

# 1 (De/Re)-Composition of Data-Parallel Computations 2 via Multi-Dimensional Homomorphisms

## 3 ANONYMOUS AUTHOR(S)

4 Data-parallel computations, such as linear algebra routines (BLAS) and stencil computations, constitute one  
5 of the most relevant classes in parallel computing, e.g., due to their importance for deep learning. Efficiently  
6 de-composing such computations for the memory and core hierarchies of modern architectures and re-  
7 composing the computed intermediate results back to the final result – we say *(de/re)-composition* for short – is  
8 key to achieve high performance for computations on, e.g., GPU and CPU. Current high-level approaches  
9 to generating data-parallel code are often restricted to a particular subclass of data-parallel computations  
10 and architectures (e.g., only linear algebra routines on only GPU, or only stencil computations), and/or  
11 the approaches rely on a user-guided optimization process for a well-performing (de/re)-composition of  
12 computations, which is complex and error prone for the user.

13 We formally introduce a systematic (de/re)-composition approach that is general enough to express a broad  
14 combination of data-parallel computations targeting different kinds of parallel architectures, inspired by  
15 the algebraic formalism for *Multi-Dimensional Homomorphisms (MDH)*. Our approach efficiently targets the  
16 memory and core hierarchies of state-of-the-art architectures, by exploiting our introduced (de/re)-composition  
17 approach for a correct-by-construction, parametrized cache blocking and parallelization strategy. We show  
18 that our approach is powerful enough to express – in the same formalism – the (de/re)-compositions of  
19 different existing state-of-the-art approaches, and we demonstrate that the parameters of our strategies enable  
20 systematically generating code that can be fully automatically optimized (auto-tuned) for the target parallel  
21 architecture and characteristics of the input and output data (e.g., their sizes and memory layouts). Our  
22 experiments confirm that our auto-tuned code achieves higher performance than the state of the art, including  
23 hand-optimized solutions provided by vendors (such as NVIDIA cuBLAS/cuDNN and Intel oneMKL/oneDNN),  
24 on real-world data sets and for a variety of data-parallel computations, including: linear algebra routines,  
25 stencil and quantum chemistry computations, data mining algorithms, and computations that recently gained  
26 high attention due to their relevance for deep learning.

## 28 1 INTRODUCTION

29 Data-parallel computations constitute one of the most relevant classes in parallel computing. Important  
30 examples of such computations include linear algebra routines (BLAS) [Whaley and Dongarra  
31 1998], various kinds of stencil computations (e.g., Jacobi method and convolutions) [Hagedorn  
32 et al. 2018], quantum chemistry computations [Kim et al. 2019], and data mining algorithms [Rasch  
33 et al. 2019b]. The success of many application areas critically depends on achieving high performance  
34 for their data-parallel building blocks on a variety of parallel architectures. For example,  
35 highly-optimized BLAS implementations combined with the computational power of modern  
36 GPUs currently enable deep learning to significantly outperform other existing machine learning  
37 approaches for speech recognition and image classification.

38 Data-parallel computations are characterized by applying the same function (a.k.a *scalar function*)  
39 to each point in a multi-dimensional grid of data (a.k.a. *array*), and combining the obtained results  
40 in the grid’s different dimensions using so-called *combine operators* [Rasch and Gorlatch 2016].

41 Figures 1 and 2 illustrate data parallelism using as examples two popular computations: i) linear  
42 algebra routine *Matrix-Vector multiplication* (MatVec), and ii) stencil computation *Jacobi* (Jacobi1D).  
43 In the case of MatVec, the grid is 2-dimensional and consists of pairs, each pointing to one element  
44 of the input matrix  $M_{i,k}$  and the vector  $v_k$ . To each pair, scalar function  $f(M_{i,k}, v_k) :=$

45 2018. 2475-1421/2018/1-ART1 \$15.00  
46  
47

48 <https://doi.org/>

$$\begin{array}{l}
 50 \\
 51 \quad \left( \begin{array}{ccc} M_{1,1} & \dots & M_{1,K} \\ \vdots & \ddots & \vdots \\ M_{I,1} & \dots & M_{I,K} \end{array} \right), \left( \begin{array}{c} v_1 \\ \vdots \\ v_K \end{array} \right) \xrightarrow{\text{MatVec}} \left( \begin{array}{ccc} f(M_{1,1}, v_1) & \dots & f(M_{1,K}, v_K) \\ \vdots & \ddots & \vdots \\ f(M_{I,1}, v_1) & \dots & f(M_{I,K}, v_K) \end{array} \right) \xrightarrow{\otimes_2} = \left( \begin{array}{c} M_{1,1} * v_1 + \dots + M_{1,K} * v_K \\ \vdots \\ M_{I,1} * v_1 + \dots + M_{I,K} * v_K \end{array} \right) = \left( \begin{array}{c} w_1 \\ \vdots \\ w_I \end{array} \right)
 \end{array}$$

54 Fig. 1. Data parallelism illustrated using the example *Matrix-Vector Multiplication* (MatVec)

55

56

57  $M_{i,k} * v_k$  (multiplication) is applied, and results in the  $i$ -dimension are combined using com-  
 58 bine operator  $\otimes_1((x_1, \dots, x_n), (y_1, \dots, y_m)) := (x_1, \dots, x_n, y_1, \dots, y_m)$  (concatenation) and in  
 59  $k$ -dimension using operator  $\otimes_2((x_1, \dots, x_n), (y_1, \dots, y_m)) := (x_1 + y_1, \dots, x_n + y_n)$  (point-wise  
 60 addition). Similarly, the scalar function of Jacobi1D is  $f(v_{i+0}, v_{i+1}, v_{i+2}) := c * (v_{i+0} + v_{i+1} + v_{i+2})$   
 61 which computes the Jacobi-specific function for an arbitrary but fixed constant  $c$ ; Jacobi1D's  
 62 combine operator  $\otimes_1$  is concatenation. We formally define scalar functions and combine operators  
 63 in Section 2 of this paper.

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

Achieving high performance for data-parallel computations is considered important in both academia and industry, but has proven to be challenging. In particular, achieving high performance that is portable (i.e., the same program code achieves a consistently high level of performance across different architectures and characteristics of input/output data, such as their size and memory layout) and in a user-productive way is identified as an ongoing, major research challenge. This is because for high performance, an efficient (de/re)-composition of computations – meaning their *de-decomposition*, *scalar computation*, and *re-composition* (explained in Figure 3 and discussed thoroughly in this paper) – is required to efficiently target the deep and complex memory and core hierarchies of state-of-the-art architectures, via efficient cache blocking and parallelization strategies [Khronos 2022b; NVIDIA 2022g; OpenMP 2022]. Moreover, to achieve performance that is portable over architectures, the programmer has to consider that architectures often differ significantly in their characteristics [Sun et al. 2019]: depth of memory and core hierarchies, automatically managed caches (as in CPUs) vs manually managed caches (as in GPUs), etc. Productivity is often also hampered: state-of-the-art programming models (such as OpenMP for CPU, CUDA for GPU, and OpenCL for multiple kinds of architectures) currently operate on a low level of abstraction. Thereby, the models require from the programmer explicitly implementing the (de/re)-composition of computations to/from architecture's different memory and core layers, which needs, for example, complex and error-prone index computations, as well as explicitly managing memory and threads on multiple layers.

Current high-level approaches to generating data-parallel code usually struggle with addressing in one combined approach all three challenges: *performance*, *portability*, and *productivity*. For example, approaches such as Apache TVM [Chen et al. 2018a], Halide [Ragan-Kelley et al. 2013], Fireiron [Hagedorn et al. 2020a] and LoopStack [Wasti et al. 2022] achieve high performance, but incorporate the user into the optimization process – by requiring from the user explicitly expressing optimizations in a so-called *scheduling language* – which is error prone and needs expert knowledge about low-level code optimizations, thus hindering user's productivity. In contrast, *polyhedral approaches* such as Facebook's TC [Vasilache et al. 2019], PPCG [Verdoolaege et al. 2013],

$$\begin{array}{l}
 92 \\
 93 \quad \left( \begin{array}{c} v_1 \\ \vdots \\ v_N \end{array} \right) \xrightarrow{\text{Jacobi1D}} \left( \begin{array}{c} f(v_1, v_2, v_3) \\ f(v_2, v_3, v_4) \\ \vdots \end{array} \right) \xrightarrow{\otimes_1} = \left( \begin{array}{c} c * (v_1 + v_2 + v_3) \\ c * (v_2 + v_3 + v_4) \\ \vdots \end{array} \right) = \left( \begin{array}{c} w_1 \\ \vdots \\ w_{N-2} \end{array} \right)
 \end{array}$$

94 Fig. 2. Data parallelism illustrated using the example *Jacobi 1D* (Jacobi1D)

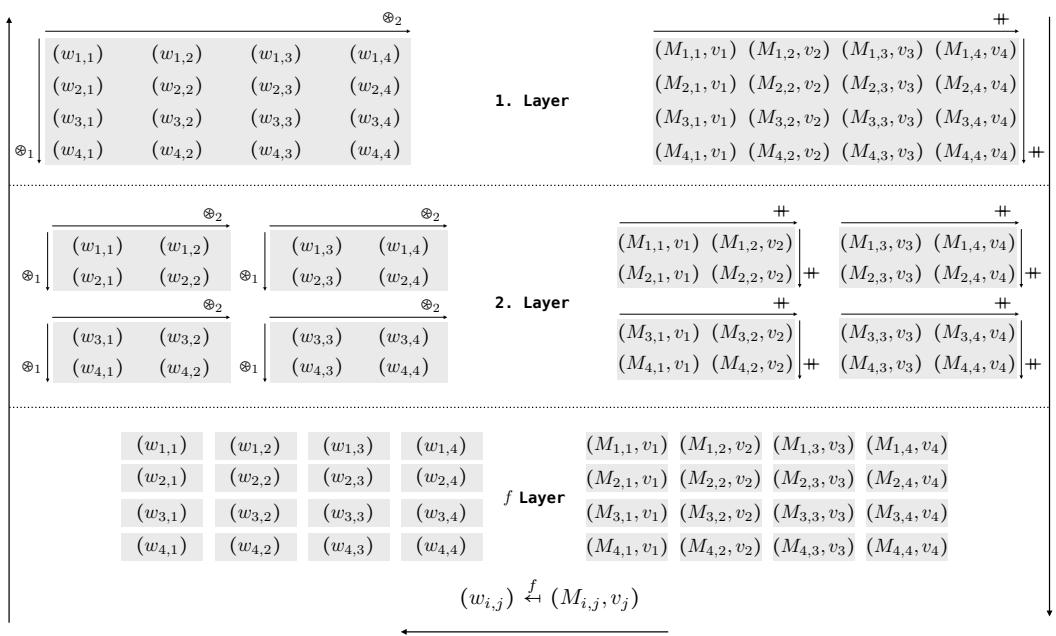
95

96

97

98

99 and Pluto [Bondhugula et al. 2008b] are often fully automatic and thus productive, but usually  
 100 specifically designed toward a particular architecture (e.g., only GPU as TC and PPCG, or only  
 101 CPU as Pluto) and thus not portable. *Functional approaches*, such as Lift [Steuwer et al. 2015],  
 102 are productive for functional programmers (e.g., with experience in Haskell [Haskell.org 2022]  
 103 programming, which is based on small, functional building blocks for expressing computations),  
 104 but the approaches often have difficulties in automatically achieving the full performance potential  
 105 of architectures [Rasch et al. 2019a]. Furthermore, many of the existing approaches are specifically  
 106 designed toward a particular subclass of data-parallel computations only, e.g., only tensor operations  
 107 (as TC and LoopStack) or only matrix multiplication (as Fireiron), or they require significant  
 108 extensions for new subclasses (as Lift for matrix multiplication [Remmelt et al. 2016] and stencil  
 109 computations [Hagedorn et al. 2018]), which further hinders the productivity of the user.



134 Fig. 3. Example (de/re)-composition of MatVec (Figure 1) on a  $4 \times 4$  input matrix  $M$  and a 4-sized vector  $v$ : i) the  
 135 *de-composition phase* (right part of the figure) partitions the concatenated input data into parts (a.k.a. *tiles*  
 136 in programming), where  $+$  denotes the concatenation operator; ii) to each part, scalar function  $f$  is applied  
 137 in the *scalar phase* (bottom part of figure), which is defined for MatVec as: multiplying matrix element  $M_{i,j}$   
 138 with vector element  $v_j$ , resulting in element  $w_{i,j}$ ; iii) the *re-composition phase* (figure's left part) combines the  
 139 computed parts to the final result, using combine operator  $\oplus_1$  for the first dimension (defined as *concatenation*  
 140 in the case of MatVec) and operator  $\oplus_2$  (*addition*) for the second dimension, correspondingly. All basic building  
 141 blocks (*scalar function*, *combine operator*, ...) and concepts (e.g. *partitioning*) are defined later in this paper,  
 142 based on algebraic concepts. For simplicity, this example presents a (de/re)-composition on 2 layers only, and  
 143 we partition the input for this example into parts that have straightforward, equal sizes. Optimized values of  
 144 semantics-preserving parameters (a.k.a. *tuning parameters*), like the number of parts or the application order  
 145 of combine operators, are crucial for achieving high performance, as we discuss in this paper. Phases are  
 146 arranged from right to left, inspired by the application order of function composition, as we also discuss later.

In this paper, we formally introduce a systematic (de/re)-composition approach for data-parallel computations targeting state-of-the-art parallel architectures. We express computations via high-level functional expressions (specifying *what* to compute), in the form of easy-to-use higher-order functions, inspired by the algebraic formalism for *Multi-Dimensional Homomorphisms* (MDHs) [Rasch and Gorlatch 2016]<sup>1</sup>. Our higher-order functions are capable of expressing various kinds of data-parallel computations (linear algebra, stencils, etc), in the same formalism and on a high level of abstraction, independently of hardware and optimization details, thereby contributing to user's productivity<sup>2</sup>. As target for our high-level expressions, we introduce functional low-level expressions (specifying *how* to compute) to formally reason about the de- and re-compositions of data-parallel computations; our low-level expressions are designed such that they can be straightforwardly transformed to executable program code (e.g., in OpenMP, CUDA, and OpenCL). To systematically lower our high-level expressions to low-level expressions, we introduce a formally sound, parameterized de-composition and re-composition approach. The parameters of our lowering process enable automatically computing low-level expressions that are optimized (auto-tuned [Balaprakash et al. 2018]) for the particular target architecture and characteristics of the input/output data, thereby achieving fully automatically high, portable performance. For example, we introduce parameters for flexibly choosing the target memory regions for de-composed and re-composed computations, and also parameters for flexibly setting an optimized data access pattern, based on a formal foundation.

We show that our high-level representation is capable of expressing various kinds of data-parallel computations, including computations that recently gained high attention due to their relevance for deep learning [Barham and Isard 2019]. For our low-level representation, we show that it can express the cache blocking and parallelization strategies of state-of-the-art parallel implementations – as generated by scheduling approach TVM and polyhedral compilers PPCG and Pluto – in one uniform formalism. Moreover, we present experimental results to confirm that via auto-tuning, our (de/re)-composition approach achieves higher performance than the state of the art, including hand-optimized implementations provided by vendors.

Summarized, we make the following three major contributions (illustrated in Figure 4):

- (1) We introduce a high-level functional representation, inspired by the algebraic formalism of Multi-Dimensional Homomorphisms (MDHs), that enables uniformly expressing data-parallel computations on a high level of abstraction.
- (2) We introduce a low-level functional representation that enables formally expressing and reasoning about (de/re)-compositions of data-parallel computations; our low-level representation is designed such that it can be straightforwardly transformed to executable program code in state-of-practice parallel programming models, including OpenMP, CUDA, and OpenCL.
- (3) We introduce a systematic process to fully automatically lower an expression in our high-level representation to a device- and data-optimized expression in our low-level representation, in a formally sound manner, based on auto-tuning.

<sup>1</sup>The existing MDH work introduces a proof-of-concept: 1) *high-level program representation* that relies on a semi-formal foundation, straightforward type system, and a higher-order function for expressing computations (but not for managing input and output data); 2) *code generator* specifically designed and optimized for the OpenCL programming model (thereby also being limited to OpenCL). We thoroughly compare to the existing MDH work in Section 6.6.

<sup>2</sup>We consider as the main users of our approach compiler engineers and library designers. Rasch et al. [2020a] show that our approach can also take straightforward, sequential code as input, which makes our approach attractive also to end users.

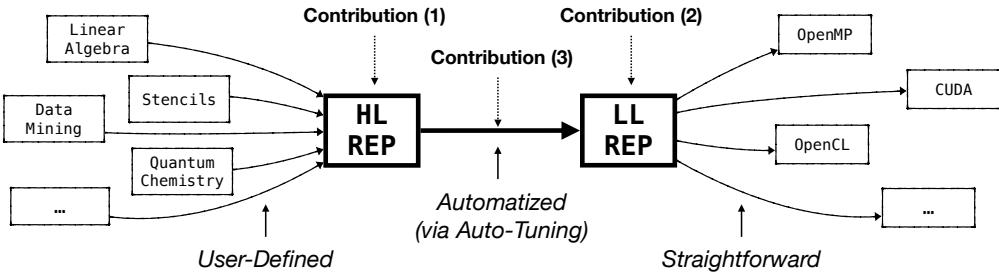


Fig. 4. Overall structure of our approach (contributions highlighted in bold)

The rest of the paper is structured as follows. We introduce our high-level functional representation (Contribution 1) in Section 2, and we show how this representation is used for expressing important data-parallel computations. In Section 3, we discuss our low-level functional representation (Contribution 2), and we confirm that it can express existing state-of-the-art low-level implementations, as generated by state-of-practice approaches TVM (scheduling based) and PPCG/Pluto (polyhedral based). Section 4 shows how we systematically lower a computation expressed in our high-level representation to an expression in our low-level representation, in a formally sound, auto-tunable manner (Contribution 3). We present experimental results in Section 5, discuss related work in Section 6, conclude in Section 7, and we present our ideas for future work in Section 8. We provide an appendix [Rasch 2022] that contains details for the interested reader that should not be required for understanding the basic ideas of this paper. In particular, our appendix contains formal details – for all the following definition, examples, and theorems in Sections 2–4 – which are simplified in the main text to better illustrate the basic concepts introduced in this paper.

## 2 HIGH-LEVEL REPRESENTATION FOR DATA-PARALLEL COMPUTATIONS

We introduce functional building blocks, in the form of higher-order functions, that express data-parallel computations on a high level of abstraction. The goal of our high-level abstraction is to express computations agnostic from hardware and optimization details, and thus in a user-productive manner, while still capturing all information relevant for generating high-performance program code. The building blocks of our abstraction are inspired by the algebraic formalism of *Multi-Dimensional Homomorphisms* (MDHs) which is an approach toward formalizing data parallelism (we compare in detail to the existing work on MDHs in Section 6.6).

Figure 5 shows a basic overview of our high-level representation. We express data-parallel computations using exactly three higher-order functions only (a.k.a. *patterns* or *skeletons* [Gorlatch and Cole 2011] in programming terminology): 1) `inp_view` transforms the domain-specific input data (e.g., a matrix and a vector in the case of matrix-vector multiplication) to a *Multi-Dimensional Array* (MDA) which is our internal data representation and defined later in this section; 2) `md_hom` performs the data-parallel computation; `md_hom` takes as input and output uniformly an MDA; 3) `out_view` transforms the computed MDA back to the domain-specific data representation.

In the following, after informally discussing an introductory example in Section 2.1, we formally define and discuss each higher-order function in detail in Section 2.2 (function `md_hom`) and Section 2.3 (functions `inp_view` and `out_view`). Note that Section 2.2 and Section 2.3 introduce and present the internals and formal details of our approach, which are not relevant for the end user of our system – the user only needs to operate on the abstraction level discussed in Section 2.1.

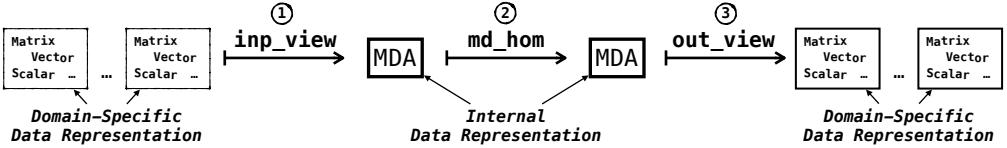


Fig. 5. High-level representation (overview)

## 2.1 Introductory Example

Figure 6 shows how our high-level representation is used for expressing the example of matrix-vector multiplication MatVec (Figure 1)<sup>3</sup>. Computation MatVec takes as input a matrix  $M \in T^{I \times K}$  and vector  $v \in T^K$  of arbitrary scalar type<sup>4</sup>  $T$  and sizes  $I \times K$  (matrix) and  $K$  (vector), for arbitrary but fixed positive natural numbers  $I, K \in \mathbb{N}^5$ . In the figure, based on index function  $(i, k) \rightarrow (i, k)$  and  $(i, k) \rightarrow (k)$ , high-level function `inp_view` computes a function that takes  $M$  and  $v$  as input and maps them to a 2-dimensional array of size  $I \times K$  (referred to as *input MDA* in the following and defined formally in the next subsection). The MDA contains at each point  $(i, k)$  the pair  $(M_{i,k}, v_k) \in T \times T$  comprising element  $M_{i,k}$  within matrix  $M$  (first component) and element  $v_k$  within vector  $v$  (second component). The input MDA is then mapped via function `md_hom` to an output MDA of size  $I \times 1$ , by applying multiplication  $*$  to each pair  $(M_{i,k}, v_k)$  within the input MDA, and combining elements within MDA's first dimension via  $+$  (concatenation – also defined formally in the next subsection) and in second dimension via  $+$  (point-wise addition). Finally, function `out_view` computes a function that straightforwardly maps the output MDA, of size  $I \times 1$ , to MatVec's result vector  $w \in T^I$ , which has scalar type  $T$  and is of size  $I$ . For the example of MatVec, the output view is trivial, but it can be used in other computations (such as matrix multiplication) to conveniently express more advanced variants of computations (e.g., computing the result matrix of matrix multiplication as transposed, as we demonstrate later).

```
273 MatVec< $T \in \text{TYPE} \mid I, K \in \mathbb{N}$ > := out_view< $T$ >( $w: (i, k) \mapsto (i)$ )  $\circ$ 
274                                     md_hom< $I, K$ >( $*, (+, +)$ )  $\circ$ 
275                                     inp_view< $T, T$ >( $M: (i, k) \mapsto (i, k)$ ,  $v: (i, k) \mapsto (k)$ )
```

Fig. 6. High-level expression for Matrix-Vector Multiplication (MatVec)<sup>6</sup>

## 2.2 Function `md_hom`

Higher-order function `md_hom` is introduced by [Rasch and Gorlatch \[2016\]](#) to express Multi-Dimensional Homomorphisms (MDHs) – a formal representation of data-parallel computations – in a convenient and structured way. In the following, we recapitulate the definition of MDHs and function `md_hom`, but in a more general and formally more precise setting than done in the original MDH work.

<sup>3</sup>The expression in Figure 6 can also be extracted from straightforward, annotated sequential code [[Rasch et al. 2020a,b](#)].

<sup>4</sup>We consider as *scalar types* integers  $\mathbb{Z}$  (a.k.a. `int` in programming), floating point numbers  $\mathbb{Q}$  (a.k.a. `float` or `double`), any fixed collection of types (a.k.a. *record* or *struct*), etc. We denote the set of scalar types as  $\text{TYPE}$  in the following.

<sup>5</sup>We denote by  $\mathbb{N}$  the set of positive natural number  $\{1, 2, \dots\}$ , and we use  $\mathbb{N}_0$  for the set of natural numbers including 0.

<sup>6</sup>Our technical implementation takes as input a representation that is equivalent to Figure 6, expressed via straightforward program code (see Appendix, Section A.4)

$$\begin{aligned}
 \mathbf{295} \quad \mathbf{296} \quad \mathbf{297} \quad \mathbf{298} \quad \mathbf{299} \quad \mathbf{300} \quad \mathbf{301} \quad \mathbf{302} \quad \mathbf{303} \quad \mathbf{304} \\
 \mathbf{305} \quad \mathbf{306} \quad \mathbf{307} \quad \mathbf{308} \quad \mathbf{309} \quad \mathbf{310} \quad \mathbf{311} \\
 \mathbf{312} \quad \mathbf{313} \quad \mathbf{314} \quad \mathbf{315} \quad \mathbf{316} \quad \mathbf{317} \quad \mathbf{318} \quad \mathbf{319} \quad \mathbf{320} \quad \mathbf{321} \\
 \mathbf{322} \quad \mathbf{323} \quad \mathbf{324} \quad \mathbf{325} \quad \mathbf{326} \quad \mathbf{327} \quad \mathbf{328} \quad \mathbf{329} \quad \mathbf{330} \quad \mathbf{331} \\
 \mathbf{332} \quad \mathbf{333} \quad \mathbf{334} \quad \mathbf{335} \quad \mathbf{336} \quad \mathbf{337} \quad \mathbf{338} \quad \mathbf{339} \quad \mathbf{340} \quad \mathbf{341} \\
 \mathbf{342} \quad \mathbf{343}
 \end{aligned}
 \begin{aligned}
 \mathbf{a} = \begin{bmatrix} \underbrace{1}_{\mathbf{a}[0,0]} & \underbrace{2}_{\mathbf{a}[0,1]} & \underbrace{3}_{\mathbf{a}[0,2]} & \underbrace{4}_{\mathbf{a}[0,3]} \\ \underbrace{5}_{\mathbf{a}[1,0]} & \underbrace{6}_{\mathbf{a}[1,1]} & \underbrace{7}_{\mathbf{a}[1,2]} & \underbrace{8}_{\mathbf{a}[1,3]} \end{bmatrix}_{\epsilon T[\mathbf{I}_1 := \{0,1\}, \mathbf{I}_2 := \{0,1,2,3\}]} \quad \mathbf{a}^{(1,1)} = \begin{bmatrix} \underbrace{1}_{\mathbf{a}[0,0]} & \underbrace{2}_{\mathbf{a}[0,1]} & \underbrace{3}_{\mathbf{a}[0,2]} & \underbrace{4}_{\mathbf{a}[0,3]} \\ \mathbf{I}_1^{(1,1)} := \{0\}, \mathbf{I}_2^{(1,1)} := \{0,1,2,3\} \end{bmatrix} \quad \mathbf{a}^{(1,2)} = \begin{bmatrix} \underbrace{5}_{\mathbf{a}[1,0]} & \underbrace{6}_{\mathbf{a}[1,1]} & \underbrace{7}_{\mathbf{a}[1,2]} & \underbrace{8}_{\mathbf{a}[1,3]} \\ \mathbf{I}_1^{(1,2)} := \{1\}, \mathbf{I}_2^{(1,2)} := \{0,1,2,3\} \end{bmatrix} \\
 \mathbf{a}^{(2,1)} = \begin{bmatrix} \underbrace{1}_{\mathbf{a}[0,0]} & \underbrace{2}_{\mathbf{a}[0,1]} \\ \mathbf{I}_1^{(2,1)} := \{0,1\}, \mathbf{I}_2^{(2,1)} := \{0,1\} \end{bmatrix} \quad \mathbf{a}^{(2,2)} = \begin{bmatrix} \underbrace{3}_{\mathbf{a}[0,2]} \\ \mathbf{I}_1^{(2,2)} := \{0,1\}, \mathbf{I}_2^{(2,2)} := \{2\} \end{bmatrix} \quad \mathbf{a}^{(2,3)} = \begin{bmatrix} \underbrace{4}_{\mathbf{a}[0,3]} \\ \mathbf{I}_1^{(2,3)} := \{0,1\}, \mathbf{I}_2^{(2,3)} := \{3\} \end{bmatrix}
 \end{aligned}$$

Fig. 7. MDA examples

In order to define Multi-Dimensional Homomorphisms (MDHs), we first need to introduce two central building blocks used in the definition of MDHs: i) *Multi-Dimensional Arrays (MDAs)* – the data type on which MDHs uniformly operate, and ii) *Combine Operators* which we use to combine elements within particular dimensions of an MDA.

## Multi-Dimensional Arrays

**Definition 1** (Multi-Dimensional Array). A *Multi-Dimensional Array (MDA)*  $\mathbf{a}$  that has *dimensionality*  $D \in \mathbb{N}$ , *size*  $N \in \mathbb{N}^D$ , *index sets*  $I_1, \dots, I_D \subset \mathbb{N}_0$ , and *scalar type*  $T \in \text{TYPE}$  is a function with the following signature:

$$\mathbf{a} : I_1 \times \dots \times I_D \rightarrow T$$

We refer to  $I_1 \times \dots \times I_D \rightarrow T$  as the *type* of MDA  $\mathbf{a}$ .

**Notation 1.** For better readability, we denote MDAs' types and accesses to them using a notation close to programming. We often write:

- $\mathbf{a} \in T[\mathbf{I}_1, \dots, \mathbf{I}_D]$  instead of  $\mathbf{a} : I_1 \times \dots \times I_D \rightarrow T$  to denote the type of MDA  $\mathbf{a}$ ;
- $\mathbf{a} \in T[\mathbf{N}_1, \dots, \mathbf{N}_D]$  instead of  $\mathbf{a} : [0, N_1)_{\mathbb{N}_0} \times \dots \times [0, N_D)_{\mathbb{N}_0} \rightarrow T$ ;
- $\mathbf{a}[i_1, \dots, i_D]$  instead of  $a(i_1, \dots, i_D)$  to access MDA  $\mathbf{a}$  at position  $(i_1, \dots, i_D)$ .

Figure 7 illustrates six example MDAs. For example, the left part of the figure shows MDA  $\mathbf{a}$  which is of type  $\mathbf{a} : I_1 \times I_2 \rightarrow T$ , for  $I_1 = \{0,1\}$ ,  $I_2 = \{0,1,2,3\}$ , and  $T = \mathbb{Z}$  (integer numbers). Note that MDAs named  $\mathbf{a}^{(1,1)}$ ,  $\mathbf{a}^{(1,2)}$ ,  $\mathbf{a}^{(2,1)}$ ,  $\mathbf{a}^{(2,2)}$ ,  $\mathbf{a}^{(2,3)}$  in Figure 7 can be considered as *parts* (a.k.a. *tiles* in programming) of MDA  $\mathbf{a}$ : the MDA named  $\mathbf{a}^{(1,1)}$  represents the first row of  $\mathbf{a}$ , MDA  $\mathbf{a}^{(2,2)}$  the third column of  $\mathbf{a}$ , etc. We formally define and use *partitionings* of MDAs in Section 3.

## Combine Operators

A central building block in our definition of MDHs is a *combine operator*. Intuitively, we use a combine operator to combine all elements within a particular dimension of an MDA. For example, in Figure 1 (matrix-vector multiplication), we combine elements of the 2-dimensional MDA via combine operator *concatenation* in MDA's first dimension and via operator *point-wise addition* in the second dimension. Technically, combine operators are functions that take as input two MDAs and yield a single MDA as their output.

We now define *combine operators* formally, and we illustrate this formal definition afterwards using the example operators *concatenation* and *point-wise combination*.

<sup>7</sup>We denote by  $[L, U)_{\mathbb{N}_0} := \{n \in \mathbb{N}_0 \mid L \leq n < U\}$  the half-open interval of natural numbers (including 0) between  $L$  (incl.) and  $U$  (excl.).

344 **Definition 2** (Combine Operator). We refer to any binary function  $\otimes$  of type

$$345 \quad \otimes : T[I_1, \dots, \underset{d}{P}, \dots, I_D] \times T[I_1, \dots, \underset{d}{Q}, \dots, I_D] \rightarrow T[I_1, \dots, \underset{d}{R}, \dots, I_D]$$

348 as *combine operator* that has *scalar type*  $T \in \text{TYPE}$ , *dimensionality*  $D \in \mathbb{N}$ , and *operating dimension*  
349  $d \in [1, D]_{\mathbb{N}}$ . We denote combine operator's type concisely as  $\text{CO}$ .

350 **Example 1** (Concatenation). We define *concatenation* (in dimension  $d$ ) as function  $+_d$  of type

$$352 \quad +_d : T[I_1, \dots, \underset{d}{P}, \dots, I_D] \times T[I_1, \dots, \underset{d}{Q}, \dots, I_D] \rightarrow T[I_1, \dots, P \cup Q, \dots, I_D]$$

354 and that is computed as:

$$356 \quad +_d(a_1, a_2)[i_1, \dots, \underset{d}{i_d}, \dots, i_D] := \begin{cases} a_1[i_1, \dots, i_d, \dots, i_D] & , i_d \in P \\ a_2[i_1, \dots, i_d, \dots, i_D] & , i_d \in Q \end{cases}$$

359 The function is well defined when  $P$  and  $Q$  are disjoint. We usually use an infix notation for  $+_d$ , i.e.,  
360 we write  $a_1 + a_2$  instead of  $+(a_1, a_2)$ , and we refrain from  $+_d$ 's subscript  $d$  when it is clear from  
361 the context.

362 **Example 2** (Point-Wise Combination). We define *point-wise combination* (in dimension  $d$ ), accord-  
363 ing to a binary function  $\oplus : T \times T \rightarrow T$  (e.g. addition), as function  $\vec{\bullet}_d$  of type

$$365 \quad \vec{\bullet}_d : \underbrace{T \times T \rightarrow T}_{\oplus} \rightarrow T[I_1, \dots, \underset{d}{\{0\}}, \dots, I_D] \times T[I_1, \dots, \underset{d}{\{0\}}, \dots, I_D] \rightarrow T[I_1, \dots, \underset{d}{\{0\}}, \dots, I_D]$$

366  $\underbrace{\qquad\qquad\qquad}_{\text{point-wise combination (according to } \oplus\text{)}}$

369 that is computed as:

$$370 \quad \vec{\bullet}_d(\oplus)(a_1, a_2)[i_1, \dots, \underset{d}{0}, \dots, i_D] := a_1[i_1, \dots, \underset{d}{0}, \dots, i_D] \oplus a_2[i_1, \dots, \underset{d}{0}, \dots, i_D]$$

373 The input MDAs are assumed to have index set  $\{0\}$  in the operating dimension  $d$ ; otherwise,  $\vec{\bullet}(\oplus)$   
374 is undefined. We refrain from  $\vec{\bullet}_d(\oplus)$ 's subscript  $d$  when it is clear from the context. For brevity,  
375 we often write  $\oplus$  only, instead of  $\vec{\bullet}_d(\oplus)$ , and we usually use an infix notation for  $\oplus$ .

## 377 Multi-Dimensional Homomorphisms

378 Now that we have defined MDAs (Definition 1) and combine operators (Definition 2), we can define  
379 *Multi-Dimensional Homomorphisms (MDH)*. Intuitively, a function  $h$  operating on MDAs is an MDH  
380 iff we can apply the function independently to parts of its input MDA and combine the obtained  
381 intermediate results to the final result using combine operators; this can be imagined as a typical  
382 divide-and-conquer pattern. Compared to classical approaches, e.g., *list homomorphisms* [Bird  
383 1989; COLE 1995; Gorlatch 1999], a major characteristic of MDH functions is that they allow  
384 (de/re)-composing computations in multiple dimensions (e.g., in Figure 1, in the concatenation as  
385 well as point-wise addition dimensions), rather than being limited to a particular dimension only  
386 (e.g., only the concatenation dimension or only point-wise addition dimension, respectively). We  
387 will see later in this paper that a multi-dimensional (de/re)-composition approach is essential to  
388 efficiently exploit the hardware of modern architectures which require fine-grained cache blocking  
389 and parallelization strategies to achieve their full performance potential.

390 Figure 8 illustrates the MDH property informally on a simple, two-dimensional input MDA. In  
391 the left part of the figure, we split the input MDA in dimension 1 (i.e., horizontally) into two parts  $a_1$   
392

393 and  $\alpha_2$ , apply MDH function  $h$  independently to each part, and combine the obtained intermediate  
 394 results to the final result using the MDH function  $h$ 's combine operator  $\otimes_1$ . Similarly, in the right  
 395 part of Figure 8, we split the input MDA in dimension 2 (i.e., vertically) into parts and combine the  
 396 results via MDH function  $h$ 's second combine operator  $\otimes_2$ .

$$h\left(\underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}}_{\alpha_1 \oplus_1 \alpha_2}\right) = h\left(\underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}}_{\alpha_1} \oplus_1 \underbrace{\begin{bmatrix} 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}}_{\alpha_2}\right)$$

$$h\left(\underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}}_{\alpha_3 \oplus_2 \alpha_4}\right) = h\left(\underbrace{\begin{bmatrix} 1 & 2 \\ 5 & 6 \\ 9 & 10 \\ 13 & 14 \end{bmatrix}}_{\alpha_3} \oplus_2 \underbrace{\begin{bmatrix} 3 & 4 \\ 7 & 8 \\ 11 & 12 \\ 15 & 16 \end{bmatrix}}_{\alpha_4}\right)$$

Fig. 8. MDH property illustrated on a two-dimensional example computation

Figure 9 shows an artificial example in which we apply the MDH property (illustrated in Figure 8) recursively. We refer in Figure 9 to the part above the horizontal dashed lines as *de-composition phase* and to the part below dashed lines as *re-composition phase*.

**Definition 3** (Multi-Dimensional Homomorphism). A function

$$h : T^{\text{INP}} \left[ I_1, \dots, I_D \right] \rightarrow T^{\text{OUT}} \left[ J_1, \dots, J_D \right]$$

is a *Multi-Dimensional Homomorphism* (*MDH*) that has *input scalar type*  $T^{\text{INP}} \in \text{TYPE}$ , *output scalar type*  $T^{\text{OUT}} \in \text{TYPE}$ , and *dimensionality*  $D \in \mathbb{N}$ , iff for each  $d \in [1, D]_{\mathbb{N}}$ , there exists a combine operator  $\oplus_d$  (Definition 2), such that for any concatenated input MDA  $a_1 \mathbin{++}_d a_2$  in dimension  $d$ , the *homomorphic property* is satisfied:

$$h(\mathfrak{a}_1 \oplus_d \mathfrak{a}_2) = h(\mathfrak{a}_1) \otimes_d h(\mathfrak{a}_2)$$

We denote the type of MDHs concisely as MDH.

MDHs are defined such that applying them to a concatenated MDA in dimension  $d$  can be computed by applying the MDH  $h$  independently to the MDA's parts  $\alpha_1$  and  $\alpha_2$  and combining the results afterwards by using its combine operator  $\otimes_d$ , as also informally discussed above.

**Example 3 (Function Mapping).** A simple example MDH is *function mapping* [González-Vélez and Leyton 2010], computed by higher-order function  $\text{map}(f)(\mathfrak{a})$ , which applies a user-defined scalar function  $f : T^{\text{INP}} \rightarrow T^{\text{OUT}}$  to each element within a  $D$ -dimensional MDA  $\mathfrak{a}$ . Function  $\text{map}(f)$  is an MDH whose combine operators are concatenation  $+$  in all of its  $D$  dimensions (Example 1).

**Example 4 (Reduction).** Another MDH function is *reduction* [González-Vélez and Leyton 2010], implemented as higher-order function  $\text{red}(\oplus)(\alpha)$ , which combines all elements within a  $D$ -dimensional MDA  $\alpha$  using a user-defined binary function  $\oplus : T \times T \rightarrow T$ . Reduction's combine operators are point-wise combination  $\vec{\bullet}(\oplus)$  in all dimensions (Example 2).

We show how Examples 3 and 4 (and particularly also more advanced examples) are expressed in our high-level representation in Section 2.4, based on higher-order functions `md_hom`, `inp_view`, and `out_view` (Figure 5) which we introduce in the following.

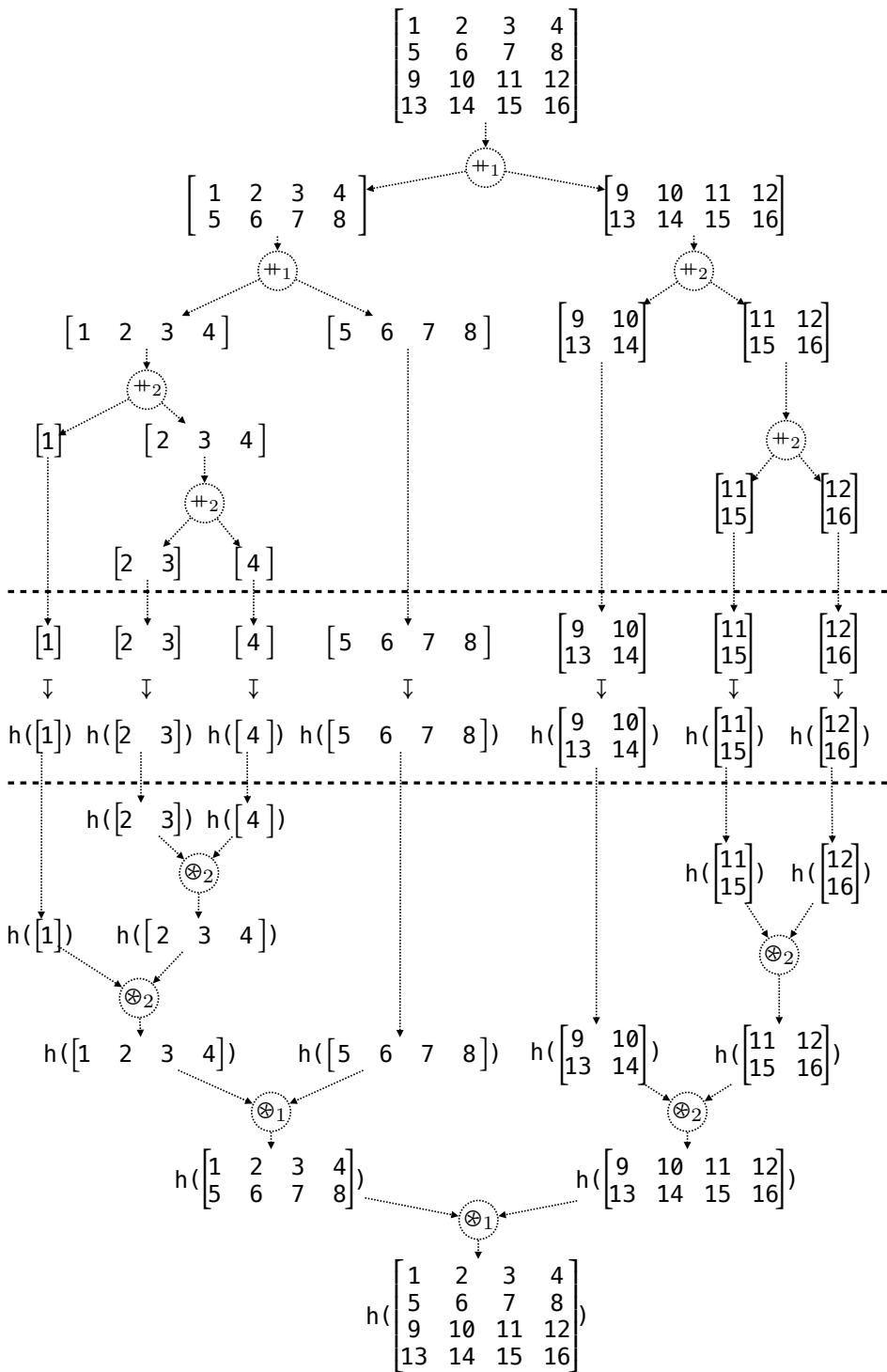


Fig. 9. MDH property recursively applied to a two-dimensional example computation

491 **Higher-Order Function `md_hom`**

492 We define higher-order function `md_hom` which conveniently expresses MDH functions in a uniform  
 493 and structured manner. For this, we exploit that any MDH function is uniquely determined by its  
 494 combine operators and its behavior on singleton MDAs, as informally illustrated in the following  
 495 figure:

$$496 \quad \begin{array}{c} \xrightarrow{\otimes_2} \\ 497 \quad h\left(\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}\right) = \begin{bmatrix} h([1]) & h([2]) & h([3]) & h([4]) \\ h([5]) & h([6]) & h([7]) & h([8]) \\ h([9]) & h([10]) & h([11]) & h([12]) \\ h([13]) & h([14]) & h([15]) & h([16]) \end{bmatrix} \end{array} \xrightarrow{\otimes_1} \begin{array}{c} \xrightarrow{\otimes_2} \\ \begin{bmatrix} f(1) & f(2) & f(3) & f(4) \\ f(5) & f(6) & f(7) & f(8) \\ f(9) & f(10) & f(11) & f(12) \\ f(13) & f(14) & f(15) & f(16) \end{bmatrix} \end{array}$$

502 Here,  $f$  is the function on scalar values that behaves the same as  $h$  when restricted to singleton  
 503 MDAs:  $f(a[i_1, \dots, i_D]) := h(a)$ , for any MDA  $a \in T[\{i_1\}, \dots, \{i_D\}]$  consisting of only one element  
 504 accessed by (arbitrary) indices  $i_1, \dots, i_D \in \mathbb{N}_0$ . For singleton MDAs, we usually use  $f$  instead of  $h$ ,  
 505 because  $f$  can be defined more conveniently by the user as  $h$  (which needs to handle MDAs of  
 506 arbitrary sizes, and not only singleton MDAs as  $f$ ). Also, since  $f$  takes as input a scalar value (rather  
 507 than a singleton MDA as  $h$ ), the type of  $f$  also becomes simpler, thereby also further simplifying  
 508 the definition of scalar functions for the user.

510 We now formally introduce function `md_hom` which uniformly expresses any MDH function, by  
 511 using only the MDH's behavior  $f$  on scalar values and its combine operators.

513 **Definition 4** (Higher-Order Function `md_hom`). The higher-order function `md_hom` is of type

$$514 \quad \text{md\_hom} : \underbrace{\text{SF}}_f \times \underbrace{(\text{CO} \times \dots \times \text{CO})}_{\otimes_1, \dots, \otimes_D} \xrightarrow{p} \underbrace{\text{MDH}}_{\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))}$$

517 where SF denotes scalar functions of type  $T^{\text{INP}} \rightarrow T^{\text{OUT}}$ . Function `md_hom` is partial (indicated  
 518 by  $\xrightarrow{p}$  instead of  $\rightarrow$ ), which we motivate after this definition. The function takes as input a  
 519 scalar function  $f$  and a tuple of  $D$ -many combine operators  $(\otimes_1, \dots, \otimes_D)$ , and it yields a function  
 520  $\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))$  which is defined as

$$521 \quad \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(a) := \underset{i_1 \in I_1}{\otimes_1} \dots \underset{i_D \in I_D}{\otimes_D} f(a[i_1, \dots, i_D])$$

523 The combine operators' underset notation denotes straightforward iteration<sup>8</sup>. For `md_hom`, we  
 524 require per definition the homomorphic property (Definition 3), i.e., for each  $d \in [1, D]_{\mathbb{N}}$ , it must  
 525 hold:

$$526 \quad \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(a_1 +_d a_2) =$$

$$527 \quad \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(a_1) \otimes_d \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(a_2)$$

530 Using Definition 4, we express any MDH function uniformly via higher-order function `md_hom`  
 531 using only the MDH's behavior  $f$  on scalar values and its combine operators  $\otimes_1, \dots, \otimes_D$ . The other  
 532 direction also holds: each function expressed via `md_hom` is an MDH function, because we require  
 533 the homomorphic property for `md_hom`.

534  
 535  
 536  
 537<sup>8</sup> We implicitly interpret the output scalar of function  $f$  as a singleton MDA, as combine operators operate on MDAs  
 538 and not on scalars (formal details provided in the Appendix, Definition 27).

540 Note that function `md_hom` is defined as partial function, because the homomorphic property  
 541 is not met for all potential combinations of combine operators, e.g.,  $\otimes_1 = +$  (point-wise addition)  
 542 and  $\otimes_2 = *$  (point-wise multiplication). However, in many real-world examples, an MDH's combine  
 543 operators are a mix of concatenations and point-wise combinations according to the same binary  
 544 function. The following lemma proves that any instance of the `md_hom` higher-order function for  
 545 such a mix of combine operators is a well-defined MDH function.

546 **Lemma 1.** Let  $\oplus : T \rightarrow T$  be an arbitrary but fixed associative and commutative binary function  
 547 on scalar type  $T \in \text{TYPE}$ . Let further  $\otimes_1, \dots, \otimes_D$  be combine operators of which any is either  
 548 concatenation (Example 1) or point-wise combination according to binary function  $\oplus$  (Example 2).  
 549

550 It holds that  $\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))$  is well defined.

551 PROOF. Proved in our Appendix, Section B.9. □

553 MDHs use as input and output uniformly an MDA. We introduce higher-order function `inp_view`  
 554 to prepare domain-specific inputs (e.g., a matrix and a vector for matrix-vector multiplication) as  
 555 an MDA, and we use function `out_view` to transform the output MDA back to the domain-specific  
 556 data requirements (like storing it as a transposed matrix in the case of matrix multiplication, or  
 557 splitting it into multiple outputs as we will see later with examples). We introduce both higher-order  
 558 functions in the following.

### 559 2.3 View Functions

561 In the following, after introducing *Buffers (BUF)* which represent domain-specific input and output  
 562 data in our approach (scalars, vectors, matrices, etc), we define, in Sections 2.3.1 and 2.3.2, the  
 563 concepts of *input views* and *output views* which are central building blocks in our approach. We  
 564 define *input views* as arbitrary functions that map a collection of user-defined BUFs to our internal  
 565 MDA data structure (Figure 5); higher-order function `inp_view` is then introduced to conveniently  
 566 compute an important class of input view functions that are relevant in real-world applications.  
 567 Correspondingly, Section 2.3.2 defines *output views* as functions that transform an MDA to a  
 568 collection of BUFs, and higher-order function `out_view` is introduced to conveniently compute  
 569 important output views.

570 Finally, we discuss in Section 2.3.3 the relationship between higher-order function `inp_view` and  
 571 `out_view`: we prove that both functions are inversely related to each other, allowing arbitrarily  
 572 switching between our internal MDA data structure and our domain-specific BUF representation  
 573 (as required for our code generation discussed later).

574 **Definition 5** (Buffer). A *Buffer (BUF)*  $b$  that has *dimensionality*  $D \in \mathbb{N}_0$ <sup>9</sup>, *size*  $N := \{N_1, \dots, N_D\} \in \mathbb{N}^D$ , and *scalar type*  $T \in \text{TYPE}$  is a function with the following signature:

$$b : [0, N_1)_{\mathbb{N}_0} \times \dots \times [0, N_D)_{\mathbb{N}_0} \rightarrow T \cup \{\perp\}$$

578 Here, we use  $\perp$  to denote the *undefined value*. We refer to  $[0, N_1)_{\mathbb{N}_0} \times \dots \times [0, N_D)_{\mathbb{N}_0} \rightarrow T \cup \{\perp\}$   
 579 as the *type* of BUF  $b$ , which we also denote as  $T^{N_1 \times \dots \times N_D}$ . Analogously to Notation 1, we write  
 580  $b[i_1, \dots, i_D]$  instead of  $b(i_1, \dots, i_D)$  to avoid a too heavy usage of parentheses.

581 In contrast to MDAs, a BUF always operates on a contiguous range of natural numbers starting  
 582 from 0, and a BUF may contain undefined values. These two differences allow us straightforwardly  
 583 transforming BUFs to data structures provided by low-level programming languages (e.g., *C arrays*  
 584 as used in OpenMP, CUDA, and OpenCL). Note that in our generated program code (discussed later  
 585 in Section 3), we implement MDAs as straightforward aliases that access BUFs, so that we do not

587 <sup>9</sup>We use the case  $D = 0$  to represent scalar values (details provided in the Appendix, Section B.11).

589 need to transform MDAs to low-level data structures and/or store them otherwise in  
 590 memory.

591 **2.3.1 Input Views.** We define *input views* as any function that compute an MDA from a collection  
 592 of user-defined BUFs. For example, in the case of `MatVec`, its input view takes as input two BUFs – a  
 593 matrix and a vector – and it yields a two-dimensional MDA containing pairs of elements taken from  
 594 the matrix and vector (illustrated in Figure 1). In contrast, the input view of `Jacobi1D` takes as input  
 595 a single BUF (vector) only, and it computes an MDA containing triples of BUF elements (Figure 2).  
 596

597 **Definition 6** (Input View). An *input view* from  $B$ -many BUFs,  $B \in \mathbb{N}$ , of arbitrary but fixed types  
 598  $T_b^{N_1^b \times \dots \times N_{D_b}^b}$ ,  $b \in [1, B]_{\mathbb{N}}$ , to an MDA of arbitrary but fixed type  $T[I_1, \dots, I_D]$  is any function  $\text{iv}$  of  
 599 type:

$$\text{iv} : \underbrace{\prod_{b=1}^B T_b^{N_1^b \times \dots \times N_{D_b}^b}}_{\text{BUFs}} \xrightarrow{p} T[\underbrace{I_1, \dots, I_D}_{\text{MDA}}]$$

600 We denote the type of  $\text{iv}$  as  $\text{IV}$ .

601 **Example 5** (Input View – `MatVec`). The input view of `MatVec` on a  $1024 \times 512$  matrix and 512-sized  
 602 vector (sizes are chosen arbitrarily) is defined as:

$$\underbrace{[M(i, k)]_{i \in [0, 1024]_{\mathbb{N}_0}, k \in [0, 512]_{\mathbb{N}_0}}}_{\text{Matrix}}, \underbrace{[v(k)]_{k \in [0, 512]_{\mathbb{N}_0}}}_{\text{Vector}} \mapsto \underbrace{[M(i, k), v(k)]_{i \in [0, 1024]_{\mathbb{N}_0}, k \in [0, 512]_{\mathbb{N}_0}}}_{\text{MDA}}$$

603 **Example 6** (Input View – `Jacobi1D`). The input view of `Jacobi1D` on a 512-sized vector is defined  
 604 as:

$$\underbrace{[v(i)]_{i \in [0, 512]_{\mathbb{N}_0}}}_{\text{Vector}} \mapsto \underbrace{[v(i+0), v(i+1), v(i+2)]_{i \in [0, 512-2]_{\mathbb{N}_0}}}_{\text{MDA}}$$

605 We introduce higher-order function `inp_view` which conveniently computes important input  
 606 views from user-defined index functions  $\text{idr}_{b,a} : \{0, 1, \dots\} \rightarrow \{0, 1, \dots\}$  for  $b \in [1, B]_{\mathbb{N}}$  and  
 607  $a \in [1, A_b]_{\mathbb{N}}$ , in a uniform, structured manner. Here,  $B \in \mathbb{N}$  represents the number of BUFs that the  
 608 computed input view will take as input, and  $A_b$  represents the number of accesses to the  $b$ -th BUF  
 609 required for computing an individual MDA element.

610 In the case of `MatVec` (Figure 1), we use  $B := 2$ , as `MatVec` has two input BUFs – a matrix  $M$   
 611 (the first input of `MatVec` and thus identified by  $b = 1$ ) and a vector  $v$  (identified by  $b = 2$ ). For the  
 612 number of accesses, we use for the matrix  $A_1 := 1$ , as one element is accessed within the matrix  $M$   
 613 to compute an individual MDA element – matrix element  $M[i, k]$  for computing MDA element  
 614 at position  $(i, k)$ ; for the vector, we use  $A_2 := 1$  – the single element  $v[k]$  is accessed. The index  
 615 functions of `MatVec` are:  $\text{idr}_{1,1}(i, k) := (i, k)$  (used to access the matrix) and  $\text{idr}_{2,1}(i, k) := (k)$  (used  
 616 for the vector).

617 In contrast, for `Jacobi1D` (Figure 2), we use  $B := 1$  (`Jacobi1D` has one input BUF  $v$ ),  $A_1 := 3$  (the  
 618 BUF is accessed three times to compute an individual MDA element at arbitrary position  $i$ : first  
 619 access  $v[i+0]$ , second access  $v[i+1]$ , and third access  $v[i+2]$ ). The index functions of `Jacobi1D`  
 620 are:  $\text{idr}_{1,1}(i) := (i+0)$ ,  $\text{idr}_{1,2}(i) := (i+1)$ , and  $\text{idr}_{1,3}(i) := (i+2)$ .

621 Figures 10 and 11 use the examples `MatVec` and `Jacobi1D` to informally illustrate how function  
 622 `inp_view` uses index functions to compute input views. In the two figures, we use domain-specific  
 623 identifiers for better clarity: in the case of `MatVec`, we use for BUFs identifiers  $M$  and  $v$  instead of  
 624  $b_1$  and  $b_2$ , as well as identifiers  $i$  and  $j$  instead of  $i_1$  and  $i_2$  for index variables; for `Jacobi1D`, we use  
 625 identifier  $v$  instead of  $b_1$ , and  $i$  instead of  $i_1$ .

638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652

653  
654

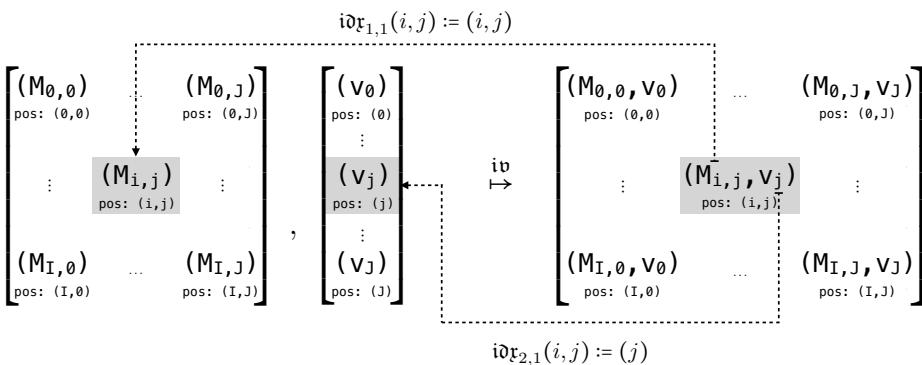


Fig. 10. Input view illustrated using the example MatVec

655  
656  
657  
658  
659  
660

661  
662

$\text{idr}_{1,3}(i) := (i + 2)$

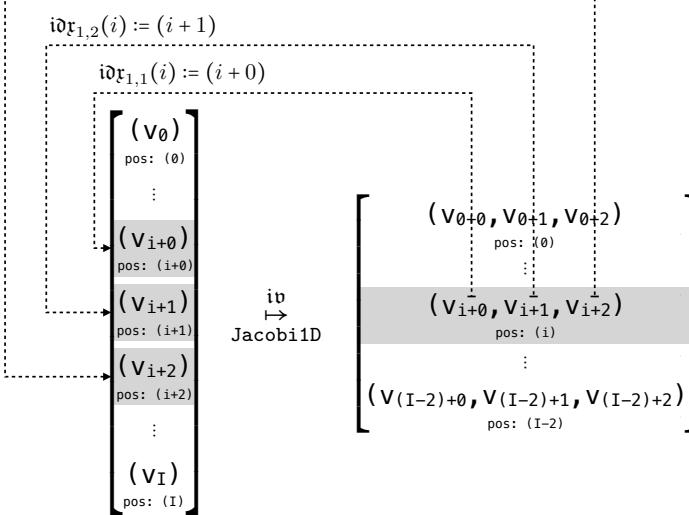


Fig. 11. Input view illustrated using the example Jacobi1D

663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679

680  
681  
682  
683  
684  
685  
686

687 **Definition 7** (Higher-Order Function `inp_view`). Function `inp_view` is of type

$$688 \quad \text{689} \quad \text{690} \quad \text{691} \quad \text{692} \quad \text{693} \quad \text{694} \quad \text{695} \quad \text{696} \quad \text{697} \quad \text{698} \quad \text{699} \quad \text{700} \quad \text{701} \quad \text{702} \quad \text{703} \quad \text{704} \quad \text{705} \quad \text{706} \quad \text{707} \quad \text{708} \quad \text{709} \quad \text{710} \quad \text{711} \quad \text{712} \quad \text{713} \quad \text{714} \quad \text{715} \quad \text{716} \quad \text{717} \quad \text{718} \quad \text{719} \quad \text{720} \quad \text{721} \quad \text{722} \quad \text{723} \quad \text{724} \quad \text{725} \quad \text{726} \quad \text{727} \quad \text{728} \quad \text{729} \quad \text{730} \quad \text{731} \quad \text{732} \quad \text{733} \quad \text{734} \quad \text{735}$$

$$\text{inp\_view} : \left( \underbrace{\begin{matrix} B \\ \times \\ b=1 \end{matrix} \quad \begin{matrix} A_b \\ \times \\ a=1 \end{matrix}}_{\text{Buffer Access}} \quad \underbrace{\text{IDX-FCT}}_{\text{Index Function: } \text{idx}_{p,a}} \right) \rightarrow \underbrace{\text{IV}}_{\text{Input View: } \text{iv}}$$

$$\underbrace{\text{Index Functions: } \text{idx}_{1,1}, \dots, \text{idx}_{B,A_B}}$$

and it is defined as:

$$\underbrace{(\text{idx}_{b,a})_{b \in [1,B]_{\mathbb{N}}, a \in [1,A_b]_{\mathbb{N}}}}_{\text{Index Functions}} \mapsto \left( \underbrace{b_1, \dots, b_B}_{\text{BUFs}} \right) \xrightarrow{\text{iv}} \underbrace{a}_{\text{MDA}}$$

$$\underbrace{\text{Input View}}$$

for

$$a[i_1, \dots, i_D] := (a_{b,a}[i_1, \dots, i_D])_{b \in [1,B]_{\mathbb{N}}, a \in [1,A_b]_{\mathbb{N}}}$$

and

$$a_{b,a}[i_1, \dots, i_D] := b_b[\text{idx}_{b,a}(i_1, \dots, i_D)]$$

Higher-order function `inp_view` takes as input a collection of index functions of types `IDX-FCT`, and it computes an input view of type `IV` (Definition 6) based on the index functions, as illustrated in Figures 10 and 11.

Note that function `inp_view` is not capable of computing every kind of input view function (Definition 6). For example, `inp_view` cannot be used for computing MDAs that are required for computations on sparse data formats [Hall 2020], because such MDAs need dynamically accessing BUFs. This limitation of `inp_view` can be relaxed by generalizing our index functions toward taking additional, dynamic input arguments, which we consider as future work (as also outlined in Section 8).

**Notation 2** (Input Views). For better readability, we use the following notation for the 2-dimensional structure of index functions taken as input by function `inp_view`, inspired by Lattner et al. [2021]:

$$\text{inp\_view}(\text{ID}_1 : \text{idx}_{1,1}, \dots, \text{idx}_{1,A_1}, \dots, \text{ID}_B : \text{idx}_{B,1}, \dots, \text{idx}_{B,A_B})$$

Here,  $\text{ID}_1, \dots, \text{ID}_B$  denote arbitrary, user-defined identifiers.

**Example 7.** Function `inp_view` is used for `MatVec` and `Jacobi1D` (in Notation 2) as follows:

$$\text{MatVec:} \quad \text{inp\_view}(\text{M: } \underbrace{(i,k) \mapsto (i,k)}_{\substack{a=1 \\ b=1}}, \text{v: } \underbrace{(i,k) \mapsto (k)}_{\substack{a=1 \\ b=2}})$$

$$\text{Jacobi1D:} \quad \text{inp\_view}(\text{v: } \underbrace{(i) \mapsto (i+0)}_{\substack{a=1 \\ b=1}}, \underbrace{(i) \mapsto (i+1)}_{\substack{a=2 \\ b=1}}, \underbrace{(i) \mapsto (i+2)}_{\substack{a=3 \\ b=1}})$$

2.3.2 *Output Views*. An *output view* is the counterpart of an input view: in contrast to an input view which maps BUFS to an MDA, an output view maps an MDA to a collection of BUFS. In the following, we define output views, and we introduce higher-order function `out_view` which computes output views in a structured manner, analogously to function `inp_view` for input views.

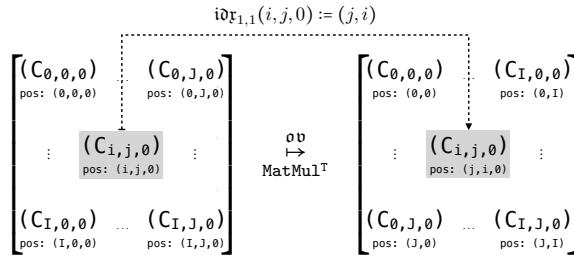


Fig. 12. Output view illustrated using the example *transposed Matrix Multiplication*

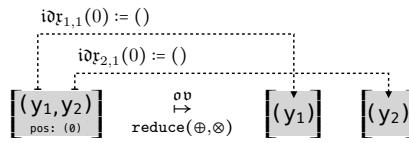


Fig. 13. Output view illustrated using the example *double reduction*

Figures 12 and 13 illustrate output views informally using the examples *transposed matrix multiplication* and *double reduction*.

In the case of transposed matrix multiplication, the computed MDA (the computation of matrix multiplication is presented later and not relevant for our following considerations) is stored via an output view as a matrix in a transposed fashion, using index function  $(i, j, 0) \mapsto (j, i)$ . Here, the MDA's third dimension (accessed via index 0) represents the so-called reduction dimension of matrix multiplication, and it contains only one element after the computation, as all elements in this dimension are combined via addition.

For double reduction, we combine the elements within the vector twice – once using operator  $\oplus$  (e.g.,  $\oplus = +$  addition) and once using operator  $\otimes$  (e.g.,  $\otimes = *$  multiplication). The final outcome of double reduction is a singleton MDA containing a pair of two elements that represent the combined vector elements (e.g., the elements' sum and product). We store this MDA via an output view as two individual scalar values, using index functions  $(0) \mapsto ()$ <sup>10</sup> for both pair elements.

**Definition 8** (Output View). An *output view* from an MDA of arbitrary but fixed type  $T[I_1, \dots, I_D]$  to  $B$ -many BUFS,  $B \in \mathbb{N}$ , of arbitrary but fixed types  $T_b^{N_1^b \times \dots \times N_{D_b}^b}$ ,  $b \in [1, B]_{\mathbb{N}}$ , is any function `ov` of type:

$$\text{ov} : \underbrace{T[I_1, \dots, I_D]}_{\text{MDA}} \rightarrow_p \underbrace{\bigtimes_{b=1}^B T_b^{N_1^b \times \dots \times N_{D_b}^b}}_{\text{BUFS}}$$

We denote the type of `ov` as `OV`.

<sup>10</sup>The empty braces denote accessing a scalar value (details provided in the Appendix, Section B.11).

785 **Example 8** (Output View – MatVec). The output view of MatVec computing a 1024-sized vector  
 786 (size is chosen arbitrarily), of integers  $\mathbb{Z}$ , is defined as:

787 
$$\underbrace{[w(i)]_{i \in [0, 1024]_{\mathbb{N}_0}, k \in \{0\}}}_{\text{MDA}} \mapsto \underbrace{[w(i)]_{i \in [0, 1024]_{\mathbb{N}_0}}}_{\text{Vector}}$$

788

789

790 **Example 9** (Output View – Jacobi1D). The output view of Jacobi1D computing a  $(512 - 2)$ -sized  
 791 vector is defined as:

792 
$$\underbrace{[w(i)]_{i \in [0, 512-2]_{\mathbb{N}_0}, k \in \{0\}}}_{\text{MDA}} \mapsto \underbrace{[w(i)]_{i \in [0, 512-2]_{\mathbb{N}_0}}}_{\text{Vector}}$$

793

794

795 We define higher-order function `out_view` formally as follows.

796 **Definition 9** (Higher-Order Function `out_view`). Function `out_view` is of type

797 
$$\text{out\_view} : \underbrace{\begin{matrix} B & A_b \\ \times & \times \\ b=1 & a=1 \end{matrix}}_{\text{Buffer Access}} \underbrace{\text{IDX-FCT}}_{\text{Index Function: } \text{id}\mathfrak{x}_{b,a}} \rightarrow \underbrace{\text{OV}}_{\text{Output View: ov}}$$

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

804 which differs from `inp_view`'s type only in mapping index functions to OV (Definition 8), rather  
 805 than IV (Definition 6). Function `out_view` is defined as:

806 
$$\underbrace{(\text{id}\mathfrak{x}_{b,a})_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_b]_{\mathbb{N}}}}_{\text{Index Functions}} \mapsto \underbrace{\begin{matrix} \text{MDA} & \xrightarrow{\text{ov}} & \underbrace{\text{a}}_{\text{BUFs}} \xrightarrow{\text{ov}} (\underbrace{\text{b}_1, \dots, \text{b}_B}_{\text{BUFs}}) \end{matrix}}_{\text{Output View}}$$

807

808

809

810

811 for

812 
$$\text{b}_b[\text{id}\mathfrak{x}_{b,a}(i_1, \dots, i_D)] := \text{a}_{b,a}[i_1, \dots, i_D]$$

813 and

814 
$$(\text{a}_{b,a}[i_1, \dots, i_D])_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_b]_{\mathbb{N}}} := \text{a}[i_1, \dots, i_D]$$

815 i.e.,  $\text{a}_{b,a}[i_1, \dots, i_D]$  is the element at point  $i_1, \dots, i_D$  within MDA  $\text{a}$  that belongs to the  $a$ -th access  
 816 of the  $b$ -th BUF. We set  $\text{b}_b[j_1, \dots, j_{D_b}] := \perp$  (symbol  $\perp$  denotes the undefined value) for all BUF  
 817 indices  $(j_1, \dots, j_{D_b}) \in [0, N_1^b]_{\mathbb{N}_0} \times \dots \times [0, N_D^b]_{\mathbb{N}_0} \setminus \bigcup_{a \in [1, A_b]_{\mathbb{N}}} \xrightarrow{\text{BUF}}_{\text{MDA } b, a}^d (I_1, \dots, I_D)$  which are not in  
 818 the function range of the index functions.

