

1 **ARCHITECTURE AND PERFORMANCE OF DEVITO, A SYSTEM**
2 **FOR AUTOMATED STENCIL COMPUTATION ***

3 FABIO LUPORINI[†], MICHAEL LANGE[‡], MATHIAS LOUBOUTIN[§], NAVJOT KUKREJA[†],
4 JAN HÜCKELHEIM[†], CHARLES YOUNT[¶], PHILIPP WITTE^{||}, PAUL H. J. KELLY[#],
5 GERARD J. GORMAN[†], AND FELIX J. HERRMANN[§]

6 **Abstract.** Stencil computations are a key part of many high-performance computing applica-
7 tions, such as image processing, convolutional neural networks, and finite-difference solvers for partial
8 differential equations. Devito is a framework capable of generating highly-optimized code given sym-
9 bolic equations expressed in *Python*, specialized in, but not limited to, affine (stencil) codes. The
10 lowering process – from mathematical equations down to C++ code – is performed by the Devito
11 compiler through a series of intermediate representations. Several performance optimizations are in-
12 troduced, including advanced common sub-expressions elimination, tiling and parallelization. Some
13 of these are obtained through well-established stencil optimizers, integrated in the back-end of the
14 Devito compiler. The architecture of the Devito compiler, as well as the performance optimiza-
15 tions that are applied when generating code, are presented. The effectiveness of such performance
16 optimizations is demonstrated using operators drawn from seismic imaging applications.

17 **Key words.** Stencil, finite difference method, symbolic processing, structured grid, compiler,
18 performance optimization

19 **AMS subject classifications.** 65N06, 68N20

20 **1. Introduction.** Developing software for high-performance computing requires
21 a considerable interdisciplinary effort, as it often involves domain knowledge from nu-
22 merous fields such as physics, numerical analysis, software engineering and low-level
23 performance optimization. The result is typically a monolithic application where
24 hardware-specific optimizations, numerical methods, and physical approximations are
25 interwoven and dispersed throughout a large number of loops, functions, files and mod-
26 ules. This frequently leads to slow innovation, high maintenance costs, and code that
27 is hard to debug and port onto new computer architectures. A powerful approach to
28 alleviate this problem is to introduce a separation of concerns and to raise the level of
29 abstraction by using domain-specific languages (DSLs). DSLs can be used to express
30 numerical methods using a syntax that closely mirrors how they are expressed math-
31 ematically, while a stack of compilers and libraries is responsible for automatically

*Submitted to SIAM Journal on Scientific Computing on July 9, 2018.

Funding: This work was supported by the Engineering and Physical Sciences Research Council through grants EP/I00677X/1, EP/L000407/1, EP/I012036/1], by the Imperial College London Department of Computing, by the Imperial College London Intel Parallel Computing Centre (IPCC), and by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Applied Mathematics and Computer Science programs under contract number DE-AC02-06CH11357.

[†]Department of Earth Science and Engineering, Imperial College London, London, UK, (f.luporini12@imperial.ac.uk, n.kukreja@imperial.ac.uk, j.hueckelheim@imperial.ac.uk, g.gorman@imperial.ac.uk)

[‡]European Centre for Medium-Range Weather Forecasts, Reading, UK, (michael.lange@ecmwf.int)

[§]Georgia Institute of Technology, School of Computational Science and Engineering, Atlanta GA, USA, (mlouboutin3@gatech.edu, felix.herrmann@gatech.edu)

[¶]Intel Corporation, (chuck.yount@intel.com)

^{||}Seismic Laboratory for Imaging and Modeling (SLIM), The University of British Columbia, Vancouver BC, CANADA, (pwitte.slim@gmail.com)

[#]Department of Computing, Imperial College London, London, SW7 2AZ, UK, (p.kelly@imperial.ac.uk)

32 creating the optimized low-level implementation in a general purpose programming
33 language such as C++. While the focus of this paper is on finite-difference (FD) based
34 codes, the DSL approach has already had remarkable success in other numerical meth-
35 ods such as the finite-element (FE) and finite-volume (FV) method, as documented
36 in Section 2.

37 This work describes the architecture of *Devito*, a system for automated stencil
38 computations from a high-level mathematical syntax. Devito was developed with an
39 emphasis on FD methods on structured grids. For this reason, Devito’s underlying
40 DSL has many features to simplify the specification of FD methods, as discussed
41 in Section 3. The original motivation was to solve large-scale partial differential
42 equations (PDEs) in the context of seismic inverse problems, where FD solvers are
43 commonly used for solving wave equations as part of complex workflows (e.g., data
44 inversion using adjoint-state methods and backpropagation). Devito is equally useful
45 as a framework for other stencil computations in general; for example, computations
46 where all array indices are affine functions of loop variables. The Devito compiler
47 is also capable of generating arbitrarily nested, possibly irregular, loops. This key
48 feature is needed to support many complex algorithms that are used in engineering and
49 scientific practice, including applications from image processing, cellular automata,
50 and machine-learning.

51 One of the design goals of Devito was to enable high-productivity, so it is fully
52 written in *Python*, with easy access to solvers, optimizers, input and output, and the
53 wide range of other libraries in the *Python* ecosystem. At the same time, Devito
54 transforms high-level symbolic input into optimized C++ code, resulting in a perfor-
55 mance that is competitive with hand-optimized implementations. While the examples
56 presented in this paper focus on using Devito from a *Python* application, exploiting
57 the full potential of on-the-fly code generation and just-in-time (JIT) compilation,
58 a practical advantage of generating C++ as an intermediate step is that it can be
59 also used to generate libraries for legacy software, thus enabling incremental code
60 modernisation.

61 Compared to other DSL frameworks that are used in practice, Devito uses com-
62 piler technology, including several layers of intermediate representations, to perform
63 optimizations in multiple passes. This allows Devito to perform more complex op-
64 timizations, and to better optimize the code for individual target platforms. The
65 fact that these optimisations are performed programmatically facilitates performance
66 portability across different computer architectures [28]. This is important, as indus-
67 trial codes are often used on a variety of platforms, including clusters with multi-core
68 CPUs, GPUs, and many-core chips spread across several compute nodes as well as
69 various cloud platforms. Devito also performs high-level transformations for floating-
70 point operation (FLOP) reduction based on symbolic manipulation, as well as loop-
71 level optimizations as implemented in Devito’s own optimizer, or using a third-party
72 stencil compiler such as YASK [40]. The Devito compiler is presented in detail in
73 Sections 4, 5 and 6.

74 After the presentation of the Devito compiler, we show test cases in Section 7
75 that are inspired by real-world seismic-imaging problems. The paper finishes with
76 directions for future work and conclusions in Sections 8 and 9.

77 **2. Related work.** The objective of maximizing productivity and performance
78 through frameworks based upon DSLs has long been pursued. In addition to well-
79 known systems such as Mathematica[®] and Matlab[®], which span broad mathematical
80 areas, there are a number of tools specialized in numerical methods for PDEs, some

81 dating back to the 1970s [6, 34, 7, 35].

82 **2.1. DSL-based frameworks for partial differential equations.** One note-
83 worthy contemporary framework centered on DSLs is FEniCS [22], which allows
84 the specification of weak variational forms, via UFL [2], and finite-element meth-
85 ods, through a high-level syntax. Firedrake [30] implements the same languages as
86 FEniCS, although it differs from it in a number of features and architectural choices.
87 Devito is heavily influenced by these two successful projects, in particular by their
88 philosophy and design. Since solving a PDE is often a small step of a larger workflow,
89 the choice of *Python* to implement these software provides access to a wide ecosystem
90 of scientific packages. Firedrake also follows the principle of graceful degradation, by
91 providing a very simple lower-level API to escape the abstraction when non-standard
92 calculations (i.e., unrelated to the finite-element formulation) are required. Likewise,
93 Devito allows injecting arbitrary expressions into the finite-difference specification;
94 this feature has been used in real-life cases, for example for interpolation in seismic
95 imaging operators. On the other hand, a major difference is that Devito lacks a for-
96 mal specification language such as UFL in FEniCS/Firedrake. This is partly because
97 there is no systematic foundation underpinning FD, as opposed to FE which relies
98 upon the theory of Hilbert spaces [5]. Yet another distinction is that, for performance
99 reasons, Devito takes control of the time-stepping loop. Other examples of embedded
100 DSLs are provided by the OpenFOAM project, with a language for FV [13], and by
101 PyFR, which targets flux reconstruction methods [36].

102 **2.2. High-level approaches to finite differences.** Due to its simplicity, the
103 FD method has been the subject of multiple research projects, chiefly targeting the
104 design of effective software abstraction and/or the generation of high performance code
105 [14, 3, 16, 21]. Devito distinguishes itself from previous work in a number of ways
106 including: support for the principle of graceful degradation for when the DSL does not
107 cover a feature required by an application; incorporation of a symbolic mathematics
108 engine; using actual compiler technology rather than template-based code generation;
109 adoption of a native *Python* interface that naturally allows composition into complex
110 workflows such as optimisation and machine-learning frameworks.

111 At a lower level of abstraction there are a number of tools targeting “stencil”
112 computation (FD codes belong to this class), whose major objective is the generation
113 of efficient code. Some of them provide a DSL [40, 31, 43, 29], whereas others are
114 compilers or user-driven code generation systems, often based upon a polyhedral
115 model, such as [4, 18]. From the Devito standpoint, the aim is to harness these
116 tools – for example by integrating them, to maximize performance portability. As a
117 proof of concept, we shall discuss the integration of one such tool, namely YASK [40],
118 with Devito.

119 **2.3. Devito and seismic imaging.** Devito is a general purpose system, not
120 restricted to specific PDEs, so it can be used for any form of the wave equation.
121 Thus, unlike software specialized in seismic exploration, like IWAVE [32] and Mada-
122 gascar [12], it suffers neither from the restriction to a small set of wave equations and
123 discretizations, nor from the lack of portability and composability typical of a pure
124 C/Fortran environment.

125 **2.4. Performance optimizations.** The Devito compiler can introduce three
126 types of performance optimizations: FLOPs reduction, data locality, and parallelism.
127 Typical FLOPs reduction transformations are common sub-expressions elimination,
128 factorization, and code motion. A thorough review is provided in [11]. To different

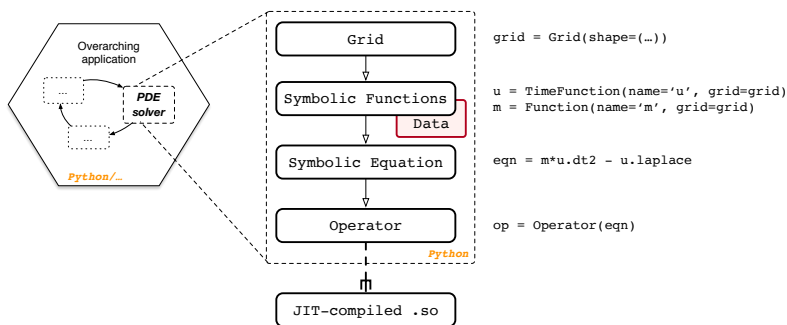


Fig. 1: The typical usage of Devito within a larger application.

129 extent, Devito applies all of these techniques (see Section 5.1). Particularly relevant
 130 for stencil computation is the search for redundancies across consecutive loop iterations
 131 [9, 10, 20]. This is at the core of the strategy described in Section 6, which
 132 essentially extends these ideas with optimizations for data locality. Typical loop trans-
 133 formations for parallelism and data locality [17] are also automatically introduced by
 134 the Devito compiler (e.g., loop blocking, vectorization); more details will be provided
 135 in Sections 5.2 and 5.3.

136 **3. Specification of a finite-difference method with Devito.** The Devito
 137 DSL allows concise expression of FD and general stencil operations using a mathe-
 138 matical notation. It uses *SymPy* [27] for the specification and manipulation of stencil
 139 expressions. In this section, we describe the use of Devito’s DSL to build PDE solvers.
 140 Although the examples used here are for FD, the DSL can describe a large class of op-
 141 erations, such as convolutions or basic linear algebra operations (e.g., chained tensor
 142 multiplications).

