



**Franz, Sebastian and Trostorff, Sebastian and Waurick, Marcus (2018)**  
**Numerical methods for changing type systems. IMA Journal of**  
**Numerical Analysis. ISSN 0272-4979 ,**  
**<http://dx.doi.org/10.1093/imanum/dry007>**

This version is available at <https://strathprints.strath.ac.uk/62949/>

**Strathprints** is designed to allow users to access the research output of the University of Strathclyde. Unless otherwise explicitly stated on the manuscript, Copyright © and Moral Rights for the papers on this site are retained by the individual authors and/or other copyright owners. Please check the manuscript for details of any other licences that may have been applied. You may not engage in further distribution of the material for any profitmaking activities or any commercial gain. You may freely distribute both the url (<https://strathprints.strath.ac.uk/>) and the content of this paper for research or private study, educational, or not-for-profit purposes without prior permission or charge.

Any correspondence concerning this service should be sent to the Strathprints administrator:  
[strathprints@strath.ac.uk](mailto:strathprints@strath.ac.uk)

The Strathprints institutional repository (<https://strathprints.strath.ac.uk>) is a digital archive of University of Strathclyde research outputs. It has been developed to disseminate open access research outputs, expose data about those outputs, and enable the management and persistent access to Strathclyde's intellectual output.

## Numerical methods for changing type systems

SEBASTIAN FRANZ\*

*Institute of Scientific Computing, Technische Universität Dresden, 01062 Dresden, Germany*

\*Corresponding author: sebastian.franz@tu-dresden.de

SASCHA TROSTORFF

*Institute of Analysis, Technische Universität Dresden, 01062 Dresden, Germany*

sascha.trostorff@tu-dresden.de

AND

MARCUS WAURICK

*Department of Mathematics, University of Strathclyde, Glasgow, UK*

marcus.waurick@strath.ac.uk

[Received on 01 August 2017; revised on 13 December 2017]

In this paper we develop a numerical method for partial differential equations with changing type. Our method is based on a unified solution theory found by Rainer Picard for several linear equations from mathematical physics. Parallel to the solution theory already developed, we frame our numerical method in a discontinuous Galerkin approach in time with certain exponentially weighted spaces combined with a finite element method in space.

*Keywords:* evolutionary equations; changing type system; discontinuous Galerkin; space-time approach.

### 1. Introduction

Following the rationale presented in the study by (Picard, 2009), most of the classical linear partial differential equations arising in mathematical physics share a common form, namely the form of an evolutionary problem. That is, we consider equations of the form

$$(\partial_t M_0 + M_1 + A)U = F, \quad (1.1)$$

where  $F$  is a given source term,  $\partial_t$  stands for the derivative with respect to time,  $M_0, M_1$  are bounded linear operators on some Hilbert space  $H$  and  $A$  is an unbounded skew-selfadjoint operator in  $H$ , which we shall identify with their canonical extensions to  $H$ -valued functions acting as abstract multiplication operators. We are seeking for a unique solution  $U$  of the above equation. We remark here that we do not impose initial conditions, since we consider the whole real line as time horizon, and hence, we implicitly assume a vanishing initial value at ‘ $-\infty$ ’. To illustrate the setting, we begin with presenting some examples.

EXAMPLE 1.1 Let  $\Omega \subseteq \mathbb{R}^n$  an open nonempty set, where  $n \in \mathbb{N}$ , but, typically  $n \in \{1, 2, 3\}$ . We define the following two differential operators

$$\nabla_0 : H_0^1(\Omega) \subseteq L^2(\Omega) \rightarrow L^2(\Omega)^n,$$

assigning each function  $u \in H_0^1(\Omega)$  its gradient, that is, the column vector of its partial derivatives in each coordinate direction. Moreover, we set

$$\operatorname{div} := -(\nabla_0)^* : D(\operatorname{div}) \subseteq L^2(\Omega)^n \rightarrow L^2(\Omega),$$

which is nothing but the operator assigning each  $L^2$  vector field its distributional divergence with maximal domain, that is,

$$D(\operatorname{div}) = \left\{ v \in L^2(\Omega)^n : \sum_{i=1}^n \partial_i v_i \in L^2(\Omega) \right\}.$$

Since both the operators  $\nabla_0$  and  $\operatorname{div}$  are closed and skew adjoints of one another, we infer that the operator

$$A := \begin{pmatrix} 0 & \operatorname{div} \\ \nabla_0 & 0 \end{pmatrix} : D(\nabla_0) \times D(\operatorname{div}) \subseteq L^2(\Omega) \times L^2(\Omega)^n \rightarrow L^2(\Omega) \times L^2(\Omega)^n$$

is skew selfadjoint on the Hilbert space  $H = L^2(\Omega) \times L^2(\Omega)^n$ . Choosing  $M_0 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$  and  $M_1 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$  in (1.1), the corresponding evolutionary problem reads as

$$\left( \partial_t \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + \begin{pmatrix} 0 & \operatorname{div} \\ \nabla_0 & 0 \end{pmatrix} \right) \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix}.$$

If  $g = 0$ , this is nothing but the wave equation. Indeed, the second line then gives  $\partial_t v = -\nabla_0 u$ , and hence, differentiating the first line with respect to time, we obtain

$$\partial_t^2 u - \operatorname{div} \nabla_0 u = \partial_{tt} u + \operatorname{div} \partial_t v = \partial_t f.$$

Note that  $\operatorname{div} \nabla_0 = \Delta_D$  is the classical Dirichlet–Laplace operator on  $L^2(\Omega)$ .

Choosing  $M_0 = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$  and  $M_1 = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$  in (1.1), the corresponding problem reads as

$$\left( \partial_t \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} + \begin{pmatrix} 0 & \operatorname{div} \\ \nabla_0 & 0 \end{pmatrix} \right) \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix}.$$

Setting again  $g = 0$ , the latter gives the heat equation. Indeed, the second line reads  $v = -\nabla_0 u$  and hence the first line yields

$$\partial_t u - \operatorname{div} \nabla_0 u = \partial_t u + \operatorname{div} v = f.$$

Finally, choosing  $M_0 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$  and  $M_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$  in (1.1), we get

$$\left( \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + \begin{pmatrix} 0 & \operatorname{div} \\ \nabla_0 & 0 \end{pmatrix} \right) \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix},$$

which in the case  $g = 0$  gives the elliptic equation

$$u - \operatorname{div} \nabla_0 u = f.$$

**REMARK 1.2** Note that we can treat the case of homogeneous Neumann boundary conditions in the same way. The only difference is that we define  $\nabla$  as the distributional gradient on  $H^1(\Omega)$  and  $\operatorname{div}_0 := -(\nabla)^*$ . Replacing now  $\nabla_0$  by  $\nabla$  and  $\operatorname{div}$  by  $\operatorname{div}_0$  yields the same hyperbolic, parabolic and elliptic type problem above, but now with homogeneous Neumann boundary conditions.

Example 1.1 shows that evolutionary problems cover all three classical types of partial differential equations, elliptic, parabolic and hyperbolic. However, also, problems of mixed type are covered as the next example shows.

**EXAMPLE 1.3** Recall the setting of Example 1.1. We decompose  $\Omega$  into three measurable, disjoint sets  $\Omega_{\text{ell}}$ ,  $\Omega_{\text{par}}$  and  $\Omega_{\text{hyp}}$  and set  $M_0 = \begin{pmatrix} \chi_{\Omega_{\text{hyp}} \cup \Omega_{\text{par}}} & 0 \\ 0 & \chi_{\Omega_{\text{hyp}}} \end{pmatrix}$  as well as  $M_1 = \begin{pmatrix} \chi_{\Omega_{\text{ell}}} & 0 \\ 0 & \chi_{\Omega_{\text{par}} \cup \Omega_{\text{ell}}} \end{pmatrix}$ . The resulting evolutionary problem then is of mixed type. More precisely, on  $\Omega_{\text{ell}}$  we get an equation of elliptic type, on  $\Omega_{\text{par}}$  the equations becomes parabolic while on  $\Omega_{\text{hyp}}$  the problem is hyperbolic.

**REMARK 1.4** The interested reader might wonder that there is no transmission condition imposed on the unknown quantities along the interfaces of  $\Omega_{\text{ell}}$ ,  $\Omega_{\text{par}}$  and  $\Omega_{\text{hyp}}$ . However, this can be implemented automatically by being in the domain of the corresponding operator sum, as can be seen, for instance, in the study by (Waurick, 2016, Remark 3.2), see also the study by (Picard *et al.*, 2013, An illustrative Example). Another example of a mixed type problem in control theory can be found in the study by (Picard *et al.*, 2016, Remark 6.2).

In the study by (Picard, 2009), the well-posedness of problems of the form (1.1) has been addressed. In fact, it was shown that these problems also cover the classical Maxwell's equations, the equations of linearized elasticity or a general class of coupled phenomena, see, for instance, the studies by (Mukhopadhyay *et al.*, 2015; Picard *et al.*, 2015; Mulholland *et al.*, 2016). All these problems are indeed well-posed (see Section 2 for the precise statement). The purpose of the present article is to provide numerical methods for such problems. In this article, for the applications to follow, we will focus, however, on problems of mixed type of the form sketched in Example 1.3. Moreover, as the spatial discretization has to be developed for each problem separately, anyway, in this work, we will put an emphasis on the time discretization. Furthermore, we want to stress that the null space of  $M_0$  in (1.1) might be infinite-dimensional. Hence, we seek to develop a numerical scheme, which in particular allows for the treatment of a certain class of (partial) differential-*algebraic* equations.

A similar approach for a unified treatment of elliptic and hyperbolic problems was already suggested by Friedrichs in 1958, see the study by (Friedrichs, 1958), where both types of problems are written as an abstract operator equation  $\mathcal{K}u = f$  with an accretive symmetric operator  $\mathcal{K}$ . These equations are known as Friedrichs systems. In particular in the studies by (Antonić *et al.*, 2013, 2014; Burazin & Erceg, 2016) time-dependent Friedrichs systems have been discussed also for the parabolic case. In these references

local operators in space are considered. Nonstationary examples of mixed type have, however, not been treated.

A drawback of Friedrichs systems is that they ignore the particular role of time in the way that they do not distinguish between the time coordinate and the spatial coordinates. Therefore, the characteristic property of time evolution, namely causality, is not considered at all. In contrast, the systems considered here are automatically causal in direction of time due to a uniform positive definiteness constraint on the operators involved (compare Remark 2.3).

Based on the framework of Friedrichs systems, a unified numerical treatment of different types of partial differential equations, also including certain equations of changing type, was studied before. We refer to the articles by (Ern & Guermond, 2006a, 2006b, 2008) and to the Ph.D. thesis by (Jensen, 2004).

We emphasize that our approach covers problems with change of type ranging from elliptic to hyperbolic but also to parabolic type problems on different spatial domains. In this sense, we obtain a natural unified treatment of a class of partial differential equations that might go beyond the consideration of Friedrichs systems. Note that in a very rough comparison the above mentioned operator  $\mathcal{K}$  is not symmetric in our situation. In particular, the operator equations considered also cover Maxwell's equations with eddy current approximation on parts of the underlying domain.

For the numerical treatment of the time derivatives we use a discontinuous Galerkin (dG) method, see also Section 3. The first dG method was published in 1973 on neutron transport (Reed & Hill, 1973). Later the methodology was developed further for classical hyperbolic, parabolic and elliptic problems, see also the survey article by (Cockburn *et al.*, 2000) and the book by (Rivière, 2008). Note that there is a strong connection between dG methods and Runge–Kutta (collocation) methods, see the study by (Akrivis *et al.*, 2011) for parabolic problems.

In Section 2, for convenience, we will recall some essentials for evolutionary equations. In particular, we recall the solution theory of problems of the type of equation (1.1). We will introduce a semidiscretized version, Equation (3.1), of Equation (1.1) at the beginning of Section 3. We will also provide a solution theory for this semidiscretized variant with general underlying (spatial) Hilbert space  $H$  (Proposition 3.2). The remainder of Section 3 is devoted to estimate difference of the exact solution of (1.1) and the approximate solution of (3.1). In Subsection 3.1, we bound the error by solely in terms of the interpolation error, which will eventually be estimated in Subsection 3.2. As our prime example, we address the full space-time discretization of Example 1.3 and derive corresponding error estimates. We verify our theoretical findings in Section 5 by means of a 1 + 1- and a 1 + 2-dimensional numerical example.

In general, we identify functions defined on a subset of  $\mathbb{R}$  with their extension to  $\mathbb{R}$  by 0.

## 2. The setting of evolutionary problems

In this section we briefly recall the well-posedness result stated in the study by (Picard, 2009). For doing so, we need to specify the functional analytic setting. Throughout, let  $H$  be a real Hilbert space.

DEFINITION 2.1 Let  $\rho > 0$  and define the space

$$H_\rho(\mathbb{R}; H) := \left\{ f : \mathbb{R} \rightarrow H; f \text{ meas.}, \int_{\mathbb{R}} |f(t)|_H^2 \exp(-2\rho t) dt < \infty \right\},$$

where we as usual identify functions which are equal almost everywhere. The space  $H_\rho(\mathbb{R}; H)$  is a Hilbert space endowed with the natural inner product given by

$$\langle f, g \rangle_\rho := \int_{\mathbb{R}} \langle f(t), g(t) \rangle_H \exp(-2\rho t) dt \quad (f, g \in H_\rho(\mathbb{R}; H)).$$

Moreover, we define  $\partial_t$  to be the closure of the operator

$$\partial_t : C_c^\infty(\mathbb{R}; H) \subseteq H_\rho(\mathbb{R}; H) \rightarrow H_\rho(\mathbb{R}; H) : \varphi \mapsto \varphi',$$

where by  $C_c^\infty(\mathbb{R}; H)$  we denote the space of infinitely differentiable  $H$ -valued functions on  $\mathbb{R}$  with compact support. We denote the domain of  $\partial_t^k$  by  $H_\rho^k(\mathbb{R}; H)$  for  $k \in \mathbb{N}$ .

