluated exactly. Later, in §7.4.5, we will discuss approximation of the right-hand side.

In view of obtaining an efficient solver, we want to replace the inverses in (7.7) while aiming to preserve quasi-optimality of the solution. To this end, let  $K_Y^{\delta} = K_Y^{\delta'} \in \mathcal{L}is(Y^{\delta'}, Y^{\delta})$  be a uniformly optimal preconditioner for  $E_Y^{\delta'}AE_Y^{\delta}$  that can be applied in linear complexity. Then, for some  $\kappa_{\Delta} \geq 1$  we have

$$\frac{((K_Y^{\underline{o}})^{-1}v)(v)}{(Av)(v)} \in [\kappa_{\Delta}^{-1}, \kappa_{\Delta}] \quad (\delta \in \Delta, v \in Y^{\underline{\delta}}).$$

Replacing  $(E_Y^{\delta'}AE_Y^{\delta})^{-1}$  by  $K_Y^{\delta}$ , we denote the solution of (7.7) again by  $u^{\delta\delta}$ . It is quasi-optimal with  $||u - u^{\delta\delta}||_X \leq \frac{\kappa_{\Delta}}{\gamma_{\Delta}} ||u - u_{\delta}||_X$ ; cf. [SW21b, Rem. 3.8].

### 7.2.2 Adaptive refinement loop

Our adaptive loop, given in Algorithm 7.1, takes the familiar <u>Solve</u>, <u>Estimate</u>, <u>Mark and refine</u> steps, and is driven by an efficient and reliable 'hierarchical basis' a posteriori error estimator.

The adaptive loop below requires a saturation assumption. Define a *partial* order on  $\Delta$  by  $\tilde{\delta} \succeq \delta$  whenever  $X^{\tilde{\delta}} \supseteq X^{\delta}$ . Let  $\delta \mapsto \underline{\delta} \succeq \delta$  be a mapping providing saturation in that for some  $\zeta < 1$ ,

(7.8) 
$$\|u - u_{\underline{\delta}}\|_X \le \zeta \|u - u_{\delta}\|_X \quad (\delta \in \Delta).$$

With this choice of  $\underline{\delta}$ , we are interested in finding  $u^{\delta} := u^{\underline{\delta}\delta} \in X^{\delta}$  that solves

(7.9) 
$$\underbrace{E_X^{\delta'}(B'E_Y^{\delta}K_Y^{\delta}E_Y^{\delta'}B + \gamma_0'\gamma_0)E_X^{\delta}}_{S^{\delta\delta}:=} u^{\delta} = \underbrace{E_X^{\delta'}(B'E_Y^{\delta}K_Y^{\delta}E_Y^{\delta'}g + \gamma_0'u_0)}_{f^{\delta}:=}.$$

Notice that (7.9) is uniquely solvable even with  $X^{\delta}$  as 'trial space', and we use this 'room' between  $X^{\delta}$  and  $X^{\delta}$  to our advantage. Expanding  $X^{\delta}$  to some intermediate space  $X^{\delta} \subset X^{\tilde{\delta}} \subset X^{\delta}$  yields a  $u^{\tilde{\delta}}$  that is a better approximation to *u* than  $u^{\delta}$ ; cf. Proposition 6.4.2. This function will be the successor of  $u^{\delta}$  in our loop, and we will show that the resulting sequence of functions converges *r*-linearly to *u*; see Algorithm 7.1 and Theorem 7.2.4.

# Solving

Instead of solving the symmetric positive definite system (7.9) exactly, we construct an approximate solution  $\hat{u}^{\delta}$  using Preconditioned Conjugate Gradients (PCG). To this end, let  $K_X^{\delta} = K_X^{\delta'} \in \mathcal{L}is(X^{\delta'}, X^{\delta})$  be a uniformly optimal preconditioner for  $S^{\delta\delta}$ . Then  $((K_X^{\delta})^{-1}w)(w) \approx ||w||_X^2 \approx ||K_X^{\delta}S^{\delta\delta}w||_X^2$  for  $w \in X^{\delta}$ . Writing  $w = K_X^{\delta}S^{\delta\delta}(u^{\delta} - v^{\delta})$  reveals that this induces an algebraic error estimator

(7.10)

$$\beta^{\delta}(v^{\delta}) := \sqrt{(f^{\delta} - S^{\underline{\delta}\delta}v^{\delta})(K_X^{\delta}(f^{\delta} - S^{\underline{\delta}\delta}v^{\delta}))} \approx \|u^{\delta} - v^{\delta}\|_X \quad (v^{\delta} \in X^{\delta}, \delta \in \Delta).$$

With  $\hat{u}_k^{\delta}$  denoting the approximant at iteration k of the PCG loop,  $\beta^{\delta}(\hat{u}_k^{\delta})$  is already available as  $\sqrt{\beta_k}$ , for  $\beta_k$  the variable used in computing the next search direction.

### **Error estimation**

Let  $\Theta_{\delta} := \{\theta_{\lambda} : \lambda \in J_{\delta}\}$  be some uniformly *X*-stable basis satisfying  $X^{\delta} \oplus \operatorname{span} \Theta_{\delta} = X^{\delta}$ , in that

(7.11) 
$$||z + \mathbf{c}^\top \Theta_{\delta}||_X^2 \approx ||z||_X^2 + ||\mathbf{c}||^2 \quad (\mathbf{c} \in \ell_2(J_{\delta}), z \in X^{\delta}, \delta \in \Delta).$$

Define the trivial embedding  $P^{\delta} : X^{\delta} \to X^{\underline{\delta}}$ . Akin to (7.9), we define  $S^{\underline{\delta}\underline{\delta}}$  and  $f^{\underline{\delta}\underline{\delta}}$ , and with it, the residual-based a posteriori error estimator  $\mathbf{r}^{\delta} : X^{\delta} \to \ell_2(J_{\delta})$ , as (7.12)

$$S^{\delta\underline{\delta}} := E_X^{\underline{\delta}'} (B' E_Y^{\underline{\delta}} K_Y^{\underline{\delta}} E_Y^{\underline{\delta}'} B + \gamma_0' \gamma_0) E_X^{\underline{\delta}}, \quad f^{\underline{\delta}\underline{\delta}} := E_X^{\underline{\delta}'} (B' E_Y^{\underline{\delta}} K_Y^{\underline{\delta}} E_Y^{\underline{\delta}'} g + \gamma_0' u_0),$$
$$\mathbf{r}^{\delta} (\hat{u}^{\delta}) := (f^{\underline{\delta}\underline{\delta}} - S^{\underline{\delta}\underline{\delta}} P^{\delta} \hat{u}^{\delta}) (\Theta_{\delta}).$$

For  $\hat{u}^{\delta}$  close to  $u^{\delta}$ , the error estimator  $\|\mathbf{r}^{\delta}(\hat{u}^{\delta})\|$  is reliable and efficient:

**Lemma 7.2.3.** Assume (7.8) and (7.11),  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} < \frac{1}{\zeta}$ , and fix some  $\xi > 0$  small enough. For  $\hat{u}^{\delta} \in X^{\delta}$  satisfying  $\beta(\hat{u}^{\delta}) \le \frac{\xi}{1-\xi} \|\mathbf{r}^{\delta}(\hat{u}^{\delta})\|$ , we have

$$\|\mathbf{r}^{\delta}(\hat{u}^{\delta})\| \approx \|u - \hat{u}^{\delta}\|_X$$
 and  $\|u - \hat{u}^{\delta}\|_X \lesssim \|u - u^{\delta}\|_X$   $(\delta \in \Delta).$ 

*Proof.* For convenience, we write  $\hat{\mathbf{r}}^{\delta} := \mathbf{r}^{\delta}(\hat{u}^{\delta})$  and  $\mathbf{r}^{\delta} := \mathbf{r}^{\delta}(u^{\delta})$ . By (7.8), (7.11) and  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} < \frac{1}{\zeta}$ , Proposition 6.4.4 shows that

(7.13) 
$$\|\mathbf{r}^{\delta}\| \approx \|u - u^{\delta}\|_{X} \quad (\delta \in \Delta).$$

From (7.11) one deduces that  $\|\mathbf{r}^{\delta} - \hat{\mathbf{r}}^{\delta}\| \lesssim \|u^{\delta} - \hat{u}^{\delta}\|_X$ ; cf. (6.25). By assumption, for  $\xi < 1$ , we find  $\beta^{\delta}(\hat{u}^{\delta}) \lesssim \xi \|\hat{\mathbf{r}}^{\delta}\|$ . Combined this reveals

(7.14) 
$$\|\mathbf{r}^{\delta} - \hat{\mathbf{r}}^{\delta}\| \lesssim \|u^{\delta} - \hat{u}^{\delta}\|_{X}^{(7.10)} \beta^{\delta}(\hat{u}^{\delta}) \lesssim \xi \|\hat{\mathbf{r}}^{\delta}\|.$$

Using this, we can show reliability of the estimator by

$$\begin{aligned} \|u - \hat{u}^{\delta}\|_{X} &\leq \|u - u^{\delta}\|_{X} + \|u^{\delta} - \hat{u}^{\delta}\|_{X} \\ \stackrel{(7.13),(7.10)}{\approx} \|\mathbf{r}^{\delta}\| + \beta^{\delta}(\hat{u}^{\delta}) &\leq \|\hat{\mathbf{r}}^{\delta}\| + \|\mathbf{r}^{\delta} - \hat{\mathbf{r}}^{\delta}\| + \beta^{\delta}(\hat{u}^{\delta}) \overset{(7.14)}{\lesssim} \|\hat{\mathbf{r}}^{\delta}\|. \end{aligned}$$

For efficiency of the estimator, we deduce

$$\begin{aligned} \|\hat{\mathbf{r}}^{\delta}\| &\lesssim \|u - u^{\delta}\|_{X} + \|\mathbf{r}^{\delta} - \hat{\mathbf{r}}^{\delta}\| \leq \|u - \hat{u}^{\delta}\|_{X} + \|u^{\delta} - \hat{u}^{\delta}\|_{X} + \|\mathbf{r}^{\delta} - \hat{\mathbf{r}}^{\delta}\| \\ &\lesssim \|u - \hat{u}^{\delta}\|_{X} + \xi \|\hat{\mathbf{r}}^{\delta}\|, \end{aligned}$$

so taking  $\xi$  sufficiently small and kicking back  $\|\hat{\mathbf{r}}^{\delta}\|$  yields

(7.15) 
$$\|\hat{\mathbf{r}}^{\delta}\| \lesssim \|u - \hat{u}^{\delta}\|_{X}.$$

Similarly, from (7.13) and (7.14) it follows that

$$\|\hat{\mathbf{r}}^{\delta}\| \lesssim \|u - u^{\delta}\|_{X}$$

We infer quasi-optimality of  $\hat{u}^{\delta}$  from

$$\|u - \hat{u}^{\delta}\|_{X}^{(7.14)} \lesssim \|u - u^{\delta}\|_{X} + \xi \|\hat{\mathbf{r}}^{\delta}\|_{\lesssim}^{(7.16)} \|u - u^{\delta}\|_{X}.$$

In the solve step, we need to iterate PCG until  $\beta^{\delta}(\hat{u}_{k}^{\delta})/\|\mathbf{r}^{\delta}(\hat{u}_{k}^{\delta})\|$  is small enough. In the algorithm below, this is ensured by the do-while loop which also avoids the (expensive) recomputation of the residual at every PCG iteration.

#### Marking and refinement

Denoting the output of the solve step by  $\hat{u}^{\delta}$ , we drive the adaptive loop by performing Dörfler marking on the residual  $\hat{\mathbf{r}}^{\delta} := \mathbf{r}^{\delta}(\hat{u}^{\delta})$ , i.e., for some  $\theta \in (0,1]$ , we mark the smallest set  $J \subset J_{\delta}$  for which  $\|\hat{\mathbf{r}}^{\delta}\|_{J} \| \geq \theta \|\hat{\mathbf{r}}^{\delta}\|$ . We then construct the smallest  $\tilde{\delta} \succeq \delta$  such that  $X^{\tilde{\delta}}$  contains span  $\Theta_{\delta}|_{J}$ .

**Theorem 7.2.4** (Theorem 6.4.9 with  $\eta = 0$ ). Assume (7.8) and (7.11). For  $\xi$  and  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1$  sufficiently small with  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1 \downarrow 0$  when  $\theta \downarrow 0$ , the sequence of approximations produced by Algorithm 7.1 converges *r*-linearly to *u*, in that after every iteration,  $||u - \hat{u}^{\delta}||_X$  decreases with a factor at least  $\rho < 1$ .

*Remark* 7.2.5. In a practical implementation, to ensure termination, Algorithm 7.1 has to be complemented by an appropriate stopping criterium; cf. Algorithm 6.30.

Algorithm 7.1: Space-time adaptive refinement loop.

**Data:**  $\theta \in (0, 1], \xi \in (0, 1), \delta := \delta_{init} \in \Delta;$   $t_{\delta} := \mathcal{E}^{\delta}(0) = \sqrt{(E_Y^{\delta'}g)(K_Y^{\delta}E_Y^{\delta'}g) + ||u_0||_H^2};$  **repeat Solve: do** Compute  $\hat{u}_*^{\delta} \in X^{\delta}$  with  $\beta^{\delta}(\hat{u}_*^{\delta}) \le t_{\delta}/2;$   $t_{\delta} := \beta^{\delta}(\hat{u}_*^{\delta});$   $e_{\delta} := ||\mathbf{r}^{\delta}(\hat{u}_*^{\delta})|| + t_{\delta};$  **while**  $t_{\delta} > \xi e_{\delta};$   $\hat{u}^{\delta} := \hat{u}_*^{\delta};$  **Estimate:** Set  $\hat{\mathbf{r}}^{\delta} := \mathbf{r}^{\delta}(\hat{u}^{\delta});$  **Mark:** Mark a smallest  $J \subset J_{\delta}$  for which  $||\hat{\mathbf{r}}^{\delta}|_J|| \ge \theta ||\hat{\mathbf{r}}^{\delta}||;$  **Refine:** Determine the smallest  $\tilde{\delta} \in \Delta$  such that  $X^{\tilde{\delta}} \supset X^{\delta} \oplus \operatorname{span} \Theta_{\delta}|_J;$  $t_{\tilde{\delta}} := e_{\delta}, \delta := \tilde{\delta};$ 

*Proof.* For convenience, we denote  $\mathbf{r}^{\delta} := \mathbf{r}^{\delta}(u^{\delta})$  and  $\hat{\mathbf{r}}^{\delta} := \mathbf{r}^{\delta}(\hat{u}^{\delta})$ . The stopping criterium of the solve step ensures that  $\beta^{\delta}(\hat{u}^{\delta}) \leq \xi(\|\hat{\mathbf{r}}^{\delta}\| + \beta^{\delta}(\hat{u}^{\delta}))$ , so for  $\xi < 1$  we are in the situation of Lemma 7.2.3.

We have

$$\|\hat{\mathbf{r}}^{\delta} - \mathbf{r}^{\delta}\| \lesssim \xi \|\hat{\mathbf{r}}^{\delta}\| \le \xi \big(\|\mathbf{r}^{\delta}\| + \|\hat{\mathbf{r}}^{\delta} - \mathbf{r}^{\delta}\|\big),$$

so taking  $\xi$  sufficiently small and kicking back  $\|\hat{\mathbf{r}}^{\delta} - \mathbf{r}^{\delta}\|$  yields

(7.17) 
$$\|\hat{\mathbf{r}}^{\delta} - \mathbf{r}^{\delta}\| \lesssim \xi \|\mathbf{r}^{\delta}\|.$$

After marking, we have  $\|\hat{\mathbf{r}}^{\delta}\| \leq \theta^{-1} \|\hat{\mathbf{r}}^{\delta}|_{J}\|$ , which shows that

$$\|\mathbf{r}^{\delta}\| \stackrel{(7.14)}{\lesssim} \|\hat{\mathbf{r}}^{\delta}\| \lesssim \|\hat{\mathbf{r}}^{\delta}|_{J}\| \le \|\mathbf{r}^{\delta}|_{J}\| + \|\mathbf{r}^{\delta} - \hat{\mathbf{r}}^{\delta}\| \stackrel{(7.17)}{\lesssim} \|\mathbf{r}^{\delta}|_{J}\| + \xi \|\mathbf{r}^{\delta}\|,$$

so for  $\xi$  small enough, kicking back  $\|\mathbf{r}^{\delta}\|$  reveals that for a  $\hat{\theta} > 0$  dependent on  $\theta$ ,

$$\|\mathbf{r}^{\delta}|_{J}\| \geq \hat{\theta} \|\mathbf{r}^{\delta}\|.$$

From Proposition 6.4.3 we now find that, for  $\frac{\kappa_{\Delta}}{\gamma_{\Delta}} - 1 \downarrow 0$  when  $\theta \downarrow 0$ , there is a  $\bar{\rho} < 1$  for which

(7.18) 
$$||u - u^{\delta}||_X \le \bar{\rho} ||u - u^{\delta}||_X.$$

Combining the results shows that

$$\begin{split} \|u - \hat{u}^{\tilde{\delta}}\|_{X} &\leq \|u - u^{\tilde{\delta}}\|_{X} + \|u^{\tilde{\delta}} - \hat{u}^{\tilde{\delta}}\|_{X} \\ &\leq (1 + \mathcal{O}(\xi))\|u - u^{\tilde{\delta}}\|_{X} \\ &\leq (1 + \mathcal{O}(\xi))\bar{\rho}\|u - u^{\delta}\|_{X} \\ &\leq (1 + \mathcal{O}(\xi))\bar{\rho}(\|u - \hat{u}^{\delta}\|_{X} + \|u^{\delta} - \hat{u}^{\delta}\|_{X}) \\ &\leq (1 + \mathcal{O}(\xi))\bar{\rho}(\|u - \hat{u}^{\delta}\|_{X} + \|u^{\delta} - \hat{u}^{\delta}\|_{X}) \\ &\leq \underbrace{(1 + \mathcal{O}(\xi))\bar{\rho}}_{=:\rho}\|u - \hat{u}^{\delta}\|_{X}, \end{split}$$

so for  $\xi$  small enough,  $\rho < 1$ , completing the proof of *r*-linear convergence.  $\Box$ 

#### 7.2.3 Adaptive trial- and test spaces

The convergence rate of our adaptive loop is determined by the approximation properties of the family  $(X^{\delta})_{\delta \in \Delta}$ . We want to construct a family that allows for *local* refinements. Here, the crucial problem is guaranteeing the inf-sup stability condition (7.6). It is known that inf-sup stability is satisfied for *full* tensor-products of (non-uniform) finite element spaces, and in [And13, Prop. 4.2], this result was generalized to families of *sparse* tensor-products. Unfortunately, neither family allows for adaptive refinements both locally in time and space.

In §7.4 we will solve this by first equipping X with a tensor-product of (infinite) bases: a wavelet basis  $\Sigma$  in time, and a hierarchical basis in space. We then construct  $X^{\delta}$  as the span of a (finite) subset of this tensor-product basis, which we grow by adding particular functions.

By imposing a *double-tree* constraint on the index set of the basis of  $X^{\delta}$ , we can apply tensor-product operators in linear complexity; see §7.3. Moreover, this constraint implies that for our model problem the inf-sup condition (7.6) is satisfied and we can construct optimal preconditioners  $K_X^{\delta}$  and  $K_X^{\delta}$ .

# 7.3 The application of linear operators in linear complexity

An efficient implementation of our adaptive loop requires the efficient application of the operators  $E_Y^{\delta'}BE_X^{\delta}$  and  $E_X^{\delta'}\gamma'_0\gamma_0E_X^{\delta}$  appearing in (7.9). Both terms are finite sums of tensor-products of operators in time and space. When we equip our trial and test spaces with tensor-products of multilevel bases, it turns out that we can evaluate these operators in linear complexity.

More precisely, this section will show the abstract result that given

• tensor-products  $\Psi := \Psi^0 \times \Psi^1$ ,  $\breve{\Psi} := \breve{\Psi}^0 \times \breve{\Psi}^1$  of multilevel bases  $\Psi^0$ ,  $\Psi^1$ ,  $\breve{\Psi}^0$ ,  $\breve{\Psi}^1$  indexed by  $\vee^0$ ,  $\vee^1$ ,  $\breve{\vee}^0$ ,  $\breve{\vee}^1$ , and

- (finite) subsets  $\Lambda \subset \vee^0 \times \vee^1$ ,  $\check{\Lambda} \subset \check{\vee}^0 \times \check{\vee}^1$  that are *double-trees*, and
- linear operators  $A_i : \operatorname{span} \Psi^i \to (\operatorname{span} \breve{\Psi}^i)'$  that are *local*  $(i \in \{0, 1\})$ ,

we can apply the matrix  $((A_0 \otimes A_1)\Psi|_{\Lambda})(\check{\Psi}|_{\check{\Lambda}})$  in  $\mathcal{O}(\#\Lambda + \#\check{\Lambda})$  operations even though this matrix is not uniformly sparse.

*Example* 7.3.1. For our model problem,  $\Psi^0$  and  $\check{\Psi}^0$  will be wavelets for  $H^1(I)$  or  $L_2(I)$  in time, and  $\Psi^1 = \check{\Psi}^1$  will be a hierarchical finite element basis for  $H^1_0(\Omega)$  in space. We will apply the result of this section to the operators  $\gamma'_0\gamma_0$  and  $B = \partial_t + A$ .

We will achieve this complexity using a variant of the *unidirectional principle*. Denote with  $I_{\Lambda}$  the extension with zeros of a vector supported on  $\Lambda$  to one on  $\vee^0 \times \vee^1$ , and with  $R_{\Lambda}$  its adjoint; define  $I_{\check{\Lambda}}$  and  $R_{\check{\Lambda}}$  analogously. Define  $\mathbf{A}_i := (A_i \Psi^i)(\check{\Psi}^i)$ . We will split  $\mathbf{A}_0$  in its upper and strictly lower triangular parts  $\mathbf{U}_0$  and  $\mathbf{L}_0$ , so that

$$R_{\check{\mathbf{A}}}(\mathbf{A}_0 \otimes \mathbf{A}_1) I_{\mathbf{A}} = R_{\check{\mathbf{A}}}(\mathbf{L}_0 \otimes \mathrm{Id}) (\mathrm{Id} \otimes \mathbf{A}_1) I_{\mathbf{A}} + R_{\check{\mathbf{A}}}(\mathbf{U}_0 \otimes \mathrm{Id}) (\mathrm{Id} \otimes \mathbf{A}_1) I_{\mathbf{A}}.$$

This in itself is not useful, as  $(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{\Lambda}}$  maps into a vector space which dimension we cannot control. However, the restriction  $R_{\mathbf{\Lambda}}$  gives us elbow room: in Theorem 7.3.13 we construct double-trees  $\Sigma, \Theta$  with  $\#\Sigma + \#\Theta \lesssim \#\mathbf{\Lambda} + \#\mathbf{\Lambda}$  s.t.

(7.19) 
$$\begin{cases} R_{\check{\mathbf{\Lambda}}}(\mathbf{L}_0 \otimes \mathrm{Id})(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{\Lambda}} = R_{\check{\mathbf{\Lambda}}}(\mathbf{L}_0 \otimes \mathrm{Id})R_{\mathbf{\Sigma}}I_{\mathbf{\Sigma}}(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{\Lambda}}, \\ R_{\check{\mathbf{\Lambda}}}(\mathbf{U}_0 \otimes \mathrm{Id})(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{\Lambda}} = R_{\check{\mathbf{\Lambda}}}(\mathbf{U}_0 \otimes \mathrm{Id})R_{\mathbf{\Theta}}I_{\mathbf{\Theta}}(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{\Lambda}}. \end{cases}$$

These right hand sides we *can* apply efficiently, and their application boils down to applications of  $\mathbf{L}_0$ ,  $\mathbf{U}_0$ , and  $\mathbf{A}_1$  in a *single* coordinate direction only. Simple matrix-vector products are inefficient though, as these matrices are again not uniformly sparse. However, by using the properties of a double-tree and the sparsity of the operator in *single scale*, we can evaluate  $\mathbf{U}_0$ ,  $\mathbf{L}_0$  and  $\mathbf{A}_1$  in linear time; see §7.3.1.

We follow the structure of [KS14, §3], which applies the aforementioned idea to *multi-trees* though with a slightly more restrictive definition of a *tree*. For readability, we defer the proofs of Theorems 7.3.7, 7.3.9, 7.3.11, and 7.3.13 to Appendix 7.A.

# 7.3.1 Evaluation of linear operators w.r.t. trees

Let  $\Psi$  be a (multilevel) collection of functions on some domain Q.

*Example* 7.3.2. In our application, Q will be either the time interval I with  $\Psi$  being a collection of wavelets, or the spatial domain  $\Omega$ , in which case  $\Psi$  is a collection of hierarchical basis functions.

Writing  $\Psi = \{\psi_{\lambda} : \lambda \in \vee\}$ , we assume that the  $\psi_{\lambda}$  are *locally supported* in the sense that with  $|\lambda| \in \mathbb{N}_0$  denoting the *level* of  $\lambda$ ,

- (7.20)  $\sup_{\lambda \in \vee} 2^{|\lambda|} \operatorname{diam} \operatorname{supp} \psi_{\lambda} < \infty,$
- (7.21)  $\sup_{\ell \in \mathbb{N}_0} \sup_{x \in Q} \#\{\lambda \in \vee : |\lambda| = \ell \wedge \operatorname{supp} \psi_{\lambda} \cap B(x; 2^{-\ell}) \neq \emptyset\} < \infty.$

We will refer to the functions  $\psi_{\lambda}$  as being *wavelets*, although not necessarily they have vanishing moments or other specific wavelet properties.

For  $\ell \in \mathbb{N}_0$ , and any  $\Lambda \subset \lor$ , we set  $\Lambda_{\ell} := \{\lambda \in \Lambda : |\lambda| = \ell\}$  and  $\Lambda_{\ell\uparrow} := \{\lambda \in \Lambda : |\lambda| \ge \ell\}$ , and write  $\Psi_{\ell} := \Psi|_{\lor_{\ell}}$ .

For  $\ell \in \mathbb{N}_0$ , we assume a collection  $\Phi_{\ell} = \{\phi_{\lambda} : \lambda \in \Delta_{\ell}\}$ , whose members will be referred to as being *scaling functions*, with

- (7.22) span  $\Phi_{\ell+1} \supseteq$  span  $\Phi_{\ell} \cup \Psi_{\ell+1}$ ,  $\Phi_0 = \Psi_0$   $(\Delta_0 := \vee_0)$ ,
- (7.23)  $\sup_{\ell \in \mathbb{N}_0} \sup_{\lambda \in \Delta_{\ell}} 2^{\ell} \operatorname{diam} \operatorname{supp} \phi_{\lambda} < \infty,$
- (7.24)  $\sup_{\ell \in \mathbb{N}_0} \sup_{x \in Q} \#\{\lambda \in \Delta_{\ell} : \operatorname{supp} \phi_{\lambda} \cap B(x; 2^{-\ell}) \neq \emptyset\} < \infty,$
- (7.25)  $\{\phi_{\lambda}|_{\Sigma}: \lambda \in \Delta_{\ell}, \phi_{\lambda}|_{\Sigma} \neq 0\}$  is independent (for all open  $\Sigma \subset Q, \ell \in \mathbb{N}_0$ ).

W.l.o.g. we assume that the index sets  $\Delta_{\ell}$  for different  $\ell$  are mutually disjoint, and set  $\Phi := \bigcup_{\ell \in \mathbb{N}_0} \Phi_{\ell}$  with index set  $\Delta := \bigcup_{\ell \in \mathbb{N}_0} \Delta_{\ell}$ . For  $\lambda \in \Delta$ , we set  $|\lambda| := \ell$ when  $\lambda \in \Delta_{\ell}$ .

Viewing  $\Psi_{\ell}$ ,  $\Phi_{\ell}$  as column vectors, the assumptions we made so far guarantee the existence of matrices  $\mathfrak{p}_{\ell}$ ,  $\mathfrak{q}_{\ell}$  such that

$$\begin{bmatrix} (\Phi_{\ell-1})^\top & (\Psi_{\ell})^\top \end{bmatrix} = (\Phi_{\ell})^\top \begin{bmatrix} \mathfrak{p}_{\ell} & \mathfrak{q}_{\ell} \end{bmatrix},$$

where the number of non-zeros per row and column of  $\mathfrak{p}_{\ell}$  and  $\mathfrak{q}_{\ell}$  is finite, uniformly in the rows and columns and in  $\ell \in \mathbb{N}$  (here also (7.25) has been used). We refer to  $\mathfrak{p}_{\ell}$  as the *prolongation matrix*. Columns of  $\mathfrak{p}_{\ell}$  contain the *mask* of the scaling functions, whereas columns of  $\mathfrak{q}_{\ell}$  contain the mask of the wavelets.

To each  $\lambda \in \forall$  with  $|\lambda| > 0$ , we associate one or more  $\mu \in \forall$  with  $|\mu| = |\lambda| - 1$ and  $|\operatorname{supp} \psi_{\lambda} \cap \operatorname{supp} \psi_{\mu}| > 0$ . We call  $\mu$  a *parent* of  $\lambda$ , and so  $\lambda$  a *child* of  $\mu$ .

To each  $\lambda \in \lor$ , we associate some neighbourhood  $S(\lambda)$  of  $\operatorname{supp} \psi_{\lambda}$ , with diameter  $\leq 2^{-|\lambda|}$ , such that for  $|\lambda| > 0$ ,  $S(\lambda) \subset \bigcup_{\mu \in \operatorname{parent}(\lambda)} S(\mu)$ .

*Remark* 7.3.3. Such a neighborhood always exists even when a child has only one parent. Indeed with  $C := \sup_{\lambda \in \vee} 2^{|\lambda|} \operatorname{diam} \operatorname{supp} \psi_{\lambda}$  and  $S(\lambda) := \{x \in Q : \operatorname{dist}(x, \operatorname{supp} \psi_{\lambda}) < C2^{-|\lambda|}\}$ , for  $\mu$  being a parent of  $\lambda$  and  $x \in S(\lambda)$ ,  $\operatorname{dist}(x, \operatorname{supp} \psi_{\mu}) \leq \operatorname{dist}(x, \operatorname{supp} \psi_{\lambda}) + \operatorname{diam} \operatorname{supp} \psi_{\lambda} < 2C2^{-|\lambda|} = C2^{-|\mu|}$ , i.e.,  $x \in S(\mu)$ .

**Definition 7.3.4** (Tree). A finite  $\Lambda \subset \lor_{\ell\uparrow}$  is called an  $\ell$ -*tree*, or simply a *tree* when  $\ell = 0$ , when for any  $\lambda \in \Lambda$  its parents in  $\lor_{\ell\uparrow}$  are in  $\Lambda$ . This is not a tree in the graph-theoretical sense, but rather one in the sense of a family history tree.

*Example* 7.3.5 (Hierarchical basis in one dimension). Figure 7.1 shows an example multilevel collection  $\Psi$  of functions defined on the interval [0, 1]. Its index set  $\vee_{\mathfrak{I}}$  with parent-child relations is shown left, with a tree  $\Lambda \subset \vee_{\mathfrak{I}}$  visualised in red. This collection is called the *hierarchical basis*. With  $S(\lambda) := \operatorname{supp} \psi_{\lambda}$  for  $\lambda \in \vee_{\mathfrak{I}}$ , the hierarchical basis satisfies conditions mentioned above.



 $\label{eq:index} \text{Index set} \vee_{\mathfrak{I}} \text{ and tree } \Lambda \subset \vee_{\mathfrak{I}} \quad \text{ Multilevel functions } \Psi \quad \text{ Scaling functions } \Phi$ 

FIGURE 7.1. Hierarchical basis for the interval [0, 1].

# A routine eval

Let  $(\Psi, \Phi)$  and  $(\check{\Psi}, \check{\Phi})$  satisfy the conditions of the previous subsection, and let *A*: span  $\Phi \to (\text{span }\check{\Phi})'$  be *local* in that  $(Au)(v) = (Au|_{\text{supp }v})(v)$ . Typically, *A* is a (partial) differential operator in variational form; e.g.  $A \in \mathcal{L}(H^1(I), L_2(I)')$ with  $(Au)(v) = \int_I \frac{\mathrm{d}u}{\mathrm{d}t} v \, \mathrm{d}t$ . For trees  $\Lambda \subset \lor$  and  $\check{\Lambda} \in \check{\lor}$ , we are interested in the efficient application of the matrix  $(A\Psi|_{\Lambda})(\check{\Psi}|_{\check{\Lambda}})$ .

Just for brevity of the following argument, assume  $\Psi = \check{\Psi}$  and  $\Phi = \check{\Phi}$ . The matrix  $(A\Psi|_{\Lambda})(\Psi|_{\Lambda})$  is not uniformly sparse, so a straight-forward matrixvector product is not of linear complexity. However, for  $\Lambda$  a uniform tree up to level  $\ell$ , i.e.  $\Lambda = \{\lambda \in \lor : |\lambda| \leq \ell\}$ , a solution is provided by the multi- to single-scale transform T characterized by  $\Psi|_{\Lambda} = T^{\top} \Phi_{\ell}$  through the equality  $(A\Psi|_{\Lambda})(\Psi|_{\Lambda}) = T^{\top}(A\Phi_{\ell})(\Phi_{\ell})T$ , as the transforms can be applied in linear complexity and the single-scale matrix is uniformly sparse.

For general trees however, we don't have dim  $\Phi_{\ell} \lesssim \dim \Psi|_{\Lambda}$  so the previous approach is not of linear complexity. Clever level-by-level multi-to-singlescale transformations and the prolongation of *only* relevant functions *does* allow applying  $(A\Psi|_{\Lambda})(\breve{\Psi}|_{\breve{\Lambda}})$  in linear complexity; see Algorithm 7.2 below.

On several places the restriction of a vector (of scalars or of functions) to its indices in some subset of the index set should be read as the vector of full length where the entries with indices outside this subset are replaced by zeros. For index sets  $\Delta$  and  $\check{\Delta}$ , matrix  $\mathfrak{m} \in \mathbb{R}^{\#\check{\Delta} \times \#\Delta}$ , and subset  $\Pi \subset \Delta$ , we write  $\operatorname{supp}(\mathfrak{m}, \Pi) \subset \check{\Delta}$  for the index set corresponding to the image of  $\mathfrak{m}$  under  $\{\mathbf{x}|_{\Pi} : \mathbf{x} \in \mathbb{R}^{\#\Delta}\}$ .

**Algorithm 7.2:** Function eval(*A*).

 $\begin{aligned} & \text{Data: } \ell \in \mathbb{N}, \breve{\Pi} \subset \breve{\Delta}_{\ell-1}, \Pi \subset \Delta_{\ell-1}, \ell\text{-trees } \breve{\Lambda} \subset \breve{\vee}_{\ell\uparrow} \text{ and } \Lambda \subset \lor_{\ell\uparrow}, \\ & \mathbf{d} \in \mathbb{R}^{\#\Pi}, \mathbf{c} \in \mathbb{R}^{\#\Lambda}. \end{aligned} \\ & \text{Result: } [\mathbf{e}, \mathbf{f}] \text{ where } \mathbf{e} = (Au)(\breve{\Phi}|_{\breve{\Pi}}), \mathbf{f} = (Au)(\breve{\Psi}|_{\breve{\Lambda}}), \text{ with} \\ & u := \mathbf{d}^{\top} \Phi|_{\Pi} + \mathbf{c}^{\top} \Psi|_{\Lambda}. \end{aligned} \\ & \text{if } \breve{\Pi} \cup \breve{\Lambda} \neq \emptyset \text{ then} \\ & \breve{\Pi}_B := \{\lambda \in \breve{\Pi} : | \operatorname{supp} \breve{\phi}_{\lambda} \cap \cup_{\mu \in \Lambda_{\ell}} S(\mu) | > 0 \}, \breve{\Pi}_A := \breve{\Pi} \setminus \breve{\Pi}_B \\ & \Pi_B := \{\lambda \in \Pi : | \operatorname{supp} \phi_{\lambda} \cap (\cup_{\mu \in \breve{\Lambda}_{\ell}} \breve{S}(\mu) \cup_{\gamma \in \breve{\Pi}_B} \operatorname{supp} \breve{\phi}_{\gamma}) | > 0 \}, \\ & \Pi_A := \Pi \setminus \Pi_B \\ & \breve{\Pi} := \operatorname{supp}(\breve{\mathfrak{p}}_{\ell}, \breve{\Pi}_B) \cup \operatorname{supp}(\breve{\mathfrak{q}}_{\ell}, \breve{\Lambda}_{\ell}) \\ & \underline{\Pi} := \operatorname{supp}(\breve{\mathfrak{p}}_{\ell}, \Pi_B) \cup \operatorname{supp}(\breve{\mathfrak{q}}_{\ell}, \Lambda_{\ell}) \\ & \underline{\mathbf{d}} := \mathfrak{p}_{\ell} \mathbf{d}|_{\Pi_B} + \mathfrak{q}_{\ell} \mathbf{c}|_{\Lambda_{\ell}} \\ & [\underline{\mathbf{e}}, \underline{\mathbf{f}}] := \operatorname{eval}(A)(\ell + 1, \breve{\Pi}, \breve{\Lambda}_{\ell+1\uparrow}, \underline{\Pi}, \Lambda_{\ell+1\uparrow}, \underline{\mathbf{d}}, \mathbf{c}|_{\Lambda_{\ell+1\uparrow}}) \\ & \mathbf{e} = \begin{bmatrix} \mathbf{e}|_{\breve{\Pi}A} \\ & \mathbf{e}|_{\breve{\Pi}B} \end{bmatrix} := \begin{bmatrix} (A\Phi|_{\Pi})(\breve{\Phi}|_{\breve{\Pi}A}) \mathbf{d} \\ & (\breve{\mathfrak{p}}_{\ell}^{\top} \mathbf{e})|_{\breve{\Pi}B} \end{bmatrix} \end{bmatrix} \\ & \mathbf{f} = \begin{bmatrix} \mathbf{f}|_{\breve{\Lambda}_{\ell}} \\ & \mathbf{f}|_{\breve{\Lambda}_{\ell+1\uparrow}} \end{bmatrix} := \begin{bmatrix} (\breve{\mathfrak{q}}_{\ell}^{\top} \mathbf{e})|_{\breve{\Lambda}_{\ell}} \\ & \underline{\mathbf{f}} \end{bmatrix} \end{aligned}$ 

*Remark* 7.3.6. Let  $\Lambda \subset \forall$ ,  $\Lambda \subset \lor$  be trees, and  $\mathbf{c} \in \ell_2(\Lambda)$ , then

$$[\mathbf{e}, \mathbf{f}] \mathrel{\mathop:}= \mathtt{eval}(A)(1, \check{\Lambda}_0, \check{\Lambda}_{1\uparrow}, \Lambda_0, \Lambda_{1\uparrow}, \mathbf{c}|_{\Lambda_0}, \mathbf{c}|_{\Lambda_{1\uparrow}}),$$

satisfies

$$(A\Psi|_{\Lambda})(\breve{\Psi}|_{\breve{\Lambda}})\mathbf{c} = \begin{bmatrix} \mathbf{e} \\ \mathbf{f} \end{bmatrix}.$$

**Theorem 7.3.7.** A call of eval yields the output as specified, at the cost of  $\mathcal{O}(\#\Pi + \#\Lambda + \#\Pi + \#\Lambda)$  operations.