819 Note that the computed output view  $\text{ov}$  is partial (indicated by  $\rightarrow_p$  in Definition 8), because for non-  
 820 injective index functions, it must hold  $\text{id}\mathfrak{x}_{b,a}(i_1, \dots, i_D) = \text{id}\mathfrak{x}_{b,a'}(i'_1, \dots, i'_D) \Rightarrow \text{a}_{b,a}[i_1, \dots, i_D] =$   
 821  $\text{a}_{b,a'}[i'_1, \dots, i'_D]$ , which may not be satisfied for each potential input MDA of the computed view.

822 **Notation 3** (Output Views). Analogously to Notation 2, we denote `out_view` for a particular  
 823 choice of index functions as:

824 
$$\text{out\_view}(\text{ID}_1 : \text{id}\mathfrak{x}_{1,1}, \dots, \text{id}\mathfrak{x}_{1,A_1}, \dots, \text{ID}_B : \text{id}\mathfrak{x}_{B,1}, \dots, \text{id}\mathfrak{x}_{B,A_B})$$

825 **Example 10.** Function `out_view` is used for MatVec and Jacobi1D (in Notation 3) as follows:

826 
$$\text{MatVec: } \text{out\_view}(\text{w: } \underbrace{(i, k) \mapsto (i)}_{\text{a=1}}) \quad \text{Jacobi1D: } \text{out\_view}(\text{w: } \underbrace{(i) \mapsto (i)}_{\text{a=1}})$$

827

828

829

830

831

832

833

2.3.3 *Relation between View Functions.* View functions transform data from their domain-specific representation to the MDA-based representation (via input views) and back (via output views). In our implementation presented later, we aim to access data uniformly in form of MDAs, thereby being independent of domain-specific data representations. However, we aim to store the data physically in the domain-specific format, as such format is usually the most efficient data representation, e.g., storing the input data of MatVec as a matrix and vector, rather than as a single MDA which contains many redundancies (e.g., each vector element once per row of the input matrix, as illustrated in Figure 10).

The following lemma proves that functions `inp_view` and `out_view` are invertible and that they are each others inverses. Consequently, the lemma shows how we can arbitrarily switch between the domain-specific data representation and the MDA representation, and consequently also that we can implicitly identify MDAs with their domain-specific data representation (represented in our formalism as BUFs, Definition 5). For example, for computing MatVec, we will specify the computations via pattern `md_hom` which operates on MDAs (see Figure 5), but we use the view functions in our implementation to implicitly forward the MDA accesses to accesses on the physically stored BUFs.

**Lemma 2.** Let

`inp_view( ID1 : idx1,1, ..., idx1,A1 , ..., IDB : idxB,1, ..., idxB,AB )`

and

`out_view( ID1 : idx1,1, ..., idx1,A1 , ..., IDB : idxB,1, ..., idxB,AB )`

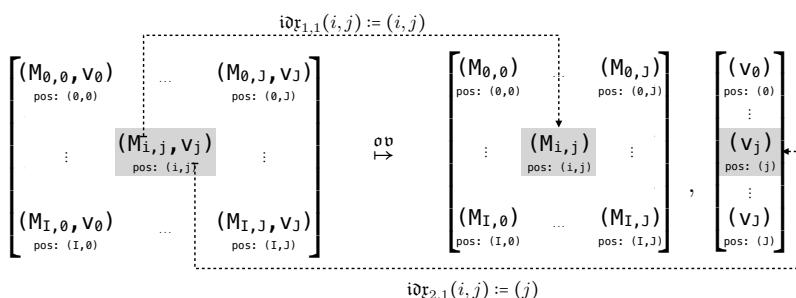
be two arbitrary instances of functions `inp_view` and `out_view` (in Notations 2 and 3), both using the same index functions  $idx_{1,1}, \dots, idx_{B,A_B}$ .

It holds (index functions omitted via ellipsis for brevity):

$inp\_view(\dots) \circ out\_view(\dots) = out\_view(\dots) \circ inp\_view(\dots) = id$

PROOF. Follows immediately from Definitions 7 and 9.  $\square$

The following figure illustrates the lemma using as example the inverse of MatVec's input view (shown in Figure 10):



## 2.4 Examples

Figure 14 shows how our high-level representation is used for expressing important data-parallel computations. For brevity, we state only the index functions, scalar function, and combine operators of the higher-order functions; an expression as in Figure 6 is then obtained by straightforwardly inserting these building blocks into the higher-order functions.

883	md_hom	f	$\oplus_1$	$\oplus_2$	$\oplus_3$	$\oplus_4$	Views	inp_view		out_view		
884	Dot	*	+					A		B	C	
885	MatVec	*	+	+				$(k) \mapsto (k)$		$(k) \mapsto (k)$	$(k) \mapsto ()$	
886	MatMul	*	+	+	+	+		$(i,k) \mapsto (i,k)$		$(i,k) \mapsto (k)$	$(i,k) \mapsto (i)$	
887	MatMul <sup>T</sup>	*	+	+	+	+		$(i,j,k) \mapsto (i,k)$		$(i,j,k) \mapsto (k,j)$	$(i,j,k) \mapsto (i,j)$	
888	bMatMul	*	+	+	+	+		$(b,i,j,k) \mapsto (b,i,k)$		$(b,i,j,k) \mapsto (b,k,j)$	$(b,i,j,k) \mapsto (b,i,j)$	
1) Linear Algebra Routines												
890	md_hom	f	$\oplus_1$	$\oplus_2$	$\oplus_3$	$\oplus_4$	$\oplus_5$	$\oplus_6$	$\oplus_7$	$\oplus_8$	$\oplus_9$	$\oplus_{10}$
891	Conv2D	*	+	+	+	+	+					
892	MCC	*	+	+	+	+	+	+	+	+		
893	MCC_Capsule	*	+	+	+	+	+	+	+	+	+	+
894	Views	inp_view						out_view				
895		I						0				
896	Conv2D	$(p,q,r,s) \mapsto (p+r, q+s)$			$(p,q,r,s) \mapsto (r,s)$			$(p,q,r,s) \mapsto (p,q)$				
897	MCC	$(n,p,...) \mapsto (n,p+r, q+s, c)$			$(n,p,...) \mapsto (k,r,s,c)$			$(n,p,...) \mapsto (n,p,q,k)$				
898	MCC_Capsule	$(n,p,...) \mapsto (n,p+r, q+s, c, mi, mk)$			$(n,p,...) \mapsto (k,r,s,c, mk, mj)$			$(n,p,...) \mapsto (n,p,q,k, mi, mj)$				
2) Convolution Stencils												
899	md_hom	f	$\oplus_1$	$\dots$	$\oplus_6$	$\oplus_7$	Views	inp_view		out_view		
900	CCSD(T)	*	+	...	+	+		A		B	C	
901	I1	$(a, \dots, g) \mapsto (g, d, a, b)$			$(a, \dots, g) \mapsto (e, f, g, c)$			$(a, \dots, g) \mapsto (a, \dots, f)$				
902	I2	$(a, \dots, g) \mapsto (g, d, a, c)$			$(a, \dots, g) \mapsto (e, f, g, b)$			$(a, \dots, g) \mapsto (a, \dots, f)$				
3) Quantum Chemistry												
904	md_hom	f	$\oplus_1$	$\dots$	$\oplus_6$	$\oplus_7$	Views	inp_view		out_view		
905	Jacobi1D	J <sub>1D</sub>	+					I		0		
906	Jacobi2D	J <sub>2D</sub>	+	+			Jacobi1D	$(i1) \mapsto (i1+0), (i1) \mapsto (i1+1), \dots$		$(i1) \mapsto (i1)$		
907	Jacobi3D	J <sub>3D</sub>	+	+	+		Jacobi2D	$(i1, i2) \mapsto (i1+0, i2+1), \dots$		$(i1, i2) \mapsto (i1, i2)$		
4) Jacobi Stencils												
910	md_hom	f	$\oplus_1$	$\dots$	$\oplus_6$	$\oplus_7$	Views	inp_view		out_view		
911	PRL	wght	+		max <sub>PRL</sub>		PRL	N	E	M		
912								$(i, j) \mapsto (i)$		$(i, j) \mapsto (j)$		
5) Probabilistic Record Linkage												
914	md_hom	f	$\oplus_1$	$\oplus_2$	Views	inp_view		out_view				
915	Histo	f <sub>Histo</sub>	+	+	Histo	Elems	Bins	Out				
916	GenHisto	f	⊕	+	GenHisto	$(e, b) \mapsto (e)$		$(e, b) \mapsto (b)$				
6) Histogram												
919	md_hom	f	$\oplus_1$	Views	inp_view		out_view					
920	map(f)	f	+	map(f)	I	$(i) \mapsto (i)$	$(i) \mapsto (i)$					
921	reduce(⊕)	id	⊕	reduce(⊕)	$(i) \mapsto (i)$	$(i) \mapsto ()$	$(i) \mapsto ()$					
922	reduce(⊕, ⊗)	$(x) \mapsto (x, x)$	$(\oplus, \otimes)$	reduce(⊕, ⊗)	$(i) \mapsto (i)$	$(i) \mapsto ()$	$(i) \mapsto ()$					
7) Map/Reduce Patterns												
925	md_hom	f	$\oplus_1$	Views	inp_view		out_view					
926	scan(⊕)	id	$\#_{\text{prefix-sum}}(\oplus)$	scan(⊕)	A	$(i) \mapsto (i)$	$(i) \mapsto (i)$					
927	MBBS	id	$\#_{\text{prefix-sum}}(+) +$	MBBS	$(i, j) \mapsto (i, j)$		$(i, j) \mapsto (i)$					
8) Prefix Sum Computations												

Fig. 14. Data-parallel computations expressed in our high-level representation

932 Subfigure 1. shows how our high-level representation is used for expressing linear algebra  
 933 routines: 1) Dot (*Dot Product*); 2) MatVec (*Matrix-Vector Multiplication*); 3) MatMul (*Matrix Multipli-*  
 934 *cation*); 4)  $\text{MatMul}^\top$  (*Transposed Matrix Multiplication*) which computes matrix multiplication on  
 935 transposed input and output matrices; 5) bMatMul (*batched Matrix Multiplication*) where multiple  
 936 matrix multiplications are computed using matrices of the same sizes.

937 We can observe from the subfigure that our high-level expressions for the routines naturally  
 938 evolve from each other. For example, the `md_hom` expression for MatVec differs from the `md_hom`  
 939 expression for Dot by only containing a further concatenation dimension `++` for its  $i$  dimension. We  
 940 consider this close relation between the high-level expressions of MatVec and Dot in our approach  
 941 as natural and favorable, as MatVec can be considered as computing multiple times Dot – one  
 942 computation of Dot for each value of MatVec’s  $i$  dimension. Similarly, the `md_hom` expression  
 943 for MatMul is very similar to the expression of MatVec, by containing the further concatenation  
 944 dimension  $j$  for MatMul’s  $j$  dimension. The same applies to bMatMul: its `md_hom` expression is the  
 945 expression of MatMul augmented with one further concatenation dimension.

946 Regarding  $\text{MatMul}^\top$ , the basic computation part of  $\text{MatMul}^\top$  and MatMul are the same, which is  
 947 exactly reflected in our formalisms: both  $\text{MatMul}^\top$  and MatMul are expressed using exactly the same  
 948 `md_hom` expression. The differences between  $\text{MatMul}^\top$  and MatMul lies only in the data accesses –  
 949 transposed accesses in the case of  $\text{MatMul}^\top$  and non-transposed accesses in the case of MatMul. Data  
 950 accesses are expressed in our formalism, in a structured way, via the view functions: for example,  
 951 for  $\text{MatMul}^\top$ , we use for its first input matrix  $A$  index function  $(i, j, k) \mapsto (k, i)$ , for transposed  
 952 access, instead of using index function  $(i, j, k) \mapsto (i, k)$  as for MatMul’s non-transposed accesses.

953 Note that all `md_hom` expressions in the subfigure are well defined according to Lemma 1.

955 Subfigure 2. shows how convolution-style stencil computations are expressed in our high-level  
 956 representation: 1) Conv2D expresses a standard convolution that uses a 2D sliding window [Podlozh-  
 957 nyuk 2007]; 2) MCC expresses a so-called *Multi-Channel Convolution* [Dumoulin and Visin 2018] – a  
 958 generalization of Conv2D that is heavily used in the area of deep learning; 3) MCC\_Capsule is a  
 959 recent generalization of MCC [Hinton et al. 2018] which attracted high attention due to its relevance  
 960 for advanced deep learning neural networks [Barham and Isard 2019].

961 While our `md_hom` expressions for convolutions are quite similar to those of linear algebra routines  
 962 (they all use multiplication `*` as scalar function, and a mix of concatenations `++` and point-wise  
 963 additions `+` as combine operators), the index functions used for the view functions of convolutions  
 964 are notably different from those used for linear algebra routines: the index functions of convolutions  
 965 contain arithmetic expressions (e.g.,  $p+r$  and  $q+s$ ), thereby access neighboring elements in their  
 966 input – a typical access pattern in stencil computations requiring special optimizations [Hagedorn  
 967 et al. 2018]. Moreover, convolution-style computations are often high-dimensional (e.g., 10 dimen-  
 968 sions in the case of MCC\_Capsule), whereas linear algebra routines usually rely on less dimensions.  
 969 Our experiments in Section 5 confirm that optimizations respecting the data access patterns and  
 970 high dimensionality of convolution computations, as in our approach, usually achieve significantly  
 971 higher performance than using optimizations chosen toward linear algebra routines as in vendor  
 972 libraries provided by NVIDIA and Intel for convolutions [Li et al. 2016].

974 Subfigure 3. shows how quantum chemistry computation *Coupled Cluster* (CCSD(T)) [Kim et al.  
 975 2019] is expressed in our high-level representation. The CCSD(T) computation notably differs from  
 976 those of linear algebra routines and convolution-style stencils, by accessing its high-dimensional  
 977 input data in sophisticated transposed fashions: for example, the view function of CCSD(T)’s  
 978 *instance one* (denoted as `I1` in the subfigure) uses indices `a` and `b` to access the last two dimensions

981 of its  $A$  input tensor (rather than the first two dimensions of the tensor, as would be the case for  
982 non-transposed accesses).

983 The subfigure presents only two CCSD(T) instances, for brevity – in our experiments in Section 5,  
984 we present experimental results for nine different real-world CCSD(T) instances.  
985

986 *Subfigures 4-6.* show computations whose scalar functions and combine operators are different  
987 from those used in Subfigures 1-3 (which are straightforward multiplications  $*$ , concatenation, and  
988 point-wise additions  $+$  only in Subfigures 1-3). For example, Jacobi stencils (Subfigure 4) use as  
989 scalar function the Jacobi-specific computation  $J_{\text{nd}}$  [Cecilia et al. 2012], and *Probabilistic Record*  
990 *Linkage (PRL)* [Christen 2012], which is used to identify duplicate entries in a data base, uses a  
991 PRL-specific both scalar function  $wght$  and combine operator  $\max_{\text{PRL}}$  (point-wise combination via  
992 the PRL-specific binary function  $\max_{\text{PRL}}$ ) [Rasch et al. 2019b]. Histograms, in their generalized  
993 version [Henriksen et al. 2020] (denoted as  $\text{GenHisto}$  in Subfigure 6), use an arbitrary, user-defined  
994 scalar function  $f$  and a user-defined associative and commutative combine operator  $\oplus$ ; the standard  
995 histogram variant  $\text{Histo}$  is then a particular instance of  $\text{GenHist}$ , for  $\oplus = +$  (point-wise addition)  
996 and  $f = f_{\text{Histo}}$ , where  $f_{\text{Histo}}(e, b) = 1$  iff  $e = b$  and  $f_{\text{Histo}}(e, b) = 0$  otherwise. Histograms are often  
997 analyzed regarding their runtime complexity [Henriksen et al. 2020]; we provide such a discussion  
998 for our MDH-based Histogram implementation in our Appendix, Section B.12, for the interested  
999 reader.  
1000

1001 *Subfigure 7.* shows how typical  $\text{map}$  and  $\text{reduce}$  patterns [González-Vélez and Leyton 2010]  
1002 are implemented in our high-level representation. Examples  $\text{map}(f)$  and  $\text{reduce}(\oplus)$  (discussed  
1003 in Examples 3 and 4) are simple and thus straightforwardly expressed in our representation. In  
1004 contrast, example  $\text{reduce}(\oplus, \otimes)$  is more complex and shows how  $\text{reduce}(\oplus)$  is extended toward  
1005 combining the input vector simultaneously twice – once combining vector elements via operator  
1006  $\oplus$  and once using operator  $\otimes$ . The outcome of  $\text{reduce}(\oplus, \otimes)$  are two scalars – one representing  
1007 the result of combination via  $\oplus$  and the other of combination via  $\otimes$  – which we map via the output  
1008 view to output elements  $0_1$  (result of  $\oplus$ ) and  $0_2$  (result of  $\otimes$ ), correspondingly; this is also illustrated  
1009 in Figure 13.  
1010

1011 *Subfigure 8.* presents *prefix-sum computations* [Blelloch 1990] which differ from the computations  
1012 in Subfigures 1-7 in terms of their combine operators: the operator used for expressing computations  
1013 in Subfigure 8 is different from concatenation (Example 1) and point-wise combinations (Example 2).  
1014 Computation  $\text{scan}(\oplus)$  uses as combine operator  $+\text{prefix-sum}(\oplus)$  which computes prefix-sum [Gor-  
1015 latch and Lengauer 1997] (formally defined in our Appendix, Section B.13) according to binary  
1016 operator  $\oplus$ , and MBBS (Maximum Bottom Box Sum) [Farzan and Nicolet 2019] uses a particular  
1017 instance of prefix-sum for  $\oplus = +$  (addition).  
1018

### 1020 3 LOW-LEVEL REPRESENTATION FOR DATA-PARALLEL COMPUTATIONS

1021 We introduce our low-level representation for expressing data-parallel computations. In contrast to  
1022 our high-level representation, our low-level representation explicitly expresses the de-composition  
1023 and re-composition of computations (informally illustrated in Figure 3). Moreover, our low-level  
1024 representation is designed such that it can be straightforwardly transformed to executable program  
1025 code, because it explicitly captures and expresses data movement and parallelization optimizations.  
1026

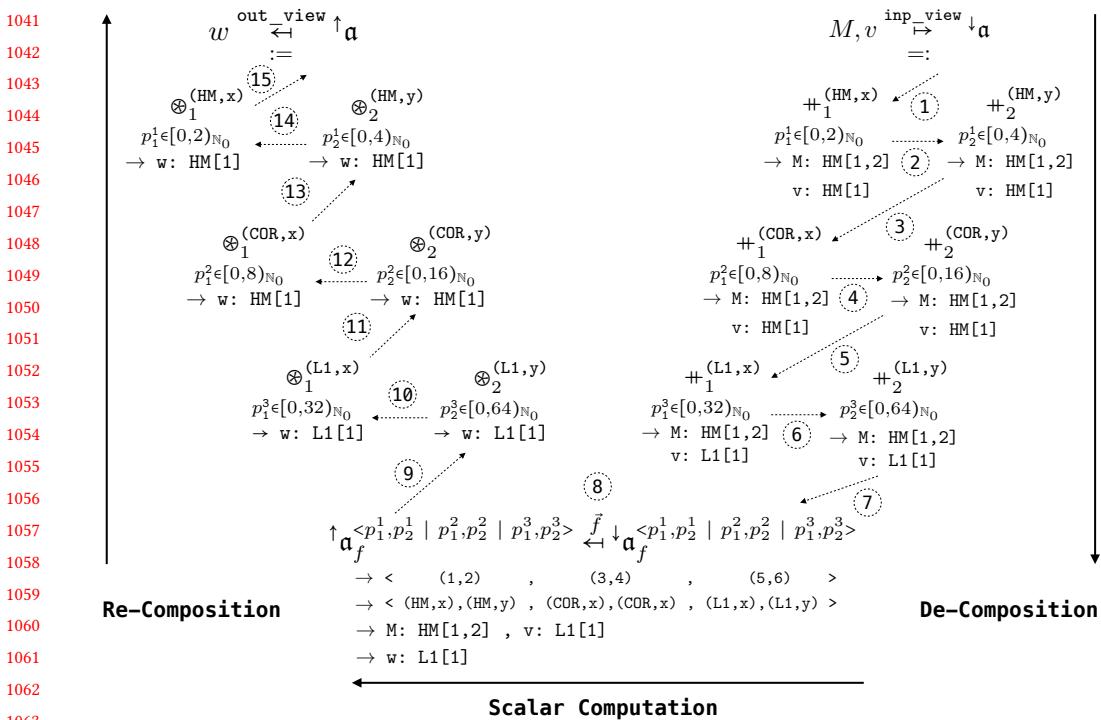
1027 In the following, after briefly discussing an introductory example in Section 3.1, we introduce in  
1028 Section 3.2 our formal representation of computer systems, to which we refer to as *Abstract System*  
1029

1030 *Model (ASM)*. Based on this model, we define *low-level MDAs*, *low-level BUFs*, and *low-level combine operators* in Section 3.3, which are basic building blocks of our low-level representation.

1032 Note that all the details and concepts discussed in this section are transparent to the end user  
 1033 of our system and therefore the user is not exposed to these details: expressions in our low-level  
 1034 representation are fully automatically generated from expressions in our high-level representation  
 1035 for the user (Figure 4), according to the methodologies presented later in Section 4 and auto-  
 1036 tuning [Rasch et al. 2021].

1038 **3.1 Introductory Example**

1040



1064 Fig. 15. Low-level expression for straightforwardly computing Matrix-Vector Multiplication (MatVec) on a  
 1065 simple, artificial architecture with two memory layers (HM and L1) and one core layer (COR). Dotted lines  
 1066 indicate data flow.

1067

1068 Figure 15 illustrates our low-level representation by showing how MatVec (Matrix-Vector Multiplication)  
 1069 is expressed in our representation. In our example, we use an input matrix  $M \in T^{512 \times 4096}$   
 1070 of size  $512 \times 4096$  (size chosen arbitrarily) that has an arbitrary but fixed scalar type  $T \in \text{TYPE}$ ; the  
 1071 input vector  $v \in T^{4096}$  is of size  $4096$ , correspondingly.

1072 For better illustration, we consider for this introductory example a straightforward, artificial  
 1073 target architecture that has only two memory layers – *Host Memory* (HM) and *Cache Memory* (L1) –  
 1074 and one *Core Layer* (COR) only; our examples presented and discussed later in this section target  
 1075 real-world architectures (e.g., CUDA-capable NVIDIA GPUs). The particular values of tuning  
 1076 parameters (discussed in detail later in this section), such as the number of threads and the order of  
 1077

1078

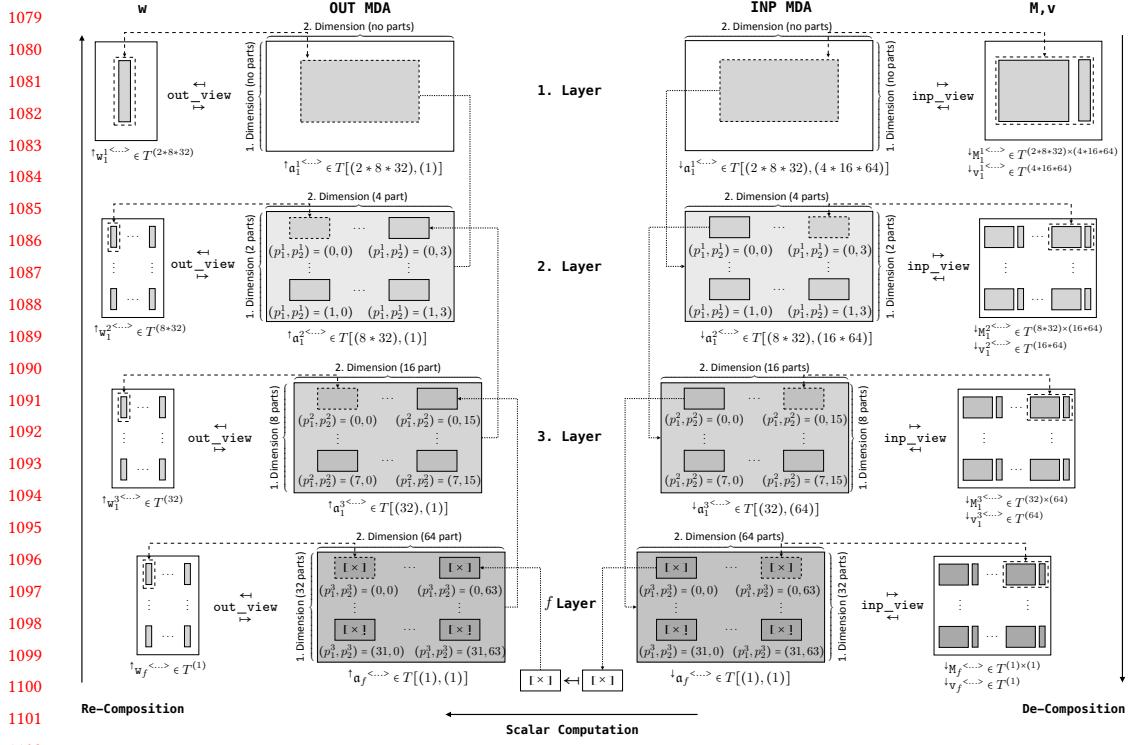


Fig. 16. Illustration of multi-layered, multi-dimensional MDA partitioning using the example MDA from Figure 15. In this example, we use three layers and two dimensions, according to Figure 15.

combine operators, are chosen by hand for this example and as straightforward for simplicity and better illustration.

Our low-level representations work in three phases: 1) *de-composition* (steps 1-7, in the right part of Figure 15), 2) *scalar* (step 8, bottom part of the figure), 3) *re-composition* (steps 9-15, left part). Steps are arranged from right to left, inspired by the application order of function composition.

**1. De-Composition Phase:** The de-composition phase (steps 1-7 in Figure 15) partitions input MDA  $\downarrow a$  (in the top right of Figure 15) to the structure  $\downarrow a_f$  (bottom right) to which we refer to as *low-level MDA* and define formally in the next subsection. The low-level MDA represents a partitioning of MDA  $\downarrow a$  (a.k.a *hierarchical, multi-dimensional tiling* in programming terminology), where each particular choice of indices  $p_1^1 \in [0, 2]_{\mathbb{N}_0}$ ,  $p_2^1 \in [0, 4]_{\mathbb{N}_0}$ ,  $p_1^2 \in [0, 8]_{\mathbb{N}_0}$ ,  $p_2^2 \in [0, 16]_{\mathbb{N}_0}$ ,  $p_1^3 \in [0, 32]_{\mathbb{N}_0}$ ,  $p_2^3 \in [0, 64]_{\mathbb{N}_0}$  refers to an MDA that represents an individual part of MDA  $\downarrow a$  (a.k.a. *tile* in programming – informally illustrated in Figure 7). The partitions are arranged on multiple layers (indicated by the  $p$ 's superscripts) and in multiple dimensions (indicated by subscripts) – as illustrated in Figure 16 – according to the memory/core layers of the target architecture and dimensions of the MDH computation: we partition for each of the target architecture's three layers (HM, L1, COR) and in each of the two dimensions of the MDH (dimensions 1 and 2, as we use example MatVec in Figure 15, which represents a two-dimensional MDH computation). Consequently, our partitioning approach allows efficiently exploiting each particular layer of the target architecture (both memory and core layers), and also optimizing for both dimensions of the target computation

1128 (in the case of MatVec, the  $i$ -dimension and also the  $k$ -dimension – see Figure 1), allowing for  
 1129 potentially fine-grained data movement and parallelization strategies, as we discuss in the following.  
 1130

1131 We compute the partitionings of MDAs by applying the concatenation operator (Example 1)  
 1132 inversely<sup>11</sup> (indicated by using  $=$ : instead of  $:=$  in the top right part of Figure 15). For example,  
 1133 we partition in Figure 15 MDA  $\downarrow a$  first via the inverse of  $+_1^{(HM, x)}$  in dimension 1 (indicated by  
 1134 the subscript 1 of  $+_1^{(HM, x)}$ ; the superscript  $(HM, x)$  is explained later) in 2 parts, as  $p_1^1$  iterates  
 1135 over interval  $[0, 2]_{\mathbb{N}_0} = \{0, 1\}$  which consists of two elements (0 and 1) – the interval is chosen  
 1136 arbitrarily for this example. Afterwards, each of the obtained parts is further partitioned, in the  
 1137 second dimension, via  $+_2^{(HM, y)}$  in 4 parts ( $p_2^1$  iterates over  $[0, 4]_{\mathbb{N}_0} = \{0, 1, 2, 3\}$  which consists of  
 1138 four elements). The  $(2 * 4)$ -many HM parts are then each further partitioned in both dimensions for  
 1139 the COR layer in  $(8 * 16)$  parts, and each individual COR part is again partitioned for the L1 layer in  
 1140  $(32 * 64)$  parts, resulting in  $(2 * 8 * 32) * (4 * 16 * 64) = 512 * 4096$  parts in total.

1141 We always use a *full partitioning* in our low-level expressions<sup>12</sup>, i.e., each particular choice of  
 1142 indices  $p_1^1, p_2^1, p_1^2, p_2^2, p_1^3, p_2^3$  refers to a part that contains a singleton MDA (in Figure 16, the individual  
 1143 MDA elements are denoted via symbol  $\times$ , in the bottom part of the figure). By relying on a full  
 1144 partitioning, we can apply scalar function  $f$  to these singleton parts in the scalar phase (described  
 1145 in the next paragraph). This is because function  $f$  is defined to take as input a single MDA element  
 1146 only (Definition 4), thereby making defining scalar functions more convenient for the user.

1147 The superscript of combine operators, e.g.,  $(COR, x)$  of operator  $+_1^{(COR, x)}$ , is a so-called *operator tag*  
 1148 (formal definition given in the next subsection). A tag indicates to our code generator  
 1149 whether its combine operator is assigned to a memory layer (and thus computed sequentially in  
 1150 our generated code) or to a core layer (and thus computed in parallel). For example, tag  $(COR, x)$   
 1151 indicates that parts processed by operator  $+_1^{(COR, x)}$  should be computed by cores COR, and thus in  
 1152 parallel; the dimension tag  $x$  indicates that COR layer's  $x$  dimension should be used for computing  
 1153 the operator (we use dimension  $x$  for our example architecture as an analogous concept to CUDA's  
 1154 thread/block dimensions  $x, y, z$  for GPU architectures [NVIDIA 2022g]), as we also discuss in the next  
 1155 subsection. In contrast, tag  $(HM, x)$  refers to a memory layer (host memory HM) and thus, operator  
 1156  $+_1^{(HM, x)}$  is computed sequentially. Since the current state-of-practice programming approaches  
 1157 (OpenMP, CUDA, OpenCL, ...) have no explicit notion of memory tiles (e.g., by offering the  
 1158 potential variables  $tileIdx.x/tileIdx.y/tileIdx.z$ , as analogous concepts to CUDA variables  
 1159  $threadIdx.x/threadIdx.y/threadIdx.z$ ), the dimensions tag  $x$  in  $(HM, x)$  is ignored by our code  
 1160 generator (which currently generates OpenMP, CUDA, or OpenCL code), because HM refers to a  
 1161 memory layer.

1162 Note that the number of parts (2 parts on layer 1 in dimension 1; 4 parts on layer 1 in dimen-  
 1163 sion 2; ...), the combine operators' tags, and our partition order (e.g. first partitioning in MDA's  
 1164 dimension 1 and afterwards in dimension 2) are chosen arbitrarily for this example and should be  
 1165 optimized (auto-tuned) for a particular target device and characteristics of the input and output  
 1166 data (size, memory layouts, etc.) to achieve high performance, which we discuss in detail later in  
 1167 this section.

1168 *2. Scalar Phase:* In the scalar phase (step 8 in Figure 15), we apply MDH's scalar function  $f$  to  
 1169 the individual MDA elements which are contained by the MDA's singleton parts

$$\downarrow a_f^{p_1^1, p_2^1 \mid p_1^2, p_2^2 \mid p_1^3, p_2^3}$$

1170 <sup>11</sup>It is easy to see that operator *concatenation* (Example 1) is invertible (formally proved in the Appendix, Section C.4).

1171 <sup>12</sup>Our future work (outlined in Section 8) aims to additionally allow coarser-grained partitioning schemas as well, e.g.,  
 1172 to target domain-specific hardware extensions (such as *NVIDIA Tensor Cores* [NVIDIA 2017] which compute  $4 \times 4$  matrices  
 1173 immediately in hardware, rather than  $1 \times 1$  matrices as obtained in the case of a full partitioning).

1177 for each particular choice of indices  $p_1^1, p_2^1, p_1^2, p_2^2, p_1^3, p_2^3$ , which results in

$$1178 \uparrow a_f^{<p_1^1, p_2^1 \mid p_1^2, p_2^2 \mid p_1^3, p_2^3>}$$

1180 In the figure,  $\vec{f}$  is the slight adaption of function  $f$  that operates on a singleton MDA, rather than a  
 1181 scalar (see Footnote 8).

1182 Annotation  $\rightarrow <(1, 2), \dots>$  indicates the application order of applying scalar function (in  
 1183 this example, first iterating over  $p_1^1$ , then over  $p_2^1$ , etc), and we use annotation  $\rightarrow <(HM, x), \dots>$  to indicate how the scalar computation is assigned to the target architecture (this is described  
 1184 in detail later in this section). Annotations  $\rightarrow M: HM, v: L1$  and  $\rightarrow w: L1$  (in the bottom part  
 1185 of Figure 15) indicate the memory regions to be used for reading and writing the input scalar of  
 1186 function  $f$  (also described later in detail).

1187 *3. Re-Composition Phase:* Finally, the re-composition phase (steps 9-15 in Figure 15) combines  
 1188 the computed parts  $\uparrow a_f^{<p_1^1, p_2^1 \mid p_1^2, p_2^2 \mid p_1^3, p_2^3>}$  (bottom left in the figure) to the final result  $\uparrow a$  (top  
 1189 left) via MDH's combine operators, which are in the case of matrix-vector multiplication  $\otimes_1 :=$   
 1190  $\text{++}$  (concatenation) and  $\otimes_2 := +$  (point-wise addition). In this example, we first combine the L1 parts  
 1191 in dimension 2 and then in dimension 1; afterwards, we combine the COR parts in both dimensions,  
 1192 and finally the HM parts. Analogously to before, this order of combine operators and their tags are  
 1193 also chosen arbitrarily for this example and should be auto-tuned for high performance.

1194 In the de- and re-composition phases, the arrow notation below combine operators allow ef-  
 1195 ficiently exploiting architecture's memory hierarchy, by indicating the memory region to read  
 1196 from (de-composition phase) or to write to (re-composition phase); the annotations also indicate  
 1197 the memory layouts to use. We exploit these memory and layout information in both: i) our code  
 1198 generation to assign combine operators' input and output data to memory regions and to chose  
 1199 memory layouts for the data (row major, column major, etc); ii) our formalism to specify constraints  
 1200 of programming models, e.g., that in CUDA, results of GPU cores can only be combined in des-  
 1201 ignated memory regions [NVIDIA 2022f]. For example, annotation  $\rightarrow M: HM[1, 2], v: L1[1]$  below  
 1202 an operator in the de-composition phase indicates to our code generator that the parts (a.k.a tiles)  
 1203 of matrix  $M$  used for this computation step should be read from host memory HM and that parts of  
 1204 vector  $v$  should be copied to and accessed from fast L1 memory. The annotation also indicates that  $M$   
 1205 should be stored using a row-major memory layout (as we use  $[1, 2]$  and not  $[2, 1]$ ). The memory  
 1206 regions and layouts are chosen arbitrarily for this example and should be chosen as optimized for  
 1207 the particular target architecture and characteristics of the input and output data (we currently rely  
 1208 on auto-tuning [Rasch et al. 2021] for choosing optimized values of performance-critical parameters,  
 1209 as we discuss in Section 5).

1210 Formally, the arrow notation of combination operators is a concise notation to hide MDAs and  
 1211 BUFs for intermediate results, which is discussed in our Appendix, Section C.5, for the interested  
 1212 reader.

## 1213 **Excursion: Code Generation**

1214 Our low-level expressions can be straightforwardly transformed to executable program code  
 1215 in imperative-style programming languages (such as OpenMP, CUDA, and OpenCL). As code  
 1216 generation is not the focus of this work, we demonstrate our code generation approach briefly  
 1217 using the example of Figure 15. Details about our code generation process are provided in Section F  
 1218 of our Appendix, and will be presented and illustrated in our future work.

1219 We implement MDAs via *preprocessor directives*. In the case of MatVec, we implement its input  
 1220 MDA, according to Definition 7, as: `#define inp_mda(i,k) M[i][k],v[k]`

1226 Combine operators are implemented as sequential or parallel loops. For example, the operator  
 1227  $+_1^{(\text{HM}, x)}$  is assigned to memory layer HM and thus implemented as a sequential loop (loop range  
 1228 indicated by  $[0, 2]_{\mathbb{N}_0}$ ), and operator  $+_1^{(\text{COR}, x)}$  is assigned to core layer COR and thus implemented as  
 1229 a parallel loop (e.g., a loop annotated with `#pragma omp parallel` for in OpenMP [OpenMP 2022],  
 1230 or variable `threadIdx.x` in CUDA [NVIDIA 2022g]).

1231 Correspondingly, our three phases (de-composition, scalar, and re-composition) each correspond  
 1232 to an individual loop nest; we generate the nests as fused when the tags of combine operators have  
 1233 the same order in phases, as in Figure 15:  $(\text{HM}, x) \rightarrow (\text{HM}, y) \rightarrow (\text{COR}, x) \rightarrow \dots \rightarrow (\text{L1}, y)$ . Note  
 1234 that our currently targeted programming models (OpenMP, CUDA, and OpenCL) have no explicit notion of *tiles*, e.g., by offering the potential variables `tileIdx.x/tileIdx.y/tileIdx.z` for managing  
 1235 tiles automatically in the programming model (similarly as variables `threadIdx.x/threadIdx.y/threadIdx.z` automatically  
 1236 manage threads in CUDA). Consequently, the dimension information  
 1237 within tags are currently ignored by our code generator when the operator tag refers to a memory  
 1238 layer (such as dimension x in tag  $(\text{HM}, x)$  which refers to memory layer HM).

1239 We implement operators' memory regions as straightforward allocations in the corresponding  
 1240 memory region (e.g., CUDA's device, shared, or register memory [NVIDIA 2022g], according  
 1241 to the arrow annotations in our low-level expression); memory layouts are also implemented via  
 1242 pre-processor directives, e.g., `#define M(i, k) M[k][i]` when storing MatVec's input matrix  $M$  as  
 1243 transposed.

1244 Code optimizations that are applied on a lower abstraction level than proposed by our representation  
 1245 in Example 15 are beyond the scope of this work and outlined in Section G of our Appendix,  
 1246 e.g., loop fusion and loop unrolling which are applied to program code on the low, loop-based  
 1247 abstraction level.

1248 In the following, we introduce in Section 3.2 our formal representation of a computer system  
 1249 (which can be a single device, but also a multi-device or a multi-node system, as we discuss  
 1250 soon), and we illustrate our formal system representation using the example architectures targeted  
 1251 by programming models OpenMP, CUDA, and OpenCL. Afterwards, in Section 3.3, we formally  
 1252 define the basic building blocks of our low-level representation – *low-level MDAs*, *low-level BUFS*,  
 1253 and *low-level combine operators* – based on our formal system representation.

### 1254 3.2 Abstract System Model (ASM)

1255 **Definition 10** (Abstract System Model). An *L-Layered Abstract System Model (ASM)*,  $L \in \mathbb{N}$ , is any  
 1256 pair of two positive natural numbers

$$1257 ( \text{NUM\_MEM\_LYRs}, \text{NUM\_COR\_LYRs} ) \in \mathbb{N} \times \mathbb{N}$$

1258 for which  $\text{NUM\_MEM\_LYRs} + \text{NUM\_COR\_LYRs} = L$ .

1259 Our ASM can model architectures with arbitrarily deep memory and core hierarchies<sup>13</sup>:  $\text{NUM\_MEM\_LYRs}$   
 1260 denotes the target architecture's number of memory layers and  $\text{NUM\_COR\_LYRs}$  the architecture's  
 1261 number of core layers, correspondingly. For example, the artificial architecture we use in Figure 15  
 1262 is represented as an ASM instance as follows (bar symbols denote set cardinality):

$$1263 \text{ASM}_{\text{artif.}} := ( |\{\text{HM}, \text{L1}\}|, |\{\text{COR}\}| ) = (2, 1)$$

1264 The instance is a pair consisting of the numbers 2 and 1 which represent the artificial architecture's  
 1265 two memory layers (HM and L1) and its single core layers (COR).

1266 <sup>13</sup>We deliberately do not model into our ASM representation an architecture's particular number of cores and sizes of  
 1267 memory regions, because our optimization process is designed to be generic in these numbers and sizes for high flexibility.

1275 **Example 11.** We show particular ASM instances that represent the device models of state-of-  
 1276 practice approaches OpenMP, CUDA, and OpenCL:

```
1277   ASMOpenMP      := ( |{MM, L2, L1}| , |{COR}| ) = (3, 1)
1278   ASMOpenMP+L3  := ( |{MM, L3, L2, L1}| , |{COR}| ) = (4, 1)
1279   ASMOpenMP+L3+SIMD := ( |{MM, L3, L2, L1}| , |{COR, SIMD}| ) = (4, 2)
1280
1281   ASMCUDA        := ( |{DM, SM, RM}| , |{SMX, CC}| ) = (3, 2)
1282   ASMCUDA+WRP   := ( |{DM, SM, RM}| , |{SMX, WRP, CC}| ) = (3, 3)
1283
1284   ASMOpenCL       := ( |{GM, LM, PM}| , |{CU, PE}| ) = (3, 2)
1285
1286
1287
```

1288 OpenMP is often used to target  $(3 + 1)$ -layered architectures which rely on 3 memory re-  
 1289 gions (main memory MM, and caches L2 and L1) and 1 core layer (COR). OpenMP-compatible ar-  
 1290 chitectures sometimes also contain the L3 memory region, and they may allow exploiting SIMD  
 1291 parallelization (a.k.a. *vectorization* [Klemm et al. 2012]), which are expressed in our ASM repres-  
 1292 entation as a further memory or core layer, respectively.

1293 CUDA's target architectures are  $(3 + 2)$ -layered: they consist of *Device Memory* (DM), *Shared Mem-  
 1294 ory* (SM), and *Register Memory* (RM), and they offer as cores so-called *Streaming Multiprocessors* (SMX)  
 1295 which themselves consist of *Cuda Cores* (CC). CUDA also has an implicit notion of so-called  
 1296 *Warps* (WRP) which are not explicitly represented in the CUDA programming model [NVIDIA  
 1297 2022g], but often exploited by programmers – via special intrinsics (e.g., *shuffle* and *tensor core  
 1298 intrinsics* [NVIDIA 2017, 2018]) – to achieve highest performance.

1299 OpenCL-compatible architectures are designed analogously to those targeted by the CUDA  
 1300 programming model; consequently, both OpenCL- and CUDA-compatible architectures are rep-  
 1301 resented by the same ASM instance in our formalism. Apart from straightforward syntactical  
 1302 differences between OpenCL and CUDA [StreamHPC 2016], we see as the main differences between  
 1303 the two programming models (from our ASM-based abstraction level) that OpenCL has no notion  
 1304 of warps, and it uses a different terminology – *Global/Local/Private Memory* (GM/LM/PM) instead of  
 1305 device/shared/register memory, and *Compute Unit* (CU) and *Processing Element* (PE), rather than SMX  
 1306 and CC.

1307 In the following, we consider memory regions and cores of ASM-represented architectures as  
 1308 arrangeable in an arbitrary number of dimensions. Programming models for such architectures  
 1309 often have native support for such arrangements. For example, in the CUDA model, memory is  
 1310 accessed via arrays which can be arbitrary-dimensional (a.k.a *multi-dimensional C arrays*), and  
 1311 cores are programmed in CUDA via threads which are arranged in CUDA's so-called dimensions x,  
 1312 y, z; further thread dimensions can be explicitly programmed in CUDA, e.g., by embedding them  
 1313 in the last dimension z.

1315 We express constraints of programming models – for example, that in CUDA, SMX can combine  
 1316 their results in DM only [NVIDIA 2022f] – via so-called *tuning parameter constraints*, which we  
 1317 discuss later in this section.

1318 Note that we call our abstraction *Abstract System Model* (rather than *Abstract Architecture Model*,  
 1319 or the like), because it can also represent systems consisting of multiple devices and/or nodes, etc.  
 1320 For example, our ASM representation of a multi-GPU system is:

```
1322   ASMMulti-GPU := ( |{HM, DM, SM, RM}| , |{GPU, SMX, CC}| ) = (4, 3)
1323
```

1324 It extends our ASM-based representation of CUDA devices (Example 11) by *Host Memory* (HM) which  
 1325 represents the memory region of the system containing the GPUs (and in which the intermediate  
 1326 of different GPUs are combined), and it introduces the further core layer GPU representing the  
 1327 system's GPUs. Analogously, our ASM representation of a multi-node, multi-GPU system is:  
 1328

$$1329 \text{ASM}_{\text{Multi-Node-Multi-GPU}} := ( |\{\text{NM, HM, DM, SM, RM}\}|, |\{\text{NOD, GPU, SMX, CC}\}| ) = (5, 4)$$

1330  
 1331 It adds to ASM  $\text{ASM}_{\text{Multi-GPU}}$  the memory layer (*Node Memory* (NM)) which represents the memory  
 1332 region of the host node, and it adds core layer (*Node* (NOD)) which represents the compute nodes.  
 1333 Our approach is currently designed for *homogeneous systems*, i.e., all devices/nodes/... are assumed  
 1334 to be identical. We aim to extend our approach toward *heterogeneous systems* (which may consist of  
 1335 different devices/nodes/...) as future work, inspired by dynamic load balancing approaches [Chen  
 1336 et al. 2010].  
 1337

### 1338 3.3 Basic Building Blocks

1339 We introduce the three main basic building blocks of our low-level representation: 1) *low-level*  
 1340 MDAs which represent multi-layered, multi-dimensionally arranged collection of ordinary MDAs  
 1341 (Definition 1) – one ordinary MDA per memory/core layer of their target ASM and for each  
 1342 dimension of the MDH computation (illustrated in Figure 16); 2) *low-level* BUFs which are a  
 1343 collection of ordinary BUFs (Definition 5) and that have a notion *memory regions* and *memory*  
 1344 *layouts*; 3) *low-level* *combine operators* which represent combine operators (Definition 2) to which  
 1345 the layer and dimension of their target ASM is assigned that are intended to be used to compute  
 1346 the operator in our generated code (e.g., a core layer to compute the operator in parallel).  
 1347

1348 **Definition 11** (Low-Level MDA). An  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned *low-level* MDA that  
 1349 has scalar type  $T$  and index sets  $I$  is any function  $a_{ll}$  of type:

$$1350 \underbrace{a_{ll}^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}}_{\text{Partitioning: Layer 1}} : \\ 1351 \underbrace{I_1^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \times \dots \times I_D^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>}}_{\text{Partitioning: Layer } L} \rightarrow T$$

1356 We use low-level MDAs in the following to represent partitionings of MDAs (as illustrated soon).  
 1357

1358 Next, we introduce *low-level* BUFs which work similarly as BUFs (Definition 5), but are tagged  
 1359 with a memory region and a memory layout. While these tags have no effect on operators' semantics,  
 1360 they indicate later to our code generator in which memory region the BUF is stored and which  
 1361 memory layout to chose for storing the BUF. Moreover, we use these tags to formally define  
 1362 constraints of programming models, e.g., that according to the CUDA specification [NVIDIA 2022f],  
 1363 SMX cores can combine their results in memory region DM only.  
 1364

1364 **Definition 12** (Low-Level BUF). An  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned *low-level* BUF that  
 1365 has scalar type  $T$  and size  $N$  is any function  $b_{ll}$  of type ( $\leftrightarrow$  denotes bijection):  
 1366

$$1367 \underbrace{b_{ll}^{<\text{MEM} \in [1, \text{NUM\_MEM\_LYRS}]_{\mathbb{N}} | \sigma: [1, D]_{\mathbb{N}} \leftrightarrow [1, D]_{\mathbb{N}} > <(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}}_{\text{Memory Region}} : \\ 1368 \underbrace{[0, N_1^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>})_{\mathbb{N}_0} \times \dots \times [0, N_D^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>})_{\mathbb{N}_0}}_{\text{Memory Layout}} \rightarrow T$$

1369  
 1370  
 1371  
 1372

1373 We refer to MEM as low-level BUF's *memory region* and to  $\sigma$  as its *memory layout*, and we refer to  
1374 the function

1381 that is defined as

**1382**  $b_{II}^{trans <MEM| \sigma > p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L >} (i_{\sigma(1)}, \dots, i_{\sigma(D)}) := b_{II}^{<MEM| \sigma > p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L >} (i_1, \dots, i_D)$   
**1383**

1384  
1385 as  $b_{ll}$ 's *transposed function representation*.

1386 Finally, we introduce local level combining.

Finally, we introduce *layered combine operators*. We define such operators to behave the same as ordinary combine operators (Definition 2), but we additionally tag them with a layer of their target ASM. Similarly as for low-level BUFs, the tag has no effect on semantics, but it is used in our code generation to assign the computation to the hardware (e.g., indicating that the operator is computed by either an SMX, WRP, or CC when targeting CUDA – see Example 11). Also, we use the tags to define model-specific constraints in our formalism (as also discussed for low-level BUFs). We also tag the combine operator with a dimension of the ASM layer, enabling later in our optimization process to express advanced data access patterns (a.k.a. *swizzles* [Phothilimthana et al. 2019]). For example, when targeting CUDA, flexibly mapping ASM dimensions on CC layer (in CUDA terminology, the dimensions are called `threadIdx.x`, `threadIdx.y`, and `threadIdx.z`) to array dimensions enables the well-performing *coalesced global memory accesses* [NVIDIA 2022f] for both transposed and non-transposed data layouts, by only using different dimension tags.

**1399** **Definition 13** (ASM Level). We call pairs  $(l_{\text{ASM}}, d_{\text{ASM}})$ , consisting of an ASM layer  $l_{\text{ASM}} \in [1, L]_{\mathbb{N}}$  and  
**1400** an ASM dimension  $d_{\text{ASM}} \in [1, D]_{\mathbb{N}}$ , an *ASM Level* (ASM-LVL<sup>14</sup>).

**1402** **Definition 14** (Low-Level Combine Operator). The *low-level representation*  $\otimes^{<(l_{ASM}, d_{ASM}) \in \text{ASM-LVL}>}$  of  
**1403** operator  $\otimes$  is a function that for each pair

$$(l_{\text{ASM}}, d_{\text{ASM}}) \in \text{ASM-L VI}$$

has the same type and semantics as  $\otimes$ :

$$\oplus^{<l_{ASM}, d_{ASM}>} \in CO \quad \quad \oplus^{<l_{ASM}, d_{ASM}>} (a, b) := \oplus(a, b)$$

<sup>1410</sup> i.e.,  $\oplus^{<l_{ASM}, d_{ASM}>}$  works exactly as combine operator  $\oplus$ , but it is enriched with a tag that captures an  
<sup>1411</sup> ASM level

1413 Note that in Figure 15, for better readability, we use domain-specific identifiers for ASM layers:  
1414  $\text{HM} := 1$  as an alias for the ASM layer that has id 1,  $\text{L1} := 2$  for layer id 2, and  $\text{COR} := 3$  for layer  
1415 id 3. For dimensions, we use aliases  $x := 1$  for ASM dimension 1 and  $y := 2$  for ASM dimension 2,  
1416 correspondingly.  
1417

<sup>1419</sup> For simplicity, we refrain from annotating identifier ASM-LVL with values  $L$  and  $D$  (e.g.,  $\text{ASM-LVL}^{<L, D>}$  ), because  
<sup>1420</sup> both values will usually be clear from the context.

## 4 LOWERING: FROM HIGH LEVEL TO LOW LEVEL

We have designed our formalism such that an expression in our high-level representation (as in Figure 6) can be systematically lowered to an expression in our low-level representation (as in Figure 15). For this, we parameterize our high-level representation, step-by-step, in tuning parameters; thereby, we obtain for concrete tuning parameter values a particular expression in our low-level representation – this is formally discussed and thoroughly demonstrated in our Appendix, Section D, for the interested reader.

Table 1 lists the tuning parameters of our low-level representation – different values of tuning parameters lead to semantically equal expressions in our low-level representation (which is proven formally in our Appendix, Section D), but the expressions will be translated to differently optimized code variants.<sup>15</sup>

In the following, we explain the 15 tuning parameters in Table 1. We give our explanations in a general, formal setting that is independent of a particular computation and programming model. Dotted lines in Table 1 separate parameters for different phases: parameters D1-D4 customize the de-composition phase, parameters S1-S6 the scalar phase, and parameters R1-R4 the re-composition phase, correspondingly; the parameter  $\theta$  impacts all three phases (separated by a straight line in the table).

Our tuning parameters in Table 1 have constraints: 1) *algorithmic constraints* which have to be satisfied by all target programming models; 2) *model constraints* which are specific for particular programming models only (CUDA-specific constraints, OpenCL-specific constraints, etc), e.g., that the results of CUDA’s thread blocks can be combined in designated memory regions only [NVIDIA 2022f]. We discuss algorithmic constraints in the following, together with our tuning parameters; model constraints are discussed in our Appendix, Section C.3, for the interested reader.

Note that our parameters do not aim to introduce novel optimization techniques, but to unify, generalize, and combine together well-proven optimizations, based on a formal foundation, toward an efficient, overall optimization process that applies to various combinations of data-parallel computations, architectures, and characteristics of input and output data (e.g., their size and memory layout).

**Definition 15** (MDH Level). We refer to pairs  $(l, d)$  – consisting of a particular layer  $l$  and dimension  $d$  – as *MDH Levels* (MDH-LVL):

$$\text{MDH-LVL} := \{ (l, d) \mid l \in [1, L]_{\mathbb{N}}, d \in [1, D]_{\mathbb{N}} \}^{16}$$

We use the pairs to say, for example, that the MDH computation is partitioned on level  $(1, 1)$  (i.e., layer  $l = 1$ , dimension  $d = 1$ ) into two parts, as in Figure 15.

*Parameter  $\theta$ :* Parameter #PRT is a function that maps pairs in MDH-LVL to natural numbers; the parameter determines *how much* data are grouped together into parts in our low-level expression, by setting the particular number of parts (a.k.a. *tiles*) used in our expression (and consequently also in our generated code later). For example, in Figure 15, we use  $\text{#PRT}(1, 1) := 2$  which causes combine operators  $+_1^{(\text{HM}, x)}$  and  $\otimes_1^{(\text{HM}, x)}$  to iterate over interval  $[0, 2]_{\mathbb{N}_0}$  (which partitions the MDH computation on level  $(1, 1)$  into two parts), and we use  $\text{#PRT}(1, 2) := 4$  to let operators  $+_2^{(\text{HM}, y)}$  and  $\otimes_2^{(\text{HM}, x)}$  iterate over interval  $[0, 4]_{\mathbb{N}_0}$  (partitioning in four parts on level  $(1, 2)$ ), etc.

<sup>15</sup>In our Appendix, Section C.2, we show that by choosing particular tuning parameter values, we can express in our formalism the (de/re)-compositions of different, existing state-of-the-art approaches, including scheduling-based approach TVM [Chen et al. 2018a] and polyhedral compilers PPCG [Verdooleaege et al. 2013] and Pluto [Bondhugula et al. 2008b].

<sup>16</sup>The same as for identifier ASM-LVL, we refrain from annotating identifier MDH-LVL with values  $L$  and  $D$ . Note that MDH-LVL and ASM-LVL (Definition 14) both refer to the same set of pairs, but we use identifier MDH-LVL when referring to MDH levels and identifier ASM-LVL when referring to ASM levels, correspondingly, for better clarity.

No.	Name	Range	Description
0	#PRT	MDH-LVL $\rightarrow \mathbb{N}$	number of parts
D1	$\sigma_{\downarrow\text{-ord}}$	MDH-LVL $\leftrightarrow$ MDH-LVL	de-composition order
D2	$\leftrightarrow_{\downarrow\text{-ass}}$	MDH-LVL $\leftrightarrow$ ASM-LVL	ASM assignment (de-composition)
D3	$\downarrow\text{-mem}^{<\text{ib}>}$	MDH-LVL $\rightarrow$ MR	memory regions of input BUFs (ib)
D4	$\sigma_{\downarrow\text{-mem}}^{<\text{ib}>}$	MDH-LVL $\rightarrow [1, \dots, D_{\text{ib}}^{\text{IB}}]_{\mathcal{S}}$	memory layouts of input BUFs (ib)
S1	$\sigma_{f\text{-ord}}$	MDH-LVL $\leftrightarrow$ MDH-LVL	scalar function order
S2	$\leftrightarrow_{f\text{-ass}}$	MDH-LVL $\leftrightarrow$ ASM-LVL	ASM assignment (scalar function)
S3	$f\downarrow\text{-mem}^{<\text{ib}>}$	MR	memory region of input BUF (ib)
S4	$\sigma_{f\downarrow\text{-mem}}^{<\text{ib}>}$	$[1, \dots, D_{\text{ib}}^{\text{IB}}]_{\mathcal{S}}$	memory layout of input BUF (ib)
S5	$f\uparrow\text{-mem}^{<\text{ob}>}$	MR	memory region of output BUF (ob)
S6	$\sigma_{f\uparrow\text{-mem}}^{<\text{ob}>}$	$[1, \dots, D_{\text{ob}}^{\text{OB}}]_{\mathcal{S}}$	memory layout of output BUF (ob)
R1	$\sigma_{\uparrow\text{-ord}}$	MDH-LVL $\leftrightarrow$ MDH-LVL	re-composition order
R2	$\leftrightarrow_{\uparrow\text{-ass}}$	MDH-LVL $\leftrightarrow$ ASM-LVL	ASM assignment (re-composition)
R3	$\uparrow\text{-mem}^{<\text{ob}>}$	MDH-LVL $\rightarrow$ MR	memory regions of output BUFs (ob)
R4	$\sigma_{\uparrow\text{-mem}}^{<\text{ob}>}$	MDH-LVL $\rightarrow [1, \dots, D_{\text{ob}}^{\text{OB}}]_{\mathcal{S}}$	memory layouts of output BUFs (ob)

Table 1. Tuning parameters of our low-level expressions

To ensure a full partitioning (so that we obtain singleton MDAs in the scalar phase to which scalar function  $f$  can be applied), we require the following algorithmic constraint for the parameter:

$$\prod_{l \in [1, L]_{\mathbb{N}}} \#PRT(l, d) = N_d, \text{ for all } d \in [1, D]_{\mathbb{N}}$$

Here,  $N_d$  is the input size in dimension  $d$ .

In our generated code, the number of parts directly translates to the number of *tiles* which are computed either sequentially (a.k.a. *cache blocking* [Lam et al. 1991]) or in parallel in our code, depending on the combine operators's tags (which are chosen via Parameters D2, S2, R2 as discussed soon). In our example from Figure 15, we process the parts belonging to combine operators tagged with HM and L1 sequentially, via for-loops, because HM and L1 correspond to ASM's memory layers (note that Parameter 0 only chooses the number of tiles; the parameter has no effect on explicitly copying data into fast memory resources, which is the purpose of Parameters D3, R3, S1, S2). The

1520 COR parts are computed in parallel in our generated code, because COR corresponds to ASM's core  
 1521 layer, and thus, the number of COR parts determines the number of threads used in our generated  
 1522 code.

1523 An optimized number of tiles is known to be essential for high performance [Bacon et al. 1994],  
 1524 e.g., due to its impact for locality-aware data accesses (number of sequentially computed tiles) and  
 1525 efficiently exploiting parallelism (number of tiles computed in parallel, which corresponds to the  
 1526 number of threads used in our generated code).

1527 *Parameters D1, S1, R1*: These three parameters are permutations on MDH-LVL(indicated by symbol  $\leftrightarrow$  in Table 1), determining *when* data is accessed and combined. The parameters specify the  
 1528 order of combine operators in the de-composition and re-composition phases (parameters D1 and  
 1529 R1), and the order of applying scalar function  $f$  to parts (parameter S1). Thereby, the parameters  
 1530 specify when parts are processed during computation.  
 1531

1532 In our generated code, combine operators are implemented as sequential/parallel loops such that  
 1533 the parameters enable optimizing loop orders (a.k.a. *loop permutation* [McKinley et al. 1996]). For  
 1534 combine operators assigned to ASM's core layer (via parameter R2 discussed in the next paragraph)  
 1535 and thus computed in parallel, parameter R1 particularly determines when the computed results  
 1536 of threads are combined: if we used in the re-composition phase of Figure 15 combine operators  
 1537 tagged with (COR, x) and (COR, y) immediately after applying scalar function  $f$  (i.e., in steps ⑩  
 1538 and ⑪, rather than steps ⑫ and ⑬), we would combine the computed intermediate results of  
 1539 threads multiple times, repeatedly after each individual computation step of threads, and using  
 1540 the two operators at the end of the re-composition phase (in steps ⑭ and ⑮) would combine the  
 1541 result of threads only once, at the end of the re-composition phase.  
 1542

1543 Note that each phase corresponds to an individual loop nest which we fuse together when  
 1544 parameters D1, S1, R1 (as well as parameters D2, S2, R2) coincide (as also outlined in our Appendix,  
 1545 Section G).

1546 *Parameters D2, S2, R2*: These parameters (symbol  $\leftrightarrow$  in the table denotes bijection) assign MDH  
 1547 levels to ASM levels, by setting the tags of low-level combine operators (Definition 14). Thereby, the  
 1548 parameters determine *by whom* data is processed (threads, for-loops, etc), similar to the concept of  
 1549 bind in scheduling languages [Apache TVM Documentation 2022a]. Consequently, the parameters  
 1550 determine which parts should be computed sequentially in our generated code and which parts in  
 1551 parallel. For example, in Figure 15, we use  $\leftrightarrow_{\downarrow\text{-ass}}(2, 1) := (\text{COR}, x)$  and  $\leftrightarrow_{\downarrow\text{-ass}}(2, 2) := (\text{COR}, y)$ ,  
 1552 thereby assigning the computation of MDA parts on layer 2 in both dimensions to ASM's COR  
 1553 layers in the de-composition phase, which causes processing the parts in parallel in our generated  
 1554 code. In CUDA-like approaches which have a notion of cores on multiple layers, the parameter  
 1555 particularly determines the thread layer to be used for parts computed in parallel (e.g., block or  
 1556 thread in CUDA).

1557 Using these parameters, we are able to flexibly set data access patterns in our generated code. In  
 1558 Figure 15, we assign parts on layer 2 to COR layers, which results in a so-called block access pattern  
 1559 of cores: we start  $8 * 16$  threads, according to the  $8 * 16$  core parts, and each thread processes a  
 1560 part of the input MDA representing a block of  $32 * 64$  MDA elements within the input data. If  
 1561 we had assigned in the figure the first computation layer to ASM's COR layer (in the figure, this  
 1562 layer is assigned to ASM's HM layer), we would start  $2 * 4$  threads and each thread would process  
 1563 MDA parts of size  $(8 * 32) * (16 * 64)$ ; assigning the last MDH layer to CORs would result in  
 1564  $(2 * 8 * 32) * (4 * 16 * 64)$  threads each processing MDA parts containing a single MDA element  
 1565 only (a.k.a. *strided access*).

1566 The parameters also enable expressing so-called *swizzle* access patterns [Phothilimthana et al.  
 1567 2019]. For example, in CUDA, processing consecutive data elements in data dimension 1 by threads  
 1568

1569 that are consecutive in thread dimension 2 (a.k.a `threadIdx.y` dimension in CUDA) can achieve  
 1570 higher performance due to the hardware design of fast memory resources in NVIDIA GPUs.  
 1571 Such swizzle patterns can be easily expressed and auto-tuned in our approach; for example, by  
 1572 interchanging in Figure 15 tags `(COR, x)` and `(COR, y)`. For memory layers (such as `HM` and `L1`),  
 1573 the dimension tags `x` and `y` currently have no effect on our generated code, as the programming  
 1574 models we target at the moment (OpenMP, CUDA, and OpenCL) have no explicit notion of tiles.  
 1575 However, this might change in the future when targeting new kinds of programming models, e.g.,  
 1576 for upcoming architectures.

1577 *Parameters D3, R3 and S3, S5:* Parameters D3 and R3 set for each BUF the memory region to  
 1578 use, thereby determining *where* data is read from or written to, respectively. In the table, we  
 1579 use  $ib \in \mathbb{N}$  to refer to a particular input BUF (e.g.,  $ib=1$  to refer to the input matrix of matrix-  
 1580 vector multiplication, and  $ib=2$  to refer to the input vector) and  $ob \in \mathbb{N}$  to refer to an output BUF,  
 1581 correspondingly. Parameter D3 specifies the memory region to read from, and parameter R3 the  
 1582 regions to write to. The set  $MR := [1, \text{NUM\_MEM\_LYRS}]_{\mathbb{N}}$  denotes the ASM’s memory regions.

1583 Similarly to parameters D3 and R3, parameters S3 and S5 set for scalar function  $f$  the memory  
 1584 regions for the input and output scalar.

1585 Exploiting fast memory resources of architectures is a fundamental optimization [Bondhugula  
 1586 2020; Hristea et al. 1997; Mei et al. 2014; Salvador Rohwedder et al. 2023], particularly due to the  
 1587 performance gap between processors’ cores and their memory system [Oliveira et al. 2021; Wilkes  
 1588 2001].

1589 *Parameters D4, R4 and S4, S6:* These parameters set the memory layouts of BUFs, thereby de-  
 1590 termining *how* data is accessed in memory; for brevity in Table 1, we denote the set of all BUF  
 1591 permutations  $[1, D]_{\mathbb{N}} \leftrightarrow [1, D]_{\mathbb{N}}$  (Definition 12) as  $[1, \dots, D]_{\mathcal{S}}$  (symbol  $\mathcal{S}$  is taken from the notation  
 1592 of *symmetric groups* [Sagan 2001]). In the case of our matrix-vector multiplication example in  
 1593 Figure 15, we use a standard memory layout for all matrices, which we express via the parameters  
 1594 by setting them to the identity function, e.g.,  $\sigma_{\text{--mem}}^{<\text{M}>}(1, 1) := id$  (Parameter D4) for the matrix read  
 1595 by operator  $+_1^{(\text{HM}, x)}$ .

1596 An optimized memory layout is important to access data in a locality-aware and thus efficient  
 1597 manner.

## 1600 5 EXPERIMENTAL RESULTS

1601 All experiments described in this section can be reproduced using [TOPLAS Artifact 2022].

1602 We experimentally evaluate our approach by comparing it to popular representatives of four  
 1603 important classes:

- 1604 (1) *scheduling approach:* TVM [Chen et al. 2018a] which generates GPU and CPU code from  
 1605 programs expressed in TVM’s high-level program representation;
- 1606 (2) *polyhedral compilers:* PPCG [Verdoolaege et al. 2013] for GPUs<sup>17</sup> and Pluto [Bondhugula et al.  
 1607 2008b] for CPUs, which automatically generate executable program code in CUDA (PPCG)  
 1608 or OpenMP (Pluto) from straightforward, unoptimized C programs;
- 1609 (3) *functional approach:* Lift [Steuwer et al. 2015] which generates OpenCL code from a Lift-  
 1610 specific, functional program representation;

1611 <sup>17</sup>We cannot compare to polyhedral compiler TC [Vasilache et al. 2019] which is optimized toward deep learning  
 1612 computations on GPUs, because TC is not under active development anymore and thus is not working for newer CUDA  
 1613 architectures [Facebook Research 2022]. Rasch et al. [2019a] show that our approach – already in its proof-of-concept  
 1614 version – achieves higher performance than TC for popular computations on real-world data sets.

1618 (4) *domain-specific libraries*: NVIDIA cuBLAS [NVIDIA 2022b] and NVIDIA cuDNN [NVIDIA  
 1619 2022e], as well as Intel oneMKL [Intel 2022c] and Intel oneDNN [Intel 2022b], which offer the  
 1620 user easy-to-use, domain-specific building blocks for programming. The libraries internally  
 1621 rely on pre-implemented assembly code that is optimized by experts for their target applica-  
 1622 tion domains: linear algebra (cuBLAS and oneMKL) or convolutions (cuDNN and oneDNN),  
 1623 respectively. To make comparison against the libraries challenging for us, we compare to  
 1624 all routines provided by the libraries. For example, the cuBLAS library offers three, seman-  
 1625 tically equal but differently optimized routines for computing MatMul: cublasSgemm (the  
 1626 default MatMul implementation in cuBLAS), cublasGemmEx which is part of the cuBLASEx  
 1627 extension of cuBLAS [NVIDIA 2022c], and the most recent cublasLtMatmul which is part  
 1628 of the cuBLASLt extension [NVIDIA 2022d]; each of these three routines may perform dif-  
 1629 ferently on different problem sizes (NVIDIA usually recommends to naively test which one  
 1630 performs best for the particular target problem). To make comparison further challenging for  
 1631 us, we exhaustively test for each routine all of its so-called cublasGemmAlgo\_t variants, and  
 1632 report the routine's runtime for the best performing variant only. In the case of oneMKL, we  
 1633 compare also to its *JIT engine* [Intel 2019] which is specifically designed and optimized for  
 1634 small problem sizes. We also compare to library *EKR* [Hentschel et al. 2008] which computes  
 1635 data mining example *PRL* (Figure 14) on CPUs – the library is implemented in the Java  
 1636 programming language and parallelized via *Java Threads*, and the library is used in practice  
 1637 by the *Epidemiological Cancer Registry* in North Rhine-Westphalia (Germany) which is the  
 1638 currently largest cancer registry in Europe.

