

# Optimizing the Memory Hierarchy by Compositing Automatic Transformations on Computations and Data

Jie Zhao

State Key Laboratory of Mathematical Engineering  
and Advanced Computing, Zhengzhou, China  
jie.zhao@inria.fr

Peng Di

Huawei Technologies Co., Ltd.  
Beijing, China  
dipeng1982@gmail.com

**Abstract**—Optimizing compilers exploit the memory hierarchy using loop tiling and fusion, but these two transformations usually interfere with each other due to the oversight of transformations on data in memories. We present a novel composition of loop tiling and fusion in this paper. Unlike existing tiling-after-fusion algorithms that only transform computation spaces, our approach first applies rectangular/parallelogram tiling to live-out computation spaces for fitting the memory hierarchy, followed by the computation of the memory footprints required by each tile. The upwards exposed data extracted from the memory footprints are used to determine the tile shapes of intermediate computation spaces, allowing the construction of arbitrary tile shapes. Finally, our technique implements a post-tiling fusion strategy for maximizing data locality without losing tilability or parallelism of live-out computation spaces, thereby enabling storage reduction and reuse, and optimizing the memory hierarchy. We demonstrate that our approach can achieve superior performance on both CPU and GPU architectures over the state of the art by experimenting on 11 benchmarks extracted from numerous domains including neural networks, image processing, sparse matrix computation and linear algebra. Also, the results of the ResNet-50 model on an AI accelerator show that our approach can obtain 16% performance improvement.

**Index Terms**—memory hierarchy, data locality, parallelism, polyhedral model, tiling, fusion

## I. INTRODUCTION

The memory hierarchy on modern heterogeneous architectures provides the programmer the illusion of unlimited, fastest memories, but it also complicates the programming issue, which is worsened by the diversity of domain-specific accelerators. Researchers from the compiler community have devised loop fusion [29], [38] and tiling [27], [58] to optimize the memory hierarchy. Loop tiling is a transformation that groups iterations of loop nests into smaller blocks, maximizing reuse along multiple loop dimensions when the block fits in registers or caches. Loop fusion is a technique that intertwines two or more loop nests while maintaining the producer-consumer relations between these loop nests, allowing more values to be allocated in faster memory and thereby enabling storage reduction.

The compositions of these two loop transformations have been proved effective in many application domains [2], [6],

[31], [32], [55], [61]. However, an optimizing compiler usually has to model tradeoffs between tilability/parallelism and locality. Worse yet, the complicated computation patterns like stencils and reductions presented frequently in applications from the domain of neural networks [13], [54] and image processing [41], [48] call for complex tiling techniques involving overlapped memory footprints between tiles [33], [41], [48], leading to the inefficient use of the memory hierarchy.

The polyhedral model is recognized for its powerful ability to composite affine loop transformations and has been integrated into numerous general-purpose optimizers [11], [21], [57] or domain-specific frameworks [6], [16], [41], [54], [59]. Typically, the model assigns a lexicographical execution date [17], [18] to each instance of the operations in programs, and schedules a new execution date by fully considering parallelism and locality. The algorithms used for computing a new lexicographical execution order are referred to as polyhedral schedulers, which are usually integrated with cost-model-based heuristics for implementing loop fusion. Loop tiling is usually implemented by expanding the dimensions of computation spaces produced from the schedulers, which implies a particular order on fusion and tiling.

Yet integrating fusion heuristics with polyhedral schedulers poses some challenges to the model. First, the fusion heuristics may heavily impact the tilability, parallelism and locality of programs. While an aggressive fusion strategy mitigates data movements between hierarchical memories at the expense of losing tilability and/or parallelism, a conservative strategy maximizes tiling possibilities by transferring data through lower-level caches or off-chip communications. Second, constructing complex tile shapes after fusion is not straightforward, since either post-pass rescheduling is required [41], which may introduce over-approximated memory footprints and therefore lead to performance degradation, or non-affine semantics have to be modeled [60], which optimize overlapped tiling without considering the producer-consumer locality. Finally, an aggressive fusion heuristic may also worsen the time complexity of polyhedral compilation.

In this paper, we do not resort to aggressive fusion heuristics before implementing tiling, but optimize the memory hierarchy

by reordering the sequence of loop fusion and tiling. Derived from a simple heuristic, our approach only applies tiling to live-out computation spaces<sup>1</sup>, which tightens the scheduling space of polyhedral schedulers and therefore reduces time complexity. The technique continues by computing the memory footprints required by each tile of live-out computation spaces, using elementary combinations of the operations on affine sets and maps. The upwards exposed data, i.e., those data used within tiled live-out computation spaces but defined in others, are extracted from the memory footprints without resorting to post-pass rescheduling algorithms nor compositing with non-affine semantics. The tile shapes of those computation spaces that produce intermediate values are then determined by compositing such upwards exposed data with access relations, allowing the construction of arbitrary tile shapes without additional efforts. Finally, we introduce a post-tiling fusion algorithm by leveraging a well-defined polyhedral representation [22], without changing the parallelism or tilability of live-out computation spaces.

Our approach models the composition of tiling and fusion in the absence of tradeoffs between parallelism, locality and recomputation that have to be faced by both domain-specific frameworks [13], [48], [54] and polyhedral optimizers [11], [21], [41], [57], [60], maximizing the utilization of the memory hierarchy on modern architectures. The extension to the well-defined polyhedral representation [22] also strengthens the power of polyhedral compilation by facilitating post-tiling fusion. Moreover, the algorithm presented in this work also moderates compilation time without restricting to special cases [52] or relaxing scheduling constraints [2].

We conduct experiments on 11 benchmarks covering application domains like neural networks, image processing, sparse matrix computation and linear algebra, validating the effectiveness of our technique by targeting both general-purpose processors and domain-specific accelerators. The experimental results demonstrate the general applicability of our approach and its portability to different architectures. We also demonstrate that our approach can achieve significant compile-time improvement over existing fusion heuristics.

The paper is organized as follows. Section II introduces the background of this work and our motivation. Section III presents our technique for constructing arbitrary tile shapes, and Section IV describes the post-tiling fusion algorithm. Section V explains the code generation strategy used by our approach, followed by the experimental results described in Section VI and related work discussed in Section VII. Finally, the conclusion is presented in Section VIII.

## II. BACKGROUND AND MOTIVATION

We now present the background knowledge of polyhedral compilation and illustrate our motivation.

<sup>1</sup>A live-out computation space writes to memory locations that will be referenced after the computation of a program.

### A. Tiling and Fusion in the Polyhedral Model

The polyhedral model is a mathematical abstraction for automatic parallelization and locality optimization. It represents a program using iteration domains, access relations, dependences and schedules. Typically, the polyhedral model uses schedules to represent both the original lexicographical order of a program and one that is generated by a scheduling algorithm. A schedule is an affine function over all statement instances, i.e., iteration domains. A scheduling algorithm has to respect the dependences relating statement instances that depend on each other, which are in turn computed on the basis of access relations. An access relation is an affine map between statement instances and memory locations.

Consider the loop nests of a 2-dimensional convolution over an input image  $A$  using kernel  $B$  in Fig. 1(a), with  $C$  representing the output image. Statement  $S_1$  represents the initialization and statement  $S_2$  performs the reduction. Statement  $S_0$  can be viewed as a quantization step and statement  $S_3$  performs a “ReLU” operation. Using the polyhedral model, the initial schedule can be expressed using a multi-dimensional affine schedule as  $[S_0(h, w) \rightarrow (0, h, w); S_1(h, w) \rightarrow (1, h, w, 0); S_2(h, w, kh, kw) \rightarrow (1, h, w, 1, kh, kw); S_3(h, w) \rightarrow (2, h, w)]$ .

The polyhedral model can compute a new, tiling-friendly schedule by integrating with different fusion heuristics. With a conservative fusion heuristic, the new schedule can be expressed as  $[S_0(h, w) \rightarrow (0, h, w); S_1(h, w) \rightarrow (1, h, w, 0, 0, 0); S_2(h, w, kh, kw) \rightarrow (1, h, w, kh, kw, 1); S_3(h, w) \rightarrow (1, h, w, KH - 1, KW - 1, 2)]$ , and we use  $(\{S_0\}, \{S_1, S_2, S_3\})$  to represent the fusion result. One can now apply rectangular tiling using tile sizes  $T_0 \times T_1$  to the first group  $\{S_0\}$  and  $T_2 \times T_3$  to the second group  $\{S_1, S_2, S_3\}$ , with the tiled code shown on the left of Fig. 1(b). We use  $ht, wt$  to represent the tile loops (iterating among tiles) and  $hp, wp$  the point loops (iterating within a tile) after loop tiling. Note that tile sizes have to be fixed integer values, but we use symbolic tile sizes throughout the paper to better explain the generality of our work.

The cost model of such a conservative heuristic is to maximize fusion without sacrificing the parallelism of the fused loops. When targeting CPUs, the compiler can add OpenMP pragmas before each group as shown in Fig. 1(b). While the tiled OpenMP code benefits from the maximal parallelism preserved by this fusion heuristic, tensor  $A$  cannot be allocated on small scratchpads but has to be stored as full buffers. When targeting GPUs, the polyhedral model generates CUDA code by mapping the parallel loops to the two-level hardware parallelism on GPUs. The GPU mapping strategy enabled by the model is shown on the right of Fig. 1(b), with the tile loops  $ht, wt$  mapped to thread blocks (red arrows) and the point loops  $hp, wp$  to threads (blue arrows). However, tensor  $A$  cannot be allocated on the shared memory.

On the contrary, an aggressive heuristic maximizes data locality by fusing all statements through the combination of loop interchange, shifting and skewing, with the generated code shown in Fig. 1(c). While this policy maximizes the producer-consumer locality, it also reduces the number of

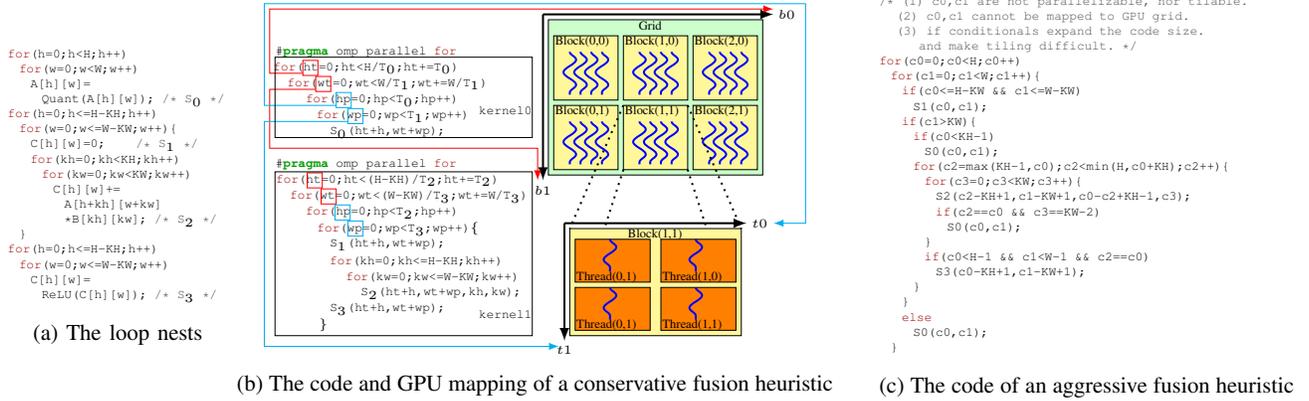


Fig. 1: A 2D convolution with its fusion results

tilable dimensions and loses the outer parallelism. The interchange transformation applied by the polyhedral model also results in another problem: the dimensions of the fused loop nests mismatch with the tile sizes specified by users due to the permutation of dimensions, especially in the case of domain-specific frameworks for neural networks.

