

Vectorizing Sparse Matrix Codes with Dependency Driven Trace Analysis

Zachary Cetinic
University of Toronto
Toronto, Canada
zachary.cetinic@mail.utoronto.ca

Kazem Cheshmi
University of Toronto
Toronto, Canada
kazem@cs.toronto.edu

Maryam Mehri Dehnavi
University of Toronto
Toronto, Canada
mmehride@cs.toronto.edu

Abstract—Sparse computations frequently appear in scientific simulations and the performance of these simulations rely heavily on the optimization of the sparse codes. The compact data structures and irregular computation patterns in sparse matrix computations introduce challenges to vectorizing these codes. Available approaches primarily vectorize regular regions of computations in the sparse code. They also reorganize data and computations, at a cost, to increase the number of regular regions. In this work, we propose a novel polyhedral model, called the partially strided codelets (PSC), that enables the vectorization of computation regions with irregular data access patterns. PSCs also improve data locality in sparse computation. Our DDF inspector-executor framework efficiently mines the memory accesses in the sparse computation, using an access function differentiation approach, to find PSC codelets. It generates vectorized code for the sparse matrix multiplication kernel (SpMV), a kernel with parallel outer loops, and for kernels with carried dependence, specifically the sparse triangular solver (SpTRSV). We demonstrate the performance of the DDF-generated code on a set of 60 large and small matrices (0.05-130M nonzeros). DDF outperforms the highly specialized library MKL with an average speedup of 1.93 and 4.5 \times for SpMV and SpTRSV, respectively. For the same matrices, DDF outperforms the state-of-the-art inspector-executor framework Sympiler [1] for the SpTRSV kernel by up to 11 \times and the work by Augustine et al [2] for the SpMV kernel by up to 12 \times .

Index Terms—Vectorization, Sparse Matrix Computations, Polyhedral Analysis

I. INTRODUCTION

Sparse matrix computations are important kernels used in a large class of scientific and machine learning applications. The performance of sparse kernels is noticeably improved if the code is *vectorized* to exploit single instruction multiple data (SIMD) capabilities of the underlying architecture. Vectorization potentially increases opportunities to optimize for locality, further increasing the performance of the sparse code. SIMD instructions can efficiently vectorize groups of operations that access consecutive data, i.e. have a regular access pattern. A computation is regular when its operations access memory addresses with a constant distance from each other, i.e. *strided access*, and otherwise irregular. However, because a sparse matrix is typically stored in a compact representation [3], the sparse matrix code accesses are often irregular, making vectorization challenging. Not vectorizing irregular segments in the sparse matrix calculation can potentially reduce locality and the performance of the sparse matrix kernel.

The most common approach to improve vectorization of irregular regions in sparse matrix codes is to reorganize data or computations so the region becomes regular and amenable to vectorization. These work typically optimize sparse kernels with no loop-carried dependencies, for example [4]–[10] improve the performance of sparse matrix-vector multiplication (SpMV) with reorganization. CSR5 [11] and CVR [12] reorganize data to create contiguous accesses and reorganize computations to increase the number of independent operations for vectorization. However, their computation reorganization typically increases the number of floating-point operations in the kernel, reducing the benefits of vectorization. Sparse storage formats such as ELL and DIA [13] are used to pad the sparse matrix with additional nonzeros. Padding increases regular regions in the matrix and hence increases the number of consecutive vector loads. In sparse matrices with very irregular patterns, the amount of padding will increase, leading to more floating-point operations. Reorganization approaches are often library-based and thus are specialized for a specific set of matrix patterns or the access patterns appearing in a specific matrix. They also need to be manually ported to new architectures. For example, work in [14] primarily optimizes matrices from power-law graph computations, CVR is for sparse graphs, and ELL is most efficient for block-structured matrices such as factors of a direct solver [15].

Compilers, such as those based on the polyhedral model, generate vectorizable and portable code for affine codes such as dense matrix computations [16]–[21]. However, compilers are limited in vectorizing sparse codes because of the existing indirection in the memory access patterns of sparse computations, caused by data compaction. This indirection introduces additional challenges to compilers when optimizing sparse kernels with loop carried dependencies, e.g. the sparse triangular solver. Recently the polyhedral method was extended to optimize sparse matrix computations via inspector-executor approaches [22]–[24]. The inspectors execute at runtime to resolve indirection and find data dependencies between operations and this information is used to transform the original source code into optimized code called the executor.

Amongst inspector-executor frameworks, Augustine et al. [2], which we refer to as the PIC framework, and Sympiler [1] inspect the sparsity pattern of input matrices to generate vectorized code for sparse matrix computations. The work in

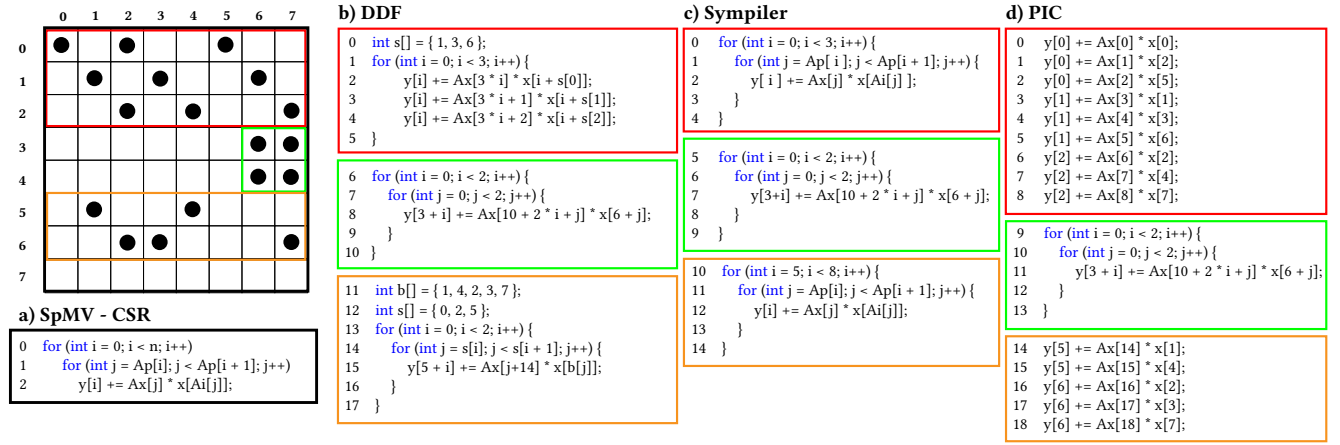


Fig. 1. Various generated codes for the given sparsity pattern (shown in the matrix) used to compute the Sparse Matrix-Vector multiplication (SpMV) kernel. **a)** Shows the code that computes SpMV using the compressed sparse row (CSR) storage format. **b)** Shows the DDF generated code. **c)** Shows the code generated from the Sympiler framework [1]. **d)** Shows the code generated using the work by Augustine *et al.* [2]

[2] proposes an inspector to mine the memory trace of the SpMV kernel to find regular regions in the computation for vectorization. The remaining regions are unrolled. Unrolling irregular regions limits the scalability of this work because of increasing code size for large matrices and for matrices with high irregularity; their framework can primarily optimize small sparse matrices (with less than 0.5 million nonzeros). Also, they do not support sparse kernels with loop-carried dependencies. Sympiler uses an inspector to resolve data dependencies in sparse kernels such as SPTRV. Its inspector finds iterations that operate on rows with similar patterns, i.e. regular regions, and applies code transformations that primarily optimize for tiling and vectorization. Sympiler adds additional nonzeros to some of the rows with non-consecutive memory accesses and thus pads the input matrix to form *row-blocks*. Row-blocks are regular regions that can be easily vectorized. To limit the overheads from padding, Sympiler only creates row-blocks when profitable, thus limiting the size of row blocks it creates. It is thus not able to efficiently optimize sparse matrices with few rows of similar sparsity patterns.

The standard approach of prior work in vectorizing sparse kernels is to either find regions with regular computations or create more regular regions in the computations with for example data reorganization. The focus of this work is to expose vectorization opportunities in computation regions with irregular computation patterns. We present a novel polyhedral model which we call the *partially strided codelet* that increases opportunities for vectorization of sparse matrix codes and improves data locality. PSCs represent computation regions that have non-consecutive, i.e. unstrided, memory accesses but can still benefit from vectorization because a subset of their accesses is strided. We also present a novel inspector-executor framework called DDF (Differentiating Data access Functions) that efficiently inspects the memory access patterns in sparse matrix codes to mine for PSCs using an *access function differentiation* approach and to generate optimized vectorized

code for the matrix computation. DDF's inspector supports sparse kernels with and without loop-carried dependencies as it is able to look for data dependencies in the computation region.

The contributions of the work are:

- A novel polyhedral model, called the partially strided codelet, that enables the vectorization of computation regions with non-consecutive memory accesses in sparse codes while improving temporal and spatial locality.
- A memory access differentiation approach, that uses the *first order partial difference (FOPD)* of access functions to distinguish strided access from non-strided accesses in the sparse matrix computation. FOPD's are used to mine for PSCs in sparse matrix codes.
- The DDF inspector-executor framework that mines the data accesses of sparse codes with or without dependencies to find PSCs. The inspector executes in parallel with low overhead and generates a compact code that can execute in parallel.
- We demonstrate the performance of DDF for SpTRSV and SpMV kernels. The speedup of the DDF-specialized code for SpTRSV is on average 1.79 \times and 4.5 \times faster than that of Sympiler and MKL respectively. For SpMV, DDF-generated code is on average 1.42 \times and 1.93 \times faster than that of CSR5 and MKL respectively. For small problems that the work from [2] supports, DDF is 10.57 \times faster than their approach.