1639 We compare to the approaches experimentally in terms of:  
 1640

- 1641 i) *performance*: via a runtime comparison of our generated code against code that is generated  
 1642 according to the related approaches;
- 1643 ii) *portability*: based on the *Pennycook Metric* [Pennycook et al. 2019] which mathematically  
 1644 defines portability<sup>18</sup> as:

$$\Phi(a, p, H) = \begin{cases} \frac{|H|}{\sum_{i \in H} \frac{1}{e_i(a, p)}} & \text{if } i \text{ is supported, } \forall i \in H \\ 0 & \text{otherwise} \end{cases}$$

1645 In words: "for a given set of platforms  $H$ , the *performance portability*  $\Phi$  of an application  $a$   
 1646 solving problem  $p$  is defined as  $\Phi(a, p, H)$ , where  $e_i(a, p)$  is the performance efficiency (i.e.  
 1647 a ratio of observed performance relative to some proven, achievable level of performance)  
 1648 of application  $a$  solving problem  $p$  on platform  $i$ ; value  $\Phi(a, p, H)$  is 0, if any platform in  $H$   
 1649 is unsupported by  $a$  running  $p$ ." [Pennycook et al. 2019]. Consequently, Pennycook defines  
 1650 portability as a real value in the interval  $[0, 1]_{\mathbb{R}}$  such that a value close to 1 indicates *high*  
 1651 portability and a value close to 0 indicates *low* portability. Here, platforms  $H$  represents a  
 1652 set of devices (CPUs, GPUs, ...), an application  $a$  is in our context a framework (such as  
 1653 TVM, a polyhedral compiler, or our approach), problems  $p$  are our case studies, and  $e_i(a, p)$   
 1654 is computed as the runtime  $a_{p,i}^{\text{best}}$  of the application that achieves the best observed runtime  
 1655 for problem  $p$  on platform  $i$ , divided by the runtime of application  $a$  for problem  $p$  running  
 1656 on platform  $i$ .

- 1657 iii) *productivity*: by intuitively arguing that our approach achieves the same/lower/higher produc-  
 1658 tivity as the related approaches, using the representative example computation *Matrix-Vector*

1659 <sup>18</sup>Pennycook's metric is actually called *Performance Portability (PP)*. Since performance portability particularly includes  
 1660 functional portability, we refer to Pennycook's PP also more generally as *Portability* only.

1667 *Multiplication* (MatVec) (Figure 6) – classical code metrics, like *Lines of Code (LOC)*, CO-  
1668 COMO [Boehm et al. 1995], McCabe’s cyclomatic complexity [McCabe 1976], and Halstead  
1669 development effort [Halstead 1977] are not meaningful for comparing the short and concise  
1670 programs in high-level languages as proposed by the related work as well as our approach.  
1671

1672 In the following, after discussing our application case studies, experimental setup, auto-tuning  
1673 system, and code generator, we compare our approach to each of the four above mentioned classes  
1674 of approaches in Sections 5.1-5.4.

## 1675 **Application Case Studies**

1676 We use for experiments in this section popular example computations from Figure 14 that belong  
1677 to different classes of computations:

- 1679 • Linear Algebra Subroutines (BLAS): *Matrix Multiplication* (MatMul) and *Matrix-Vector Multi-  
1680 plication* (MatVec);
- 1681 • Stencil Computations: *Jacobi Computation* (Jacobi3D) and *Gaussian Convolution* (Conv2D)  
1682 which differ from linear algebra routines by accessing neighboring elements in their input  
1683 data;
- 1684 • Quantum Chemistry: *Coupled Cluster* (CCSD(T)) computations which differ from linear  
1685 algebra routines and stencil computations by accessing their high-dimensional input data in  
1686 complex, transposed fashions;
- 1687 • Data Mining: *Probabilistic Record Linkage* (PRL) which differs from the previous computations  
1688 by relying on a PRL-specific combine operator and scalar function, instead of straightforward  
1689 additions or multiplications as previous computations;
- 1690 • Deep Learning: the most time-intensive computations within the popular neural networks  
1691 ResNet-50 [He et al. 2015], VGG-16 [Simonyan and Zisserman 2014], and MobileNet [Howard  
1692 et al. 2017], according to their TensorFlow implementations [TensorFlow 2022a,b,c]. Deep  
1693 learning computations rely on advanced variants of linear algebra routines and stencil  
1694 computations, e.g., MCC and MCC\_Capsule for computing convolution-like stencils, instead  
1695 of the classical Conv2D variant of convolution (Figure 14) – the deep learning variants are  
1696 considered as significantly more challenging to optimize than their classical variants [Barham  
1697 and Isard 2019].  
1698

1699 We use for experiments this subset of computations from Figure 14 to make experimenting chal-  
1700 lenging for us: state-of-practice approaches, such as TVM and vendor libraries from NVIDIA and  
1701 Intel, are specifically designed and optimized toward these computations. To make experimenting  
1702 further challenging for us, we consider data sizes and characteristics either taken from real-world  
1703 computations (e.g., from the TCCG benchmark suite [Springer and Bientinesi 2016] for quantum  
1704 chemistry computations) or sizes that are preferable for our competitors, e.g., powers of two for  
1705 which many competitors are highly optimized, e.g., vendor libraries. For the deep learning case  
1706 studies, we use data characteristics (sizes, strides, padding strategy, image/filter formats, etc.)  
1707 taken from the particular implementations of the neural networks when computing the popular  
1708 *ImageNet* [Krizhevsky et al. 2012] data set (the particular characteristics are listed in our Appendix,  
1709 Section E.1, for the interested reader). For all experiments, we use single precision floating point  
1710 numbers (a.k.a. float or fp32), as such precision is the default in TensorFlow and many other  
1711 frameworks.

## 1712 **Experimental Setup**

1713 We run our experiments on a cluster containing two different kinds of GPUs and CPUs:

- NVIDIA Ampere GPU A100-PCIE-40GB
- NVIDIA Volta GPU V100-SXM2-16GB
- Intel Xeon Broadwell CPU E5-2683 v4 @ 2.10GHz
- Intel Xeon Skylake CPU Gold-6140 @ 2.30GHz

We represent the two CUDA GPUs in our formalism using model  $\text{ASM}_{\text{CUDA+WRP}}$  (Example 11). We rely on model  $\text{ASM}_{\text{CUDA+WRP}}$ , rather than the CUDA's standard model  $\text{ASM}_{\text{CUDA}}$  (also in Example 11), in order to exploit CUDA's (implicit) warp level for a fair comparison to competitors: warp-level optimizations are exploited by our competitors, e.g., for *shuffle operations* [NVIDIA 2018] which combine the results of threads within a warp with high performance. To fairly compare to TVM and PPCG, we avoid exploiting warps' *tensor core intrinsics* [NVIDIA 2017] in all experiments, which compute the multiplication of small matrices with high performance [Feng et al. 2022], because these intrinsics are not used in the TVM- and PPCG-generated CUDA code. For the two CPUs, we rely on model  $\text{ASM}_{\text{OpenCL}}$  (Example 11) for generating OpenCL code. The same as our approach, TVM also generates OpenCL code for CPUs; Pluto relies on the OpenMP approach for CPUs.

For all experiments, we use the currently newest versions of frameworks, libraries, and compilers, as follows. We compile our generated GPU code using library CUDA NVRTC [NVIDIA 2022h] from CUDA Toolkit 11.4, and we use Intel's OpenCL runtime version 18.1.0.0920 for compiling CPU code. For both compilers, we do not set any flags so that they run in their default modes. For the related approaches, we use the following versions of frameworks, libraries, and compilers:

- TVM [Apache 2022] version 0.8.0 which also uses our system's CUDA Toolkit version 11.4 for GPU computations and Intel's runtime version 18.1.0.0920 for computations on CPU;
- PPCG [Michael Kruse 2022] version 0.08.04 using flag `--target=cuda` for generating CUDA code, rather than OpenCL, as CUDA is usually better performing than OpenCL on NVIDIA GPUs, and we use flag `--sizes` followed by auto-tuned tile sizes – we rely on the *Auto-Tuning Framework (ATF)* [Rasch et al. 2021] for choosing optimized tile size values (as we discuss in the next subsection);
- Pluto [Uday Bondhugula 2022] commit 12e075a using flag `--parallel` for generating OpenMP-parallelized C code (rather than sequential C), as well as flag `--tile` to use ATF-tuned tile sizes for Pluto; the Pluto-generated OpenMP code is compiled via Intel's `icx` compiler version 2022.0.0 using the Pluto-recommended optimization flags `-O3 -fopenmp`;
- NVIDIA cuBLAS [NVIDIA 2022b] from CUDA Toolkit 11.4, using the NVIDIA-recommended compiler flags `-fast -O3 -DNDEBUG`;
- NVIDIA cuDNN [NVIDIA 2022e] from CUDA Toolkit 11.4, using the NVIDIA-recommended compiler flags `-fast -O3 -DNDEBUG`;
- Intel oneMKL [Intel 2022c] compiled with Intel's `icpx` compiler version 2022.0.0, using flags `-DMKL_ILP64 -qmk1=parallel -L${MKLROOT}/lib/intel64 -liomp5 -lpthread -lm -ldl`, as recommended for oneMKL by Intel's *Link Line Advisor* tool [Intel 2022a], as well as standard flags `-O3 -NDEBUG`;
- Intel oneDNN [Intel 2022b] also compiled with Intel's `icpx` compiler version 2022.0.0, using flags `-I${DNNLROOT}/include -L${DNNLROOT}/lib -ldnnl`, according to oneDNN's documentation, as well as standard flags `-O3 -NDEBUG`;
- EKR [Hentschel et al. 2008] executed via Java SE 1.8.0 Update 281.

We profile runtimes of CUDA and OpenCL programs using the corresponding, event-based profiling APIs provided by CUDA and OpenCL. For Pluto which generates OpenMP-annotated C code, we measure runtimes via system call `clock_gettime` [GNU/Linux 2022]. In the case of C++ libraries Intel oneMKL and Intel oneDNN, we use C++'s chrono library [C++ reference 2022].

1765 Libraries NVIDIA cuBLAS and NVIDIA cuDNN are also based on the CUDA programming model; thus,  
1766 we profile them also via CUDA events. To measure the runtimes of the EKR Java library, we use  
1767 `Java function System.currentTimeMillis()`.

1768 All measurements of CUDA and OpenCL programs contain the pure program runtime only (a.k.a.  
1769 *kernel runtime*). The runtime of *host code*<sup>19</sup> is not included in the reported runtimes, as performance  
1770 of host code is not relevant for this work and the same for all approaches.

1771 In all experiments, we collect measurements until the 99% confidence interval was within 5% of  
1772 our reported means, according to the guidelines for *scientific benchmarking of parallel computing*  
1773 *systems* by [Hoefler and Belli \[2015\]](#).

1774

## 1775 Auto-Tuning

1776 The auto-tuning process of our approach relies completely on the generic *Auto Tuning Frame-  
1777 work (ATF)* [[Rasch et al. 2021](#)]. The ATF framework has proven to be efficient for exploring large  
1778 search spaces that are based on constrained tuning parameters (as our space introduced in Section 4).  
1779 We use ATF, out of the box, exactly as described by [Rasch et al. \[2021\]](#): 1) we straightforwardly  
1780 represent in ATF our search space (Table 1) via *tuning parameters* which express the parameters  
1781 in the table and their constraints; 2) we use ATF’s pre-implemented cost functions for CUDA and  
1782 OpenCL to measure the cost of our generated OpenCL and CUDA codes (in this work, we consider  
1783 as cost program’s runtime, rather than its energy consumption, etc); 3) we start the tuning process  
1784 using ATF’s default search technique (*AUC bandit* [[Ansel et al. 2014](#)]). ATF then fully automatically  
1785 determines a well-performing tuning parameter configuration for the particular combination of a  
1786 case study, architecture, and input/output characteristics (size, memory layout, etc).

1787 For scheduling approach TVM, we use its Ansor [[Zheng et al. 2020a](#)] optimization engine which  
1788 is specifically designed and optimized toward generating optimized TVM schedules. Polyhedral  
1789 compilers PPCG and Pluto do not provide own auto-tuning systems; thus, we use for them also the  
1790 ATF framework for auto-tuning, the same as for our approach. For both compilers, we additionally  
1791 also report their runtimes when relying on their internal heuristics, rather than on auto-tuning, to  
1792 fairly compare to them.

1793 To achieve the best possible performance results for TVM, PPCG, and Pluto, we auto-tune each  
1794 of these frameworks individually for each particular combination of case study, architecture, and  
1795 input/output characteristics, the same as for our approach. For example, we start for TVM one tuning  
1796 run when auto-tuning case study MatMul for architecture GPU on one input size, and another, new  
1797 tuning run for a new input size, etc.

1798 Hand-optimized libraries NVIDIA cuBLAS/cuDNN and Intel oneMKL/oneDNN rely on heuristics  
1799 provided by experts, rather than auto-tuning. By relying on heuristics, the libraries avoid the  
1800 time-intensive process of auto-tuning. However, auto-tuning is well amortized in many application  
1801 areas (e.g., deep learning), because the auto-tuned implementations are re-used in many program  
1802 runs. Moreover, auto-tuning avoids the complex and costly process of hand optimization by experts,  
1803 and it often achieves higher performance than hand optimizations, as confirmed later by our  
1804 experiments.

1805 For a fair comparison, we use for each tuning run uniformly the same tuning time of 12h. Even  
1806 though for many computations well-performing tuning results could often also be found in less  
1807 than 12h for our approach as well as for other frameworks, we use such generous tuning time  
1808 for all frameworks to avoid auto-tuning issues in our reported results – analyzing, improving,  
1809 and accelerating the auto-tuning process is beyond the scope of this work and intended for our  
1810

1811 <sup>19</sup>*Host code* is required in approaches CUDA and OpenCL for program execution – it compiles the CUDA and OpenCL  
1812 programs, performs data transfers between host and device, etc.

1813

1814 future work (as also outlined in Section 8). In particular, TVM’s Ansor optimizer was often able  
 1815 to find well performing optimization decisions in 6h of tuning time or less. This is because Ansor  
 1816 explores a small search space that is specifically designed and optimized toward deep learning  
 1817 computations – Ansor’s space is a proper subset of our space, as our space aims at capturing general  
 1818 optimizations that apply to a broad class of data-parallel computations. Thereby, Ansor usually  
 1819 struggles with achieving high performance for computations not taken from the deep learning  
 1820 area, as we confirm in our experiments later.

1821 To improve the auto-tuning efficiency for our implementations, we rely on a straightforward  
 1822 cost model that shrinks our search space in Table 1 before starting our ATF-based auto-tuning pro-  
 1823 cess: i) we always use the same values for Parameters D1, S1, R1 as well as for Parameters D2, S2, R2,  
 1824 thereby generating the same loop structure for all three phases (de-composition, scalar, and re-  
 1825 composition) such that the structures can be generated as a fused loop nest; ii) we restrict Pa-  
 1826 rameters D2, S2, R2 to two values – one value that let threads process outer parts (a.k.a. *blocked*  
 1827 *access* or *outer parallelism*, respectively) and one to let threads process inner parts (*strided access*  
 1828 or *inner parallelism*); all other permutations are currently ignored for simplicity or because they  
 1829 have no effect on the generated code (e.g., permutations of Parameters D2, S2, R2 that only differ in  
 1830 dimension tags belonging to memory layers, as discussed in the previous sections); iii) we restrict  
 1831 Parameters D3, S3, S5, R3 such that each parameter is invariant under different values of  $d$  of its  
 1832 input pairs  $(l, d) \in \text{MDH-LVL}$ , i.e., we always copy full tiles in memory regions (and not a full tile  
 1833 of one input buffer and a half tile of another input buffer, which sometimes might achieve higher  
 1834 performance when memory is a limited resource). Our cost model is straightforward and might  
 1835 filter out configurations of our search space that achieve potentially higher performance than we  
 1836 report for our approach in Sections 5.1-5.4. We aim to substantially improve our naive cost model  
 1837 in future work, based on *operational semantics* for our low-level representation, in order to improve  
 1838 the auto-tuning quality and to reduce (or even avoid) tuning time.

## 1839 **Code Generator**

1840 We provide a straightforward, proof-of-concept code generator, implemented in C++. Our generator  
 1841 takes as input the high-level representation of the target computation (Figure 14), in the form of a  
 1842 straightforward text file (see Appendix, Section A.4), and it fully automatically generates auto-tuned  
 1843 program code, based on the concepts and methodologies introduced and discussed in this paper  
 1844 and the ATF auto-tuning framework.

1845 In our future work, we aim to integrate our code generation approach into the *MLIR* compiler  
 1846 framework [Lattner et al. 2021], building on work-in-progress results [Google SIG MLIR Open  
 1847 Design Meeting 2020], thereby making our work better accessible for the community.

## 1848 **5.1 Scheduling Approaches**

1849 *Performance.* Figures 17-22 report the performance of the TVM-generated code, which is in CUDA  
 1850 for GPUs and in OpenCL for CPUs. We observe that we usually achieve the high performance  
 1851 of TVM and often perform even better. For example, in Figure 21, we achieve a speedup  $> 2\times$   
 1852 over TVM on NVIDIA Ampere GPU for matrix multiplications as used in the inference phase of the  
 1853 ResNet-50 neural network – an actually favorable example for TVM which is designed and optimized  
 1854 toward deep learning computations executed on modern GPUs. Our performance advantage over  
 1855 TVM is because we parallelize and optimize more efficiently reduction-like computations – in the  
 1856 case of MatMul (Figure 14), its 3rd-dimension (a.k.a.  $k$ -dimension). The difficulties of TVM with  
 1857 reduction computations becomes particularly obvious when computing dot products (Dot) on GPUs  
 1858 (Figure 17): the Dot’s main computation part is a reduction computation (via point-wise addition,  
 1859 see Figure 14), thus requiring reduction-focussed optimization, in particular when targeting the  
 1860

1863 highly-parallel architecture of GPUs: in the case of Dot (Figure 17), our MDH-generated CUDA  
1864 code exploits parallelization over CUDA blocks, whereas the Ansor-generated TVM code exploits  
1865 parallelization over threads within in a single block only, because TVM currently cannot use blocks for  
1866 parallelizing reduction computations [Apache TVM Community 2022a]. Furthermore, while TVM's  
1867 Ansor rigidly parallelizes outer dimensions [Zheng et al. 2020a], our ATF-based tuning process has  
1868 auto-tuned our tuning parameters D2, S2, R2 in Table 1 to exploit parallelism for inner dimensions,  
1869 which achieves higher performance for this particular MatMul example used in ResNet-50. Also,  
1870 for MatMul-like computations, Ansor always caches parts of the input in GPU's shared memory,  
1871 and it computes these cached parts always in register memory. In contrast, our caching strategy  
1872 is auto-tunable (via parameters D3, S3 S5, R3 in Table 1), and ATF has determined to not cache  
1873 the input matrices into fast memory resources for the MatMul example in ResNet-50. Surprisingly,  
1874 Ansor does not exploit fast memory resources for Jacobi stencils (Figure 18), as required to achieve  
1875 high performance for them: our MDH-generated and ATF-tuned CUDA kernel for Jacobi uses  
1876 register memory for both inputs (image buffer and filter) when targeting NVIDIA Ampere GPU (small  
1877 input size), thereby achieving a speedup over TVM+Ansor of 1.93 $\times$  for Jacobi. Most likely, Ansor  
1878 fails to foresee the potential of exploiting fast memory resources for Jacobi stencils, because the  
1879 Jacobi's index functions used for memory accesses (Figure 14) are injective. For the MatMul example  
1880 of ResNet-50's training phase (Figure 21), we achieve a speedup over TVM on NVIDIA Ampere GPU  
1881 of 1.26 $\times$ , because auto-tuning determined to store parts of input matrix A as transposed into fast  
1882 memory (via parameter D4 in Table 1). Storing parts of the input/output data as transposed is not  
1883 considered by Ansor as optimization, perhaps because such optimization must be expressed in  
1884 TVM's high-level language, rather than scheduling language [Apache TVM Community 2022c]. For  
1885 MatVec on NVIDIA Ampere GPU (Figure 17), we achieve a speedup over TVM of 1.22 $\times$  for the small  
1886 input size, by exploiting a so-called *swizzle pattern* [Phothilimthana et al. 2019]: our ATF tuner  
1887 has determined to assign threads that are consecutive in CUDA's x-dimension to the second MDA  
1888 dimension (via parameters D2, S2, R2 in Table 1), thereby accessing the input matrix in a GPU-  
1889 efficient manner (a.k.a *coalesced global memory accesses* [NVIDIA 2022f]). In contrast, for MatVec  
1890 computations, Ansor assigns threads with consecutive x-ids always to the first data dimension, in a  
1891 non-tunable manner, causing lower performance.

1892 Our positive speedups over TVM on CPU are for the same reasons as discussed above for GPU. For  
1893 example, we achieve a speedup of  $> 3\times$  over TVM on Intel Skylake CPU for MCC (Figure 22) as used  
1894 in the training phase of the MobileNet neural network, because we exploit fast memory resources  
1895 more efficiently than TVM: our auto-tuning process has determined to use register memory for the  
1896 MCC's second input (the filter buffer F, see Table 14) and using no fast memory for the first input  
1897 (image buffer I), whereas Ansor uses shared memory rigidly for both inputs of MCC. Moreover,  
1898 our auto-tuning process has determined to parallelize the inner dimensions of MCC, while Ansor  
1899 always parallelizes outer dimensions. We achieve the best speedup over TVM for MCC on an input size  
1900 taken from TVM's own tutorials [Apache TVM Documentation 2022b] (Figure 18), rather than from  
1901 neural networks (as in Figures 21 and 22). This is because TVM's MCC size includes large reduction  
1902 computations, which are not efficiently optimized by TVM (as discussed above).

1903 The TVM compiler achieves higher performance than our approach for some examples in Figures  
1904 17-22. However, in most cases, this has a technical reason only: TVM uses the NVCC compiler for  
1905 compiling CUDA code, while our proof-of-concept code generator relies on NVIDIA's NVRTC library  
1906 which surprisingly generates less efficient CUDA assembly than NVCC. In three cases, the higher  
1907 performance of TVM over our approach is because our ATF auto-tuning framework was not able  
1908 to find a better performing tuning configuration than TVM's Ansor optimization engine; the three  
1909 cases are: 1) MCC (capsule variant) from ResNet-50's training phase on Intel Skylake CPU, 2) MCC  
1910 from VGG-16's inference phase on NVIDIA Ampere GPU, and 3) MCC (capsule variant) from VGG-16's  
1911

1912 training phase on NVIDIA Ampere GPU. However, when we manually set the Ansor-found tuning  
 1913 configuration also for our approach (analogously as done in Section C.2), instead of using the  
 1914 ATF-found configuration, we achieve for these three cases exactly the same high performance as  
 1915 TVM+Ansor, i.e., the well-performing configurations are contained in our MDH-based search space.  
 1916 Most likely, Ansor was able to find this well-performing configuration automatically, because it  
 1917 explores a smaller search space that is particularly designed for deep learning computations. To  
 1918 avoid such tuning issues in our approach, we aim to substantially improve our auto-tuning process  
 1919 in future work: we plan to introduce an analytical cost model that assists (or even replaces) our  
 1920 auto-tuner, as we also outline in Section 8.

1921 Note that the TVM compiler crashes for our data mining example PRL, because TVM has difficulties  
 1922 with computations relying on user-defined combine operators [Apache TVM Community 2022d].

1923 *Portability.* Figure 23 reports the portability of the TVM compiler over our GPU and CPU ar-  
 1924 chitectures. The portability measurements are based on the Pennycook metric where a value  
 1925 close to 1 indicates high portability and a value close to 0 low portability, correspondingly. We  
 1926 observe that except for the example of transposed matrix multiplication  $\text{GEMM}^\top$ , we always achieve  
 1927 higher portability than TVM. The higher portability of TVM for  $\text{GEMM}^\top$  is because TVM achieves for  
 1928 this example higher performance than our approach on NVIDIA Volta GPU. However, the higher  
 1929 performance of TVM is only due to the fact that TVM uses NVIDIA’s NVCC for compiling CUDA code,  
 1930 while we currently rely on NVIDIA’s NVRTC library which surprisingly generates less efficient  
 1931 CUDA assembly, as discussed above.

1932 *Productivity.* Listing 1 shows how matrix-vector multiplication (MatVec) is implemented in TVM’s  
 1933 high-level program representation which is embedded into the Python programming language.  
 1934 In line 1, the input size  $(I, K) \in \mathbb{N} \times \mathbb{N}$  of matrix  $M \in T^{I \times K}$  (line 2) and vector  $v \in T^K$  (line 3) are  
 1935 declared, in the form of function parameters; the matrix and vector are named  $M$  and  $v$  and both are  
 1936 assumed to contain elements of scalar type  $T = \text{float32}$  (floating point numbers). Line 5 defines  
 1937 a so-called *reduction axis* in TVM in which all values are combined in line 8 via `te.sum` (addition).  
 1938 The basic computation part of MatVec – multiplying matrix element  $M[i, k]$  with vector element  
 1939  $v[k]$  – is also specified in line 8.

```
1941
1942 1 def MatVec(I, K):
1943 2     M = te.placeholder((I, K), name='M', dtype='float32')
1944 3     v = te.placeholder((K,), name='v', dtype='float32')
1945 4
1946 5     k = te.reduce_axis((0, K), name='k')
1947 6     w = te.compute(
1948 7         (I,),
1949 8         lambda i: te.sum(M[i, k] * v[k], axis=k)
1950 9     )
1951 10    return [M, v, w]
```

1951 Listing 1. TVM program expressing Matrix-Vector Multiplication (MatVec)

1952
 1953
 1954 While we see the MatVec implementations of TVM (Listing 1) and our approach (Figure 6) basi-  
 1955 cally on the same level of abstraction, we consider our approach in general as more expressive.  
 1956 This is because our approach supports multiple reduction dimensions that may rely on different  
 1957 combine operators, e.g., as required for expressing the MBBS example in Figure 14 with which TVM  
 1958 is struggling – adding support for multiple, different reduction dimensions is considered in the  
 1959 TVM community as a non-trivial extension of TVM [Apache TVM Community 2020, 2022b]. Also,  
 1960

1961 we consider our approach as slightly less error-prone: we automatically compute the expected  
1962 sizes of matrix  $M$  (as  $I \times K$ ) and vector  $v$  (as  $K$ ), based on the user-defined input size  $(I, K)$  in line 1  
1963 and index functions  $(i, k) \mapsto (i, k)$  for the matrix and  $(i, k) \mapsto (k)$  for the vector in line 8 (see  
1964 Definition 7). In contrast, TVM redundantly requests these matrix and vector sizes from the user  
1965 for computing the function specification of its generated MatVec code (once in lines 2 and 3 of  
1966 Listing 1, and again in lines 5 and 7), which lets TVM generate incorrect low-level code – without  
1967 issuing an error message – when the user sets sizes different from  $I \times K$  for the matrix and  $K$  for  
1968 the vector in lines 2 and 3 [Apache TVM Community 2022f].

## 1969 5.2 Polyhedral Compilers

1970 *Performance.* Figures 17-22 report for our application case studies the performance of the PPCG-  
1971 generated CUDA code for GPUs and of the OpenMP-annotated C code generated by polyhedral  
1972 compiler Pluto for CPUs. For a fair comparison, we report for both polyhedral compilers their  
1973 performance achieved for ATF-tuned tile sizes (denoted as PPCG+ATF/Pluto+ATF in Figures 17-22),  
1974 as well as the performance of the two compilers when relying on their internal heuristics instead  
1975 of auto-tuning (denoted as PPCG and Pluto in the figures). In some cases, PPCG’s heuristic crashed  
1976 with error too many resources requested for launch, because the heuristic seems to not take  
1977 into account device-specific constraints, e.g., limited availability of GPUs’ fast memory resources.  
1978

1979 We observe that our MDH-based approach achieves better performance than PPCG and Pluto in  
1980 all cases – sometimes by multiple orders of magnitude – in particular for deep learning computations  
1981 (Figures 21 and 22). This is caused by the rigid optimization goals of PPCG and Pluto, e.g., always  
1982 parallelizing outer dimensions, which causes severe performance losses. For example, we achieve a  
1983 speedup over PPCG of  $> 13\times$  on NVIDIA Ampere GPU and of  $> 60\times$  over Pluto on Intel Skylake CPU  
1984 for MCC as used in the inference phase of the real-world ResNet-50 neural network. Compared to  
1985 PPCG, our better performance for this MCC example is because PPCG has difficulties with efficiently  
1986 parallelizing computations relying on more than 3 dimension; most likely, this is because CUDA  
1987 offers per default 3 dimensions for parallelization (called x, y, z dimension in CUDA). However,  
1988 MCC relies on 7 parallelizable dimensions (as shown in Figure 14); exploiting the parallelization  
1989 opportunities of the 4 further dimensions (as done in our generated CUDA code) is essential to  
1990 achieve high performance for this MCC example from ResNet-50. Our performance advantage over  
1991 Pluto for the MCC example is because Pluto parallelizes the outer dimensions of MCC only; however,  
1992 the dimension has a size of only 1 for this real-world example, resulting in starting only 1 thread in  
1993 the Pluto-generated OpenMP code.

1994 For dot products Dot (Figure 17), we can observe that PPCG fails to generate parallel CUDA code,  
1995 because PPCG cannot parallelize and optimize computations which rely solely on combine operators  
1996 different from concatenation, as we also discuss in Section 6.2. In Section 6.2, we also discuss that  
1997 we do not consider the performance issues of PPCG and Pluto as weaknesses of the polyhedral  
1998 approach in general, but of the particular polyhedral transformations used in PPCG and Pluto.

1999 Note that Pluto crashes for our data mining examples (Figure 20), with Error extracting  
2000 polyhedra from source file, which is due to the complex scalar function of the example,  
2001 which involves if-statements. Moreover, Intel’s icx compiler struggles with compiling the Pluto-  
2002 generated OpenMP code for quantum chemistry computations (Figure 19): we aborted icx’s  
2003 compilation process after 24h compilation time. The icx’s issues with the Pluto-generated code  
2004 are most likely because of too aggressive loop unrolling of Pluto – the Pluto-generated OpenMP  
2005 code has often a size  $> 50\text{MB}$  for our real-world quantum chemistry examples.

2006  
2007 *Portability.* Since PPCG and Pluto are each designed for particular architectures only, they achieve  
2008 the lowest portability of 0 for all our studies, according to the Pennycook metric. To simplify for  
2009

2010 PPCG and Pluto the portability comparison with our approach, we compute the Pennycook metric  
 2011 additionally also for two restricted sets of devices: only GPUs to make comparison against our  
 2012 approach easier for PPCG, and only CPUs for Pluto.

2013 Figures 24-28 report the portability of PPCG when considering only GPUs, as well as the portability  
 2014 of Pluto for only CPUs. We observe that we achieve higher portability for all our studies, as we  
 2015 constantly achieve higher performance than the two polyhedral compilers for the studies.

2016 Note that even when restricting our set of devices to only GPUs for PPCG or only CPUs for Pluto,  
 2017 the two polyhedral compilers still achieve a portability of 0 for some examples, because they fail to  
 2018 generate code for them (as discussed above).

2019 *Productivity.* Listing 2 shows the input program of polyhedral compilers PPCG and Pluto for  
 2020 MatVec. Both take as input easy-to-implement, straightforward, sequential C code. We consider  
 2021 these two polyhedral compilers as more productive than our approach (as well as scheduling/func-  
 2022 tional approaches and polyhedral compilers relying on a DSL, rather than sequential programs, as  
 2023 TC [Vasilache et al. 2019]), because both compilers fully automatically generate optimized parallel  
 2024 code from unoptimized, sequential programs.

2025 Rasch et al. [2020a,b] show that our approach can achieve the same, high user productivity as  
 2026 polyhedral compilers, by using a polyhedral frontend for our approach: we can alternatively take  
 2027 as input the same sequential user programs as PPCG and Pluto, instead of programs implemented  
 2028 in our high-level program representation (as in Figure 6). The sequential input program is then  
 2029 transformed via polyhedral tool *pet* [Verdoolaege and Grosser 2012] to its polyhedral representation  
 2030 which is then automatically transformed to our high-level program representation, according to  
 2031 the methodology presented by Rasch et al. [2020a,b].

```
2033 1 for( int i = 0 ; i < M ; ++i )
2034 2   for( int k = 0 ; k < K ; ++k )
2035 3     w[i] += M[i][k] * v[k];
```

2036 Listing 2. PPCG/Pluto program expressing Matrix-Vector Multiplication (MatVec)

### 2039 5.3 Functional Approaches

2041 Our previous work [Rasch et al. 2019a] already shows that while functional approaches provide  
 2042 a solid formal foundation for computations, they typically have performance and portability  
 2043 issues; for this, our previous work used the state-of-the-art Lift [Steuwer et al. 2015] framework  
 2044 as running example (which, to the best of our knowledge, has so far not been improved toward  
 2045 higher performance and/or better portability). Therefore, we refrain from a further performance and  
 2046 portability evaluation of Lift and focus in the following on analyzing and discussing the productivity  
 2047 potentials of functional approaches, using again the state-of-the-art Lift approach as running  
 2048 example. We discuss the performance and portability issues of functional approaches from a general  
 2049 perspective thoroughly in Section 6.3.

2050 *Performance/Portability.* Already experimentally evaluated in previous work [Rasch et al. 2019a]  
 2051 and discussed in general terms in Section 6.3.

2053 *Productivity.* Listing 3 shows the implementation of MatVec in the Lift approach. In line 1, type  
 2054 parameters *n* and *m* are declared via Lift building block *nFun* – type parameters are limited to  
 2055 natural numbers in the Lift formalism. Line 2 declares a function in Lift that takes as input a  
 2056 matrix of size *m* × *n* and a vector of size *n*, correspondingly; the matrix and vector are both assumed  
 2057 to consist of float numbers (floating point numbers). The computation of MatVec is specified in  
 2058

2059 lines 3 and 4: the map function of Lift in line 3 iterates over all rows of the matrix, each row is  
 2060 pair-wise combined with the input vector via Lift pattern zip, multiplication \* is applied to each  
 2061 pair in the combined vector via Lift's map pattern in line 4, and the obtained products are finally  
 2062 combined via addition + using Lift's reduce pattern.

2063

```
2064 1 nFun(n => nFun(m =>
2065 2   fun(matrix: [[float]]n)m => fun(xs: [float]n =>
2066 3     matrix :>> map(fun(row =>
2067 4       zip(xs, row) :>> map(*) :>> reduce(+, 0)
2068 5     )))) )
```

2069 Listing 3. Lift program expressing Matrix-Vector Multiplication (MatVec)

2070

2071 Already for expressing MatVec, we can observe that Lift relies on a vast set of small, functional  
 2072 building blocks (five building blocks for MatVec: nFun, fun, map, zip, and reduce), and the blocks  
 2073 have to be composed and nested in complex ways for expressing computations. Consequently,  
 2074 we consider programming in Lift-like approaches as complex and their productivity for the user  
 2075 as limited. Moreover, Lift-like approaches often need fundamental extension for targeting new  
 2076 kinds of computations, e.g., so-called *macro-rules* which had to be added to Lift for efficiently  
 2077 targeting matrix multiplications [Remmeli et al. 2016] and primitives slide and pad together with  
 2078 optimization *overlapped tiling* for expressing stencil computations [Hagedorn et al. 2018]. This  
 2079 need for extensions limits the expressiveness of the Lift language and thus hinders productivity.

2080 In contrast to Lift, our approach relies on exactly three higher-order functions (Figure 5)  
 2081 to express various kinds of data-parallel computations (Figure 14): 1) *inp\_view* (Definition 7)  
 2082 which prepares the input data; our *inp\_view* function is designed as general enough to sub-  
 2083 sume – in a structured manner – the subset of all Lift patterns intended to change the view on  
 2084 input data, including patterns zip, pad, and slide; 2) *md\_hom* (Definition 3) expresses the actual  
 2085 computation and subsumes all Lift patterns performing actual computations (fun, map, reduce,  
 2086 ...); 3) *out\_view* (Definition 9) expresses the view on output data and is designed to work similarly  
 2087 as function *inp\_view* (Lemma 2). Our three functions are always composed straightforwardly in  
 2088 the same, fixed order (Figure 5), and they do not rely on complex function nesting for expressing  
 2089 computations.

2090 Note that even though our language is designed as minimalist, it should cover the expressivity  
 2091 of the Lift language<sup>20</sup> and beyond: for example, we are currently not aware of any Lift program  
 2092 being able to express the prefix-sum examples in Figure 14. For the above reasons, we consider  
 2093 programming in our high-level language as more productive for the user than programming in  
 2094 Lift-like functional languages. Furthermore, as discussed in Section 5.2, our approach can take as  
 2095 input also straightforward, sequential program code, which further contributes to the productivity  
 2096 of our approach.

2097

2098

2099

2100

2101

2102

2103

2104

<sup>20</sup>This work is focussed on dense computations. Lift supports sparse computations [Pizzuti et al. 2020] which we consider as future work (as also outlined in Section 8). We consider Lift's approach, based on their so-called *position dependent arrays*, as a great inspiration for our future goal.

2105

2106

2107

2108 2109 2110 2111 2112 2113 2114 2115 2116 2117 2118	Linear Algebra	NVIDIA Ampere GPU							
		Dot		MatVec		MatMul		MatMult	bMatMul
		$2^{24}$	$10^7$	4096, 4096	8192, 8192	10,500, 64	1024, 1024, 1024	10,500, 64	16,10,500, 64
TVM+Ansor		172.48	128.22	1.74	1.23	1.00	1.00	1.00	1.17
PPCG		-	-	5.44	2.95	2.20	2.73	3.40	162.92
PPCG+ATF		-	-	4.22	2.77	1.20	1.87	1.32	3.06
cuBLAS		1.10	1.11	1.14	1.01	1.40	0.92	1.60	1.50
cuBLASEx		-	-	-	-	1.20	0.91	1.60	1.33
cuBLASL <sub>t</sub>		-	-	-	-	1.20	0.88	1.60	-
2119									
2119 2120 2121 2122 2123 2124 2125 2126 2127 2128 2129	Linear Algebra	NVIDIA Volta GPU							
		Dot		MatVec		MatMul		MatMult	bMatMul
		$2^{24}$	$10^7$	4096, 4096	8192, 8192	10,500, 64	1024, 1024, 1024	10,500, 64	16,10,500, 64
TVM+Ansor		82.28	67.97	1.06	1.04	1.00	1.08	0.80	1.00
PPCG		-	-	2.67	1.71	1.40	3.07	2.60	111.98
PPCG+ATF		-	-	2.44	2.24	1.00	2.16	1.20	2.83
cuBLAS		1.06	1.09	1.10	1.07	2.60	1.11	1.80	1.83
cuBLASEx		-	-	-	-	1.80	0.30	1.40	1.17
cuBLASL <sub>t</sub>		-	-	-	-	1.20	0.96	1.40	-
2130									
2130 2131 2132 2133 2134 2135 2136 2137 2138 2139 2140	Linear Algebra	Intel Skylake CPU							
		Dot		MatVec		MatMul		MatMult	bMatMul
		$2^{24}$	$10^7$	4096, 4096	8192, 8192	10,500, 64	1024, 1024, 1024	10,500, 64	16,10,500, 64
TVM+Ansor		5.07	6.14	1.03	3.39	1.06	1.15	1.02	1.10
Pluto		5.40	6.48	2.49	6.24	3.21	12.25	5.45	14.30
Pluto+ATF		5.39	6.01	1.43	3.38	2.98	4.78	4.79	2.14
oneMKL		0.64	0.57	0.42	3.83	6.27	0.69	3.42	0.98
oneMKL(JIT)		-	-	-	-	0.65	-	1.13	-
2141									
2141 2142 2143 2144 2145 2146 2147 2148 2149 2150	Linear Algebra	Intel Broadwell CPU							
		Dot		MatVec		MatMul		MatMult	bMatMul
		$2^{24}$	$10^7$	4096, 4096	8192, 8192	10,500, 64	1024, 1024, 1024	10,500, 64	16,10,500, 64
TVM+Ansor		5.60	8.46	1.21	1.63	1.20	1.11	1.11	1.00
Pluto		4.78	6.73	3.01	1.28	4.89	5.26	6.74	11.97
Pluto+ATF		4.75	6.72	2.91	1.21	1.94	2.85	3.46	1.23
oneMKL		1.03	0.41	0.57	0.59	2.00	0.66	1.98	0.84
oneMKL(JIT)		-	-	-	-	1.03	-	1.30	-

2151 Fig. 17. Speedup (higher is better) of our approach for linear algebra routines on GPUs and CPUs over:  
2152 i) scheduling approach TVM, ii) polyhedral compilers PPCG (GPU) and Pluto (CPU), as well as iii) hand-  
2153 optimized libraries provided by vendors. Dash symbol "-" means this framework does not support this  
2154 combination of architecture, computation, and data characteristic.

NVIDIA Ampere GPU					
Stencils	Jacobi3D		Conv2D		MCC
	256,256,256	512,512,512	224,224,5,5	4096,4096,5,5	1,512,7,7,512,3,3
TVM+Ansor	1.93	2.04	1.00	2.32	1.63
PPCG	4.19	5.27	1.58	2.36	-
PPCG+ATF	1.08	1.02	1.22	1.38	9.37
cuDNN	-	-	2.20	5.29	2.44
NVIDIA Volta GPU					
Stencils	Jacobi3D		Conv2D		MCC
	256,256,256	512,512,512	224,224,5,5	4096,4096,5,5	1,512,7,7,512,3,3
TVM+Ansor	2.05	1.86	1.00	2.00	1.50
PPCG	7.01	13.87	1.45	1.75	-
PPCG+ATF	1.03	1.00	1.23	1.34	8.28
cuDNN	-	-	2.60	3.58	4.42
Intel Skylake CPU					
Stencils	Jacobi3D		Conv2D		MCC
	256,256,256	512,512,512	224,224,5,5	4096,4096,5,5	1,512,7,7,512,3,3
TVM+Ansor	2.30	1.64	1.59	2.46	2.76
Pluto	3.65	2.66	2.39	1.38	143.80
Pluto+ATF	1.81	1.38	2.09	1.06	61.47
oneDNN	3.92	2.60	6.47	2.83	3.91
Intel Broadwell CPU					
Stencils	Jacobi3D		Conv2D		MCC
	256,256,256	512,512,512	224,224,5,5	4096,4096,5,5	1,512,7,7,512,3,3
TVM+Ansor	2.21	1.78	3.14	3.98	3.99
Pluto	2.10	1.67	2.29	2.17	74.48
Pluto+ATF	1.29	1.05	1.74	1.25	74.47
oneDNN	16.09	15.02	7.29	16.42	7.69

Fig. 18. Speedup (higher is better) of our approach for stencil computations on GPUs and CPUs over: i) scheduling approach TVM, ii) polyhedral compilers PPCG (GPU) and Pluto (CPU), as well as iii) hand-optimized libraries provided by vendors. Dash symbol "-" means this framework does not support this combination of architecture, computation, and data characteristic.

NVIDIA Ampere GPU								
	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
TVM+Ansor	1.15	1.07	1.25	1.00	1.36	1.05	1.00	1.15
PPCG	10585.85	10579.40	9819.81	11211.57	10181.14	10482.81	11693.21	10585.85
PPCG+ATF	11.19	15.60	14.06	11.45	11.81	12.06	11.72	11.19

NVIDIA Volta GPU								
	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
TVM+Ansor	1.09	0.93	1.04	1.03	1.01	1.11	1.01	1.09
PPCG	6466.22	6019.64	6300.31	6468.40	6608.80	5256.49	6602.22	6466.22
PPCG+ATF	8.28	9.61	9.38	7.21	6.60	5.14	7.77	8.28

Intel Skylake CPU								
	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
TVM+Ansor	1.60	1.50	2.06	1.70	1.20	2.12	1.56	1.60
Pluto	147.45	151.55	206.60	162.58	157.43	145.17	321.66	147.45
Pluto+ATF	1.89	2.01	1.89	1.80	1.82	1.92	1.84	1.89

Intel Broadwell CPU								
	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
TVM+Ansor	1.06	1.28	1.16	1.15	1.29	1.13	2.07	1.06
Pluto	-	-	-	-	-	-	-	-
Pluto+ATF	-	-	-	-	-	-	-	-

Fig. 19. Speedup (higher is better) of our approach for quantum chemistry computations Coupled Cluster (CCSD(T)) on GPUs and CPUs over: i) scheduling approach TVM, and ii) polyhedral compilers PPCG (GPU) and Pluto (CPU). Dash symbol "-" means this framework does not support this combination of architecture, computation, and data characteristic.

Data Mining		NVIDIA Ampere GPU					
		2 <sup>15</sup>	2 <sup>16</sup>	2 <sup>17</sup>	2 <sup>18</sup>	2 <sup>19</sup>	2 <sup>20</sup>
TVM+Ansor	-	-	-	-	-	-	-
PPCG	1.49	1.05	1.12	1.22	1.37	1.56	
PPCG+ATF	1.40	1.22	1.50	1.63	1.83	2.12	

Data Mining		NVIDIA Volta GPU					
		2 <sup>15</sup>	2 <sup>16</sup>	2 <sup>17</sup>	2 <sup>18</sup>	2 <sup>19</sup>	2 <sup>20</sup>
TVM+Ansor	-	-	-	-	-	-	-
PPCG	1.11	1.15	1.10	1.30	1.51	1.82	
PPCG+ATF	1.26	1.37	1.47	1.77	2.07	2.48	

Data Mining		Intel Skylake CPU					
		2 <sup>15</sup>	2 <sup>16</sup>	2 <sup>17</sup>	2 <sup>18</sup>	2 <sup>19</sup>	2 <sup>20</sup>
TVM+Ansor	-	-	-	-	-	-	-
Pluto	-	-	-	-	-	-	-
Pluto+ATF	-	-	-	-	-	-	-
EKR	6.18	5.39	9.62	19.87	26.42	24.78	

Data Mining		Intel Broadwell CPU					
		2 <sup>15</sup>	2 <sup>16</sup>	2 <sup>17</sup>	2 <sup>18</sup>	2 <sup>19</sup>	2 <sup>20</sup>
TVM+Ansor	-	-	-	-	-	-	-
Pluto	-	-	-	-	-	-	-
Pluto+ATF	-	-	-	-	-	-	-
EKR	8.01	9.17	23.58	66.90	119.33	167.19	

2294 Fig. 20. Speedup (higher is better) of our approach for data mining algorithm Probabilistic Record Linkage  
2295 (PRL) on GPUs and CPUs over: i) scheduling approach TVM, and ii) polyhedral compilers PPCG (GPU) and  
2296 Pluto (CPU), as well as the iii) hand-implemented Java CPU implementation used by EKR – the largest  
2297 cancer registry in Europa. Dash symbol “-” means this framework does not support this combination of  
2298 architecture, computation, and data characteristic.

2304	2305	2306	2307	2308	2309	2310	2311	2312	2313	2314	NVIDIA Ampere GPU													
											ResNet-50				VGG-16				MobileNet					
											Training		Inference		Training		Inference		Training		Inference			
MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	0.93	1.42	0.88	1.14	0.94	1.00	0.94	1.00	MCC	MCC	0.94	1.00		
TVM+Ansor	1.00	1.26	1.05	2.22	—	7.89	1661.14	7.06	5.77	5.08	2254.67	—	7.55	—	—	—	—	—	—	—	—	—		
PPCG	3456.16	8.26	—	—	3.28	13.76	5.44	4.26	3.92	9.46	3.73	3.31	10.71	—	—	—	—	—	—	—	—	—		
PPCG+ATF	—	—	—	—	cuDNN	0.92	—	1.85	—	1.22	—	1.94	—	1.81	2.14	—	—	—	—	—	—	—		
cuBLAS	—	1.58	—	2.67	—	—	—	—	0.93	—	—	1.04	—	—	—	—	—	—	—	—	—	—		
cuBLASEx	—	1.47	—	2.56	—	—	—	—	0.92	—	—	1.02	—	—	—	—	—	—	—	—	—	—		
cuBLASLt	—	—	1.26	—	—	1.22	—	—	0.91	—	—	1.01	—	—	—	—	—	—	—	—	—	—		
2315														NVIDIA Volta GPU										
2316	2317	2318	2319	2320	2321	2322	2323	2324	2325	2326	ResNet-50				VGG-16				MobileNet					
											Training	Inference	Training	Inference	Training	Inference	Training	Inference	Training	Inference	Training	Inference		
MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	0.75	1.21	0.72	1.79	1.00	1.11	1.06	1.00	1.00	1.00	1.00	1.00		
TVM+Ansor	0.75	1.21	0.72	1.79	PPCG	1976.38	5.88	—	5.64	994.16	3.41	8.21	2.51	1411.92	—	7.26	—	—	—	—	—	—	—	
PPCG	1976.38	5.88	—	5.64	PPCG+ATF	3.43	3.54	3.42	4.93	3.85	3.15	8.13	2.05	3.49	3.56	—	—	—	—	—	—	—	—	
PPCG+ATF	3.43	3.54	3.42	4.93	cuDNN	1.21	—	1.29	—	2.80	—	3.50	—	2.32	3.14	—	—	—	—	—	—	—	—	
cuBLAS	—	1.33	—	1.14	cuBLASEx	—	1.21	—	1.07	—	1.04	—	1.03	—	—	—	—	—	—	—	—	—	—	
cuBLASEx	—	1.21	—	1.07	cuBLASLt	—	1.00	—	1.07	—	1.04	—	1.02	—	—	—	—	—	—	—	—	—	—	
2327														NVIDIA Ampere GPU										
2328	2329	2330	2331	2332	2333	2334	2335	2336	2337	2338	ResNet-50				VGG-16				MobileNet					
											Training	Inference	Training	Inference	Training	Inference	Training	Inference	Training	Inference	Training	Inference		
MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	TVM+Ansor	0.96	1.00	0.79	1.02	0.88	—	0.99	—	—	—	—	—	—	—	—	—	—	—	
PPCG	4642.24	—	—	—	PPCG+ATF	25.98	85.33	4.41	13.64	8.89	—	4017.74	—	—	—	—	—	—	—	—	—	—	—	
PPCG+ATF	25.98	85.33	4.41	13.64	cuDNN	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	
2339														NVIDIA Volta GPU										
2340	2341	2342	2343	2344	2345	2346	2347	2348	2349	2350	ResNet-50				VGG-16				MobileNet					
											Training	Inference	Training	Inference	Training	Inference	Training	Inference	Training	Inference	Training	Inference		
MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	TVM+Ansor	0.95	1.01	1.05	0.97	1.04	—	0.87	—	—	—	—	—	—	—	—	—	—	—	—
PPCG	2935.40	—	—	945.16	PPCG+ATF	19.24	19.68	8.28	12.29	8.84	—	2885.90	—	—	—	—	—	—	—	—	—	—	—	
PPCG+ATF	19.24	19.68	8.28	12.29	cuDNN	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	

Fig. 21. Speedup (higher is better) of our approach for the most time-intensive computations used in deep learning neural networks ResNet-50, VGG-16, and MobileNet on GPUs over: i) scheduling approach TVM, ii) polyhedral compilers PPCG (GPU), as well as iii) hand-optimized libraries provided by vendors. Dash symbol “—” means this framework does not support this combination of architecture, computation, and data characteristic.

2353	2354	2355	Deep Learning	Intel Skylake CPU							
				ResNet-50				VGG-16			
				Training		Inference		Training		Inference	
				MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MCC
2356	TVM+Ansor			1.53	1.05	1.14	1.20	1.97	1.14	2.38	1.27
2357	Pluto			355.81	49.57	364.43	13.93	130.80	93.21	186.25	36.30
2358	Pluto+ATF			13.08	19.70	170.69	6.57	3.11	6.29	53.61	8.29
2359	oneDNN			0.39	—	5.07	—	1.22	—	9.01	—
2360	oneMKL			—	0.44	—	1.09	—	0.88	—	0.53
2361	oneMKL(JIT)			—	6.43	—	8.33	—	27.09	—	9.78
2362											
2363											
2364											
2365											
2366	2367	2368	Deep Learning	Intel Broadwell CPU							
				ResNet-50				VGG-16			
				Training	Inference	Training	Inference	Training	Inference	Training	Inference
				MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MCC
2369	TVM+Ansor			1.53	1.60	1.29	1.53	1.32	1.00	1.27	1.02
2370	Pluto			4349.20	40.41	137.21	15.96	1865.07	53.57	113.40	24.10
2371	Pluto+ATF			6.43	8.93	61.60	6.91	5.07	4.38	42.63	4.45
2372	oneDNN			1.30	—	1.81	—	2.94	—	2.85	—
2373	oneMKL			—	1.45	—	1.36	—	1.35	—	0.50
2374	oneMKL(JIT)			—	19.78	—	9.77	—	50.58	—	10.70
2375											
2376											
2377											
2378											
2379	2380	2381	Deep Learning (Capsule)	Intel Skylake CPU							
				ResNet-50				VGG-16			
				Training	Inference	Training	Inference	Training	Inference	Training	Inference
				MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule
2382	TVM+Ansor			0.94	1.14	3.50	1.18	2.94	1.59		
2383	Pluto			209.36	265.77	—	166.45	160.49	159.34		
2384	Pluto+ATF			14.33	265.77	3.33	60.66	4.40	57.21		
2385	oneDNN			—	—	—	—	—	—	—	—
2386											
2387											
2388	2389	2390	Deep Learning (Capsule)	Intel Broadwell CPU							
				ResNet-50				VGG-16			
				Training	Inference	Training	Inference	Training	Inference	Training	Inference
				MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule
2391	TVM+Ansor			2.61	1.30	3.55	1.00	1.32	2.24		
2392	Pluto			—	—	—	—	—	—	—	—
2393	Pluto+ATF			4418.82	56.17	75.77	2173.72	202.34	158.52		
2394	oneDNN			—	—	—	—	—	—	—	—
2395											
2396											

2397 Fig. 22. Speedup (higher is better) of our approach for the most time-intensive computations used in deep  
 2398 learning neural networks ResNet-50, VGG-16, and MobileNet on CPUs over: i) scheduling approach TVM,  
 2399 ii) polyhedral compilers Pluto (CPU), as well as iii) hand-optimized libraries provided by vendors. Dash  
 2400 symbol “—” means this framework does not support this combination of architecture, computation, and data  
 2401 characteristic.

Pennycook Metric								
Linear Algebra		Dot		MatVec		MatMul		MatMult
		2 <sup>24</sup>	10 <sup>7</sup>	4096, 4096	8192, 8192	10, 500, 64	1024, 1024, 1024	10, 500, 64
MDH+ATF		0.88	0.64	0.65	0.85	0.88	0.54	0.94
TVM+Ansor		0.01	0.02	0.54	0.47	0.83	0.50	0.97
								0.89
Pennycook Metric								
Stencils		Jacobi3D		Conv2D		MCC		
		256, 256, 256	512, 512, 512	224, 224, 5, 5	4096, 4096, 5, 5	1, 512, 7, 7, 512, 3, 3		
MDH+ATF		1.00	1.00	1.00	1.00	1.00	1.00	1.00
TVM+Ansor		0.47	0.55	0.59	0.37	0.41		
Pennycook Metric								
Quantum Chemistry		abcdefg-gdab-efgc	abcdefg-gdac-efgb	abcdefg-gdbc-efga	abcdefg-geab-dfgc	abcdefg-geac-dfgb	abcdefg-gebc-dfga	abcdefg-gfab-degc
		1.00	0.98	1.00	1.00	1.00	1.00	1.00
MDH+ATF		0.82	0.82	0.73	0.82	0.82	0.74	0.71
TVM+Ansor		0.82	0.82	0.73	0.82	0.82	0.74	0.84
Pennycook Metric								
Data Mining		2 <sup>15</sup>	2 <sup>16</sup>	2 <sup>17</sup>	2 <sup>18</sup>	2 <sup>19</sup>	2 <sup>20</sup>	
		1.00	1.00	1.00	1.00	1.00	1.00	
MDH+ATF		0.00	0.00	0.00	0.00	0.00	0.00	
TVM+Ansor		0.00	0.00	0.00	0.00	0.00	0.00	
Pennycook Metric								
Deep Learning		ResNet-50			VGG-16			MobileNet
		Training	Inference		Training	Inference		Training Inference
		MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC MCC
MDH+ATF		0.67	0.76	0.91	1.00	0.98	0.95	0.97 0.68
TVM+Ansor		0.53	0.62	0.89	0.59	0.76	0.81	0.70 0.61
								0.98 1.00
Pennycook Metric								
Deep Learning (Capsule)		ResNet-50		VGG-16		MobileNet		
		Training	Inference	Training	Inference	Training	Inference	
		MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule
MDH+ATF		0.96	1.00	0.94	0.99	0.97	0.96	
TVM+Ansor		0.71	0.90	0.44	0.95	0.63	0.69	

Fig. 23. Portability (higher is better), according to Pennycook metric, of our approach and TVM over GPUs and CPUs for case studies. Polyhedral compilers PPCG/Pluto and vendor libraries by NVIDIA and Intel are not listed: due to their limitation to certain architectures, all of them achieve the lowest portability of 0 only.

2451	2452	2453	2454	2455	2456	2457	2458	2459	2460	Pennycook Metric (GPUs only)								
										Dot		MatVec		MatMul		MatMul <sup>T</sup>		bMatMul
										2 <sup>24</sup>	10 <sup>7</sup>	4096, 4096	8192, 8192	10, 500, 64	1024, 1024, 1024	10, 500, 64	16, 10, 500, 64	1.00
MDH+ATF	1.00	1.00	1.00	1.00	1.00	0.45	0.89	1.00	1.00	0.78	0.48	0.44	0.32	0.79	0.67	1.00	0.90	
TVM+Ansor	0.01	0.01	0.71	0.88	1.00	0.42	1.00	0.42	1.00	0.15	0.30	0.30	0.00	0.91	0.30	0.01	0.92	
PPCG	0.00	0.00	0.25	0.43	0.56	0.15	0.30	0.21	0.71	0.40	0.40	0.30	0.30	0.91	0.34	0.01	0.34	
PPCG+ATF	0.00	0.00	0.30	0.40	0.91	0.21	0.71	0.21	0.71	0.93	0.91	0.90	0.00	0.89	0.52	0.60	0.60	
cuBLAS	0.93	0.91	0.89	0.96	0.50	0.42	0.52	0.42	0.60	0.00	0.00	0.00	0.00	0.50	0.52	0.60	0.60	
cuBLASEx	0.00	0.00	0.00	0.00	0.67	0.98	0.60	0.60	0.00	0.00	0.00	0.00	0.00	0.98	0.00	0.00	0.00	
cuBLASL <sub>t</sub>	0.00	0.00	0.00	0.00	0.83	0.48	0.60	0.60	0.00	0.00	0.00	0.00	0.00	0.48	0.60	0.60	0.60	
2461	2462	2463	2464	2465	2466	2467	2468	2469	2470	Pennycook Metric (CPUs only)								
										Dot		MatVec		MatMul		MatMul <sup>T</sup>	bMatMul	
										2 <sup>24</sup>	10 <sup>7</sup>	4096, 4096	8192, 8192	10, 500, 64	1024, 1024, 1024	10, 500, 64	16, 10, 500, 64	1.00
MDH+ATF	0.78	0.48	0.48	0.74	0.79	0.67	1.00	0.67	1.00	0.78	0.48	0.44	0.32	0.71	0.60	0.94	0.86	
TVM+Ansor	0.15	0.06	0.44	0.24	0.20	0.08	0.16	0.08	0.07	0.15	0.07	0.23	0.37	0.31	0.18	0.24	0.55	
Pluto	0.15	0.07	0.18	0.24	0.20	0.08	0.07	0.08	0.07	0.15	0.07	0.23	0.41	0.17	1.00	0.37	1.00	
Pluto+ATF	0.15	0.07	0.23	0.37	0.31	0.18	0.24	0.24	0.24	0.15	0.00	0.00	0.00	0.98	0.00	0.83	0.00	
oneMKL	0.99	1.00	1.00	0.41	0.17	1.00	0.37	1.00	1.00	0.99	0.00	0.00	0.00	0.98	0.00	0.83	0.00	
oneMKL(JIT)	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	

Fig. 24. Portability (higher is better), according to Pennycook metric, for linear algebra routines computed on only GPUs or CPUs, respectively. The restriction simplifies for frameworks with limited architectural support (such as polyhedral compilers and vendor libraries) the portability comparisons against our approach.

2471	2472	2473	2474	2475	2476	2477	2478	2479	2480	2481	2482	2483	2484	Pennycook Metric (GPUs only)					
														Jacobi3D		Conv2D		MCC	
														256, 256, 256	512, 512, 512	224, 224, 5, 5	4096, 4096	1, 512, 7, 7, 512, 3, 3	
MDH+ATF	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
TVM+Ansor	0.50	0.51	1.00	0.51	0.46	0.64	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	
PPCG	0.18	0.10	0.66	0.10	0.49	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
PPCG+ATF	0.95	0.99	0.82	0.99	0.74	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	0.11	
cuDNN	0.00	0.00	0.42	0.00	0.23	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	0.29	
2485	2486	2487	2488	2489	2490	2491	2492	2493	2494	2495	2496	2497	2498	2499	Pennycook Metric (CPUs only)				
															Jacobi3D		Conv2D		MCC
															256, 256, 256	512, 512, 512	224, 224, 5, 5	4096, 4096	1, 512, 7, 7, 512, 3, 3
MDH+ATF	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
TVM+Ansor	0.44	0.58	0.42	0.42	0.31	0.30	0.30	0.30	0.30	0.30	0.30	0.30	0.30	0.30	0.30	0.30	0.30	0.30	
Pluto	0.35	0.46	0.43	0.43	0.56	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	
Pluto+ATF	0.65	0.83	0.52	0.52	0.86	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	
oneDNN	0.10	0.11	0.15	0.15	0.10	0.17	0.17	0.17	0.17	0.17	0.17	0.17	0.17	0.17	0.17	0.17	0.17	0.17	

Fig. 25. Portability (higher is better), according to Pennycook metric, for stencil computations computed on only GPUs or CPUs, respectively. The restriction simplifies for frameworks with limited architectural support (such as polyhedral compilers and vendor libraries) the portability comparisons against our approach.

Pennycook Metric (GPUs only)									
	Quantum Chemistry	abcdefg-gdab-efgc	abcdefg-gdac-efgb	abcdefg-gdbc-efga	abcdefg-geab-dfgc	abcdefg-geac-dfgb	abcdefg-gebc-dfga	abcdefg-gfab-degc	abcdefg-gfbc-dega
2500	MDH+ATF	1.00	0.96	1.00	1.00	1.00	1.00	1.00	1.00
2501	TVM+Ansor	0.90	0.96	0.87	0.99	0.84	0.93	0.99	1.00
2502	PPCG	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
2503	PPCG+ATF	0.10	0.08	0.09	0.11	0.11	0.12	0.10	0.15
2504	Pennycook Metric (CPUs only)								
2505	Quantum Chemistry	abcdefg-gdab-efgc	abcdefg-gdac-efgb	abcdefg-gdbc-efga	abcdefg-geab-dfgc	abcdefg-geac-dfgb	abcdefg-gebc-dfga	abcdefg-gfab-degc	abcdefg-gfbc-dega
2506	MDH+ATF	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
2507	TVM+Ansor	0.75	0.72	0.62	0.70	0.80	0.62	0.55	0.72
2508	Pluto	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
2509	Pluto+ATF	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

2507 Fig. 26. Portability (higher is better), according to Pennycook metric, for quantum chemistry computation  
 2508 Coupled Cluster (CCSD(T)) computed on only GPUs or CPUs, respectively. The restriction simplifies for  
 2509 frameworks with limited architectural support (such as polyhedral compilers and vendor libraries) the  
 2510 portability comparisons against our approach.

Pennycook Metric (GPUs only)							
	Data Mining	2 <sup>15</sup>	2 <sup>16</sup>	2 <sup>17</sup>	2 <sup>18</sup>	2 <sup>19</sup>	
2519	MDH+ATF	1.00	1.00	1.00	1.00	1.00	
2520	TVM+Ansor	0.00	0.00	0.00	0.00	0.00	
2521	PPCG	0.77	0.91	0.90	0.80	0.69	
2522	PPCG+ATF	0.75	0.77	0.67	0.59	0.51	
2523	Pennycook Metric (CPUs only)						
	Data Mining	2 <sup>15</sup>	2 <sup>16</sup>	2 <sup>17</sup>	2 <sup>18</sup>	2 <sup>19</sup>	
2524	MDH+ATF	1.00	1.00	1.00	1.00	1.00	
2525	TVM+Ansor	0.00	0.00	0.00	0.00	0.00	
2526	Pluto	0.00	0.00	0.00	0.00	0.00	
2527	Pluto+ATF	0.00	0.00	0.00	0.00	0.00	
2528	EKR	0.14	0.14	0.06	0.02	0.01	
2529	EKR						

2529 Fig. 27. Portability (higher is better), according to Pennycook metric, for data mining algorithm Probabilistic  
 2530 Record Linkage (PRL) computed on only GPUs or CPUs, respectively. The restriction simplifies for  
 2531 frameworks with limited architectural support (such as polyhedral compilers and vendor libraries) the portability  
 2532 comparisons against our approach.

Deep Learning		Pennycook Metric (GPUs only)									
		ResNet-50				VGG-16				MobileNet	
		Training		Inference		Training		Inference		Training	Inference
MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MCC
MDH+ATF	0.82	1.00	0.84	1.00	0.96	0.95	0.94	1.00	0.97	1.00	
TVM+Ansor	0.96	0.81	0.98	0.50	1.00	0.75	0.97	0.93	1.00	1.00	
PPCG	0.00	0.14	0.00	0.15	0.00	0.18	0.14	0.26	0.00	0.13	
PPCG+ATF	0.24	0.33	0.11	0.19	0.24	0.27	0.11	0.35	0.28	0.14	
cuBLAS	0.76	0.00	0.55	0.00	0.48	0.00	0.35	0.00	0.47	0.38	
cuBLASEx	0.00	0.69	0.00	0.53	0.00	0.95	0.00	0.96	0.00	0.00	
cuBLASLt	0.00	0.75	0.00	0.55	0.00	0.97	0.00	0.97	0.00	0.00	
cuDNN	0.00	0.88	0.00	0.87	0.00	0.98	0.00	0.98	0.00	0.00	
2561											
Deep Learning		Pennycook Metric (CPUs only)									
		ResNet-50				VGG-16				MobileNet	
MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MCC
MDH+ATF	0.56	0.61	1.00	1.00	1.00	0.94	1.00	0.51	1.00	1.00	
TVM+Ansor	0.37	0.50	0.82	0.73	0.61	0.87	0.55	0.45	0.37	0.60	
Pluto	0.00	0.01	0.00	0.07	0.00	0.01	0.01	0.02	0.00	0.02	
Pluto+ATF	0.05	0.04	0.01	0.15	0.24	0.17	0.02	0.08	0.20	0.04	
oneMKL	0.87	0.00	0.29	0.00	0.48	0.00	0.17	0.00	0.69	0.23	
oneMKL(JIT)	0.00	0.82	0.00	0.82	0.00	0.85	0.00	1.00	0.00	0.00	
oneDNN	0.00	0.06	0.00	0.11	0.00	0.02	0.00	0.05	0.00	0.00	
2573											
Deep Learning (Capsule)		Pennycook Metric (GPUs only)									
		ResNet-50		VGG-16		MobileNet					
MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule
MDH+ATF	0.95	1.00		0.88		0.98		0.94		0.93	
TVM+Ansor	1.00		0.99	0.98		0.99		0.98		1.00	
PPCG	0.00		0.00	0.00		0.00		0.00		0.00	
PPCG+ATF	0.04		0.02	0.14		0.08		0.11		0.07	
cuDNN	0.00		0.00	0.00		0.00		0.00		0.00	
2583											
Deep Learning (Capsule)		Pennycook Metric (CPUs only)									
		ResNet-50		VGG-16		MobileNet					
MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule
MDH+ATF	0.97	1.00		1.00		1.00		1.00		1.00	
TVM+Ansor	0.55		0.82	0.28		0.92		0.47		0.52	
Pluto	0.00		0.00	0.00		0.00		0.00		0.00	
Pluto+ATF	0.00		0.01	0.03		0.00		0.01		0.01	
oneDNN	0.00		0.00	0.00		0.00		0.00		0.00	

Fig. 28. Portability (higher is better), according to Pennycook metric, for deep learning computations computed on only GPUs or CPUs, respectively. The restriction simplifies for frameworks with limited architectural support (such as polyhedral compilers and vendor libraries) the portability comparisons against our approach.