Proof. See Appendix 7.A.

#### Routines evalupp and evallow

Let  $A: \operatorname{span} \Phi \to (\operatorname{span} \Phi)'$  be local and *linear*. Set

$$\mathbf{A} := (A\Psi)(\check{\Psi}) = [(A\psi_{\mu})(\check{\psi}_{\lambda})]_{(\lambda,\mu)\in\check{\vee}\times\check{\vee}}$$

as well as  $\mathbf{U} := [(A\psi_{\mu})(\check{\psi}_{\lambda})]_{|\lambda| \le |\mu|}$  and  $\mathbf{L} := [(A\psi_{\mu})(\check{\psi}_{\lambda})]_{|\lambda| > |\mu|}$  so that  $\mathbf{A} = \mathbf{L} + \mathbf{U}$ . As sketched in the introduction of this section, this splitting is going to be necessary for the application of system matrices in the tensor-product setting; see also (7.19). Algorithms 7.3 and 7.4 below can be used to evaluate  $\mathbf{U}$  and  $\mathbf{L}$  in linear complexity.

 $\begin{aligned} & \textbf{Algorithm 7.3: Function evalupp}(A). \\ & \textbf{Data: } \ell \in \mathbb{N}, \Bar{\Pi} \subset \check{\Delta}_{\ell-1}, \Bar{\Pi} \subset \Delta_{\ell-1}, \ell\text{-trees }\check{\Lambda} \subset \check{\vee}_{\ell\uparrow} \text{ and } \Lambda \subset \vee_{\ell\uparrow}, \\ & \textbf{d} \in \mathbb{R}^{\#\Pi}, \textbf{c} \in \mathbb{R}^{\#\Lambda}. \\ & \textbf{Result: } [\textbf{e}, \textbf{f}] \text{ where } \textbf{e} = (Au)(\check{\Phi}|_{\check{\Pi}}), \textbf{f} = \textbf{U}|_{\check{\Lambda} \times \Lambda} \textbf{c}, \text{ with} \\ & u := \textbf{d}^{\top} \Phi|_{\Pi} + \textbf{c}^{\top} \Psi|_{\Lambda}. \\ & \textbf{if } \Bar{\Pi} \cup \check{\Lambda} \neq \emptyset \text{ then} \\ & \Bar{\Pi}_{B} := \{\lambda \in \Bar{\Pi} : | \sup \breve{\phi}_{\lambda} \cap \cup_{\mu \in \Lambda_{\ell}} S(\mu) | > 0 \}, \Bar{\Pi}_{A} := \Bar{\Pi} \setminus \Bar{\Pi}_{B} \\ & \Bar{\Pi} := \sup p(\check{\mathfrak{p}}_{\ell}, \Bar{\Pi}_{B}) \cup \sup p(\check{\mathfrak{q}}_{\ell}, \Aar{\Lambda}_{\ell}) \\ & \Bar{\Pi} := \sup p(\check{\mathfrak{p}}_{\ell}, \Bar{\Pi}_{B}) \cup \sup p(\check{\mathfrak{q}}_{\ell}, \Aar{\Lambda}_{\ell}) \\ & \Bar{\Pi} := \sup p(\mathfrak{q}_{\ell}, \Lambda_{\ell}) \\ & \Bar{\Pi} := evalupp(A)(\ell+1, \Bar{\Pi}, \Aar{\Lambda}_{\ell+1\uparrow}, \Bar{\Pi}, \Lambda_{\ell+1\uparrow}, \Bar{\Pi}, \mathbf{c}|_{\Lambda_{\ell+1\uparrow}}) \\ & \textbf{e} = \begin{bmatrix} \textbf{e}|_{\Bar{\Pi}_{B}} \\ & \textbf{e}|_{\Bar{\Pi}_{B}} \end{bmatrix} := \begin{bmatrix} (A\Phi|_{\Bar{\Pi}})(\breve{\Phi}|_{\Bar{\Pi}_{B}}) \textbf{d} \\ & (A\Phi|_{\Bar{\Pi}})(\breve{\Phi}|_{\Bar{\Pi}_{B}}) \textbf{d} + (\breve{\mathfrak{p}}_{\ell}^{\top} \textbf{e})|_{\Bar{\Pi}_{B}} \end{bmatrix} \\ & \textbf{f} = \begin{bmatrix} \textbf{f}|_{\Bar{\Lambda}_{\ell}} \\ & \textbf{f}|_{\Bar{\Lambda}_{\ell+1\uparrow}} \end{bmatrix} := \begin{bmatrix} (\breve{\mathfrak{q}}_{\ell}^{\top} \textbf{e})|_{\Bar{\Lambda}_{\ell}} \\ & \textbf{f} \end{bmatrix} \end{aligned}$ 

*Remark* 7.3.8. Let  $\Lambda \subset \forall$ ,  $\Lambda \subset \lor$  be trees, and  $\mathbf{c} \in \ell_2(\Lambda)$ , then

$$[\mathbf{e},\,\mathbf{f}] := \mathtt{evalupp}(A)(1,\check{\Lambda}_0,\check{\Lambda}_{1\uparrow},\Lambda_0,\Lambda_{1\uparrow},\mathbf{c}|_{\Lambda_0},\mathbf{c}|_{\Lambda_{1\uparrow}}),$$

satisfies

$$\left. \mathbf{U} \right|_{\breve{\Lambda} imes \Lambda} \mathbf{c} = egin{bmatrix} \mathbf{e} \ \mathbf{f} \end{bmatrix}$$

**Theorem 7.3.9.** A call of evalupp yields the output as specified, at the cost of  $\mathcal{O}(\#\Pi + \#\Lambda + \#\Pi + \#\Lambda)$  operations.

Proof. See Appendix 7.A.

**Algorithm 7.4:** Function evallow(A).

 $\begin{aligned} & \operatorname{Data:} \ \ell \in \mathbb{N}, \Pi \subset \Delta_{\ell-1}, \ell\text{-trees } \check{\Lambda} \subset \check{\vee}_{\ell\uparrow} \text{ and } \Lambda \subset \vee_{\ell\uparrow}, \mathbf{d} \in \mathbb{R}^{\#\Pi}, \\ & \mathbf{c} \in \mathbb{R}^{\#\Lambda}. \end{aligned}$   $\begin{aligned} & \operatorname{Result:} \mathbf{f} = (A\Phi|_{\Pi})(\check{\Psi}|_{\check{\Lambda}})\mathbf{d} + \mathbf{L}|_{\check{\Lambda} \times \Lambda} \mathbf{c}. \end{aligned}$   $& \operatorname{if } \check{\Pi} \cup \check{\Lambda} \neq \emptyset \text{ then} \\ & \Pi_B := \{\lambda \in \Pi : |\operatorname{supp} \phi_{\lambda} \cap \cup_{\mu \in \check{\Lambda}_{\ell}} \check{S}(\mu)| > 0\}, \\ & \underline{\Pi} := \operatorname{supp}(\mathfrak{p}_{\ell}, \Pi_B) \cup \operatorname{supp}(\mathfrak{q}_{\ell}, \Lambda_{\ell}) \\ & \underline{\Pi}_B := \operatorname{supp}(\mathfrak{p}_{\ell}, \Pi_B) \cup \operatorname{supp}(\mathfrak{q}_{\ell}, \Lambda_{\ell}) \\ & \underline{\Pi}_B := \operatorname{supp}(\check{\mathfrak{p}}_{\ell}, \check{\Lambda}_{\ell}) \\ & \underline{\mathbf{d}} := \mathfrak{p}_{\ell} \mathbf{d}|_{\Pi_B} + \mathfrak{q}_{\ell} \mathbf{c}|_{\Lambda_{\ell}} \\ & \underline{\mathbf{e}} := (A\Phi|_{\underline{\Pi}_B})(\check{\Phi}|_{\underline{\Pi}})\mathfrak{p}_{\ell} \mathbf{d}|_{\Pi_B} \\ & \mathbf{f} = \begin{bmatrix} \mathbf{f}|_{\check{\Lambda}_{\ell}} \\ & \mathbf{f}|_{\check{\Lambda}_{\ell+1\uparrow}} \end{bmatrix} := \begin{bmatrix} (\check{\mathfrak{q}}_{\ell}^{\top}\underline{\mathbf{e}})|_{\check{\Lambda}_{\ell}} \\ & \mathbf{evallow}(A)(\ell+1, \check{\Lambda}_{\ell+1\uparrow}, \underline{\Pi}, \Lambda_{\ell+1\uparrow}, \underline{\mathbf{d}}, \mathbf{c}|_{\Lambda_{\ell+1\uparrow}}) \end{bmatrix} \end{aligned}$ 

*Remark* 7.3.10. Let  $\Lambda \subset \forall$ ,  $\Lambda \subset \lor$  be trees, and  $\mathbf{c} \in \ell_2(\Lambda)$ , then

$$\mathbf{L}|_{\check{\Lambda}\times\Lambda}\mathbf{c} = \texttt{evallow}(A)(1,\Lambda_{1\uparrow},\Lambda_0,\Lambda_{1\uparrow},\mathbf{c}|_{\Lambda_0},\mathbf{c}|_{\Lambda_{1\uparrow}}).$$

**Theorem 7.3.11.** A call of evallow yields the output as specified, at the cost of  $\mathcal{O}(\#\breve{\Lambda} + \#\Pi + \#\Lambda)$  operations.

Proof. See Appendix 7.A.

### 7.3.2 Application of tensor-product operators w.r.t. double-trees

For  $i \in \{0, 1\}$ , let  $A_i$ : span  $\Phi_i \to \text{span } \check{\Phi}'_i$  be local and linear and let

$$\mathbf{A}_i = (A\Psi_i)(\breve{\Psi}_i) = [(A\psi^i_{\mu})(\breve{\psi}^i_{\lambda})]_{\lambda \in \breve{\vee}^i, \mu \in \vee^i} = \mathbf{L}_i + \mathbf{U}_i.$$

where  $\mathbf{U}_i := [(\mathbf{A}_i)_{\lambda,\mu}]_{|\lambda| \leq |\mu|}$  and  $\mathbf{L}_i := [(\mathbf{A}_i)_{\lambda,\mu}]_{|\lambda| > |\mu|}$ . For  $i \in \{0,1\}$ , let  $\neg i := 1 - i$ .

**Definition 7.3.12** (Double-tree). Define the coordinate projector  $P_i(b_0, b_1) := b_i$ . We call  $\mathbf{\Lambda} \subset \{ \breve{\vee}^0 \times \breve{\vee}^1, \lor^0 \times \breve{\vee}^1, \breve{\vee}^0 \times \lor^1, \lor^0 \times \lor^1 \}$ , a *double-tree* when for  $i \in \{0, 1\}$  and any  $\mu \in P_{\neg i} \mathbf{\Lambda}$ , the *fiber* 

$$\mathbf{\Lambda}_{i,\mu} := P_i (P_{\neg i}|_{\mathbf{\Lambda}})^{-1} \{\mu\}$$

is a tree (in  $\check{\lor}^i$  or  $\lor^i$ ), i.e.,  $\Lambda$  is a double-tree when 'frozen' in each of its coordinates, at any value of that coordinate, it is a tree in the remaining coordinate.



FIGURE 7.2. With  $\forall_{\mathfrak{I}}$  from Figure 7.1:  $\forall_{\mathfrak{I}} \times \forall_{\mathfrak{I}}$  in black; a double-tree  $\Lambda \subset \forall_{\mathfrak{I}} \times \forall_{\mathfrak{I}}$  in red; the projection  $P_0\Lambda$  in gray, and a fiber  $\Lambda_{0,\mu}$  for  $\mu \in P_1\Lambda$  in brown.

From  $\Lambda = \bigcup_{\mu \in P_{\neg i} \Lambda} (P_{\neg i}|_{\Lambda})^{-1} \{\mu\}$ , we have  $P_i \Lambda = \bigcup_{\mu \in P_{\neg i} \Lambda} \Lambda_{i,\mu}$ , which, being a union of trees, is a tree itself. See also Figure 7.2.

For a subset  $\lhd$  of a (double) index set  $\Diamond$ , let  $I_{\lhd}^{\Diamond}$  denote the extension operator with zeros of a vector supported on  $\lhd$  to one on  $\Diamond$ , and let  $R_{\lhd}^{\Diamond}$  denotes its (formal) adjoint, being the restriction operator of a vector supported on  $\Diamond$  to one on  $\lhd$ . Since the set  $\Diamond$  will always be clear from the context, we will denote these operators simply by  $I_{\lhd}$  and  $R_{\lhd}$ .

As sketched in the introduction of this section, the pieces are now in place to apply  $R_{\check{\Lambda}}(\mathbf{A}_0 \otimes \mathbf{A}_1)I_{\Lambda}$  in linear complexity.

**Theorem 7.3.13.** Let  $\check{\Lambda} \subset \check{\vee}^0 \times \check{\vee}^1$ ,  $\Lambda \subset \vee^0 \times \vee^1$  be finite double-trees. Then

$$\begin{split} \boldsymbol{\Sigma} &\coloneqq \bigcup_{\boldsymbol{\lambda} \in P_0 \boldsymbol{\Lambda}} \Big( \{\boldsymbol{\lambda}\} \times \bigcup_{\substack{\{\mu \in P_0 \check{\boldsymbol{\Lambda}} : |\mu| = |\boldsymbol{\lambda}| + 1, \ |\check{S}^0(\mu) \cap S^0(\boldsymbol{\lambda})| > 0\}}} \check{\boldsymbol{\Lambda}}_{1,\mu} \Big), \\ \boldsymbol{\Theta} &\coloneqq \bigcup_{\boldsymbol{\lambda} \in P_1 \boldsymbol{\Lambda}} \Big( \{\mu \in P_0 \check{\boldsymbol{\Lambda}} : \exists \gamma \in \boldsymbol{\Lambda}_{0,\boldsymbol{\lambda}} \ s.t. \ |\gamma| = |\mu|, \ |\check{S}^0(\mu) \cap S^0(\gamma)| > 0\} \times \{\boldsymbol{\lambda}\} \Big), \end{split}$$

are double-trees with  $\#\Sigma \lesssim \#\check{\Lambda}$  and  $\#\Theta \lesssim \#\Lambda$ , and

$$R_{\check{\mathbf{\Lambda}}}(\mathbf{A}_0 \otimes \mathbf{A}_1) I_{\mathbf{\Lambda}} = R_{\check{\mathbf{\Lambda}}}(\mathbf{L}_0 \otimes \mathrm{Id}) I_{\mathbf{\Sigma}} R_{\mathbf{\Sigma}}(\mathrm{Id} \otimes \mathbf{A}_1) I_{\mathbf{\Lambda}} + R_{\check{\mathbf{\Lambda}}}(\mathrm{Id} \otimes \mathbf{A}_1) I_{\mathbf{\Theta}} R_{\mathbf{\Theta}}(\mathbf{U}_0 \otimes \mathrm{Id}) I_{\mathbf{\Lambda}}.$$

Proof. See Appendix 7.A.

The application of  $R_{\check{\Lambda}_{0,\mu}} \mathbf{L}_0 \otimes \mathrm{Id} I_{\Sigma}$  boils down to the application of  $R_{\check{\Lambda}_{0,\mu}} \mathbf{L}_0 I_{\Sigma_{0,\mu}}$ for every  $\mu \in P_1 \Sigma \cap P_1 \check{\Lambda}$ . Such an application can be performed in  $\mathcal{O}(\#\check{\Lambda}_{0,\mu} + \#\Sigma_{0,\mu})$  operations by means of a call of evallow( $A_0$ ); see also Algorithm 7.9.

Since  $\sum_{\mu \in \check{\vee}_1} \# \check{\Lambda}_{0,\mu} + \# \Sigma_{0,\mu} = \# \check{\Lambda} + \# \Sigma$ , we conclude that the application of  $R_{\check{\Lambda}}(\mathbf{L}_0 \otimes \mathrm{Id}) I_{\Sigma}$  can be performed in  $\mathcal{O}(\# \check{\Lambda} + \# \Sigma)$  operations.

Similarly, applications of  $R_{\Sigma}(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{\Lambda}}$ ,  $R_{\check{\mathbf{\Lambda}}}(\mathrm{Id} \otimes \mathbf{A}_1)I_{\Theta}$ , and  $R_{\Theta}(\mathbf{U}_0 \otimes \mathrm{Id})I_{\mathbf{\Lambda}}$  using calls of eval $(A_1)$ , eval $(A_1)$ , and evalupp $(A_0)$  respectively, can be done in  $\mathcal{O}(\#\Sigma + \#\Lambda)$ ,  $\mathcal{O}(\#\check{\mathbf{\Lambda}} + \#\Theta)$ , and  $\mathcal{O}(\#\Theta + \#\Lambda)$  operations. From  $\#\Sigma \lesssim \#\check{\mathbf{\Lambda}}$  and  $\#\Theta \lesssim \#\Lambda$  we conclude the following.

**Corollary 7.3.14.** Let  $\check{\Lambda} \subset \check{\vee}^0 \times \check{\vee}^1$ ,  $\Lambda \subset \vee^0 \times \vee^1$  be finite double-trees, then  $R_{\check{\Lambda}}(\mathbf{A}_0 \otimes \mathbf{A}_1)I_{\mathbf{\Lambda}}$  can be applied in  $\mathcal{O}(\#\check{\Lambda} + \#\Lambda)$  operations.

# 7.4 The heat equation and practical realization

In this section, we consider the numerical approximation of the *heat equation* 

(7.26) 
$$\begin{cases} \frac{\mathrm{d}u}{\mathrm{d}t}(t) - (\Delta_{\mathbf{x}}u)(t) &= g(t) \quad (t \in I), \\ u(0) &= u_0. \end{cases}$$

For some bounded domain  $\Omega \subset \mathbb{R}^2$ , we take  $H := L_2(\Omega)$  and  $V := H_0^1(\Omega)$ , so that  $X = L_2(I; H_0^1(\Omega)) \cap H^1(I; H^{-1}(\Omega))$  and  $Y = L_2(I; H_0^1(\Omega))$ . We define

$$a(t;\eta,\zeta) := \int_{\Omega} \nabla \eta \cdot \nabla \zeta \, \mathrm{d}\mathbf{x},$$

and aim to solve the parabolic initial value problem (7.1) numerically. The bilinear forms present in our variational formulation (7.4) satisfy

$$A = M_t \otimes A_{\mathbf{x}}, \quad B = D_t \otimes M_{\mathbf{x}} + A, \quad \text{and} \quad \gamma'_0 \gamma_0 = G_t \otimes M_{\mathbf{x}}$$

where

(7.27) 
$$(M_t v)(w) := \int_I vw \, \mathrm{d}t, \quad (D_t v)(w) := \int_I v'w \, \mathrm{d}t, \quad (G_t v)(w) := v(0)w(0), \\ (A_\mathbf{x}\eta)(\zeta) := \int_\Omega \nabla\eta \cdot \nabla\zeta \, \mathrm{d}\mathbf{x}, \quad (M_\mathbf{x}\eta)(\zeta) := \int_\Omega \eta\zeta \, \mathrm{d}\mathbf{x}.$$

In this section, we first construct suitable tensor-product bases for X and Y which functions are wavelets in time and hierarchical finite element functions in space. We then build our discrete 'trial' and 'test' spaces  $(X^{\delta}, Y^{\delta})_{\delta \in \Delta}$  as the span of subsets of these tensor-product bases. We finish with concrete uniformly optimal preconditioners  $K_X^{\delta}$  and  $K_Y^{\delta}$ , the basis necessary for error estimation in the adaptive loop, and evaluation of the right-hand side of (7.9) using interpolants.

#### 7.4.1 Wavelets in time

We construct piecewise linear wavelet bases  $\Sigma$  for  $H^1(I)$  and  $\Xi$  for  $L_2(I)$ .

#### Basis on the trial side

For  $\Sigma$ , we choose the three-point wavelet basis from [Ste98]; for completeness, we include its construction. For  $\ell \geq 0$ , define the scaling functions as the nodal continuous piecewise linears w.r.t. a uniform partition into  $2^{\ell}$  subintervals, that is  $\Phi_{\ell}^{\Sigma} := \{\phi_{(\ell,n)} : 0 \leq n \leq 2^{\ell}\}$  with  $\phi_{(\ell,n)}(k2^{-\ell}) = \delta_{kn}$  for  $0 \leq k \leq 2^{\ell}$ . Define  $\Sigma_0 := \Phi_0^{\Sigma}$ , and for  $\ell \geq 1$ , define  $\Sigma_{\ell} := \{\sigma_{\lambda} : \lambda := (\ell, n) \text{ with } 0 \leq n < 2^{\ell-1}\}$  with  $\sigma_{\lambda} = \sigma_{(\ell,n)}$  as in the right of Figure 7.3. Note that each  $\sigma_{\lambda}$  is a linear combination of three nodal functions from  $\Phi_{\ell}^{\Sigma}$ , hence the name *three-point wavelet*.

By imposing the parent-child structure

(7.28) 
$$\tilde{\lambda} \triangleleft_{\Sigma} \lambda \iff |\tilde{\lambda}| + 1 = |\lambda| \text{ and } |\operatorname{supp} \sigma_{\lambda} \cap \operatorname{supp} \sigma_{\tilde{\lambda}}| > 0,$$

on any two indices  $\lambda$ ,  $\lambda$ , we get the tree shown left in Figure 7.3.

Define  $\Sigma := \bigcup_{\ell \ge 0} \Sigma_{\ell}, \forall_{\Sigma} := \{\lambda : \sigma_{\lambda} \in \Sigma\}$ , and  $S(\sigma_{\lambda}) := \operatorname{supp} \sigma_{\lambda}$ . We see that  $\Sigma$  satisfies (7.20)–(7.21) and that the  $\Phi_{\ell}^{\Sigma}$  satisfy (7.22)–(7.25). Moreover, one can show that  $\Sigma$  is a Riesz basis for  $L_2(I)$  (cf. [Ste98, Thm. 4.2]), and that  $\{2^{-|\lambda|}\sigma_{\lambda}\}$  is a Riesz basis for  $H^1(I)$  (cf. [Ste98, Thm. 4.3]).



FIGURE 7.3. Left: three-point wavelet index set  $\vee_{\Sigma}$  with parent-child relations; right: three-point wavelets.

#### Basis on the test side

We construct an  $L_2(I)$ -orthonormal basis  $\Xi$ .

For  $\ell \geq 0$ , define the (discontinuous) piecewise linear scaling functions w.r.t. a uniform partition into  $2^{\ell}$  subintervals by  $\Phi_{\ell}^{\Xi} := \{\phi_{(\ell,n)} : 0 \leq n < 2^{\ell+1}\}$ where  $\phi_{(0,0)}(t) := \mathbb{1}_{[0,1]}(t)$  and  $\phi_{(0,1)}(t) := \sqrt{3}(2t-1)\mathbb{1}_{[0,1]}$ , and for  $\ell \geq 1$ ,  $\phi_{(\ell,2k)}(t) := \phi_{(0,0)}(2^{\ell}t - k)$  and  $\phi_{(\ell,2k+1)}(t) := \phi_{(0,1)}(2^{\ell}t - k)$ . Let  $\Xi_0 := \Phi_0^{\Xi}$ , and define  $\Xi_1 := \{\xi_{(1,0)}, \xi_{(1,1)}\}$  as in the right of Figure 7.4. For  $\ell \geq 2$ , we take  $\Xi_{\ell} := \{\xi_{(\ell,n)} : 0 \leq n < 2^{\ell}\}$  with

$$\xi_{(\ell,2k)}(t) := 2^{(\ell-1)/2} \xi_{(1,0)}(2^{\ell-1}t-k), \quad \xi_{(\ell,2k+1)}(t) := 2^{(\ell-1)/2} \xi_{(1,1)}(2^{\ell-1}t-k).$$

The resulting  $\Xi := \bigcup_{\ell \ge 0} \Xi_{\ell}$  is an *orthonormal* basis for  $L_2(I)$ , and together with its scaling functions  $\bigcup_{\ell} \Phi_{\ell}^{\Xi}$ , the conditions from §7.3.1 are satisfied with  $S(\xi_{\mu}) := \operatorname{supp} \xi_{\mu}$ . We impose a parent-child relation analogously to (7.28); see the left of Figure 7.4.



FIGURE 7.4. Left: orthonormal wavelet index set  $\lor_{\Xi}$  with parent-child relations; right: the wavelets at levels 0 and 1.

## 7.4.2 Finite elements in space

Let  $\mathbb{T}$  be the family of all *conforming* partitions of  $\Omega$  into triangles that can be created by Newest Vertex Bisection from some given conforming initial triangulation  $\mathcal{T}_{\perp}$  with an assignment of newest vertices satisfying the matching condition; cf. [Ste08b].

Define  $\mathfrak{T} := \bigcup_{\mathcal{T} \in \mathbb{T}} \{T : T \in \mathcal{T}\}$ . For  $T \in \mathfrak{T}$ , set gen(T) as the number of bisections needed to create T from its 'ancestor'  $T' \in \mathcal{T}_{\perp}$ . With  $\mathfrak{N}$  the set of all vertices of all  $T \in \mathfrak{T}$ , for  $\nu \in \mathfrak{N}$  we set  $gen(\nu) := \min\{gen(T) :$  $\nu$  is a vertex of  $T \in \mathfrak{T}\}$ .

Any  $\nu \in \mathfrak{N}$  with  $gen(\nu) > 0$  is the midpoint of an edge  $e_{\nu}$  of one or two  $T \in \mathfrak{T}$  with  $gen(T) = gen(\nu) - 1$ . The set of newest vertices  $\tilde{\nu}$  of these T, so those vertices of T with  $|\tilde{\nu}| = gen(\nu) - 1$ , are defined as the parents of  $\nu$ , denoted  $\tilde{\nu} \triangleleft_{\mathfrak{N}} \nu$ . The set of *godparents* of  $\nu$ , denoted  $gp(\nu)$ , are defined as the two endpoints of  $e_{\nu}$ . Vertices with  $gen(\nu) = 0$  have no parents or godparents.

*Example* 7.4.1. In Figure 7.5, the parents of  $\nu_4$  are  $\nu_1$  and  $\nu_3$  and its godparents are  $\nu_0$ ,  $\nu_2$ ; the sole parent of  $\nu_5$  is  $\nu_4$ , and its godparents are  $\nu_0$  and  $\nu_3$ .

**Proposition 7.4.2** ([DKS16]). An (essentially) non-overlapping partition  $\mathcal{T}$  of  $\overline{\Omega}$  into triangles is in  $\mathbb{T}$  if and only if the set  $N_{\mathcal{T}}$  of vertices of all  $T \in \mathcal{T}$  forms a tree in the sense of §7.3.1, meaning that it contains every vertex of generation zero as well as all parents of any  $\nu \in N_{\mathcal{T}}$ ; see also Figure 7.5.

Let  $\mathcal{O}$  be the collection of spaces  $W_{\mathcal{T}}$  of continuous piecewise linears w.r.t.  $\mathcal{T} \in \mathbb{T}$  vanishing on  $\partial \Omega$ . For  $\nu \in \mathfrak{N}$ , we set  $\psi_{\nu}$  as that continuous piecewise linear function on the *uniform partition*  $\mathcal{T}_{\nu} := \{T \in \mathfrak{T} : \text{gen}(T) = \text{gen}(\nu)\} \in \mathbb{T}$ 



FIGURE 7.5. Vertex tree  $N_{\mathcal{T}}$  and its triangulation  $\mathcal{T}$  shown level-by-level.

for which  $\psi_{\nu}(\tilde{\nu}) = \delta_{\nu\tilde{\nu}}$  for  $\tilde{\nu} \in \mathcal{T}_{\nu}$ . Setting  $\mathfrak{N}_{0} := \mathfrak{N} \setminus \partial\Omega$ , the collection  $\{\psi_{\nu} : \nu \in \mathfrak{N}_{0}\}$  is known as the *hierarchical basis*. For  $\mathcal{T} \in \mathbb{T}$ , write  $N_{\mathcal{T},0} := N_{\mathcal{T}} \setminus \partial\Omega$  and  $\Psi_{\mathcal{T}} := \{\psi_{\nu} : \nu \in N_{\mathcal{T},0}\}$ ; it holds that  $W_{\mathcal{T}} = \operatorname{span} \Psi_{\mathcal{T}}$ .

#### **Applying stiffness matrices**

The hierarchical basis satisfies conditions (7.20) and (7.21), and so, the application of stiffness matrices  $(A\Psi_{\mathcal{T}})(\Psi_{\mathcal{T}})$  for  $A \in \{A_{\mathbf{x}}, M_{\mathbf{x}}\}$  can be done through  $\operatorname{eval}(A)$ .<sup>1</sup> However, the computation in Theorem 7.3.13 does not involve the lower and upper parts of A. This crucial insight allows for a faster and easier approach using standard finite element techniques:  $\operatorname{span}\Psi_{\mathcal{T}}$  is a continuous piecewise linear finite element space, so it has a *canonical* single-scale basis  $\Phi_{\mathcal{T}} := \operatorname{span}\{\phi_{\mathcal{T},\nu}\}$  characterized by  $\phi_{\mathcal{T},\nu}(\tilde{\nu}) = \delta_{\nu\tilde{\nu}}$  for  $\tilde{\nu} \in N_{\mathcal{T},0}$ , for which the application of  $(A\Phi_{\mathcal{T}})(\Phi_{\mathcal{T}})$  at linear cost using local element matrices is standard. This is different from the general setting in §7.3.1, in that  $\dim \Phi_{\mathcal{T}} = \dim \Psi_{\mathcal{T}}$  also for locally refined triangulations. Let T be the transformation characterized by  $\Psi_{\mathcal{T}} = T^{\top} \Phi_{\mathcal{T}}$ , we find

(7.29) 
$$(A\Psi_{\mathcal{T}})(\Psi_{\mathcal{T}}) = T^{\top}(A\Phi_{\mathcal{T}})(\Phi_{\mathcal{T}})T.$$

We can apply T in linear complexity by iterating over the vertices bottom-up while applying elementary local transformations in which not parent-child, but *godparent*-child relations play a role.

## 7.4.3 Inf-sup stable family of trial- and test spaces

With  $\underline{\Sigma}$  and  $\underline{\Xi}$  from §7.4.1 and  $\underline{\Psi}_{\mathfrak{N}_0} := \{\psi_{\nu} : \nu \in \mathfrak{N}_0\}$  from §7.4.2, we find that  $X = \overline{\operatorname{span}(\Sigma \otimes \Psi_{\mathfrak{N}_0})}$  and  $Y = \overline{\operatorname{span}(\Xi \otimes \Psi_{\mathfrak{N}_0})}$ . We now turn to the construction of  $X^{\delta}$  and  $Y^{\delta}$ .

<sup>&</sup>lt;sup>1</sup>This would require the definition of a suitable single-scale basis.

**Definition 7.4.3.** For a double-tree  $\Lambda^{\delta} \subset \vee_{\Sigma} \times \mathfrak{N}$ , define  $\Lambda_{0}^{\delta} := \Lambda^{\delta} \setminus \vee_{\Sigma} \times \partial \Omega$ . We construct our 'trial' space as

$$X^{\delta} := \operatorname{span} \{ \sigma_{\lambda} \otimes \psi_{\nu} : (\lambda, \nu) \in \mathbf{\Lambda}_{0}^{\delta} \}.$$

Defining the double-tree  $\Lambda_{Y,0}^{\delta} \subset \vee_{\Xi} \times \mathfrak{N}_0$  as

$$\mathbf{\Lambda}^{\delta}_{Y,0} := \{(\mu,\nu) : \exists (\lambda,\nu) \in \mathbf{\Lambda}^{\delta}_{0}, \, \mu \in \vee_{\Xi}, \, |\mu| = |\lambda|, \, |\operatorname{supp} \xi_{\mu} \cap \operatorname{supp} \sigma_{\lambda}| > 0\},$$

we construct our 'test' space as  $Y^{\delta} = Y^{\delta}(X^{\delta}) := \operatorname{span}\{\xi_{\mu} \otimes \psi_{\nu} : (\mu, \nu) \in \mathbf{\Lambda}_{Y,0}^{\delta}\}.$ 

**Theorem 7.4.4.** Define  $\Delta := \{\delta : \Lambda^{\delta} \subset \bigvee_{\Sigma} \times \mathfrak{N} \text{ is a double-tree}\}$  equipped with the partial ordering  $\delta \preceq \tilde{\delta} \iff \Lambda^{\delta} \subseteq \Lambda^{\tilde{\delta}}$ . With  $X^{\delta}$  and  $Y^{\delta}$  as above, uniform inf-sup stability holds; cf. (7.6).

Proof. See Proposition 6.5.2.

**Definition 7.4.5.** Given a double-tree  $\Lambda^{\delta} \subset \vee_{\Sigma} \times \mathfrak{N}$ , we define  $\Lambda^{\delta} \supset \Lambda^{\delta}$  by adding, for  $(\lambda, \nu) \in \Lambda^{\delta}$  and any child  $\tilde{\lambda}$  of  $\lambda$  and descendant  $\tilde{\nu}$  of  $\nu$  up to generation 2, all pairs  $(\tilde{\lambda}, \nu)$  and  $(\lambda, \tilde{\nu})$ . We expect this choice of  $X^{\delta}$  to provide saturation; cf. (7.8).

## 7.4.4 Preconditioners

We follow §6.5.6 for the construction of optimal preconditioners  $K_Y^{\delta}$  for  $E_Y^{\delta'} A E_Y^{\delta}$ and  $K_X^{\delta}$  for  $S^{\delta\delta}$  necessary for solving (7.9). With notation from Definition 7.3.12, we equip  $X^{\delta}$  and  $Y^{\delta}$  with bases

$$\begin{cases} \bigcup_{\lambda \in P_0 \mathbf{\Lambda}_0^{\delta}} \sigma_{\lambda} \otimes \Psi_{\lambda}^{\delta} \quad \text{with} \quad \Psi_{\lambda}^{\delta} \coloneqq \{\psi_{\nu} : \nu \in (\mathbf{\Lambda}_0^{\delta})_{1,\lambda}\}, \\ \bigcup_{\mu \in P_0 \mathbf{\Lambda}_{Y,0}^{\delta}} \xi_{\mu} \otimes \Psi_{\mu}^{\delta} \quad \text{with} \quad \Psi_{\mu}^{\delta} \coloneqq \{\psi_{\nu} : \nu \in (\mathbf{\Lambda}_{Y,0}^{\delta})_{1,\mu}\}.\end{cases}$$

Matrix representations of preconditioners from §6.5.6 are then given by

$$\begin{cases} \mathbf{K}_{Y}^{\delta} \coloneqq \operatorname{blockdiag}[\mathbf{K}_{\mu}^{\delta}]_{\mu \in P_{0} \mathbf{\Lambda}_{Y,0}^{\delta}} & \text{where} \quad \mathbf{K}_{\mu}^{\delta} \eqsim (\mathbf{A}_{\mu}^{\delta})^{-1}, \\ \mathbf{K}_{X}^{\delta} \coloneqq \operatorname{blockdiag}[\mathbf{K}_{\lambda}^{\delta} \mathbf{A}_{\lambda}^{\delta} \mathbf{K}_{\lambda}^{\delta}]_{\lambda \in P_{0} \mathbf{\Lambda}_{0}^{\delta}} & \text{where} \quad \mathbf{K}_{\lambda}^{\delta} \eqsim (\mathbf{A}_{\lambda}^{\delta} + 2^{|\lambda|} \mathbf{M}_{\lambda}^{\delta})^{-1} \end{cases}$$

with  $\mathbf{A}^{\delta}_{\mu} := (A_{\mathbf{x}} \Psi^{\delta}_{\mu})(\Psi^{\delta}_{\mu})$ ,  $\mathbf{A}^{\delta}_{\lambda} := (A_{\mathbf{x}} \Psi^{\delta}_{\lambda})(\Psi^{\delta}_{\lambda})$ , and  $\mathbf{M}^{\delta}_{\lambda} := (M_{\mathbf{x}} \Psi^{\delta}_{\lambda})(\Psi^{\delta}_{\lambda})$ . Suitable spatial preconditioners  $\mathbf{K}^{\delta}_{\mu}$  are provided by multigrid methods. In [OR00] it was shown that for quasi-uniform triangulations, satisfying a 'full-regularity' assumption, a multiplicative multigrid method yields suitable  $\mathbf{K}^{\delta}_{\lambda}$ , and we assume these results to hold for our locally refined triangulations  $\mathcal{T} \in \mathbb{T}$  as well. In §7.5.1 below, we detail our linear-complexity multigrid implementation following [WZ17].