Within the setting introduced, we can formulate the well-posedness for evolutionary equations of the form (1.1).

**THEOREM 2.2** (Picard, 2009, Solution Theory) Let  $M_0, M_1 : H \rightarrow H$  be bounded linear operators,  $M_0$  selfadjoint and  $A : D(A) \subseteq H \rightarrow H$  skew selfadjoint. Moreover, assume that there is some  $\rho_0 > 0$  such that

$$\exists \gamma > 0 \forall \rho \geq \rho_0, x \in H : \langle (\rho M_0 + M_1)x, x \rangle_H \geq \gamma \langle x, x \rangle_H.$$

Then, for each  $\rho \geq \rho_0$  and each  $F \in H_\rho(\mathbb{R}; H)$  there exists a unique  $U \in H_\rho(\mathbb{R}; H)$  such that

$$\overline{(\partial_t M_0 + M_1 + A)U} = F, \quad (2.1)$$

where the closure is taken in  $H_\rho(\mathbb{R}; H)$ . Moreover, the following continuity estimate holds

$$\|U\|_\rho \leq \frac{1}{\gamma} \|F\|_\rho.$$

If  $F \in H_\rho^k(\mathbb{R}; H)$  for  $k \in \mathbb{N}$ , then so is  $U$  and we can omit the closure bar in (2.1).

**REMARK 2.3**

- (a) Note that the positive definiteness condition in the latter theorem especially implies  $\langle M_0 x, x \rangle_H \geq 0$  for each  $x \in H$ .
- (b) We remark that  $H_\rho^1(\mathbb{R}; H) \hookrightarrow C_\rho(\mathbb{R}; H)$  by a variant of the Sobolev embedding theorem (Picard & McGhee, 2011, Lemma 3.1.59) or (Kalauch *et al.*, 2014, Lemma 5.2). Here,

$$C_\rho(\mathbb{R}; H) := \left\{ f : \mathbb{R} \rightarrow H; f \text{ cont.}, \sup_{t \in \mathbb{R}} |f(t)| \exp(-\rho t) < \infty \right\}.$$

- (c) If  $F \in H_\rho^1(\mathbb{R}; H)$  then  $U \in H_\rho^1(\mathbb{R}; H)$ , and hence

$$AU = F - \partial_t M_0 U - M_1 U \in H_\rho(\mathbb{R}; H),$$

which yields that  $U(t) \in D(A)$  for almost every  $t \in \mathbb{R}$ . If even  $F, U \in H_\rho^2(\mathbb{R}; H)$  the latter gives  $AU \in H_\rho^1(\mathbb{R}; H)$  and hence, using the Sobolev embedding result (see part (b)),  $U \in C_\rho(\mathbb{R}; D(A))$ .

- (d) We note that the constant  $\gamma$  in the positive definiteness constraint above is chosen uniformly in  $\rho \geq \rho_0$ . This uniformity yields the causality of the solution operator, see e.g., the study by (Picard, 2009).

- (e) The original result in the study by (Picard, 2009) treats a general class of time-translation invariant coefficients. We refer to the studies by (Picard *et al.*, 2013; Waurick, 2015) for nonautonomous variants as well as to the studies by (Trostorff, 2012; Trostorff & Wehowski, 2014) for nonautonomous and/or nonlinear versions of Theorem 2.2.
- (f) In order to incorporate initial value problems, one extends the solution theory to certain distributional right-hand sides (Picard & McGhee, 2011, Section 6.2.5). Indeed, an initial condition of the form  $M_0 U(t_0+) = M_0 x_0$  for some  $t_0 \in \mathbb{R}, x_0 \in H$  can be implemented by solving the equation

$$(\partial_t M_0 + M_1 + A)U = F + \delta_{t_0} M_0 x_0$$

for some  $F \in H_\rho(\mathbb{R}; H)$  supported on  $[t_0, \infty)$ . Employing causality, one indeed obtains the solution  $U$  satisfies the asserted initial condition (see again the study by Picard & McGhee, 2011, Section 6.2.5).

We note that the equations treated in Example 1.1 and Example 1.3 satisfy the conditions of the previous theorem, and hence are well-posed.

### 3. Semidiscretization in time

In this section, we discretize (1.1) with respect to time and do the *a priori* analysis. We assume that  $A, M_0, M_1$  satisfy the assumptions of Theorem 2.2. Let  $\rho \geq \rho_0$  and fix  $T > 0$  and consider the time interval  $[0, T]$  instead of the whole real line. We partition the time interval  $[0, T]$  into subintervals  $I_m = (t_{m-1}, t_m]$  of length  $\tau_m$  for  $m \in \{1, 2, \dots, M\}$  with  $t_0 = 0$  and  $t_M = T$ . Let  $q \in \mathbb{N}$ . We define the space

$$\mathcal{V}^\tau := \{u \in H_\rho(\mathbb{R}; H) : \forall m \in \{1, \dots, M\} : u|_{I_m} \in \mathcal{P}_q(I_m; H)\},$$

where we denote by

$$\mathcal{P}_q(I_m; H) := \text{lin} \{I_m \ni t \mapsto t^k \zeta \in H; k \in \{0, \dots, q\}, \zeta \in H\}$$

the space of  $H$ -valued polynomials of degree at most  $q$  defined on  $I_m$ . We endow  $\mathcal{P}_q(I_m; H)$  with the scalar product

$$\langle p, q \rangle_{\rho, m} := \int_{t_{m-1}}^{t_m} \langle p(t), q(t) \rangle_H \exp(-2\rho(t - t_{m-1})) dt$$

turning the space  $\mathcal{P}_q(I_m; H)$  into a Hilbert space.

The time integrals have to be evaluated numerically. We choose on each time interval  $I_m$  a right-sided weighted Gauß–Radau quadrature formula. To this end, denote by  $\omega_i^m$  and  $\hat{t}_i^m$ ,  $i \in \{0, \dots, q\}$ , the weights and nodes of the weighted Gauß–Radau formula with  $q + 1$  nodes on the reference time interval  $\hat{I} = (-1, 1]$ , such that

$$\int_{\hat{I}} e^{-\rho\tau_m(t+1)} p(t) dt = \sum_{i=0}^q \omega_i^m p(\hat{t}_i^m)$$

holds for all polynomials  $p$  of degree at most  $2q$ . Note that the weights and nodes can always be numerically computed as shown for instance in (Press *et al.*, 2007, Chapter 4.6), see also the technical

paper (Trostorff & Waurick, 2016) for some basic facts on the Gauß–Radau quadrature. With the standard linear transformation  $T_m : \widehat{T} \rightarrow I_m$  and the transformed Gauß–Radau points  $t_{m,i} := T_m(\hat{t}_i^m)$ ,  $i \in \{0, \dots, q\}$ , we define by

$$Q_m[v] := \frac{\tau_m}{2} \sum_{i=0}^q \omega_i^m v(t_{m,i})$$

a quadrature formula on  $I_m$ . Note that for all polynomials of degree at most  $2q$  we have  $Q_m[p] = \langle p, 1 \rangle_{\rho, m}$ .

Using

$$Q_m[a, b]_{\rho} := Q_m[\langle a, b \rangle_H]$$

instead of the scalar products  $\langle a, b \rangle_{\rho}$  we employ the following discrete **quadrature formulation**:

For given  $F \in \mathcal{V}^{\tau}$  and  $x_0 \in H$ , find  $U^{\tau} \in \mathcal{V}^{\tau}$ , such that for all  $\Phi \in \mathcal{V}^{\tau}$  and  $m \in \{1, 2, \dots, M\}$  it holds

$$Q_m[(\partial_t M_0 + M_1 + A)U^{\tau}, \Phi]_{\rho} + \langle M_0 \llbracket U^{\tau} \rrbracket_{m-1}^{x_0}, \Phi_{m-1}^+ \rangle_H = Q_m[F, \Phi]_{\rho}. \quad (3.1)$$

Here, we denote by

$$\llbracket U^{\tau} \rrbracket_{m-1}^{x_0} := \begin{cases} U^{\tau}(t_{m-1}+) - U(t_{m-1}-), & m \in \{2, \dots, M\} \\ U^{\tau}(t_0+) - x_0, & m = 1, \end{cases}$$

the jump at  $t_{m-1}$  and by  $\Phi_{m-1}^+ := \Phi(t_{m-1}+)$ .

**REMARK 3.1** We shall briefly comment on the derivation of (3.1). We start with the formulation of an initial value problem given in Remark 2.3 (f), i.e., with the equation

$$(\partial_t M_0 + M_1 + A)U = F + \delta_0 M_0 x_0.$$

We shall do so for the example case  $M = 1$ , i.e., only one interval, first. Testing the latter equation with  $\Phi \in \mathcal{V}^{\tau}$  we obtain

$$\langle (\partial_t M_0 + M_1 + A)U, \Phi \rangle_{\rho} = \langle F, \Phi \rangle_{\rho} + \langle M_0 x_0, \Phi(0+) \rangle_H.$$

By a standard penalization technique we impose the initial condition weakly in the following way

$$\langle (\partial_t M_0 + M_1 + A)U, \Phi \rangle_{\rho} + \langle M_0 U(0+) - M_0 x_0, \Phi(0+) \rangle_H = \langle F, \Phi \rangle_{\rho}.$$

Using the quadrature formula instead of the inner products, we derive (3.1) for  $m = M = 1$ .

For the general case, we repeat this argument and take  $U(t_{m-1}+)$  as the new initial value and hence obtain (3.1) for each interval.

**PROPOSITION 3.2** Let  $F \in \mathcal{V}^{\tau}$ ,  $x_0 \in H$ . Then there exists a unique solution of (3.1).

*Proof.* Let  $m \in \{1, \dots, M\}$  and recall that  $\mathcal{P}_q(I_m; H)$  is a Hilbert space with the aforementioned scalar product. We note that

$$\partial_t : \mathcal{P}_q(I_m; H) \rightarrow \mathcal{P}_q(I_m; H) : p \mapsto p'$$



and

$$\delta_{m-1} : \mathcal{P}_q(I_m; H) \rightarrow \mathbb{R} : p \mapsto p(t_{m-1}+)$$

are bounded linear operators. Consequently, the mapping

$$\mathcal{P}_q(I_m; H) \rightarrow \mathbb{R} : p \mapsto \langle x, \delta_{m-1} p \rangle_H$$

is linear and bounded for each  $x \in H$  and thus, by the Riesz representation theorem, there is a unique  $\Psi(x) \in \mathcal{P}_q(I_m; H)$  such that

$$\langle \Psi(x), p \rangle_{\rho, m} = \langle x, \delta_{m-1} p \rangle_H.$$

Moreover, the mapping  $\Psi : H \rightarrow \mathcal{P}_q(I_m; H)$  is linear and bounded, since

$$|\Psi(x)|_{\rho, m}^2 = \langle \Psi(x), \Psi(x) \rangle_{\rho, m} = \langle x, \delta_{m-1} \Psi(x) \rangle_H \leq |x|_H \|\delta_{m-1}\| |\Psi(x)|_{\rho, m} \quad (x \in H).$$

We now prove that for each  $g \in \mathcal{P}_q(I_m; H)$  there is a unique  $u \in \mathcal{P}_q(I_m; D(A))$  such that

$$(\partial_t M_0 + M_1 + A)u + \Psi M_0 \delta_{m-1} u = g. \quad (3.2)$$

For doing so, we first compute using integration by parts

$$\begin{aligned} \langle \partial_t M_0 v, v \rangle_{\rho, m} &= \frac{1}{2} \langle \partial_t M_0 v, v \rangle_{\rho, m} + \frac{1}{2} \int_{t_{m-1}}^{t_m} \langle v(t), M_0 v'(t) \rangle_H \exp(-2\rho(t - t_{m-1})) dt \\ &= \frac{1}{2} \langle \partial_t M_0 v, v \rangle_{\rho, m} - \frac{1}{2} \int_{t_{m-1}}^{t_m} \langle M_0 v'(t), v(t) \rangle_H \exp(-2\rho(t - t_{m-1})) dt \\ &\quad + \rho \int_{t_{m-1}}^{t_m} \langle M_0 v(t), v(t) \rangle_H \exp(-2\rho(t - t_{m-1})) dt + \frac{1}{2} \langle M_0 v(t_m), v(t_m) \rangle_H \exp(-2\rho\tau_m) \\ &\quad - \frac{1}{2} \langle M_0 v(t_{m-1}), v(t_{m-1}) \rangle_H \\ &\geq \rho \langle M_0 v, v \rangle_{\rho, m} - \frac{1}{2} \langle \Psi M_0 \delta_{m-1} v, v \rangle_{\rho, m} \end{aligned}$$

for each  $v \in \mathcal{P}_q(I_m; H)$ . Next, from  $A^* = -A$  it follows  $\langle Ax, x \rangle_H = 0$  for each  $x \in D(A)$ . Therefore, for all  $u \in \mathcal{P}_q(I_m; D(A))$ , we get

$$\begin{aligned} \langle (\partial_t M_0 + M_1 + A)u + \Psi M_0 \delta_{m-1} u, u \rangle_{\rho, m} &= \langle \partial_t M_0 u, u \rangle_{\rho, m} + \langle M_1 u, u \rangle_{\rho, m} + \langle \Psi M_0 \delta_{m-1} u, u \rangle_{\rho, m} \quad (3.3) \\ &\geq \langle (\rho M_0 + M_1)u, u \rangle_{\rho, m} + \frac{1}{2} \langle \Psi M_0 \delta_{m-1} u, u \rangle_{\rho, m} \\ &\geq \gamma \langle u, u \rangle_{\rho, m}, \end{aligned}$$

where we have used

$$\langle \Psi M_0 \delta_{m-1} u, u \rangle_{\rho, m} = \langle M_0 u(t_{m-1}+), u(t_{m-1}+) \rangle_H \geq 0.$$

