

# A Domain-Specific Optimizing Compiler for Solving Partial Differential Equations

Fabio Luporini\*, Florian Rathgeber\*, Ana Lucia Varbanescu<sup>†</sup>, Gheorghe-Teodor Bercea\*,  
Paul H. J. Kelly\*, and David Ham\*

\*Imperial College London, Department of Computing,

<sup>†</sup>TU Delft, Faculty of Engineering, Mathematics and Computer Science (EWI)

\*{f.luporini12|f.rathgeber10|gheorghe-teodor.bercea08|p.kelly|d.ham@imperial.ac.uk}@imperial.ac.uk,

<sup>†</sup>a.l.varbanescu@uva.nl

**Abstract**—The abstract goes here.

## I. INTRODUCTION

In many fields, like computational fluid dynamics, computational electromagnetics, and structural mechanics, programs model phenomena by means of partial differential equations (PDEs). Numerical techniques, like Finite Volume Method and Finite Element Method, are widely employed in real-world applications to approximate solutions of PDEs. Unstructured meshes are often used to discretize the equation domain, since their geometric flexibility allows solvers to be extremely effective. The solution is sought in each cell of the discretized domain by applying suitable numerical kernels. As the number of cells can be of the order of millions, a major issue is the time required to execute the computation, which can be hours or days. To address this problem, domain-specific languages (DSLs) have been developed. The successful porting of Hydra, a computational fluid dynamics industrial application devised by Rolls Royce for turbomachinery design (based on Finite Volume Method, roughly 50000 lines of code and mesh sizes that can be over 100 millions edges), to OP2 [3], demonstrates the effectiveness of the DSL approach for implementing PDEs solvers [?].

OP2 adopts a kernel-oriented programming model, in which the computation semantics is expressed through self-contained functions (“kernels”). A kernel is applied to all elements in a set of mesh components (e.g. edges, vertices, elements), with an implicit synchronization between the application of two consecutive kernels. On commodity multi-cores, a kernel is executed sequentially by a thread, while parallelism is achieved partitioning the mesh and assigning each partition to a thread. Similar programming and execution models are adopted by [?], [?], [2]. Kernel optimization is one of the major concerns in unstructured mesh applications. In this paper, we tackle this problem by proposing a domain-driven optimization strategy for a class of kernels used in Finite Element Methods.

We focus on Local Assembly (“assembly”, in the following), a fundamental step of a Finite Element Method that covers an important fraction of the overall computation runtime, typically in the range 30%-60%. During the assembly phase, the solution of the PDE is approximated by executing a suitable kernel over all elements in the discretized domain. A kernel’s working set is usually small enough to fit the L1 cache; it might need L2 cache when high-order methods are employed

to improve the accuracy of the solution. However, we do not consider the latter case. An assembly kernel is characterized by the presence of an affine, usually non-perfect loop nest, where individual loops are rather small (the trip count rarely exceeds 40, with a minimum value of 3, depending on the order of the method). With such small kernels, we focus on aspects like minimization of floating-point operations, register allocation and instruction-level parallelism, especially in the form of SIMD vectorisation. Our study is conducted in the context of Firedrake, a system for solving PDEs through the Finite Element Method based on the OP2 abstraction [?].

Optimization of assembly kernels is non-trivial. Polyhedral compilation could be used to improve performance; however, given the exceptionally small size of the kernels, effective register locality can be obtained only through an in-depth analysis of the data space, rather than by reasoning on the iteration space (our attempts to use PoCC’s optimizations [?], especially unrolling and register tiling, did not improve performance over *icc* compilation). BLAS routines could be employed as well, although fairly complicated code manipulation would be required, for a performance gain potentially far from the expectation. The major issue would concern automating identification and extraction of matrix-matrix multiplications from the code, which requires a non-trivial analysis of control- and data-flow. BLAS libraries are also known to perform distant from peak performance when the dimension of the matrices is small [5]. In addition, register locality would be affected, since the same matrix could appear in a number of products. Hand-made BLAS-based implementations of the Helmholtz assembly kernel (illustrated later), varying the polynomial order from  $p = 1$  to  $p = 4$ , has showed that the performance is worse than that achieved with our optimization strategy.

