

# Temporal blocking of finite-difference stencil operators with sparse “off-the-grid” sources

George Bisbas  
Imperial College London  
London, UK  
g.bisbas18@imperial.ac.uk

Fabio Luporini  
Devito Codes  
London, UK  
fabio@devitocodes.com

Mathias Louboutin  
Georgia Institute of Technology  
Atlanta, GA  
mlouboutin3@gatech.edu

Rhodri Nelson  
Imperial College London  
London, UK  
rnelson@imperial.ac.uk

Gerard J. Gorman  
Imperial College London  
London, UK  
g.gorman@imperial.ac.uk

Paul H.J. Kelly  
Imperial College London  
London, UK  
p.kelly@imperial.ac.uk

**Abstract**— Stencil kernels dominate a range of scientific applications, including seismic and medical imaging, image processing, and neural networks. Temporal blocking is a performance optimization that aims to reduce the required memory bandwidth of stencil computations by re-using data from the cache for multiple time steps. It has already been shown to be beneficial for this class of algorithms. However, applying temporal blocking to practical applications’ stencils remains challenging. These computations often consist of sparsely located operators not aligned with the computational grid (“off-the-grid”). Our work is motivated by modelling problems in which source injections result in wavefields that must then be measured at receivers by interpolation from the grided wavefield. The resulting data dependencies make the adoption of temporal blocking much more challenging. We propose a methodology to inspect these data dependencies and reorder the computation, leading to performance gains in stencil codes where temporal blocking has not been applicable. We implement this novel scheme in the Devito domain-specific compiler toolchain. Devito implements a domain-specific language embedded in Python to generate optimized partial differential equation solvers using the finite-difference method from high-level symbolic problem definitions. We evaluate our scheme using isotropic acoustic, anisotropic acoustic, and isotropic elastic wave propagators of industrial significance. After auto-tuning, performance evaluation shows that this enables substantial performance improvement through temporal blocking over highly-optimized vectorized spatially-blocked code of up to 1.6x.

**Index Terms**—temporal blocking, stencil computations, code generation, partial differential equations, seismic imaging, domain-specific languages, wave-propagation

## I. INTRODUCTION

Stencils are commonly encountered in scientific applications such as image processing [1], convolutional neural networks, weather forecasting [2], computational fluid dynamics, seismic [3], [4], and medical imaging [5]. We present a scheme that enables the application of temporal blocking, a common technique to enhance cache-locality [6]–[8], to a class of finite-difference (FD) [9], [10] stencil kernels where this is challenging. This class consists of additional sparsely-located in the computational grid (off-the-grid) operators. Typical stencil kernels are computational patterns that are usually

functions of the nearest neighboring point values. In a more general context, a stencil defines the iterative computation of an element in an  $n$ -dimensional spatial grid at time  $t$  as a function of neighboring grid elements (space dependencies) at time  $t - 1, \dots, t - k$  (time dependencies). A typical stencil update in a scientific simulation has a 3-dimensional spatial iteration space and 1-dimensional temporal iteration space. Figure 1 illustrates a 1D stencil and its flow dependences. Each point is updated using values from the previous timestep and the right and left neighbors. Arrows illustrate the flow dependencies. Halo points (grey) are used to extend the computational domain by the stencil radius size. Wider stencils in 3D and their respective data dependencies are illustrated in Figure 2.

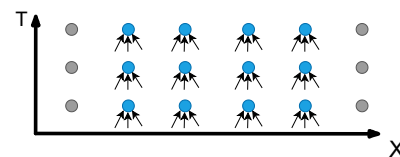


Fig. 1: A 1D-3pt Jacobi stencil update. Arrows show the data flow dependencies, grey points indicate the read-only halo area.

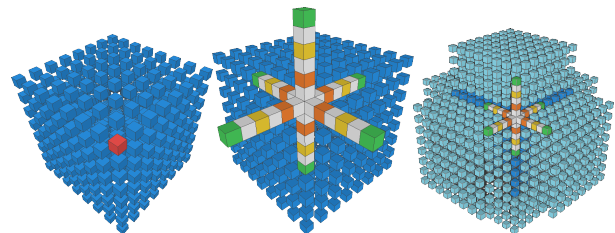


Fig. 2: A 6th-order ( $O(1,8)$ ) 3D-19pt stencil update. A point (red) at the edge of a block (blue) depends on a four-deep halo of neighbouring points which extends outside the block.

In addition to data dependencies of the kind illustrated in Figure 1, applications such as seismic wave-modeling carry additional dependencies owing to the interpolation of data not directly associated with grid points into the model (e.g., source injection). These positions, that are not aligned with the grid points are sparsely-distributed off-the-grid [11], [12] positions as shown in Figure 3a. They may also include receivers that interpolate neighboring values to take measurements, as shown in Figure 3b. Sources and receivers are sets of sparsely-distributed off-the-grid points. We iterate over these sparsely-located sets through indirections applying their effect to the grid points after iterating the 3D grid for stencil updates for each timestep. A loop nest structure illustrating this computation pattern is shown in Listing 1, where  $nt$  is the number of time steps;  $nx$ ,  $ny$ ,  $nz$  are the number of grid points along the  $x$ ,  $y$ ,  $z$  dimensions respectively,  $A(t, x, y, z)$  is the stencil kernel update and  $so$  is the space discretisation order. The `src` is of size  $nt \times len(sources)$  holding the wavefield for each timestep for every source where `sources` is the structure holding the information for modeling source injection. Variable `np` accounts for the number of points affected from a source, and  $f$  is the function defining the type of interpolation (e.g., bilinear, trilinear). An example of bilinear interpolation is shown in Figure 3, where 4 points are affected in 2D space. The `sources` set shown in Listing 1 provides the sparse off-the-grid coordinates for the injection. We iterate this set of coordinates that determine the affected neighboring points. The wave amplitude is scattered to these affected points.

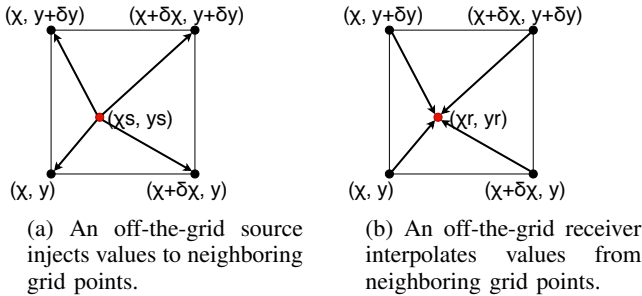


Fig. 3: A source injection and a receiver measurement interpolation at off-the-grid positions in a 2-D FD-grid. We assume linear interpolation.

#### A. Problem overview: a running example

Space blocking [13], [14] can be applied in computational patterns similar to Listing 1, but applying temporal blocking is challenging as we illustrate with a 1-D example in Figure 4. Red diamonds indicate off-the-grid coordinates where sparse operators are applied. Sparse operators are applied after a time iteration for the whole domain is finished. Consequently, separating the FD grid into blocks does not violate any data dependencies.

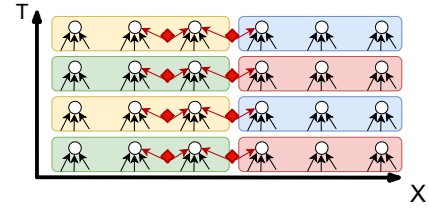
In contrast to space blocking, temporal blocking cannot be applied. When a sparse operator is located at an off-the-grid

**Listing 1:** A typical time-stepping loop nest structure for a stencil update with source injection. This stencil has one temporal and three spatial dimensions.