In particular, both  $B := (\partial_t M_0 + M_1) + \Psi M_0 \delta_{m-1}$  and  $B + A$  are strictly positive definite. Moreover, since  $B$  is bounded,  $B^*$  is strictly positive definite, as well. Hence, from

$$(B + A)^* = B^* - A$$

we read off that  $(B + A)^*$  is strictly positive definite as well. Thus, for each  $g \in \mathcal{P}_q(I_m; H)$  there is a unique  $u \in \mathcal{P}_q(I_m; D(A)) = D(A + B)$  such that

$$(\partial_t M_0 + M_1 + A)u + \Psi M_0 \delta_{m-1} u = g. \quad (3.4)$$

Thus, we are in the position to define a solution for (3.1) by induction on  $m$ . For this, we put  $U(t_0-) := x_0$ . Next, assume we have solved (3.1) for  $U^\tau$  on  $I_{m-1}$  for some  $m \in \{1, \dots, M\}$  ( $I_0 := \{t_0\}$  and the equation is void). Then, let  $u \in \mathcal{P}_q(I_m; D(A))$  be such that (3.4) holds for  $g = F|_{I_m} - \Psi M_0 U^\tau(t_{m-1}-)$ . We put  $U^\tau|_{I_m} := u$ . The thus defined function  $U^\tau$  solves (3.1): we observe

$$\begin{aligned} & \langle (\partial_t M_0 + M_1 + A)U^\tau, \Phi \rangle_{\rho, m} + \langle \Psi M_0 \delta_{m-1} U^\tau, \Phi \rangle_{\rho, m} \\ &= \langle F - \Psi M_0 U^\tau(t_{m-1}-), \Phi \rangle_{\rho, m} = \langle F, \Phi \rangle_{\rho, m} + \langle \Psi M_0 U^\tau(t_{m-1}-), \Phi \rangle_{\rho, m}, \end{aligned}$$

by definition for all  $\Phi \in \mathcal{V}^\tau$  and  $m \in \{1, \dots, M\}$ . The latter is the same as saying

$$\begin{aligned} & \langle (\partial_t M_0 + M_1 + A)U^\tau, \Phi \rangle_{\rho, m} + \langle M_0 U^\tau(t_{m-1}+), \Phi(t_{m-1}+) \rangle_H \\ &= \langle F, \Phi \rangle_{\rho, m} + \langle M_0 U^\tau(t_{m-1}-), \Phi(t_{m-1}+) \rangle_H. \end{aligned}$$

But, since the quadrature is exact for polynomials up to degree  $2q$ , the latter equation in turn is equivalent to

$$Q_m \left[ (\partial_t M_0 + M_1 + A)U^\tau, \Phi \right]_\rho + \langle M_0 \llbracket U^\tau \rrbracket_{m-1}^{x_0}, \Phi_{m-1}^+ \rangle_H = Q_m [F, \Phi]_\rho,$$

which yields existence of  $U^\tau$ . Uniqueness follows from the uniqueness of  $u$  satisfying (3.4).  $\square$

**REMARK 3.3 (Dissipation of energy)** To use dG methods for time stepping is known to be slightly dissipative in the following sense. Let us consider the weak formulation of (1.1) with  $U$  as test function and time integration over  $(-\infty, T)$ , i.e.,

$$\langle (\partial_t M_0 + M_1 + A)U, U \rangle_{\rho, (-\infty, T)} = \langle F, U \rangle_{\rho, (-\infty, T)}.$$

Assuming a nonzero value of  $U(0)$  and  $F(t) = 0$  for  $t \geq 0$ , we obtain for each  $t > 0$  the following conservation law of the energy:

$$e^{-2\rho t} \langle M_0 U(t), U(t) \rangle_H + 2 \langle (\rho M_0 + M_1)U, U \rangle_{\rho, (0, t)} = \langle M_0 U(0), U(0) \rangle_H.$$

For our dG method we obtain in the discrete points  $t_i > 0$  the conservation law of the discrete energy

$$\begin{aligned} e^{-2\rho t_i} \langle M_0 U_i^\tau, U_i^\tau \rangle_H + 2 \langle (\rho M_0 + M_1) U^\tau, U^\tau \rangle_{\rho, (0, t_i)} \\ + \sum_{m=1}^i e^{-2\rho t_{m-1}} \langle M_0 [[U^\tau]]_{m-1}^{x_0}, [[U^\tau]]_{m-1}^{x_0} \rangle_H = \langle M_0 U_0^\tau, U_0^\tau \rangle_H. \end{aligned}$$

From the two conservation laws we observe a reduction of the discrete energy compared to the original energy over time due to the jumps. There are time-stepping methods, using different ansatz and test spaces, without such a dissipation. We will consider them in a future publication. Here we analyse the dissipative dG method because its ansatz and trial spaces are the same. It therefore fits better in the theoretical framework of (Picard, 2009).

**REMARK 3.4** In the proof of Proposition 3.2 the importance of the introduction of the exponential weight—also for the finite time regime—becomes apparent. Indeed, the exponential weight serves to ensure strict positive definiteness of the operator appearing in (3.2) in space time, see also (3.3). If on the other hand one uses an unweighted  $L^2$  space in Proposition 3.2 the same result holds due to the equivalence of the corresponding norms on finite time intervals.

### 3.1 On some a priori error estimates in time

After having proved the unique solvability of (3.1), we address the error estimates in the following. In our analysis we will use the discretized norms

$$\|v\|_{Q, \rho, m}^2 := Q_m[v, v]_\rho \quad \text{and} \quad \|v\|_{Q, \rho}^2 := \sum_{m=1}^M Q_m[v, v]_\rho e^{-2\rho t_{m-1}}$$

as approximations of  $\|v\|_{\rho, m}^2 := \int_{t_m} |v(t)|_H^2 \exp(-2\rho(t - t_{m-1})) dt$  and  $|v|_\rho^2$ . Note that for  $v \in \mathcal{V}^\tau$  the approximation is exact.

Let us start by defining an interpolation operator into  $\mathcal{V}^\tau$  and define by  $\varphi_{m,i}$  with  $i \in \{0, \dots, q\}$  the associated Lagrange basis functions to the nodes  $t_{m,i}$ . Then we obtain for a function  $v \in C([0, T], H)$  by

$$(Pv)(0) = v(0), \quad (Pv)|_{I_m}(t) = \sum_{i=0}^q v(t_{m,i}) \varphi_{m,i}(t), \quad m \in \{1, \dots, M\}, \quad (3.5)$$

an interpolation operator in time.

In the analysis to follow, we will consider the problem (2.1). In particular, we emphasize that we assume that the hypotheses of Theorem 2.2 are in effect. Furthermore, we fix a right-hand side  $F \in H_\rho^2(\mathbb{R}; H)$ . Thus, by Theorem 2.2 (and Remark 2.3(c)) there exists a unique solution

$$U \in H_\rho^2(\mathbb{R}; H) \text{ with } (\partial_t M_0 + M_1 + A)U = F. \quad (3.6)$$

Also, by Remark 2.3(c), we obtain  $F \in C_\rho(\mathbb{R}; H)$  and  $U \in C_\rho(\mathbb{R}; D(A)) \cap C_\rho^1(\mathbb{R}; H)$ . Moreover, we set  $U^\tau \in \mathcal{V}^\tau$  to satisfy (3.1) for the right-hand side  $PF \in \mathcal{V}^\tau$  and  $x_0 := U(0+)$ . We consider the following splitting:

$$U^\tau - U = \xi - \eta, \quad \text{where } \xi = U^\tau - PU \in \mathcal{V}^\tau \quad \text{and} \quad \eta = U - PU.$$

Note that due to the better regularity of  $U$  mentioned above, the equation holds pointwise, that is, for every  $t \in [0, T]$  we have that

$$(\partial_t M_0 + M_1 + A)U(t) = F(t)$$

and thus,

$$\langle (\partial_t M_0 + M_1 + A)U(t), \Phi(t) \rangle_H = \langle F(t), \Phi(t) \rangle_H$$

for each  $\Phi \in \mathcal{V}^\tau$  and every  $t \in [0, T]$ , which gives

$$Q_m [(\partial_t M_0 + M_1 + A)U, \Phi]_\rho + \langle M_0[[U]]_{m-1}^{x_0}, \Phi_{m-1}^+ \rangle_H = Q_m [F, \Phi]_\rho = Q_m [PF, \Phi]_\rho,$$

where we have used  $M_0[[U]]_{m-1}^{x_0} = M_0[[U]]_{m-1}^{U(0+)} = 0$ , due to the continuity of  $U$  and  $Q_m [F, \Phi]_\rho = Q_m [PF, \Phi]_\rho$ , since  $PF$  interpolates at the Gauß–Radau points used in the quadrature. Hence,  $U$  (formally) solves the same semidiscretized problem as  $U^\tau$ . Thus, we obtain with  $\chi \in \mathcal{V}^\tau$  as test function the **error equation**

$$\begin{aligned} & Q_m [(\partial_t M_0 + M_1 + A)\xi, \chi]_\rho + \langle M_0[[\xi]]_{m-1}^0, \chi_{m-1}^+ \rangle_H \\ & = Q_m [(\partial_t M_0 + M_1 + A)\eta, \chi]_\rho + \langle M_0[[\eta]]_{m-1}^0, \chi_{m-1}^+ \rangle_H. \end{aligned} \tag{3.7}$$

For the special case  $\chi = \xi$  (use  $A = -A^*$ ) we obtain

$$\begin{aligned} E_d^m & := Q_m [(\partial_t M_0 + M_1)\xi, \xi]_\rho + \langle M_0[[\xi]]_{m-1}^0, \xi_{m-1}^+ \rangle_H \\ & = Q_m [(\partial_t M_0 + M_1 + A)\eta, \xi]_\rho + \langle M_0[[\eta]]_{m-1}^0, \xi_{m-1}^+ \rangle_H =: E_i^m \end{aligned} \tag{3.8}$$

for all  $m \in \{1, \dots, M\}$ , where the subscripts d and i should remind of discretization and interpolation, respectively.

LEMMA 3.5 For all  $m \in \{1, \dots, M\}$ , we have

$$E_d^m \geq \gamma \|\xi\|_{Q,\rho,m}^2 + \frac{1}{2} \left[ \langle M_0 \xi_m^-, \xi_m^- \rangle_H e^{-2\rho\tau_m} - \langle M_0 \xi_{m-1}^-, \xi_{m-1}^- \rangle_H + \langle M_0 [[\xi]]_{m-1}^0, [[\xi]]_{m-1}^0 \rangle_H \right],$$

where  $\xi_m^- := \xi(t_m^-)$  and  $\xi_0^- := 0$ .

*Proof.* Let  $m \in \{1, \dots, M\}$ . Since  $\xi$  is a (piecewise) polynomial of order  $q$  in time, we obtain

$$\begin{aligned} Q_m [\partial_t M_0 \xi, \xi]_\rho &= \langle \partial_t M_0 \xi, \xi \rangle_{\rho, m} \\ &= \frac{1}{2} \int_{I_m} e^{-2\rho(t-t_{m-1})} \partial_t \langle M_0 \xi, \xi \rangle_H dt \\ &= \frac{1}{2} \left[ \langle M_0 \xi_m^-, \xi_m^- \rangle_H e^{-2\rho\tau_m} - \langle M_0 \xi_{m-1}^+, \xi_{m-1}^+ \rangle_H \right] + \rho \langle M_0 \xi, \xi \rangle_{\rho, m}. \end{aligned}$$

Further, we compute

$$\langle M_0 \llbracket \xi \rrbracket_{m-1}^0, \xi_{m-1}^+ \rangle_H = \frac{1}{2} \left[ \langle M_0 \xi_{m-1}^+, \xi_{m-1}^+ \rangle_H - \langle M_0 \xi_{m-1}^-, \xi_{m-1}^- \rangle_H + \langle M_0 \llbracket \xi \rrbracket_{m-1}^0, \llbracket \xi \rrbracket_{m-1}^0 \rangle_H \right].$$

Therefore, we have

$$\begin{aligned} E_d^m &= Q_m [(\partial_t M_0 + M_1) \xi, \xi]_\rho + \langle M_0 \llbracket \xi \rrbracket_{m-1}^0, \xi_{m-1}^+ \rangle_H \\ &= \frac{1}{2} \left[ \langle M_0 \xi_m^-, \xi_m^- \rangle_H e^{-2\rho\tau_m} - \langle M_0 \xi_{m-1}^-, \xi_{m-1}^- \rangle_H + \langle M_0 \llbracket \xi \rrbracket_{m-1}^0, \llbracket \xi \rrbracket_{m-1}^0 \rangle_H \right] + \langle (\rho M_0 + M_1) \xi, \xi \rangle_{\rho, m}. \end{aligned}$$

Together with

$$\langle (\rho M_0 + M_1) \xi, \xi \rangle_{\rho, m} \geq \gamma \|\xi\|_{\rho, m}^2 = \gamma \|\xi\|_{Q, \rho, m}^2$$

the lemma is proved.  $\square$

In order to analyse  $E_i^m$  we introduce another interpolation operator that enables us to estimate the time derivative of the interpolation error with a higher order. This operator utilizes  $t_{m,-1} := t_{m-1}$  in addition to  $t_{m,i}$ ,  $i \in \{0, \dots, q\}$  as interpolation points. Denoting the associated Lagrange basis functions by  $\psi_{m,i}$ ,  $i \in \{-1, 0, \dots, q\}$ , this interpolation operator is given by

$$(\widehat{P}v)|_{I_m}(t) := \sum_{i=-1}^q v(t_{m,i}) \psi_{m,i}(t), \quad m \in \{1, \dots, M\}. \quad (3.9)$$

Note the  $\widehat{P}$  maps to functions that are continuous in time (recall that  $t_{m,q} = t_m$ ) while the image of  $P$  is allowed to be discontinuous at the time mesh points.