143 **3.1. Symbolic types.** The key steps to implement a numerical kernel with De-
 144 vito are shown in Figure 1. We describe this workflow, as well as fundamental features
 145 of the Devito API, using the acoustic wave equation, also known as d’Alembertian or
 146 Box operator. Its continuous form is given by:

$$\begin{aligned}
 147 \quad (3.1) \quad m(x, y, z) \frac{d^2 u(x, y, z, t)}{dt^2} - \nabla^2 u(x, y, z, t) &= q_s, \\
 u(x, y, z, 0) &= 0, \\
 \frac{du(x, y, z, t)}{dt} \Big|_{t=0} &= 0,
 \end{aligned}$$

149 where the variables of this expression are defined as follows:

- 150 • $m(x, y, z) = \frac{1}{c(x, y, z)^2}$, is the parametrization of the subsurface with $c(x, y, z)$ being
 151 the speed of sound as a function of the three space coordinates (x, y, z) ;
- 152 • $u(x, y, z, t)$, is the spatially varying acoustic wavefield, with the additional dimen-
 153 sion of time t ;
- 154 • q_s is the source term, which is a point source in this case.

155 The first step towards solving this equation is the definition of a discrete computa-
 156 tional grid, on which the model parameters, wavefields and source are defined. The
 157 computational grid is defined as a `Grid(shape)` object, where `shape` is the number

158 of grid points in each spatial dimension. Optional arguments for instantiating a `Grid`
 159 are `extent`, which defines the extent in physical units, and `origin`, the origin of the
 160 coordinate system, with respect to which all other coordinates are defined.

161 The next step is the symbolic definition of the squared slowness, wavefield and
 162 source. For this, we introduce some fundamental types.

- 163 • `Function` represents a discrete spatially varying function, such as the velocity. A
 164 `Function` is instantiated for a defined `name` and a given `Grid`.
- 165 • `TimeFunction` represents a discrete function that is both spatially varying and
 166 time dependent, such as wavefields. Again, a `TimeFunction` object is defined on
 167 an existing `Grid` and is identified by its `name`.
- 168 • `SparseFunction` and `SparseTimeFunction` represent sparse functions, that is
 169 functions that are only defined over a subset of the grid, such as a seismic point
 170 source. The corresponding object is defined on a `Grid`, identified by a `name`, and
 171 also requires the `coordinates` defining the location of the sparse points.

172 Apart from the grid information, these objects carry their respective FD dis-
 173 cretization information in space and time. They also have a `data` field that contains
 174 values of the respective function at the defined grid points. By default, `data` is ini-
 175 tialized with zeros and therefore automatically satisfies the initial conditions from
 176 equation 3.1. The initialization of the fields to solve the wave equation over a one-
 177 dimensional grid is displayed in Listing 1.

Listing 1 Setup Functions to express and solve the acoustic wave equation.

```

1 >>> from devito import Grid, TimeFunction, Function, SparseTimeFunction
2 >>> g = Grid(shape=(nx,), origin=(ox,), extent=(sx,))
3 >>> u = TimeFunction(name="u", grid=g, space_order=2, time_order=2) # Wavefield
4 >>> m = Function(name="m", grid=g) # Physical parameter
5 >>> q = SparseTimeFunction(name="q", grid=g, coordinates=coordinates) # Source

```

178 **3.2. Discretization.** With symbolic objects that represent the discrete velocity
 179 model, wavefields and source function, we can now define the full discretized wave
 180 equation. As mentioned earlier, one of the main features of Devito is the possibility
 181 to formulate stencil computations as concise mathematical expressions. To do so, we
 182 provide shortcuts to classic FD stencils, as well as the functions to define arbitrary
 183 stencils. The shortcuts are accessed as object properties and are supported by `Time-`
 184 `Function` and `Function` objects. For example, we can take spatial and temporal
 185 derivatives of the wavefield `u` via the shorthand expressions `u.dx` and `u.dt` (Listing 2).

Listing 2 Example of spatial and temporal FD stencil creation.

```

1 >>> u.dx
2 -u(t, x - h_x)/(2*h_x) + u(t, x + h_x)/(2*h_x)
3 >>> u.dt
4 -u(t - dt, x)/(2*dt) + u(t + dt, x)/(2*dt)
5 >>> u.dt2
6 -2*u(t, x)/dt**2 + u(t - dt, x)/dt**2 + u(t + dt, x)/dt**2

```

186 Furthermore, Devito provides shortcuts for common differential operations such
 187 as the Laplacian via `u.laplace`. The full discrete wave equation can then be imple-
 188 mented in a single line of *Python* (Listing 3).

189 To solve the time-dependent wave equation with an explicit time-stepping scheme,
 190 the symbolic expression representing our PDE has to be rearranged such that it yields
 191 an update rule for the wavefield u at the next time step: $u(t + dt) = f(u(t), u(t -$

Listing 3 Expressing the wave equation.

```
1 >>> wave_equation = m * u.dt2 - u.laplace
2 >>> wave_equation
3 (-2*u(t, x)/dt**2 + u(t - dt, x)/dt**2 + u(t + dt, x)/dt**2)*m(x) + 2*u(t, x)/h_x
   **2 - u(t, x - h_x)/h_x**2 - u(t, x + h_x)/h_x**2
```

192 `dt`)). `Devito` allows to rearrange the PDE expression automatically using the `solve`
193 function, as shown in Listing 4.

Listing 4 Time-stepping scheme for the acoustic wave equation. `region=INTERIOR`
ensures that the Dirichlet boundary conditions at the edges of the Grid are satisfied.

```
1 >>> from devito import Eq, INTERIOR, solve
2 >>> stencil = Eq(u.forward, solve(wave_equation, u.forward), region=INTERIOR)
3 >>> stencil
4 Eq(u(t + dt, x), -2*dt**2*u(t, x)/(h_x**2*m(x)) + dt**2*u(t, x - h_x)/(h_x**2*m(x))
   ) + dt**2*u(t, x + h_x)/(h_x**2*m(x)) + 2*u(t, x) - u(t - dt, x))
```

194 Note that the `stencil` expression in Listing 4 does not yet contain the point
195 source `q`. This could be included as a regular `Function` which has zeros all over the
196 grid except for a few points; this, however, would obviously be wasteful. Instead,
197 `SparseFunctions` allow to perform operations, such as injecting a source or sampling
198 the wavefield, at a subset of grid points determined by `coordinates`. In the case in
199 which coordinates do not coincide with grid points, bilinear (for 2D) or trilinear (for
200 3D) interpolation are employed. To inject a point source into the `stencil` expression,
201 we use the `inject` function of the `SparseTimeFunction` object that represents our
202 seismic source (Listing 5).¹

Listing 5 Expressing the injection of a source into a field.

```
1 >>> injection = q.inject(field=u.forward, expr=dt**2 * q / m)
2 >>> injection
3 [Eq(u[t + 1, INT(floor((-o_x + q_coords[p_q, 0])/h_x))], dt**2*(1 - FLOAT(-h_x*INT
   (floor((-o_x + q_coords[p_q, 0])/h_x)) - o_x + q_coords[p_q, 0])/h_x)*q[time,
   p_q]/m[INT(floor((-o_x + q_coords[p_q, 0])/h_x))] + u[t + 1, INT(floor((-o_x +
   q_coords[p_q, 0])/h_x))],
4 Eq(u[t + 1, INT(floor((-o_x + q_coords[p_q, 0])/h_x)) + 1], dt**2*FLOAT(-h_x*INT(
   floor((-o_x + q_coords[p_q, 0])/h_x)) - o_x + q_coords[p_q, 0])*q[time, p_q
   ]/(h_x*m[INT(floor((-o_x + q_coords[p_q, 0])/h_x)) + 1]) + u[t + 1, INT(floor
   ((-o_x + q_coords[p_q, 0])/h_x)) + 1]])
```

203 The `inject` function takes the field being updated as an input argument (in this
204 case `u.forward`), while `expr=dt**2 * q / m` is the expression being injected. The
205 result of the `inject` function is a list of symbolic expressions, similar to the `stencil`
206 expression we defined earlier. As we shall see, these expressions are eventually joined
207 together and used to create an `Operator` object – the solver of our PDE.

208 **3.3. Boundary conditions.** Simple boundary conditions (BCs), such as Dirich-
209 let BCs, can be imposed on individual equations through special keywords (see List-
210 ing 4). For more exotic schemes, instead, the BCs need to be explicitly written (e.g.,
211 Higdon BCs [15]), just like any of the symbolic expressions defined in the Listings

¹More complicated interpolation schemes can be defined by precomputing the grid points cor-
responding to each sparse point, and their respective coefficients. The result can then be used to
create a `PrecomputedSparseFunction`, which behaves like a `SparseFunction` at the symbolic level.

212 above. For reasons of space, this aspect is not elaborated further; the interested
213 reader may refer to [26].

214 **3.4. Control flow.** By default, the extent of a `TimeFunction` in the time dimen-
215 sion is limited by its time order. Hence, the shape of u in Listing 1 is $(time_order +$
216 $1, nx) = (3, nx)$. The iterative method will then access u via modulo iteration, that
217 is $u[t\%3, \dots]$. In many scenarios, however, the entire time history, or at least periodic
218 time slices, should be saved (e.g., for inversion algorithms). Listing 6 expands our
219 running example with an equation that saves the content of u every 4 iterations, up
220 to a maximum of $save = 100$ time slices.

Listing 6 Implementation of time sub-sampling.

```
1 >>> from devito import ConditionalDimension
2 >>> ts = ConditionalDimension('ts', parent=g.time_dim, factor=4)
3 >>> us = TimeFunction(name='us', grid=g, save=100, time_dim=ts)
4 >>> save = Eq(us, u)
```

221 In general, all equations that access `Functions` (or `TimeFunctions`) employing
222 one or more `ConditionalDimensions` will be conditionally executed. The condition
223 may be a number indicating how many iterations should pass between two executions
224 of the same equation, or even an arbitrarily complex expression.

225 **3.5. Domain, halo, and padding regions.** A `Function` internally distin-
226 guishes between three regions of points.

227 **Domain** Represents the *computational domain* of the `Function` and is inferred from
228 the input `Grid`. This includes any elements added to the *physical domain*
229 purely for computational purposes, e.g. absorbing boundary layers.

230 **Halo** The grid points surrounding the domain region, i.e. “ghost” points that are
231 accessed by the stencil when iterating in proximity of the domain boundary.

232 **Padding** The grid points surrounding the halo region, which are allocated for perfor-
233 mance optimizations, such as data alignment. Normally this region should be
234 of no interest to a user of Devito, except for precise measurement of memory
235 allocated for each `Function`.

236 **4. The Devito compiler.** In Devito, an `Operator` carries out three fundamen-
237 tal tasks: generation of low-level code, JIT compilation, and execution. The `Operator`
238 input consists of one or more symbolic equations. In the generated code, these equa-
239 tions are scheduled within loop nests of suitable depth and extent. The `Operator` also
240 accepts substitution rules (to replace symbols with constant values) and optimization
241 levels for the Devito Symbolic Engine (DSE) and the Devito Loop Engine (DLE). By
242 default, all DSE and DLE optimizations that are known to unconditionally improve
243 performance are automatically applied. The same `Operator` may be reused with dif-
244 ferent input data; JIT-compilation occurs only once, triggered by the first execution.
245 Overall, this lowering process – from high-level equations to dynamically compiled
246 and executable code – consists of multiple compiler passes, summarized in Figure 2
247 and discussed in the following sections (a minimal background in data dependence
248 analysis is recommended; the unfamiliar reader may refer to a classic textbook such
249 as [1]).

250 **4.1. Equations lowering.** In this pass, three main tasks are carried out: *in-*
251 *dexification*, *substitution*, and *domain-alignment*.

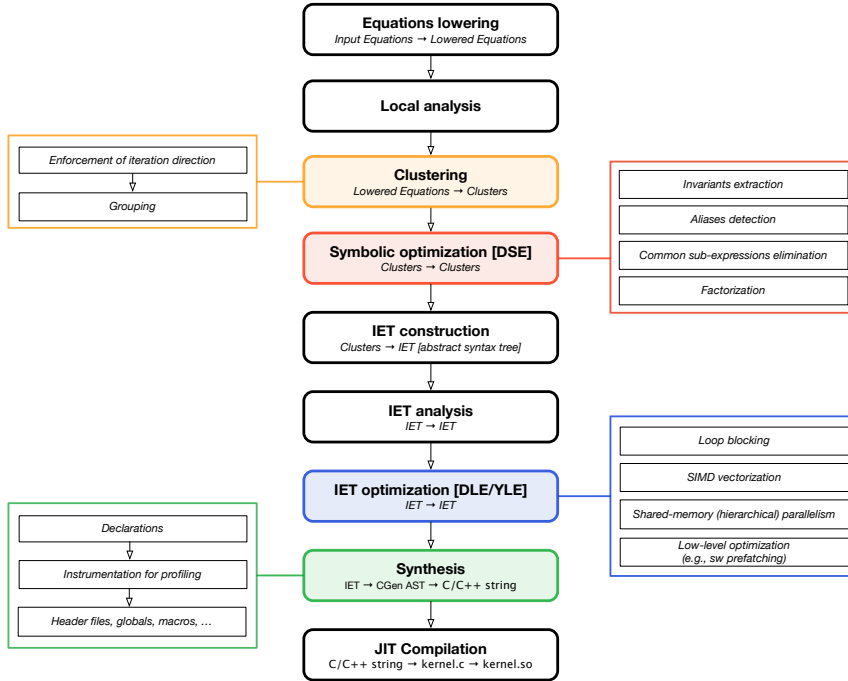


Fig. 2: Compiler passes to lower symbolic equations into shared objects through an Operator.