We therefore argue that the tiling-after-fusion strategy in existing polyhedral optimizers cannot fully exploit the memory hierarchy and intend to take another way to avoid such tradeoffs between tilability/parallelism and locality by reordering the sequence of tiling and fusion. To implement this reordering, we choose to leverage schedule trees [22].

### B. Schedule Trees

The polyhedral model uses multi-dimensional affine schedules for representing the lexicographical order of a program, but this representation cannot be easily extended to automate memory managements on GPUs, e.g., the automatic insertion of thread-level synchronizations for CUDA code. Affine schedules can be explicitly encoded using a tree structure [22], which can simplify the modeling of automatic memory managements in polyhedral compilers.

Building an initial schedule tree for a program on the basis of a multi-dimensional affine schedule is straightforward. A schedule tree starts with a so-called *domain* node containing all statement instances, i.e., iteration domains, expressed using Presburger formulas [46]. A *sequence* node is introduced to explicitly express the scalar dimensions used in multi-dimensional affine schedules, defining a particular order on its children. Each child of a sequence node has to be a *filter* node, which collects a subset of statement instances introduced by its outer domain or filter node. A *band* node is used to encode the variable and/or constant dimensions of multi-dimensional affine schedules, in the form of a piecewise multi-dimensional affine function over the iteration domain. The polyhedral code generators like [8], [22], [53] take as input the iteration domain and the new schedule produced from the polyhedral model for generating imperative code. The schedule tree representation uses band nodes and sequence nodes for encoding the clas-

sical multi-dimensional affine schedules, but also introduces a domain node to represent the iteration domain. Encoding iteration domains and schedules together makes it possible to generate code by only scanning schedule trees.

Taking Fig. 1(a) as an example, it is easy to obtain its initial schedule tree representation, depicted in Fig. 2(a). The domain node can be expressed using a Presburger set  $\{S_0(h, w) : 0 \leq h < H \wedge 0 \leq w < W; S_1(h, w) : 0 \leq h \leq H - KH \wedge 0 \leq w \leq W - KW; S_2(h, w, kh, kw) : 0 \leq h \leq H - KH \wedge 0 \leq w \leq W - KW \wedge 0 \leq kh < KH \wedge 0 \leq kw < KW; S_3(h, w) : 0 \leq h \leq H - KH \wedge 0 \leq w \leq W - KW\}$ , which is omitted for the sake of conciseness. An initial schedule tree can be automatically transformed into a new one, with the information about the parallelism and tilability of the loop nest attached in band nodes, as shown in Fig. 2(b) (the attached information is not shown).

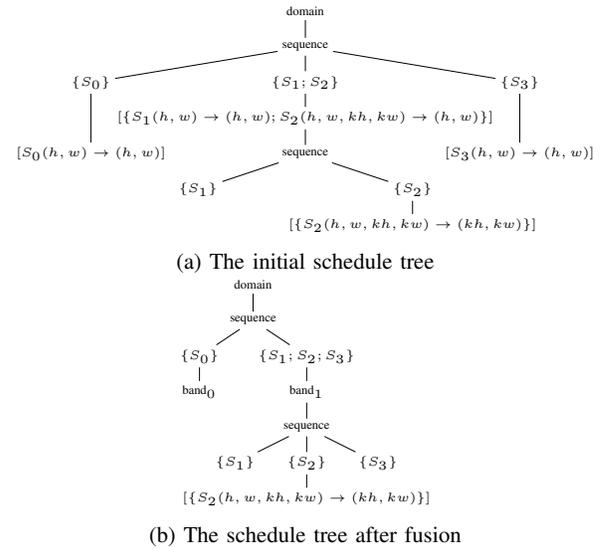


Fig. 2: Schedule trees of the code in Fig. 1(a)

A band node is used to represent a loop nest, and the attached information is used to guide the compiler to apply transformations. The content of the band node  $band_0$  is an

affine function  $\{\{S_0(h, w) \rightarrow (h, w)\}\}$  attached with two attributes, a boolean value *permutable* and a vector *coincident*, where *permutable* is used to indicate whether the loop nest is tilable, and each component of *coincident* is used for expressing the parallelism of a single loop (1 for parallelizable and 0 otherwise). In Fig. 2(b), *permutable* is 1 and *coincident* should be [1, 1] for  $band_0$ .  $band_1$  can be written as  $\{\{S_1(h, w) \rightarrow (h, w); S_2(h, w, kh, kw) \rightarrow (h, w); S_3(h, w) \rightarrow (h, w)\}\}$ , with the attributes represented as 1 and [1, 1]. It means the two loop nests in Fig. 2(b) are both tilable and possess 2D parallelism.

We only present the node types that will be used throughout this work. The readers are invited to refer to the work of Grosser et. al [22] for a detailed description of schedule trees. In particular, we introduce the *extension* node that defines an affine function over its outer schedule dimensions to statement instances or array/scalar elements. An extension node can be used to automate memory managements by retaining itself to a suitable position in schedule trees. We will leverage extension nodes to implement our post-tiling fusion algorithm.

### III. CONSTRUCTING TILE SHAPES

The drawback of a conservative fusion heuristic is that it may result in different tile shapes on the separated computation spaces it produces, which may lead to the mismatch of memory footprints required by these computation spaces.

Still consider the example shown in Fig. 1(a) which is composed of three loop nests. For the sake of clarity, we assume  $H = W = 6$  and  $KH = KW = 3$ . A fusion strategy  $(\{S_0\}, \{S_1, S_2, S_3\})$  obtained by a conservative heuristic results in two groups for the example, as shown in Fig. 1(b). A follow-up rectangular tiling can be applied using a piecewise affine relation:

$$\begin{aligned} & \{\{S_0(h, w) \rightarrow (h/T_0, w/T_1, h, w)\}, \{S_1(h, w) \rightarrow (h/T_2, w/T_3, h, w); \\ & \quad S_2(h, w, kh, kw) \rightarrow (h/T_2, w/T_3, h, w, kh, kw); \} (1) \\ & \quad S_3(h, w) \rightarrow (h/T_2, w/T_3, h, w)\} \end{aligned}$$

which tiles the first group with tile sizes  $T_0 \times T_1$  and the second with  $T_2 \times T_3$ . Each piece of the affine relation (1) represents a tiling schedule, expressed within a pair of braces. We use *tiling schedules* to refer to such piecewise affine relations.

The polyhedral model constructs one computation space for each fusion group. We use a *quantization space* to refer to the first group and a *reduction space* the second. The tiled computation spaces are shown at the bottom of Fig. 3 when given tiles sizes  $T_0 = T_1 = 4$  and  $T_2 = T_3 = 2$ . The quantization space is separated into four blocks, with one full tile (the yellow tile) and three partial tiles. The reduction space is composed of four full tiles. Full tiles are entirely covered by the computation space; partial tiles are on the boundaries of a computation space and interleave with the later [30].

There exists a dependence caused by tensor  $A$  between the two computation spaces, which is written by  $S_0$  and read by  $S_2$ . We show the data space of tensor  $A$  on the top of Fig. 3. Let us focus on the red tiles in both computation spaces. The access relations between each red tile and tensor  $A$  are represented using dotted lines and dashed lines respectively. While the red tile of the quantization space writes to 4 points of

tensor  $A$ , the red tile of the reduction space requires 16 points. The conflict between the memory footprints requires a gather-scatter communication of tensor  $A$  between two computation spaces, which prevents the fusion of the two red tiles.

Such conflict is due to the tiling-after-fusion strategy implemented in existing polyhedral compilation frameworks without considering the transformations on data spaces. If we only apply loop tiling to the reduction space, one can obtain the memory footprints of tensor  $A$  required by each tile of the reduction space, which can then be used to determine the tile shapes of the quantization space in conjunction with the access relation between  $S_0$  and tensor  $A$ . The tiles with the same color from different computation spaces can be fused as the mismatch of memory footprints does not exist. Constructing the tile shape of the quantization space by taking into account the transformations on the data space can also reduce the magnitude of the tile size space, since users only need to specify the tile sizes for the reduction space.

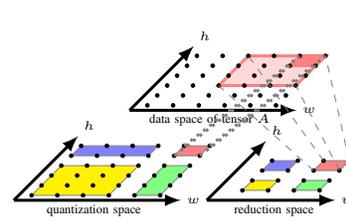


Fig. 3: Tiling computation spaces individually

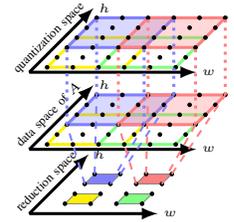


Fig. 4: Constructing tile shapes via upwards exposed data

We thus first use a conservative fusion heuristic to enforce separated computation spaces by setting a proper compilation option provided by *isl* [56], an integer set library used by most polyhedral tools [11], [41], [54], [57] for making decisions on loop fusion. The next step is to apply rectangular/parallelogram tiling only to the live-out computation spaces, which will be used to compute the memory footprints.

#### A. Extracting Upwards Exposed Data

The data accessed within the reduction space can be divided into read and write accesses. We use *upwards exposed data* to refer to those data read by the reduction space but defined in the quantization space. One can easily construct the relation between the computation tiles and the upwards exposed data by assembling the dependence relations and the read access relation of the reduction space. The data space of tensor  $A$  is depicted in the middle of Fig. 4, with each data tile represented using the same color with its corresponding computation tile in the reduction space.

Let us go through the deduction in the example. We only discuss statement  $S_2$  that reads tensor  $A$  in the reduction space for the sake of simplicity. The tiling transformation applied to the reduction space can be expressed using the second piece of

the tiling schedule (1). Meanwhile, one can extract the affine relation over statement  $S_2$  to tile dimensions from (1):

$$\{\{S_2(h, w, kh, kw) \rightarrow (o_0 = h/T_2, o_1 = w/T_3)\}\} \quad (2)$$

where  $o_0$  and  $o_1$  represent the tile dimensions. The access relation over  $S_2$  to the upwards exposed data is written as:

$$\{S_2(h, w, kh, kw) \rightarrow A(h + kh, w + kw) : 0 \leq h \leq H - KH \wedge 0 \leq w \leq W - KW \wedge 0 \leq kh < KH \wedge 0 \leq kw < KW\} \quad (3)$$

The relation between the tile dimensions ( $o_0, o_1$ ) and the upwards exposed data tensor  $A$  can be constructed by intersecting the reverse of (2) with (3):

$$\{(o_0, o_1) \rightarrow A(h', w') : 0 \leq o_0 < \lceil (H - KH + 1)/T_2 \rceil \wedge 0 \leq o_1 < \lceil (W - KW + 1)/T_3 \rceil \wedge T_2 \cdot o_0 \leq h' < T_2 \cdot o_0 + KH + T_2 - 1 \wedge T_3 \cdot o_1 \leq w' < T_3 \cdot o_1 + KW + T_3 - 1\} \quad (4)$$

which represents the dashed lines between the reduction space and the data space in Fig. 4. It allows the overlapped memory footprints between two consecutive computation tiles.