#### 7.4.5 Right-hand side

We follow §6.6.4. For  $g \in C(\overline{I \times \Omega})$ ,  $u_0 \in C(\overline{\Omega})$ , we can approximate the right-hand side of (7.9) by interpolants, avoiding quadrature issues.

The procedure of §7.4.2 for constructing the hierarchical basis  $\Psi_{\mathfrak{N}} := \{\psi_{\nu} : \nu \in \mathfrak{N}\}\$  can be applied in time as well, yielding the basis  $\{\psi_{\lambda} : \lambda \in \vee_{\mathfrak{I}}\}\$  from Figure 7.1 which index set  $\vee_{\mathfrak{I}}$  coincides with  $\vee_{\Sigma}$ . We construct  $\{\tilde{\psi}_{\nu} : \nu \in \mathfrak{N}\} \subset C(\overline{\Omega})'$  biorthogonal to  $\Psi_{\mathfrak{N}}$ , with  $\tilde{\psi}_{\nu} := \delta_{\nu} - \sum_{\tilde{\nu} \in \mathrm{gp}(\nu)} \delta_{\tilde{\nu}}/2$ . In time, define  $\{\tilde{\psi}_{\lambda} : \lambda \in \vee_{\mathfrak{I}}\} \subset C(\overline{I})'$  analogously. Define the vectors  $\mathbf{g} := [(\tilde{\psi}_{\lambda} \otimes \tilde{\psi}_{\nu})(g)]_{(\lambda,\nu) \in \mathbf{\Lambda}^{\delta}}$  and  $\mathbf{u}_{0} := [\tilde{\psi}_{\nu}(u_{0})]_{\nu \in P_{1}} \mathbf{\Lambda}^{\delta}$ . Upon replacing  $(g, u_{0})$  in (7.9) by the interpolants

$${}^{\delta}g := \sum_{(\lambda,\nu)\in\mathbf{\Lambda}^{\delta}} \mathbf{g}_{(\lambda,\nu)}\psi_{\lambda}\otimes\psi_{\nu}, \quad {}^{\delta}u_{0} := \sum_{\nu\in P_{1}\mathbf{\Lambda}^{\delta}} \mathbf{u}_{0,\nu}\psi_{\nu},$$

we can evaluate its right-hand side in linear complexity by computing the quantities

$$\begin{split} [\langle \xi_{\mu} \otimes \psi_{\nu}, {}^{\delta}g \rangle_{L_{2}(I \times \Omega)}]_{(\mu,\nu) \in \mathbf{\Lambda}_{Y,0}^{\delta}} &= R_{\mathbf{\Lambda}_{Y,0}^{\delta}}(M_{t} \otimes M_{\mathbf{x}})I_{\mathbf{\Lambda}^{\delta}}\mathbf{g}, \\ [\sigma_{\lambda}(0) \langle \psi_{\nu}, {}^{\delta}u_{0} \rangle_{L_{2}(\Omega)}]_{(\lambda,\nu) \in \mathbf{\Lambda}_{0}^{\delta}} &= [\sigma_{\lambda}(0)\mathbf{w}_{\nu}]_{(\lambda,\nu) \in \mathbf{\Lambda}_{0}^{\delta}} \\ \mathbf{w}here \quad \mathbf{w} &= (M_{\mathbf{x}}\Psi_{\mathfrak{N}}|_{P_{1}\mathbf{\Lambda}^{\delta}})(\Psi_{\mathfrak{N}}|_{P_{1}\mathbf{\Lambda}^{\delta}})\mathbf{u}_{0} \end{split}$$

### 7.4.6 Two-level basis

We now discuss the construction of a uniformly X-stable basis  $\Theta_{\delta}$ , needed in the local error estimator  $\mathbf{r}^{\delta}$  of (7.12). Following §6.6.3, define a *modified hierarchical basis* { $\hat{\psi}_{\nu} : \nu \in \mathfrak{N}_0$ } by

$$\hat{\psi}_{\nu} = \psi_{\nu} \text{ when } \operatorname{gen}(\nu) = 0, \quad \text{else } \quad \hat{\psi}_{\nu} \coloneqq \psi_{\nu} - \frac{\sum_{\{\tilde{\nu} \in \mathfrak{N} : \tilde{\nu} \triangleleft_{\mathfrak{N}}\nu\}} \frac{\int_{\Omega} \psi_{\nu} \, \mathrm{d}\mathbf{x}}{\int_{\Omega} \psi_{\bar{\nu}} \, \mathrm{d}\mathbf{x}} \psi_{\bar{\nu}}}{\#\{\tilde{\nu} \in \mathfrak{N} : \tilde{\nu} \triangleleft_{\mathfrak{N}}\nu\}}.$$

For any  $\mathcal{T} \in \mathbb{T}$ ,  $W_{\mathcal{T}} = \operatorname{span}\{\hat{\psi}_{\nu} : \nu \in N_{\mathcal{T},0}\} = \operatorname{span}\Psi_{\mathcal{T}}$  and the transformation from modified to unmodified hierarchical basis can be performed in linear complexity. For  $\underline{\mathcal{T}} \succeq \mathcal{T} \in \mathbb{T}$ ,  $\mathbf{d} \in \ell_2(N_{\underline{\mathcal{T}},0} \setminus N_{\mathcal{T},0})$  and  $v \in W_{\mathcal{T}}$ , Lemma 6.6.7 shows that

(7.30) 
$$\begin{cases} \|v + \sum_{\nu} \mathbf{d}_{\nu} \hat{\psi}_{\nu}\|_{H^{1}(\Omega)}^{2} \approx \|v\|_{H^{1}(\Omega)}^{2} + \|\mathbf{d}\|^{2}, \\ \|v + \sum_{\nu} \mathbf{d}_{\nu} \hat{\psi}_{\nu}\|_{H^{-1}(\Omega)}^{2} \approx \|v\|_{H^{-1}(\Omega)}^{2} + \sum_{\nu} 4^{-\operatorname{gen}(\nu)} |\mathbf{d}_{\nu}|^{2}, \end{cases}$$

with the constants in the  $\approx$ -symbols dependent on  $\max_{\{\underline{T} \ni \underline{T} \subset T \in \mathcal{T}\}} \{\operatorname{gen}(\underline{T}) - \operatorname{gen}(T)\}$  only. We then construct a basis for  $X^{\underline{\delta}} \ominus X^{\delta}$  as

$$\Theta_{\delta} := \{ e_{\lambda\nu} \sigma_{\lambda} \otimes \hat{\psi}_{\nu} : (\lambda, \nu) \in \mathbf{\Lambda}_{0}^{\underline{\delta}} \setminus \mathbf{\Lambda}_{0}^{\underline{\delta}} \} \quad \text{where} \quad \frac{1}{e_{\lambda\nu}} = \sqrt{1 + 4^{|\lambda| - \text{gen}(\nu)}}.$$

Define the gradedness of a double-tree  $\Lambda^{\delta} \subset \bigvee_{\Sigma} \times \mathfrak{N}$  as the smallest  $L_{\delta} \in \mathbb{N}$ for which every  $(\lambda, \nu) \in \Lambda^{\delta}$  with  $\tilde{\nu}$  an ancestor of  $\nu$  with  $gen(\nu) - gen(\tilde{\nu}) = L_{\delta}$ , it holds that  $(\check{\lambda}, \tilde{\nu}) \in \Lambda^{\delta}$  for all  $\check{\lambda} \triangleleft_{\Sigma} \lambda$ . Thanks to  $\Sigma$  being a (scaled) Riesz basis for  $L_2(I)$  and  $H^1(I)$ , together with the  $H^1(\Omega)$ - and  $H^{-1}(\Omega)$ -stable splittings of (7.30), it holds that

$$\|z + \mathbf{c}^{ op} \Theta_{\delta}\|_X^2 pprox \|z\|_X^2 + \|\mathbf{c}\|^2 \quad (\mathbf{c} \in \ell_2(\mathbf{\Lambda}_0^{\delta} \setminus \mathbf{\Lambda}_0^{\delta}), z \in X^{\delta}),$$

with the constant in the  $\approx$ -symbol dependent on  $L_{\delta}$  only, so when  $L_{\delta}$  is uniformly bounded, condition (7.11) is satisfied.

# 7.5 Implementation

A tree-based implementation of the aforementioned adaptive algorithm in C++ can be found at [vVW21d]. In this section, we describe our design choices for a linear complexity implementation.

## 7.5.1 Trees and linear operators in one axis

In §7.3, we consider an abstract multilevel collection  $\Psi$  indexed on  $\vee_{\Psi}$ . Endowed with a parent-child relation,  $\vee_{\Psi}$  has a tree-like structure that we call a *mother tree*; see also Figures 7.3 and 7.4.

In our applications, the support of a wavelet  $\psi_{\lambda}$  is a union of simplices of generation  $|\lambda|$ . In time, these simplices are subintervals of *I* found by dyadic refinement. In space, they are elements of  $\mathfrak{T}$ , the collection of all triangles found by newest vertex bisection. Endowed with the natural parent-child relation, both collections of simplices have a tree structure we call the *domain mother tree*. Every wavelet  $\psi_{\lambda}$  stores references to the simplices *T* of generation  $|\lambda|$  that make up its support; conversely, every *T* stores a reference to  $\psi_{\lambda}$ .

Every mother tree  $\lor$  is stored once in memory, and every node  $\lambda \in \lor$  stores references to its parents, children, and siblings. We treat the mother tree as infinite by *lazy initialization*, constructing new nodes as they are needed.

#### Trees

We store a tree  $\Lambda \subset \lor$  using the parent-child relation, and additionally, at each  $\lambda \in \Lambda$  store a reference to the corresponding node in  $\lor$ . This allows us to compare different trees subject to the same mother tree. This tree-like representation does not allow direct access of arbitrary nodes: in any operation, we traverse  $\Lambda$  from its roots in breadth-first, or level-wise, order.

#### **Tree operations**

One important operation is the union of one tree  $\Lambda$  into another  $\Lambda$ . This can be implemented by traversing both trees simultaneously in breadth-first order.
The union allows us to easily perform high-level operations, such as vector addition: given two vectors  $\mathbf{c} \in \ell_2(\Lambda)$ ,  $\mathbf{d} \in \ell_2(\check{\Lambda})$  on the same mother tree  $\lor$ , we use the union to perform  $\mathbf{c} := \mathbf{c} + \mathbf{d}$ . See Figure 7.6 for an example.



FIGURE 7.6. Left:  $\mathbf{c} \in \ell_2(\Lambda)$  for  $\Lambda \subset \vee_{\mathfrak{I}}$ ; Middle:  $\mathbf{d} \in \ell_2(\Lambda)$  for  $\Lambda \subset \vee_{\mathfrak{I}}$ ; Right: in-place sum  $\mathbf{c} := \mathbf{c} + \mathbf{d}$ .

#### Tree operations in time

The routines eval, evalupp, and evallow from §7.3.1 involve various level-wise index sets (represented as arrays of references into their mother trees). One example is  $\breve{\Pi}_B = \{\lambda \in \breve{\Pi} : | \operatorname{supp} \breve{\phi}_{\lambda} \cap \cup_{\mu \in \Lambda_{\ell}} S(\mu) | > 0 \}$ , which we constructed efficiently using the domain mother tree; see Algorithm 7.5.

```
Algorithm 7.5: The construction of \breve{\Pi}_B.Data: \ell \in \mathbb{N}, \breve{\Pi} \subset \check{\Delta}_{\ell-1}, \Lambda_\ell \subset \vee_\ell.Result: [\breve{\Pi}_A, \breve{\Pi}_B] where \breve{\Pi}_A = \breve{\Pi} \setminus \breve{\Pi}_B,\breve{\Pi}_B = \{\lambda \in \breve{\Pi} : | \operatorname{supp} \breve{\phi}_\lambda \cap \cup_{\mu \in \Lambda_\ell} S(\mu) | > 0 \}.\breve{\Pi}_A := \emptyset;\breve{\Pi}_B := \emptyset;for \mu \in \Lambda_\ell dofor T \in \psi_\mu.support doT.parent.marked := true;for \lambda \in \breve{\Pi} doif \exists T \in \breve{\phi}_\lambda.support with T.marked = true then\breve{\Pi}_B.insert(\lambda);else
```

 $\Pi_A.\operatorname{insert}(\lambda);$ for  $\mu \in \Lambda_\ell$  do
for  $T \in \psi_\mu.\operatorname{support}$  do  $T.\operatorname{parent.marked} := \operatorname{false};$ 

We can apply the linear operators appearing in the routines of §7.3.1 efficiently by again traversing the domain mother tree; for example, Algorithm 7.6 details a matrix-free application of  $(A\Phi|_{\Pi})(\check{\Phi}|_{\check{\Pi}})$ .

Algorithm 7.6: The computation of  $\mathbf{e} = (A\Phi|_{\Pi})(\check{\Phi}|_{\check{\Pi}})\mathbf{d}$ .

**Data:** Index sets  $\Pi \subset \Delta_{\ell}$ ,  $\breve{\Pi} \subset \breve{\Delta}_{\ell}$ ,  $\mathbf{d} \in \ell_2(\Pi)$ , local and linear  $A: \operatorname{span} \Phi \to \operatorname{span} \breve{\Phi}'$ . **Result:**  $\mathbf{e} = (A\Phi|_{\Pi})(\breve{\Phi}|_{\breve{\Pi}})\mathbf{d}$ for  $\lambda \in \Pi$  do  $\phi_{\lambda}$ .data :=  $\mathbf{d}_{\lambda}$ ; for  $\mu \in \breve{\Pi}$  do  $\mathbf{e}_{\lambda} := 0$ ; for  $T \in \breve{\phi}_{\mu}$ .support do for  $\phi_{\lambda} \in T$ .functions $(\Delta_{\ell})$  do  $// \{\phi_{\lambda} : \lambda \in \Delta_{\ell}, |\operatorname{supp} \phi_{\lambda} \cap T| > 0\}$   $\mathbf{e}_{\lambda} := \mathbf{e}_{\lambda} + A(\phi_{\lambda})(\breve{\phi}_{\mu}|_{T}) \cdot \phi_{\lambda}$ .data; for  $\lambda \in \Pi$  do  $\phi_{\lambda}$ .data := 0;

#### **Operations in space**

We can construct a triangulation  $\mathcal{T}$  from a vertex tree  $N_{\mathcal{T}}$  in linear complexity. First mark every  $\nu \in N_{\mathcal{T}}$  in its mother tree, then traverse the domain mother tree  $\mathfrak{T}$ . A triangle T visited in this traversal is in  $\mathcal{T}$  exactly when the newest vertex of its children is not marked.

For the preconditioners  $\mathbf{K}_{\mu}^{\delta}$  and  $\mathbf{K}_{\lambda}^{\delta}$  from §7.4.4 we use multigrid. We apply multiplicative V-cycle multigrid, in each cycle applying one pre- and one post Gauss-Seidel smoother with reversed ordering of the unknowns.

In view of obtaining a linear complexity algorithm, at level k we restrict smoothing to the vertices of generation k as well as their godparents, cf. [WZ17]. For  $\mathcal{T} \in \mathbb{T}$  we consider  $W_{\mathcal{T}}$ , the space of continuous piecewise linears w.r.t.  $\mathcal{T}$ , zero on  $\partial\Omega$ , now equipped with the single-scale basis  $\Phi_{\mathcal{T}}$ . Set  $L = L(\mathcal{T}) := \max_{T \in \mathcal{T}} gen(T)$ , and define the sequence

$$\mathcal{T}_{\perp} = \mathcal{T}_0 \prec \mathcal{T}_1 \prec \cdots \prec \mathcal{T}_L = \mathcal{T} \subset \mathbb{T}$$

where  $\mathcal{T}_{k-1}$  is constructed from  $\mathcal{T}_k$  by removing all vertices  $\nu \in N_{\mathcal{T}_k}$  for which  $gen(\nu) = k$ . For  $1 \leq k \leq L$ , let  $M_k$  be the set of new vertices and their godparents, i.e.,  $M_k := \bigcup_{\nu \in N_{\mathcal{T}_k} \setminus N_{\mathcal{T}_{k-1}}} \{\nu\} \cup gp(\nu)$ , and let  $M_{k,0} := M_k \setminus \partial\Omega$  be the vertices not on the boundary. We consider the multilevel decomposition, cf. [WZ17],

(7.31) 
$$W_{\mathcal{T}_L} = W_{\mathcal{T}_0} + \sum_{k=1}^{L} \sum_{\nu \in M_{k,0}} \operatorname{span} \phi_{k,\nu}, \quad \text{where} \quad \phi_{k,\nu} := \phi_{\mathcal{T}_k,\nu}.$$

For  $1 \leq k \leq L$ , let  $\mathbf{P}_k$  be the prolongation matrix, i.e., the matrix representation of the embedding  $W_{\mathcal{T}_{k-1}} \to W_{\mathcal{T}_k}$ , and enumerate the vertices  $M_{k,0}$  as  $(\nu_k^i)_{i=1}^{n_k}$ . Algorithm 7.7 details a (non-recursive) implementation of a single multiplicative V-cycle for the multilevel decomposition (7.31) using Gauss-Seidel smoothing. We assume the availability of an efficient coarse-grid solver; in our application, a direct solve suffices. For linear complexity, we use in-place vector updates restricted to non-zeros.

Note that this multigrid method is given in terms of the single-scale basis  $\Phi_{\tau}$ ; it can be transformed to the hierarchical basis  $\Psi_{\tau}$  similarly to (7.29). Multiple V-cycles are done by setting  $u_0 := 0$  and iterating  $u_k := MG(A, f - MG(A, f))$  $Au_{k-1}$ ).

	Algorithm 7.7	7: Single mul	tiplicative V-c	vcle multigrid	MG(A, f)
--	---------------	---------------	-----------------	----------------	----------

**Data:** Some  $f \in W'_{\mathcal{T}}$  and a linear operator  $A \colon W_{\mathcal{T}} \to W'_{\mathcal{T}}$ . **Result:**  $u = \mathbf{u}^{\top} \Phi_{\mathcal{T}} \in W_{\mathcal{T}}$ , the result of a single V-cycle applied to f.  $\mathbf{r} := f(\Phi_{\mathcal{T}});$ for  $L \ge k \ge 1$  do for  $\nu = \nu_k^1, \ldots, \nu_k^{n_k}$  do  $r_{k,\nu} := \mathbf{r}_{\nu};$  $e_{k,\nu} := r_{k,\nu} / (A\phi_{k,\nu})(\phi_{k,\nu});$  $\mathbf{r} := \mathbf{r} - e_{k,\nu} (A\phi_{k,\nu}) (\Phi_{\mathcal{T}_k});$  $\mathbf{r} := \mathbf{P}_{k}^{\top}\mathbf{r};$ Solve  $(A\Phi_{\mathcal{T}_0})(\Phi_{\mathcal{T}_0})\mathbf{u} = \mathbf{r};$ for 1 < k < L do  $\mathbf{u} := \mathbf{P}_k \mathbf{u};$ for  $\nu = \nu_k^{n_k}, \ldots, \nu_k^1$  do  $\mathbf{u}_{\nu} := \mathbf{u}_{\nu} + e_{k,\nu};$  $\mathbf{u}_{\nu} := \mathbf{u}_{\nu} + (r_{k,\nu} - (A\phi_{k,\nu})(\mathbf{u}^{\top}\Phi_{\mathcal{T}_{k}}))/(A\phi_{k,\nu})(\phi_{k,\nu});$ 

#### 7.5.2 Double-trees and tensor-product operators

For every node in a double-tree  $\Lambda \subset \vee^0 \times \vee^1$ , we store a reference to the underlying pair of nodes in their mother trees. This allows growing doubletrees intuitively, and allows comparing different double-trees over the same pair of mother trees. C++ templates allow us to re-use much of the tree code without runtime performance loss.

In §7.3.2 we saw how to apply a tensor-product operator. For this, we first construct the double-trees  $\Sigma$  and  $\Theta$ ; construction of  $\Sigma$  is illustrated in Algorithm 7.8. Evaluation of the operator then reduces to the four simple steps of Algorithm 7.9.

Algorithm 7.8: Function GenerateSigma( $\Lambda$ ,  $\Lambda$ ).

**Data:**  $\check{\Lambda} \subset \check{\lor}^0 \times \check{\lor}^1, \Lambda \subset \lor^0 \times \lor^1$  **Result:**  $\Sigma$  for application of Theorem 7.3.13 with  $\check{\Lambda}$  and  $\Lambda$ .  $\Sigma := P_0 \Lambda \times \{\nu \in P_1 \check{\Lambda} : |\nu| = 0\};$ for  $\lambda \in \Sigma$ .project(0) do for  $T \in \phi_{\lambda}$ .support do for  $\mu \in T$ .functions $(\check{\lor}^0_{|\lambda|})$  do  $\Sigma$ .fiber $(1, \lambda)$ .union $(\check{\Lambda}$ .fiber $(1, \mu)$ );

**Algorithm 7.9:** Algorithm to evaluate  $\mathbf{d} = R_{\check{\mathbf{\Lambda}}} (\mathbf{A}_0 \otimes \mathbf{A}_1) I_{\mathbf{\Lambda}} \mathbf{c}$ .

 $\begin{array}{l} \textbf{Data:} \ \boldsymbol{\Lambda} \subset \vee^0 \times \vee^1, \ \boldsymbol{\check{\Lambda}} \subset \breve{\vee}^0 \times \breve{\vee}^1, \ \mathbf{c} \in \ell_2(\boldsymbol{\Lambda}), \ \mathbf{d} \in \ell_2(\breve{\boldsymbol{\Lambda}}).\\ \boldsymbol{\Sigma} := \texttt{GenerateSigma}(\breve{\boldsymbol{\Lambda}}, \boldsymbol{\Lambda});\\ \boldsymbol{\Theta} := \texttt{GenerateTheta}(\breve{\boldsymbol{\Lambda}}, \boldsymbol{\Lambda});\\ \mathbf{s} := \mathbf{0} \in \ell_2(\boldsymbol{\Sigma});\\ \mathbf{t} := \mathbf{0} \in \ell_2(\boldsymbol{\Theta});\\ \mathbf{l} := \mathbf{0} \in \ell_2(\breve{\boldsymbol{\Theta}});\\ \mathbf{l} := \mathbf{0} \in \ell_2(\breve{\boldsymbol{\Lambda}});\\ \texttt{for } \boldsymbol{\lambda} \in \texttt{s.project}(0) \ \mathbf{do} \ \texttt{eval}(A_1)(\texttt{s.fiber}(1, \boldsymbol{\lambda}), \texttt{c.fiber}(1, \boldsymbol{\lambda}));\\ \texttt{for } \boldsymbol{\mu} \in \texttt{l.project}(1) \ \mathbf{do} \ \texttt{evallow}(A_0)(\texttt{l.fiber}(0, \boldsymbol{\mu}), \texttt{s.fiber}(0, \boldsymbol{\mu}));\\ \texttt{for } \boldsymbol{\mu} \in \texttt{t.project}(1) \ \mathbf{do} \ \texttt{evalupp}(A_0)(\texttt{t.fiber}(0, \boldsymbol{\mu}), \texttt{c.fiber}(0, \boldsymbol{\mu}));\\ \texttt{for } \boldsymbol{\lambda} \in \texttt{d.project}(0) \ \mathbf{do} \ \texttt{eval}(A_1)(\texttt{d.fiber}(1, \boldsymbol{\lambda}), \texttt{t.fiber}(1, \boldsymbol{\lambda}));\\ \mathbf{d} := \mathbf{d} + \mathbf{l}; \end{aligned}$ 

#### Memory optimizations

As the memory consumption of a double-tree is significant, at around 280 bytes per node, we want to have as few double-trees in memory as possible. By storing the nodes of  $\Lambda$  in a persistent container, every node is uniquely identified with its index in the container. This induces a mapping  $\mathbb{R}^{\#\Lambda} \leftrightarrow \ell_2(\Lambda)$  and allows us to overlay multiple vectors on the same underlying double-tree in a memory-friendly way.

The  $\Sigma$  generated by Algorithm 7.8 for the application of a tensor-product operator can play the role of  $\Theta$  necessary for the application of its transpose operator (and vice versa). This allows tensor-product operators and their transposes to share the double-trees  $\Sigma$  and  $\Theta$ .

With these insights, our implementation of the heat equation has at most 5 different double-trees in memory.

# 7.5.3 The adaptive loop

In the refine step of the adaptive loop, we first mark a set *J* of nodes in  $\Lambda^{\delta} \setminus \Lambda^{\delta}$  using Dörfler marking (possible in linear complexity; cf. [PP20]). We then refine  $\Lambda^{\delta}$  to the smallest double-tree containing *J*:

- 1. mark all nodes in  $\Lambda^{\delta}$  that are also present in  $\Lambda^{\delta}$  ((ii) in Fig. 7.7);
- 2. traverse  $\Lambda^{\delta}$  from every node in *J*, top-down in level-wise order, until hitting a previously marked node. Mark all nodes along the way ((iii–iv) in Fig. 7.7);
- 3. union the marked nodes of  $\Lambda^{\delta}$  into  $\Lambda^{\delta}$  ((v) in Fig. 7.7).

As  $\#\Lambda^{\delta} \lesssim \#\Lambda^{\delta}$  and we visit every node of  $\Lambda^{\delta}$  at most twice, the traversal is linear in  $\#\Lambda^{\delta}$ . See also Figure 7.7.



FIGURE 7.7. Adaptive refinement of a double-tree with underlying *unary* mother trees. Left to right: (i)  $\Lambda^{\delta}$ ; (ii)  $\Lambda^{\delta}$  with nodes in  $\Lambda^{\delta} \setminus \Lambda^{\delta}$  in white; (iii) nodes in J marked in red; (iv) nodes marked in the top-down traversal; (v) refined  $\Lambda^{\delta}$ .

# 7.6 Numerical experiments

We consider the heat equation (7.26), and assess our implementation of the adaptive Algorithm 7.1 for its numerical solution. Complementing the convergence results gathered in §6.7, here we provide results on the practical performance of the adaptive loop. Results were gathered on a multi-core 2.2 GHz machine, provided by the Dutch national e-infrastructure with the support of SURF Cooperative.

# 7.6.1 The adaptive loop

We summarize the main results from §6.7. We run Algorithm 7.1 with  $\theta = \frac{1}{2}$  and  $\xi = \frac{1}{2}$ . We consider four problems.

In the *smooth problem*, we select  $\Omega := [0, 1]^2$  and prescribe the solution

$$u(t, x, y) := (1 + t^2)x(1 - x)y(1 - y).$$

In the *moving peak* problem, we again select  $\Omega := [0,1]^2$  with prescribed solution

$$u(t, x, y) := x(1 - x)y(1 - y)\exp(-100[(x - t)^{2} + (y - t)^{2}]);$$



FIGURE 7.8. Error convergence and peak memory usage of the adaptive loop for the four problems of §7.6.1.

here, u is essentially zero outside a small strip along the diagonal (0,0,0) to (1,1,1).

In the *cylinder* problem, we select  $\Omega := [-1, 1]^2 \setminus [-1, 0]^2$  with data

$$u_0 \equiv 0$$
, and  $g(t, x, y) := t \cdot \mathbb{1}_{\{x^2 + y^2 < 1/4\}}$ .

The solution has singularities in the re-entrant corner and along the wall of the cylinder  $\{(t, x, y) : x^2 + y^2 = 1/4\}$ .

In the *singular* problem, we select  $\Omega := [-1, 1]^2 \setminus [-1, 0]^2$  with data  $u_0 \equiv 1$  and  $g \equiv 0$ ; the solution then has singularities along  $\{0\} \times \partial \Omega$  and  $I \times \{(0, 0)\}$ .

#### Convergence

To estimate the error  $||u - \hat{u}^{\delta}||_X$ , we measure the residual error estimator  $||\mathbf{r}^{\delta}(\hat{u}^{\delta})||$  from (7.12); see also Lemma 7.2.3. In the left pane of Figure 7.8, for the first three problems, we observe a convergence rate of 1/2, which is the best that can be expected from our family of trial spaces  $(X^{\delta})_{\delta \in \Delta}$ . For the singular problem, the reduced rate 0.4 is found; it is unknown if a better rate can be expected.

#### Memory

The right pane of Figure 7.8 shows the peak memory consumption after every iteration of the adaptive algorithm. We see that the peak memory is linear in dim  $X^{\delta}$ , stabilizing to around 15kB per degree of freedom. This is relatively high due to our implementation based on double-trees. In fact, the double-trees together make up around 85% of the total memory. For the singular problem, the largest double-tree  $\Lambda_Y^{\delta}$  occupies around 40% of the total memory.



FIGURE 7.9. Time (in ms) per DoF of bilinear form evaluations in time.

### 7.6.2 Linearity of operations

The majority of our runtime is spent in the application of bilinear forms. In this section, we measure the application times to assert their linear complexity.

#### In time

We select three sequences  $\{\Lambda_U\}$ ,  $\{\Lambda_L\}$ ,  $\{\Lambda_R\}$  of trees in  $\vee_{\Sigma}$ , one uniformly refined and two graded towards the left and right respectively. For each such tree  $\Lambda \subset \vee_{\Sigma}$ , we define a corresponding tree  $\check{\Lambda} := \{\mu \in \vee_{\Xi} : \exists \lambda \in \Lambda, |\lambda| = |\mu|, |\operatorname{supp} \xi_{\mu} \cap \operatorname{supp} \sigma_{\lambda}| > 0\} \subset \vee_{\Xi}$ .

We select the bilinear forms  $M_t$  and  $D_t$  from (7.27), and run the algorithms from §7.3.1. We see in Figure 7.9 that the runtime per degree of freedom stabilizes to  $10^{-3}$  ms, essentially independent of the bilinear form and the trees. We suspect the increase until  $10^7$  degrees of freedom has to do with cache locality.

#### In space

On the L-shaped domain  $\Omega := [-1,1]^2 \setminus [-1,0]^2$ , we select two sequences of hierarchical basis trees, one uniformly refined and the other refined by a standard adaptive loop on  $-\Delta u = 1$ ,  $u|_{\partial\Omega} = 0$ .

For a hierarchical basis tree  $\Psi_{\mathcal{T}} = \{\psi_{\nu} : \nu \in N_{\mathcal{T},0}\}$ , we denote the stiffness matrix  $\langle \nabla \Psi_{\mathcal{T}}, \nabla \Psi_{\mathcal{T}} \rangle_{L_2(\Omega)}$  as  $\mathbf{A}_{\mathcal{T}}$ . We measure the runtime of the conversion from vertex tree  $N_{\mathcal{T}}$  to triangulation  $\mathcal{T}$  (cf. §7.5.1), the application time of  $\mathbf{A}_{\mathcal{T}}$  through (7.29), and that of multigrid on  $\mathbf{A}_{\mathcal{T}}$  (with 1 and 3 V-cycles) through Algorithm 7.7. Figure 7.10 confirms that the relative runtime of every operation is essentially independent of the refinement strategy. Interesting is again the increase until  $10^5$  degrees of freedom.



FIGURE 7.10. Time (in ms) per DoF of important operations in space, for uniform and adaptive refinements.



FIGURE 7.11. Time (in ms) per DoF of the four bilinear forms applied in the solve step of the adaptive algorithm.

#### In space-time

Solving (7.9) using PCG requires the application of the four linear operators  $E_Y^{\delta}BE_X^{\delta}$ ,  $E_X^{\delta'}\gamma'_0\gamma_0 E_X^{\delta}$ ,  $K_X^{\delta}$ , and  $K_Y^{\delta}$ . For the first two, Corollary 7.3.14 asserts that their application time is of linear complexity, while for the preconditioners  $K_X^{\delta}$  and  $K_Y^{\delta}$ , this follows from the block-diagonal structure of their matrix representation.

We run the adaptive algorithm on the four problems of §7.6.1. Figure 7.11 shows that the application time of the aforementioned operators is essentially independent of the problem, even though the underlying double-trees are vastly different. We again see an increase in relative runtime until  $10^6$  degrees of freedom.

Figure 7.12 shows the runtimes of the solve, estimate, mark and refine steps of the adaptive loop. We confirm that each step is of linear complexity, and that the total runtime is governed by the solve and estimate steps.



FIGURE 7.12. Time (in ms) per DoF of the steps in the adaptive loop.



FIGURE 7.13. Speedup and time (in ms) per DoF of the solve step in the adaptive loop, for different number of parallel processors.

### 7.6.3 Shared-memory parallelism

Most of our execution time is spent applying the linear operators from Figure 7.11. We can obtain a significant speedup with multithreading. In Algorithm 7.9, all fibers inside each of the four for-loops are disjoint, and we can easily parallelize each loop using OpenMP.

We run the parallel code on the smooth and singular problems. The right pane of Figure 7.13 shows decent parallel performance for the singular problem, with  $10 \times$  speedup at 16 cores. The left pane however reveals a load balancing issue: as *u* is smooth, the two fibers  $(\Lambda_0^{\delta})_{1,\lambda}$  with  $|\lambda| = 0$  contain the majority of the degrees of freedom. This results in poor parallel efficiency for the first and fourth loop in Algorithm 7.9.

# 7.7 Conclusion

We discussed an implementation of an adaptive solver for a space-time variational formulation of parabolic evolution equations where every step is of linear complexity.

We constructed a family of trial spaces spanned by tensor-products of wavelets in time and hierarchical basis functions in space. The resulting adaptive loop is able to resolve singularities locally in space and time, and we proved its *r*-linear convergence.

After imposing a *double-tree* constraint on the index set of the trial spaces, we devised an abstract algorithm that is able to apply the system matrices in linear complexity. We achieve this complexity in practice by a *tree-based* implementation. The numerical results show high performance of the adaptive loop as a whole.

# 7.A Proofs of Theorems in §7.3

**Theorem 7.3.7.** A call of eval yields the output as specified, at the cost of  $\mathcal{O}(\# \Pi + \# \Lambda + \# \Pi + \# \Lambda)$  operations.

*Proof.* By locality of the collections  $\check{\Phi}$  and  $\check{\Psi}$ , and sparsity of the matrices  $\check{\mathfrak{p}}_{\ell}$  and  $\check{\mathfrak{q}}_{\ell}$ , we see that  $\#\underline{\Pi} \lesssim \#\Pi_B + \#\check{\Lambda}_{\ell} \lesssim \#\Lambda_{\ell} + \#\check{\Lambda}_{\ell}$ . So after sufficiently many recursive calls, the current set  $\Pi \cup \check{\Lambda}$  will be empty. For use later, we note that similarly  $\#\underline{\Pi} \lesssim \#\Pi_B + \#\Lambda_{\ell} \lesssim \#\check{\Lambda}_{\ell} + \#\check{\Pi}_B + \#\Lambda_{\ell} \lesssim \#\Lambda_{\ell} + \#\check{\Lambda}_{\ell}$ .

For  $\Pi \cup \Lambda = \emptyset$ , the call produces nothing, which is correct.

Now let  $\Pi \cup \Lambda \neq \emptyset$ . From  $\Lambda$  being an  $\ell$ -tree, the definitions of  $S(\cdot)$  and  $\Pi_A$ , and the locality of A, one has

$$\mathbf{e}|_{\breve{\Pi}_A} = (Au)(\breve{\Phi}|_{\breve{\Pi}_A}) = (A(\mathbf{d}^{\top}\Phi|_{\Pi}))(\breve{\Phi}|_{\breve{\Pi}_A}).$$

By choice of  $\underline{\Pi}$  we have

$$\underline{u} := \underline{\mathbf{d}}^{\top} \Phi |_{\underline{\Pi}} + \mathbf{c} |_{\Lambda_{\ell+1\uparrow}}^{\top} \Psi |_{\Lambda_{\ell+1\uparrow}} = (\mathbf{d} |_{\Pi_B})^{\top} \Phi |_{\Pi_B} + \mathbf{c}^{\top} \Psi |_{\Lambda} = u - (\mathbf{d} |_{\Pi_A})^{\top} \Phi |_{\Pi_A}.$$

By induction the recursive call yields  $\underline{\mathbf{e}} = (A\underline{u})(\check{\Phi}|_{\underline{\check{\Pi}}})$ , and  $\underline{\mathbf{f}} = (A\underline{u})(\check{\Psi}|_{\check{\Lambda}_{\ell+1\uparrow}})$ . From  $\check{\Lambda}$  being an  $\ell$ -tree, the definitions of  $\check{S}(\cdot)$  and  $\Pi_A$ , and the locality of A, we have

$$(Au)(\check{\Psi}|_{\check{\Lambda}_{\ell\uparrow}}) = (A\underline{u})(\check{\Psi}|_{\check{\Lambda}_{\ell\uparrow}}),$$

and so in particular  $\mathbf{f}|_{\Lambda_{\ell+1,\uparrow}} = \mathbf{\underline{f}}$ .