LEMMA 3.6 For  $m \in \{1, \dots, M\}$  and  $\psi \in \mathcal{V}^\tau$ , we have

$$Q_m [\partial_t M_0 \eta, \psi]_\rho + \langle M_0 \llbracket \eta \rrbracket_{m-1}^0, \psi_{m-1}^+ \rangle_H = Q_m [\partial_t M_0 (U - \widehat{P}U), \psi]_\rho + R(U, \psi),$$

where

$$|R(U, \psi)| \leq C\alpha\tau_m |M_0 \eta_{m-1}^+|_H^2 + \beta \|\psi\|_{Q, \rho, m}^2$$

for all  $\alpha, \beta > 0$  satisfying  $\alpha\beta = 1/4$  and with  $C \geq 0$  depending on  $T$  (the finite time horizon) and  $\rho$  only.

*Proof.* With  $U$  being continuous in time, we only have to consider the discrete part. Using the exactness of the quadrature rule for polynomials of degree  $2q$  and integration by parts, we obtain for  $m \in \{1, \dots, M\}$

$$\begin{aligned}
 & Q_m [\partial_t M_0 P U, \psi]_\rho + \underbrace{\langle M_0 \llbracket P U \rrbracket_{m-1}^{x_0}, \psi_{m-1}^+ \rangle_H}_{=:a} \\
 &= \langle \partial_t M_0 P U, \psi \rangle_{\rho, m} + a \\
 &= -\langle M_0 P U, \partial_t \psi \rangle_{\rho, m} + 2\rho \langle M_0 P U, \psi \rangle_{\rho, m} + \underbrace{\langle e^{-2\rho(t-t_{m-1})} M_0 P U, \psi \rangle_H \Big|_{t_{m-1}}^{t_m}}_{=:b} + a \\
 &= - \underbrace{Q_m [M_0 P U, \partial_t \psi]_\rho}_{=Q_m [M_0 \widehat{P} U, \partial_t \psi]_\rho = \langle M_0 \widehat{P} U, \partial_t \psi \rangle_{\rho, m}} + 2\rho \langle M_0 P U, \psi \rangle_{\rho, m} + a + b \\
 &= \langle \partial_t M_0 \widehat{P} U, \psi \rangle_{\rho, m} + 2\rho (\langle M_0 P U, \psi \rangle_{\rho, m} - \langle M_0 \widehat{P} U, \psi \rangle_{\rho, m}) + a + b - \underbrace{\langle e^{-2\rho(t-t_{m-1})} M_0 \widehat{P} U, \psi \rangle_H \Big|_{t_{m-1}}^{t_m}}_{=:c}.
 \end{aligned}$$

Using  $(P U)_{m-1}^- = (\widehat{P} U)_{m-1}^+$  ( $m \geq 2$ ),  $(\widehat{P} U)_0^+ = U(0+) = x_0$  and  $(P U)_m^- = (\widehat{P} U)_m^-$  ( $m \geq 1$ ), we have

$$a + b - c = 0.$$

Furthermore, it holds

$$\begin{aligned}
 2\rho (\langle M_0 P U, \psi \rangle_{\rho, m} - \langle M_0 \widehat{P} U, \psi \rangle_{\rho, m}) &= 2\rho \langle M_0 ((P - \widehat{P}) U)(t_{m-1}^+) \chi, \psi \rangle_{\rho, m} \\
 &= 2\rho \langle M_0 (P U - U)(t_{m-1}^+) \chi, \psi \rangle_{\rho, m} =: R(U, \psi),
 \end{aligned}$$

where  $\chi \in \mathcal{P}_{q+1}(I_m)$  with  $\chi(t_{m-1}) = 1$  and  $\chi(t_{m,i}) = 0$ ,  $i \in \{0, \dots, q\}$ . By (Trostorff & Waurick, 2016, Corollary 1.4) for  $K = T$  (note that  $0 < \tau_m = |I_m| \leq T$ ), we obtain  $\|\chi\|_{\rho, m}^2 \leq C\tau_m$  for some  $C \geq 0$ . Thus, we get

$$\begin{aligned}
 |R(U, \psi)| &\leq 2\rho \|M_0 (P U - U)(t_{m-1}^+)\|_H \|\chi\|_{\rho, m} \|\psi\|_{\rho, m} \\
 &\leq C^2 (2\rho)^2 \alpha \tau_m |M_0 (P U - U)(t_{m-1}^+)|^2 + \beta \|\psi\|_{Q, \rho, m}^2,
 \end{aligned}$$

where  $\alpha\beta = 1/4$ . Combining above transformations we are done.  $\square$

LEMMA 3.7 For all  $m \in \{1, \dots, M\}$ , we have for all  $\psi \in \mathcal{V}^\tau$

$$Q_m [M_1 \eta, \psi]_\rho = 0 = Q_m [A \eta, \psi]_\rho.$$

*Proof.* These equalities follow from the fact that  $\eta(t_{m,i}) = P U(t_{m,i}) - U(t_{m,i}) = 0$  for each  $i \in \{0, \dots, q\}$  and  $M_1, A$  are purely spatial operators.  $\square$

Combining the previous lemmas gives the first result.

**THEOREM 3.8** There exists a  $C \geq 0$  depending on  $T$ ,  $\rho$  and  $\gamma$ , only, such that

$$\langle M_0 \xi_M^-, \xi_M^- \rangle_H + e^{2\rho T} \|\xi\|_{Q,\rho}^2 \leq C e^{2\rho T} \left( \|\partial_t M_0(U - \widehat{P}U)\|_{Q,\rho}^2 + T \max_{1 \leq m \leq M} \left\{ |M_0 \eta_{m-1}^+|_H^2 e^{-2\rho t_{m-1}} \right\} \right) =: g(U).$$

*Proof.* From Lemma 3.5 we obtain upon summation with the weights  $e^{-2\rho t_{m-1}}$  for  $m \in \{1, \dots, M\}$  due to cancellation

$$\sum_{m=1}^M e^{-2\rho t_{m-1}} |E_d^m| \geq \frac{1}{2} \langle M_0 \xi_M^-, \xi_M^- \rangle_H e^{-2\rho T} + \gamma \|\xi\|_{Q,\rho}^2,$$

by  $\xi_0^- = 0$  and neglecting the positive jump contributions. Combining Lemmas 3.6 and 3.7 for  $\psi = \xi$  we have for some  $C \geq 1$  depending on  $T$  and  $\rho$  only

$$|E_i^m| \leq C\alpha \left( \|\partial_t M_0(U - \widehat{P}U)\|_{Q,\rho,m}^2 + \tau_m |M_0 \eta_{m-1}^+|_H^2 \right) + 2\beta \|\xi\|_{Q,\rho,m}^2,$$

which yields upon the same weighted summation

$$\sum_{m=1}^M e^{-2\rho t_{m-1}} |E_i^m| \leq C\alpha \left( \|\partial_t M_0(U - \widehat{P}U)\|_{Q,\rho}^2 + T \max_{1 \leq m \leq M} \left\{ |M_0 \eta_{m-1}^+|_H^2 e^{-2\rho t_{m-1}} \right\} \right) + 2\beta \|\xi\|_{Q,\rho}^2.$$

Thus for  $\beta < \gamma/2$  the result is proved upon the equality  $E_i^m = E_d^m$ .  $\square$

**REMARK 3.9** Let  $m \in \{1, \dots, M\}$ . Note that the estimate in Theorem 3.8 remains valid, if one replaces  $T$  by  $t_m$ ,  $\xi_M^-$  by  $\xi_m^-$ ,  $\|\xi\|_{Q,\rho}$  by  $\|\xi \chi_{\bar{I}}\|_{Q,\rho}$  and  $\|\partial_t M_0(U - \widehat{P}U)\|_{Q,\rho}$  by  $\|\partial_t M_0(U - \widehat{P}U) \chi_{\bar{I}}\|_{Q,\rho}$ , with  $\chi_{\bar{I}}$  being the characteristic function of  $\bar{I} = \bigcup_{k=1}^m I_k$ .

In the following, we want to improve Theorem 3.8. In order to do so, we will need the following technical lemmas. The first is an adaptation of the study by (Akrivis & Makridakis, 2004, Lemma 2.1). For the upcoming result and the corresponding proof, we recall for polynomials  $a, b \in \mathcal{P}_q(0, 1; H)$

$$\langle a, b \rangle_\rho = \int_0^1 \langle a(t), b(t) \rangle_H e^{-2\rho t} dt$$

and the corresponding integration by parts formula

$$\langle a', b \rangle_\rho = -\langle a, b' \rangle_\rho + 2\rho \langle a, b \rangle_\rho + e^{-2\rho t} \langle a, b \rangle_H \Big|_0^1. \quad (3.10)$$

**LEMMA 3.10** Let  $t_i, w_i, i \in \{0, \dots, q\}$  be the points and weights of the right-sided Gauß–Radau quadrature rule of order  $q$  on  $(0, 1]$  with weighting function  $t \mapsto e^{-2\rho t}$ .

Let  $p \in \mathcal{P}_q(0, 1; H)$  and  $\tilde{p}$  the Lagrange interpolant w.r.t.  $(t_i)_{i \in \{0, \dots, q\}}$  of  $\varphi: (0, 1] \ni t \mapsto p(t)/t$ . Then

$$\langle p', \tilde{p} \rangle_\rho + \langle p(0), \tilde{p}(0) \rangle_H \geq \frac{1}{2} \left( |p(1)|_H^2 e^{-2\rho} + \langle \tilde{p}, \tilde{p} \rangle_\rho \right) + \rho \langle p, p \rangle_\rho.$$

*Proof.* We can basically follow the proof of the study by (Akrivis & Makridakis, 2004, Lemma 2.1) step by step. The only difference lies with the weighted scalar product and (3.10). We will sketch the proof.

As in the study by (Akrivis & Makridakis, 2004) define  $v \in \mathcal{P}_{q-1}((0, 1); H)$  by  $v(t) = (p(t) - p(0))/t$  and  $\Lambda \in \mathcal{P}_q[0, 1]$  by  $\Lambda(t_i) = 1/t_i$ ,  $i \in \{0, \dots, q\}$ . Then

$$\langle p', \tilde{p} \rangle_\rho = \langle v, v \rangle_\rho + \langle mv', v \rangle_\rho + \langle v, p(0)\Lambda \rangle_\rho + \langle v', p(0)m\Lambda \rangle_\rho,$$

where we denote by  $m$  the multiplication with the argument, that is,  $(mf)(t) := tf(t)$ . With (3.10) we obtain for the second term

$$\langle mv', v \rangle_\rho = \langle v', mv \rangle_\rho = \frac{1}{2} \left( e^{-2\rho} |v(1)|_H^2 + 2\rho \langle mv, v \rangle_\rho - \langle v, v \rangle_\rho \right).$$

From  $mv' \Lambda \in \mathcal{P}_{2q-1}$  and  $m\Lambda' \Lambda \in \mathcal{P}_{2q}$  together with the exactness of the quadrature rule it follows that

$$\begin{aligned} \langle p', \tilde{p} \rangle_\rho + \langle p(0), \tilde{p}(0) \rangle_H &= \frac{1}{2} \left( e^{-2\rho} |p(1)|_H^2 + \langle v, v \rangle_\rho + 2\rho \langle mv, v \rangle_\rho \right) \\ &\quad + \langle v, p(0)\Lambda \rangle_\rho + 2\rho \langle v, p(0) \rangle_\rho + |p(0)|_H^2 \left( \Lambda(0) - \frac{e^{-2\rho}}{2} \right). \end{aligned} \quad (3.11)$$

Next  $\langle \Lambda', m\Lambda \rangle_\rho = 2\rho \langle \Lambda, 1 \rangle_\rho + e^{-2\rho} - \Lambda(0)$  and (3.10) yield

$$\Lambda(0) = 2\rho \langle \Lambda, 1 \rangle_\rho - \rho \langle \Lambda, m\Lambda \rangle_\rho + \frac{1}{2} \left( e^{-2\rho} + \langle \Lambda, \Lambda \rangle_\rho \right),$$

which can be substituted into (3.11). With

$$\frac{1}{2} \langle v, v \rangle_\rho + \langle v, p(0)\Lambda \rangle_\rho + \frac{1}{2} |p(0)|_H^2 \langle \Lambda, \Lambda \rangle_\rho = \frac{1}{2} \langle v + p(0)\Lambda, v + p(0)\Lambda \rangle_\rho = \frac{1}{2} \langle \tilde{p}, \tilde{p} \rangle_\rho$$

and

$$\langle mv, v \rangle_\rho + 2\langle v, p(0) \rangle_\rho + |p(0)|^2 \langle \Lambda, 1 \rangle_\rho = \langle p, \tilde{p} \rangle_\rho$$

it follows

$$\langle p', \tilde{p} \rangle_\rho + \langle p(0), \tilde{p}(0) \rangle_H = \frac{1}{2} \left( e^{-2\rho} |p(1)|_H^2 + \langle \tilde{p}, \tilde{p} \rangle_\rho \right) + \rho \left( \langle p, \tilde{p} \rangle_\rho + |p(0)|_H^2 \langle \Lambda, 1 - m\Lambda \rangle_\rho \right).$$

Using  $\langle \Lambda, 1 - m\Lambda \rangle_\rho \geq 0$ , which we provide in Lemma 3.11, the exactness of the quadrature rule and  $t_i^{-1} > 1$  for  $\langle p, \tilde{p} \rangle_\rho$  the result is proved.  $\square$

**LEMMA 3.11** Let  $\Lambda \in \mathcal{P}_q[0, 1]$  such that  $\Lambda(t_i) = \frac{1}{t_i}$  for  $i \in \{0, \dots, q\}$ , where  $t_i$  is chosen as in Lemma 3.10. Then

$$\langle \Lambda, 1 - m\Lambda \rangle_\rho \geq 0,$$

where  $(m\Lambda)(t) := t\Lambda(t)$ .