We continue by focusing on the blue tile and the red tile of the reduction space. Without loss of generality, we assume  $T_2 = T_3 = 2$ , and the blue tile can be represented using a coordinate ( $o_0 = 1, o_1 = 0$ ) in the space of tile dimensions while the red tile can be represented using ( $o_0 = 1, o_1 = 1$ ). One can apply (4) to these tiles and obtain their memory footprints, which can be written as  $\{A(h', w') : 2 \leq h' \leq 5 \wedge 0 \leq w' \leq 3\}$  and  $\{A(h', w') : 2 \leq h' \leq 5 \wedge 2 \leq w' \leq 5\}$  respectively. In other words, their intersection is accessed by both tiles, which represents the interleaved region between the blue and red tiles in the data space.

### B. Tiling Intermediate Computation Spaces

The memory footprints obtained by relation (4) can be used to construct the tile shape of the quantization space, which writes intermediate values to tensor  $A$ . With the polyhedral model, one can determine the tile shape of the quantization space using elementary operations on affine maps.

The polyhedral model can provide the write access relation of the quantization space over tensor  $A$ . An affine function mapping tensor  $A$  to the statement  $S_0$  can be computed by reversing this write access relation:

$$\{A(h, w) \rightarrow S_0(h, w) : 0 \leq h < H \wedge 0 \leq w < W\} \quad (5)$$

This relation is represented using the dotted lines between the data space and quantization space in Fig. 4. Intersecting (4) with (5) generates another affine relation:

$$\{(o_0, o_1) \rightarrow S_0(h, w) : 0 \leq o_0 < \lceil (H - KH + 1)/T_2 \rceil \wedge 0 \leq o_1 < \lceil (W - KW + 1)/T_3 \rceil \wedge T_2 \cdot o_0 \leq h < T_2 \cdot o_0 + KH + T_2 - 1 \wedge T_3 \cdot o_1 \leq w < T_3 \cdot o_1 + KW + T_3 - 1\} \quad (6)$$

which represents a function mapping the tile dimensions ( $o_0, o_1$ ) of the reduction space to a set of  $S_0$ 's instances. In other words, the tile dimensions of the reduction space divides the statement instances of  $S_0$  into multiple subsets/tiles, which implements the tiling of the quantization space without using the first piece of the tiling schedule (1). We use *extension schedules* to represent affine relations like (6), since they will be used by an extension node in Section IV.

Intersecting (4) with (5) can be interpreted as the conjunction of the dashed lines and dotted lines in Fig. 4. All of the statement instances within the blue tile of the quantization space can thus be represented as  $\{S_0(h, w) : 2 \leq h \leq 5 \wedge 0 \leq w \leq 3\}$ , and those of the red tile should be  $\{S_0(h, w) : 2 \leq h \leq 5 \wedge 2 \leq w \leq 5\}$ . It means that the tile shape of the quantization space computed using the extension schedule (6) can overlap with each other, without modeling non-affine expressions or refining scheduling algorithms. Note that applying loop tiling using an extension schedule is not possible in existing polyhedral compilation frameworks [6], [11], [21], [54], [57].

### C. The Tiling Algorithm

To summarize our approach to constructing arbitrary tile shapes, we assume that there are only one live-out computation space and multiple intermediate computation spaces after the start-up fusion. Algorithm 1 formally describes the construction of arbitrary tile shapes.

---

#### Algorithm 1: Construct arbitrary tile shapes

---

**Input:** *Spaces*—A group of affine sets for computation spaces  
1 *liveout*  $\leftarrow$  The live-out computation space of *Spaces*;  
2 *Spaces*  $\leftarrow$  *Spaces* - *liveout*; *Untiled*  $\leftarrow$   $\emptyset$ ;  
**if** *liveout* is tilable **then**  
3     *Mixed Schedules*  $\leftarrow$  A tiling schedule like (1) for *liveout*;  
4     *m*  $\leftarrow$  Number of parallelizable loops of *Mixed Schedules*;  
5     *data*  $\leftarrow$  Upwards exposed data of *liveout*;  
6     *f*  $\leftarrow$  A function of *liveout* like (4);  
**foreach** *S* **in** *Spaces* **do**  
7     *n*  $\leftarrow$  Number of parallelizable loops of *S*;  
8     **if** *m* > *n* **then**  
9         | *Untiled*  $\leftarrow$  *Untiled*  $\cup$  *S*; **continue**;  
9     *stmts*  $\leftarrow$  The set of statements in *S*; *h*  $\leftarrow$   $\emptyset$ ;  
10     **while** *stmts*  $\neq$   $\emptyset$  **do**  
11         | *dep*  $\leftarrow$  Dependences caused by *data*;  
11         | *s*  $\leftarrow$  Source statements of *dep*  $\cap$  *stmts*;  
12         | *s*  $\leftarrow$  *s* - the ranges of each map in *h*;  
13         | *S*  $\leftarrow$  The reverse of the write accesses of *s*;  
14         | *h*  $\leftarrow$  *h*  $\cup$  *f*·*S*; *stmts*  $\leftarrow$  *stmts* - *s*;  
15         | *data*  $\leftarrow$  *data*  $\cup$  upwards exposed data of *s*;  
16     *Mixed Schedules*  $\leftarrow$  *Mixed Schedules*  $\cup$  *h*;  
17     **if** *Untiled*  $\neq$   $\emptyset$  **then**  
18         | *Mixed Schedules*  $\leftarrow$  *Mixed Schedules*  $\cup$  Apply Algorithm 1 to *Untiled*;  
**else**  
18     *Mixed Schedules*  $\leftarrow$  *Mixed Schedules*  $\cup$  Apply Algorithm 1 to *Spaces*;  
**Output:** *Mixed Schedules*

---

The algorithm takes a group of affine sets *Spaces* that represent the computation spaces produced by a conservative heuristic as input, and obtains the live-out computation space *liveout* (line 1), which is subtracted from *Spaces* (line 2) and the latter becomes the set of intermediate computation spaces. When *liveout* is tilable, the algorithm first constructs a tiling schedule of *liveout* for simple rectangular/parallelogram tiling (line 3), which is used to initialize the output *Mixed Schedules*, a union of tiling schedules and extension schedules. The Pluto scheduler [11] or its variants can be used to construct rectangular/parallelogram tile shapes. The relation between upwards exposed data to the tile dimensions of *liveout* is computed at lines 5-6.

The construction of tile shapes for each intermediate computation space  $S$  is implemented between lines 7 and 17. The algorithm compares the numbers of parallelizable loops  $n$  of  $S$  with the parallelizable dimensions  $m$  of *liveout*.  $S$  should not be tiled using upwards exposed data but be added to another set *Untiled* that collects all affine sets like  $S$  (line 8) when  $m$  is greater than  $n$ , which means the live-out computation space has more parallelizable dimensions than that of  $S$ . One may obtain an incorrect tiled version of  $S$  without this condition.

Note that each  $S$  will be either considered to be fused with *liveout* (lines 9-16) or added to the *Untiled* set. The union of  $S$  and *liveout* will be considered as the live-out computation space. This guarantees that the visiting order of affine sets in *Spaces* will not impact the correctness of post-tiling fusion.

The comparison between  $m$  and  $n$  is used to guarantee the correctness and effectiveness of Algorithm 1, as well as those of the post-tiling fusion as we will introduce in Section IV-B. In particular, the condition  $m > n$  promises that an intermediate computation space with fewer parallel loops will not be fused with a live-out computation space with more parallel dimensions. Recall that we use two attributes, *permutable* and *coincident*, of a band node to represent its tilability and parallelism. A modern polyhedral scheduler always prefers outer parallelism. It means that the  $n$  parallelizable loops always appear at the outermost  $n$  levels of a multi-dimensional loop nest. In practice, the parallelizable dimensions of a live-out computation space may be greater than  $m$ . For example, a live-out computation space has a 3D parallelizable band. One can force  $m$  to be equal to 1, i.e., only the outermost loop to be parallelizable, when targeting CPUs because OpenMP code only provides 1D parallelism. One can also let  $m$  be equal to 2 when generating CUDA code for GPUs, which can allow more aggressive fusion strategies without losing the two-level hardware parallelism.

Comparing  $m$  and  $n$  preserves the parallelism of the live-out computation space, but it may lose the parallelism of a fused intermediate computation space. In the worst case when  $n > m = 0$ , the parallelism of an intermediate computation is completely lost. In such cases, we only assume a live-out computation space is tilable if  $m$  is greater than 0 when targeting CPUs, or if  $m$  is greater than 1 when targeting GPUs.

Lines 10-15 are used to compute an extension schedule like (6). Lines 12-16 perform the elementary operations like intersection and reverse, as mentioned in Sections III-A and III-B.  $h$  is a union of extension schedules, which is used to model the tiling of an intermediate computation space with multiple statements. Each  $S$  should be tiled with  $h$  and added to *Mixed Schedules*. The algorithm is recursively applied to *Untiled* when the latter is not empty (line 17), or *Spaces* when *liveout* is not tilable (line 18).

Our tiling algorithm can construct overlapped tile shapes without refining the scheduling algorithms like the PolyMage framework [41] does. The algorithm constructs overlapped tile shapes of the intermediate computation spaces through the tiles of upwards exposed data, avoiding the introduction of complicated constraints to the affine relations implemented

in [60]. More importantly, unlike existing tiling techniques for constructing complex tile shapes [9], [19], [33], [41], [60], Algorithm 1 provides the ability to construct arbitrary tile shapes and the general applicability to more application domains due to the consideration of transformations on data in memories. The tile shapes are determined by the access manner of upwards exposed data. For example, one can convert the example shown in Fig. 1(a) into matrix multiplication code by fine-tuning the  $kh, kw$  loops and the corresponding subscripts. The readers will find that our tiling technique can still apply to the code by constructing rectangular tile shapes.

#### IV. POST-TILING FUSION

While used to construct arbitrary tile shapes, Algorithm 1 also implies aggressive fusion strategies. The number of fusion groups suggested by Algorithm 1 is equal to its number of invocation times. *Mixed Schedules* should be the union of the second piece of the tiling schedule (1) and the extension schedule (6) for the example in Fig. 1(a). The tiling schedule maps each statement instance of the reduction space to a lexicographical execution date; the extension schedule defines a mapping over the tile dimensions  $(o_0, o_1)$  to the statement instances of  $S_0$ , which has to be integrated with post-tiling fusion to apply loop tiling to the quantization space.

##### A. Facilitating Fusion using Schedule Trees

The schedule tree shown in Fig. 2(b) is the result of a conservative heuristic, with the fusion strategy  $(\{S_0\}, \{S_1, S_2, S_3\})$  represented using the children of the top sequence node.

According to Algorithm 1, we first apply rectangular tiling using the second piece of the tiling schedule (1) to the reduction space, which is represented as the second child of the top sequence node. The original multi-dimensional affine schedule of the  $band_1$  node in Fig. 2(b) should be replaced with this tiling schedule, which is in turn split into two parts, with one represented as  $\{\{S_1(h, w) \rightarrow (h/T_2, w/T_3); S_2(h, w, kh, kw) \rightarrow (h/T_2, w/T_3); S_3(h, w) \rightarrow (h/T_2, w/T_3)\}\}$  and the other  $\{\{S_1(h, w) \rightarrow (h, w); S_2(h, w, kh, kw) \rightarrow (h, w, kh, kw); S_3(h, w) \rightarrow (h, w)\}\}$ . This splitting operation isolates the tile dimensions and thus makes it possible to implement tile-wise fusion between computation spaces.

We use *tile\_band* and *point\_band* to represent these affine relations after splitting the original  $band_1$  node, as shown in Fig. 5. *tile\_band* represents the dimensions iterating among tiles and *point\_band* represents the dimensions iterating within tiles. We will explain the meaning of each introduced node between *tile\_band* and *point\_band* soon. Note that we have not yet applied tiling to the quantization space; the band node  $band_0$  of the first filter node  $\{S_0(h, w)\}$  is unchanged.