2598 **5.4 Domain-Specific Approaches**

2599 *Performance.* Figures 17–22 report for completeness also performance results achieved by domain-  
 2600 specific approaches. Since domain-specific approaches are specifically designed and optimized  
 2601 toward particular applications domains and often also architectures (e.g., only linear algebra routines  
 2602 on only GPU), we consider comparing to them as challenging for us: our approach is designed and  
 2603 optimized toward data-parallel computations in general, from arbitrary application domains (the  
 2604 same as also TVM, polyhedral compilers, and many functional approaches), and our approach is  
 2605 also flexible in the target parallel architecture.

2606 We observe in Figures 17–22 that the domain-specific libraries NVIDIA cuBLAS/cuDNN (for lin-  
 2607 ear algebra routines and convolutions on GPUs) and Intel oneMKL/oneDNN (for linear algebra  
 2608 routines and convolutions on CPUs) sometimes perform better and sometimes worse than our  
 2609 approach. The better performance of libraries over our approach is most likely<sup>21</sup> because the  
 2610 libraries internally rely on assembly-level optimizations, while we currently focus on the higher  
 2611 CUDA/OpenCL abstraction level which offers less optimization opportunities [Goto and Geijn  
 2612 2008; Lai and Seznec 2013]. The cuBLASEx extension of cuBLAS achieves in one case – MatMul  
 2613 on NVIDIA Volta GPU for square  $1024 \times 1024$  input matrices – significantly higher performance  
 2614 than our approach. The high performance is achieved by cuBLASEx when using its algorithm  
 2615 variant CUBLAS\_GEMM\_ALG01\_TENSOR\_OP which casts the float-typed inputs implicitly to the  
 2616 half precision type (a.k.a. half or fp16), allowing cuBLASEx exploiting the GPU’s tensor core  
 2617 extension [NVIDIA 2017]. Thereby, cuBLASEx achieves significantly higher performance than our  
 2618 approach, because tensor cores compute small matrix multiplication immediately in hardware; how-  
 2619 ever, at the cost of a significant precision loss: the half scalar type achieves only half the accuracy  
 2620 achieved by scalar type float. When using cuBLASEx’s default algorithm CUBLAS\_GEMM\_DEFAULT  
 2621 (rather than algorithm CUBLAS\_GEMM\_ALG01\_TENSOR\_OP), which retains the float type and thus  
 2622 meets the accuracy expected from the computation, we achieve a speedup of  $1.11 \times$  over cuBLASEx.  
 2623 For the interested reader, we report in our Appendix, Section E.2, the runtime of cuBLASEx for all  
 2624 its algorithm variants, including reports for the accuracy achieved by the different variants.

2625 The reason for the better performance of our approach over NVIDIA and Intel libraries is most  
 2626 likely because our approach allows generating code that is also optimized (auto-tuned) for data  
 2627 characteristics, which is important for high performance [Tillet and Cox 2017]. In contrast, the  
 2628 vendor libraries usually rely for each computation on pre-implemented implementations each  
 2629 optimized toward only average high performance for a range of data characteristics (size, memory  
 2630 layout, etc). By relying on these fixed, pre-implemented implementations, the libraries avoid the  
 2631 auto-tuning overhead. However, auto-tuning is often amortized, particularly for deep learning  
 2632 computations – the main target of libraries NVIDIA cuDNN and Intel oneDNN – because the auto-  
 2633 tuned implementations are re-used in many program runs. Moreover, we achieve better performance  
 2634 for convolutions (Figure 18), because the libraries re-use optimizations for these computations  
 2635 originally intended for linear algebra routines [Li et al. 2016], while our optimization space (Table 1)  
 2636 is designed as general and not oriented toward linear algebra.

2637 Compared to the EKR library (Figure 20), we achieve higher performance, because the EKR’s Java  
 2638 implementation inefficiently handles memory: the library is implemented using Java’s ArrayList  
 2639 data structure which is convenient to use for the Java programmer, but inefficient in terms of  
 2640 performance, because the structure internally performs costly memory reallocations.

2641 *Portability.* Similarly to polyhedral compilers PPCG and Pluto, the domain-specific approaches  
 2642 work for particular architectures only. The domain-specific approaches are also restricted to a

2643 <sup>21</sup>Since the Intel and NVIDIA libraries are not open source, we cannot explain their performance behavior with certainty.

2647 narrow set of studies, e.g., only linear algebra routines as NVIDIA cuBLAS and Intel oneMKL or  
 2648 only data mining example PRL as EKR. Consequently, the approaches achieve for these unsupported  
 2649 studies the lowest portability of only 0 in Figure 23 as well as Figures 24-28. For their target  
 2650 studies, domain-specific approaches can achieve high portability. This is because the approaches are  
 2651 specifically designed and optimized toward these studies, e.g., via application-specific assembly-level  
 2652 optimizations which are currently beyond the scope of our work.

2653  
 2654 *Productivity.* Listing 4 shows the implementation of MatVec in domain-specific approach NVIDIA  
 2655 cuBLAS; the implementation of MatVec in other domain-specific approaches, e.g., Intel oneMKL,  
 2656 is analogous to the implementation in Listing 4. We consider domain-specific approaches as  
 2657 most productive for their target domain: in the case of MatVec, the user simply calls the high-level  
 2658 function `cublasSgemv` and passes to it the input matrices (omitted via ellipsis in the listing) together  
 2659 with some meta information (memory layout of matrices, etc); cuBLAS then automatically starts  
 2660 the GPU computation for MatVec.

2661 Besides the fact that domain-specific approaches typically target only particular target architec-  
 2662 tures, a further fundamental productivity issue of domain-specific approaches is that they can only  
 2663 be used for a narrow class of computations, e.g., only linear algebra routines as NVIDIA cuBLAS and  
 2664 Intel oneMKL. Moreover, in the case of domain-specific libraries from NVIDIA and Intel, it is often  
 2665 up to the user to manually choose among different, semantically equal but differently performing  
 2666 implementations for high performance. For example, the cuBLAS library offers three different  
 2667 routines for computing matrix multiplications – routines `cublasSgemm` (part of standard cuBLAS),  
 2668 `cublasGemmEx` (part of the cuBLASEx extension of cuBLAS), and routine `cublasLtMatmul` (part of  
 2669 the cuBLASLt extension) – and the routines often also offer different, so-called *algorithms* (e.g., 42  
 2670 algorithm variants in the case cuBLASEx) which impact the internal optimization process. When  
 2671 striving for the highest performance potentials of libraries, the user is in charge of naively testing  
 2672 each possible combination of routine and algorithm variant (as we have done in Figures 17-22  
 2673 to make experimenting challenging for us). In addition, the user must be aware that different  
 2674 combinations of routines and algorithms can produce results of reduced accuracy (as discussed  
 2675 above), which can be critical for accuracy-sensitive use cases.

2676  
 2677 1 `cublasSgemv( /* ... */ );`

2678  
 2679 Listing 4. cuBLAS program expressing Matrix-Vector Multiplication (MatVec)

2680  
 2681  
 2682 **6 RELATED WORK**  
 2683  
 2684 Three major classes of approaches currently focus on code generation and optimization for data-  
 2685 parallel computations: 1) scheduling, 2) polyhedral, and 3) functional. In the following, we compare  
 2686 in Sections 6.1-6.3 our approach to each of these three classes – in terms of *performance*, *portability*,  
 2687 and *productivity*. In contrast to Section 5, which has compared our approach against these classes  
 2688 experimentally, this section is focussed on discussions in a more general, non-experimental context.  
 2689 Afterwards, we outline domain-specific approaches in Section 6.4, which are specifically designed  
 2690 and optimized toward their target application domains. In Section 6.5, we briefly outline approaches  
 2691 focussing on optimizations at the algorithmic level of abstraction which we consider as higher-level  
 2692 optimizations than proposed by our approach and as greatly combinable with our work. Finally,  
 2693 we discuss in Section 6.6 the differences between our approach introduced in this paper and the  
 2694 already existing work on MDHs.

## 2696 6.1 Scheduling Approaches

2697 Popular examples of scheduling approaches include TVM [Chen et al. 2018a], Halide [Ragan-Kelley  
 2698 et al. 2013], Elevate [Hagedorn et al. 2020b], DaCe [Ben-Nun et al. 2019], Tiramisu [Baghdadi et al.  
 2699 2019], CHill [Chen et al. 2008; Khan et al. 2013], Clay [Bagnères et al. 2016], UTF [Kelly and Pugh  
 2700 1998], URUK [Girbal et al. 2006], Fireiron [Hagedorn et al. 2020a], DISTAL [Yadav et al. 2022], and  
 2701 LoopStack [Wasti et al. 2022]. While scheduling approaches usually achieve high performance, they  
 2702 often have difficulties with achieving portability and productivity, as we discuss in the following.<sup>22</sup>  
 2703

2704 *Performance.* Scheduling approaches usually achieve high performance. For this, the approaches  
 2705 incorporate human expert knowledge into their optimization process which is based on two  
 2706 major steps: 1) a human expert implements an optimization program (a.k.a *schedule*) in a so-called  
 2707 *scheduling language* – the program specifies the basic optimizations to perform, such as tiling  
 2708 and parallelization; 2) an auto-tuning system (or, alternatively, a human hardware expert) chooses  
 2709 performance-critical parameter values of the optimizations implemented in the schedule, e.g.,  
 2710 particular values of tile sizes and concrete numbers of threads.

2711 Our experiments in Section 5 show that compared to scheduling approach TVM (using its recent  
 2712 Ansor optimizer [Zheng et al. 2020a] for schedule generation), our approach achieves competitive  
 2713 and sometimes even better performance, e.g., speedups up to 2.22× on GPU and 3.55× on CPU over  
 2714 TVM+Ansor for computations taken from TVM’s favorable application domain (deep learning).  
 2715 Section 5 discusses that our better performance is due to the design and structure of our general  
 2716 optimization space (Table 1) which can be efficiently explored fully automatically using state-of-  
 2717 the-art auto tuning techniques [Rasch et al. 2021]. We focus on TVM in our experiments (rather  
 2718 than, e.g. Halide) to make experimenting challenging for us: TVM+Ansor has proved to achieve  
 2719 higher performance on GPUs and CPUs than popular state-of-practice approaches [Zheng et al.  
 2720 2020a], including Halide, pyTorch [Paszke et al. 2019], and the recent FlexTensor optimizer [Zheng  
 2721 et al. 2020b].

2722 Recent approach TensorIR [Feng et al. 2022] is a compiler for deep learning computations that  
 2723 achieves higher performance than TVM on NVIDIA GPUs. However, this performance gain over  
 2724 TVM is mainly achieved by exploiting the domain-specific *tensor core* [NVIDIA 2017] extensions  
 2725 of NVIDIA GPUs, which compute in hardware the multiplications of small, low-precision 4 × 4  
 2726 matrices. For this, TensorIR introduces the concept of *blocks* which represent sub-computations,  
 2727 e.g., for computing matrix multiplication on 4 × 4 sub-matrices. These blocks are then mapped by  
 2728 TensorIR to domain-specific hardware extensions, such as NVIDIA’s tensor cores, leading to high  
 2729 performance.

2730 While domain-specific hardware extensions are not targeted by this paper, we can naturally ex-  
 2731 ploit them in our approach, similar to TensorIR, as we plan for our future work: the sub-computations  
 2732 targeted by the current hardware extensions, like matrix multiplication on 4 × 4 matrices, can be  
 2733 straightforwardly expressed in our approach (Figure 14). Thus, we can match these sub-expressions  
 2734 in our low-level representation and map them to hardware extensions in our generated code. For  
 2735 this, instead of relying on a full partitioning in our low-level representation (as in Figure 15) such  
 2736 that we can apply scalar function  $f$  to the fully de-composed data (consisting of a single data  
 2737 element only in the case of a full partitioning), we plan to rely on a coarser-grained partitioning  
 2738 schema, e.g., down to only 4 × 4 matrices (rather than 1 × 1 matrices, as in the case of a full parti-  
 2739 tioning). This allows us replacing scalar function  $f$  (which in the case of matrix multiplication is

2740 <sup>22</sup>Rasch et al. [2023] introduce (optionally) a scheduling language for MDH to incorporate expert knowledge into MDH’s  
 2741 optimization process, e.g., to achieve 1) better optimization, as an auto-tuning system might not always make the same  
 2742 high-quality optimization decisions as a human expert, or 2) faster auto-tuning, as some (or even all) optimization decisions  
 2743 might be made by the expert user and thus are not left to the costly auto-tuner.

2745 a simple scalar multiplication  $*$ ) with the operation supported by the hardware extension, such  
2746 as matrix multiplication on  $4 \times 4$  matrices. We expect for our future work to achieve the same  
2747 advantages over TensorIR as over TVM, because apart from supporting domain-specific hardware  
2748 extensions, TensorIR is very similar to TVM.

2749  
2750 *Portability.* While scheduling approaches achieve high performance, they tend to struggle with  
2751 achieving portability. This is because even though the approaches often provide different, pre-  
2752 implemented backends (e.g., a CUDA backend to target NVIDIA GPUs and an OpenCL backend for  
2753 CPUs), they do not propose any structured methodology about how new backends can be added, e.g.,  
2754 for potentially upcoming architectures, with potentially deeper memory and core hierarchies than  
2755 GPUs and CPUs. This might be particularly critical (or requiring significant development effort)  
2756 for the application area of deep learning which is the main target of many scheduling approaches,  
2757 e.g., TVM and TensorIR, and for which new architectures are arising continuously [Hennessy and  
2758 Patterson 2019].

2759 In contrast, we introduce in this paper a formally precise recipe for correct-by-construction code  
2760 generation in different backends (including OpenMP, CUDA, and OpenCL), generically in the target  
2761 architecture: we introduce an architecture-agnostic low-level representation (Section 3) as target  
2762 for our high-level programs (Section 2), and we describe formally how our high-level programs are  
2763 automatically lowered to our low-level representation (Section 4), based on the architecture-agnostic  
2764 optimization space in Table 1. Our Appendix, Section F, outlines how executable, imperative-style  
2765 program code is straightforwardly generated from low-level expressions, which we plan to discuss  
2766 and illustrate in detail in our future work.

2767  
2768 *Productivity.* Scheduling approaches rely on a two-step optimization process, as discussed  
2769 above: implementing a schedule (first step) and choosing optimized values of performance-critical  
2770 parameters within that schedule (second step). While the second step often can be easily automa-  
2771 tized, e.g., via auto-tuning [Chen et al. 2018b], the first step – implementing a schedule – usually  
2772 has to be conducted manually by the user for high performance, which requires expert knowledge  
2773 and thus hinders productivity. The lack of formal foundation of many scheduling approaches  
2774 further complicates implementing schedules for the user, as implementation becomes error prone  
2775 and hardly predictable. For example, Fireiron’s schedules can achieve high performance, close to  
2776 GPUs’ peak, but schedules in Fireiron can easily generate incorrect low-level code: Fireiron cannot  
2777 guarantee that optimizations expressed in its scheduling language are semantics preserving, e.g.,  
2778 based on a formal foundation as done in this work, making programming Fireiron’s schedules error  
2779 prone and complex for the user. Similarly, TVM is sometimes unable to detect user errors in both  
2780 its high-level language (as discussed in Section 5.1) as well as scheduling language [Apache TVM  
2781 Community 2022e]. Safety in parallel programming is an ongoing major demand, in particular  
2782 from industry [Khronos 2022a].

2783 Auto schedulers, such as Halide’s optimization engine [Mullapudi et al. 2016] and TVM’s recent  
2784 Ansor [Zheng et al. 2020a], aim at automatically generating well-performing, correct schedules  
2785 for the user. However, a major flaw of the current auto schedulers is that even though they work  
2786 well for some computations (e.g., from deep learning, as TVM’s Ansor), they may perform worse  
2787 for others. For example, our approach achieves a speedup over TVM+Ansor of  $> 100\times$  already  
2788 for straightforward dot products. This is because Ansor does not exploit multiple thread blocks  
2789 and uses only a small number of threads for reduction computations, which is usually beneficial  
2790 for reductions as computed within convolutions and matrix multiplications used in deep learning  
2791 applications (because parallelization can be better exploited for outer loops of these computations),  
2792 but not for pure reductions.

To avoid the productivity issues of scheduling approaches, we have designed our optimization process as fully auto-tunable, thereby freeing the user from the burden and complexity of making complex optimization decisions. Our optimization space (Table 1) is designed as agnostic of a target application area and hardware architecture, thereby achieving high performance for various combinations of applications and architectures (Section 5). Correctness of optimizations is ensured in our approach by introducing a formal foundation that enables mathematical reasoning about correctness. Our optimization process is designed as *correct-by-construction*, meaning that any valid optimization decisions (i.e., a particular choice of tuning parameters in Table 1 that satisfy the constraints) leads to a correct expression in our low-level expression (as in Figure 15). In contrast, approaches such as introduced by Clément and Cohen [2022] formally validate optimization decisions of scheduling approaches in already generated low-level code. Thereby, such approaches work potentially for arbitrary scheduling approaches (Halide, TVM, ...), but they cannot save the user at the high abstraction level from implementing incorrect optimizations (e.g., via easy-to-understand, high-level error messages indicating that an invalid optimization decisions is made) or restricting the optimization space otherwise to valid decisions only, e.g., for an efficient auto-tuning process, because the approaches check already generated low-level program code.

Scheduling approaches often also suffer from expressivity issues. For example, Fireiron is restricted to computing only matrix multiplications on only NVIDIA GPUs, and TVM does not support computations that rely on multiple combine operators different from concatenation [Apache TVM Community 2020, 2022b], e.g., as required for expressing the *Maximum Bottom Box Sum* example in Figure 14. Also, TVM has difficulties with user-defined combine operators [Apache TVM Community 2022d] and thus crashes for example *Probabilistic Record Linkage* in Figure 14. In contrast to TVM, we introduce a formal methodology about of how to manage different kinds of arbitrary, user-defined combine operators (Section 3), which is considered as challenging [Apache TVM Community 2020].

## 6.2 Polyhedral Approaches

Polyhedral approaches such as TC [Vasilache et al. 2019], PPCC [Verdoolaege et al. 2013], Pluto [Bondhugula et al. 2008b], Polly [Grosser et al. 2012], and the recent AKG [Bastoul et al. 2022] rely on a formal, geometrically-inspired representation, called *polyhedral model*. Polyhedral approaches often achieve high user productivity, e.g., by automatically parallelizing and optimizing straightforward sequential code. However, the approaches tend to have problems with achieving high performance and portability when used for generating low-level code, as we outline in the following. In Section 6.5, we revisit the polyhedral approach as a potential frontend for our approach, as polyhedral transformations have proven to be efficient when used for high-level code optimizations (e.g., *loop skewing* [Wolf and Lam 1991]), rather than low-level code generation.

*Performance.* Polyhedral compilers tend to struggle with achieving their full performance potential. We argue that this performance issue of polyhedral compilers is mainly caused by the following two major reasons.

While we consider the set of polyhedral transformation (so-called *affine transformation*) as broad, expressive, and powerful, each polyhedral compiler implements a subset of expert-chosen transformations. This subset of transformations, as well as the application order of transformations, are usually fixed in a particular polyhedral compiler and chosen toward specific optimization goals only, e.g., coarse-grained parallelization and locality-aware data accesses (a.k.a. *Pluto algorithm* [Bondhugula et al. 2008a]), causing the search spaces of polyhedral compilers to be a proper subset of our space in Table 1 only. Consequently, computations that require for high performance other subsets of polyhedral transformations and/or application orders of transformations (e.g.,

2843 transformations toward fine-grained parallelization) might not achieve their full performance  
2844 potential when compiled with a particular polyhedral compiler.

2845 In contrast to the currently existing polyhedral compilers, we have designed our optimization  
2846 process as generic in goals: for example, our space is designed such that the degree of parallelization  
2847 (coarse, fine, ...) is completely auto-tunable for the particular combination of target architecture  
2848 and computation to optimize. We consider it as an interesting future work to investigate the strength  
2849 and weaknesses of the polyhedral model for expressing our generic optimization space.

2850 We see the second reason for potential performance issues in polyhedral compilers in their  
2851 difficulties with reduction-like computations. This is mainly caused by the fact that the polyhedral  
2852 model captures less semantic information than the high-level program representation introduced  
2853 in Section 2 of this paper: combine operators which are used to combine the intermediate results of  
2854 computations (e.g., operator + from Example 2 for combining the intermediate results of the dot  
2855 products within matrix multiplication) are not explicitly represented in the polyhedral model; the  
2856 polyhedral model is rather focussed on modeling memory accesses and their relative order only.  
2857 Most likely, these semantic information are missing in the polyhedral model, because polyhedral  
2858 approaches were originally intended to fully automatically optimize sequential code (such as Pluto  
2859 and PPCG) – extracting combine operators automatically from sequential code is challenging and  
2860 often even impossible (Rice’s theorem).

2861 In contrast, our proposed high-level representation explicitly captures combine operators (Figure  
2862 14), by requesting these operators explicitly from the user. This is important, because the  
2863 operators are often required for generating code that fully utilizes the highly parallel hardware of  
2864 state-of-the-art parallel architectures (GPUs, etc), as discussed in Section 5. Similarly to our ap-  
2865 proach, polyhedral compiler TC also requests combine operators explicitly from the user. However,  
2866 TC is restricted to operators + (addition), \* (multiplication), min (minimum), and max (maximum)  
2867 only, thereby TC is not able to express important examples in Figure 14, e.g., PRL which is popular  
2868 in data mining. Moreover, TC outsources the computation of its combine operators to the NVIDIA  
2869 CUB library [NVIDIA 2022a]; most likely as a workaround, because TC relies on the polyhedral  
2870 model which is not designed to capture and exploit semantic information about combine opera-  
2871 tors for optimization. Thereby, TC is dependent on external approaches for computing combine  
2872 operators, which might not always be available (e.g., for upcoming architectures).

2873 Workarounds have been proposed by the polyhedral community to target reduction-like compu-  
2874 tations [Doerfert et al. 2015; Reddy et al. 2016]. However, these approaches are limited to a subset  
2875 of computations, e.g., by not supporting user-defined scalar types [Doerfert et al. 2015] (as required  
2876 for our PRL example in Figure 14), or by being limited to GPUs only [Reddy et al. 2016]. Comparing  
2877 the semantic information captured in the polyhedral model vs our MDH-based representation have  
2878 been the focus of discussions between polyhedral experts and MDH developers [Google SIG MLIR  
2879 Open Design Meeting 2020].

2880 *Portability.* The polyhedral approach, in its general form, is a framework offering transformation  
2881 rules (affine transformations), and each individual polyhedral compiler implements a set of such  
2882 transformations which are then instantiated (e.g., with particular tile sizes) and applied when  
2883 compiling a particular application. However, individual polyhedral compilers (e.g., PPCG and Pluto)  
2884 apply a fixed set of affine transformations, thereby directly generating a schedule that is optimized  
2885 for a particular target architecture only, e.g., only GPU (as PPCG) or only CPU (as Pluto), and  
2886 it remains open which affine transformations have to be used and how for other architectures,  
2887 e.g., upcoming accelerators for deep learning computations [Hennessy and Patterson 2019] with  
2888 potentially more complex memory and core hierarchies than GPUs and CPUs. Moreover, while  
2889 we introduce an explicit low-level representation (Section 3), the polyhedral approach does not

2892 introduce representations on different abstraction levels: the model relies on one representation  
 2893 that is transformed via affine transformations. Apart from the ability of our low-level representation  
 2894 to handle combine operators (which we consider as complex and important), we see the advantages  
 2895 of our explicit low-level representation in, for example, explicitly representing memory regions,  
 2896 which allows formally defining important correctness constraints, e.g., that GPU architectures  
 2897 allow combining the results of threads in designated memory regions only – shared and device  
 2898 memory, but not registers. Furthermore, our low-level representation also allows straightforwardly  
 2899 generating executable code from it (focus of Section F in our Appendix, and planned to be discussed  
 2900 thoroughly in future work). In contrast, code generation from the polyhedral model has proven  
 2901 challenging [Bastoul et al. 2022; Grosser et al. 2015; Vasilache et al. 2022].

2902 *Productivity.* Most polyhedral compilers achieve high user productivity, by fully automatically  
 2903 parallelizing and optimizing straightforward sequential code (as Pluto and PPCG). Our approach  
 2904 currently relies on a DSL (Domain-Specific Language) for expressing computations, as discussed  
 2905 in Section 2; thus, our approach can be considered as less productive than many polyhedral  
 2906 compilers. However, Rasch et al. [2020a,b] confirm that DSL programs in our approach can be  
 2907 automatically generated from sequential code (optionally annotated with simple, OpenMP-like  
 2908 directives for expressing combine operators, enabling advanced optimizations), by using polyhedral  
 2909 tool pet [Verdoolaege and Grosser 2012] as a frontend for our approach. Thereby, we are able to  
 2910 achieve the same, high user productivity as polyhedral compilers. We consider this direction –  
 2911 combining the polyhedral model with our approach – as promising, as it enables benefitting from  
 2912 the advantages of both directions: optimizing sequential programs and making them parallelizable  
 2913 using polyhedral techniques (like *loop skewing*, as also outlined in Section 6.5), and mapping the  
 2914 optimized and parallelizable code eventually to parallel architectures based on the concepts and  
 2915 methodologies introduced in this paper.

### 2917 6.3 Functional Approaches

2918 Functional approaches map data-parallel computations that are expressed via small, formally  
 2919 defined building blocks (a.k.a. patterns [Gorlatch and Cole 2011]), such as map and reduce, to the  
 2920 memory and core hierarchies of parallel architectures, based on a strong formal foundation. Notable  
 2921 functional approaches include Accelerate [Chakravarty et al. 2011], Obsidian [Svensson et al. 2011],  
 2922 so-called *skeleton libraries* [Aldinucci et al. 2017; Enmyren and Kessler 2010; Ernstsson et al. 2018;  
 2923 Steuwer et al. 2011], and the modern Lift approach [Steuwer et al. 2015] (recently also known as  
 2924 RISE [Steuwer et al. 2022]).

2925 In the following, as functional approaches usually follow the same basic concepts and methodologies,  
 2926 we focus on comparing to Lift, because Lift is more recent than, e.g., Accelerate and  
 2927 Obsidian.

2928 *Performance.* Functional approaches tend to struggle with achieving their full performance  
 2929 potential, often caused by the design of their optimization spaces. For example, analogously to our  
 2930 approach, functional approach Lift relies on an internal low-level representation [Steuwer et al.  
 2931 2017] that is used as target for Lift’s high-level programs. However, Lift’s transformation process,  
 2932 from high level to low level, turned out to be challenging: Lift’s lowering process relies on an  
 2933 infinitely large optimization space – identifying a well-performing configuration within that space  
 2934 is too complex to be done automatically, in general, due to the space’s large and complex structure.  
 2935 As a workaround, Lift currently uses approach Elevate [Hagedorn et al. 2020b] to incorporate  
 2936 user knowledge into the optimization process; however, at the cost of productivity, as manually  
 2937 expressing optimization is challenging, particularly for non-expert users.

2938

2941 In contrast, our optimization process is designed as auto-tunable (Table 1), thereby achieving fully  
2942 automatically high performance, as confirmed in our experiments (Section 5), without involving the  
2943 user for optimization decisions. In particular, our previous work already showed that our approach –  
2944 even in its original version [Rasch and Gorlatch 2016; Rasch et al. 2019a] – can significantly  
2945 outperform Lift on GPU and CPU [Rasch et al. 2019a]. Our performance advantage over Lift is  
2946 mainly caused by the design of our optimization process: relying on formally defined tuning  
2947 parameters (Table 1) – rather than on formal transformation rules that span a too large and complex  
2948 search space, as in Lift – thereby contributing to a simpler, fully auto-tunable optimization process.  
2949

2950 *Portability.* The current functional approaches usually are designed and optimized toward code  
2951 generation in a particular programming model only. For example, Lift inherently relies on the  
2952 OpenCL programming model, because OpenCL works for multiple kinds of architectures: NVIDIA  
2953 GPU, Intel CPU, etc. However, we see two major disadvantages in addressing the portability issue via  
2954 OpenCL only: 1) GPU-specific optimizations (such as *shuffle operations* [NVIDIA 2018]) are available  
2955 only in the CUDA programming model, but not OpenCL; 2) the set of OpenCL-compatible devices is  
2956 broad but still limited; in particular, in the *new golden age for computer architectures* [Hennessy and  
2957 Patterson 2019], upcoming architectures are arising continuously and may not support the OpenCL  
2958 standard. We consider targeting new programming models as challenging for Lift, as its formal  
2959 low-level representation is inherently designed for OpenCL [Steuwer et al. 2017]; targeting further  
2960 programming models with Lift would require the design and implementation of new low-level  
2961 representations, which we do not consider as straightforward.

2962 To allow easily targeting new programming models with our approach, we have designed our  
2963 formalism as generic in the target model: our low-level representation (Figure 15) and optimization  
2964 space (Table 1) are designed and optimized toward an abstract system model (Definition 10) which  
2965 is capable of representing the device models of important programming approaches, including  
2966 OpenMP, CUDA, and OpenCL (Example 11). Furthermore, we have designed our high- and low-  
2967 level representations as minimalistic (Figures 6 and 15), e.g., by relying on three higher-order  
2968 functions only for expressing programs at the high abstraction level, which simplifies and reduces  
2969 the development effort for implementing code generators for programming models.

2970 In addition, we believe that compared to our approach, the following basic design decisions  
2971 of Lift (and similar functional approaches) complicate the process of code generation for them  
2972 and increase the development effort for implementing code generators: 1) relying on a vast set  
2973 of small patterns for expressing computations, rather than aiming at a minimalistic design as  
2974 we do (also discussed in Section 5.3); 2) relying on complex function nestings and compositions  
2975 for expressing computations, rather than avoiding nesting and relying on a fixed composition  
2976 structure of functions, as in our approach (Figure 5); 3) requiring new patterns for targeting new  
2977 classes of data-parallel computations (such as patterns *slide* and *pad* for stencils [Hagedorn et al.  
2978 2018]), which have to be non-trivially integrated into Lift’s type and optimization system (often via  
2979 extensions of the systems [Hagedorn et al. 2018; Remmelg et al. 2016]), instead of relying on a fixed  
2980 set of expressive patterns (Figure 6) and generalized optimizations (Table 1) that work for various  
2981 kinds of data-parallel computations (Figure 14); 4) expressing high-level and low-level concepts in  
2982 the same language, instead of separating high-level and low-level concepts for a more structured  
2983 and thus simpler code generation process (Figure 4). We consider these four design decisions as  
2984 disadvantageous for code generation, because they require from a code generator handling various  
2985 kinds of patterns (decision 1), and the patterns need to be translated to significantly different code  
2986 variants, depending on their nesting level and composition order (decision 2). Moreover, each  
2987 extension of patterns (decision 3) might affect code generation also for the already supported  
2988 patterns, because the existing patterns need to be combined with the new ones via composition and  
2989

nesting (decision 2). We consider mixing up high-level and low-level concepts in the same language (decision 4) as further complicating the code generation process, because code generators cannot be implemented in clear, distinct stages: *high-level language* → *low-level language* → *executable program code*.

*Productivity.* Functional approaches are expressive frameworks – to the best of our knowledge, the majority of these approaches should be able to express (possibly after some extension) many of the high-level programs that can also be expressed via our high-level representation (e.g., those presented in Figure 14).

A main difference we see between the high-level representations of existing functional approaches and the representation introduced by our approach is that the existing approaches rely on a vast set of higher-order functions for expressing computations; these functions have to be functionally composed and nested in complex ways for expressing computations. For example, expressing matrix multiplication in Lift requires also involving Lift’s pattern transpose (also when operating on non-transposed input matrices) [Remmelg et al. 2016], as per design in Lift, multi-dimensional data is considered as an array of arrays (rather than a multi-dimensional array, as in our approach as well as polyhedral approaches). In contrast, we aim to keep our high-level language minimalistic, by expressing data-parallel computations using exactly three higher-order functions and which are always used in the same, fixed order (shown in Figure 5). Rasch et al. [2020a,b] confirm that due to the minimalistic and structured design of our high-level representation, programs in our representation can even be systematically generated from straightforward, sequential program code.

Functional approaches also tend to require extension when targeting new application areas, which hinders the expressivity of the frameworks and thus also their productivity. For example, functional approach Lift [Steuwer et al. 2015] required notable extension for targeting, e.g., matrix multiplications (so-called *macro-rules* had to be added to Lift [Remmelg et al. 2016]) and stencil computations (primitives *slide* and *pad* were added and Lift’s tiling optimization had to be extended toward *overlapped tiling* [Hagedorn et al. 2018]). In contrast, the generality of our approach allows expressing matrix multiplications and stencils out of the box, without relying on domain-specific building blocks.

#### 6.4 Domain-Specific Approaches

Many approaches focus on code generation and optimization for particular domains. A popular domain-specific approach is *ATLAS* [Whaley and Dongarra 1998] which offers a convenient user interface for automatically generating and optimizing CPU code for the domain of linear algebra<sup>23</sup>. Similarly to *ATLAS*, approach *FFTW* [Frigo and Johnson 1998] is specifically designed and optimized for *Fast Fourier Transform (FFT)*, and *SPIRAL* [Puschel et al. 2005] targets the domain of *Digital Signal Processing (DSP)*.

Nowadays, the best performing, state-of-practice domain-specific approaches are often provided by vendors and specifically designed and optimized toward their target application domain and also architecture. For example, the popular vendor library *NVIDIA cuBLAS* [NVIDIA 2022b] is optimized by hand, on the assembly level, toward computing linear algebra routines on NVIDIA GPUs – cuBLAS is considered in the community as gold standard for computing linear algebra routines on GPUs. Similarly, Intel’s oneMKL library [Intel 2022c] computes with high performance linear algebra routines on Intel CPUs, and libraries *NVIDIA cuDNN* [NVIDIA 2022e] and Intel

<sup>23</sup>Previous work [Rasch et al. 2021] shows that MDH – already in its original, proof-of-concept implementation [Rasch and Gorlatch 2016; Rasch et al. 2019a] – achieves higher performance than *ATLAS*.

3039 oneDNN [Intel 2022b] work well for convolution computations on either NVIDIA GPU (cuDNN) or  
3040 Intel CPU (oneDNN), respectively.

3041 In the following, we discuss domain-specific approaches in terms of *performance*, *portability*, and  
3042 *productivity*.

3043  
3044 *Performance.* Domain-specific approaches, such as cuBLAS and cuDNN, usually achieve high  
3045 performance. This is because the approaches are hand-optimized by performance experts – on  
3046 the assembly level – to exploit the full performance potential of their target architecture. In our  
3047 experiments, we show that our approach often achieves competitive and sometimes even better  
3048 performance than domain-specific approaches provided by NVIDIA and Intel, which is mainly due  
3049 to their portability issues over different data characteristics, as we discuss in the next paragraph.  
3050

3051  
3052 *Portability.* Domain-specific approaches usually struggle with achieving portability over different  
3053 architectures. This is because the approaches are often implemented in architecture-specific  
3054 assembly code to achieve high performance. The domain-specific approaches often also struggle  
3055 with achieving portability over different characteristics of their input and output data: they usually  
3056 rely on a set of pre-implemented implementations that are each designed and optimized toward  
3057 average high performance over a range of input characteristic (e.g., their sizes). In contrast, our  
3058 approach (as well as many scheduling and polyhedral approaches) allow automatically optimizing  
3059 (auto-tuning) kernels for particular data characteristics, which is important for performance [Tillet  
3060 and Cox 2017], thereby often achieving higher performance than domain-specific approaches  
3061 for advanced data characteristics (small, uneven, irregularly shaped, . . .), e.g., as used in deep  
3062 learning. The costly time for auto-tuning is well amortized in many application areas, because the  
3063 auto-tuned implementations are re-used in many program runs. Moreover, auto-tuning avoids the  
3064 time-intensive and costly process of hand optimization by human experts.  
3065

3066  
3067 *Productivity.* Domain-specific approaches usually achieve highest productivity for their target  
3068 domain (e.g., linear algebra), by providing easy to use high-level abstractions. However, the  
3069 approaches suffer from significant expressivity issues, because – per design – they are inherently  
3070 restricted to their target application domain only. Also, the approaches are often inherently bound  
3071 to only particular architectures, e.g., only GPU (as NVIDIA cuBLAS and cuDNN) or only CPU  
3072 (as Intel oneMKL and oneDNN). Domain-specific vendor libraries, such as NVIDIA cuBLAS and  
3073 Intel oneMKL, also tend to offer the user differently performing variants of computations; the  
3074 variants have to be naively tested by the user when striving for the full performance potentials of  
3075 approaches (as discussed in Section 5.4), which is cumbersome for the user.  
3076

## 3077 6.5 Higher-Level Approaches

3078 There is a broad range of existing work that is focused on higher-level optimizations than proposed  
3079 by this work. We consider such higher-level approaches as greatly combinable with our approach.  
3080 For example, the polyhedral approach is capable of expressing algorithmic-level optimizations, like  
3081 *loop skewing* [Wolf and Lam 1991], to make programs parallelizable; such optimizations are beyond  
3082 the scope of this work, but they can be combined with our approach as demonstrated by Rasch et al.  
3083 [2020a,b]. Similarly, we consider approaches introduced by Farzan and Nicolet [2019]; Frigo et al.  
3084 [1999]; Gunnels et al. [2001]; Yang et al. [2021], which also focus on algorithmic-level optimizations,  
3085 as greatly combinable with our approach: algorithmically optimizing user code according to the  
3086 approaches’ techniques, and using our methodologies to eventually map the optimized code to  
3087 low-level code for parallel architectures.

3088 Futhark [Henriksen et al. 2017], Dex [Paszke et al. 2021], and ATL [Liu et al. 2022] are further  
 3089 approaches focussed on high-level program transformations, like advanced *flattening* mechanisms  
 3090 [Henriksen et al. 2019], thereby optimizing programs at the algorithmic level of abstraction.  
 3091 We consider using our work as backend for these approaches as promising: the three approaches  
 3092 often struggle with mapping their algorithmically optimized program variants eventually to the  
 3093 multi-layered memory and core hierarchies of state-of-the-art parallel architectures, which is  
 3094 exactly the focus of this work.

3095

## 3096 6.6 Existing Work on MDH

3097 Our work is inspired by the algebraic approach of Multi-Dimensional Homomorphisms (MDHs)  
 3098 which is introduced in the work-in-progress paper [Rasch and Gorlatch 2016]. The MDH approach,  
 3099 as presented in the previous work, relies on a semi-formal foundation and focuses on code generation  
 3100 for the OpenCL programming model only [Rasch et al. 2019a]. This work makes major contributions  
 3101 over the existing work on MDHs and its OpenCL code generation approach.

3102 We introduce a full formalization of MDH’s high-level program representation. In our new  
 3103 formalism, we rely on expressive typing: for example, we encode MDHs’ data sizes into our type  
 3104 system, e.g., by introducing *index sets* for MDAs (Definition 1), and we respect and maintain these  
 3105 sets thoroughly during MDH computations. Our expressive typing significantly contributes to  
 3106 correct and simplified code generation, as all relevant type and data size information are contained  
 3107 in our formal, low-level program representation (Figure 15) from which we eventually generate exe-  
 3108 cutable program code (Section 3). In contrast, the existing MDH work considers multi-dimensional  
 3109 arrays (MDAs) of arbitrary sizes and dimensionalities to be all of the same, straightforward type,  
 3110 which has greatly simplified the design of the proof-of-concept MDH formalism introduced by Rasch  
 3111 and Gorlatch [2016] (in particular, the definition and usage of combine operators), but at the cost  
 3112 of significantly harder and error-prone code generation: all the missing, type-relevant information  
 3113 need to be elaborated by the implementer of the code generator in the existing MDH work, e.g.,  
 3114 allocation sizes of fast memory resources used for caching input data or for storing computed  
 3115 intermediate results. Furthermore, while the original MDH work [Rasch and Gorlatch 2016] is  
 3116 focused on introducing higher-order function `md_hom` only, this work in particular also introduces  
 3117 higher order functions `inp_view` and `out_view` (Section 2.3) which express input and output  
 3118 views in a formally structured and concise manner, and which are central building blocks for  
 3119 expressing computations (Figure 14). Also, by introducing and exploiting the index set concept  
 3120 for MDAs, we have improved the definition of the concatenation operator `++` (Example 1) toward  
 3121 commutativity, which is required for important optimizations. e.g., loop permutations (expressed  
 3122 via Parameters  $D1, S1, R1$  in Table 1).

3123 A further substantial improvement is the introduction of our low-level representation (Section 3).  
 3124 It relies on a novel combination of tuning parameters (Table 1) that enhance, generalize, and extend  
 3125 the existing, proof-of-concept MDH parameters which capture a subset of OpenCL-orientated  
 3126 features only [Rasch et al. 2019a]. Moreover, while the existing MDH work introduces formally only  
 3127 parameters for flexibly choosing numbers of threads [Rasch and Gorlatch 2016] (which corresponds  
 3128 to a very limited variant of our tuning parameter  $\theta$  in Table 1, because our parameter  $\theta$  also chooses  
 3129 numbers of memory tiles and is not restricted to OpenCL); the other OpenCL parameters are  
 3130 introduced and discussed by Rasch et al. [2019a] only informally, from a technical perspective.  
 3131 With our novel parameter set, we are able to target various kinds of programming models (e.g., also  
 3132 CUDA, as in Section 5) and also to express important optimizations that are beyond the existing  
 3133 work on MDH, e.g., optimizing the memory access pattern of MDH computations: for example,  
 3134 we achieve speedups  $> 2 \times$  over existing MDH for the deep learning computations discussed  
 3135 in Section 5. Our new tuning parameters are expressive enough to represent state-of-the-art,  
 3136

3137 data-parallel implementations, e.g., as generated by scheduling and polyhedral approaches (see  
3138 Appendix, Figures 32-35), and our experiments in Section 5 confirm that auto-tuning our parameters  
3139 enables performance beyond the state of the art, including hand-optimized solutions provided by  
3140 vendors, which is not possible when using the existing MDH approach. The expressivity of our  
3141 parameters in particular also enables comparing significantly differently optimized implementations,  
3142 based on the values of formally specified tuning parameters, which we consider as promising for  
3143 structured performance analysis in future work. Moreover, our new low-level representation targets  
3144 architectures that may have arbitrarily deep memory and core hierarchies, by having optimized  
3145 our representation toward an *Abstract System Model* (Definition 10). In contrast, the existing MDH  
3146 work is focused on OpenCL-compatible architectures only.

3147 Our experimental evaluation extends previous MDH experiments by comparing also to the popular  
3148 state-of-practice approach TVM which is attracting increasing attention from both academia [[Apache](#)  
3149 [Software Foundation 2021](#)] and industry [[OctoML 2022](#)]. Also, we compare to the popular poly-  
3150 hedral compilers PPCG and Pluto, as well as the currently newest versions of hand-optimized  
3151 high-performance libraries provided by vendors. Furthermore, we have included a real-world  
3152 case study in our experiments, considering the most time-intensive computations within the three  
3153 popular deep learning neural networks ResNet-50, VGG-16, and MobileNet; the study also includes  
3154 Capsule-style convolution computations, which are considered as challenging to optimize [[Barham](#)  
3155 and [Isard 2019](#)]. Moreover, Table 14 analyzes MDH’s expressivity using new examples: it shows  
3156 that MDH – based on the new contributions of this work (e.g., view functions) – is capable of  
3157 expressing computations bMatMul, MCC\_Capsule, Histo, scan, and MBBS, which have not been  
3158 expressed via MDH in previous work. Our experiments confirm that we achieve high performance  
3159 for bMatMul and MCC\_Capsule on GPUs and CPUs, and our future work aims to thoroughly analyze  
3160 our approach for computations Histo, scan, and MBBS in terms of performance, portability, and  
3161 productivity.

3162

3163

## 3164 7 CONCLUSION

3165 We introduce a formal (de/re)-composition approach for data-parallel computations targeting  
3166 state-of-the-art parallel architectures. Our approach aims to combine three major advantages  
3167 over related approaches – performance, portability, and productivity – by introducing formal  
3168 program representations on both: 1) *high level*, for conveniently expressing – in one uniform  
3169 formalism – various kinds of data-parallel computations (including linear algebra routines, stencil  
3170 computations, data mining algorithms, and quantum chemistry computations), agnostic from  
3171 hardware and optimization details, while still capturing all information relevant for generating high-  
3172 performance program code; 2) *low level*, which allows uniformly reasoning – in the same formalism –  
3173 about optimized (de/re)-compositions of data-parallel computations targeting different kinds of  
3174 parallel architectures (GPUs, CPUs, etc). We *lower* our high-level representation to our low-level  
3175 representation, in a formally sound manner, by introducing a generic search space that is based on  
3176 performance-critical parameters. The parameters of our lowering process enable fully automatically  
3177 optimizing (auto-tuning) our low-level representations for a particular target architecture and  
3178 characteristics of the input and output data, and our low-level representation is designed such  
3179 that it can be straightforwardly transformed to executable program code in imperative-style  
3180 programming approaches (including OpenMP, CUDA, and OpenCL). Our experiments confirm  
3181 that due to the design and structure of our generic search space in combination with auto-tuning,  
3182 our approach achieves higher performance on GPUs and CPUs than popular state-of-practice  
3183 approaches, including hand-optimized libraries provided by vendors.

3184

3185

3186 **8 FUTURE WORK**

3187 We consider this work as a promising starting point for future directions. A major future goal is  
 3188 to extend our approach toward expressing and optimizing simultaneously multiple data-parallel  
 3189 computations (e.g., matrix multiplication followed by convolution), rather than optimizing computa-  
 3190 tions individually and thus independently from each other only (e.g., only matrix multiplication or  
 3191 only convolution). Such extension enables optimizations, such as *kernel fusion*, which is important  
 3192 for the overall application performance and considered as challenging [Fukuhara and Takimoto  
 3193 2022; Li et al. 2022; Wahib and Maruyama 2014]. We see this work as a promising foundation for  
 3194 our future goal, because it enables expressing and reasoning about different computations in the  
 3195 same formal framework. Targeting computations on sparse input/output data formats, inspired  
 3196 by Ben-Nun et al. [2017]; Hall [2020]; Kjolstad et al. [2017]; Pizzuti et al. [2020], is a further major  
 3197 goal, which requires extending our approach toward irregularly-shaped input and output data,  
 3198 similarly as done by Pizzuti et al. [2020]. Regarding our optimization process, we aim to introduce  
 3199 an analytical cost model for computations expressed in our formalism – based on operational  
 3200 semantics – thereby accelerating (or even avoiding) the auto-tuning overhead, similarly as done  
 3201 by Li et al. [2021]; Muller and Hoffmann [2021]. Moreover, we aim to incorporate machine-learning  
 3202 based methods into our optimization process [Leather et al. 2014], instead of relying on empirical  
 3203 auto-tuning methods only. To make our work better accessible for the community, we aim to  
 3204 implement our approach into *MLIR* [Lattner et al. 2021] which offers a reusable compiler infras-  
 3205 tructure. The contributions of this work give a precise, formal recipe of how to implement our  
 3206 introduced methods into approaches like *MLIR*. Moreover, relying on the *MLIR* framework will  
 3207 contribute to a structured code generation process in assembly-level programming models, like  
 3208 *LLVM* [Lattner and Adve 2004] and *NVIDIA PTX* [NVIDIA 2022i]. We consider targeting assembly  
 3209 languages as important for our future work: assembly code offers further, low-level optimization  
 3210 opportunities [Goto and Geijn 2008; Lai and Seznec 2013], thereby enabling our approach to poten-  
 3211 tially achieve higher performance than presented in this work for our MDH-generated CUDA and  
 3212 OpenCL code. Also, we aim to extend our approach toward distributed multi-device systems that  
 3213 are heterogeneous, inspired by dynamic load balancing approaches [Chen et al. 2010] and advanced  
 3214 data distributions techniques [Yadav et al. 2022]. Targeting domain-specific hardware extensions,  
 3215 such as *NVIDIA Tensor Cores* [NVIDIA 2017], is also an important goal for our future work, as such  
 3216 extensions allow significantly accelerating computations for the target of the extensions (e.g., deep  
 3217 learning [Markidis et al. 2018]).

3218

3219 **REFERENCES**

- 3220 Marco Aldinucci, Marco Danelutto, Peter Kilpatrick, and Massimo Torquati. 2017. Fastflow: high-level and efficient streaming  
 3221 on multi-core. *Programming multi-core and many-core computing systems, parallel and distributed computing* (2017).
- 3222 Jason Ansel, Shoaib Kamil, Kalyan Veeramachaneni, Jonathan Ragan-Kelley, Jeffrey Bosboom, Una-May O'Reilly, and  
 3223 Saman Amarasinghe. 2014. OpenTuner: An Extensible Framework for Program Autotuning. In *Proceedings of the 23rd  
 3224 International Conference on Parallel Architectures and Compilation (PACT '14)*. Association for Computing Machinery,  
 3225 New York, NY, USA, 303–316. <https://doi.org/10.1145/2628071.2628092>
- 3226 Apache. 2022. TVM: Open Deep Learning Compiler Stack. <https://github.com/apache/tvm>.
- 3227 Apache Software Foundation. 2021. TVM and Open Source ML Acceleration Conference. <https://www.tvmcon.org>.
- 3228 Apache TVM Community. 2020. Non top-level reductions in compute statements. <https://discuss.tvm.apache.org/t/non-top-level-reductions-in-compute-statements/5693>.
- 3229 Apache TVM Community. 2022a. Bind reduce axis to blocks. <https://discuss.tvm.apache.org/t/bind-reduce-axis-to-blocks/2907>.
- 3230 Apache TVM Community. 2022b. Expressing nested reduce operations. <https://discuss.tvm.apache.org/t/expressing-nested-reduce-operations/8784>.
- 3231 Apache TVM Community. 2022c. Implementing Array Packing via cache\_read. <https://discuss.tvm.apache.org/t/implementing-array-packing-via-cache-read/13360>.

3232

- 3235 Apache TVM Community. 2022d. Invalid comm\_reducer. <https://discuss.tvm.apache.org/t/invalid-comm-reducer/12788>.  
3236 Apache TVM Community. 2022e. Undetected parallelization issue. <https://discuss.tvm.apache.org/t/undetected-parallelization-issue/13224>.  
3237 Apache TVM Community. 2022f. Undetected type issue. <https://discuss.tvm.apache.org/t/undetected-type-issue/13223>.  
3238 Apache TVM Documentation. 2022a. Bind ivar to thread index thread\_ivar. <https://tvm.apache.org/docs/reference/api/python/te.html?highlight=bind#tvm.te.Stage.bind>.  
3239 Apache TVM Documentation. 2022b. Tuning High Performance Convolution on NVIDIA GPUs. [https://tvm.apache.org/docs/how\\_to/tune\\_with\\_autotvm/tune\\_conv2d\\_cuda.html](https://tvm.apache.org/docs/how_to/tune_with_autotvm/tune_conv2d_cuda.html).  
3240 David F. Bacon, Susan L. Graham, and Oliver J. Sharp. 1994. Compiler Transformations for High-Performance Computing. *ACM Comput. Surv.* 26, 4 (dec 1994), 345–420. <https://doi.org/10.1145/197405.197406>  
3241 Riyadh Baghdadi, Jessica Ray, Malek Ben Romdhane, Emanuele Del Sozzo, Abdurrahman Akkas, Yunming Zhang, Patricia Suriana, Shoail Kamil, and Saman Amarasinghe. 2019. Tiramisu: A Polyhedral Compiler for Expressing Fast and Portable Code. In *2019 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 193–205. <https://doi.org/10.1109/CGO.2019.8661197>  
3242 Lénaïc Bagnères, Oleksandr Zinenko, Stéphane Huot, and Cédric Bastoul. 2016. Opening polyhedral compiler’s black box. In *Proceedings of the 2016 International Symposium on Code Generation and Optimization (CGO ’16)*. Association for Computing Machinery, New York, NY, USA, 128–138. <https://doi.org/10.1145/2854038.2854048>  
3243 Prasanna Balaprakash, Jack Dongarra, Todd Gamblin, Mary Hall, Jeffrey K. Hollingsworth, Boyana Norris, and Richard Vuduc. 2018. Autotuning in High-Performance Computing Applications. *Proc. IEEE* 106, 11 (2018), 2068–2083. <https://doi.org/10.1109/JPROC.2018.2841200>  
3244 Paul Barham and Michael Isard. 2019. Machine Learning Systems Are Stuck in a Rut. In *Proceedings of the Workshop on Hot Topics in Operating Systems (HotOS ’19)*. Association for Computing Machinery, New York, NY, USA, 177–183. <https://doi.org/10.1145/3317550.3321441>  
3245 Cedric Bastoul, Zhen Zhang, Harenome Razanajato, Nelson Lossing, Adilla Susungi, Javier de Juan, Etienne Filhol, Baptiste Jarry, Gianpietro Consolaro, and Renwei Zhang. 2022. Optimizing GPU Deep Learning Operators with Polyhedral Scheduling Constraint Injection. In *2022 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 313–324. <https://doi.org/10.1109/CGO53902.2022.9741260>  
3246 Tal Ben-Nun, Johannes de Fine Licht, Alexandros Nikolaos Ziogas, Timo Schneider, and Torsten Hoefer. 2019. Stateful Dataflow Multigraphs: A Data-Centric Model for Performance Portability on Heterogeneous Architectures. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC ’19)*.  
3247 Tal Ben-Nun, Michael Sutton, Sreepathi Pai, and Keshav Pingali. 2017. Groute: An Asynchronous Multi-GPU Programming Model for Irregular Computations. *SIGPLAN Not.* 52, 8 (Jan. 2017), 235–248. <https://doi.org/10.1145/3155284.3018756>  
3248 Richard S. Bird. 1989. Lectures on Constructive Functional Programming. In *Constructive Methods in Computing Science*, Manfred Broy (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 151–217.  
3249 Guy E. Blelloch. 1990. *Prefix Sums and Their Applications*. Technical Report CMU-CS-90-190. School of Computer Science, Carnegie Mellon University.  
3250 Barry Boehm, Bradford Clark, Ellis Horowitz, Chris Westland, Ray Madachy, and Richard Selby. 1995. Cost models for future software life cycle processes: COCOMO 2.0. *Annals of Software Engineering* 1, 1 (1995), 57–94. <https://doi.org/10.1007/BF02249046>  
3251 Uday Bondhugula. 2020. High Performance Code Generation in MLIR: An Early Case Study with GEMM. [arXiv:cs/2003.00532](https://arxiv.org/abs/cs/2003.00532)  
3252 Uday Bondhugula, Muthu Baskaran, Sriram Krishnamoorthy, J. Ramanujam, Atanas Rountev, and P. Sadayappan. 2008a. Automatic Transformations for Communication-Minimized Parallelization and Locality Optimization in the Polyhedral Model. In *Compiler Construction*, Laurie Hendren (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 132–146.  
3253 Uday Bondhugula, A Hartono, J Ramanujam, and P Sadayappan. 2008b. Pluto: A practical and fully automatic polyhedral program optimization system. In *Proceedings of the ACM SIGPLAN 2008 Conference on Programming Language Design and Implementation (PLDI 08)*, Tucson, AZ (June 2008). Citeseer.  
3254 C++ reference. 2022. Date and time utilities. <https://en.cppreference.com/w/cpp/chrono>.  
3255 José María Cecilia, José Manuel García, and Manuel Ujaldón. 2012. CUDA 2D Stencil Computations for the Jacobi Method. In *Applied Parallel and Scientific Computing*, Kristján Jónasson (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 173–183.  
3256 Manuel M.T. Chakravarty, Gabriele Keller, Sean Lee, Trevor L. McDonell, and Vinod Grover. 2011. Accelerating Haskell Array Codes with Multicore GPUs. In *Proceedings of the Sixth Workshop on Declarative Aspects of Multicore Programming (DAMP ’11)*. Association for Computing Machinery, New York, NY, USA, 3–14. <https://doi.org/10.1145/1926354.1926358>  
3257 Chun Chen, Jacqueline Chame, and Mary Hall. 2008. *CHiLL: A framework for composing high-level loop transformations*. Technical Report. Technical Report 08-897, U. of Southern California.  
3258  
3259  
3260  
3261  
3262  
3263  
3264  
3265  
3266  
3267  
3268  
3269  
3270  
3271  
3272  
3273  
3274  
3275  
3276  
3277  
3278  
3279  
3280  
3281  
3282  
3283

- 3284 Long Chen, Oreste Villa, Sriram Krishnamoorthy, and Guang R. Gao. 2010. Dynamic load balancing on single- and  
 3285 multi-GPU systems. In *2010 IEEE International Symposium on Parallel Distributed Processing (IPDPS)*. 1–12. <https://doi.org/10.1109/IPDPS.2010.5470413>
- 3286 Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang,  
 3287 Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. 2018a. TVM: An Automated End-to-End Optimizing  
 3288 Compiler for Deep Learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*.  
 3289 USENIX Association, Carlsbad, CA, 578–594. <https://www.usenix.org/conference/osdi18/presentation/chen>
- 3290 Tianqi Chen, Lianmin Zheng, Eddie Yan, Ziheng Jiang, Thierry Moreau, Luis Ceze, Carlos Guestrin, and Arvind Kr-  
 3291 ishnamurthy. 2018b. Learning to Optimize Tensor Programs. In *Advances in Neural Information Processing Systems*,  
 3292 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc.  
 3293 <https://proceedings.neurips.cc/paper/2018/file/8b5700012be65c9da25f49408d959ca0-Paper.pdf>
- 3294 Peter Christen. 2012. *Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*.  
 3295 Springer Publishing Company, Incorporated.
- 3296 Krzysztof Ciesielski. 1997. *Set theory for the working mathematician*. Number 39. Cambridge University Press.
- 3297 Basile Clément and Albert Cohen. 2022. End-to-End Translation Validation for the Halide Language. In *OOPSLA 2022 - Conference on Object-Oriented Programming Systems, Languages, and Applications (Proceedings of the ACM on Programming Languages (PACMPL))*. Vol. 6. Auckland, New Zealand. <https://doi.org/10.1145/3527328>
- 3298 MURRAY I. COLE. 1995. PARALLEL PROGRAMMING WITH LIST HOMOMORPHISMS. *Parallel Processing Letters* 05, 02  
 3299 (1995), 191–203. <https://doi.org/10.1142/S0129626495000175> arXiv:<https://doi.org/10.1142/S0129626495000175>
- 3300 Haskell B. Curry. 1980. Some Philosophical Aspects of Combinatory Logic. In *The Kleene Symposium*, Jon Barwise, H. Jerome  
 3301 Keisler, and Kenneth Kunen (Eds.). Studies in Logic and the Foundations of Mathematics, Vol. 101. Elsevier, 85–101.  
 3302 [https://doi.org/10.1016/S0049-237X\(08\)71254-0](https://doi.org/10.1016/S0049-237X(08)71254-0)
- 3303 Johannes Doerfert, Kevin Streit, Sebastian Hack, and Zino Benaissa. 2015. Polly’s Polyhedral Scheduling in the Presence of  
 3304 Reductions. *CoRR* abs/1505.07716 (2015). arXiv:1505.07716 <http://arxiv.org/abs/1505.07716>
- 3305 Vincent Dumoulin and Francesco Visin. 2018. A guide to convolution arithmetic for deep learning. arXiv:stat.ML/1603.07285
- 3306 Johan Enmyren and Christoph W. Kessler. 2010. SkePU: A Multi-Backend Skeleton Programming Library for Multi-GPU  
 3307 Systems. In *Proceedings of the Fourth International Workshop on High-Level Parallel Programming and Applications (HLPP  
 3308 ’10)*. Association for Computing Machinery, New York, NY, USA, 5–14. <https://doi.org/10.1145/1863482.1863487>
- 3309 August Ernstsson, Lu Li, and Christoph Kessler. 2018. SkePU 2: Flexible and Type-Safe Skeleton Programming for Hetero-  
 3310 geneous Parallel Systems. *International Journal of Parallel Programming* 46, 1 (2018), 62–80. <https://doi.org/10.1007/s10766-017-0490-5>
- 3311 Facebook Research. 2022. Tensor Comprehensions. <https://github.com/facebookresearch/TensorComprehensions>.
- 3312 Azadeh Farzan and Victor Nicolet. 2019. Modular Divide-and-Conquer Parallelization of Nested Loops. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2019)*. Association for  
 3313 Computing Machinery, New York, NY, USA, 610–624. <https://doi.org/10.1145/3314221.3314612>
- 3314 Siyuan Feng, Bohan Hou, Hongyi Jin, Wuwei Lin, Junru Shao, Ruihang Lai, Zihao Ye, Lianmin Zheng, Cody Hao Yu,  
 3315 Yong Yu, and Tianqi Chen. 2022. TensorIR: An Abstraction for Automatic Tensorized Program Optimization. <https://doi.org/10.48550/ARXIV.2207.04296>
- 3316 M. Frigo and S.G. Johnson. 1998. FFTW: an adaptive software architecture for the FFT. In *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP ’98 (Cat. No.98CH36181)*, Vol. 3. 1381–1384  
 3317 vol.3. <https://doi.org/10.1109/ICASSP.1998.681704>
- 3318 M. Frigo, C.E. Leiserson, H. Prokop, and S. Ramachandran. 1999. Cache-oblivious algorithms. In *40th Annual Symposium on Foundations of Computer Science (Cat. No.99CB37039)*. 285–297. <https://doi.org/10.1109/SFCS.1999.814600>
- 3319 Junji Fukuhara and Munehiro Takimoto. 2022. Automated Kernel Fusion for GPU Based on Code Motion. In *Proceedings of the 23rd ACM SIGPLAN/SIGBED International Conference on Languages, Compilers, and Tools for Embedded Systems (LCTES 2022)*. Association for Computing Machinery, New York, NY, USA, 151–161. <https://doi.org/10.1145/3519941.3535078>
- 3320 Sylvain Girbal, Nicolas Vasilache, Cédric Bastoul, Albert Cohen, David Parello, Marc Sigler, and Olivier Temam. 2006.  
 3321 Semi-Automatic Composition of Loop Transformations for Deep Parallelism and Memory Hierarchies. *International Journal of Parallel Programming* 34, 3 (2006), 261–317. <https://doi.org/10.1007/s10766-006-0012-3>
- 3322 GNU/Linux. 2022. `clock_gettime(3)` – Linux man page. [https://linux.die.net/man/3/clock\\_gettime](https://linux.die.net/man/3/clock_gettime).
- 3323 Horacio González-Vélez and Mario Leyton. 2010. A survey of algorithmic skeleton frameworks: high-level structured parallel  
 3324 programming enablers. *Software: Practice and Experience* 40, 12 (2010), 1135–1160. <https://doi.org/10.1002/spe.1026>  
 3325 arXiv:<https://onlinelibrary.wiley.com/doi/10.1002/spe.1026>
- 3326 Google SIG MLIR Open Design Meeting. 2020. Using MLIR for Multi-Dimensional Homomorphisms. <https://drive.google.com/file/d/1bS4vapyzf7705wWj7t3WzwWkcbxZyF6/view>
- 3327 Sergei Gorlatch. 1999. Extracting and implementing list homomorphisms in parallel program development. *Science of Computer Programming* 33, 1 (1999), 1–27. [https://doi.org/10.1016/S0167-6423\(97\)00014-2](https://doi.org/10.1016/S0167-6423(97)00014-2)
- 3328

- 3333 3334 Sergei Gorlatch and Murray Cole. 2011. Parallel skeletons. In *Encyclopedia of parallel computing*. Springer-Verlag GmbH, 1417–1422.
- 3335 3336 S. Gorlatch and C. Lengauer. 1997. (De) composition rules for parallel scan and reduction. In *Proceedings. Third Working Conference on Massively Parallel Programming Models (Cat. No.97TB100228)*, 23–32. <https://doi.org/10.1109/MPPM.1997.715958>
- 3337 3338 Kazushige Goto and Robert A. van de Geijn. 2008. Anatomy of High-Performance Matrix Multiplication. *ACM Trans. Math. Softw.* 34, 3, Article 12 (may 2008), 25 pages. <https://doi.org/10.1145/1356052.1356053>
- 3339 3340 Tobias Grosser, Armin Groesslinger, and Christian Lengauer. 2012. Polly - Performing Polyhedral Optimizations on a Low-level Intermediate Representation. *Parallel Processing Letters* 22, 04 (2012), 1250010. <https://doi.org/10.1142/S0129626412500107> arXiv:<https://doi.org/10.1142/S0129626412500107>
- 3341 3342 Tobias Grosser, Sven Verdoolaege, and Albert Cohen. 2015. Polyhedral AST Generation Is More Than Scanning Polyhedra. *ACM Trans. Program. Lang. Syst.* 37, 4, Article 12 (jul 2015), 50 pages. <https://doi.org/10.1145/2743016>
- 3343 3344 John A. Gunnels, Fred G. Gustavson, Greg M. Henry, and Robert A. van de Geijn. 2001. FLAME: Formal Linear Algebra Methods Environment. *ACM Trans. Math. Softw.* 27, 4 (dec 2001), 422–455. <https://doi.org/10.1145/504210.504213>
- 3345 3346 Bastian Hagedorn, Archibald Samuel Elliott, Henrik Barthels, Rastislav Bodik, and Vinod Grover. 2020a. Fireiron: A Data-Movement-Aware Scheduling Language for GPUs. In *Proceedings of the ACM International Conference on Parallel Architectures and Compilation Techniques (PACT '20)*. Association for Computing Machinery, New York, NY, USA, 71–82. <https://doi.org/10.1145/3410463.3414632>
- 3347 3348 Bastian Hagedorn, Johannes Lenfers, Thomas Kundefinedhler, Xueying Qin, Sergei Gorlatch, and Michel Steuwer. 2020b. Achieving High-Performance the Functional Way: A Functional Pearl on Expressing High-Performance Optimizations as Rewrite Strategies. *Proc. ACM Program. Lang.* 4, ICFP, Article 92 (Aug. 2020), 29 pages. <https://doi.org/10.1145/3408974>
- 3349 3350 Bastian Hagedorn, Larisa Stoltzfus, Michel Steuwer, Sergei Gorlatch, and Christophe Dubach. 2018. High Performance Stencil Code Generation with Lift. In *Proceedings of the 2018 International Symposium on Code Generation and Optimization (CGO 2018)*. Association for Computing Machinery, New York, NY, USA, 100–112. <https://doi.org/10.1145/3168824>
- 3351 3352 Mary Hall. 2020. Research Challenges in Compiler Technology for Sparse Tensors. In *2020 IEEE/ACM 10th Workshop on Irregular Applications: Architectures and Algorithms (IA3)*. viii–viii. <https://doi.org/10.1109/IA351965.2020.00006>
- 3353 Maurice H Halstead. 1977. *Elements of Software Science (Operating and programming systems series)*. Elsevier Science Inc.
- 3354 3355 Mark Harris et al. 2007. Optimizing Parallel Reduction in CUDA. *NVIDIA Developer Technology* (2007).
- 3356 Haskell Wiki. 2013. Parameter Order. [https://wiki.haskell.org/Parameter\\_order](https://wiki.haskell.org/Parameter_order)
- 3357 Haskell.org. 2022. Haskell: An advanced, purely functional programming language. <https://www.haskell.org>.
- 3358 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep Residual Learning for Image Recognition. *CoRR* abs/1512.03385 (2015). arXiv:1512.03385 <http://arxiv.org/abs/1512.03385>
- 3359 3360 John L. Hennessy and David A. Patterson. 2019. A New Golden Age for Computer Architecture. *Commun. ACM* 62, 2 (jan 2019), 48–60. <https://doi.org/10.1145/3282307>
- 3361 3362 Troels Henriksen, Sune Hellfrtzsch, Ponnuswamy Sadayappan, and Cosmin Oancea. 2020. Compiling Generalized Histograms for GPU. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. 1–14. <https://doi.org/10.1109/SC41405.2020.00101>
- 3363 3364 Troels Henriksen, Niels G. W. Serup, Martin Elsman, Fritz Henglein, and Cosmin E. Oancea. 2017. Futhark: Purely Functional GPU-Programming with Nested Parallelism and in-Place Array Updates. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2017)*. Association for Computing Machinery, New York, NY, USA, 556–571. <https://doi.org/10.1145/3062341.3062354>
- 3365 3366 Troels Henriksen, Frederik Thorøe, Martin Elsman, and Cosmin Oancea. 2019. Incremental Flattening for Nested Data Parallelism. In *Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming (PPoPP '19)*. Association for Computing Machinery, New York, NY, USA, 53–67. <https://doi.org/10.1145/3293883.3295707>
- 3367 K Hentschel et al. 2008. Das Krebsregister-Manual der Gesellschaft der epidemiologischen Krebsregister in Deutschland e.V. Zuckschwerdt Verlag.
- 3368 3369 Geoffrey E Hinton, Sara Sabour, and Nicholas Frosst. 2018. Matrix capsules with EM routing. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=HJWLfGWRb>
- 3370 3371 Torsten Hoefler and Roberto Belli. 2015. Scientific Benchmarking of Parallel Computing Systems: Twelve Ways to Tell the Masses When Reporting Performance Results. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '15)*. Association for Computing Machinery, New York, NY, USA, Article 73, 12 pages. <https://doi.org/10.1145/2807591.2807644>
- 3372 3373 Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. 2017. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. *CoRR* abs/1704.04861 (2017). arXiv:1704.04861 <http://arxiv.org/abs/1704.04861>
- 3374 3375 Cristina Hristea, Daniel Lenoski, and John Keen. 1997. Measuring Memory Hierarchy Performance of Cache-Coherent Multiprocessors Using Micro Benchmarks. In *Proceedings of the 1997 ACM/IEEE Conference on Supercomputing (SC '97)*.
- 3376 3377 3378 3379 3380 3381