The definition of  $\underline{\Pi}$  shows that

$$\check{\Phi}|_{\check{\Pi}_B} = (\check{\mathfrak{p}}_{\ell}^{\top}\check{\Phi}|_{\underline{\check{\Pi}}})|_{\check{\Pi}_B}, \quad \check{\Psi}|_{\check{\Lambda}_{\ell}} = (\check{\mathfrak{q}}_{\ell}^{\top}\check{\Phi}|_{\underline{\check{\Pi}}})|_{\check{\Lambda}_{\ell}}.$$

We conclude that

$$\mathbf{f}|_{\check{\Lambda}_{\ell}} = (Au)(\check{\Psi}|_{\check{\Lambda}_{\ell}}) = (A\underline{u})(\check{\Psi}|_{\check{\Lambda}_{\ell}}) = \left(\check{\mathfrak{q}}_{\ell}^{\top}\underline{\mathbf{e}}\right)|_{\check{\Lambda}_{\ell}},$$

and from  $|\operatorname{supp} \phi_{\lambda} \cap \operatorname{supp} \check{\phi}_{\mu}| = 0$  for  $(\lambda, \mu) \in \Pi_A \times \check{\Pi}_B$ , that

$$\mathbf{e}|_{\breve{\Pi}_B} = (Au)(\breve{\Phi}|_{\breve{\Pi}_B}) = (A\underline{u})(\breve{\Phi}|_{\breve{\Pi}_B}) = (\breve{\mathfrak{p}}_{\ell}^{\top}\underline{\mathbf{e}})|_{\breve{\Pi}_B}.$$

From the assumptions on the collections  $\Phi$ ,  $\check{\Phi}$ ,  $\check{\Psi}$ , and  $\Psi$ , and their consequences on the sparsity of the matrices  $\mathfrak{p}_{\ell}$ ,  $\check{\mathfrak{p}}_{\ell}$ ,  $\mathfrak{q}_{\ell}$ , and  $\check{\mathfrak{q}}_{\ell}$ , one infers that the total cost of the evaluations of the statements in eval is  $\mathcal{O}(\#\check{\Pi} + \#\check{\Lambda}_{\ell} + \#\Pi + \#\Lambda_{\ell})$ plus the cost of the recursive call. Using  $\#\check{\Pi} + \#\Pi \leq \#\check{\Lambda}_{\ell} + \#\Lambda_{\ell}$  and induction, we conclude the second statement of the theorem.

**Theorem 7.3.9.** A call of evalupp yields the output as specified, at the cost of  $\mathcal{O}(\#\Pi + \#\Lambda + \#\Pi + \#\Lambda)$  operations.

*Proof.* By locality of the collections  $\check{\Phi}$  and  $\check{\Psi}$ , and sparsity of the matrices  $\check{\mathfrak{p}}_{\ell}$  and  $\check{\mathfrak{q}}_{\ell}$ , we see that  $\#\underline{\check{\Pi}} \lesssim \#\check{\Pi}_B + \#\check{\Lambda}_{\ell} \lesssim \#\Lambda_{\ell} + \#\check{\Lambda}_{\ell}$ . So after sufficiently many recursive calls, the current set  $\breve{\Pi} \cup \check{\Lambda}$  will be empty. Notice that  $\#\underline{\Pi} \lesssim \#\Lambda_{\ell}$ .

For  $\Pi \cup \Lambda = \emptyset$ , the call produces nothing, which is correct.

Now let  $\Pi \cup \Lambda \neq \emptyset$ . From  $\Lambda$  being an  $\ell$ -tree, the definitions of  $S(\cdot)$  and  $\Pi_A$ , and the locality of A, one has

$$\mathbf{e}|_{\breve{\Pi}_A} = (Au)(\breve{\Phi}|_{\breve{\Pi}_A}) = (A(\mathbf{d}^{\top}\Phi|_{\Pi}))(\breve{\Phi}|_{\breve{\Pi}_A}).$$

By definition of  $\underline{\Pi}$  we have

$$\underline{u} := \underline{\mathbf{d}}^{\top} \Phi |_{\underline{\Pi}} + \mathbf{c} |_{\Lambda_{\ell+1\uparrow}}^{\top} \Psi |_{\Lambda_{\ell+1\uparrow}} = \mathbf{c}^{\top} \Psi |_{\Lambda} = u - \mathbf{d}^{\top} \Phi |_{\Pi}.$$

By induction the recursive call yields

$$\underline{\mathbf{e}} = (A\underline{u})(\check{\Phi}|_{\underline{\check{\Pi}}}), \quad \underline{\mathbf{f}} = \mathbf{U}_{\check{\Lambda}_{\ell+1\uparrow} \times \Lambda_{\ell+1\uparrow}} c|_{\Lambda_{\ell+1\uparrow}} = \mathbf{f}|_{\check{\Lambda}_{\ell+1\uparrow}}.$$

The definition of  $\underline{\Pi}$  shows that

$$\check{\Phi}|_{\check{\Pi}_B} = (\check{\mathfrak{p}}_{\ell}^{\top}\check{\Phi}|_{\underline{\check{\Pi}}})|_{\check{\Pi}_B}, \quad \check{\Psi}|_{\check{\Lambda}_{\ell}} = (\check{\mathfrak{q}}_{\ell}^{\top}\check{\Phi}|_{\underline{\check{\Pi}}})|_{\check{\Lambda}_{\ell}}.$$

We conclude that

$$\mathbf{f}|_{\check{\Lambda}_{\ell}} = (A(\mathbf{c}^{\top}\Psi|_{\Lambda}))(\check{\Psi}|_{\check{\Lambda}_{\ell}}) = (A\underline{u})(\check{\Psi}|_{\check{\Lambda}_{\ell}}) = \left(\check{\mathfrak{q}}_{\ell}^{\top}\underline{\mathbf{e}}\right)|_{\check{\Lambda}_{\ell}},$$

and

$$\mathbf{e}|_{\breve{\Pi}_B} = (Au)(\breve{\Phi}|_{\breve{\Pi}_B}) = (A\underline{u})(\breve{\Phi}|_{\breve{\Pi}_B}) + (A(\mathbf{d}^{\top}\Phi|_{\Pi}))(\breve{\Phi}|_{\breve{\Pi}_B}) \\ = (\mathbf{p}_{\ell}^{\top}\underline{\mathbf{e}})|_{\breve{\Pi}_B} + (A(\mathbf{d}^{\top}\Phi|_{\Pi}))(\breve{\Phi}|_{\breve{\Pi}_B}).$$

From the assumptions on the collections  $\Phi$ ,  $\check{\Phi}$ ,  $\check{\Psi}$ , and  $\Psi$ , and their consequences on the sparsity of the matrices  $\mathfrak{p}_{\ell}$ ,  $\check{\mathfrak{p}}_{\ell}$ ,  $\mathfrak{q}_{\ell}$ , and  $\check{\mathfrak{q}}_{\ell}$ , one infers that the total cost of the evaluations of the statements in eval is  $\mathcal{O}(\#\check{\Pi} + \#\check{\Lambda}_{\ell} + \#\Pi + \#\Lambda_{\ell})$ plus the cost of the recursive call. Using  $\#\check{\Pi} + \#\underline{\Pi} \lesssim \#\check{\Lambda}_{\ell} + \#\Lambda_{\ell}$  and induction, we conclude the second statement of the theorem.

**Theorem 7.3.11.** A call of evallow yields the output as specified, at the cost of  $\mathcal{O}(\#\breve{\Lambda} + \#\Pi + \#\Lambda)$  operations.

*Proof.* Notice that  $\#\underline{\Pi} \lesssim \#\Lambda_{\ell} + \#\Pi_B \lesssim \#\Lambda_{\ell} + \#\breve{\Lambda}_{\ell}$ .

For  $\Pi \cup \Lambda = \emptyset$ , the call produces nothing, which is correct. Now let  $\Pi \cup \Lambda \neq \emptyset$ . The definitions of  $\underline{\Pi}$  and  $\underline{\Pi}_B$  show that

$$\begin{aligned} \mathbf{f}|_{\check{\Lambda}_{\ell}} &= (A\Phi|_{\Pi})(\check{\Psi}|_{\check{\Lambda}_{\ell}})\mathbf{d} = (A\Phi|_{\Pi})(\check{\Psi}|_{\check{\Lambda}_{\ell}})\mathbf{d}|_{\Pi_{B}} \\ &= (\check{\mathfrak{q}}_{\ell}^{\top}(A\Phi|_{\underline{\Pi}_{B}})(\check{\Phi}|_{\underline{\check{\Pi}}})\mathfrak{p}_{\ell}\mathbf{d}|_{\Pi_{B}})|_{\check{\Lambda}_{\ell}} = (\check{\mathfrak{q}}_{\ell}^{\top}\underline{\mathbf{e}})|_{\check{\Lambda}_{\ell}} \end{aligned}$$

From  $\Lambda$  being an  $\ell$ -tree, the definitions of  $\check{S}(\cdot)$  and  $\Pi_B$ , and the locality of a, and for the third equality, the definition of  $\underline{\Pi}$ , one has

$$\begin{split} f|_{\breve{\Lambda}_{\ell+1\uparrow}} &= a(\breve{\Psi}|_{\breve{\Lambda}_{\ell+1\uparrow}}, \Phi|_{\Pi})\mathbf{d} + \mathbf{L}|_{\breve{\Lambda}_{\ell+1\uparrow} \times \Lambda_{\ell}} \mathbf{c}|_{\Lambda_{\ell}} + \mathbf{L}|_{\breve{\Lambda}_{\ell+1\uparrow} \times \Lambda_{\ell+1\uparrow}} \mathbf{c}|_{\Lambda_{\ell+1\uparrow}} \\ &= (A\Phi|_{\Pi})(\breve{\Psi}|_{\breve{\Lambda}_{\ell+1\uparrow}})\mathbf{d}|_{\Pi_{B}} + (A\Psi|_{\Lambda_{\ell}})(\breve{\Psi}|_{\breve{\Lambda}_{\ell+1\uparrow}})\mathbf{c}|_{\Lambda_{\ell}} + \mathbf{L}|_{\breve{\Lambda}_{\ell+1\uparrow} \times \Lambda_{\ell+1\uparrow}} \mathbf{c}|_{\Lambda_{\ell+1\uparrow}} \\ &= (A\Phi|_{\underline{\Pi}})(\breve{\Psi}|_{\breve{\Lambda}_{\ell+1\uparrow}})\underline{\mathbf{d}} + \mathbf{L}|_{\breve{\Lambda}_{\ell+1\uparrow} \times \Lambda_{\ell+1\uparrow}} \mathbf{c}|_{\Lambda_{\ell+1\uparrow}} \\ &= \mathrm{evallow}(A)(\ell+1,\breve{\Lambda}_{\ell+1\uparrow},\underline{\Pi},\Lambda_{\ell+1\uparrow},\underline{\mathbf{d}},\mathbf{c}|_{\Lambda_{\ell+1\uparrow}}) \end{split}$$

by induction.

From the assumptions on the collections  $\Phi$ ,  $\Psi$ , and  $\Psi$ , and their consequences on the sparsity of the matrices  $\mathfrak{p}_{\ell}$ ,  $\mathfrak{q}_{\ell}$ , and  $\check{\mathfrak{q}}_{\ell}$ , one easily infers that the total cost of the evaluations of the statements in evallow is  $\mathcal{O}(\#\check{\Lambda}_{\ell} + \#\Pi + \#\Lambda_{\ell})$ plus the cost of the recursive call. Using  $\#\underline{\Pi} \lesssim \#\check{\Lambda}_{\ell} + \#\Lambda_{\ell}$  and induction, we conclude the second statement of the theorem.

**Theorem 7.3.13.** Let  $\check{\Lambda} \subset \check{\vee}^0 \times \check{\vee}^1$ ,  $\Lambda \subset \vee^0 \times \vee^1$  be finite double-trees. Then

$$\begin{split} \boldsymbol{\Sigma} &\coloneqq \bigcup_{\boldsymbol{\lambda} \in P_0 \boldsymbol{\Lambda}} \Big( \{\boldsymbol{\lambda}\} \times \bigcup_{\substack{\{\mu \in P_0 \check{\boldsymbol{\Lambda}} : |\mu| = |\boldsymbol{\lambda}| + 1, \ |\check{S}^0(\mu) \cap S^0(\boldsymbol{\lambda})| > 0\}}} \check{\boldsymbol{\Lambda}}_{1,\mu} \Big), \\ \boldsymbol{\Theta} &\coloneqq \bigcup_{\boldsymbol{\lambda} \in P_1 \boldsymbol{\Lambda}} \Big( \{\mu \in P_0 \check{\boldsymbol{\Lambda}} : \exists \gamma \in \boldsymbol{\Lambda}_{0,\boldsymbol{\lambda}} \ s.t. \ |\gamma| = |\mu|, \ |\check{S}^0(\mu) \cap S^0(\gamma)| > 0\} \times \{\boldsymbol{\lambda}\} \Big), \end{split}$$

are double-trees with  $\#\Sigma \lesssim \#\breve{\Lambda}$  and  $\#\Theta \lesssim \#\Lambda$ , and

$$R_{\check{\mathbf{\Lambda}}}(\mathbf{A}_0 \otimes \mathbf{A}_1) I_{\mathbf{\Lambda}} = R_{\check{\mathbf{\Lambda}}}(\mathbf{L}_0 \otimes \mathrm{Id}) I_{\mathbf{\Sigma}} R_{\mathbf{\Sigma}}(\mathrm{Id} \otimes \mathbf{A}_1) I_{\mathbf{\Lambda}} + R_{\check{\mathbf{\Lambda}}}(\mathrm{Id} \otimes \mathbf{A}_1) I_{\mathbf{\Theta}} R_{\mathbf{\Theta}}(\mathbf{U}_0 \otimes \mathrm{Id}) I_{\mathbf{\Lambda}}.$$

Proof. We write

(7.32)  

$$R_{\check{\mathbf{\Lambda}}}(\mathbf{A}_{0} \otimes \mathbf{A}_{1})I_{\mathbf{\Lambda}} = R_{\check{\mathbf{\Lambda}}}((\mathbf{L}_{0} + \mathbf{U}_{0}) \otimes \mathbf{A}_{1})I_{\mathbf{\Lambda}}$$

$$= R_{\check{\mathbf{\Lambda}}}(\mathbf{L}_{0} \otimes \mathrm{Id})(\mathrm{Id} \otimes \mathbf{A}_{1})I_{\mathbf{\Lambda}} + R_{\check{\mathbf{\Lambda}}}(\mathrm{Id} \otimes \mathbf{A}_{1})(\mathbf{U}_{0} \otimes \mathrm{Id})I_{\mathbf{\Lambda}}.$$

Considering (7.32), the range of  $(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{\Lambda}}$  consists of vectors whose entries with first index outside  $P_0\mathbf{\Lambda}$  are zero. In view of the subsequent application of  $\mathbf{L}_0 \otimes \mathrm{Id}$ , furthermore only those indices  $(\lambda, \gamma) \in P_0\mathbf{\Lambda} \times \check{\nabla}^1$ of these vectors might be relevant for which  $\exists (\mu, \gamma) \in \check{\mathbf{\Lambda}}$ , i.e.  $\gamma \in \mathbf{\Lambda}_{1,\mu}$ , with  $|\mu| > |\lambda|$  and  $|\check{S}^0(\mu) \cap S^0(\lambda)| > 0$ . Indeed  $|\check{S}^0(\mu) \cap S^0(\lambda)| = 0$  implies  $|\sup p \check{\psi}^0_\mu \cap \sup p \psi^0_\lambda| = 0$ , and so  $A_0(\check{\psi}^0_\mu, \psi^0_\lambda) = 0$ . If for given  $(\lambda, \gamma)$  such a pair  $(\mu, \gamma)$  exists for  $|\mu| > |\lambda|$ , then such a pair exists for  $|\mu| = |\lambda| + 1$  as well, because  $\check{\mathbf{\Lambda}}_{0,\gamma}$  is a tree, and  $\check{S}^0(\mu') \supset \check{S}^0(\mu)$  for any ancestor  $\mu'$  of  $\mu$ . In order words, the condition  $|\mu| > |\lambda|$  can be read as  $|\mu| = |\lambda| + 1$ . The set of  $(\lambda, \gamma)$  that we just described is given by the set  $\Sigma$ , and so we infer that

$$R_{\check{\mathbf{A}}}(\mathbf{L}_0 \otimes \mathrm{Id})(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{A}} = R_{\check{\mathbf{A}}}(\mathbf{L}_0 \otimes \mathrm{Id})I_{\mathbf{\Sigma}}R_{\mathbf{\Sigma}}(\mathrm{Id} \otimes \mathbf{A}_1)I_{\mathbf{A}}$$

Now let  $(\lambda, \gamma) \in \Sigma$ . Using that  $P_0 \Lambda$  is a tree, and  $S^0(\lambda) \subset S^0(\lambda')$  for any ancestor  $\lambda'$  of  $\lambda$ , we infer that  $(\lambda', \gamma) \in \Sigma$ . Using that for any  $\mu \in P_0 \check{\Lambda}$ ,  $\check{\Lambda}_{1,\mu}$  is a tree, we infer that for any ancestor  $\gamma'$  of  $\gamma$ ,  $(\lambda, \gamma') \in \Sigma$ , so that  $\Sigma$  is a double-tree.

For any  $\mu \in \check{\vee}^0$ , the number of  $\lambda \in \vee^0$  with  $|\mu| = |\lambda| + 1$  and  $|\check{S}^0(\mu) \cap S^0(\lambda)| > 0$  is uniformly bounded, from which we infer that  $\#\Sigma \lesssim \sum_{\mu \in P_0\check{\Lambda}} \#\check{\Lambda}_{1,\mu} = \#\check{\Lambda}$ .

Considering (7.33), the range of  $(\mathbf{U}_0 \otimes \mathrm{Id})I_{\mathbf{\Lambda}}$  consists of vectors that can only have non-zero entries for indices  $(\mu, \lambda) \in \check{\vee}^0 \times P_1 \mathbf{\Lambda}$  for which there exists a  $\gamma \in \mathbf{\Lambda}_{0,\lambda}$  with  $|\gamma| \geq |\mu|$  and  $|\check{S}^0(\mu) \cap S^0(\gamma)| > 0$ . Since  $\mathbf{\Lambda}_{0,\lambda}$  is a tree, and  $S^0(\gamma') \supset S^0(\gamma)$  for any ancestor  $\gamma'$  of  $\gamma$ , equivalently  $|\gamma| \geq |\mu|$  can be read as  $|\gamma| = |\mu|$ . Furthermore, in view of the subsequent application of  $R_{\mathbf{\Lambda}}(\mathrm{Id} \otimes \mathbf{A}_1)$ , it suffices to consider those indices  $(\mu, \lambda)$  with  $\mu \in P_0 \mathbf{\Lambda}$ . The set of  $(\mu, \lambda)$  that we just described is given by the set  $\mathbf{\Theta}$ , and so we infer that

$$R_{\check{\mathbf{\Lambda}}}(\mathrm{Id}\otimes\mathbf{A}_1)(\mathbf{U}_0\otimes\mathrm{Id})I_{\mathbf{\Lambda}}=R_{\check{\mathbf{\Lambda}}}(\mathrm{Id}\otimes\mathbf{A}_1)I_{\mathbf{\Theta}}R_{\mathbf{\Theta}}(\mathbf{U}_0\otimes\mathrm{Id})I_{\mathbf{\Lambda}}.$$

Now let  $(\mu, \lambda) \in \Theta$ . If  $\lambda'$  is an ancestor of  $\lambda$ , then from  $P_0 \Lambda$  being a tree, and  $\Lambda_{0,\lambda} \subset \Lambda_{0,\lambda'}$ , we have  $(\mu, \lambda') \in \Theta$ . If  $\mu'$  is an ancestor of  $\mu$ , then from  $P_0 \check{\Lambda}$  being a tree, and  $\check{S}^0(\mu') \supset \check{S}^0(\mu)$ , we infer that  $(\mu', \lambda) \in \Theta$ , and thus that  $\Theta$  is a double-tree.

For any  $\gamma \in \vee^0$ , the number of  $\mu \in \check{\vee}^0$  with  $|\mu| = |\gamma|$  and  $|\check{S}^0(\mu) \cap S^0(\gamma)| > 0$ is uniformly bounded, from which we infer that  $\#\Theta \lesssim \sum_{\lambda \in P_1 \Lambda} \#\Lambda_{0,\lambda} =$  $\#\Lambda$ .

# 8.1 Introduction

This chapter deals with solving parabolic evolution equations in a time-parallel fashion using tensor-product discretizations. Time-parallel algorithms for solving parabolic evolution equations have become a focal point following the enormous increase in parallel computing power. Spatial parallelism is a ubiquitous component in large-scale computations, but when spatial parallelism is exhausted, parallelization of the time axis is of interest.

Time-stepping methods first discretize the problem in space, and then solve the arising system of coupled ODEs sequentially, immediately revealing a primary source of difficulty for time-parallel computation.

Alternatively, one can solve simultaneously in space *and* time. Originally introduced in [BJ89, BJ90], these space-time methods are very flexible: some can guarantee quasi-best approximations, meaning that their error is proportional to that of the best approximation from the trial space [And13, DS18, FK21, SZ20], or drive adaptive routines [SY18, RS19]. Many are especially well-suited for time-parallel computation [GN16, NS19]. Since the first significant contribution to time-parallel algorithms [Nie64] in 1964, many methods suitable for parallel computation have surfaced; see the review [Gan15].

#### Parallel complexity

The (serial) complexity of an algorithm measures asymptotic runtime on a single processor in terms of the input size. *Parallel complexity* measures asymptotic runtime given *sufficiently many* parallel processors having access to a shared memory, i.e., assuming there are no communication costs.

In the current context of tensor-product discretizations of parabolic PDEs, we denote with  $N_t$  and  $N_x$  the number of unknowns in time and space respectively.

The parareal method [LMT01] aims at time-parallelism by alternating a serial coarse-grid solve with fine-grid computations in parallel. This way, each

iteration has a time-parallel complexity of  $\mathcal{O}(\sqrt{N_t}N_x)$ , and combined with parallel multigrid in space, a parallel complexity of  $\mathcal{O}(\sqrt{N_t}\log N_x)$ . The popular MGRIT algorithm extends these ideas to multiple levels in time; cf. [FFK<sup>+</sup>14].

Recently, Neumüller and Smears proposed an iterative algorithm that uses a Fast Fourier Transform in time. Each iteration runs serially in  $\mathcal{O}(N_t \log(N_t)N_x)$ and parallel in time, in  $\mathcal{O}(\log(N_t)N_x)$ . By also incorporating parallel multigrid in space, its parallel runtime may be reduced to  $\mathcal{O}(\log N_t + \log N_x)$ .

#### **Our contribution**

We study a variational formulation introduced in [SW21b] which was based on work by Andreev [And13, And16]. Recently in [SvVW21, vVW21b], we studied this formulation in the context of space-time adaptivity and its efficient implementation in serial and on shared-memory parallel computers. The current chapter instead focuses on its massively parallel implementation and time-parallel performance.

Our method has remarkable similarities with the approach of [NS19], and the most essential difference is the substitution of their Fast Fourier Transform by our Fast Wavelet Transform. The strengths of both methods include a solid inf-sup theory that enables quasi-optimal approximate solutions from the trial space, ease of implementation, and excellent parallel performance in practice.

Our method has another strength: based on a wavelet transform, for fixed algebraic tolerance it runs serially in linear complexity. Parallel in time, it runs in complexity  $O(\log(N_t)N_x)$ ; parallel in *space and time*, in  $O(\log(N_tN_x))$ . Moreover, when solving to an algebraic error proportional to the discretization error, incorporating a *nested iteration* (cf. [Hac85, Ch. 5]) results in complexities  $O(N_tN_x)$ ,  $O(\log(N_t)N_x)$ , and  $O(\log^2(N_tN_x))$  respectively. This is on par with best-known results on parallel complexity for elliptic problems; see also [Bra81].

#### Organization of this chapter

In §8.2, we formally introduce the problem, derive a saddle-point formulation, and provide sufficient conditions for quasi-optimality of discrete solutions. In §8.3, we detail on the efficient computation of these discrete solutions. In §8.4 we take a concrete example—the reaction-diffusion equation—and analyze the serial and parallel complexity of our algorithm. In §8.5, we test these theoretical findings in practice. We conclude in §8.6.

### Notations

For normed linear spaces U and V, in this work for convenience over  $\mathbb{R}$ ,  $\mathcal{L}(U, V)$  will denote the space of bounded linear mappings  $U \to V$  endowed with the operator norm  $\|\cdot\|_{\mathcal{L}(U,V)}$ . The subset of invertible operators in  $\mathcal{L}(U, V)$  with inverses in  $\mathcal{L}(V, U)$  will be denoted as  $\mathcal{L}$ is(U, V).

Given a finite-dimensional subspace  $U^{\delta}$  of a normed linear space U, we denote the trivial embedding  $U^{\delta} \rightarrow U$  by  $E_U^{\delta}$ . For a basis  $\Phi^{\delta}$ —viewed formally as a column vector—of  $U^{\delta}$ , we define the *synthesis operator* as

$$\mathcal{F}_{\Phi^{\delta}}: \mathbb{R}^{\dim U^{\delta}} \to U^{\delta}: \boldsymbol{c} \mapsto \boldsymbol{c}^{\top} \Phi^{\delta} =: \sum_{\phi \in \Phi^{\delta}} c_{\phi} \phi.$$

Equip  $\mathbb{R}^{\dim U^{\delta}}$  with the Euclidean inner product and identify  $(\mathbb{R}^{\dim U^{\delta}})'$  with  $\mathbb{R}^{\dim U^{\delta}}$  using the corresponding Riesz map. We find the adjoint of  $\mathcal{F}_{\Phi^{\delta}}$ , the *analysis operator*, to satisfy

$$(\mathcal{F}_{\Phi^{\delta}})' : (U^{\delta})' \to \mathbb{R}^{\dim U^{\delta}} : f \mapsto f(\Phi^{\delta}) := [f(\phi)]_{\phi \in \Phi^{\delta}}.$$

For quantities f and g, by  $f \leq g$ , we mean that  $f \leq C \cdot g$  with a constant that does not depend on parameters that f and g may depend on. By f = g, we mean that  $f \leq g$  and  $g \leq f$ . For matrices A and  $B \in \mathbb{R}^{N \times N}$ , by A = B we will denote *spectral equivalence*, i.e.  $x^{\top}Ax = x^{\top}Bx$  for all  $x \in \mathbb{R}^{N}$ .

# 8.2 Quasi-optimal approximations to the parabolic problem

Let V, H be separable Hilbert spaces of functions on some spatial domain such that V is continuously embedded in H, i.e.  $V \hookrightarrow H$ , with dense compact embedding. Identifying H with its dual yields the Gelfand triple  $V \hookrightarrow H \simeq H' \hookrightarrow V'$ .

For a.e.

$$t \in I := (0, T),$$

let  $a(t; \cdot, \cdot)$  denote a bilinear form on  $V \times V$  so that for any  $\eta, \zeta \in V, t \mapsto a(t; \eta, \zeta)$  is measurable on I, and such that for a.e.  $t \in I$ ,

$$\begin{aligned} |a(t;\eta,\zeta)| &\lesssim \|\eta\|_V \|\zeta\|_V \quad (\eta,\zeta \in V) \quad (boundedness) \\ a(t;\eta,\eta) &\gtrsim \|\eta\|_V^2 \qquad (\eta \in V) \quad (coercivity). \end{aligned}$$

With  $(A(t)\cdot)(\cdot) := a(t; \cdot, \cdot) \in \mathcal{L}is(V, V')$ , given a forcing function g and initial value  $u_0$ , we want to solve the *parabolic initial value problem* of

(8.1) finding 
$$u: I \to V$$
 such that 
$$\begin{cases} \frac{\mathrm{d}u}{\mathrm{d}t}(t) + A(t)u(t) &= g(t) \quad (t \in I), \\ u(0) &= u_0. \end{cases}$$

### 8.2.1 An equivalent self-adjoint saddle-point system

In a simultaneous space-time variational formulation, the parabolic problem reads as finding u from a suitable space of functions of time and space s.t.

$$(Bw)(v) := \int_{I} \langle \frac{\mathrm{d}w}{\mathrm{d}t}(t), v(t) \rangle_{H} + a(t; w(t), v(t)) \mathrm{d}t = \int_{I} \langle g(t), v(t) \rangle_{H} =: g(v)$$

for all v from another suitable space of functions of time and space. One possibility to enforce the initial condition is by testing against additional test functions.

**Theorem 8.2.1 ([SS09]).** With  $X := L_2(I; V) \cap H^1(I; V')$ ,  $Y := L_2(I; V)$ , we have

$$\begin{bmatrix} B\\ \gamma_0 \end{bmatrix} \in \mathcal{L}\mathrm{is}(X, Y' \times H),$$

where for  $t \in \overline{I}$ ,  $\gamma_t : u \mapsto u(t, \cdot)$  denotes the trace map. In other words,

(8.2) finding 
$$u \in X$$
 s.t.  $(Bu, \gamma_0 u) = (g, u_0)$  given  $(g, u_0) \in Y' \times H$ 

is a well-posed simultaneous space-time variational formulation of (8.1).

We define  $A \in \mathcal{L}is(Y, Y')$  and  $\partial_t \in \mathcal{L}is(X, Y')$  as

$$(Au)(v) := \int_I a(t; u(t), v(t)) dt$$
, and  $\partial_t := B - A$ .

Following [SW21b], we assume that *A* is *symmetric*. We can reformulate (8.2) as the self-adjoint saddle point problem

(8.3) finding 
$$(v, \sigma, u) \in Y \times H \times X$$
 s.t. 
$$\begin{bmatrix} A & 0 & B \\ 0 & \text{Id} & \gamma_0 \\ B' & \gamma'_0 & 0 \end{bmatrix} \begin{bmatrix} v \\ \sigma \\ u \end{bmatrix} = \begin{bmatrix} g \\ u_0 \\ 0 \end{bmatrix}.$$

By taking a Schur complement w.r.t. the *H*-block, we can reformulate this as

(8.4) finding 
$$(v, u) \in Y \times X$$
 s.t.  $\begin{bmatrix} A & B \\ B' & -\gamma'_0 \gamma_0 \end{bmatrix} \begin{bmatrix} v \\ u \end{bmatrix} = \begin{bmatrix} g \\ -\gamma'_0 u_0 \end{bmatrix}$ .

We equip *Y* and *X* with 'energy'-norms

$$\|\cdot\|_{Y}^{2} := (A \cdot)(\cdot), \quad \|\cdot\|_{X}^{2} := \|\partial_{t} \cdot\|_{Y'}^{2} + \|\cdot\|_{Y}^{2} + \|\gamma_{T} \cdot\|_{H^{2}}^{2}$$

which are equivalent to the canonical norms on Y and X.

## 8.2.2 Uniformly quasi-optimal Galerkin discretizations

Our numerical approximations will be based on the saddle-point formulation (8.4). Let  $(Y^{\delta}, X^{\delta})_{\delta \in \Delta}$  be a collection of closed subspaces of  $Y \times X$  satisfying

(8.5) 
$$X^{\delta} \subset Y^{\delta}, \quad \partial_t X^{\delta} \subset Y^{\delta} \quad (\delta \in \Delta),$$

and

(8.6) 
$$1 \ge \gamma_{\Delta} := \inf_{\delta \in \Delta} \inf_{0 \ne u \in X^{\delta}} \sup_{0 \ne v \in Y^{\delta}} \frac{(\partial_t u)(v)}{\|\partial_t u\|_{Y'} \|v\|_Y} > 0.$$

*Remark* 8.2.2. In [SW21b, §4], these conditions were verified for  $X^{\delta}$  and  $Y^{\delta}$  being tensor-products of (locally refined) finite element spaces in time and space. In Chapter 6, we relax these conditions to  $X_t^{\delta}$  and  $Y^{\delta}$  being *adaptive sparse grids*, allowing adaptive refinement locally in space *and* time simultaneously.

For  $\delta \in \Delta$ , let  $(v^{\delta}, \overline{u}^{\delta}) \in Y^{\delta} \times X^{\delta}$  solve the Galerkin discretization of (8.4):

(8.7) 
$$\begin{bmatrix} E_Y^{\delta'} A E_Y^{\delta} & E_Y^{\delta'} B E_X^{\delta} \\ E_X^{\delta'} B' E_Y^{\delta} & -E_X^{\delta'} \gamma_0' \gamma_0 E_X^{\delta} \end{bmatrix} \begin{bmatrix} v^{\delta} \\ \overline{u}^{\delta} \end{bmatrix} = \begin{bmatrix} E_Y^{\delta'} g \\ -E_X^{\delta'} \gamma_0' u_0 \end{bmatrix}$$

The solution  $(v^{\delta}, \overline{u}^{\delta})$  of (8.7) exists uniquely, and exhibits *uniform quasi-optimality* in that  $||u - \overline{u}^{\delta}||_X \leq \gamma_{\Delta}^{-1} \inf_{u_{\delta} \in X^{\delta}} ||u - u_{\delta}||_X$  for all  $\delta \in \Delta$ .

Instead of solving a matrix representation of (8.7) using e.g. preconditioned MINRES, we will opt for a computationally more attractive method. By taking the Schur complement w.r.t. the  $Y^{\delta}$ -block in (8.7), and replacing  $(E_Y^{\delta} A E_Y^{\delta})^{-1}$  in the resulting formulation by a *preconditioner*  $K_Y^{\delta}$  that can be applied cheaply, we arrive at the *Schur complement formulation* of finding  $u^{\delta} \in X^{\delta}$  s.t.

(8.8) 
$$\underbrace{E_X^{\delta'}(B'E_Y^{\delta}K_Y^{\delta}E_Y^{\delta'}B+\gamma_0'\gamma_0)E_X^{\delta}}_{=:S^{\delta}}u^{\delta} = \underbrace{E_X^{\delta'}(B'E_Y^{\delta}K_Y^{\delta}E_Y^{\delta'}g+\gamma_0'u_0)}_{=:f^{\delta}}$$

The resulting operator  $S^{\delta} \in \mathcal{L}$ is $(X^{\delta}, X^{\delta'})$  is self-adjoint and elliptic. Given a self-adjoint operator  $K_Y^{\delta} \in \mathcal{L}(Y^{\delta'}, Y^{\delta})$  satisfying, for some  $\kappa_{\Delta} \geq 1$ ,

(8.9) 
$$\frac{\left((K_Y^{\delta})^{-1}v\right)(v)}{(Av)(v)} \in [\kappa_{\Delta}^{-1}, \kappa_{\Delta}] \quad (\delta \in \Delta, \ v \in Y^{\delta}),$$

the solution  $u^{\delta}$  of (8.8) exists uniquely as well. In fact, the following holds.

**Theorem 8.2.3** ([SW21b, Rem. 3.8]). *Take*  $(Y^{\delta} \times X^{\delta})_{\delta \in \Delta}$  *satisfying* (8.5)–(8.6), and  $K_{Y}^{\delta}$  satisfying (8.9). Solutions  $u^{\delta} \in X^{\delta}$  of (8.8) are uniformly quasi-optimal, i.e.

$$\|u - u^{\delta}\|_{X} \le \frac{\kappa_{\Delta}}{\gamma_{\Delta}} \inf_{u_{\delta} \in X^{\delta}} \|u - u_{\delta}\|_{X} \quad (\delta \in \Delta).$$

# 8.3 Solving efficiently on tensor-product discretizations

From now on, we assume that  $X^{\delta} := X_t^{\delta} \otimes X_x^{\delta}$  and  $Y^{\delta} := Y_t^{\delta} \otimes Y_x^{\delta}$  are *tensor-products*, and for ease of presentation, we assume that the spatial discretizations on  $X^{\delta}$  and  $Y^{\delta}$  coincide, i.e.  $X_x^{\delta} = Y_x^{\delta}$ , reducing (8.5) to  $X_t^{\delta} \subset Y_t^{\delta}$  and  $\frac{\mathrm{d}}{\mathrm{d}t} X_t^{\delta} \subset Y_t^{\delta}$ .

We equip  $X_t^{\delta}$  with a basis  $\Phi_t^{\delta}$ ,  $X_{\mathbf{x}}^{\delta}$  with  $\Phi_{\mathbf{x}}^{\delta}$ , and  $Y_t^{\delta}$  with  $\Xi^{\delta}$ .

## 8.3.1 Construction of $K_Y^{\delta}$

Define  $O := \langle \Xi^{\delta}, \Xi^{\delta} \rangle_{L_2(I)}$  and  $A_{\mathbf{x}} := \langle \Phi_{\mathbf{x}}^{\delta}, \Phi_{\mathbf{x}}^{\delta} \rangle_V$ . Given  $K_{\mathbf{x}} = A_{\mathbf{x}}^{-1}$  uniformly in  $\delta \in \Delta$ , define

$$K_Y := O^{-1} \otimes K_x.$$

Then, the preconditioner  $K_Y^{\delta} := \mathcal{F}_{\Xi^{\delta} \otimes \Phi_{\mathbf{x}}^{\delta}} \mathbf{K}_Y(\mathcal{F}_{\Xi^{\delta} \otimes \Phi_{\mathbf{x}}^{\delta}})' \in \mathcal{L}(Y^{\delta'}, Y^{\delta})$  satisfies (8.9); cf. §6.5.6.