$S_0$  can be viewed as foreign to the subtree rooted under the filter node  $\{S_1(h, w); S_2(h, w, kh, kw); S_3(h, w)\}$ . We mentioned in Section II-B that we would use *extension* nodes to facilitate post-tiling fusion. An extension node is originally designed to add additional statements that are not covered by the domain node of a schedule tree using an affine relation [22]. We

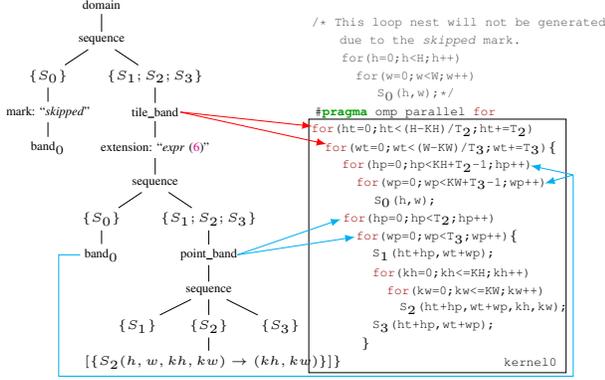


Fig. 5: The schedule tree of post-tiling fusion

extend the expressiveness of an expansion node to introduce additional statements under a filter node in our case.

As shown on the left of Fig. 5, the extension node is inserted underneath the *tile\_band* node, with its affine relation assigned with (6). This simple manipulation on the schedule tree implements overlapped tiling of  $S_0$  and the tile-wise fusion of the original two computation spaces. Note that such an extension to the expressiveness of the schedule tree representation makes it possible to implement post-tiling fusion in the polyhedral model, which is not possible in existing polyhedral compilation frameworks [54], [57] that also use schedule trees.

One can now schedule the statement instances of  $S_0$  under the extension node. As we expect to implement tile-wise fusion, a sequence node has to be introduced underneath the extension node. The first child should be the filter node of the original quantization space, i.e.,  $\{S_0(h, w)\}$ , while the second should be the original reduction space. The subtree rooted at *point\_band* is attached under the new introduced filter node  $\{S_1(h, w); S_2(h, w, kh, kw); S_3(h, w)\}$ . The original subtree *band\_0* of the extended filter node  $\{S_0(h, w)\}$  is also introduced to instruct how the extended filter node is scheduled. Duplicating such subtrees guarantees that the multiple statements encompassed by the extended filter will be scheduled together. Introducing a sequence node under the *tile\_band* also benefits the intra-tile distribution transformation [63], which is used to exploit the spatial locality on small scratchpads.

As we have fused the statement  $S_0$  using an extension node, its original schedule, i.e., the left child of the top sequence node should be ignored. This can be implemented by introducing a *mark* node in schedule trees, as depicted in Fig. 5. A mark node is used to attach information to the schedule tree, providing more flexibilities to the polyhedral model for code generation. We attach a string “skipped” to the mark node, instructing the code generator to bypass the subtree below.

A code generator can generate code as shown on the right of Fig. 5 by scanning this schedule tree. The red and blue arrows represent the relations between band nodes with the loops they represent. Unlike the code shown in Fig. 1(b),

this code fuses all three loop nests into one group, allowing tensor  $A$  to be allocated on small scratchpads. In addition, the post-tiling fusion strategy does not lose the parallelism of the fused dimensions, and one can add an OpenMP pragma before the outermost loop when targeting CPUs. When generating CUDA code for GPUs, the entire loop nest can be executed by launching a single kernel, with  $ht, wt$  mapped to thread blocks and each pair of  $hp, wp$  mapped to threads, and tensor  $A$  is allowed to be declared on the share memory.

## B. The Fusion Algorithm

Algorithm 2 formally describes the post-tiling fusion strategy. It requires two inputs: one is a schedule tree built from a multi-dimensional affine schedule obtained by a polyhedral scheduler, and the other is the output of Algorithm 1.

### Algorithm 2: The post-tiling fusion algorithm

**Input:** 1) *Schedule Tree*—a schedule tree before tiling, and 2) *Mixed Schedules*—Output of Algorithm 1

```

foreach Schedule in Mixed Schedules do
  if Schedule is a tiling schedule then
    1  band ← Replace the original band node using Tiling;
    2  m ← Number of parallelizable loops in Schedule;
    3  tile_band, point_band ← Split band into tile dimensions and
      point dimensions;
    4  Intermediates ← All intermediate computation spaces that are
      to be fused;
    foreach I in Intermediates do
      5  Schedule ← Extract the extension schedule over I like
        (6) from Mixed Schedules;
      6  n ← Number of parallelizable loops in Schedule;
      if m > n then
        7  Replace the extension schedule over I with a tiling
          schedule in Mixed Schedules;
        continue;
      8  Insert an extension node to Schedule Tree;
      9  Insert sequence and filter nodes to Schedule Tree;
      10 Mark the original subtree of I as “skipped”;
  Output: Schedule Tree—a schedule tree after tiling and fusion

```

The number of tiling schedules in *Mixed Schedules* is exactly the number of fusion groups suggested by Algorithm 1. For each group, Algorithm 2 first replaces the original band node using *Schedule* (line 1), and splits it into two parts as described in Section IV-A (line 3). The inner loop (lines 5-10) iterates over the intermediate computation spaces that to be fused with the current live-out computation space. An extension schedule of  $I$  should not be fused when  $m > n$  (line 7), with  $m$  and  $n$  representing the numbers of parallelizable loops of a live-out computation space and  $I$ , respectively. The purpose of comparing  $m$  and  $n$  has been explained in Section III-C. Lines 8-10 perform the manipulations on schedule trees.

Algorithm 2 returns a fusion strategy of  $(\{S_0, S_1, S_2, S_3\})$  for the illustrative example, and the tiled and fused code is shown in Fig. 5. Our post-tiling fusion algorithm does not resort to tedious aggressive fusion heuristics used by existing optimizers [11], [21], [54], [57] to maximize data locality. More importantly, unlike the code shown in Fig. 1(c), the post-tiling fusion algorithm does not lose the parallelism of the program, which guarantees the high performance of the

generated code on various architectures by optimizing the memory hierarchy.

### C. Generalization

So far, we have always assumed that there exists only one live-out computation space. We discuss the case of multiple live-out computation spaces below.

Without loss of generality, we assume that there exist two live-out computation spaces,  $liveout_0$  and  $liveout_1$ . The intermediate computation spaces can be divided into three categories: the first is composed of all intermediate computation spaces of  $liveout_0$ , the second is a collection of those used by  $liveout_1$ , and the third consists of all intermediate computation spaces that used by both. The difficulty is how to handle an intermediate computation space of the third category.

Consider the scene shown in Fig. 6(a) where the values defined by  $op_0$  are used by both  $op_1$  and  $op_2$ . We use  $op'_0$  to represent the subset of  $op_0$  that computes the values used by  $op_1$ , and  $op''_0$  the subset that writes to the values read by  $op_2$ . Existing heuristics [10], [28], [40] do not apply fusion when the third category is present due to the possible redundant computations. We observe that fusion is still possible in such cases. With the post-tiling fusion strategy, one can still apply loop fusion as shown in Fig. 6(b) when  $op'_0$  and  $op''_0$  do not intersect with each other, which will not result in redundant computations. We do not allow fusion when the intersection of  $op'_0$  and  $op''_0$  is not empty which will produce redundancy. We also prevent fusion when  $op_0$  cannot be fused to either of its uses since the generated code cannot benefit from aggressive memory optimizations. In summary, our fusion strategy never introduces redundancy to the code.

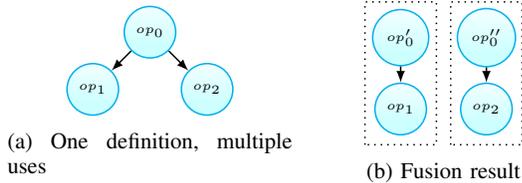


Fig. 6: Fusion strategy for multiple uses

Algorithm 3 describes our approach to compositing tiling and fusion. The algorithm takes a schedule tree constructed from a multi-dimensional affine schedule obtained by the polyhedral model as input, and generates a tiled and fused schedule tree as output. It implements a novel composition of tiling and fusion in three steps.

First, each live-out computation space and its intermediate computation spaces are extracted from the iteration domain of the input schedule tree and saved in  $Spaces$ , to which Algorithm 1 is then applied (line 3). This prevents the fusions between live-out computation spaces, and it indeed makes sense because live-out values do not necessarily need to be allocated on small scratchpads.  $GroupsSet$  is a set of  $Groups$ , i.e., the output of Algorithm 1 for each  $Spaces$ .

---

### Algorithm 3: Reorder the sequence of tiling and fusion

---

**Input:** *Schedule Tree*—a schedule tree before tiling  
1  $Domain \leftarrow$  Domain node of *Schedule Tree*;  $GroupsSet \leftarrow \emptyset$ ;  
**foreach** *liveout* **in** *Domain* **do**  
2  $Spaces \leftarrow$  *liveout* and its intermediate computation spaces;  
3  $GroupsSet \leftarrow GroupsSet \cup$  Apply Algorithm 1 to  $Spaces$ ;  
**foreach** *SharedSpace* **in**  $GroupsSet$  **do**  
4  $Intersect \leftarrow$  The intersection of all uses of *SharedSpace*;  
**if**  $Intersect \neq \emptyset$  **then**  
5 Replace each extension schedule of *SharedSpace* with a tiling schedule;  
**foreach** *Groups* and *SharedSpace* **in**  $GroupsSet$  **do**  
6  $Schedule Tree \leftarrow$  Apply Algorithm 2 to *Groups*;  
**if** *SharedSpace* cannot be fused with its uses **then**  
7 Remove the introduced nodes related to *SharedSpace* in *Schedule Tree*;  
**Output:** *Schedule Tree*—after tiling and fusion

---

The second step of is to handle each intermediate computation space *SharedSpace* that is used by multiple live-out computations. Line 4 is computing the intersection of all extension schedules with respect to *SharedSpace*. If the intersection is not empty, the algorithm replaces all extension nodes of *SharedSpace* using tiling schedules (line 5). This means that *SharedSpace* will not be fused.

The final step is the last loop in the algorithm. It first applies Algorithm 2 to *Groups* (line 7) for constructing schedule trees after tiling and fusion. If *SharedSpace* cannot be fused to any of its uses, the algorithm will remove the introduced nodes related to *SharedSpace* in *Schedule Tree*, which prevents the possible fusion and avoids redundant computations.

Note that Algorithm 3 can also implement dead code elimination in some extreme cases. Suppose that the union of the tiles computed using (6) is a strict subset of the iteration domain of  $S_0$ . These tiles can still be fused with the subtree of the reduction space, while the original subtree representing  $S_0$  are skipped. The post-tiling algorithm in this case eliminates dead stores of  $S_0$  while maintaining the semantic of the program. This fine-grained dead code elimination was not considered by existing polyhedral optimizers [11], [21], [56].

### D. General Applicability

Now we discuss the general applicability of our approach. First, the tile shapes constructed by Algorithm 1 can be rectangular/parallelogram or in an overlapped form. In the latter case, our algorithm enables tile-wise concurrent start [33] by minimizing the recomputations required by overlapped tiling. Algorithm 1 boils down to classical rectangular/parallelogram tiling when no extension schedules are found.