II. MOTIVATION

In this section, we use the matrix in Figure 1 to explain the approach used in inspector-executor frameworks, i.e. Sympiler and the work in [2] (which we call PIC), to optimize different computation regions in the SpMV kernel and compare to the approach proposed in this work. The example uses the compressed sparse row (CSR) format to store matrix A and then computes $y = A \times x$ where x is a vector. The region with

$$I = [i_0, i_1, n] \left\{ \begin{array}{l} \text{Iteration Space} \\ i_0 \geq 0 \wedge i_0 < n \\ i_1 \geq \text{Ap}[i_0] \wedge i_1 < \text{Ap}[i_0 + 1] \end{array} \right\} \quad \begin{array}{l} \text{Access Functions} \\ y[*]: f[i_0, i_1] = i_0 \\ \text{Ax}[*]: g[i_0, i_1] = i_1 \\ x[*]: h[i_0, i_1] = \text{Ai}[i_1] \end{array}$$

Fig. 2. Polyhedral representation of the SpMV kernel shown in Figure 1a

the green border is a regular computation region in the SpMV computation and the red and orange colors show regions with non-strided accesses.

All tools can efficiently vectorize the green region since all of its accesses are consecutive and thus strided. BLAS routines are typically used to optimize such regions. They enable vectorization, and in part, also improve locality by reusing the entries in the input vector (e.g. $x[6]$, $x[7]$) and the output vector (e.g. $y[3]$, $y[4]$).

For the non-strided segments, Sympiler uses a padding strategy, when profitable, to create regular regions that can be easily vectorized. However, Sympiler will not pad the red and yellow regions, since they are too irregular, and padding will introduce a large number of additional operations. Instead, Sympiler will fall back to the original non-affine code and miss any additional vectorization opportunities. PIC will generate fully unrolled code for the non-strided regions, as shown in Figure 1d. While this reduces the number of instructions, it increases code size and as a result, increases the number of instruction cache misses, especially for large matrices, limiting PIC's scalability.

For irregular regions, DDF uses partially strided codelets to improve locality and to provide a better vectorization efficiency. The SpMV code for the red region has non-strided accesses in x , while its accesses to Ax and y are strided. DDF generates the PSC codelet bordered with red in Figure 1b. The code stores the non-strided addresses in s to be reused across all three iterations of i instead of loading indices for every operation. This also enables the reuse of index $\text{Ap}[i]$ in Ax so it does not have to be loaded per each outermost iteration. Similarly, for the yellow region, accesses to both Ax and x are non-strided, but vectorizing computations of the innermost loop improves spatial locality due to contiguous accesses in Ax .

III. PARTIALLY STRIDED CODELET

This section introduces partially strided codelets and discusses their classification as well as efficiency in optimizing computation regions with irregular access patterns. We also present a novel cost model and a differentiation-based PSC detection strategy, both of which are used in the DDF framework to efficiently find PSCs in sparse matrix computations.

A. Definitions

Polyhedral model and data access functions. A loop nest that contains a set of statements is represented with a polyhedral model through an integer polyhedron sets \mathcal{I} and relations f . A statement is made of a data space described with \mathcal{D} which is a disjoint set containing $\mathcal{D}_0, \dots, \mathcal{D}_n$. An integer polyhedral

$$\begin{array}{l} \text{a)} \\ y[0] += \text{Ax}[0] * x[0] \\ y[0] += \text{Ax}[1] * x[2] \\ y[0] += \text{Ax}[2] * x[5] \\ y[1] += \text{Ax}[3] * x[1] \\ y[1] += \text{Ax}[4] * x[3] \\ y[1] += \text{Ax}[5] * x[6] \\ y[2] += \text{Ax}[6] * x[2] \\ y[2] += \text{Ax}[7] * x[4] \\ y[2] += \text{Ax}[8] * x[7] \end{array} \quad \begin{array}{l} \text{b)} \\ I = [i_0, i_1] \left\{ \begin{array}{l} i_0 \geq 0 \wedge i_0 < 3 \\ i_1 \geq 0 \wedge i_1 < 3 \end{array} \right\} \\ s = [0, 2, 5] \\ y[*]: f[i_0, i_1] = i_0 \\ \text{Ax}[*]: g[i_0, i_1] = 3 * i_0 + i_1 \\ x[*]: h[i_0, i_1] = i_0 + s[i_1] \end{array}$$

Fig. 3. a) Shows the set of SpMV operations from the red outlined region of the matrix in Figure 1. b) Shows the iteration space and access functions of a codelet that maps to the data-space seen in section a) of the figure.

set $\mathcal{I} = [i_0, \dots, i_n]$ is a collection of inequalities that create bounds for each dimension inside $i \in \mathcal{I}$. For each $\mathcal{D}_d \in \mathcal{D}$ a *data access function* f is used to describe how the data space \mathcal{D}_d is accessed by the iteration space of \mathcal{I} . In other words, a data access function maps an iteration space to a data space, i.e. $f_{\mathcal{I} \rightarrow \mathcal{D}}$. The SpMV code in Figure 1a has one statement. That statement has three data spaces y , Ax , and x as well as three data access functions. Figure 2 shows the polyhedron sets for $\mathcal{I} = [i_0, i_1]$ and also the access functions corresponding to each data space.

Codelet. A polyhedral model that has a convex integer polyhedron with no flow dependencies, i.e. read after write dependency between access functions, is a codelet. A codelet only has one statement and the operation in that statement is a SIMD-supported operation. A set of SpMV operations is shown in Figure 3a and its polyhedral representation including iteration space and data access functions are shown in Figure 3b. All operations in this model are independent and have only one statement ($y[*] += \text{Ax}[*] * x[*]$), with a multiply-add operation which is supported by Intel and AMD SIMD units. These properties satisfy the criteria for being a codelet.

Strided data access function. A function that can be expressed with a linear combination of induction variables in \mathcal{I} is a strided access function. The data access function f and g in Figure 3b are both strided and can be expressed as a linear combination of $[i_0, i_1]$.

B. Partially Strided Codelet Classification

An efficient way to vectorize a codelet is to find operations with strided accesses across different iterations. This enables the vectorization of more operations and potentially increases data reuse between different iterations. However, current approaches are primarily limited to vectorizing codelets that all of their access functions are strided. An efficient BLAS [25] implementation is used for this purpose and thus, we call these codelets types BLAS codelets. In this work we define a novel set of codelets called *Partially Strided Codelets* that have at most $n-1$ strided access functions and at least one non-strided access function. These codelets can benefit from vectorization because they have one or more strided access functions.

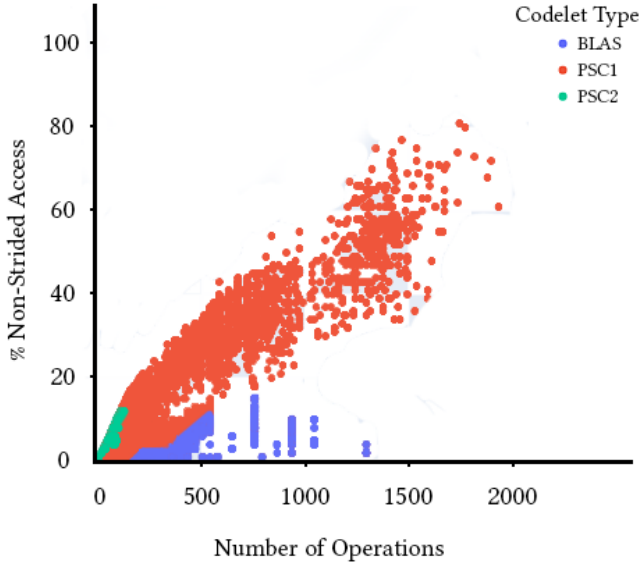


Fig. 4. Shows best performing codelet type on a wide variety (150K total points) of computational regions found in SpMV/SpTRSV.

PSCs are classified into different types based on the number of access functions that are strided. For example in SpMV and SpTRSV that have codelets with three access functions, $n = 3$, two types of PSCs can be defined. The PSC type I codelet is used when two of access functions are strided, and the PSC type II codelet is used when only one access function is strided. Figure 1b shows two types of PSCs, outlined in yellow and red, and also the BLAS codelet, outlined in green.

C. PSC Efficiency and Cost Model

Depending on the number of operations and strided accesses in a computation region, either PSC or BLAS codelets can be profitable. To illustrate this, we extracted over 150,000 computation regions, with a different number of operations and strided accesses, from the SpMV and SpTRSV kernels on 60 matrices in the Suitesparse repository [26]. Each region is optimized using one or several codelets of the same type and the best performing codelet type for that region is inserted as a point in Figure 4. The red, green, and blue colors are used for PSC I, PSC II, or BLAS codelets respectively. The figure shows that some regions perform best with PSCs and some with BLAS codelets, e.g. PSC I provides the fastest execution time over all codelets for regions with a large percentage of non-strided accesses.