Given the constraints on polyhedral compilation and linear-algebra-specialized libraries, we address assembly optimization by studying a set of domain-specific code transformations, applicable to a wide class of problems. We have developed an optimizing compiler and we have integrated it with Firedrake. This allows us to evaluate our code transformations in a range of real-world problems, varying parameters that impact both solution accuracy and kernel cost, namely the polynomial order of the method (from  $p = 1$  to  $p = 4$ ) and the mesh type (2D, 3D, 3D hybrid structured-unstructured).

Performance results show that Firedrake-generated code for non-trivial assembly kernels, if not adequately optimized, is sub-optimal. A cost-model-driven sequence of source-to-

```

1 void helmholtz(double A[3][3], double **coords) {
2   // K, det = Compute Jacobian (coords)
3
4   static const double W3[3] = {...}
5   static const double X_D10[3][3] = {...}
6   static const double X_D01[3][3] = {...}
7
8   for (int i = 0; i<3; i++)
9     for (int j = 0; j<3; j++)
10      for (int k = 0; k<3; k++)
11        A[j][k] += ((Y[i][k]*Y[i][j]+
12          +((K[1]*X_D10[i][k]+K[3]*X_D01[i][k])*
13            *(K[1]*X_D10[i][j]+K[3]*X_D01[i][j]))+
14            +((K[0]*X_D10[i][k]+K[2]*X_D01[i][k])*
15              *(K[0]*X_D10[i][j]+K[2]*X_D01[i][j])))*
16            *det*W3[i]);
17 }

```

**Algorithm 1:** Local assembly code generated by Firedrake for a Helmholtz problem (2D mesh, Lagrange  $p = 1$  elements).

source code transformations, aimed at improving register allocation and register data locality, can result in performance improvements up to 33% over “softly-optimized” code (i.e. where only basic transformations are performed, such as auto-vectorizable loop-invariant code motion, padding and data alignment), and up to 77% over the original kernels. The contribution of this paper is therefore twofold:

- An optimisation strategy for a class of kernels widely-used in scientific applications, namely Local Assembly in the context of the Finite Element Method. Our approach exploits domain knowledge to go beyond the limits of both vendor (e.g. *icc*) and research (e.g. those based on the polyhedral model) compilers.
- Automation of such code optimizations in Firedrake. This enables performing an in-depth performance evaluation of the proposed optimization strategy in real-world simulations.

The paper is organized as follows. ...

## II. BACKGROUND

Local assembly consists of evaluating so called element stiffness matrix and vector; in this work, we focus on computation of matrices, which is the costly part of the process. A stiffness matrix can be intuitively thought as an approximated representation of the PDE solution in a specific cell of the discretized domain. Numerical integration algorithms are widely-used in assembly codes to implement evaluation of stiffness matrices [4], [?].

Given a mathematical description of the input problem, expressed through the domain-specific Unified Form Language [?], Firedrake generates C-code kernels implementing assembly using a numerical integration algorithm. It then triggers compilation of such kernels using an available vendor compiler, and eventually manages parallel execution over the mesh. In our work, we have enhanced this execution model by adding an optimization stage prior to the generation of C code. The code transformations described next are also generalizable to non-Firedrake assembly kernels, provided that numerical

```

1 void burgers(double A[12][12], double **c, double **w) {
2   // K, det = Compute Jacobian (c)
3
4   static const double W5[5] = {...}
5   static const double X1_D001[5][12] = {...}
6   static const double X2_D001[5][12] = {...}
7   //It follows 11 other basis functions definitions.
8   ...
9
10  for (int i = 0; i<5; i++) {
11    double F0 = 0.0;
12    //It follows 10 other declarations (F1, F2,...)
13    ...
14    for (int r = 0; r<12; r++) {
15      F0 += (w[r][0]*X1_D001[ip][r]);
16      //It follows 10 analogous statements (F1, F2, ...)
17    }
18    ...
19    for (int j = 0; j<12; j++)
20      for (int k = 0; k<12; k++)
21        A[j][k] += (.(K[5]*F9)+(K[8]*F10))*Y1[i][j])+
22          +(((K[0]*X1_D100[i][k])+(K[3]*X1_D010[i][k])+
23            +(K[6]*X1_D001[i][k]))*Y2[i][j]))*F11)+
24          +(.((K[2]*X2_D100[i][k])+...+(K[8]*X2_D001[i][k])*
25            *(K[2]*X2_D100[i][j])+...+(K[8]*X2_D001[i][j]))..)+
26            + <roughly a hundred of sum/muls go here>..)*
27            *det*W5[i]);
28  }
29 }

```

**Algorithm 2:** Local assembly code generated by FFC for a Burgers problem (3D mesh, Lagrange  $p = 1$  elements).

integration is used. Notable examples are [?], [?], [?], in which the kernel’s loop nest and memory access pattern are reducible to that generated by Firedrake.