- 252 • As explained in Section 3, the input equations typically involve one or more indexed
253 **Functions**. The *indexification* consists of converting such objects into actual ar-
254 rays. An array always keeps a reference to its originating **Function**. For instance,
255 all accesses to u such as $u[t, x + 1]$ and $u[t + 1, x - 2]$ would store a pointer to the
256 same, user-defined **Function** $u(t, x)$. This metadata is exploited throughout the
257 various compilation passes.
- 258 • During *substitution*, the user-provided substitution rules are applied. These may
259 be given for any literal appearing in the input equations, such as the grid spacing
260 symbols. Applying a substitution rule increases the chances of constant folding,
261 but it makes the **Operator** less generic. The values of symbols for which no
262 substitution rule is available are provided at execution time.
- 263 • The *domain-alignment* step shifts the array accesses deriving from **Functions** hav-
264 ing non-empty halo and padding regions. Thus, the array accesses become logically
265 aligned to the equation’s natural domain. For instance, given the usual **Function**
266 $u(t, x)$ having two points on each side of the x halo region, the array accesses $u[t, x]$
267 and $u[t, x + 2]$ are transformed, respectively, into $u[t, x + 2]$ and $u[t, x + 4]$. When
268 $x = 0$, therefore, the values $u[t, 2]$ and $u[t, 4]$ are fetched, representing the first and
269 third points in the computational domain.

270 **4.2. Local analysis.** The lowered equations are analyzed to collect information
271 relevant for the **Operator** construction and execution. In this pass, an equation is
272 inspected “in isolation”, ignoring its relationship with the rest of the input. The
273 following metadata are retrieved and/or computed:

- 274 • input and output **Functions**;
- 275 • **Dimensions**, which are topologically ordered based on how they appear in the
- 276 various array index functions; and
- 277 • two notable **Spaces**: the iteration space, **ISpace**, and the data space, **DSpace**.

278 A **Space** is a collection of points given by the product of n compact intervals on
 279 \mathbb{Z} . With the notation $d[o_m, o_M]$ we indicate the compact interval $[d_m + o_m, d_M +$
 280 $o_M]$ over the **Dimension** d , in which d_m and d_M are parameters (specialized only
 281 at runtime), while o_m and o_M are known integers. For instance, $[x[0, 0], y[-1, 1]]$
 282 describes a rectangular two-dimensional space over x and y , whose points are given
 283 by the Cartesian product $[x_m, x_M] \times [y_m - 1, y_M + 1]$. The **ISpace** and **DSpace** are
 284 two special types of **Space**. They usually span different sets of **Dimensions**. A **DSpace**
 285 may have **Dimensions** that do not appear in an **ISpace**, in particular those that are
 286 accessed only via integer indices. Likewise, an **ISpace** may have **Dimensions** that are
 287 not part of the **DSpace**, such as a reduction axis. Further, an **ISpace** also carries, for
 288 each **Dimension**, its iteration direction.

289 As an example, consider the equation *stencil* in Listing 4. Immediately we see
 290 that $\text{input} = [u, m]$, $\text{output} = [u]$, $\text{Dimensions} = [t, x]$. The compiler constructs the
 291 **ISpace** $[t[0, 0]^+, x[0, 0]^*]$. The first entry $t[0, 0]^+$ indicates that, along t , the equation
 292 should run between $t_m + 0$ and $t_M + 0$ (extremes included) in the *forward* direction,
 293 as indicated by the symbol $+$. This is due to the fact that there is a flow dependence
 294 in t , so only a unitary positive stepping increment (i.e., $t = t + 1$) allows a correct
 295 propagation of information across consecutive iterations. The only difference along
 296 x is that the iteration direction is now arbitrary, as indicated by $*$. The **DSpace** is
 297 $[t[0, 1], x[0, 0]]$; intuitively, the entry $t[0, 1]$ is used right before running an **Operator**
 298 to provide a default value for $t_M - 1$ - in particular, t_M will be set to the largest possible
 299 value that does not cause out-of-domain accesses (i.e., out-of-bounds array accesses).

300 **4.3. Clustering.** A **Cluster** is a sequence of equations having (i) same **ISpace**,
 301 (ii) same control flow (i.e., same **ConditionalDimensions**), and (iii) no dimension-
 302 carried “true” anti-dependences among them.

303 As an example, consider again the setup in Section 3. The equation *stencil* cannot
 304 be “clusterized” with the equations in the *injection* list as their **ISpaces** are different.
 305 On the other hand, the equations in *injection* can be grouped together in the same
 306 **Cluster** as (i) they have same **ISpace** $[t[0, 0]^*, p_q[0, 0]^*]$, (ii) same control flow, and
 307 (iii) there are no true anti-dependences among them (note that the second equation
 308 in *injection* does write to $u[t + 1, \dots]$, but as explained later this is in fact a reduction,
 309 that is a “false” anti-dependence).

310 **4.3.1. Iteration direction.** First, each equation is assigned a new **ISpace**,
 311 based upon a *global* analysis. Any of the iteration directions that had been marked as
 312 “arbitrary” ($*$) during local analysis may now be enforced to *forward* ($+$) or *backward*
 313 ($-$). This process exploits data dependence analysis.

314 For instance, consider the flow dependence between *stencil* and the *injection* equa-
 315 tions. If we want u to be up-to-date when evaluating *injection*, then we eventually
 316 need all equations to be scheduled sequentially within the t loop. For this, the **ISpaces**
 317 of the *injection* equations are specialized by enforcing the direction *forward* along the
 318 **Dimension** t . The new **ISpace** is $[t[0, 0]^+, p_q[0, 0]^*]$.

319 Algorithm 1 illustrates how the enforcement of iteration directions is achieved in
 320 general. Whenever a clash is detected (i.e., two equations with **ISpace** $[d[0, 0]^+, \dots]$
 321 and $[d[0, 0]^-, \dots]$), the original direction determined by the local analysis pass is kept
 322 (lines 11 and 13), which will eventually lead to generating different loops.

Algorithm 1: Clustering: enforcement of iteration directions (pseudocode).

Input: A sequence of equations \mathcal{E} .**Output:** A sequence of equations \mathcal{E}' with altered ISpace.

// Map each dimension to a set of expected iteration directions

```
1 mapper ← DETECT_FLOW_DIRECTIONS( $\mathcal{E}$ );
2 for  $e$  in  $\mathcal{E}$  do
3   for  $dim, directions$  in mapper do
4     if  $len(directions) == 1$  then
5       // No ambiguity
6       forced[ $dim$ ] ← directions.pop();
7     else if  $len(directions) == 2$  then
8       // No ambiguity as long as one of the two items is /Any/
9       try
10        directions.remove(Any);
11        forced[ $dim$ ] ← directions.pop();
12      except
13        forced[ $dim$ ] ← e.directions[ $dim$ ];
14      else
15        forced[ $dim$ ] ← e.directions[ $dim$ ];
16      end if
17    end for
18   $\mathcal{E}'$ .append( $e$ .rebuild(directions=forced))
19 end for
20 return  $\mathcal{E}'$ 
```

323 **4.3.2. Grouping.** This step performs the actual clustering, checking ISpaces
324 and anti-dependences, as well as handling control flow. The procedure is shown in
325 Algorithm 2; some explanations follow.

Algorithm 2: Clustering: grouping expressions into Clusters (pseudocode)

Input: A sequence of equations \mathcal{E} .**Output:** A sequence of clusters \mathcal{C} .

```
1  $\mathcal{C}$  ← ClusterGroup();
2 for  $e$  in  $\mathcal{E}$  do
3   grouped ← false;
4   for  $c$  in reversed( $\mathcal{C}$ ) do
5     anti, flow ← GET_DEPENDENCES( $c, e$ );
6     if  $e.ispace == c.ispace$  and  $anti.carried$  is empty then
7       c.add( $e$ );
8       grouped ← true;
9       break;
10    else if  $anti.carried$  is not empty then
11      c.atomics.update( $anti.carried.cause$ );
12      break;
13    else if  $flow.cause.intersection(c.atomics)$  then
14      // cannot search across earlier clusters
15      break;
16    end for
17  if not grouped then
18     $\mathcal{C}$ .append(Cluster( $e$ ));
19  end if
20 end for
21  $\mathcal{C}$  ← CONTROL_FLOW( $\mathcal{C}$ );
22 return  $\mathcal{C}$ 
```

326 • Robust data dependence analysis, capable of tracking flow-, anti-, and output-

dependencies at the level of array accesses, is necessary. In particular, it must be able to tell whether two generic *array accesses* induce a dependence or not. The data dependence analysis performed is conservative; that is, a dependence is always assumed when a test is inconclusive. Dependence testing is based on the standard Lamport test [1]. In Algorithm 2, data dependence analysis is carried out by the function GET_DEPENDENCES.

- If an anti-dependence is detected along a Dimension i , then i is marked as *atomic* – meaning that no further clustering can occur along i . This information is also exploited by later Operator passes (see Section 4.5).
- Reductions, and in particular increments, are treated specially. They represent a special form of anti-dependence, as they do not break clustering. GET_DEPENDENCES detects reductions and removes them from the set of anti-dependencies.
- Given the sequence of equations $[E_1, E_2, E_3]$, it is possible that E_3 can be grouped with E_1 , but not with its immediate predecessor E_2 (e.g., due to a different ISpace). However, this can only happen when there are no flow or anti-dependencies between E_2 and E_3 ; i.e. when the if commands at lines 10 and 13 are not entered, thus allowing the search to proceed with the next equation. This optimization was originally motivated by gradient operators in seismic imaging kernels.
- The routine CONTROL_FLOW, omitted for brevity, creates additional Clusters if one or more ConditionalDimensions are encountered. These are tracked in a special Cluster field, *guards*, as also required by later passes (see Section 4.5).

4.4. Symbolic optimization. The DSE – Devito Symbolic Engine – is a macro-pass reducing the *arithmetic strength* of Clusters (e.g., their operation count). It consists of a series of passes, ranging from standard common sub-expression elimination (CSE) to more advanced rewrite procedures, applied individually to each Cluster. The DSE output is a new ordered sequence of Clusters: there may be more or fewer Clusters than in the input, and both the overall number of equations as well as the sequence of arithmetic operations might differ. The DSE passes are discussed in Section 5.1. We remark that the DSE only operates on Clusters (i.e., on collections of equations); there is no concept of “loop” at this stage yet. However, by altering Clusters, the DSE has an indirect impact on the final loop-nest structure.

4.5. IET construction. In this pass, the intermediate representation is lowered to an Iteration/Expression Tree (IET). An IET is an abstract syntax tree in which Iterations and Expressions – two special node types – are the main actors. Equations are wrapped within Expressions, while Iterations represent loops. Loop nests embedding such Expressions are constructed by suitably nesting Iterations. Each Cluster is eventually placed in its own loop (Iteration) nest, although some (outer) loops may be shared by multiple Clusters.

Consider again our running acoustic wave equation example. There are three Clusters in total: C_1 for *stencil*, C_2 for *save*, and C_3 for the equations in *injection*. We use Algorithm 3 – an excerpt of the actual cluster scheduling algorithm – to explain how this sequence of Clusters is turned into an IET. Initially, the *schedule* list is empty, so when C_1 is handled two nested Iterations are created (line 15), respectively for the Dimensions t and x . Subsequently, C_2 ’s ISpace and the current *schedule* are compared (line 5). It turns out that t appears among C_2 ’s guards, hence the for loop is exited at line 12 without inspecting the second and last iteration. Thus, $index = 1$, and the previously built Iteration over t is reused. Finally, when processing C_3 , the for loop is exited at the second iteration due to line 6, since $p_q! = x$. Again, the t Iteration is reused, while a new Iteration is constructed for

Algorithm 3: An excerpt of the cluster scheduling algorithm, turning a list (of Clusters) into a tree (IET). Here, the fact that different Clusters may eventually share some outer Iterations is highlighted.