The performance of codelets differs based on the number of memory accesses that occur from running the codelets. Per Section III-A, an access function is described as $f(\mathcal{I}) = x + s[i_0, \dots, i_n]$ where x is an integer and s is an n -dimensional array. If f is a strided access function, an integer $s[i]$ for each $i \in \mathcal{I}$ needs to be loaded from memory before the operation in the codelets executes. Each integer, $s[i]$, is the i^{th} coefficient in the linear combination of f . However, for a non-strided access function, the number of memory accesses in the codelet

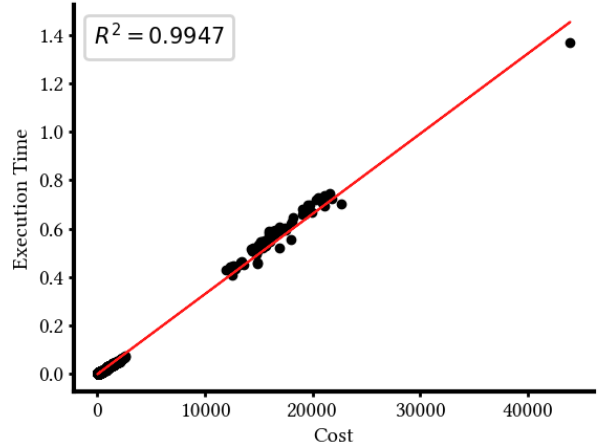


Fig. 5. Correlation between predicted cost of codelets and average execution time of each codelet. Averages were obtained using 1,000,000 executions per codelet.

is equal to the size of $s[i]$ for $i \in \mathcal{I}$, which is larger than one. Thus, compared to BLAS codelets, PSCs access more memory addresses. But this does not imply that PSC codelets are inefficient because the number of codelets used to cover a computational region also affects the total memory accesses. For example, the computation region in Figure 3b is a PSC I with one non-strided access function $h(i_0, i_1) = i_0 + s[i_1]$ where $s = \{0, 2, 5\}$ and thus requires loading 4 memory addresses. If BLAS is selected for the same region, a minimum of 3 codelets will need to be used leading to 9 loads.

To efficiently determine the types of codelets for a computation region, we first define a cost model to compute the cost of a single codelet. The cost model uses the number of memory accesses of each memory access function, i.e. $|s_d| + 1$ and the number of operations in a codelet ($|p|$) to estimate the execution time of the codelet. The cost of a codelet l is c_l defined as:

$$c_l = |p| + m + \sum_{d=0}^m |s_d| \quad (1)$$

If a computation region is described with k codelet types then the overall cost of the codelets for that region is $C = \sum_{l=0}^k c_l$. Figure 5 shows the correlation between the codelet cost model and the execution time of codelets. We use a set of codelets with different sizes from SpMV and SpTRSV on different computation regions of the sparse matrices in Suitesparse [27]. The X-axis shows the cost of each codelet, and the Y-axis shows their corresponding execution time in seconds. As shown, the cost model predicts the performance of the codelets with a correlation of ($r^2 = 0.99$).

D. Differentiation Based PSC Detection

In this section, we explain how the *first order partial difference* (FOPD) of the access functions in a computation region can be used to detect a codelet type.

$$\begin{aligned}
& \begin{array}{ccc} i_0=0 & i_0=1 & i_0=2 \\ f(i_0, i_1) = & \boxed{[[0,0,0], [1,1,1], [2,2,2]]} \\ g(i_0, i_1) = & \boxed{[[0,1,2], [3,4,5], [6,7,8]]} \\ \Delta_{i_0} h(i_0, i_1) = & \begin{array}{c} \begin{array}{ccc} [[1,1,1], [1,1,1]] \\ \diagdown \quad \diagup \quad \diagdown \quad \diagup \\ \diagdown \quad \diagup \quad \diagdown \quad \diagup \\ \diagdown \quad \diagup \quad \diagdown \quad \diagup \end{array} \\ = FOPD_{i_0} \\ h(i_0, i_1) = & \boxed{[[0,2,5], [1,3,6], [2,4,7]]} \\ \Delta_{i_1} h(i_0, i_1) = & \begin{array}{ccc} \begin{array}{c} \sqcup \sqcup \\ \sqcup \sqcup \\ \sqcup \sqcup \end{array} & \begin{array}{c} \sqcup \sqcup \\ \sqcup \sqcup \\ \sqcup \sqcup \end{array} & \begin{array}{c} \sqcup \sqcup \\ \sqcup \sqcup \\ \sqcup \sqcup \end{array} \\ [[2,3], [2,3], [2,3]] & = FOPD_{i_1} \end{array} \end{array}
\end{aligned}$$

Fig. 6. Illustrates the process of taking the derivative of the access function h with respect to i_0 and i_1 .

First Order Partial Differentiation (FOPD _{i}). Given the data access function f with iteration space of $\mathcal{I} = [i_0, i_1]$, the first order partial difference of f with respect to $i_1 \in \mathcal{I}$ is computed as $\frac{\Delta}{\Delta i_1} f(\mathcal{I}) = \Delta_{i_1} f = f(i_0, i_1 + 1) - f(i_0, i_1)$. FOPD shows if accesses to a data space are strided with respect to the induction variable i_1 . Figure 6 illustrates the process of computing the FOPD for the access function h given the computation region shown in Figure 3. For example, the FOPD of h evaluated at $i_0 = 1, i_1 = 1$ is $\Delta_{i_1} h(1, 1) = 3$.

FOPD of access functions are used to distinguish types of codelets by finding strided access functions. Given a codelet with three access functions and the iteration space of $\mathcal{I} = [i_0, i_1]$, an access function f is strided if its FOPDs with respect to \mathcal{I} are equal in the entire iteration space. In other words, f is strided if the elements in $\Delta_{i_0} f(i_0, i_1)$ are equal to each other and similarly for elements in $\Delta_{i_1} f(i_0, i_1)$. With the strided definition above, all codelet types can be defined per definitions in Section III-B. For example, the first three accesses in function $g(i_0, i_1)$ in Figure 6 are to consecutive locations (0, 1, 2) in $\mathbb{A}\mathbb{X}$. The FOPD of the first two accesses, wrt. i_1 , is $FOPD_{i_1}(0, 0) : g(0, 1) - g(0, 0) = 1$ and for the second and third accesses is $FOPD_{i_1}(0, 1) : g(0, 2) - g(0, 1) = 1$. Similarly the FOPD of the accesses for i_0 are $FOPD_{i_0}(0, 0) : g(1, 0) - g(0, 0) = 3$ and $FOPD_{i_0}(1, 0) : g(2, 0) - g(1, 0) = 3$. Since $FOPD_{i_1}(0, 0)$ and $FOPD_{i_1}(0, 1)$ are equal, and $FOPD_{i_0}(0, 0)$ and $FOPD_{i_0}(1, 0)$ are equal, the function in the iteration space $\mathcal{I} = \{i = 0 \wedge 0 \leq j \leq 3\}$ is strided. Because function f is also strided and h is not strided, the codelet is categorized as a PSC type I.

```

1 #include "DDF.h"
2 int main() {
3     Kernel SpMV;
4     SpMV.generate_c("spmv.h", Arch::x86);
5     return 1;
6 }

```

Listing 1. The input code to DDF.

```

1 #include "DDF.h"
2 #include "spmv.h"
3 int main() {

```

Algorithm 1: DDF Inspector

```

Input :  $\mathcal{K}, \mathcal{P}$ 
Output:  $cList$ 
/* 1) Grouping independent memory accesses */
1  $f_0, f_1, f_2, \mathcal{I} \leftarrow \text{computeAccessFunctions}(\mathcal{K}, \mathcal{P});$ 
2  $\mathcal{G} \leftarrow \text{findDependencies}(f_0, f_1, f_2, \mathcal{I});$ 
3  $\mathcal{S} \leftarrow \text{partitionIterationSpace}(\mathcal{G}, \mathcal{I});$ 
/* 2) Codelet Creation */
4 for  $p \in \mathcal{S}$  do
5     FOPDs  $\leftarrow \text{ComputeFOPD}(p, f_0, f_1, f_2);$ 
6      $\mathcal{R} \leftarrow \text{GetConsecutiveIterations}(p);$ 
7     for  $r \in \mathcal{R}$  do
8          $clb \leftarrow \text{BLAS\_first}(r, \text{FOPDs});$ 
9          $clp1 \leftarrow \text{PSCI\_first}(r, \text{FOPDs});$ 
10         $clp2 \leftarrow \text{PSCII\_first}(r, \text{FOPDs});$ 
11         $\text{min\_cl} \leftarrow \text{min\_cost}(clb, clp1, clp2);$ 
12         $cList.append(\text{min\_cl});$ 

```

```

4 Matrix A("input.mtx");
5 Vector x("x.mtx"), y("y.mtx");
6 // ----- Inspector ----- //
7 vector<codelet *> clist=DDF_Inspector(A.pattern(),
8     SpMV);
9 // ----- Executor ----- //
9 SpMV_Vectorized(clist, A, x, y);
10 return 1;
11 }

```

Listing 2. A sample driver code to call the DDF-generated code.