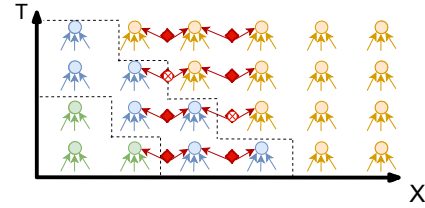
```

1 for t = 1 to nt do
2   for x = 1 to nx do
3     for y = 1 to ny do
4       for z = 1 to nz do
5          $A(t, x, y, z) \equiv u[t, x, y, z] = u[t-1, x, y, z] + \sum_{r=1}^{so/2} w_r ($ 
            $u[t-1, x-r, y, z] + u[t-1, x+r, y, z] + u[t-1, x, y-r, z] +$ 
            $u[t-1, x, y+r, z] + u[t-1, x, y, z-r] + u[t-1, x, y, z+r] );$ 
6       foreach s in sources do // For every source
7         for i = 1 to np do // Get the points affected
8            $xs, ys, zs = map(s, i)$  // through indirection
9            $u[t, xs, ys, zs] += f(src(t, s))$  // add their impact
           on the field

```



(a) Rectangular space blocking. All grid points can be updated in parallel at a specific time-step. Sparse operators fit within space blocking as their effect is imposed after all points have been updated.



(b) Skewed/Wave-front temporal blocking. Grid points are updated in waves. During a wave-front update, we compute grid point values for multiple timesteps. Applying sparse operators in the space boundary may lead to erroneous updates since source injection may precede the stencil update for a particular timestep. We have a data dependency violation.

Fig. 4: Sparse operators do not violate data dependencies in space blocking, Fig. 4a in contrast to temporal blocking, Fig. 4b.

position, among points that belong to different space blocks, data dependencies are violated, thus yielding incorrect results. The violation occurs because updates in space may pause for a particular time-step, and computation will proceed in time rather than space. Consequently, a sparse operator update may be computed, and points that have not yet been updated through the stencil kernel updates may be affected. Similarly, a point may be erroneously updated due to a forward move in time but may miss injection from a neighboring off-the-grid operator due to space-time block constraint. Data dependencies are violated, and similar violations are raised in other variants of temporal blocking, such as wave-front temporal blocking [15], [16] diamond temporal blocking [17], [18] and others.

We aim to overcome this limitation through our contributions in this paper.

Because the set of sources is sparse, the loops generated in Listing 1 by modelling source injection consist of non-affine accesses as illustrated in Listing 1. While polyhedral tools such as PLUTO [19], [20], Polly [21], Loopy [22], and CLoG [23] manage to deal with the first uniform stencil update, they are not capable of dealing with the non-affine nature of the source injection loop nests.

The sparse off-the-grid nature of the source operator combined with the non-affine nature of our loop structure illustrated in Listing 1 is blocking the application of time-tiling as described in subsection I-A. The methodology presented in Section II aims to overcome this limitation.

## B. Related work

1) *Improving stencil performance*: stencils offer good parallelisation opportunities, ranging from Instruction Level Parallelism (ILP), SIMD register-level parallelism (SSE, AVX) to shared-memory (OpenMP, OpenACC) and task parallelism. Furthermore, distributed-memory parallelism is often employed.

In most applications of interest, stencil kernels have low operational intensity, having few floating-point operations per byte of data accessed, and are therefore memory bandwidth-bound. Considerable effort has been put into caches for improving stencil performance. Rescheduling the order of computations towards increased cached memory reuse can increase throughput by alleviating memory bandwidth boundedness issues.

a) *Spatial cache blocking*: as illustrated in [13], [14], [24]–[30] FD grids can easily be decomposed into tiles and benefit from cache blocking. The same idea is applicable to unstructured grids, though this requires more sophisticated algorithms to create efficient schedules [31]. Space tiling techniques have also been implemented for execution on GPUs. Related work includes automated split tiling with trapezoids [32], hybrid hexagonal tiling [33] and automated HPC GPU code [34], [35], [36] as well as hybrid spatial/temporal blocking on FPGAs [37].

b) *Temporal cache blocking*: extending cache reuse in the time dimension led to the development of temporal blocking algorithms. To further reduce cache misses, we utilize computed values in a block to update values in the next timestep where possible. While one or more timesteps for a given block’s values are stored in the cache, we start computing the next time step for this block, not depending on the requirement to compute all the blocks of a given grid for previous timesteps. Spatial and temporal reuse are often fused into hybrid models (equidistant locality) to harness the advantages of both methods [37]. Plenty of research has been conducted in designing and evaluating temporal blocking schemes ranging from simple skewing [6]–[8], [16], [38], [39] and wave-front [15] to more sophisticated such as diamond [17], [18], [40]–[43]. The technique presented in this paper

enables such schedules to be used in applications with off-the-grid operators. While narrow stencil kernels exhibit good temporal locality, temporal blocking gains decrease when space-order increases. Higher space order problems limit temporal locality as more space updates are required to update one value in time.

2) *Domain Specific Languages*: improving performance is essential but usually comes with the price of error-prone hand-optimization. Finding ways to automate HPC code generation led to the birth of several domain-specific Languages (DSLs) like Devito [3], [4], which is used in this paper. Several DSL/compiler frameworks are working towards the automated generation of PDE solvers, such as FEniCS [44], and Firedrake [45]. Halide [1], implemented as an internal DSL in C++. OpenSBLI [46] is another framework that generated C code from Python-based high-level abstractions targeting equations written in Einstein notation. Halide is a language targeting code generation for digital image processing featuring memory locality and vectorized computation optimizations and ported to multi-core CPUs and GPUs.

Other automatic code-generation frameworks are OPS [47] for GPU code generation and YASK [48], a DSL to create high-performance FD-stencil code. Lift [49], achieves performance portability on parallel accelerators by combining high-level functional data-parallel language with rewrite rules which encode algorithmic and hardware-specific optimization choices. Stella [2] and GridTools [50] are DSLs embedded in C++, focusing on weather and climate HPC simulations.

## C. Contributions

Our contributions are:

- We propose a scheme that precomputes the off-the-grid sparse operators’ effect, allowing to reorder the computations for FD wave propagators, thus enabling the application of temporal blocking to stencil codes consisting of sparse operators such as source injection and measurement interpolation. Our scheme is cost-efficient, adding a negligible overhead compared to the measured gains.
- We implement the algorithm directly on top of the Devito DSL, harnessing the power of automated code generation, thus providing a pathway to express any similar operator in a form that exploits the benefits of time tiling with only minimal coding effort.
- We evaluate our scheme using 3D stencils encountered in wave propagation applications (isotropic acoustic, isotropic elastic, and anisotropic acoustic (TTI)), each having different memory and compute requirements.
- We achieve performance gains ranging from 15% to 60% for space order 4 and 8 for isotropic acoustic and elastic and anisotropic acoustic as well as 5% to 10% gains for elastic and TTI cases at space order 12.

Our work is mainly motivated by the need to enable and automate these optimizations for a class of seismic and medical imaging problems. Characteristic examples of such applications include full-waveform inversion (FWI) [51] and

reverse time migration (RTM) [52]. As future work, we aim to deliver these optimizations as a fully automated workflow.

The rest of the paper is organized as follows: Section II presents the approach followed to solve our problem. Section III introduces the wave-propagation kernels to be evaluated, and Section IV presents an experimental evaluation of the applicability and impact of the approach. Finally, in Section V, we discuss and summarise our work and briefly refer to our plans for future work.