Input: A sequence of Clusters \mathcal{C} .

Output: An Iteration/Expression Tree.

```

1 schedule ← list();
2 for c in C do
3   root ← None;
4   index ← 0;
5   for i0, i1 in zip(c.ispace, schedule) do
6     if i0 != i1 or i0.dimension in c.atomics then
7       break;
8     end if
9     root ← schedule[i1];
10    index ← index + 1;
11    if i0.dim in c.guards then
12      break;
13    end if
14  end for
15  <build as many Iterations as Dimensions in c.ispace[index:] and nest them inside root>;
16  <update schedule>;
17  <...>
18 end for

```

376 the Dimension p_q . Eventually, the constructed IET is as in Listing 7.

Listing 7 Graphical representation of the IET produced by the cluster scheduling algorithm for the running example.

```

1 for t = t_m to t_M:
2   |-- for x = x_m to x_M:
3     |-- <Eq(u[t+1,x], ...) >
4     |
5     |-- if t % 4 == 0
6       |-- for x = x_m to x_M:
7         |-- <Eq(us[t/4, x], ...) >
8         |
9       |-- for p_q = p_q_m to p_q_M:
10        |-- <Eq(u[t+1,f(p_q)], ...) >
11        |-- <Eq(u[t+1,g(p_q)], ...) >

```

377 **4.6. IET analysis.** The newly constructed IET is analyzed to determine Itera-
378 tion properties such as **sequential**, **parallel**, and **vectorizable**, which are then
379 attached to the relevant nodes in the IET. These properties are used for loop optimiza-
380 tion, but only by a later pass (see Section 4.7). To determine whether an Iteration
381 is **parallel** or **sequential**, a fundamental result from compiler theory is used –
382 the i -th Iteration in a nest comprising n Iterations is parallel if for all dependen-
383 ces D , expressed as distance vectors $D = (d_0, \dots, d_{n-1})$, either $(d_1, \dots, d_{i-1}) > 0$ or
384 $(d_1, \dots, d_i) = 0$ [1].

385 **4.7. IET optimization.** This macro-pass transforms the IET for performance
386 optimization. Apart from runtime performance, this pass also optimizes for rapid
387 JIT compilation with the underlying C compiler. A number of loop optimizations are
388 introduced, including loop blocking, minimization of remainder loops, SIMD vector-
389 ization, shared-memory (hierarchical) parallelism via OpenMP, software prefetching.
390 These will be detailed in Section 5. A *backend* (see Section 4.9) might provide its own

391 loop optimization engine.

392 **4.8. Synthesis, dynamic compilation, and execution.** Finally, the IET
393 adds variable declarations and header files, as well as instrumentation for performance
394 profiling, in particular, to collect execution times of specific code regions. Declara-
395 tions are injected into the IET, ensuring they appear as close as possible to the scope
396 in which the relative variables are used, while honoring the OpenMP semantics of pri-
397 vate and shared variables. To generate C code, a suitable tree visitor inspects the IET
398 and incrementally builds a *CGen* tree [19], which is ultimately translated into a string
399 and written to a file. Such files are stored in a software cache of Devito-generated
400 **Operators**, JIT-compiled into a shared object, and eventually loaded into the *Python*
401 environment. The compiled code has a default entry point (a special function), which
402 is called directly from *Python* at **Operator** application time.

403 **4.9. Operator specialization through backends.** In Devito, a *backend* is
404 a mechanism to specialize data types as well as **Operator** passes, while preserving
405 software modularity (inspired by [25]).

406 One of the main objectives of the backend infrastructure is promoting software
407 composability. As explained in Section 2, there exist a significant number of interest-
408 ing tools for stencil optimization, which we may want to integrate with Devito. For
409 example, one of the future goals is to support GPUs, and this might be achieved by
410 writing a new backend implementing the interface between Devito and third-party
411 software specialized for this particular architecture.

412 Currently, two backends exist:

413 **core** the default backend, which relies on the DLE for loop optimization.
414 **yask** an alternative backend using the YASK stencil compiler to generate optimized
415 C++ code for Intel® Xeon® and Intel® Xeon Phi™ architectures [40].
416 Devito transforms the IET into a format suitable for YASK, and uses its API
417 for data management, JIT-compilation, and execution. Loop optimization is
418 performed by YASK through the YASK Loop Engine (YLE).

419 The *core* and *yask* backends share the compilation pipeline in Figure 2 until the loop
420 optimization stage.

421 **5. Automated performance optimizations.** As discussed in Section 4, De-
422 vito performs symbolic optimizations to reduce the arithmetic strength of the expres-
423 sions, as well as loop transformations for data locality and parallelism. The former are
424 implemented as a series of compiler passes in the DSE, while for the latter there cur-
425 rently are two alternatives, namely the DLE and the YLE (depending on the chosen
426 execution backend).

427 Devito abstracts away the single optimizations passes by providing users with a
428 certain number of optimization levels, called “modes“, which trigger pre-established
429 sequences of optimizations – analogous to what general-purpose compilers do with,
430 for example, -O2 and -O3. In Sections 5.1, 5.2, and 5.3 we describe the individual
431 passes provided by the DSE, DLE, and YLE respectively, while in Section 7.1 we
432 explain how these are composed into modes.

433 **5.1. DSE - Devito Symbolic Engine.** The DSE passes attempt to reduce the
434 arithmetic strength of the expressions through FLOP-reducing transformations [11].
435 They are illustrated in Listings 8-11, which derive from the running example used
436 throughout the article. A detailed description follows.

437 • **Common sub-expression elimination (CSE).** Two implementations are avail-
438 able: one based upon *SymPy*'s *cse* routine and one built on top of more basic

439 *SymPy* routines, such as `xreplace`. The former is more powerful, being aware
 440 of key arithmetic properties such as associativity; hence it can discover more re-
 441 dundancies. The latter is simpler, but avoids a few critical issues: (i) it has a
 442 much quicker turnaround time; (ii) it does not capture integer index expressions
 443 (for increased quality of the generated code); and (iii) it tries not to break factor-
 444 ization opportunities. A generalized common sub-expressions elimination routine
 445 retaining the features and avoiding the drawbacks of both implementations is still
 446 under development. By default, the latter implementation is used when the CSE
 447 pass is selected.

Listing 8 An example of common sub-expressions elimination.

```
1 >>> 9.0*dt*dt*u[t, x + 1] - 18.0*dt*dt*u[t][x + 2] + 9.0*dt*dt*u[t, x + 3]
2 temp0 = dt*dt
3 9.0*temp0*u[t, x + 1] - 18.0*temp0*u[t][x + 2] + 9.0*temp0*u[t, x + 3]
```

448 • **Factorization.** This pass visits each expression tree and tries to factorize FD
 449 weights. Factorization is applied without altering the expression structure (e.g.,
 450 without expanding products) and without performing any heuristic search across
 451 groups of expressions. This choice is based on the observation that a more ag-
 452 gressive approach is only rarely helpful (never in the test cases in Section 7),
 453 while the increase in symbolic processing time could otherwise be significant. The
 454 implementation exploits the *SymPy* `collect` routine. However, while `collect`
 455 only searches for common factors across the immediate children of a single node,
 456 the DSE implementation recursively applies `collect` to each `Add` node (i.e., an
 457 addition) in the expression tree, until the leaves are reached.

Listing 9 An example of FD weights factorization.

```
1 >>> 9.0*temp0*u[t, x + 1] - 18.0*temp0*u[t][x + 2] + 9.0*temp0*u[t, x + 3]
2 9.0*temp0*(u[t, x + 1] + u[t, x + 3]) - 18.0*temp0*u[t][x + 2]
```

458 • **Extraction.** The name stems from the fact that sub-expressions matching a
 459 certain condition are pulled out of a larger expression, and their values are stored
 460 into suitable scalar or tensor temporaries. For example, a condition could be
 461 “*extract all time-varying sub-expressions whose operation count is larger than a*
 462 *given threshold*”. A tensor temporary may be preferred over a scalar temporary if
 463 the intention is to let the *IET construction* pass (see Section 4.5) place the pulled
 464 sub-expressions within an outer loop nest. Obviously, this comes at the price of
 465 additional storage. This peculiar effect – trading operations for memory – will be
 466 thoroughly analyzed in Sections 6 and 7.

Listing 10 An example of time-varying sub-expressions extraction. Only sub-
 expressions performing at least one floating-point operation are extracted.

```
1 >>> 9.0*temp0*(u[t, x + 1] + u[t, x + 3]) - 18.0*temp0*u[t][x + 2]
2 temp1[x] = u[t, x + 1] + u[t, x + 3]
3 9.0*temp0*temp1[x] - 18.0*temp0*u[t][x + 2]
```

467 • **Detection of aliases.** The Alias-Detection Algorithm implements the most ad-
 468 vanced DSE pass. In essence, an alias is a sub-expression that is redundantly com-
 469 puted at multiple iteration points. Because of its key role in the Cross-Iteration

470 Redundancy-Elimination algorithm, the formalization of the Alias-Detection Al-
471 gorithm is postponed until Section 6.

Listing 11 An example of alias detection.

```
1 >>> 9.0*temp0*u[t, x + 1] - 18.0*temp0*u[t][x + 2] + 9.0*temp0*u[t, x + 3]
2 temp[x] = 9.0*temp0*u[t, x]
3 temp[x + 1] - 18.0*temp0*u[t][x + 2] + temp[x + 3]
```

472 **5.2. DLE - Devito Loop Engine.** The DLE transforms the IET via classic
473 loop optimizations for parallelism and data locality [17]. These are summarized below.

474 • **SIMD Vectorization.** Implemented by enforcing compiler auto-vectorization
475 via special `pragmas` from the OpenMP 4.0 language. With this approach, the
476 DLE aims to be performance-portable across different architectures. However,
477 this strategy causes a significant fraction of vector loads/stores to be unaligned
478 to cache boundaries, due to the stencil offsets. As we shall see, this is a primary
479 cause of performance loss.

480 • **Loop Blocking.** Also known as “tiling”, this technique implemented by replacing
481 `Iteration` trees in the IET. The current implementation only supports blocking
482 over fully-parallel `Iterations`. Blocking over dimensions characterized by flow- or
483 anti-dependences, such as the time dimension in typical explicit finite difference
484 schemes, is instead work in progress (this would require a preliminary pass known
485 as loop skewing; see Section 8 for more details). On the other hand, a feature of
486 the present implementation is the capability of blocking across particular *sequences*
487 of loop nests. This is exploited by the Cross-Iteration Redundancy-Elimination
488 algorithm, as shown in Section 6.3. To determine an optimal block shape, an
489 `Operator` resorts to empirical auto-tuning.

490 • **Parallelism.** Shared-memory parallelism is introduced by decorating `Iterations`
491 with suitable OpenMP `pragmas`. The OpenMP `static` scheduling is used. Nor-
492 mally, only the outermost fully-parallel `Iteration` is annotated with the parallel
493 `pragma`. However, heuristically nested fully-parallel `Iterations` are collapsed if
494 the core count is greater than a certain threshold. This pass also ensures that all
495 array temporaries allocated in the scope of the parallel `Iteration` are declared as
496 `private` and that storage is allocated where appropriate (stack, heap).

497 Summarizing, the DLE applies a sequence of typical stencil optimizations, aiming
498 to reach a minimum level of performance across different architectures. As we shall
499 see, the effectiveness of this approach, based on simple transformations, deteriorates
500 on architectures strongly conceived for hierarchical parallelism. This is one of the main
501 reasons behind the development of the `yask` backend (see Section 4.9), described in
502 the following section.