IV. DDF: DIFFERENTIATING DATA ACCESS FUNCTIONS TO MINE PSC CODELETS

DDF is an inspector-executor framework that finds partially strided codelets from an input code and the matrix sparsity pattern using an inspector and generates a vectorized code as its executor. The DDF input specification is shown in Listing 1. It takes the input kernel, SpMV or SpTRSV, and the target architecture type and generates a driver code shown in Listing 2 and a vectorized code in `spmv.h`. The vectorized code is generated when `generate_c` is called in line 4 of Listing 1. In lines 4–5 of the Listing 2, the input matrix A and vectors x and y are loaded. DDF's inspector then creates a list of codelets in line 7, and the vectorized code executes in line 9.

A. The DDF Inspector

DDF's inspector uses Algorithm 1 to generate efficient code for the input sparse kernel. The algorithm first creates groups of independent iterations and then for the computation regions in each group finds a best codelet combination to generate.

Inputs and output: The inputs to the DDF inspector algorithm are a kernel code \mathcal{K} , i.e., SpMV or SpTRSV, and the pattern of the input matrix \mathcal{P} . The algorithm generates a list of codelets $cList$ to be used in the vectorized code. In DDF, the kernel code \mathcal{K} is represented with an abstract syntax tree (AST) that represents loops and operations. The pattern \mathcal{P} is stored in a compressed sparse row (CSR) storage.

1) *Grouping independent memory accesses*: The first step of the inspector algorithm computes memory access functions of the input code and then generates a graph \mathcal{G} and uses it to create groups of operations that are independent to support parallelism. Function *computeAccessFunctions* in line 1 of Algorithm 1 uses kernel \mathcal{K} and goes over \mathcal{P} to compute access functions and the iteration space of the input kernel. Then *findDependencies* in line 2 uses access functions to compute a graph \mathcal{G} . The vertices in this graph represent iterations of the outermost loop $i \in \mathcal{I} = [i, j]$ and its edges are the flow dependencies across iterations of i . In kernels with loop carried dependencies, \mathcal{G} is passed to the LBC algorithm [28] to create independent parallel workloads. For sparse kernels with parallel loops, \mathcal{G} does not have any edges, thus instead, the function groups iterations based on the number of operations. Finally, iterations of the outermost loop are grouped and the groups are stored in \mathcal{S} using *partitionIterationSpace* in line 3.

2) *Codelet creation*: For every group of iterations $p \in \mathcal{S}$, the algorithm generates the FOPD of the access functions in line 5. Function *GetConsecutiveIterations* in line 6 groups operations across t consecutive iterations and creates a number of computation regions stored in \mathcal{R} . To reduce the overhead of mining, the algorithm limits the search window in line 6 to t consecutive iterations. Our experiments show that this does not lead to slowdowns because consecutive iterations in sparse codes typically operate on consecutive rows of a matrix, and thus their vectorization can potentially improve spatial locality.

Lines 8-12 in the algorithm inspect each computation region r and finds the best codelet combination for that region using three different strategies, i.e. *BLAS_first*, *PSCI_first*, and *PSCII_first*. Since the strategies mine for codelets, they take as input the FOPDs computed from line 5 along with r . The BLAS-first strategy prioritizes finding BLAS codelets and then mines for PSCs. The PSC I-first strategy mines for PSC type I codelets first, and the PSC II-strategy looks for PSC type II codelets and then other codelets. To reduce the execution time of each strategy, we store the list of codelets for each mined region and look up this list when executing the other strategies. The costs of codelets from each approach are stored in constants *clb*, *clp1*, and *clp2* and then in line 11 the combination with the lowest cost is chosen. Line 12 appends the most efficient codelet combination found for the current computation region to the output list of codelets in the Algorithm, i.e. *cList*.

B. DDF Executor

The DDF driver runs the executor code, as shown in line 9 of listing 2. To efficiently use the instruction cache and to eliminate the need for recomputing a new executor code per matrix pattern, DDF generates a parametric code. The parametric code for the three types of codelets, BLAS, PSC I, and PSC II are shown in Figure 7. As shown each generic code of a codelet class takes the data access functions as input and can vectorize any codelet of the same class. These generic codelets are called inside a switch statement in the executor

```

a) void BLAS_codelet(int m, int n, Fn f, Fn g, Fn h) {
    for (int i = 0; i < m; i++)
        for (int j = 0; j < n; j++)
            y[ f.x + f.s[0][0] * i + f.s[1][0] * j ] +=
            Ax[ g.x+ g.s[0][0] * i + g.s[1][0] * j ] *
            x[ h.x+ h.s[0][0] * i + h.s[1][0] * j ];
}

b) void PSCI_codelet(int m, int n, Fn f, Fn g, Fn h) {
    for (int i = 0; i < m; i++)
        for (int j = 0; j < n; j++)
            y[ f.x + f.s[0][0] * i + f.s[1][0] * j ] +=
            Ax[ g.x+ g.s[0][0] * i + g.s[1][0] * j ] *
            x[ h.x+ h.s[0][0] * i + h.s[1][j] ];
}

c) void PSCII_codelet(int m, int n, Fn f, Fn g, Fn h) {
    for (int i = 0; i < m; i++)
        for (int j = 0; j < n; j++)
            y[ f.x + f.s[0][i] + f.s[1][0] * j ] +=
            Ax[ g.x+ g.s[0][0] * i + g.s[1][0] * j ] *
            x[ h.x+ h.s[0][0] * i + h.s[1][j] ];
}

```

Fig. 7. Shows generic codes used to evaluate different codelet types; all codelets are parameterized by m, n which bound the outer and inner dimensions respectively. Each type also uses f, g and h which are special structs containing offset information used to process each codelet.

```

1 #include "Codeletsx86.h"
2
3 void vectorized_code(vector<Codlet*> clsit, CSR* A,
4     double *x, double *y){
5     for(int i = 0; i < clsit.partitions(); i++){
6         #pragma omp parallel
7         for(int j = 0; j < clsit[i].size(); j++){
8             switch(clsit[i][j].type){
9                 case BLAS:
10                    BLAS_codelet();
11                    break;
12                 case PSCI:
13                    PSCI_codelet();
14                    break;
15                 case PSCII:
16                    PSCII_codelet();
17                    break;
18             }
19         }
20     }
21 }

```

Listing 3. The vectorized code for SpMV generated by DDF in *spmv.h*.

code, shown in Listing 3. In the executor, each mined codelet in the codelet list is mapped to a generic codelet using this switch statement. A change in the input sparsity pattern results in a different codelet list, however, the same parametric code is still used. For matrices of different size, the codelets list size changes, however, the same executor code in Listing 3 can be used. Thus, the code size is invariant to the size of matrix.

V. RESULTS

We evaluate the performance of DDF using two kernels, sparse triangular solver (SpTRSV) and sparse matrix-vector multiplication (SpMV). DDF is compared to two other inspector-executor approaches, i.e. Sympiler and PIC [2]. We also compare DDF to libraries CSR5 [11], MKL [29], and our implementation of SpMV using the ELLPACK [13] storage format (referred to as ELL in the figures). Sympiler does not support SpMV, and SpTRV is not supported by the CSR5, PIC, and ELL implementations. Hence those tools are omitted from the respective figures.

The set of symmetric positive definite (SPD) matrices is used for evaluation. These matrices are selected from the SPD symmetric matrices in the SuiteSparse repository [26] and are diverse in both size and sparsity pattern. For the SpTRSV evaluations, we only use the lower triangular half of the SPD symmetric matrices. Matrices with ID 0=30 are the thirty largest matrices in the repository, Matrix IDs 31-49 are L-factors [30] of the largest matrices, and 10 small matrices (ID = 50-60) with 50-350K nonzeros are also included. L-factors matrices are included to compare the performance of DDF to Sympiler because Sympiler is best suited to optimize these matrices. Small matrices are included to compare with PIC as it does not scale to larger problems. The test-bed architectures are an Intel(R) Xeon(R) Gold 5115 CPU (2.8GHz, 14080K L3 Cache) with 20 cores and 64GB of main memory (Intel) and an AMD Ryzen 3900x CPU (3.8GHz, 64MB L3 Cache) with 12 cores and 32GB of main memory (AMD). The performance of kernels is reported on both architectures but the analysis is only shown for the Intel machine. All generated code, implementations of different approaches, and library drivers are compiled with GCC v.7.2.0 compiler and with the `-O3` flag. All benchmarks are executed 5 times and the median value of all runs is reported. To report the performance in GFLOP/s, we compute the theoretical floating-point operations for each kernel and matrix and divide it by execution times of each tool.

The sequential implementation of SpMV CSR and SpTRSV CSR are used as a baseline. Since the code for PIC is not publicly available, we created an in-house inspector-executor implementation of their approach with feedback from the authors of [2]. PIC was originally developed for a single thread, however, we extended its support to be parallel and report the best performance between the two implementations in all figures. A timeout of 4 hours was used for all runs (including the inspector and the executor time). Thus, PIC does not have data points in respective figures for large matrices. Also, our implementation of ELL which is based on [13], has missing points in some figures because the number of non-zero fill-ins exceeds the main memory of our test-bed architecture.

A. SpMV Performance

The performance of the SpMV kernel for DDF, MKL, CSR5, ELLPACK, baseline, and PIC is shown in Figure 8. As shown, DDF is faster than all other tools for over 88% of the matrices tested. DDF's code is on average 1.93x, 1.42x,

10.57x, and 7.19x faster than in order MKL, CSR5, ELLPACK, and PIC respectively on the Intel architecture. DDF's code is on average 1.22x, 1.29x, 8.54x, and 10.93x faster than in order MKL, CSR5, ELLPACK, and PIC respectively on the AMD architecture.