When  $\Xi^{\delta}$  is orthogonal, O is diagonal and can be inverted exactly. For standard finite element bases  $\Phi_{\mathbf{x}}^{\delta}$ , suitable  $\mathbf{K}_{\mathbf{x}}$  that can be applied efficiently (at cost linear in the discretization size) are provided by symmetric multigrid methods.

### 8.3.2 Preconditioning the Schur complement formulation

We will solve a matrix representation of (8.8) with an iterative solver, thus requiring a preconditioner. Inspired by the constructions of [And16, NS19], we build an *optimal* self-adjoint coercive preconditioner  $K_X^{\delta} \in \mathcal{L}(X^{\delta'}, X^{\delta})$  as a wavelet-in-time block-diagonal matrix with multigrid-in-space blocks.

Let *U* be a separable Hilbert space of functions over some domain. A given collection  $\Psi = {\psi_{\lambda}}_{\lambda \in \vee_{\Psi}}$  is a *Riesz basis* for *U* when

$$\overline{\operatorname{span}\Psi} = U$$
, and  $\|\boldsymbol{c}\|_{\ell_2(\vee_{\Psi})} \approx \|\boldsymbol{c}^{\top}\Psi\|_U$  for all  $\boldsymbol{c} \in \ell_2(\vee_{\Psi})$ .

Thinking of  $\Psi$  being a basis of wavelet-type, for indices  $\lambda \in \vee_{\Psi}$ , its *level* is denoted  $|\lambda| \in \mathbb{N}_0$ . We call  $\Psi$  *uniformly local* when for all  $\lambda \in \vee_{\Psi}$ ,

diam $(\operatorname{supp} \psi_{\lambda}) \lesssim 2^{-|\lambda|}$  and  $\#\{\mu \in \vee_{\Psi} : |\mu| = |\lambda|, |\operatorname{supp} \psi_{\mu} \cap \operatorname{supp} \psi_{\lambda}| > 0\} \lesssim 1.$ 

Assume  $\Sigma := \{\sigma_{\lambda} : \lambda \in \vee_{\Sigma}\}$  is a uniformly local Riesz basis for  $L_2(I)$  with  $\{2^{-|\lambda|}\sigma_{\lambda} : \lambda \in \vee_{\Sigma}\}$  Riesz for  $H^1(I)$ . Writing  $w \in X$  as  $\sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes w_{\lambda}$  for some  $w_{\lambda} \in V$ , we define the bounded, symmetric, and coercive bilinear form

$$(D_X \sum_{\lambda \in \vee_{\Sigma}} \sigma_{\lambda} \otimes w_{\lambda}) (\sum_{\mu \in \vee_{\Sigma}} \sigma_{\mu} \otimes v_{\mu}) \coloneqq \sum_{\lambda \in \vee_{\Sigma}} \langle w_{\lambda}, v_{\lambda} \rangle_{V} + 4^{|\lambda|} \langle w_{\lambda}, v_{\lambda} \rangle_{V'}.$$

The operator  $D_X^{\delta} := E_X^{\delta'} D_X E_X^{\delta}$  is in  $\mathcal{L}is(X^{\delta}, X^{\delta'})$ . Its norm and that of its inverse are bounded uniformly in  $\delta \in \Delta$ . When  $X^{\delta} = \operatorname{span} \Sigma^{\delta} \otimes \Phi_{\mathbf{x}}^{\delta}$  for some  $\Sigma^{\delta} := \{\sigma_{\lambda} : \lambda \in \vee_{\Sigma^{\delta}}\} \subset \Sigma$ , the matrix representation of  $D_X^{\delta}$  w.r.t.  $\Sigma^{\delta} \otimes \Phi_{\mathbf{x}}^{\delta}$  is

$$(\mathcal{F}_{\Sigma^{\delta}\otimes\Phi^{\delta}})'D_{X}^{\delta}\mathcal{F}_{\Sigma^{\delta}\otimes\Phi^{\delta}} =: \mathbf{D}_{X}^{\delta} = \text{blockdiag}[\mathbf{A}_{\mathbf{x}} + 4^{|\lambda|}\langle \Phi_{\mathbf{x}}^{\delta}, \Phi_{\mathbf{x}}^{\delta} \rangle_{V'}]_{\lambda \in \vee_{\Sigma^{\delta}}}.$$

**Theorem 8.3.1** (§6.5.6). Define  $M_{\mathbf{x}} := \langle \Phi_{\mathbf{x}}^{\delta}, \Phi_{\mathbf{x}}^{\delta} \rangle_{H}$ . When we have matrices  $K_{j} = (\mathbf{A}_{\mathbf{x}} + 2^{j} \mathbf{M}_{\mathbf{x}})^{-1}$  uniformly in  $\delta \in \Delta$  and  $j \in \mathbb{N}_{0}$ , it follows that

$$\mathbf{D}_X^{-1} = \mathbf{K}_X := \operatorname{blockdiag}[\mathbf{K}_{|\lambda|} \mathbf{A}_{\mathbf{x}} \mathbf{K}_{|\lambda|}]_{\lambda \in \vee_{\Sigma^{\delta}}}$$

This yields an optimal preconditioner  $K_X^{\delta} := \mathcal{F}_{\Sigma^{\delta} \otimes \Phi^{\delta}} \mathbf{K}_X (\mathcal{F}_{\Sigma^{\delta} \otimes \Phi^{\delta}})' \in \mathcal{L}is(X^{\delta'}, X^{\delta}).$ 

In [OR00] it was shown that under a 'full-regularity' assumption, for quasiuniform meshes, a multiplicative multigrid method yields  $K_j$  satisfying the conditions of Thm. 8.3.1, which can moreover be applied in linear time.

### 8.3.3 Wavelets in time

The preconditioner  $K_X$  requires  $X_t^{\delta}$  to be equipped with a *wavelet* basis  $\Sigma^{\delta}$ , whereas one typically uses a different (single-scale) basis  $\Phi_t^{\delta}$  on  $X_t^{\delta}$ . To bridge this gap, a basis transformation from  $\Sigma^{\delta}$  to  $\Phi_t^{\delta}$  is required. We define the wavelet transform as  $W_t := (\mathcal{F}_{\Phi^{\delta}})^{-1} \mathcal{F}_{\Sigma^{\delta}}$ .<sup>1</sup>

Define  $V_j := \operatorname{span}\{\sigma_{\lambda} \in \Sigma : |\lambda| \leq j\}$ . Equip each  $V_j$  with a (single-scale) basis  $\Phi_j$ , and assume that  $\Phi_t^{\delta} := \Phi_J$  for some J, so that  $X_t^{\delta} := V_J$ . Since  $V_{j+1} = V_j \oplus \operatorname{span} \Sigma_j$  where  $\Sigma_j := \{\sigma_{\lambda} : |\lambda| = j\}$ , there exist matrices  $P_j$  and  $Q_j$  such that  $\Phi_i^{\top} = \Phi_{j+1}^{\top} P_j$  and  $\Psi_j^{\top} = \Phi_{j+1}^{\top} Q_j$ , with  $M_j := [P_j|Q_j]$  invertible.

Writing  $v \in V_J$  in both forms  $v = c_0^\top \Phi_0 + \sum_{j=0}^{J-1} d_j^\top \Psi_j$  and  $v = c_J^\top \Phi_J$ , the basis transformation  $W_t := W_J$  mapping wavelet coordinates  $(c_0^\top, d_0^\top, \dots, d_{J-1}^\top)$  to single-scale coordinates  $c_J$  satisfies

(8.10) 
$$W_J = M_{J-1} \begin{bmatrix} W_{J-1} & \mathbf{0} \\ \mathbf{0} & \mathrm{Id} \end{bmatrix}$$
, and  $W_0 := \mathrm{Id}$ .

Uniform locality of  $\Sigma$  implies *uniform sparsity* of the  $M_j$ , i.e. with  $\mathcal{O}(1)$  nonzeros per row and column. Then, assuming a geometrical increase in dim  $V_j$  in terms of j, which is true in the concrete setting below, matrix-vector products  $\boldsymbol{x} \mapsto \boldsymbol{W}_t \boldsymbol{x}$  can be performed (serially) in linear complexity; cf. [Ste03b].

### 8.3.4 Solving the system

The matrix representation of  $S^{\delta}$  and  $f^{\delta}$  from (8.8) w.r.t. a basis  $\Phi_t^{\delta} \otimes \Phi_x^{\delta}$  of  $X^{\delta}$  is

$$oldsymbol{S} \mathrel{\mathop:}= (\mathcal{F}_{\Phi^\delta_t\otimes\Phi^\delta_{\mathbf{x}}})'S^\delta\mathcal{F}_{\Phi^\delta_t\otimes\Phi^\delta_{\mathbf{x}}} \quad ext{and} \quad oldsymbol{f} \mathrel{\mathop:}= (\mathcal{F}_{\Phi^\delta_t\otimes\Phi^\delta_{\mathbf{x}}})'f^\delta.$$

Envisioning an iterative solver, using §8.3.2 we have a preconditioner in terms of the wavelet-in-time basis  $\Sigma^{\delta} \otimes \Phi_{\mathbf{x}}^{\delta}$ , with which their matrix representation is

(8.11) 
$$\hat{\boldsymbol{S}} := (\mathcal{F}_{\Sigma^{\delta} \otimes \Phi_{\mathbf{x}}^{\delta}})' S^{\delta} \mathcal{F}_{\pm^{\delta} \otimes \oplus_{\mathbf{x}}^{\delta}} \text{ and } \hat{\boldsymbol{f}} := (\mathcal{F}_{\Sigma^{\delta} \otimes \Phi_{\mathbf{x}}^{\delta}})' f^{\delta}.$$

These two forms are related: with the wavelet transform  $W := W_t \otimes \text{Id}_x$ , we have  $\hat{S} = W^{\top}SW$  and  $\hat{f} = W^{\top}f$ , and the matrix representation of (8.8) becomes

(8.12) finding  $\boldsymbol{w}$  s.t.  $\hat{\boldsymbol{S}}\boldsymbol{w} = \hat{\boldsymbol{f}}$ .

<sup>&</sup>lt;sup>1</sup>In literature, this transform is typically called an *inverse wavelet transform*.

We can then recover the solution in single-scale coordinates as u = Ww.

We use Preconditioned Conjugate Gradients (PCG), with preconditioner  $K_X$ , to solve (8.12). Given an algebraic error tolerance  $\varepsilon > 0$  and current guess  $w_k$ , we monitor  $r_k^\top K_X r_k \leq \varepsilon^2$  where  $r_k := \hat{f} - \hat{S} w_k$ . This data is available within PCG, and constitutes a stopping criterium: with  $u_k^{\delta} := \mathcal{F}_{\Sigma^{\delta} \otimes \Phi_x^{\delta}} w_k \in X^{\delta}$ , we see

(8.13) 
$$\boldsymbol{r}_{k}^{\top}\boldsymbol{K}_{X}\boldsymbol{r}_{k} = (f^{\delta} - S^{\delta}u_{k}^{\delta})(K_{X}^{\delta}(f^{\delta} - S^{\delta}u_{k}^{\delta})) \approx \|u^{\delta} - u_{k}^{\delta}\|_{X}^{2}$$

where  $\approx$  follows from (6.34), so that the algebraic error satisfies  $||u^{\delta} - u_k^{\delta}||_X \lesssim \varepsilon$ .

### 8.4 A concrete setting: the reaction-diffusion equation

On a bounded Lipschitz domain  $\Omega \subset \mathbb{R}^d$ , take  $H := L_2(\Omega)$ ,  $V := H_0^1(\Omega)$ , and

$$a(t;\eta,\zeta) := \int_{\Omega} \boldsymbol{D} \nabla \eta \cdot \nabla \zeta + c \eta \zeta \mathrm{d} \mathbf{x}$$

where  $\boldsymbol{D} = \boldsymbol{D}^{\top} \in \mathbb{R}^{d \times d}$  is positive definite, and  $c \ge 0.^2$  We note that A(t) is symmetric and coercive. W.l.o.g. we take I := (0, 1), i.e. T := 1.

Fix  $p_t, p_x \in \mathbb{N}$ . With  $\{\mathcal{T}_I\}$  the family of quasi-uniform partitions of I into subintervals, and  $\{\mathcal{T}_\Omega\}$  that of conforming quasi-uniform triangulations of  $\Omega$ , we define  $\Delta$  as the collection of pairs  $(\mathcal{T}_I, \mathcal{I}_\Omega)$ . We construct our trial- and test spaces as

$$X^{\delta} := X^{\delta}_t \otimes X^{\delta}_{\mathbf{x}}, \quad Y^{\delta} := Y^{\delta}_t \otimes X^{\delta}_{\mathbf{x}},$$

where, with  $\mathbb{P}_p^{-1}(\mathcal{T})$  denoting the space of piecewise degree-p polynomials on  $\mathcal{T}$  ,

$$X_t^{\delta} := H^1(I) \cap \mathbb{P}_{p_t}^{-1}(\mathcal{T}_I), \quad X_{\mathbf{x}}^{\delta} := H_0^1(\Omega) \cap \mathbb{P}_{p_{\mathbf{x}}}^{-1}(\mathcal{T}_\Omega), \quad Y_t^{\delta} := \mathbb{P}_{p_t}^{-1}(\mathcal{T}_I).$$

These spaces satisfy condition (8.5), with coinciding spatial discretizations on  $X^{\delta}$  and  $Y^{\delta}$ . For this choice of  $\Delta$ , inf-sup condition (8.6) follows from [SW21b, Thm. 4.3].

For  $X_t^{\delta}$ , we choose  $\Phi_t^{\delta}$  to be the Lagrange basis of degree  $p_t$  on  $\mathcal{T}_I$ ; for  $X_{\mathbf{x}}^{\delta}$ , we choose  $\Phi_{\mathbf{x}}^{\delta}$  to be that of degree  $p_{\mathbf{x}}$  on  $\mathcal{T}_{\Omega}$ . An orthogonal basis  $\Xi^{\delta}$  for  $Y_t^{\delta}$  may be built as piecewise shifted Legendre polynomials of degree  $p_t$  w.r.t.  $\mathcal{T}_I$ .

For  $p_t = 1$ , one finds a suitable wavelet basis  $\Sigma$  in [Ste98]. For  $p_t > 1$ , one can either split the system into lowest- and higher-order parts and perform the transform on the lowest-order part only, or construct higher-order wavelets directly; cf. [Dij09].

Owing to the tensor-product structure of  $X^{\delta}$  and  $Y^{\delta}$  and of the operators A and  $\partial_t$ , the matrix representation of our formulation becomes remarkably simple.

<sup>&</sup>lt;sup>2</sup>This is easily generalized to variable coefficients, but notation becomes more obtuse.

**Lemma 8.4.1.** Define  $g := (\mathcal{F}_{\Xi^{\delta} \otimes \Phi^{\delta}_{\mathbf{x}}})'g$ ,  $u_0 := \Phi^{\delta}_t(0) \otimes \langle u_0, \Phi^{\delta}_{\mathbf{x}} \rangle_{L_2(\Omega)}$ , and

$$\begin{split} \boldsymbol{T} &:= \langle \frac{\mathrm{d}}{\mathrm{d}t} \boldsymbol{\Phi}_t^{\delta}, \Xi^{\delta} \rangle_{L_2(I)}, & \boldsymbol{N} &:= \langle \boldsymbol{\Phi}_t^{\delta}, \Xi^{\delta} \rangle_{L_2(I)}, \\ \boldsymbol{\Gamma}_0 &:= \boldsymbol{\Phi}_t^{\delta}(0) [\boldsymbol{\Phi}_t^{\delta}(0)]^{\top}, & \boldsymbol{M}_{\mathbf{x}} &:= \langle \boldsymbol{\Phi}_{\mathbf{x}}^{\delta}, \boldsymbol{\Phi}_{\mathbf{x}}^{\delta} \rangle_{L_2(\Omega)}, \\ \boldsymbol{A}_{\mathbf{x}} &:= \langle \boldsymbol{D} \nabla \boldsymbol{\Phi}_{\mathbf{x}}^{\delta}, \nabla \boldsymbol{\Phi}_{\mathbf{x}}^{\delta} \rangle_{L_2(\Omega)} + c \boldsymbol{M}_{\mathbf{x}}, & \boldsymbol{B} &:= \boldsymbol{T} \otimes \boldsymbol{M}_{\mathbf{x}} + \boldsymbol{N} \otimes \boldsymbol{A}_{\mathbf{x}}. \end{split}$$

With  $K_Y := O^{-1} \otimes K_x$  from §8.3.1, we can write S and f from §8.3.4 as

$$oldsymbol{S} = oldsymbol{B}^ op oldsymbol{K}_Y oldsymbol{B} + oldsymbol{\Gamma}_0 \otimes oldsymbol{M}_{\mathbf{x}}, \quad oldsymbol{f} = oldsymbol{B}^ op oldsymbol{K}_Y oldsymbol{g} + oldsymbol{u}_0.$$

Note that N and T are non-square,  $\Gamma_0$  is very sparse, and T is bidiagonal.

In fact, assumption (8.5) allows us to write S in an even simpler form.

**Lemma 8.4.2.** *The matrix S can be written as* 

$$egin{aligned} oldsymbol{S} &= oldsymbol{A}_t \otimes (oldsymbol{M}_{\mathbf{x}}oldsymbol{K}_{\mathbf{x}}oldsymbol{M}_{\mathbf{x}}) + oldsymbol{M}_t \otimes (oldsymbol{A}_{\mathbf{x}}oldsymbol{K}_{\mathbf{x}}oldsymbol{A}_{\mathbf{x}}) + oldsymbol{L} \otimes (oldsymbol{A}_{\mathbf{x}}oldsymbol{K}_{\mathbf{x}}oldsymbol{M}_{\mathbf{x}}) + oldsymbol{L} \otimes (oldsymbol{A}_{\mathbf{x}}oldsymbol{K}_{\mathbf{x}}oldsymbol{M}_{\mathbf{x}}) + oldsymbol{\Gamma}_0 \otimes oldsymbol{M}_{\mathbf{x}} \end{aligned}$$

where

$$\boldsymbol{L} := \langle \frac{\mathrm{d}}{\mathrm{d}t} \Phi_t^{\delta}, \Phi_t^{\delta} \rangle_{L_2(I)}, \quad \boldsymbol{M}_t := \langle \Phi_t^{\delta}, \Phi_t^{\delta} \rangle_{L_2(I)}, \quad \boldsymbol{A}_t := \langle \frac{\mathrm{d}}{\mathrm{d}t} \Phi_t^{\delta}, \frac{\mathrm{d}}{\mathrm{d}t} \Phi_t^{\delta} \rangle_{L_2(I)},$$

This matrix representation does not depend on  $Y_t^{\delta}$  or  $\Xi^{\delta}$  at all.

*Proof.* The expansion of  $B := T \otimes M_x + N \otimes A_x$  in S yields a sum of five Kronecker products, one of which is

$$(oldsymbol{T}^{ op}\otimes oldsymbol{M}_{\mathbf{x}})oldsymbol{K}_{Y}(oldsymbol{T}\otimes oldsymbol{A}_{\mathbf{x}})=(oldsymbol{T}^{ op}oldsymbol{O}^{-1}oldsymbol{N})\otimes(oldsymbol{M}_{\mathbf{x}}oldsymbol{K}_{\mathbf{x}}oldsymbol{A}_{\mathbf{x}}).$$

We will show that  $T^{\top}O^{-1}N = L^{\top}$ ; similar arguments hold for the other terms. Thanks to  $X_t^{\delta} \subset Y_t^{\delta}$ , we can define the trivial embedding  $F_t^{\delta} : X_t^{\delta} \to Y_t^{\delta}$ . Defining

$$T^{\delta} \colon X_t^{\delta} \to Y_t^{\delta'}, \quad (T^{\delta}u)(v) := \langle \frac{\mathrm{d}}{\mathrm{d}t}u, v \rangle_{L_2(I)}, M^{\delta} \colon Y_t^{\delta} \to Y_t^{\delta'}, \quad (M^{\delta}u)(v) := \langle u, v \rangle_{L_2(I)},$$

we find  $\boldsymbol{O} = (\mathcal{F}_{\Xi^{\delta}})' M^{\delta} \mathcal{F}_{\Xi^{\delta}}$ ,  $\boldsymbol{N} = (\mathcal{F}_{\Xi^{\delta}})' M^{\delta} \mathcal{F}_{t}^{\delta} \mathcal{F}_{\Phi_{t}^{\delta}}$  and  $\boldsymbol{T} = (\mathcal{F}_{\Xi^{\delta}})' T^{\delta} \mathcal{F}_{\Phi_{t}^{\delta}}$ , so

$$\boldsymbol{T}^{\top}\boldsymbol{O}^{-1}\boldsymbol{N} = (\mathcal{F}_{\Phi_t^{\delta}})'T^{\delta'}F_t^{\delta}\mathcal{F}_{\Phi_t^{\delta}} = \langle \Phi_t, \frac{\mathrm{d}}{\mathrm{d}t}\Phi_t \rangle_{L_2(I)} = \boldsymbol{L}^{\top}.$$

#### 8.4.1 Parallel complexity

The *parallel complexity* of our algorithm is the asymptotic runtime of solving (8.12) for  $u \in \mathbb{R}^{N_t N_x}$  in terms of  $N_t := \dim X_t^{\delta}$  and  $N_x := \dim X_x^{\delta}$ , given sufficiently many parallel processors and assuming no communication cost.

We understand the serial (resp. parallel) cost of a matrix B, denoted  $C_B^s$  (resp.  $C_B^p$ ), as the asymptotic runtime of performing  $x \mapsto Bx \in \mathbb{R}^N$  in terms of N, on a single (resp. sufficiently many) processors at no communication cost. For *uniformly sparse* matrices, i.e. with  $\mathcal{O}(1)$  nonzeros per row and column, the serial cost is  $\mathcal{O}(N)$ , and the parallel cost is  $\mathcal{O}(1)$  by computing each cell of the output concurrently.

From Theorem 8.3.1, we see that  $K_X$  is such that  $\kappa_2(K_X \hat{S}) \lesssim 1$  uniformly in  $\delta \in \Delta$ . Therefore, for a given algebraic error tolerance  $\varepsilon$ , we require  $\mathcal{O}(\log \varepsilon^{-1})$  PCG iterations. Assuming that the parallel cost of matrices dominates that of vector addition and inner products, the parallel complexity of a single PCG iteration is dominated by the cost of applying  $K_X$  and  $\hat{S}$ . As  $\hat{S} = W^{\top}SW$ , our algorithm runs in complexity

(8.14) 
$$\mathcal{O}(\log \varepsilon^{-1}[C^{\circ}_{\boldsymbol{K}_{\boldsymbol{X}}} + C^{\circ}_{\boldsymbol{W}^{\top}} + C^{\circ}_{\boldsymbol{S}} + C^{\circ}_{\boldsymbol{W}}]) \quad (\circ \in \{s, p\}).$$

**Theorem 8.4.3.** For fixed algebraic error tolerance  $\varepsilon > 0$ , our algorithm runs in

- serial complexity  $\mathcal{O}(N_t N_{\mathbf{x}})$ ;
- time-parallel complexity  $\mathcal{O}(\log(N_t)N_{\mathbf{x}})$ ;
- space-time-parallel complexity  $\mathcal{O}(\log(\mathcal{N}_{\sqcup}\mathcal{N}_{\mathbf{x}}))$ .

*Proof.* We absorb the constant factor  $\log \varepsilon^{-1}$  of (8.14) into  $\mathcal{O}$ . We analyse the cost of every matrix separately.

#### The (inverse) wavelet transform

As  $W = W_t \otimes \operatorname{Id}_{\mathbf{x}}$ , its serial cost equals  $\mathcal{O}(C^s_{W_t}N_{\mathbf{x}})$ . The choice of wavelet allows performing  $\mathbf{x} \mapsto W_t \mathbf{x}$  at linear serial cost (cf. §8.3.3), so that  $C^s_W = \mathcal{O}(N_t N_{\mathbf{x}})$ .

Using (8.10), we write  $W_t$  as the composition of J matrices, each uniformly sparse and hence at parallel cost  $\mathcal{O}(1)$ . Because the mesh in time is quasiuniform, we have  $J = \log N_t$ . We find that  $C_{W_t}^p = \mathcal{O}(J) = \mathcal{O}(\log N_t)$ , so that the time-parallel cost of W equals  $\mathcal{O}(\log(N_t)N_x)$ . By exploiting spatial parallelism as well, we find  $C_W^p = \mathcal{O}(\log N_t)$ . Analogous arguments hold for  $W_t^{\top}$  and  $W^{\top}$ .

#### The preconditioner

Recall that  $K_X := \text{blockdiag}[K_{|\lambda|}A_xK_{|\lambda|}]_{\lambda}$ . Since the cost of  $K_j$  is independent of j, we see that

$$C_{\boldsymbol{K}_{X}}^{s} = \mathcal{O}\left(N_{t} \cdot \left(2C_{\boldsymbol{K}_{j}}^{s} + C_{\boldsymbol{A}_{\mathbf{x}}}^{s}\right)\right) = \mathcal{O}(2N_{t}C_{\boldsymbol{K}_{j}}^{s} + N_{t}N_{\mathbf{x}}).$$

Implementing the  $K_j$  as typical multiplicative multigrid solvers with linear serial cost, we find  $C^s_{K_x} = O(N_t N_x)$ .

Through temporal parallelism, we can apply each block of  $K_X$  concurrently, resulting in a time-parallel cost of  $\mathcal{O}(2C_{K_i}^s + C_{A_x}^s) = \mathcal{O}(N_x)$ .

By parallelizing in space as well, we reduce the cost of the uniformly sparse  $A_{\mathbf{x}}$  to  $\mathcal{O}(1)$ . The parallel cost of multiplicative multigrid on quasi-uniform triangulations is  $\mathcal{O}(\log N_{\mathbf{x}})$ ; cf. [MFL+91]. It follows that  $C_{K_{\mathbf{x}}}^{p} = \mathcal{O}(\log N_{\mathbf{x}})$ .

#### The Schur matrix

Using Lemma 8.4.1, we write  $S = B^{\top} K_Y B + \Gamma_0 \otimes M_x$  where  $B = T \otimes M_x + N \otimes A_x$ , which immediately reveals that

$$C_{\boldsymbol{S}}^{s} = C_{\boldsymbol{B}^{\top}}^{s} + C_{\boldsymbol{K}_{Y}}^{s} + C_{\boldsymbol{B}}^{s} + C_{\boldsymbol{\Gamma}_{0}}^{s} \cdot C_{\boldsymbol{M}}^{s} = \mathcal{O}(N_{t}N_{\mathbf{x}} + C_{\boldsymbol{K}_{Y}}^{s}), \quad \text{and}$$
$$C_{\boldsymbol{S}}^{p} = \max\left\{C_{\boldsymbol{B}^{\top}}^{p} + C_{\boldsymbol{K}_{Y}}^{p} + C_{\boldsymbol{B}}^{p}, \ C_{\boldsymbol{\Gamma}_{0}}^{p} \cdot C_{\boldsymbol{M}}^{p}\right\} = \mathcal{O}(C_{\boldsymbol{K}_{Y}}^{p})$$

because every matrix except  $K_Y$  is uniformly sparse. With arguments similar to the previous paragraph, we see that  $K_Y$  (and hence S) has serial cost  $\mathcal{O}(N_tN_x)$ , time-parallel cost  $\mathcal{O}(N_x)$ , and space-time-parallel cost  $\mathcal{O}(\log N_x)$ .

### 8.4.2 Solving to higher accuracy

Instead of *fixing* the algebraic error tolerance, maybe more realistic is is to desire a solution  $\tilde{u}^{\delta} \in X^{\delta}$  for which the error is proportional to the discretization error, i.e.  $\|u - \tilde{u}^{\delta}\|_X \lesssim \inf_{u_{\delta} \in X^{\delta}} \|u - u_{\delta}\|_X$ .

Assuming that this error decays with a (problem-dependent) rate s > 0, i.e.  $\inf_{u_{\delta} \in X^{\delta}} ||u - u_{\delta}||_X \leq (N_t N_x)^{-s}$ , then the same holds for the solution  $u^{\delta}$  of (8.8); cf. Thm. 8.2.3. When the algebraic error tolerance decays as  $\varepsilon \leq (N_t N_x)^{-s}$ , a triangle inequality and (8.13) show that the error of our solution  $\tilde{u}^{\delta}$  obtained by PCG decays at rate *s* too.

In this case,  $\log \epsilon^{-1} = \mathcal{O}(\log(N_t N_{\mathbf{x}}))$ . From (8.14) and the proof of Theorem 8.4.3, we find our algorithm to run in superlinear serial complexity  $\mathcal{O}(N_t N_{\mathbf{x}} \log(N_t N_{\mathbf{x}}))$ , time-parallel complexity  $\mathcal{O}(\log^2(N_t) \log(N_{\mathbf{x}})N_{\mathbf{x}})$ , and polylogarithmic complexity  $\mathcal{O}(\log^2(N_t N_{\mathbf{x}}))$  parallel in space and time.

For elliptic PDEs, algorithms are available that offer quasi-optimal solutions, serially in linear complexity  $O(N_x)$ —the cost of a serial solve to *fixed* algebraic error—and in parallel in  $O(\log^2 N_x)$ , by combining a *nested iteration* with parallel multigrid; cf. [Hac85, Ch. 5] and [Bra81].

In [HVW95], the question is posed whether "good serial algorithms for parabolic PDEs are intrinsically as parallel as good serial algorithms for elliptic PDEs", basically asking if the lower bound of  $O(\log^2(N_t N_x))$  can be attained by an algorithm that runs serially in  $O(N_t N_x)$ ; see [Wor91, §2.2] for a formal discussion.

Nested iteration drives down the serial complexity of our algorithm to a linear  $\mathcal{O}(N_t N_{\mathbf{x}})$ , and also improves the time-parallel complexity to  $\mathcal{O}(\log(N_t)N_{\mathbf{x}})$ .<sup>3</sup> This is on par with the best-known results for elliptic problems, so we answer the question posed in [HVW95] in the affirmative.

### 8.5 Numerical experiments

We take the simple heat equation, i.e.  $D = \text{Id}_x$  and c = 0. We select  $p_t = p_x = 1$ , i.e. lowest order finite elements in space and time. We will use the 3-point wavelet introduced in [Ste98].

We implemented our algorithm in Python using the open source finite element library NGSolve [Sch14] for meshing and discretization of the bilinear forms in space and time, MPI through mpi4py [DPS05] for distributed computations, and SciPy [Vir20] for the sparse matrix-vector computations. The source code is available at [vVW21c].

#### 8.5.1 Preconditioner calibration on a 2D problem

Our preconditioner is optimal, meaning that  $\kappa_2(K_X \hat{S}) \lesssim 1$ . Here we will investigate this condition number quantitatively.

As a model problem, we partition the temporal interval I uniformly into  $2^J$  subintervals. We consider the domain  $\Omega := [0, 1]^2$ , and triangulate it uniformly into  $4^K$  triangles. We set  $N_t := \dim X_t^{\delta} = 2^J + 1$  and  $N_{\mathbf{x}} := \dim X_{\mathbf{x}}^{\delta} = (2^K - 1)^2$ .

We start by using direct inverses  $K_j = (A_x + 2^j M_x)^{-1}$  and  $K_x = A_x^{-1}$  to determine the best possible condition numbers. We found that replacing  $K_j$  by  $K_j^{\alpha} = (\alpha A_x + 2^j M_x)^{-1}$  for  $\alpha = 0.3$  gave better conditioning; see also the left of Table 8.1. At the right of Table 8.1, we see that the condition numbers are very robust with respect to spatial refinements, but less so for refinements in time. Still, at  $N_t = 16129$ , we observe a modest  $\kappa_2(K_X \hat{S})$  of 8.74.

Replacing the direct inverses with multigrid solvers, we found a good balance between speed and conditioning at 2 V-cycles with 3 Gauss-Seidel smoothing steps per grid. We decided to use these for our experiments.

<sup>&</sup>lt;sup>3</sup>Interestingly, nested iteration offers no improvements parallel in space *and* time, with complexity still  $O(\log^2(N_t N_x))$ .

50 -	/	$N_t$	= 65	129	257	513	1025	2049	4097	8 193
40 -		$N_{\mathbf{x}} = 49$	6.34	7.05	7.53	7.89	8.15	8.37	8.60	8.78
20		225	6.33	6.89	7.55	7.91	8.14	8.38	8.57	8.73
30 -		961	6.14	6.89	7.55	7.93	8.15	8.38	8.57	8.74
20 -		3 969	6.14	7.07	7.56	7.87	8.16	8.38	8.57	8.74
10		16129	6.14	6.52	7.55	7.86	8.16	8.37	8.57	8.74
10 1										
0.0	0.2 0.4 0.6 0.8 1.0									

TABLE 8.1. Computed condition numbers  $\kappa_2(\mathbf{K}_X \hat{\mathbf{S}})$ . Left: fixed  $N_t = 1025$ ,  $N_x = 961$  for varying  $\alpha$ . Right: fixed  $\alpha = 0.3$  for varying  $N_t$  and  $N_x$ .

### 8.5.2 Time-parallel results

α

We perform computations on Cartesius, the Dutch supercomputer. Each Cartesius node has 64GB of memory and 12 cores (at 2 threads per core) running at 2.6GHz. Using the preconditioner detailed above, we iterate PCG on (8.12) with *S* computed as in Lemma 8.4.2, until achieving an algebraic error of  $\varepsilon = 10^{-6}$ ; see also §8.3.4. For the spatial multigrid solvers, we use 2 V-cycles with 3 Gauss-Seidel smoothing steps per grid.

#### Memory-efficient time-parallel implementation

For  $X \in \mathbb{R}^{N_x \times N_t}$ , we define  $(X) \in \mathbb{R}^{N_t N_x}$  as the vector obtained by stacking columns of X vertically. For memory efficency, we do not build matrices of the form  $B_t \otimes B_x$  appearing in Lemma 8.4.2 directly, but instead perform matrix-vector products using the identity

$$(8.15) \qquad (\boldsymbol{B}_t \otimes \boldsymbol{B}_{\mathbf{x}})(\boldsymbol{X}) = (\boldsymbol{B}_{\mathbf{x}}(\boldsymbol{B}_t \boldsymbol{X}^{\top})^{\top}) = (\mathrm{Id}_t \otimes \boldsymbol{B}_{\mathbf{x}})(\boldsymbol{B}_t \boldsymbol{X}^{\top}).$$

Each parallel processor stores only a subset of the temporal degrees of freedom, e.g. a subset of columns of X. When  $B_t$  is uniformly sparse, which holds true for all of our temporal matrices, using (8.15) we can evaluate  $(B_t \otimes B_x)(X)$ in  $\mathcal{O}(C^s_{B_x})$  operations parallel in time: on each parallel processor, we compute 'our' columns of  $Y := B_t X^{\top}$  by receiving the necessary columns of X from neighbouring processors, and then compute  $B_x Y^{\top}$  without communication.

The preconditioner  $K_X$  is block-diagonal, making its time-parallel application trivial. Representing the wavelet transform of §8.3.3 as the composition of J Kronecker products allows a time-parallel implementation using the above.

#### 2D problem

We select  $\Omega := [0,1]^2$  with a uniform triangulation  $\mathcal{T}_{\Omega}$ , and we triangulate I uniformly into  $\mathcal{T}_I$ . We prescribe the smooth solution  $u(t, x, y) := \exp(-2\pi^2 t) \sin(\pi x) \sin(\pi y)$ , so the problem has vanishing forcing data g.

Table 8.2 details the strong scaling results, i.e. fixing the problem size and increasing the number of processors P. We triangulate I into  $2^{14}$  time slabs, yielding  $N_t = 16\,385$  temporal degrees of freedom, and  $\Omega$  into  $4^8$  triangles, yielding a  $X_x^{\delta}$  of dimension  $N_x = 65\,025$ . The resulting system contains  $1\,065\,434\,625$  degrees of freedom and our solver reaches the algebraic error tolerance after 16 iterations. In perfect strong scaling, the total number of CPU-hours remains constant. Even at 2 048 processors, we observe a parallel efficiency of around 92.9%, solving this system in a modest 11.7 CPU-hours. Acquiring strong scaling results on a single node was not possible due to memory limitations.

Table 8.3 details the weak scaling results, i.e. fixing the problem size per processor and increasing the number of processors. In perfect weak scaling, the time per iteration should remain constant. We observe a slight increase in time per iteration on a single node, but when scaling to multiple nodes, we observe a near-perfect parallel efficiency of around 96.7%, solving the final system with 4 278 467 585 degrees of freedom in a mere 109 seconds.

#### 3D problem

We select  $\Omega := [0,1]^3$  with  $u(t, x, y, z) := \exp(-3\pi^2 t) \sin(\pi x) \sin(\pi y) \sin(\pi z)$ , so the problem has vanishing forcing data *g*.

Table 8.4 shows the strong scaling results. We triangulate *I* uniformly into  $2^{14}$  time slabs, and  $\Omega$  uniformly into  $8^6$  tetrahedra. The arising system has  $N = 4\,097\,020\,095$  unknowns, which we solve to tolerance in 18 iterations. The results are comparable to those in two dimensions, albeit a factor two slower at similar problem sizes.

Table 8.5 shows the weak scaling results for the 3D problem. As in the two-dimensional case, we observe excellent scaling properties, and see that the time per iteration is nearly constant.