503 **5.3. YLE - YASK Loop Engine.** “YASK” (Yet Another Stencil Kernel) is an
504 open-source C++ software framework for generating high-performance implementa-
505 tions of stencil codes for Intel[®] Xeon[®] and Intel[®] Xeon Phi[™] processors. Previous
506 publications on YASK have discussed its overall structure [40] and its application to
507 the Intel[®] Xeon Phi[™] x100 family (code-named Knights Corner) [37] and Intel[®]
508 Xeon Phi[™] x200 family (code-named Knights Landing) [38, 33] many-core CPUs.
509 Unlike Devito, it does not expose a symbolic language to the programmer or create
510 stencils from finite-difference approximations of differential equations. Rather, the
511 programmer provides simple declarative descriptions of the stencil equations using
512 a C++ or Python API. Thus, Devito operates at a level of abstraction higher than

513 that of YASK, while YASK provides performance portability across Intel architectures
514 and is more focused on low-level optimizations. Following is a sample of some of the
515 optimizations provided by YASK:²

- 516 • **Vector-folding.** In traditional SIMD vectorization, such as that provided by
517 a vectorizing compiler, the vector elements are arranged sequentially along the
518 unit-stride dimension of the grid, which must also be the dimension iterated over
519 in the inner-most loop of the stencil application. Vector-folding is an alternative
520 data-layout method whereby neighboring elements are arranged in small *multi-*
521 *dimensional* tiles. Figure 3 illustrates three ways to pack eight double-precision
522 floating-point values into a 512-bit SIMD register. Figure 3a shows a traditional
523 1D “in-line” layout, and 3b and 3c show alternative 2D and 3D “folded” lay-
524 outs. Furthermore, these tiles may be ordered in memory in a dimension indepen-
525 dent of the dimensions used in vectorization [37]. The combination of these two
526 techniques can significantly increase overlap and reuse between successive stencil-
527 application iterations, reducing the memory-bandwidth demand. For stencils that
528 are bandwidth-bound, this can provide significant performance gains [37, 33].
- 529 • **Software prefetching.** Many high-order or staggered-grid stencils require many
530 streams of data to be read from memory, which can overwhelm the hardware
531 prefetchers. YASK can be directed to automatically generate software prefetch
532 instructions to improve the cache hit rates, especially on Xeon Phi CPUs.
- 533 • **Hierarchical parallelism.** Dividing the spatial domain into tiles to increase tem-
534 poral cache locality is a common stencil optimization as discussed earlier. When
535 implementing this technique, sometimes called “cache-blocking”, it is typical to
536 assign each thread to one or more small rectilinear subsets of the domain in which
537 to apply the stencil(s). However, if these threads share caches, one thread’s data
538 will often evict data needed later by another thread, reducing the effective capacity
539 of the cache. YASK addresses this by employing two levels of OpenMP paralleliza-
540 tion: the outer level of parallel loops are applied across the cache-blocks, and an
541 inner level is applied across sub-blocks within those tiles. In the case of the Xeon
542 Phi, the eight hyper-threads that share each L2 cache can now cooperate on filling
543 and reusing the data in the cache, rather than evicting each other’s data.

544 YASK also provides other optimizations, such as temporal wave-front tiling, as
545 well as MPI support. These features, however, are not exploited by Devito yet. The
546 interested reader may refer to [38, 39].

547 To obtain the best of both tools, we have integrated the YASK framework into
548 the Devito package. In essence, the Devito `yask` backend exploits the intermediate
549 representation of an `Operator` to generate YASK kernels. This process is based upon
550 sophisticated compiler technology. In *Devito v3.1*, roughly 70% of the Devito API is
551 supported by the `yask` backend³.

552 **6. The Cross-Iteration Redundancy-Elimination Algorithm.** Aliases, or
553 “cross-iteration redundancies” (informally introduced in Section 5.1), in FD operators
554 depend on the differential operators used in the PDE(s) and the chosen discretization
555 scheme. From a performance viewpoint, the presence of aliases is a non-issue as long
556 as the operator is memory-bound, while it becomes relevant in kernels with a high
557 arithmetic intensity. In Devito, the Cross-Iteration Redundancy-Elimination (CIRE)
558 algorithm attempts to remove aliases with the goal of reducing the operation count. As

²Not all YASK features are currently used by Devito.

³At the time of writing, reaching feature-completeness is one the major on-going development efforts

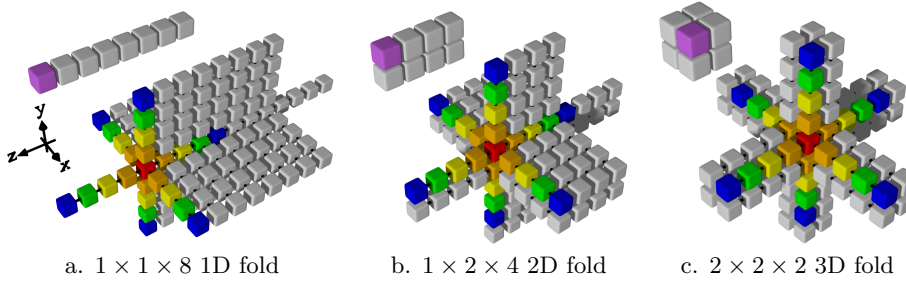


Fig. 3: Various folds of 8 elements [37]. The smaller diagram in the upper-left of each sub-figure illustrates a single SIMD layout, and the larger diagram shows the input values needed for a typical 25-point stencil, as from an 8th-order finite-difference approximation of an isotropic acoustic wave. Note that the $1 \times 1 \times 8$ 1D fold corresponds to the traditional in-line vectorization.

559 shown in Section 7, the CIRE algorithm has considerable impact in seismic imaging
 560 kernels. The algorithm is implemented through the orchestration of multiple DSE
 561 and DLE/YLE passes, namely *extraction of candidate expressions (DSE)*, *detection*
 562 *of aliases (DSE)*, *loop blocking (DLE/YLE)*.

563 **6.1. Extraction of candidate expressions.** The criteria for extraction of candi-
 564 date sub-expressions are:

- 565 • Any *maximal time-invariant* whose operation count is greater than $Thr_0 = 10$
 566 (floating point arithmetic only). The term “maximal” means that the expression
 567 is not embedded within a larger time-invariant. The default value $Thr_0 = 10$,
 568 determined empirically, provides systematic performance improvements in a series
 569 of seismic imaging kernels. Transcendental functions are given a weight in the
 570 order of tens of operations, again determined empirically.
- 571 • Any *maximal time-varying* whose operation count is greater than $Thr_1 = 10$. Such
 572 expressions often lead to aliases, since they typically result from taking spatial and
 573 time derivatives on `TimeFunctions`. In particular, cross-derivatives are a major
 574 cause of aliases.

575 This pass leverages the *extraction* routine described in Section 5.1.

576 **6.2. Detection of aliases.** To define the concept of aliasing expressions, we
 577 first need to formalize the notion of *translated operands*. Here, an operand is regarded
 578 as the arithmetic product of a scalar value (or “coefficient”) and one or more indexed
 579 objects. An indexed object is characterized by a label (i.e., its name), a vector of n
 580 dimensions, and a vector of n displacements (one for each dimension). We say that
 581 an operand o_1 is translated with respect to an operand o_0 if o_0 and o_1 have same
 582 coefficient, label, and dimensions, and if their displacement vectors are such that one
 583 is the translation of the other (in the classic geometric sense). For example, the
 584 operand $2 * u[x, y, z]$ is translated with respect to the operand $2 * u[x + 1, y + 2, z + 3]$
 585 since they have same coefficient (2), label (u), and dimensions ($[x, y, z]$), while the
 586 displacement vectors $[0, 0, 0]$ and $[1, 2, 3]$ are expressible by means of a translation.

587 Now consider two expressions e_0 and e_1 in fully-expanded form (i.e., a non-nested
 588 sum-of-operands). We say that e_0 is an alias of e_1 if the following conditions hold:

- 589 • the operands in e_0 (e_1) are expressible as a translation of the operands in e_1 (e_0);
- 590 • the same arithmetic operators are applied to the involved operands.

591 For example, consider $e = u[x] + v[x]$, having two operands $u[x]$ and $v[x]$; then:
592 • $\mathbf{u}[\mathbf{x}-1] + \mathbf{v}[\mathbf{y}-1]$ is *not* an alias of e , due to a different dimension vector.
593 • $\mathbf{u}[\mathbf{x}] + \mathbf{w}[\mathbf{x}]$ is *not* an alias of e , due to a different label.
594 • $\mathbf{u}[\mathbf{x}+2] + \mathbf{v}[\mathbf{x}]$ is *not* an alias of e , since a translation cannot be determined.
595 • $\mathbf{u}[\mathbf{x}+2] + \mathbf{v}[\mathbf{x}+2]$ is an alias of e , as the operands $u[x+2]$ and $v[x+2]$ can be
596 expressed as a translation of $u[x]$ and $v[x]$ respectively, with $T(o_d) = o_d + \mathbf{2}$ and
597 o_d representing the displacement vector of an operand.

598 The relation “ e_0 is an alias of e_1 ” is an equivalence relation, as it is at the same
599 time reflexive, symmetric, and transitive. Thanks to these properties, the turnaround
600 times of the Alias-Detection Algorithm are extremely quick (less than 2 seconds run-
601 ning on an Intel[®] Xeon[®] E5-2620 v4 for the challenging `tti` test case with `so=16`,
602 described in Section 7.2), despite the $O(n^2)$ computational complexity (with n repre-
603 senting the number of candidate expressions, see Section 6.1).

604 Algorithm 4 highlights the fundamental steps of the Alias-Detection Algorithm.
605 In the worst case scenario, all pairs of candidate expressions are compared by ap-
606 plying the aliasing definition given above. Aggressive pruning, however, is applied
607 to minimize the cost of the search. The algorithm uses some auxiliary functions:
608 (i) `CALCULATE_DISPLACEMENTS` returns a mapper associating, to each candidate, its
609 displacement vectors (one for each indexed object); (ii) `COMPARE_OPS`(e_1, e_2) eval-
610 uates to true if e_1 and e_2 perform the same operations on the same operands; (iii)
611 `IS_TRANSLATED`(d_1, d_2) evaluates to true if the displacement vectors in d_2 are pairwise-
612 translated with respect to the vectors in d_1 by the same factor. Together, (ii) and
613 (iii) are used to establish whether two expressions alias each other (line 8).

614 Eventually, m sets of aliasing expressions are determined. For each of these sets
615 G_0, \dots, G_{m-1} , a *pivot* – a special aliasing expression – is constructed. This is the key
616 for operation count reduction: the pivot p_i of $G_i = \{e_0, \dots, e_{k-1}\}$ will be used in place
617 of e_0, \dots, e_{k-1} (thus obtaining a reduction proportional to k). A simple example is
618 illustrated in Listing 11.

Algorithm 4: The Alias-Detection Algorithm (pseudocode).

Input: A sequence of expressions \mathcal{E} .

Output: A sequence of Alias objects \mathcal{A} .

```

1 displacements ← CALCULATE_DISPLACEMENTS( $\mathcal{E}$ );
2  $\mathcal{A}$  ← list();
3 unseen ← list( $\mathcal{E}$ );
4 while unseen is not empty do
5     top ← unseen.pop();
6     G = Alias(top);
7     for e in unseen do
8         if COMPARE_OPS(top, e) and IS_TRANSLATED(displacements[top], displacements[e])
9             then
10                G.append(e);
11                unseen.remove(e);
12            end if
13        end for
14         $\mathcal{A}$ .append(G)
15    end while
16 return  $\mathcal{A}$ 

```

619 Several optimizations for data locality, not shown in Algorithm 4, are also applied.
620 The interested reader may refer to the documentation and the examples of *Devito v3.1*
621 for more details; below, we only mention the underlying ideas.

- 622 • The pivot of G_i is *constructed*, rather than selected out of e_0, \dots, e_{k-1} , so that
623 it could coexist with as many other pivots as possible within the same `Cluster`.
624 For example, consider again Listing 11: there are infinite possible pivots `temp[x`
625 `+ s] = 9.0*temp0*u[t, x + s]`, and the one with $s = 0$ is chosen. However,
626 this choice is not random: the Alias-Detection Algorithm chooses pivots based
627 on a global optimization strategy, which takes into account all of the m sets of
628 aliasing expressions. The objective function consists of choosing s so that multiple
629 pivots will have identical `ISpace`, and thus be scheduled to the same `Cluster`
630 (and, eventually, to the same loop nest).
- 631 • Conservatively, the chosen pivots are assigned to array variables. A second opti-
632 mization pass, called *index bumping and array contraction* in *Devito v3.1*, attempts
633 to turn these arrays into scalar variables, thus reducing memory consumption.
634 This pass is based on data dependence analysis, which essentially checks whether
635 a given pivot is required only within its `Cluster` or by later `Clusters` as well. In
636 the former case, the optimization is applied.