- 3382      Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/509593.509638>
- 3383      Intel. 2019. Math Kernel Library Improved Small Matrix Performance Using Just-in-Time (JIT) Code Generation for Matrix Multiplication (GEMM). <https://www.intel.com/content/www/us/en/developer/articles/technical/onemkl-improved-small-matrix-performance-using-just-in-time-jit-code.html>.
- 3384      Intel. 2022a. oneAPI Math Kernel Library Link Line Advisor. <https://www.intel.com/content/www/us/en/developer/tools/oneapi/onemkl-link-line-advisor.html>.
- 3385      Intel. 2022b. oneDNN. [https://oneapi-src.github.io/oneDNN/group\\_dnnl\\_api.html](https://oneapi-src.github.io/oneDNN/group_dnnl_api.html).
- 3386      Intel. 2022c. oneMKL. <https://www.intel.com/content/www/us/en/develop/documentation/oneapi-programming-guide/top-api-based-programming/intel-oneapi-math-kernel-library-onemkl.html>.
- 3387      Wayne Kelly and William Pugh. 1998. *A framework for unifying reordering transformations*. Technical Report. Technical Report UMIACS-TR-92-126.1.
- 3388      Malik Khan, Protonu Basu, Gabe Rudy, Mary Hall, Chun Chen, and Jacqueline Chame. 2013. A Script-Based Autotuning Compiler System to Generate High-Performance CUDA Code. *ACM Trans. Archit. Code Optim.* 9, 4, Article 31 (jan 2013), 25 pages. <https://doi.org/10.1145/2400682.2400690>
- 3389      Khronos. 2022a. Khronos Releases Vulkan SC 1.0 Open Standard for Safety-Critical Accelerated Graphics and Compute. <https://www.khronos.org/news/press/khronos-releases-vulkan-safety-critical-1.0-specification-to-deliver-safety-critical-graphics-compute>.
- 3390      Khronos. 2022b. OpenCL: Open Standard For Parallel Programming of Heterogeneous Systems. <https://www.khronos.org/opencl/>.
- 3391      Jinsung Kim, Aravind Sukumaran-Rajam, Vineeth Thumma, Sriram Krishnamoorthy, Ajay Panyala, Louis-Noël Pouchet, Atanas Rountev, and P. Sadayappan. 2019. A Code Generator for High-Performance Tensor Contractions on GPUs. In *2019 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 85–95. <https://doi.org/10.1109/CGO.2019.8661182>
- 3392      Fredrik Kjolstad, Shoaib Kamil, Stephen Chou, David Lugato, and Saman Amarasinghe. 2017. The Tensor Algebra Compiler. *Proc. ACM Program. Lang.* 1, OOPSLA, Article 77 (oct 2017), 29 pages. <https://doi.org/10.1145/3133901>
- 3393      Michael Klemm, Alejandro Duran, Xinmin Tian, Hideki Saito, Diego Caballero, and Xavier Martorell. 2012. Extending OpenMP\* with Vector Constructs for Modern Multicore SIMD Architectures. In *Proceedings of the 8th International Conference on OpenMP in a Heterogeneous World (IWOMP'12)*. Springer-Verlag, Berlin, Heidelberg, 59–72. [https://doi.org/10.1007/978-3-642-30961-8\\_5](https://doi.org/10.1007/978-3-642-30961-8_5)
- 3394      Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*, F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger (Eds.), Vol. 25. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>
- 3395      Junjie Lai and André Seznec. 2013. Performance upper bound analysis and optimization of SGEMM on Fermi and Kepler GPUs. In *Proceedings of the 2013 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 1–10. <https://doi.org/10.1109/CGO.2013.6494986>
- 3396      Monica D. Lam, Edward E. Rothberg, and Michael E. Wolf. 1991. The Cache Performance and Optimizations of Blocked Algorithms. In *Proceedings of the Fourth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS IV)*. Association for Computing Machinery, New York, NY, USA, 63–74. <https://doi.org/10.1145/106972.106981>
- 3397      C. Lattner and V. Adve. 2004. LLVM: a compilation framework for lifelong program analysis and transformation. In *International Symposium on Code Generation and Optimization, 2004. CGO 2004*. 75–86. <https://doi.org/10.1109/CGO.2004.1281665>
- 3398      Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilescu, and Oleksandr Zinenko. 2021. MLIR: Scaling Compiler Infrastructure for Domain Specific Computation. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 2–14. <https://doi.org/10.1109/CGO51591.2021.9370308>
- 3399      Hugh Leather, Edwin Bonilla, and Michael O’boyle. 2014. Automatic Feature Generation for Machine Learning-Based Optimising Compilation. *ACM Trans. Archit. Code Optim.* 11, 1, Article 14 (Feb. 2014), 32 pages. <https://doi.org/10.1145/2536688>
- 3400      Ao Li, Bojian Zheng, Gennady Pekhimenko, and Fan Long. 2022. Automatic Horizontal Fusion for GPU Kernels. In *2022 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 14–27. <https://doi.org/10.1109/CGO53902.2022.9741270>
- 3401      Rui Li, Yufan Xu, Aravind Sukumaran-Rajam, Atanas Rountev, and P. Sadayappan. 2021. Analytical Characterization and Design Space Exploration for Optimization of CNNs. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS 2021)*. Association for Computing Machinery, New York, NY, USA, 928–942. <https://doi.org/10.1145/3445814.3446759>

- 3431 Xiaqing Li, Guangyan Zhang, H. Howie Huang, Zhufan Wang, and Weimin Zheng. 2016. Performance Analysis of  
3432 GPU-Based Convolutional Neural Networks. In *2016 45th International Conference on Parallel Processing (ICPP)*. 67–76.  
3433 <https://doi.org/10.1109/ICPP.2016.15>
- 3434 Amanda Liu, Gilbert Louis Bernstein, Adam Chlipala, and Jonathan Ragan-Kelley. 2022. Verified Tensor-Program Op-  
3435 timization via High-Level Scheduling Rewrites. *Proc. ACM Program. Lang.* 6, POPL, Article 55 (jan 2022), 28 pages.  
3436 <https://doi.org/10.1145/3498717>
- 3437 Stefano Markidis, Steven Wei Der Chien, Erwin Laure, Ivy Bo Peng, and Jeffrey S. Vetter. 2018. NVIDIA Tensor Core  
3438 Programmability, Performance & Precision. In *2018 IEEE International Parallel and Distributed Processing Symposium  
3439 Workshops (IPDPSW)*. 522–531. <https://doi.org/10.1109/IPDPSW.2018.00091>
- 3440 T.J. McCabe. 1976. A Complexity Measure. *IEEE Transactions on Software Engineering* SE-2, 4 (1976), 308–320. <https://doi.org/10.1109/TSE.1976.233837>
- 3441 Kathryn S. McKinley, Steve Carr, and Chau-Wen Tseng. 1996. Improving Data Locality with Loop Transformations. *ACM  
3442 Trans. Program. Lang. Syst.* 18, 4 (jul 1996), 424–453. <https://doi.org/10.1145/233561.233564>
- 3443 Xinxin Mei, Kaiyong Zhao, Chengjian Liu, and Xiaowen Chu. 2014. Benchmarking the Memory Hierarchy of Modern  
3444 GPUs. In *Network and Parallel Computing*, Ching-Hsien Hsu, Xuanhua Shi, and Valentina Salapura (Eds.). Springer Berlin  
Heidelberg, Berlin, Heidelberg, 144–156.
- 3445 Michael Kruse. 2022. Polyhedral Parallel Code Generation. <https://github.com/Meinersbur/ppcg>, commit = 8a74e46, date =  
19.11.2020.
- 3446 Ravi Teja Mullapudi, Andrew Adams, Dillon Sharlet, Jonathan Ragan-Kelley, and Kayvon Fatahalian. 2016. Automatically  
3447 Scheduling Halide Image Processing Pipelines. *ACM Trans. Graph.* 35, 4, Article 83 (jul 2016), 11 pages. <https://doi.org/10.1145/2897824.2925952>
- 3448 Stefan K. Muller and Jan Hoffmann. 2021. Modeling and Analyzing Evaluation Cost of CUDA Kernels. *Proc. ACM Program.  
3449 Lang.* 5, POPL, Article 25 (Jan. 2021), 31 pages. <https://doi.org/10.1145/3434306>
- 3450 NVIDIA. 2017. Programming Tensor Cores in CUDA 9. [https://developer.nvidia.com/blog/3451\\_programming-tensor-cores-cuda-9/](https://developer.nvidia.com/blog/3451_programming-tensor-cores-cuda-9/)
- 3452 NVIDIA. 2018. Warp-level Primitives. <https://developer.nvidia.com/blog/using-cuda-warp-level-primitives/>
- 3453 NVIDIA. 2022a. CUB. <https://docs.nvidia.com/cuda/cub/>.
- 3454 NVIDIA. 2022b. cuBLAS. <https://developer.nvidia.com/cublas>.
- 3455 NVIDIA. 2022c. cuBLAS – BLAS-like Extension. <https://docs.nvidia.com/cuda/cublas/index.html#blas-like-extension>
- 3456 NVIDIA. 2022d. cuBLAS – Using the cuBLASLT API. <https://docs.nvidia.com/cuda/cublas/index.html#using-the-cublaslt-api>
- 3457 NVIDIA. 2022e. CUDA Deep Neural Network library. <https://developer.nvidia.com/cudnn>
- 3458 NVIDIA. 2022f. CUDA Programming Guide. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/>.
- 3459 NVIDIA. 2022g. CUDA Toolkit Documentation. <https://docs.nvidia.com/cuda/>.
- 3460 NVIDIA. 2022h. NVRTC. <https://docs.nvidia.com/cuda/nvrtc>.
- 3461 NVIDIA. 2022i. Parallel Thread Execution ISA. <https://docs.nvidia.com/cuda/parallel-thread-execution>.
- 3462 OctoML. 2022. Accelerated Machine Learning Deployment. <https://octoml.ai>.
- 3463 Geraldo F. Oliveira, Juan Gómez-Luna, Lois Orosa, Saugata Ghose, Nandita Vijaykumar, Ivan Fernandez, Mohammad  
3464 Sadrosadati, and Onur Mutlu. 2021. DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement  
3465 Bottlenecks. *IEEE Access* 9 (2021), 134457–134502. <https://doi.org/10.1109/ACCESS.2021.3110993>
- 3466 OpenMP. 2022. The OpenMP API Specification for Parallel Programming. <https://www.openmp.org>.
- 3467 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin,  
3468 Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison,  
3469 Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An  
3470 Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems*,  
3471 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc.  
3472 <https://proceedings.neurips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf>
- 3473 Adam Paszke, Daniel D. Johnson, David Duvenaud, Dimitrios Vytiniotis, Alexey Radul, Matthew J. Johnson, Jonathan Ragan-  
3474 Kelley, and Dougal Maclaurin. 2021. Getting to the Point: Index Sets and Parallelism-Preserving Autodiff for Pointful  
3475 Array Programming. *Proc. ACM Program. Lang.* 5, ICFP, Article 88 (aug 2021), 29 pages. <https://doi.org/10.1145/3473593>
- 3476 S.J. Pennycook, J.D. Sewall, and V.W. Lee. 2019. Implications of a metric for performance portability. *Future Generation  
3477 Computer Systems* 92 (2019), 947–958. <https://doi.org/10.1016/j.future.2017.08.007>
- 3478 Phitchaya Mangpo Phothilimthana, Archibald Samuel Elliott, An Wang, Abhinav Jangda, Bastian Hagedorn, Henrik Barthels,  
3479 Samuel J. Kaufman, Vinod Grover, Emina Torlak, and Rastislav Bodik. 2019. Swizzle Inventor: Data Movement Synthesis  
3480 for GPU Kernels. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming  
3481 Languages and Operating Systems (ASPLOS '19)*. Association for Computing Machinery, New York, NY, USA, 65–78.  
3482 <https://doi.org/10.1145/3297858.3304059>

- 3480 Federico Pizzuti, Michel Steuwer, and Christophe Dubach. 2020. Generating Fast Sparse Matrix Vector Multiplication from a  
 3481 High Level Generic Functional IR. In *Proceedings of the 29th International Conference on Compiler Construction (CC 2020)*.  
 3482 Association for Computing Machinery, New York, NY, USA, 85–95. <https://doi.org/10.1145/3377555.3377896>
- 3483 Victor Podlozhnyuk. 2007. Image Convolution with CUDA. *NVIDIA Corporation White Paper* (2007).
- 3484 M. Puschel, J.M.F. Moura, J.R. Johnson, D. Padua, M.M. Veloso, B.W. Singer, Jianxin Xiong, F. Franchetti, A. Gacic, Y.  
 3485 Voronenko, K. Chen, R.W. Johnson, and N. Rizzolo. 2005. SPIRAL: Code Generation for DSP Transforms. *Proc. IEEE* 93, 2  
 3486 (2005), 232–275. <https://doi.org/10.1109/JPROC.2004.840306>
- 3487 Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédéric Durand, and Saman Amarasinghe. 2013.  
 3488 Halide: A Language and Compiler for Optimizing Parallelism, Locality, and Recomputation in Image Processing Pipelines.  
 3489 In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '13)*.  
 3490 Association for Computing Machinery, New York, NY, USA, 519–530. <https://doi.org/10.1145/2491956.2462176>
- 3491 Ari Rasch. 2022. (De/Re)-Composition of Data-Parallel Computations  
 3492 via Multi-Dimensional Homomorphisms – A Deep Dive into the MDH Formalism. (2022). <http://arxiv.org/abs/1911.11313>
- 3493 Ari Rasch and Sergei Gorlatch. 2016. Multi-Dimensional Homomorphisms and Their Implementation in OpenCL. In  
 3494 *International Workshop on High-Level Parallel Programming and Applications (HLPP)*. 101–119.
- 3495 Ari Rasch, Richard Schulze, and Sergei Gorlatch. 2019a. Generating Portable High-Performance Code via Multi-Dimensional  
 3496 Homomorphisms. In *2019 28th International Conference on Parallel Architectures and Compilation Techniques (PACT)*.  
 3497 354–369. <https://doi.org/10.1109/PACT.2019.00035>
- 3498 Ari Rasch, Richard Schulze, and Sergei Gorlatch. 2020a. `md_poly`: A Performance-Portable Polyhedral Compiler based on  
 3499 Multi-Dimensional Homomorphisms. In *Proceedings of the International Workshop on Polyhedral Compilation Techniques  
 3500 (IMPACT'20)*. 1–4.
- 3501 Ari Rasch, Richard Schulze, and Sergei Gorlatch. 2020b. `md_poly`: A Performance-Portable Polyhedral Compiler based on  
 3502 Multi-Dimensional Homomorphisms. In *ACM SRC Grand Finals Candidates, 2019 - 2020*. 1–5.
- 3503 Ari Rasch, Richard Schulze, Waldemar Gorus, Jan Hiller, Sebastian Bartholomäus, and Sergei Gorlatch. 2019b. High-  
 3504 Performance Probabilistic Record Linkage via Multi-Dimensional Homomorphisms. In *Proceedings of the 34th ACM/SI-  
 3505 GAPP Symposium on Applied Computing (SAC '19)*. Association for Computing Machinery, New York, NY, USA, 526–533.  
 3506 <https://doi.org/10.1145/3297280.3297330>
- 3507 Ari Rasch, Richard Schulze, Denys Shabalin, Anne Elster, Sergei Gorlatch, and Mary Hall. 2023. (De/Re)-Compositions  
 3508 Expressed Systematically via MDH-Based Schedules. In *Proceedings of the 32nd ACM SIGPLAN International Conference  
 3509 on Compiler Construction (CC 2023)*. Association for Computing Machinery, New York, NY, USA, 61–72. <https://doi.org/10.1145/3578360.3580269>
- 3510 Ari Rasch, Richard Schulze, Michel Steuwer, and Sergei Gorlatch. 2021. Efficient Auto-Tuning of Parallel Programs with  
 3511 Interdependent Tuning Parameters via Auto-Tuning Framework (ATF). *ACM Trans. Archit. Code Optim.* 18, 1, Article 1  
 3512 (Jan. 2021), 26 pages. <https://doi.org/10.1145/3427093>
- 3513 Chandan Reddy, Michael Kruse, and Albert Cohen. 2016. Reduction Drawing: Language Constructs and Polyhedral  
 3514 Compilation for Reductions on GPU. In *Proceedings of the 2016 International Conference on Parallel Architectures and  
 3515 Compilation (PACT '16)*. Association for Computing Machinery, New York, NY, USA, 87–97. <https://doi.org/10.1145/2967938.2967950>
- 3516 Toomas Remmelg, Thibaut Lutz, Michel Steuwer, and Christophe Dubach. 2016. Performance Portable GPU Code Generation  
 3517 for Matrix Multiplication. In *Proceedings of the 9th Annual Workshop on General Purpose Processing Using Graphics  
 3518 Processing Unit (GPGPU '16)*. Association for Computing Machinery, New York, NY, USA, 22–31. <https://doi.org/10.1145/2884045.2884046>
- 3519 Bertrand Russell. 2020. *The principles of mathematics*. Routledge.
- 3520 Bruce Sagan. 2001. *The symmetric group: representations, combinatorial algorithms, and symmetric functions*. Vol. 203. Springer  
 3521 Science & Business Media.
- 3522 Caio Salvador Rohwedder, Nathan Henderson, João P. L. De Carvalho, Yufei Chen, and José Nelson Amaral. 2023. To Pack  
 3523 or Not to Pack: A Generalized Packing Analysis and Transformation. In *Proceedings of the 21st ACM/IEEE International  
 3524 Symposium on Code Generation and Optimization (CGO 2023)*. Association for Computing Machinery, New York, NY,  
 3525 USA, 14–27. <https://doi.org/10.1145/3579990.3580024>
- 3526 Richard Schulze. 2024. MDH Python Compiler. <https://github.com/mdh-lang>.
- 3527 Karen Simonyan and Andrew Zisserman. 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition.  
 3528 <https://doi.org/10.48550/ARXIV.1409.1556>
- 3529 Paul Springer and Paolo Bientinesi. 2016. Design of a high-performance GEMM-like Tensor-Tensor Multiplication. *CoRR*  
 3530 (2016). arXiv:quant-ph/1607.00145 <http://arxiv.org/abs/1607.00145>
- 3531 Michel Steuwer, Christian Fensch, Sam Lindley, and Christophe Dubach. 2015. Generating Performance Portable Code  
 3532 Using Rewrite Rules: From High-Level Functional Expressions to High-Performance OpenCL Code. In *Proceedings of  
 3533 the 20th ACM SIGPLAN International Conference on Functional Programming (ICFP 2015)*. Association for Computing  
 3534

- 3529 Machinery, New York, NY, USA, 205–217. <https://doi.org/10.1145/2784731.2784754>
- 3530 Michel Steuwer, Philipp Kegel, and Sergei Gorlatch. 2011. SkelCL - A Portable Skeleton Library for High-Level GPU  
3531 Programming. In *2011 IEEE International Symposium on Parallel and Distributed Processing Workshops and Phd Forum*.  
1176–1182. <https://doi.org/10.1109/IPDPS.2011.269>
- 3532 Michel Steuwer, Thomas Koehler, Bastian Köpcke, and Federico Pizzuti. 2022. RISE & Shine: Language-Oriented Compiler  
3533 Design. *CoRR* abs/2201.03611 (2022). arXiv:2201.03611 <https://arxiv.org/abs/2201.03611>
- 3534 Michel Steuwer, Toomas Remmelg, and Christophe Dubach. 2017. LIFT: A functional data-parallel IR for high-performance  
3535 GPU code generation. In *2017 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 74–85.  
<https://doi.org/10.1109/CGO.2017.7863730>
- 3536 StreamHPC. 2016. Comparing Syntax for CUDA, OpenCL and HiP. <https://streamhpc.com/blog/2016-04-05-comparing-syntax-cuda-opencl-hip/>.
- 3537 Yifan Sun, Nicolas Bohm Agostini, Shi Dong, and David R. Kaeli. 2019. Summarizing CPU and GPU Design Trends with  
3538 Product Data. *CoRR* abs/1911.11313 (2019). arXiv:1911.11313 <http://arxiv.org/abs/1911.11313>
- 3539 Joel Svensson, Mary Sheeran, and Koen Claessen. 2011. Obsidian: A Domain Specific Embedded Language for Parallel  
3540 Programming of Graphics Processors. In *Implementation and Application of Functional Languages*, Sven-Bodo Scholz and  
3541 Olaf Chitil (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 156–173.
- 3542 Walid Taha and Tim Sheard. 1997. Multi-Stage Programming with Explicit Annotations. In *Proceedings of the 1997 ACM  
3543 SIGPLAN Symposium on Partial Evaluation and Semantics-Based Program Manipulation (PEPM '97)*. Association for  
3544 Computing Machinery, New York, NY, USA, 203–217. <https://doi.org/10.1145/258993.259019>
- 3545 TensorFlow. 2022a. MobileNet v1 models for Keras. <https://github.com/keras-team/keras/blob/master/keras/applications/mobilenet.py>.
- 3546 TensorFlow. 2022b. ResNet models for Keras. <https://github.com/keras-team/keras/blob/master/keras/applications/resnet.py>.
- 3547 TensorFlow. 2022c. VGG16 model for Keras. <https://github.com/keras-team/keras/blob/master/keras/applications/vgg16.py>.
- 3548 Philippe Tillet and David Cox. 2017. Input-Aware Auto-Tuning of Compute-Bound HPC Kernels. In *Proceedings of the  
3549 International Conference for High Performance Computing, Networking, Storage and Analysis (SC '17)*. Association for  
3550 Computing Machinery, New York, NY, USA, Article 43, 12 pages. <https://doi.org/10.1145/3126908.3126939>
- 3551 TIOBE. 2023. The Software Quality Company. <https://www.tiobe.com/tiobe-index/>.
- 3552 TOPLAS Artifact. 2022. [https://gitlab.com/mdh-project/toplas22\\_artifact](https://gitlab.com/mdh-project/toplas22_artifact).
- 3553 Uday Bondhugula. 2022. Pluto: An automatic polyhedral parallelizer and locality optimizer. <https://github.com/bondhugula/pluto>, commit = 12e075a, date = 31.10.2021.
- 3554 Nicolas Vasilache, Oleksandr Zinenko, Aart J. C. Bik, Mahesh Ravishankar, Thomas Raoux, Alexander Belyaev, Matthias  
3555 Springer, Tobias Gysi, Diego Caballero, Stephan Herhut, Stella Laurenzo, and Albert Cohen. 2022. Composable  
3556 and Modular Code Generation in MLIR: A Structured and Retargetable Approach to Tensor Compiler Construction.  
arXiv:cs.PL/2202.03293
- 3557 Nicolas Vasilache, Oleksandr Zinenko, Theodoros Theodoridis, Priya Goyal, Zachary Devito, William S. Moses, Sven  
3558 Verdoolaege, Andrew Adams, and Albert Cohen. 2019. The Next 700 Accelerated Layers: From Mathematical Expressions  
3559 of Network Computation Graphs to Accelerated GPU Kernels, Automatically. *ACM Trans. Archit. Code Optim.* 16, 4,  
3560 Article 38 (Oct. 2019), 26 pages. <https://doi.org/10.1145/3355606>
- 3561 Sven Verdoolaege, Juan Carlos Juega, Albert Cohen, José Ignacio Gómez, Christian Tenllado, and Francky Catthoor. 2013.  
3562 Polyhedral Parallel Code Generation for CUDA. *ACM Trans. Archit. Code Optim.* 9, 4, Article 54 (Jan. 2013), 23 pages.  
<https://doi.org/10.1145/2400682.2400713>
- 3563 Sven Verdoolaege and Tobias Grosser. 2012. Polyhedral Extraction Tool. In *International Workshop on Polyhedral Compilation  
3564 Techniques (IMPACT'12), Paris, France*, Vol. 141.
- 3565 J. von Neumann. 1925. Eine Axiomatisierung der Mengenlehre. 1925, 154 (1925), 219–240. <https://doi.org/doi:10.1515/crll.1925.154.219>
- 3566 Mohamed Wahib and Naoya Maruyama. 2014. Scalable Kernel Fusion for Memory-Bound GPU Applications. In *SC '14:  
3567 Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 191–202.  
3568 <https://doi.org/10.1109/SC.2014.21>
- 3569 Bram Wasti, José Pablo Cambronero, Benoit Steiner, Hugh Leather, and Aleksandar Zlateski. 2022. LoopStack: a Lightweight  
3570 Tensor Algebra Compiler Stack. <https://doi.org/10.48550/ARXIV.2205.00618>
- 3571 R.C. Whaley and J.J. Dongarra. 1998. Automatically Tuned Linear Algebra Software. In *SC '98: Proceedings of the 1998  
3572 ACM/IEEE Conference on Supercomputing*. 38–38. <https://doi.org/10.1109/SC.1998.10004>
- 3573 Maurice V Wilkes. 2001. The memory gap and the future of high performance memories. *ACM SIGARCH Computer  
3574 Architecture News* 29, 1 (2001), 2–7.
- 3575 M.E. Wolf and M.S. Lam. 1991. A loop transformation theory and an algorithm to maximize parallelism. *IEEE Transactions  
3576 on Parallel and Distributed Systems* 2, 4 (1991), 452–471. <https://doi.org/10.1109/71.97902>
- 3577

- 3578 Hongwei Xi and Frank Pfenning. 1999. Dependent Types in Practical Programming. In *Proceedings of the 26th ACM*  
3579 *SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '99)*. Association for Computing Machinery,  
3580 New York, NY, USA, 214–227. <https://doi.org/10.1145/292540.292560>
- 3581 Rohan Yadav, Alex Aiken, and Fredrik Kjolstad. 2022. DISTAL: The Distributed Tensor Algebra Compiler. In *Proceedings*  
3582 *of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI 2022)*.  
3583 Association for Computing Machinery, New York, NY, USA, 286–300. <https://doi.org/10.1145/3519939.3523437>
- 3584 Cambridge Yang, Eric Atkinson, and Michael Carbin. 2021. Simplifying Dependent Reductions in the Polyhedral Model.  
3585 *Proc. ACM Program. Lang.* 5, POPL, Article 20 (Jan. 2021), 33 pages. <https://doi.org/10.1145/3434301>
- 3586 Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo,  
3587 Koushik Sen, Joseph E. Gonzalez, and Ion Stoica. 2020a. AnsoR: Generating High-Performance Tensor Programs for Deep  
3588 Learning. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. USENIX Association,  
3589 863–879. <https://www.usenix.org/conference/osdi20/presentation/zheng>
- 3590 Size Zheng, Yun Liang, Shuo Wang, Renze Chen, and Kaiwen Sheng. 2020b. *FlexTensor: An Automatic Schedule Exploration*  
3591 *and Optimization Framework for Tensor Computation on Heterogeneous System*. Association for Computing Machinery,  
3592 New York, NY, USA, 859–873. <https://doi.org/10.1145/3373376.3378508>
- 3593
- 3594
- 3595
- 3596
- 3597
- 3598
- 3599
- 3600
- 3601
- 3602
- 3603
- 3604
- 3605
- 3606
- 3607
- 3608
- 3609
- 3610
- 3611
- 3612
- 3613
- 3614
- 3615
- 3616
- 3617
- 3618
- 3619
- 3620
- 3621
- 3622
- 3623
- 3624
- 3625
- 3626

3627 **APPENDIX**

3628 Our appendix provides details for the interested reader, which should not be required for under-  
 3629 standing the basic ideas presented in this paper.

3630

3631 **A MATHEMATICAL FOUNDATION**

3632 We rely on a set theoretical foundation, based on ZFC set theory [Ciesielski 1997]. We avoid class  
 3633 theory, such as NBG [von Neumann 1925], by assuming, for example, that our universe of types  
 3634 contains all relevant representatives (int, float, struct, etc), but is not the "class of all types".  
 3635 Thereby, we avoid fundamental issues [Russell 2020] which are not relevant for this work.

3636

3637 **A.1 Family**

3638 **Definition 16** (Family). Let  $I$  and  $A$  be two sets. A *family*  $F$  from  $I$  to  $A$  is any set

3639

$$F := \{ (i, a) \mid i \in I \wedge a \in A \}$$

3640 such that the following two properties are satisfied:

3641

- *left-total*:  $\forall i \in I : \exists a \in A : (i, a) \in F$
- *right-unique*:  $(i, a) \in F \wedge (i, a') \in F \Rightarrow a = a'$

3642 We refer to  $I$  also as *index set* of family  $F$  and to  $A$  as  $F$ 's *image set*. If  $I$  has a strict total order  $<$ , we  
 3643 refer to  $F$  also as *ordered family*.

3644

3645 **Notation 4** (Family). Let  $F$  be a family from  $I$  to  $A$ .

3646

3647 We write:

- $F_i$  for the unique  $a \in A$  such that  $(i, a) \in F$ ;
- $(F_i)_{i \in I}$  instead of  $F$  to explicitly state  $F$ 's index and image sets in our notation;
- $(F_{i_1, \dots, i_n})_{i_1 \in I_1, \dots, i_n \in I_n}$  instead of  $(\dots (F_{i_1, \dots, i_n})_{i_n \in I_n} \dots)_{i_1 \in I_1}$ .

3648 Alternatively, depending on the context, we use the following notation:

3649

- $F^{<i>}$  instead of  $F_i$ ;
- $(F_i)^{<i \in I>}$  instead of  $(F_i)_{i \in I}$ ;
- $(F_{i_1, \dots, i_n})^{<i_1 \in I_1 | \dots | i_n \in I_n>}$  instead of  $(F_{i_1, \dots, i_n})_{i_1 \in I_1, \dots, i_n \in I_n}$ .

3650

3651 For nested families, each index set  $I_k$  may depend on the earlier-defined values  $i_1, \dots, i_{k-1}$  (not  
 3652 explicitly stated above for brevity).

3653

3654 **Definition 17** (Tuple). We identify  $n$ -tuples as families that have index set  $[1, n]_{\mathbb{N}}$ .

3655

3656 **Example 12** (Tuple). A 2-tuple  $(a, b)$  (a.k.a *pair*) is a family  $(F_i)_{i \in I := [1, 2]_{\mathbb{N}}}$  for which  $F_1 = a$  and  
 3657  $F_2 = b$ .

3658

3659 **A.2 Scalar Types**

3660

3661 We denote by

3662

$$\text{TYPE} := \{ \text{int}, \text{int8}, \text{int16}, \dots, \text{float}, \text{double}, \dots, \text{struct}, \dots \}$$

3663

3664 our set of *scalar types*, where *int8* and *int16* represent 8-bit/16-bit integer numbers, *float* and  
 3665 *double* are the types of single/double precision floating point numbers (IEEE 754 standards),  
 3666 *structs* contain a fixed set of other scalar types, etc. For simplicity, we interpret integer types  
 3667 (*int*, *int8*, *int16*, ...) uniquely as integers  $\mathbb{Z}$ , floating point number types (*float* and *double*) as  
 3668 rationale numbers  $\mathbb{Q}$ , etc.

3669

3670 For high flexibility, we avoid fixing *TYPE* to a particular set of scalar types, i.e., we assume that  
 3671 *TYPE* contains all practice-relevant types. This is legal, because our formalism makes no assumptions  
 3672 on the number and kinds of scalar types.

3673

We consider operations on scalar types (addition, multiplication, etc) to be: 1) *atomic*: we do not aim at parallelizing or otherwise optimizing operations on scalar values in this work; 2) *size preserving*: we assume that all values of a scalar type have the same arbitrary but fixed size.

Note that we can potentially also define, for example, the set of arbitrarily sized matrices  $\{T^{m \times n} \mid m, n \in \mathbb{N}, T \in \text{TYPE}\}$  as scalar type in our approach. However, this would prevent any kind of formal reasoning about performance and type correctness of matrix-related operations (e.g., matrix multiplication), like parallelization (due to our atomic assumption above) or type correctness (e.g., assuring in matrix multiplication that number of columns of the first input matrix coincides with and number of rows of the second matrix: due to our size preservation assumption above, we would not be able to distinguish matrices based on their sizes).

### A.3 Functions

**Definition 18** (Function). Let  $A \in \text{TYPE}$  and  $B \in \text{TYPE}$  be two scalar types.

A (total) function  $f$  is a tuple of the form

$$f \in \{ (\underbrace{(A, B)}_{\text{function type}}, \underbrace{G_f}_{\text{function graph}}) \mid G_f \subseteq \{ (a, b) \mid a \in A \wedge b \in B \} \}$$

that satisfies the following two properties:

- *left-total*:  $\forall a \in A : \exists b \in B : (a, b) \in G_f$ ;
- *right-unique*:  $(a, b) \in G_f \wedge (a, b') \in G_f \Rightarrow b = b'$ .

We write  $f(a)$  for the unique  $b \in B$  such that  $(a, b) \in G_f$ . Moreover, we denote  $f$ 's function type as  $A \rightarrow B$ , and we write  $f : A \rightarrow B$  to state that  $f$  has function type  $A \rightarrow B$ .

We refer to:

- $A$  as the *domain* of  $f$
- $B$  as the *co-domain* (or *range*) of  $f$
- $(A, B)$  as the *type* of  $f$
- $G_f$  as the *graph* of  $f$

If  $f$  does not satisfy the left total property, we say  $f$  is *partial*, and we denote  $f$ 's type as  $f : A \rightarrow_p B$  (where  $\rightarrow$  is replaced by  $\rightarrow_p$ ).

We allow functions to have so-called *dependent types* [Xi and Pfenning 1999] for expressive typing. For example, dependent types enable encoding the sizes of families into the type system, which contributes to better error checking. We refer to dependently typed functions as *meta-functions*, as outlined in the following.

**Definition 19** (Meta-Function). We refer to any family of functions

$$(f^{<i>} : A^{<i>} \rightarrow B^{<i>})^{<i \in I>}$$

as *meta-function*. In the context of meta-functions, we refer to index  $i \in I$  also as *meta-parameter*, to index set  $I$  as *meta-parameter type*, to  $A^{<i \in I>}$  and  $B^{<i \in I>}$  as *meta-types* (as both are generic in meta-parameter  $i \in I$ ), and to  $A^{<i>} \rightarrow B^{<i>}$  for concrete  $i$  as *meta-function*  $f$ 's ordinary function type.

In the following, we often write:

- $f^{<i \in I>} : A^{<i>} \rightarrow B^{<i>}$  instead of  $(f^{<i>} : A^{<i>} \rightarrow B^{<i>})^{<i \in I>}$ ;
- $f^{<i>} : A' \rightarrow B^{<i>}$  (or  $f^{<i>} : A^{<i>} \rightarrow B'$ ) iff  $A^{<i>} = A'$  (or  $B^{<i>} = B'$ ) for all  $i \in I$ .

3725 We use *multi-stage meta-functions* as a concept analogous to *multi staging* [Taha and Sheard  
 3726 1997] in programming and similar to *currying* in mathematics.

3727 **Definition 20** (Multi-Stage Meta-Function). A *multi-stage meta-function* is a nested family of  
 3728 functions:

$$3730 \quad \overbrace{f^{<i_1 \in I_1^{<>} > \dots <i_S \in I_S^{<i_1, \dots, i_{S-1}>} >}}^{\text{stage 1}} : \underbrace{A^{<i_1, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_S>}}_{\text{stage } S} \\ 3731 \quad \text{function instance}$$

3734 Here,  $I_s^{<i_1, \dots, i_{s-1}>}$ ,  $s \in [1, S]_{\mathbb{N}}$ , is the meta-parameter type on stage  $s$ , which may depend on all  
 3735 meta-parameters of the previous stages  $i_1, \dots, i_{s-1}$ . We refer to such meta-functions also as *S-stage*  
 3736 *meta-functions*, and we denote their type also as

$$3737 \quad f^{<i_1 \in I_1 | \dots | i_S \in I_S>} : A^{<i_1, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_S>}$$

3739 and access to them as

$$3740 \quad f^{<i_1 | \dots | i_S>} (x)$$

3741 where different stages are separated by vertical bars.

3743 We allow partially applying parameters (meta and ordinary) of meta-functions.

3745 **Definition 21** (Partial Meta-Function Application). Let

$$3746 \quad f^{<i_1 \in I_1 | \dots | i_S \in I_S>} : A^{<i_1, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_S>}$$

3748 be a meta-function (meta-parameters of meta-types  $I_1, \dots, I_S$  omitted for brevity).

- 3749 • The *partial application* of meta-function  $f$  on stage  $s$  to meta-parameter  $\hat{i}_s$  is the meta-function

$$3750 \quad f'^{<i_1 \in \hat{I}_1 | \dots | i_{s-1} \in \hat{I}_{s-1} | i_{s+1} \in I_{s+1} | \dots | i_S \in I_S>} : A^{<i_1, \dots, i_{s-1}, \hat{i}_s, i_{s+1}, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_{s-1}, \hat{i}_s, i_{s+1}, \dots, i_S>}$$

3751 where  $\hat{I}_1 \subseteq I_1, \dots, \hat{I}_{s-1} \subseteq I_{s-1}$  are the largest sets such that  $\hat{i}_s \in I_s^{<i_1, \dots, i_{s-1}>}$  for all  $i_1 \in \hat{I}_1, \dots, i_{s-1} \in \hat{I}_{s-1}$ . The function is defined as:

$$3754 \quad f'^{<i_1 | \dots | i_{s-1} | i_{s+1} | \dots | i_S>} (x) := f^{<i_1 | \dots | i_{s-1} | \hat{i}_s | i_{s+1} | \dots | i_S>} (x)$$

3756 We write for  $f'$ 's type also

$$3757 \quad f'^{<i_1 \in \hat{I}_1 | \dots | i_{s-1} \in \hat{I}_{s-1} | \hat{i}_s | i_{s+1} \in I_{s+1} | \dots | i_S \in I_S>} : A^{<i_1, \dots, i_{s-1}, \hat{i}_s, i_{s+1}, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_{s-1}, \hat{i}_s, i_{s+1}, \dots, i_S>}$$

3759 where  $f'$  is replaced by  $f$  and  $i_s \in I_s$  is replaced by the concrete value  $\hat{i}_s$ .

- 3760 • The *partial application* of meta-function  $f$  to ordinary parameter  $x$  is the meta-function

$$3761 \quad f'^{<i_1 \in \hat{I}_1 | \dots | i_S \in \hat{I}_S>} : \underbrace{B_1^{<i_1, \dots, i_S>} \rightarrow B_2^{<i_1, \dots, i_S>}}_{:= B^{<i_1, \dots, i_S>}}$$

3764 where  $\hat{I}_1 \subseteq I_1, \dots, \hat{I}_S \subseteq I_S$  are the largest sets such that  $x \in A^{<i_1, \dots, i_S>}$  for all  $i_1 \in \hat{I}_1, \dots, i_S \in \hat{I}_S$ .

3765 The function is defined as:

$$3766 \quad f'^{<i_1 | \dots | i_S>} (x') := f^{<i_1 | \dots | i_S>} (x)(x')$$

3768 We allow *generalizing* meta-parameters. For example, by generalizing meta-parameters for input  
 3769 sizes, we can use functions on arbitrarily sized inputs (a.k.a. *dynamic size* in programming). In our  
 3770 generated code, meta-parameters are available at compile time such that concrete meta-parameter  
 3771 values can be exploited for generating well-performing code (e.g., for static loop boundaries).  
 3772 Consequently, generalization increases expressivity of the generated code, e.g., by being able to  
 3773

3774 process differently sized inputs without requiring re-compilation for unseen input sizes, but usually  
 3775 at the cost of performance.

3776 **Definition 22** (Generalized Meta-Parameters). Let

$$3778 f^{<i_1 \in I_1 | \dots | i_s \in I_s | \dots | i_S \in I_S>} : A^{<i_1, \dots, i_s, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_s, \dots, i_S>}$$

3779 be a meta-function (meta-parameters of  $I_1, \dots, I_S$  omitted for brevity) such that

$$3780 f^{<i_1 | \dots | i_s | \dots | i_S>} (x) = f^{<i_1 | \dots | i'_s | \dots | i_S>} (x)$$

3782 i.e.,  $f$ 's behavior is invariant under different values of meta-parameter  $i_s$  in stage  $s$ .

3783 The *generalization* of  $f$  in meta-parameter  $s \in [1, S]_{\mathbb{N}}$  is the meta-function

$$3784 f'^{<i_1 \in I_1 | \dots | i_{s-1} \in I_{s-1} | i_{s+1} \in I_{s+1} | \dots | i_S \in I_S>} :$$

$$3786 \bigcup_{i_s \in I_s^{<i_1, \dots, i_{s-1}>}} A^{<i_1 | \dots | i_{s-1} | i_s | i_{s+1} | \dots | i_S>} \rightarrow \bigcup_{i_s \in I_s^{<i_1, \dots, i_{s-1}>}} B^{<i_1 | \dots | i_{s-1} | i_s | i_{s+1} | \dots | i_S>}$$

3788 which is defined as:

$$3789 f'^{<i_1 | \dots | i_{s-1} | i_{s+1} | \dots | i_S>} (x) := f^{<i_1 | \dots | i_s | i_{s+1} | \dots | i_S>} (x)$$

3791 for an arbitrary  $i_s \in I_s$  such that  $x \in A^{<i_1 | \dots | i_{s-1} | i_s | i_{s+1} | \dots | i_S>}$ .

3792 We write for  $f$ 's type also

$$3794 f^{<i_1 \in I_1 | \dots | i_{s-1} \in I_{s-1} | * \in I_s | i_{s+1} \in I_{s+1} | \dots | i_S \in I_S>} : A^{<i_1, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_S>}$$

3795 where  $i_s$  is replaced by  $*$ , and for access to  $f$

$$3797 f^{<i_1 | \dots | i_{s-1} | * | i_{s+1} | \dots | i_S>} (x)$$

3798 We use *postponed meta-parameters* to change the order of meta-parameters of already defined  
 3799 meta-functions. For example, we use postponed meta-parameters in Definition 7 to compute the  
 3800 values of meta-parameters based on the particular meta-parameter values of later stages.

3802 **Definition 23** (Postponed Meta-Parameters). Let

$$3803 f^{<i_1 \in I_1 | \dots | i_s \in I_s | \dots | i_S \in I_S>} : A^{<i_1, \dots, i_s, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_s, \dots, i_S>}$$

3804 be a meta-function (meta-parameters of  $I_1, \dots, I_S$  omitted via ellipsis for brevity) such that for each  
 3805  $k \in (s, S]_{\mathbb{N}}$ , it holds:

$$3807 I_k^{<i_1 | \dots | i_s | \dots | i_{k-1}>} = I_k^{<i_1 | \dots | i'_s | \dots | i_{k-1}>}$$

3808 i.e., the  $I_k$  are invariant under different values of meta-parameter  $i_s$  in stage  $s$ , such that  $i_s$  can be  
 3809 ignored in the parameter list of  $I_k$ .

3810 Function  $f'$  is function  $f$  *postponed* on stage  $s$  to meta-type

$$3812 \hat{I}_s^{<i_1 | \dots | i_{s-1} | i_{s+1} | \dots | i_S>} \subseteq I_s^{<i_1 | \dots | i_{s-1}>}$$

3813 which, in contrast to  $I_s$ , may also depend on meta-parameter values  $i_{s+1}, \dots, i_S$ , iff  $f'$  is of type

$$3815 f'^{<i_1 \in I_1^{<\dots>} | \dots | i_{s-1} \in I_{s-1}^{<\dots>} | i_{s+1} \in I_{s+1}^{<\dots>} | \dots | i_S \in I_S^{<\dots>} > < i_s \in \hat{I}_s^{<i_1, \dots, i_S>} >} : A^{<i_1, \dots, i_s, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_s, \dots, i_S>}$$

3816 and defined as:

$$3817 f'^{<i_1 | \dots | i_{s-1} | i_{s+1} | \dots | i_S > < i_s >} (a) = f^{<i_1 | \dots | i_{s-1} | i_s | i_{s+1} | \dots | i_S>} (a)$$

3819 We write for  $f'$ 's type also

$$3820 f^{<i_1 \in I_1^{<\dots>} | \dots | i_{s-1} \in I_{s-1}^{<\dots>} | \rightarrow | i_{s+1} \in I_{s+1}^{<\dots>} | \dots | i_S \in I_S^{<\dots>} > < i_s \in \hat{I}_s^{<\dots>} >} : A^{<i_1, \dots, i_s, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_s, \dots, i_S>}$$

3822

3823 where  $f'$  is replaced by  $f$  and  $i_s$  by symbol " $\rightarrow$ ". For access to  $f'$ , we write

$$3824 \quad 3825 \quad f^{<i_1 | \dots | i_{s-1} | \rightarrow | i_{s+1} | \dots | i_s>} (x)$$

3826 When using a binary function for combining a family of elements, we often use the following  
 3827 notation.

3829 **Notation 5** (Iterative Function Application). Let  $\otimes : T \times T \rightarrow T$  be an arbitrary associative and  
 3830 commutative binary function on scalar type  $T \in \text{TYPE}$ . Let further  $x$  be an arbitrary family that has  
 3831 index set  $I := \{i_1, \dots, i_N\}$  and image set  $\{x_i\}_{i \in I} \subseteq T$ .

3832 We write  $\otimes_{i \in I} x_i$  instead of  $x_{i_1} \otimes \dots \otimes x_{i_N}$  (infix notation).

#### 3834 A.4 MatVec Expressed in MDH DSL

3836 Our MatVec example from Figure 6 is expressed in MDH's high-level *Domain-Specific Language (DSL)*,  
 3837 used as input by our proof-of-concept MDH compiler, as follows:

```
3839
3840 1 MatVec<T in TYPE | I,K in IN> := 
3841 2     out_view<T>( w:(i,k)->(i) ) o
3842 3     md_hom<I,K>( *, (++,+) ) o
3843 4     inp_view<T,T>( M:(i,k)->(i,k) ,
3844 5             v:(i,k)->(k) )
```

3844 Listing 5. MatVec expressed in MDH's DSL

3845 Our compiler takes an expression as in Listing 5 as input, and it generates auto-tunable code from  
 3846 it, according to the methods presented in this paper (particularly in Section F). Our compile times  
 3847 currently vary for our examples in Figure 14 from < 1 Minute for low-dimensional examples (e.g.,  
 3848 linear algebra routines) to 5 Minutes for high-dimensional examples (such as quantum chemistry  
 3849 computations). We aim to significantly reduce our compile times and to describe our code generation  
 3850 process in our future work, as code generation is not the focus of this paper.

3851 In addition, we also offer an implementation of our MDH approach and its high-level DSL  
 3852 embedded in the *Python* programming language [Schulze 2024], as Python is becoming increasingly  
 3853 popular in both academia and industry [TIOBE 2023].

## 3854 B ADDENDUM SECTION 2

### 3855 B.1 Multi-Dimensional Array

3856 **Definition 24** (Multi-Dimensional Array). Let  $\text{MDA-IDX-SETS} := \{I \subset \mathbb{N}_0 \mid |I| < \infty\}$  be the set of all  
 3857 finite subsets of natural numbers, to which we also refer to as set of *MDA index sets* in the context of  
 3858 MDAs. Let further  $T \in \text{TYPE}$  be an arbitrary scalar type,  $D \in \mathbb{N}$  a natural number,  $I := (I_1, \dots, I_D) \in$   
 3859  $\text{MDA-IDX-SETS}^D$  a tuple of  $D$ -many MDA index sets, and  $N := (N_1, \dots, N_D) := (|I_1|, \dots, |I_D|)$  the  
 3860 tuple of index sets' sizes.

3861 A *Multi-Dimensional Array (MDA)*  $\alpha$  that has *dimensionality D*, *size N*, *index sets I*, and *scalar*  
 3862 *type T* is a function with the following signature:

$$3863 \quad \alpha : I_1 \times \dots \times I_D \rightarrow T$$

3864 We refer to  $I_1 \times \dots \times I_D \rightarrow T$  as the *type* of MDA  $\alpha$ .

## 3872 B.2 Combine Operators

Technically, combine operators are functions that take as input two MDAs and yield a single MDA as their output (formal definition follows soon). Per definition, we require that the input MDAs' index sets coincide in all dimensions except in the dimension in which we combine; thereby, we catch undefined cases already at the type level, e.g., trying to concatenate improperly sized MDAs:

$$\underbrace{\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}}_{3 \times 3\text{-many elements}} + \underbrace{\begin{bmatrix} 11 & 12 \\ 13 & 14 \\ 15 & 16 \end{bmatrix}}_{3 \times 2\text{-many elements}} = \underbrace{\begin{bmatrix} 1 & 2 & 3 & 11 & 12 \\ 4 & 5 & 6 & 13 & 14 \\ 7 & 8 & 9 & 15 & 16 \end{bmatrix}}_{3 \times 5\text{-many elements}}$$

3886 Here, on the left, we can reasonably define the concatenation of MDAs that contain  $3 \times 3$ -many  
3887 elements and  $3 \times 2$  elements. However, as indicated in the right part of the figure, it is not possible  
3888 to intuitively concatenate MDAs of sizes  $3 \times 3$  and  $2 \times 2$ , as the MDAs do not match in their number  
3889 of elements in any of the two dimensions.

Figure 29 illustrates combine operators informally using the example operators *concatenation* (left part of the figure) and *point-wise addition* (right part). We illustrate concatenation using the example MDAs  $a^{(1,1)}$  and  $a^{(1,2)}$  from Figure 7; for point-wise addition, we use MDAs  $a^{(2,2)}$  and  $a^{(2,3)}$  from Figure 7 (all MDAs are chosen arbitrarily, and the example works the same for other MDAs).

In the case of concatenation (left part of Figure 29), MDAs  $\alpha^{(1,1)}$  and  $\alpha^{(1,2)}$  coincide in their second dimension  $I_2 := \{0, 1, 2, 3\}$ , which is important, because we concatenate in the first dimension, thus requiring coinciding index sets in all other dimensions (as motivated above). In the case of the point-wise addition example (right part of Figure 29), the example MDAs  $\alpha^{(2,2)}$  and  $\alpha^{(2,3)}$  coincide in their first dimension  $I_1 := \{0, 1\}$ , as required for combining the MDAs in the second dimension. The varying index sets of the four MDAs are denoted as  $P$  and  $Q$  in the figure, which are in the case of the concatenation example, index sets in the first dimension of MDAs  $\alpha^{(1,1)}$  and  $\alpha^{(1,2)}$ ; in the

$\mathfrak{a}^{(1,1)} = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$ $\in T[P := \{0\}, I_2 := \{0,1,2,3\}]$	$\mathfrak{a}^{(1,2)} = \begin{bmatrix} 5 & 6 & 7 & 8 \end{bmatrix}$ $\in T[Q := \{1\}, I_2 := \{0,1,2,3\}]$	$\mathfrak{a}^{(2,2)} = \begin{bmatrix} 3 \\ 7 \end{bmatrix}$ $\in T[I_1 := \{0,1\}, P := \{2\}]$	$\mathfrak{a}^{(2,3)} = \begin{bmatrix} 4 \\ 8 \end{bmatrix}$ $\in T[I_1 := \{0,1\}, Q := \{3\}]$
<b>index set function</b>		<b>index set function</b>	
$\mathfrak{a}_{\downarrow \oplus_1}^{(1,1)} := \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$ $\in T[id(P = \{0\}) = \{0\}, I_2 = \{0,1,2,3\}]$	$\mathfrak{a}_{\downarrow \oplus_1}^{(1,2)} := \begin{bmatrix} 5 & 6 & 7 & 8 \end{bmatrix}$ $\in T[id(Q = \{1\}) = \{1\}, I_2 = \{0,1,2,3\}]$	$\mathfrak{a}_{\rightarrow \oplus_2}^{(2,2)} := \begin{bmatrix} 3 \\ 7 \end{bmatrix}$ $\in T[I_1 = \{0,1\}, 0_f(P = \{2\}) = \{0\}]$	$\mathfrak{a}_{\rightarrow \oplus_2}^{(2,3)} := \begin{bmatrix} 4 \\ 8 \end{bmatrix}$ $\in T[I_1 = \{0,1\}, 0_f(Q = \{3\}) = \{0\}]$
<b>concatenation</b>		<b>point-wise addition</b>	
$\dots \rightarrow \mathfrak{a}^{(1)} := \mathfrak{a}_{\downarrow \oplus_1}^{(1,1)} \oplus_1 \mathfrak{a}_{\downarrow \oplus_1}^{(1,2)} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}$ $\in T[id(P \cup Q = \{0,1\}) = \{0,1\}, I_2 = \{0,1,2,3\}]$		$\dots \rightarrow \mathfrak{a}^{(2)} := \mathfrak{a}_{\rightarrow \oplus_2}^{(2,2)} \oplus_2 \mathfrak{a}_{\rightarrow \oplus_2}^{(2,3)} = \begin{bmatrix} 7 \\ 15 \end{bmatrix}$ $\in T[I_1 = \{0,1\}, 0_f(P \cup Q = \{2,3\}) = \{0\}]$	

Fig. 29. Illustration of *combine operators* using the examples *concatenation* (left) and *point-wise addition* (right)

3921 case of the point-wise addition example, the varying index sets of MDAs  $\alpha^{(2,1)}$  and  $\alpha^{(2,2)}$  belong to  
 3922 the second dimensions.

3923 In the following, we assume w.l.o.g. that that the varying index sets  $P$  and  $Q$  of MDAs to combine  
 3924 are disjoint. Our assumption will not be a limitation for us, because we will apply combine operators  
 3925 always to MDA parts that belong to the same MDA, causing the index sets of the parts to be disjoint  
 3926 by construction. For example, in the case of the concatenation example in Figure 29, the MDA parts  
 3927  $\alpha^{(1,1)}$  and  $\alpha^{(1,2)}$  correspond to the first and second row of the same MDA  $\alpha$  in Figure 7, and in the  
 3928 case of the point-wise addition example, the parts  $\alpha^{(2,2)}$  and  $\alpha^{(2,3)}$  represent the third and fourth  
 3929 column of MDA  $\alpha$ .

3930 We define combine operators based on *index set functions* (also introduced formally soon). Index  
 3931 set functions precisely describe, on the type level, the index set of the combined output MDA  
 3932 and thus how an MDA's index set evolves during combination. For this, an index set function  
 3933 takes as input the input MDA's index set in the dimension to combine, and the function yields  
 3934 as its output the index set of the output MDA which is combined in this dimension. In the case  
 3935 of the concatenation example in Figure 29, the index set function is identity  $id$  and thus trivial.  
 3936 However, in the case of point-wise addition, the corresponding index set function is constant  
 3937 function  $0_f$  which maps any index set to the singleton set  $\{0\}$  containing index 0 only. This is  
 3938 because when combining via point-wise addition, the MDA shrinks in the combined dimension to  
 3939 only one element which we aim to uniformly access via MDA index 0. In Figure 29, we denote MDAs  
 3940  $\alpha^{(1,1)}, \alpha^{(1,2)}, \alpha^{(2,2)}, \alpha^{(2,3)}$  after applying the corresponding index set function as  $\alpha_{\downarrow \otimes_1}^{(1,1)}, \alpha_{\downarrow \otimes_1}^{(1,2)}, \alpha_{\rightarrow \otimes_2}^{(2,2)},$   
 3941  $\alpha_{\rightarrow \otimes_2}^{(2,3)}$ ; the combined MDAs are denoted as  $\alpha^{(1)}$  and  $\alpha^{(2)}$  in the figure. The concatenation operator is  
 3942 denoted in the figure generically as  $\otimes_1$ , and point-wise addition is denoted as  $\otimes_2$ , correspondingly.  
 3943

3944 We now define *combine operators* formally, and we illustrate this formal definition afterwards  
 3945 using the example operators *concatenation* and *point-wise combination*. For the interested reader,  
 3946 details on some technical design decisions of combine operators are outlined in Section B.6.  
 3947

3948 **Definition 25** (Combine Operator). Let  $\text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} := \{ (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} \mid P \cap Q = \emptyset \}$  denote the set of all pairs of MDA index sets that are disjoint. Let  
 3949 further  $\Rightarrow^{\text{MDA}}_{\text{MDA}} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$  be a function on MDA index sets,  $T \in \text{TYPE}$  a scalar  
 3950 type,  $D \in \mathbb{N}$  an MDA dimensionality, and  $d \in [1, D]_{\mathbb{N}}$  an MDA dimension.  
 3951

3952 We refer to any binary function  $\otimes$  of type (parameters in angle brackets are type parameters)  
 3953

3954  $\otimes <(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS}> :$

$$3955 \otimes : \underbrace{T[I_1, \dots, \overbrace{\Rightarrow^{\text{MDA}}_{\text{MDA}}(P), \dots, I_D}^d, \dots, I_D] \times T[I_1, \dots, \overbrace{\Rightarrow^{\text{MDA}}_{\text{MDA}}(Q), \dots, I_D}^d, \dots, I_D]}_{d} \rightarrow T[I_1, \dots, \overbrace{\Rightarrow^{\text{MDA}}_{\text{MDA}}(P \cup Q), \dots, I_D}^d, \dots, I_D]$$

3960 as *combine operator* that has *index set function*  $\Rightarrow^{\text{MDA}}_{\text{MDA}}$ , *scalar type*  $T$ , *dimensionality*  $D$ , and *operating*  
 3961 *dimension*  $d$ . We denote combine operator's type concisely as  $\text{CO}^{<\Rightarrow^{\text{MDA}}_{\text{MDA}} \mid T \mid D \mid d>}$ .  
 3962

3963 Since function  $\otimes$ 's ordinary function type  $T[\dots] \times T[\dots] \rightarrow T[\dots]$  is generic in parameters  
 3964  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D)$  and  $(P, Q)$  (these type parameters are denoted in angle brackets in  
 3965 Definition 25), we refer to function  $\otimes$  as *meta-function*, to the type parameters in angle brackets  
 3966 as *meta-parameters*, and we say *meta-types* to  $T[I_1, \dots, \overbrace{\Rightarrow^{\text{MDA}}_{\text{MDA}}(P), \dots, I_D}^d, \dots, I_D]$  (first input MDA),  
 3967  $T[I_1, \dots, \overbrace{\Rightarrow^{\text{MDA}}_{\text{MDA}}(Q), \dots, I_D}^d, \dots, I_D]$  (second input MDA), and  $T[I_1, \dots, \overbrace{\Rightarrow^{\text{MDA}}_{\text{MDA}}(P \cup Q), \dots, I_D}^d, \dots, I_D]$  (output  
 3968

3970 MDA), as these types are generic in meta-parameters. Formal definitions and details about our  
 3971 meta-parameter concept are provided in Section A for the interested reader.

3972 We use meta-functions as an analogous concept to *metaprogramming* in programming language  
 3973 theory to achieve high generality. For example, by defining combine operators as meta-functions,  
 3974 we can use the operators on input MDAs that operate on arbitrary index sets while still guaranteeing  
 3975 correctness, e.g., that index sets of the two input MDAs match in all dimensions except in the  
 3976 dimension to combine (as discussed above). For simplicity, we often refrain from explicitly stating  
 3977 meta-parameters when they are clear from the context; for example, when they can be deduced  
 3978 from the types of their particular inputs (a.k.a. *type deduction* in programming).

3979 We now formally discuss the example operators *concatenation* and *point-wise combination*. For  
 3980 high flexibility, we define both operators generically in the scalar type  $T$  of their input and output  
 3981 MDAs, the MDAs' dimensionality  $D$ , as well as in the dimension  $d$  to combine.  
 3982

3983 **Example 13** (Concatenation). We define *concatenation* as function  $+$  of type

$$3984 + : \langle T \in \text{TYPE} \mid D \in \mathbb{N} \mid d \in [1, D]_{\mathbb{N}} \mid (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} \rangle :$$

$$3985 + : T[I_1, \dots, \underbrace{id(P), \dots, I_D}_d] \times T[I_1, \dots, \underbrace{id(Q), \dots, I_D}_d] \rightarrow T[I_1, \dots, \underbrace{id(P \cup Q), \dots, I_D}_d]$$

3986 where  $id : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$  is the identity function on MDA index sets. The  
 3987 function is computed as:

$$3988 + : \langle T \mid D \mid d \mid (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D), (P, Q) \rangle (\alpha_1, \alpha_2)[i_1, \dots, i_d, \dots, i_D]$$

$$3989 := \begin{cases} \alpha_1[i_1, \dots, i_d, \dots, i_D] & , i_d \in P \\ \alpha_2[i_1, \dots, i_d, \dots, i_D] & , i_d \in Q \end{cases}$$

3990 The function is well defined, because  $P$  and  $Q$  are disjoint. We usually use an infix notation for  
 3991  $+^{<\dots>}$  (meta-parameters omitted via ellipsis), i.e., we write  $\alpha_1 +^{<\dots>} \alpha_2$  instead of  $+^{<\dots>}(\alpha_1, \alpha_2)$ .  
 3992

3993 The vertical bar in the superscript of  $+$  denotes that function  $+$  can be partially evaluated  
 3994 (a.k.a. *Currying* [Curry 1980] in math and *multi staging* [Taha and Sheard 1997] in programming)  
 3995 for particular values of meta-parameters:  $T \in \text{TYPE}$  (first stage),  $D \in \mathbb{N}$  (second stage), etc. Partial  
 3996 evaluation (formally defined in Definition 20) enables both: 1) expressive typing and thus better error  
 3997 elimination: for example, parameter  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}$  can depend on  
 3998 meta-parameter  $D \in \mathbb{N}$ , because  $D$  is defined in an earlier stage, which allows precisely limiting  
 3999 the length of the tuple  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D)$  to  $D - 1$  index sets; 2) generality: for example, we  
 4000 can instantiate  $+$  to  $+^{<T>}$  which is specific for a particular scalar type  $T \in \text{TYPE}$ , but still generic  
 4001 in meta-parameters  $D \in \mathbb{N}$ ,  $d \in [1, D]_{\mathbb{N}}$ , ..., as these meta-parameters are defined in later stages.  
 4002 We specify stages and their order according to the recommendations in Haskell Wiki [2013], e.g.,  
 4003 using earlier stages for meta-parameters that are expected to change less frequently than other  
 4004 meta-parameters.  
 4005

4006 It is easy to see that  $+^{<T \mid D \mid d>}$  is a combine operator of type  $\text{CO}^{<id \mid T \mid D \mid d>}$  for any particular  
 4007 choice of meta-parameters  $T \in \text{TYPE}$ ,  $D \in \mathbb{N}$ , and  $d \in [1, D]_{\mathbb{N}}$ .  
 4008

4009 **Example 14** (Point-Wise Combination). We define *point-wise combination*, according to a binary  
 4010 function  $\oplus : T \times T \rightarrow T$  (e.g. addition), as function  $\overrightarrow{\oplus}$  of type  
 4011

$$4012 \overrightarrow{\oplus} : \langle T \in \text{TYPE} \mid D \in \mathbb{N} \mid d \in [1, D]_{\mathbb{N}} \mid (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} \rangle :$$

4013

where  $0_f : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$ ,  $I \mapsto \{0\}$  is the constant MDA index set function that maps any index set  $I$  to the index set containing MDA index 0 only. The function is computed as:

$$\bullet^{< T | D | d |} (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D), (P, Q) > (\oplus) (\mathfrak{a}_1, \mathfrak{a}_2) [i_1, \dots, \overset{\uparrow}{0}, \dots, i_D] := \\ \mathfrak{a}_1 [i_1, \dots, \overset{\uparrow}{0}, \dots, i_D] \oplus \mathfrak{a}_2 [i_1, \dots, \overset{\uparrow}{0}, \dots, i_D]$$

We often write  $\oplus$  only, instead of  $\vec{\bullet}(\oplus)$ , and we usually use an infix notation for  $\oplus$ .

Function  $\overrightarrow{\bullet}^{<T|D|d>}(\oplus)$  (meaning:  $\overrightarrow{\bullet}$  is partially applied to ordinary function parameter  $\oplus$  and thus still generic in parameters  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D)$  and  $(P, Q)$  – formal details provided in Definition 21) is a combine operator of type  $\text{CO}^{<0_f|T|D|d>}$  for any binary operator  $\oplus : T \times T \rightarrow T$ .

### B.3 Multi-Dimensional Homomorphisms

**Definition 26** (Multi-Dimensional Homomorphism). Let  $T^{\text{INP}}, T^{\text{OUT}} \in \text{TYPE}$  be two arbitrary scalar types,  $D \in \mathbb{N}$  a natural number, and  $\stackrel{1}{\Rightarrow}_{\text{MDA}}^{\text{MDA}}, \dots, \stackrel{D}{\Rightarrow}_{\text{MDA}}^{\text{MDA}} : \text{MDA-IDX-SETs} \rightarrow \text{MDA-IDX-SETs}$  functions on MDA index sets. Let further  $+_d := +^{<T^{\text{INP}}|D|d} \in \text{CO}^{<\text{id}|T^{\text{INP}}|D|d}$  be concatenation (Definition 13) in dimension  $d \in [1, D]_{\mathbb{N}}$  on  $D$ -dimensional MDAs that have scalar type  $T^{\text{INP}}$ .

## A function

$$h^{< I_1, \dots, I_D \in \text{MDA-IDX-SETS} >} : T^{\text{INP}}[ I_1, \dots, I_D ] \rightarrow T^{\text{OUT}}[ \overset{1}{\Rightarrow}_{\text{MDA}}(I_1), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(I_D) ]$$

is a *Multi-Dimensional Homomorphism (MDH)* that has *input scalar type*  $T^{\text{INP}}$ , *output scalar type*  $T^{\text{OUT}}$ , *dimensionality*  $D$ , and *index set functions*  $\Rightarrow^1_{\text{MDA}}, \dots, \Rightarrow^D_{\text{MDA}}$ , iff for each  $d \in [1, D]_{\mathbb{N}}$ , there exists a combine operator  $\otimes_d \in \text{CO}^{<\Rightarrow^d_{\text{MDA}} | T^{\text{OUT}} | D | d >}$  (Definition 25), such that for any concatenated input MDA  $a_1 \mathbin{++}_d a_2$  in dimension  $d$ , the *homomorphic property* is satisfied:

$$h(\mathfrak{q}_1 \#_d \mathfrak{q}_2) \equiv h(\mathfrak{q}_1) \otimes_d h(\mathfrak{q}_2)$$

We denote the type of MDHs concisely as  $\text{MDH}^{<T^{\text{INP}}, T^{\text{OUT}} | D |}(\stackrel{d}{\text{MDA}})_{d \in [1, D]}_{\mathbb{N}}$

MDHs are defined such that applying them to a concatenated MDA in dimension  $d$  can be computed by applying the MDH  $h$  independently to the MDA's parts  $\alpha_1$  and  $\alpha_2$  and combining the results afterwards by using its combine operator  $\otimes_d$ , as also discussed above. Note that per definition of MDHs, their combine operators are associative and commutative (which follows from the associativity and commutativity of  $+_d$ ). Note further that for simplicity, Definition 26 is specialized to MDHs whose input algebraic structure relies on concatenation, as such kinds of MDHs already cover the currently practice-relevant data-parallel computations (as we will see later). We provide a generalized definition of MDHs in Section B.7, for the algebraically interested reader.

**Example 15** (Function Mapping). Function  $\text{map}^{<T^{\text{INP}}, T^{\text{OUT}} | D | (I_1, \dots, I_D)}(f)$  maps a function  $f : T^{\text{INP}} \rightarrow T^{\text{OUT}}$  to each element of an MDA that has scalar type  $T^{\text{INP}} \in \text{TYPE}$ , dimensionality  $D \in \mathbb{N}$ , and index sets  $I := (I_1, \dots, I_D) \in \text{MDA-IDX-SETs}^D$ .