## II. METHODOLOGY AND IMPLEMENTATION

In this section, we describe our approach that enables temporal blocking for wave-propagators with sparse operators. We describe the individual steps and present the details of our implementation. The whole precomputation workflow benefits from the power of the *Devito* DSL [3], [4] to automatically generate code and the data structures required by our scheme in its DSL. Afterward, we manually transform the generated loops to implement wave-front temporal blocking (WTB) [8], [53], [54], a representative temporal blocking schedule.

### A. Source injection precomputation

The source injection modeling consists of the following parameters: the number of sources, their coordinates, and their wavelet time-series. This data is enough to precompute their effect on an empty grid. We assume that the sources' coordinates are constant across our models' time-domain though this may not always be the case. However, *Devito*'s API can support the moving sources' case, and our algorithm is independent of it.

1) *Iterate sources' coordinates and store indices of affected points*: initially, we iterate over each source and inject to an empty grid for one timestep, assuming the wavefield is not zero at the first timestep. If the wavefield is zero at the first timestep, we may inject for more timesteps. Our experiments use wavefields with non-zero values at the first timesteps. The pseudocode is illustrated in 2. We use *Devito* to automatically generate code for this step. Our scheme is independent of the injection and interpolation type (e.g., non-linear injection). Then, we store the non-zero grid point coordinates.

**Listing 2:** Source injection over an empty grid. No PDE stencil update is happening.

```

1 for t = 1 to 2 do
2   foreach s in sources do
3     for i = 1 to np do
4       |xs, ys, zs = map(s, i);
5       |u[t, xs, ys, zs] += f(src(t, s))

```

2) *Generate sparse binary mask and unique IDs*: using the nonzero indices, we populate two arrays. The first array (Fig. 5b) is a binary integer mask of our grid with 1s at indices where  $u$  is nonzero. Ones are shown as filled bullet circles, with a green background, Fig. 5b. The second one (Fig. 5c) has the same shape and is populated with unique ascending values for each unique point affected. It is quite common to encounter points being affected by more than one source. Figures show an x-y plane (z-slice) of the 3D grid.

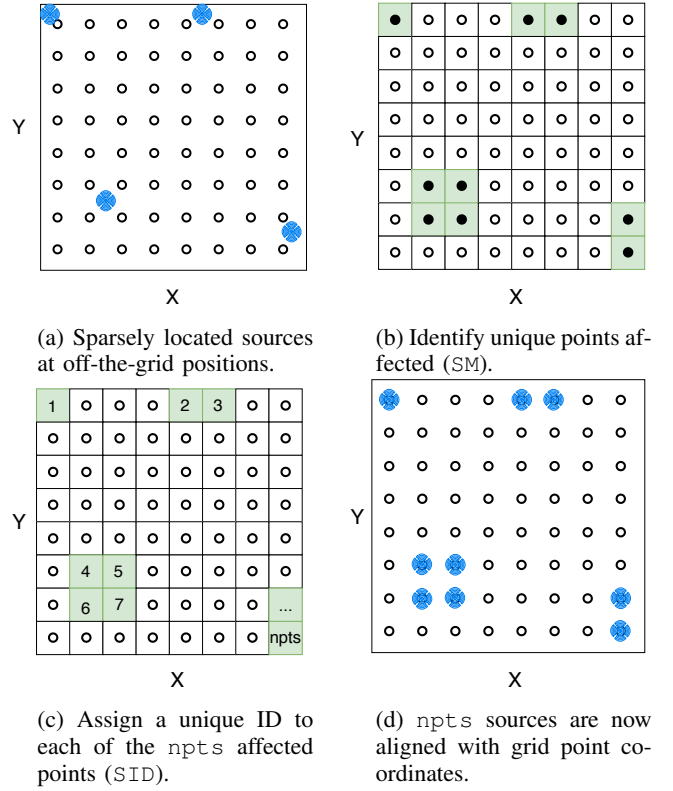


Fig. 5: Illustration of the four steps through which source impact is aligned to the computational grid. The figures show an x-y plane slice of the 3D grid.

3) *Decompose wavefields*: knowing the unique positions affected and their coordinates, we now use *Devito*'s source injection mechanism to decompose the off-the-grid positioned wavefields to grid-aligned point wavefields. Using the *SID* structure, we perform an indirection and decomposition of the sources' wavefields to per-affected-point wavefields. The pseudocode for that workflow is presented in Listing 3. *src\_dcmp* now replaces *src* in our source injection computations. Instead of having sources at off-the-grid positions (Fig. 5a), we now have decomposed, aligned to the grid point sources (Fig. 5d).

**Listing 3:** Decomposing the source injection wavefields.

```

1 for t = 1 to nt do
2   foreach s in sources do
3     for i = 1 to np do
4       |xs, ys, zs = map(s, i);
5       |src_dcmp[t, SID[xs, ys, zs]] += f(src(t, s));

```

4) *Fuse iteration spaces*: using the aligned structure *src\_dcmp*, we can now fuse the source injection loop inside the kernel update iteration space. There is no *sources* loop as sparse data can be expressed in 3D coordinates. We fuse the source injection loop at the same loop level as the stencil update  $z$  loop. The source mask *SM* acts as a binary mask and is used to add (if 1) or not (if 0) the source impact while *SID* is used to access the impact values indirectly as we iterate over

the grid dimensions. The resulting loop structure is illustrated in the following pseudocode in Listing 4 and also offers SIMD vectorization opportunities over the  $z2$  loop.

**Listing 4:** Stencil kernel update with fused source injection.

```

1 for t = 1 to nt do
2   for x = 1 to nx do
3     for y = 1 to ny do
4       for z = 1 to nz do
5         |A(t, x, y, z, s);
6         for z2 = 1 to nz do
7           |u[t, x, y, z2] += SM[x, y, z2] * src_dcnp[t, SID[x, y, z2]];

```

5) *Reducing the iteration space size:* the 3D structures that are used to iterate through sources (SM and SID) in the  $z2$  loop are, in the general case, massively sparse. Multiplications by zero are dominant. Only the necessary iterations in  $z$  dimension need to be performed to alleviate this issue. We aggregate nonzero occurrences along the  $z$ -axis of SM recording them in a structure named `nnz_mask`. We reduce the size of SID cutting off  $z$ -slices where all elements are zero. For naming convention, we use `Sp_SID` for the new structure. These structures reduce the iteration space size of  $z2$  to perform the necessary computation. Pseudocode for the new structure is illustrated in Listing 5. The opportunity to reduce the iteration space generally applies to the majority of problems in seismic. However, we show that benefits are not limited in cases where the reduction is small (see subsection IV-E).

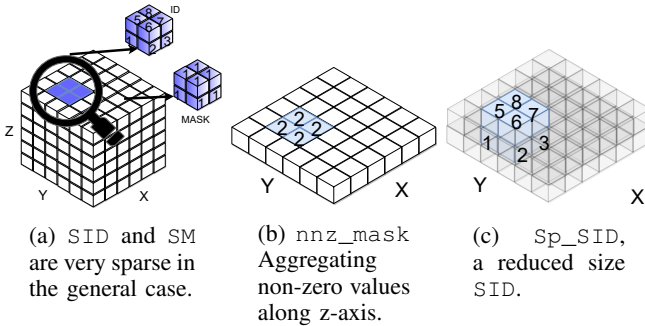


Fig. 6: We aggregate nonzero occurrences along the  $z$ -axis, keeping count of them. The size of SID is reduced by cutting off  $z$ -slices where all elements are zero.