### 8.6 Conclusion

We have presented a framework for solving linear parabolic evolution equations massively in parallel. Based on earlier ideas [And16, NS19, SW21b], we found a remarkably simple symmetric Schur-complement equation. With a tensor-product discretization of the space-time cylinder using standard finite elements in time and space together with a wavelet-in-time multigrid-in-space preconditioner, we were able to solve the arising systems to fixed accuracy in a uniformly bounded number of PCG steps.

We found that our algorithm runs in linear complexity on a single processor. Moreover, when *sufficiently many* parallel processors are available and communication is free, its runtime scales *logarithmically* in the discretization size. These complexity results translate to a highly efficient algorithm in practice.

P	$N_t$	$N_{\mathbf{x}}$	$N=N_tN_{\mathbf{x}}$	its	time (s)	time/it (s)	CPU-hrs
1–16	16385	65 0 25	1065434625		ou	t of memory -	
32	16385	65 0 25	1065434625	16	1224.85	76.55	10.89
64	16385	65 0 25	1065434625	16	615.73	38.48	10.95
128	16385	65 0 25	1065434625	16	309.81	19.36	11.02
256	16385	65 0 25	1065434625	16	163.20	10.20	11.61
512	16385	65 0 25	1065434625	16	96.54	6.03	13.73
512	16385	65 0 25	1065434625	16	96.50	6.03	13.72
1024	16385	65 0 25	1065434625	16	45.27	2.83	12.88
2048	16385	65025	1065434625	16	20.59	1.29	11.72

TABLE 8.2. Strong scaling results for the 2D problem.

	P	$N_t$	$N_{\mathbf{x}}$	$N = N_t N_{\mathbf{x}}$	its	time (s)	time/it (s)	CPU-hrs
e	1	9	261 121	2 350 089	8	33.36	4.17	0.01
por	2	17	261 121	4439057	11	46.66	4.24	0.03
ler	4	33	261 121	8616993	12	54.60	4.55	0.06
ng	8	65	261 121	16972865	13	65.52	5.04	0.15
si	16	129	261 121	33 684 609	13	86.94	6.69	0.39
s	32	257	261 121	67 108 097	14	93.56	6.68	0.83
de	64	513	261 121	133 955 073	14	94.45	6.75	1.68
no	128	1025	261 121	267 649 025	14	93.85	6.70	3.34
ole	256	2049	261 121	535 036 929	15	101.81	6.79	7.24
ltip	512	4097	261 121	1 069 812 737	15	101.71	6.78	14.47
ทท	1024	8193	261 121	2 1 39 364 353	16	108.32	6.77	30.81
Ч	2048	16385	261 121	4278467585	16	109.59	6.85	62.34

TABLE 8.3. Weak scaling results for the 2D problem.

P	$N_t$	$N_{\mathbf{x}}$	$N=N_tN_{\mathbf{x}}$	its	time (s)	time/it (s)	CPU-hrs
1–64	16385	250 047	4097020095		out	t of memory -	
128	16385	250047	4097020095	18	3 308.49	174.13	117.64
256	16385	250 047	4097020095	18	1 655.92	87.15	117.75
512	16385	250 047	4097020095	18	895.01	47.11	127.29
1024	16385	250 047	4097020095	18	451.59	23.77	128.45
2048	16385	250047	4097020095	18	221.12	12.28	125.80

TABLE 8.4. Strong scaling results for the 3D problem.

P	$N_t$	$N_{\mathbf{x}}$	$N = N_t N_{\mathbf{x}}$	its	time (s)	time/it (s)	CPU-hrs
16	129	250 047	32 256 063	15	183.65	12.24	0.82
32	257	250 047	64 262 079	16	196.26	12.27	1.74
64	513	250 047	128274111	16	197.55	12.35	3.51
128	1025	250047	256 298 175	17	210.21	12.37	7.47
256	2049	250047	512 346 303	17	209.56	12.33	14.90
512	4097	250047	1024442559	17	210.14	12.36	29.89
1024	8 1 9 3	250047	2048635071	18	221.77	12.32	63.08
2048	16385	250047	4097020095	18	221.12	12.28	125.80

TABLE 8.5. Weak scaling results for the 3D problem.

The numerical experiments serve as a showcase for the described spacetime method, and exhibit its excellent time-parallelism by solving a linear system with over 4 billion unknowns in just 109 seconds, using just over 2 thousand parallel processors. By incorporating spatial parallelism as well, we expect these results to scale well to much larger problems.

Although performed in the rather restrictive setting of the heat equation discretized using piecewise linear polynomials on uniform triangulations, the parallel framework already allows solving more general linear parabolic PDEs using polynomials of varying degree on locally refined (tensor-product) meshes. In this more general setting, we envision load balancing to become the main hurdle in achieving good scaling results.

# 9.1 Introduction

This chapter is about the adaptive numerical approximation of the heat equation using a simultaneous space-time boundary element method (BEM). In the last years, there has been a growing interest in space-time BEM for the heat equation [CS13, MST14, MST15, HT18, CR19, DNS19, DZO<sup>+</sup>19, Tau19, ZWOM21]. In contrast to the differential operator based variational formulation on the space-time cylinder, the variational formulation corresponding to space-time BEM is coercive [AN87, Cos90] so that the discretized version always has a unique solution regardless of the chosen trial space which is even quasi-optimal in the natural energy norm. Moreover, it is naturally applicable on unbounded domains and only requires a mesh of the lateral boundary of the space-time cylinder resulting in a dimension reduction. The potential disadvantage that discretizations lead to dense matrices due to the nonlocality of the boundary integral operators has been tackled, e.g., in [MST14, MST15, HT18] via the fast multipole method and  $\mathcal{H}$ -matrices.

Two often mentioned advantages of simultaneous space-time methods are their potential for massive parallelization as well as their potential for fully adaptive refinement to resolve singularities local in both space and time. While the first advantage has been investigated in, e.g., [DZO<sup>+</sup>19, ZWOM21], the latter requires suitable *a posteriori* computable error estimators, which have not been developed yet for the heat equation. Indeed, concerning *a posteriori* error estimation as well as adaptive refinement for BEM for time-dependent problems, we are only aware of the works [Glä12, GÖSS20] for the wave equation in two and three space dimensions, respectively.

In the present manuscript, we generalize the results [Fae00, Fae02] from Faermann for stationary PDEs to the heat equation: Let  $\Omega \subset \mathbb{R}^d$ , d = 2, 3, be a Lipschitz domain with boundary  $\Gamma := \partial \Omega$  and T > 0 a given end time point with corresponding time interval I := (0, T). We abbreviate the space-time cylinder  $Q := I \times \Omega$  with lateral boundary  $\Sigma := I \times \Gamma$  and corresponding outer normal vector  $\boldsymbol{n} \in \mathbb{R}^d$ . With the heat kernel

$$G(t, \boldsymbol{x}) := \begin{cases} \frac{1}{(4\pi t)^{d/2}} e^{-\frac{|\boldsymbol{x}|^2}{4t}} & \text{for } (t, \boldsymbol{x}) \in (0, \infty) \times \mathbb{R}^d, \\ 0 & \text{else,} \end{cases}$$

and a given function  $f: \Sigma \to \mathbb{R}$ , we consider the boundary integral equation,

(9.1) 
$$(\mathscr{V}\phi)(t,\boldsymbol{x}) := \int_{\Sigma} G(t-s,\boldsymbol{x}-\boldsymbol{y})\phi(t-s,\boldsymbol{x}-\boldsymbol{y})\,\mathrm{d}\boldsymbol{y}\,\mathrm{d}\boldsymbol{s} = f(t,\boldsymbol{x}),$$

for a.e.  $(t, \mathbf{x}) \in \Sigma$ . Here,  $\mathscr{V}$  is the single-layer operator. For given initial condition  $u_0 : \Omega \to \mathbb{R}$  and Dirichlet data  $u_D : \Sigma \to \mathbb{R}$ , such equations arise from the heat equation

(9.2) 
$$\begin{array}{rcl} \partial_t u - \Delta u &= 0 & \text{ on } Q, \\ u &= u_D & \text{ on } \Sigma, \\ u(0, \cdot) &= u_0 & \text{ on } \Omega. \end{array}$$

Let  $\mathcal{P}$  be a mesh of the space-time boundary  $\Sigma$  consisting of prismatic elements  $J \times K$  with  $J \subseteq \overline{I}$  and  $K \subseteq \Gamma$ , and let  $\Phi$  be an associated approximation of  $\phi$ . Typically,  $\Phi$  is a piecewise polynomial with respect to  $\mathcal{P}$ . As  $\mathscr{V}$  is an isomorphism from the dual space  $H^{-1/2,-1/4}(\Sigma) := H^{1/2,1/4}(\Sigma)'$  to the anisotropic Sobolev space  $H^{1/2,1/4}(\Sigma)$ , the discretization error  $\|\phi - \Phi\|_{H^{-1/2,-1/4}(\Sigma)}$  is equivalent to the norm of the residual  $\|f - \mathscr{V}\Phi\|_{H^{1/2,1/4}(\Sigma)}$ . We show that the residual norm can be localized up to weighted  $L_2$ -terms, i.e.,

$$\sum_{J \times K \in \mathcal{P}} \eta_{\mathcal{P}}(\Phi, J \times K)^2 \lesssim \|f - \mathscr{V}\Phi\|_{H^{1/2, 1/4}(\Sigma)}^2 \lesssim \sum_{J \times K \in \mathcal{P}} \eta_{\mathcal{P}}(\Phi, J \times K)^2 + \zeta_{\mathcal{P}}(\Phi, J \times K)^2,$$

where  $\eta_{\mathcal{P}}(\Phi, J \times K)^2$  measures the  $H^{1/2,1/4}$ -seminorm of the residual in a neighborhood of  $J \times K$  and  $\zeta_{\mathcal{P}}(\Phi) := (\operatorname{diam}(K)^{-1} + |J|^{-1/2}) \|f - \mathscr{V}\Phi\|_{L_2(J \times K)}^2$ . The hidden constants depend only on the regularity of the of the meshes found by fixing either the temporal or the spatial coordinate in  $\mathcal{P}$ . In particular, we do not require any assumption on the relation between the spatial and temporal size of the mesh elements, making anisotropically refined meshes possible.

If the elements satisfy the scaling  $|J| \approx \operatorname{diam}(K)^2$  and if  $\Phi$  is the Galerkin approximation of  $\phi$  in a discrete space  $\mathcal{X}$  that contains at least all  $\mathcal{P}$ -piecewise constant functions, then we can additionally prove that

$$\zeta_{\mathcal{P}}(\Phi, J \times K) \lesssim \eta_{\mathcal{P}}(\Phi, J \times K).$$

Indeed, numerical experiments (with d = 2) suggest that this is not the case in general: If the scaling condition is not enforced, we observe situations where the  $L_2$ -terms  $\zeta$  do not decay under mesh-refinement.

That being said, the estimator  $\eta$  does not only behave efficiently but also reliably in all considered examples. Moreover, anisotropic refinement steered by the space- and time-components of the estimator always yield the optimal algebraic convergence rate of both the estimator and the error.

# Outline

The remainder of this chapter is organized as follows: Section 9.2 summarizes the general principles of the space-time boundary element method for the heat equation. Section 9.3 recalls the localization argument of [Fae00, Fae02] and applies it to anisotropic Sobolev spaces (Theorem 9.3.3). This result is then invoked in Corollary 9.3.5 for the residual, resulting in efficient and reliable *a posteriori* computable error bounds. In particular, a Poincaré-type inequality (Lemma 9.3.4) allows to estimate the weighted  $L_2$ -terms that are still present in the upper bound from Theorem 9.3.3. Finally, Section 9.4 introduces an adaptive algorithm for d = 2 which is based on the derived error estimator. Different marking and refinement strategies are presented. The adaptive algorithm is subsequently applied to several concrete examples with typical singularities in space and time.

# 9.2 Preliminaries

# 9.2.1 General notation

Throughout and without any ambiguity,  $|\cdot|$  denotes the absolute value of scalars, the Euclidean norm of vectors in  $\mathbb{R}^n$ , or the the measure of a set in  $\mathbb{R}^n$ , e.g., the length of an interval or the area of a surface in  $\mathbb{R}^3$ . We write  $A \leq B$  to abbreviate  $A \leq CB$  with some generic constant C > 0, which is clear from the context. Moreover,  $A \approx B$  abbreviates  $A \leq B \leq A$ .

# 9.2.2 Anisotropic Sobolev spaces

For *d*-dimensional  $\omega \subseteq \Omega$  or (d-1)-dimensional  $\omega \subseteq \Gamma$ , and  $\mu \in (0, 1]$ , we first recall the Sobolev space  $H^{\mu}(\omega) := \{v \in L_2(\omega) : \|v\|_{H^{\mu}(\omega)} < \infty\}$  associated with the Sobolev–Slobodeckij norm  $\|v\|_{H^{\mu}(\omega)}^2 := \|v\|_{L_2(\omega)}^2 + |v|_{H^{\mu}(\omega)}^2$ , with

$$|v|_{H^{\mu}(\omega)}^{2} \coloneqq \begin{cases} \int_{\omega} \int_{\omega} \frac{|v(\boldsymbol{x}) - v(\boldsymbol{y})|^{2}}{|\boldsymbol{x} - \boldsymbol{y}|^{\dim(\omega) + 2\mu}} \, \mathrm{d}\boldsymbol{y} \, \mathrm{d}\boldsymbol{x} & \text{if } \mu \in (0, 1), \\ \|\nabla_{\omega} v\|_{L_{2}(\omega)}^{2} & \text{if } \mu = 1, \end{cases}$$

where dim( $\omega$ ) denotes the dimension of  $\omega$ , i.e., d or d - 1, and  $\nabla_{\omega}$  denotes the (weak) gradient on  $\omega$ , i.e., the standard gradient or the surface gradient.

Moreover, we define for any subinterval  $J \subseteq \overline{I}$ ,  $\nu \in (0, 1]$ , and any Banach space *X*,

$$H^{\nu}(J;X) := \left\{ v \in L_2(J;X) : \|v\|_{H^{\nu}(J;X)} < \infty \right\}$$

associated with the norm  $\|v\|_{H^{\nu}(J;X)}^2 := \|v\|_{L_2(J;X)}^2 + |v|_{H^{\nu}(J;X)'}^2$  with

$$|v|_{H^{\nu}(J;X)}^{2} := \begin{cases} \int_{J} \int_{J} \frac{\|v(t) - v(s)\|_{X}^{2}}{|t - s|^{1 + 2\nu}} \, \mathrm{d}s \, \mathrm{d}t & \text{ if } \nu \in (0,1), \\ \|\partial_{t}v\|_{L_{2}(\omega)}^{2} & \text{ if } \nu = 1, \end{cases}$$

where  $\partial_t$  denotes the (weak) time derivative. If  $X = \mathbb{R}$ , we simply write  $H^{\nu}(J)$ ,  $||v||_{H^{\nu}(J)}$ , and  $|v|_{H^{\nu}(J)}$ . Finally, we recall the anisotropic Sobolev space

$$H^{\mu,\nu}(J\times\omega) := L_2(J; H^\mu(\omega)) \cap H^\nu(J; L_2(\omega))$$

with corresponding norm

$$\|v\|_{H^{\mu,\nu}(J\times\omega)}^2 := \|v\|_{L_2(J;H^{\mu}(\omega))}^2 + \|v\|_{H^{\nu}(J;L_2(\omega))}^2 \quad (v \in H^{\mu,\nu}(J\times\omega)).$$

We will sometimes use the abbreviation

$$|v|^2_{L_2(J;H^{\mu}(\omega))} := \int_J |v(t,\cdot)|^2_{H^{\mu}(\omega)} dt \quad (v \in L_2(J;H^{\mu}(\omega))).$$

For  $\omega \in {\Omega, \Gamma}$ , we denote  $H^{-\mu, -\nu}(I \times \omega)$  for the dual of  $H^{\mu, \nu}(I \times \omega)$  with duality pairing  $\langle \cdot, \cdot \rangle_{I \times \omega}$ . We view  $L_2(I \times \omega)$  as subspace of  $H^{-\mu, -\nu}(I \times \omega)$  via

$$\langle v, \psi \rangle_{I \times \omega} := \int_{I} \int_{\omega} v(t, \boldsymbol{x}) \psi(t, \boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} \, \mathrm{d}t \quad \big( v \in H^{\mu, \nu}(I \times \omega), \psi \in L_2(I \times \omega) \big).$$

### 9.2.3 Boundary integral equations

It is well-known that for  $u_0 \in L^2(\Omega)$  and  $u_D \in H^{1/2,1/4}(\Sigma)$ , the heat equation (9.2) admits a unique solution  $u \in H^{1,1/2}(Q)$ . With the normal derivative  $\phi_N := \partial_n u \in H^{-1/2,-1/4}(\Sigma)$ , u satisfies the representation formula

(9.3) 
$$u = \widetilde{\mathcal{M}}_0 u_0 + \widetilde{\mathcal{V}} \phi_N - \widetilde{\mathcal{K}} u_D,$$

where

(9.4) 
$$(\widetilde{\mathscr{M}}_0 u_0)(t, \boldsymbol{x}) := \int_{\Omega} G(t, \boldsymbol{x} - \boldsymbol{y}) u_0(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y} \quad ((t, \boldsymbol{x}) \in Q)$$

denotes the initial potential,

(9.5) 
$$(\widetilde{\mathscr{V}}\phi_N)(t,\boldsymbol{x}) := \int_{\Sigma} G(t-s,\boldsymbol{x}-\boldsymbol{y})\phi_N(\boldsymbol{y}) \,\mathrm{d}\boldsymbol{y} \,\mathrm{d}s \quad ((t,\boldsymbol{x})\in Q)$$

denotes the single-layer potential, and

(9.6) 
$$(\widetilde{\mathscr{H}}u_D)(t, \boldsymbol{x}) := \int_{\Sigma} \partial_{\boldsymbol{n}(\boldsymbol{y})} G(t-s, \boldsymbol{x}-\boldsymbol{y}) u_D(\boldsymbol{y}) \, \mathrm{d}\boldsymbol{y} \, \mathrm{d}s \quad ((t, \boldsymbol{x}) \in Q)$$

denotes the double-layer potential. These linear operators satisfy the mapping properties  $\widetilde{\mathscr{M}_0}: L^2(\Omega) \to H^{1,1/2}(Q), \widetilde{\mathscr{V}_0}: H^{-1/2,-1/4}(\Sigma) \to H^{1,1/2}(Q)$ , and  $\widetilde{\mathscr{K}_0}: H^{1/2,1/4}(\Sigma) \to H^{1,1/2}(Q)$ . The lateral trace  $(\cdot)|_{\Sigma}$  of these potentials is given by

$$(\widetilde{\mathscr{M}}_0 u_0)|_{\Sigma} = \mathscr{M}_0 u_0, \quad (\widetilde{\mathscr{V}}\phi_N)|_{\Sigma} = \mathscr{V}\phi_N, \quad (\widetilde{\mathscr{K}}u_D)|_{\Sigma} = (\mathscr{K} - 1/2)u_D,$$

where the initial operator  $\mathcal{M}_0$ , the single-layer operator  $\mathcal{V}$ , and the doublelayer operator  $\mathcal{K}$  are defined as in (9.4)–(9.6) for  $(t, x) \in \Sigma$ . Applying the lateral trace to (9.3) thus results in

(9.7) 
$$\mathscr{V}\phi_N = (\mathscr{K} + 1/2)u_D - \mathscr{M}_0 u_0,$$

i.e., (9.1) with  $f := (\mathscr{K} + 1/2)u_D - \mathscr{M}_0 u_0$ . As the single-layer operator  $\mathscr{V}$  is also coercive, i.e.,

(9.8) 
$$\langle \mathscr{V}\psi,\psi\rangle_{\Sigma} \ge c_{\text{coe}} \|\psi\|_{H^{-1/2,-1/4}(\Sigma)}^2 \quad \left(\psi\in H^{-1/2,-1/4}(\Sigma)\right)$$

with some constant  $c_{\text{coe}} > 0$ , (9.7) is uniquely solvable and the solution  $\phi_N$  is just the missing normal derivative  $\partial_n u$  to compute u via the representation formula (9.3).

Alternatively, one can make the ansatz  $u = \widetilde{\mathcal{M}}_0 u_0 + \widetilde{\mathcal{V}} \phi$ . Indeed, both  $\widetilde{\mathcal{M}}_0 u_0$ and  $\widetilde{\mathcal{V}} \phi$  satisfy the heat equation, where  $\widetilde{\mathcal{M}}_0 u_0$  restricted to  $\{0\} \times \Omega$  coincides with  $u_0$  and  $\widetilde{\mathcal{V}} \phi$  vanishes there. To satisfy the Dirichlet boundary conditions, one has to solve

(9.9) 
$$\mathscr{V}\phi = u_D - \mathscr{M}_0 u_0,$$

i.e., (9.1) with  $f := u_D - \mathcal{M}_0 u_0$ . While (9.7) is called direct method as it directly provides the physically relevant quantity  $\phi_N = \partial_n u$ , (9.9) is called indirect method.

For more details and proofs, we refer to the seminal works [AN87, Noo88, Cos90], which considered  $u_0 = 0$ , and to [DNS19, Doh19] for the general case.

#### 9.2.4 Boundary meshes

Throughout this work, we consider prismatic meshes  $\mathcal{P}$  of  $\Sigma$ :

- $\mathcal{P}$  is a finite set of prisms of the form  $P = J \times K$ , where  $J \subseteq \overline{I} = [0, T]$  is some non-empty compact interval and  $K \subseteq \Gamma$  is the image of some compact Lipschitz domain<sup>1</sup>  $\hat{K} \subset \mathbb{R}^{d-1}$  under some bi-Lipschitz mapping;
- for all  $P, \tilde{P} \in \mathcal{P}$  with  $P \neq \tilde{P}$ , the intersection has measure zero on  $\Sigma$ ;
- $\mathcal{P}$  is a partition of  $\Sigma$ , i.e.,  $\Sigma = \bigcup_{P \in \mathcal{P}} P$ .

For arbitrary  $t \in \overline{I}$  and  $x \in \Gamma$ , we abbreviate the induced sets

$$\mathcal{P}|_t := \left\{ K \subseteq \Gamma \ : \ (\{t\} \times \Gamma) \cap (J \times K) \neq \emptyset \text{ for some } J \times K \in \mathcal{P} \right\}$$

<sup>&</sup>lt;sup>1</sup>A compact Lipschitz domain is the closure of a bounded Lipschitz domain. For d = 2, it is a compact interval with non-empty interior.

and

$$\mathcal{P}|_{\boldsymbol{x}} := \{ J \subseteq \overline{I} : (\overline{I} \times \{\boldsymbol{x}\}) \cap (J \times K) \neq \emptyset \text{ for some } J \times K \in \mathcal{P} \}.$$

For almost all  $t \in \overline{I}$ ,  $\mathcal{P}|_t$  is a mesh of  $\Gamma$ , i.e., a partition of  $\Gamma$  into finitely many compact Lipschitz domains such that the intersection of two distinct elements has measure zero on  $\Gamma$ . Similarly, for almost all  $x \in \Gamma$ ,  $\mathcal{P}|_x$  is a mesh of  $\overline{I}$ , i.e., a partition of  $\overline{I}$  into finitely many non-empty compact intervals such that the intersection of two different intervals is at most a point. Note that for one fixed prismatic mesh  $\mathcal{P}$  there exist constants  $C_{\text{nei}} \geq 1$ ,  $C_{\text{dist}} \geq 1$ ,  $C_{\text{shape}} \geq 1$ , and  $C_{\text{lqu}} \geq 1$  such that:

• for almost all  $t \in \overline{I}$ , the number of neighbors of an element in  $\mathcal{P}|_t$  is bounded, i.e.,

(9.10) 
$$\#\{\tilde{K}\in\mathcal{P}|_t:K\cap\tilde{K}\neq\emptyset\}\leq C_{\text{nei}}\quad\text{for all }K\in\mathcal{P}|_t.$$

• for almost all  $t \in \overline{I}$ , the elements of  $\mathcal{P}|_t$  are uniformly away from nonneighboring elements, i.e.,

(9.11) diam(K)  $\leq C_{\text{dist}} \operatorname{dist}(K, \tilde{K})$  for all  $K, \tilde{K} \in \mathcal{P}|_t$  with  $K \cap \tilde{K} = \emptyset$ ;

• for almost all  $t \in \overline{I}$ , the elements of  $\mathcal{P}|_t$  are shape-regular, i.e.,

$$(9.12) C_{\text{shape}}^{-1}|K|^{d-1} \le \operatorname{diam}(K)^{d-1} \le C_{\text{shape}}|K| for all K \in \mathcal{P}|_t;$$

• for almost all  $x \in \Gamma$ ,  $\mathcal{P}|_x$  is locally quasi-uniform, i.e.,

(9.13) 
$$|J| \leq C_{lqu} |\tilde{J}|$$
 for all  $J, \tilde{J} \in \mathcal{P}|_{\boldsymbol{x}}$  with  $J \cap \tilde{J} \neq \emptyset$ .

In the remainder of this work, we will always indicate the dependence of estimates on these particular constants.

*Remark* 9.2.1. If, for d = 2, the meshes  $\mathcal{P}|_t$  are found by iteratively bisecting some initial mesh and the level difference of neighboring elements is bounded by 1, then the constants from (9.10)–(9.12) depend only on the initial mesh; cf. [AFF+13]. For d = 3, the same holds true if the initial mesh is for instance a conforming (curvilinear) triangulation of  $\Gamma$  and one iteratively applies newest vertex bisection. The arguments for (9.10)–(9.11) are found in [AFF+17, Section 2.3 and 4.1].

# 9.2.5 Boundary element method

Given a prismatic boundary mesh  $\mathcal{P}$  and an associated finite-dimensional trial space  $\mathcal{X} \subset H^{-1/2,-1/4}(\Sigma)$ , e.g., the space of all  $\mathcal{P}$ -piecewise polynomials of

some fixed degree in space and time, let  $\Phi \in \mathcal{X}$  denote the Galerkin discretization of the solution  $\phi$  of the boundary integral equation (9.1), i.e.,

$$\langle \mathscr{V}\Phi, \Psi \rangle_{\Sigma} = \langle f, \Psi \rangle_{\Sigma} \quad (\Psi \in \mathcal{X}),$$

which is equivalent to the Galerkin orthogonality

(9.14) 
$$\langle \mathscr{V}(\phi - \Phi), \Psi \rangle_{\Sigma} = 0 \quad (\Psi \in \mathcal{X}).$$

Note that coercivity (9.8) guarantees unique solvability of the latter equations, and the Céa lemma applies

$$\|\phi - \Phi\|_{H^{-1/2, -1/4}(\Sigma)} \le \frac{C_{\text{cont}}}{c_{\text{coe}}} \min_{\Psi \in \mathcal{X}} \|\phi - \Psi\|_{H^{-1/2, -1/4}(\Sigma)},$$

where  $C_{\text{cont}}$  is the operator norm of  $\mathscr{V}: H^{-1/2,-1/4}(\Sigma) \to H^{1/2,1/4}(\Sigma)$ .

Suppose  $\mathcal{P} = \{\overline{J} \times K : J \in \mathcal{P}_{\overline{I}}, K \in \mathcal{P}_{\Gamma}\}$  is a full tensor-mesh corresponding to a mesh  $\mathcal{P}_{\Gamma}$  of  $\Gamma$  with uniform mesh-size  $h_{\boldsymbol{x}} = \operatorname{diam}(K)$  for all  $K \in \mathcal{P}_{\Gamma}$ and a mesh  $\mathcal{P}_{\overline{I}}$  of  $\overline{I}$  with uniform step-size  $h_t = h_{\boldsymbol{x}}^{\sigma}$  for some  $\sigma > 0$ . Using  $\mathcal{P}$ -piecewise polynomials of some degree  $p_{\boldsymbol{x}} \in \mathbb{N}_0$  in space- and some degree  $p_t \in \mathbb{N}_0$  in time-direction as trial space  $\mathcal{X}$ , then gives the error decay rate

(9.15) 
$$\min_{\Psi \in \mathcal{X}} \|\phi - \Psi\|_{H^{-1/2, -1/4}(\Sigma)} \lesssim N^{-\frac{\min\{p_{x}+3/2, (p_{t}+5/4)\sigma\}}{d-1+\sigma}} \quad \text{for all smooth } \phi;$$

see [CR19, Theorem 3.3]. Here,  $N \approx h_{\boldsymbol{x}}^{-(d-1)} h_t^{-1} = h_{\boldsymbol{x}}^{d-1+\sigma}$  denotes the number of degrees of freedom. The optimal grading parameter is thus given by  $\sigma = (p_{\boldsymbol{x}} + \frac{3}{2})/(p_t + \frac{5}{4})$  with resulting rate  $\mathcal{O}(N^{-\frac{p_{\boldsymbol{x}}+3/2}{d-1+\sigma}})$ .

### 9.3 A posteriori error estimation

As  $\mathscr{V}$  is an isomorphism, it holds that

(9.16) 
$$\|\phi - \Phi\|_{H^{-1/2, -1/4}(\Sigma)} = \|f - \mathscr{V}\Phi\|_{H^{1/2, 1/4}(\Sigma)}.$$

Here,  $\Phi \in H^{-1/2,-1/4}(\Sigma)$  can be an arbitrary approximation of the solution  $\phi$  of (9.1). While the right-hand side is in principle *a posteriori* computable, the computation of the Sobolev–Slobodeckij norm over the full space-time boundary  $\Sigma$  is expensive, and it does not provide any information on where to locally refine the given mesh to increase the accuracy of the approximation. According to (9.16), it is sufficient to derive suitable estimate for the residual  $f - \mathscr{V}\Phi$  in the  $H^{1/2,1/4}(\Sigma)$ -norm. Recall that this term is  $L_2(\Sigma)$ -orthogonal to all functions  $\Psi \in \mathcal{X}$  provided that  $\Phi$  is the Galerkin approximation of  $\phi$  in  $\mathcal{X}$ ; see (9.14).
#### 9.3.1 Localization of the anisotropic Sobolev–Slobodeckij norm

The following proposition provides the key argument for our *a posteriori* error estimation. While the first inequality is trivial, the original version of the second one already goes back to [Fae00, Fae02]. We make use of the slightly generalized version from [GP20, Lemma 4.5]; see [Gan17, Lemma 5.3.2] for a detailed proof.

**Proposition 9.3.1.** Let  $\mu \in (0,1)$  and  $\mathcal{P}_{\Gamma}$  be a mesh of  $\Gamma$ . Then, there exist constants  $C_1, C_2 > 0$  such that for all  $v \in H^{\mu}(\Gamma)$ , there holds that

$$\begin{split} C_1^{-1} \sum_{K \in \mathcal{P}_{\Gamma}} \sum_{\substack{\tilde{K} \in \mathcal{P}_{\Gamma} \\ K \cap \tilde{K} \neq \emptyset}} |v|_{H^{\mu}(K \cup \tilde{K})}^2 &\leq \|v\|_{H^{\mu}(\Gamma)}^2 \leq \sum_{K \in \mathcal{P}_{\Gamma}} \sum_{\substack{\tilde{K} \in \mathcal{P}_{\Gamma} \\ K \cap \tilde{K} \neq \emptyset}} |v|_{H^{\mu}(K \cup \tilde{K})}^2 \\ &+ C_2 \sum_{K \in \mathcal{P}_{\Gamma}} \operatorname{diam}(K)^{-2\mu} \|v\|_{L_2(K)}^2. \end{split}$$

The constant  $C_1$  is given as  $C_1 = 2(C_{\text{nei}} + 1)^2$  with  $C_{\text{nei}}$  from (9.10) (with  $\mathcal{P}|_t$  replaced by  $\mathcal{P}_{\Gamma}$ ), and  $C_2$  depends only on the dimension  $d, \mu, \Gamma$ , and the constant  $C_{\text{dist}}$  from (9.11) (with  $\mathcal{P}|_t$  replaced by  $\mathcal{P}_{\Gamma}$ ).

Note that local quasi-uniformity (9.13) (with  $\mathcal{P}|_{\boldsymbol{x}}$  replaced by  $\mathcal{P}_{\overline{I}}$ ) of a time mesh  $\mathcal{P}_{\overline{I}}$  is actually equivalent to

diam
$$(J) = |J| \leq C_{lqu} \operatorname{dist}(J, \tilde{J})$$
 for all  $J, \tilde{J} \in \mathcal{P}_{\overline{I}}$  with  $J \cap \tilde{J} = \emptyset$ .

Moreover, for any element  $J \in \mathcal{P}_{\overline{I}}$ , there are at most three  $\tilde{J} \in \mathcal{P}_{\overline{I}}$  with  $J \cap \tilde{J} \neq \emptyset$ . In particular, the same reference as before applies and we also obtain the following proposition.

**Proposition 9.3.2.** Let  $\nu \in (0,1)$  and  $\mathcal{P}_{\overline{I}}$  be a mesh of  $\overline{I}$ . Then, there exist constants  $C_1, C_2 > 0$  such that for all  $v \in H^{\nu}(I)$ , there holds that

$$\begin{split} C_1^{-1} \sum_{J \in \mathcal{P}_{\overline{T}}} \sum_{\substack{\tilde{J} \in \mathcal{P}_{\overline{T}} \\ J \cap J \neq \emptyset}} |v|_{H^{\nu}(J \cup \tilde{J})}^2 \leq \|v\|_{H^{\nu}(I)}^2 \leq \sum_{J \in \mathcal{P}_{\overline{T}}} \sum_{\substack{\tilde{J} \in \mathcal{P}_{\overline{T}} \\ J \cap J \neq \emptyset}} |v|_{H^{\nu}(J \cup \tilde{J})}^2 \\ &+ C_2 \sum_{J \in \mathcal{P}_{\overline{T}}} |J|^{-2\nu} \|v\|_{L_2(J)}^2. \end{split}$$

*The constant*  $C_1$  *is given as*  $C_1 = 32$ *, and*  $C_2$  *depends only on*  $\nu$ *,* |I|*, and the constant*  $C_{lqu}$  *from* (9.13) *(with*  $\mathcal{P}|_{\boldsymbol{x}}$  *replaced by*  $\mathcal{P}_{\overline{I}}$ *).* 

The latter two propositions allow to derive the following *a posteriori error estimation*, which can be employed for arbitrary approximations  $\Phi$ .

**Theorem 9.3.3.** Let  $\mu, \nu \in (0, 1)$  and  $\mathcal{P}$  be a prismatic mesh of  $\Sigma$ . Then, there exist constants  $C'_{\text{eff}}, C''_{\text{rel}} > 0$  such that for all  $v \in H^{\mu,\nu}(\Sigma)$ , there holds that

$$\sum_{J\times K\in\mathcal{P}} \left(\sum_{\substack{\tilde{J}\times\tilde{K}\in\mathcal{P}\\|J\cap\tilde{J}|>0\\K\cap\tilde{K}\neq\emptyset}} |v|^2_{L_2(J\cap\tilde{J};H^{\mu}(K\cup\tilde{K}))} + \sum_{\substack{\tilde{J}\times\tilde{K}\in\mathcal{P}\\|\tilde{J}\cap\tilde{J}\neq\emptyset\\|K\cap\tilde{K}|>0}} |v|^2_{H^{\nu}(J\cup\tilde{J};L_2(K\cap\tilde{K}))}\right) \leq (C'_{\text{eff}})^2 \|v\|^2_{H^{\mu,\nu}(\Sigma)}$$

as well as

$$(C_{\rm rel}')^{-2} \|v\|_{H^{\mu,\nu}(\Sigma)}^2 \leq \sum_{J \times K \in \mathcal{P}} \left( \sum_{\substack{J \times \bar{K} \in \mathcal{P} \\ |J \cap \bar{J}| > 0 \\ K \cap \bar{K} \neq \emptyset}} |v|_{L_2(J \cap \tilde{J}; H^{\mu}(K \cup \tilde{K}))}^2 + \sum_{\substack{J \times \bar{K} \in \mathcal{P} \\ |J \cap \bar{J}| > 0 \\ K \cap \bar{K} \neq \emptyset}} |v|_{K \cap \bar{K}| > 0}^2 \right)$$

$$(9.17) \qquad + \sum_{J \times K \in \mathcal{P}} \left( \operatorname{diam}(K)^{-2\mu} + |J|^{-2\nu} \right) \|v\|_{L_2(J \times K)}^2.$$

The constant  $C'_{\text{eff}}$  is given as  $C'_{\text{eff}} = \max(2(C_{\text{nei}} + 1)^2, 32)$  with  $C_{\text{nei}}$  from (9.10), and  $C'_{\text{rel}}$  depends only on d,  $\mu$ ,  $\nu$ ,  $\Gamma$ , |I| and the constants  $C_{\text{dist}}$  from (9.11) as well as  $C_{\text{lgu}}$  from (9.13).

*Proof.* We split the proof into four steps.