4068 The function is of type

4069  $\text{map}^{<T^{\text{INP}}, T^{\text{OUT}} \in \text{TYPE} | D \in \mathbb{N} | (I_1, \dots, I_D) \in \text{MDA-IDX-SETS}^D} :$

4070

$$\underbrace{T^{\text{INP}} \rightarrow T^{\text{OUT}}}_{f} \rightarrow \underbrace{T^{\text{INP}}[I_1, \dots, I_D]}_{\text{map}^{<T^{\text{INP}}, T^{\text{OUT}} | D | (I_1, \dots, I_D)}(f)} \rightarrow \underbrace{T^{\text{OUT}}[I_1, \dots, I_D]}_{\text{map}^{<T^{\text{INP}}, T^{\text{OUT}} | D | (I_1, \dots, I_D)}(f)}$$

4071 and it is computed as:

4072

$$\text{a} \xrightarrow{\text{map}^{<T^{\text{INP}}, T^{\text{OUT}} | D | (I_1, \dots, I_D)}(f)} \underset{i_1 \in I_1}{\underset{i_D \in I_D}{\underset{\text{++}_1 \dots \text{++}_D}{\vec{f}_{\text{map}}(\text{a}|_{\{i_1\} \times \dots \times \{i_D\}})}}}$$

4073 Here,  $\text{++}_d := \text{++}^{<T^{\text{OUT}} | D | d}$  denotes concatenation (Example 13) in dimension  $d$ , MDA  $\text{a}|_{\{i_1\} \times \dots \times \{i_D\}}$  is the restriction of  $\text{a}$  to the single element accessed via indices  $(i_1, \dots, i_D)$ , and  $\vec{f}_{\text{map}}$  denotes the adaption of function  $f$  to operate on MDAs comprising a single value only: it is of type

4074  $\vec{f}_{\text{map}}^{<i_1, \dots, i_D \in \mathbb{N}} : T^{\text{INP}}[\{i_1\}, \dots, \{i_D\}] \rightarrow T^{\text{OUT}}[\{i_1\}, \dots, \{i_D\}]$

4075 and defined as

4076  $\vec{f}_{\text{map}}(x)[i_1, \dots, i_D] := f(x[i_1, \dots, i_D])$

4077 It is easy to see that  $\text{map}^{<T^{\text{INP}}, T^{\text{OUT}} | D >}(f)$  is an MDH of type  $\text{MDH}^{<T^{\text{INP}}, T^{\text{OUT}} | D | id, \dots, id >}$  whose combine operators are concatenation  $\text{++}_d \in \text{CO}^{<id | T^{\text{OUT}} | D | d >}$  in all dimensions  $d \in [1, D]_{\mathbb{N}}$ .

4078 We have chosen  $\text{map}$  function's order of stages –  $T^{\text{INP}}, T^{\text{OUT}} \in \text{TYPE}$  (stage 1),  $D \in \mathbb{N}$  (stage 2), and  $(I_1, \dots, I_D) \in \text{MDA-IDX-SETS}^D$  (stage 3) – according to the recommendations in [Haskell Wiki \[2013\]](#), i.e., earlier stages (such as the scalar types  $T^{\text{INP}}, T^{\text{OUT}}$ ) are expected to change less frequently than later stages (e.g., the MDAs' index sets  $I_1, \dots, I_D$ ).

4079 **Example 16** (Reduction). Function  $\text{red}^{<T | D | (I_1, \dots, I_D)}(\oplus)$  combines all elements within an MDA that has scalar type  $T \in \text{TYPE}$ , dimensionality  $D \in \mathbb{N}$ , and index sets  $I := (I_1, \dots, I_D) \in \text{MDA-IDX-SETS}^D$ , using an associative and commutative binary function  $\oplus : T \times T \rightarrow T$ .

4080 The function is of type

4081  $\text{red}^{<T \in \text{TYPE} | D \in \mathbb{N} | (I_1, \dots, I_D) \in \text{MDA-IDX-SETS}^D} : \underbrace{T \times T \rightarrow T}_{\oplus} \rightarrow \underbrace{T[I_1, \dots, I_D]}_{\text{red}^{<T | D | (I_1, \dots, I_D)}(\oplus)} \rightarrow T[1, \dots, 1]$

4082 and it is computed as:

4083  $\text{a} \xrightarrow{\text{red}^{<T | D | (I_1, \dots, I_D)}(\oplus)} \underset{i_1 \in I_1}{\underset{i_D \in I_D}{\underset{\vec{\bullet}_1(\oplus) \dots \vec{\bullet}_D(\oplus)}{\vec{f}_{\text{red}}(\text{a}|_{\{i_1\} \times \dots \times \{i_D\}})}}}$

4084 Here,  $\vec{\bullet}_d(\oplus) := \vec{\bullet}^{<T | D | d >}(\oplus)$  denotes point-wise combination (Example 14) in dimension  $d$ , MDA  $\text{a}|_{\{i_1\} \times \dots \times \{i_D\}}$  is defined as above, and  $\vec{f}_{\text{red}}$  is the function of type

4085  $\vec{f}_{\text{red}}^{<i_1, \dots, i_D \in \mathbb{N}} : T^{\text{INP}}[\{i_1\}, \dots, \{i_D\}] \rightarrow T^{\text{OUT}}[\{0\}, \dots, \{0\}]$

4086 that is defined as

4087  $\vec{f}_{\text{red}}(x)[0, \dots, 0] := x[i_1, \dots, i_D]$

4088 It is easy to see that  $\text{red}^{<T | D >}(\oplus)$  is an MDH of type  $\text{MDH}^{<T, T | D | 0_f, \dots, 0_f >}$  whose combine operators are point-wise addition  $\vec{\bullet}_d(\oplus) \in \text{CO}^{<id | T | D | d >}$  in all dimensions  $d \in [1, D]_{\mathbb{N}}$ . The same as for function  $\text{map}$ , function  $\text{red}$ 's order of meta-parameter stages are chosen according to [\[Haskell Wiki 2013\]](#).

4117 **Definition 27** (Higher-Order Function `md_hom`). The higher-order function `md_hom` is of type

4118  $\text{md\_hom}^{<T^{\text{INP}}, T^{\text{OUT}} \in \text{TYPE} \mid D \in \mathbb{N} \mid (\xrightarrow{d}^{\text{MDA}}_{\text{MDA}} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS})_{d \in [1, D]_{\mathbb{N}} >}$  :

4120 
$$\underbrace{\text{SF}^{<T^{\text{INP}}, T^{\text{OUT}} >}}_f \times \underbrace{(\text{CO}^{<\xrightarrow{1}^{\text{MDA}}_{\text{MDA}} | T^{\text{OUT}} | D | 1 >} \times \dots \times \text{CO}^{<\xrightarrow{D}^{\text{MDA}}_{\text{MDA}} | T^{\text{OUT}} | D | D >})}_{\otimes_1, \dots, \otimes_D} \rightarrow_p \underbrace{\text{MDH}^{<T^{\text{INP}}, T^{\text{OUT}} \mid D \mid (\xrightarrow{d}^{\text{MDA}}_{\text{MDA}})_{d \in [1, D]_{\mathbb{N}} >}}}_{\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))}$$

4126 where  $\text{SF}^{<T^{\text{INP}}, T^{\text{OUT}} >}$  denotes scalar functions of type  $T^{\text{INP}} \rightarrow T^{\text{OUT}}$ . Function `md_hom` is partial (indicated by  $\rightarrow_p$  instead of  $\rightarrow$ ), which we motivate after this definition. The function takes as input a scalar function  $f$  and a tuple of  $D$ -many combine operators  $(\otimes_1, \dots, \otimes_D)$ , and it yields a function  $\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))$  which is defined as

4131 
$$\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\alpha) := \otimes_{i_1 \in I_1} \dots \otimes_{i_D \in I_D} \vec{f}(\alpha|_{\{i_1\} \times \dots \times \{i_D\}})$$

4133 The combine operators' underset notation denotes straightforward iteration (explained formally in  
4134 Notation 5), and the MDA  $\alpha|_{\{i_1\} \times \dots \times \{i_D\}}$  is the restriction of  $\alpha$  to the MDA containing the single  
4135 element accessed via MDA indices  $(i_1, \dots, i_D)$ . Function  $\vec{f}$  behaves like scalar function  $f$ , but  
4136 operates on singleton MDAs (rather than scalars): the function is of type

4137 
$$\vec{f}^{<i_1, \dots, i_D \in \mathbb{N}_0 >} : T^{\text{INP}}[\{i_1\}, \dots, \{i_D\}] \rightarrow T^{\text{OUT}}[\xrightarrow{1}^{\text{MDA}}(\{i_1\}), \dots, \xrightarrow{D}^{\text{MDA}}(\{i_D\})]$$

4139 and defined as

4140 
$$\vec{f}(x)[j_1, \dots, j_D] := f(x[i_1, \dots, i_D])^{24}$$

4141 For  $\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))$ , we require per definition the homomorphic property (Definition 3),  
4142 i.e., for each  $d \in [1, D]_{\mathbb{N}}$ , it must hold:

4144 
$$\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\alpha_{1+d} \alpha_2) =$$

4145 
$$\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\alpha_1) \otimes_d \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\alpha_2)$$

4147 Note that we can potentially allow in Definition 27 the case  $D = 0$  in which we would define the  
4148 `md_hom` instance equal to the scalar function  $f$ :

4149 
$$\text{md\_hom}(f, ()) := f$$

#### 4151 B.4 View Functions

4152 **Definition 28** (Buffer). Let  $T \in \text{TYPE}$  be an arbitrary scalar type,  $D \in \mathbb{N}_0$  a natural number<sup>25</sup>, and  
4153  $N := (N_1, \dots, N_D) \in \mathbb{N}^D$  a tuple of natural numbers.

4154 A *Buffer (BUF)*  $b$  that has *dimensionality D*, *size N*, and *scalar type T* is a function with the  
4155 following signature:

4157 
$$b : [0, N_1]_{\mathbb{N}_0} \times \dots \times [0, N_D]_{\mathbb{N}_0} \rightarrow T \cup \{\perp\}$$

4158 Here, we use  $\perp$  to denote the *undefined value*. We refer to  $[0, N_1]_{\mathbb{N}_0} \times \dots \times [0, N_D]_{\mathbb{N}_0} \rightarrow T \cup \{\perp\}$   
4159 as the *type* of BUF  $b$ , which we also denote as  $T^{N_1 \times \dots \times N_D}$ , and we refer to the set  $\text{BUF-IDX-SETS} :=$   
4160  $\{[0, N]_{\mathbb{N}_0} \mid N \in \mathbb{N}\}$  as *BUF index sets*. Analogously to Notation 1, we write  $b[i_1, \dots, i_D]$  instead of  
4161  $b(i_1, \dots, i_D)$  to avoid a too heavy usage of parentheses.

4162 <sup>24</sup>We assume in this work that MDH functions when used on singleton MDAs return a singleton MDA, as such MDHs  
4163 already cover all real-world cases we currently are aware of.

4164 <sup>25</sup>We use the case  $D = 0$  to represent scalar values, as we discuss later.

In our approach, *index functions* are used by the user to conveniently instantiate higher-order functions `inp_view` and `out_view` (e.g., index function  $(i, k) \mapsto (i, k)$  and  $(i, k) \mapsto (k)$  used in Figure 6 to instantiate `inp_view`, and  $(i, k) \mapsto (i)$  used for `out_view`).

**Definition 29** (Index Function). Let  $D \in \mathbb{N}$  be a natural number (representing an MDA's dimensionality) and  $D_b \in \mathbb{N}_0$  (representing a BUF's dimensionality).

An *index function*  $\text{idx}$  from  $D$ -dimensional MDA indices to  $D_b$ -dimensional BUF indices is any meta-function of type

$$\text{idx}^{<I_1^{\text{MDA}}, \dots, I_D^{\text{MDA}} \in \text{MDA-IDX-SETS}^D>} : I_1^{\text{MDA}} \times \dots \times I_D^{\text{MDA}} \rightarrow I_1^{\text{BUF}} \times \dots \times I_{D_b}^{\text{BUF}}$$

for  $(I_1^{\text{BUF}}, \dots, I_{D_b}^{\text{BUF}}) := \Rightarrow_{\text{BUF}}^{\text{MDA}}(I_1^{\text{MDA}}, \dots, I_D^{\text{MDA}})$  where  $\Rightarrow_{\text{BUF}}^{\text{MDA}} : \text{MDA-IDX-SETS}^D \rightarrow \text{BUF-IDX-SETS}^{D_b}$  is an arbitrary but fixed function that maps  $D$ -many MDA index sets to  $D_b$ -many BUF index sets. We denote the type of index functions as  $\text{MDA-IDX-to-BUF-IDX}^{< D, D_b } \mid \Rightarrow_{\text{BUF}}^{\text{MDA}}$ .

In words: Index functions have to be capable of operating on any potential MDA index set. This generality will be required later for using index functions also on parts of MDAs whose index sets are subsets of the original MDA's index sets.

We will use index functions to access BUFS. For example, in the case of MatVec (Figure 1), we access its input matrix using index function  $(i, k) \mapsto (i, k)$  which is of type

MDA-IDX-to-BUF-IDX  $\leftarrow D := 2, D_b := 2 \mid \Rightarrow_{\text{BUF}}^{\text{MDA}} (I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_1^{\text{MDA}})]_{N_0}, [0, \max(I_2^{\text{MDA}})]_{N_0} \right>$

and we use index function  $(i, k) \mapsto (k)$  to access MatVec's input vector, which is of type

MDA-TDX-to-BUE-TDX < $D := 2, D_b := 1$  |  $\Rightarrow_{\text{BUF}}^{\text{MDA}} (I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_2^{\text{MDA}})]_{\mathbb{N}_0}$ >

Further examples of index function, e.g., for stencil computation `Jacobi1D`, are presented in Section [B.10](#), for the interested reader.

**Definition 30** (Input View). An *input view* from  $B$ -many BUFs,  $B \in \mathbb{N}$ , of arbitrary but fixed types  $T_i^{N_1^b \times \dots \times N_D^b}$ ,  $b \in [1, B]_{\mathbb{N}}$ , to an MDA of arbitrary but fixed type  $T[I_1, \dots, I_D]$  is any function  $\text{iv}$  of type:

$$\text{iv: } \underbrace{\prod_{b=1}^B T_b^{N_1^b \times \dots \times N_{D_b}^b}}_{\text{BUFs}} \xrightarrow{p} \underbrace{T[\, I_1, \dots, I_D \,]}_{\text{MDA}}$$

We denote the type of  $\mathbf{iv}$  as  $\mathbf{IV}^{ $b$ } | (D_b)_{b \in [1, B]_{\mathbb{N}}} | (N_1^b, \dots, N_{D_b}^b)_{b \in [1, B]_{\mathbb{N}}} | (T_b)_{b \in [1, B]_{\mathbb{N}}} | D | I_1, \dots, I_D | T >$ .

**Example 17** (Input View – MatVec). The input view of MatVec on a  $1024 \times 512$  matrix and 512-sized vector (sizes are chosen arbitrarily), both of integers  $\mathbb{Z}$ , is of type

$$\mathcal{IV}^{<B=2 \mid D_1=2, D_2=1 \mid (N_1^1=1024, N_2^1=512), (N_1^2=512) \mid T_1=\mathbb{Z}, T_2=\mathbb{Z} \mid D=2 \mid I_1=[0, 1024]_{\mathbb{N}_0}, I_2=[0, 512]_{\mathbb{N}_0} \mid T=\mathbb{Z} \times \mathbb{Z} >}$$

and defined as

$$\underbrace{[ M(i, k) ]_{i \in [0, 1024]_{\mathbb{N}_0}, k \in [0, 512]_{\mathbb{N}_0}}}_{\text{Matrix}}, \underbrace{[ v(k) ]_{k \in [0, 512]_{\mathbb{N}_0}}}_{\text{Vector}} \mapsto \underbrace{[ M(i, k), v(k) ]_{i \in [0, 1024]_{\mathbb{N}_0}, k \in [0, 512]_{\mathbb{N}_0}}}_{\text{MDA}}$$

Here, the BUFS' meta-parameters are as follows:  $B = 2$  is the number of BUFS (matrix and vector);  $D_1 = 2$  is dimensionality of the matrix and  $D_2 = 1$  the vector's dimensionality;  $(N_1^1, N_2^1) =$

4215 (1024, 512) is the matrix size and  $N_1^2 = 512$  the vector's size;  $T_1, T_2 = \mathbb{Z}$  are the scalar types of  
 4216 matrix and vector. The MDA's meta-parameters are:  $D = 2$  is the computed MDA's dimensionality;  
 4217  $I_1, I_2$  are the MDA's index sets; parameter  $T = \mathbb{Z} \times \mathbb{Z}$  is MDA's scalar type (pairs of matrix/vector  
 4218 elements – see Figure 1).

4219 **Example 18** (Input View – Jacobi1D). The input view of Jacobi1D on a 512-sized vector of  
 4220 integers is of type  
 4221

4222  $\overbrace{\text{IV}^{<B=1 \mid D=1 \mid (N_1^1=512) \mid T_1=\mathbb{Z} \mid D=1 \mid I_1=[0,512-2]_{\mathbb{N}_0} \mid T=\mathbb{Z} \times \mathbb{Z} \times \mathbb{Z}>}}^{\substack{\text{BUFs' meta-parameters} \\ \text{MDA's meta-parameters}}}$

4223

4224 and defined as

4225  $\underbrace{[v(i)]_{i \in [0,512]_{\mathbb{N}_0}}}_{\text{Vector}} \mapsto \underbrace{[v(i+0), v(i+1), v(i+2)]_{i \in [0,512-2]_{\mathbb{N}_0}}}_{\text{MDA}}$

4226

4227 Higher-order function `inp_view` uses the index functions  $\text{idx}_{b,a}$  to compute an input view that  
 4228 maps BUFs  $b_1, \dots, b_B$  to an MDA  $a$  that contains at position  $i_1, \dots, i_D$  the following element:  
 4229

4230  $a[i_1, \dots, i_D] := (( \underbrace{b_1[\text{idx}_{1,1}(i_1, \dots, i_D)] \in T_1}_{a=1}, \dots, \underbrace{b_1[\text{idx}_{1,A_1}(i_1, \dots, i_D)] \in T_1}_{a=A_1} ),$

4231  $\underbrace{\vdots}_{b=1}$

4232  $( \underbrace{b_B[\text{idx}_{B,1}(i_1, \dots, i_D)] \in T_B}_{a=1}, \dots, \underbrace{b_B[\text{idx}_{B,A_B}(i_1, \dots, i_D)] \in T_B}_{a=A_B} ) )$

4233

4234 The element consists of  $B$ -many tuples – one per BUF – and each such tuple contains  $A_b$ -many  
 4235 elements – one element per access to the  $b$ -th BUF. For MatVec, the element is of the form  
 4236

4237  $a[i_1, i_2] := (( \underbrace{b_1[\text{idx}_{1,1}(i_1, i_2)] := (i_1, i_2)}_{a=1} \in T_1 ), ( \underbrace{b_2[\text{idx}_{2,1}(i_1, i_2)] := (i_2)}_{a=1} \in T_2 ))$

4238  $\underbrace{\vdots}_{b=1} \quad \underbrace{\vdots}_{b=2}$

4239 and for Jacobi1D, the element is

4240  $a[i_1] := (( \underbrace{b_1[\text{idx}_{1,1}(i_1)] := (i_1+0)}_{a=1} \in T_1, \underbrace{b_1[\text{idx}_{1,2}(i_1)] := (i_1+1)}_{a=2} \in T_1, \underbrace{b_1[\text{idx}_{1,3}(i_1)] := (i_1+2)}_{a=3} \in T_1 ))$

4241  $\underbrace{\vdots}_{b=1}$

4242 We now formally define higher-order function `inp_view`. For high flexibility and formal correctness,  
 4243 function `inp_view` relies on a type that involves many meta-parameters. The high number of  
 4244 meta-parameters, and the resulting complex type of function `inp_view`, might appear daunting to  
 4245 the user. However, Notation 6 confirms that despite the complex type of function `inp_view`, the function  
 4246 can be conveniently used by the user (as also illustrated in Figure 6), because meta-parameters  
 4247 can be automatically deduced from `inp_view`'s input parameters.

4264 **Definition 31** (Higher-Order Function `inp_view`). Function `inp_view` is of type

<code>inp_view</code>	$B \in \mathbb{N}$	$A_1, \dots, A_B \in \mathbb{N}$	$D_1, \dots, D_B \in \mathbb{N}_0$	$D \in \mathbb{N}$
	$\underbrace{\phantom{A_1, \dots, A_B \in \mathbb{N}}}_{\text{Number BUFS}}$	$\underbrace{\phantom{D_1, \dots, D_B \in \mathbb{N}_0}}_{\text{Number BUFS' Accesses}}$	$\underbrace{\phantom{D_1, \dots, D_B \in \mathbb{N}_0}}_{\text{BUFS' Dimensionalities}}$	$\underbrace{\phantom{D \in \mathbb{N}}}_{\text{MDA Dimensionality}}$
	$(\Rightarrow_{\substack{\text{MDA } b, a \\ \text{BUF}}}^{\text{MDA-IDX-SETS}^D \rightarrow \text{BUF-IDX-SETS}^{D_b}})_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_b]_{\mathbb{N}}} > :$			

The diagram illustrates the mapping of Buffer Access to MDA's Meta-Parameters and Postponed Parameters. On the left, 'Buffer Access' is shown with indices  $b=1$  and  $a=1$ . An arrow labeled 'MDA-IDX-to-BUF-IDX' maps this to an 'Index Function:  $\text{idx}_{b,a}$ '. This function is then mapped to 'MDA's Meta-Parameters' (a set of parameters  $D, D_b \rightarrow T, T_b$ ) and 'Postponed Parameters' (a set of parameters  $D, I_1, \dots, I_D \rightarrow N_1^1, \dots, N_{D_B}^B$ ). The overall mapping is labeled 'Input View:  $\text{iv}$ '.

for  $N_d^b := 1 + \max(\bigcup_{a \in [1, A_b]_{\mathbb{N}}} \xrightarrow[\text{BUF}]{d_{\text{MDA}, b, a}} (I_1, \dots, I_D))$  and  $T := \times_{b=1}^B \times_{a=1}^{A_b} T_b$ , and it is defined as:

$$\underbrace{(\mathsf{idx}_{b,a})_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_b]_{\mathbb{N}}}}_{\text{Index Functions}} \mapsto \underbrace{(\underbrace{\mathsf{b}_1, \dots, \mathsf{b}_B}_{\text{BUFs}})}_{\text{MDA}} \xrightarrow{\text{iv}} \underbrace{\mathsf{a}}_{\text{Input View}}$$

for

$$\mathfrak{a}[i_1, \dots, i_D] := (\mathfrak{a}_{b,a}[i_1, \dots, i_D])_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_b]_{\mathbb{N}}}$$

and

$$\mathfrak{a}_{b,a}[i_1, \dots, i_D] := \mathfrak{b}_b[\mathfrak{id}\mathfrak{x}_{b,a}(i_1, \dots, i_D)]$$

Higher-order function `inp_view` takes as input a collection of index functions that are of types `MDA-IDX-to-BUF-IDX` (Definition 29), and it computes an input view of type IV (Definition 30) based on the index functions, as illustrated in Figures 10 and 11.

We use for type MDA-IDX-to-BUF-IDX as concrete meta-parameter values (listed in angle brackets) straightforwardly the values of meta-parameters passed to function `inp_view`. Similarly, for type IV's meta-parameters  $B$ ,  $D_1, \dots, D_B$ , and  $D$ , we use again straightforwardly the particular meta-parameter values of function `inp_view`.

To be able using the computed input view on arbitrarily typed input buffers and letting the input view compute MDAs that have arbitrary index sets, we keep IV's meta-parameters  $T_1, \dots, T_B$  (scalar types of the computed view's input buffers) and  $I_1, \dots, I_D$  (index sets of the view's returned MDA) flexible. Being flexible in the BUFs' scalar types and MDA's index sets is important for convenience – for example, in the case of MatVec, this flexibility allows using the computed input view generically for matrices and vectors that have arbitrary scalar types (e.g., either `int` or `float`) and sizes  $(I, J)$  (matrix) and  $J$  (vector), for arbitrary  $I, J \in \mathbb{N}$ , without needing to re-compute a new input view every time again when BUFs' scalar types and/or sizes change.

We automatically compute the sizes  $N_d^b$  of BUFs in IV's meta-parameter list (e.g., in the case of MatVec, the size of the input matrix  $(I, J)$  and vector size  $J$ ), according to the formula in Definition 31, based on the flexible MDA's index sets (e.g., sets  $[0, I]_{\mathbb{N}_0}$  and  $[0, J]_{\mathbb{N}_0}$  for MatVec). By computing BUF sizes from MDA index sets (rather than requesting the sizes explicitly from the user), we achieve strong error checking: for example, for MatVec, we can ensure – already on

4313 the type level – that the number of columns of its input matrix and the size of its input vector  
 4314 match. To compute the BUF sizes, we *postpone* via  $\rightarrow$  (defined formally in Definition 23) the sizes  
 4315 in IV's meta-parameter list to later meta-parameter stages; this is because the sizes are defined in  
 4316 early stages and thus have no access to the MDA's index sets which are defined in later stages. Our  
 4317 formula in Definition 31 then works as follows: for each BUF  $b$ , its size  $N_d^b$  has to be well-defined in  
 4318 each of its dimensions  $d$ , for all accesses  $a$  via the BUF's index functions when used for all indices  
 4319  $I_1, \dots, I_D$  within the MDA index sets. Here, in the computation of  $N_d^b$ , function  $\Rightarrow_{\text{BUF}}^{MDA, b, a}$  computes  
 4320 the  $d$ -th component of the  $D_b$ -sized output tuple of  $\Rightarrow_{\text{BUF}}^{MDA, b, a}$  (the computed component is the index  
 4321 set of BUF  $b$  in dimension  $d$  for the  $a$ -th index function used to access the BUF).  
 4322

4323 We automatically compute also MDA's scalar type  $T$  using the formula presented in Definition 31.  
 4324 The formula computes  $T$  as a tuple that consists of the BUFs' scalar types, as each MDA element  
 4325 consist of BUF elements (illustrated in Figures 10 and 11). Postponing  $T$  in IV's meta-parameter  
 4326 list is done (analogously as for  $N_d^b$ ), but is actually not required, because the BUFs' scalar types  
 4327  $T_1, \dots, T_B$  are already defined in earlier meta-parameter stages than  $T$ . However, we will see that  
 4328 postponing  $T$  is required later in the Definition 33 of higher-order function *out\_view*; therefore,  
 4329 we postpone  $T$  also in our definition of *inp\_view* to increase consistency between our definitions  
 4330 of *inp\_view* (Definition 31) and *out\_view* (Definition 33).  
 4331

4331 **Notation 6** (Input Views). Let  $\text{inp\_view}^{<\dots>}((\text{id}x_{b,a})_{b \in [1,B]_{\mathbb{N}}, a \in [1,A_b]_{\mathbb{N}}})$  be a particular instance  
 4332 of higher-order function *inp\_view* (meta-parameters omitted via ellipsis for simplicity) for an  
 4333 arbitrary but fixed choice of index functions. Let further  $ID_1, \dots, ID_B \in \Sigma^*$  be arbitrary, user-defined  
 4334 BUF identifiers (e.g.,  $ID_1 = "M"$  and  $ID_2 = "v"$  in the case of MatVec), for an arbitrary, fixed collection  
 4335 of letters  $\Sigma = \{A, B, C, \dots, a, b, c, \dots, 1, 2, 3, \dots\}$ .  
 4336

4336 For better readability, we use the following notation for the 2-dimensional structure of index  
 4337 functions taken as input by function *inp\_view*, inspired by Lattner et al. [2021]:  
 4338

4339  $\text{inp\_view}(ID_1 : \text{id}x_{1,1}, \dots, \text{id}x_{1,A_1}, \dots, ID_B : \text{id}x_{B,1}, \dots, \text{id}x_{B,A_B})$   
 4340

4341 We refrain from stating *inp\_view*'s meta-parameters in our notation, as the parameters can be  
 4342 automatically deduced from the number and types of index functions.  
 4343

4343 **Definition 32** (Output View). An *output view* from an MDA of arbitrary but fixed type  $T[I_1, \dots, I_D]$   
 4344 to  $B$ -many BUFs,  $B \in \mathbb{N}$ , of arbitrary but fixed types  $T_b^{N_1^b \times \dots \times N_{D_b}^b}$ ,  $b \in [1, B]_{\mathbb{N}}$ , is any function  $\text{ov}$  of  
 4345 type:  
 4346

$$4347 \text{ov} : \underbrace{T[I_1, \dots, I_D]}_{\text{MDA}} \rightarrow_p \underbrace{\bigtimes_{b=1}^B T_b^{N_1^b \times \dots \times N_{D_b}^b}}_{\text{BUFs}}$$

$$4348 \underbrace{\text{MDA's Meta-Parameters}}_{\text{MDA's Meta-Parameters}} \quad \underbrace{\text{BUFs' Meta-Parameters}}_{\text{BUFs' Meta-Parameters}}$$

4349 We denote the type of  $\text{ov}$  as  $\text{OV}^{< D \mid I_1, \dots, I_D \mid T \mid B \mid (D_b)_{b \in [1, B]_{\mathbb{N}}} \mid (N_1^b, \dots, N_{D_b}^b)_{b \in [1, B]_{\mathbb{N}}} \mid (T_b)_{b \in [1, B]_{\mathbb{N}}} >}$ .  
 4350

4351 **Example 19** (Output View – MatVec). The output view of MatVec computing a 1024-sized vector  
 4352 (size is chosen arbitrarily), of integers  $\mathbb{Z}$ , is of type  
 4353

$$4354 \underbrace{\text{MDA's meta-parameters}}_{\text{MDA's meta-parameters}} \quad \underbrace{\text{BUFs' meta-parameters}}_{\text{BUFs' meta-parameters}}$$

$$4355 \text{OV}^{< D=2 \mid I_1=[0, 1024]_{\mathbb{N}_0}, I_2=\{0\} \mid T=\mathbb{Z} \mid B=1 \mid D_1=1 \mid (N_1^1=1024) \mid T_1=\mathbb{Z} >}$$

4362 and defined as

$$\underbrace{[ w(i) ]_{i \in [0, 1024]_{\mathbb{N}_0}, k \in \{0\}}}_{\text{MDA}} \mapsto \underbrace{[ w(i) ]_{i \in [0, 1024]_{\mathbb{N}_0}}}_{\text{Vector}}$$

4366 **Example 20** (Output View – Jacobi1D). The output view of Jacobi1D computing a  $(512-2)$ -sized  
4367 vector of integers is of type  
4368

$$\text{OV}^{<D=1 | I_1=[0, 512-2]_{\mathbb{N}_0} | T=\mathbb{Z} | B=1 | D_1=1 | (N_1^1=(512-2)) | T_1=\mathbb{Z}>}$$

4372 and defined as

$$\underbrace{[ w(i) ]_{i \in [0, 512-2]_{\mathbb{N}_0}, k \in \{0\}}}_{\text{MDA}} \mapsto \underbrace{[ w(i) ]_{i \in [0, 512-2]_{\mathbb{N}_0}}}_{\text{Vector}}$$

4376 We define higher-order function `out_view` formally as follows.  
4377

4378 **Definition 33** (Higher-Order Function `out_view`). Function `out_view` is of type

$$\begin{aligned} \text{out\_view}^< & \underbrace{B \in \mathbb{N}}_{\substack{\text{Number of BUFS} \\ \text{Buffer Access}}} \quad | \quad \underbrace{A_1, \dots, A_B \in \mathbb{N}}_{\substack{\text{Number BUFS' Accesses} \\ \text{Index Function: } \text{idr}_{b,a}}} \quad | \quad \underbrace{D_1, \dots, D_B \in \mathbb{N}_0}_{\substack{\text{BUFS' Dimensionalities} \\ \text{Index Set Functions (MDA indices to BUF indices)}}} \quad | \quad \underbrace{D \in \mathbb{N}}_{\substack{\text{MDA Dimensionality} \\ \text{MDA's Meta-Parameters}}} \quad | \\ & (\Rightarrow_{\text{BUF}}^{\text{MDA } b, a : \text{MDA-IDX-SETS}^D \rightarrow \text{BUF-IDX-SETS}^{D_b}})_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_b]_{\mathbb{N}}} : \\ & \underbrace{\text{Index Functions: } \text{idr}_{1,1}, \dots, \text{idr}_{B, A_B}}_{\text{Index Functions: } \text{idr}_{b,a}} \quad \underbrace{\text{OV}^{<D | I_1, \dots, I_D \in \text{MDA-IDX-SETS} | \rightarrow | B | D_1, \dots, D_B | \rightarrow | T_1, \dots, T_B \in \text{TYPE} > < N_1^1, \dots, N_{D_B}^B | T >}}_{\substack{\text{MDA's Meta-Parameters} \\ \text{BUF's Meta-Parameters} \\ \text{Postponed Parameters}}} \end{aligned}$$

4392 which differs from `inp_view`'s type only in mapping index functions to OV (Definition 32), rather  
4393 than IV (Definition 30). Function `out_view` is defined as:  
4394

$$\underbrace{(\text{idr}_{b,a})_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_b]_{\mathbb{N}}}}_{\text{Index Functions}} \mapsto \underbrace{\underbrace{a}_{\substack{\text{MDA} \\ \text{Output View}}} \mapsto \underbrace{(\text{b}_1, \dots, \text{b}_B)}_{\substack{\text{BUFs} \\ \text{Output View}}}}_{\text{Output View: } \text{ov}}$$

4399 for

$$\text{b}_b[\text{idr}_{b,a}(i_1, \dots, i_D)] := \text{a}_{b,a}[i_1, \dots, i_D]$$

4402 and

$$(\text{a}_{b,a}[i_1, \dots, i_D])_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_b]_{\mathbb{N}}} := \text{a}[i_1, \dots, i_D]$$

4404 i.e.,  $\text{a}_{b,a}[i_1, \dots, i_D]$  is the element at point  $i_1, \dots, i_D$  within MDA  $\text{a}$  that belongs to the  $a$ -th access  
4406 of the  $b$ -th BUF. We set  $\text{b}_b[j_1, \dots, j_{D_b}] := \perp$  (symbol  $\perp$  denotes the undefined value) for all BUF  
4407 indices  $(j_1, \dots, j_{D_b}) \in [0, N_1^b]_{\mathbb{N}_0} \times \dots \times [0, N_D^b]_{\mathbb{N}_0} \setminus \bigcup_{a \in [1, A_b]_{\mathbb{N}}} \Rightarrow_{\text{BUF}}^{\text{MDA } b, a}(I_1, \dots, I_D)$  which are not in  
4409 the function range of the index functions.  
4410

4411 Note that the computed output view  $\mathbf{v}$  is partial (indicated by  $\rightarrow_p$  in Definition 32), because for  
 4412 non-injective index functions, it must hold  $\mathbf{idx}_{b,a}(i_1, \dots, i_D) = \mathbf{idx}_{b,a'}(i'_1, \dots, i'_D) \Rightarrow \mathbf{a}_{b,a}[i_1, \dots, i_D] =$   
 4413  $\mathbf{a}_{b,a'}[i'_1, \dots, i'_D]$ , which may not be satisfied for each potential input MDA of the computed view.  
 4414

## 4415 B.5 Generic High-Level Expression

4416 Figure 30 shows a generic expression in our high-level representation – consisting of higher-order  
 4417 functions `inp_view`, `md_hom`, and `out_view` (Figure 5) – for an arbitrary but fixed choice of index  
 4418 functions, scalar function, and combine operators. We express data-parallel computations using a  
 4419 particular instance of this generic expression in Figure 30.

4420 Note that all meta-parameters of higher-order function `inp_view`, `out_view`, and `md_hom` are  
 4421 automatically deduced from the particular numbers and/or types of the three functions' inputs  
 4422 (index functions in the case of `inp_view` and `out_view`, and scalar function and combine operators  
 4423 for `md_hom`).

4424 The concrete instance of `md_hom`(...) (i.e., the MDH function returned by `md_hom` for the particular  
 4425 input scalar function and combine operators) has as meta-parameter the MDH function's  
 4426 index sets (see Definition 26) for high flexibility. We use as index sets straightforwardly the input  
 4427 size  $(N_1, \dots, N_D)$  (which abbreviates  $([0, N_1]_{\mathbb{N}_0}, \dots, [0, N_D]_{\mathbb{N}_0})$  - see Notation 1)<sup>26</sup>. Instances of  
 4428 `inp_view`(...) and `out_view`(...) (i.e., the input and output view returned by `inp_view` and  
 4429 `out_view` for concrete index functions) have as meta-parameters the MDA's index sets and the  
 4430 scalar types of BUFS. We explicitly state only the meta-parameter for the BUFS' scalar types in  
 4431 our generic high-level expression (Figure 30), and we avoid explicitly stating the MDA's index  
 4432 sets for simplicity and to avoid redundancies, because the sets can be taken from the `md_hom`'s  
 4433 meta-parameter list.

4434 Note that for better readability of our high-level expressions (Figure 30), we list meta-parameters  
 4435 before parentheses, i.e., instead of writing `inp_view`(...)<sup><...></sup>, `out_view`(...)<sup><...></sup>, and `md_hom`(...)<sup><...></sup>  
 4436 for the particular instances of higher-order functions, where meta-parameters are listed at the end,  
 4437 we write `inp_view`<...>(...), `out_view`<...>(...), and `md_hom`<...>(...).

4439  $\text{out\_view}^{OB} < T_1^{OB}, \dots, T_{B^{OB}}^{OB} > ( OB_1: \mathbf{idx}_{1,1}^{OUT}, \dots, \mathbf{idx}_{1,A_1^{OB}}^{OUT}, \dots, OB_{B^{OB}}: \mathbf{idx}_{B^{OB},1}^{OUT}, \dots, \mathbf{idx}_{B^{OB},A_{B^{OB}}^{OB}}^{OUT} ) \circ$   
 4440  
 4441  $\text{md\_hom} < N_1, \dots, N_D > ( f, (\oplus_1, \dots, \oplus_D) ) \circ$   
 4442  
 4443  $\text{inp\_view}^{IB} < T_1^{IB}, \dots, T_{B^{IB}}^{IB} > ( IB_1: \mathbf{idx}_{1,1}^{INP}, \dots, \mathbf{idx}_{1,A_1^{IB}}^{INP}, \dots, IB_{B^{IB}}: \mathbf{idx}_{B^{IB},1}^{INP}, \dots, \mathbf{idx}_{B^{IB},A_{B^{IB}}^{IB}}^{INP} )$   
 4444  
 4445

4446 Fig. 30. Generic high-level expression for data-parallel computations

## 4447 B.6 Design Decisions: Combine Operators

4448 We list some design decisions for combine operators (Definition 25).

### 4449 Note 1.

- 4450 • We deliberately restrict index set function  $\Rightarrow_{MDA}^{MDA}$  to compute the index set in the particular  
 4451 dimension  $d$  only, and not of all  $D$  Dimensions (i.e., the function's output is in MDA-IDX-SETS  
 4452 and not MDA-IDX-SETS $^D$ ), because this enables applying combine operator  $\otimes$  iteratively:

4453  $\dots ( (\mathbf{a}_1 \otimes^{(P,Q)} \mathbf{a}_2) \otimes^{(P \cup Q, R)} \mathbf{a}_3 ) \otimes^{(P \cup Q \cup R, \dots)} \dots$

4454  
 4455 <sup>26</sup>Our formalism allows dynamic shapes, by using symbol  $*$  instead of a particular natural number for  $N_i$  (formal details  
 4456 provided in Definition 22), which we aim to discuss thoroughly in future work.

4460 for MDAs  $a_1, a_2, a_3, \dots$  that have index sets  $\Rightarrow_{MDA}^{MDA}(P), \Rightarrow_{MDA}^{MDA}(Q), \Rightarrow_{MDA}^{MDA}(R), \dots$  in dimension  $d$ .  
 4461 This is because the index set of the output MDA changes only in dimension  $d$ , to the new  
 4462 index set  $\Rightarrow_{MDA}^{MDA}(P \cup Q), \Rightarrow_{MDA}^{MDA}(\Rightarrow_{MDA}^{MDA}(P \cup Q) \cup R), \dots$ , so that the output MDA can be used as  
 4463 input for a new application of  $\otimes$ .

- 4464 • It is a design decision that a combine operator's index set function  $\Rightarrow_{MDA}^{MDA}$  takes as input the  
 4465 MDA index set  $P$  or  $Q$  in the particular dimension  $d$  only, rather than the all sets  $(I_1, \dots, I_D)$ . Our approach can be easily extended to index set functions  $\Rightarrow_{MDA}^{MDA} : MDA\text{-IDX-SETS}^D \rightarrow MDA\text{-IDX-SETS}$  that take the entire MDA's index sets as input. However, we avoid this additional complexity, because we are currently not aware of any real-world application that would benefit from such extension.
- 4466 • For better convenience, we could potentially define the meta-type of combine operators (Definition 25) such that meta-parameter  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D)$  is separated from parameter  $(P, Q)$  in a distinct, earlier stage (Definition 20). This would allow automatically deducing  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D)$  from the input MDAs' types, whereas for meta-parameter  $(P, Q)$ , automatic deduction is usually not possible: function  $\Rightarrow_{MDA}^{MDA}$  has to be either invertible for automatically deducing  $P$  and  $Q$  from the input MDAs or invariant under different values of  $P$  and  $Q$ . Consequently, separating parameter  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D)$  in a distinct, earlier stage would allow avoiding explicitly stating this parameter, by deducing it from the input MDAs' type, and only explicitly stating parameter  $(P, Q)$ , e.g.,  $\otimes_2^{<(P,Q)>}(a, b)$  instead of  $\otimes_2^{<(I_1)|(P,Q)>}(a, b)$  for  $a \in T[I_1, \Rightarrow_{MDA}^{MDA}(P)]$  and  $b \in T[I_1, \Rightarrow_{MDA}^{MDA}(Q)]$ .

4467 We avoid separating  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D)$  and  $(P, Q)$  in this work, as we focus on concatenation (Example 13), prefix-sum (Example 23), and point-wise combination (Example 14) only, which have invertible or  $P/Q$ -invariant index set functions, respectively. Consequently, for the practice-relevant combine operators considered in this work, we can deduce all meta-parameters automatically.

## 4486 B.7 Generalized Notion of MDHs

4487 The MDH Definition 26 can be generalized to have an arbitrary algebraic structure as input.

4488 **Definition 34** (Multi-Dimensional Homomorphism). Let

$$4489 \mathcal{A}^\downarrow := (T^{INP}[\stackrel{1}{\Rightarrow_{MDA}^{MDA}}(*), \dots, \stackrel{D}{\Rightarrow_{MDA}^{MDA}}(*)], (\downarrow \otimes_d)_{d \in [1, D]_{\mathbb{N}}})$$

4490 and

$$4491 \mathcal{A}^\uparrow := (T^{OUT}[\stackrel{1}{\Rightarrow_{MDA}^{MDA}}(*), \dots, \stackrel{D}{\Rightarrow_{MDA}^{MDA}}(*)], (\uparrow \otimes_d)_{d \in [1, D]_{\mathbb{N}}})$$

4492 be two algebraic structures, where

$$4493 (\downarrow \otimes_d \in CO^{<\stackrel{d}{\Rightarrow_{MDA}^{MDA}} \mid T^{INP} \mid D \mid d>} )_{d \in [1, D]_{\mathbb{N}}}$$

4494 and

$$4495 (\uparrow \otimes_d \in CO^{<\stackrel{d}{\Rightarrow_{MDA}^{MDA}} \mid T^{OUT} \mid D \mid d>} )_{d \in [1, D]_{\mathbb{N}}}$$

4500 are tuples of combine operators, for  $D \in \mathbb{N}$ ,  $T^{INP}, T^{OUT} \in \text{TYPE}$ ,  $\stackrel{d}{\Rightarrow_{MDA}^{MDA}}, \stackrel{d}{\Rightarrow_{MDA}^{MDA}} : MDA\text{-IDX-SETS} \rightarrow MDA\text{-IDX-SETS}$ , and the two structures' carrier sets

$$4501 T^{INP}[\stackrel{1}{\Rightarrow_{MDA}^{MDA}}(*), \dots, \stackrel{D}{\Rightarrow_{MDA}^{MDA}}(*)]$$

4502 and

$$4503 T^{OUT}[\stackrel{1}{\Rightarrow_{MDA}^{MDA}}(*), \dots, \stackrel{D}{\Rightarrow_{MDA}^{MDA}}(*)]$$

4509 denote the set of MDAs that are in the function domain of combine operators (the star symbol is  
 4510 used for indicating the function range of index functions).

4511 A *Multi-Dimensional Homomorphism* (MDH) from the algebraic structure  $\mathcal{A}^\downarrow$  to the structure  $\mathcal{A}^\uparrow$   
 4512 is any function

$$4513 \quad h^{<I_1, \dots, I_D \in \text{MDA-IDX-SETS}>} : T^{\text{INP}}[\overset{1}{\Rightarrow}_{\text{MDA}}(I_1), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(I_D)] \rightarrow T^{\text{OUT}}[\overset{1}{\Rightarrow}_{\text{MDA}}(I_1), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(I_D)]$$

4515 that satisfies the *homomorphic property*:

$$4517 \quad h(\alpha_1 \downarrow \otimes_d \alpha_2) = h(\alpha_1) \uparrow \otimes_d h(\alpha_2)$$

4519 The MDH Definition 26 is a special case of our generalized MDH Definition 34 above, for  
 4520  $\downarrow \otimes_d = +^{<T^{\text{INP}} | D | d>} \text{ (Example 13).}$

4521 Higher-order function `md_hom` (originally introduced in Definition 27) is defined for the generalized  
 4522 MDH Definition 34 as follows.

4523 **Definition 35** (Higher-Order Function `md_hom`). The higher-order function `md_hom` is of type

$$4525 \quad \text{md\_hom}^{<T^{\text{INP}}, T^{\text{OUT}} \in \text{TYPE} \mid D \in \mathbb{N} \mid (\overset{d}{\Rightarrow}_{\text{MDA}} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS})_{d \in [1, D]_{\mathbb{N}}},}$$

$$4527 \quad (\overset{d}{\Rightarrow}_{\text{MDA}} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS})_{d \in [1, D]_{\mathbb{N}}} > :$$

$$4529 \quad \underbrace{(\text{CO}^{<\overset{1}{\Rightarrow}_{\text{MDA}} \mid T^{\text{OUT}} \mid D \mid 1>} \times \dots \times \text{CO}^{<\overset{D}{\Rightarrow}_{\text{MDA}} \mid T^{\text{OUT}} \mid D \mid D>})}_{\downarrow \otimes_1, \dots, \downarrow \otimes_D} \times$$

$$4532 \quad \underbrace{\text{SF}^{<T^{\text{INP}}, T^{\text{OUT}}>}}_f \times$$

$$4536 \quad \underbrace{(\text{CO}^{<\overset{1}{\Rightarrow}_{\text{MDA}} \mid T^{\text{OUT}} \mid D \mid 1>} \times \dots \times \text{CO}^{<\overset{D}{\Rightarrow}_{\text{MDA}} \mid T^{\text{OUT}} \mid D \mid D>})}_{\uparrow \otimes_1, \dots, \uparrow \otimes_D}$$

$$4540 \quad \rightarrow_p \quad \underbrace{\text{MDH}^{<T^{\text{INP}}, T^{\text{OUT}} \mid D \mid (\overset{d}{\Rightarrow}_{\text{MDA}})_{d \in [1, D]_{\mathbb{N}}}, (\overset{d}{\Rightarrow}_{\text{MDA}})_{d \in [1, D]_{\mathbb{N}}}>}}_{\text{md\_hom}((\otimes_1^\downarrow, \dots, \otimes_D^\downarrow), f, (\otimes_1^\uparrow, \dots, \otimes_D^\uparrow))}$$

4543 The function takes as input a scalar function  $f$  and two tuples of  $D$ -many combine operators  
 4544  $(\downarrow \otimes_1, \dots, \downarrow \otimes_D)$  and  $(\uparrow \otimes_1, \dots, \uparrow \otimes_D)$ , and it yields a function  $\text{md\_hom}((\downarrow \otimes_1, \dots, \downarrow \otimes_D), f, (\uparrow \otimes_1, \dots, \uparrow \otimes_D))$   
 4545 which is defined as:

$$4547 \quad \downarrow \alpha \in T^{\text{INP}}[\overset{1}{\Rightarrow}_{\text{MDA}}(I_1), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(I_D)]$$

4548  $=:$

$$4550 \quad \downarrow \otimes_1 \dots \downarrow \otimes_D \downarrow \alpha^{<i_1, \dots, i_d>} \in T^{\text{INP}}[\overset{1}{\Rightarrow}_{\text{MDA}}(\{i_1\}), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(\{i_D\})]$$

4552  $\mapsto$

$$4553 \quad +_1 \dots +_D \downarrow \alpha_f^{<i_1, \dots, i_d>} \in T^{\text{INP}}[\{i_1\}, \dots, \{i_D\}]$$

4555  $\vec{f}$

$$\begin{aligned}
 4558 \quad & \underset{i_1 \in I_1}{+} \dots \underset{i_D \in I_D}{+} \underset{\mathfrak{a}_f}{\uparrow} f^{<i_1, \dots, i_D>} \in T^{\text{OUT}}[\{i_1\}, \dots, \{i_D\}] \\
 4559 \quad & \mapsto \\
 4560 \quad & \underset{i_1 \in I_1}{\uparrow} \underset{i_D \in I_D}{\uparrow} \underset{\mathfrak{a}}{\uparrow} f^{<i_1, \dots, i_D>} \in T^{\text{OUT}}[\underset{\text{MDA}}{\overset{1}{\Rightarrow}}(\{i_1\}), \dots, \underset{\text{MDA}}{\overset{D}{\Rightarrow}}(\{i_D\})] \\
 4561 \quad & =: \\
 4562 \quad & \underset{\mathfrak{a}}{\uparrow} \in T^{\text{OUT}}[\underset{\text{MDA}}{\overset{1}{\Rightarrow}}(I_1), \dots, \underset{\text{MDA}}{\overset{D}{\Rightarrow}}(I_D)]
 \end{aligned}$$

4566 Here,  $\vec{f}$  denotes the adaption of function  $f$  to operate on MDAs comprising a single value only: it  
 4567 is of type

$$4569 \quad \vec{f}^{<i_1, \dots, i_D \in \mathbb{N}>} : T^{\text{INP}}[\{i_1\}, \dots, \{i_D\}] \rightarrow T^{\text{OUT}}[\{i_1\}, \dots, \{i_D\}]$$

4570 and defined as

$$4572 \quad \vec{f}(x)[i_1, \dots, i_D] := f(x[i_1, \dots, i_D])$$

4573 We refer to the first application of  $\mapsto$  as *de-composition*, to the application of  $\mapsto$  as *scalar function*  
 4574 *application*, and to the second application of  $\mapsto$  as *re-composition*.

4575 For  $\text{md\_hom}((\downarrow \otimes_1, \dots, \downarrow \otimes_d), f, (\uparrow \otimes_1, \dots, \uparrow \otimes_d))$ , we require per definition the homomorphic property  
 4576 (Definition 34), i.e., for each  $d \in [1, D]_{\mathbb{N}}$ , it must hold:

$$\begin{aligned}
 4578 \quad & \text{md\_hom}((\downarrow \otimes_1, \dots, \downarrow \otimes_d), f, (\uparrow \otimes_1, \dots, \uparrow \otimes_d))(\mathfrak{a}_1 \downarrow \otimes_d \mathfrak{a}_2) = \\
 4579 \quad & \quad \text{md\_hom}((\downarrow \otimes_1, \dots, \downarrow \otimes_d), f, (\uparrow \otimes_1, \dots, \uparrow \otimes_d))(\mathfrak{a}_1) \\
 4580 \quad & \quad \quad \uparrow \otimes_d \\
 4581 \quad & \quad \quad \text{md\_hom}((\downarrow \otimes_1, \dots, \downarrow \otimes_d), f, (\uparrow \otimes_1, \dots, \uparrow \otimes_d))(\mathfrak{a}_2)
 \end{aligned}$$

## 4585 B.8 Design Decisions: `md_hom`

4586 We list some design decisions for higher-order function `md_hom` (Definition 27).

4587 **Note 2.** For some MDHs (such as Mandelbrot), the scalar function  $f$  (Definition 27) is dependent  
 4588 on the position in the input MDA, i.e., it takes as arguments, in addition to  $\mathfrak{a}[i_1, \dots, i_D]$ , also the  
 4589 indices  $i_1, \dots, i_D$ . Such MDHs can be easily expressed via `md_hom` after a straightforward type  
 4590 adjustment: type  $\text{SF}^{<T^{\text{INP}}, T^{\text{OUT}}>}$  has to be defined as the set of functions  $f : T^{\text{INP}} \times \text{MDA-IDX-SETS}^D \rightarrow$   
 4591  $T^{\text{OUT}}$  (rather than of functions  $f : T^{\text{INP}} \rightarrow T^{\text{OUT}}$ , as in Definition 27).

4592 Since we do not aim at forcing scalar functions to always take MDA indices as input arguments –  
 4593 for expressing most computations, this is not required (Figure 14) and only causes additional  
 4594 complexity – we assume in the following two different definitions of pattern `md_hom`: one variant  
 4595 exactly as in Definition 27, and one variant with the adjusted type for scalar functions and that  
 4596 passes automatically indices  $i_1, \dots, i_D$  to  $f$ . The two variants can be easily differentiated, via an  
 4597 additional, boolean meta-parameter `USE_MDA_INDICES`: first variant iff `USE_MDA_INDICES = false`  
 4598 and second variant iff `USE_MDA_INDICES = true`.

4599 For simplicity, we focus in this paper on the first variant (as in Definition 27), because, in practice,  
 4600 it is the more common variant, and because all insights presented in this work apply to both  
 4602 variants.

## 4603 B.9 Proof `md_hom` Lemma

4604 We prove Lemma 1.

4605

4607 PROOF. Let  $\mathbf{a}_1 \in T[I_1^{\alpha_1}, \dots, I_D^{\alpha_2}]$  and  $\mathbf{a}_2 \in T[I_1^{\beta_1}, \dots, I_D^{\beta_2}]$  be two arbitrary MDAs for which  $I_d^{\alpha_1} \cap I_d^{\beta_2} = \emptyset$ ,  
 4608 for  $d \in [1, D]_{\mathbb{N}}$  arbitrary but fixed, i.e., the two MDAs are concatenable in dimension  $d$ .

4609 According to Definition 27, we have to show that

4610  $\text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\mathbf{a}_1 \mathbf{a}_2) =$

4612  $\text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\mathbf{a}_1) \oplus_d \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\mathbf{a}_2)$

4613 For this, we first show for arbitrary  $k \in [1, D]_{\mathbb{N}}$  that

$$4615 \dots \underset{i_k \in I_k}{\otimes_k} \underset{i_{k+1} \in I_{k+1}}{\otimes_{k+1}} \dots x|_{\dots, \{i_k\}, \{i_{k+1}\}, \dots} = \dots \underset{i_{k+1} \in I_{k+1}}{\otimes_{k+1}} \underset{i_k \in I_k}{\otimes_k} \dots x|_{\dots, \{i_k\}, \{i_{k+1}\}, \dots}$$

4616 from which follows

$$4617 \underset{i_1 \in I_1}{\otimes_1} \dots \underset{i_D \in I_D}{\otimes_D} x|_{\{i_1\}, \dots, \{i_D\}} = \underset{i_{\sigma(1)} \in I_{\sigma(1)}}{\otimes_{\sigma(1)}} \dots \underset{i_{\sigma(D)} \in I_{\sigma(D)}}{\otimes_{\sigma(D)}} x|_{\{i_1\}, \dots, \{i_D\}}$$

4618 for any permutation  $\sigma : \{1, \dots, D\} \leftrightarrow \{1, \dots, D\}$ . Afterwards, in our assumption above, we can  
 4619 assume w.l.o.g. that  $d = 1$ .

4620 Case 1:  $[\otimes_k = \otimes_{k+1}]$  Follows immediately from the commutativity of  $+$  or  $\vec{\bullet}(\oplus)$  for com-  
 4621 mutative  $\oplus$ , respectively. ✓

4622 Case 2:  $[\otimes_k \neq \otimes_{k+1}]$  Trivial, as it is either  $\otimes_k = +$  or  $\otimes_{k+1} = +$ , and

$$4623 \underset{i_d \in I_d}{\oplus_d} x|_{\dots, \{i_d\}, \dots} [i_1, \dots, i_D] = (x|_{\dots, \{i_d\}, \dots}) [i_1, \dots, i_D]$$

4624 according to the definition of MDA concatenation  $+$  (Example 13). ✓

4625 Let now be  $d = 1$  (see assumption above), it holds:

$$4626 \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\mathbf{a}_1 \mathbf{a}_2) \\ 4627 = \underset{i_1 \in I_1}{\otimes_1} \dots \underset{i_D \in I_D}{\otimes_D} \vec{f}((\mathbf{a}_1 \mathbf{a}_2)|_{\{i_1\} \times \dots \times \{i_D\}}) \\ 4628 = \underset{i_1 \in I_1^{\alpha_1}}{\otimes_1} \dots \underset{i_D \in I_D^{\alpha_D}}{\otimes_D} \vec{f}(\mathbf{a}_1|_{\{i_1\} \times \dots \times \{i_D\}}) \oplus_1 \underset{i_1 \in I_1^{\beta_2}}{\otimes_1} \dots \underset{i_D \in I_D^{\beta_D}}{\otimes_D} \vec{f}(\mathbf{a}_2|_{\{i_1\} \times \dots \times \{i_D\}}) \\ 4629 = \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\mathbf{a}_1) \oplus_1 \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\mathbf{a}_2) \quad \checkmark$$

4630 □

## 4631 B.10 Examples of Index Functions

4632 We present examples of index functions (Definition 29).

4633 **Example 21** (Matrix-Vector Multiplication). The index functions we use for expressing Matrix-  
 4634 Vector Multiplication (MatVec) are:

- 4635 • Input Matrix:

4636  $\text{id}_{\mathbf{X}}(i, k) := (i, k) \in \text{MDA-IDX-to-BUF-IDX}^{< D=2, D_1=2 \Rightarrow \text{BUF}^{\text{MDA}}}$

4637 for  $\Rightarrow_{\text{BUF}}^{\text{MDA}}(I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_1^{\text{MDA}})]_{\mathbb{N}_0}, [0, \max(I_2^{\text{MDA}})]_{\mathbb{N}_0}$

- 4638 • Input Vector:

4639  $\text{id}_{\mathbf{X}}(i, k) := (k) \in \text{MDA-IDX-to-BUF-IDX}^{< D=2, D_1=1 \Rightarrow \text{BUF}^{\text{MDA}}}$

4640 for  $\Rightarrow_{\text{BUF}}^{\text{MDA}}(I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_2^{\text{MDA}})]_{\mathbb{N}_0}$

- Output Vector:

`idx(i, k) := (i) ∈ MDA-IDX-to-BUF-IDX<D=2, D1=1| ⇒BUFMDA`

for  $\Rightarrow_{\text{BUF}}^{\text{MDA}} (I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_1^{\text{MDA}})]_{\mathbb{N}_0}$

**Example 22** (Jacobi 1D). The index functions we use for expressing Jacobi 1D (Jacobi1D) are:

- Input Buffer, 1. Access:

`idx(i) := (i + 0) ∈ MDA-IDX-to-BUF-IDX<D=1, D1=1| ⇒MDA BUF`

for  $\Rightarrow_{\text{BUF}}^{\text{MDA}} (I_1^{\text{MDA}}) := [0, \max(I_1^{\text{MDA}}) + 0]_{\mathbb{N}_0}$

- Input Buffer, 2. Access:

$\text{idx}(i) := (i + 1) \in \text{MDA-IDX-to-BUF-IDX}^{< D=1, D_1=1 | \Rightarrow \text{MDA-BUF} }$

for  $\Rightarrow_{\text{BLIE}}^{\text{MDA}}(I_1^{\text{MDA}}) := \lceil 0, \max(I_1^{\text{MDA}}) + 1 \rceil_{\mathbb{N}_0}$

- Input Buffer, 3. Access:

`iidx(i) := (i + 2) ∈ MDA-TDX-to-BUF<D=1,D1=1|⇒MDA-BUF`

for  $\Rightarrow_{\text{MDA}}^{\text{MDA}}(J_1^{\text{MDA}}) \equiv [0, \max(J_1^{\text{MDA}}) + 2]_{\mathbb{N}}$

- Output Buffer:

$\text{isr}(i) : (i) \in \text{MDA\_IDY} \rightarrow \text{RUE\_IDY}^{< D=1, D_1=1} \Rightarrow^{\text{MDA}}_{\text{RUE}}$

for  $\rightarrow^{\text{MDA}}(I^{\text{MDA}}) := [0 \text{ } \max(I^{\text{MDA}})]_{\text{v}}$

## B.11 Representation of Scalar Values

Scalar values can be considered as 0-dimensional BUFs (Definition 28). Consequently, in Definition 28, the cartesian product  $[0, N_1]_{\mathbb{N}_0} \times \dots \times [0, N_D]_{\mathbb{N}_0}$  is empty for  $D = 0$ , and thus results in the neutral element of the cartesian product. As any singleton set can be considered as neutral element of cartesian product (up to bijection), we define the set  $\{\epsilon\}$  containing the dedicated symbol epsilon only, as the uniquely determined neutral element of cartesian product (inspired by the notation of the *empty word*).

We often refrain from explicitly stating symbol  $\epsilon$ , e.g., by writing  $b$  instead of  $b[\epsilon]$  for accessing a BUF, or  $(i_1, \dots, i_D) \rightarrow ()$  instead of  $(i_1, \dots, i_D) \rightarrow (\epsilon)$  for index functions.

Note that alternatively, scalar values can be considered as any multi-dimensional BUF containing a single element only. For example, a scalar value  $s$  can be represented as 1-dimensional BUF  $b_{1D}[0] := s$ , or a 2-dimensional BUF  $b_{2D}[0, 0] := s$ , or a 3-dimensional BUF  $b_{3D}[0, 0, 0] := s$ , etc. However, this results in an ambiguous representation of scalar values, which we aim to avoid by considering scalars as 0-dimensional BUFs, as described above.

## B.12 Runtime Complexity of Histograms

Our implementation of Histograms (Subfigure 5 in Figure 14) has a *work complexity* of  $\mathcal{O}(E * B)$ , where  $E$  is the number of elements to check and  $B$  the number of bins, i.e., our MDH Histogram implementation is not work efficient. However, our Histograms' *step complexity* [Harris et al. 2007] is  $\mathcal{O}(\log(E))$ : step complexity is often used for parallel algorithms and assumes an infinite number of cores, i.e., we can ignore in our implementation of Histogram the concatenation dimension

4705  $B$  (which has a step complexity of  $\mathcal{O}(1)$ ) and take into account its reduction dimension  $B$  only,  
 4706 which has a step complexity of  $\log(B)$  (parallel reduction [Harris et al. 2007]). In contrast, related  
 4707 approaches [Henriksen et al. 2020] are often work efficient, by having a work complexity of  
 4708  $\mathcal{O}(B)$ ; however, their high work efficiency is at the cost of their step complexity which is also  $\mathcal{O}(B)$ ,  
 4709 rather than  $\mathcal{O}(\log(B))$  as for our implementation in Subfigure 5, thereby being asymptotically less  
 4710 efficient for parallel machines consisting of many cores. Our future work will show that the work-  
 4711 efficient Histogram implementation introduced in Henriksen et al. [2020] can also be expressed in  
 4712 our approach, by using for scalar function  $f$  an optimized micro kernel for Histogram computation,  
 4713 similarly as done in the related work.

## B.13 Combine Operator of Prefix-Sum Computations

We define *prefix-sum* which is the combine operator of compute pattern `scan` and example MBBS in Section 2.4.

**Example 23 (Prefix-Sum).** We define *prefix-sum*, according to a binary function  $\oplus : T \times T \rightarrow T$  (e.g. addition), as function  $\otimes_{\text{prefix-sum}}$  of type

where  $id : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$  is the identity function on MDA index sets. The function is computed as (w.l.o.g., we assume  $\max(P) < \max(O)$  for commutativity):

$$\begin{aligned} & \oplus_{\text{prefix-sum}}^{<T|D|d| (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D), (P, Q)>} (\oplus)(\mathfrak{a}_1, \mathfrak{a}_2)[i_1, \dots, i_d, \dots, i_D] \\ & := \begin{cases} \mathfrak{a}_1[i_1, \dots, i_d, \dots, i_D] & , i_d \in P \\ \mathfrak{a}_1[i_1, \dots, \max(P), \dots, i_D] \oplus \mathfrak{a}_2[i_1, \dots, i_d, \dots, i_D] & , i_d \in Q \end{cases} \end{aligned}$$

Function  $\oplus^{< T \mid D \mid d >}$  (meaning:  $\oplus_{\text{prefix-sum}}$  is partially applied to ordinary function parameter  $\oplus$ ; formal details provided in Definition 21) is a combine operator of type  $\text{CO}^{< id \mid T \mid D \mid d >}$  for any binary operator  $\oplus : T \times T \rightarrow T$ .

## C ADDENDUM SECTION 3

## Basic Building Blocks

**Definition 36** (Low-Level MDA). Let be  $L \in \mathbb{N}$  (representing an ASM's number of layers) and  $D \in \mathbb{N}$  (representing an MDH's number of dimensions). Let further be  $P = ((P_1^1, \dots, P_D^1), \dots, (P_1^L, \dots, P_D^L)) \in \mathbb{N}^{L \times D}$  an arbitrary tuple of  $L$ -many  $D$ -tuples of positive natural numbers,  $T \in \text{TYPE}$  a scalar type, and  $I := ((I_d^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \in \text{MDA-IDX-SETS})_{d \in [1, D]_{\mathbb{N}}})^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}$  an arbitrary collection of  $D$ -many MDA index sets (Definition 24) for each particular choice of indices  $(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1, \dots, (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L$  (illustrated in Figure 16).

<sup>27</sup> Analogously to Notation 1, we identify each  $P^l, l \in \mathbb{N}$  implicitly also with the interval  $[0, P^l]_{\mathbb{N}_0}$  (inspired by set theory).

4754 An  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned *low-level MDA* that has scalar type  $T$  and index sets  $I$   
 4755 is any function  $a_{ll}$  of type:

4756

$$4757 \quad \overbrace{a_{ll}^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}}^{\text{Partitioning: Layer 1}} : \\ 4758 \quad \overbrace{I_1^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \times \dots \times I_D^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>}}^{\text{Partitioning: Layer } L} \rightarrow T \\ 4759 \\ 4760 \\ 4761$$

4762 **Definition 37** (Low-Level BUF). Let be  $L \in \mathbb{N}$  (representing an ASM's number of layers) and  
 4763  $D \in \mathbb{N}$  (representing an MDH's number of dimensions). Let further  $P = ((P_1^1, \dots, P_D^1), \dots,$   
 4764  $(P_1^L, \dots, P_D^L)) \in \mathbb{N}^{L \times D}$  be an arbitrary tuple of  $L$ -many  $D$ -tuples of positive natural numbers,  $T \in$   
 4765 TYPE a scalar type, and  $N := ((N_d^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \in \mathbb{N})_{d \in [1, D]_{\mathbb{N}}})^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}$  be a BUF's size (Definition 28) for each particular choice of  $p_1^1, \dots, p_D^L$ .

4766 An  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned *low-level BUF* that has scalar type  $T$  and size  $N$  is  
 4767 any function  $b_{ll}$  of type ( $\leftrightarrow$  denotes bijection):

4768

$$4769 \quad \overbrace{b_{ll}^{<\text{MEM} \in [1, \text{NUM\_MEM\_LYRS}]_{\mathbb{N}} | \sigma: [1, D]_{\mathbb{N}} \leftrightarrow [1, D]_{\mathbb{N}} > <(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}}^{\text{Memory Region}} : \\ 4770 \quad \overbrace{[0, N_1^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>}]_{\mathbb{N}_0} \times \dots \times [0, N_D^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>}]_{\mathbb{N}_0}}^{\text{Memory Layout}} \rightarrow T \\ 4771 \\ 4772 \\ 4773 \\ 4774 \\ 4775$$

4776 We refer to *MEM* as low-level BUF's *memory region* and to  $\sigma$  as its *memory layout*, and we refer to  
 4777 the function

4778

$$4779 \quad \overbrace{b_{ll}^{\text{trans} <\text{MEM} \in [1, \text{NUM\_MEM\_LYRS}]_{\mathbb{N}} | \sigma: [1, D]_{\mathbb{N}} \leftrightarrow [1, D]_{\mathbb{N}} > <(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}}^{\text{Memory Region}} : \\ 4780 \quad [0, N_{\sigma(1)}^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>}]_{\mathbb{N}_0} \times \dots \times [0, N_{\sigma(D)}^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>}]_{\mathbb{N}_0} \rightarrow T \\ 4781 \\ 4782 \\ 4783$$

4784 that is defined as

4785

$$4786 \quad b_{ll}^{\text{trans} <\text{MEM} | \sigma > <p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} (i_{\sigma(1)}, \dots, i_{\sigma(D)}) := b_{ll}^{<\text{MEM} | \sigma > <p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} (i_1, \dots, i_D)$$

4787 as  $b_{ll}$ 's *transposed function representation*.

4788 **Definition 38** (Low-Level Combine Operator). Let be  $L \in \mathbb{N}$  (representing an ASM's number of layers)  
 4789 and  $D \in \mathbb{N}$  (representing an MDH's number of dimensions). Let further be  $\otimes \in \text{CO}^{<\Rightarrow_{\text{MDA}}^{\text{MDA}} | T | D | d >}$   
 4790 an arbitrary  $D$ -dimensional combine operator (Definition 25).

4791 The *low-level representation*  $\otimes^{<(l_{\text{ASM}}, d_{\text{ASM}}) \in \text{ASM-LVL}>}$  of operator  $\otimes$  is a function that for each pair

4792

$$4793 \quad (l_{\text{ASM}}, d_{\text{ASM}}) \in \{ (l, d) \mid l \in [1, L]_{\mathbb{N}}, d \in [1, D]_{\mathbb{N}} \} \subset \text{ASM-LVL}$$

4794 has the same type and semantics as  $\otimes$ :

4795

$$4796 \quad \otimes^{<l_{\text{ASM}}, d_{\text{ASM}}>} \in \text{CO}^{<\Rightarrow_{\text{MDA}}^{\text{MDA}} | T | D | d >}, \quad \otimes^{<l_{\text{ASM}}, d_{\text{ASM}}>} (a, b) := \otimes (a, b)$$

4797 i.e.,  $\otimes^{<l_{\text{ASM}}, d_{\text{ASM}}>}$  works exactly as combine operator  $\otimes$ , but its type is enriched with a meta-parameter  
 4798 that captures the notation of an ASM layer  $l_{\text{ASM}} \in [1, L]_{\mathbb{N}}$  and dimension  $d_{\text{ASM}} \in [1, D]_{\mathbb{N}}$ .