The complexity of Firedrake-generated kernels depends on the mathematical problem being solved. In simpler cases, the loop nest is perfect, it has short trip counts (in the range 3-15), and the computation reduces to a summation of a few products. An example is provided in Figure 1, which shows an assembly kernel for a Helmholtz problem, using Lagrange basis functions on 2D elements with polynomial order  $p = 1$ . The stiffness matrix in the code is called  $A$ . In other scenarios, for instance when solving a non-linear problem like Burgers (see Figure 2), the number of arrays involved in the computation of  $A$  can be much larger: in this case, 14 unique arrays are accessed, and the same array can be referenced multiple times within the expression. Also, constants evaluated in outer loops (called  $F$  in the code), acting as scaling factors of arrays, may be required; trip counts can be larger (proportionally to the order of the method); arrays may be block-sparse. Note that in addition to a larger number of operations to compute the stiffness matrix, the Burgers case shows a register pressure higher than that in Helmholtz. Despite assembly kernels being problem-dependent, and so infinite is the space of codes that Firedrake can generate, it is still possible to identify common domain-specific traits, which can be exploited for effective code transformations and SIMD vectorization.

The class of kernels we are considering has, in particular, some peculiarities. 1) The computation of the Jacobian, which is the first step of the assembly, is independent of the loop

```

1 void helmholtz(double A[3][4], double **coords) {
2   // K, det = Compute Jacobian (coords)
3
4   static const double W3[3] = {...}
5   static const double X_D10[3][4] = {...}
6   static const double X_D01[3][4] = {...}
7
8   for (int i = 0; i<3; i++) {
9     double LI_0[4];
10    double LI_1[4];
11    for (int r = 0; r<4; r++) {
12      LI_0[r] = ((K[1]*X_D10[ip][r])+(K[3]*X_D01[ip][r]));
13      LI_1[r] = ((K[0]*X_D10[ip][r])+(K[2]*X_D01[ip][r]));
14    }
15    for (int j = 0; j<3; j++)
16      for (int k = 0; k<4; k++)
17        A[j][k] += (Y[i][k]*Y[i][j]+LI_0[k]*LI_0[j]+
18                  +LI_1[k]*LI_1[j])*det*W3[i];
19  }
20 }

```

**Algorithm 3:** Local assembly code generated by Firedrake for a Helmholtz problem (2D mesh, Lagrange  $p = 1$  elements). The padding, data alignment and *licm* optimizations are applied. Data alignment and padding relate to an AVX machine. In this specific case, sub-expressions invariant to  $j$  are identical to those invariant to  $k$ ; in general, this could not be the case.

nest; this is not true in general, since the unstructured mesh might use bended elements that would require the Jacobian be re-computed every  $i$  iteration. 2) Memory accesses along the three loop dimensions are always 1-stride. 3) The  $j$  and  $k$  loops are interchangeable, whereas permutation of  $i$  is might be subjected to pre-computation of values and introduction of temporaries (e.g. the  $F$  values in Burgers). 4) So called basis functions arrays (denoted by  $X$ ,  $Y$ , ...) are constants, and most of the sub-expressions on the right hand side of the stiffness matrix computation depend on just two loops (either  $i$ - $j$  or  $i$ - $k$ ). In Section III we show how to exploit these observations to define a set of systematic, composable optimizations.