637 **6.3. Loop blocking for working-set minimization.** In essence, the CIRE
638 algorithm trades operation for memory – the (array) temporaries to store the aliases.
639 From a run-time performance viewpoint, this is convenient only in arithmetic-intensive
640 kernels. Unsurprisingly, we observed that storing temporary arrays spanning the
641 entire grid rarely provides benefits (e.g., only when the operation count reductions
642 are exceptionally high). We then considered the following options.

- 643 1. Capturing redundancies arising along the innermost dimension only. Thus,
644 only scalar temporaries would be necessary. This approach presents three
645 main issues, however: (i) only a small percentage of all redundancies are
646 captured; (ii) the implementation is non-trivial, due to the need for circular
647 buffers in the generated code; (iii) SIMD vectorization is affected, since inner
648 loop iterations are practically serialised. Some previous articles followed this
649 path [9, 10].
- 650 2. A generalization of the previous approach: using both scalar and array tempo-
651 raries, without searching for redundancies across the outermost loop(s). This
652 mitigates issue (i), although the memory pressure is still severely affected.
653 Issue (iii) is also unsolved. This strategy was discussed in [20].
- 654 3. Using loop blocking. Redundancies are sought and captured along all avail-
655 able dimensions, although they are now assigned to array temporaries whose
656 size is a function of the block shape. A first loop nest produces the array tempo-
657 raries, while a subsequent loop nest consumes them, to compute the actual
658 output values. The block shape should be chosen so that writes and reads to
659 the temporary arrays do not cause high latency accesses to the DRAM. An
660 illustrative example is shown in Listing 12.

661 The CIRE algorithm uses the third approach, based on cross-loop-nest blocking.
662 This pass is carried out by the DLE, which can introduce blocking over sequences of
663 loops (see Section 5.2).

664 **7. Performance evaluation.** We outline in Section 7.1 the compiler setup,
665 computer architectures, and measurement procedure that we used for our performance
666 experiments. Following that, we outline the physical model and numerical setup that
667 define the problem being solved in Section 7.2. This leads to performance results,
668 presented in Sections 7.3 and 7.4.

Listing 12 The loop nest produced by the CIRE algorithm for the example in Listing 11. Note that the block loop (line 2) wraps both the producer (line 3) and consumer (line 5) loops. For ease of read, unnecessary information are omitted.

```

1 for t = t_m to t_M:
2   for xb = x_m to x_M, xb += blocksize:
3     for x = xb to xb + blocksize + 3, x += 1
4       temp[x] = 9.0*temp0*u[t, x]
5       for x = xb to xb + blocksize; x += 1:
6         u[t+1,x,y] = ... + temp[x + 1] - 18.0*temp0*u[t][x + 2] + temp[x + 3] + ...

```

669 **7.1. Compiler and system setup.** We analyse the performance of generated
670 code using enriched roofline plots. Since the DSE transformations may alter the
671 operation count by allocating extra memory, only by looking at GFlops/s performance
672 and runtime jointly can a quality measure of code syntheses be derived.

673 For the roofline plots, Stream TRIAD was used to determine the attainable mem-
674 ory bandwidth of the node. Two peaks for the maximum floating-point performance
675 are shown: the ideal peak, calculated as

$$676 \quad \#[\text{cores}] \cdot \#[\text{avx units}] \cdot \#[\text{vector lanes}] \cdot \#[\text{FMA ports}] \cdot [\text{ISA base frequency}]$$

678 and a more realistic one, given by the LINPACK benchmark. The reported runtimes
679 are the minimum of three runs (the variance was negligible). The model used to
680 calculate the operational intensity assumes that the time-invariant Functions are
681 reloaded at each time iteration. This is a more realistic setting than a “compulsory-
682 traffic-only” model (i.e., an infinite cache).

683 We had exclusive access to two architectures: an Intel[®] Xeon[®] Platinum 8180
684 (formerly code-named Skylake) and an Intel[®] Xeon Phi[™] 7250 (formerly code-
685 named Knights Landing), which will be referred to as `sk18180` and `kn17250`. Thread
686 pinning was enabled with the program `numactl`. The Intel[®] compiler `icc version`
687 `18.0` was used to compile the generated code. The experiments were run with *Devito*
688 *v3.1* [41]. The experimentation framework with instructions for reproducibility is
689 available at [42]. All floating point operations are performed in single precision, which
690 is typical for seismic imaging applications.

691 Any arbitrary sequence of DSE and DLE/YLE transformations is applicable to an
692 Operator. Devito, provides three preset optimization sequences, or “modes”, which
693 vary in aggressiveness and affect code generation in three major ways:

- 694 • the time required by the Devito compiler to generate the code,
- 695 • the potential reduction in operation count, and
- 696 • the potential amount of additional memory that might be allocated to store (scalar,
697 tensor) temporaries.

698 A more aggressive mode might obtain a better operation count reduction than a non-
699 aggressive one, although this does not necessarily imply a better time to solution as
700 the memory pressure might also increase. The three optimization modes – `basic`,
701 `advanced`, and `aggressive`– apply the same sequence of DLE/YLE transformations,
702 which includes OpenMP parallelism, SIMD vectorization, and loop blocking. How-
703 ever, they vary in the number, type, and order of DSE transformations. In particular,
704 `basic` enables common sub-expressions elimination only;
705 `advanced` enables `basic`, then factorization, extraction of time-invariant aliases;
706 `aggressive` enables `advanced`, then extraction of time-varying aliases.
707 Thus, `aggressive` triggers the full-fledged CIRE algorithm, while `advanced` uses
708 only a relaxed version (based on time invariants). All runs used loop tiling with a

709 block shape that was determined individually for each case using auto-tuning. The
 710 auto-tuning phase, however, was not included in the measured experiment runtime.
 711 Likewise, the code generation phase is not included in the reported runtime.

712 **7.2. Test case setup.** In the following sections, we benchmark the performance
 713 of operators modeling the propagation of acoustic waves in two different models:
 714 isotropic and Tilted Transverse Isotropy (TTI, [44]), henceforth `isotropic` and `tti`,
 715 respectively. These operators were chosen for their relevance in seismic imaging tech-
 716 niques [44].

717 Acoustic `isotropic` modeling is the most commonly used technique for seismic
 718 inverse problems, due to the simplicity of its implementation, as well as the compar-
 719 atively low computational cost in terms of FLOPs. The `tti` wave equation provides
 720 a more realistic simulation of wave propagation and accounts for local directional
 721 dependency of the wave speed, but comes with increased computational cost and
 722 mathematical complexity. For our numerical tests, we use the `tti` wave equation as
 723 defined in [44]. The full specification of the equation as well as the finite difference
 724 schemes and its implementation using Devito are provided in [24, 23]. Essentially,
 725 the `tti` wave equation consists of two coupled acoustic wave equations, in which
 726 the Laplacians are constructed from spatially rotated first derivative operators. As
 727 indicated by Figure 4, these spatially rotated Laplacians have a significantly larger
 728 number of stencil coefficients in comparison to its isotropic equivalent which comes
 729 with an increased operational intensity.

730 The `tti` and `isotropic` equations are discretized with second order in time and
 731 varying space orders of 4, 8, 12 and 16. For both test cases, we use zero initial condi-
 732 tions, Dirichlet boundary conditions and absorbing boundaries with a 10 point mask
 733 (Section 3.5). The waves are excited by injecting a time-dependent, but spatially-
 734 localized seismic source wavelet into the subsurface model, using Devito’s sparse point
 735 interpolation and injection as described in Section 3.1. We carry out performance mea-
 736 surements for two velocity models of 512^3 and 768^3 grid points with a grid spacing
 737 of 20 m. Wave propagation is modeled for 1000 ms, resulting in 327 time steps for
 738 `isotropic`, and 415 time steps for `tti`. The time-stepping interval is chosen accord-
 739 ing to the Courant-Friedrichs-Lewy (CFL) condition [8], which guarantees stability
 740 of the explicit time-marching scheme and is determined by the highest velocity of the
 741 subsurface model and the grid spacing.

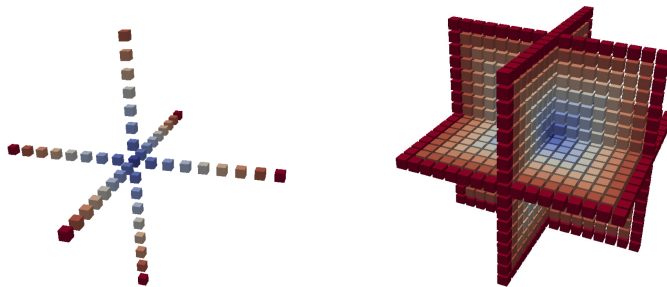


Fig. 4: Stencils of the acoustic Laplacian for the `isotropic` (left) and `tti` (right) wave equations and `so=16`. The anisotropic Laplacian corresponds to a spatially rotated version of the isotropic Laplacian. The color indicates the distance from the central coefficient.

742 **7.3. Performance: acoustic wave in isotropic model.** This section illus-
743 trates the performance of `isotropic` with the `core` and `yask` backends. To relieve
744 the exposition, we show results for the DSE in advanced mode only; the `aggressive`
745 has no impact on `isotropic`, due to the memory-bound nature of the code [23].

746 The performance of `core` on `skl8180`, illustrated in Figure 5a (`yask` uses slightly
747 smaller grids than `core` due to a flaw in the API of *Devito v3.1*, which will be fixed
748 in *Devito v3.2*), degrades as the space order (henceforth, `so`) increases. In particular,
749 it drops from 59% of the attainable machine peak to 36% in the case of `so=16`. This
750 is the result of multiple issues. As `so` increases, the number of streams of unaligned
751 virtual addresses also increases, causing more pressure on the memory system. Intel[®]
752 VTune[™] revealed that the lack of split registers to efficiently handle split loads was
753 a major source of performance degradation. Another major issue for `isotropic`
754 on `core` concerns the quality of the generated SIMD code. The in-line vectorization
755 performed by the auto-vectorizer produces a large number of pack/unpack instructions
756 to move data between vector registers, which introduces substantial overhead. Intel[®]
757 VTune[™] also confirmed that, unsurprisingly, `isotropic` is a memory-bound kernel.
758 Indeed, switching off the DSE basically did not impact the runtime, although it did
759 increase the operational intensity of the four test cases.

760 The performance of `core` on `kn17250` is not as good as that on `skl8180`. Figure 5b
761 shows an analogous trend to that on `skl8180`, with the attainable machine peak
762 systematically dropping as `so` increases. The issue is that here the distance from the
763 peak is even larger. This simply suggests that `core` is failing at exploiting the various
764 levels of parallelism available on `kn17250`.

765 The `yask` backend overcomes all major limitations to which `core` is subjected.
766 On both `skl8180` and `kn17250`, `yask` outperforms `core`, essentially since it does not
767 suffer from the issues presented above. Vector folding minimizes unaligned accesses;
768 software prefetching helps especially for larger values of `so`; hierarchical OpenMP
769 parallelism is fundamental to leverage shared caches. The speed-up on `kn17250` is
770 remarkable, since even in the best scenario for `core` (`so=4`), `yask` is roughly 3×
771 faster, and more than 4× faster when `so=12`.

772 **7.4. Performance: acoustic wave in tilted transverse isotropy model.**

773 This sections illustrates the performance of `tti` with the `core` backend. `tti` cannot
774 be run on the `yask` backend in *Devito v3.1* as some fundamental features are still
775 missing; this is part of our future work (more details in Section 8).

776 Unlike `isotropic`, `tti` significantly benefits from different levels of DSE optimiza-
777 tions, which play a key role in reducing the operation count as well as the register
778 pressure. Figure 6 displays the performance of `tti` for the usual range of space orders
779 on `skl8180` and `kn17250`, for two different cubic grids.

780 Generally, `tti` does not reach the same level of performance as `isotropic`. This
781 is not surprising given the complexity of the PDEs (e.g., in terms of differential oper-
782 ators), which translates into code with much higher arithmetic intensity. In `tti`, the
783 memory system is stressed by a considerably larger number of loads per loop iteration
784 than in `isotropic`. On `skl8180`, we ran some profiling with Intel[®] VTune[™]. We
785 determined that one of the major issues is the pressure on both L1 cache (lack of split
786 registers, unavailability of “fill buffers” to handle requests to the other levels of the
787 hierarchy) and DRAM (bandwidth and latency). Clearly, this is only a summary from
788 some sample kernels – the actual situation varies depending on the DSE optimizations
789 as well as the `so` employed.