Second, compositing tiling and fusion using extension schedules requires the presence of producer-consumer relations across loop nests. This makes our approach fail to construct complex tile shapes for a single loop nest of stencils that was studied by existing techniques [9], [19], [33]. However, our approach is well suited for applications with multiple loop nests, like neural networks, image processing pipelines, finite element methods and linear algebra. One can unroll the time dimension of a stencil kernel to transform it into multiple loop

nests, to which our work can apply due to the introduced producer-consumer relations across loop nests.

Finally, our approach is very useful for the cases when maximizing fusion does not lose the parallelism of intermediate computation spaces, but it may not be a good choice for programs with multiple consecutive reduction operations, where the parallelism of intermediate reductions cannot be preserved.

## V. CODE GENERATION

We implement our approach using the *isl* library [56] due to its ability to generate AST by scanning schedule trees. This allows us to generate codes for different architectures by first generating AST and then converting the AST to imperative codes using a pretty-print scheme. The PPCG [57] compiler is a polyhedral code generator that wraps *isl* for manipulating integer sets/maps and generating AST; it finally converts the AST generated by *isl* to OpenMP C code or CUDA code. We implement our algorithms in the PPCG compiler to generate OpenMP code for CPUs and CUDA code for GPUs.

A weakness of PPCG’s OpenMP backend is that it cannot exploit automatic vectorization. We identify the innermost parallel loop and add an *ivdep* directive for indicating the absence of loop-carried dependences. The compilers used to compile the generated OpenMP code like Intel *icc* can thus execute the innermost loop with SIMD instructions.

Generating CUDA code for GPUs requires the ability to map parallel loops to thread blocks and threads on GPUs. The CUDA backend of the PPCG compiler models the GPU mapping by leveraging mark nodes of schedule trees. The outermost parallel band node *tile\_band* in Fig. 5 is marked using a “*kernel*” string, which instructs the code generator to map the *ht* and *wt* loops represented by *tile\_band* to GPU thread blocks. A “*thread*” mark is introduced before the *point\_band* node and the *band<sub>0</sub>* node, which tells the code generator to map the *hp* and *wp* loops to GPU threads.

### A. Domain-Specific Code Generation

We also integrate our approach into the *akg* (auto kernel generator) compiler, a wrapper of TVM [13] to generate code for neural networks on domain-specific accelerators. The *akg* compiler is publicly accessible at <https://gitee.com/mindspore/akg>, and it has been integrated into the Huawei MindSpore framework [1] which can provide a DSL to its users for expressing algorithms without taking into consideration the details of underlying architectures. The DSL will then be transformed into Halide IR (intermediate representation), which can be optimized and scheduled either by an expert that is familiar with the target architecture using the schedule primitives provided by the framework, or by the *akg* compiler. For example, one can apply loop tiling using the *tile* primitive or loop fusion with the *fuse* primitive.

The AI platform we target is a dedicated accelerator to boost neural networks—Huawei Ascend 910 chip, of which the DaVinci architecture [36] is depicted in Fig. 7. Cube Unit is a specialized execution unit for performing tensor/matrix

operations by taking as input the data in LOA and LOB, of which the output is stored into LOC. LOA/LOB can fetch data from L1 Buffer. The data in LOC can also be transferred to Vector Unit. Vector/Scalar Units are designed for executing vector/scalar operations. They are allowed to read/write data from/to Unified Buffer, which in turn can exchange data with LOC. L1 Buffer and Unified Buffer are serving as on-chip lower-level caches and used for exchanging data with external memory which is not shown in the figure. Data exchange between L1 Buffer and Unified Buffer is permitted.

The programming model of the accelerator is designed by fully considering the domain-specific properties of applications and the underlying architecture. For example, a convolution operator can be implemented by emitting a single vector instruction using the programming model. The generate CCE code will be compiled with native compilers on the chip, with the same compilation options set for all the versions used in the experiments.

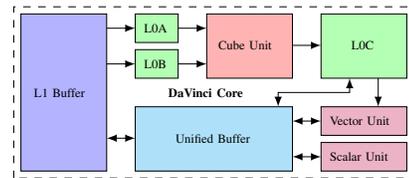


Fig. 7: Overview of the DaVinci architecture

Instead of refining the polyhedral scheduling algorithms to model more tile shapes and/or aggressive fusion strategies, the purpose of this work is to allow more combinations of fusion and tiling that were missed before. The post-tiling fusion is very useful when optimizing neural network applications. The *akg* project indeed implements a more accurate cost model by fully considering the DaVinci architecture, which is enabled when experimenting different versions in the evaluation for a fair comparison. In addition, we also implement a technique to handle parametric tile sizes in the *akg* project, but this property is disabled in the experiment. Explaining these features are out of the scope of this work. Please refer to the open accessible website of *akg* for the detailed implementation.

Also, rather than resorting to the manual scheduling approach, we introduce another pass in the *akg* project by lowering the Halide IR to the schedule tree representation, which will be optimized using our approach. The output schedule tree will be transformed back to the Halide IR for the follow-up code generation. The target imperative code of our AI accelerator, which we refer to as CCE code, will then be generated from the automatically optimized Halide IR.

### B. Aggressive Memory Optimizations

Our approach maximizes data locality without hampering tiling or parallelism, but compilation optimizations may not be very effective without storage reductions for the memory hierarchy due to the streaming nature of applications.

The values produced by an intermediate computation space are only used within a tile, and can thus be discarded after

the computation of a tile [28], [41], [60]. We automatically allocate such values on small scratchpads when generating OpenMP code. The indexing expressions are determined using the range of the affine relations generated by Algorithm 1.

The CUDA backend of PPCG provides us a software-controlled scheme to use the shared/private memory on GPU effectively. PPCG computes an over-approximated rectangular box for complex tile shapes that access non-rectangular blocks of data and therefore enables the allocation of the intermediate values on the shared/private memory. This strategy is also used by existing compilation techniques [19], [60].

We also automate the memory promotion to higher-level caches of the DaVinci architecture using schedule trees. Generally speaking, we resort to mark nodes for managing the data flow between different computation units and extension nodes for the memory optimization and allocation, like the TensorComprehensions framework [54] does. The integration of *isl* into TVM implements the deployment of a neural network on Ascend 910 by only manipulating schedule trees, allowing the framework work with both inference and training of a network.

## VI. EXPERIMENTAL EVALUATION

We select benchmarks from PolyBench [45], PolyMage benchmarks [42], SPEC CPU2000 [26] by considering the following criteria. First, a benchmark is composed of multiple loop nests so that loop fusion can make sense. Second, a polyhedral compiler may find different fusion strategies using an aggressive fusion heuristic; otherwise, our approach may generate the same code as the default fusion heuristic.

The platform for evaluating OpenMP code is a 32-core, dual-socket workstation, of which each CPU is a 2.10GHz 16-core Intel Xeon(R) E5-2683 v4. The OpenMP code is compiled with Intel *icc* compiler 18.0.1 with options *-gopenmp -ipo -O3* enabled. CUDA code is compiled by the NVIDIA CUDA toolkit version 9.1 with *-O3* flag, of which the executable is run on an NVIDIA Quadro P6000 GPU. Each benchmark is executed 11 times with the first run used as a warm-up execution and discarded. We report the average of the remaining 10 execution of each benchmark. By considering 7 possible tiles sizes including 8, 16, 32, 64, 128, 256 and 512 for each dimension, the PolyMage framework uses an auto-tuning strategy for tile size selection. Such auto-tuned tile sizes are also listed in Table I.

### A. Performance on CPU

**Image Processing Pipelines.** An image processing pipeline performs a given task on input images, using a variety of operations like stencils and complex reductions. We use six image processing pipelines extracted from the PolyMage benchmarks which vary widely in structure and complexity. Table I lists the PolyMage benchmarks [42]. We compare the performance with a domain-specific compiler, PolyMage [41], and Halide’s [48] manual schedule written by experts. The PolyMage compiler generates both naïve and optimized OpenMP codes by taking a DSL as input. The sequential code

of a naïve version is used as the baseline and also the input of PPCG that implements our approach, as a naïve version is generated by PolyMage without applying tiling or fusion. On the contrary, an optimized version is generated by fully exploiting fusion and overlapped tiling opportunities.

The parameters like tile sizes, vector lengths and unroll factors of both PolyMage and Halide have been tuned for our platform. We use the same auto-tuned tile sizes and keep the code generation parameters identical with PolyMage for a fair comparison. This isolates the effect of the fusion strategies and tile shapes. The tile sizes and the execution times of different versions are also reported. One can obtain the speedups of over PolyMage and Halide using such numbers, which are shown in Fig. 8. On average, our approach provides 20% and 33% improvements over PolyMage and Halide.

Our approach produces more aggressive fusion strategies for Bilateral Grid [12], [44], Multiscale Interpolation, Local Laplacian Filter [5], [43] and Unsharp Mask than both PolyMage and Halide. Aggressive fusion strategies found by our approach imply more intermediate stages when coupled with overlapped tiling, and more intermediate values can be allocated on small scratchpads, leading to better performance. PolyMage implements the tiling-after-fusion policy; Halide only provides schedule primitives for computation spaces without considering the transformations on data spaces. Neither of them can construct an extension schedule like (6) and thus fails to find the same fusion results as our approach.

PolyMage and our work can also automatically apply an inlining transformation to Harris Corner Detection [23], which was missed by the manual schedule of Halide. This inlining transformation results in 2 remaining stages, and our work generates the same code as PolyMage and thus obtains the same result, outperforming Halide’s manual schedule by 2×.

For Camera Pipeline, our approach generates the same fusion result as PolyMage, which is more aggressive than that of Halide, but our approach can construct tighter overlapped tile shapes which are determined by the memory footprints. Conversely, PolyMage applies overlapped tiling by only transforming computation spaces, leading to over-approximated recomputations and performance degradation.

**Finite Element Method.** *equake* [7] is a benchmark extracted from SPEC CPU2000. It performs a finite element method using a 3D sparse matrix-vector (SpMV) computation. The 3D SpMV computation updates an unstructured mesh using a reduction array, which is followed by a group of affine loop nests performing elementary operations on the mesh. The imperfect loop nest of the 3D SpMV computation consists of three components, with the first one initializing the reduction array, the second performing reduction using a `while` loop, and the third gathering the reduction variables to update the global mesh. The reduction step involves a dynamic condition along the second dimension due to the use of a `while` loop, which is handled by PPCG as a block box.