### III. DOMAIN-DRIVEN CODE TRANSFORMATIONS

#### A. Padding and Data Alignment

Auto-vectorization of assembly code computing the stiffness matrix can be less effective if data are not aligned and if the length of the innermost loop is smaller than the vector length ( $vl$ ). Data alignment is enforced in two steps. Initially, both arrays and matrices are allocated to addresses that are multiples of  $vl$ . Then, matrices are padded by rounding the number of columns to the nearest multiple of  $vl$ . For example, assume the original size of a matrix is  $3 \times 3$  and that the underlying machine possesses AVX, which implies  $vl = 4$  since a vector register is 256 bits long and our kernels use 64-bits double-precision floating-point values. Then, a padded matrix on this architecture will have size  $3 \times 4$ . The compiler is explicitly informed about data alignment using a suitable pragma. Padding of all matrices involved in the evaluation of the stiffness matrix allows us to safely round the loop trip count to the nearest multiple of  $vl$ . This avoids the introduction of a remainder (scalar) loop from the compiler, which would be responsible for inefficient vectorization.

```

1 void helmholtz(double A[3][4], double **coords) {
2   // K, det = Compute Jacobian (coords)
3   // Declaration of basis function matrices
4
5   for (int i = 0; i<3; i++) {
6     double LI_0[4];
7     double LI_1[4];
8     for (int r = 0; r<4; r++) {
9       LI_0[r] = ((K[1]*X_D10[ip][r])+(K[3]*X_D01[ip][r]));
10      LI_1[r] = ((K[0]*X_D10[ip][r])+(K[2]*X_D01[ip][r]));
11    }
12    for (int j = 0; j<3; j++)
13      for (int k = 0; k<4; k++)
14        A[j][k] += (Y[i][k]*Y[i][j]+LI_0[k]*LI_0[j])*det*W3[i];
15    for (int j = 0; j<3; j++)
16      for (int k = 0; k<4; k++)
17        A[j][k] += LI_1[k]*LI_1[j]*det*W3[i];
18  }
19 }

```

**Algorithm 4:** Local assembly code generated by Firedrake for a Helmholtz problem (2D mesh, Lagrange  $p = 1$  elements). The padding, data alignment, *licm* and *split* optimizations are applied. Data alignment and padding relate to an AVX machine. In this specific case, the split factor is 2.

#### B. Auto-vectorizable Loop-invariant Code Motion

From inspection of the codes in Figures 1 and 2, it can be noticed that the computation of  $A$  involves evaluating many sub-expressions that depend on two iteration variables only. Since symbols in most of these sub-expressions are read-only variables, there is ample space for loop-invariant code motion. Vendor compilers apply this technique, although not in the systematic way we need for our assembly kernels. We want to overcome two deficiencies that both *icc* and *gcc* have. First, these compilers can identify sub-expressions that are invariant with respect to the innermost loop only. This is an issue for sub-expressions depending on  $i$ - $k$ , which are not automatically lifted. Second, the hoisted code is scalar, i.e. it is not subjected to auto-vectorization. We work around these limitations with source-level loop-invariant code motion. In particular, we pre-compute all values that an invariant sub-expression assumes along the fastest varying dimension. This is implemented by introducing a temporary array (per invariant sub-expression) and by adding a new loop to the nest. At the price of extra memory for storing temporaries, the gain is that lifted terms can be auto-vectorized, because part of an inner loop. Note that given the short trip counts of our loops, it is important to achieve auto-vectorization of all hoisted terms in order to minimize the percentage of scalar instructions, which could be otherwise significant.

Figure 3 shows the Helmholtz assembly code after the application of loop-invariant code motion, padding, and data alignment.

#### C. Optimizing Register Management

One notable problem assembly kernels are exposed to concerns register allocation and register locality. The critical situation occurs when loop trip counts and accessed variables are such that the vector-registers pressure is high. Due to a kernel's working set fitting the L1 cache, it is remarkably

important to optimize register management in order to maximize performance. Canonical optimizations, such as loop interchange, unroll, unroll-and-jam and register tiling are typically employed to deal with this problem. In the compiler we have developed, we support these optimizations either by means of explicit source code transformation (interchange, unroll-and-jam) or indirectly advising the compiler through standard pragmas (unroll). In this context, register tiling is a critical optimization. A way of obtaining register tiles consists of slicing the innermost loop into rectangular blocks of identical size (except the “reminder” block), and applying unrolling to individual blocks in order to expose ILP. This is automated by some vendor compilers [?]. Register tiles along both inner dimensions can be determined as well. Based on the observation that the evaluation of the stiffness matrix can be reduced to a “summation of outer products” along the  $j$  and  $k$  dimensions, we implement a domain-specific vector-register tiling strategy that leads to larger tiles, enhancing reuse. If we consider the code snippet in Figure 3 (Helmholtz after loop-invariant code motion), we can notice that the computation of  $A$  is abstractly expressible as

$$A_{jk} = \sum_{\substack{x \in B' \subseteq B \\ y \in B'' \subseteq B}} x_j \cdot y_k \quad j, k = 0, \dots, 4 \quad (1)$$