790 It is remarkable that on both `skl8180` and `kn17250`, and on both grids, the

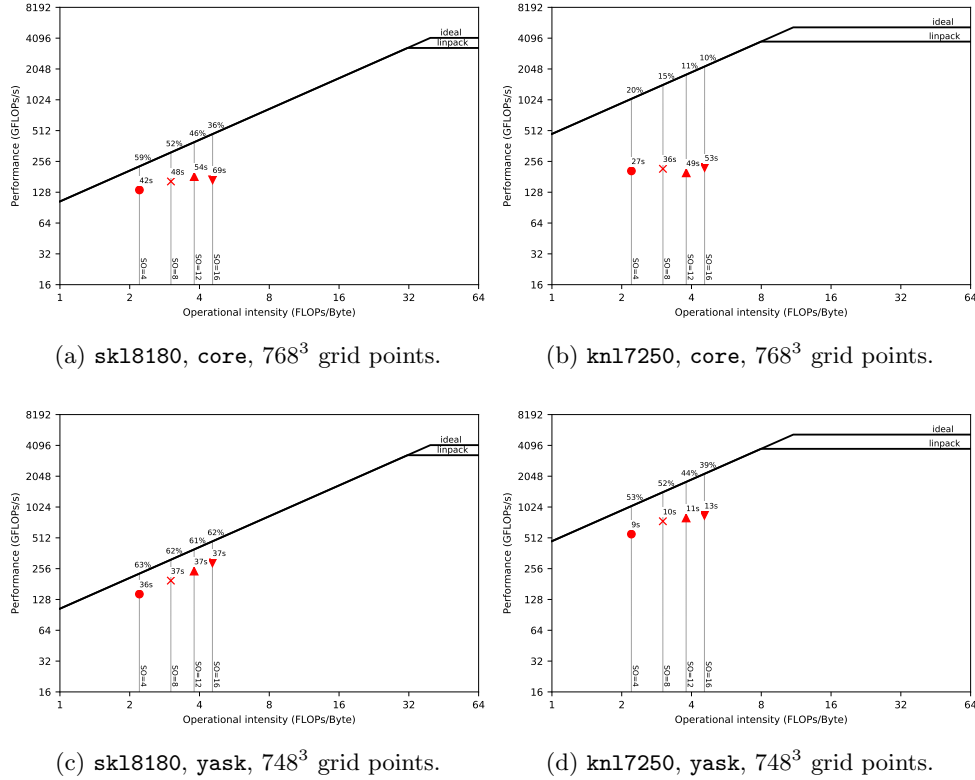


Fig. 5: Performance of `isotropic` on multiple Devito backends and architectures.

791 cutoff point beyond which `advanced` results in worse runtimes than `aggressive` is
 792 `so=8`. One issue with `aggressive` is that to avoid redundant computation, not only
 793 additional memory is required, but also more data communication may occur through
 794 caches, rather than through registers. In Figure 12, for example, we can easily deduce
 795 that `temp` is first stored, and then reloaded in the subsequent loop nest. This is an
 796 overhead that `advanced` does not pay, since temporaries are communicated through
 797 registers, for as much as possible. Beyond `so=8`, however, this overhead is overtaken
 798 by the reduction in operation count, which grows almost quadratically with `so`, as
 799 reported in Table 1.

Table 1: Operation counts for different DSE modes in `tti`

<code>so</code>	<code>basic</code>	<code>advanced</code>	<code>aggressive</code>
4	299	260	102
8	857	707	168
12	1703	1370	234
16	2837	2249	300

800 The performance on `kn17250` is overall disappointing. This is unfortunately
 801 caused by multiple factors – some of which already discussed in the previous sec-

802 tions. These results, and more in general, the need for performance portability across
 803 future (Intel[®] or non-Intel[®]) architectures, motivated the ongoing *yask* project. Here,
 804 the overarching issue is the inability to exploit the multiple levels of parallelism typ-
 805 ical of architectures such as *kn17250*. Approximately 17% of the attainable peak
 806 obtained when *so=4* with *advanced* (best runtime out of the three DSE modes for
 807 the given space order). This occurs when using 512^3 points per grid, which allows
 808 the working set to completely fit in MCDRAM (our calculations estimated a size of
 809 roughly 7.5GB). With the larger grid size (Figure 6d), the working set increases up
 810 to 25.5GB, which exceeds the MCDRAM capacity. This partly accounts for the $5\times$
 811 slow down in runtime (from 34s to 173s) in spite of only a $3\times$ increase in number of
 812 grid points computed per time iteration.

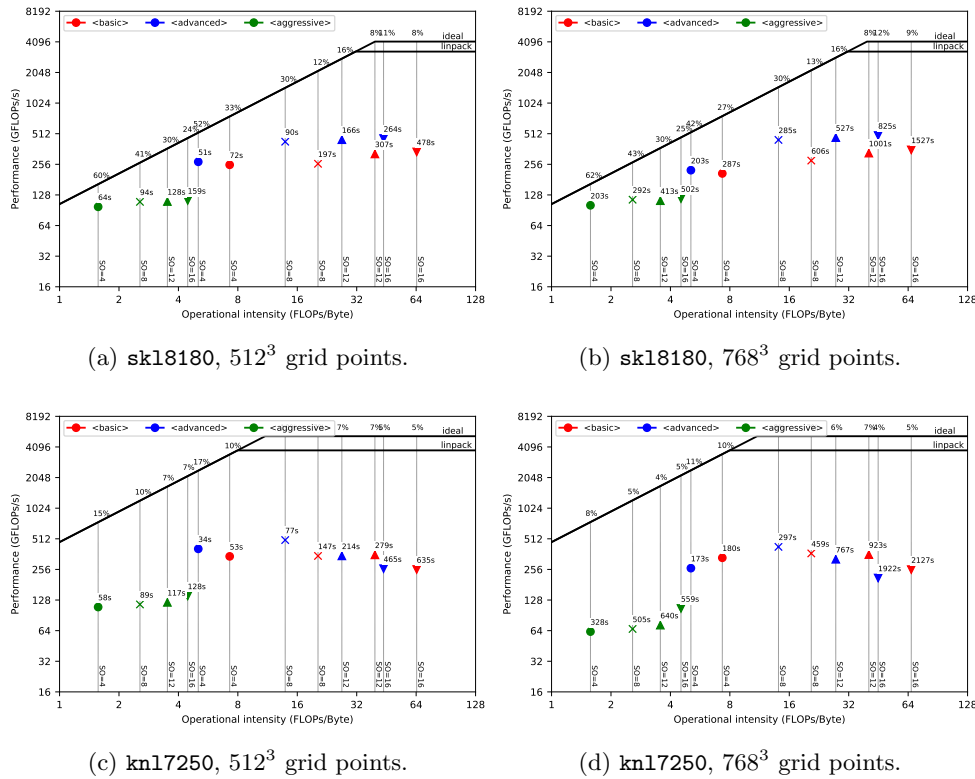


Fig. 6: Performance of *tti* on *core* for different architectures and grids.

813 **8. Further work.** While many simulation and inversion problems such as full-
 814 waveform inversion only require the solver to run on a single shared memory node,
 815 many other applications require support for distributed memory parallelism (typically
 816 via MPI) so that the solver can run across multiple compute nodes. The immediate
 817 plan is to leverage *yask*'s MPI support, and perhaps to include MPI support into *core*
 818 at a later stage. Another important feature is staggered grids, which are necessary for
 819 a wide range of FD discretization methods (e.g. modelling elastic wave propagation).
 820 Basic support for staggered grids is already included in *Devito v3.1*, but currently

821 only through a low-level API – the principle of graceful degradation in action. We
822 plan to make the use of this feature more convenient.

823 As discussed in Section 7.4, the `yask` backend is not feature-complete yet; in
824 particular, it cannot run the `tti` equations in the presence of array temporaries.
825 As `tti` is among the most advanced models for wave propagation used in industry,
826 extending Devito in this direction has high priority.

827 There also is a range of advanced performance optimization techniques that we
828 want to implement, such as “time tiling” (i.e., loop blocking across the time dimen-
829 sion), on-the-fly data compression, and mixed-precision arithmetic exploiting appli-
830 cation knowledge. Finally, there is an on-going effort towards adding an `ops` [31]
831 backend, which will enable code generation for GPUs and also supports distributed
832 memory parallelism via MPI.

833 **9. Conclusions.** Devito is a system to automate high-performance stencil com-
834 putations. While Devito provides a *Python*-based syntax to easily express FD ap-
835 proximations of PDEs, it is not limited to finite differences. A Devito `Operator`
836 can implement arbitrary loop nests, and can evaluate arbitrarily long sequences of
837 heterogeneous expressions such as those arising in FD solvers, linear algebra, or in-
838 terpolation. The compiler technology builds upon years of experience from other
839 DSL-based systems such as FEniCS and Firedrake, and wherever possible Devito
840 uses existing software components including *SymPy* and *NumPy*, and YASK. The
841 experiments in this article show that Devito can generate production-level code with
842 compelling performance on state-of-the-art architectures.

843

REFERENCES

- 844 [1] *Compilers: principles, techniques, and tools*, Pearson/Addison Wesley, Boston, MA, USA,
845 second ed., 2007, <http://www.loc.gov/catdir/toc/ecip0618/2006024333.html>.
- 846 [2] M. S. ALNÆS, A. LOGG, K. B. ØLGAARD, M. E. ROGNES, AND G. N. WELLS, *Unified Form Lan-*
847 *guage: a domain-specific language for weak formulations of partial differential equations*,
848 *ACM Transactions on Mathematical Software (TOMS)*, 40 (2014), p. 9.
- 849 [3] A. ARBONA, B. MIÑANO, A. RIGO, C. BONA, C. PALENZUELA, A. ARTIGUES, C. BONA-CASAS,
850 AND J. MASSÓ, *Simflowny 2: An upgraded platform for scientific modeling and simulation*,
851 arXiv preprint arXiv:1702.04715, (2017).
- 852 [4] U. BONDHUGULA, A. HARTONO, J. RAMANUJAM, AND P. SADAYAPPAN, *A practical automatic*
853 *polyhedral parallelizer and locality optimizer*, in Proceedings of the 2008 ACM SIGPLAN
854 Conference on Programming Language Design and Implementation, PLDI '08, New York,
855 NY, USA, 2008, ACM, pp. 101–113, <https://doi.org/10.1145/1375581.1375595>, <http://doi.acm.org/10.1145/1375581.1375595>.
- 856 [5] S. C. BRENNER AND L. R. SCOTT, *The Mathematical Theory of Finite Element Methods*, vol. 15,
857 Springer New York, New York, NY, 2008, <https://doi.org/10.1007/978-0-387-75934-0>,
858 <http://dx.doi.org/10.1007/978-0-387-75934-0>.
- 859 [6] A. F. CÁRDENAS AND W. J. KARPLUS, *Pdel—a language for partial differential equations*,
860 *Communications of the ACM*, 13 (1970), pp. 184–191.
- 861 [7] G. O. COOK JR, *Alpal: A tool for the development of large-scale simulation codes*, tech. report,
862 Lawrence Livermore National Lab., CA (USA), 1988.
- 863 [8] R. COURANT, K. FRIEDRICHS, AND H. LEWY, *On the partial difference equations of mathemati-*
864 *cal physics*, *International Business Machines (IBM) Journal of Research and Development*,
865 11 (1967), pp. 215–234, <https://doi.org/10.1147/rd.112.0215>.
- 866 [9] K. DATTA, S. WILLIAMS, V. VOLKOV, J. CARTER, L. OLIKER, J. SHALF, AND K. YELICK,
867 *Auto-tuning the 27-point stencil for multicore*, in In Proc. iWAPT2009: The Fourth
868 International Workshop on Automatic Performance Tuning, 2009.
- 869 [10] S. J. DEITZ, B. L. CHAMBERLAIN, AND L. SNYDER, *Eliminating redundancies in sum-of-product*
870 *array computations*, in Proceedings of the 15th International Conference on Supercomput-
871 ing, ICS '01, New York, NY, USA, 2001, ACM, pp. 65–77, <https://doi.org/10.1145/377792.377807>,
872 <http://doi.acm.org/10.1145/377792.377807>.
- 873