*Proof.* We rewrite the scalar product as a quadrature error:

$$\langle \Lambda, 1 - m\Lambda \rangle_\rho = \sum_{i=0}^q w_i \frac{1}{t_i} - \int_0^1 e^{-2\rho t} t \Lambda^2(t) dt = Q[f] - I[f]$$

for  $f$  given by  $f(t) = t\Lambda^2(t)$ , where  $Q[a] = \sum_{i=0}^q w_i a(t_i)$  and  $I[a] = \int_0^1 e^{-2\rho t} a(t) dt$  for a suitable function  $a$ . There exists a constant  $\alpha \in \mathbb{R}$  and a polynomial  $w_0 \in \mathcal{P}_{q-1}[0, 1]$ , such that  $\Lambda(t) = \alpha t^q + w_0(t)$ , which implies  $f(t) = \alpha^2 t^{2q+1} + w_1(t)$ , where  $w_1 \in \mathcal{P}_{2q}[0, 1]$ . Thus, setting  $a(t) = t^{2q+1}$ , we have that

$$\langle \Lambda, 1 - m\Lambda \rangle_\rho = \alpha^2 (Q[a] - I[a]),$$

due to the exactness of the quadrature rule for polynomials of degree  $2q$ .

Let  $\Pi w \in \mathcal{P}_{2q}[0, 1]$  be an Hermite interpolant of a given function  $w$  satisfying

$$\begin{aligned} \Pi w(t_i) &= w(t_i), \quad i \in \{0, \dots, q\}, \\ (\Pi w)'(t_i) &= w'(t_i), \quad i \in \{0, \dots, q-1\}. \end{aligned}$$

Then it follows

$$Q[a] = \sum_{i=0}^q w_i t_i^{2q+1} = \sum_{i=0}^q w_i (\Pi a)(t_i^{2q+1}) = Q[\Pi a] = I[\Pi a].$$

Using that for each  $t \in [0, 1]$  there is  $\zeta \in (0, 1)$  such that

$$(\Pi a)(t) - a(t) = -\frac{a^{(2q+1)}(\zeta)}{(2q+1)!} (t-1) \prod_{i=0}^{q-1} (t-t_i)^2 = (1-t) \prod_{i=0}^{q-1} (t-t_i)^2,$$

see, for instance, the study by (Stoer *et al.*, 2002, Section 2.1.5), we infer that

$$\langle \Lambda, 1 - m\Lambda \rangle_\rho = \alpha^2 I[\Pi a - a] = \alpha^2 \int_0^1 e^{-2\rho t} \left( \prod_{i=0}^{q-1} (t-t_i)^2 \right) (1-t) dt \geq 0.$$

□

Now we are able to improve Theorem 3.8 following the studies by (Akrivis & Makridakis, 2004, Corollary 2.1) and (Vlasak & Roos, 2014).

**THEOREM 3.12** There exists  $C \geq 0$  depending on  $T$ ,  $q$ ,  $\|M_0\|$ ,  $\|M_1\|$ ,  $\gamma$  and  $\rho$  such that

$$\sup_{t \in [0, T]} \langle M_0 \xi(t), \xi(t) \rangle_H \leq Cg(U),$$

with  $g(U)$  defined as in Theorem 3.8.

*Proof.* For the discrete error  $\xi = U^\tau - PU \in \mathcal{V}^\tau$  we define  $\varphi$  by

$$\varphi|_{I_m} = P \left( t \mapsto \frac{\tau_m}{t - t_{m-1}} \xi(t) \right) \quad (m \in \{1, \dots, M\}).$$

Then for all  $m \in \{1, \dots, M\}$  and  $i \in \{0, \dots, q\}$  we have

$$\langle M_0 \varphi(t_{m,i}), \varphi(t_{m,i}) \rangle_H = \frac{\tau_m^2}{(t_{m,i} - t_{m-1})^2} \langle M_0 \xi(t_{m,i}), \xi(t_{m,i}) \rangle_H \geq \langle M_0 \xi(t_{m,i}), \xi(t_{m,i}) \rangle_H$$

and by Lemma 3.10 (applied to  $p = \sqrt{M_0} \xi$  and  $\tilde{p} = \sqrt{M_0} \varphi$  rescaled on  $[0, 1]$ , which is valid due to the nonnegativity and self-adjointness of  $M_0$ )

$$Q_m [\partial_t M_0 \xi, 2\varphi]_\rho + \langle M_0 \xi_{m-1}^+, 2\varphi_{m-1}^+ \rangle_H \geq \frac{1}{\tau_m} Q_m [M_0 \varphi, \varphi]_\rho \geq \frac{1}{\tau_m} Q_m [M_0 \xi, \xi]_\rho.$$

By the equivalence of norms on  $\mathcal{P}_q([0, 1])$ , there exists  $K_1 \geq 0$  depending on  $q$  only, such that

$$\sup_{t \in [0, 1]} |p(t)| \leq K_1 \int_0^1 |p(t)| dt \quad (p \in \mathcal{P}_q([0, 1])).$$

Consequently, we obtain for all  $m \in \{1, \dots, M\}$

$$\sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H \leq \frac{K_1}{\tau_m} e^{2\rho\tau_m} Q_m [M_0 \xi, \xi]_\rho \leq \frac{K}{\tau_m} Q_m [M_0 \xi, \xi]_\rho$$

where  $K := K_1 e^{2\rho T} \geq \max_{m \in \{1, \dots, M\}} \{e^{2\rho\tau_m}\} K_1$ . Moreover, we have

$$\begin{aligned} Q_m [A\xi, 2\varphi]_\rho &= \frac{\tau_m}{2} \sum_{i=0}^q \omega_i^m \langle A\xi(t_{m,i}), 2\varphi(t_{m,i}) \rangle_H \\ &= \frac{\tau_m}{2} \sum_{i=0}^q \omega_i^m \frac{2\tau_m}{t_{m,i} - t_{m-1}} \langle A\xi(t_{m,i}), \xi(t_{m,i}) \rangle_H = 0 \end{aligned}$$

upon  $A = -A^*$ . Together, it follows for all  $m \in \{1, \dots, M\}$

$$\begin{aligned} \sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H &\leq K \left( Q_m [(\partial_t M_0 + M_1 + A)\xi, 2\varphi]_\rho + \langle M_0 \xi_{m-1}^+, 2\varphi_{m-1}^+ \rangle_H - Q_m [M_1 \xi, 2\varphi]_\rho \right) \\ &= K \left( Q_m [(\partial_t M_0 + M_1 + A)\xi, 2\varphi]_\rho + \langle M_0 \|\xi\|_{m-1}^0, 2\varphi_{m-1}^+ \rangle_H \right. \\ &\quad \left. + \langle M_0 \xi_{m-1}^-, 2\varphi_{m-1}^+ \rangle_H - Q_m [M_1 \xi, 2\varphi]_\rho \right). \end{aligned}$$

Using the error equation (3.7) with  $\chi = 2\varphi$  (recall  $\eta = U - PU$ ), we obtain

$$\begin{aligned} \sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H &\leq K \left( Q_m [(\partial_t M_0 + M_1 + A)\eta, 2\varphi]_\rho + \langle M_0 \llbracket \eta \rrbracket_{m-1}^0, 2\varphi_{m-1}^+ \rangle_H \right. \\ &\quad \left. + \langle M_0 \xi_{m-1}^-, 2\varphi_{m-1}^+ \rangle_H - Q_m [M_1 \xi, 2\varphi]_\rho \right). \end{aligned}$$

Using Lemma 3.6, Lemma 3.7 with  $\psi = 2\varphi$  and Theorem 3.8, we estimate further with some  $C_1 \geq 1$  depending on  $q, T$  and  $\rho$  such that

$$\begin{aligned} \sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H &\leq K \left( Q_m [\partial_t M_0 (U - \widehat{P}U), 2\varphi]_\rho + R(U, 2\varphi) \right. \\ &\quad \left. + \langle M_0 \xi_{m-1}^-, 2\varphi_{m-1}^+ \rangle_H - Q_m [M_1 \xi, 2\varphi]_\rho \right) \\ &\leq C_1 \alpha_1 \left( \|\partial_t M_0 (U - \widehat{P}U)\|_{Q, \rho, m}^2 + |M_0 (PU - U)(t_{m-1}^+)|^2 \right) \\ &\quad + C_1 \alpha_2 \langle M_0 \xi_{m-1}^-, \xi_{m-1}^- \rangle_H + C_1 \alpha_1 \|M_1\|^2 \|\xi\|_{Q, \rho, m}^2 \\ &\quad + 3\beta_1 Q_m [2\varphi, 2\varphi]_\rho + \beta_2 \langle M_0 \varphi_{m-1}^+, 2\varphi_{m-1}^+ \rangle_H, \end{aligned}$$

where  $\alpha_i \beta_i = \frac{1}{4}$ ,  $i \in \{1, 2\}$  and we used that

$$\langle M_0 u, v \rangle_H = \langle \sqrt{M_0} u, \sqrt{M_0} v \rangle_H \leq \langle M_0 u, u \rangle_H \langle M_0 v, v \rangle_H$$

for all  $u, v \in H$ , by the non-negativity and selfadjointness of  $M_0$ . Using Theorem 3.8 (and Remark 3.9), we, thus, get

$$\sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H \leq C(\alpha_1 + \alpha_2)g(U) + 12\beta_1 Q_m [\varphi, \varphi]_\rho + 4\beta_2 \langle M_0 \varphi_{m-1}^+, \varphi_{m-1}^+ \rangle_H \quad (3.12)$$

for some  $C \geq 1$  depending on  $q, T, \rho$  and  $\|M_1\|$ , where  $g(U)$  is defined in Theorem 3.8. Next, by (Trostorff & Waurick, 2016, Corollary 1.5), we find  $c > 0$  depending on  $\rho$  and  $T$  only such that

$$\frac{\tau_m}{t_{m,i} - t_{m-1}} \leq \frac{\tau_m}{t_{m,0} - t_{m-1}} \leq \frac{1}{c} \quad (m \in \{1, \dots, M\}).$$

Hence, for all  $m \in \{1, \dots, M\}$ ,

$$Q_m [\varphi, \varphi]_\rho \leq \frac{1}{c^2} \|\xi\|_{Q, \rho, m}^2 \quad \text{and} \quad \langle M_0 \varphi_{m-1}^+, 1\varphi_{m-1}^+ \rangle_H \leq \frac{1}{c^2} \sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H.$$

Next, we choose  $\beta_2 = \frac{c^2}{8}$ . Thus, appealing to (3.12), we obtain for all  $m \in \{1, \dots, M\}$

$$\begin{aligned} \frac{1}{2} \sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H &= \sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H - \frac{1}{2} \sup_{t \in I_m} \langle M_0 \xi(t), \xi(t) \rangle_H \\ &\leq C \left( \alpha_1 + \frac{2}{c^2} \right) g(U) + \frac{12}{c^2} \beta_1 \|\xi\|_{Q,\rho,m}^2, \end{aligned}$$

using Theorem 3.8 (i.e., Remark 3.9) again for the second term on the right-hand side and computing the supremum over  $m \in \{1, \dots, M\}$  in the latter inequality, we obtain the assertion.  $\square$

### 3.2 Estimating the interpolation error in time

In the previous section we showed that the discrete error is bounded in terms of the interpolation errors. We finalize the error estimates in time in this section focusing on the interpolation error. The aim and, thus, main theorem of this section is Theorem 3.16, where we estimate the difference between the exact solution  $U$  of (3.6) and the solution  $U^\tau$  of the quadrature formulation (3.1) with right-hand side  $PF$  and initial value  $U(0+)$ . We use the same notation as in the previous section. In addition, we set

$$\tau := \max\{\tau_m : m \in \{1, \dots, M\}\}.$$

Moreover, we shall further assume that the hypotheses of Theorem 2.2 are in effect.