To demonstrate the effect of partially strided codelets on the performance of DDF, in Figure 9 we use a stacked bar for DDF. The stacks show the GFlop/s for the baseline code (input to DDF), for DDF when it only mines for BLAS codelets (DDF-baseline+BLAS and refer to DDF BLAS_only), and from running the entire DDF algorithm, i.e. Algorithm 1, that mines for BLAS and PSC (shown with DDF or DDF-baseline+BLAS+PSC in the figure). As shown, DDF is on average 10 \times faster than DDF BLAS_only which demonstrates the importance of using PSC codelets. Figure 9 shows the percentage of operations that are vectorized in the generated code from DDF with PSC I, PSC II, and BLAS codelets, obtained by averaging over all matrices in the benchmark. As shown, over 80% of the operations in SpMV are vectorized with PSC codelets.

The PSC codelets improve data locality in DDF's generated code. Figure 10 shows the relation between DDF's performance and the performance of the parallel (OpenMP) version of the baseline code. Average memory access latency [31] is used as a measure for locality and is computed by gathering the number of misses and accesses to L1, L2, and LLC caches using the PAPI [32] performance counters. Figure 10 shows the coefficient of determination or R^2 is 0.83 which indicates a good correlation between speedup and the memory access latency.

To further explain why DDF is faster than other tools, we conduct a few experiments and report the resulting average over all matrices in this paragraph. For MKL, we compute its average memory access latency, which is 4.36x slower than that of DDF, contributing to the worse performance compared to DDF. To compare to CSR5, we count the number of instructions. CSR5 executes 1.73x more instructions compared to DDF. This is potentially due to the overhead of the segmented sum calculations used in their approach to improving vectorization and load balance. DDF is on average 10.57x faster than ELL, because the percentage of fill-in in the ELL implementation is 800% relative to the non-zero elements in the matrix, significantly increasing the number of operations in their method. To conclude, the SpMV code of DDF is faster than existing implementations because it improves data locality and/or reduces the number of instructions via vectorization.

B. SpTRSV Performance

Figure 11 compares the performance of SpTRSV using DDF, MKL, and Sympiler. The stacked bar of DDF shows the additional performance gained from mining partially strided codelets; the trend is similar to that in Figure 8 for SpMV. On average DDF is faster than MKL, Sympiler, and the baseline 4.5 \times , 1.79 \times , and 3.62 \times respectively for the intel architecture. On the AMD architecture, DDF is faster than MKL and Sympiler 2.66, 1.38 \times respectively.

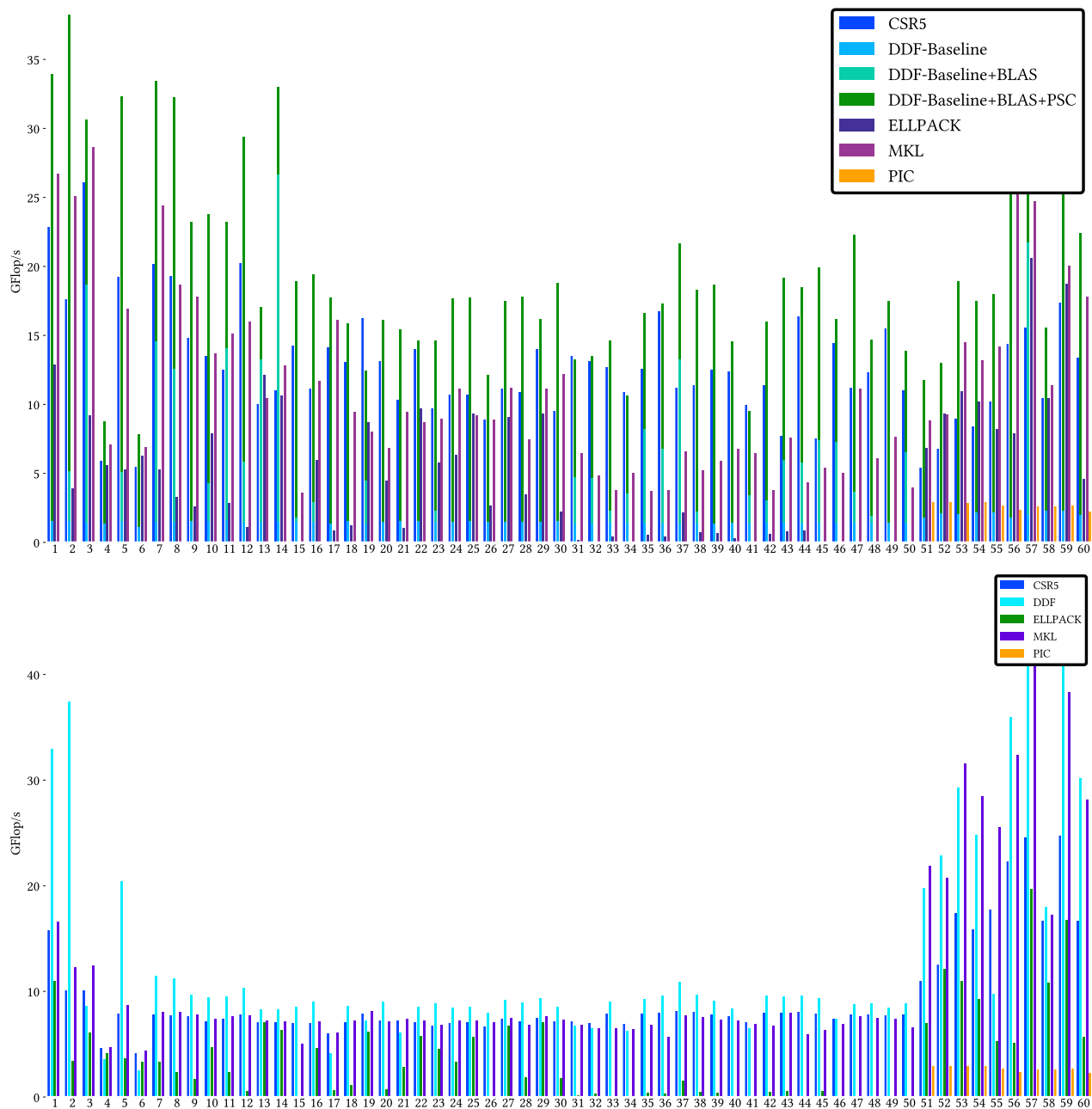


Fig. 8. *Top*: Shows DDF speedups for the SpMV kernel over the single threaded baseline SpMV CSR algorithm ran on an Intel architecture. Speedups also reported for parallel MKL, parallel CSR5, parallel ELLPACK and single threaded PIC [2] over the single threaded baseline SpMV CSR algorithm. *Bottom*: Shows DDF, MKL, CSR5, ELLPACK and PIC speedups for the SpMV kernel using an AMD architecture.

As shown in Figure 11, partially strided codelets are the main contributors to the overall performance of the DDF's SpTRSV code. On average, DDF is $3.9\times$ faster compared to when only BLAS codelets are generated. Similar to SpMV, PSCs contribute to optimizing 78% of the operations over all matrices (Figure 6b). However, the number of PSC II codelets has increased from 18% in SpMV to 53% in SpTRSV. The number of computational regions with more than one strided access function is small, due to the existing dependencies in the SpTRSV kernel, thus, more PSC type II codelets are

generated. Similar to the SpMV kernel, the mined PSCs in DDF improve locality. The correlation coefficient between the speedup and relative memory cycle is 0.67 which is consistent with the trend in SpMV.

While MKL provides an efficient and vectorized implementation for single-threaded SPTRV executions, it's not optimized to execute on parallel processors; the performance of MKL's parallel code is similar to its serial implementation. Sympiler performs well for matrices that contain row-blocks or can be padded with up to 30% nonzeros to create row-

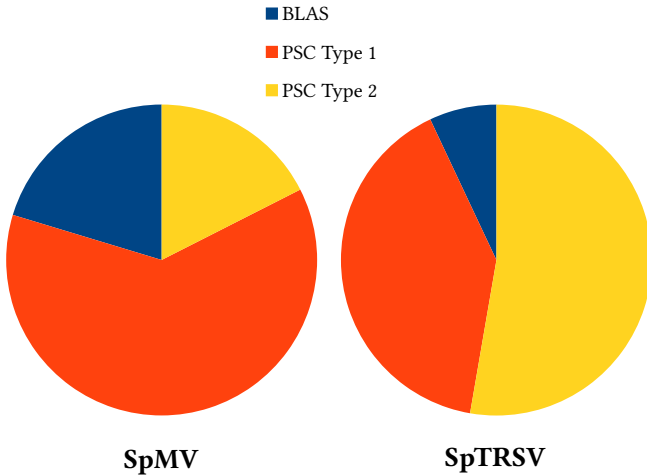


Fig. 9. Shows percentage of points for all matrices processed by BLAS, PSC type 1 or PSC type 2 codelets for SpMV and SpTRSV.