**Step 1:** In this step, we bound  $||v||_{L_2(I;H^{\mu}(\Gamma))}$  from below. Proposition 9.3.1 gives that

$$\|v\|_{L_2(I;H^{\mu}(\Gamma))}^2 = \int_I \|v(t,\cdot)\|_{H^{\mu}(\Gamma)}^2 \,\mathrm{d}t \gtrsim \int_I \sum_{K\in\mathcal{P}|_t} \sum_{\substack{\tilde{K}\in\mathcal{P}|_t\\K\cap\tilde{K}\neq\emptyset}} |v(t,\cdot)|_{H^{\mu}(K\cup\tilde{K})}^2 \,\mathrm{d}t.$$

Note that  $K \in \mathcal{P}|_t$  is equivalent to  $J \times K \in \mathcal{P}$  for some J with  $t \in J$ . With the indicator function  $\mathbb{1}_S$  of a set S, the last term thus is equal to

$$\begin{split} \int_{I} \sum_{K \in \mathcal{P}|_{t}} \sum_{\tilde{K} \in \mathcal{P}|_{t}} |v(t, \cdot)|^{2}_{H^{\mu}(K \cup \tilde{K})} \, \mathrm{d}t &= \int_{I} \sum_{J \times K \in \mathcal{P}} \mathbb{1}_{J}(t) \sum_{\tilde{J} \times \tilde{K} \in \mathcal{P} \atop K \cap \tilde{K} \neq \emptyset} \mathbb{1}_{\tilde{J}}(t) |v(t, \cdot)|^{2}_{H^{\mu}(K \cup \tilde{K})} \, \mathrm{d}t \\ &= \sum_{J \times K \in \mathcal{P}} \sum_{J \times \tilde{K} \in \mathcal{P} \atop |J \cap \tilde{J}| > 0 \atop K \cap \tilde{K} \neq \emptyset} |v|^{2}_{L_{2}(J \cap \tilde{J}; H^{\mu}(K \cup \tilde{K}))}. \end{split}$$

**Step 2:** In this step, we bound  $||v||_{L_2(I;H^{\mu}(\Gamma))}$  from above. Proposition 9.3.1 gives that

$$\begin{split} \|v\|_{L_{2}(I;H^{\mu}(\Gamma))}^{2} &= \int_{I} \|v(t,\cdot)\|_{H^{\mu}(\Gamma)}^{2} \,\mathrm{d}t \\ \lesssim \int_{I} \sum_{K \in \mathcal{P}|_{t}} \sum_{\substack{\tilde{K} \in \mathcal{P}|_{t} \\ K \cap \tilde{K} \neq \emptyset}} |v(t,\cdot)|_{H^{\mu}(K \cup \tilde{K})}^{2} + \sum_{K \in \mathcal{P}|_{t}} \operatorname{diam}(K)^{-2\mu} \|v(t,\cdot)\|_{L_{2}(K)}^{2} \,\mathrm{d}t. \end{split}$$

The first term in this expression has already been treated in Step 1. As  $K \in \mathcal{P}|_t$  is equivalent to  $J \times K \in \mathcal{P}$  for some J with  $t \in J$ , the second term reads

$$\begin{split} \int_{I} \sum_{K \in \mathcal{P}|_{t}} \operatorname{diam}(K)^{-2\mu} \| v(t, \cdot) \|_{L_{2}(K)}^{2} \, \mathrm{d}t &= \int_{I} \sum_{J \times K \in \mathcal{P}} \mathbb{1}_{J}(t) \, \operatorname{diam}(K)^{-2\mu} \| v(t, \cdot) \|_{L_{2}(K)}^{2} \, \mathrm{d}t \\ &= \sum_{J \times K \in \mathcal{P}} \operatorname{diam}(K)^{-2\mu} \| v \|_{L_{2}(J \times K)}^{2}. \end{split}$$

**Step 3:** In this step, we bound  $||v||_{H^{\nu}(I;L_2(\Gamma))}$  from below. The Fubini theorem, Proposition 9.3.2, and the same argument as in Step 1 give that

$$\begin{split} \|v\|_{H^{\nu}(I;L_{2}(\Gamma))}^{2} \gtrsim & \int_{\Gamma} \sum_{J \in \mathcal{P}|_{\boldsymbol{x}}} \sum_{\substack{\tilde{J} \in \mathcal{P}|_{\boldsymbol{x}} \\ J \cap \tilde{J} \neq \emptyset}} |v(\cdot,\boldsymbol{x})|_{H^{\nu}(J \cup \tilde{J})}^{2} \, \mathrm{d}\boldsymbol{x} \\ &= \int_{\Gamma} \sum_{J \times K \in \mathcal{P}} \mathbbm{1}_{K}(\boldsymbol{x}) \sum_{\substack{\tilde{J} \times \tilde{K} \in \mathcal{P} \\ J \cap \tilde{J} \neq \emptyset}} \mathbbm{1}_{\tilde{K}}(\boldsymbol{x}) |v(\cdot,\boldsymbol{x})|_{H^{\nu}(J \cup \tilde{J})}^{2} \, \mathrm{d}\boldsymbol{x} \\ &= \sum_{\substack{J \times K \in \mathcal{P} \\ J \cap \tilde{J} \neq \emptyset \\ |K \cap \tilde{K}| > 0}} \sum_{\|v\|_{H^{\nu}(J \cup \tilde{J}; L_{2}(K \cap \tilde{K}))}^{2} \cdot \mathbbm{1}_{\tilde{K}}(\tilde{K})} \|v\|_{H^{\nu}(J \cup \tilde{J}; L_{2}(K \cap \tilde{K}))}^{2} \cdot \mathbbm{1}_{\tilde{K}}(\tilde{K})} \|v\|_{H^{\nu}(J \cup \tilde{J}; L_{2}(K \cap \tilde{K}))}^{2} \cdot \mathbbm{1}_{\tilde{K}}(\tilde{K}) \|v\|_{H^{\nu}(J \cup \tilde{J}; L_{2}(K \cap \tilde{K}))}^{2} \cdot \mathbbm{1}_{\tilde{K}}(\tilde{K})} \|v\|_{H^{\nu}(J \cup \tilde{J}; L_{2}(K \cap \tilde{K}))}^{2} \cdot \mathbbm{1}_{\tilde{K}}(\tilde{K}) \|v\|_{H^{\nu}(J \cup \tilde{J}; L_{2}(K \cap \tilde{K})}^{2} \cdot \mathbbm{1}_{\tilde{K}}(\tilde{K}) \|v\|_{H^{\nu}(J \cup \tilde{J}; L_{2}(K \cap \tilde{K}))}^{2} \cdot \mathbbm{1}_{\tilde{K}}(\tilde{K}) \|v\|_{H^{\nu}(J \cup \tilde{J}; L_{2}(K \cap \tilde{K})}^{2} \cdot \mathbbm{1}_{\tilde{K}}(\tilde{K}) \|v\|_{H^{\nu}(J \cup \tilde{K})}^$$

**Step 4:** In this step, we bound  $||v||_{H^{\nu}(I;L_2(\Gamma))}$  from above. The Fubini theorem and Proposition 9.3.2 give that

$$\begin{split} \|v\|_{H^{\nu}(I;L_{2}(\Gamma))}^{2} &= \int_{\Gamma} \|v(\cdot,\boldsymbol{x})\|_{H^{\nu}(I)}^{2} \,\mathrm{d}\boldsymbol{x} \\ &\lesssim \int_{\Gamma} \sum_{J \in \mathcal{P}|_{\boldsymbol{x}}} \sum_{\substack{\tilde{J} \in \mathcal{P}|_{\boldsymbol{x}} \\ J \cap \tilde{J} \neq \emptyset}} |v(\cdot,\boldsymbol{x})|_{H^{\nu}(J \cup \tilde{J})}^{2} + \sum_{J \in \mathcal{P}|_{\boldsymbol{x}}} |J|^{-2\nu} \|v(\cdot,\boldsymbol{x})\|_{L_{2}(J)}^{2} \,\mathrm{d}\boldsymbol{x}. \end{split}$$

The first term has already been treated in Step 3. The second term reads

$$\int_{\Gamma} \sum_{J \in \mathcal{P}|_{x}} |J|^{-2\nu} \|v(\cdot, x)\|_{L_{2}(J)}^{2} dx = \int_{\Gamma} \sum_{J \times K \in \mathcal{P}} \mathbb{1}_{K}(x) |J|^{-2\nu} \|v(\cdot, x)\|_{L_{2}(J)}^{2} dx$$
$$= \sum_{J \times K \in \mathcal{P}} |J|^{-2\nu} \|v\|_{L_{2}(J \times K)}^{2}.$$

This concludes the proof.

#### 9.3.2 Poincaré-type inequality

Assuming the grading  $|J| \approx \operatorname{diam}(K)^{\mu/\nu}$  as well as  $L_2(\Sigma)$ -orthogonality of v to piecewise constants, the following local Poincaré-type inequality allows to get rid of the weighted  $L_2$ -terms in (9.17). The proof works essentially as in [Cos90, Proposition 5.3], where a global version on uniform meshes is considered.

**Lemma 9.3.4.** Let  $\mu, \nu \in (0, 1)$  and  $\mathcal{P}$  be a prismatic mesh of  $\Sigma$ . Then, there holds for all  $v \in H^{\mu,\nu}(\Sigma)$  and all  $J \times K \in \mathcal{P}$  with  $\langle v, 1 \rangle_{L_2(J \times K)} = 0$  that

$$\|v\|_{L_2(J\times K)}^2 \le C_{\text{shape}} \left( \operatorname{diam}(K)^{2\mu} |v|_{L_2(J;H^{\mu}(K))}^2 + |J|^{2\nu} |v|_{H^{\nu}(J;L_2(K))}^2 \right)$$

*Here*,  $C_{\text{shape}} \ge 1$  *is the constant from* (9.12).

*Proof.* Let  $\Pi_J$ ,  $\Pi_K$ , and  $\Pi_{J \times K}$  denote the  $L_2$ -orthogonal projection onto the space of constants on J, K, and  $J \times K$ , respectively. Note that  $\Pi_{J \times K} = \Pi_J \otimes \Pi_K$  and thus

$$\begin{aligned} \|v\|_{L_2(J\times K)} &= \|(1-\Pi_{J\times K})v\|_{L_2(J\times K)} \\ &\leq \|(1-\Pi_J\otimes\operatorname{Id})v\|_{L_2(J\times K)} + \|(\Pi_J\otimes\operatorname{Id}-\Pi_J\otimes\Pi_K)v\|_{L_2(J\times K)}. \end{aligned}$$

As  $\Pi_J$  has operator norm 1, a standard Poincaré-type inequality, see, e.g., [Fae02, Lemma 3.4] for the elementary proof, shows for the second term that

$$\begin{split} \|(\Pi_J \otimes \mathrm{Id} - \Pi_J \otimes \Pi_K) v\|_{L_2(J \times K)}^2 &\leq \|(1 - \mathrm{Id} \otimes \Pi_K) v\|_{L_2(J \times K)}^2 \\ &= \int_J \|(1 - \Pi_K) v(t, \cdot)\|_{L_2(K)}^2 \,\mathrm{d}t \\ &\leq \frac{\mathrm{diam}(K)^{d-1+2\mu}}{2|K|} \int_J |v(t, \cdot)|_{H^{\mu}(K)}^2 \,\mathrm{d}t \\ &\leq \frac{C_{\mathrm{shape}}}{2} \,\mathrm{diam}(K)^{2\mu} |v|_{L_2(J;H^{\mu}(K)}^2. \end{split}$$

The first term can be estimated similarly

$$\|(1 - \Pi_J \otimes \mathrm{Id})v\|_{L_2(J \times K)}^2 \le \frac{1}{2} |J|^{2\nu} |v|_{H^{\nu}(J;L_2(K))}^2,$$

which concludes the proof.

#### 9.3.3 A posteriori error estimators

For arbitrary prismatic meshes  $\mathcal{P}$  of  $\Sigma$  with some associated trial space  $\mathcal{X} \subset H^{-1/2,-1/4}(\Sigma)$  and  $\Phi \in \mathcal{X}$ , we define the following error indicators for all  $J \times K \in \mathcal{P}$ ,

$$\begin{split} \eta_{\mathcal{P}}^{\boldsymbol{x}}(\Phi, J \times K) &\coloneqq \sum_{\substack{J \times \tilde{K} \in \mathcal{P} \\ |J \cap \tilde{J}| > 0 \\ K \cap \tilde{K} \neq \emptyset}} |f - \mathscr{V}\Phi|^2_{L_2(J \cap \tilde{J}; H^{1/2}(K \cup \tilde{K}))}, \\ \eta_{\mathcal{P}}^t(\Phi, J \times K)^2 &\coloneqq \sum_{\substack{J \times \tilde{K} \in \mathcal{P} \\ J \cap \tilde{J} \neq \emptyset \\ |K \cap \tilde{K}| > 0}} |f - \mathscr{V}\Phi|^2_{H^{1/4}(J \cup \tilde{J}; L_2(K \cap \tilde{K}))}, \\ \zeta_{\mathcal{P}}^{\boldsymbol{x}}(\Phi, J \times K) &\coloneqq \operatorname{diam}(K)^{-1} \|f - \mathscr{V}\Phi\|^2_{L_2(J \times K)}, \\ \zeta_{\mathcal{P}}^t(\Phi, J \times K)^2 &\coloneqq |J|^{-1/2} \|f - \mathscr{V}\Phi\|^2_{L_2(J \times K)}. \end{split}$$

The corresponding error estimators read as

$$\begin{split} \eta_{\mathcal{P}}(\Phi, J \times K)^2 &:= \eta_{\mathcal{P}}^{\boldsymbol{x}}(\Phi, J \times K)^2 + \eta_{\mathcal{P}}^t(\Phi, J \times K)^2, \\ \zeta_{\mathcal{P}}(\Phi, J \times K)^2 &:= \zeta_{\mathcal{P}}^{\boldsymbol{x}}(\Phi, J \times K)^2 + \zeta_{\mathcal{P}}^t(\Phi, J \times K)^2, \\ \eta_{\mathcal{P}}(\Phi)^2 &:= \sum_{J \times K \in \mathcal{P}} \eta_{\mathcal{P}}(\Phi, J \times K)^2, \quad \zeta_{\mathcal{P}}(\Phi)^2 &:= \sum_{J \times K \in \mathcal{P}} \zeta_{\mathcal{P}}(\Phi, J \times K)^2. \end{split}$$

With (9.16), we overall obtain the following *a posteriori* estimates.

**Corollary 9.3.5.** Let  $\phi$  be the solution of (9.1) and  $\mathcal{P}$  be a prismatic mesh of  $\Sigma$  with some associated discrete trial space  $\mathcal{X} \subset H^{-1/2,-1/4}(\Sigma)$ . Then, there exist constants  $C_{\text{eff}}, \tilde{C}_{\text{rel}} > 0$  such that for arbitrary  $\Phi \in \mathcal{X}$ , there holds that

$$C_{\text{eff}}^{-1}\eta_{\mathcal{P}}(\Phi) \le \|\phi - \Phi\|_{H^{-1/2, -1/4}(\Sigma)} \le \tilde{C}_{\text{rel}}(\eta_{\mathcal{P}}(\Phi)^2 + \zeta_{\mathcal{P}}(\Phi)^2)^{1/2}.$$

*If the space* X *contains all* P*-piecewise constant functions and*  $\Phi \in X$  *is the Galerkin approximation of*  $\phi$ *, there further holds that* 

$$\zeta_{\mathcal{P}}(\Phi, J \times K)^{2} \leq C_{\text{shape}} \Big( \operatorname{diam}(K)^{-1} + |J|^{-1/2} \Big) \Big( \operatorname{diam}(K) + |J|^{1/2} \Big) \eta_{\mathcal{P}}(\Phi, J \times K)^{2}$$

for all  $J \times K \in \mathcal{P}$ . If  $C_{\text{grad}}^{-1} \operatorname{diam}(K) \leq |J|^{1/2} \leq C_{\text{grad}} \operatorname{diam}(K)$  is satisfied for all  $J \times K \in \mathcal{P}$  and a uniform constant  $C_{\text{grad}} \geq 1$ , this implies the existence of a constant  $C_{\text{rel}} > 0$  such that

$$\|\phi - \Phi\|_{H^{-1/2, -1/4}(\Sigma)} \le C_{\operatorname{rel}} \eta_{\mathcal{P}}(\Phi).$$

The constants  $C_{\text{eff}}$  and  $\tilde{C}_{\text{rel}}$  are given as  $C_{\text{eff}} = C'_{\text{eff}} C_{\text{cont}}$  and  $C_{\text{rel}} = C'_{\text{rel}}/c_{\text{coe}}$  with  $C'_{\text{eff}}$  and  $C'_{\text{rel}}$  from Theorem 9.3.3, the operator norm  $C_{\text{cont}}$  of  $\mathscr{V}$ , and  $c_{\text{coe}}$  from (9.8). The constant  $C_{\text{rel}}$  is given as  $C_{\text{rel}} = \tilde{C}_{\text{rel}} \sqrt{2C_{\text{shape}}(1 + C_{\text{grad}})}$ .

*Remark* 9.3.6. According to (9.15), the required scaling  $|J| \approx \operatorname{diam}(K)^2$ , i.e.,  $\sigma = 2$ , is the optimal scaling for approximating smooth solutions  $\phi$  if the polynomial degrees of  $\mathcal{X}$  satisfy  $p_x = 2p_t + 1$ .

#### 9.4 Numerical experiments

In this section, we employ the error estimator  $\eta$  within an adaptive algorithm using different refinement strategies, and investigate the resulting convergence rates. We restrict ourselves to the case d = 2, with  $\Gamma = \partial \Omega$  being the boundary of a polygonal domain  $\Omega \subset \mathbb{R}^2$ , and set the time domain to be I = (0, 1).

For  $\mathcal{P}$  a prismatic mesh of the space-time boundary, i.e., a quadrilateral mesh as d = 2, we consider the trial space  $\mathcal{X}$  of piecewise constants with respect to  $\mathcal{P}$ . In particular, this allows us to perform integration in time analytically for all integrals that are involved in the computation of the Galerkin matrix and the evaluation of the single-layer operator  $\mathscr{V}$ ; see, e.g., [Cos90]. The remaining

integrals over  $\Gamma$  have a logarithmic singularity, for which we use the quadrature rules from [Smi00]. For the computation of the Sobolev–Slobodeckij seminorm in the Faermann estimator  $\eta_{\mathcal{P}}(\Phi)$ , we use Duffy transformations and Gauss quadrature for the regularized integrands.

### 9.4.1 Adaptive algorithm

In our numerical experiments below, we employ the following adaptive algorithm with  $\theta = 0.9$ .

**Algorithm.** Let  $0 < \theta \leq 1$  be a marking parameter and  $\mathcal{P} = \{J \times K : J \in \mathcal{P}_{\overline{I}}, K \in \mathcal{P}_{\Gamma}\}$  be an initial tensor-mesh corresponding to a mesh  $\mathcal{P}_{\Gamma}$  of  $\Gamma$  and a mesh  $\mathcal{P}_{\overline{I}}$  of  $\overline{I} = [0, T]$ . For each  $\ell = 0, 1, 2, ...$ , iterate the following steps:

- (i) Compute Galerkin approximation  $\Phi_{\ell}$  of  $\phi$  in the space  $\mathcal{X}_{\ell}$  of all  $\mathcal{P}_{\ell}$ -piecewise constant functions on  $\Sigma$ .
- (ii) Compute indicators  $\eta_{\mathcal{P}_{\ell}}^{\boldsymbol{x}}(\Phi_{\ell}, J \times K)$  and  $\eta_{\mathcal{P}_{\ell}}^{t}(\Phi_{\ell}, J \times K)$  for all elements  $J \times K \in \mathcal{P}_{\ell}$ .
- (iii) Determine two minimal sets of marked elements  $\mathcal{M}^{\boldsymbol{x}}_{\ell}, \mathcal{M}^{t}_{\ell} \subseteq \mathcal{P}_{\ell}$  such that

$$(9.18) \ \theta^2 \eta_{\mathcal{P}_{\ell}}(\Phi_{\ell})^2 \leq \sum_{J \times K \in \mathcal{M}_{\ell}^{\mathfrak{a}}} \eta_{\mathcal{P}_{\ell}}^{\mathfrak{a}}(\Phi_{\ell}, J \times K)^2 + \sum_{J \times K \in \mathcal{M}_{\ell}^{t}} \eta_{\mathcal{P}_{\ell}}^{t}(\Phi_{\ell}, J \times K)^2.$$

(iv) Refine at least all marked elements to obtain a new mesh  $\mathcal{P}_{\ell+1}$ .

We will focus on *isotropic* and *anisotropic* adaptive strategies:

- In isotropic refinement, we require  $\mathcal{M}_{\ell}^{\boldsymbol{x}} = \mathcal{M}_{\ell}^{t}$  in the marking step (iii), so that (9.18) simplifies to  $\theta^{2}\eta_{\mathcal{P}_{\ell}}(\Phi_{\ell})^{2} \leq \sum_{J \times K \in \mathcal{M}_{\ell}} \eta_{\mathcal{P}_{\ell}}(\Phi_{\ell}, J \times K)^{2}$ . In the refinement step (iv), we iteratively mark additional elements to ensure that, after subdividing all marked elements into four congruent rectangles, the new mesh  $\mathcal{P}_{\ell+1}$  has only one hanging node per edge.
- In anisotropic refinement, we bisect the elements M<sup>t</sup><sub>ℓ</sub> \ M<sup>t</sup><sub>ℓ</sub> in space, bisect the elements M<sup>t</sup><sub>ℓ</sub> \ M<sup>t</sup><sub>ℓ</sub> in time, and subdivide all elements M<sup>t</sup><sub>ℓ</sub> ∩ M<sup>t</sup><sub>ℓ</sub> into four congruent rectangles. Then, we iteratively bisect additional elements in space and/or time to ensure that the level difference in space and in time between elements sharing an edge in the new mesh P<sub>ℓ+1</sub> is bounded by 1. Here, the level in space and the level in time of elements are defined as the number of bisections in space and time, respectively, to obtain the element from the initial mesh P<sub>0</sub>.

For comparison, we also include uniform refinement, where  $\mathcal{P}_{\ell+1}$  is obtained from  $\mathcal{P}_{\ell}$  by subdividing each element into four congruent rectangles. For all considered refinement strategies, it is easy to see that the mesh constants from (9.10)–(9.13) for to  $(\mathcal{P}_{\ell})_{\ell \in \mathbb{N}_0}$  depend only on the initial mesh  $\mathcal{P}_0$ .

#### 9.4.2 Reference for exact error

As the exact error  $\|\phi - \Phi\|_{H^{-1/2,-1/4}(\Sigma)}$  cannot be readily computed in the examples below, we compare the error estimator  $\eta_{\mathcal{P}}$  and the weighted  $L_2$ -terms  $\zeta_{\mathcal{P}}$  from Section 9.3.3 with the following (h - h/2)-estimator: For a mesh  $\mathcal{P}$ , define the uniformly refined mesh as  $\widehat{\mathcal{P}}$ . With the the Galerkin approximation  $\widehat{\Phi}$  from the refined trial space, we define the (h - h/2)-estimator as

$$\|\Phi - \widehat{\Phi}\|_{\mathscr{V}}^2 := \langle \mathscr{V}(\Phi - \widehat{\Phi}), \, \Phi - \widehat{\Phi} \rangle_{\Sigma}.$$

Under the saturation assumption  $\|\phi - \widehat{\Phi}\|_{\mathscr{V}} \leq q_{\text{sat}} \|\phi - \Phi\|_{\mathscr{V}}$ , the triangle inequality shows that this estimator is equivalent to  $\|\phi - \Phi\|_{\mathscr{V}}$ , and therefore to the error  $\|\phi - \Phi\|_{H^{-1/2,-1/4}(\Sigma)}$  by coercivity of  $\mathscr{V}$ . Note that the saturation assumption is indeed satisfied under the realistic (asymptotic) assumption that  $\|\phi - \Phi\|_{\mathscr{V}} = \mathcal{O}\left((\#\mathcal{P})^{-s}\right)$  for some arbitrary rate s > 0.

#### 9.4.3 Smooth problem

Let  $\Omega = (0, 1)^2$  and prescribe  $u(t, x_1, x_2) := \exp(-2\pi^2 t) \sin(\pi x_1) \sin(\pi x_2)$  with initial condition  $u_0(x_1, x_2) := \sin(\pi x_1) \sin(\pi x_2)$  and Dirichlet data  $u_D \equiv 0$ . We choose  $\mathcal{P}_0 := \{[0, 1] \times K : K \in \mathcal{P}_{\Gamma}\}$  with the uniform mesh  $\mathcal{P}_{\Gamma}$  of  $\Gamma$  being aligned with the corners and consisting of four elements, as initial mesh of the space-time boundary  $\Sigma$ .

Figure 9.1 displays the results in double-logarithmic plots so that the slopes of the lines indicate the corresponding convergence rates. With the number of degrees of freedom  $N = \#\mathcal{P}$ , we see the expected rate  $\mathcal{O}(N^{-5/8}) = \mathcal{O}(N^{-0.625})$ from (9.15) for both uniform refinement and isotropic refinement (with still slightly worse rate for the  $L_2$ -terms  $\zeta_{\mathcal{P}}(\Phi)$  for uniform refinement), albeit adaptive isotropic refinement offers quantitively better results. For anisotropic refinement refinement, the rate is improved to  $\mathcal{O}(N^{-15/22}) \approx \mathcal{O}(N^{-0.68})$ . According to (9.15), this coincides with the best possible rate that can be achieved with uniform tensor-meshes, where the optimal scaling parameter in  $h_t \approx h_x^{\sigma}$ is given by  $\sigma = 6/5$ . Note that we do not require setting an explicit scaling in our anisotropic adaptive algorithm, it recovers the optimal rate automatically.

#### 9.4.4 Mildly singular problem

Let  $\Omega = (0, 1)^2$ , with initial condition  $u_0 \equiv 0$  and Dirichlet data  $u_D(t, x_1, x_2) := t^2$ . We expect the solution here to be only singular in the four corners of the unit square as the initial condition is compatible with the Dirichlet data. The initial mesh  $\mathcal{P}_0$  is chosen as in Section 9.4.3. Figure 9.2 displays the results. The assymptotic decay rate for all estimators under uniform refinement seems to be  $\mathcal{O}(N^{-1/3})$ , which is improved to  $\mathcal{O}(N^{-1/2})$  for isotropic refinement, and finally, under anistriopic refinement this becomes the optimal rate  $\mathcal{O}(N^{-15/22})$ .



FIGURE 9.1. Error estimators for the smooth problem of Section 9.4.3 plotted double-logarithmically over the degrees of freedom N = #P: uniform refinement (left), isotropic refinement (middle), and anisotropic refinement (right).



FIGURE 9.2. Error estimators for the mildly singular problem of Section 9.4.4 plotted double-logarithmically over the degrees of freedom  $N = \#\mathcal{P}$ : uniform refinement (left), isotropic refinement (middle), and anisotropic refinement (right).

#### 9.4.5 Singular problem

Let  $\Omega = (0,1)^2$  with initial condition  $u_0 \equiv 0$  and Dirichlet data  $u_D \equiv 1$ . The solution to this problem is known to have a strong singularity for t = 0 due to the incompatibility of initial and boundary conditions, in addition to singularities in the four corners of the unit square. The initial mesh  $\mathcal{P}_0$  is chosen as in Section 9.4.3.

Figure 9.3 displays the results. The Faermann estimator  $\eta_{\mathcal{P}}(\Phi)$  and the (h - h/2)-estimator  $\|\Phi - \widehat{\Phi}\|_{\mathscr{V}}$  show both the same sensible convergence behavior for this problem. For uniform refinement, they display the rate  $\mathcal{O}(N^{-1/8})$ , which is then improved by isotropic refinement to  $\mathcal{O}(N^{-1/4})$ . Finally, for anistropic



FIGURE 9.3. Error estimators for the singular problem of Section 9.4.5 plotted double-logarithmically over the degrees of freedom N = #P: uniform refinement (left), isotropic refinement (middle), and anisotropic refinement (right).

refinement, they achieve the best possible rate  $\mathcal{O}(N^{-15/22})$ , recovering the rate for a smooth problem. Looking at Figure 9.4, we see strong anisotropic refinement towards t = 0 with elements of size  $h_x = 1$ ,  $h_t = 2^{-18}$ , and some mild refinement towards the corners of the unit square.

On the other hand, the weighted  $L_2$ -terms  $\zeta_{\mathcal{P}}(\Phi)$  do not seem to decay for uniform or isotropic refinement, and seem to degenerate for anisotropic refinement. This is problematic for the reliability bound in Corollary 9.3.5. Further inspection suggests that this is a theoretical shortcoming rather than a practical one. This is hinted by the (h - h/2)-estimator, which one generally assumes to be reliable. Note that this does not contradict the theoretical results from Corollary 9.3.5, which states  $\zeta_{\mathcal{P}}(\Phi) \lesssim \eta_{\mathcal{P}}(\Phi) \lesssim \|\phi - \Phi\|_{H^{-1/2, -1/4}(\Sigma)}$  only under the additional parabolic scaling assumption  $h_t \approx h_x^2$  for all space-time elements.

Under this parabolic scaling assumption, the optimal error decay rate for smooth problems becomes  $\mathcal{O}(N^{-1/2})$ ; see (9.15). Figure 9.5 displays the results of uniform and adaptive refinement, with meshes that satisfy this scaling constraint<sup>2</sup>, providing convergence rates  $\mathcal{O}(N^{-0.18})$  and  $\mathcal{O}(N^{-0.4})$ , respectively, for *all* considered estimators.

#### 9.4.6 Singular L-shape problem

We consider the L-shaped domain  $\Omega := (-1,1)^2 \setminus [-1,0]^2$  with data  $u_0 \equiv 1$ and  $u_D \equiv 0$ . The solution has a strong singularity for t = 0, in addition to

<sup>&</sup>lt;sup>2</sup> For uniform refinement, all elements are bisected once in space-direction and three times in time-direction. For adaptive refinement, we assume  $\mathcal{M}_{\ell}^x = \mathcal{M}_{\ell}^t$ . All these marked elements are subdivided into four congruent rectangles, where we use additional bisections in space and/or time to guarantee that the level differences between elements sharing an edge is bounded by 1 and that  $\frac{1}{2}h_x^2 \leq h_t \leq 2h_x^2$ .





Figure 9.4. Mesh bv anisotropic Section 9.4.5.

FIGURE 9.5. Error estimators for the singuwith lar problem of Section 9.4.5 plotted double-N = 1391 elements, generated logarithmically over the degrees of freedom refinement  $N = \# \mathcal{P}$ : uniform refinement (left) and for the singular problem of adaptive refinement (right) with parabolic scaling  $h_t \equiv h_x^2$ .



FIGURE 9.6. Error estimators for the singular L-shape problem of Section 9.4.6 plotted double-logarithmically over the degrees of freedom  $N = \#\mathcal{P}$ : uniform refinement (left), isotropic refinement (middle), and anisotropic refinement (right).

a singularity at the re-entrant corner (0,0). We choose  $\mathcal{P}_0 := \{[0,1] \times K :$  $K \in \mathcal{P}_{\Gamma}$ , with the uniform mesh  $\mathcal{P}_{\Gamma}$  of  $\Gamma$  being aligned with the corners and consisting of eight elements, as initial mesh of the space-time boundary  $\Sigma$ . Figure 9.6 displays the results, which are similar to those of Section 9.4.5 with a better behavior of the  $L_2$ -terms  $\zeta_{\mathcal{P}}(\Phi)$  for anisotropic refinement.

Enforcing the parabolic scaling  $h_t \approx h_x^2$  as in Section 9.4.5, *all* estimators converge again with the same rates, being  $\mathcal{O}(N^{-0.18})$  for uniform refinement and  $\mathcal{O}(N^{-0.45})$  for adaptive refinement (not displayed).

## Summary

Partial differential equations are used for the modeling of many (natural) phenomena such as fluid flow, heat dissipation, sound propagation, chemical reactions and the weather. Typically, it is impossible to find closed-form solutions to these differential equations. Instead, one may use numerical methods to construct approximations of the solution to these differential equations. Ideally, these numerical methods provide us with *high quality* approximations and do so *quickly*.

In this thesis we look at such numerical methods for approximating linear partial differential equations. This thesis is structured into two parts. In the first part we study *preconditioning*, a technique used to accelerate iterative solvers for systems of linear equations arising in our approximation schemes. In the second part we focus on (adaptive) numerical methods for parabolic evolution equations in a simultaneous *space-time* approach.

In Part I we study the construction of preconditioners for discretized operators using the concept of *operator preconditioning*. The idea is to precondition the discretized operator by a discretized operator of opposite order. It turns out that in order to get a uniformly well-conditioned system, as well as a preconditioner that can be implemented efficiently, the second discretization has to be carefully chosen dependent on the first one.

In Chapter 2 we apply this idea to construct uniform preconditioners for operators of *negative order* discretized by (dis)continuous piecewise polynomials, of any degree, with respect to some mesh. The application cost of the preconditioner is the cost of the opposite order operator discretized by continuous piecewise linears on the same mesh, plus minor cost that scales linearly in the number of mesh cells. Compared to earlier proposals, our approach has the advantages that it does not require a (barycentric) refinement of the mesh, nor the inversion of a non-diagonal matrix, and that it applies without any mildly grading assumption on the mesh.

In Chapter 3 we propose a multi-level type operator that both fulfills the role of the opposite order operator *and* can be applied in optimal *linear complexity*. So, when it is used within the framework from Chapter 2, it provides a uniform

preconditioner for operators of negative order that can be applied in linear complexity.

In Chapter 4 we investigate the operator preconditioning framework for *positive order* operators, discretized by continuous piecewise polynomials with respect to some mesh. We obtain uniform preconditioners with similar advantages as those found in Chapter 2. That is, we construct uniform preconditioners whose application cost equals the cost of the opposite order operator discretized with discontinuous or continuous finite elements on the same mesh, plus minor cost of linear complexity.

In Chapter 5 we study the preconditioning framework in a more restrictive setting. We build uniform preconditioners for operators discretized by continuous piecewise polynomials as the composition of an opposite order operator, discretized on the same ansatz space, and two identical diagonal scaling operators, whose matrix representation is the lumped mass matrix.

In Part II of this thesis, we discuss the (adaptive) numerical solution of parabolic evolution equations, e.g. the heat equation, written in a simultaneous spacetime variational formulation. The 'natural' weak formulation of these problems is not coercive, making it hard to discretize the problem. We consider a minimal residual discretization, which leads to quasi-optimal approximations under an inf-sup stability condition on the pairs of trial- and test spaces.

In Chapter 6 we propose an *r*-linearly converging adaptive method for parabolic evolution equations that is able to resolve singularities locally in space and time. We achieve this by using trial- and test spaces that are given as the spans of wavelets-in-time tensorized with (locally refined) finite element spaces-in-space. Thanks to this tensor product ansatz, we are able to solve the whole time evolution at the cost of solving the corresponding stationary problem. We also introduce optimal preconditioners, allowing one to solve the discrete problem efficiently. Numerical results illustrate our theoretical findings.

In Chapter 7 we discuss an implementation of the aforementioned adaptive method in which every step is of linear complexity. In particular, we propose a matrix-free algorithm that can apply the system matrices in linear complexity, even though their matrix representation is not sparsely populated.

In Chapter 8 we consider a (time-)parallel algorithm for parabolic evolution equations using trial- and test spaces that are tensor products of temporal and spatial spaces. We present large parallel computations showing the effective-ness of the method in practice.

Finally, in Chapter 9 we propose an adaptive space-time boundary element method for the solution of the homogeneous heat equation with prescribed initial condition and Dirichlet data. We introduce an a posteriori error estimator, and use it to drive an adaptive loop with anisotropic refinement. In all our numerical experiments, the adaptive algorithm recovers the optimal error decay rate.

## Samenvatting

Partiële differentiaalvergelijkingen worden gebruikt voor het modelleren van veel (natuurlijke) verschijnselen, zoals de stroming van vloeistoffen, de verspreiding van warmte, de voortplanting van geluid, het verloop van chemische reacties en het weer. In veel realistische gevallen is het onmogelijk om een expliciet functievoorschrift voor de oplossing van de differentiaalvergelijking te vinden. Daarom gebruikt men numerieke methoden om benaderingen van de oplossing uit te rekenen. Idealiter geven deze numerieke methoden *snel* een *goede* benadering.

In dit proefschrift doen we onderzoek naar dergelijke numerieke methoden. Het proefschrift is opgebouwd uit twee delen. In het eerste deel bestuderen we *preconditionering*, een techniek voor het versnellen van iteratieve methoden die gebruikt worden om de stelsels lineaire vergelijkingen in de numerieke methoden op te lossen. In het tweede deel bekijken we het (adaptief) benaderen van parabolische evolutieproblemen, zoals de warmtevergelijking, met ruimte-tijd methoden.

In deel I bekijken we *operatorpreconditionering*: een methode om een gediscretiseerde operator te preconditioneren met een gediscretiseerde operator van tegenovergestelde orde. Om te zorgen dat dit tot uniform goed-geconditioneerde systemen leidt en dat de preconditioneerder efficiënt kan worden toegepast, moet de tweede discretisatieruimte zorgvuldig gekozen worden.

In hoofdstuk 2 passen we dit concept toe om uniforme preconditioneerders te maken voor operatoren van negatieve orde gediscretiseerd met (dis)continue stuksgewijs polynomen, van een willekeurige graad, op een gegeven rooster. De kosten van het toepassen van de preconditioneerder zijn gelijk aan de kosten van de gediscretiseerde operator van tegenovergestelde orde op hetzelfde rooster, plus lage kosten die lineair schalen in het aantal roosterelementen. In vergelijking met methoden uit de literatuur, heeft onze aanpak de voordelen dat er geen (barycentrisch) verfijnd rooster nodig is, dat er geen niet-diagonaal matrix geïnverteerd moet worden en we geen aanname nodig hebben over dat het verschil in grootte van naastliggende roosterelementen voldoende klein is.

In hoofdstuk 3 construeren we een multischaaloperator die de rol van de tegenovergestelde operator vervult én die in optimale (lineaire) complexiteit

kan worden toegepast. In combinatie met de bevindingen uit hoofdstuk 2 geeft dit een uniforme preconditioneerder voor operators van negatieve orde die in lineaire complexiteit kan worden toegepast.