One can manually permute the `while` loop into the innermost dimension to create fusion opportunities for PPCG that provides three different fusion heuristics. The default

TABLE I: Results of the PolyMage Benchmarks

Benchmark	stages	Tile size	GPU grid	CPU execution time (ms)				GPU execution time (ms)			Compilation time (s)			
				parameter	naïve (1 core)	PolyMage (32 cores)	Halide (32 cores)	Our work (32 cores)	PPCG (minfuse)	Halide	Our work	minfuse	smartfuse	maxfuse
Bilateral Grid	7	8×128	8×64	66.01	5.57	4.23	4.11	5.07	3.79	4.09	0.15	120	>24h	0.86
Camera Pipeline	32	64×256	16×32	116.32	4.68	4.76	4.40	3.51	2.47	2.38	18.3	20.9	>24h	4560
Harris Corner Detection	11	32×256	16×32	246.88	5.10	10.71	5.10	1.79	1.68	1.60	0.03	0.06	0.12	435
Local Laplacian Filter	99	8×256	8×64	480.48	35.35	29.12	27.08	16.73	12.53	11.12	6.94	90.8	>24h	89.3
Multiscale Interpolation	49	32×128	32×16	209.10	16.44	20.07	14.87	15.75	25.65	13.37	0.68	1.40	>24h	3.30
Unsharp Mask	4	8×512	8×32×3	142.16	5.01	5.02	3.68	2.03	1.94	2.01	0.06	0.08	0.10	0.05

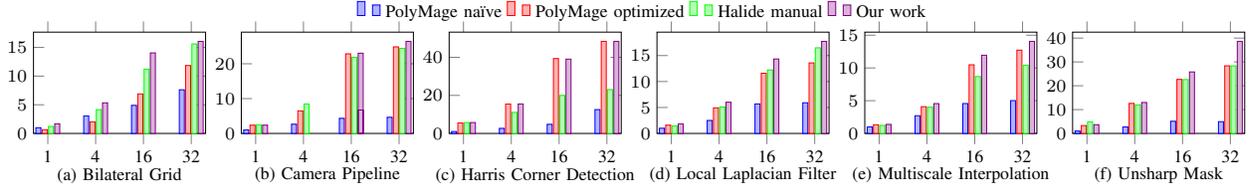


Fig. 8: Performance of PolyMage benchmarks on CPU (x axis: # of threads; y axis: speedup over sequential code)

heuristic that is represented as *smartfuse* tries to maximize fusion without hampering the parallelism or tilability. A more conservative strategy, *minfuse*, does not fuse any loop nests. Our implementation in PPCG is derived from *minfuse*. The most aggressive heuristic maximizes fusion regardless of the parallelism or tilability, represented as *maxfuse*.

*smartfuse* fuses the three components of the 3D SpMV computation together. On the contrary, *maxfuse* fuses the gathering component with the follow-up affine loop nests. The fusion strategy found by our approach is identical with that of *maxfuse*. The speedups over the baseline version of different fusion heuristics are shown in Fig. 9, with the *x* axis representing the problem sizes. As only the outermost loop is tilable, all versions did not apply loop tiling. Algorithm 1 returns an extension schedule with an empty domain, allowing the fusion without loop tiling. This also validates that our post-tiling fusion scheme is applicable without tiling.

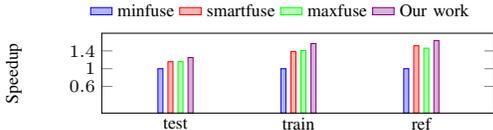


Fig. 9: Performance of quake on CPU (32 cores)

Without the manual permutation of the `while` loop, PPCG cannot exploit loop fusion due to the dynamic condition introduced by the `while` loop. However, this permutation transformation is harmful to data locality, which makes the performance of PPCG’s fusion heuristics fall behind our approach. Our approach does not require such manual permutation as a preprocessing step.

**Linear Algebra and Data Mining.** The benchmarks extracted from PolyBench [45] are summarized in Table II. PolyBench is a collection of micro kernels for linear algebra, stencil computation and physical simulation. 20 out of the 30

PolyBench benchmarks are excluded since they do not require loop fusion due to the structure of perfectly nested loops. Our approach generates the same fusion results as *smartfuse* for 3 of the remaining 10, which are not considered in the evaluation. This validates that our algorithm can fall back to *smartfuse* in the worst case. Due to the page limit, we only choose 3 representative kernels that our approach generates different fusion results from *smartfuse*. The others perform similarly to those shown here.

TABLE II: CPU execution time of the PolyBench benchmarks

threads	2mm			gemver			covariance		
	1	8	32	1	8	32	1	8	32
sequential	4.9	-	-	0.07	-	-	8.1	-	-
icc	2.8	2.5	2.4	0.06	0.07	0.21	8.4	8.5	8.9
minfuse	15.3	2.6	1.1	0.08	0.03	0.03	10.3	2.7	1.1
smartfuse	15.3	2.6	1.1	0.08	0.03	0.03	10.3	2.7	1.1
maxfuse	15.4	2.5	1.0	0.46	8.86	16.88	10.7	3.6	4.6
hybridfuse	9.1	1.6	0.7	0.11	0.03	0.03	×	×	×
Our work	15.3	2.5	1.1	0.08	0.03	0.03	10.5	2.7	1.1

We still compare the performance with different fusion heuristics of PPCG. Besides, we also compare with *hybridfuse*, the hybrid fusion heuristic used by the Pluto optimizer [11] which fuses outer loop dimensions using a conservative heuristic but maximizes the fusion at the inner level of a loop nest. We use same tile sizes (32×32) default enabled by these compilers for each benchmark; tuning the tile sizes does not impact the execution time too much.

*2mm* is a kernel performing 2 matrix-matrix multiplications. We did not observe significant variations in its execution time when using different fusion heuristics of PPCG or our approach, since the parallelism/tilability is preserved by each fusion heuristic. *hybridfuse* achieves the best performance, since maximizing fusion at the innermost level benefits the automatic vectorization of the *icc* compiler. Integrating with a hybrid fusion heuristic may be an interesting direction for our approach to follow.

*gemver* is composed of 4 loop nests performing vector multiplications, additions and matrix-vector multiplications. *covariance* is used to compute the covariance of data samples from different populations in data mining. One can observe that *maxfuse* suffers from significant performance degradation due to the lose of parallelism for these two benchmarks. Our approach enables rectangular/parallelogram tile shapes for these benchmarks. The fusion strategy found by our approach is more aggressive than that of *smartfuse*, but we did not lose parallelism or tilability. *hybridfuse* generates a segmentation fault (represented as  $\times$ ) for *covariance*.

We did not apply aggressive storage optimizations for these micro kernels. This demonstrates that the composition of tiling and fusion exploited by our approach can also improve the performance of programs by maximizing data locality.

### B. Performance on GPU

We now evaluate the performance on GPU. The performance of those benchmarks extracted from the PolyBench benchmark suites follows the same trend with that of the CPU case; we thus do not discuss them here.

**Image Processing Pipelines.** As PolyMage does not target GPUs, we only compare the performance with Halide’s manual schedule. The baseline version is generated by PPCG without our approach, which implements rectangular/parallelogram tiling and the *minfuse* heuristic. The auto-tuned tile sizes are identical with those shown in Table I, with the auto-tuned GPU grid parameters listed in the fifth column.

The results are shown in Fig. 10. The numbers of *smartfuse* and *maxfuse* are missing for some of the benchmarks because they cannot terminate within a reasonable amount of time. We will explain the time complexity issue in Section VI-D.

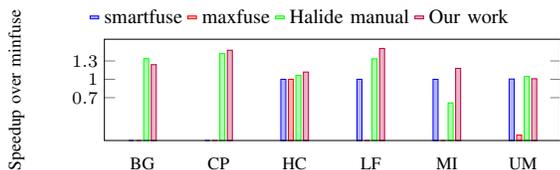


Fig. 10: Performance of PolyMage benchmarks on GPU

*minfuse* does not fuse any of the 4 stages of Unsharp Mask, failing to benefit from the shared/private memory. While *smartfuse* exploits 3D parallelism by fusing the 4 stages into 2 groups, *maxfuse* groups all of the stages together but boils down to 2D parallelism and using  $128 \times 3$  as the GPU grid parameters. *maxfuse* suffers from performance degradation due to the lose of parallelism. None of the fusion heuristics of PPCG applies fusion to Harris Corner Detection which is prevented by overlapped memory footprints. Our approach obtains superior performances because our CUDA code maximizes the utilization of the shared/private memory due to the aggressive fusion results and overlapped tile shapes. Besides, our technique does not lose parallelism.

Halide outperforms our approach slightly for Bilateral Grid and Unsharp Mask since it manually applies unrolling trans-

formations to the channel dimension of the input images after tiling. This can benefit the instruction level pipelined parallelism of the benchmarks, and is an interesting direction for our approach to follow in the future. Our approach provides a mean performance improvement of 17% over Halide.

**Finite Element Method.** PPCG cannot generate effective CUDA code for *equake* due to the presence of the `while` loop. Its enhancement [61] converts the `while` loop into a so-called dynamic counted loop using a preprocessing step, which allows the exploration of loop tiling and fusion in the polyhedral model. The code generation algorithm then introduces a `goto` statement to eliminate the over-approximated iterations caused by the preprocessing. This may achieve  $2.3\times$  speedup over the default setting of PPCG. However, the fusion strategy is exploited by hand in [61]. Our approach achieves the same result as the manual approach but automates the composition of tiling and fusion, which benefits the performance by taking advantage of faster memory on GPU.

### C. Performance on AI Accelerator

We use the ResNet-50 model [25], a 50-layer deep network that is trained on more than one million images from the ImageNet database [15], to conduct experiments on the AI accelerator. The ResNet-50 model used in our experiment is composed of a variety of operators including forward/backward 2D convolutions, batch normalizations, and ReLUs. The *smartfuse* heuristic of *isl* failed to fuse convolutions and batch normalizations. Our approach uses *smartfuse* as the start-up heuristic and enables overlapped tiling and allows the fusion of each forward convolution with its following batch normalization. *minfuse* may prevent the vectorization of the CCE code by separating the initialization and reduction statements of a convolution operator and we therefore do not compare with it. The execution time on our AI accelerator is shown in Table III. The numbers of *maxfuse* are missing due to the tedious compilation time of the heuristic.

TABLE III: Results of the ResNet-50 model

	Execution time (ms)			Compilation time (s)	
	smart	Our work	Speedup	smart	Our work
fwd conv+batchnorm	11.50	6.69	$1.72\times$	-	-
entire workload	35.03	30.25	$1.16\times$	736	487

The tile sizes are specified by experts in the DSL and we did not use the auto-tuner of the framework. The network is trained with the requirement of no less than 76% validation accuracy and the execution time is reported for a single training epoch. We first compare the execution time of all forward convolution and batch normalization operations. This isolates the effect of our approach from other unrelated optimizations of the entire workload. Our technique can obtain 72% performance improvement over *smartfuse*. This is because the off-chip memory latency is very expensive on Ascend 910 chips, and our approach avoids such data transfers. We also show the execution time of the entire workload of Resnet-50, with an improvement of 16% obtained for the entire network.

Optimizing the fusion of backward convolution with other operators is our future work.

#### D. Comparison of Time Complexity

Our approach can also benefit the compilation time of the polyhedral model. We mainly discuss the image processing pipelines and the ResNet-50 model which challenge the scalability of aggressive fusion heuristics. The compilation overheads for the remaining benchmarks are lightweight and we do not discuss them here. We report the compilation time for generating OpenMP code; the compilation overhead for generating CUDA code follows a similar trend with the CPU case. The data of the ResNet-50 workload is collected when generating CCE code on the Ascend 910 chip. The comparison of compilation time is shown in Table I and III.

*maxfuse* cannot finish within one day for most of the image processing pipelines, including Bilateral Grid, Camera Pipeline, Multiscale Interpolation and Local Laplacian Filter. *smartfuse* also suffers from the tedious compilation time for two of them. Our approach can always terminate within 8 minutes (except Camera Pipeline).

The exceptional case is Harris Corner Detection, for which our approach takes longer time than the heuristics. This is because the very complex access pattern presented in the benchmark will lead to an extremely tedious computation of upwards exposed data when using *isl* [56].

Our algorithm exploits an aggressive fusion strategy than *smartfuse* when targeting the Ascend 910 chip, generating fewer computation spaces that need to be scanned by the code generator. This reduces the time of code generation for ResNet-50 and therefore moderates the compilation time of the ResNet-50 model on our AI accelerator.