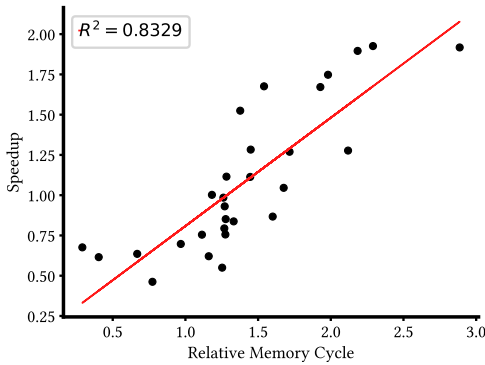


Fig. 10. Shows the correlation between speedup over the baseline while profiling the code with PAPI [32] and the relative memory cycle for matrices used in the test set for SpMV. Relative memory cycle is the average memory cycle of the baseline over the average memory cycle of the DDF framework codes.

blocks. These row-blocks are converted to BLAS calls and thus improve locality. L-factors are generated from a Cholesky factorization of the original matrices in the SuiteSparse collection. Since L-factors are created with a supernodal factorization process [33], the resulting matrices are well-structured and have many row-blocks. As shown, Sympiler’s performance on L-factor matrices (IDs 31-50) is close to that of DDF’s with an average speedup of $1.18\times$, while for all other matrices DDF’s performs on average time $2.1\times$ better than Sympiler.

C. The Inspector Overhead

We compare the inspection time of DDF to that of other inspector-executors/code-generation frameworks, i.e. Sympiler and PIC. We compute the number of executor runs (NER) that amortize the cost of the inspector using

$$\frac{\text{Inspector Time}}{\text{Baseline Time} - \text{Executor Time}}$$
The *baseline* time is obtained by running sequential implementation of the kernel. The PIC

inspector would timeout for over 83.3% of matrices because of its large inspection overhead and code compilation time. For small matrices that PIC would execute, more than one million executor runs are needed to amortize the cost of the inspection. DDF’s inspection time is on average 0.5 seconds with an average NER of 15 for SpMV. The inspection time of DDF and Sympiler are similar with an average NER of less than 100 for both tools for SPTRV. The largest matrices in our benchmark are inspected in less than 7 seconds with DDF. Sparse kernels such as SpMV and SpTRSV are typically used in iterative solvers, for example, to compute a residual in each iteration or to apply a preconditioner per iteration. Even with preconditioning, these solvers typically converge to a solution after tens of thousands of iterations [34]–[36] and hence inspector-executor frameworks such as DDF and Sympiler lead to noticeable speedups as their inspection time overheads are amortized after a few initial iterations of the solver.

VI. RELATED WORK

Numerous hand-optimized libraries [4], [37] and implementations [30], [38]–[40] exist that optimize the performance of sparse matrix computations for different parallel architectures and also optimize vectorization on a single core. Libraries such as MKL and Eigen as well as implementations in [14], [41]–[44] optimize the performance of SpMV on shared memory architectures and improve SIMD vectorizability. A number of library implementations such as [11]–[13], [45]–[48] reorganize data and computation to increase opportunities for vectorization. A class of these libraries implement and optimize sparse kernels based on available storage formats; for example [49] optimizes SpMV based on ELLPACK. Other libraries work best for matrices arising from specific applications such as [14], which optimizes SpMV for large matrices from graph analytics, or KLU [50] which works best for circuit simulation problems.

Domain-specific compilers use domain information to enable the application of code transformations and optimizations such as vectorization. For example domain-specific compilers such as [20], [51], [52] optimize signal processing applications, stencil computations are optimized in [51]–[53], and matrix assembly and mesh analysis is optimized in [54], [55]. While these compilers provide highly-optimized code for the applications that they accelerate, they do not generate specialized code for the input pattern and the computations do not have indirect accesses.

Inspector-Executor approaches inspect the irregular access patterns of sparse matrix computations at run-time to enable the automatic optimization of sparse codes [23], [56]–[63]. The index array accesses of the sparse code is analyzed using an inspector and the information is used at run-time to execute the code efficiently. The sparse polyhedral framework [22], [64], [65] uses uninterpreted function symbols to express regular and irregular segments of sparse codes. As a result it is able to automatically generate inspector executors at compile time that can resolve data dependencies in sparse

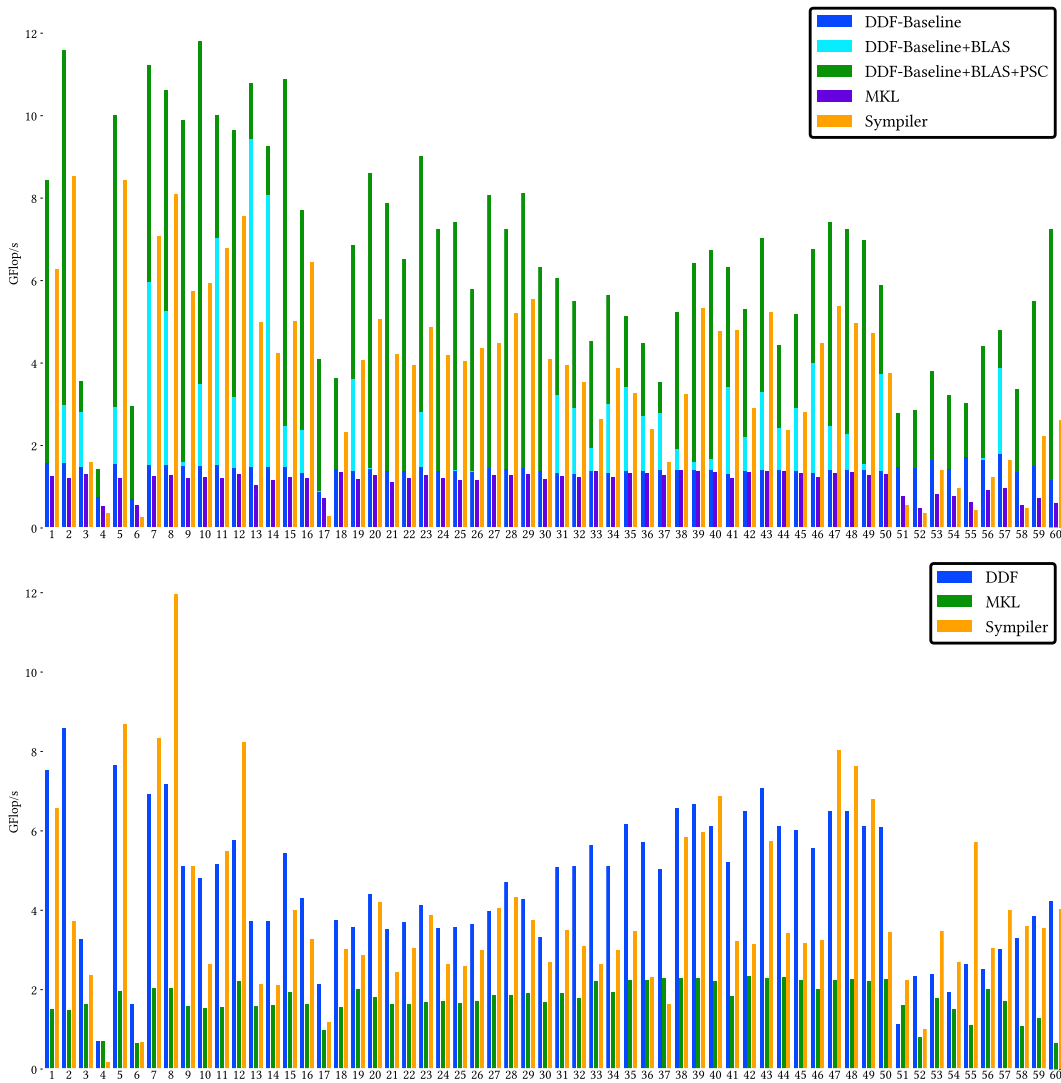


Fig. 11. *Top*: Shows breakdowns for DDF speedups, when enabling BLAS and PSC codelet mining, MKL and Sympiler speedups over the single threaded baseline SpTRSV CSR algorithm on Intel architecture. *Bottom*: Shows DDF, MKL and Sympiler speedups on AMD architecture.

computations. These approaches do not generate code that is specialized for the sparsity of the input matrix. Sympiler [1] and ParSy [28] are amongst the inspector executor frameworks that inspect the matrix sparsity pattern and as a result generate vectorized and parallel code specialized for the input sparsity. Their optimizations for tiling and vactorization are based on detecting row-blocks that primarily exist in matrices obtained from numerical factorizations.

Augustine et al. [2] proposed an approach based on the Trace Reconstruction Engine [66], [67] where polyhedral representations are built by inspecting the sequence of addresses being accessed in the sparse matrix vector multiplication. A followup to this work, proposes to use program guided optimization for better vectorization [68]. These approaches lead to generating code that is specialized for the sparsity pattern of the input matrix and improves SIMD vectorization in SpMV. However, their work can only support small matrices

(below 0.5M nonzeros) because of inspector overheads and because the size of the generated code increases with the matrix size. Sparse computations with loop-carried dependencies such as the sparse triangular solver are also not supported. DDF inspects memory address accesses to find SIMD vectorizable code for sparse triangular solve and SpMV and supports both small and large matrices.

VII. CONCLUSION

In this work, we present partially strided codelets that enable the vectorization of computation regions with unstrided memory accesses in sparse matrix codes. We demonstrate how these codelets increase opportunities for vectorization in sparse codes and also improve data locality in their computation. A novel inspector-executor framework called DDF is proposed. DDF uses an efficient inspector to mine for PSCs with a memory access differentiation approach and as

a result, generates highly efficient code for sparse kernels. The performance of the DDF-generated code is compared to state-of-the-art library implementations and other inspector-executor frameworks for the sparse matrix-vector multiply and the sparse triangular solve kernels. We also demonstrate that the inspection overhead of DDF is negligible.