LEMMA 3.13 There exists  $C \geq 0$  depending on  $q$  and  $T$  such that for all  $V \in H_\rho^{q+3}(\mathbb{R}; H)$

$$\|\partial_t(V - \widehat{P}V)\|_{Q,\rho} \leq C\tau^{q+1} \sup_{t \in [0,T]} |\partial_t^{p+2}V(t)|_H.$$

*Proof.* First we note that  $H_\rho^{q+3}(\mathbb{R}; H) \hookrightarrow C_\rho^{q+2}(\mathbb{R}; H)$  by the Sobolev-embedding theorem. By the definition of  $\|\cdot\|_{Q,\rho}$  we have that

$$\|\partial_t(V - \widehat{P}V)\|_{Q,\rho}^2 = \sum_{m=1}^M \frac{\tau_m}{2} \sum_{i=0}^q \omega_i^m |\partial_t(V - \widehat{P}V)(t_{m,i})|_H^2 e^{-2\rho t_{m,i}}.$$

Using the standard result from interpolation theory, see, for instance, the study by (Stoer *et al.*, 2002, Section 2.1.4)

$$\sup_{t \in I_m} |(v - \widehat{P}v)'(t)| \leq C\tau_m^{q+1} \sup_{t \in I_m} |v^{(q+2)}(t)|,$$

for all  $v \in W^{q+2,\infty}(0, T)$  we obtain

$$\begin{aligned} \|\partial_t(V - \widehat{P}V)\|_{Q,\rho}^2 &\leq C^2 \sum_{m=1}^M \frac{\tau_m}{2} \tau_m^{2(q+1)} \sum_{i=0}^q \omega_i^m \sup_{t \in I_m} |\partial_t^{p+2}V(t)|_H^2 e^{-2\rho t_{m,i}} \\ &\leq C^2 \tau^{2(q+1)} \sup_{t \in [0,T]} |\partial_t^{p+2}V(t)|_H^2, \end{aligned}$$

which yields the assertion.  $\square$

For the next two lemmas, we recall the standard result from interpolation theory

$$\sup_{t \in I_m} |(v - Pv)(t)| \leq C \tau_m^{q+1} \sup_{t \in I_m} |v^{(q+1)}(t)|, \quad (3.13)$$

for all  $v \in W^{q+1, \infty}(0, T)$ , see, for instance, the study by (Stoer *et al.*, 2002, Section 2.1.4).

LEMMA 3.14 There exists  $C \geq 0$  depending on  $q$  and  $T$  such that for all  $V \in H_\rho^{q+2}(\mathbb{R}; H)$

$$|(V - PV)(t_{m-1}^+)|_H \leq C \tau^{q+1} \sup_{t \in [0, T]} |\partial_t^{p+1} V(t)|_H.$$

*Proof.* This follows directly from (3.13) and the Sobolev-embedding theorem.  $\square$

With the previous lemmas we can already estimate  $g(U)$ . Now let us estimate the final term needed to estimate the error  $U - U^\tau$ .

LEMMA 3.15 There exists  $C \geq 0$  depending on  $\|M_0\|$ ,  $q$  and  $T$  such that for all  $U \in H_\rho^{q+2}(\mathbb{R}; H)$

$$\sup_{t \in [0, T]} \langle M_0(U - PU)(t), (U - PU)(t) \rangle_H \leq C \tau^{2(q+1)} \sup_{t \in [0, T]} |\partial_t^{p+1} U(t)|_H^2.$$

*Proof.* The result follows by applying Cauchy–Schwarz and Young inequality, and (3.13) with  $v = M_0 U$  and  $v = U$ .  $\square$

Combining the previous lemmas, Theorem 3.8 and Theorem 3.12, we can bound the discrete error in time.

THEOREM 3.16 Assume that  $U \in H_\rho^{q+3}(\mathbb{R}; H)$ . Then there exists  $C \geq 0$  depending on  $\|M_0\|$ ,  $\|M_1\|$ ,  $\rho$ ,  $T$ ,  $\gamma$ ,  $q$  such that

$$\sup_{t \in [0, T]} \langle M_0(U - U^\tau)(t), (U - U^\tau)(t) \rangle_H + e^{2\rho T} \|U - U^\tau\|_{Q, \rho}^2 \leq C e^{2\rho T} \tau^{2(q+1)} \sup_{t \in [0, T]} |\partial_t^{p+2} U(t)|_H^2.$$

*Proof.* By Lemma 3.13 and Lemma 3.14 applied to  $V = M_0 U$  we have that

$$g(U) \leq C_1 e^{2\rho T} \tau^{2(q+1)} \sup_{t \in [0, T]} |\partial_t^{p+2} U(t)|_H^2$$

for some  $C_1 \geq 0$ . We note that  $\|U - U^\tau\|_{Q, \rho} \leq \|\eta\|_{Q, \rho} + \|\xi\|_{Q, \rho} = \|\xi\|_{Q, \rho}$  and hence, by Theorem 3.8 we obtain

$$\|U - U^\tau\|_{Q, \rho}^2 \leq g(U).$$

Moreover,

$$\langle M_0(U - U^\tau)(t), (U - U^\tau)(t) \rangle_H \leq 2 \langle M_0 \eta(t), \eta(t) \rangle_H + 2 \langle M_0 \xi(t), \xi(t) \rangle_H$$

and thus, by Theorem 3.12 and Lemma 3.15 we infer

$$\sup_{t \in [0, T]} \langle M_0(U - U^\tau)(t), (U - U^\tau)(t) \rangle_H \leq C \left( \tau^{2(q+1)} \sup_{t \in [0, T]} |\partial_t^{q+1} U(t)|_H^2 + g(U) \right).$$

for some  $C \geq 0$ . Combining these estimates, the claim follows.  $\square$

**REMARK 3.17** In the above lemmas the regularity assumptions on  $U$  are higher than actually needed. Instead of  $U \in H_\rho^{q+3}(\mathbb{R}, H)$  it would be sufficient to assume  $U \in W_\rho^{q+2, \infty}(\mathbb{R}, H)$ . But in order to prove that claim from conditions on the right-hand side the easiest way is by proving  $U \in H_\rho^{q+3}(\mathbb{R}, H)$  and using the Sobolev embedding.

**REMARK 3.18** The above analysis holds for all evolutionary problems and the above theorem gives error bounds for the semidiscrete solution of order  $q + 1$ , assuming enough regularity in time. In the case of less regularity Theorems 3.8 and 3.12 still hold. For a fully discrete method, a spatial discretization has to be defined too. This step, however, has to be done for each problem considered separately.

#### 4. Full discretization for Example 1.3

Let us assume a regular, quasi uniform and shape-regular triangulation  $\Omega_h$  of  $\Omega$  into triangular open cells  $\sigma$  with maximal cell diameter  $h$ . Moreover, we assume that the interfaces between  $\Omega_{\text{ell}}$ ,  $\Omega_{\text{par}}$  and  $\Omega_{\text{hyp}}$  are polygonal such that the triangulation  $\Omega_h$  fits to these interfaces.

As the whole article is mainly concerned with the correct time discretization, in this section, we will employ the custom of the ‘generic constant’  $C \geq 0$  that may vary from line to line, which, however, depends on  $T$ ,  $\rho$ ,  $\|M_1\|$ ,  $\|M_0\|$ ,  $q$  and  $\gamma$  and on  $k$ , the order of the assumed spatial regularity, only.

Then the fully discretized counterpart  $\mathcal{V}_h^\tau$  to  $\mathcal{V}$  is given by

$$\mathcal{V}_h^\tau := \{(u_h, v_h) \in \mathcal{V}^\tau : u_h|_{I_m} \in \mathcal{P}_q(I_m, V_1(\Omega)), v_h|_{I_m} \in \mathcal{P}_q(I_m, V_2(\Omega)), m \in \{1, \dots, M\}\},$$

where the spatial spaces are

$$V_1(\Omega) := \left\{ v \in H_0^1(\Omega); \forall \sigma : v|_\sigma \in \mathcal{P}_k(\sigma) \right\},$$

$$V_2(\Omega) := \{w \in H(\text{div}, \Omega); \forall \sigma : w|_\sigma \in RT_{k-1}(\sigma)\}.$$

Here  $\mathcal{P}_k(\sigma)$  is the space of polynomials of degree up to  $k$  on the cell  $\sigma$  and  $RT_{k-1}(\sigma)$  is the Raviart–Thomas space, defined by

$$RT_{k-1}(\sigma) = (\mathcal{P}_{k-1}(\sigma))^n + \mathbf{x}\mathcal{P}_{k-1}(\sigma).$$

Note that

$$(\mathcal{P}_{k-1}(\sigma))^n \subset RT_{k-1}(\sigma) \subset (\mathcal{P}_k(\sigma))^n,$$

$$\text{div}(RT_{k-1}(\sigma)) \subset \mathcal{P}_{k-1}(\sigma) \quad \text{and}$$

$$RT_{k-1}(\sigma) \cdot \mathbf{n}|_{\partial\sigma} \subset \mathcal{P}_{k-1}(\partial\sigma).$$

Furthermore, if the mesh consists of quadrilateral or hexahedral cells, in the above definitions and statements the polynomial space  $\mathcal{P}_k(\sigma)$  can be replaced by a mapped  $\mathcal{Q}_k$  space, including all polynomials of total degree  $k$  over a reference element and then mapped onto  $\sigma$ . If the mesh is a combination of both types of cells, a combination of spaces also works with a suitable mapping ensuring the continuities.

**REMARK 4.1** (Solvability of the fully discrete system) We can apply the general existence theory that was also used in Proposition 3.2. More precisely, the positive definiteness still holds, since the triangulation fits to the interfaces and hence, the uniqueness of the system is warranted. However, since the problem is finite-dimensional, the uniqueness implies the existence of a solution of the fully discretized problem.

Let us come to the interpolation operator  $I = (I_1, I_2)$ . For  $I_1 : C(\Omega) \rightarrow V_1$  we use the Scott–Zhang interpolant on each cell  $\sigma$ , see the study by (Scott & Zhang, 1990) for a precise definition, that is patched together continuously. Here local interpolation error estimates can be given using  $L^2$  norms also in 3d, which is not possible for standard Lagrange interpolation. For  $I_2 : H(\text{div}, \Omega) \cap (L^s(\Omega))^n \rightarrow V_2$  with  $s > 2$  we also use the standard interpolator, defined via moments, see the study by (Brezzi & Fortin, 1991). Note that in the following, in order to avoid a cluttered notation as much as possible, we will not explicitly keep track on the number of components of the  $L^2(\Omega)$  or  $H^k(\Omega)$  spaces under consideration, as it will be obvious from the context.

Standard local interpolation error estimates yield for all  $v \in H_0^1(\Omega) \cap H^r(\Omega)$ ,

$$\|v - I_1 v\|_{0,\Omega} \leq Ch^r \|v\|_{r,\Omega}, \quad \|\nabla(v - I_1 v)\|_{0,\Omega} \leq Ch^{r-1} \|v\|_{r,\Omega}, \quad (4.1)$$

where  $1 \leq r \leq k + 1$ , see the study by (Scott & Zhang, 1990), and for all  $q \in H^s(\Omega)$  such that  $\text{div } q \in H^s(\Omega)$

$$\|q - I_2 q\|_{0,\Omega} \leq Ch^s \|q\|_{s,\Omega}, \quad \|\text{div}(q - I_2 q)\|_{0,\Omega} \leq Ch^s \|\text{div } q\|_{s,\Omega}, \quad (4.2)$$

where  $1 \leq s \leq k$ , see the study by (Brezzi & Fortin, 1991).

Let  $U_h^\tau \in \mathcal{V}_h^\tau$  be the solution of the fully discretized system and  $PIU \in \mathcal{V}_h^\tau$  be the interpolated solution of (1.1) for the operators  $M_0, M_1$  given in Example 1.3 and  $A$  given as in Example 1.1. Then we obtain analogously to the derivation of the errors of the semidiscretization

$$\begin{aligned} & \sup_{t \in [0, T]} \langle M_0(PIU - U_h^\tau)(t), (PIU - U_h^\tau)(t) \rangle_H + e^{2\rho T} \|PIU - U_h^\tau\|_{\mathcal{Q},\rho}^2 \\ & \leq C e^{2\rho T} \left( \|\partial_t M_0(U - \widehat{PIU})\|_{\mathcal{Q},\rho}^2 + \|M_1(U - PIU)\|_{\mathcal{Q},\rho}^2 + \|A(U - PIU)\|_{\mathcal{Q},\rho}^2 \right. \\ & \quad \left. + T \max_{1 \leq m \leq M} \left\{ |M_0(PIU - IU)(t_{m-1}^+)|_H^2 e^{-2\rho t_{m-1}} \right\} \right), \end{aligned} \quad (4.3)$$

where we remark that in contrast to Theorem 3.8 the terms  $\|M_1(U - PIU)\|_{\mathcal{Q},\rho}^2$  and  $\|A(U - PIU)\|_{\mathcal{Q},\rho}^2$  do not vanish, since we also interpolate with respect to space. In the following group of lemmas we estimate the terms on the right-hand side of (4.3) and start with a term particularly needed for the final

convergence estimate in Theorem 4.7. Beforehand, let us introduce

$$\|u\|_{Q,\rho,k,D}^2 = \sum_{m=1}^M Q_m \left[ |u|_{k,D}^2 \right] e^{-2\rho t_{m-1}},$$

where  $D \subseteq \Omega$  is measurable.

