

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

3 A Deep Dive into the MDH Formalism

4  
5  
6 ANONYMOUS AUTHOR(S)

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

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

## 29 1 INTRODUCTION

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

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

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

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

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-composition*, *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], and

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

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

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

110  
111  
112  
113 In this paper, we formally introduce a systematic (de/re)-composition approach for data-parallel  
114 computations targeting state-of-the-art parallel architectures. We express computations via high-  
115 level functional expressions (specifying *what* to compute), in the form of easy-to-use higher-order  
116 functions, inspired by the algebraic formalism for *Multi-Dimensional Homomorphisms (MDHs)* [Rasch  
117 and Gorlatch 2016]<sup>1</sup>. Our higher-order functions are capable of expressing various kinds of data-  
118 parallel computations (linear algebra, stencils, etc), in the same formalism and on a high level of  
119 abstraction, independently of hardware and optimization details, thereby contributing to user's pro-  
120 ductivity<sup>2</sup>. As target for our high-level expressions, we introduce functional low-level expressions  
121 (specifying *how* to compute) to formally reason about the de- and re-compositions of data-parallel  
122 computations; our low-level expressions are designed such that they can be straightforwardly trans-  
123 formed to executable program code (e.g., in OpenMP, CUDA, and OpenCL). To systematically lower  
124 our high-level expressions to low-level expressions, we introduce a formally sound, parameterized  
125 de-composition and re-composition approach. The parameters of our lowering process enable  
126 automatically computing low-level expressions that are optimized (auto-tuned [Balaprakash et al.  
127 2018]) for the particular target architecture and characteristics of the input/output data, thereby  
128 achieving fully automatically high, portable performance. For example, we introduce parameters for  
129 flexibly choosing the target memory regions for de-composed and re-composed computations, and  
130 also parameters for flexibly setting an optimized data access pattern, based on a formal foundation.

131 We show that our high-level representation is capable of expressing various kinds of data-  
132 parallel computations, including computations that recently gained high attention due to their  
133 relevance for deep learning [Barham and Isard 2019]. For our low-level representation, we show  
134 that it can express the cache blocking and parallelization strategies of state-of-the-art parallel  
135 implementations – as generated by scheduling approach TVM and polyhedral compilers PPCG and  
136 Pluto – in one uniform formalism. Moreover, we present experimental results to confirm that via

141  
142 <sup>1</sup> The existing MDH work introduces a proof-of-concept: 1) *high-level program representation* that relies on a semi-formal  
143 foundation, straightforward type system, and a higher-order function for expressing computations (but not for managing  
144 input and output data); 2) *code generator* specifically designed and optimized for the OpenCL programming model (thereby  
145 also being limited to OpenCL). We thoroughly compare to the existing MDH work in Section 6.6.

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

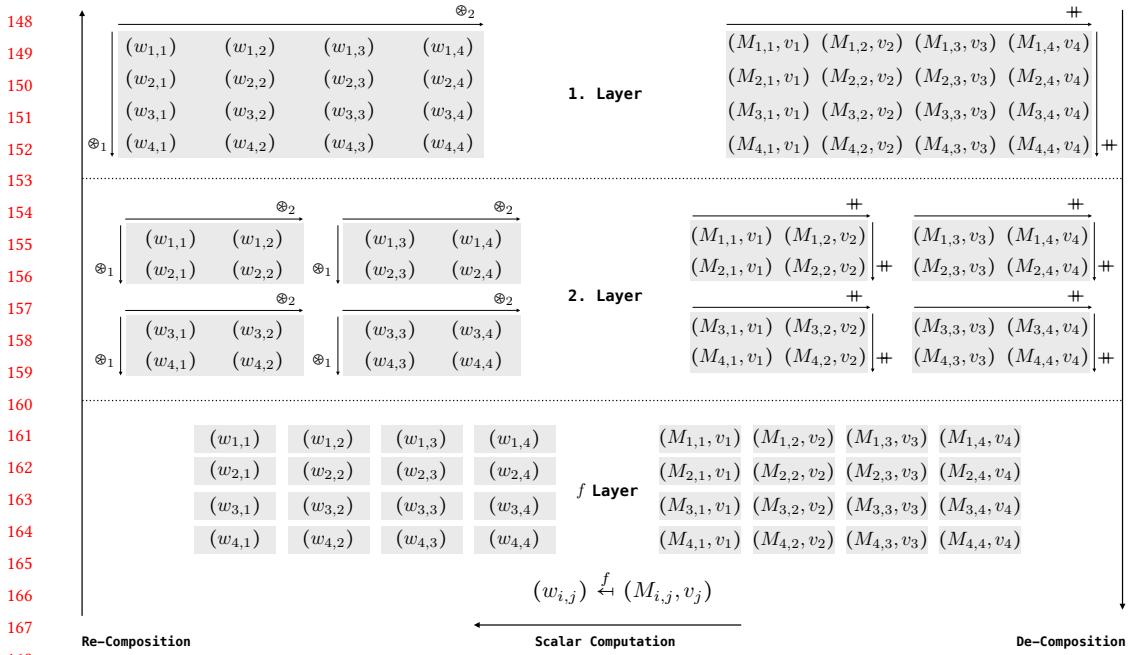


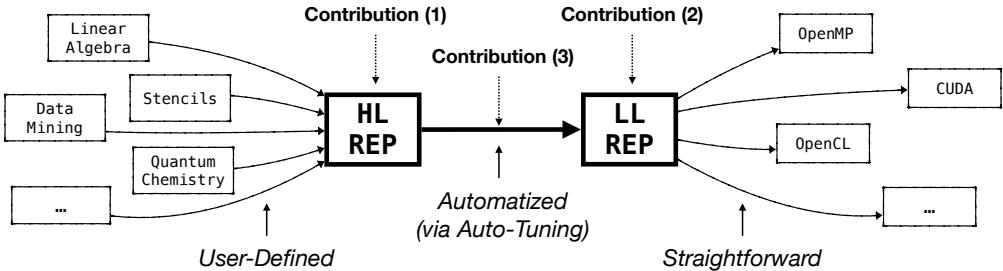
Fig. 3. Example (de/re)-composition of MatVec (Figure 1) on a  $4 \times 4$  input matrix  $M$  and a 4-sized vector  $v$ : i) the *de-composition phase* (right part of the figure) partitions the concatenated input data into parts (a.k.a. *tiles* in programming), where  $+$  denotes the concatenation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

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.

- 197 (3) We introduce a systematic process to fully automatically lower an expression in our high-level  
 198 representation to a device- and data-optimized expression in our low-level representation, in  
 199 a formally sound manner, based on auto-tuning.



210 Fig. 4. Overall structure of our approach (contributions highlighted in bold)  
 211

212 The rest of the paper is structured as follows. We introduce our high-level functional representation  
 213 (Contribution 1) in Section 2, and we show how this representation is used for expressing  
 214 important data-parallel computations. In Section 3, we discuss our low-level functional representation  
 215 (Contribution 2), and we confirm that it can express existing state-of-the-art low-level  
 216 implementations, as generated by state-of-practice approaches TVM (scheduling based) and PPCG/-  
 217 Pluto (polyhedral based). Section 4 shows how we systematically lower a computation expressed  
 218 in our high-level representation to an expression in our low-level representation, in a formally  
 219 sound, auto-tunable manner (Contribution 3). We present experimental results in Section 5, discuss  
 220 related work in Section 6 (including a thorough comparison to previous work on MDHs), conclude  
 221 in Section 7, and we present our ideas for future work in Section 8. Our Appendix, in Sections A-E,  
 222 provides details for the interested reader that should not be required for understanding the basic  
 223 concepts introduced in this paper.

## 224 2 HIGH-LEVEL REPRESENTATION FOR DATA-PARALLEL COMPUTATIONS

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

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

240 In the following, after informally discussing an introductory example in Section 2.1, we formally  
 241 define and discuss each higher-order function in detail in Section 2.2 (function *md\_hom*) and Section  
 242 2.3 (functions *inp\_view* and *out\_view*). Note that Section 2.2 and Section 2.3 introduce and  
 243 present the internals and formal details of our approach, which are not relevant for the end user of  
 244 our system – the user only needs to operate on the abstraction level discussed in Section 2.1.

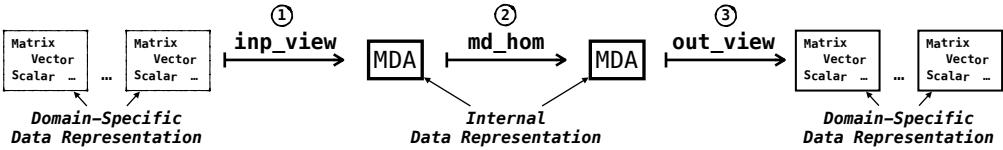


Fig. 5. High-level representation (overview)