VIII. ACKNOWLEDGMENTS

This work was supported in part by NSERC Discovery Grants (RGPIN-06516, DGECR00303), the Canada Research Chairs program, and U.S. NSF awards NSF CCF-1814888, NSF CCF-1657175; used the Extreme Science and Engineering Discovery Environment (XSEDE) [Townes et al. 2014] which is supported by NSF grant number ACI-1548562; and was enabled in part by Compute Canada and Scinet¹

REFERENCES

- [1] K. Cheshmi, S. Kamil, M. M. Strout, and M. M. Dehnavi, "Sympiler: Transforming sparse matrix codes by decoupling symbolic analysis," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2017*. New York, NY, USA: Association for Computing Machinery, Inc, nov 2017, pp. 1–13. [Online]. Available: <https://dl.acm.org/doi/10.1145/3126908.3126936>
- [2] T. Augustine, L. N. Pouchet, J. Sarma, and G. Rodríguez, "Generating piecewise-regular code from irregular structures," in *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. New York, NY, USA: Association for Computing Machinery, jun 2019, pp. 625–639. [Online]. Available: <https://dl.acm.org/doi/10.1145/3314221.3314615>
- [3] A. Buluç, J. T. Fineman, M. Frigo, J. R. Gilbert, and C. E. Leiserson, "Parallel sparse matrix-vector and matrix-transpose-vector multiplication using compressed sparse blocks," in *Proceedings of the twenty-first annual symposium on Parallelism in algorithms and architectures*, 2009, pp. 233–244.
- [4] *Intel Math Kernel Library. Reference Manual*. Intel Corporation, 2009, santa Clara, USA. ISBN 630813-054US.
- [5] M. Kreutzer, G. Hager, G. Wellein, H. Fehske, and A. R. Bishop, "A unified sparse matrix data format for efficient general sparse matrix-vector multiply on modern processors with wide SIMD units," *SIAM Journal on Scientific Computing*, vol. 36, no. 5, pp. C401–C423, jul 2013. [Online]. Available: <http://arxiv.org/abs/1307.6209http://dx.doi.org/10.1137/130930352>
- [6] X. Liu, M. Smelyanskiy, E. Chow, and P. Dubey, "Efficient sparse matrix-vector multiplication on x86-based many-core processors," in *Proceedings of the International Conference on Supercomputing*. New York, New York, USA: ACM Press, 2013, pp. 273–282. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=2464996.2465013>
- [7] W. Liu and B. Vinter, "CSR5: An Efficient Storage Format for Cross-Platform Sparse Matrix-Vector Multiplication," *Proceedings of the International Conference on Supercomputing*, vol. 2015-June, pp. 339–350, mar 2015. [Online]. Available: <http://arxiv.org/abs/1503.05032>
- [8] W. T. Tang, R. Zhao, M. Lu, Y. Liang, H. P. Huyng, X. Li, and R. S. M. Goh, "Optimizing and auto-tuning scale-free sparse matrix-vector multiplication on Intel Xeon Phi," in *Proceedings of the 2015 IEEE/ACM International Symposium on Code Generation and Optimization, CGO 2015*. Institute of Electrical and Electronics Engineers Inc., mar 2015, pp. 136–145.
- [9] L. Chen, P. Jiang, and G. Agrawal, "Exploiting recent SIMD architectural advances for irregular applications," in *Proceedings of the 14th International Symposium on Code Generation and Optimization, CGO 2016*. New York, NY, USA: Association for Computing Machinery, Inc, feb 2016, pp. 47–58. [Online]. Available: <https://dl.acm.org/doi/10.1145/2854038.2854046>
- [10] B. Xie, W. Gao, J. Zhan, Z. Jia, L. Zhang, X. Liu, and X. He, "CVR: Efficient Vectorization of SpMV on X86 Processors," in *CGO 2018 - Proceedings of the 2018 International Symposium on Code Generation and Optimization*, vol. 2018-February. New York, NY, USA: Association for Computing Machinery, Inc, feb 2018, pp. 149–162. [Online]. Available: <https://dl.acm.org/doi/10.1145/3168818>
- [11] W. Liu and B. Vinter, "Csr5: An efficient storage format for cross-platform sparse matrix-vector multiplication," in *Proceedings of the 29th ACM on International Conference on Supercomputing*, 2015, pp. 339–350.
- [12] B. Xie, J. Zhan, X. Liu, W. Gao, Z. Jia, X. He, and L. Zhang, "Cvr: Efficient vectorization of spmv on x86 processors," in *Proceedings of the 2018 International Symposium on Code Generation and Optimization*, 2018, pp. 149–162.
- [13] Y. Saad, *Iterative methods for sparse linear systems*. SIAM, 2003.
- [14] S. Yesil, A. Heidarshenas, A. Morrison, and J. Torrellas, "Speeding up spmv for power-law graph analytics by enhancing locality vectorization," in *International Conference for High Performance Computing, Networking, Storage and Analysis, SC*, vol. 2020-November. IEEE Computer Society, nov 2020, pp. 1–15.
- [15] F. Vazquez, G. Ortega, J.-J. Fernández, and E. M. Garzón, "Improving the performance of the sparse matrix vector product with gpus," in *2010 10th IEEE International Conference on Computer and Information Technology*. IEEE, 2010, pp. 1146–1151.
- [16] P. Boulet and P. Feautrier, "Scanning polyhedra without Do-loops," in *Parallel Architectures and Compilation Techniques - Conference Proceedings, PACT*. Institute of Electrical and Electronics Engineers Inc., 1998, pp. 4–11.
- [17] C. Chen, "Polyhedra scanning revisited," in *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. New York, New York, USA: ACM Press, 2012, pp. 499–508. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=2254064.2254123>
- [18] D. S. Kershaw, "The incomplete Cholesky-conjugate gradient method for the iterative solution of systems of linear equations," pp. 43–65, jan 1978.
- [19] M. Püschel, J. M. Moura, J. R. Johnson, D. Padua, M. M. Veloso, B. W. Singer, J. Xiong, F. Franchetti, A. Gačić, Y. Voronenko, K. Chen, R. W. Johnson, and N. Rizzolo, "SPIRAL: Code generation for DSP transforms," in *Proceedings of the IEEE*, vol. 93, no. 2. Institute of Electrical and Electronics Engineers Inc., 2005, pp. 232–273. [Online]. Available: <https://experts.illinois.edu/en/publications/spiral-code-generation-for-dsp-transforms>
- [20] D. G. Spampinato and M. Püschel, "A basic linear algebra compiler," in *Proceedings of the 12th ACM/IEEE International Symposium on Code Generation and Optimization, CGO 2014*. New York, NY, USA: Association for Computing Machinery, feb 2014, pp. 23–32. [Online]. Available: <http://dl.acm.org/doi/10.1145/2544137.2544155>
- [21] A. Tiwari, C. Chen, J. Chame, M. Hall, and J. K. Hollingsworth, "A scalable auto-tuning framework for compiler optimization," in *IPDPS 2009 - Proceedings of the 2009 IEEE International Parallel and Distributed Processing Symposium*, 2009.
- [22] M. M. Strout, G. Georg, and C. Olschanowsky, "Set and relation manipulation for the sparse polyhedral framework," in *International Workshop on Languages and Compilers for Parallel Computing*. Springer, 2012, pp. 61–75.
- [23] A. Venkat, M. Hall, and M. Strout, "Loop and data transformations for sparse matrix code," in *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, vol. 2015-June. New York, NY, USA: Association for Computing Machinery, jun 2015, pp. 521–532. [Online]. Available: <https://dl.acm.org/doi/10.1145/2737924.2738003>
- [24] M. M. Strout, M. Hall, and C. Olschanowsky, "The sparse polyhedral framework: Composing compiler-generated inspector-executor code," *Proceedings of the IEEE*, vol. 106, no. 11, pp. 1921–1934, 2018.
- [25] L. S. Blackford, A. Petitet, R. Pozo, K. Remington, R. C. Whaley, J. Demmel, J. Dongarra, I. Duff, S. Hammarling, G. Henry et al., "An updated set of basic linear algebra subprograms (blas)," *ACM Transactions on Mathematical Software*, vol. 28, no. 2, pp. 135–151, 2002.
- [26] T. A. Davis and Y. Hu, "The university of florida sparse matrix collection," *ACM Transactions on Mathematical Software (TOMS)*, vol. 38, no. 1, pp. 1–25, 2011.