4803	$b_1^{\text{IB}}, \dots, b_{B^{\text{IB}}}^{\text{IB}}$	$\xrightarrow{\text{inp\_view}}$	$a$	$=:$
4804	$\#_1^{\leftrightarrow_{\text{prt-ass}}(1,1)}$	$\dots$	$\#_D^{\leftrightarrow_{\text{prt-ass}}(1,D)}$	
4805	$p_1^1 \in \#PRT(1,1)$		$p_D^1 \in \#PRT(1,D)$	
4806	$\rightarrow b_1^{\text{IB}} : \downarrow\text{-mem}(1,1)[\sigma_{\downarrow\text{-mem}}(1,1)(1), \dots, \sigma_{\downarrow\text{-mem}}(1,1)(D_1^{\text{IB}})]$	$\rightarrow b_1^{\text{IB}} : \downarrow\text{-mem}(1,D)[\sigma_{\downarrow\text{-mem}}(1,D)(1), \dots, \sigma_{\downarrow\text{-mem}}(1,D)(D_1^{\text{IB}})]$		
4807	$\vdots$		$\vdots$	
4808	$b_{B^{\text{IB}}}^{\text{IB}} : \downarrow\text{-mem}(1,1)[\sigma_{\downarrow\text{-mem}}(1,1)(1), \dots, \sigma_{\downarrow\text{-mem}}(1,1)(D_{B^{\text{IB}}}^{\text{IB}})]$	$b_{B^{\text{IB}}}^{\text{IB}} : \downarrow\text{-mem}(1,D)[\sigma_{\downarrow\text{-mem}}(1,D)(1), \dots, \sigma_{\downarrow\text{-mem}}(1,D)(D_{B^{\text{IB}}}^{\text{IB}})]$	$\hookrightarrow \sigma_{\downarrow\text{-ord}}(1,D)$	
4809	$\hookrightarrow \sigma_{\downarrow\text{-ord}}(1,1)$			
4810	$\vdots$		$\vdots$	
4811	$\vdots$		$\vdots$	
4812	$\#_1^{\leftrightarrow_{\text{prt-ass}}(L,1)}$	$\dots$	$\#_D^{\leftrightarrow_{\text{prt-ass}}(L,D)}$	
4813	$p_1^1 \in \#PRT(L,1)$		$p_D^1 \in \#PRT(L,D)$	
4814	$\rightarrow b_1^{\text{IB}} : \downarrow\text{-mem}(L,1)[\sigma_{\downarrow\text{-mem}}(L,1)(1), \dots, \sigma_{\downarrow\text{-mem}}(L,1)(D_1^{\text{IB}})]$	$\rightarrow b_1^{\text{IB}} : \downarrow\text{-mem}(L,D)[\sigma_{\downarrow\text{-mem}}(L,D)(1), \dots, \sigma_{\downarrow\text{-mem}}(L,D)(D_1^{\text{IB}})]$		
4815	$\vdots$		$\vdots$	
4816	$b_{B^{\text{IB}}}^{\text{IB}} : \downarrow\text{-mem}(L,1)[\sigma_{\downarrow\text{-mem}}(L,1)(1), \dots, \sigma_{\downarrow\text{-mem}}(L,1)(D_{B^{\text{IB}}}^{\text{IB}})]$	$b_{B^{\text{IB}}}^{\text{IB}} : \downarrow\text{-mem}(L,D)[\sigma_{\downarrow\text{-mem}}(L,D)(1), \dots, \sigma_{\downarrow\text{-mem}}(L,D)(D_{B^{\text{IB}}}^{\text{IB}})]$	$\hookrightarrow \sigma_{\downarrow\text{-ord}}(L,D)$	
4817	$\hookrightarrow \sigma_{\downarrow\text{-ord}}(L,1)$			
4818	$\downarrow a_f^{< p_1^1, \dots, p_D^1   \dots   p_1^L, \dots, p_D^L >}$			
4819				
4820				

### (a) De-Composition Phase

```

4821
4822  $\downarrow \mathfrak{a}_f^{<p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L>} \xrightarrow{\vec{f}} \uparrow \mathfrak{a}_f^{<p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L>}$ 
4823  $\rightarrow <\sigma_{f\text{-ord}}(1,1), \dots, \sigma_{f\text{-ord}}(L,D)>$ 
4824  $\rightarrow \leftrightarrow_{f\text{-ass}}(1,1), \dots, \leftrightarrow_{f\text{-ord}}(L,D)>$ 
4825  $\rightarrow \mathfrak{b}_1^{\text{IB}} : f^\dagger\text{-mem}[\sigma_{f^\dagger\text{-mem}}(1), \dots, \sigma_{f^\dagger\text{-mem}}(D_1^{\text{IB}})], \dots, \mathfrak{b}_{B^{\text{IB}}}^{\text{IB}} : f^\dagger\text{-mem}[\sigma_{f^\dagger\text{-mem}}(1), \dots, \sigma_{f^\dagger\text{-mem}}(D_{B^{\text{IB}}}^{\text{IB}})]$ 
4826  $\rightarrow \mathfrak{b}_1^{\text{OB}} : f^\dagger\text{-mem}[\sigma_{f^\dagger\text{-mem}}(1), \dots, \sigma_{f^\dagger\text{-mem}}(D_1^{\text{OB}})], \dots, \mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} : f^\dagger\text{-mem}[\sigma_{f^\dagger\text{-mem}}(1), \dots, \sigma_{f^\dagger\text{-mem}}(D_{B^{\text{OB}}}^{\text{OB}})]$ 
4827
4828 (b) Scalar Phase
4829
4830  $\mathfrak{b}_1^{\text{OB}}, \dots, \mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} \xrightarrow{\text{out-view}} \uparrow \mathfrak{a} :=$ 
4831  $\oplus_1^{\leftrightarrow_{\text{prt-ass}}(1,1)} \dots \oplus_D^{\leftrightarrow_{\text{prt-ass}}(1,D)}$ 
4832  $p_1^1 \in \#PRT(1,1) \dots p_D^1 \in \#PRT(1,D)$ 
4833  $\rightarrow \mathfrak{b}_1^{\text{OB}} : \uparrow\text{-mem}(1,1)[\sigma_{\uparrow\text{-mem}}(1,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,1)(D_1^{\text{OB}})] \rightarrow \mathfrak{b}_1^{\text{OB}} : \uparrow\text{-mem}(1,D)[\sigma_{\uparrow\text{-mem}}(1,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,D)(D_1^{\text{OB}})]$ 
4834  $\vdots \vdots$ 
4835  $\mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(1,1)[\sigma_{\uparrow\text{-mem}}(1,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,1)(D_{B^{\text{OB}}}^{\text{OB}})] \mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(1,D)[\sigma_{\uparrow\text{-mem}}(1,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,D)(D_{B^{\text{OB}}}^{\text{OB}})]$ 
4836  $\hookrightarrow \sigma_{\uparrow\text{-ord}}(1,1) \hookrightarrow \sigma_{\uparrow\text{-ord}}(1,D)$ 

```

### (b) Scalar Phase

4829	$b_1^{\text{OB}}, \dots, b_{B^{\text{OB}}}^{\text{OB}}$	$\xrightarrow{\text{out\_view}} \uparrow a :=$
4830		
4831	$\oplus_1^{\leftrightarrow_{\text{prt-ass}}(1,1)}$	$\dots$
4832	$p_1^i \in \#PRT(1,1)$	$\oplus_D^{\leftrightarrow_{\text{prt-ass}}(1,D)}$
4833	$\rightarrow b_1^{\text{OB}} : \uparrow\text{-mem}(1,1)[\sigma_{\uparrow\text{-mem}}(1,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,1)(D_1^{\text{OB}})]$	$\rightarrow b_1^{\text{OB}} : \uparrow\text{-mem}(1,D)[\sigma_{\uparrow\text{-mem}}(1,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,D)(D_1^{\text{OB}})]$
4834	$\vdots$	$\vdots$
4835	$b_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(1,1)[\sigma_{\uparrow\text{-mem}}(1,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,1)(D_{B^{\text{OB}}}^{\text{OB}})]$	$b_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(1,D)[\sigma_{\uparrow\text{-mem}}(1,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,D)(D_{B^{\text{OB}}}^{\text{OB}})]$
4836	$\underbrace{\leftrightarrow \sigma_{\uparrow\text{-ord}}(1,1)}$	$\underbrace{\leftrightarrow \sigma_{\uparrow\text{-ord}}(1,D)}$
4837	$\vdots$	$\vdots$
4838	$\vdots$	$\vdots$
4839	$\oplus_1^{\leftrightarrow_{\text{prt-ass}}(L,1)}$	$\dots$
4840	$p_1^L \in \#PRT(L,1)$	$\oplus_D^{\leftrightarrow_{\text{prt-ass}}(L,D)}$
4841	$\rightarrow b_1^{\text{OB}} : \uparrow\text{-mem}(L,1)[\sigma_{\uparrow\text{-mem}}(L,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(L,1)(D_1^{\text{OB}})]$	$\rightarrow b_1^{\text{OB}} : \uparrow\text{-mem}(L,D)[\sigma_{\uparrow\text{-mem}}(L,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(L,D)(D_1^{\text{OB}})]$
4842	$\vdots$	$\vdots$
4843	$b_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(L,1)[\sigma_{\uparrow\text{-mem}}(L,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(L,1)(D_{B^{\text{OB}}}^{\text{OB}})]$	$b_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(L,D)[\sigma_{\uparrow\text{-mem}}(L,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(L,D)(D_{B^{\text{OB}}}^{\text{OB}})]$
4844	$\underbrace{\leftrightarrow \sigma_{\uparrow\text{-ord}}(L,1)}$	$\underbrace{\leftrightarrow \sigma_{\uparrow\text{-ord}}(L,D)}$
4845	$\uparrow a \leftarrow p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L$	

### (c) Re-Composition Phase

Fig. 31. Generic low-level expression for data-parallel computations

4852 **C.1 Generic Low-Level Expression**

4853 Figure 31 shows a generic expression in our low-level representation: it targets an arbitrary but  
 4854 fixed  $L$ -layered ASM instance, and it implements – on low level – the generic instance of our  
 4855 high-level expression in Figure 30. Inserting into the low-level expression a particular value for  
 4856 ASM’s numbers of layer  $L$ , as well as particular values for the generic parameters of the high-level  
 4857 expression in Figure 30 (dimensionality  $D$ , combine operators  $\otimes_1, \dots, \otimes_d$ , and input/output views)  
 4858 results in an instance of the expression in Figure 31 that remains generic in tuning parameters  
 4859 only; this auto-tunable instance will be the focus of our discussion in the remainder of this section.  
 4860

4861 In Section 4, we show that we fully automatically compute the auto-tunable low-level expression  
 4862 for a concrete ASM instance and high-level expression, and we automatically optimize this tunable  
 4863 expression for a particular target device and characteristics of the input and output data using  
 4864 auto-tuning [Rasch et al. 2021]. Our final outcome is a concrete (non-generic) low-level expression  
 4865 (as in Figure 15) that is auto-tuned for the particular target device (represented via an ASM  
 4866 instance, e.g., ASM instance  $\text{ASM}_{\text{CUDA}}$  when targeting an NVIDIA Ampere GPU) and high-level  
 4867 MDH expression; from this auto-tuned low-level expression, we can straightforwardly generate  
 4868 executable program code, because all the major optimization decisions have already been made in  
 4869 the previous auto-tuning step. Our overall approach is illustrated in Figure 4.  
 4870

4871 **C.2 Examples**

4872 Figures 32-35 show how our low-level representation is used for expressing the (de/re)-compositions  
 4873 of concrete, state-of-the-art implementations, using the popular example of matrix multiplication  
 4874 (MatMul) on a real-world input size taken from the ResNet-50 [He et al. 2015] deep learning  
 4875 neural network (training phase). To challenge our formalism: i) we express implementations gener-  
 4876 ated and optimized according to notably different approaches: scheduling approach TVM using its  
 4877 recent Ansor [Zheng et al. 2020a] optimization engine which is specifically designed and optimized  
 4878 toward deep learning computations; polyhedral compilers PPCG and Pluto with auto-tuned tile  
 4879 sizes; ii) we consider optimizations for two fundamentally different kinds of architectures: NVIDIA  
 4880 Ampere GPU and Intel Skylake CPU. We consider our study as challenging for our formalism,  
 4881 because it needs to express – in the same formal framework – the (de/re)-compositions of imple-  
 4882 ments generated and optimized according to notably different approaches (scheduling-based  
 4883 and polyhedral-based) and for significantly different kinds of architectures (GPU and CPU). Experi-  
 4884 mental results for TVM, PPCG, and Pluto (including the MatMul study used in this section) are  
 4885 presented and discussed in Section 5, as the focus of this section is on analyzing and discussing the  
 4886 expressivity of our low-level representation, rather than its experimental evaluation.  
 4887

4888 In Figures 32-35, we list our low-level representation’s particular tuning parameter values for  
 4889 expressing the TVM- and PPCG/Pluto-generated implementations, which concisely describe the  
 4890 (de/re)-composition strategies used by TVM, PPCG and Pluto for MatMul on GPU or CPU for  
 4891 the ResNet-50’s input size; inserting these tuning parameter values into our generic low-level  
 4892 expression in Figure 31 results in the formal representation of the (de/re)-composition strategies  
 4893 used by TVM, PPCG and Pluto.

4894 In the following, we describe the columns of the tables in Figures 32-35, each of which listing  
 4895 particular values of tuning parameters in Table 1: column 0 in Figures 32-35 lists values of tuning  
 4896 parameter 0 in Table 1, column D1 of tuning parameter D1, etc. As all four tables follow the  
 4897 same structure, we focus on describing the particular example table in Figure 32 (example chosen  
 4898 arbitrarily), which shows the (de/re)-composition used in TVM’s generated CUDA code for MatMul  
 4899 on NVIDIA Ampere GPU using input matrices of sizes  $16 \times 2048$  and  $2048 \times 1000$  taken from  
 4900

4901 ResNet-50<sup>28</sup>. Note that for clarity, we use in the figures domain-specific aliases, instead of numerical  
 4902 values, to better differentiate between different ASM layers and memory regions. For example,  
 4903 we use in Figure 32 as aliases DEV := 1, SHR := 2, and REG := 3 to refer to CUDA’s three memory  
 4904 layers (device memory layer DEV, shared memory layer SHR, and register memory layer REG), and  
 4905 we use DM := 1, SM := 2, and RM := 3 to refer to CUDA’s memory regions device DM, shared SM, and  
 4906 register RM; aliases BLK := 4 and THR := 5 refer to CUDA’s core layers which are programmed via  
 4907 blocks and threads in CUDA. We differentiate between memory layers and memory regions for  
 4908 the following reason: for example, using tuning parameter 0 in Table 1, we partition input data  
 4909 hierarchically for each particular memory layer of the target architecture (sometime possibly into  
 4910 one part only, which is equivalent to not partitioning). However, depending on the value of tuning  
 4911 parameter D3, we do not necessarily copy the input’s parts always into the corresponding memory  
 4912 regions (e.g., a part on SHR layer is not necessarily copied into shared memory SM), for example,  
 4913 when the architecture provides automatically managed memory regions (as caches in CPUs) or  
 4914 when only some parts of the input are accessed multiple times (e.g., the input vector in the case of  
 4915 matrix-vector multiplication, but not the input matrix), etc.

4916 *Column 0.* The column lists the particular number of parts (a.k.a. *tiles*) used in TVM’s multi-  
 4917 layered, multi-dimensional partitioning strategy for MatMul on the ResNet-50’s input matrices of  
 4918 sizes  $16 \times 2048$  and  $2048 \times 1000$ . We can observe from this column that the input MDA, which is  
 4919 initially of size  $(16, 1000, 2048)$  for the ResNet-50’s input matrices (follows from MatMul’s particular  
 4920 input view functions shown in Figure 14), is partitioned into  $(2 * 50 * 1)$ -many parts (indicated by  
 4921 the first three rows in column 1): 2 parts in the first dimension, 50 parts in the second dimension,  
 4922 and 1 part in the third dimension. Each of these parts is then further partitioned into  $(2 * 1 * 8)$ -many  
 4923 parts (rows 4,5,6), and these parts are again partitioned into  $(4 * 20 * 1)$ -many further parts (rows  
 4924 7,8,9), etc.

4925 *Columns D1, S1, R1.* These 3 columns describe the order in which parts are processed in the  
 4926 different phases: de-composition (column D1), scalar phase (column S1), and re-composition (col-  
 4927 umn R1). For example, we can observe from column R1 that TVM’s CUDA code first combines parts  
 4928 on layer 1 in dimensions 1, 2, 3 (indicated via 1, 2, 3 in rows 1, 2, 3 of column R1); afterwards, parts  
 4929 on layer 3 in dimensions 1, 2, 3 are processed (indicated via 1, 2, 3 in rows 7, 8, 9 of column R1), etc.

4930 Note that TVM uses the same order in all three phases (i.e., columns D1, S1, R1 coincide). Most  
 4931 likely this is because in CUDA, iteration over memory tiles are programmed via for-loops such that  
 4932 columns D1, S1, R1 represent loop orders – using the same order of loops in columns D1, S1, R1  
 4933 consequently allows TVM to generate the loop nests of the three different phases as fused, in one  
 4934 single loop nest (rather than three individual nests), which usually achieves high performance in  
 4935 CUDA.

4936 Note further that the order of parts that are processed in parallel (columns D2, S2, R2, described  
 4937 in the next paragraph, determine if parts are processed in parallel or not) effects when results of  
 4938 blocks and threads are combined (a.k.a. *parallel reduction* [Harris et al. 2007]), e.g., early in the  
 4939 computation and thus often (but probably at the cost of low memory consumption) or late and thus  
 4940 only once (but requiring possibly more memory), etc.

4941 *Columns D2, S2, R2.* The columns determine how computations are assigned to the target archi-  
 4942 tecture. In our example in Figure 32, the  $(2 * 50 * 1)$ -many parts in the first partitioning layer (rows  
 4943 1-3) are assigned to CUDA blocks (BLK); consequently, each such part is computed by an individual  
 4944 block (i.e., TVM uses a so-called *grid size* of 2, 50, 1 in its generated CUDA code for MatMul). The  
 4945  $(4 * 20 * 1)$ -many parts in the third partitioning layer (rows 7-9) are processed by CUDA threads

4946 <sup>28</sup>For the interested reader, TVM’s corresponding, Ansor-generated scheduling program is presented in Section C.10.

4950 (THR) (i.e., the CUDA *block size* in the TVM-generated code is 4, 20, 1). All other parts, e.g., those  
 4951 belonging to the  $(2 * 1 * 8)$ -many parts in the second partitioning layer (rows 4-6), are assigned to  
 4952 CUDA's memory layers (denoted as DEV, SHR, REG in Figure 32) and thus processed sequentially,  
 4953 via for-loops.

4954     *Columns* D3, S3, S5, R3. While column 0 sets the multi-layered, multi-dimensional partitioning  
 4955 strategy used in TVM's CUDA code (according to the CUDA model's multi-layered memory and  
 4956 core hierarchies, shown in Example 11), column 0 does not indicate how CUDA's fast memory  
 4957 regions are exploited in the TVM-generated CUDA code for MatMul – column 0 only describes  
 4958 TVM's partitioning of the input/output computations such that parts of the input and output data  
 4959 can potentially fit into fast memory resources.  
 4960

4961     The actual assignment of parts to memory regions is set via columns D3 (memory regions to be  
 4962 used for input data), columns S3 and S5 (memory regions to be used for the scalar computations),  
 4963 and column R3 (memory regions to be used for storing the computed output data). For example,  
 4964 column D3 indicates that in TVM's CUDA code for MatMul, parts of the *A* and *B* input matrices are  
 4965 stored in CUDA's fast shared memory SM (rows 10-15), and column R3 indicates that each thread  
 4966 computes its results within CUDA's faster register memory RM (rows 10-15).  
 4967

4968     Our flexibility of separating the tiling strategy (Parameter 0) from the actual usage of memory  
 4969 regions (columns D3, S3, S5, R3) allows us, for example, to store parts belonging to one input buffer  
 4970 into fast memory resources (e.g., the input vector of matrix-vector multiplication, whose values  
 4971 are accessed multiple times), but not parts of other buffers (e.g., the input matrix of matrix-vector  
 4972 multiplication, whose values are accessed only once) or only subparts of buffers, etc.  
 4973

4974     Note that in the case of Figure 35 which shows Pluto's (de/re)-composition for OpenMP code,  
 4975 the memory tags in columns D3, S3, S5, R3 have no effect on the generated code: OpenMP relies  
 4976 on its target device's implicit memory management mechanism (e.g., CPU caches), rather than  
 4977 exposing the memory hierarchy explicitly to the programmer. Consequently, the memory tags are  
 4978 ignored by our OpenMP code generator and only emphasize the implementer's intention, e.g., that  
 4979 in Figure 35, each of the  $(2 * 962 * 218)$ -many tiles (rows 7-9) are intended by the implementer  
 4980 to be processed in L2 memory (even though this decision is eventually made automatically in the  
 4981 CPU hardware's automatic cache engine).  
 4982

4983     *Columns* D4, S4, S6, R4. These columns set the memory layout to be used for memory allocations  
 4984 in the CUDA code. TVM chooses in all cases CUDA's standard transpositions layout (indicated by  
 4985  $[1, 2]$ , called *row-major* layout, instead of  $[2, 1]$  which is known as *column-major*). Since the  
 4986 same layout and memory region is used on consecutive layers, the same memory allocation is  
 4987 re-used in the CUDA code. For example, only one memory region is allocated in shared memory  
 4988 for input buffer *A*; the region is accessed in the computations of all SHR, REG tiles as well as CC tiles  
 4989 in dimensions *x* and *z*. Similarly, only one region is allocated in register memory for computing  
 4990 the results of SHR, REG, and DEV tiles.  
 4991

4992  
 4993  
 4994  
 4995  
 4996  
 4997  
 4998

TVM's (de/re)-composition for MatMul in CUDA on NVIDIA Ampere GPU																																					
5000	5001	5002	5003	5004	5005	5006	5007	5008	5009	5010	5011	5012	5013	5014	5015	5016	5017	5018	Part.				De-Comp. Phase				Scalar Phase						Re-Comp. Phase				
																			P0	D1	D2	D3	D4	S1	S2	S3	S4	S5	S6	R1	R2	R3	R4				
																					A	B	A,B		A	B	A,B	C	C			C	C				
4999																			2	1	BLK,y	DM	DM	[1,2]	1	BLK,y							1	BLK,y	DM	[1,2]	
5000																			50	2	BLK,x	DM	DM	[1,2]	2	BLK,x							2	BLK,x	DM	[1,2]	
5001																			1	3	BLK,z	DM	DM	[1,2]	3	BLK,z							3	BLK,z	DM	[1,2]	
5002																			2	14	DEV,x	SM	SM	[1,2]	14	DEV,x							14	DEV,x	RM	[1,2]	
5003																			5004	1	15	DEV,y	SM	SM	[1,2]	15	DEV,y							15	DEV,y	RM	[1,2]
5004																			5005	8	7	DEV,z	DM	DM	[1,2]	7	DEV,z							7	DEV,z	RM	[1,2]
5005																			5006	4	4	THR,y	DM	DM	[1,2]	4	THR,y							4	THR,y	DM	[1,2]
5006																			5007	20	5	THR,x	DM	DM	[1,2]	5	THR,x	RM	RM	[1,2]	RM	[1,2]		5	THR,x	DM	[1,2]
5007																			5008	1	6	THR,z	DM	DM	[1,2]	6	THR,z							6	THR,z	DM	[1,2]
5008																			5009	128	8	SHR,z	SM	SM	[1,2]	8	SHR,z							9	SHR,x	RM	[1,2]
5009																			5010	1	12	REG,x	SM	SM	[1,2]	12	REG,x							10	SHR,y	RM	[1,2]
5010																			5011	1	13	REG,y	SM	SM	[1,2]	13	REG,y							8	SHR,z	RM	[1,2]
5011																			5012	2	11	REG,z	SM	SM	[1,2]	11	REG,z							12	REG,x	RM	[1,2]
5012																			5013														13	REG,y	RM	[1,2]	
5013																			5014														11	REG,z	RM	[1,2]	
5014																			5015																		
5015																			5016																		
5016																			5017																		
5017																			5018																		

Fig. 32. TVM's (de/re)-composition for MatMul in CUDA on GPU expressed in our low-level representation

TVM's (de/re)-composition for MatMul in OpenCL on Intel Skylake CPU																																					
5021	5022	5023	5024	5025	5026	5027	5028	5029	5030	5031	5032	5033	5034	5035	5036	5037	5038	5039	5040	Part.				De-Comp. Phase				Scalar Phase						Re-Comp. Phase			
																				P0	D1	D2	D3	D4	S1	S2	S3	S4	S5	S6	R1	R2	R3	R4			
																						A	B	A,B		A	B	A,B	C	C			C	C			
5019																			1	1	WG,1	GM	GM	[1,2]	1	WG,1							1	WG,1	GM	[1,2]	
5020																			125	2	WG,0	GM	GM	[1,2]	2	WG,0							2	WG,0	GM	[1,2]	
5021																			1	3	WG,2	GM	GM	[1,2]	3	WG,2							3	WG,2	GM	[1,2]	
5022																			1	14	GLB,0	LM	LM	[1,2]	14	GLB,0							14	GLB,0	PM	[1,2]	
5023																			1	15	GLB,1	LM	LM	[1,2]	15	GLB,1							15	GLB,1	PM	[1,2]	
5024																			16	7	GLB,2	GM	GM	[1,2]	7	GLB,2							7	GLB,2	PM	[1,2]	
5025																			16	4	WI,1	GM	GM	[1,2]	4	WI,1							4	WI,1	GM	[1,2]	
5026																			5027	5	WI,0	GM	GM	[1,2]	5	WI,0	PM	PM	[1,2]	PM	[1,2]		5	WI,0	GM	[1,2]	
5027																			1	6	WI,2	GM	GM	[1,2]	6	WI,2							6	WI,2	GM	[1,2]	
5028																			1	9	LCL,0	LM	LM	[1,2]	9	LCL,0							9	LCL,0	PM	[1,2]	
5029																			1	10	LCL,1	LM	LM	[1,2]	10	LCL,1							10	LCL,1	PM	[1,2]	
5030																			128	8	LCL,2	LM	LM	[1,2]	8	LCL,2							8	LCL,2	PM	[1,2]	
5031																			128	12	PRV,0	LM	LM	[1,2]	12	PRV,0							12	PRV,0	PM	[1,2]	
5032																			5033	8	13	PRV,1	LM	LM	[1,2]	13	PRV,1							13	PRV,1	PM	[1,2]
5033																			5034	1	11	PRV,2	LM	LM	[1,2]	11	PRV,2							11	PRV,2	PM	[1,2]
5034																			5035																		
5035																			5036																		
5036																			5037																		
5037																			5038																		
5038																			5039																		
5039																			5040																		
5040																																					

Fig. 33. TVM's (de/re)-composition for MatMul in OpenCL on CPU expressed in our low-level representation

PPCG's (de/re)-composition for MatMul in CUDA on NVIDIA Ampere GPU																
Part.	De-Comp. Phase				Scalar Phase						Re-Comp. Phase					
	P0	D1	D2	D3	D4	S1	S2	S3	S4	S5	S6	R1	R2	R3	R4	
				A	B	A,B		A	B	A,B	C	C		C	C	
16	1	BLK,y	DM	DM	[1,2]	1	BLK,y					1	BLK,y	DM	[1,2]	
8	2	BLK,x	DM	DM	[1,2]	2	BLK,x					2	BLK,x	DM	[1,2]	
1	3	BLK,z	DM	DM	[1,2]	3	BLK,z					3	BLK,z	DM	[1,2]	
1	7	DEV,x	DM	DM	[1,2]	7	DEV,x					7	DEV,x	DM	[1,2]	
1	8	DEV,y	DM	DM	[1,2]	8	DEV,y					8	DEV,y	DM	[1,2]	
1	9	DEV,z	DM	DM	[1,2]	9	DEV,z					9	DEV,z	DM	[1,2]	
1	4	THR,y	DM	DM	[1,2]	4	THR,y					4	THR,y	RM	[1,2]	
125	5	THR,x	DM	DM	[1,2]	5	THR,x	RM	RM	[1,2]	RM	1,2]	5	THR,x	RM	[1,2]
1	6	THR,z	DM	DM	[1,2]	6	THR,z					6	THR,z	RM	[1,2]	
1	10	SHR,x	SM	DM	[1,2]	10	SHR,x					10	SHR,x	RM	[1,2]	
1	11	SHR,y	SM	DM	[1,2]	11	SHR,y					11	SHR,y	RM	[1,2]	
1181	12	SHR,z	SM	DM	[1,2]	12	SHR,z					12	SHR,z	RM	[1,2]	
1	13	REG,x	SM	DM	[1,2]	13	REG,x					13	REG,x	RM	[1,2]	
1	14	REG,y	SM	DM	[1,2]	14	REG,y					14	REG,y	RM	[1,2]	
1	15	REG,z	SM	DM	[1,2]	15	REG,z					15	REG,z	RM	[1,2]	

Fig. 34. PPCG's (de/re)-composition for MatMul in CUDA on GPU expressed in our low-level representation

Pluto's (de/re)-composition for MatMul in OpenMP on Intel Skylake CPU																
Part.	De-Comp. Phase				Scalar Phase						Re-Comp. Phase					
	P0	D1	D2	D3	D4	S1	S2	S3	S4	S5	S6	R1	R2	R3	R4	
				A	B	A,B		A	B	A,B	C	C		C	C	
8	1	COR,0	MM	MM	[1,2]	1	COR,0					1	COR,0	MM	[1,2]	
1	2	COR,1	MM	MM	[1,2]	2	COR,1					2	COR,1	MM	[1,2]	
1	3	COR,2	MM	MM	[1,2]	3	COR,2					3	COR,2	MM	[1,2]	
1	4	MM,0	MM	MM	[1,2]	4	MM,0					4	MM,0	MM	[1,2]	
1	5	MM,1	MM	MM	[1,2]	5	MM,1					5	MM,1	MM	[1,2]	
9	6	MM,2	MM	MM	[1,2]	6	MM,2					6	MM,2	MM	[1,2]	
2	7	L2,0	L2	L2	[1,2]	7	L2,0					7	L2,0	L2	[1,2]	
962	8	L2,1	L2	L2	[1,2]	8	L2,1	L1	L1	[1,2]	L1	1,2]	8	L2,1	L2	[1,2]
218	9	L2,2	L2	L2	[1,2]	9	L2,2					9	L2,2	L2	[1,2]	
1	10	L1,0	L1	L1	[1,2]	10	L1,0					10	L1,0	L1	[1,2]	
1	11	L1,1	L1	L1	[1,2]	11	L1,1					11	L1,1	L1	[1,2]	
1	12	L1,2	L1	L1	[1,2]	12	L1,2					12	L1,2	L1	[1,2]	

Fig. 35. Pluto's (de/re)-composition for MatMul in OpenMP on CPU expressed in our low-level representation

### C.3 Constraints of Programming Models

Constraints of programming models can be expressed in our formalism; we demonstrate this using the example models CUDA and OpenCL. For this, we add to the general, model-unspecific constraints (described in Section 4) the new, model-specific constraints listed in Tables 2 or 3 for

5097 CUDA or the constraints in Table 4 for OpenCL, respectively. For brevity, we use in the following:

$$(l_{\text{MDH}}, d_{\text{MDH}}) := \leftrightarrow_{\bullet\text{-ass}}^{-1} (l_{\text{ASM}}, d_{\text{ASM}}), \bullet \in \{\downarrow, f, \uparrow\}$$

5101

5102 No. Constraint

5104	5105 0	$\prod_{d \in [1, D]_{\mathbb{N}}} \#PRT(CC, d) \leq 1024$	(Number of CCs limited)
5106	5107 R3	$\#PRT(BLK, d) > 1 \wedge \otimes_{\downarrow\text{-MDH}} \neq \oplus_{\uparrow\text{-MDH}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\text{BLK}, d) \in \{\text{DM}\}$	(SMXs combine in DM)
5109	5110	$\#PRT(CC, d) > 1 \wedge \otimes_{\downarrow\text{-MDH}} \neq \oplus_{\uparrow\text{-MDH}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\text{CC}, d) \in \{\text{DM, SM}\}$	(CCs combine in DM/SM)

5113 Table 2. CUDA model constraints on tuning parameters

5114

5115

5116

5117 No. Constraint

5118	5119 0	$\prod_{d \in [1, D]_{\mathbb{N}}} \#PRT(CC, d) \leq 1024$	(Number of CCs limited)
5121	5122 R3	$\#PRT(BLK, d) > 1 \wedge \otimes_{\downarrow\text{-MDH}} \neq \oplus_{\uparrow\text{-MDH}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\text{BLK}, d) \in \{\text{DM}\}$	(SMXs combine in DM)
5125	5126	$\#PRT(WRP, d) > 1 \wedge \otimes_{\downarrow\text{-MDH}} \neq \oplus_{\uparrow\text{-MDH}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\text{WRP}, d) \in \{\text{DM, SM}\}$	(WRPs combine in DM/SM)

5128 Table 3. CUDA+WRP model constraints on tuning parameters

5129

5130

5131

5132 No. Constraint

5133	5134 0	$\prod_{d \in [1, D]_{\mathbb{N}}} \#PRT(WI, d) \leq C_{\text{DEV}}$	(Number of PEs limited)
5136	5137 R3	$\#PRT(WG, d) > 1 \wedge \otimes_{\downarrow\text{-MDH}} \neq \oplus_{\uparrow\text{-MDH}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\text{WG}, d) \in \{\text{GM}\}$	(CUs combine in GM)
5139	5140	$\#PRT(WI, d) > 1 \wedge \otimes_{\downarrow\text{-MDH}} \neq \oplus_{\uparrow\text{-MDH}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\text{WI}, d) \in \{\text{GM, LM}\}$	(PEs combine in GM/LM)

5143 Table 4. OpenCL model constraints on tuning parameters

5144

5145

In Tables 2 and 3 for CUDA, the constraint No. 1 (which constrains tuning parameter No. 0 in Table 1) limits the number of cuda cores (CC) to 1024, according to the CUDA specification [NVIDIA 2022g]. The constraints on tuning parameters D3 and R3 specify that the results of SMX can be combined in device memory (DM) only, and that of CCs/WRPs in only device memory (DM) or shared memory (SM). Note that in the case of Table 3, CCs are not constrained in parameter D3 and R3, as CCs within a WRP have access to all CUDA memory regions: DM, SM, as well as RM (via warp shuffles [NVIDIA 2018]).

In Table 4 for OpenCL, the constraints are similar to the CUDA’s constraints in Tables 2: they limit the number of PEs to  $C_{DEV}$  (which is a device-specific constant in OpenCL), and they specify the valid memory regions for combining the results of cores, according to the OpenCL specification [Khronos 2022b].

Note that the tables present some important example constraints only and are not complete: for example, CUDA and OpenCL devices are also constrained regarding their memory sizes (shared/private memory), which is not considered in the tables.

## C.4 Inverse Concatenation

**Definition 39** (Inverse Concatenation). The inverse of operator *concatenation* (Example 13) is function  $+_{-1}$  which is of type

$\#_{-1}^{< T \in \text{TYPE} \mid D \in \mathbb{N} \mid d \in [1, D]_{\mathbb{N}} \mid (I_1, \dots, I_{d-1}, I_d, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} } :$

$$T[\underbrace{I_1, \dots, id(P \cup Q)}_{\uparrow d}, \dots, I_D] \rightarrow T[\underbrace{I_1, \dots, id(P)}_{\uparrow d}, \dots, I_D] \times T[\underbrace{I_1, \dots, id(Q)}_{\uparrow d}, \dots, I_D]$$

where  $id : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$  is the identity function on MDA index sets. The function is computed as:

$$\#^{-1 < T |D| d |(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D), (P, Q) >} (\mathfrak{a}) \equiv (\mathfrak{a}_1, \mathfrak{a}_2)$$

for

$$q_1[i_1, \dots, i_d, \dots, i_D] \equiv q[i_1, \dots, i_d, \dots, i_D], \quad i_d \in P$$

and

$$\mathfrak{g}_2[i_1, \dots, i_d, \dots, i_D] \equiv \mathfrak{g}[i_1, \dots, i_d, \dots, i_D], \quad i_d \in Q$$

i.e.,  $q_1$  and  $q_2$  behave exactly as MDA  $q$  on their restricted index sets  $P$  or  $Q$ , respectively.

We often write for  $(\alpha_1, \alpha_2) := +^{-1}^{<\dots>}(\alpha)$  (meta-parameters omitted via ellipsis) also  $\alpha = \alpha_1 +^{<\dots>} \alpha_2$ . Our notation is justified by the fact that the inverse of MDA  $\alpha$  is uniquely determined as the two MDAs  $\alpha_1$  and  $\alpha_2$  which are equal to  $\alpha$  when concatenating them (follows immediately from the definition of inverse functions).

## C 5 Example 15 in Verbose Math Notation

Figures 36-38 show our low-level representation from Example 15 in verbose math notation. The symbols  $\blacksquare, \dots, \blacksquare_{\ell}$  used in the figures are a textual abbreviation for:

$\blacksquare_{\perp}$	$\coloneqq$	$*, *$	$ $	$*, *$	$ $	$*, *$
$\blacksquare_1^1$	$\coloneqq$	$*, *$	$ $	$*, *$	$ $	$*, *$
$\blacksquare_2^1$	$\coloneqq$	$p_1^1, *$	$ $	$*, *$	$ $	$*, *$
$\blacksquare_1^2$	$\coloneqq$	$p_1^1, p_2^1$	$ $	$*, *$	$ $	$*, *$
$\blacksquare_2^2$	$\coloneqq$	$p_1^1, p_2^1$	$ $	$p_2^2, *$	$ $	$*, *$

5195	$\blacksquare_1^3$	$::=$	$p_1^1, p_2^1$	$ $	$p_1^2, p_2^2$	$ $	$*, *$
5196	$\blacksquare_2^3$	$::=$	$p_1^1, p_2^1$	$ $	$p_1^2, p_2^2$	$ $	$p_1^3, *$
5197	$\blacksquare_f$	$::=$	$p_1^1, p_2^1$	$ $	$p_1^2, p_2^2$	$ $	$p_1^3, p_2^3$

5199 where symbol  $*$  indicates generalization in meta-parameters (Definition 22).

5200 In Example 15, the arrow annotation of combine operators is formally an abbreviation. For  
 5201 example, operator  $\text{++}_2^{(\text{COR}, y)}$  in Figure 15 is annotated with  $\rightarrow \text{M: HM}[1, 2], v: \text{HM}[1]$  which  
 5202 abbreviates

$$5204 \quad \dots \quad \downarrow a_2^{2 < p_1^1, p_2^1 | p_1^2, p_2^2 := * | p_1^3 := *, p_2^3 := * >} =: \text{++}_2^{(\text{COR}, y)} \dots \\ 5205 \quad \quad \quad p_2^2 \in [0, 16]_{\mathbb{N}_0}$$

5206 Here,  $\downarrow a_2^2$  represents the low-level MDA (Definition 36) that is already partitioned for layer 1 in  
 5207 dimensions 1 and 2, and for layer 2 in dimension 1 (because in Figure 15, operators  $\text{++}_1^{(\text{HM}, x)}, \text{++}_2^{(\text{HM}, y)},$   
 5208  $\text{++}_1^{(\text{COR}, x)}$  appear before operator  $\text{++}_2^{(\text{COR}, y)}$ ), but not yet for layer 2 in dimension 2 as well as for  
 5209 layer 3 in both dimensions (indicated by symbol  $*$  which is described formally in Definition 22). In  
 5210 our generated code (discussed in Section F), we store low-level MDAs, like  $\downarrow a_2^2$ , using their domain-  
 5211 specific data representation, as the domain-specific representation is usually more efficient: in the  
 5212 case of MatVec, we physically store matrix  $M$  and vector  $v$  for the input MDA, and vector  $w$  for  
 5213 the output MDA. For example, low-level MDA  
 5214

$$5215 \quad \downarrow a_2^{2 < p_1^1, p_2^1 | p_1^2, p_2^2 := * | p_1^3 := *, p_2^3 := * >}$$

5216 can be transformed via view functions (Definitions 31 and 33) to *low-level BUFs* (Definition 37)

$$5218 \quad M_2^{< \text{HM} | id > < p_1^1, p_2^1 | p_1^2, p_2^2 := * | p_1^3 := *, p_2^3 := * >}$$

5219 and

$$5221 \quad v_2^{2 < \text{HM} | id > < p_1^1, p_2^1 | p_1^2, p_2^2 := * | p_1^3 := *, p_2^3 := * >}$$

5222 and back (Lemma 2). Similarly as for data structures in low-level programming models (e.g., C  
 5223 arrays as in OpenMP, CUDA, and OpenCL), low-level BUFs are defined to have an explicit notion  
 5224 of memory regions and memory layouts.

5225 In Figure 36, we de-compose the input MDA  $\downarrow a$ , step by step, for the MDH levels  $(1, 1), \dots, (3, 2)$ :

$$5228 \quad \downarrow a =: \downarrow a_{\perp}^{< \blacksquare_{\perp} >} \rightarrow \downarrow a_1^{1 < \blacksquare_1^1 >} \rightarrow \downarrow a_2^{1 < \blacksquare_2^1 >} \rightarrow \downarrow a_1^{2 < \blacksquare_1^2 >} \rightarrow \downarrow a_2^{2 < \blacksquare_2^2 >} \rightarrow \downarrow a_1^{3 < \blacksquare_1^3 >} \rightarrow \downarrow a_2^{3 < \blacksquare_2^3 >} \rightarrow \downarrow a_f^{< \blacksquare_f >}$$

5229 The input MDAs  $(\downarrow a_d^l)_{l \in [1, 3]_{\mathbb{N}}, d \in [1, 2]_{\mathbb{N}}}$  and  $\downarrow a_f$  are all low-level MDA representations (Definition 36).  
 5230 We use as partitioning schema  $P$  (Definition 36)

$$5232 \quad P := ( (P_1^1, P_2^1), (P_1^2, P_2^2), (P_1^3, P_2^3) ) = ( (2, 4), (8, 16), (32, 64) )$$

5233 and we use the index sets  $I_d$  from Definition 40 (which define a uniform index set partitioning) as  
 5234 follows:

$$5236 \quad ( \xrightarrow[d]{\text{MDA}} (I_d^{< p_1^1, p_2^1 | p_1^2, p_2^2 | p_1^3, p_2^3 >} )^{< (p_1^1, p_2^1) \in P_1^1 \times P_2^1 | (p_1^2, p_2^2) \in P_1^2 \times P_2^2 | (p_1^3, p_2^3) \in P_1^3 \times P_2^3 >} )_{d \in [1, D]_{\mathbb{N}}}$$

5238 Here,  $\xrightarrow[d]{\text{MDA}}$  denotes the index set function of combine operator concatenation (Example 13), which  
 5239 is the identity function and explicitly stated for the sake of completeness only. Note that in Figure 36,  
 5240 we access low-level MDAs  $\downarrow a_d^l$  as generalized in some partition sizes, via  $*$  (Definition 22), according  
 5241 to the definitions of the  $\blacksquare_d^l$ .  
 5242

5244 Each MDA  $\downarrow \mathbf{a}$  can be transformed to its domain-specific data representation matrix  $\downarrow M$  and  
 5245 vector  $\downarrow v$  and vice versa, using the view functions, as discussed above.

5246 Figure 37 shows our scalar phase, which is formally trivial.

5247 In Figure 38, we re-compose the computed data  $\uparrow \mathbf{a}_f^{<\blacksquare_f>}$ , step by step, to the final result  $\uparrow \mathbf{a}$ :

$$5248 \quad \uparrow \mathbf{a}_f^{<\blacksquare_f>} \rightarrow \uparrow \mathbf{a}_1^{<\blacksquare_1^1>} \rightarrow \uparrow \mathbf{a}_2^{<\blacksquare_2^1>} \rightarrow \uparrow \mathbf{a}_2^{<\blacksquare_2^2>} \rightarrow \uparrow \mathbf{a}_1^{<\blacksquare_1^3>} \rightarrow \uparrow \mathbf{a}_2^{<\blacksquare_2^3>} \rightarrow \uparrow \mathbf{a}_\perp^{<\blacksquare_\perp>} =: \uparrow \mathbf{a}$$

5250 Analogously to the de-composition phase, each output MDA  $(\uparrow \mathbf{a}_d^l)_{l \in [1,3]_{\mathbb{N}}, d \in [1,2]_{\mathbb{N}}}$  and  $\uparrow \mathbf{a}_f$  is a  
 5251 low-level MDA representation, for  $P$  as defined above and index sets

$$5252 \quad \left( \underset{\otimes_d}{\overset{d}{\Rightarrow}}_{MDA} (I_d^{<p_1^1, p_2^1 | p_1^2, p_2^2 | p_1^3, p_2^3>} )^{<(p_1^1, p_2^1) \in P_1^1 \times P_2^1 | (p_1^2, p_2^2) \in P_1^2 \times P_2^2 | (p_1^3, p_2^3) \in P_1^3 \times P_2^3>} \right)_{d \in [1, D]_{\mathbb{N}}}$$

5255 where  $\underset{\otimes_d}{\overset{d}{\Rightarrow}}_{MDA}$  are the index set functions of the combine operators (Definition 25) used in the re-  
 5256 composition phase. The same as in the de-composition phase, we access the output low-level MDAs  
 5257 as generalized in some partition sizes, according to our definitions of the  $\blacksquare_d^l$ , and we identify each  
 5258 MDA with its domain-specific data representation (the output vector  $w$ ).  
 5259

5260

5261

5262

5263

5264

5265

5266

5267

5268

5269

5270

5271

5272

5273

5274

5275

5276

5277

5278

5279

5280

5281

5282

5283

5284

5285

5286

5287

5288

5289

5290

5291

5292

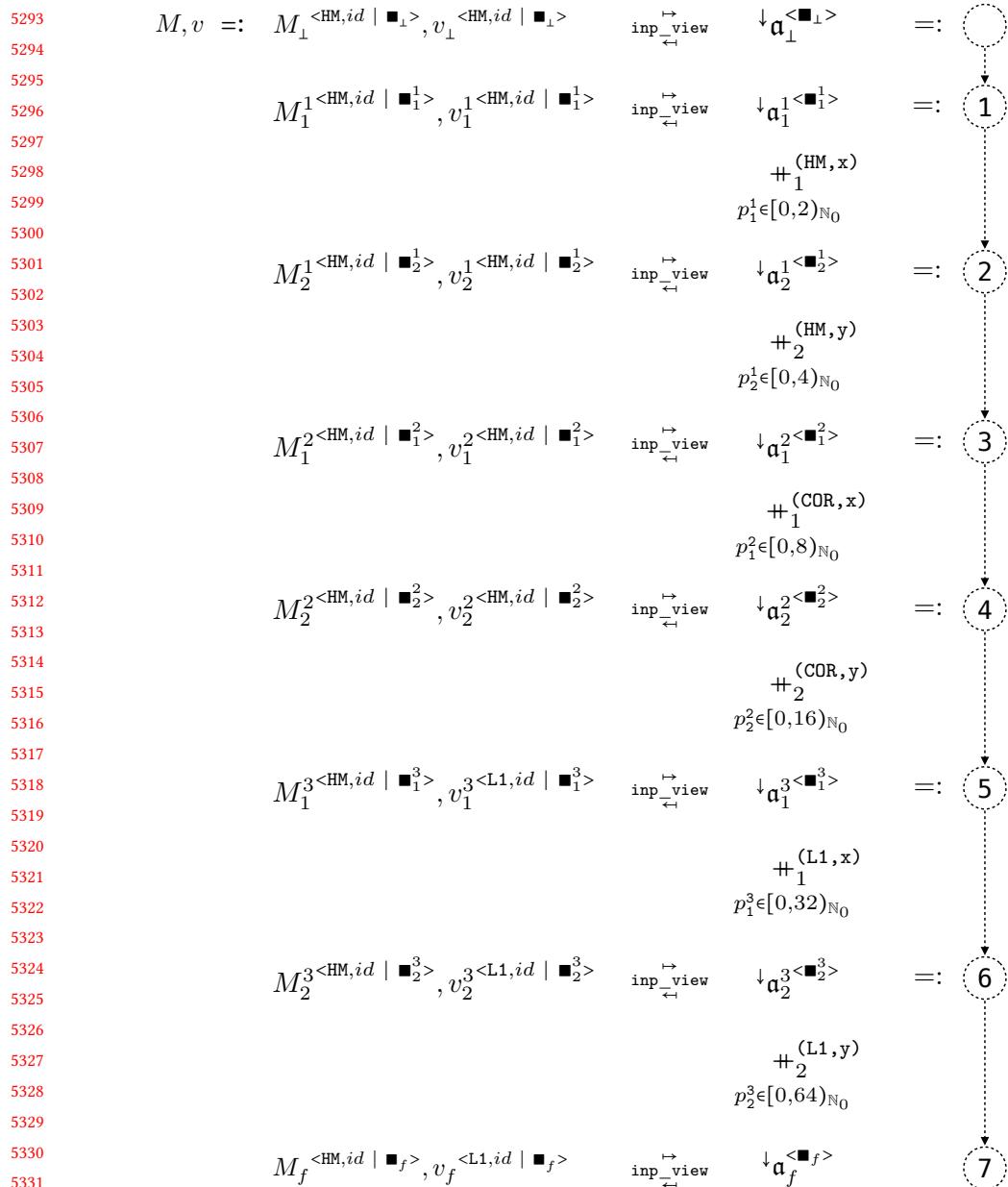


Fig. 36. De-composition phase of Example 15 in verbose math notation.

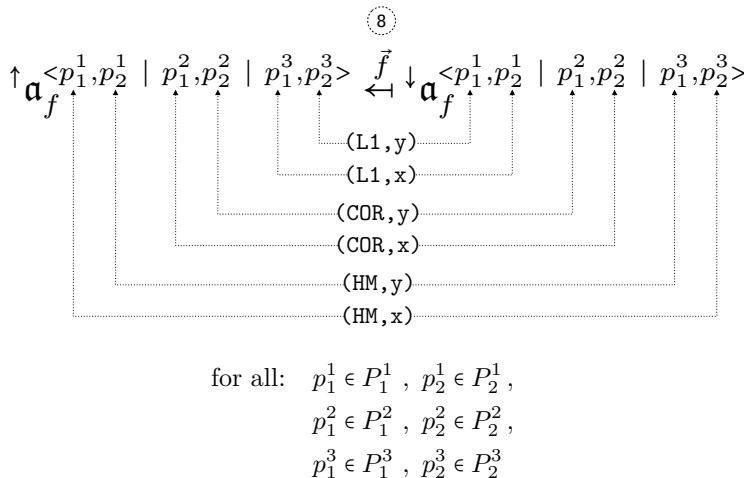


Fig. 37. Scalar phase of Example 15 in verbose math notation.

5391	$w := w_{\perp}^{<\text{HM}, id   \blacksquare_{\perp}>}$	$\text{out}_{\perp \rightarrow}^{\leftrightarrow \text{view}}$	$\uparrow \mathfrak{a}_{\perp}^{<\blacksquare_{\perp}>}$	$:=$	15
5392					
5393					
5394	$w_1^1^{<\text{HM}, id   \blacksquare_1^1>}$	$\text{out}_{\perp \rightarrow}^{\leftrightarrow \text{view}}$	$\uparrow \mathfrak{a}_1^1^{<\blacksquare_1^1>}$	$:=$	15
5395					
5396					
5397			$\oplus_1^{(\text{HM}, x)}$		
5398			$p_1^1 \in [0, 2)_{\mathbb{N}_0}$		
5399					
5400	$w_2^1^{<\text{HM}, id   \blacksquare_2^1>}$	$\text{out}_{\perp \rightarrow}^{\leftrightarrow \text{view}}$	$\uparrow \mathfrak{a}_2^1^{<\blacksquare_2^1>}$	$:=$	14
5401					
5402			$\oplus_2^{(\text{HM}, y)}$		
5403			$p_2^1 \in [0, 4)_{\mathbb{N}_0}$		
5404					
5405	$w_1^2^{<\text{HM}, id   \blacksquare_1^2>}$	$\text{out}_{\perp \rightarrow}^{\leftrightarrow \text{view}}$	$\uparrow \mathfrak{a}_1^2^{<\blacksquare_1^2>}$	$:=$	13
5406					
5407			$\oplus_1^{(\text{COR}, x)}$		
5408			$p_1^2 \in [0, 8)_{\mathbb{N}_0}$		
5409					
5410					
5411	$w_2^2^{<\text{HM}, id   \blacksquare_2^2>}$	$\text{out}_{\perp \rightarrow}^{\leftrightarrow \text{view}}$	$\uparrow \mathfrak{a}_2^2^{<\blacksquare_2^2>}$	$:=$	12
5412					
5413			$\oplus_2^{(\text{COR}, y)}$		
5414			$p_2^2 \in [0, 16)_{\mathbb{N}_0}$		
5415					
5416					
5417	$w_1^3^{<\text{L1}, id   \blacksquare_1^3>}$	$\text{out}_{\perp \rightarrow}^{\leftrightarrow \text{view}}$	$\uparrow \mathfrak{a}_1^3^{<\blacksquare_1^3>}$	$:=$	11
5418					
5419			$\oplus_1^{(\text{L1}, x)}$		
5420			$p_1^3 \in [0, 32)_{\mathbb{N}_0}$		
5421					
5422					
5423	$w_2^3^{<\text{L1}, id   \blacksquare_2^3>}$	$\text{out}_{\perp \rightarrow}^{\leftrightarrow \text{view}}$	$\uparrow \mathfrak{a}_2^3^{<\blacksquare_2^3>}$	$:=$	10
5424					
5425			$\oplus_2^{(\text{L1}, y)}$		
5426			$p_2^3 \in [0, 64)_{\mathbb{N}_0}$		
5427					
5428					
5429	$w_f^{<\text{L1}, id   \blacksquare_f>}$	$\text{out}_{\perp \rightarrow}^{\leftrightarrow \text{view}}$	$\uparrow \mathfrak{a}_f^{<\blacksquare_f>}$	$:=$	9
5430					

Fig. 38. Re-composition phase of Example 15 in verbose math notation.

5440 **C.6 Multi-Dimensional ASM Arrangements**

5441 We demonstrate how we arrange memory regions and cores of ASM-represented systems (Section  
 5442 3.2) in multiple dimensions using the example of CUDA.

5443     *Cores (COR)*: In CUDA, SMX cores are programmed via so-called *CUDA Blocks*, and CUDA's  
 5444 CC cores are programmed via *CUDA Threads*. CUDA has native support for arranging its blocks  
 5445 and threads in up to three dimensions which are called x, y, and z in CUDA [NVIDIA 2022f].  
 5446 Consequently, even though the original CUDA specification [NVIDIA 2022g] introduces SMX and  
 5447 CC without having an order, the CUDA programmer benefits from imagining SMX and CC as three-  
 5448 dimensionally arranged.

5449     Additional dimensions can be explicitly programmed in CUDA. For example, to add a fourth  
 5450 dimension to CUDA, we can embed the additional dimension in the CUDA's z dimension, thereby  
 5451 splitting CUDA dimension z in the explicitly programmed dimensions z\_1 (third dimension) and  
 5452 z\_2 (fourth dimension), as follows:

5453                  $z_1 := z \% Z_1$  and  $z_2 := z / Z_1$

5454     Here,  $Z_1$  represents the number of threads in the additional dimension, and symbol % the modulo  
 5455 operator.

5456     *Memory (MEM)*: In CUDA, memory is managed via *C arrays* which may be multi-dimensional: to  
 5457 arrange  $(\text{DIM}_1 \times \dots \times \text{DIM}_D)$ -many memory regions, each of size N, we use a CUDA array of the  
 5458 following type (pseudocode):

5459                 array[ DIM\_1 ]...[ DIM\_D ][ N ]

5460     Note that CUDA implicitly arranges its *shared* and *private* memory allocations in multiple  
 5461 dimensions, depending on the number of blocks and threads: a shared memory array of type  
 5462 `shared_array[ DIM_1 ]...[ DIM_D ][ N ]` is internally managed in CUDA as `shared_array[`  
 5463 `blockIdx.x ][ blockIdx.y ][ blockIdx.z ][ DIM_1 ]...[ DIM_D ][ N ]`, i.e., each CUDA  
 5464 block has its own shared memory region. Analogously, a private memory array `private_array[`  
 5465 `DIM_1 ]...[ DIM_D ][ N ]` is managed in CUDA as `private_array[ blockIdx.x ][ blockIdx.y ][`  
 5466 `blockIdx.z ][ threadIdx.x ][ threadIdx.y ][ threadIdx.z ][ DIM_1 ]...[ DIM_D ][ N ]`, correspondingly.  
 5467 Our arrangement methodology continues the CUDA's approach by explicitly  
 5468 programming the additional arrangement dimensions  $\text{DIM}_1, \dots, \text{DIM}_D$ .

5469     Figure 39 illustrates our multi-dimensional core and memory arrangement using the example of  
 5470 CUDA for  $D = 2$  (two-dimensional arrangement).

5471 **C.7 ASM Levels**

5472     ASM levels are pairs  $(l_{\text{ASM}}, d_{\text{ASM}})$  consisting of an ASM layer  $l_{\text{ASM}} \in \mathbb{N}$  and ASM dimension  $d_{\text{ASM}} \in \mathbb{N}$ .

5473     Figure 40 illustrates ASM levels using the example of CUDA's thread hierarchy. The figure shows  
 5474 that thread hierarchies can be considered as a tree in which each level is uniquely determined by a  
 5475 particular combination of a layer (block or thread in the case of CUDA) and dimension (x, y, or  
 5476 z). In the figure, we use lvl as an abbreviation for *level*, l for *layer*, and d for *dimension*.

5477     For ASM layers and dimensions, we usually use their domain-specific identifiers, e.g., BLK/CC  
 5478 and x/y/z as aliases for numerical values of layers and dimensions.

5479 **C.8 MDH Levels**

5480     MDH levels are pairs  $(l_{\text{MDH}}, d_{\text{MDH}})$  consisting of an MDH layer  $l_{\text{MDH}} \in \mathbb{N}$  and MDH dimension  $d_{\text{MDH}} \in \mathbb{N}$ .

5481

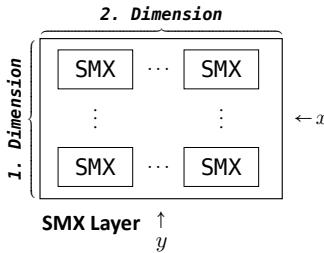
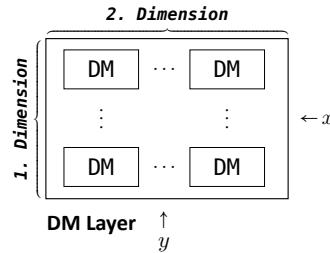
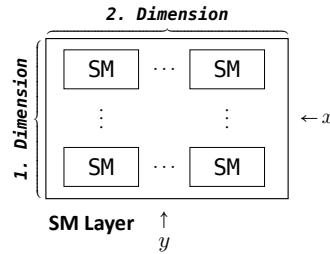
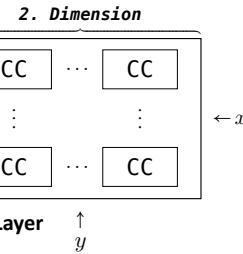
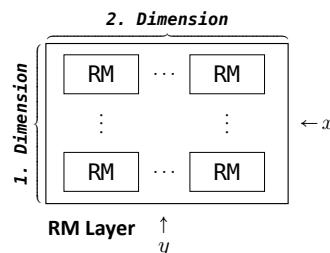
5489  
5490  
5491  
5492  
5493  
5494  
5495  
5496  
5497  
5498**COR Hierarchy****MEM Hierarchy**5499  
5500  
5501  
5502  
5503  
5504  
5505  
5506  
5507  
5508  
5509  
5510  
5511  
5512  
5513  
5514  
5515  
55165517  
5518Fig. 39. Multi-dimensional ASM arrangement illustrated using CUDA for the case  $D = 2$  (two dimensions)5519  
5520  
5521  
5522

Figure 41 illustrates MDH levels using as example the de-composition phase in Figure 15. The levels ( $l_{MDH}, d_{MDH}$ ) can be derived from the super- and subscripts of combine operators' variables  $p_{d_{MDH}}^{l_{MDH}}$ .

5523  
5524  
5525  
5526  
5527  
5528  
5529  
5530  
5531  
5532  
5533  
5534  
5535  
5536  
5537

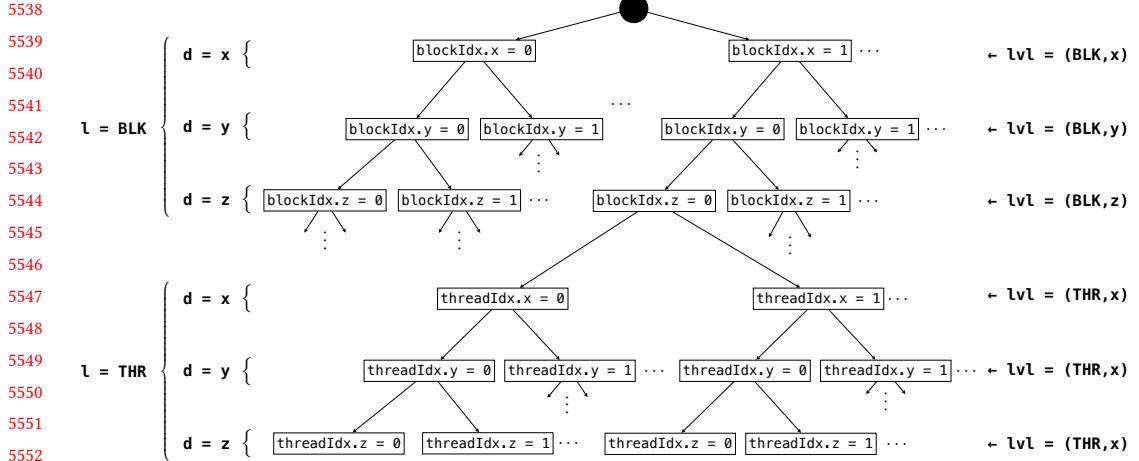


Fig. 40. ASM levels illustrated using CUDA's thread hierarchy

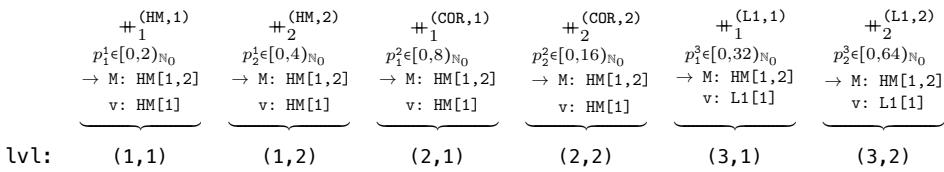


Fig. 41. MDH levels illustrated using as example the de-composition phase in Figure 15

## C.9 MDA Partitioning

We demonstrate how we partition MDAs into equally sized parts (a.k.a. *uniform partitioning*).

**Definition 40** (MDA Partitioning). Let  $\alpha \in T[I_1, \dots, I_D]$  be an arbitrary MDA that has scalar type  $T \in \text{TYPE}$ , dimensionality  $D \in \mathbb{N}$ , index sets  $I = (I_1, \dots, I_D) \in \text{MDA-IDX-SETS}^D$ , and size  $N = \{|I_1|, \dots, |I_D|\} \in \mathbb{N}^D$ . We consider  $I_d = \{i_1^d, \dots, i_{N_d}^d\}$ ,  $d \in [1, D]_{\mathbb{N}}$ , such that  $i_1^d < \dots < i_{N_d}^d$  represents a sorted enumeration of the elements in  $I_d$ . Let further  $P = ((P_1^1, \dots, P_D^1), \dots, (P_1^L, \dots, P_D^L))$  be an arbitrary tuple of  $L$ -many  $D$ -tuples of positive natural numbers such that  $\prod_{l \in [1, L]_{\mathbb{N}}} P_d^l$  divides  $N_d$  (the number of indices of MDA  $\alpha$  in dimension  $d$ ), for each  $d \in \{1, \dots, D\}$ .

The  $L$ -layered,  $D$ -dimensional,  $P$ -partitioning of MDA  $\alpha$  is the  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned low-level MDA  $\alpha_{\text{prt}}$  (Definition 36) that has scalar type  $T$  and index sets

$$I_d^{< p_d^1, \dots, p_d^L >} :=$$

$$\{i_j \in I_d \mid j = OS + j', \text{ for } OS := \sum_{l \in [1, L]_{\mathbb{N}}} p_d^l * \frac{N_d}{\prod_{l' \in [1, l]_{\mathbb{N}}} P_d^{l'}} \text{ and } j' \in PS := \frac{N_d}{\prod_{l' \in [1, L]_{\mathbb{N}}} P_d^{l'}}\}$$

i.e., set  $I_d^{< p_d^1, \dots, p_d^L >}$  denotes for each choice of parameters  $p_d^1, \dots, p_d^L$  a part of the uniform partitioning of the ordered index set  $I_d$  ( $OS$  in the formula above represents the OffSet to the part, and  $PS$

5587 the Part's Size). The partitioned MDA  $\alpha_{\text{prt}}$  is defined as:

$$5588 \quad \alpha =: \underbrace{\begin{array}{c} +_1 \dots +_D \\ p_1^1 \in P_1^1 \quad p_D^1 \in P_D^1 \end{array}}_{\text{Layer 1}} \dots \underbrace{\begin{array}{c} +_1 \dots +_D \\ p_1^L \in P_1^L \quad p_D^L \in P_D^L \end{array}}_{\text{Layer } L} \alpha_{\text{prt}}^{< p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L >}$$

5593 i.e., the parts  $\alpha_{\text{prt}}^{< p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L >}$  are defined such that concatenating them results in the  
5594 original MDA  $\alpha$ .

## 5595 C.10 TVM Schedule for MatMul

5597 Listing 6 shows TVM's Ansor-generated schedule program for MatMul on input matrices of sizes  
5598  $16 \times 2048$  and  $2048 \times 1000$  taken from ResNet-50's training phase, discussed in Section C.2. Code  
5599 formatting, like names of variables and comments, have been shortened and adapted in the listing  
5600 for brevity.