## 2.1 Introductory Example

Figure 6 shows how our high-level representation is used for expressing the simple but illustrative example of matrix-vector multiplication MatVec<sup>3</sup> (Figure 1). 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< $T, T$ >( $M: (i, k) \mapsto (i, k)$ ,  $v: (i, k) \mapsto (k)$ )
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294

```

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 *TYPE* and provide details on them in the Appendix, Section A.2, for the interested reader.

<sup>5</sup> We denote by  $\mathbb{N}$  the set of positive natural numbers  $\{1, 2, \dots\}$ , 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)

$$\begin{aligned}
& \mathfrak{a} = \begin{bmatrix} \underbrace{1}_{\mathfrak{a}[0,0]} & \underbrace{2}_{\mathfrak{a}[0,1]} & \underbrace{3}_{\mathfrak{a}[0,2]} & \underbrace{4}_{\mathfrak{a}[0,3]} \\ \underbrace{5}_{\mathfrak{a}[1,0]} & \underbrace{6}_{\mathfrak{a}[1,1]} & \underbrace{7}_{\mathfrak{a}[1,2]} & \underbrace{8}_{\mathfrak{a}[1,3]} \end{bmatrix} \quad \epsilon T[\ I_1 := \{0,1\}, I_2 := \{0,1,2,3\}] \\
& \mathfrak{a}^{(1,1)} = \begin{bmatrix} \underbrace{1}_{\mathfrak{a}[0,0]} & \underbrace{2}_{\mathfrak{a}[0,1]} & \underbrace{3}_{\mathfrak{a}[0,2]} & \underbrace{4}_{\mathfrak{a}[0,3]} \\ \in T[\ I_1^{(1,1)} := \{0\}, I_2^{(1,1)} := \{0,1,2,3\}] \end{bmatrix} \quad \mathfrak{a}^{(1,2)} = \begin{bmatrix} \underbrace{5}_{\mathfrak{a}[1,0]} & \underbrace{6}_{\mathfrak{a}[1,1]} & \underbrace{7}_{\mathfrak{a}[1,2]} & \underbrace{8}_{\mathfrak{a}[1,3]} \\ \in T[\ I_1^{(1,2)} := \{1\}, I_2^{(1,2)} := \{0,1,2,3\}] \end{bmatrix} \\
& \mathfrak{a}^{(2,1)} = \begin{bmatrix} \underbrace{1}_{\mathfrak{a}[0,0]} & \underbrace{2}_{\mathfrak{a}[0,1]} \\ \underbrace{5}_{\mathfrak{a}[1,0]} & \underbrace{6}_{\mathfrak{a}[1,1]} \\ \in T[\ I_1^{(2,1)} := \{0,1\}, I_2^{(2,1)} := \{0,1\}] \end{bmatrix} \quad \mathfrak{a}^{(2,2)} = \begin{bmatrix} \underbrace{3}_{\mathfrak{a}[0,2]} \\ \underbrace{7}_{\mathfrak{a}[1,2]} \\ \in T[\ I_1^{(2,2)} := \{0,1\}, I_2^{(2,2)} := \{2\}] \end{bmatrix} \quad \mathfrak{a}^{(2,3)} = \begin{bmatrix} \underbrace{4}_{\mathfrak{a}[0,3]} \\ \underbrace{8}_{\mathfrak{a}[1,3]} \\ \in T[\ I_1^{(2,3)} := \{0,1\}, I_2^{(2,3)} := \{3\}] \end{bmatrix}
\end{aligned}$$

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

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

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

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

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

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

Figure 7 illustrates six example MDAs. The left part of the figure shows MDA  $\mathfrak{a}$  which is of type  $\mathfrak{a} : I_1 \times I_2 \rightarrow T$ , for  $I_1 = \{0, 1\}$ ,  $I_2 = \{0, 1, 2, 3\}$ , and  $T = \mathbb{Z}$  (integer numbers). On the right side, five MDAs are shown, named  $\mathfrak{a}^{(1,1)}$ ,  $\mathfrak{a}^{(1,2)}$ ,  $\mathfrak{a}^{(2,1)}$ ,  $\mathfrak{a}^{(2,2)}$ ,  $\mathfrak{a}^{(2,3)}$  – the superscripts are part of the names and represent a two-dimensional numbering of the five MDAs. The MDAs  $\mathfrak{a}^{(1,1)}$  and  $\mathfrak{a}^{(1,2)}$  are of types  $\mathfrak{a}^{(1,1)} : I_1^{(1,1)} \times I_2^{(1,1)} \rightarrow T$  and  $\mathfrak{a}^{(1,2)} : I_1^{(1,2)} \times I_2^{(1,2)} \rightarrow T$ , for  $I_1^{(1,1)} = \{0\}$  and  $I_1^{(1,2)} = \{1\}$ , and coinciding second dimensions  $I_2^{(1,1)} = I_2^{(1,2)} = \{0, 1, 2, 3\}$ . The MDAs  $\mathfrak{a}^{(2,1)}$ ,  $\mathfrak{a}^{(2,2)}$ , and  $\mathfrak{a}^{(2,3)}$  are of types  $\mathfrak{a}^{(2,1)} : I_1^{(2,1)} \times I_2^{(2,1)} \rightarrow T$ ,  $\mathfrak{a}^{(2,2)} : I_1^{(2,2)} \times I_2^{(2,2)} \rightarrow T$ , and  $\mathfrak{a}^{(2,3)} : I_1^{(2,3)} \times I_2^{(2,3)} \rightarrow T$ , and they coincide in the first dimension  $I_1^{(2,1)} = I_1^{(2,2)} = I_1^{(2,3)} = \{0, 1\}$ ; their second dimensions are  $I_2^{(2,1)} = \{0, 1\}$ ,  $I_2^{(2,2)} = \{2\}$ , and  $I_2^{(2,3)} = \{3\}$ . Note that MDAs  $\mathfrak{a}^{(1,1)}$ ,  $\mathfrak{a}^{(1,2)}$ ,  $\mathfrak{a}^{(2,1)}$ ,  $\mathfrak{a}^{(2,2)}$ ,  $\mathfrak{a}^{(2,3)}$  can be

<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 considered as *parts* (a.k.a. *tiles* in programming) of MDA  $\alpha$ . We formally define and use *partitionings*  
 345 of MDAs in Section 3.

346

### 347 Combine Operators

348 A central building block in our definition of MDHs is a *combine operator*. Intuitively, we use a  
 349 combine operator to combine all elements within a particular dimension of an MDA. For example,  
 350 in Figure 1 (matrix–vector multiplication), we combine elements of the 2-dimensional MDA via  
 351 combine operator *concatenation* in MDA’s first dimension and via operator *point-wise addition* in  
 352 the second dimension.

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

358

$$\begin{array}{c}
 \begin{array}{ccc}
 1 & 2 & 3 \\
 4 & 5 & 6 \\
 7 & 8 & 9
 \end{array}
 \underbrace{+ \begin{array}{cc}
 11 & 12 \\
 13 & 14 \\
 15 & 16
 \end{array}}_{\substack{3 \times 3\text{-many} \\ \text{elements}}} = \underbrace{\begin{array}{ccccc}
 1 & 2 & 3 & 11 & 12 \\
 4 & 5 & 6 & 13 & 14 \\
 7 & 8 & 9 & 15 & 16
 \end{array}}_{\substack{3 \times 5\text{-many} \\ \text{elements}}} \\
 \text{well defined}
 \end{array}
 \qquad
 \begin{array}{c}
 \begin{array}{ccc}
 1 & 2 & 3 \\
 4 & 5 & 6 \\
 7 & 8 & 9
 \end{array}
 \underbrace{+ \begin{array}{cc}
 11 & 12 \\
 13 & 14
 \end{array}}_{\substack{2 \times 2\text{-many} \\ \text{elements}}} = ?
 \end{array}
 \qquad
 \text{\textit{\$ undefined}}$$

365

366

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

371 Figure 8 illustrates combine operators informally using the example operators *concatenation*  
 372 (left part of the figure) and *point-wise addition* (right part). We illustrate concatenation using the  
 373 example MDAs  $\alpha^{(1,1)}$  and  $\alpha^{(1,2)}$  from Figure 7; for point-wise addition, we use MDAs  $\alpha^{(2,2)}$  and  
 374  $\alpha^{(2,3)}$ .

375

$$\begin{array}{cccc}
 \alpha^{(1,1)} = \underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}}_{\in T[\ P:=\{0\}, I_2:=\{0,1,2,3\} ]} & \alpha^{(1,2)} = \underbrace{\begin{bmatrix} 5 & 6 & 7 & 8 \end{bmatrix}}_{\in T[\ Q:=\{1\}, I_2:=\{0,1,2,3\} ]} & \alpha^{(2,2)} = \underbrace{\begin{bmatrix} 3 \\ 7 \end{bmatrix}}_{\in T[\ I_1:=\{0,1\}, P:=\{2\} ]} & \alpha^{(2,3)} = \underbrace{\begin{bmatrix} 4 \\ 8 \end{bmatrix}}_{\in T[\ I_1:=\{0,1\}, Q:=\{3\} ]} \\
 & \text{index set function} & & \text{index set function} \\
 \alpha_{\downarrow \otimes_1}^{(1,1)} := \underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}}_{\in T[\ id(P=\{0\})=\{0\}, I_2=\{0,1,2,3\} ]} & \alpha_{\downarrow \otimes_1}^{(1,2)} := \underbrace{\begin{bmatrix} 5 & 6 & 7 & 8 \end{bmatrix}}_{\in T[\ id(Q=\{1\})=\{1\}, I_2=\{0,1,2,3\} ]} & \alpha_{\rightarrow \otimes_2}^{(2,2)} := \underbrace{\begin{bmatrix} 3 \\ 7 \end{bmatrix}}_{\in T[\ I_1=\{0,1\}, 0_f(P=\{2\})=\{0\} ]} & \alpha_{\rightarrow \otimes_2}^{(2,3)} := \underbrace{\begin{bmatrix} 4 \\ 8 \end{bmatrix}}_{\in T[\ I_1=\{0,1\}, 0_f(Q=\{3\})=\{0\} ]} \\
 & & & \\
 \xrightarrow{\quad} \alpha^{(1)} := \alpha_{\downarrow \otimes_1}^{(1,1)} \otimes_1 \alpha_{\downarrow \otimes_1}^{(1,2)} = \underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}}_{\in T[\ id(P \cup Q=\{0,1\})=\{0,1\}, I_2=\{0,1,2,3\} ]} & \xrightarrow{\quad} \alpha^{(2)} := \alpha_{\rightarrow \otimes_2}^{(2,2)} \otimes_2 \alpha_{\rightarrow \otimes_2}^{(2,3)} = \underbrace{\begin{bmatrix} 7 \\ 15 \end{bmatrix}}_{\in T[\ I_1=\{0,1\}, 0_f(P \cup Q=\{2,3\})=\{0\} ]} \\
 & \text{concatenation} & & \text{point-wise} \\
 & & & \text{addition}
 \end{array}$$

390

391

392 Fig. 8. Illustration of *combine operators* using the examples *concatenation* (left) and *point-wise addition* (right)

<sup>393</sup>  $\alpha^{(2,3)}$  from Figure 7 (all MDAs are chosen arbitrarily, and the example works the same for other  
<sup>394</sup> MDAs).

<sup>395</sup> In the case of concatenation (left part of Figure 8), MDAs  $\alpha^{(1,1)}$  and  $\alpha^{(1,2)}$  coincide in their second  
<sup>396</sup> dimension  $I_2 := \{0, 1, 2, 3\}$ , which is important, because we concatenate in the first dimension, thus  
<sup>397</sup> requiring coinciding index sets in all other dimensions (as motivated above). In the case of the  
<sup>398</sup> point-wise addition example (right part of Figure 8), the example MDAs  $\alpha^{(2,2)}$  and  $\alpha^{(2,3)}$  coincide  
<sup>399</sup> in their first dimension  $I_1 := \{0, 1\}$ , as required for combining the MDAs in the second dimension.  
<sup>400</sup> The varying index sets of the four MDAs are denoted as  $P$  and  $Q$  in the figure, which are in the case  
<sup>401</sup> of the concatenation example, index sets in the first dimension of MDAs  $\alpha^{(1,1)}$  and  $\alpha^{(1,2)}$ ; in the  
<sup>402</sup> case of the point-wise addition example, the varying index sets of MDAs  $\alpha^{(2,1)}$  and  $\alpha^{(2,2)}$  belong to  
<sup>403</sup> the second dimensions.

<sup>404</sup> In the following, we assume w.l.o.g. that that the varying index sets  $P$  and  $Q$  of MDAs to combine  
<sup>405</sup> are disjoint. Our assumption will not be a limitation for us, because we will apply combine operators  
<sup>406</sup> always to MDA parts that belong to the same MDA, causing the index sets of the parts to be disjoint  
<sup>407</sup> by construction. For example, in the case of the concatenation example in Figure 8, the MDA parts  
<sup>408</sup>  $\alpha^{(1,1)}$  and  $\alpha^{(1,2)}$  correspond to the first and second row of the same MDA  $\alpha$  in Figure 7, and in the  
<sup>409</sup> case of the point-wise addition example, the parts  $\alpha^{(2,1)}$  and  $\alpha^{(2,2)}$  represent the third and fourth  
<sup>410</sup> column of MDA  $\alpha$ .

<sup>411</sup> We define combine operators based on *index set functions* (also introduced formally soon). Index  
<sup>412</sup> set functions precisely describe, on the type level, the index set of the combined output MDA and  
<sup>413</sup> thus how an MDA's index set evolves during combination. For this, an index set function takes  
<sup>414</sup> as input the input MDA's index set in the dimension to combine, and the function yields as its  
<sup>415</sup> output the index set of the output MDA which is combined in this dimension. In the case of the  
<sup>416</sup> concatenation example in Figure 8, the index set function is identity  $id$  and thus trivial. However, in  
<sup>417</sup> the case of point-wise addition, the corresponding index set function is constant function  $0_f$  which  
<sup>418</sup> maps any index set to the singleton set  $\{0\}$  containing index 0 only. This is because when combining  
<sup>419</sup> via point-wise addition, the MDA shrinks in the combined dimension to only one element which we  
<sup>420</sup> aim to uniformly access via MDA index 0. In Figure 8, we denote MDAs  $\alpha^{(1,1)}$ ,  $\alpha^{(1,2)}$ ,  $\alpha^{(2,2)}$ ,  $\alpha^{(2,3)}$   
<sup>421</sup> after applying the corresponding index set function as  $\alpha_{\downarrow \otimes_1}^{(1,1)}$ ,  $\alpha_{\downarrow \otimes_1}^{(1,2)}$ ,  $\alpha_{\rightarrow \otimes_2}^{(2,2)}$ ,  $\alpha_{\rightarrow \otimes_2}^{(2,3)}$ ; the combined  
<sup>422</sup> MDAs are denoted as  $\alpha^{(1)}$  and  $\alpha^{(2)}$  in the figure. The concatenation operator is denoted in the  
<sup>423</sup> figure generically as  $\otimes_1$ , and point-wise addition is denoted as  $\otimes_2$ , correspondingly.

<sup>425</sup> We now define *combine operators* formally, and we illustrate this formal definition afterwards  
<sup>426</sup> using the example operators *concatenation* and *point-wise combination*. For the interested reader,  
<sup>427</sup> details on some technical design decisions of combine operators are outlined in the Appendix,  
<sup>428</sup> Section B.1.

<sup>430</sup> **Definition 2** (Combine Operator). Let  $\text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} := \{(P, Q) \in \text{MDA-IDX-SETS}$   
<sup>431</sup>  $\times \text{MDA-IDX-SETS} \mid P \cap Q = \emptyset\}$  denote the set of all pairs of MDA index sets that are disjoint. Let  
<sup>432</sup> further  $\Rightarrow^{\text{MDA}} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$  be a function on MDA index sets,  $T \in \text{TYPE}$  a scalar  
<sup>433</sup> type,  $D \in \mathbb{N}$  an MDA dimensionality, and  $d \in [1, D]_{\mathbb{N}}$  an MDA dimension.  
<sup>434</sup>

<sup>435</sup> We refer to any binary function  $\otimes$  of type (parameters in angle brackets are type parameters)

<sup>436</sup>  $\otimes : \langle (I_1, \dots, I_{d-1}, I_d, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} \rangle :$

$$\text{437} \quad \text{438} \quad T[I_1, \dots, \underbrace{\Rightarrow^{\text{MDA}}(P), \dots, I_D}_{d}, \times T[I_1, \dots, \underbrace{\Rightarrow^{\text{MDA}}(Q), \dots, I_D}_{d}]$$

$$\rightarrow T[ I_1, \dots, \underbrace{\Rightarrow_{MDA}^{MDA}(P \cup Q), \dots, I_D]}_d ]$$

as *combine operator* that has *index set function*  $\Rightarrow_{MDA}^{MDA}$ , *scalar type*  $T$ , *dimensionality*  $D$ , and *operating dimension*  $d$ . We denote combine operator's type concisely as  $CO^{<\Rightarrow_{MDA}^{MDA}|T|D|d>}$ .

Since function  $\otimes$ 's ordinary function type  $T[\dots] \times T[\dots] \rightarrow T[\dots]$  is generic in parameters  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D)$  and  $(P, Q)$  (these type parameters are denoted in angle brackets in Definition 2), we refer to function  $\otimes$  as *meta-function*, to the type parameters in angle brackets as *meta-parameters*, and we say *meta-types* to  $T[ I_1, \dots, \Rightarrow_{MDA}^{MDA}(P), \dots, I_D ]$  (first input MDA),  $T[ I_1, \dots, \Rightarrow_{MDA}^{MDA}(Q), \dots, I_D ]$  (second input MDA), and  $T[ I_1, \dots, \Rightarrow_{MDA}^{MDA}(P \cup Q), \dots, I_D ]$  (output MDA), as these types are generic in meta-parameters. Formal definitions and details about our meta-parameter concept are provided in Section A of our Appendix for the interested reader.

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