### VII. RELATED WORK

As we presented in the previous sections, our approach exploits a novel combination of loop tiling and fusion to maximize the utilization of the memory hierarchy by introducing a post-tiling fusion pass. Loop fusion [37], [47], [49] was revisited widely in the past few years for optimizing locality on modern domain-specific chips. The fusion heuristics were well studied for both general-purpose optimizers and domain-specific frameworks, but a well-defined cost model for an optimal solution to different architectures is still not found. The underlying principle is due to the conflict demands of parallelism and locality, as demonstrated in this work. Designing aggressive fusion heuristics cannot avoid this difficulty, and we thus implemented a post-tiling fusion strategy which makes no tradeoffs between parallelism, locality and recomputation.

Tiling [27], [58] was unified into the polyhedral model using affine relations [11], followed by numerous publications on complex tile shapes [9], [14], [19], [20], [33], [50], [62] and parameterized tile sizes [24], [30], [39], [51]. The tile size selection issue was partially addressed by auto-tuning tools [3], [4], [13], [54] that can be used as a complementary optimization for our approach, but complex tile shapes that benefit aggressive storage optimizations like overlapped tiling

were insufficiently integrated within polyhedral frameworks. PolyMage [41] constructs looser overlapped tile shapes by refining its scheduling algorithm and code generator, leading to rather insufficient parallelism of the applications it targets; a schedule-tree-based enhancement of overlapped tiling [60] minimized the recomputation caused by such shapes but missed the exploration of optimized fusion strategies. Our approach does not introduce redundant recomputation when compared with the PolyMage framework [41], but also goes beyond the work of tighter tile shapes [60] by designing a two-level fusion strategy.

Halide [3], [48] was proposed to describe and optimize image processing pipelines in an easier way by isolating schedules from algorithms. Such idea of isolation was also adapted by later deep learning compilers like TVM [13]. Part of such frameworks [13], [41], [48] only provide users with schedule primitives for transforming computations. This makes such frameworks generate potentially similar fusion result as that shown in Fig. 1(b) for the example in Fig. 1(a) due to the conflict in data spaces. As the schedule primitives cannot transform data spaces, these frameworks are not able to compute extension schedules like (6) and therefore fail to fuse the two computation spaces. This problem also arises in the XLA compiler [35]. Our approach overcomes this weakness. The polyhedral model was also integrated into MLIR [34] which provides a mechanism to transform both computation and data spaces. We plan to implement our approach in the MLIR project in the future.

### VIII. CONCLUSION

Compiler optimizations are designed to make effective use of the memory structure on architectures, but some transformations like tiling and fusion usually interfere with each other and therefore fail to maximize the utilization of the memory hierarchy. We proposed a new composition of tiling and fusion in the context of the polyhedral model, enabling aggressive storage optimizations and optimizing the memory hierarchy by constructing arbitrary tile shapes and facilitating aggressive fusion strategies. We have validated the effectiveness of our technique on various architectures by conducting experiments on the benchmarks extracted from numerous application domains. In particular, we considered the portability of our approach to domain-specific accelerators like Huawei Ascend 910 chips. Our approach also provided a reasonable improvement of compilation overhead over existing aggressive heuristics.

### ACKNOWLEDGMENT

The authors are grateful to the MICRO 2020 reviewers for detailed and valuable comments that have improved the paper. This work was partly supported by the National Natural Science Foundation of China under Grant No. 61702546. Both authors are the corresponding authors. The authors would also like to acknowledge the input and discussions from the Huawei MindSpore team.

## REFERENCES

- [1] “Mindspore,” URL: <https://www.mindspore.cn/en>, 2020.
- [2] A. Acharya, U. Bondhugula, and A. Cohen, “Polyhedral auto-transformation with no integer linear programming,” in *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI 2018. New York, NY, USA: ACM, 2018, pp. 529–542. [Online]. Available: <http://doi.acm.org/10.1145/3192366.3192401>
- [3] A. Adams, K. Ma, L. Anderson, R. Baghdadi, T.-M. Li, M. Gharbi, B. Steiner, S. Johnson, K. Fatahalian, F. Durand, and J. Ragan-Kelley, “Learning to optimize halide with tree search and random programs,” *ACM Trans. Graph.*, vol. 38, no. 4, pp. 121:1–121:12, Jul. 2019. [Online]. Available: <http://doi.acm.org/10.1145/3306346.3322967>
- [4] J. Ansel, S. Kamil, K. Veeramachaneni, J. Ragan-Kelley, J. Bosboom, U.-M. O’Reilly, and S. Amarasinghe, “Opentuner: An extensible framework for program autotuning,” in *Proc. of the 23rd Intl. Conf. on Parallel Architectures and Compilation*, ser. PACT ’14. New York, NY, USA: ACM, 2014, pp. 303–316. [Online]. Available: <http://doi.acm.org/10.1145/2628071.2628092>
- [5] M. Aubry, S. Paris, S. W. Hasinoff, J. Kautz, and F. Durand, “Fast local laplacian filters: Theory and applications,” *ACM Trans. Graph.*, vol. 33, no. 5, pp. 167:1–167:14, Sep. 2014. [Online]. Available: <http://doi.acm.org/10.1145/2629645>
- [6] R. Baghdadi, J. Ray, M. B. Romdhane, E. Del Sozzo, A. Akkas, Y. Zhang, P. Suriana, S. Kamil, and S. Amarasinghe, “Tiramisu: A polyhedral compiler for expressing fast and portable code,” in *Proceedings of the 2019 IEEE/ACM International Symposium on Code Generation and Optimization*, ser. CGO 2019. Piscataway, NJ, USA: IEEE Press, 2019, pp. 193–205. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3314872.3314896>
- [7] H. Bao, J. Bielak, O. Ghattas, L. F. Kallivokas, D. R. O’Hallaron, J. R. Shewchuk, and J. Xu, “Large-scale simulation of elastic wave propagation in heterogeneous media on parallel computers,” *Computer Methods in Applied Mechanics and Engineering*, vol. 152, no. 1, pp. 85 – 102, 1998, containing papers presented at the Symposium on Advances in Computational Mechanics. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0045782597001837>
- [8] C. Bastoul, “Code generation in the polyhedral model is easier than you think,” in *Proceedings of the 13th International Conference on Parallel Architectures and Compilation Techniques*, ser. PACT ’04. Washington, DC, USA: IEEE Computer Society, 2004, pp. 7–16. [Online]. Available: <https://doi.org/10.1109/PACT.2004.11>
- [9] U. Bondhugula, V. Bandishti, and I. Pananilath, “Diamond tiling: Tiling techniques to maximize parallelism for stencil computations,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 28, no. 5, pp. 1285–1298, Oct. 2017.
- [10] U. Bondhugula, O. Gunluk, S. Dash, and L. Renganarayanan, “A model for fusion and code motion in an automatic parallelizing compiler,” in *Proceedings of the 19th International Conference on Parallel Architectures and Compilation Techniques*, ser. PACT ’10. New York, NY, USA: ACM, 2010, pp. 343–352. [Online]. Available: <http://doi.acm.org/10.1145/1854273.1854317>
- [11] U. Bondhugula, A. Hartono, J. Ramanujam, and P. Sadayappan, “A practical automatic polyhedral parallelizer and locality optimizer,” in *Proceedings of the 29th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI ’08. New York, NY, USA: ACM, 2008, pp. 101–113. [Online]. Available: <http://doi.acm.org/10.1145/1375581.1375595>
- [12] J. Chen, S. Paris, and F. Durand, “Real-time edge-aware image processing with the bilateral grid,” in *ACM SIGGRAPH 2007 Papers*, ser. SIGGRAPH ’07. New York, NY, USA: ACM, 2007. [Online]. Available: <http://doi.acm.org/10.1145/1275808.1276506>
- [13] T. Chen, T. Moreau, Z. Jiang, L. Zheng, E. Yan, M. Cowan, H. Shen, L. Wang, Y. Hu, L. Ceze, C. Guestrin, and A. Krishnamurthy, “Tvm: An automated end-to-end optimizing compiler for deep learning,” in *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation*, ser. OSDI’18. Berkeley, CA, USA: USENIX Association, 2018, pp. 579–594. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3291168.3291211>
- [14] E. C. Davis, M. M. Strout, and C. Olschanowsky, “Transforming loop chains via macro dataflow graphs,” in *Proc. of the 2018 Intl. Symp. on Code Generation and Optimization*, ser. CGO 2018. New York, NY, USA: ACM, 2018, pp. 265–277. [Online]. Available: <http://doi.acm.org/10.1145/3168832>
- [15] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, June 2009, pp. 248–255.
- [16] V. Elango, N. Rubin, M. Ravishankar, H. Sandanagobalane, and V. Grover, “Diesel: Dsl for linear algebra and neural net computations on gpus,” in *Proceedings of the 2nd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, ser. MAPL 2018. New York, NY, USA: ACM, 2018, pp. 42–51. [Online]. Available: <http://doi.acm.org/10.1145/3211346.3211354>
- [17] P. Feautrier, “Some efficient solutions to the affine scheduling problem. part i. one-dimensional time,” *International journal of parallel programming*, vol. 21, no. 5, pp. 313–347, 1992.
- [18] P. Feautrier, “Some efficient solutions to the affine scheduling problem. part ii. multidimensional time,” *International journal of parallel programming*, vol. 21, no. 6, pp. 389–420, 1992.
- [19] T. Grosser, A. Cohen, J. Holewinski, P. Sadayappan, and S. Verdoolaege, “Hybrid hexagonal/classical tiling for gpus,” in *Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization*, ser. CGO ’14. New York, NY, USA: ACM, 2014, pp. 66:66–66:75. [Online]. Available: <http://doi.acm.org/10.1145/2544137.2544160>
- [20] T. Grosser, A. Cohen, P. H. J. Kelly, J. Ramanujam, P. Sadayappan, and S. Verdoolaege, “Split tiling for gpus: Automatic parallelization using trapezoidal tiles,” in *Proc. of the 6th Workshop on General Purpose Processor Using Graphics Processing Units*, ser. GPGPU-6. New York, NY, USA: ACM, 2013, pp. 24–31. [Online]. Available: <http://doi.acm.org/10.1145/2458523.2458526>
- [21] T. Grosser, A. Groesslinger, and C. Lengauer, “Polly—performing polyhedral optimizations on a low-level intermediate representation,” *Parallel Processing Letters*, vol. 22, no. 04, p. 1250010, 2012.
- [22] T. Grosser, S. Verdoolaege, and A. Cohen, “Polyhedral ast generation is more than scanning polyhedra,” *ACM Trans. Program. Lang. Syst.*, vol. 37, no. 4, pp. 12:1–12:50, Jul. 2015. [Online]. Available: <http://doi.acm.org/10.1145/2743016>
- [23] C. Harris and M. Stephens, “A combined corner and edge detector,” in *Alvey vision conference*, vol. 15, no. 50, 1988, pp. 10–5244.
- [24] A. Hartono, M. M. Baskaran, C. Bastoul, A. Cohen, S. Krishnamoorthy, B. Norris, J. Ramanujam, and P. Sadayappan, “Parametric multi-level tiling of imperfectly nested loops,” in *Proceedings of the 23rd International Conference on Supercomputing*, ser. ICS ’09. New York, NY, USA: ACM, 2009, pp. 147–157. [Online]. Available: <http://doi.acm.org/10.1145/1542275.1542301>
- [25] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 770–778. [Online]. Available: <https://www.computer.org/10.1109/CVPR.2016.90>
- [26] J. L. Henning, “Spec cpu2000: measuring cpu performance in the new millennium,” *Computer*, vol. 33, no. 7, pp. 28–35, July 2000.
- [27] F. Irigoien and R. Triolet, “Supernode partitioning,” in *Proc. of the 15th ACM SIGPLAN-SIGACT Symp. on Principles of Programming Languages*, ser. POPL ’88. New York, NY, USA: ACM, 1988, pp. 319–329. [Online]. Available: <http://doi.acm.org/10.1145/73560.73588>
- [28] A. Jangda and U. Bondhugula, “An effective fusion and tile size model for optimizing image processing pipelines,” in *Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPOPP ’18. New York, NY, USA: ACM, 2018, pp. 261–275. [Online]. Available: <http://doi.acm.org/10.1145/3178487.3178507>
- [29] K. Kennedy and K. S. McKinley, “Maximizing loop parallelism and improving data locality via loop fusion and distribution,” in *Proceedings of the 6th International Workshop on Languages and Compilers for Parallel Computing*. Berlin, Heidelberg: Springer-Verlag, 1993, pp. 301–320.
- [30] D. Kim, L. Renganarayanan, D. Rostron, S. Rajopadhye, and M. M. Strout, “Multi-level tiling: M for the price of one,” in *Proceedings of the 2007 ACM/IEEE Conference on Supercomputing*, ser. SC ’07. New York, NY, USA: ACM, 2007, pp. 51:1–51:12. [Online]. Available: <http://doi.acm.org/10.1145/1362622.1362691>
- [31] M. Kong and L.-N. Pouchet, “Model-driven transformations for multi- and many-core cpus,” in *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI 2019. New York, NY, USA: Association for Computing