**Listing 5:** Stencil kernel update with reduced size iteration space for source injection.

```

1 for t = 1 to nt do
2   for x = 1 to nx do
3     for y = 1 to ny do
4       for z = 1 to nz do
5         |A(t, x, y, z, s);
6         for z2 = 1 to nnz_mask[x][y] do
7           |I(t, x, y, z) ≡ { zind = Sp_SID[x, y, z2];
8           |u[t, x, y, z2] += src_dcnp[t, SID[x, y, zind]];

```

Finally, we present a methodology that aligns the source injection impact to the grid points, thus enabling the application of TB to stencil operators with sparse off-the-grid points.

## B. Applying wave-front temporal blocking

We present the loop transformations to apply WTB to the loop structure in Listing 5. In temporal blocking, we extend space blocking so that multiple timesteps are evaluated in a subset of the overall problem domain. In WTB space-time, wave-fronts traverse our domain computing grid point values. For naming convention, as used in [15], we are going to use the term “block” for spatial-only grouping and “tile” when multiple temporal updates are allowed. Fig. 7 shows grid point updates as they happen in WTB. The green point is updated using orange values. This stencil kernel has a radius of size two. Thus a margin of 2 points is required to preserve data dependencies. Only two timesteps are kept in memory (for time order one problems), so the green value substitutes the yellow one in the buffer. The stencil radius affects the wave-front angle (the ratio of spatial indices needed to update one point in the next time step). The angle gets steeper with a higher stencil radius. WTB can also be applied to staggered grids. In this case, two or more grids may be updated, often having inter-dependencies [15]. It is then necessary to shift the wave-front angle (allow more margin to preserve dependencies) by an amount, depending on the stencil radius of data dependencies in each loop as shown in Figure 8b.

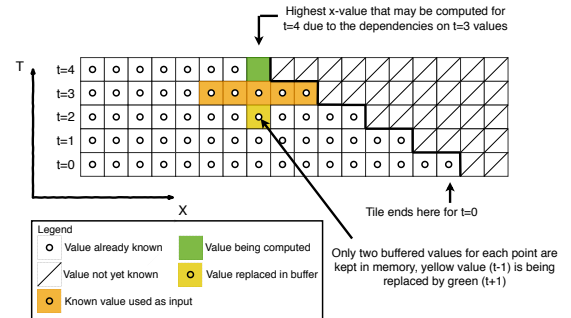
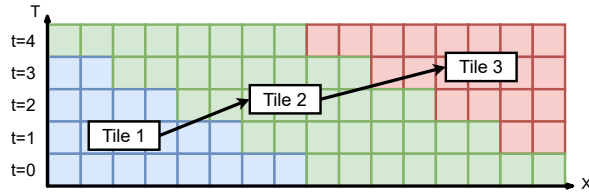


Fig. 7: Illustration of stencil kernel update in WTB. The green point is updated using the orange values. This stencil kernel has a space order of 4, thus allowing a margin of 2 points on the right in order not to violate data dependencies. Only two timesteps are kept in memory (for time order one problems), so the green value substitutes the yellow one. Figure partially adapted from [15].

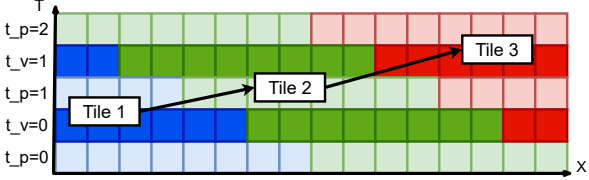
After precomputing source injection, data dependencies are now aligned with the computational grid points. Applying temporal blocking is now feasible. We split the time-space iteration space into tiles as shown in Figure 8a. Each tile is then partitioned into space blocks. By applying the transformations required from wave-front temporal blocking to Listing 4 we now have the loop structure 6. This structure is a time-tiled wave-front scheme over a stencil kernel update with source injection.

The next section provides details about the evaluated kernels, their data dependencies, and their inherent loop structure.





(a) The figure shows multiple wave-front tiles evaluated sequentially, partially adapted from [15].



(b) The figure shows multiple wave-front tiles evaluated sequentially in multigrid stencil codes.

Fig. 8: Wave-front updates for single- and multi-grid stencil updates.

**Listing 6:** The figure shows the loop structure after applying our proposed scheme.

```

1 for t_tile in time_tiles do
2   for xtile in xtiles do
3     for ytile in ytiles do
4       for t in t_tile do
5         OpenMP parallelism
6           for xblk in xtile do
7             for yblk in ytile do
8               for x in xblk do
9                 for y in yblk do
10                  SIMD vectorization
11                    for z = 1 to nz do
12                      A(t, x - time, y - time, z);
13                      for z2 = 1 to nnz_mask[x][y] do
14                        I(t, x - time, y - time, z2);

```

### III. STRUCTURE OF WAVE-PROPAGATION KERNELS

To illustrate our technique, we selected three representative kernels implementing explicit FD methods for wave propagation. The chosen kernels significantly differ in the operational intensity and working set size [55]. The kernels are implemented and validated in the *Devito* framework. The *Devito* compiler generates a C implementation for each kernel given a symbolic specification expressed with the *Devito* DSL.

#### A. Isotropic acoustic

The first equation we consider is the most straightforward and generally known wave-equation in an anisotropic acoustic medium. This equation is a single scalar PDE with a Jacobi-like stencil. The acoustic wave equation for the square slowness  $m$ , defined as  $m = \frac{1}{c^2}$ , where  $c$  is the speed of sound in the given physical media, and a source  $q$  is given by:

$$\begin{cases} m(x) \frac{\partial^2 u(t, x)}{\partial t^2} - \Delta u(t, x) = \delta(x_s) q(t) \\ u(0, \cdot) = \frac{\partial u(t, x)}{\partial t}(0, \cdot) = 0 \\ d(t, \cdot) = u(t, x_r). \end{cases} \quad (1)$$

where  $u(t, x)$  is the pressure wavefield,  $x_s$  is the point source position,  $q(t)$  is the source time signature,  $d(t, \cdot)$  is the measured data at positions  $x_r$  and  $m(x)$  is the squared slowness. This equation writes in few lines with the *Devito* symbolic API as follows:

---

```

from devito import solve, Eq, Operator
eq = m * u.dt2 - u.laplace
update = Eq(u.forward, solve(eq, u.forward))
src_eqns = s.inject(u.forward, expr=s*dt**2/m)
d_eqns = d.interpolate(u)

```

---

Listing 1: Wave-equation symbolic definition

The discretized acoustic wave-equation is generally memory-bound due to the low computational count of the standard Laplacian [55], [56].

#### B. Anisotropic acoustic

The second wave-equation kernel we consider is the most commonly used in industrial applications for subsurface imaging (RTM, FWI) [57]–[61]. This equation is a pseudo-acoustic anisotropic equation that consists of a coupled system of two scalar PDEs. Unlike the most simple acoustic isotropic equation, this formulation considers direction-dependent propagation speeds that translate into the discretized equation into a rotated anisotropic laplacian. Such a kernel increases the operation count drastically [55]. For example, the first dimension  $x$  component of the Laplacian is defined as:

$$\begin{aligned} G_{\bar{x}\bar{x}} &= D_{\bar{x}}^T D_{\bar{x}} \\ D_{\bar{x}} &= \cos(\theta) \cos(\phi) \frac{\partial}{\partial x} + \cos(\theta) \sin(\phi) \frac{\partial}{\partial y} - \sin(\theta) \frac{\partial}{\partial z}. \end{aligned} \quad (2)$$

where  $\theta$  is the (spatially dependent) tilt angle (rotation around  $z$ ),  $\phi$  is the (spatially dependent) azimuth angle (rotation around  $y$ ). A more detailed description of the physics and discretization can be found in [57], [58].