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

**Example 1** (Concatenation). We define *concatenation* as function  $+$  of type

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

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

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

$$+ < T \mid D \mid (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D), (P, Q) > (a_1, a_2)[ i_1, \dots, i_d, \dots, i_D ]$$

$$:= \begin{cases} a_1[ i_1, \dots, i_d, \dots, i_D ] & , i_d \in P \\ a_2[ i_1, \dots, i_d, \dots, i_D ] & , i_d \in Q \end{cases}$$

The function is well defined, because  $P$  and  $Q$  are disjoint. We usually use an infix notation for  $+^{<\dots>}$  (meta-parameters omitted via ellipsis), i.e., we write  $a_1 +^{<\dots>} a_2$  instead of  $+^{<\dots>}(a_1, a_2)$ .

The vertical bar in the superscript of  $+$  denotes that function  $+$  can be partially evaluated (a.k.a. *Currying* [Curry 1980] in math and *multi staging* [Taha and Sheard 1997] in programming) for particular values of meta-parameters:  $T \in \text{TYPE}$  (first stage),  $D \in \mathbb{N}$  (second stage), etc. Partial evaluation (formally defined in the Appendix, Definition 21) enables both: 1) expressive typing and thus better error elimination: for example, parameter  $(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}$

491 can depend on meta-parameter  $D \in \mathbb{N}$ , because  $D$  is defined in an earlier stage, which allows precisely  
 492 limiting 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  
 493 example, we can instantiate  $\text{++}$  to  $\text{++}^{<T>}$  which is specific for a particular scalar type  $T \in \text{TYPE}$ , but  
 494 still generic in meta-parameters  $D \in \mathbb{N}$ ,  $d \in [1, D]_{\mathbb{N}}, \dots$ , as these meta-parameters are defined in  
 495 later stages. We specify stages and their order according to the recommendations in [Haskell Wiki](#)  
 496 [[2013](#)], e.g., using earlier stages for meta-parameters that are expected to change less frequently  
 497 than other meta-parameters.

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

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

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

$$\underbrace{T \times T \rightarrow T}_{\oplus} \rightarrow T[I_1, \dots, \underbrace{0_f(P), \dots, I_D}_{\substack{\uparrow \\ d}}] \times T[I_1, \dots, \underbrace{0_f(Q), \dots, I_D}_{\substack{\uparrow \\ d}}] \rightarrow T[I_1, \dots, \underbrace{0_f(P \cup Q), \dots, I_D}_{\substack{\uparrow \\ d}}]$$

$\underbrace{\hspace{10cm}}_{\text{point-wise combination (according to } \oplus\text{)}}$

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:

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

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 the Appendix, Definition 22) is a combine operator of type  $\text{CO}^{<0_f|T|D|d>}$  for any binary operator  $\oplus : T \times T \rightarrow T$ .

## Multi-Dimensional Homomorphisms

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

Figure 9 illustrates the MDH property informally on a simple, two-dimensional input MDA. In the left part of the figure, we split the input MDA in dimension 1 (i.e., horizontally) into two parts  $a_1$

and  $\alpha_2$ , apply MDH function  $h$  independently to each part, and combine the obtained intermediate results to the final result using the MDH function  $h$ 's combine operator  $\otimes_1$ . Similarly, in the right part of Figure 9, we split the input MDA in dimension 2 (i.e., vertically) into parts and combine the results via MDH function  $h$ 's second combine operator  $\otimes_2$ .

$$h\left(\underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}}_{\mathfrak{a}_1 \leftrightarrow_1 \mathfrak{a}_2}\right) = h\left(\underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}}_{\mathfrak{a}_2} \circledast_1 \underbrace{\begin{bmatrix} 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}}_{\mathfrak{a}_1}\right)$$

$$h\left(\underbrace{\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}}_{\mathfrak{a}_3 \leftrightarrow_2 \mathfrak{a}_4}\right) = h\left(\underbrace{\begin{bmatrix} 1 & 2 \\ 5 & 6 \\ 9 & 10 \\ 13 & 14 \end{bmatrix}}_{\mathfrak{a}_3} \circledast_2 h\left(\underbrace{\begin{bmatrix} 3 & 4 \\ 7 & 8 \\ 11 & 12 \\ 15 & 16 \end{bmatrix}}_{\mathfrak{a}_4}\right)\right)$$

Fig. 9. MDH property illustrated on a two-dimensional example computation

Figure 10 shows an artificial example in which we apply the MDH property (illustrated in Figure 9) recursively. We refer in Figure 10 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). 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  $\text{++}_d := \text{++}^{< T^{\text{INP}} | D | d >} \in \text{CO}^{< id | T^{\text{INP}} | D | d >}$  be concatenation (Definition 1) in dimension  $d \in [1, D]_{\mathbb{N}}$  on  $D$ -dimensional MDAs that have scalar type  $T^{\text{INP}}$ .

## A function

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

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*  $\overset{1}{\Rightarrow}_{\text{MDA}}, \dots, \overset{D}{\Rightarrow}_{\text{MDA}}$ , iff for each  $d \in [1, D]_{\mathbb{N}}$ , there exists a combine operator  $\otimes_d \in \text{CO}^{<\overset{d}{\Rightarrow}_{\text{MDA}} | T^{\text{OUT}} | D | d >}$  (Definition 2), such that for any concatenated input MDA  $\alpha_1 \text{++}_d \alpha_2$  in dimension  $d$ , the *homomorphic property* is satisfied:

$$h(\alpha_1 \oplus_d \alpha_2) = h(\alpha_1) \otimes_d h(\alpha_2)$$

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

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. Note that per definition of MDHs, their combine operators are associative and commutative (which follows from the associativity and commutativity of  $+$ ). Note further that for simplicity, Definition 3 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.2 of our Appendix, for the algebraically interested reader.

**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  $\text{map}(f)$

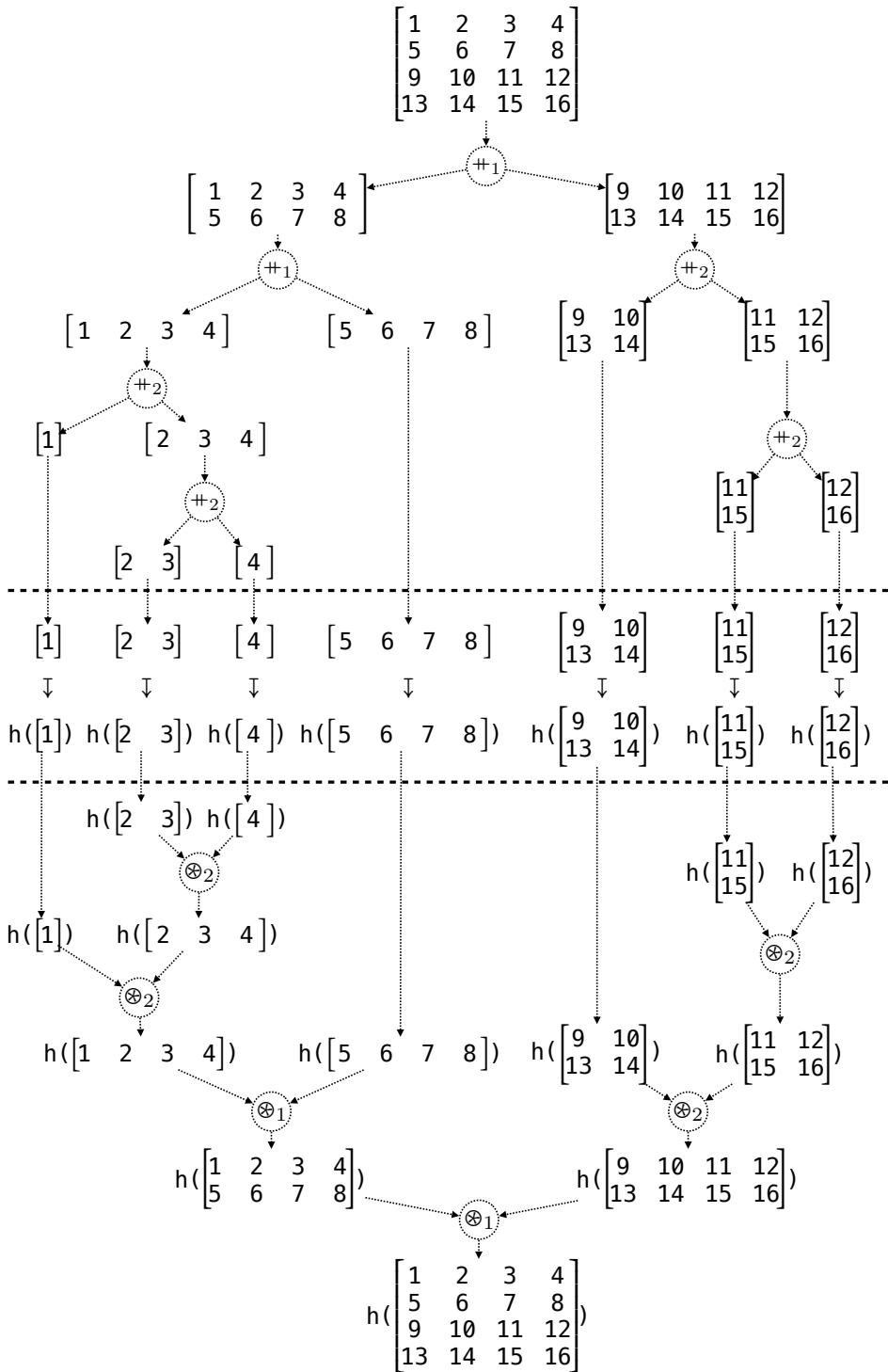


Fig. 10. MDH property recursively applied to a two-dimensional example computation

638 is an MDH of type  $\text{MDH}^{<T^{\text{INP}}, T^{\text{OUT}} | D | id, \dots, id>}$  whose combine operators are concatenation  $+$  in all  
 639 of its  $D$  dimensions (Example 1). Function  $id$  is the index set function of  $+$  (see Example 1) and  
 640 consequently also of MDH  $\text{map}(f)$ . Formal details and definitions for *function mapping* can be  
 641 found in the Appendix, Section B.3.

642 **Example 4** (Reduction). Another MDH function is *reduction* [González-Vélez and Leyton 2010],  
 643 implemented as higher-order function  $\text{red}(\oplus)(\alpha)$ , which combines all elements within a  $D$ -  
 644 dimensional MDA  $\alpha$  using a user-defined binary function  $\oplus : T \times T \rightarrow T$ . Reduction's MDH  
 645 type is  $\text{MDH}^{<T, T | D | 0_f, \dots, 0_f>}$ , and its combine operators are point-wise combination  $\rightarrow$  ( $\oplus$ ) in all  
 646 dimensions (Example 2) which have  $0_f$  index set function. Formal details and definitions for *reduction*  
 647 can be found in the Appendix, Section B.3.

648 We show how Examples 3 and 4 (and particularly also more advanced examples) are expressed  
 649 in our high-level representation in Section 2.5, based on higher-order functions `md_hom`, `inp_view`,  
 650 and `out_view` (Figure 5) which we introduce in the following.

### 653 Higher-Order Function `md_hom`

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

$$658 \quad \begin{aligned} 659 \quad h\left(\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}\right) &= \overbrace{\begin{bmatrix} h([1]) & h([2]) & h([3]) & h([4]) \\ h([5]) & h([6]) & h([7]) & h([8]) \\ h([9]) & h([10]) & h([11]) & h([12]) \\ h([13]) & h([14]) & h([15]) & h([16]) \end{bmatrix}}^{\oplus_2} \\ 660 &= \overbrace{\begin{bmatrix} f(1) & f(2) & f(3) & f(4) \\ f(5) & f(6) & f(7) & f(8) \\ f(9) & f(10) & f(11) & f(12) \\ f(13) & f(14) & f(15) & f(16) \end{bmatrix}}^{\oplus_1} \end{aligned}$$

661 Here,  $f$  is the function on scalar values that behaves the same as  $h$  when restricted to singleton  
 662 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  
 663 accessed by (arbitrary) indices  $i_1, \dots, i_D \in \mathbb{N}_0$ . For singleton MDAs, we usually use  $f$  instead of  $h$ ,  
 664 because  $f$  can be defined more conveniently by the user as  $h$  (which needs to handle MDAs of  
 665 arbitrary sizes, and not only singleton MDAs as  $f$ ). Also, since  $f$  takes as input a scalar value (rather  
 666 than a singleton MDA as  $h$ ), the type of  $f$  also becomes simpler, thereby also further simplifying  
 667 the definition of scalar functions for the user.

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

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

$$671 \quad \text{md\_hom}^{<T^{\text{INP}}, T^{\text{OUT}} \in \text{TYPE} \mid D \in \mathbb{N} \mid (\xrightarrow{D_{\text{MDA}}}^{\text{MDA}}_{\text{MDA}} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS})_{d \in [1, D] \in \mathbb{N}}>} :$$

$$672 \quad \underbrace{\text{SF}^{<T^{\text{INP}}, T^{\text{OUT}}>}}_f \times \underbrace{(\text{CO}^{<\xrightarrow{1_{\text{MDA}}}^{\text{MDA}}_{\text{MDA}} | T^{\text{OUT}} | D | 1>} \times \dots \times \text{CO}^{<\xrightarrow{D_{\text{MDA}}}^{\text{MDA}} | T^{\text{OUT}} | D | D>})}_{\oplus_1, \dots, \oplus_D} \rightarrow_p \underbrace{\text{MDH}^{<T^{\text{INP}}, T^{\text{OUT}} | D | (\xrightarrow{D_{\text{MDA}}}^{\text{MDA}}_{\text{MDA}})_{d \in [1, D] \in \mathbb{N}}>}}_{\text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))}$$

673 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  
 674 by  $\rightarrow_p$  instead of  $\rightarrow$ ), which we motivate after this definition. The function takes as input a

687 scalar function  $f$  and a tuple of  $D$ -many combine operators  $(\otimes_1, \dots, \otimes_D)$ , and it yields a function  
 688  $\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))$  which is defined as

$$689 \quad \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\alpha) := \bigotimes_{i_1 \in I_1} \dots \bigotimes_{i_D \in I_D} \vec{f}(\alpha|_{\{i_1\} \times \dots \times \{i_D\}})$$

691 The combine operators' underset notation denotes straightforward iteration (explained formally in  
 692 the Appendix, Notation 5), and the MDA  $\alpha|_{\{i_1\} \times \dots \times \{i_D\}}$  is the restriction of  $\alpha$  to the MDA containing  
 693 the single element accessed via MDA indices  $(i_1, \dots, i_D)$ . Function  $\vec{f}$  behaves like scalar function  
 694  $f$ , but operates on singleton MDAs (rather than scalars): the function is of type  
 695

$$696 \quad \vec{f}^{< i_1, \dots, i_D \in \mathbb{N}_0 >} : T^{\text{INP}}[\{i_1\}, \dots, \{i_D\}] \rightarrow T^{\text{OUT}}[\Rightarrow_{\text{MDA}}^1(\{i_1\}), \dots, \Rightarrow_{\text{MDA}}^D(\{i_D\})]$$

697 and defined as

$$699 \quad \vec{f}(x)[j_1, \dots, j_D] := f(x[i_1, \dots, i_D])^8$$

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

$$703 \quad \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\alpha_1 +_d \alpha_2) =$$

$$704 \quad \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\alpha_1) \otimes_d \text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\alpha_2)$$

706 Using Definition 4, we express any MDH function uniformly via higher-order function  $\text{md\_hom}$   
 707 using only the MDH's behavior  $f$  on scalar values<sup>9</sup> and its combine operators  $\otimes_1, \dots, \otimes_D$ . The  
 708 other direction also holds: each function expressed via  $\text{md\_hom}$  is an MDH function, because we  
 709 require the homomorphic property for  $\text{md\_hom}$ .

710 Note that we can potentially allow in Definition 4 the case  $D = 0$  in which we would define the  
 711  $\text{md\_hom}$  instance equal to the scalar function  $f$ :

$$712 \quad \text{md\_hom}(f, ()) := f$$

714 Note further that function  $\text{md\_hom}$  is defined as partial function, because the homomorphic  
 715 property is not met for all potential combinations of combine operators, e.g.,  $\otimes_1 = +$  (point-  
 716 wise addition) and  $\otimes_2 = *$  (point-wise multiplication). However, in many real-world examples, an  
 717 MDH's combine operators are a mix of concatenations and point-wise combinations according to  
 718 the same binary function. The following lemma proves that any instance of the  $\text{md\_hom}$  higher-order  
 719 function for such a mix of combine operators is a well-defined MDH function.

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

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

725 PROOF. See Appendix, Section B.5. □

726 MDHs use as input and output uniformly an MDA. We introduce higher-order function  $\text{inp\_view}$   
 727 to prepare domain-specific inputs (e.g., a matrix and a vector for matrix-vector multiplication) as  
 728 an MDA, and we use function  $\text{out\_view}$  to transform the output MDA back to the domain-specific  
 729 data requirements (like storing it as a transposed matrix in the case of matrix multiplication, or

731 <sup>8</sup> We assume in this work that MDH functions when used on singleton MDAs return a singleton MDA, as such MDHs  
 732 already cover all real-world cases we currently are aware of.

733 <sup>9</sup> For simplicity, we ignore that the scalar functions of some MDHs (such as Mandelbrot) also take as input MDA indices,  
 734 which requires slight, straightforward extension of function  $\text{md\_hom}$ , as outlined in the Appendix, Section B.4.

736 splitting it into multiple outputs as we will see later with examples). We introduce both higher-order  
 737 functions in the following.

### 738 2.3 View Functions

740 We start, in Section 2.3.1, by formally introducing *Buffers (BUF)* and *Index Functions* – both concepts  
 741 are central building blocks in our definition of higher-order functions `inp_view` and `out_view`. In  
 742 our approach, we use BUFs to represent domain-specific input and output data (scalars, vectors,  