- Machinery, 2019, pp. 469–484. [Online]. Available: <https://doi.org/10.1145/3314221.3314653>
- [32] M. Kong, R. Veras, K. Stock, F. Franchetti, L.-N. Pouchet, and P. Sadayappan, “When polyhedral transformations meet simd code generation,” in *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI ’13. New York, NY, USA: Association for Computing Machinery, 2013, pp. 127–138. [Online]. Available: <https://doi.org/10.1145/2491956.2462187>
- [33] S. Krishnamoorthy, M. Baskaran, U. Bondhugula, J. Ramanujam, A. Rountev, and P. Sadayappan, “Effective automatic parallelization of stencil computations,” in *Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI ’07. New York, NY, USA: ACM, 2007, pp. 235–244. [Online]. Available: <http://doi.acm.org/10.1145/1250734.1250761>
- [34] C. Lattner, J. Pienaar, M. Amini, U. Bondhugula, R. Riddle, A. Cohen, T. Shpeisman, A. Davis, N. Vasilache, and O. Zinenko, “Mlir: A compiler infrastructure for the end of moore’s law,” *arXiv preprint arXiv:2002.11054*, 2020.
- [35] C. Leary and T. Wang, “Xla: Tensorflow, compiled,” *TensorFlow Dev Summit*, 2017.
- [36] H. Liao, J. Tu, J. Xia, and X. Zhou, “Davinci: A scalable architecture for neural network computing,” in *2019 IEEE Hot Chips 31 Symposium (HCS)*. IEEE, 2019, pp. 1–44.
- [37] G. Long, J. Yang, K. Zhu, and W. Lin, “Fusionstitching: Deep fusion and code generation for tensorflow computations on gpus,” *arXiv preprint arXiv:1811.05213*, 2018.
- [38] K. S. McKinley, S. Carr, and C.-W. Tseng, “Improving data locality with loop transformations,” *ACM Trans. Program. Lang. Syst.*, vol. 18, no. 4, pp. 424–453, Jul. 1996. [Online]. Available: <https://doi.org/10.1145/233561.233564>
- [39] S. Mehta, G. Beeraka, and P.-C. Yew, “Tile size selection revisited,” *ACM Trans. Archit. Code Optim.*, vol. 10, no. 4, pp. 35:1–35:27, Dec. 2013. [Online]. Available: <http://doi.acm.org/10.1145/2541228.2555292>
- [40] S. Mehta, P.-H. Lin, and P.-C. Yew, “Revisiting loop fusion in the polyhedral framework,” in *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPoPP ’14. New York, NY, USA: ACM, 2014, pp. 233–246. [Online]. Available: <http://doi.acm.org/10.1145/2555243.2555250>
- [41] R. T. Mullapudi, V. Vasista, and U. Bondhugula, “Polymage: Automatic optimization for image processing pipelines,” in *Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS ’15. New York, NY, USA: ACM, 2015, pp. 429–443. [Online]. Available: <http://doi.acm.org/10.1145/2694344.2694364>
- [42] R. T. Mullapudi, V. Vasista, and U. Bondhugula, “Polymage benchmarks,” URL: <https://github.com/bondhugula/polymage-benchmarks> (commit d20264ef), 2017.
- [43] S. Paris, S. W. Hasinoff, and J. Kautz, “Local laplacian filters: Edge-aware image processing with a laplacian pyramid,” *Commun. ACM*, vol. 58, no. 3, pp. 81–91, Feb. 2015. [Online]. Available: <http://doi.acm.org/10.1145/2723694>
- [44] S. Paris, P. Kornprobst, and J. Tumblin, *Bilateral Filtering*. Hanover, MA, USA: Now Publishers Inc., 2009.
- [45] L.-N. Pouchet *et al.*, “Polybench: The polyhedral benchmark suite,” URL: <http://www.cs.ucla.edu/pouchet/software/polybench>, 2012.
- [46] W. Pugh and D. Wonnacott, “Static analysis of upper and lower bounds on dependences and parallelism,” *ACM Trans. Program. Lang. Syst.*, vol. 16, no. 4, pp. 1248–1278, Jul. 1994. [Online]. Available: <http://doi.acm.org/10.1145/183432.183525>
- [47] B. Qiao, O. Reiche, F. Hannig, and J. Teich, “From loop fusion to kernel fusion: A domain-specific approach to locality optimization,” in *Proceedings of the 2019 IEEE/ACM International Symposium on Code Generation and Optimization*, ser. CGO 2019. Piscataway, NJ, USA: IEEE Press, 2019, pp. 242–253. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3314872.3314901>
- [48] J. Ragan-Kelley, C. Barnes, A. Adams, S. Paris, F. Durand, and S. Amarasinghe, “Halide: A language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines,” in *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI ’13. New York, NY, USA: ACM, 2013, pp. 519–530. [Online]. Available: <http://doi.acm.org/10.1145/2491956.2462176>
- [49] S. Rajbhandari, J. Kim, S. Krishnamoorthy, L.-N. Pouchet, F. Rastello, R. J. Harrison, and P. Sadayappan, “On fusing recursive traversals of k-d trees,” in *Proceedings of the 25th International Conference on Compiler Construction*, ser. CC 2016. New York, NY, USA: ACM, 2016, pp. 152–162. [Online]. Available: <http://doi.acm.org/10.1145/2892208.2892228>
- [50] P. S. Rawat, C. Hong, M. Ravishankar, V. Grover, L.-N. Pouchet, A. Rountev, and P. Sadayappan, “Resource conscious reuse-driven tiling for gpus,” in *Proceedings of the 2016 International Conference on Parallel Architectures and Compilation*, ser. PACT ’16. New York, NY, USA: ACM, 2016, pp. 99–111. [Online]. Available: <http://doi.acm.org/10.1145/2967938.2967967>
- [51] L. Renganarayanan, D. Kim, M. M. Strout, and S. Rajopadhye, “Parameterized loop tiling,” *ACM Trans. Program. Lang. Syst.*, vol. 34, no. 1, pp. 3:1–3:41, May 2012. [Online]. Available: <http://doi.acm.org/10.1145/2160910.2160912>
- [52] R. Upadrasta and A. Cohen, “Sub-polyhedral scheduling using (unit-/two-variable-per-inequality polyhedra),” in *Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, ser. POPL ’13. New York, NY, USA: ACM, 2013, pp. 483–496. [Online]. Available: <http://doi.acm.org/10.1145/2429069.2429127>
- [53] N. Vasilache, C. Bastoul, and A. Cohen, “Polyhedral code generation in the real world,” in *Compiler Construction*, ser. CC 2006, A. Mycroft and A. Zeller, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 185–201.
- [54] N. Vasilache, O. Zinenko, T. Theodoridis, P. Goyal, Z. Devito, W. S. Moses, S. Verdoolaege, A. Adams, and A. Cohen, “The next 700 accelerated layers: From mathematical expressions of network computation graphs to accelerated gpu kernels, automatically,” *ACM Trans. Archit. Code Optim.*, vol. 16, no. 4, Oct. 2019. [Online]. Available: <https://doi.org/10.1145/3355606>
- [55] A. Venkat, M. Hall, and M. Strout, “Loop and data transformations for sparse matrix code,” *SIGPLAN Not.*, vol. 50, no. 6, pp. 521–532, Jun. 2015. [Online]. Available: <https://doi.org/10.1145/2813885.2738003>
- [56] S. Verdoolaege, “Isl: An integer set library for the polyhedral model,” in *Proceedings of the Third International Congress Conference on Mathematical Software*, ser. ICMS’10. Berlin, Heidelberg: Springer-Verlag, 2010, pp. 299–302. [Online]. Available: [https://doi.org/10.1007/978-3-642-15582-6\\_49](https://doi.org/10.1007/978-3-642-15582-6_49)
- [57] S. Verdoolaege, J. Carlos Juega, A. Cohen, J. Ignacio Gómez, C. Tenllado, and F. Catthoor, “Polyhedral parallel code generation for cuda,” *ACM Trans. Archit. Code Optim.*, vol. 9, no. 4, pp. 54:1–54:23, Jan. 2013. [Online]. Available: <http://doi.acm.org/10.1145/2400682.2400713>
- [58] J. Xue, *Loop tiling for parallelism*. Springer Science & Business Media, 2012, vol. 575.
- [59] T. Zerrell and J. Bruestle, “Stripe: Tensor compilation via the nested polyhedral model,” *arXiv preprint arXiv:1903.06498*, 2019.
- [60] J. Zhao and A. Cohen, “Flexextended tiles: A flexible extension of overlapped tiles for polyhedral compilation,” *ACM Trans. Archit. Code Optim.*, vol. 16, no. 4, Dec. 2019. [Online]. Available: <https://doi.org/10.1145/3369382>
- [61] J. Zhao, M. Kruse, and A. Cohen, “A polyhedral compilation framework for loops with dynamic data-dependent bounds,” in *Proceedings of the 27th International Conference on Compiler Construction*, ser. CC 2018. New York, NY, USA: Association for Computing Machinery, 2018, pp. 14–24. [Online]. Available: <https://doi.org/10.1145/3178372.3179509>
- [62] X. Zhou, J.-P. Giacalone, M. J. Garzarán, R. H. Kuhn, Y. Ni, and D. Padua, “Hierarchical overlapped tiling,” in *Proc. of the Tenth Intl. Symp. on Code Generation and Optimization*, ser. CGO ’12. New York, NY, USA: ACM, 2012, pp. 207–218. [Online]. Available: <http://doi.acm.org/10.1145/2259016.2259044>
- [63] O. Zinenko, S. Verdoolaege, C. Reddy, J. Shirako, T. Grosser, V. Sarkar, and A. Cohen, “Modeling the conflicting demands of parallelism and temporal/spatial locality in affine scheduling,” in *Proceedings of the 27th International Conference on Compiler Construction*, ser. CC 2018. New York, NY, USA: ACM, 2018, pp. 3–13.