#### C. Isotropic elastic

Finally, we consider the isotropic elastic equation. Unlike the two previous acoustic approximations, this equation has two significant properties. First, this is a first-order system in time, which allows us to extend our work to a smaller range of local data dependency over time. Consequently, we demonstrate that the benefits of time-blocking and our implementation of it are not limited to a single pattern along the time dimension. Second, this equation is a coupled system of a vectorial and a tensorial PDE, which increases the data movement drastically (one or two versus nine state parameters) on the wavefield and contains non-scalar expressions of the source and receiver expressions that involve multiple wavefields.

The isotropic elastic wave-equation, parametrized by the Lamé parameters  $\lambda, \mu$  and the density  $\rho$ , is defined as [62]:

$$\begin{aligned} \frac{1}{\rho} \frac{\partial v}{\partial t} &= \nabla \cdot \tau \\ \frac{\partial \tau}{\partial t} &= \lambda \text{tr}(\nabla v) \mathbf{I} + \mu (\nabla v + (\nabla v)^T) \end{aligned} \quad (3)$$

where  $v$  is a vector-valued function with one component per cartesian direction, and the stress  $\tau$  is a symmetric second-order tensor-valued function.

In the following section, we consider these three wave-equations for varying spatial discretization orders to verify and analyze our temporal blocking method.

#### IV. EXPERIMENTAL EVALUATION

We outline in subsection IV-A the experimental setup followed for performance evaluation. We aim to demonstrate the performance improvement achieved by our approach, illustrate its potential for impact on key applications, and probe its applicability limits.

##### A. Compiler and system setup

To evaluate our scheme, we used Virtual Machines in Azure with two architectures: Intel® Xeon® Processor E5 v4 Family (formerly called Broadwell) and Intel® Xeon® Scalable Processors (formerly called Skylake 8171M). For our experiments, access was granted on VMs (on Microsoft Azure) called `Standard_E16s_v3` and `Standard_E32s_v3`, running Ubuntu 18.04.4. The first system, `E16s_v3`, has a single socket 8-core Intel Broadwell E5-2673 v4 CPUs with AVX2 support. Each Intel Broadwell CPU has three cache levels: L1 (32KB) and L2 (256KB) caches private to each core and a 50MB L3 cache shared per socket. The second system has single-socket 16-core Intel Skylake Platinum 8171M CPUs with AVX512 support. Each Intel Skylake CPU has three cache levels: L1 (32KB) and L2 (1MB) caches private to each core and a 35.75MB L3 cache shared per socket. Compilers used were GCC 7.5.0\* and ICC 2021.1. We used OpenMP shared-memory parallelism with dynamic scheduling and SIMD vectorization. Thread pinning was enabled using the environment variables `OMP_PROC_BIND` (for GCC) and `KMP_AFFINITY` (for ICC). Experiments were built with Devito v.4.2.3. The experimentation framework and instructions on reproducibility are available in subsection V-A.

##### B. Test case setup

We evaluate the performance of operators relevant to seismic imaging. We model the propagation of waves for three different models: *isotropic acoustic*, *anisotropic acoustic (TTI)* and *isotropic elastic*. The isotropic acoustic and TTI wave equations are discretized with second order in time while isotropic elastic with first-order, and we study varying space orders of 4, 8, 12. For all test cases, we use zero initial conditions and damping fields with absorbing boundary layers. Waves are injected using one time-dependent, spatially localized seismic source wavelet into the subsurface model. We benchmark velocity models of  $512^3$  grid points, with

a grid spacing of 10 for isotropic and elastic and 20 for TTI. Wave propagation is modeled in single precision for 512ms, resulting in 228 time-steps for isotropic acoustic, 436 for isotropic elastic, and 587 for anisotropic acoustic. The time-stepping interval is selected regarding the Courant-Friedrichs-Lewy (CFL) condition [63], ensuring the explicit time-stepping scheme’s stability determined by the highest velocity of the subsurface model and the grid spacing.

##### C. Autotuning temporally blocked code

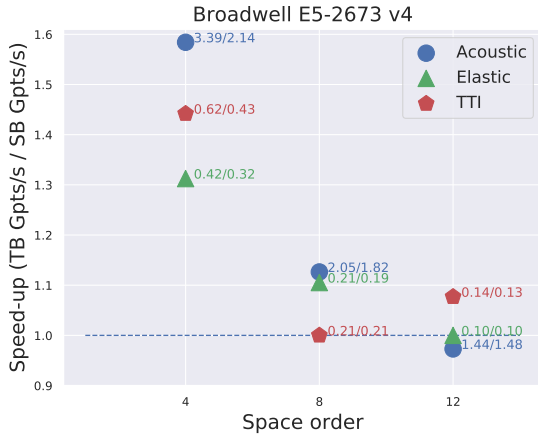
It should be noted that the parameter space for temporal blocking schemes is extensive. We report results obtained from guidance and experience from the state-of-the-art codes and literature [64]. Simulation codes are hard to generalize in terms of performance as multiple configurations may be used from case to case. An operator’s performance depends upon many factors, such as grid shape, discretization space order, tile and block shapes, number of other fields, number of timesteps, platforms, and others. To tune our C code for the underlying hardware, we swept over the whole parameter space to find the global performance maxima. We executed our experiments using the best-performing tile and block sizes, ensuring a fair comparison versus Devito’s aggressively tuned optimized spatially-blocked and vectorized code. The best performing tile sizes for temporal blocking are reported in Table I. Figure 9 illustrates the throughput (GPoints/s) speedup achieved for each evaluated model for space order discretizations of 4, 8, and 12.

Problem	$tile_x, tile_y, block_x, block_y$	
	Broadwell	Skylake
Acoustic O(2,4)	32, 32, 8, 8	64, 64, 8, 8
Acoustic O(2,8)	64, 64, 8, 8	64, 64, 8, 8
Acoustic O(2,12)	256, 256, 8, 8	128, 128, 8, 8
Elastic O(1,4)	32, 32, 8, 8	32, 32, 8, 8
Elastic O(1,8)	32, 32, 8, 8	64, 56, 8, 12
Elastic O(1,12)	256, 256, 8, 8	256, 256, 8, 8
TTI O(2,4)	40, 32, 4, 4	48, 48, 8, 8
TTI O(2,8)	32, 32, 8, 8	64, 64, 8, 8
TTI O(2,12)	256, 256, 8, 8	256, 256, 8, 8

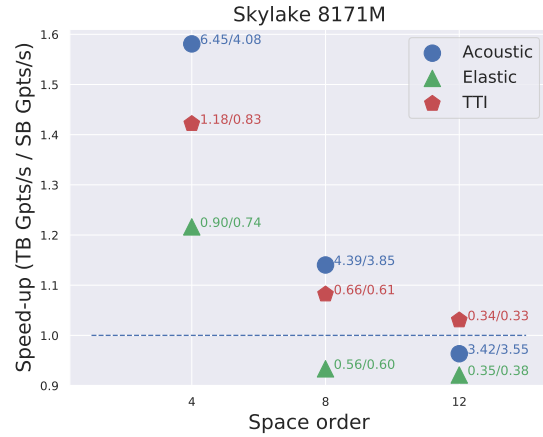
TABLE I: Optimal tile-block shapes after tuning WTB

##### D. Results discussion

Figure 9 illustrates the speedup achieved for each of our wave propagation models. All of our models show speedup for space order four discretization on both platforms. Acoustic benefits the most with around 1.6x and TTI follows with around 1.44x. Elastic wave propagation is accelerated by 1.3x on Broadwell and 1.22x on Skylake. Concerning space order 8, a commonly used practice, we observe speedups of 1.13x or more for acoustic, elastic on Broadwell and acoustic, TTI on Skylake. No significant performance gains are observed for space order 12, excluding some gains of around 5% on Broadwell with isotropic elastic and TTI. Figure 11 shows the isotropic acoustic kernels’ roofline performance for the Broadwell microarchitecture. The roofline is a cache-aware roofline



(a) Throughput speed-up of kernels for Broadwell.