 743 matrices, etc), and *index functions* are used by the user to conveniently instantiate higher-order  
 744 functions `inp_view` and `out_view` (e.g., index function  $(i, k) \mapsto (i, k)$  and  $(i, k) \mapsto (k)$  used in  
 745 Figure 6 to instantiate `inp_view`, and  $(i, k) \mapsto (i)$  used for `out_view`).

746 Sections 2.3.2 and 2.3.3 introduce *input views* and *output views* which are central concepts in our  
 747 approach. We define *input views* as arbitrary functions that map a collection of BUFs to an MDA  
 748 (Figure 5); higher-order function `inp_view` is then defined to conveniently compute an important  
 749 class of input view functions relevant in real-world applications. Correspondingly, Section 2.3.3  
 750 defines *output views* as functions that transform an MDA to a collection of BUFs, and higher-order  
 751 function `out_view` is defined to conveniently compute important output views.

752 Finally, we discuss in Section 2.3.4 the relationship between higher-order function `inp_view` and  
 753 `out_view`: we prove that both functions are inversely related to each other, allowing arbitrarily  
 754 switching between our internal MDA and the domain-specific BUF representation (as required for  
 755 our code generation discussed later).

756 2.3.1 *Preparation.* We formally introduce *Buffers (BUF)* and *Index Functions* in the following.

758 **Definition 5** (Buffer). Let  $T \in \text{TYPE}$  be an arbitrary scalar type,  $D \in \mathbb{N}_0$  a natural number<sup>10</sup>, and  
 759  $N := (N_1, \dots, N_D) \in \mathbb{N}^D$  a tuple of natural numbers.

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

$$762 \quad b : [0, N_1)_{\mathbb{N}_0} \times \dots \times [0, N_D)_{\mathbb{N}_0} \rightarrow T \cup \{\perp\}$$

764 Here, we use  $\perp$  to denote the *undefined value*. We refer to  $[0, N_1)_{\mathbb{N}_0} \times \dots \times [0, N_D)_{\mathbb{N}_0} \rightarrow T \cup \{\perp\}$   
 765 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  $\text{BUF-IDX-SETS} :=$   
 766  $\{[0, N)_{\mathbb{N}_0} \mid N \in \mathbb{N}\}$  as *BUF index sets*. Analogously to Notation 1, we write  $b[i_1, \dots, i_D]$  instead of  
 767  $b(i_1, \dots, i_D)$  to avoid a too heavy usage of parentheses.

768 In contrast to MDAs, a BUF always operates on a contiguous range of natural numbers starting  
 769 from 0, and a BUF may contain undefined values. These two differences allow us straightforwardly  
 770 transforming BUFs to data structures provided by low-level programming languages (e.g., *C arrays*  
 771 as used in OpenMP, CUDA, and OpenCL). Note that in our generated program code (discussed later  
 772 in Section 3), we implement MDAs as straightforward aliases that access BUFs, on top of BUFs, so  
 773 that we do not need to transform MDAs to low-level data structures and/or store them otherwise  
 774 physically in memory.

776 **Definition 6** (Index Function). Let  $D \in \mathbb{N}$  be a natural number (representing an MDA’s dimension-  
 777 ality) and  $D_b \in \mathbb{N}_0$  (representing a BUF’s dimensionality).

778 An *index function*  $\text{id}\mathbf{x}$  from  $D$ -dimensional MDA indices to  $D_b$ -dimensional BUF indices is any  
 779 meta-function of type

$$780 \quad \text{id}\mathbf{x}^{<I_1^{\text{MDA}}, \dots, I_D^{\text{MDA}} \in \text{MDA-IDX-SETS}^D>} : I_1^{\text{MDA}} \times \dots \times I_D^{\text{MDA}} \rightarrow I_1^{\text{BUF}} \times \dots \times I_{D_b}^{\text{BUF}}$$

783 <sup>10</sup> We use the case  $D = 0$  to represent scalar values, as we discuss later.

785 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  
 786 arbitrary but fixed function that maps  $D$ -many MDA index sets to  $D_b$ -many BUF index sets. We  
 787 denote the type of index functions as  $\text{MDA-IDX-to-BUF-IDX}^{<D, D_b} \mid \Rightarrow_{\text{BUF}}^{\text{MDA}}$ .

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

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

$$794 \text{MDA-IDX-to-BUF-IDX}^{<D:=2, D_b:=2} \mid \Rightarrow_{\text{BUF}}^{\text{MDA}}(I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_1^{\text{MDA}})]_{\mathbb{N}_0}, [0, \max(I_2^{\text{MDA}})]_{\mathbb{N}_0} >$$

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

$$797 \text{MDA-IDX-to-BUF-IDX}^{<D:=2, D_b:=1} \mid \Rightarrow_{\text{BUF}}^{\text{MDA}}(I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_2^{\text{MDA}})]_{\mathbb{N}_0} >$$

798 Further examples of index function, e.g., for stencil computation Jacobi1D, are presented in the  
 799 Appendix, Section B.6, for the interested reader.

800 2.3.2 *Input Views*. We define *input views* as any function that compute an MDA from a collection  
 801 of user-defined BUFs. For example, in the case of MatVec, its input view takes as input two BUFs – a  
 802 matrix and a vector – and it yields a two-dimensional MDA containing pairs of elements taken from  
 803 the matrix and vector (illustrated in Figure 1). In contrast, the input view of Jacobi1D takes as input  
 804 a single BUF (vector) only, and it computes an MDA containing triples of BUF elements (Figure 2).

805 **Definition 7** (Input View). An *input view* from  $B$ -many BUFs,  $B \in \mathbb{N}$ , of arbitrary but fixed types  
 806  $T_b^{N_1^b \times \dots \times N_{D_b}^b}$ ,  $b \in [1, B]_{\mathbb{N}}$ , to an MDA of arbitrary but fixed type  $T[I_1, \dots, I_D]$  is any function  $\text{iv}$  of  
 807 type:

$$808 \text{iv} : \underbrace{\bigtimes_{b=1}^B T_b^{N_1^b \times \dots \times N_{D_b}^b}}_{\text{BUFs}} \xrightarrow{p} \underbrace{T[I_1, \dots, I_D]}_{\text{MDA}}$$

809  $\overbrace{\text{BUFs' Meta-Parameters}}^{\text{BUFs' Meta-Parameters}} \overbrace{\text{MDA's Meta-Parameters}}^{\text{MDA's Meta-Parameters}}$

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

811 **Example 5** (Input View – MatVec). The input view of MatVec on a  $1024 \times 512$  matrix and 512-sized  
 812 vector (sizes are chosen arbitrarily), both of integers  $\mathbb{Z}$ , is of type

$$813 \text{IV}^{<B=2 \mid D_1=2, D_2=1 \mid (N_1^1=1024, N_2^1=512), (N_1^2=512) \mid T_1=\mathbb{Z}, T_2=\mathbb{Z} \mid D=2 \mid I_1=[0, 1024]_{\mathbb{N}_0}, I_2=[0, 512]_{\mathbb{N}_0} \mid T=\mathbb{Z} \times \mathbb{Z} >}$$

814 and defined as

$$815 \underbrace{[M(i, k)]_{i \in [0, 1024]_{\mathbb{N}_0}, k \in [0, 512]_{\mathbb{N}_0}}}_{\text{Matrix}} \underbrace{[v(k)]_{k \in [0, 512]_{\mathbb{N}_0}}}_{\text{Vector}} \mapsto \underbrace{[M(i, k), v(k)]_{i \in [0, 1024]_{\mathbb{N}_0}, k \in [0, 512]_{\mathbb{N}_0}}}_{\text{MDA}}$$

816 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;  $(N_1^1, N_2^1) = (1024, 512)$  is the matrix size and  $N_1^2 = 512$  the vector's size;  $T_1, T_2 = \mathbb{Z}$  are the scalar types of matrix and vector. The MDA's meta-parameters are:  $D = 2$  is the computed MDA's dimensionality;  $I_1, I_2$  are the MDA's index sets; parameter  $T = \mathbb{Z} \times \mathbb{Z}$  is MDA's scalar type (pairs of matrix/vector elements – see Figure 1).

834 **Example 6** (Input View – Jacobi1D). The input view of Jacobi1D on a 512-sized vector of integers  
 835 is of type

836  $\overbrace{\text{IV}^{<B=1 \mid D_1=1 \mid (N_1^1=512) \mid T_1=\mathbb{Z} \mid D=1 \mid I_1=[0,512-2)_{\mathbb{N}_0} \mid T=\mathbb{Z} \times \mathbb{Z} \times \mathbb{Z}>}}^{\substack{\text{BUFs' meta-parameters} \\ \text{MDA's meta-parameters}}}$

837

838

839 and defined as

840

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

842

843

844 We introduce higher-order function `inp_view` which computes important input views conveniently and in a structured manner from user-defined index functions  $(\text{idx}_{b,a})_{b \in [1,B]_{\mathbb{N}}, a \in [1,A_b]_{\mathbb{N}}}$  (Definition 6). Here,  $B \in \mathbb{N}$  represents the number of BUFs that the computed input view will take as input, and  $A_b$  represents the number of accesses to the  $b$ -th BUF required for computing an individual MDA element.

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

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

855 More generally, higher-order function `inp_view` uses the index functions  $\text{idx}_{b,a}$  to compute an  
 856 input view that maps BUFs  $b_1, \dots, b_B$  to an MDA  $a$  that contains at position  $i_1, \dots, i_D$  the following  
 857 element:

858  $a[i_1, \dots, i_D] := \left( \left( \underbrace{\text{b}_1[\text{idx}_{1,1}(i_1, \dots, i_D)] \in T_1}_{a=1}, \dots, \underbrace{\text{b}_1[\text{idx}_{1,A_1}(i_1, \dots, i_D)] \in T_1}_{a=A_1} \right) \right.$

859  $\left. \underbrace{\vdots}_{b=1} \right)$

860  $\left( \left( \underbrace{\text{b}_B[\text{idx}_{B,1}(i_1, \dots, i_D)] \in T_B}_{a=1}, \dots, \underbrace{\text{b}_B[\text{idx}_{B,A_B}(i_1, \dots, i_D)] \in T_B}_{a=A_B} \right) \right)$

861  $\left. \underbrace{\vdots}_{b=B} \right)$

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

864  $a[i_1, i_2] := \left( \left( \underbrace{\text{b}_1[\text{idx}_{1,1}(i_1, i_2)] := (i_1, i_2)}_{a=1} \in T_1 \right), \left( \underbrace{\text{b}_2[\text{idx}_{2,1}(i_1, i_2)] := (i_2)}_{a=1} \in T_2 \right) \right)$

865  $\left. \underbrace{\vdots}_{b=1} \right)$

866  $\left. \underbrace{\vdots}_{b=2} \right)$

883 and for Jacobi1D, the element is

884  $a[i_1] :=$

885  
 886 
$$(( \underbrace{b_1[\text{id}x_{1,1}(i_1) := (i_1 + 0)] \in T_1}_{a=1}, \underbrace{b_1[\text{id}x_{1,2}(i_1) := (i_1 + 1)] \in T_1}_{a=2}, \underbrace{b_1[\text{id}x_{1,3}(i_1) := (i_1 + 2)] \in T_1}_{a=3} ) )$$
  
 887  
 888 
$$\underbrace{\quad\quad\quad}_{b=1}$$
  
 889

890 In the following, we introduce higher-order function `inp_view` which computes important  
 891 input views conveniently and in a uniform, structured manner. Function `inp_view` takes as input  
 892 a collection of index functions (Definition 6), and it uses these index functions to compute a  
 893 corresponding input view (Definition 7), as follows.

894 Figures 11 and 12 use the examples `MatVec` and `Jacobi1D` to informally illustrate how function  
 895 `inp_view` uses index functions to compute input views. In the two figures, we use domain-specific  
 896 identifiers for better clarity: in the case of `MatVec`, we use for `BUFs` identifiers  $M$  and  $v$  instead of  
 897  $b_1$  and  $b_2$ , as well as identifiers  $i$  and  $j$  instead of  $i_1$  and  $i_2$  for index variables; for `Jacobi1D`, we use  
 898 identifier  $v$  instead of  $b_1$ , and  $i$  instead of  $i_1$ .

899 We now formally define higher-order function `inp_view`. For high flexibility and formal correctness,  
 900 function `inp_view` relies on a type that involves many meta-parameters. The high number of  
 901 meta-parameters, and the resulting complex type of function `inp_view`, might appear daunting to  
 902 the user. However, Notation 2 confirms that despite the complex type of function `inp_view`, the  
 903 function can be conveniently used by the user (as also illustrated in Figure 6), because meta-parameters  
 904 can be automatically deduced from `inp_view`'s input parameters.

905 **Definition 8** (Higher-Order Function `inp_view`). Function `inp_view` is of type

906  
 907 
$$\text{inp\_view}^c \quad \underbrace{B \in \mathbb{N}}_{\text{Number BUFS}} \quad | \quad \underbrace{A_1, \dots, A_B \in \mathbb{N}}_{\text{Number BUFS' Accesses}} \quad | \quad \underbrace{D_1, \dots, D_B \in \mathbb{N}_0}_{\text{BUFS' Dimensionalities}} \quad | \quad \underbrace{D \in \mathbb{N}}_{\text{MDA Dimensionality}}$$
  
 908  
 909  
 910  
 911 
$$(\Rightarrow_{\text{BUF}}^{\text{MDA } b, a : \text{MDA-IDX-SETS}^D \rightarrow \text{BUF-IDX-SETS}^{D_B}})_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_B]_{\mathbb{N}}} > :$$
  
 912  
 913 
$$\underbrace{\quad\quad\quad}_{\text{Index Set Functions (MDA indices to BUF indices)}}$$
  
 914  
 915 
$$\underbrace{\begin{array}{c} B \\ \times \\ b=1 \end{array} \quad \underbrace{\begin{array}{c} A_B \\ \times \\ a=1 \end{array}}_{\text{Index Function: } \text{id}x_{b, a}}}_{\text{Buffer Access}} \quad \underbrace{\text{MDA-IDX-to-BUF-IDX}^{<D, D_B | \Rightarrow_{\text{BUF}}^{\text{MDA } b, a}}}_{\text{Index Function: } \text{id}x_{b, a}} \rightarrow \text{IV}^{<B | D_1, \dots, D_B | \rightarrow | T_1, \dots, T_B \in \text{TYPE} |}_{\text{Input View: } \text{iv}} \quad \underbrace{\begin{array}{c} D | I_1, \dots, I_D \in \text{MDA-IDX-SETS} | \rightarrow > < N_1^1, \dots, N_{D_B}^B | T > \\ \text{MDA's Meta-Parameters} \quad \underbrace{\quad\quad\quad}_{\text{Postponed Parameters}} \end{array}}_{\text{Input View: } \text{iv}}$$
  
 916  
 917  
 918  
 919  
 920

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

922  
 923  
 924 
$$\underbrace{(\text{id}x_{b, a})_{b \in [1, B]_{\mathbb{N}}, a \in [1, A_B]_{\mathbb{N}}}}_{\text{Index Functions}} \mapsto \underbrace{(\underbrace{b_1, \dots, b_B}_{\text{BUFs}}) \xrightarrow{\text{id}v} \underbrace{a}_{\text{MDA}}}_{\text{Input View}}$$
  
 925  
 926  
 927  
 928 for

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

930

931

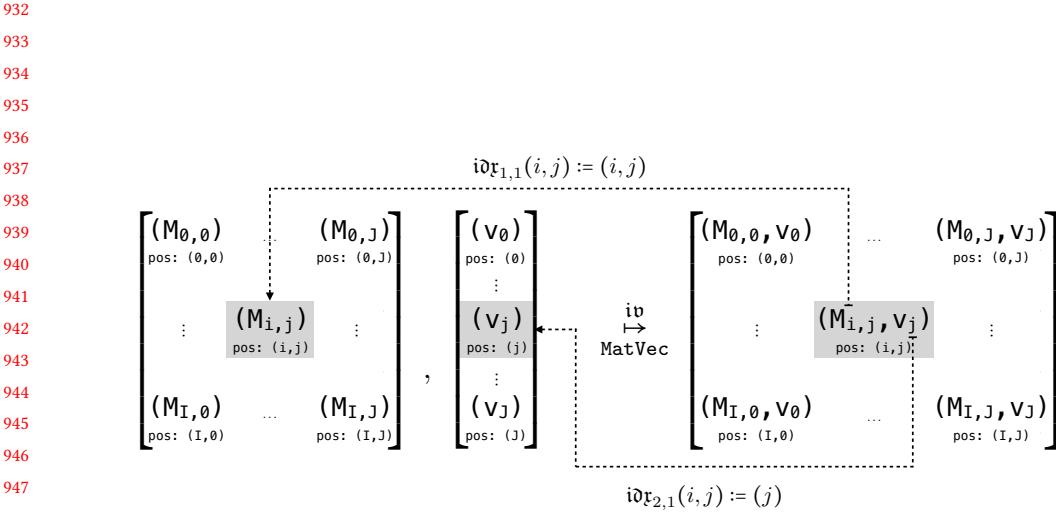


Fig. 11. Input view illustrated using the example MatVec

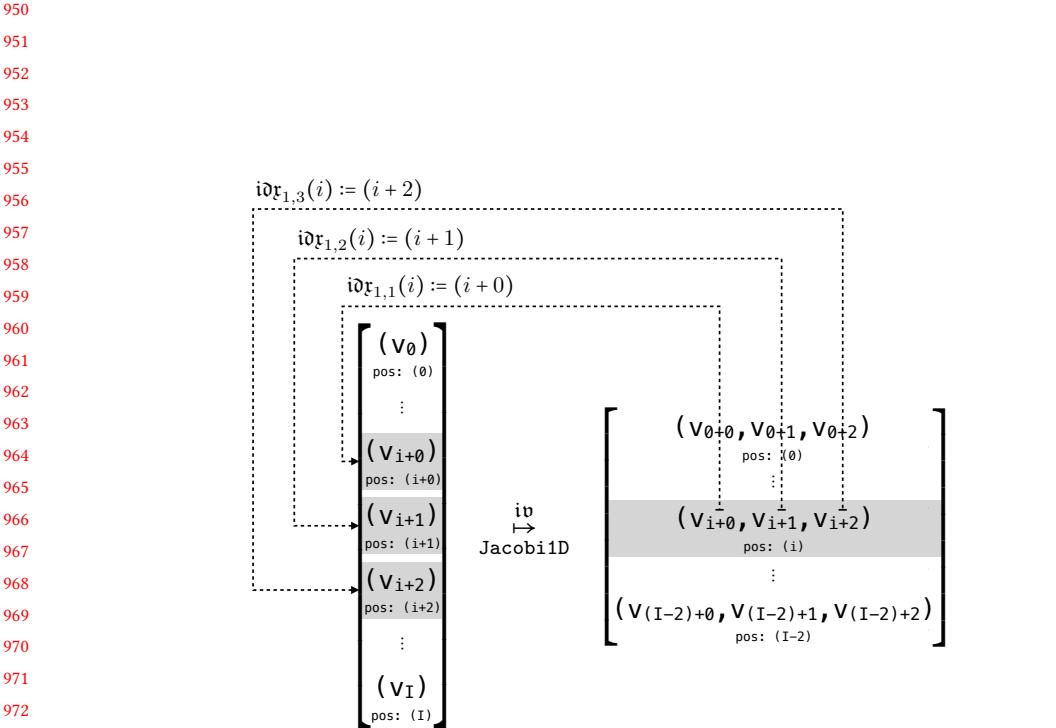


Fig. 12. Input view illustrated using the example Jacobi1D

981 and

$$982 \quad 983 \quad \mathbf{a}_{b,a}[i_1, \dots, i_D] := \mathbf{b}_b[\mathbf{idx}_{b,a}(i_1, \dots, i_D)]$$

984 Higher-order function `inp_view` takes as input a collection of index functions that are of types  
 985 MDA-IDX-to-BUF-IDX (Definition 6), and it computes an input view of type IV (Definition 7) based  
 986 on the index functions, as illustrated in Figures 11 and 12.

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

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

999 We automatically compute the sizes  $N_d^b$  of BUFs in IV's meta-parameter list (e.g., in the case of  
 1000 `MatVec`, the size of the input matrix  $(I, J)$  and vector size  $J$ ), according to the formula in Definition 8,  
 1001 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  
 1002 sizes from MDA index sets (rather than requesting the sizes explicitly from the user), we achieve  
 1003 strong error checking: for example, for `MatVec`, we can ensure – already on the type level – that  