where  $B$  is the set of all matrices (or temporaries) accessed in the kernel, whereas  $B'$  and  $B''$  are generic problem-dependent subsets. Without loss of generality, the presence of constants or other variables independent of both  $j$  and  $k$  can be momentarily neglected. Note that regardless of the specific input problem, the stiffness matrix computation is always reducible to this kind of form. Figure III-C illustrates how we can evaluate 16 elements ( $j, k = 0, \dots, 4$ ) of the stiffness matrix using just 2 vector registers (a  $4 \times 4$  tile), assuming  $|B'| = |B''| = 1$ . Values in a register are shuffled each time a product is performed. Standard compiler auto-vectorization (*gcc* and *icc*), instead, executes 4 broadcast operations (i.e. “splat” of a value over all of the register locations) along the outer dimension to perform the calculation, and would also need to keep between  $f = 1$  and  $f = 3$  extra registers to perform the same 16 evaluations if unroll-and-jam is found to be useful, with  $f$  being the unroll-and-jam factor.

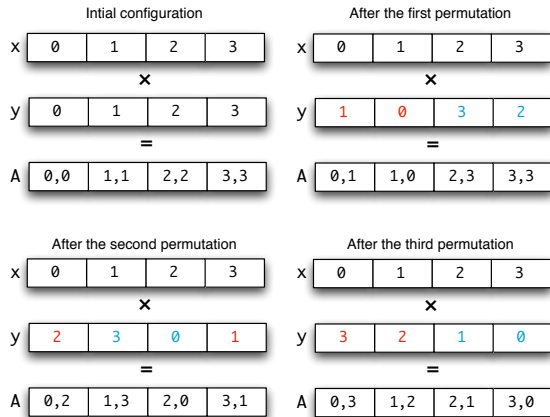


Fig. 1. Outer-product vectorization by permuting values in a vector register.

The storage layout of  $A$ , however, is incorrect after the

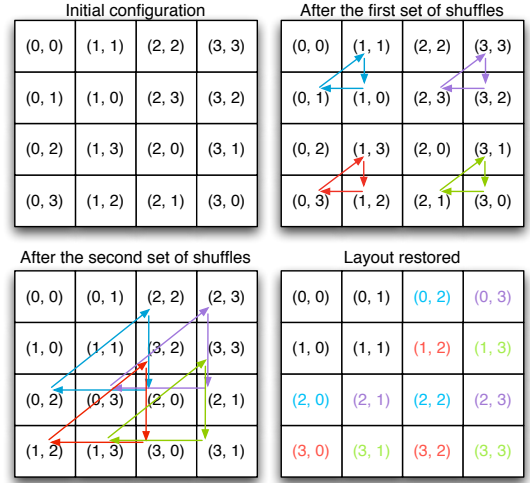


Fig. 2. Restoring the storage layout after *op-vect*. The figure shows how  $4 \times 4$  elements in the top-left block of the stiffness matrix  $A$  can be brought to their correct positions.

application of this “outer-product vectorization” (*op-vect*). We efficiently restore it with a sequence of vector shuffle following the pattern highlighted in Figure III-C, executed once outside of the  $ijk$  loop nest. The generated pseudo-code for the simple Helmholtz problem when using *op-vect* is shown in Figure 5.

#### D. Assembly splitting

In complex kernels, like Burgers in Figure 2, and on certain architectures, achieving effective register allocation can be challenging. If the number of variables independent of the innermost-loop dimension, i.e. basis function matrices, temporaries introduced by loop-invariant code motion and constants, is even close to the number of available CPU registers, it is likely to obtain poor register reuse. For example, applying loop-invariant code motion to Burgers on a 3D mesh needs 33 temporaries for the  $ijk$  loop order, and compiler’s hoisting of invariant loads out of the  $k$  loop can be inefficient on architectures with a relatively low number of registers. One potential solution to this problem consists of suitably “splitting” the computation of the stiffness matrix  $A$  into multiple sub-expressions; an example, for the simpler Helmholtz problem, is given in Figure 4. Splitting an expression has, in general, several drawbacks. Firstly, it increases the number of accesses to  $A$  proportionally to the “split factor”, which is the number of sub-expressions produced. Secondly, depending on how the split is executed, it can lead to redundant computation (e.g. the product  $det * W3[i]$  is performed times number of sub-expressions in the code of Figure 4). Finally, it might affect register locality, although this is not the case of the Helmholtz example: for instance, the same matrix could be accessed in different sub-expressions, requiring a proportional number of loads be performed. Nevertheless, as shown in Section V, the performance gain from improved register reuse along inner dimensions can still be greater, especially if the split factor and the splitting itself use heuristics to minimize the aforementioned issues.