(b) Throughput speed-up of kernels for Skylake.

Fig. 9: Throughput speed-up of temporal blocking kernels versus highly-optimised vectorized spatially-blocked code in Intel Xeon architectures, Broadwell and Skylake.

model representing cumulative (L1+L2+LLC+DRAM) traffic-based Arithmetic Intensity for application kernels <sup>1</sup>. We showcase improvement for the acoustic model breaking the ceiling of the L3 cache.

### E. Corner cases

Although our test cases use a single source, it is interesting to explore how our model performs with the presence of more off-the-grid operators. Each source is decomposed into its surrounding grid points, so the overhead increases due to the number of initial sources and the number of grid points affected. We evaluate the overhead induced for two cases a) an increasing number of sparsely located sources: in this case, we have an increasing number of sources located at an x-y plane slice of the 3D grid, a scenario which can be of practical interest and b) an increasing number of sources densely and uniformly located all over the 3D grid. Figure 10 shows that for isotropic acoustic wave propagation, the increasing number of sources is not affecting performance gains except with really densely located sources where our scheme is not taking advantage of the structure sparsity. Still, though, we observe gains of around 1.4x compared to 1.55x previously.

## V. CONCLUSIONS

This paper introduced a mechanism to enable temporal blocking in stencil computations involving sparse off-the-grid operators as encountered, for example, with sources and receivers in seismic inversion problems. We applied wave-front temporal blocking to wave-propagators ranging from isotropic acoustic to more advanced, such as isotropic elastic and anisotropic acoustic (TTI). Experimental evaluation of the improved kernels on Broadwell and Skylake microarchitectures showed compelling evidence of substantial acceleration

<sup>1</sup><https://software.intel.com/content/www/us/en/develop/articles/integrated-roofline-model-with-intel-advisor.html>

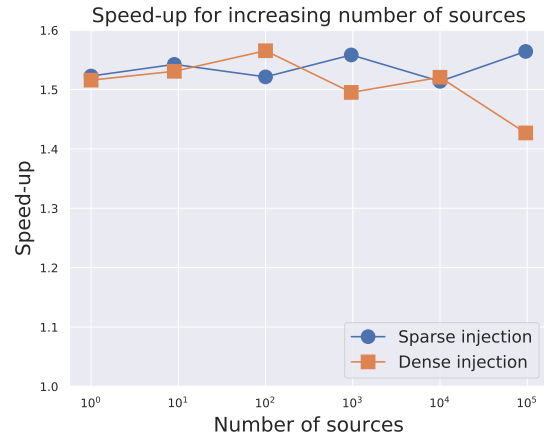


Fig. 10: Throughput speed-up for an isotropic acoustic operator for space order 4 over an increasing number of sources, sparsely and densely located.

of at least 1.5x for low and at least 1.1x for medium space order wave-propagation kernels.

### A. Code availability

An implementation of the methods described in this paper is available in a [Devito](#) fork repository under the MIT open-source license at [georgebisbas/devito v0.9-alpha](#). See the [README.md](#) for instructions on how to reproduce the results in the paper.

### B. Future work

Achieving performance improvement with high-space order kernels requires further research work. Methods like stencil retiming [65] have shown promise in alleviating this performance bottleneck, and a possible combination with temporal blocking may be promising. Another possible solution can



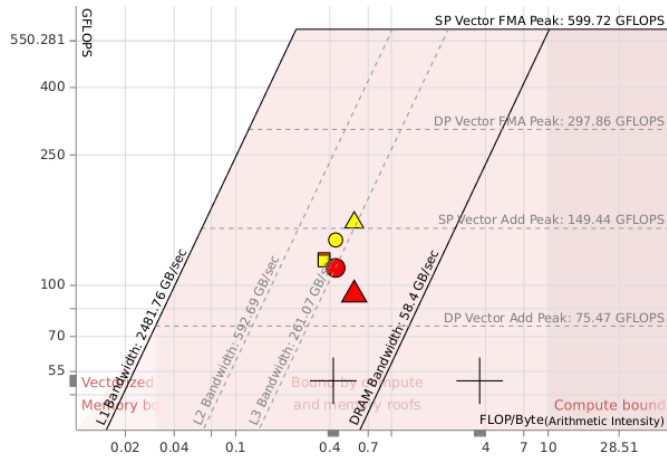


Fig. 11: Cache-aware roofline model on Broadwell for isotropic acoustic model space order 4 (triangles), 8 (circles), and 12 (squares). Red markers show the performance of spatially blocked vectorized kernels, while yellow ones show our temporal blocking scheme's performance.

be data layout transformations [48]. Near-term plans include evaluating our scheme on more diverse architectures (e.g., ARM) and accelerators (e.g., GPUs). The next step is the full automation and integration in the *Devito* DSL [3]. We aim to deliver automated, scalable optimizations on generated code beyond our kernels' current roofline performance limit. This paper's evaluation results are mainly motivated by the seismic imaging domain; however, the target applications are not limited to this scope.

#### ACKNOWLEDGMENTS

This research is funded by the Engineering and Physical Sciences Research Council (EPSRC) grants EP/L016796/1, EP/R029423/1, EP/V001493/1, and HiPEDS Center for Doctoral Training. The author thanks Richard M. Veras, Navjot Kukreja, John Washbourne, and Giacomo La Scala for the fruitful discussions as well as the whole *Devito* community.

#### REFERENCES

- [1] J. Ragan-Kelley, A. Adams, D. Sharlet, C. Barnes, S. Paris, M. Levoy, S. Amarasinghe, and F. Durand, "Halide: decoupling algorithms from schedules for high-performance image processing," *Communications of the ACM*, vol. 61, no. 1, pp. 106–115, 2017.
- [2] T. Gysi, C. Osuna, O. Fuhrer, M. Bianco, and T. C. Schulthess, "Stella: A domain-specific tool for structured grid methods in weather and climate models," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '15. New York, NY, USA: Association for Computing Machinery, 2015.
- [3] F. Luporini, M. Louboutin, M. Lange, N. Kukreja, P. Witte, J. Hüchelheim, C. Yount, P. H. J. Kelly, F. J. Herrmann, and G. J. Gorman, "Architecture and performance of devito, a system for automated stencil computation," *ACM Trans. Math. Softw.*, vol. 46, no. 1, Apr. 2020.
- [4] M. Louboutin, M. Lange, F. Luporini, N. Kukreja, P. A. Witte, F. J. Herrmann, P. Velesko, and G. J. Gorman, "Devito (v3.1.0): an embedded domain-specific language for finite differences and geophysical exploration," *Geoscientific Model Development*, vol. 12, no. 3, pp. 1165–1187, 2019.