 1004 the number of columns of its input matrix and the size of its input vector match. To compute the  
 1005 BUF sizes, we *postpone* via  $\rightarrow$  (defined formally in the Appendix, Definition 24) the sizes in IV's  
 1006 meta-parameter list to later meta-parameter stages; this is because the sizes are defined in early  
 1007 stages and thus have no access to the MDA's index sets which are defined in later stages. Our  
 1008 formula in Definition 8 then works as follows: for each BUF  $b$ , its size  $N_d^b$  has to be well-defined in  
 1009 each of its dimensions  $d$ , for all accesses  $a$  via the BUF's index functions when used for all indices  
 1010  $I_1, \dots, I_D$  within the MDA index sets. Here, in the computation of  $N_d^b$ , function  $\Rightarrow_{\text{BUF}}^{MDAb,a}$  computes  
 1011 the  $d$ -th component of the  $D_b$ -sized output tuple of  $\Rightarrow_{\text{BUF}}^{MDAb,a}$  (the computed component is the index  
 1012 set of BUF  $b$  in dimension  $d$  for the  $a$ -th index function used to access the BUF).

1013 We automatically compute also MDA's scalar type  $T$  using the formula presented in Definition 8.  
 1014 The formula computes  $T$  as a tuple that consists of the BUFs' scalar types, as each MDA element  
 1015 consist of BUF elements (illustrated in Figures 11 and 12). Postponing  $T$  in IV's meta-parameter  
 1016 list is done (analogously as for  $N_d^b$ ), but is actually not required, because the BUFs' scalar types  
 1017  $T_1, \dots, T_B$  are already defined in earlier meta-parameter stages than  $T$ . However, we will see that  
 1018 postponing  $T$  is required later in the Definition 10 of higher-order function `out_view`; therefore,  
 1019 we postpone  $T$  also in our definition of `inp_view` to increase consistency between our definitions  
 1020 of `inp_view` (Definition 8) and `out_view` (Definition 10).

1021 Note that function `inp_view` is not capable of computing every kind of input view function  
 1022 (Definition 7). For example, `inp_view` cannot be used for computing MDAs that are required for  
 1023 computations on sparse data formats [Hall 2020], because such MDAs need dynamically accessing  
 1024 BUFs. This limitation of `inp_view` can be relaxed by generalizing our index functions toward  
 1025 taking additional, dynamic input arguments, which we consider as future work (as also outlined in  
 1026 Section 8).

**Notation 2** (Input Views). Let  $\text{inp\_view}^{<\dots>}((\text{id}_{x,b,a})_{b \in [1, B]_N, a \in [1, A_b]_N})$  be a particular instance of higher-order function  $\text{inp\_view}$  (meta-parameters omitted via ellipsis for simplicity) for an arbitrary but fixed choice of index functions. Let further  $\text{ID}_1, \dots, \text{ID}_B \in \Sigma^*$  be arbitrary, user-defined BUF identifiers (e.g.,  $\text{ID}_1 = "M"$  and  $\text{ID}_2 = "v"$  in the case of MatVec), for an arbitrary, fixed collection of letters  $\Sigma = \{A, B, C, \dots, a, b, c, \dots, 1, 2, 3, \dots\}$ .

1035 For better readability, we use the following notation for the 2-dimensional structure of index  
1036 functions taken as input by function `inp_view`, inspired by [Lattner et al. \[2021\]](#):

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

1039 We refrain from stating `inp_view`'s meta-parameters in our notation, as the parameters can be  
1040 automatically deduced from the number and types of index functions.

**1041 Example 7.** Function `inp_view` is used for `MatVec` and `Jacobi1D` (in Notation 2) as follows:

<u>MatVec:</u>	$\text{inp\_view}(\underbrace{M: (i, k) \mapsto (i, k)}_{\substack{a=1 \\ b=1}}, \underbrace{v: (i, k) \mapsto (k)}_{\substack{a=1 \\ b=2}})$
<u>Jacobi1D:</u>	$\text{inp\_view}(\underbrace{v: (i) \mapsto (i + 0)}_{\substack{a=1 \\ h=1}}, \underbrace{(i) \mapsto (i + 1)}_{a=2}, \underbrace{(i) \mapsto (i + 2)}_{a=3})$

2.3.3 *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.

The diagram illustrates the computation of the element  $(j, i)$  in the result matrix. It shows two parallel operations: one for the row  $C_0$  and one for the column  $C_I$ . The row operation involves multiplying the row vector  $C_{0,0,0}$  by the column vector  $C_{0,j,0}$ . The column operation involves multiplying the row vector  $C_{I,0,0}$  by the column vector  $C_{I,j,0}$ . The result of the row multiplication is  $(j, i)$  and the result of the column multiplication is  $(i, j)$ .

Fig. 13. Output view illustrated using the example *transposed Matrix Multiplication*

1073 Figures 13 and 14 illustrate output views informally using the examples *transposed matrix*  
1074 *multiplication* and *double reduction*.

1075 In the case of transposed matrix multiplication, the computed MDA (the computation of matrix  
1076 multiplication is presented later and not relevant for our following considerations) is stored via  
1077 an output view as a matrix in a transposed fashion, using index function  $(i, j, 0) \mapsto (j, i)$ . Here,

Diagram illustrating the execution of a reduce operation. A dashed box encloses the assignment of  $\text{idr}_{2,1}(0) := ()$  and the reduce operation. The reduce operation is labeled with  $\text{reduce}(\oplus, \otimes)$  and  $\text{ov}$  above it. Inputs are  $(y_1, y_2)$  and  $\text{pos: } (0)$ . The output is a list of two elements:  $(y_1)$  and  $(y_2)$ .

Fig. 14. Output view illustrated using the example *double reduction*

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  $(0) \mapsto ()^{11}$  for both pair elements.

**Definition 9** (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  $\text{ov}$  of type:

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

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

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

$$\text{QV}^{\underbrace{\text{MDA's meta-parameters}}_{D=2 \mid I_1=[0, 1024]_{\mathbb{N}_0}, I_2=\{0\} \mid T=\mathbb{Z}}, \underbrace{\text{BUF's meta-parameters}}_{B=1 \mid D_1=1 \mid (N_1^1=1024) \mid T_1=\mathbb{Z}}}$$

and defined as

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

**Example 9** (Output View – Jacobi1D). The output view of Jacobi1D computing a  $(512 - 2)$ -sized vector of integers is of type

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

$$\underbrace{\text{Ov}^{<D=1 \mid I_1=[0, 512-2]_{\mathbb{N}_0} \mid T=\mathbb{Z} \mid B=1 \mid D_1=1 \mid (N_1^1=(512-2)) \mid T_1=\mathbb{Z} \geq 0}}_{\text{MDA's meta-parameters}}$$

and defined as

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

---

<sup>11</sup> The empty braces denote accessing a scalar value (formal details provided in the Appendix, Section B.7).

1128 We define higher-order function `out_view` formally as follows.

1129 **Definition 10** (Higher-Order Function `out_view`). Function `out_view` is of type

1131  $\text{out\_view} \in \{ B \in \mathbb{N} \mid A_1, \dots, A_B \in \mathbb{N} \mid D_1, \dots, D_B \in \mathbb{N}_0 \mid D \in \mathbb{N} \mid$

1132     Number of BUFs     Number BUFs' Accesses     BUFs' Dimensionalities     MDA Dimensionality

1134      $(\Rightarrow_{\text{BUF}}^{\text{MDA } b, a, \text{MDA-IDX-SETS}^D \rightarrow \text{BUF-IDX-SETS}^{D_b}})_{b \in [1, B] \mathbb{N}, a \in [1, A_b] \mathbb{N}} \dots$

1136     Index Set Functions (MDA indices to BUF indices)

1137     MDA's Meta-Parameters

1138      $\underbrace{B \times A_b}_{\substack{b=1 \\ a=1}} \underbrace{\text{MDA-IDX-to-BUF-IDX}^{<D, D_b | \Rightarrow_{\text{BUF}}^{\text{MDA } b, a}}}_{\substack{\text{Index Function: } \text{id}\mathfrak{x}_{b, a} \\ \text{Buffer Access}}} \rightarrow \underbrace{\text{OV}^{<D | I_1, \dots, I_D \in \text{MDA-IDX-SETS}^D \rightarrow |}}_{\substack{\text{BUFs' Meta-Parameters} \\ \text{Postponed Parameters}}} \underbrace{| B | D_1, \dots, D_B | \rightarrow | T_1, \dots, T_B \in \text{TYPE} > < N_1^1, \dots, N_{D_B}^B | T >}_{\substack{\text{MDA's Meta-Parameters} \\ \text{Postponed Parameters}}}$

1142     Index Functions:  $\text{id}\mathfrak{x}_{1,1}, \dots, \text{id}\mathfrak{x}_{B, A_B}$

1143     Output View:  $\text{ov}$

1144 which differs from `inp_view`'s type only in mapping index functions to `OV` (Definition 9), rather  
1145 than `IV` (Definition 7). Function `out_view` is defined as:

1146

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

1151 for

1153  $\mathfrak{b}_b[\text{id}\mathfrak{x}_{b, a}(i_1, \dots, i_D)] := \mathfrak{a}_{b, a}[i_1, \dots, i_D]$

1154 and

1156  $(\mathfrak{a}_{b, a}[i_1, \dots, i_D])_{b \in [1, B] \mathbb{N}, a \in [1, A_b] \mathbb{N}} := \mathfrak{a}[i_1, \dots, i_D]$

1157

1158 i.e.,  $\mathfrak{a}_{b, a}[i_1, \dots, i_D]$  is the element at point  $i_1, \dots, i_D$  within MDA  $\mathfrak{a}$  that belongs to the  $a$ -th access  
1159 of the  $b$ -th BUF. We set  $\mathfrak{b}_b[j_1, \dots, j_{D_b}] := \perp$  (symbol  $\perp$  denotes the undefined value) for all BUF  
1160 indices  $(j_1, \dots, j_{D_b}) \in [0, N_1^b] \mathbb{N}_0 \times \dots \times [0, N_{D_b}^b] \mathbb{N}_0 \setminus \bigcup_{a \in [1, A_b] \mathbb{N}} \xrightarrow{d_{\text{MDA } b, a}} (I_1, \dots, I_D)$  which are not in  
1161 the function range of the index functions.

1162 Note that the computed output view  $\text{ov}$  is partial (indicated by  $\rightarrow_p$  in Definition 9), because for non-  
1163 injective index functions, it must hold  $\text{id}\mathfrak{x}_{b, a}(i_1, \dots, i_D) = \text{id}\mathfrak{x}_{b, a'}(i'_1, \dots, i'_D) \Rightarrow \mathfrak{a}_{b, a}[i_1, \dots, i_D] =$   
1164  $\mathfrak{a}_{b, a'}[i'_1, \dots, i'_D]$ , which may not be satisfied for each potential input MDA of the computed view.  
1165

1166 **Notation 3** (Output Views). Analogously to Notation 2, we denote `out_view` for a particular  
1167 choice of index functions as:

1168  $\text{out\_view}(\text{ID}_1 : \text{id}\mathfrak{x}_{1,1}, \dots, \text{id}\mathfrak{x}_{1, A_1}, \dots, \text{ID}_B : \text{id}\mathfrak{x}_{B,1}, \dots, \text{id}\mathfrak{x}_{B, A_B})$

1170 **Example 10.** Function `out_view` is used for `MatVec` and `Jacobi1D` (in Notation 3) as follows:

1172  $\text{MatVec: } \text{out\_view}(\text{w}: (i, k) \mapsto (i)) \quad \text{Jacobi1D: } \text{out\_view}(\text{w}: (i) \mapsto (i))$

1173      $\underbrace{\text{a=1}}_{\text{a=1}}$

1174      $\underbrace{\text{b=1}}_{\text{b=1}}$

1175      $\underbrace{\text{a=1}}_{\text{a=1}}$

1176      $\underbrace{\text{b=1}}_{\text{b=1}}$

1177 2.3.4 *Relation between View Functions.* View functions transform data from their domain-specific  
 1178 representation to the MDA-based representation (via input views) and back (via output views).  
 1179 In our implementation presented later, we aim to access data uniformly in the form of MDAs,  
 1180 thereby being independent of domain-specific data representations. However, we aim to store the  
 1181 data physically in the domain-specific format, as such format is usually the most efficient data  
 1182 representation, e.g., storing the input data of MatVec as a matrix and vector, rather than as a single  
 1183 MDA which contains many redundancies (e.g., each vector element once per row of the input  
 1184 matrix, as illustrated in Figure 11).

1185 The following lemma proves that functions `inp_view` and `out_view` are invertible and that they  
 1186 are each others inverses. Consequently, the lemma shows how we can arbitrarily switch between  
 1187 the domain-specific data representation and the MDA representation, and consequently also that  
 1188 we can implicitly identify MDAs with their domain-specific data representation (represented  
 1189 in our formalism as BUFs, Definition 5). For example, for computing MatVec, we will specify  
 1190 the computations via pattern `md_hom` which operates on MDAs (see Figure 5), but we use the  
 1191 view functions in our implementation to implicitly forward the MDA accesses to accesses on the  
 1192 physically stored BUFs.

1193 **Lemma 2.** Let

1195  $\text{inp\_view}(\text{ID}_1 : \text{id}\mathfrak{x}_{1,1}, \dots, \text{id}\mathfrak{x}_{1,A_1}, \dots, \text{ID}_B : \text{id}\mathfrak{x}_{B,1}, \dots, \text{id}\mathfrak{x}_{B,A_B})$

1196 and

1198  $\text{out\_view}(\text{ID}_1 : \text{id}\mathfrak{x}_{1,1}, \dots, \text{id}\mathfrak{x}_{1,A_1}, \dots, \text{ID}_B : \text{id}\mathfrak{x}_{B,1}, \dots, \text{id}\mathfrak{x}_{B,A_B})$

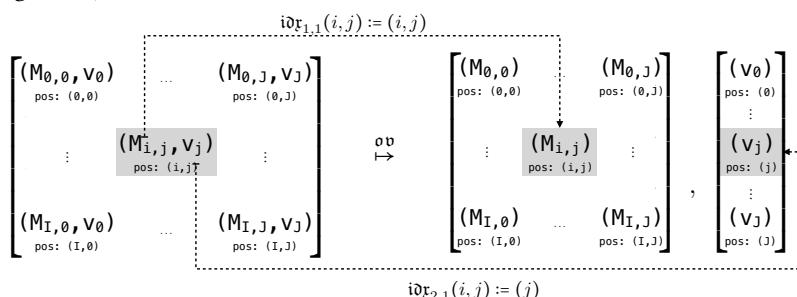
1199 be two arbitrary instances of functions `inp_view` and `out_view` (in Notations 2 and 3), both using  
 1200 the same index functions  $\text{id}\mathfrak{x}_{1,1}, \dots, \text{id}\mathfrak{x}_{B,A_B}$ .

1201 It holds (index functions omitted via ellipsis for brevity):

1202  $\text{inp\_view}(\dots) \circ \text{out\_view}(\dots) = \text{out\_view}(\dots) \circ \text{inp\_view}(\dots) = id$

1204 PROOF. Follows immediately from Definitions 8 and 10.  $\square$

1205 The following figure illustrates the lemma using as example the inverse of MatVec's input view  
 1206 (shown in Figure 11):



## 2.4 Generic High-Level Expression

1219 Figure 15 shows a generic expression in our high-level representation – consisting of higher-order  
 1220 functions `inp_view`, `md_hom`, and `out_view` (Figure 5) – for an arbitrary but fixed choice of index  
 1221 functions, scalar function, and combine operators. We express data-parallel computations using a  
 1222 particular instance of this generic expression in Figure 15.

1223 Note that all meta-parameters of higher-order function `inp_view`, `out_view`, and `md_hom` are  
 1224 automatically deduced from the particular numbers and/or types of the three functions' inputs

1226 (index functions in the case of `inp_view` and `out_view`, and scalar function and combine operators  
 1227 for `md_hom`).

1228 The concrete instance of `md_hom(...)` (i.e., the MDH function returned by `md_hom` for the par-  
 1229 ticular input scalar function and combine operators) has as meta-parameter the MDH function's  
 1230 index sets (see Definition 3) for high flexibility. We use as index sets straightforwardly the input  
 1231 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>12</sup>. Instances of  
 1232 `inp_view(...)` and `out_view(...)` (i.e., the input and output view returned by `inp_view` and  