LEMMA 4.2 It holds for  $U = (u, v) \in H_\rho(\mathbb{R}; H^k(\Omega) \times H^k(\Omega))$

$$\|U - PIU\|_{Q,\rho} \leq Ch^k \left( \|u\|_{Q,\rho,k,\Omega} + \|v\|_{Q,\rho,k,\Omega} \right).$$

Moreover, if  $U = (u, v) \in H_\rho(\mathbb{R}; D(A))$  such that  $AU \in H_\rho(\mathbb{R}; H^k(\Omega) \times H^k(\Omega))$ , then

$$\|A(U - PIU)\|_{Q,\rho} \leq Ch^k \left( \|u\|_{Q,\rho,k+1,\Omega} + \|\operatorname{div} v\|_{Q,\rho,k,\Omega} \right).$$

*Proof.* By the definition of  $Q_m[\cdot]_\rho$  we have with (4.1) ( $r = k$ ) and (4.2) ( $s = k$ )

$$\begin{aligned} \|U - PIU\|_{Q,\rho}^2 &= \|U - IU\|_{Q,\rho}^2 = \sum_{m=1}^M Q_m \left[ \|u - I_1 u\|_{0,\Omega}^2 + \|v - I_2 v\|_{0,\Omega}^2 \right] e^{-2\rho t_{m-1}} \\ &\leq C \sum_{m=1}^M Q_m \left[ h^{2k} |u(\cdot)|_{k,\Omega}^2 + h^{2k} |v(\cdot)|_{k,\Omega}^2 \right] e^{-2\rho t_{m-1}} \\ &= Ch^{2k} \left( \|u\|_{Q,\rho,k,\Omega}^2 + \|v\|_{Q,\rho,k,\Omega}^2 \right). \end{aligned}$$

Very similarly we have for the second norm using (4.1) ( $r = k + 1$ ) and (4.2) ( $s = k$ )

$$\begin{aligned} \|A(U - PIU)\|_{Q,\rho}^2 &= \sum_{m=1}^M Q_m \left[ \|\nabla(u - I_1 u)\|_{0,\Omega}^2 + \|\operatorname{div}(v - I_2 v)\|_{0,\Omega}^2 \right] e^{-2\rho t_{m-1}} \\ &\leq Ch^{2k} \left( \|u\|_{Q,\rho,k+1,\Omega}^2 + \|\operatorname{div} v\|_{Q,\rho,k,\Omega}^2 \right). \end{aligned}$$

□

LEMMA 4.3 It holds for  $U = (u, v) \in H_\rho(\mathbb{R}; H^k(\Omega) \times H^k(\Omega))$

$$\|M_1(U - PIU)\|_{Q,\rho} \leq Ch^k \left( \|u\|_{Q,\rho,k,\Omega} + \|v\|_{Q,\rho,k,\Omega} \right).$$

*Proof.* The assertion follows from Lemma 4.2 and the boundedness of  $M_1$ .

□



LEMMA 4.4 For  $U = (u, v) \in H_\rho^1(\mathbb{R}; H^k(\Omega) \times H^k(\Omega)) \cap H_\rho^{q+2}(\mathbb{R}; L^2(\Omega) \times L^2(\Omega))$  we have

$$\begin{aligned} & \sup_{t \in [0, T]} \langle M_0(U - PIU)(t), (U - PIU)(t) \rangle_H \\ & \leq C \left( h^{2k} \sup_{t \in [0, T]} (|u(t)|_{k, \Omega} + |v(t)|_{k, \Omega})^2 + \tau^{2(q+1)} \sup_{t \in [0, T]} |\partial_t^{q+1} IU(t)|_H^2 \right). \end{aligned}$$

*Proof.* The operator  $M_0$  is selfadjoint and non-negative. Thus it follows that

$$\begin{aligned} \langle M_0(U - PIU)(t), (U - PIU)(t) \rangle_H &= |\sqrt{M_0}(U - PIU)(t)|_H^2 \\ &\leq 2 \left( |\sqrt{M_0}(U - IU)(t)|_H^2 + |\sqrt{M_0}(IU - PIU)(t)|_H^2 \right) \end{aligned}$$

for each  $t \in [0, T]$ . The second term can be estimated by

$$|\sqrt{M_0}(IU - PIU)(t)|_H^2 \leq C \tau^{2(q+1)} \sup_{t \in [0, T]} |\partial_t^{q+1} IU(t)|_H^2$$

according to Lemma 3.15, while the first term can be estimated by

$$|\sqrt{M_0}(U - IU)(t)|_{L^2(\Omega)}^2 \leq Ch^{2k} \left( |u(t)|_{k, \Omega}^2 + |v(t)|_{k, \Omega}^2 \right),$$

due to the boundedness of  $\sqrt{M_0}$ . Hence, the assertion follows.  $\square$

LEMMA 4.5 For  $U = (u, v) \in H_\rho^1(\mathbb{R}; H^k(\Omega) \times H^k(\Omega)) \cap H_\rho^{q+3}(\mathbb{R}; L^2(\Omega) \times L^2(\Omega))$ , we get

$$\|\partial_t M_0(U - \widehat{PIU})\|_{Q, \rho}^2 \leq C \left( h^{2k} (\|\partial_t u\|_{Q, \rho, k, \Omega} + \|\partial_t v\|_{Q, \rho, k, \Omega})^2 + \tau^{2(q+1)} \sup_{t \in [0, T]} |\partial_t^{q+2} IU(t)|_H^2 \right).$$

*Proof.* We have that

$$\begin{aligned} \|\partial_t M_0(U - \widehat{PIU})\|_{Q, \rho} &\leq \|\partial_t M_0(U - IU)\|_{Q, \rho} + \|\partial_t(M_0 IU - \widehat{PI} M_0 IU)\|_{Q, \rho} \\ &\leq \|M_0(\partial_t U - I\partial_t U)\|_{Q, \rho} + C \tau^{q+1} \sup_{t \in [0, T]} |\partial_t^{q+2} IU(t)|_H, \end{aligned}$$

by Lemma 3.13. For the first term we have by Lemma 4.2

$$\|M_0(\partial_t U - I\partial_t U)\|_{Q, \rho}^2 \leq Ch^{2k} \left( \|\partial_t u\|_{Q, \rho, k, \Omega}^2 + \|\partial_t v\|_{Q, \rho, k, \Omega}^2 \right).$$

$\square$

LEMMA 4.6 It holds for  $U = (u, v) \in H_\rho^{q+2}(\mathbb{R}; L^2(\Omega) \times L^2(\Omega))$

$$\max_{1 \leq m \leq M} \left\{ \|M_0(PIU - IU)(t_{m-1}^+) \|_{L^2(\Omega)} e^{-\rho t_{m-1}} \right\} \leq C \tau^{q+1} \sup_{t \in [0, T]} |\partial_t^{q+1} IU(t)|_H.$$

*Proof.* This is a direct consequence of Lemma 3.14.  $\square$

Lemmas 4.2 to 4.6 give us all needed estimates for the final convergence result for Example 1.3.

**THEOREM 4.7** We assume for the solution  $U = (u, v)$  of Example 1.3 the regularity

$$U \in H_\rho^1(\mathbb{R}; H^k(\Omega) \times H^k(\Omega)) \cap H_\rho^{q+3}(\mathbb{R}; L^2(\Omega) \times L^2(\Omega))$$

as well as

$$AU \in H_\rho(\mathbb{R}; H^k(\Omega) \times H^k(\Omega)).$$

Then we have for the error of the numerical solution by (3.1)

$$\sup_{t \in [0, T]} \langle M_0(U - U_h^\tau)(t), (U - U_h^\tau)(t) \rangle_H + e^{2\rho T} \|U - U_h^\tau\|_{Q, \rho}^2 \leq C e^{2\rho T} (\tau^{2(q+1)} + Th^{2k}).$$

## 5. Numerical examples

In the following section we consider some examples to verify numerically our theoretical findings. In all of them we use homogeneous initial conditions, i.e., we set  $x_0 = 0$ . Using other values would also be possible. Our computations were done with the finite-element framework `SOFE` developed by L. Ludwig, see [github.com/SOFE-Developers/SOFE](https://github.com/SOFE-Developers/SOFE).

### 5.1 Changing type system—one space dimension

Let  $\Omega = \left(-\frac{3\pi}{2}, \frac{3\pi}{2}\right)$ ,  $\Omega_{\text{hyp}} = \left(-\frac{3\pi}{2}, 0\right)$ ,  $\Omega_{\text{par}} = \left(0, \frac{3\pi}{2}\right)$ . The problem is given on  $\mathbb{R} \times \Omega$  by

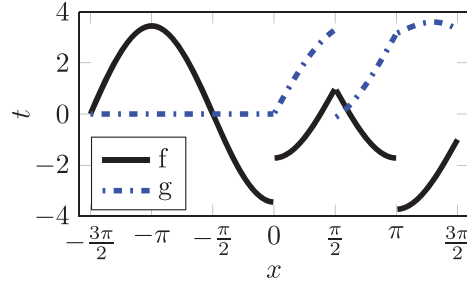
$$\left( \partial_t \begin{pmatrix} 1 & 0 \\ 0 & \chi_{\Omega_{\text{hyp}}} \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & \chi_{\Omega_{\text{par}}} \end{pmatrix} + \begin{pmatrix} 0 & \partial_x \\ \partial_x & 0 \end{pmatrix} \right) \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix} \quad (5.1a)$$

with  $u(t, -\frac{3\pi}{2}) = u(t, \frac{3\pi}{2}) = 0$  and

$$\begin{aligned} f(t, x) = \chi_{\mathbb{R}_{\geq 0}}(t) & \left( - (2e^t - t - 1) \chi_{(-\frac{\pi}{2}, 0)}(x) \cos(x) + \chi_{(0, \pi)}(x) - \chi_{\left(\pi, \frac{3\pi}{2}\right)}(x) \right. \\ & \left. + e^t \left( \chi_{\left(\frac{\pi}{2}, \frac{3\pi}{2}\right)}(x) - \chi_{\left(0, \frac{\pi}{2}\right)}(x) \right) \cos(x) \right), \end{aligned} \quad (5.1b)$$

$$\begin{aligned} g(t, x) = \chi_{\mathbb{R}_{\geq 0}}(t) & \left( \chi_{(0, \pi)}(x)x + \chi_{\left(\pi, \frac{3\pi}{2}\right)}(x)(2\pi - x) \right. \\ & \left. - (e^t - 1) \left( \chi_{\left(\frac{\pi}{2}, \frac{3\pi}{2}\right)}(x) - \chi_{\left(0, \frac{\pi}{2}\right)}(x) \right) \sin(x) \right). \end{aligned} \quad (5.1c)$$

Note that the right-hand side ( $f$  and  $g$ ) are only in  $L^2$ ; Fig. 1 shows them for  $t = 1$ .

FIG. 1. Right-hand sides  $f$  and  $g$  of problem (5.1).

The solution can be derived as

$$u(t, x) = \chi_{\mathbb{R}_{\geq 0}}(t)(e^t - 1)(\chi_{(\frac{\pi}{2}, \frac{3\pi}{2})}(x) - \chi_{(-\frac{3\pi}{2}, \frac{\pi}{2})}(x)) \cos(x),$$

$$v(t, x) = \chi_{\mathbb{R}_{\geq 0}}(t) \left( -(e^t - t - 1) \chi_{(-\frac{3\pi}{2}, 0)}(x) \sin(x) + \chi_{(0, \pi)}(x)x + \chi_{(\pi, \frac{3\pi}{2})}(x)(2\pi - x) \right).$$

We observe that  $u$  and  $v$  are nondifferentiable, but piecewise smooth. Figure 2 shows the solutions for  $t \in [0, 1]$ . Note that *a priori*, we impose no transmission condition. However, as in (Waurick, 2016, Remark 3.2), they can be derived for  $u$  satisfying (5.1) as

$$u(t, 0+) = u(t, 0-), \quad \partial_x u(t, 0+) = \int_0^t \partial_x u(s, 0-) ds.$$

For the numerical solution we use  $T = 1$ , an equidistant mesh of  $M$  cells in time and an equidistant mesh of  $N$  cells in space, thus  $\tau = 1/M$  and  $h = 3\pi/N$ . In order to capture the jumps of  $f$  and  $g$ , and to resolve the boundary  $S = \overline{\Omega}_h \cap \overline{\Omega}_p = \{0\}$  we use an equidistant mesh in space with the number of cells  $N$  divisible by 6. Note that we can use  $\rho = 1$  for the given solution  $u$ .

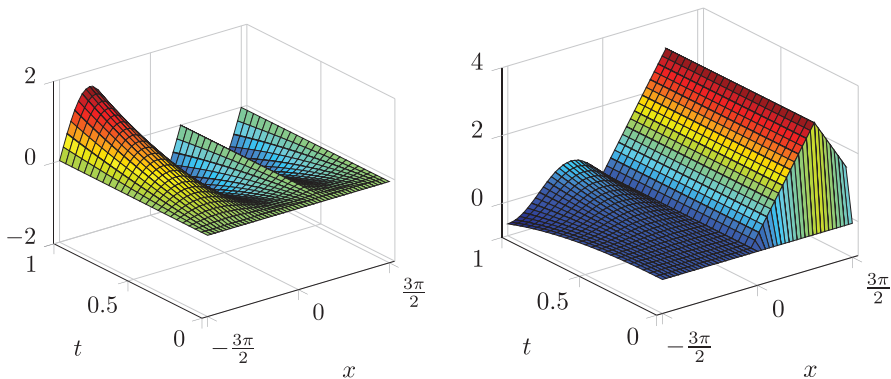
FIG. 2. Solution  $u$  (left) and  $v$  (right) of problem (5.1).