- [5] L. Guasch, O. C. Agudo, M.-X. Tang, P. Nachev, and M. Warner, "Full-waveform inversion imaging of the human brain," *NPJ digital medicine*, vol. 3, no. 1, pp. 1–12, 2020.
- [6] D. Wonnacott, "Using time skewing to eliminate idle time due to memory bandwidth and network limitations," *Proceedings 14th International Parallel and Distributed Processing Symposium. IPDPS 2000*, pp. 171–180, May 2000.
- [7] Guohua Jin, J. Mellor-Crummey, and R. Fowler, "Increasing temporal locality with skewing and recursive blocking," in *SC '01: Proceedings of the 2001 ACM/IEEE Conference on Supercomputing*, 2001, pp. 57–57.
- [8] D. Wonnacott, "Achieving scalable locality with time skewing," *International Journal of Parallel Programming*, vol. 30, pp. 181–221, 2004.
- [9] R. J. LeVeque, *Finite difference methods for ordinary and partial differential equations: steady-state and time-dependent problems*. Siam, 2007, vol. 98.
- [10] G. Smith, *Numerical Solution of Partial Differential Equations: Finite Difference Methods*. Clarendon, 1985.
- [11] H. Yan, J. Xu, and X. Zhang, "Compressed sensing radar imaging of off-grid sparse targets," in *2015 IEEE Radar Conference (RadarCon)*, 2015, pp. 0690–0693.
- [12] C. Poon and G. Peyré, "Degrees of freedom for off-the-grid sparse estimation," *arXiv preprint arXiv:1911.03577*, 2019.
- [13] M. Wolfe, "More iteration space tiling," *Proceedings of the 1989 ACM/IEEE conference on Supercomputing - Supercomputing '89*, pp. 655–664, 1989.
- [14] M. E. Wolf and M. S. Lam, "A data locality optimizing algorithm," *ACM SIGPLAN Notices*, vol. 26, no. 6, pp. 30–44, 1991.
- [15] C. Yount and A. Duran, "Effective use of large high-bandwidth memory caches in hpc stencil computation via temporal wave-front tiling," in *2016 7th International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS)*, 2016, pp. 65–75.
- [16] G. Wellein, G. Hager, T. Zeiser, M. Wittmann, and H. Fehske, "Efficient temporal blocking for stencil computations by multicore-aware wave-front parallelization," *Proceedings - International Computer Software and Applications Conference*, vol. 1, pp. 579–586, 2009.
- [17] I. J. Bertolacci, C. Olschanowsky, B. Harshbarger, B. L. Chamberlain, D. G. Wonnacott, and M. M. Strout, "Parameterized diamond tiling for stencil computations with chapel parallel iterators," in *Proceedings of the 29th ACM International Conference on Supercomputing*, ser. ICS '15. New York, NY, USA: ACM, 2015, pp. 197–206.
- [18] T. Malas, G. Hager, H. Ltaief, H. Stengel, G. Wellein, and D. Keyes, "Multicore-optimized wavefront diamond blocking for optimizing stencil updates," *SIAM Journal on Scientific Computing*, vol. 37, no. 4, pp. C439–C464, 2015.
- [19] U. Bondhugula, A. Hartono, J. Ramanujam, and P. Sadayappan, "A practical automatic polyhedral parallelizer and locality optimizer," in *ACM SIGPLAN Notices*, vol. 43, no. 6. ACM, 2008, pp. 101–113.
- [20] U. Bondhugula, "Compiling affine loop nests for distributed-memory parallel architectures," in *SC'13: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. IEEE, 2013, pp. 1–12.
- [21] T. Gresser, A. GröBlinger, and C. Lengauer, "Polly - performing polyhedral optimizations on a low-level intermediate representation," *Parallel Process. Lett.*, vol. 22, 2012.
- [22] A. Klöckner, "Loo.py: transformation-based code generation for GPUs and CPUs," in *Proceedings of ARRAY '14: ACM SIGPLAN Workshop on Libraries, Languages, and Compilers for Array Programming*. Edinburgh, Scotland: Association for Computing Machinery, 2014.
- [23] C. Bastoul, "Code generation in the polyhedral model is easier than you think," in *Proceedings of the 13th International Conference on Parallel Architectures and Compilation Techniques*. IEEE Computer Society, 2004, pp. 7–16.
- [24] J. Ramanujam, "Iteration Spaces for Multicomputers," *Computer*, vol. i, 1992.
- [25] J. Xue, "On tiling as a loop transformation," *Parallel Processing Letters*, vol. 7, no. 04, pp. 409–424, 1997.
- [26] M. Frigo and V. Strumpfen, "Cache oblivious stencil computations," in *Proceedings of the 19th Annual International Conference on Supercomputing*, ser. ICS '05. New York, NY, USA: ACM, 2005, pp. 361–366.
- [27] —, "The cache complexity of multithreaded cache oblivious algorithms," in *Proceedings of the Eighteenth Annual ACM Symposium on Parallelism in Algorithms and Architectures*, ser. SPAA '06. New York, NY, USA: Association for Computing Machinery, 2006, p. 271–280.