 1233 `out_view` for concrete index functions) have as meta-parameters the MDA's index sets and the  
 1234 scalar types of BUFS. We explicitly state only the meta-parameter for the BUFS' scalar types in  
 1235 our generic high-level expression (Figure 15), and we avoid explicitly stating the MDA's index  
 1236 sets for simplicity and to avoid redundancies, because the sets can be taken from the `md_hom`'s  
 1237 meta-parameter list.

1238 Note that for better readability of our high-level expressions (Figure 15), we list meta-parameters  
 1239 before parentheses, i.e., instead of writing `inp_view(...)`, `out_view(...)`, and `md_hom(...)`,  
 1240 for the particular instances of higher-order functions, where meta-parameters are listed at the end,  
 1241 we write `inp_view<...>(...)`, `out_view<...>(...)`, and `md_hom<...>(...)`.

1242

1243  
 1244  $\text{out\_view} < T_1^{\text{OB}}, \dots, T_{B^{\text{OB}}}^{\text{OB}} > ( \text{OB}_1 : \text{id}x_{1,1}^{\text{OUT}}, \dots, \text{id}x_{1,A_1^{\text{OB}}}^{\text{OUT}}, \dots, \text{OB}_{B^{\text{OB}}} : \text{id}x_{B^{\text{OB}},1}^{\text{OUT}}, \dots, \text{id}x_{B^{\text{OB}},A_{B^{\text{OB}}}^{\text{OB}}}^{\text{OUT}} ) \circ$   
 1245  
 1246  $\text{md\_hom} < N_1, \dots, N_D > ( f, (\oplus_1, \dots, \oplus_D) ) \circ$   
 1247  
 1248  $\text{inp\_view} < T_1^{\text{IB}}, \dots, T_{B^{\text{IB}}}^{\text{IB}} > ( \text{IB}_1 : \text{id}x_{1,1}^{\text{INP}}, \dots, \text{id}x_{1,A_1^{\text{IB}}}^{\text{INP}}, \dots, \text{IB}_{B^{\text{IB}}} : \text{id}x_{B^{\text{IB}},1}^{\text{INP}}, \dots, \text{id}x_{B^{\text{IB}},A_{B^{\text{IB}}}^{\text{IB}}}^{\text{INP}} )$   
 1249

1250 Fig. 15. Generic high-level expression for data-parallel computations  
 1251  
 1252

## 1253 2.5 Examples

1254 Figure 16 shows how our high-level representation is used for expressing important data-parallel  
 1255 computations. For brevity, we state only the index functions, scalar function, and combine operators  
 1256 of the higher-order functions; the expression in Figure 15 is then obtained by straightforwardly  
 1257 inserting these building blocks into the higher-order functions.  
 1258

1259 *Subfigure 1.* shows how our high-level representation is used for expressing linear algebra  
 1260 routines: 1) Dot (*Dot Product*); 2) MatVec (*Matrix-Vector Multiplication*); 3) MatMul (*Matrix Multipli-*  
 1261 *cation*); 4)  $\text{MatMul}^{\top}$  (*Transposed Matrix Multiplication*) which computes matrix multiplication on  
 1262 transposed input and output matrices; 5) bMatMul (*batched Matrix Multiplication*) where multiple  
 1263 matrix multiplications are computed using matrices of the same sizes.  
 1264

1265 We can observe from the subfigure that our high-level expressions for the routines naturally  
 1266 evolve from each other. For example, the `md_hom` expression for MatVec differs from the `md_hom`  
 1267 expression for Dot by only containing a further concatenation dimension `++` for its  $i$  dimension. We  
 1268 consider this close relation between the high-level expressions of MatVec and Dot in our approach  
 1269 as natural and favorable, as MatVec can be considered as computing multiple times Dot – one  
 1270 computation of Dot for each value of MatVec's  $i$  dimension. Similarly, the `md_hom` expression  
 1271 for MatMul is very similar to the expression of MatVec, by containing the further concatenation

1272 <sup>12</sup> Our formalism allows dynamic shapes, by using symbol  $*$  instead of a particular natural number for  $N_i$  (formal details  
 1273 provided in the Appendix, Definition 23), which we aim to discuss thoroughly in future work.

1274

1275	md_hom	f	$\oplus_1$	$\oplus_2$	$\oplus_3$	$\oplus_4$	Views	inp_view		out_view					
1276	Dot	*	+					A		B					
1277	MatVec	*	+	+				$(k) \mapsto (k)$		$(k) \mapsto (k)$					
1278	MatMul	*	+	+	+	+		$(i,k) \mapsto (i,k)$		$(i,k) \mapsto (i,k)$					
1279	MatMul <sup>T</sup>	*	+	+	+	+		$(i,j,k) \mapsto (i,j,k)$		$(i,j,k) \mapsto (i,j,k)$					
1280	bMatMul	*	+	+	+	+		$(b,i,j,k) \mapsto (b,i,j,k)$		$(b,i,j,k) \mapsto (b,i,j,k)$					
1281	1) Linear Algebra Routines														
1282	md_hom	f	$\oplus_1$	$\oplus_2$	$\oplus_3$	$\oplus_4$	$\oplus_5$	$\oplus_6$	$\oplus_7$	$\oplus_8$	$\oplus_9$				
1283	Conv2D	*	+	+	+	+	+								
1284	MCC	*	+	+	+	+	+	+	+						
1285	MCC_Capsule	*	+	+	+	+	+	+	+	+	+				
1286	Views	inp_view						out_view							
1287	Conv2D	I						F							
1288	MCC	$(p,q,r,s) \mapsto (p+r,q+s)$						$(p,q,r,s) \mapsto (r,s)$							
1289	MCC_Capsule	$(n,p,...) \mapsto (n,p+r,q+s,c)$						$(n,p,...) \mapsto (k,r,s,c)$							
1290	2) Convolution Stencils														
1291	md_hom	f	$\oplus_1$	$\dots$	$\oplus_6$	$\oplus_7$	Views	inp_view		out_view					
1292	CCSD(T)	*	+	...	+	+		A		B					
1293	I1							$(a,\dots,g) \mapsto (g,d,a,b)$		$(a,\dots,g) \mapsto (e,f,g,c)$					
1294	I2							$(a,\dots,g) \mapsto (g,d,a,c)$		$(a,\dots,g) \mapsto (e,f,g,b)$					
1295	3) Quantum Chemistry														
1296	md_hom	f	$\oplus_1$	$\oplus_2$	$\oplus_3$		Views	inp_view		out_view					
1297	Jacobi1D							I		0					
1298	Jacobi2D							$(i1) \mapsto (i1+0)$		$(i1) \mapsto (i1+1)$					
1299	Jacobi3D							$(i1,i2) \mapsto (i1+0,i2+1)$		$(i1,i2) \mapsto (i1,i2)$					
1300	4) Jacobi Stencils														
1301	md_hom	f	$\oplus_1$	$\oplus_2$	$\oplus_3$		Views	inp_view		out_view					
1302	PRL							I		0					
1303	PRL							$(i,j) \mapsto (i)$		$(i,j) \mapsto (j)$					
1304	5) Probabilistic Record Linkage														
1305	md_hom	f	$\oplus_1$	$\oplus_2$		Views	inp_view		out_view						
1306	Histo							N		M					
1307	GenHisto							$(i,j) \mapsto (i)$		$(i,j) \mapsto (j)$					
1308	PRL							$(i,j) \mapsto (i)$		$(i,j) \mapsto (i)$					
1309	6) Histogram														
1310	md_hom	f	$\oplus_1$		Views	inp_view		out_view							
1311	map(f)							I		0 <sub>1</sub>					
1312	reduce( $\oplus$ )							$(i) \mapsto (i)$		$(i) \mapsto (i)$					
1313	reduce( $\oplus$ )							$(i) \mapsto (i)$		$(i) \mapsto ()$					
1314	reduce( $\oplus, \otimes$ )							$(i) \mapsto (i)$		$(i) \mapsto ()$					
1315	7) Map/Reduce Patterns														
1316	md_hom	f	$\oplus_1$		Views	inp_view		out_view							
1317	scan( $\oplus$ )							A		Out					
1318	MBBS							$(i) \mapsto (i)$		$(i) \mapsto (i)$					
1319	MBBS							$(i,j) \mapsto (i,j)$		$(i,j) \mapsto (i,j)$					
1320	8) Prefix Sum Computations														
1321															

Fig. 16. Data-parallel computations expressed in our high-level representation

1324 dimension  $j$  for MatMul's  $j$  dimension. The same applies to bMatMul: its `md_hom` expression is the  
 1325 expression of MatMul augmented with one further concatenation dimension.

1326 Regarding  $\text{MatMul}^\top$ , the basic computation part of  $\text{MatMul}^\top$  and MatMul are the same, which is  
 1327 exactly reflected in our formalisms: both  $\text{MatMul}^\top$  and MatMul are expressed using exactly the same  
 1328 `md_hom` expression. The differences between  $\text{MatMul}^\top$  and MatMul lies only in the data accesses –  
 1329 transposed accesses in the case of  $\text{MatMul}^\top$  and non-transposed accesses in the case of MatMul. Data  
 1330 accesses are expressed in our formalism, in a structured way, via the view functions: for example,  
 1331 for  $\text{MatMul}^\top$ , we use for its first input matrix  $A$  index function  $(i, j, k) \mapsto (k, i)$ , for transposed  
 1332 access, instead of using index function  $(i, j, k) \mapsto (i, k)$  as for MatMul's non-transposed accesses.

1333 Note that all `md_hom` expressions in the subfigure are well defined according to Lemma 1.

1334  
 1335 *Subfigure 2.* shows how convolution-style stencil computations are expressed in our high-level  
 1336 representation: 1) Conv2D expresses a standard convolution that uses a 2D sliding window [Podlozh-  
 1337 nyuk 2007]; 2) MCC expresses a so-called *Multi-Channel Convolution* [Dumoulin and Visin 2018] – a  
 1338 generalization of Conv2D that is heavily used in the area of deep learning; 3) MCC\_Capsule is a  
 1339 recent generalization of MCC [Hinton et al. 2018] which attracted high attention due to its relevance  
 1340 for advanced deep learning neural networks [Barham and Isard 2019].

1341 While our `md_hom` expressions for convolutions are quite similar to those of linear algebra rou-  
 1342 tines (they all use multiplication  $*$  as scalar function, and a mix of concatenations  $+$  and point-wise  
 1343 additions  $+$  as combine operators), the index functions used for the view functions of convolutions  
 1344 are notably different from those used for linear algebra routines: the index functions of convolutions  
 1345 contain arithmetic expressions (e.g.,  $p+r$  and  $q+s$ ), thereby access neighboring elements in their  
 1346 input – a typical access pattern in stencil computations requiring special optimizations [Hagedorn  
 1347 et al. 2018]. Moreover, convolution-style computations are often high-dimensional (e.g., 10 dimen-  
 1348 sions in the case of MCC\_Capsule), whereas linear algebra routines usually rely on less dimensions.  
 1349 Our experiments in Section 5 confirm that optimizations respecting the data access patterns and  
 1350 high dimensionality of convolution computations, as in our approach, usually achieve significantly  
 1351 higher performance than using optimizations chosen toward linear algebra routines as in vendor  
 1352 libraries provided by NVIDIA and Intel for convolutions [Li et al. 2016].

1353 *Subfigure 3.* shows how quantum chemistry computation *Coupled Cluster* (CCSD(T)) [Kim et al.  
 1354 2019] is expressed in our high-level representation. The CCSD(T) computation notably differs from  
 1355 those of linear algebra routines and convolution-style stencils, by accessing its high-dimensional  
 1356 input data in sophisticated transposed fashions: for example, the view function of CCSD(T)'s  
 1357 *instance one* (denoted as I1 in the subfigure) uses indices  $a$  and  $b$  to access the last two dimensions  
 1358 of its  $A$  input tensor (rather than the first two dimensions of the tensor, as would be the case for  
 1359 non-transposed accesses).

1360 The subfigure presents only two CCSD(T) instances, for brevity – in our experiments in Section 5,  
 1361 we present experimental results for nine different real-world CCSD(T) instances.

1362  
 1363 *Subfigures 4-6.* show computations whose scalar functions and combine operators are different  
 1364 from those used in Subfigures 1-3 (which are straightforward multiplications  $*$ , concatenation, and  
 1365 point-wise additions  $+$  only in Subfigures 1-3). For example, Jacobi stencils (Subfigure 4) use as  
 1366 scalar function the Jacobi-specific computation  $J_{\text{nd}}$  [Cecilia et al. 2012], and *Probabilistic Record*  
 1367 *Linkage* (PRL) [Christen 2012], which is used to identify duplicate entries in a data base, uses a  
 1368 PRL-specific both scalar function  $\text{wght}$  and combine operator  $\text{max}_{\text{PRL}}$  (point-wise combination via  
 1369 the PRL-specific binary function  $\text{max}_{\text{PRL}}$ ) [Rasch et al. 2019b]. Histograms, in their generalized  
 1370 version [Henriksen et al. 2020] (denoted as GenHisto in Subfigure 6), use an arbitrary, user-defined  
 1371 scalar function  $f$  and a user-defined associative and commutative combine operator  $\oplus$ ; the standard  
 1372

1373 histogram variant  $\text{Histo}$  is then a particular instance of  $\text{GenHist}$ , for  $\oplus = +$  (point-wise addition)  
 1374 and  $f = f_{\text{Histo}}$ , where  $f_{\text{Histo}}(e, b) = 1$  iff  $e = b$  and  $f_{\text{Histo}}(e, b) = 0$  otherwise. Histogram's are often  
 1375 analyzed regarding their runtime complexity [Henriksen et al. 2020]; we provide such a discussion  
 1376 for our MDH-based Histogram implementation in the Appendix, Section B.8, for the interested  
 1377 reader.

1378  
 1379 *Subfigure 7.* shows how typical  $\text{map}$  and  $\text{reduce}$  patterns [González-Vélez and Leyton 2010]  
 1380 are implemented in our high-level representation. Examples  $\text{map}(f)$  and  $\text{reduce}(\oplus)$  (discussed  
 1381 in Examples 3 and 4) are simple and thus straightforwardly expressed in our representation. In  
 1382 contrast, example  $\text{reduce}(\oplus, \otimes)$  is more complex and shows how  $\text{reduce}(\oplus)$  is extended toward  
 1383 combining the input vector simultaneously twice – once combining vector elements via operator  
 1384  $\oplus$  and once using operator  $\otimes$ . The outcome of  $\text{reduce}(\oplus, \otimes)$  are two scalars – one representing  
 1385 the result of combination via  $\oplus$  and the other of combination via  $\otimes$  – which we map via the output  
 1386 view to output elements  $O_1$  (result of  $\oplus$ ) and  $O_2$  (result of  $\otimes$ ), correspondingly; this is also illustrated  
 1387 in Figure 14.

1388  
 1389 *Subfigure 8.* presents *prefix-sum computations* [Blelloch 1990] which differ from the computations  
 1390 in Subfigures 1-7 in terms of their combine operators: the operator used for expressing computations  
 1391 in Subfigure 8 is different from concatenation (Example 1) and point-wise combinations  
 1392 (Example 2). Computation  $\text{scan}(\oplus)$  uses as combine operator  $+\text{prefix-sum}(\oplus)$  which computes  
 1393 prefix-sum [Gorlatch and Lengauer 1997] (formally defined in the Appendix, Section B.9) according  
 1394 to binary operator  $\oplus$ , and MBBS (Maximum Bottom Box Sum) [Farzan and Nicolet 2019] uses a  
 1395 particular instance of prefix-sum for  $\oplus = +$  (addition).

### 1396 3 LOW-LEVEL REPRESENTATION FOR DATA-PARALLEL COMPUTATIONS

1397 We introduce our low-level representation for expressing data-parallel computations. In contrast to  
 1398 our high-level representation, our low-level representation explicitly expresses the de-composition  
 1399 and re-composition of computations (informally illustrated in Figure 3). Moreover, our low-level  
 1400 representation is designed such that it can be straightforwardly transformed to executable program  
 1401 code, because it explicitly captures and expresses data movement and parallelization optimizations.  
 1402