¹www.computecanada.ca

- [27] —, “The University of Florida Sparse Matrix Collection,” *ACM Transactions on Mathematical Software*, vol. 38, no. 1, pp. 1–25, nov 2011. [Online]. Available: <https://dl.acm.org/doi/10.1145/2049662.2049663>
- [28] K. Cheshmi, S. Kamil, M. M. Strout, and M. M. Dehnavi, “ParSy: Inspection and transformation of sparse matrix computations for parallelism,” in *Proceedings - International Conference for High Performance Computing, Networking, Storage, and Analysis, SC 2018*. Institute of Electrical and Electronics Engineers Inc., mar 2019, pp. 779–793.
- [29] E. Wang, Q. Zhang, B. Shen, G. Zhang, X. Lu, Q. Wu, and Y. Wang, *Intel Math Kernel Library*, 05 2014, pp. 167–188.
- [30] T. A. Davis and W. Hager, “Cholmod: supernodal sparse cholesky factorization and update/downdate,” 2005.
- [31] J. Hennessy and D. Patterson, *Computer Architecture - A Quantitative Approach*, 01 2007.
- [32] D. Terpstra, H. Jagode, H. You, and J. Dongarra, “Collecting performance data with papi-c,” in *Tools for High Performance Computing 2009*. Springer, 2010, pp. 157–173.
- [33] Y. Chen, T. A. Davis, W. W. Hager, and S. Rajamanickam, “Algorithm 887: Cholmod, supernodal sparse cholesky factorization and update/downdate,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 35, no. 3, pp. 1–14, 2008.
- [34] M. Benzi, J. K. Cullum, and M. Tuma, “Robust approximate inverse preconditioning for the conjugate gradient method,” *SIAM Journal on Scientific Computing*, vol. 22, no. 4, pp. 1318–1332, 2000.
- [35] D. S. Kershaw, “The incomplete cholesky-conjugate gradient method for the iterative solution of systems of linear equations,” *Journal of computational physics*, vol. 26, no. 1, pp. 43–65, 1978.
- [36] M. Papadrakakis and N. Bitoulas, “Accuracy and effectiveness of preconditioned conjugate gradient algorithms for large and ill-conditioned problems,” *Computer methods in applied mechanics and engineering*, vol. 109, no. 3-4, pp. 219–232, 1993.
- [37] G. Guennebaud, B. Jacob *et al.*, “Eigen v3,” <http://eigen.tuxfamily.org>, 2010.
- [38] X. S. Li and J. W. Demmel, “SuperLU_DIST: A scalable distributed-memory sparse direct solver for unsymmetric linear systems,” *ACM Transactions on Mathematical Software*, vol. 29, no. 2, pp. 110–140, jun 2003. [Online]. Available: <https://dl.acm.org/doi/10.1145/779359.779361>
- [39] S. A. Hutchinson, J. N. Shadid, and R. S. Tuminaro, “Aztec User’s Guide Version 1.1,” Tech. Rep., 1995.
- [40] T. A. Davis, “Algorithm 832: Umfpack v4. 3—an unsymmetric-pattern multifrontal method,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 30, no. 2, pp. 196–199, 2004.
- [41] R. Li and Y. Saad, “Gpu-accelerated preconditioned iterative linear solvers,” *The Journal of Supercomputing*, vol. 63, no. 2, pp. 443–466, 2013.
- [42] S. Williams, L. Oliker, R. Vuduc, J. Shalf, K. Yelick, and J. Demmel, “Optimization of sparse matrix-vector multiplication on emerging multicore platforms,” in *SC’07: Proceedings of the 2007 ACM/IEEE Conference on Supercomputing*. IEEE, 2007, pp. 1–12.
- [43] S. Kamin, M. J. Garzarán, B. Aktumur, D. Xu, B. Yılmaz, and Z. Chen, “Optimization by runtime specialization for sparse matrix-vector multiplication,” in *Proceedings of the 2014 International Conference on Generative Programming: Concepts and Experiences*, 2014, pp. 93–102.
- [44] D. Merrill and M. Garland, “Merge-based parallel sparse matrix-vector multiplication,” in *SC’16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 2016, pp. 678–689.
- [45] R. W. Vuduc, *Automatic performance tuning of sparse matrix kernels*. University of California, Berkeley, 2003.
- [46] R. W. Vuduc and H.-J. Moon, “Fast sparse matrix-vector multiplication by exploiting variable block structure,” in *International Conference on High Performance Computing and Communications*. Springer, 2005, pp. 807–816.
- [47] Y. Li, P. Xie, X. Chen, J. Liu, B. Yang, S. Li, C. Gong, X. Gan, and H. Xu, “Vbsf: a new storage format for simd sparse matrix-vector multiplication on modern processors,” *The Journal of Supercomputing*, vol. 76, no. 3, pp. 2063–2081, 2020.
- [48] W. T. Tang, R. Zhao, M. Lu, Y. Liang, H. P. Huyng, X. Li, and R. S. M. Goh, “Optimizing and auto-tuning scale-free sparse matrix-vector multiplication on intel xeon phi,” in *2015 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. IEEE, 2015, pp. 136–145.
- [49] M. Kreutzer, G. Hager, G. Wellein, H. Fehske, and A. R. Bishop, “A unified sparse matrix data format for efficient general sparse matrix-vector multiplication on modern processors with wide simd units,” *SIAM Journal on Scientific Computing*, vol. 36, no. 5, pp. C401–C423, 2014.
- [50] T. A. Davis and E. Palamadai Natarajan, “Algorithm 907: Klu, a direct sparse solver for circuit simulation problems,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 37, no. 3, pp. 1–17, 2010.
- [51] N. D. Jones, C. K. Gomard, and P. Sestoft, *Partial evaluation and automatic program generation*. Peter Sestoft, 1993.
- [52] F. Quilleré, S. Rajopadhye, and D. Wilde, “Generation of efficient nested loops from polyhedra,” *International journal of parallel programming*, vol. 28, no. 5, pp. 469–498, 2000.
- [53] D. G. Spampinato and M. Püschel, “A basic linear algebra compiler,” in *Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization*, 2014, pp. 23–32.
- [54] M. S. Alnaes, A. Logg, K. B. Ølgaard, M. E. Rognes, and G. N. Wells, “Unified form language: A domain-specific language for weak formulations of partial differential equations,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 40, no. 2, pp. 1–37, 2014.
- [55] F. Luporini, D. A. Ham, and P. H. Kelly, “An algorithm for the optimization of finite element integration loops,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 44, no. 1, pp. 1–26, 2017.
- [56] J. Shin, M. W. Hall, J. Chame, C. Chen, P. F. Fischer, and P. D. Hovland, “Speeding up Nek5000 with autotuning and specialization,” in *Proceedings of the International Conference on Supercomputing*. New York, New York, USA: ACM Press, 2010, pp. 253–262. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1810085.1810120>
- [57] Y. Tang, R. A. Chowdhury, B. C. Kuszmaul, C. K. Luk, and C. E. Leiserson, “The pochoir stencil compiler,” in *Annual ACM Symposium on Parallelism in Algorithms and Architectures*. New York, New York, USA: ACM Press, 2011, pp. 117–128. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1989493.1989508>
- [58] H. L. A. van der Spek and H. A. G. Wijshoff, “Sublimation: Expanding data structures to enable data instance specific optimizations,” in *Languages and Compilers for Parallel Computing*, K. Cooper, J. Mellor-Crummey, and V. Sarkar, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 106–120.
- [59] N. Vasilache, C. Bastoul, and A. Cohen, “Polyhedral code generation in the real world,” in *International Conference on Compiler Construction*. Springer, 2006, pp. 185–201.
- [60] G. Agrawal, J. Saltz, and R. Das, “Interprocedural partial redundancy elimination and its application to distributed memory compilation,” *ACM SIGPLAN Notices*, vol. 30, no. 6, pp. 258–269, 1995.
- [61] R. Das, P. Havlak, J. Saltz, and K. Kennedy, “Index array flattening through program transformation,” in *Proceedings of the 1995 ACM/IEEE conference on Supercomputing*, 1995, pp. 70–es.
- [62] G. Rodríguez and L. Pouchet, “Polyhedral Modeling of Immutable Sparse Matrices,” *impact.gforge.inria.fr*. [Online]. Available: <http://impact.gforge.inria.fr/impact2018/papers/modeling-immutable-sparsemat.pdf>
- [63] J. Saltz, K. Crowley, R. Michandaney, and H. Berryman, “Run-time scheduling and execution of loops on message passing machines,” *Journal of Parallel and Distributed Computing*, vol. 8, no. 4, pp. 303–312, 1990.
- [64] A. LaMille and M. M. Strout, “Enabling code generation within the sparse polyhedral framework,” *Technical report, Technical Report CS-10-102*, 2010.
- [65] A. Venkat, M. S. Mohammadi, J. Park, H. Rong, R. Barik, M. M. Strout, and M. Hall, “Automating wavefront parallelization for sparse matrix computations,” in *SC’16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 2016, pp. 480–491.
- [66] G. Rodríguez, J. M. Andión, M. T. Kandemir, and J. Touriño, “Trace-based affine reconstruction of codes,” in *Proceedings of the 14th International Symposium on Code Generation and Optimization, CGO 2016*. New York, NY, USA: Association for Computing Machinery, Inc, feb 2016, pp. 139–149. [Online]. Available: <https://dl.acm.org/doi/10.1145/2854038.2854056>
- [67] A. Ketterlin and P. Clauss, “Prediction and trace compression of data access addresses through nested loop recognition,” in *Proceedings of the 2008 CGO - Sixth International Symposium on Code Generation and Optimization*. New York, New York, USA: ACM Press, 2008, pp. 94–103. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1356058.1356071>

- [68] M. Selva, F. Gruber, D. Sampaio, C. Guillon, L. N. Pouchet, and F. Rastello, "Building a polyhedral representation from an instrumented execution: Making dynamic analyses of nonaffine programs scalable," *ACM Transactions on Architecture and Code Optimization*, vol. 16, no. 4, pp. 1–26, dec 2019. [Online]. Available: <https://dl.acm.org/doi/10.1145/3363785>