TABLE 1 *Convergence results for  $U - U_h$  of problem (5.1)*

$N = M$	$E_{\text{sup}}(U - U_h)$		$\ U - U_h\ _{Q,\rho}$		$\ U - U_h\ _\rho$	
$p = 2, q = 1$						
96	3.577e-04		6.400e-05		6.637e-05	
192	9.010e-05	1.99	1.601e-05	2.00	1.660e-05	2.00
384	2.261e-05	1.99	4.002e-06	2.00	4.150e-06	2.00
768	5.662e-06	2.00	1.001e-06	2.00	1.037e-06	2.00
$p = 3, q = 2$						
96	6.981e-08		7.500e-10		1.468e-08	
192	8.726e-09	3.00	2.343e-11	5.00	1.833e-09	3.00
384	1.091e-09	3.00	7.329e-13	5.00	2.291e-10	3.00
768	1.363e-10	3.00	2.474e-14	4.89	2.864e-11	3.00

TABLE 2 *Estimated convergence rates for  $E(U - U_h)$  of problem (5.1) and several polynomial orders*

$p \setminus q$	1	2	3	4	5
1	2	3	3	3	3
2	2	2	2	2	2
3	2	3	5	5	5
4	2	3	4	4	4
5	2	3	4	7	7

Defining

$$E_{\text{sup}}(v)^2 := \sup_{t \in [0, T]} \langle M_0 v(t), v(t) \rangle_H, \quad E(v)^2 := \sup_{t \in [0, T]} \langle M_0 v(t), v(t) \rangle_H + \|v\|_{Q,\rho}^2$$

we consider in Table 1 the convergence behaviour of  $U_h$  for  $N = M$  and polynomial degrees  $q = p + 1 = 2$  and  $q = p + 1 = 3$ . Note that we approximate the supremum by a maximum over a large number of evaluations and also show the norm  $\|U - U_h\|_\rho$  estimated by a refined quadrature rule in the last columns. The estimated rates of convergence support our theoretical result in Theorem 4.7 that the error  $E$  is of order  $\min\{p, q + 1\}$ . For odd polynomial degrees  $p$  the component  $\|U - U_h\|_{Q,\rho}$  shows a higher convergence order, hinting at a superconvergence property. In Table 2 the estimated convergence rates for all combinations of polynomial degrees  $\{p, q\} \subseteq \{1, \dots, 5\}$  are given. Clearly the rates for even  $p$  follow the predicted  $\min\{p, q + 1\}$ , while for odd  $p$  the rates are larger. Thus there might be a superconvergence phenomenon.

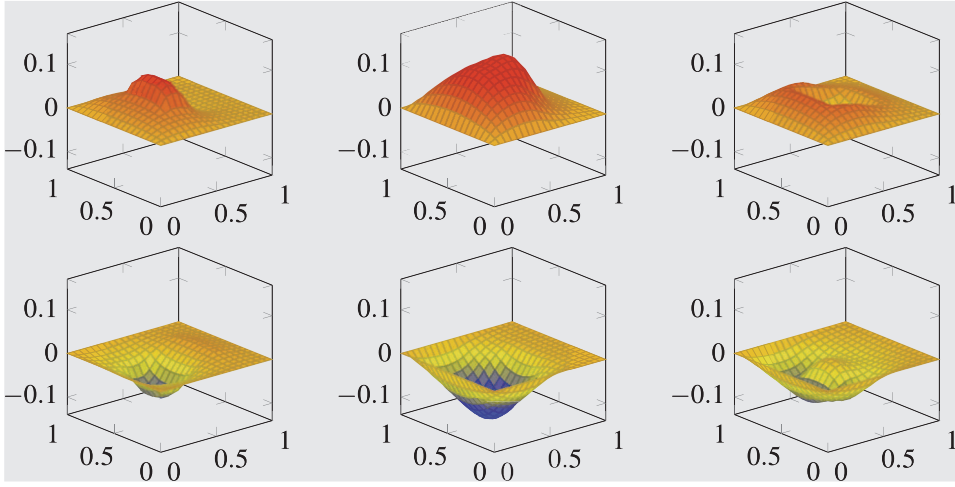


FIG. 3. Solution  $u$  at times  $t = 5k/16$  for  $k \in \{1, \dots, 6\}$  (top left to bottom right) of problem (5.2) for  $T = 1.875$ .

TABLE 3 Convergence results for  $\tilde{U} - U_h$  of problem (5.2)

$N = M$	$E_{\text{sup}}(\tilde{U} - U_h)$		$\ \tilde{U} - U_h\ _{Q,\rho}$		$\ \tilde{U} - U_h\ _\rho$	
$p = 2, q = 1$						
16	1.666e-03		7.445e-04		8.517e-04	
32	5.260e-04	1.66	2.790e-04	1.42	3.012e-04	1.50
64	1.926e-04	1.45	1.300e-04	1.10	1.331e-04	1.18
$p = 3, q = 2$						
16	4.015e-04		2.414e-04		2.419e-04	
32	1.430e-04	1.49	1.175e-04	1.04	1.175e-04	1.04
64	5.245e-05	1.45	5.075e-05	1.21	5.072e-05	1.21

## 5.2 Changing type system—two space dimensions

As a final example we consider a problem with an unknown solution. Let  $\Omega = (0, 1)^2 \subset \mathbb{R}^2$ ,  $\Omega_{\text{hyp}} = \left(\frac{1}{4}, \frac{3}{4}\right)^2$ ,  $\Omega_{\text{ell}} = \Omega \setminus \bar{\Omega}_{\text{hyp}}$  and  $\Omega_{\text{par}} = \emptyset$ . The problem is given on  $\mathbb{R} \times \Omega$  by

$$\left( \partial_t \begin{pmatrix} \chi_{\Omega_{\text{hyp}}} & 0 \\ 0 & \chi_{\Omega_{\text{hyp}}} \end{pmatrix} + \begin{pmatrix} \chi_{\Omega_{\text{ell}}} & 0 \\ 0 & \chi_{\Omega_{\text{ell}}} \end{pmatrix} + \begin{pmatrix} 0 & \text{div} \\ \nabla_0 & 0 \end{pmatrix} \right) \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} f \\ 0 \end{pmatrix}, \quad (5.2)$$

where

$$f(t, x) = 2 \sin(\pi t) \chi_{\mathbb{R}_{<1/2} \times \mathbb{R}}(x).$$

For  $T = 1.875$  Fig. 3 shows some snapshots of the component  $u$  of the solution  $U$ , approximated by a numerical simulation.

In order to investigate the error behaviour upon refinement of the discretization, we use a numerically computed reference solution  $\tilde{U}$  instead of the real one  $U$ . For this we set  $T = 1$  and use an equidistant mesh of  $128 \times 128$  rectangular cells in space and 128 cells in time, and polynomial degrees  $p = 3$  and  $q = 2$ . Thus  $u$  is approximated in space by piecewise  $\mathcal{Q}_3$  elements,  $v$  by  $RT_2$  elements and both in time by  $\mathcal{P}_2$  elements. In Table 3 we see the results of our numerical simulation for two pairs of polynomial order. We observe that the error rates are independent of the polynomial order and furthermore less than the optimal orders given in Theorem 4.7. The reason for this decrease in convergence order lies in the reduced regularity of the solution to this given problem. The interior boundaries where the type of the problem changes introduce corners, where it is very likely for singular solution components to arise.

### Acknowledgements

M. W. carried out this work with financial support of the EPSRC grant EP/L018802/2: ‘Mathematical foundations of metamaterials: homogenization, dissipation and operator theory’. This is gratefully acknowledged.

### REFERENCES

- AKRIVIS, G. & MAKRIDAKIS, C. (2004) Galerkin time-stepping methods for nonlinear parabolic equations. *ESAIM: M2AN*, **38**, 261–289.
- AKRIVIS, G., MAKRIDAKIS, C. & NOCHETTO, R. (2011) Galerkin and Runge–Kutta methods: unified formulation, a posteriori error estimates and nodal superconvergence. *Numerische Mathematik*, **118**, 429–456.
- ANTONIĆ, N., BURAZIN, K. S. & VRDOLJAK, M. (2013) Heat equation as a Friedrichs system. *J. Math. Anal. Appl.*, **404**, 537–553.
- ANTONIĆ, N., BURAZIN, K. S. & VRDOLJAK, M. (2014) Second-order equations as Friedrichs systems. *Nonlinear Anal. Real World Appl.*, **15**, 290–305.
- BREZZI, F. & FORTIN, M. (1991) *Mixed and Hybrid Finite Element Methods*. Springer Series in Computational Mathematics, vol. **15**. New York: Springer.
- BURAZIN, K. S. & ERCEG, M. (2016) Non-stationary abstract Friedrichs systems. *Mediterr. J. Math.*, **13**, 3777–3796.
- COCKBURN, B., KARNIADAKIS, G. E. & SHU, C.-W. (2000) *The Development of Discontinuous Galerkin Methods*. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 3–50.
- ERN, A. & GUERMOND, J.-L. (2006a) Discontinuous Galerkin methods for Friedrichs’ systems. I: General theory. *SIAM J. Numer. Anal.*, **44**, 753–778.
- ERN, A. & GUERMOND, J.-L. (2006b) Discontinuous Galerkin methods for Friedrichs’ systems. II: Second-order elliptic PDEs. *SIAM J. Numer. Anal.*, **44**, 2363–2388.
- ERN, A. & GUERMOND, J.-L. (2008) Discontinuous Galerkin methods for Friedrichs’ systems. Part III. Multifield theories with partial coercivity. *SIAM J. Numer. Anal.*, **46**, 776–804.
- FRIEDRICHS, K. (1958) Symmetric positive linear differential equations. *Commun. Pure Appl. Math.*, **11**, 333–418.
- JENSEN, M. (2004) Discontinuous Galerkin methods for Friedrichs systems with irregular solutions. *Ph.D. Thesis*, Oxford: University of Oxford.
- KALAUCH, A., PICARD, R., SIEGMUND, S., TROSTORFF, S. & WAURICK, M. (2014) A Hilbert space perspective on ordinary differential equations with memory term. *J. Dyn. Differ. Equations*, **26**, 369–399.
- MUKHOPADHYAY, S., PICARD, R., TROSTORFF, S. & WAURICK, M. (2017) A Note on a Two-Temperature Model in Linear Thermoelasticity. *Math. Mech. Solids*, **22**, 905–918.
- MULHOLLAND, A. J., PICARD, R., TROSTORFF, S. & WAURICK, M. (2016) On well-posedness for some thermopiezoelectric coupling models. *Math. Methods Appl. Sci.*, **39**, 4375–4384.
- PICARD, R. (2009) A structural observation for linear material laws in classical mathematical physics. *Math. Methods Appl. Sci.*, **32**, 1768–1803.

- PICARD, R. & MCGHEE, D. (2011) *Partial Differential Equations. A Unified Hilbert Space Approach*. de Gruyter Expositions in Mathematics 55. Berlin: de Gruyter. xviii.
- PICARD, R., TROSTORFF, S., WAURICK, M. & WEHOWSKI, M. (2013) On non-autonomous evolutionary problems. *J. Evol. Equ.*, **13**, 751–776.
- PICARD, R., TROSTORFF, S. & WAURICK, M. (2015) On some models for elastic solids with micro-structure. *ZAMM, Z. Angew. Math. Mech.*, **95**, 664–689.
- PICARD, R., TROSTORFF, S. & WAURICK, M. (2016) On a comprehensive class of linear control problems. *IMA J. Math. Control Inf.*, **33**, 257–291.
- PRESS, W. H., TEUKOLSKY, S. A., VETTERLING, W. T. & FLANNERY, B. P. (2007) *Numerical Recipes* 3rd Edition: *The Art of Scientific Computing*, 3rd edn. New York: Cambridge University Press.
- REED, W. H. & HILL, T. R. (1973) Triangular mesh methods for the neutron transport equation. Submitted to American Nuclear Society Topical Meeting on Mathematical Models and Computational Techniques for Analysis of Nuclear Systems, Los Alamos Laboratory.
- RIVIÈRE, B. (2008) *Discontinuous Galerkin Methods for Solving Elliptic and Parabolic Equations*. Society for Industrial and Applied Mathematics.
- SCOTT, L. R. & ZHANG, S. (1990) Finite element interpolation of nonsmooth functions satisfying boundary conditions. *Math. Comp.*, **54**, 483–493.
- STOER, J., BARTELS, R., GAUTSCHI, W., BULIRSCH, R. & WITZGALL, C. (2002) *Introduction to Numerical Analysis. Texts in Applied Mathematics*. New York: Springer New York.
- TROSTORFF, S. (2012) An alternative approach to well-posedness of a class of differential inclusions in Hilbert spaces. *Nonlinear Anal., Theory Methods Appl., Ser. A, Theory Methods*, **75**, 5851–5865.
- TROSTORFF, S. & WAURICK, M. (2016) On the weighted Gauss-Radau Quadrature. Preprint arXiv:1610.09016.
- TROSTORFF, S. & WEHOWSKI, M. (2014) Well-posedness of non-autonomous evolutionary inclusions. *Nonlinear Anal., Theory Methods Appl., Ser. A, Theory Methods*, **101**, 47–65.
- VLASAK, M. & ROOS, H.-G. (2014) An optimal uniform *a priori* error estimate for an unsteady singularly perturbed problem. *Int. J. Numer. Anal. Model.*, **11**, 24–33.
- WAURICK, M. (2015) On non-autonomous integro-differential-algebraic evolutionary problems. *Math. Methods Appl. Sci.*, **38**, 665–676.
- WAURICK, M. (2016) Stabilization via homogenization. *Appl. Math. Lett.*, **60**, 101–107.