In hoofdstuk 4 bekijken we nogmaals het concept van operatorpreconditionering, ditmaal voor operatoren van positieve orde gediscretiseerd door continue stuksgewijs polynomen op een rooster. We construeren uniforme preconditioneerders met dezelfde voordelen als die uit hoofdstuk 2.

In hoofdstuk 5 bekijken we het concept van operatorpreconditionering in een beperkte setting. We maken uniforme preconditioneerders voor operatoren die gediscretiseerd zijn met continue stuksgewijze polynomen. De gevonden preconditioneerders bestaan uit de tegenovergestelde operator, gediscretiseerd op dezelfde deelruimte, en twee gelijke diagonale schaaloperatoren.

In deel II van dit proefschrift kijken we naar parabolische evolutieproblemen met ruimte-tijd methoden. De natuurlijke zwakke formulering die hoort bij zulke vergelijkingen is niet coercief, wat het lastig maakt om deze problemen te discretiseren. We bekijken een residu-minimalisatie formulering van het probleem die leidt tot quasi-optimale benaderingen onder een inf-sup aanname op de paren van zoek- en testruimten.

In hoofdstuk 6 introduceren we een adaptieve numerieke methode voor parabolische problemen die *r*-linear convergeert en die singulariteiten lokaal in plaats en tijd kan benaderen. We bereiken dit door gebruik te maken van zoeken testruimten die opgespannen zijn door het tensorproduct van wavelets in tijd en (lokaal verfijnde) eindige elementenruimten in ruimte. Door deze ruimten te gebruiken, kunnen we de gehele tijdsevolutie oplossen met dezelfde kosten als het oplossen van het bijbehorende tijdsonafhankelijke probleem. We introduceren ook optimale preconditioneerders, zodat het gediscretiseerde probleem efficiënt opgelost kan worden. We illustreren onze theoretische bevindingen met numerieke experimenten.

In hoofdstuk 7 bestuderen we een implementatie van de bovengenoemde adaptieve methode, waarin elke stap in lineaire complexiteit kan worden uitgevoerd. In het bijzonder geven we een algoritme om de systeemmatrices in lineaire complexiteit toe te toepassen, hoewel deze matrices niet ijl zijn.

In hoofdstuk 8 geven we een parallel (in tijd) algoritme voor parabolische vergelijkingen met zoek- en testruimten die het tensorproduct zijn van ruimten in tijd en ruimte. We laten de effectiviteit van de methode in de praktijk zien middels grootschalige berekeningen op een supercomputer.

Tot slot, in hoofdstuk 9 introduceren we een adaptieve ruimte-tijd randelementenmethode voor de homogene warmtevergelijking met gegeven beginen randvoorwaarden. We introduceren een a posteriori foutschatter, die we vervolgens gebruiken om een adaptief algoritme te sturen met anisotrope verfijning. In de door ons bekeken experimenten, behaalt deze methode de optimale convergentiesnelheid, ook voor problemen die singulier in tijd en ruimte zijn.

## Dankwoord

Allereerst wil ik mijn promotor Rob Stevenson bedanken. Ik heb veel van je geleerd. Je hebt mij iedere keer opnieuw verbaasd met jouw (wiskundig) inzicht. Dank voor je geduld, advies, humor, het becommentariëren van talloze stukken tekst en de ruimte die je mij hebt geboden om een half jaar stage te lopen. Je bent een goede mentor en ik ben bevoorrecht dat jij mijn promotor was.

Ik ben dankbaar dat ik dit promotietraject samen met jou, Jan Westerdiep, heb mogen doorlopen. De band die we hebben opgebouwd als medestudenten, collega's en vrienden is bijzonder. Je was er altijd om advies te geven of mijn gedachtestroom aan te horen op zowel inhoudelijk vlak als daarbuiten. De samenwerking met jou heeft mij scherp gehouden, zeker in de coronaperiode. Je ongekende enthousiasme, behulpzaamheid en positiviteit maken je een mooi mens. Dank voor alles.

De collega's van het Korteweg-de Vries instituut bedank ik voor het creëren van een prettige werksfeer. Gregor, je was een fijne kamergenoot. Dank voor je belangstelling, de vele discussies, het kritisch lezen van mijn introductie en je toewijding aan het paper dat wij samen hebben gemaakt. Jan Brandts bedank ik voor zijn rol als copromotor en Marieke voor het mogelijk maken van werken op kantoor tijdens de coronapandemie. Mees, Ruben en Jan, de programmeerwedstrijden waar wij als team aan hebben meegedaan waren vermakelijke maar vaak ook frustrerende uitdagingen waar ik veel van heb geleerd, dank hiervoor.

Ik spreek ook mijn dank uit aan alle vrienden om mij heen. Ontspanning in de vorm van sporten, borrelen, vakantie vieren of bordspellen spelen heeft mij veel energie gegeven. In het bijzonder dank ik Floris en Rick voor alle etentjes, Maarten voor de squashpartijen, Danny en Niki voor de fietstochten en Roos voor het zijn van een fijne huisgenoot.

Mijn familie bedank ik voor alle steun door de jaren heen, met name mijn moeder Monique voor haar betrokkenheid en de aanmoedigende woorden wanneer deze nodig waren.

Tot slot bedank ik mijn vriendin Ingrid. Ik waardeer jouw liefdevolle steun enorm, dankjewel.

# Bibliography

[AFF+13]	M. Aurada, M. Feischl, T. Führer, M. Karkulik, and D. Praetorius. Efficiency and optimality of some weighted-residual error estimator for adaptive 2D boundary
[AFF <sup>+</sup> 17]	<ul> <li>element methods. Comput. Methods Appl. Math., 13(3):305–332, 2013.</li> <li>M. Aurada, M. Feischl, T. Führer, M. Karkulik, J.M. Melenk, and D. Praetorius. Local inverse estimates for non-local boundary integral operators. Math. Comp., 86(308):2651–2686, 2017.</li> </ul>
[Alp93]	B.K. Alpert. A class of bases in $L^2$ for the sparse representation of integral operators. <i>SIAM J. Math. Anal.</i> , 24:246–262, 1993.
[AN87]	D.N. Arnold and P.J. Noon. Boundary integral equations of the first kind for the heat equation. In <i>Boundary elements IX, Vol. 3 (Stuttgart, 1987),</i> pages 213–229. Comput. Mech., Southampton, 1987.
[And13]	R. Andreev. Stability of sparse space-time finite element discretizations of linear parabolic evolution equations. <i>IMA J. Numer. Anal.</i> , 33(1):242–260, 2013.
[And16]	R. Andreev. Wavelet-in-time multigrid-in-space preconditioning of parabolic evolution equations. <i>SIAM J. Sci. Comput.</i> , 38(1):A216–A242, 2016.
[BC07]	A. Buffa and S.H. Christiansen. A dual finite element complex on the barycentric refinement. <i>Math. Comp.</i> , 76(260):1743–1769, 2007.
[BD04]	P. Binev and R. DeVore. Fast computation in adaptive tree approximation. <i>Numer. Math.</i> , 97(2):193 – 217, 2004.
[BDDP02]	P. Binev, W. Dahmen, R. DeVore, and P. Petruchev. Approximation classes for adaptive methods. <i>Serdica Math. J.</i> , 28:391–416, 2002.
[Beb00]	M. Bebendorf. Approximation of boundary element matrices. <i>Numer. Math.</i> , 86(4):565–589, 2000.
[BFV19]	P. Binev, F. Fierro, and A. Veeser. Near-best adaptive approximation on conforming meshes. <i>Preprint</i> , 2019, arXiv:1912.13437.
[BJ89]	I. Babuška and T. Janik. The h-p version of the finite element method for parabolic equations. Part I. The p-version in time. <i>Numer. Methods Partial Differential Equations</i> , 5(4):363–399, 1989.
[BJ90]	I. Babuška and T. Janik. The h-p version of the finite element method for parabolic equations. II. The h-p version in time. <i>Numer. Methods Partial Differential Equations</i> , 6(4):343–369, 1990.
[BPV00]	J.H. Bramble, J.E. Pasciak, and P.S. Vassilevski. Computational scales of Sobolev norms with application to preconditioning. <i>Math. Comp.</i> , 69(230):463–480, 2000.
[Bra81]	A. Brandt. Multigrid solvers on parallel computers. In <i>Elliptic Problem Solvers</i> , pages 39–83. Elsevier, 1 1981.
[Bra01]	D. Braess. Finite Elements. Cambridge University Press, 2001. Second edition.
[BRU20]	N. Beranek, M.A. Reinhold, and K. Urban. A space-time variational method for optimal control problems. <i>Preprint</i> , 2020, arXiv:2010.00345.

[BY14]	R.E. Bank and H. Yserentant. On the $H^1$ -stability of the $L_2$ -projection onto finite element spaces. <i>Numer. Math.</i> , 126(2):361–381, 2014.
[BZ96]	R. Balder and C. Zenger. The solution of multidimensional real Helmholtz equations on sparse grids. <i>SIAM J. Sci. Comput.</i> , 17(3):631–646, 1996.
[Car02]	C. Carstensen. Merging the Bramble-Pasciak-Steinbach and the Crouzeix-Thomée criterion for $H^1$ -stability of the $L^2$ -projection onto finite element spaces. <i>Math. Comp.</i> , 71(237):157–163, 2002.
[Car04]	C. Carstensen. An adaptive mesh-refining algorithm allowing for an $H^1$ stable $L^2$ projection onto Courant finite element spaces. <i>Constr. Approx.</i> , 20(4):549–564, 2004.
[CDD01]	A. Cohen, W. Dahmen, and R. DeVore. Adaptive wavelet methods for elliptic operator equations – Convergence rates. <i>Math. Comp</i> , 70:27–75, 2001.
[CN00]	S.H. Christiansen and JC. Nédélec. Des préconditionneurs pour la résolution numérique des équations intégrales de frontière de l'acoustique. <i>C. R. Acad. Sci. Paris Sér. I Math.</i> , 330(7):617–622, 2000.
[Cos90]	M. Costabel. Boundary integral operators for the heat equation. <i>Integral Equations Operator Theory</i> , 13(4):498–552, 1990.
[CR19]	A. Chernov and A. Reinarz. Sparse grid approximation spaces for space-time bound- ary integral formulations of the heat equation. <i>Comput. Math. Appl.</i> , 78(11):3605–3619, 2019.
[CS11]	N.G. Chegini and R.P. Stevenson. Adaptive wavelets schemes for parabolic problems: Sparse matrices and numerical results. <i>SIAM J. Numer. Anal.</i> , 49(1):182–212, 2011.
[CS13]	A. Chernov and Ch. Schwab. Sparse space-time Galerkin BEM for the nonstationary heat equation. <i>ZAMM Z. Angew. Math. Mech.</i> , 93(6–7):403–413, 2013.
[Dev20]	D. Devaud. Petrov-Galerkin space-time $hp$ -approximation of parabolic equations in $H^{1/2}$ . <i>IMA J. Numer. Anal.</i> , 40(4):2717–2745, 2020.
[DFG <sup>+</sup> 04]	W. Dahmen, B. Faermann, I.G. Graham, W. Hackbusch, and S.A. Sauter. Inverse inequalities on non-quasiuniform meshes and application to the mortar element method. <i>Math. Comp.</i> , 73:1107–1138, 2004.
[DGVdZ18]	R. Dyja, B. Ganapathysubramanian, and K.G. Van der Zee. Parallel-in-space-time, adaptive finite element framework for nonlinear parabolic equations. <i>SIAM J. Sci. Comput.</i> , 40(3):C283–C304, 2018.
[Dij09]	T.J. Dijkema. Adaptive tensor product wavelet methods for the solution of PDEs. PhD
[DKS16]	thesis, Utrecht University, 2009.
	<ul><li>thesis, Utrecht University, 2009.</li><li>L. Diening, Ch. Kreuzer, and R.P. Stevenson. Instance optimality of the adaptive maximum strategy. <i>Found. Comput. Math.</i>, 16(1):33–68, 2016.</li></ul>
[DL92]	<ul> <li>thesis, Utrecht University, 2009.</li> <li>L. Diening, Ch. Kreuzer, and R.P. Stevenson. Instance optimality of the adaptive maximum strategy. <i>Found. Comput. Math.</i>, 16(1):33–68, 2016.</li> <li>R. Dautray and JL. Lions. <i>Mathematical analysis and numerical methods for science and technology. Vol. 5.</i> Springer-Verlag, Berlin, 1992.</li> </ul>
[DL92] [DNS19]	<ul> <li>thesis, Utrecht University, 2009.</li> <li>L. Diening, Ch. Kreuzer, and R.P. Stevenson. Instance optimality of the adaptive maximum strategy. <i>Found. Comput. Math.</i>, 16(1):33–68, 2016.</li> <li>R. Dautray and JL. Lions. <i>Mathematical analysis and numerical methods for science and technology. Vol. 5.</i> Springer-Verlag, Berlin, 1992.</li> <li>S. Dohr, K. Niino, and O. Steinbach. Space-time boundary element methods for the heat equation. In <i>Space-Time Methods: Applications to Partial Differential Equations</i>, pages 1–60. De Gruyter, 2019.</li> </ul>
[DL92] [DNS19] [Doh19]	<ul> <li>thesis, Utrecht University, 2009.</li> <li>L. Diening, Ch. Kreuzer, and R.P. Stevenson. Instance optimality of the adaptive maximum strategy. <i>Found. Comput. Math.</i>, 16(1):33–68, 2016.</li> <li>R. Dautray and JL. Lions. <i>Mathematical analysis and numerical methods for science and technology. Vol. 5.</i> Springer-Verlag, Berlin, 1992.</li> <li>S. Dohr, K. Niino, and O. Steinbach. Space-time boundary element methods for the heat equation. In <i>Space-Time Methods: Applications to Partial Differential Equations</i>, pages 1–60. De Gruyter, 2019.</li> <li>S. Dohr. <i>Distributed and Preconditioned Space-Time Boundary Element Methods for the Heat Equation.</i> PhD thesis, TU Graz, 2019.</li> </ul>
[DL92] [DNS19] [Doh19] [DPS05]	<ul> <li>thesis, Utrecht University, 2009.</li> <li>L. Diening, Ch. Kreuzer, and R.P. Stevenson. Instance optimality of the adaptive maximum strategy. <i>Found. Comput. Math.</i>, 16(1):33–68, 2016.</li> <li>R. Dautray and JL. Lions. <i>Mathematical analysis and numerical methods for science and technology. Vol. 5.</i> Springer-Verlag, Berlin, 1992.</li> <li>S. Dohr, K. Niino, and O. Steinbach. Space-time boundary element methods for the heat equation. In <i>Space-Time Methods: Applications to Partial Differential Equations</i>, pages 1–60. De Gruyter, 2019.</li> <li>S. Dohr. <i>Distributed and Preconditioned Space–Time Boundary Element Methods for the Heat Equation.</i> PhD thesis, TU Graz, 2019.</li> <li>L. Dalcín, R. Paz, and M. Storti. MPI for Python. <i>J. Parallel Distrib. Comput.</i>, 65(9):1108–1115, 9 2005.</li> </ul>
[DL92] [DNS19] [Doh19] [DPS05] [DS18]	<ul> <li>thesis, Utrecht University, 2009.</li> <li>L. Diening, Ch. Kreuzer, and R.P. Stevenson. Instance optimality of the adaptive maximum strategy. <i>Found. Comput. Math.</i>, 16(1):33–68, 2016.</li> <li>R. Dautray and JL. Lions. <i>Mathematical analysis and numerical methods for science and technology. Vol. 5.</i> Springer-Verlag, Berlin, 1992.</li> <li>S. Dohr, K. Niino, and O. Steinbach. Space-time boundary element methods for the heat equation. In <i>Space-Time Methods: Applications to Partial Differential Equations</i>, pages 1–60. De Gruyter, 2019.</li> <li>S. Dohr. <i>Distributed and Preconditioned Space-Time Boundary Element Methods for the Heat Equation.</i> PhD thesis, TU Graz, 2019.</li> <li>L. Dalcín, R. Paz, and M. Storti. MPI for Python. <i>J. Parallel Distrib. Comput.</i>, 65(9):1108–1115, 9 2005.</li> <li>D. Devaud and Ch. Schwab. Space-time <i>hp</i>-approximation of parabolic equations. <i>Calcolo</i>, 55(3):Paper No. 35, 23, 2018.</li> </ul>
[DL92] [DNS19] [Doh19] [DPS05] [DS18] [DS20]	<ul> <li>thesis, Utrecht University, 2009.</li> <li>L. Diening, Ch. Kreuzer, and R.P. Stevenson. Instance optimality of the adaptive maximum strategy. <i>Found. Comput. Math.</i>, 16(1):33–68, 2016.</li> <li>R. Dautray and JL. Lions. <i>Mathematical analysis and numerical methods for science and technology. Vol. 5.</i> Springer-Verlag, Berlin, 1992.</li> <li>S. Dohr, K. Niino, and O. Steinbach. Space-time boundary element methods for the heat equation. In <i>Space-Time Methods: Applications to Partial Differential Equations</i>, pages 1–60. De Gruyter, 2019.</li> <li>S. Dohr. <i>Distributed and Preconditioned Space-Time Boundary Element Methods for the Heat Equation</i>. PhD thesis, TU Graz, 2019.</li> <li>L. Dalcín, R. Paz, and M. Storti. MPI for Python. <i>J. Parallel Distrib. Comput.</i>, 65(9):1108–1115, 9 2005.</li> <li>D. Devaud and Ch. Schwab. Space-time <i>hp</i>-approximation of parabolic equations. <i>Calcolo</i>, 55(3):Paper No. 35, 23, 2018.</li> <li>L. Diening and J. Storn. A Space-Time DPG Method for the Heat Equation. <i>Preprint</i>, 2020, arXiv:2012.13229.</li> </ul>

- [DSW21] W. Dahmen, R.P. Stevenson, and J. Westerdiep. Accuracy controlled data assimilation for parabolic problems. *Preprint*, 2021, arXiv:2105.05836.
- [DZO<sup>+</sup>19] S. Dohr, J. Zapletal, G. Of, M. Merta, and M. Kravčenko. A parallel space-time boundary element method for the heat equation. *Comput. Math. Appl.*, 78(9):2852– 2866, 2019.
- [EG04] A. Ern and J.-L. Guermond. Theory and practice of finite elements, volume 159 of Applied Mathematical Sciences. Springer, New York, 2004.
- [ESV17] A. Ern, I. Smears, and M. Vohralík. Guaranteed, locally space-time efficient, and polynomial-degree robust a posteriori error estimates for high-order discretizations of parabolic problems. *SIAM J. Numer. Anal.*, 55(6):2811–2834, 2017.
- [Fae00] B. Faermann. Localization of the Aronszajn-Slobodeckij norm and application to adaptive boundary elements methods. Part I. The two-dimensional case. IMA J. Numer. Anal., 20(2):203–234, 2000.
- [Fae02] B. Faermann. Localization of the Aronszajn-Slobodeckij norm and application to adaptive boundary element methods. Part II. The three-dimensional case. *Numer. Math.*, 92(3):467–499, 2002.
- [FFK<sup>+</sup>14] R. D. Falgout, S. Friedhoff, Tz. V. Kolev, S. P. MacLachlan, and J. B. Schroder. Parallel Time Integration with Multigrid. *SIAM J. Sci. Comput.*, 36(6):C635–C661, 1 2014.
- [FH19] T. Führer and N. Heuer. Optimal quasi-diagonal preconditioners for pseudodifferential operators of order minus two. J. Sci. Comput., 79(2):1161–1181, 2019.
- [FHPS19] T. Führer, A. Haberl, D. Praetorius, and S. Schimanko. Adaptive BEM with inexact PCG solver yields almost optimal computational costs. *Numer. Math.*, 141(4):967–1008, 2019.
- [FK21] T. Führer and M. Karkulik. Space-time least-squares finite elements for parabolic equations. *Comput. Math. Appl.*, 92:27–36, 6 2021.
- [For77] M. Fortin. An analysis of the convergence of mixed finite element methods. RAIRO Anal. Numér., 11(4):341–354, iii, 1977.
- [Gan15] M.J. Gander. 50 years of time parallel time integration. In *Multiple shooting and time domain decomposition methods*, volume 9 of *Contrib. Math. Comput. Sci.*, pages 69–113. Springer, Cham, 2015.
- [Gan17] G. Gantner. *Optimal adaptivity for splines in finite and boundary element methods*. PhD thesis, TU Wien, 2017.
- [GHS05] I.G. Graham, W. Hackbusch, and S.A. Sauter. Finite elements on degenerate meshes: inverse-type inequalities and applications. IMA J. Numer. Anal., 25(2):379–407, 2005.
- [GHS16] F. D. Gaspoz, C.-J. Heine, and K. G. Siebert. Optimal grading of the newest vertex bisection and H<sup>1</sup>-stability of the L<sub>2</sub>-projection. IMA J. Numer. Anal., 36(3):1217–1241, 2016.
- [GK11] M.D. Gunzburger and A. Kunoth. Space-time adaptive wavelet methods for control problems constrained by parabolic evolution equations. SIAM J. Contr. Optim., 49(3):1150–1170, 2011.
- [Glä12] M. Gläfke. *Adaptive methods for time domain boundary integral equations*. PhD thesis, Brunel University London, 2012.
- [GN16] M.J. Gander and M. Neumüller. Analysis of a new space-time parallel multigrid algorithm for parabolic problems. *SIAM J. Sci. Comput.*, 38(4):A2173–A2208, 2016.
- [GO07] M. Griebel and D. Oeltz. A sparse grid space-time discretization scheme for parabolic problems. *Computing*, 81(1):1–34, 2007.
- [GÖSS20] H. Gimperlein, C. Özdemir, D. Stark, and E.P. Stephan. A residual a posteriori error estimate for the time-domain boundary element method. *Numer. Math.*, 146(2):239– 280, 2020.
- [GP20] G. Gantner and D. Praetorius. Adaptive BEM for elliptic PDE systems, part I: abstract framework, for weakly-singular integral equations. *Appl. Anal.*, published online:1– 34, 2020.

[GS19]	H. Gimperlein and J. Stocek. Space-time adaptive finite elements for nonlocal parabolic variational inequalities. <i>Comput. Methods Appl. Mech. Engrg.</i> , 352:137–171, 2019.
[GS21]	G. Gantner and R.P. Stevenson. Further results on a space-time FOSLS formulation of parabolic PDEs. <i>ESAIM Math. Model. Numer. Anal.</i> , 55:283–299, 2021.

- [GvV21] G. Gantner and R. van Venetië. Adaptive space-time BEM for the heat equation. *Preprint*, 2021, arXiv:2108.03055.
- [Hac85] W. Hackbusch. Multi-Grid Methods and Applications, volume 4 of Springer Series in Computational Mathematics. Springer Berlin Heidelberg, Berlin, Heidelberg, 1985.
- [Hac99] W. Hackbusch. A sparse matrix arithmetic based on *H*-matrices. Part i: Introduction to *H*-matrices. *Computing*, 62(2):89–108, 1999.
- [Hac12] W. Hackbusch. Tensor Spaces and Numerical Tensor Calculus, volume 42 of Springer Series in Computational Mathematics. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.
- [Hip06] R. Hiptmair. Operator preconditioning. Comput. Math. Appl., 52(5):699–706, 2006.
- [HJHUT18] R. Hiptmair, C. Jerez-Hanckes, and C. Urzúa-Torres. Closed-form inverses of the weakly singular and hypersingular operators on disks. *Integral Equations Operator Theory*, 90(1):Art. 4, 14, 2018.
- [HJHUT20] R. Hiptmair, C. Jerez-Hanckes, and C. Urzúa-Torres. Optimal operator preconditioning for Galerkin boundary element methods on 3-dimensional screens. SIAM J. Numer. Anal., 58(1):834–857, 2020.
- [HK12] R. Hiptmair and L. Kielhorn. BETL A generic boundary element template library. Technical Report 2012-36, Seminar for Applied Mathematics, ETH Zürich, Switzerland, 2012.
- [HLNS19] Ch. Hofer, U. Langer, M. Neumüller, and R. Schneckenleitner. Parallel and robust preconditioning for space-time isogeometric analysis of parabolic evolution problems. *SIAM J. Sci. Comput.*, 41(3):A1793–A1821, 2019.
- [HT18] H. Harbrecht and J. Tausch. A fast sparse grid based space–time boundary element method for the nonstationary heat equation. *Numer. Math.*, 140(1):1–26, 2018.
- [HUT16] R. Hiptmair and C. Urzúa-Torres. Dual mesh operator preconditioning on 3d screens: Low-order boundary element discretization. Technical Report 2016-14, Seminar for Applied Mathematics, ETH Zürich, Switzerland, 2016.
- [HVW95] G. Horton, S. Vandewalle, and P. Worley. An Algorithm with Polylog Parallel Complexity for Solving Parabolic Partial Differential Equations. SIAM J. Sci. Comput., 16(3):531–541, 5 1995.
- [Kat60] T. Kato. Estimation of iterated matrices, with application to the von Neumann condition. *Numer. Math.*, 2:22–29, 1960.
- [KS08] Y. Kondratyuk and R.P. Stevenson. An optimal adaptive finite element method for the Stokes problem. SIAM J. Numer. Anal., 46(2):747–775, 2008.
- [KS14] S. Kestler and R.P. Stevenson. Fast evaluation of system matrices w.r.t. multi-tree collections of tensor product refinable basis functions. J. Comput. Appl. Math., 260:103– 116, 2014.
- [KSU16] S. Kestler, K. Steih, and K. Urban. An efficient space-time adaptive wavelet Galerkin method for time-periodic parabolic partial differential equations. *Math. Comp.*, 85(299):1309–1333, 2016.
- [LM17] S. Larsson and M. Molteni. Numerical solution of parabolic problems based on a weak space-time formulation. *Comput. Methods Appl. Math.*, 17(1):65–84, 2017.
- [LMN16] U. Langer, S.E. Moore, and M. Neumüller. Space-time isogeometric analysis of parabolic evolution problems. *Comput. Methods Appl. Mech. Engrg.*, 306:342–363, 2016.

- [LMT01] J.-L. Lions, Y. Maday, and G. Turinici. Résolution d'EDP par un schéma en temps «pararéel ». Comptes Rendus de l'Académie des Sciences - Series I - Mathematics, 332(7):661–668, 4 2001.
- [LS20] U. Langer and A. Schafelner. Adaptive Space-Time Finite Element Methods for Nonautonomous Parabolic Problems with Distributional Sources. Comput. Methods Appl. Math., 20(4):677–693, 2020.
- [Mau95] J. Maubach. Local bisection refinement for n-simplicial grids generated by reflection. SIAM J. Sci. Comput., 16(1):210–227, 1995.
- [MFL+91] O.A McBryan, P.O Frederickson, J. Lindenand, A. Schüller, K. Solchenbach, K. Stüben, C.-A. Thole, and U. Trottenberg. Multigrid methods on parallel computers—a survey of recent developments. *Impact Comput. Sci. Engrg.*, 3(1):1–75, 1991.
- [MST14] M. Messner, M. Schanz, and J. Tausch. A fast Galerkin method for parabolic spacetime boundary integral equations. J. Comput. Phys., 258:15–30, 2014.
- [MST15] M. Messner, M. Schanz, and J. Tausch. An efficient Galerkin boundary element method for the transient heat equation. SIAM J. Sci. Comput., 37(3):A1554–A1576, 2015.
- [Nie64] J. Nievergelt. Parallel methods for integrating ordinary differential equations. *Comm. ACM*, 7:731–733, 1964.
- [Nit06] P.-A. Nitsche. Best N-term approximation spaces for tensor product wavelet bases. Constr. Approx., 24(1):49–70, 2006.
- [Noo88] P.J. Noon. *The single layer heat potential and Galerkin boundary element methods for the heat equation.* PhD thesis, University of Maryland, 1988.
- [NS19] M. Neumüller and I. Smears. Time-parallel iterative solvers for parabolic evolution equations. *SIAM J. Sci. Comput.*, 41(1):C28–C51, 2019.
- [OR00] M.A. Olshanskii and A. Reusken. On the convergence of a multigrid method for linear reaction-diffusion problems. *Computing*, 65(3):193–202, 2000.
- [Osw94] P. Oswald. *Multilevel finite element approximation: Theory and applications*. B.G. Teubner, Stuttgart, 1994.
- [Pab15] R. Pabel. Adaptive Wavelet Methods for Variational Formulations of Nonlinear Elliptic PDEs on Tensor-Product Domains. PhD thesis, Universität zu Köln, 2015.
- [Pf110] D. Pflüger. Spatially Adaptive Sparse Grids for High-Dimensional Problems. PhD thesis, Technische Universität München, 2010.
- [PP20] C.-M. Pfeiler and D. Praetorius. Dörfler marking with minimal cardinality is a linear complexity problem. *Math. Comp.*, 89(326):2735–2752, 2020.
- [PW12] J. W. Pearson and A. J. Wathen. A new approximation of the Schur complement in preconditioners for PDE-constrained optimization. *Numer. Linear Algebra Appl.*, 19(5):816–829, 2012.
- [Rek18] N. Rekatsinas. Optimal adaptive wavelet methods for solving first order system least squares. PhD thesis, University of Amsterdam, 2018.
- [RS18] N. Rekatsinas and R.P. Stevenson. A quadratic finite element wavelet Riesz basis. Int. J. Wavelets Multiresolut. Inf. Process., 16(4):1850033, 17, 2018.
- [RS19] N. Rekatsinas and R. Stevenson. An optimal adaptive tensor product wavelet solver of a space-time FOSLS formulation of parabolic evolution problems. *Adv. Comput. Math.*, 45(2):1031–1066, 4 2019.
- [SBA+15] W. Śmigaj, T. Betcke, S. Arridge, J. Phillips, and M. Schweiger. Solving boundary integral problems with BEM++. ACM Trans. Math. Software, 41(2):Art. 6, 40, 2015.
- [Sch14] Joachim Schöberl. C++11 Implementation of Finite Elements in NGSolve. Technical report, Institute for Analysis and Scientific Computing, Vienna University of Technology, 2014.
- [Smi00] R.N.L. Smith. Direct Gauss quadrature formulae for logarithmic singularities on isoparametric elements. *Eng. Anal. Bound. Elem.*, 24(2):161–167, 2000.

[SS09]	Ch. Schwab and R.P. Stevenson. Space-time adaptive wavelet methods for parabolic evolution problems. <i>Math. Comp.</i> , 78(267):1293–1318, 2009.
[Ste96]	R.P. Stevenson. The frequency decomposition multi-level method: A robust additive hierarchical basis preconditioner. <i>Math. Comp.</i> , 65(215):983–997, July 1996.
[Ste98]	R.P. Stevenson. Stable three-point wavelet bases on general meshes. <i>Numer. Math.</i> , 80(1):131–158, 1998.
[Ste02]	O. Steinbach. On a generalized $L_2$ projection and some related stability estimates in Sobolev spaces. <i>Numer. Math.</i> , 90(4):775–786, 2002.
[Ste03a]	O. Steinbach. <i>Stability estimates for hybrid coupled domain decomposition methods</i> , volume 1809 of <i>Lecture Notes in Mathematics</i> . Springer-Verlag, Berlin, 2003.
[Ste03b]	R.P. Stevenson. Locally supported, piecewise polynomial biorthogonal wavelets on nonuniform meshes. <i>Constr. Approx.</i> , 19(4):477–508, 2003.
[Ste08a]	Olaf Steinbach. <i>Numerical approximation methods for elliptic boundary value problems</i> . Springer, New York, 2008. Finite and boundary elements, Translated from the 2003 German original.
[Ste08b]	R.P. Stevenson. The completion of locally refined simplicial partitions created by bisection. <i>Math. Comp.</i> , 77:227–241, 2008.
[Ste15]	O. Steinbach. Space-Time Finite Element Methods for Parabolic Problems. <i>Comput. Methods Appl. Math.</i> , 15(4):551–566, 2015.
[SU09]	W. Sickel and T. Ullrich. Tensor products of Sobolev-Besov spaces and applications to approximation from the hyperbolic cross. <i>J. Approx. Theory</i> , 161:748–786, 2009.
[SvV20a]	R.P. Stevenson and R. van Venetië. Uniform preconditioners for problems of negative order. <i>Math. Comp.</i> , 89(322):645–674, 2020.
[SvV20b]	R.P. Stevenson and R. van Venetië. Uniform preconditioners for problems of positive order. <i>Comput. Math. Appl.</i> , 79(12):3516–3530, 2020.
[SvV21a]	R.P. Stevenson and R. van Venetië. Uniform Preconditioners of Linear Complexity for Problems of Negative Order. <i>Comput. Methods Appl. Math.</i> , 21(2):469–478, 2021.
[SvV21b]	R.P. Stevenson and R. van Venetië. Operator preconditioning: the simplest case. <i>Preprint</i> , 2021, arXiv:2102.11951.
[SvVW21]	R.P. Stevenson, R. van Venetië, and J. Westerdiep. A wavelet-in-time, finite element-in- space adaptive method for parabolic evolution equations. <i>Preprint</i> , arXiv:2101.03956, 2021.
[SW98]	O. Steinbach and W. L. Wendland. The construction of some efficient preconditioners in the boundary element method. <i>Adv. Comput. Math.</i> , 9(1-2):191–216, 1998.
[SW21a]	R.P. Stevenson and J. Westerdiep. Minimal residual space-time discretizations of parabolic equations: Asymmetric spatial operators. <i>Preprint</i> , 2021, arXiv:2106.01090.
[SW21b]	R.P. Stevenson and J. Westerdiep. Stability of Galerkin discretizations of a mixed space-time variational formulation of parabolic evolution equations. <i>IMA J. Numer. Anal.</i> , 41(1):28–47, 2021.
[SY18]	O. Steinbach and H. Yang. Comparison of algebraic multigrid methods for an adap- tive space-time finite-element discretization of the heat equation in 3D and 4D. <i>Numer.</i> <i>Linear Algebra Appl.</i> , 25(3):e2143, 17, 2018.
[SZ90]	L. R. Scott and S. Zhang. Finite element interpolation of nonsmooth functions satisfying boundary conditions. <i>Math. Comp.</i> , 54(190):483–493, 1990.
[SZ20]	O. Steinbach and M. Zank. Coercive space-time finite element methods for initial boundary value problems. <i>Electron. Trans. Numer. Anal.</i> , 52:154–194, 2020.
[Szy06]	D.B. Szyld. The many proofs of an identity on the norm of oblique projections. <i>Numer. Algorithms</i> , 42(3-4):309–323, 2006.
[Tau19]	J. Tausch. Nyström method for BEM of the heat equation with moving boundaries. <i>Adv. Comput. Math.</i> , 45(5):2953–2968, 2019.

- [Tra97] C.T. Traxler. An algorithm for adaptive mesh refinement in *n* dimensions. *Computing*, 59(2):115–137, 1997.
- [Vir20] P. Virtanen. SciPy 1.0: fundamental algorithms for scientific computing in Python. Nat. Methods, 17(3):261–272, 3 2020.
- [VS18] F. J. Vermolen and A. Segal. On an integration rule for products of barycentric coordinates over simplexes in  $\mathbb{R}^n$ . J. Comput. Appl. Math., 330:289–294, 2018.
- [vVW21a] R. van Venetië and J. Westerdiep. A parallel algorithm for solving linear parabolic evolution equations. In 9th Parallel-in-Time Workshop, 2021, arXiv:2009.08875.
- [vVW21b] R. van Venetië and J. Westerdiep. Efficient space-time adaptivity for parabolic evolution equations using wavelets in time and finite elements in space. *Preprint*, 2021, arXiv:2104.08143.
- [vVW21c] R. van Venetië and J. Westerdiep. Implementation of: A parallel algorithm for solving linear parabolic evolution equations, 2021, zenodo:4475959.
- [vVW21d] R. van Venetië and J. Westerdiep. Implementation of: Efficient space-time adaptivity for parabolic evolution equations using wavelets in time and finite elements in space, 2021, zenodo:4697250.
- [Wat87] A. J. Wathen. Realistic eigenvalue bounds for the Galerkin mass matrix. IMA Journal of Numerical Analysis, 7(4):449–457, 1987.
- [Wlo82] J. Wloka. *Partielle Differentialgleichungen*. B. G. Teubner, Stuttgart, 1982. Sobolevräume und Randwertaufgaben.
- [Wor91] P.H. Worley. Limits on parallelism in the numerical solution of linear partial differential equations. SIAM J. Sci. Statist. Comput., 12(1):1–35, 1991.
- [WZ17] J. Wu and H. Zheng. Uniform convergence of multigrid methods for adaptive meshes. *Appl. Numer. Math.*, 113:109–123, 2017.
- [XZ03] J. Xu and L. Zikatanov. Some observations on Babuška and Brezzi theories. Numer. Math., 94(1):195–202, 2003.
- [ZMD<sup>+</sup>11] J. Zitelli, I. Muga, L. Demkowicz, J. Gopalakrishnan, D. Pardo, and V. M. Calo. A class of discontinuous Petrov-Galerkin methods. Part IV: the optimal test norm and time-harmonic wave propagation in 1D. J. Comput. Phys., 230(7):2406–2432, 2011.
- [ZWOM21] J. Zapletal, R. Watschinger, G. Of, and M. Merta. Semi-analytic integration for a parallel space-time boundary element method modeling the heat equation. *Preprint*, arXiv:2102.09811, 2021.