5601  
5602  
5603  
5604  
5605  
5606  
5607  
5608  
5609  
5610  
5611  
5612  
5613  
5614  
5615  
5616  
5617  
5618  
5619  
5620  
5621  
5622  
5623  
5624  
5625  
5626  
5627  
5628  
5629  
5630  
5631  
5632  
5633  
5634  
5635

```

5636 1 # exploiting fast memory resources for computed results
5637 2 matmul_local, = s.cache_write([matmul], "local")
5638 3 matmul_1, matmul_2, matmul_3 = tuple(matmul_local.op.axis) + tuple(matmul_local
5639 4 .op.reduce_axis)
5640 5 SHR_1, REG_1 = s[matmul_local].split(matmul_1, factor=1)
5641 6 THR_1, SHR_1 = s[matmul_local].split(SHR_1, factor=1)
5642 7 DEV_1, THR_1 = s[matmul_local].split(THR_1, factor=4)
5643 8 BLK_1, DEV_1 = s[matmul_local].split(DEV_1, factor=2)
5644 9 SHR_2, REG_2 = s[matmul_local].split(matmul_2, factor=1)
5645 10 THR_2, SHR_2 = s[matmul_local].split(SHR_2, factor=1)
5646 11 DEV_2, THR_2 = s[matmul_local].split(THR_2, factor=20)
5647 12 BLK_2, DEV_2 = s[matmul_local].split(DEV_2, factor=1)
5648 13 SHR_3, REG_3 = s[matmul_local].split(matmul_3, factor=2)
5649 14 DEV_3, SHR_3 = s[matmul_local].split(SHR_3, factor=128)
5650 15 s[matmul_local].reorder(BLK_1, BLK_2, DEV_1, DEV_2, THR_1, THR_2, DEV_3, SHR_3,
5651 16 SHR_1, SHR_2, REG_3, REG_1, REG_2)
5652 17 # low-level optimizations:
5653 18 s[matmul_local].pragma(BLK_1, "auto_unroll_max_step", 512)
5654 19 s[matmul_local].pragma(BLK_1, "unroll_explicit", True)
5655 20 # tiling
5656 21 matmul_1, matmul_2, matmul_3 = tuple(matmul.op.axis) + tuple(matmul.op.
5657 22 .reduce_axis)
5658 23 THR_1, SHR_REG_1 = s[matmul].split(matmul_1, factor=1)
5659 24 DEV_1, THR_1 = s[matmul].split(THR_1, factor=4)
5660 25 BLK_1, DEV_1 = s[matmul].split(DEV_1, factor=2)
5661 26 THR_2, SHR_REG_2 = s[matmul].split(matmul_2, factor=1)
5662 27 DEV_2, THR_2 = s[matmul].split(THR_2, factor=20)
5663 28 BLK_2, DEV_2 = s[matmul].split(DEV_2, factor=1)
5664 29 s[matmul].reorder(BLK_1, BLK_2, DEV_1, DEV_2, THR_1, THR_2, SHR_REG_1,
5665 30 SHR_REG_2)
5666 31 s[matmul_local].compute_at(s[matmul], THR_2)
5667 32 # block/thread assignments:
5668 33 BLK_fused = s[matmul].fuse(BLK_1, BLK_2)
5669 34 s[matmul].bind(BLK_fused, te.thread_axis("blockIdx.x"))
5670 35 DEV_fused = s[matmul].fuse(DEV_1, DEV_2)
5671 36 s[matmul].bind(DEV_fused, te.thread_axis("vthread"))
5672 37 THR_fused = s[matmul].fuse(THR_1, THR_2)
5673 38 s[matmul].bind(THR_fused, te.thread_axis("threadIdx.x"))
5674 39 # exploiting fast memory resources for first input matrix:
5675 40 A_shared = s.cache_read(A, "shared", [matmul_local])
5676 41 A_shared_ax0, A_shared_ax1 = tuple(A_shared.op.axis)
5677 42 A_shared_ax0_ax1_fused = s[A_shared].fuse(A_shared_ax0, A_shared_ax1)
5678 43 A_shared_ax0_ax1_fused_o, A_shared_ax0_ax1_fused_i = s[A_shared].split(
5679 44 A_shared_ax0_ax1_fused, factor=1)
5680 45 s[A_shared].vectorize(A_shared_ax0_ax1_fused_i)
5681 46 A_shared_ax0_ax1_fused_o_o, A_shared_ax0_ax1_fused_o_i = s[A_shared].split(
5682 47 A_shared_ax0_ax1_fused_o, factor=80)
5683 48 s[A_shared].bind(A_shared_ax0_ax1_fused_o_i, te.thread_axis("threadIdx.x"))
5684 49 s[A_shared].compute_at(s[matmul_local], DEV_3)
5685 50 # exploiting fast memory resources for second input matrix:
5686 51 # ... (analogous to lines 40 - 47)

```

5683 Listing 6. TVM schedule for Matrix Multiplication on NVIDIA Ampere GPU (variable names shortened for  
5684 brevity)

Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

## **D ADDENDUM SECTION 4**

We have designed our formalism such that an expression in our high-level representation (such as in Figure 6) can be systematically lowered to an expression in our low-level representation (as in Figure 15). We confirm this by parameterizing the high-level expression in Figure 30 – step-by-step – in the tuning parameters listed in Table 1, in a formally sound manner, which results exactly in the low-level expression in Figure 31.

PARAMETER  $\theta$ . According to Definition 36, we use the  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned low-level representation  ${}^{\downarrow}\alpha_f$  of  $\alpha$  for

$$P := \left( \underbrace{(\#PRT(1,1), \dots, \#PRT(1,D))}_{\text{Dimension 1}}, \dots, \underbrace{(\#PRT(L,1), \dots, \#PRT(L,D))}_{\text{Dimension } D} \right) \underbrace{\qquad\qquad\qquad}_{\text{Layer 1}} \qquad\qquad \underbrace{\qquad\qquad\qquad}_{\text{Layer } L}$$

where  $\#PRT$  denotes the number of partitions (Parameter 0 in Table 1).

Applying the homomorphic property (Definition 27)  $L * D$  times, we get:

`md_hom( f , (⊗1, ..., ⊗D) )( a )` =

`md_hom( f , ( ⊕_1, ..., ⊕_D ) )( ↴ a_f^{p_1^1, ..., p_D^1} | ... | p_1^L, ..., p_D^L )`

where  $\downarrow a_f^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \text{ is defined as:}$

$$\downarrow \mathfrak{a}_f =: \underbrace{\begin{array}{ccccccccc} \mathfrak{a}_f^{+1} & \dots & \mathfrak{a}_f^{+D} & \dots & \mathfrak{a}_f^{+1} & \dots & \mathfrak{a}_f^{+D} & \downarrow \mathfrak{a}_f^{< p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L >} \\ \underbrace{p_1^1 \in P_1} & & \underbrace{p_D^1 \in P_D} & & \underbrace{p_1^L \in P_1} & & \underbrace{p_D^L \in P_D} & \end{array}}_{\begin{array}{cc} \text{Dimension 1} & \text{Dimension } D \end{array}} \quad \underbrace{\begin{array}{ccccccccc} \mathfrak{a}_f^{-1} & \dots & \mathfrak{a}_f^{-D} & \dots & \mathfrak{a}_f^{-1} & \dots & \mathfrak{a}_f^{-D} & \downarrow \mathfrak{a}_f^{< p_1^{-1}, \dots, p_D^{-1} | \dots | p_1^{-L}, \dots, p_D^{-L} >} \\ \underbrace{p_1^{-1} \in P_1} & & \underbrace{p_D^{-1} \in P_D} & & \underbrace{p_1^{-L} \in P_1} & & \underbrace{p_D^{-L} \in P_D} & \end{array}}_{\begin{array}{cc} \text{Dimension 1} & \text{Dimension } D \end{array}}$$

Since each part  $\downarrow a_f^{p_1^1, \dots, p_D^1} \mid \dots \mid p_1^L, \dots, p_D^L \rangle$  contains a single scalar value only (according to the algorithmic constraint of Parameter  $\theta$ , discussed in Section 4), it holds

$$\text{md\_hom}( f, (\oplus_1, \dots, \oplus_D) )(\downarrow_{\mathfrak{a}_f}^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} ) = \vec{f}(\downarrow_{\mathfrak{a}_f}^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} )$$

for  $\vec{f}$  defined as in Definition 27

PARAMETER D1. Parameter D1 reorders concatenation operators  $+_1, \dots, +_D$  (Example 13) which are of types  $+_d \in \text{CO}^{<id \mid T \mid D \mid d>}$  for arbitrary but fixed  $T \in \text{TYPE}$  and  $D \in \mathbb{N}$ . We prove our assumption w.l.o.g. for the case  $D = 2$ ; the general case  $D \in \mathbb{N}$  follows analogously.

Let  $++_{d_1} \in \text{CO}^{<id|T|D|d_1>}$  and  $++_{d_2} \in \text{CO}^{<id|T|D|d_2>}$  be two arbitrary concatenation operators that coincide in meta-parameters  $T$  and  $D$ , but may differ in their operating dimensions  $d_1$  and  $d_2$ . We have to show

$$(a_1 \#_{d_1} a_2) \#_{d_2} (a_3 \#_{d_1} a_4) = (a_1 \#_{d_2} a_3) \#_{d_1} (a_2 \#_{d_2} a_4)$$

which follows from the definition of the concatenation operator  $\#$  in Example 13.

5734 PARAMETERS D2, S2, R2. These parameters replaces combine operators (Definition 25) by low-  
 5735 level combine operators (Definition 38) which has no effect on semantics.  $\square$

5736 PARAMETERS D3, S3, S5, R3. These parameters set the memory tags of low-level BUFs (Definition  
 5737 37), which have no effect on semantics.  $\square$

5738 5739 PARAMETERS D4, S4, S6, R4. The parameters change the memory layout of low-level BUFs (Definition  
 5740 37), which does not affect extensional equality.  $\square$

5741 5742 PARAMETERS S1. This parameter sets the order in which function  $f$  is applied to parts, which is  
 5743 trivially sound for any order.  $\square$

5744 5745 PARAMETERS R1. Similarly to parameter D1, parameter R1 reorders combine operators  $\otimes_1, \dots, \otimes_D$ .  
 5746 In contrast to parameter D1, the combine operators reordered by parameter R1 are not restricted to  
 5747 be concatenation.

5748 We prove our assumption by exploiting the MDH property (Definition 26) together with the  
 5749 proof of parameter D1, as follows:

$$\begin{aligned}
 & (a_1 \otimes_{d_1} a_2) \otimes_{d_2} (a_3 \otimes_{d_1} a_4) \\
 = & \text{md\_hom}(\dots)(a_1 \#_{d_1} a_2) \#_{d_2} (a_3 \#_{d_1} a_4) \\
 = & \text{md\_hom}(\dots)(a_1 \#_{d_2} a_3) \#_{d_1} (a_2 \#_{d_2} a_4) \\
 = & (a_1 \otimes_{d_2} a_3) \otimes_{d_1} (a_2 \otimes_{d_2} a_4) \quad \checkmark
 \end{aligned}$$

$\square$

## 5756 E ADDENDUM SECTION 5

### 5757 E.1 Data Characteristics used in Deep Neural Networks

5758 Figure 42 shows the data characteristics used for the deep learning experiments in Figures 21 and 22  
 5759 of Section 5. We use real-world characteristics taken from the neural networks ResNet-50, VGG-16,  
 5760 and MobileNet. For each network, we consider computations MCC and MatMul (Table 14), because  
 5761 these are the networks' most time-intensive building blocks. Each computation is called in each  
 5762 network on different data characteristics – we use for each combination of network and computation  
 5763 the two most time-intensive characteristics. Note that the MobileNet network does not use MatMul  
 5764 in its implementation.

5765 The capsule variants MCC\_Capsule in Figures 21 and 22 of Section 5 have the same characteristics  
 5766 as those listed for MCCs in Figure 42 – the only difference is that MCC\_Capsule, in addition to the  
 5767 dimensions N, H, W, K, R, S, C, uses three additional dimensions MI, MJ, MK, each with a fixed size of 4.  
 5768 This is because MCC\_Capsule operates on  $4 \times 4$  matrices, rather than scalars as MCC does.

### 5771 E.2 Runtime and Accuracy of cuBLASEx

5772 Listing 7 shows the runtime of cuBLASEx for its different *algorithm* variants. For demonstration, we  
 5773 use the example of matrix multiplication MatMul on NVIDIA Volta GPU for square input matrices  
 5774 of sizes  $1024 \times 1024$ . For each algorithm variant, we list both: 1) the runtime achieved by cuBLASEx  
 5775 (in nanoseconds ns), as well as 2) the maximum absolute deviation ( $\delta_{\max}$  values) compared  
 5776 to a straightforward, sequential CPU computation. For example, the  $\delta_{\max}$  value of algorithm  
 5777 CUBLAS\_GEMM\_DEFAULT is  $3.14713e-05$ , i.e., at least one value  $c_{i,j}^{\text{GPU}}$  in the GPU-computed output  
 5778 matrix deviates by  $3.14713e-05$  from its corresponding, sequentially computed value  $c_{i,j}^{\text{seq}}$  such that  
 5779  $|c_{i,j}^{\text{GPU}}| = |c_{i,j}^{\text{seq}}| + 3.14713e-05$  (bar symbols  $|\dots|$  denote absolute value). All other GPU-computed  
 5780 values  $c_{i',j'}^{\text{GPU}}$  deviate from their sequentially computed CPU-variant by  $3.14713e-05$  or less.

5781 5782

5783	Network	Phase	N	H	W	K	R	S	C	Stride H	Stride W	Padding	P	Q	Image Format	Filter Format	Output Format	
5784	ResNet-50	Training	16	230	230	64	7	7	3		2	2	VALID	112	112	NHWC	KRSC	NPQK
		Inference	1	230	230	64	7	7	3		2	2	VALID	112	112	NHWC	KRSC	NPQK
5785	VGG-16	Training	16	224	224	64	3	3	3		1	1	VALID	224	224	NHWC	KRSC	NPQK
		Inference	1	224	224	64	3	3	3		1	1	VALID	224	224	NHWC	KRSC	NPQK
5786	MobileNet	Training	16	225	225	32	3	3	3		2	2	VALID	112	112	NHWC	KRSC	NPQK
		Inference	1	225	225	32	3	3	3		2	2	VALID	112	112	NHWC	KRSC	NPQK

(a) Data characteristics used for MCC experiments

5788	Network	Phase	M	N	K	Transposition
5789	ResNet-50	Training	16	1000	2048	NN
		Inference	1	1000	2048	NN
5790	VGG-16	Training	16	4096	25088	NN
		Inference	1	4096	25088	NN

(b) Data characteristics used for MatMul experiments

Fig. 42. Data characteristics used for experiments in Section 5

Note that cuBLASEx offers 42 algorithm variants, but not all of them are supported for all potential characteristics of the input and output data (size, memory layout, . . .). For our MatMul example, the list of unsupported variants includes: CUBLAS\_GEMM\_ALG01, CUBLAS\_GEMM\_ALG012, etc.

5801

5802

5803

5804

5805

5806

5807

5808

5809

5810

5811

5812

5813

5814

5815

5816

5817

5818

5819

5820

5821

5822

5823

5824

5825

5826

5827

5828

5829

5830

5831

```

5832
5833 CUBLAS_GEMM_DEFAULT: 188416ns (delta_max: 3.14713e-05)
5834 CUBLAS_GEMM_ALGO2: 190464ns (delta_max: 6.86646e-05)
5835 CUBLAS_GEMM_ALGO3: 186368ns (delta_max: 6.86646e-05)
5836 CUBLAS_GEMM_ALGO4: 185344ns (delta_max: 6.86646e-05)
5837 CUBLAS_GEMM_ALGO5: 181248ns (delta_max: 6.86646e-05)
5838 CUBLAS_GEMM_ALGO6: 181248ns (delta_max: 6.86646e-05)
5839 CUBLAS_GEMM_ALGO7: 178176ns (delta_max: 4.1008e-05)
5840 CUBLAS_GEMM_ALGO8: 189440ns (delta_max: 4.1008e-05)
5841 CUBLAS_GEMM_ALGO9: 171008ns (delta_max: 4.1008e-05)
5842 CUBLAS_GEMM_ALGO10: 188416ns (delta_max: 4.1008e-05)
5843 CUBLAS_GEMM_ALGO11: 191488ns (delta_max: 4.1008e-05)
5844 CUBLAS_GEMM_ALGO18: 185344ns (delta_max: 2.67029e-05)
5845 CUBLAS_GEMM_ALGO19: 172032ns (delta_max: 2.67029e-05)
5846 CUBLAS_GEMM_ALGO20: 192512ns (delta_max: 2.67029e-05)
5847 CUBLAS_GEMM_ALGO21: 201728ns (delta_max: 1.90735e-05)
5848 CUBLAS_GEMM_ALGO22: 177152ns (delta_max: 1.90735e-05)
5849 CUBLAS_GEMM_ALGO23: 194560ns (delta_max: 1.90735e-05)
5850 CUBLAS_GEMM_DEFAULT_TENSOR_OP: 184320ns (delta_max: 3.14713e-05)
5851 CUBLAS_GEMM_ALGO0_TENSOR_OP: 62464ns (delta_max: 0.0131454)
5852 CUBLAS_GEMM_ALGO1_TENSOR_OP: 52224ns (delta_max: 0.0131454)
5853 CUBLAS_GEMM_ALGO2_TENSOR_OP: 190464ns (delta_max: 3.14713e-05)
5854 CUBLAS_GEMM_ALGO3_TENSOR_OP: 189440ns (delta_max: 3.14713e-05)
5855 CUBLAS_GEMM_ALGO4_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5856 CUBLAS_GEMM_ALGO5_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5857 CUBLAS_GEMM_ALGO6_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5858 CUBLAS_GEMM_ALGO7_TENSOR_OP: 189440ns (delta_max: 3.14713e-05)
5859 CUBLAS_GEMM_ALGO8_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5860 CUBLAS_GEMM_ALGO9_TENSOR_OP: 189440ns (delta_max: 3.14713e-05)
5861 CUBLAS_GEMM_ALGO10_TENSOR_OP: 188416ns (delta_max: 3.14713e-05)
5862 CUBLAS_GEMM_ALGO11_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5863 CUBLAS_GEMM_ALGO12_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5864 CUBLAS_GEMM_ALGO13_TENSOR_OP: 188416ns (delta_max: 3.14713e-05)
5865 CUBLAS_GEMM_ALGO14_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5866 CUBLAS_GEMM_ALGO15_TENSOR_OP: 189440ns (delta_max: 3.14713e-05)

```

Listing 7. Runtime of cuBLASEx for its different *algorithm* variants on NVIDIA Volta GPU when computing MatMul on square  $1024 \times 1024$  input matrices

```

5867
5868
5869
5870
5871
5872
5873
5874
5875
5876
5877
5878
5879
5880

```

## 5881 F CODE GENERATION

5882 This section outlines how imperative-style pseudocode is generated from our low-level program  
 5883 representation in Section 3. Code optimizations that are below the abstraction level of our low-  
 5884 level representation (like loop fusion and loop unrolling) are considered in this work as *code-level*  
 5885 *optimizations* and briefly outlined in Section G. We aim to discuss and illustrate our code generation  
 5886 approach in detail in future work.

5887 In the following, we highlight tuning parameters gray in our pseudocode which are substituted  
 5888 by concrete, optimized values in our executable program code. Static parameters, such as scalar  
 5889 types and the number of input/output buffers, are denoted in math font and also substituted by  
 5890 concrete values in our executable code. We list meta-parameters in angle brackets  $\langle \dots \rangle$ , and  
 5891 other static function annotations in double angle brackets  $\langle\langle \dots \rangle\rangle$ , e.g.,  $\text{idx} \langle\langle \text{OUT} \rangle\rangle \langle\langle 1, 1 \rangle\rangle$  for  
 5892 denoting in our pseudocode index function  $\text{idx}_{1,1}^{\text{OUT}}$  used in Figure 30.  
 5893

## 5894 Overall Structure

5895 Listing 8 shows the overall structure of our generated code. We implement a particular expression  
 5896 in our low-level representation (Figure 31) as a compute kernel that is structured in the follow-  
 5897 ing phases: 0) preparation (Section F.0), 1) de-composition phase (Section F.1), 2) scalar phase  
 5898 (Section F.2), 3) re-composition phase (Section F.3).  
 5899

```
5900 1 kernel mdh(
5901 2    $T_1^{\text{IB}}$  trans_ll_IB  $\langle\langle \perp \rangle\rangle \langle\langle 1 \rangle\rangle \langle\langle \star, \dots, \star \rangle\rangle, \dots, T_1^{\text{IB}}$  trans_ll_IB  $\langle\langle \perp \rangle\rangle \langle\langle B^{\text{IB}} \rangle\rangle \langle\langle \star, \dots, \star \rangle\rangle,$ 
5902 3    $T_1^{\text{OB}}$  trans_ll_OB  $\langle\langle \perp \rangle\rangle \langle\langle 1 \rangle\rangle \langle\langle \star, \dots, \star \rangle\rangle, \dots, T_1^{\text{OB}}$  trans_ll_OB  $\langle\langle \perp \rangle\rangle \langle\langle B^{\text{OB}} \rangle\rangle \langle\langle \star, \dots, \star \rangle\rangle$ 
5903 4 )
5904 5
5905 6   // 0. preparation
5906 7   ...
5907 8   // 2. de-composition phase
5908 9   ...
5909 10  // 3. scalar phase
5910 11  ...
5911 12  // 4. re-composition phase
5912 13  ...
5913 14 }
```

5912 Listing 8. Overall structure of our generated code  
 5913

## 5914 F.0 Preparation

5915 Listing 9 shows the preparation phase. It prepares in five sub-phases the basic building blocks used  
 5916 in our low-level representation: 1) md\_hom (Section F.0.1), 2) inp\_view (Section F.0.2), 3) out\_view  
 5917 (Section F.0.3), 4) BUFs (Section F.0.4), 5) MDAs (Section F.0.5).  
 5918

```
5919 1 // 0. preparation
5920 2   // 0.1. md_hom
5921 3   ...
5922 4   // 0.2. inp_view
5923 5   ...
5924 6   // 0.3. out_view
5925 7   ...
5926 8   // 0.4. BUFs
5927 9   ...
5928 10  // 0.5. MDAs
5929 11  ...
```

## Listing 9. Preparation Phase

## F.0.1 md\_hom.

Listing 10 shows the user-defined scalar function and low-level combine operators (Definition 38) which are both provided by the user via higher-order function `md_hom` (Definition 27).

Listing 11 shows how we pre-implement for the user the two combine operators *concatenation* (Example 13) and *point-wise combination* (Example 14).

Listing 12 shows how we pre-implement the *inverse of concatenation* (Definition C.4), which we will use in the de-composition phase (via Definition 40).

```

5940 1 // 0.1. md_hom
5941 2
5942 3     // 0.1.1. scalar function
5943 4     f(  $T^{INP}$  inp ) ->  $T^{OUT}$  out
5944 5     {
5945 6         // ... (user defined)
5946 7     }
5947 8
5948 9     // 0.1.2. combine operators
5949 10     $\forall d \in [1, D]_{\mathbb{N}}$ :
5950 11         $co << d >> < I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in MDA\text{-}IDX\text{-}SETs, (P, Q) \in MDA\text{-}IDX\text{-}SETs \times MDA\text{-}IDX\text{-}SETs > ($ 
5951 12             $T^{OUT}[I_1, \dots, I_{d-1}, \xrightarrow{d}^{MDA}(P), I_{d+1}, \dots, I_D] \text{ lhs} ,$ 
5952 13             $T^{OUT}[I_1, \dots, I_{d-1}, \xrightarrow{d}^{MDA}(Q), I_{d+1}, \dots, I_D] \text{ rhs } ) \rightarrow T^{OUT}[I_1, \dots, I_{d-1}, \xrightarrow{d}^{MDA}(P \cup Q), I_{d+1}, \dots, I_D] \text{ res}$ 
5953 14        {
5954 15            // ... (user defined)
5955 16        }
```

## Listing 10. Scalar Function &amp; Combine Operators

```

5956 1 // 0.1.2. combine operators
5957 2
5958 3     // pre-implemented combine operators
5959 4
5960 5     // concatenation
5961 6      $\forall d \in \mathbb{N}$ :
5962 7         $cc << d >> < I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in MDA\text{-}IDX\text{-}SETs, (P, Q) \in MDA\text{-}IDX\text{-}SETs \times MDA\text{-}IDX\text{-}SETs > ($ 
5963 8             $T^{OUT}[I_1, \dots, I_{d-1}, id(P), I_{d+1}, \dots, I_D] \text{ lhs} ,$ 
5964 9             $T^{OUT}[I_1, \dots, I_{d-1}, id(Q), I_{d+1}, \dots, I_D] \text{ rhs } ) \rightarrow T^{OUT}[I_1, \dots, I_{d-1}, id(P \cup Q), I_{d+1}, \dots, I_D] \text{ res}$ 
5965 10    {
5966 11        int i_1  $\in I_1$ 
5967 12        ..
5968 13        int i_{d-1}  $\in I_{d-1}$ 
5969 14        int i_{d+1}  $\in I_{d+1}$ 
5970 15        ..
5971 16        int i_D  $\in I_D$ 
5972 17        {
5973 18            int i_d  $\in P$ 
5974 19            res[ i_1, \dots, i_d, \dots, i_D ] := lhs[ i_1, \dots, i_d, \dots, i_D ];
5975 20            int i_d  $\in Q$ 
5976 21            res[ i_1, \dots, i_d, \dots, i_D ] := rhs[ i_1, \dots, i_d, \dots, i_D ];
5977 22        }
5978 23    }
5979 24
5980 25     // point-wise combination
```

```

5979 26    $\forall d \in \mathbb{N} : \text{pw} << d >> < I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in \text{MDA-IDX-SETS}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} > ($ 
5980 27    $\oplus : T^{\text{OUT}} \times T^{\text{OUT}} \rightarrow T^{\text{OUT}}) ( \begin{array}{l} T^{\text{OUT}}[I_1, \dots, I_{d-1}, 0_f(P), I_{d+1}, \dots, I_D] \text{ lhs} , \\ T^{\text{OUT}}[I_1, \dots, I_{d-1}, 0_f(Q), I_{d+1}, \dots, I_D] \text{ rhs } \end{array} )$ 
5981 28    $-> T^{\text{OUT}}[I_1, \dots, I_{d-1}, 0_f(P \cup Q), I_{d+1}, \dots, I_D] \text{ res}$ 
5982 29    $\{$ 
5983 30      $\text{int } i\_1 \in I_1$ 
5984 31      $\ddots$ 
5985 32      $\text{int } i_{\{d-1\}} \in I_{d-1}$ 
5986 33      $\text{int } i_{\{d+1\}} \in I_{d+1}$ 
5987 34      $\ddots$ 
5988 35      $\text{int } i\_D \in I_D$ 
5989 36      $\{$ 
5990 37        $\text{res}[i\_1, \dots, i_{\{d-1\}}, 0, i_{\{d+1\}}, \dots, i\_D]$ 
5991 38        $:= \text{lhs}[i\_1, \dots, i_{\{d-1\}}, 0, i_{\{d+1\}}, \dots, i\_D]$ 
5992 39        $\text{atomic}(\oplus) \text{ rhs}[i\_1, \dots, i_{\{d-1\}}, 0, i_{\{d+1\}}, \dots, i\_D];$ 
5993 40      $\}$ 
5994 41    $\}$ 

```

Listing 11. Pre-Implemented Combine Operators

```

5996
5997 1 // 0.1.2. combine operators
5998 2
5999 3 // pre-implemented combine operators
6000 4
6001 5 // inverse concatenation
6002 6  $\forall d \in \mathbb{N}:$ 
6003 7  $cc\_inv << d >> < I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in \text{MDA-IDX-SETS}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} > ($ 
6004 8  $T^{\text{INP}}[I_1, \dots, I_{d-1}, id(P \cup Q), I_{d+1}, \dots, I_D] \text{ res} ) \rightarrow ( T^{\text{INP}}[I_1, \dots, I_{d-1}, id(P), I_{d+1}, \dots, I_D] \text{ lhs} ,$ 
6005 9  $T^{\text{INP}}[I_1, \dots, I_{d-1}, id(Q), I_{d+1}, \dots, I_D] \text{ rhs} )$ 
6006 10 {
6007 11     int i_1  $\in I_1$ 
6008 12     ...
6009 13     int i_{d-1}  $\in I_{d-1}$ 
6010 14     int i_{d+1}  $\in I_{d+1}$ 
6011 15     ...
6012 16     int i_D  $\in I_D$ 
6013 17     {
6014 18         int i_d  $\in P$ 
6015 19         res[ i_1, \dots, i_d, \dots, i_D ] =: lhs[ i_1, \dots, i_d, \dots, i_D ];
6016 20         int i_d  $\in Q$ 
6017         res[ i_1, \dots, i_d, \dots, i_D ] =: rhs[ i_1, \dots, i_d, \dots, i_D ];
6018 21     }
6019 22 }
6020 23 }
```

### Listing 12 Pre-Implemented Combine Operators

6018      *F.0.2*    `inp_view`.  
6019      Listing 13 shows the user-defined index functions provided by the user via higher-order function  
6020      `inp_view` (Definition 31).

```
6021
6022 1 // 0.2. inp_view
6023 2
6024 3 // index functions
6025 4  $\forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}, a \in [1, A_b^{\text{IB}}]_{\mathbb{N}}: \forall d \in [1, D_b^{\text{IB}}]_{\mathbb{N}}:$ 
6026 5 static
6027 6 idx<<INP>><<b,a>><<d>>( int i_MDA_1 , . . . , i_MDA_D ) -> int i_BUF_d
```

```
6028 7  {
6029 8  // ... (user defined)
6030 9 }
```

Listing 13. Index Functions (input)

6032

6033 *F.0.3 out\_view.*

6034 Listing 14 shows the user-defined index functions provided by the user via higher-order function  
 6035 *out\_view* (Definition 33).

```
6036 1 // 0.3. out_view
6037 2
6038 3 // index functions
6039 4  $\forall b \in [1, B^{OB}]_{\mathbb{N}}, a \in [1, A_b^{OB}]_{\mathbb{N}}: \forall d \in [1, D_b^{OB}]_{\mathbb{N}}:$ 
6040 5 static
6041 6  $\text{idx} << \text{OUT} >> << b, a >> << d >> ( \text{int } i_{\text{MDA\_1}}, \dots, i_{\text{MDA\_D}} ) \rightarrow \text{int } i_{\text{BUF\_d}}$ 
6042 7
6043 8 // ... (user defined)
6044 9 }
```

Listing 14. Index Functions (output)

6045

6046

6047 *F.0.4 BUFs.*

6048 Listing 15 shows our implementation of low-level BUFs (Definition 37). We compute BUFs' sizes  
 6049 using the ranges of their index functions (Definitions 31 and 33). Moreover, we partially evaluate  
 6050 BUFs' meta-parameters *MEM* (memory region) and  $\sigma$  (memory layout) immediately, as the same  
 6051 values are re-used for them during program runtime.

6052 The BUFs in lines 30 and 45 as well as in lines 69 and 84 represent the BUFs' transposed function  
 6053 representation (Definition 37), and the BUFs in lines 23, 37, and 52 as well as in lines 62, 76, and 91  
 6054 are the transposed BUFs' ordinary low-level BUF representation.

```
6055 1 // 0.4. BUFs
6056 2
6057 3 // 0.4.1. compute BUF sizes
6058 4  $\forall IO \in \{IB, OB\}: \forall b \in [1, B^{IO}]_{\mathbb{N}} \forall d \in [1, D_b^{IO}]_{\mathbb{N}}:$ 
6059 5 static  $N << IO >> << b >> << d >> ( \text{mda\_idx\_set } I_1, \dots, I_D ) \rightarrow \text{int } N_{b\_d}$ 
6060 6
6061 7  $N_{b\_d} := 0;$ 
6062 8
6063 9  $i_1 \in I_1$ 
6064 10  $\dots$ 
6065 11  $i_D \in I_D$ 
6066 12 {
6067 13      $\forall a \in [1, A_b^{IB}]_{\mathbb{N}}:$ 
6068 14      $N_{b\_d} :=_{\max} 1 + \text{idx} << IO >> << b, a >> << d >> ( i_1, \dots, i_D );$ 
6069 15 }
6070 16
6071 17
6072 18 // 0.4.2. input BUFs
6073 19
6074 20 // initial BUFs
6075 21  $\forall b \in [1, B^{IB}]_{\mathbb{N}}:$ 
6076 22 static  $l1_{IB} << \perp >> << b >> < \nabla_1^{(\perp)} \in \#PRT(1, 1), \dots, \nabla_D^{(\perp)} \in \#PRT(L, D) >> ( \text{int } i_1, \dots,$ 
6077 23      $\text{int } i_{D_b^{IB}} ) \rightarrow T_b^{IB} a$ 
6078 24 }
```

6079

```

6077 25     a := trans_ll_IB<<_1>><<b>><<  $\nabla_1^{(1)}, \dots, \nabla_D^{(1)} >[ i_1, \dots, i_D ];$ 
6078 26   }
6079 27
6080 28   // de-composition BUFS
6081 29    $\forall (l, d) \in \text{MDH-LVL}: \forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}:$ 
6082 30   auto trans_ll_IB<<l, d>><<b>><<  $\nabla_1^{(l, d)} \in \#PRT(1, 1), \dots, \nabla_D^{(l, d)} \in \#PRT(L, D) >$ 
6083 31   :=  $\downarrow \text{-mem}^{<b>} (l, d) T_b^{\text{IB}} [ N << \text{INP} >> <<b>><< \sigma_{\downarrow \text{-mem}}^{<b>} (l, d)(1) >> ( \frac{d}{\# \text{MDA}} (N_d) )_{d \in [1, D]_{\mathbb{N}} } ) ,$ 
6084 32   :
6085 33    $N << \text{INP} >> <<b>><< \sigma_{\downarrow \text{-mem}}^{<b>} (l, d)(D_b^{\text{IB}}) >> ( \frac{d}{\# \text{MDA}} (N_d) )_{d \in [1, D]_{\mathbb{N}} } ) ];$ 
6086 34
6087 35    $\forall (l, d) \in \text{MDH-LVL}: \forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}:$ 
6088 36   static ll_IB<<l, d>><<b>><<  $\nabla_1^{(l, d)} \in \#PRT(1, 1), \dots, \nabla_D^{(l, d)} \in \#PRT(L, D) >> ( \text{int } i_1, \dots,$ 
6089 37    $\text{int } i_D ) \rightarrow T_b^{\text{IB}} a$ 
6090 38   {
6091 39   a := trans_ll_IB<<l, d>><<b>><<  $\nabla_1^{(l, d)}, \dots, \nabla_D^{(l, d)} >[ i_1 \sigma_{\downarrow \text{-mem}}^{<b>} (l, d)(1) , \dots ,$ 
6092 40    $i_1 \sigma_{\downarrow \text{-mem}}^{<b>} (l, d)(D_b^{\text{IB}}) ];$ 
6093 41   }
6094 42
6095 43   // scalar BUFS
6096 44    $\forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}:$ 
6097 45   auto trans_ll_IB<f>><<b>><<  $\nabla_1^{(f)} \in \#PRT(1, 1), \dots, \nabla_D^{(f)} \in \#PRT(L, D) >$ 
6098 46   :=  $\text{f} \downarrow \text{-mem}^{<b>} T_b^{\text{IB}} [ N << \text{INP} >> <<b>><< \sigma_{\text{f} \downarrow \text{-mem}}^{<b>} (1) >> ( \frac{d}{\# \text{MDA}} (N_d) )_{d \in [1, D]_{\mathbb{N}} } ) ,$ 
6099 47   :
6100 48    $N << \text{INP} >> <<b>><< \sigma_{\text{f} \downarrow \text{-mem}}^{<b>} (D_b^{\text{IB}}) >> ( \frac{d}{\# \text{MDA}} (N_d) )_{d \in [1, D]_{\mathbb{N}} } ) ];$ 
6101 49
6102 50    $\forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}:$ 
6103 51   static ll_IB<<f>><<b>><<  $\nabla_1^{(f)} \in \#PRT(1, 1), \dots, \nabla_D^{(f)} \in \#PRT(L, D) >> ( \text{int } i_1, \dots,$ 
6104 52    $\text{int } i_D ) \rightarrow T_b^{\text{IB}} a$ 
6105 53   {
6106 54   a := trans_ll_IB<<f>><<b>><<  $\nabla_1^{(f)}, \dots, \nabla_D^{(f)} >[ i_1 \sigma_{\text{f} \downarrow \text{-mem}}^{<b>} (1) , \dots , i_1 \sigma_{\text{f} \downarrow \text{-mem}}^{<b>} (D_b^{\text{IB}}) ];$ 
6107 55   }
6108 56
6109 57   // 0.4.3. output BUFS
6110 58
6111 59   // initial BUFS
6112 60    $\forall b \in [1, B^{\text{OB}}]_{\mathbb{N}}:$ 
6113 61   static ll_OB<<_1>><<b>><<  $\blacktriangle_1^{(1)} \in \#PRT(1, 1), \dots, \blacktriangle_D^{(1)} \in \#PRT(L, D) >> ( \text{int } i_1, \dots,$ 
6114 62    $\text{int } i_D ) \rightarrow T_b^{\text{OB}} a$ 
6115 63   {
6116 64   a := trans_ll_OB<<_1>><<b>><<  $\blacktriangle_1^{(1)}, \dots, \blacktriangle_D^{(1)} >[ i_1 , \dots , i_D ];$ 
6117 65   }
6118 66
6119 67   // re-composition BUFS
6120 68    $\forall (l, d) \in \text{MDH-LVL}: \forall b \in [1, B^{\text{OB}}]_{\mathbb{N}}:$ 
6121 69   auto trans_ll_OB<<l, d>><<b>><<  $\blacktriangle_1^{(l, d)} \in \#PRT(1, 1), \dots, \blacktriangle_D^{(l, d)} \in \#PRT(L, D) >$ 
6122 70
6123 71
6124 72
6125 73

```

```

6126 70   :=  $\uparrow \text{-mem}^{<\flat>} (l, d)$   $T_b^{\text{OB}} [$   $\text{N}^{<<\text{OUT}>><<b>><<}$   $\sigma_{\uparrow\text{-mem}}^{<\flat>} (l, d)(1)$   $>> ($   $(\xrightarrow[\oplus]{d_{\text{MDA}}} (N_d))_{d \in [1, D]_{\mathbb{N}}}$   $)$  ,
6127 71   :
6128 72    $\text{N}^{<<\text{OUT}>><<b>><<}$   $\sigma_{\uparrow\text{-mem}}^{<\flat>} (l, d)(D_b^{\text{OB}})$   $>> ($   $(\xrightarrow[\oplus]{d_{\text{MDA}}} (N_d))_{d \in [1, D]_{\mathbb{N}}}$   $)$  ];
6129 73
6130 74    $\forall (l, d) \in \text{MDH-LVL} : \forall b \in [1, B^{\text{OB}}]_{\mathbb{N}} :$ 
6131 75   static  $\text{ll\_OB}^{<<l, d>><<b>><<}$   $\Delta_1^{(l, d)} \in \#PRT(1, 1)$ , ...,  $\Delta_D^{(l, d)} \in \#PRT(L, D)$   $> ($  int  $i_1, \dots,$ 
6132 76   int  $i_D$   $\rightarrow T_b^{\text{OB}}$  a
6133 77   {
6134 78   a :=  $\text{trans\_ll\_OB}^{<<l, d>><<b>><<}$   $\Delta_1^{(l, d)}, \dots, \Delta_D^{(l, d)}$   $> [$   $i_{-} \sigma_{\uparrow\text{-mem}}^{<\flat>} (l, d)(1)$  , ..., ,
6135 79    $i_{-} \sigma_{\uparrow\text{-mem}}^{<\flat>} (l, d)(D_b^{\text{OB}})$  ];
6136 80   }
6137 81
6138 82   // scalar BUFS
6139 83    $\forall b \in [1, B^{\text{OB}}]_{\mathbb{N}} :$ 
6140 84   auto  $\text{trans\_ll\_OB}^{<<\text{f}>><<b>><<}$   $\Delta_1^{(f)} \in \#PRT(1, 1)$ , ...,  $\Delta_D^{(f)} \in \#PRT(L, D)$  >
6141 85   :=  $\text{f}^{\uparrow\text{-mem}}^{<\flat>} T_b^{\text{OB}} [$   $\text{N}^{<<\text{OUT}>><<b>><<}$   $\sigma_{\uparrow\text{-mem}}^{<\flat>} (1)$   $>> ($   $(\xrightarrow[\oplus]{d_{\text{MDA}}} (N_d))_{d \in [1, D]_{\mathbb{N}}}$   $)$  ,
6142 86   :
6143 87    $\text{N}^{<<\text{OUT}>><<b>><<}$   $\sigma_{\uparrow\text{-mem}}^{<\flat>} (D_b^{\text{OB}})$   $>> ($   $(\xrightarrow[\oplus]{d_{\text{MDA}}} (N_d))_{d \in [1, D]_{\mathbb{N}}}$   $)$  ];
6144 88
6145 89    $\forall b \in [1, B^{\text{OB}}]_{\mathbb{N}} :$ 
6146 90   static  $\text{ll\_OB}^{<<\text{f}>><<b>><<}$   $\Delta_1^{(f)} \in \#PRT(1, 1)$ , ...,  $\Delta_D^{(f)} \in \#PRT(L, D)$   $> ($  int  $i_1, \dots,$ 
6147 91   int  $i_D$   $\rightarrow T_b^{\text{OB}}$  a
6148 92   {
6149 93   a :=  $\text{trans\_ll\_OB}^{<<\text{f}>><<b>><<}$   $\Delta_1^{(f)}, \dots, \Delta_D^{(f)}$   $> [$   $i_{-} \sigma_{\uparrow\text{-mem}}^{<\flat>} (1)$  , ...,  $i_{-} \sigma_{\uparrow\text{-mem}}^{<\flat>} (D_b^{\text{OB}})$  ];
6150 94   }

```

Listing 15. Low-Level BUFS

6154  
6155  
6156 where  $\Delta_{\flat}^{(\bullet)}$  and  $\Delta_{\flat}^{(\bullet)}$ , for  $\bullet \in \{\perp\} \cup \text{MDH-LVL} \cup \{\text{f}\}$ , are textually replaced by:

6158  
6159  
6160  
6161  
6162

$$\Delta_{\flat}^{(\bullet)} = \begin{cases} p_{\flat}^{\perp} & : \sigma_{\downarrow\text{-ord}} (l, d) < \bullet \\ * & : \text{else} \end{cases}$$

6163  
6164  
6165  
6166

$$\Delta_{\flat}^{(\bullet)} = \begin{cases} p_{\flat}^{\perp} & : \sigma_{\uparrow\text{-ord}} (l, d) < \bullet \\ * & : \text{else} \end{cases}$$

6167 (symbol \* is taken from Definition 22) where < is defined according to the lexicographical order on  
6168  $\text{MDH-LVL} = [1, L]_{\mathbb{N}} \times [1, D]_{\mathbb{N}}$ , and:

6169  
6170

$$\forall (l, d) \in \text{MDH-LVL} : \perp < (l, d) < \text{f}$$

6171 Functions

6172  
6173

$$\xrightarrow[\oplus]{1}^{\text{MDA}}, \dots, \xrightarrow[\oplus]{D}^{\text{MDA}}$$

6174

6175 are the index set functions  $id$  of combine operator concatenation  $+$  (Example 13), and functions

$$\begin{array}{c} 1 \text{ MDA} \\ \overrightarrow{\oplus} \text{ MDA} , \dots , \overrightarrow{\oplus} \text{ MDA} \\ \overrightarrow{\oplus} \end{array}$$

6178 are the index set functions of combine operators  $\otimes_1, \dots, \otimes_D$ .

6179 Note that we use generous BUFs sizes (lines 31-33, 46-48, 70-72, 85-87), as imperative-style  
6180 programming models usually struggle with non-contiguous index ranges. We discuss optimizations  
6181 targeting BUF sizes in Section G.

6182 Note further that we do not need to initialize output buffers with neutral elements of combine  
6183 operators in lines 64, 79, and 93 of Listing 15, as the buffers are initialized implicitly in the re-  
6184 composition phase (Section F.3).

#### 6186 F.0.5 MDAs.

6187 Listing 16 shows our implementation of low-level MDAs (Definitions 36 and 40).

6188 Note that for a particular choice of meta-parameters, low-level BUFs (Definition 37) are ordinary  
6189 BUFs (Definition 28), as required by the types of functions `inp_view` and `out_view` (Definitions 31  
6190 and 33).

```

6191 1 // 0.5. MDAs
6192 2
6193 3 // 0.5.1. partitioned index sets
6194 4  $\forall d \in [1, D]_{\mathbb{N}}$ :
6195 5 static  $I << d >> < p_d^1 \in \#PRT(1, d) , \dots , p_d^L \in \#PRT(L, d) > ( \text{int } j' ) \rightarrow \text{int } i_j$ 
6196 6 {
6197 7      $i_j := ( p_d^1 * ( N_d / ( \#PRT(1, d) ) ) + \dots + ( p_d^L * ( N_d / ( \#PRT(1, d) * \dots * \#PRT(L, d) ) ) + j' )$ ;
6198 8
6199 9
6200 10
6201 11
6202 12 // 0.5.2. input MDAs
6203 13  $\forall \bullet \in \{\perp\} \cup \text{MDH-LVL} \cup \{f\}$ :
6204 14 static  $ll\_inp\_mda << \bullet >> < \overset{(\bullet)}{\nabla}_1^1 \in \#PRT(1, 1) , \dots , \overset{(\bullet)}{\nabla}_D^L \in \#PRT(L, D) > ( \text{int } i_1 , \dots ,$ 
6205 15      $\text{int } i_D ) \rightarrow T_b^{IB} a$ 
6206 16 {
6207 17      $\forall b \in [1, B^{IB}]_{\mathbb{N}}, a \in [1, A_b^{IB}]_{\mathbb{N}}$ :
6208 18      $a := ll\_INP << \bullet >> << b >> < \overset{(\bullet)}{\nabla}_1^1 , \dots , \overset{(\bullet)}{\nabla}_D^L > ( \text{idx} << \text{INP} >> << b , a >> << 1 >> ( i_1 , \dots , i_D ) ,$ 
6209 19
6210 20      $\vdots$ 
6211 21      $\text{idx} << \text{INP} >> << b , a >> << D_b^{IB} >> ( i_1 , \dots , i_D ) )$ ;
6212 22
6213 23 // 0.5.3. output MDAs
6214 24  $\forall \bullet \in \{\perp\} \cup \text{MDH-LVL} \cup \{f\}$ :
6215 25 static  $ll\_out\_mda << \bullet >> < \overset{(\bullet)}{\Delta}_1^1 \in \#PRT(1, 1) , \dots , \overset{(\bullet)}{\Delta}_D^L \in \#PRT(L, D) > ( \text{int } i_1 , \dots ,$ 
6216 26      $\text{int } i_D ) \rightarrow T_b^{OB} a$ 
6217 27 {
6218 28      $\forall b \in [1, B^{OB}]_{\mathbb{N}}, a \in [1, A_b^{OB}]_{\mathbb{N}}$ :
6219 29      $a := ll\_OUT << \bullet >> << b >> < \overset{(\bullet)}{\Delta}_1^1 , \dots , \overset{(\bullet)}{\Delta}_D^L > ( \text{idx} << \text{OUT} >> << b , a >> << 1 >> ( i_1 , \dots , i_D ) ,$ 
6220 30
6221 31      $\vdots$ 
6222 32      $\text{idx} << \text{OUT} >> << b , a >> << D_b^{OB} >> ( i_1 , \dots , i_D ) )$ ;
```

6222 Listing 16. Low-Level MDAs

6224 For computing the partitioned index sets (lines 3-10), we exploit the following proposition.

6225 **Proposition 1.** Let  $\alpha \in T[N_1, \dots, N_D]$  be an arbitrary MDA that operates on contiguous index  
6226 sets  $[1, N_d]_{\mathbb{N}}$ ,  $d \in [1, D]_{\mathbb{N}}$ . Let further be

$$6228 \quad \alpha^{<(p_1^1, \dots, p_d^1) \in P_1^1 \times \dots \times P_D^1 \mid \dots \mid (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>} : I_1^{<p_1^1 \dots p_1^L>} \times \dots \times I_D^{<p_D^1 \dots p_D^L>} \rightarrow T$$

6229 an arbitrary  $L$ -layered,  $D$ -dimensional,  $P$ -partitioning of MDA  $\alpha$ .

6230 It holds that  $j$ -th element within an MDA's part is accessed via index  $j$ :

$$6232 \quad I_d^{<p_d^1 \dots p_d^L>} = \left\{ \underbrace{\sum_{l \in [1, L]_{\mathbb{N}}} p_d^l * \frac{N_d}{\prod_{l' \in [1, l]_{\mathbb{N}}} P_d^{l'}} + 0,}_{\text{OS}} \underbrace{\sum_{l \in [1, L]_{\mathbb{N}}} p_d^l * \frac{N_d}{\prod_{l' \in [1, l]_{\mathbb{N}}} P_d^{l'}} + 1, \dots}_{\text{OS}} \right\}$$

6236 PROOF. Since MDA  $\alpha$ 's index sets are contiguous ranges of natural numbers, it holds the  $i_j$  – the  
6237 index to access the  $j$ -th element within an MDA's part (Definition 40) – is equal to  $j$  itself.  $\square$

## 6239 F.1 De-Composition Phase

6240 Listing 17 shows our implementation of the de-composition phase (Figure 31).

```
6241 1 // 1. de-composition phase
6242 2
6243 3 // 1.1. initialization
6244 4 ll_inp_mda<<1>> =: ll_inp_mda<<σ↓-ord(1, 1)>>
6245 5
6246 6 // 1.2. main
6247 7 int p_ σ↓-ord(1, 1) ∈ <↔-ass (1, 1) > #PRT ( σ↓-ord(1, 1) )
6248 8 {
6249 9     ll_inp_mda<<σ↓-ord(1, 1)>> =: cc< σ↓-ord(1, 1) > inp_mda<<σ↓-ord(1, 2)>>;
6250 10    int p_ σ↓-ord(1, 2) ∈ <↔-ass (1, 2) > #PRT ( σ↓-ord(1, 2) )
6251 11    {
6252 12        ll_inp_mda<<σ↓-ord(1, 2)>> =: cc< σ↓-ord(1, 2) > inp_mda<<σ↓-ord(1, 3)>>;
6253 13        ..
6254 14        int p_ σ↓-ord(L, D) ∈ <↔-ass (L, D) > #PRT ( σ↓-ord(L, D) )
6255 15        {
6256 16            ll_inp_mda<<σ↓-ord(L, D)>> =: cc< σ↓-ord(L, D) > inp_mda<<f>>;
6257 17        }
6258 18        ..
6259 19    }
6260 20 }
```

6264 Listing 17. De-Composition Phase

6265 where

$$6266 \quad ll\_inp\_mda<<l, d>> =: cc<l, d> ll\_inp\_mda<<l', d'>>$$

6268 abbreviates

$$6269 \quad ll\_inp\_mda<<l', d'>><\nabla_1^{(l', d')}, \dots, \nabla_D^{(l', d')}>, ll\_inp\_mda<<l, d>><\nabla_1^{(l, d)}, \dots, \nabla_D^{(l, d)}>
6270 \quad := cc\_inv<<d>><\frac{1}{MDA} MDA ( I <<1>> <\blacksquare_1^{(l, d)}, \dots, \blacksquare_1^{(l, d)} >> (0) ), // I_1
6271$$

```

6273      :
6274       $\stackrel{D}{\Rightarrow}_{\#MDA}^{MDA}( I << D >> < \underset{D}{\blacksquare}_1^{(l,d)}, \dots, \underset{D}{\blacksquare}_L^{(l,d)} > (0) ), // I_D$ 
6275
6276
6277       $\stackrel{d}{\Rightarrow}_{\#MDA}^{MDA}( I << d >> < \underset{d}{\boxplus}_1^{(l,d)}, \dots, \underset{d}{\boxplus}_L^{(l,d)} > (0) ), // P$ 
6278
6279       $\stackrel{d}{\Rightarrow}_{\#MDA}^{MDA}( I << d >> < \underset{d}{\boxtimes}_1^{(l,d)}, \dots, \underset{d}{\boxtimes}_L^{(l,d)} > (0) ) // Q$ 
6280      > ( 11_inp_mda << l, d >> < \underset{1}{\blacktriangledown}^{(l,d)}, \dots, \underset{D}{\blacktriangledown}^{(l,d)} > )
6281

```

6282 Here, functions  $\stackrel{1}{\#MDA}, \dots, \stackrel{D}{\#MDA}$  are the index set functions  $id$  of combine operator concatenation  $\#_1, \dots, \#_D$  (Example 13), and  $\underset{d}{\blacksquare}_1^{(\bullet)}, \underset{d}{\boxplus}_1^{(\bullet)}, \underset{d}{\boxtimes}_1^{(\bullet)}$ , for  $\bullet \in \text{MDH-LVL}$ , are textually replaced by:

$$\underset{d}{\blacksquare}_1^{(\bullet)} := \begin{cases} p_-(l, d) & : \sigma_{\downarrow\text{-ord}}(l, d) < \bullet \\ [p_-(l, d), \#PRT(l, d)]_{\mathbb{N}_0} & : (l, d) = \bullet \\ [0, \#PRT(l, d)]_{\mathbb{N}_0} & : \sigma_{\downarrow\text{-ord}}(l, d) > \bullet \end{cases}$$

$$\underset{d}{\boxplus}_1^{(\bullet)} := \begin{cases} p_-(l, d) & : \sigma_{\downarrow\text{-ord}}(l, d) < \bullet \\ p_-(l, d) & : (l, d) = \bullet \\ [0, \#PRT(l, d)]_{\mathbb{N}_0} & : \sigma_{\downarrow\text{-ord}}(l, d) > \bullet \end{cases}$$

$$\underset{d}{\boxtimes}_1^{(\bullet)} := \begin{cases} p_-(l, d) & : \sigma_{\downarrow\text{-ord}}(l, d) < \bullet \\ (p_-(l, d), \#PRT(l, d)]_{\mathbb{N}_0} & : (l, d) = \bullet \\ [0, \#PRT(l, d)]_{\mathbb{N}_0} & : \sigma_{\downarrow\text{-ord}}(l, d) > \bullet \end{cases}$$

6305 where  $<$  is defined as lexicographical order, according to Section F.0.4.

6306 Note that we re-use `inp_mda<<l, d>>` for the intermediate results given by different iterations of 6307 variable  $p_-(l, d)$ . Correctness is ensured, as it holds:

$$6309 A \subseteq B \Rightarrow \stackrel{d}{\#MDA}^{MDA}(A) \subseteq \stackrel{d}{\#MDA}^{MDA}(B)$$

6311 MDA `inp_mda<<l, d>>` has the following type when used for the intermediate result in a particular 6312 iteration of  $p_-(l, d)$ :

$$6314 \stackrel{1}{\#MDA}^{MDA}( I_1^{< \underset{1}{\blacksquare}_1^{(l,d)}, \dots, \underset{D}{\blacksquare}_1^{(l,d)} | \dots | \underset{1}{\blacksquare}_L^{(l,d)}, \dots, \underset{D}{\blacksquare}_L^{(l,d)} >} ) \times \dots \times \stackrel{D}{\#MDA}^{MDA}( I_D^{< \underset{1}{\blacksquare}_1^{(l,d)}, \dots, \underset{D}{\blacksquare}_1^{(l,d)} | \dots | \underset{1}{\blacksquare}_L^{(l,d)}, \dots, \underset{D}{\blacksquare}_L^{(l,d)} >} ) \rightarrow T^{\text{INP}}$$

6316 Here, for a set  $P \subseteq [0, \#PRT(l, d)]_{\mathbb{N}_0}$ , index set  $I_d^{<\dots| \dots P \dots | \dots >}$  denotes  $\bigcup_{p_d^l \in P} I_d^{<\dots| \dots p_d^l \dots | \dots >}$ .

## F.2 Scalar Phase

6320 Listing 18 shows our implementation of the scalar phase (Figure 31).

```

6322 1 // 2. scalar phase
6323 2 | int p_  $\sigma_{f\text{-}ord}(1, 1)$   $\in \leftarrow_{f\text{-}ass}(1, 1)$  #PRT(  $\sigma_{f\text{-}ord}(1, 1)$  )
6324 3 | `.
6325 4 | int p_  $\sigma_{f\text{-}ord}(L, D)$   $\in \leftarrow_{f\text{-}ass}(L, D)$  #PRT(  $\sigma_{f\text{-}ord}(L, D)$  )
6326 5 |
6327 6 {
6328 7 | (
6329 8 |     ll_out_mda<<f>><<
6330 9 |         p_(1, 1) , . . . , p_(1, D) ,
6331 10 |         . . .
6332 11 |         p_(L, 1) , . . . , p_(L, D)>><<b, a>>(
6333 12 |              $\stackrel{1}{\Rightarrow}_{MDA}$ ( I<<1>><p_(1, 1) , . . . , p_(L, 1)>(0) ) ,
6334 13 |             . . .
6335 14 |              $\stackrel{D}{\Rightarrow}_{MDA}$ ( I<<D>><p_(1, D) , . . . , p_(L, D)>(0) ) )
6336 15 |         )b ∈ [1, BOB]N, a ∈ [1, AbOB]N := f( ( ll_inp_mda<<f>><< p_(1, 1) , . . . , p_(1, D) ,
6337 16 |             . . .
6338 17 |             p_(L, 1) , . . . , p_(L, D)>><<b, a>>(
6339 18 |                  $\stackrel{d}{\Rightarrow}_{MDA}$ ( I<<1>><p_(1, 1) , . . . , p_(L, 1)>(0) ) ,
6340 19 |                 . . .
6341 20 |                  $\stackrel{d}{\Rightarrow}_{MDA}$ ( I<<D>><p_(1, D) , . . . , p_(L, D)>(0) ) )
6342 21 |         )b ∈ [1, BIB]N, a ∈ [1, AbIB]N )
6343
6344
6345

```

Listing 18. Scalar Phase

```

6346
6347
6348
6349
6350
6351
6352
6353
6354
6355
6356
6357
6358
6359
6360
6361
6362
6363
6364
6365
6366
6367
6368
6369
6370

```

### 6371 F.3 Re-Composition Phase

6372 Listing 19 shows our implementation of the re-composition phase (Figure 31).

6373

```

6374 1 // 3. re-composition phase
6375 2
6376 3 // 3.1. main
6377 4 int p_  $\sigma_{\uparrow\text{-ord}}(1, 1)$   $\in \leftrightarrow_{\uparrow\text{-ass}}(1, 1)$  > #PRT(  $\sigma_{\uparrow\text{-ord}}(1, 1)$  )
6378 5 {
6379 6     int p_  $\sigma_{\uparrow\text{-ord}}(1, 2)$   $\in \leftrightarrow_{\uparrow\text{-ass}}(1, 2)$  > #PRT(  $\sigma_{\uparrow\text{-ord}}(1, 2)$  )
6380 7     {
6381 8         ..
6382 9             int p_  $\sigma_{\uparrow\text{-ord}}(L, D)$   $\in \leftrightarrow_{\uparrow\text{-ass}}(L, D)$  > #PRT(  $\sigma_{\uparrow\text{-ord}}(L, D)$  )
6383 10            {
6384 11                ll_out_mda<<  $\sigma_{\uparrow\text{-ord}}(L, D)$  >> :=  $_{\text{co}< \sigma_{\uparrow\text{-ord}}(L, D) >}$  out_mda<<f>>;
6385 12            }
6386 13            ..
6387 14                ll_out_mda<<  $\sigma_{\uparrow\text{-ord}}(1, 2)$  >> :=  $_{\text{co}< \sigma_{\uparrow\text{-ord}}(1, 2) >}$  out_mda<<  $\sigma_{\uparrow\text{-ord}}(1, 3)$  >>;
6388 15            }
6389 16                ll_out_mda<<  $\sigma_{\uparrow\text{-ord}}(1, 1)$  >> :=  $_{\text{co}< \sigma_{\uparrow\text{-ord}}(1, 1) >}$  out_mda<<  $\sigma_{\uparrow\text{-ord}}(1, 2)$  >>;
6390 17            }
6391 18
6392 19 // 3.2. finalization
6393 20 ll_out_mda<<1>> := ll_out_mda<<  $\sigma_{\uparrow\text{-ord}}(1, 1)$  >>
6394
6395
6396

```

6397 Listing 19. Re-Composition Phase

6398

6399 where

6400

6401  $\text{ll\_out\_mda}<<l, d>> :=_{\text{co}<l, d>} \text{ll\_out\_mda}<<l', d'>>$

6402

6403 abbreviates

6404

```

6405      $\text{ll\_out\_mda}<<l, d>>< \overset{(l, d)}{\blacktriangle}_1, \dots, \overset{(l, d)}{\blacktriangle}_D >$ 
6406     :=  $\text{co}<<\text{d}>>< \overset{1}{\text{MDA}} \otimes \overset{D}{\text{MDA}}( \text{I} << 1 >> < \overset{(l, d)}{\blacksquare}_1, \dots, \overset{(l, d)}{\blacksquare}_D > (0) ), // I_1$ 
6407     :
6408     //  $\dots, I_{d-1}, I_{d+1}, \dots$ 
6409      $\overset{D}{\text{MDA}} \otimes \overset{D}{\text{MDA}}( \text{I} << D >> < \overset{(l, d)}{\blacksquare}_D, \dots, \overset{(l, d)}{\blacksquare}_D > (0) ), // I_D$ 
6410
6411      $\overset{d}{\text{MDA}} \otimes \overset{d}{\text{MDA}}( \text{I} << d >> < \overset{(l, d)}{\boxplus}_d, \dots, \overset{(l, d)}{\boxplus}_d > (0) ), // P$ 
6412      $\overset{d}{\text{MDA}} \otimes \overset{d}{\text{MDA}}( \text{I} << d >> < \overset{(l, d)}{\boxtimes}_d, \dots, \overset{(l, d)}{\boxtimes}_d > (0) ) // Q$ 
6413     > (  $\text{ll\_out\_mda}<<l, d>>< \overset{(l, d)}{\blacktriangle}_1, \dots, \overset{(l, d)}{\blacktriangle}_D >$  ,
6414      $\text{ll\_out\_mda}<<l', d'>>< \overset{(l', d')}{\blacktriangle}_1, \dots, \overset{(l', d')}{\blacktriangle}_D >$  )
6415
6416
6417
6418
6419

```

6420 Here, functions  $\Rightarrow_{\otimes \text{MDA}}^1, \dots, \Rightarrow_{\otimes \text{MDA}}^D$  are the index set function of combine operators  $\otimes_1, \dots, \otimes_D$  (Defi-  
 6421  
 6422 nition 25), and  $\blacksquare_{\mathfrak{d}}^{\bullet}, \boxplus_{\mathfrak{d}}^{\bullet}, \boxtimes_{\mathfrak{d}}^{\bullet}$ , for  $\bullet \in \text{MDH-LVL}$ , are textually replaced by:  
 6423

$$\blacksquare_{\mathfrak{d}}^{\bullet} := \begin{cases} p_{-}(\mathfrak{l}, \mathfrak{d}) & : \sigma_{\uparrow\text{-ord}}(\mathfrak{l}, \mathfrak{d}) < \bullet \\ [0, p_{-}(\mathfrak{l}, \mathfrak{d})]_{N_0} & : (\mathfrak{l}, \mathfrak{d}) = \bullet \\ [0, \#PRT(\mathfrak{l}, \mathfrak{d})]_{N_0} & : \sigma_{\uparrow\text{-ord}}(\mathfrak{l}, \mathfrak{d}) > \bullet \end{cases}$$

$$\boxplus_{\mathfrak{d}}^{\bullet} := \begin{cases} p_{-}(\mathfrak{l}, \mathfrak{d}) & : \sigma_{\uparrow\text{-ord}}(\mathfrak{l}, \mathfrak{d}) < \bullet \\ [0, p_{-}(\mathfrak{l}, \mathfrak{d})]_{N_0} & : (\mathfrak{l}, \mathfrak{d}) = \bullet \\ [0, \#PRT(\mathfrak{l}, \mathfrak{d})]_{N_0} & : \sigma_{\uparrow\text{-ord}}(\mathfrak{l}, \mathfrak{d}) > \bullet \end{cases}$$

$$\boxtimes_{\mathfrak{d}}^{\bullet} := \begin{cases} p_{-}(\mathfrak{l}, \mathfrak{d}) & : \sigma_{\uparrow\text{-ord}}(\mathfrak{l}, \mathfrak{d}) < \bullet \\ p_{-}(\mathfrak{l}, \mathfrak{d}) & : (\mathfrak{l}, \mathfrak{d}) = \bullet \\ [0, \#PRT(\mathfrak{l}, \mathfrak{d})]_{N_0} & : \sigma_{\uparrow\text{-ord}}(\mathfrak{l}, \mathfrak{d}) > \bullet \end{cases}$$

6435 where  $<$  is defined as lexicographical order, according to Section F.0.4.

6440 Note that we assume for index set functions  $\Rightarrow_{\otimes \text{MDA}}^d$  that

$$A \subseteq B \Rightarrow \Rightarrow_{\otimes \text{MDA}}^d(A) \subseteq \Rightarrow_{\otimes \text{MDA}}^d(B)$$

6445 (which holds for all kinds of index set functions used in this paper, e.g., in Examples 13 and 14)  
 6450 so that we can re-use `out_mda<<l, d>>` for the intermediate results given by different iterations  
 6451 of `p_{-}(l, d)`. `MDA inp_mda<<l, d>>` has the following type when used for the intermediate result  
 6452 in a particular iteration of variable `p_{-}(l, d)`:

$$T^{\text{INP}} \left[ \Rightarrow_{\otimes \text{MDA}}^1(I_1^{< \blacksquare_1^{(l, d)_1}, \dots, \blacksquare_D^{(l, d)_D} | \dots | \blacksquare_1^{(l, d)_L}, \dots, \blacksquare_D^{(l, d)_L} >}), \dots, \Rightarrow_{\otimes \text{MDA}}^D(I_D^{< \blacksquare_1^{(l, d)_1}, \dots, \blacksquare_D^{(l, d)_D} | \dots | \blacksquare_1^{(l, d)_L}, \dots, \blacksquare_D^{(l, d)_L} >}) \right]$$

6453 Note that in line 11 of Listing 19, we implicitly override the uninitialized value in `out_mda<<l, d>>`  
 6454 (not explicitly stated in the listing for brevity), thereby avoiding initializing output buffers with  
 6455 neutral elements of combine operators.

## G CODE-LEVEL OPTIMIZATIONS

6456 We consider optimizations that are below the abstraction level of our low-level representation  
 6457 (Section 3) as *code-level optimizations*. For some of these code-level optimizations, like loop fusion,  
 6458 we do not want to rely on the underlying compiler (e.g., the OpenMP/CUDA/OpenCL compiler): we  
 6459 exactly know the structure of our code presented in Section F, thus being able to implement  
 6460 code-level optimizations without requiring complex compiler analyses for optimizations.

6461 Since code-level optimizations are not the focus of this work, we give a brief, basic overview of  
 6462 our these optimizations which our systems performs for the underlying compiler (OpenMP, CUDA,  
 6463

6469 OpenCL, etc). We will thoroughly discuss our code-level optimizations and how we apply them  
6470 to our code in future work. Further code-level optimizations, like loop unrolling, are currently  
6471 mostly left to the underlying compiler, e.g., the OpenMP, CUDA, or OpenCL compiler. In our future  
6472 work, we aim to incorporate code-level optimizations, such as loop unrolling, into our auto-tunable  
6473 optimization process.

6474 *Loop Fusion.* In Listings 17, 18, 19, the lines containing symbol " $\epsilon$ " are mapped to for-loops. These  
6475 loops can often be fused; for example, when parameters 3, 6, 11 in Table 1 coincide (as in Figure 15).  
6476 Besides reducing the overhead caused by loop control structures, loop fusion in particular allows  
6477 significantly reducing the memory footprint: we can re-use the same memory region for each BUF  
6478 partition (Definition 37), rather than allocating memory for all these partition.  
6479

6480 *Buffer Elimination.* In Listing 15, we allocate BUFs for each combination of a layer and dimension.  
6481 However, when memory regions and memory layouts of BUFs coincide, we can avoid a new BUF  
6482 allocation, by re-using the BUF of the upper level, thereby again reducing memory footprint.  
6483

6484 *Buffer Size Reduction.* We can reduce the sizes of BUFs when specific classes of index functions  
6485 are used for views (Definitions 31 and 33). For example, in the case of *Dot Product* (DOT) (Figure 14),  
6486 when accessing its input in a strided fashion – via index function  $k \mapsto (2 * k)$ , instead of function  
6487  $k \mapsto (k)$  (as in Figure 14) – we would have to allocate BUFs (Listing 15, lines 30 and 45) of size  
6488  $2 * K$  for an input size of  $K \in \mathbb{N}$ ; in these BUFs, each second value (accessed via indices  $2 * k + 1$ )  
6489 would be undefined. We avoid this waste of memory, by using index function  $k \mapsto (k)$  instead  
6490 of  $k \mapsto (2 * k)$  for allocated BUFs (Listing 16, lines 18-20 and 29-31, case " $\bullet \neq \perp$ "), which avoids  
6491 index functions' leading factors and potential constant additions. Thereby, we reduce the memory  
6492 footprint from  $2 * K$  to  $K$ . Furthermore, according to our partitioning strategy (Listing 16, line 5,  
6493 and Listings 17, 18, 19), we often access BUFs via offsets:  $k \mapsto (2 * k)$  for  $k \in \{OS + k' \mid k' \in \mathbb{N}\}$  and  
6494 offset  $OS \in \mathbb{N}$ . We avoid such offset by using  $k \mapsto (2 * (k - OS))$ , thereby further reducing the  
6495 memory footprint.  
6496

6496 *Memory Operation Minimization.* In our code, we access BUFs uniformly via MDAs (Listing 16),  
6497 which may cause unnecessary memory operations. For example, in the de-composition phase  
6498 (Listing 17) of, e.g., matrix multiplication (MatMul) (Figure 14), we iterate over all dimensions of the  
6499 input (i.e., the  $i, j, k$  dimensions) for de-composition (Listing 12). However, the  $A$  input matrix of  
6500 MatMul is accessed via MDA indices  $i$  and  $k$  only (Figure 14). We avoid these unnecessary memory  
6501 operations ( $J$ -many in the case of an input MDA of size  $J$  in dimension  $j$ ) by using index 0 only  
6502 in dimension  $j$  for the de-composition of the  $A$  matrix. Analogously, we use index 0 only in the  $i$   
6503 dimension for the de-composition of MatMul's  $B$  matrix which is accessed via MDA indices  $k$  and  $j$   
6504 only. Moreover, we exploit all available parallelism for memory copy operations. For example, for  
6505 MatMul, we use also the threads intended for the  $j$  dimension when de-composing the  $A$  matrix,  
6506 and we use the threads in the  $i$  dimension for the  $B$  matrix. For this, we flatten the thread ids over  
6507 all dimensions  $i, j, k$  and re-structure them only in dimensions  $i, k$  (for the  $A$  matrix) or  $k, j$  (for the  
6508  $B$  matrix).  
6509

6510 *Constant Substitution.* We use constants whenever possible. For example, in CUDA, variable  
6511 `threadIdx.x` returns the thread id in dimension  $x$ . However, in our code, we use constant  $0$  instead  
6512 of `threadIdx.x` when only one thread is started in dimension  $x$ , enabling the CUDA compiler to  
6513 significantly simplify arithmetic expressions.  
6514  
6515  
6516  
6517