Table I summarizes the code transformations described so far. Given that many of these transformations depend on some

```

1 void helmholtz(double A[8][8], double **coords) {
2   // K, det = Compute Jacobian (coords)
3   // Declaration of basis function matrices
4
5   for (int i = 0; i<6; i++) {
6     // Do loop-invariant code motion
7     for (int j = 0; j<4; j+=4)
8       for (int k = 0; k<8; k+=4) {
9         // Call Load and set intrinsics
10        // Compute A[1,1],A[2,2],A[3,3],A[4,4]
11        // One permute_pd intrinsics per k-loop load
12        // Compute A[1,2],A[2,1],A[3,4],A[4,3]
13        // One permute2f128_pd intrinsics per k-loop load
14        // ...
15      }
16    // Do Remainder loop (from j = 4 to j = 6)
17  }
18  // Restore the storage layout:
19  for (int j = 0; j<4; j+=4) {
20    __m256d r0, r1, r2, r3, r4, r5, r6, r7;
21    for (int k = 0; k<8; k+=4) {
22      r0 = __mm256_load_pd (&A[j+0][k]);
23      // Load A[j+1][k], A[j+2][k], A[j+3][k]
24      r4 = __mm256_unpackhi_pd (r1, r0);
25      r5 = __mm256_unpacklo_pd (r0, r1);
26      r6 = __mm256_unpackhi_pd (r2, r3);
27      r7 = __mm256_unpacklo_pd (r3, r2);
28      r0 = __mm256_permute2f128_pd (r5, r7, 32);
29      r1 = __mm256_permute2f128_pd (r4, r6, 32);
30      r2 = __mm256_permute2f128_pd (r7, r5, 49);
31      r3 = __mm256_permute2f128_pd (r6, r4, 49);
32      __mm256_store_pd (&A[j+0][k], r0);
33      // Store A[j+1][k], A[j+2][k], A[j+3][k]
34    }
35  }
36 }

```

**Algorithm 5:** Local assembly code generated by Firedrake for a Helmholtz problem (2D mesh, Lagrange  $p = 2$  elements). The padding, data alignment, *licm* and *op-vect* optimizations are applied. Data alignment and padding relate to an AVX machine. The original size of the  $j$ - $k$  iteration space (i.e. before padding was applied) was  $6 \times 6$ . In this example, the unroll-and-jam factor is 1.

parameters (e.g. tile size), we need a mechanism to prune such a large space of optimization. This aspect is treated in Section IV

Name (Abbreviation)	Parameter
Auto-vectorizable Loop-invariant code motion ( <i>licm</i> )	
Padding	
Data Alignment	
Loop interchange	loops
Loop unrolling	unroll factor
Register tiling	tile size
Outer-product vectorization ( <i>op-vect</i> )	tile size
Assembly splitting ( <i>split</i> )	split point, split factor

TABLE I. OVERVIEW OF CODE TRANSFORMATIONS FOR FIREDRAKE-GENERATED ASSEMBLY KERNELS.

#### IV. THE PYOP2 COMPILER

Firedrake provides users with the Unified Form Language to write problems in a notation resembling mathematical equations. This high-level specification is translated by the Fenics Form Compiler [?] into an Abstract Syntax Tree representation

of a Finite Element assembly kernel. ASTs are then passed to PyOP2 [?], which lies at the core of Firedrake, where parallel execution over the unstructured mesh is managed. Our compiler, capable of applying the transformations described in Section III, is integrated with PyOP2: it receives FFC’s ASTs as input, it introduces optimizations, and it generates C code as output, which is eventually just-in-time compiled on the underlying architecture.

TODO: say there is a very huge space of implementation opportunities; in this section, we say how to prune it (without autotuning, which would be prohibitive for our jit-compilation).