- 874 [11] Y. DING AND X. SHEN, *Glore: Generalized loop redundancy elimination upon ler-notation*,
875 Proc. ACM Program. Lang., 1 (2017), pp. 74:1–74:28, <https://doi.org/10.1145/3133898>,
876 <http://doi.acm.org/10.1145/3133898>.
- 877 [12] S. FOMEL, P. SAVA, I. VLAD, Y. LIU, AND V. BASHKARDIN, *Madagascar: open-source software*
878 *project for multidimensional data analysis and reproducible computational experiments*,
879 Journal of Open Research Software, 1 (2013), p. e8, <https://doi.org/http://dx.doi.org/10.5334/jors.ag>.
- 880 [13] T. O. FOUNDATION, *OpenFOAM v5 User Guide*, <https://cfd.direct/openfoam/user-guide/>.
- 881 [14] K. A. HAWICK AND D. P. PLAYNE, *Simulation software generation using a domain-specific lan-*
882 *guage for partial differential field equations*, in 11th International Conference on Software
883 Engineering Research and Practice (SERP'13), no. CSTN-187, Las Vegas, USA, 22-25 July
884 2013, WorldComp, p. SER3829.
- 885 [15] R. L. HIGDON, *Numerical absorbing boundary conditions for the wave equation*, Mathematics
886 of Computation, 49 (1987), pp. 65–90, <http://www.jstor.org/stable/2008250>.
- 887 [16] C. T. JACOBS, S. P. JAMMY, AND N. D. SANDHAM, *Opensbli: A framework for the automated*
888 *derivation and parallel execution of finite difference solvers on a range of computer archi-*
889 *tectures*, CoRR, abs/1609.01277 (2016), <http://arxiv.org/abs/1609.01277>.
- 890 [17] J. JEFFERS AND J. REINDERS, *High Performance Parallelism Pearls Volume Two: Multicore*
891 *and Many-core Programming Approaches*, Morgan Kaufmann Publishers Inc., San Fran-
892 cisco, CA, USA, 1st ed., 2015.
- 893 [18] A. KLÖCKNER, *Loo.py: transformation-based code generation for GPUs and CPUs*, in Proceed-
894 ings of ARRAY '14: ACM SIGPLAN Workshop on Libraries, Languages, and Compilers for
895 Array Programming, Edinburgh, Scotland., 2014, Association for Computing Machinery,
896 <https://doi.org/10.1145/2627373.2627387>.
- 897 [19] A. KLCKNER, *Cgen - c/c++ source generation from an ast*. <https://github.com/inducer/cgen>,
898 2016.
- 899 [20] S. KRONAWITTER, S. KUCKUK, AND C. LENGAUER, *Redundancy elimination in the exastencils*
900 *code generator*, in ICA3PP Workshops, 2016.
- 901 [21] C. LENGAUER, S. APEL, M. BOLTEN, A. GRÖSSLINGER, F. HANNIG, H. KÖSTLER, U. RÜDE,
902 J. TEICH, A. GREBHahn, S. KRONAWITTER, S. KUCKUK, H. RITTICH, AND C. SCHMITT,
903 *Exastencils: Advanced stencil-code engineering*, in Euro-Par 2014: Parallel Process-
904 ing Workshops - Euro-Par 2014 International Workshops, Porto, Portugal, August 25-
905 26, 2014, Revised Selected Papers, Part II, 2014, pp. 553–564, https://doi.org/10.1007/978-3-319-14313-2_47,
906 https://doi.org/10.1007/978-3-319-14313-2_47.
- 907 [22] A. LOGG, K.-A. MARDAL, G. N. WELLS, ET AL., *Automated Solution of Differential Equations*
908 *by the Finite Element Method*, Springer, 2012, <https://doi.org/10.1007/978-3-642-23099-8>.
- 909 [23] M. LOUBOUTIN, M. LANGE, F. J. HERRMANN, N. KUKREJA, AND G. GORMAN, *Per-*
910 *formance prediction of finite-difference solvers for different computer architectures*,
911 Computers & Geosciences, 105 (2017), pp. 148–157, <https://doi.org/https://doi.org/10.1016/j.cageo.2017.04.014>,
912 <https://www.slim.eos.ubc.ca/Publications/Public/Journals/ComputersAndGeosciences/2016/louboutin2016ppf/louboutin2016ppf.pdf>. (Computers &
913 Geosciences).
- 914 [24] M. LOUBOUTIN, M. LANGE, F. LUPORINI, N. KUKREJA, P. A. WITTE, P. VELESKO, G. GORMAN,
915 AND F. J. HERRMANN, *Devito: A portable and flexible mathematical api for geophysical*
916 *applications*. 2018.
- 917 [25] G. R. MARKALL, F. RATHGEBER, L. MITCHELL, N. LORIAN, C. BERTOLLI, D. A. HAM, AND
918 P. H. J. KELLY, *Performance-portable finite element assembly using pyop2 and fenics*, in
919 28th International Supercomputing Conference, ISC, Proceedings, J. M. Kunkel, T. Lud-
920 wig, and H. W. Meuer, eds., vol. 7905 of Lecture Notes in Computer Science, Springer,
921 2013, pp. 279–289, https://doi.org/10.1007/978-3-642-38750-0_21, http://dx.doi.org/10.1007/978-3-642-38750-0_21.
- 922 [26] MATHIAS LOUBOUTIN, FABIO LUPORINI, *Boundary conditions in Devito*, in preparation (2018).
- 923 [27] A. MEURER, C. P. SMITH, M. PAPROCKI, O. ČERTÍK, S. B. KIRPICHEV, M. ROCKLIN, A. KU-
924 MAR, S. IVANOV, J. K. MOORE, S. SINGH, T. RATHNAYAKE, S. VIG, B. E. GRANGER, R. P.
925 MULLER, F. BONAZZI, H. GUPTA, S. VATS, F. JOHANSSON, F. PEDREGOSA, M. J. CURRY,
926 A. R. TERREL, v. ROUČKA, A. SABOO, I. FERNANDO, S. KULAL, R. CIMRMAN, AND A. SCO-
927 PATZ, *Sympy: symbolic computing in python*, PeerJ Computer Science, 3 (2017), p. e103,
928 <https://doi.org/10.7717/peerj-cs.103>, <https://doi.org/10.7717/peerj-cs.103>.
- 929 [28] S. J. PENNYCOOK, J. SEWALL, AND V. LEE, *A metric for performance portability*, arXiv preprint
930 arXiv:1611.07409, (2016).
- 931 [29] J. RAGAN-KELLEY, C. BARNES, A. ADAMS, S. PARIS, F. DURAND, AND S. AMARASINGHE,
932 *Halide: A language and compiler for optimizing parallelism, locality, and recomputation*

- 936 *in image processing pipelines*, in Proceedings of the 34th ACM SIGPLAN Conference
937 on Programming Language Design and Implementation, PLDI '13, New York, NY, USA,
938 2013, ACM, pp. 519–530, <https://doi.org/10.1145/2491956.2462176>, <http://doi.acm.org/10.1145/2491956.2462176>.
- 940 [30] F. RATHGEBER, D. A. HAM, L. MITCHELL, M. LANGE, F. LUPORINI, A. T. T. MCRAE, G.-T.
941 BERCEA, G. R. MARKALL, AND P. H. J. KELLY, *Firedrake: Automating the finite element
942 method by composing abstractions*, ACM Trans. Math. Softw., 43 (2016), pp. 24:1–24:27,
943 <https://doi.org/10.1145/2998441>, <http://doi.acm.org/10.1145/2998441>.
- 944 [31] I. Z. REGULY, G. R. MUDALIGE, M. B. GILES, D. CURRAN, AND S. MCINTOSH-SMITH, *The
945 ops domain specific abstraction for multi-block structured grid computations*, in Pro-
946 ceedings of the Fourth International Workshop on Domain-Specific Languages and High-
947 Level Frameworks for High Performance Computing, WOLFHPC '14, Piscataway, NJ,
948 USA, 2014, IEEE Press, pp. 58–67, <https://doi.org/10.1109/WOLFHPC.2014.7>, <http://dx.doi.org/10.1109/WOLFHPC.2014.7>.
- 950 [32] W. W. SYMES, D. SUN, AND M. ENRIQUEZ, *From modelling to inversion: designing a well-
951 adapted simulator*, Geophysical Prospecting, 59 (2011), pp. 814–833, <https://doi.org/10.1111/j.1365-2478.2011.00977.x>, <http://dx.doi.org/10.1111/j.1365-2478.2011.00977.x>.
- 953 [33] J. TOBIN, A. BREUER, A. HEINECKE, C. YOUNT, AND Y. CUI, *Accelerating seismic simulations
954 using the intel xeon phi knights landing processor*, in Proceedings of ISC High Performance
955 2017 (ISC17), to appear 2017.
- 956 [34] Y. UMETANI, *Deqsol a numerical simulation language for vector/parallel processors*, Proc. IFIP
957 TC2/WG22, 1985, 5 (1985), pp. 147–164.
- 958 [35] R. VAN ENGELEN, L. WOLTERS, AND G. CATS, *Ctadel: A generator of multi-platform high
959 performance codes for pde-based scientific applications*, in Proceedings of the 10th inter-
960 national conference on Supercomputing, ACM, 1996, pp. 86–93.
- 961 [36] F. WITHERDEN, A. FARRINGTON, AND P. VINCENT, *Pyfr: An open source frame-
962 work for solving advectiondiffusion type problems on streaming architectures using
963 the flux reconstruction approach*, Computer Physics Communications, 185 (2014),
964 pp. 3028 – 3040, <https://doi.org/https://doi.org/10.1016/j.cpc.2014.07.011>, <http://www.sciencedirect.com/science/article/pii/S0010465514002549>.
- 966 [37] C. YOUNT, *Vector folding: Improving stencil performance via multi-dimensional simd-vector
967 representation*, in Proceedings of the IEEE 17th International Conference on High Perfor-
968 mance Computing and Communications (HPCC), Aug 2015, pp. 865–870, <https://doi.org/10.1109/HPCC-CSS-ICCESS.2015.27>.
- 970 [38] C. YOUNT AND A. DURAN, *Effective use of large high-bandwidth memory caches in HPC stencil
971 computation via temporal wave-front tiling*, in Proceedings of the 7th International
972 Workshop in Performance Modeling, Benchmarking and Simulation of High Performance
973 Computer Systems held as part of ACM/IEEE Supercomputing 2016 (SC16), PMBS'16,
974 Nov 2016.
- 975 [39] C. YOUNT, A. DURAN, AND J. TOBIN, *Multi-level spatial and temporal tiling for efficient hpc
976 stencil computation on many-core processors with large shared caches*, Future Generation
977 Computer Systems, (2017), <https://doi.org/https://doi.org/10.1016/j.future.2017.10.041>,
978 <http://www.sciencedirect.com/science/article/pii/S0167739X17304648>.
- 979 [40] C. YOUNT, J. TOBIN, A. BREUER, AND A. DURAN, *Yask—yet another stencil kernel: a framework
980 for hpc stencil code-generation and tuning*, in Proceedings of the 6th International Work-
981 shop on Domain-Specific Languages and High-Level Frameworks for High Performance
982 Computing held as part of ACM/IEEE Supercomputing 2016 (SC16), WOLFHPC'16, Nov
983 2016, <https://doi.org/10.1109/WOLFHPC.2016.08>.
- 984 [41] ZENODO/DEVITO, *Devito v3.1*, October 2017, <https://doi.org/10.5281/zenodo.836688>.
- 985 [42] ZENODO/DEVITO-PERFORMANCE, *Devito Experimentation Framework*, July 2018, <https://doi.org/TODO>.
- 986 [43] Y. ZHANG AND F. MUELLER, *Auto-generation and auto-tuning of 3d stencil codes on gpu
987 clusters*, in Proceedings of the Tenth International Symposium on Code Generation
988 and Optimization, CGO '12, New York, NY, USA, 2012, ACM, pp. 155–164, <https://doi.org/10.1145/2259016.2259037>, <http://doi.acm.org/10.1145/2259016.2259037>.
- 990 [44] Y. ZHANG, H. ZHANG, AND G. ZHANG, *A stable tti reverse time migration and its
991 implementation*, GEOPHYSICS, 76 (2011), pp. WA3–WA11, <https://doi.org/10.1190/1.3554411>, <https://doi.org/10.1190/1.3554411>, <https://arxiv.org/abs/https://doi.org/10.1190/1.3554411>.