1403 In the following, after briefly discussing an introductory example in Section 3.1, we introduce in  
 1404 Section 3.2 our formal representation of computer systems, to which we refer to as *Abstract System  
 1405 Model (ASM)*. Based on this model, we define *low-level MDAs*, *low-level BUFS*, and *low-level combine  
 1406 operators* in Section 3.3, which are basic building blocks of our low-level representation.

1407 Note that all the details and concepts discussed in this section are transparent to the end user  
 1408 of our system and therefore the user is not exposed to these details: expressions in our low-level  
 1409 representation are fully automatically generated from expressions in our high-level representation  
 1410 for the user (Figure 4), according to the methodologies presented later in Section 4 and auto-  
 1411 tuning [Rasch et al. 2021].

#### 1412 3.1 Introductory Example

1413 Figure 17 illustrates our low-level representation by showing how *MatVec* (matrix-vector multipli-  
 1414 cation) is expressed in our representation. In our example, we use an input matrix  $M \in T^{512 \times 4096}$  of  
 1415 size  $512 \times 4096$  (size chosen arbitrarily) that has an arbitrary but fixed scalar type  $T \in \text{TYPE}$ ; the  
 1416 input vector  $v \in T^{4096}$  is of size 4096, correspondingly.

1417 For better illustration, we consider for this introductory example a straightforward, artificial  
 1418 target architecture that has only two memory layers – *Host Memory (HM)* and *Cache Memory (L1)* –  
 1419 and one *Core Layer (COR)* only; our examples presented and discussed later in this section target  
 1420

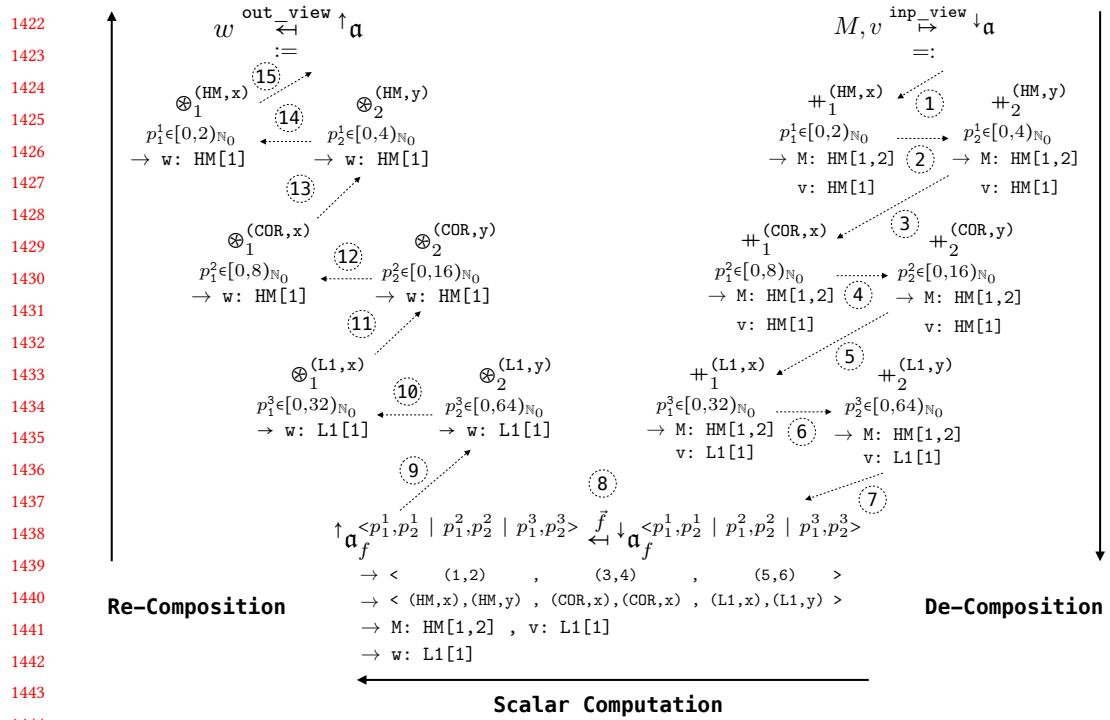


Fig. 17. Low-level expression for straightforwardly computing Matrix-Vector Multiplication (MatVec) on a simple, artificial architecture with two memory layers (HM and L1) and one core layer (COR). Dotted lines indicate data flow.

real-world architectures (e.g., CUDA-capable NVIDIA GPUs). The particular values of tuning parameters (discussed in detail later in this section), such as the number of threads and the order of 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 17), 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 17) partitions input MDA  $\downarrow a$  (in the top right of Figure 17) to the structure  $\downarrow a_f^{<\dots>}$  (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 18 – 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 17, which represents a two-dimensional MDH computation). Consequently, our

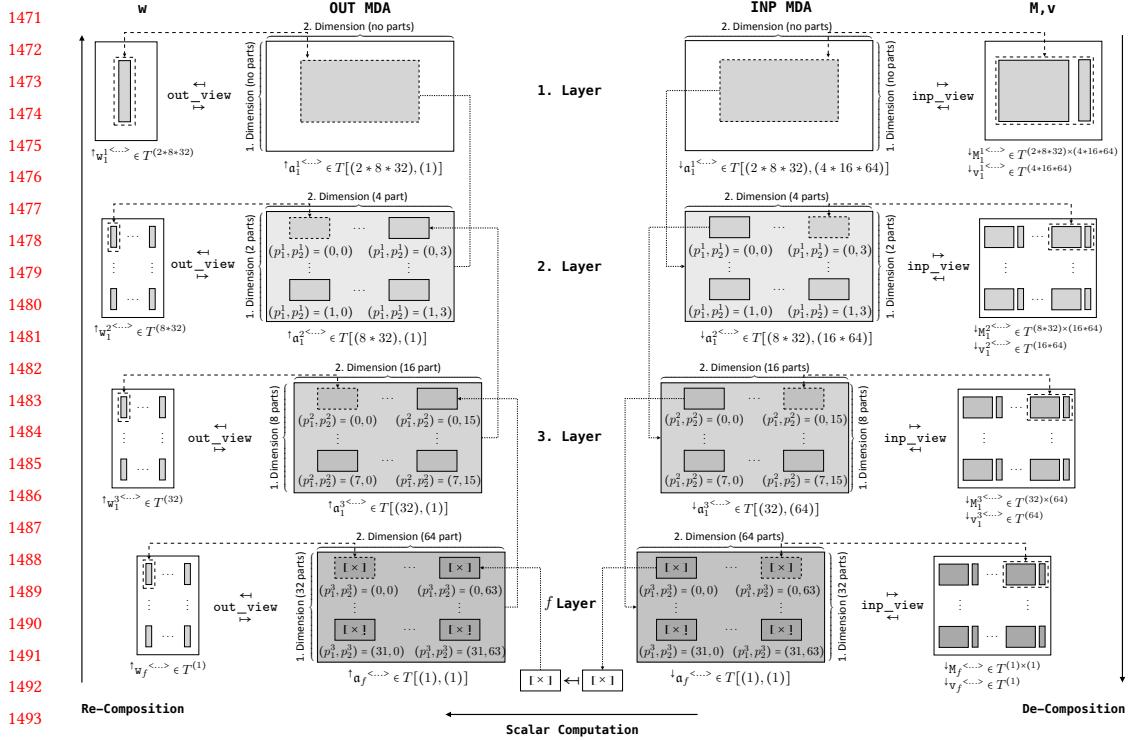


Fig. 18. Illustration of multi-layered, multi-dimensional MDA partitioning using the example MDA from Figure 17. In this example, we use three layers and two dimensions, according to Figure 17.

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

We compute the partitionings of MDAs by applying the concatenation operator (Example 1) inversely<sup>13</sup> (indicated by using  $=:$  instead of  $:=$  in the top right part of Figure 17). For example, we partition in Figure 17 MDA  $\mathbf{a}$  first via the inverse of  $+_1^{(HM, x)}$  in dimension 1 (indicated by the subscript 1 of  $+_1^{(HM, x)}$ ; the superscript  $(HM, x)$  is explained later) in 2 parts, as  $p_1^1$  iterates over interval  $[0, 2]_{N_0} = \{0, 1\}$  which consists of two elements (0 and 1) – the interval is chosen arbitrarily for this example. Afterwards, each of the obtained parts is further partitioned, in the second dimension, via  $+_2^{(HM, y)}$  in 4 parts ( $p_2^1$  iterates over  $[0, 4]_{N_0} = \{0, 1, 2, 3\}$  which consists of four elements). The  $(2 * 4)$ -many HM parts are then each further partitioned in both dimensions for the COR layer in  $(8 * 16)$  parts, and each individual COR part is again partitioned for the L1 layer in  $(32 * 64)$  parts, resulting in  $(2 * 8 * 32) * (4 * 16 * 64) = 512 * 4096$  parts in total.

<sup>13</sup>It is easy to see that operator *concatenation* (Example 1) is invertible for any particular choice of meta-parameters (formally proved in the Appendix, Section C.2).

1520 We always use a *full partitioning* in our low-level expressions<sup>14</sup>, i.e., each particular choice of  
 1521 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 18, the individual  
 1522 MDA elements are denoted via symbol  $\times$ , in the bottom part of the figure). By relying on a full  
 1523 partitioning, we can apply scalar function  $f$  to these singleton parts in the scalar phase (described  
 1524 in the next paragraph). This is because function  $f$  is defined to take as input a single MDA element  
 1525 only (Definition 4), thereby making defining scalar functions more convenient for the user.

1526 The superscript of combine operators, e.g.,  $(\text{COR}, x)$  of operator  $+_1^{(\text{COR}, x)}$ , is a so-called *operator tag*  
 1527 (formal definition given in the next subsection). A tag indicates to our code generator  
 1528 whether its combine operator is assigned to a memory layer (and thus computed sequentially in  
 1529 our generated code) or to a core layer (and thus computed in parallel). For example, tag  $(\text{COR}, x)$   
 1530 indicates that parts processed by operator  $+_1^{(\text{COR}, x)}$  should be computed by cores COR, and thus in  
 1531 parallel; the dimension tag  $x$  indicates that COR layer's  $x$  dimension should be used for computing  
 1532 the operator (we use dimension  $x$  for our example architecture as an analogous concept to CUDA's  
 1533 thread/block dimensions  $x, y, z$  for GPU architectures [NVIDIA 2022g]), as we also discuss in the next  
 1534 subsection. In contrast, tag  $(\text{HM}, x)$  refers to a memory layer (host memory HM) and thus, operator  
 1535  $+_1^{(\text{HM}, x)}$  is computed sequentially. Since the current state-of-practice programming approaches  
 1536 (OpenMP, CUDA, OpenCL, ...) have no explicit notion of memory tiles (e.g., by offering the  
 1537 potential variables  $\text{tileIdx}.x/\text{tileIdx}.y/\text{tileIdx}.z$ , as analogous concepts to CUDA variables  
 1538  $\text{threadIdx}.x/\text{threadIdx}.y/\text{threadIdx}.z$ ), the dimensions tag  $x$  in  $(\text{HM}, x)$  is ignored by our code  
 1539 generator (which currently generates OpenMP, CUDA, or OpenCL code), because HM refers to a  
 1540 memory layer.

1541 Note that the number of parts (2 parts on layer 1 in dimension 1; 4 parts on layer 1 in dimension  
 1542 2; ...), the combine operators' tags, and our partition order (e.g. first partitioning in MDA's  
 1543 dimension 1 and afterwards in dimension 2) are chosen arbitrarily for this example and should be  
 1544 optimized (auto-tuned) for a particular target device and characteristics of the input and output  
 1545 data (size, memory layouts, etc.) to achieve high performance, which we discuss in detail later in  
 1546 this section.

1547 *2. Scalar Phase:* In the scalar phase (step 8 in Figure 17), we apply MDH's scalar function  $f$  to  
 1548 the individual MDA elements which are contained by the MDA's singleton parts

$$\downarrow_{\mathfrak{a}_f}^{< p_1^1, p_2^1 \mid p_1^2, p_2^2 \mid p_1^3, p_2^3 >}$$

1549 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

$$\uparrow_{\mathfrak{a}_f}^{< p_1^1, p_2^1 \mid p_1^2, p_2^3 \mid p_1^3, p_2^2 >}$$

1550 In the figure,  $\vec{f}$  (introduced in Definition 4) is the slight adaption of function  $f$  that operates on a  
 1551 singleton MDA, rather than a scalar.

1552 Annotation  $\rightarrow < (1, 2), \dots >$  indicates the application order of applying scalar function (in  
 1553 this example, first iterating over  $p_1^1$ , then over  $p_2^1$ , etc), and we use annotation  $\rightarrow < (\text{HM}, x), \dots$   
 1554  $>$  to indicate how the scalar computation is assigned to the target architecture (this is described  
 1555 in detail later in this section). Annotations  $\rightarrow M: \text{HM}$ ,  $v: \text{L1}$  and  $\rightarrow w: \text{L1}$  (in the bottom part  
 1556 of Figure 17) indicate the memory regions to be used for reading and writing the input scalar of  
 1557 function  $f$  (also described later in detail).

1558  
 1559 <sup>14</sup> Our future work (outlined in Section 8) aims to additionally allow coarser-grained partitioning schemas as well, e.g., to  
 1560 target domain-specific hardware extensions (such as NVIDIA Tensor Cores [NVIDIA 2017] which compute  $4 \times 4$  matrices  
 1561 immediately in hardware, rather than  $1 \times 1$  matrices as obtained in the case of a full partitioning).

1569     3. *Re-Composition Phase*: Finally, the re-composition phase (steps 9-15 in Figure 17) combines  
 1570     the computed parts  $\uparrow a_f^{<p_1^1, p_1^1 \mid p_1^2, p_2^2 \mid p_1^3, p_2^3>}$  (bottom left in the figure) to the final result  $\uparrow a$  (top  
 1571     left) via MDH's combine operators, which are in the case of matrix-vector multiplication  $\otimes_1 :=$   
 1572      $\text{++}$  (concatenation) and  $\otimes_2 := +$  (point-wise addition). In this example, we first combine the L1 parts  
 1573     in dimension 2 and then in dimension 1; afterwards, we combine the COR parts in both dimensions,  
 1574     and finally the HM parts. Analogously to before, this order of combine operators and their tags are  
 1575     also chosen arbitrarily for this example and should be auto-tuned for high performance.  
 1576

1577     In the de- and re-composition phases, the arrow notation below combine operators allow ef-  
 1578     ficiently exploiting architecture's memory hierarchy, by indicating the memory region to read  
 1579     from (de-composition phase) or to write to (re-composition phase); the annotations also indicate  
 1580     the memory layouts to use. We exploit these memory and layout information in both: i) our code  
 1581     generation to assign combine operators' input and output data to memory regions and to chose  
 1582     memory layouts for the data (row major, column major, etc); ii) our formalism to specify constraints  
 1583     of programming models, e.g., that in CUDA, results of GPU cores can only be combined in desig-  
 1584     nated memory regions [NVIDIA 2022f]. For example, annotation  $\rightarrow M: \text{HM}[1, 2], v: \text{L1}[1]$  below  
 1585     an operator in the de-composition phase indicates to our code generator that the parts (a.k.a tiles)  
 1586     of matrix  $M$  used for this computation step should be read from host memory HM and that parts of  
 1587     vector  $v$  should be copied to and accessed from fast L1 memory. The annotation also indicates that  $M$   
 1588     should be stored using a row-major memory layout (as we use  $[1, 2]$  and not  $[2, 1]$ ). The memory  
 1589     regions and layouts are chosen arbitrarily for this example and should be chosen as optimized for  
 1590     the particular target architecture and characteristics of the input and output data (we currently rely  
 1591     on auto-tuning [Rasch et al. 2021] for choosing optimized values of performance-critical parameters,  
 1592     as we discuss in Section 5).

1593     Formally, the arrow notation of combination operators is a concise notation to hide MDAs and  
 1594     BUFs for intermediate results, which is discussed in the Appendix, Section C.3, for the interested  
 1595     reader.

## 1596     Excuse: Code Generation

1597     Our low-level expressions can be straightforwardly transformed to executable program code  
 1598     in imperative-style programming languages (such as OpenMP, CUDA, and OpenCL). As code  
 1599     generation is not the focus of this work, we demonstrate our code generation approach briefly  
 1600     using the example of Figure 17. Details about our code generation process are provided in Section E  
 1601     of our Appendix, and will be presented and illustrated in detail in our future work.

1602     We implement MDAs via *preprocessor directives*. In the case of MatVec, we implement its input  
 1603     MDA, according to Definition 8, as: `#define inp_mda(i,k) M[i][k],v[k]`

1604     Combine operators are implemented as sequential or parallel loops. For example, the operator  
 1605      $+_1^{(\text{HM}, x)}$  is assigned to memory layer HM and thus implemented as a sequential loop (loop range  
 1606     indicated by  $[0, 2]_{\mathbb{N}_0}$ ), and operator  $+_1^{(\text{COR}, x)}$  is assigned to core layer COR and thus implemented as  
 1607     a parallel loop (e.g., a loop annotated with `#pragma omp parallel` for in OpenMP [OpenMP 2022],  
 1608     or variable `threadIdx.x` in CUDA [NVIDIA 2022g]).

1609     Correspondingly, our three phases (de-composition, scalar, and re-composition) each correspond  
 1610     to an individual loop nest; we generate the nests as fused when the tags of combine operators have  
 1611     the same order in phases, as in Figure 17:  $(\text{HM}, x) \rightarrow (\text{HM}, y) \rightarrow (\text{COR}, x) \rightarrow \dots \rightarrow (\text{L1}, y)$ . Note  
 1612     that our currently targeted programming models (OpenMP, CUDA, and OpenCL) have no explicit no-  
 1613     tion of *tiles*, e.g., by offering the potential variables `tileIdx.x/tileIdx.y/tileIdx.z` for managing  
 1614     tiles automatically in the programming model (similarly as variables `threadIdx.x/threadIdx.y/`  
 1615     `threadIdx.z` automatically manage threads in CUDA). Consequently, the dimension information  
 1616

1618 within tags are currently ignored by our code generator when the operator tag refers to a memory  
 1619 layer (such as dimension  $x$  in tag  $(HM, x)$  which refers to memory layer  $HM$ ).

1620 We implement operators' memory regions as straightforward allocations in the corresponding  
 1621 memory region (e.g., CUDA's device, shared, or register memory [NVIDIA 2022g], according  
 1622 to the arrow annotations in our low-level expression); memory layouts are also implemented via  
 1623 pre-processor directives, e.g., `#define M(i,k) M[k][i]` when storing MatVec's input matrix  $M$  as  
 1624 transposed.

1625 Code optimizations that are applied on a lower abstraction level than proposed by our represen-  
 1626 tation in Example 17 are beyond the scope of this work and outlined in Section E.3 of our Appendix,  