- [28] K. Datta, M. Murphy, V. Volkov, S. Williams, J. Carter, L. Oliker, D. Patterson, J. Shalf, and K. Yelick, "Stencil computation optimization and auto-tuning on state-of-the-art multicore architectures," in *Proceedings of the 2008 ACM/IEEE conference on Supercomputing*. IEEE Press, 2008, p. 4.
- [29] Y. Tang, R. A. Chowdhury, B. C. Kuszmaul, C.-K. Luk, and C. E. Leiserson, "The Pochair stencil compiler," *Proceedings of the 23rd ACM symposium on Parallelism in algorithms and architectures - SPAA '11*, p. 117, 2011.
- [30] M. Frigo, C. E. Leiserson, H. Prokop, and S. Ramachandran, "Cache-oblivious algorithms," *ACM Trans. Algorithms*, vol. 8, no. 1, Jan. 2012.
- [31] F. Luporini, M. Lange, C. T. Jacobs, G. J. Gorman, J. Ramanujam, and P. H. Kelly, "Automated tiling of unstructured mesh computations with application to seismological modeling," *ACM Transactions on Mathematical Software (TOMS)*, vol. 45, no. 2, p. 17, 2019.
- [32] T. Grosser, A. Cohen, P. H. J. Kelly, J. Ramanujam, P. Sadayappan, and S. Verdoolaege, "Split tiling for gpus: Automatic parallelization using trapezoidal tiles," in *Proceedings of the 6th Workshop on General Purpose Processor Using Graphics Processing Units*, ser. GPGPU-6. New York, NY, USA: Association for Computing Machinery, 2013, p. 24–31.
- [33] T. Grosser, A. Cohen, J. Holewinski, P. Sadayappan, and S. Verdoolaege, "Hybrid hexagonal/classical tiling for gpus," in *Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization*, ser. CGO '14. New York, NY, USA: Association for Computing Machinery, 2014, p. 66–75.
- [34] A. Schäfer and D. Fey, "High performance stencil code algorithms for GPGPUs," *Procedia Computer Science*, vol. 4, pp. 2027–2036, 2011.
- [35] J. Holewinski, L.-N. Pouchet, and P. Sadayappan, "High-performance code generation for stencil computations on GPU architectures," *Proceedings of the 26th ACM international conference on Supercomputing - ICS '12*, p. 311, 2012.
- [36] P. S. Rawat, M. Vaidya, A. Sukumaran-Rajam, M. Ravishankar, V. Grover, A. Rountev, L. N. Pouchet, and P. Sadayappan, "Domain-Specific Optimization and Generation of High-Performance GPU Code for Stencil Computations," *Proceedings of the IEEE*, vol. 106, no. 11, pp. 1902–1920, 2018.
- [37] H. R. Zohouri, A. Podobas, and S. Matsuoka, "Combined spatial and temporal blocking for high-performance stencil computation on fpgas using openc1," in *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, ser. FPGA '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 153–162.
- [38] M. Wolfe, "Loop skewing: The wavefront method revisited," *Int. J. Parallel Program.*, vol. 15, no. 4, p. 279–293, Oct. 1986.
- [39] R. Strzodka, M. Shaheen, D. Pajak, and H.-P. Seidel, "Cache accurate time skewing in iterative stencil computations," in *Proceedings of the 2011 International Conference on Parallel Processing*, ser. ICPP '11. Washington, DC, USA: IEEE Computer Society, 2011, pp. 571–581.
- [40] V. Bandishti, I. Pananilath, and U. Bondhugula, "Tiling stencil computations to maximize parallelism," in *SC'12: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. IEEE, 2012, pp. 1–11.
- [41] T. Muranushi and J. Makino, "Optimal temporal blocking for stencil computation," *Procedia Computer Science*, vol. 51, pp. 1303–1312, 2015.
- [42] V. Levchenko and A. Perepelkina, "The diamondtetris algorithm for maximum performance vectorized stencil computation," in *Parallel Computing Technologies*, V. Malyshev, Ed. Cham: Springer International Publishing, 2017, pp. 124–135.
- [43] K. Akbudak, H. Ltaief, V. Etienne, R. Abdelkhalak, T. Tonellot, and D. Keyes, "Asynchronous computations for solving the acoustic wave propagation equation," *The International Journal of High Performance Computing Applications*, vol. 34, pp. 377 – 393, 2020.
- [44] A. Logg, K.-A. Mardal, and G. Wells, *Automated solution of differential equations by the finite element method: The FEniCS book*. Springer Science & Business Media, 2012, vol. 84.
- [45] F. Rathgeber, D. A. Ham, L. Mitchell, M. Lange, F. Luporini, A. T. McRae, G.-T. Bercea, G. R. Markall, and P. H. Kelly, "Firedrake: automating the finite element method by composing abstractions," *ACM Transactions on Mathematical Software (TOMS)*, vol. 43, no. 3, p. 24, 2017.
- [46] C. T. Jacobs, S. P. Jammy, and N. Sandham, "Opensbli: A framework for the automated derivation and parallel execution of finite difference solvers on a range of computer architectures," *J. Comput. Sci.*, vol. 18, pp. 12–23, 2017.
- [47] I. Z. Reguly, G. R. Mudalige, M. B. Giles, D. Curran, and S. McIntosh-Smith, "The ops domain specific abstraction for multi-block structured grid computations," in *2014 Fourth International Workshop on Domain-Specific Languages and High-Level Frameworks for High Performance Computing*. IEEE, 2014, pp. 58–67.
- [48] C. Yount, "Vector folding: improving stencil performance via multi-dimensional simd-vector representation," in *2015 IEEE 17th International Conference on High Performance Computing and Communications, 2015 IEEE 7th International Symposium on CyberSpace Safety and Security, and 2015 IEEE 12th International Conference on Embedded Software and Systems*. IEEE, 2015, pp. 865–870.
- [49] B. Hagedorn, L. Stoltzfus, M. Steuer, S. Gorlatch, and C. Dubach, "High performance stencil code generation with lift," in *Proceedings of the 2018 International Symposium on Code Generation and Optimization*. ACM, 2018, pp. 100–112.
- [50] F. Thaler, S. Moosbrugger, C. Osuna, M. Bianco, H. Vogt, A. Afanasyev, L. Mosimann, O. Fuhrer, T. C. Schulthess, and T. Hoefler, "Porting the cosmo weather model to manycore cpus," in *Proceedings of the Platform for Advanced Scientific Computing Conference*, ser. PASC '19. New York, NY, USA: Association for Computing Machinery, 2019.
- [51] J. Virieux and S. Operto, "An overview of full-waveform inversion in exploration geophysics," *Geophysics*, vol. 74, no. 6, pp. WCC1–WCC26, 2009.
- [52] E. Baysal, D. D. Kosloff, and J. W. Sherwood, "Reverse time migration," *Geophysics*, vol. 48, no. 11, pp. 1514–1524, 1983.
- [53] J. Ramanujam and P. Sadayappan, "Tiling multidimensional iteration spaces for nonshared memory machines," in *Supercomputing '91: Proceedings of the 1991 ACM/IEEE Conference on Supercomputing*, 1991, pp. 111–120.
- [54] L. Lamport, "The parallel execution of do loops," *Commun. ACM*, vol. 17, no. 2, p. 83–93, Feb. 1974.
- [55] M. Louboutin, M. Lange, F. J. Herrmann, N. Kukreja, and G. Gorman, "Performance prediction of finite-difference solvers for different computer architectures," *Computers and Geosciences*, vol. 105, pp. 148–157, 2017.
- [56] S. Williams, A. Waterman, and D. Patterson, "Roofline: an insightful visual performance model for multicore architectures," *Communications of the ACM*, vol. 52, no. 4, p. 65, 2009.
- [57] Y. Zhang, H. Zhang, and G. Zhang, "A stable tti reverse time migration and its implementation," *Geophysics*, vol. 76, no. 3, pp. WA3–WA11, 2011.
- [58] M. Louboutin, P. Witte, and F. J. Herrmann, "Effects of wrong adjoints for rtm in tti media," in *SEG Technical Program Expanded Abstracts 2018*. Society of Exploration Geophysicists, 2018, pp. 331–335.
- [59] E. Duvencek and P. M. Bakker, "Stable p-wave modeling for reverse-time migration in tilted ti media," *GEOPHYSICS*, vol. 76, no. 2, pp. S65–S75, 2011.
- [60] T. Alkhalifah, "An acoustic wave equation for anisotropic media," *Geophysics*, vol. 65, pp. 1239–1250, 2000.
- [61] K. Bube, J. Washbourne, R. Ergas, and T. Nemeth, "Self-adjoint, energy-conserving second-order pseudoacoustic systems for vti and tti media for reverse time migration and full-waveform inversion," in *SEG Technical Program Expanded Abstracts 2016*. Society of Exploration Geophysicists, 2016, pp. 1110–1114.
- [62] J. Virieux, "P-sv wave propagation in heterogeneous media: Velocity-stress finite-difference method," *Geophysics*, vol. 51, no. 4, pp. 889–901, 1986.
- [63] R. Courant, K. Friedrichs, and H. Lewy, "On the partial difference equations of mathematical physics," *IBM Journal of Research and Development*, vol. 11, no. 2, pp. 215–234, 1967.
- [64] M. Wittmann, G. Hager, and G. Wellein, "Multicore-aware parallel temporal blocking of stencil codes for shared and distributed memory," in *2010 IEEE International Symposium on Parallel Distributed Processing, Workshops and Phd Forum (IPDPSW)*, 2010, pp. 1–7.
- [65] K. Stock, M. Kong, T. Grosser, L.-N. Pouchet, F. Rastello, J. Ramanujam, and P. Sadayappan, "A framework for enhancing data reuse via associative reordering," in *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI '14. New York, NY, USA: Association for Computing Machinery, 2014, p. 65–76.