The compiler structure is outlined in Figure 6. Initially, an AST is inspected, looking for the presence of iteration spaces and other domain-specific information provided by the higher layer. If the kernel lacks an iteration space, then so-called inter-kernel vectorization, in which the non-affine loop over mesh elements is vectorized, can be applied. This feature, currently under development, has been proved to be useful in several Finite-Volume-based applications [?]. The second transformation step is applied if the backend is a manycore machine, like a GPU: the compiler tries to extract parallelism from inside the kernel, by partitioning loop iterations among different threads, if these are found to be independent [?]. Then, a sequence of optimization steps are executed. A cost model is queried to estimate the potential gain if the *split* and *op-vect* transformations were applied. Optimization passes must follow a specific order. Application of *licm* must precede padding and data alignment, due to the introduction of temporary arrays. Then, *split* might take place, depending on the outcome of the cost model. The *op-vect* optimization is introduced lastly. Note that the compiler can handle all possible corner cases: for example, if *op-vect* is to be applied but the size of the iteration space is not a multiple of the vector length, then a reminder loop, amenable to auto-vectorization, is inserted.

Loop unroll and unroll-and-jam of outer loops are fundamental to expose ILP and data reuse, and so tuning critical parameters such as the unroll factor becomes of great importance. Autotuning might be used, although we avoid it to minimize the overhead of Firedrake’s just-in-time compilation. It is our experience (inspection of assembly code, comparison with other hand-made implementations), however, that for our assembly kernels, where the loop nest is affine, bounds are known at compile-time, and memory accesses are unit-stride, newest versions of a vendor compiler like *icc* employ cost models capable of estimating close-to-optimal values for such parameters. We therefore leave the backend compiler in charge to select unroll and unroll-and-jam factors. This choice also simplifies the PyOP2 compiler’s cost model. The only situation in which we explicitly unroll-and-jam a loop is when *op-vect* is used, since the transformed code seems to prevent the *icc* compiler from applying any other unroll-based optimizations.

All loops are interchangeable provided that temporaries are introduced if the nest is not perfect. It would be possible, therefore, to permute loops in order to improve data movement, although we discovered that a specific loop order often leads to the best run-time, for the storage layout we employ. The loop order  $ijk$  is likely to maximize performance. Conceptually, this is motivated by the following considerations: 1) all of

## 1 The PyOP2 Compiler

**Input:** ast, wrapper, isa

**Output:** code

```
2 // Analyze ast and build optimization plan
3 it_space = analyze(ast)
4 if not it_space then
5     ast.apply_inter_kernel_vectorization(wrapper)
6     return wrapper + ast.from_ast_to_c()
7 endif
8 if isa.backend == gpu then
9     if it_space then
10         ast.extract_iteration_space(wrapper)
11     endif
12     return wrapper + ast.from_ast_to_c()
13 endif
13 plan = cost_model(it_space.n_inner_arrays, isa.n_regs)
14 // Optimize ast based on plan
15 ast.apply_licm()
16 ast.padding()
17 ast.data_align()
18 if plan.sz_split then
19     ast.split(plan.sz_split)
20 endif
21 if plan.uaj_factor then
22     ast.vr_tile(plan.uaj_factor)
23 endif
24 return wrapper + ast.from_ast_to_c()
```

**Algorithm 6:** The PyOP2 compiler.

the accessed arrays depend on the  $i$  loop, and having it in one of the inner positions would mean a significantly higher number of load operations, with a performance loss greater than the reduction of store operations into the stiffness matrix; 2) if loop-invariant code motion is applied, then the number of loads along the  $j$  and  $k$  dimensions becomes identical (this property derives from the nature of assembly), which implies the loop order  $ijk$  is performance-wise identical to  $ikj$ . These observations, which allow us to fix a specific loop order, contribute to relieve the decisions the PyOP2 compiler has to make.

The cost model ..... TODO

## 1 Algorithm: ApplyCostModel

**Input:** n\_inner\_arrays, n\_regs

**Output:** uaj\_factor, sz\_split

```
2 sz_split = 0
3 // Compute splitting factor
4 while n_inner_arrays > n_regs do
5     n_inner_arrays = n_inner_arrays / 2
6     sz_split = sz_split + 1
7 endw
8 n_inner_regs = n_regs / 2
9 if n_inner_arrays > n_inner_regs or sz_split > 0 then
10     // Rely on autovectorization, as there might be not
11     // enough registers to improve performance by tiling
12     return <sz_split, 0>
13 endif
13 // Compute unroll-and-jam factor for vector-register
14 // tiling
14 n_regs_avail = n_regs - n_inner_arrays
15 uaj_factor = ⌈n_regs_avail / n_inner_arrays⌉
16 return <sz_split, uaj_factor>
```

**Algorithm 7:** Procedure to estimate the most suitable unroll-and-jam factor and/or split size.

## V. PERFORMANCE EVALUATION

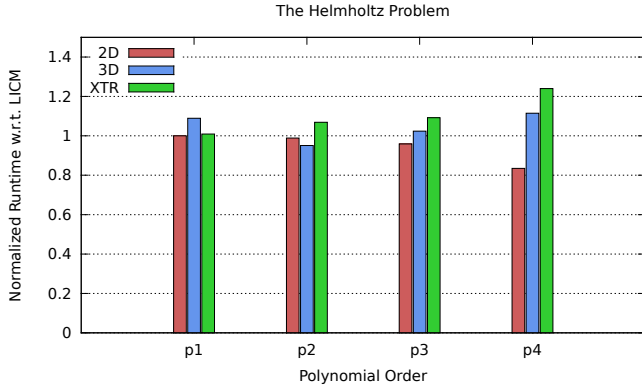


Fig. 3. Helmholtz

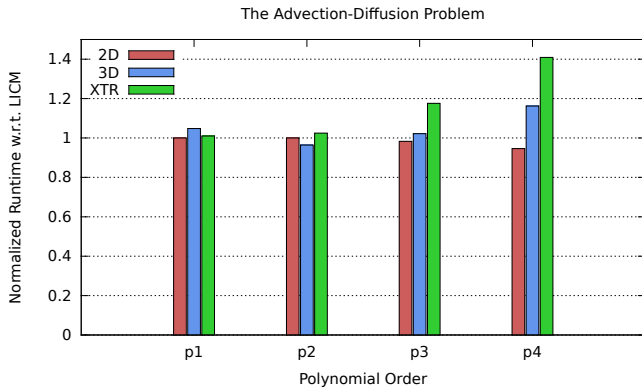


Fig. 4. Advection-Diffusion

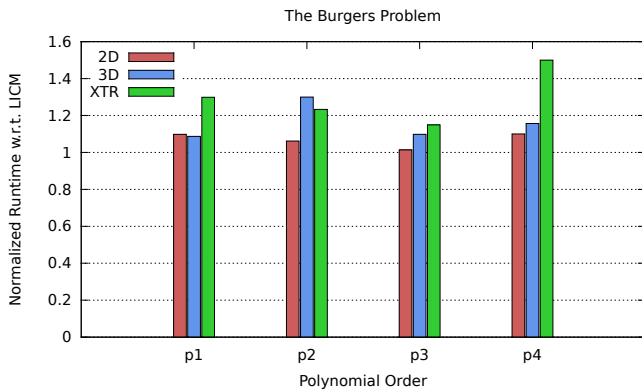


Fig. 5. Burgers

## VI. RELATED WORK

## VII. CONCLUSIONS

## ACKNOWLEDGMENT

The authors would like to thank...

## REFERENCES

- [1] The firedrake project, 2014.

- [2] Zachary DeVito, Niels Joubert, Francisco Palacios, Stephen Oakley, Montserrat Medina, Mike Barrientos, Erich Elsen, Frank Ham, Alex Aiken, Karthik Duraisamy, Eric Darve, Juan Alonso, and Pat Hanrahan. Liszt: A domain specific language for building portable mesh-based pde solvers. In *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis, SC '11*, pages 9:1–9:12, New York, NY, USA, 2011. ACM.
- [3] G. R. Markall, F. Rathgeber, L. Mitchell, N. Lorient, C. Bertolli, D. A. Ham, and P. H. J. Kelly. Performance portable finite element assembly using PyOP2 and FEniCS. In *Proceedings of the International Supercomputing Conference (ISC) '13*, volume 7905 of *Lecture Notes in Computer Science*, June 2013. In press.
- [4] Kristian B. Olgaard and Garth N. Wells. Optimizations for quadrature representations of finite element tensors through automated code generation. *ACM Trans. Math. Softw.*, 37(1):8:1–8:23, January 2010.
- [5] Jaewook Shin, Mary W. Hall, Jacqueline Chame, Chun Chen, Paul F. Fischer, and Paul D. Hovland. Speeding up nek5000 with autotuning and specialization. In *Proceedings of the 24th ACM International Conference on Supercomputing, ICS '10*, pages 253–262, New York, NY, USA, 2010. ACM.