 1627 e.g., loop fusion and loop unrolling which are applied to program code on the low, loop-based  
 1628 abstraction level.

1629  
 1630 In the following, we introduce in Section 3.2 our formal representation of a computer sys-  
 1631 tem (which can be a single device, but also a multi-device or a multi-node system, as we discuss  
 1632 soon), and we illustrate our formal system representation using the example architectures targeted  
 1633 by programming models OpenMP, CUDA, and OpenCL. Afterwards, in Section 3.3, we formally  
 1634 define the basic building blocks of our low-level representation – *low-level MDAs*, *low-level BUFs*,  
 1635 and *low-level combine operators* – based on our formal system representation.

### 1636 3.2 Abstract System Model (ASM)

1637 **Definition 11** (Abstract System Model). An *L-Layered Abstract System Model (ASM)*,  $L \in \mathbb{N}$ , is any  
 1638 pair of two positive natural numbers

$$1640 \quad ( \text{NUM\_MEM\_LYRs}, \text{NUM\_COR\_LYRs} ) \in \mathbb{N} \times \mathbb{N}$$

1641 for which  $\text{NUM\_MEM\_LYRs} + \text{NUM\_COR\_LYRs} = L$ .

1642 Our ASM can model architectures with arbitrarily deep memory and core hierarchies<sup>15</sup>:  $\text{NUM\_MEM\_LYRs}$   
 1643 denotes the target architecture's number of memory layers and  $\text{NUM\_COR\_LYRs}$  the architecture's  
 1644 number of core layers, correspondingly. For example, the artificial architecture we use in Figure 17  
 1645 is represented as an ASM instance as follows (bar symbols denote set cardinality):

$$1646 \quad \text{ASM}_{\text{artif.}} := ( \bar{|\{HM, L1\}|}, \bar{|\{COR\}|} ) = (2, 1)$$

1647 The instance is a pair consisting of the numbers 2 and 1 which represent the artificial architecture's  
 1648 two memory layers ( $HM$  and  $L1$ ) and its single core layers ( $COR$ ).

1649 **Example 11.** We show particular ASM instances that represent the device models of state-of-  
 1650 practice approaches OpenMP, CUDA, and OpenCL:

$$\begin{aligned} 1651 \quad \text{ASM}_{\text{OpenMP}} &:= ( \bar{|\{MM, L2, L1\}|}, \bar{|\{COR\}|} ) = (3, 1) \\ 1652 \quad \text{ASM}_{\text{OpenMP+L3}} &:= ( \bar{|\{MM, L3, L2, L1\}|}, \bar{|\{COR\}|} ) = (4, 1) \\ 1653 \quad \text{ASM}_{\text{OpenMP+L3+SIMD}} &:= ( \bar{|\{MM, L3, L2, L1\}|}, \bar{|\{COR, SIMD\}|} ) = (4, 2) \\ 1654 \quad \text{ASM}_{\text{CUDA}} &:= ( \bar{|\{DM, SM, RM\}|}, \bar{|\{SMX, CC\}|} ) = (3, 2) \\ 1655 \quad \text{ASM}_{\text{CUDA+WRP}} &:= ( \bar{|\{DM, SM, RM\}|}, \bar{|\{SMX, WRP, CC\}|} ) = (3, 3) \\ 1656 \quad \text{ASM}_{\text{OpenCL}} &:= ( \bar{|\{GM, LM, PM\}|}, \bar{|\{CU, PE\}|} ) = (3, 2) \end{aligned}$$

1664 <sup>15</sup> We deliberately do not model into our ASM representation an architecture's particular number of cores and sizes of  
 1665 memory regions, because our optimization process is designed to be generic in these numbers and sizes for high flexibility.

1667 OpenMP is often used to target  $(3 + 1)$ -layered architectures which rely on 3 memory regions (main memory MM, and caches L2 and L1) and 1 core layer (COR). OpenMP-compatible architectures sometimes also contain the L3 memory region, and they may allow exploiting SIMD parallelization (a.k.a. *vectorization* [Klemm et al. 2012]), which are expressed in our ASM representation as a further memory or core layer, respectively.

1672 CUDA's target architectures are  $(3 + 2)$ -layered: they consist of *Device Memory* (DM), *Shared Memory* (SM), and *Register Memory* (RM), and they offer as cores so-called *Streaming Multiprocessors* (SMX) which themselves consist of *Cuda Cores* (CC). CUDA also has an implicit notion of so-called *Warps* (WRP) which are not explicitly represented in the CUDA programming model [NVIDIA 2022g], but often exploited by programmers – via special intrinsics (e.g., *shuffle* and *tensor core intrinsics* [NVIDIA 2017, 2018]) – to achieve highest performance.

1678 OpenCL-compatible architectures are designed analogously to those targeted by the CUDA programming model; consequently, both OpenCL- and CUDA-compatible architectures are represented by the same ASM instance in our formalism. Apart from straightforward syntactical differences between OpenCL and CUDA [StreamHPC 2016], we see as the main differences between the two programming models (from our ASM-based abstraction level) that OpenCL has no notion of warps, and it uses a different terminology – *Global/Local/Private Memory* (GM/LM/PM) instead of device/shared/register memory, and *Compute Unit* (CU) and *Processing Element* (PE), rather than SMX and CC.

1686 In the following, we consider memory regions and cores of ASM-represented architectures as  
1687 arrangeable in an arbitrary number of dimensions. Programming models for such architectures  
1688 often have native support for such arrangements. For example, in the CUDA model, memory is  
1689 accessed via arrays which can be arbitrary-dimensional (a.k.a *multi-dimensional C arrays*), and  
1690 cores are programmed in CUDA via threads which are arranged in CUDA's so-called dimensions x,  
1691 y, z; further thread dimensions can be explicitly programmed in CUDA, e.g., by embedding them  
1692 in the last dimension z. Details on our arrangements of memory and cores are provided in the  
1693 Appendix, Section C.4.

1695 We express constraints of programming models – for example, that in CUDA, SMX can combine  
1696 their results in DM only [NVIDIA 2022f] – via so-called *tuning parameter constraints*, which we  
1697 discuss later in this section.

1698 Note that we call our abstraction *Abstract System Model* (rather than *Abstract Architecture Model*,  
1699 or the like), because it can also represent systems consisting of multiple devices and/or nodes, etc.  
1700 For example, our ASM representation of a multi-GPU system is:

$$1702 \text{ASM}_{\text{Multi-GPU}} := ( |\{\text{HM, DM, SM, RM}\}|, |\{\text{GPU, SMX, CC}\}| ) = (4, 3)$$

1703 It extends our ASM-based representation of CUDA devices (Example 11) by *Host Memory* (HM) which  
1704 represents the memory region of the system containing the GPUs (and in which the intermediate  
1705 results of different GPUs are combined), and it introduces the further core layer GPU representing  
1706 the system's GPUs. Analogously, our ASM representation of a multi-node, multi-GPU system is:

$$1708 \text{ASM}_{\text{Multi-Node-Multi-GPU}} := ( |\{\text{NM, HM, DM, SM, RM}\}|, |\{\text{NOD, GPU, SMX, CC}\}| ) = (5, 4)$$

1709 It adds to ASM  $\text{ASM}_{\text{Multi-GPU}}$  the memory layer (*Node Memory* (NM)) which represents the memory  
1710 region of the host node, and it adds core layer (*Node* (NOD)) which represents the compute nodes.  
1711 Our approach is currently designed for *homogeneous systems*, i.e., all devices/nodes/... are assumed  
1712 to be identical. We aim to extend our approach toward *heterogeneous systems* (which may consist of  
1713 different devices/nodes/...) as future work, inspired by dynamic load balancing approaches [Chen  
1714 et al. 2010].

1716 **3.3 Basic Building Blocks**

1717 We introduce the three main basic building blocks of our low-level representation: 1) *low-level*  
 1718 *MDAs* which represent multi-layered, multi-dimensionally arranged collection of ordinary MDAs  
 1719 (Definition 1) – one ordinary MDA per memory/core layer of their target ASM and for each  
 1720 dimension of the MDA computation (illustrated in Figure 18); 2) *low-level BUFs* which are a collection  
 1721 of ordinary BUFs (Definition 5) and that have a notion *memory regions* and *memory layouts*; 3) *low-*  
 1722 *level combine operators* which represent combine operators (Definition 2) to which the layer and  
 1723 dimension of their target ASM is assigned that are intended to be used to compute the operator in  
 1724 our generated code (e.g., a core layer to compute the operator in parallel).

1725  
 1726 **Definition 12** (Low-Level MDA). Let be  $L \in \mathbb{N}$  (representing an ASM's number of layers) and  
 1727  $D \in \mathbb{N}$  (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,  $T \in \text{TYPE}$   
 1728 a scalar type, and  $I := ((I_d^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \in \text{MDA-IDX-SETs})_{d \in [1, D]_{\mathbb{N}}})^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}$   
 1729 an arbitrary collection of  $D$ -many MDA index sets (Definition 1) for each  
 1730 particular choice of indices  $(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1, \dots, (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L$ <sup>16</sup>  
 1731 (illustrated in Figure 18).

1732 An  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned *low-level MDA* that has scalar type  $T$  and index sets  $I$   
 1733 is any function  $a_{ll}$  of type:

$$1734 \underbrace{a_{ll}^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>} :}_{1735} \underbrace{I_1^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \times \dots \times I_D^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \rightarrow T}_{1736}$$

1737 We use low-level MDAs in the following to represent partitionings of MDAs (as illustrated soon  
 1738 and formally discussed the Appendix, Section C.7).

1739 Next, we introduce *low-level BUFs* which work similarly as BUFs (Definition 5), but are tagged  
 1740 with a memory region and a memory layout. While these tags have no effect on operators' semantics,  
 1741 they indicate later to our code generator in which memory region the BUF is stored and which  
 1742 memory layout to chose for storing the BUF. Moreover, we use these tags to formally define  
 1743 constraints of programming models, e.g., that according to the CUDA specification [NVIDIA 2022f],  
 1744 SMX cores can combine their results in memory region DM only.

1745  
 1746 **Definition 13** (Low-Level BUF). Let be  $L \in \mathbb{N}$  (representing an ASM's number of layers) and  
 1747  $D \in \mathbb{N}$  (representing an MDH's number of dimensions). Let further  $P = ((P_1^1, \dots, P_D^1), \dots, (P_1^L, \dots, P_D^L)) \in \mathbb{N}^{L \times D}$  be an arbitrary tuple of  $L$ -many  $D$ -tuples of positive natural numbers,  $T \in$   
 1748  $\text{TYPE}$  a scalar type, and  $N := ((N_d^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \in \mathbb{N})_{d \in [1, D]_{\mathbb{N}}})^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 | \dots | (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}$   
 1749 be a BUF's size (Definition 5) for each particular choice of  $p_1^1, \dots, p_D^L$ .

1750 An  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned *low-level BUF* that has scalar type  $T$  and size  $N$  is  
 1751 any function  $b_{ll}$  of type ( $\leftrightarrow$  denotes bijection):

$$1752 \underbrace{b_{ll}^{<\text{Memory Region} | \text{Memory Layout} | \text{Partitioning: Layer 1} | \text{Partitioning: Layer } L>} :}_{1753} \underbrace{[0, N_1^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>})_{\mathbb{N}_0} \times \dots \times [0, N_D^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>})_{\mathbb{N}_0} \rightarrow T}_{1754}$$

1755  
 1756 <sup>16</sup> Analogously to Notation 1, we identify each  $P_d^l \in \mathbb{N}$  implicitly also with the interval  $[0, P_d^l]_{\mathbb{N}_0}$  (inspired by set theory).

1765 We refer to MEM as low-level BUF's *memory region* and to  $\sigma$  as its *memory layout*, and we refer to  
1766 the function

1773 that is defined as

$$\mathfrak{b}_{II}^{\text{trans} < \text{MEM} | \sigma > p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L >} (i_{\sigma(1)}, \dots, i_{\sigma(D)}) := \mathfrak{b}_{II}^{< \text{MEM} | \sigma > p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L >} (i_1, \dots, i_D)$$

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

1778 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.

**1791** **Definition 14** (ASM Level). We call pairs  $(l_{\text{ASM}}, d_{\text{ASM}})$ , consisting of an ASM layer  $l_{\text{ASM}} \in [1, L]_{\mathbb{N}}$   
**1792** and an ASM dimension  $d_{\text{ASM}} \in [1, D]_{\mathbb{N}}$ , an *ASM Level* (ASM-LVL<sup>17</sup>) (terminology motivated in the  
**1793** Appendix, Section C.5).

**1795 Definition 15** (Low-Level Combine Operator). Let be  $L \in \mathbb{N}$  (representing an ASM's number of layers)  
**1796** and  $D \in \mathbb{N}$  (representing an MDH's number of dimensions). Let further be  $\otimes \in \text{CO}^{<\Rightarrow \text{MDA}} | T | D | d$   
**1797** an arbitrary  $D$ -dimensional combine operator (Definition 2).

1798 The *low-level representation*  $\otimes_{<(l_{\text{ASM}}, d_{\text{ASM}})>\text{-ASM-LVL}}$  of operator  $\otimes$  is a function that for each pair

$$(l_{\text{ASM}}, d_{\text{ASM}}) \in \{ (l, d) \mid l \in [1, L]_{\mathbb{N}}, d \in [1, D]_{\mathbb{N}} \} \subset \text{ASM-LVU}$$

1802 has the same type and semantics as  $\otimes$ :

$$\bigoplus_{\langle l_{ASM}, d_{ASM} \rangle} \in \bigcap_{MDA} \Rightarrow_{MDA} |T| |D| |d\rangle \quad \quad \bigoplus_{\langle l_{ASM}, d_{ASM} \rangle} (a, b) := \bigoplus (a, b)$$

1805 i.e.,  $\oplus^{<l_{\text{ASM}}, d_{\text{ASM}}>}$  works exactly as combine operator  $\oplus$ , but its type is enriched with a meta-parameter  
 1806 that captures an ASM layer  $l_{\text{ASM}} \in [1, L]_{\mathbb{N}}$  and dimension  $d_{\text{ASM}} \in [1, D]_{\mathbb{N}}$ .  
 1807

1808  
1809 Note that in Figure 17, for better readability, we use domain-specific identifiers for ASM lay-  
1810 ers: HM:=1 as an alias for the ASM layer that has id 1, L1:=2 for layer id 2, and COR:=3 for layer  
1811 id 3. For dimensions, we use aliases  $x := 1$  for ASM dimension 1 and  $y := 2$  for ASM dimension 2,  
1812 correspondingly.

1814	$\mathfrak{b}_1^{\text{IB}}, \dots, \mathfrak{b}_{B^{\text{IB}}}^{\text{IB}}$	$\xrightarrow{\text{inp\_view}}$	$\uparrow \mathfrak{a} :=$
1815	$\oplus_1^{\leftrightarrow_{\text{prt-ass}}(1,1)}$	$\dots$	$\oplus_D^{\leftrightarrow_{\text{prt-ass}}(1,D)}$
1816	$p_1^1 \in \#PRT(1,1)$		$p_D^1 \in \#PRT(1,D)$
1817	$\rightarrow \mathfrak{b}_1^{\text{IB}} : \downarrow\text{-mem}(1,1)[\sigma_{\downarrow\text{-mem}}(1,1)(1), \dots, \sigma_{\downarrow\text{-mem}}(1,1)(D_1^{\text{IB}})]$	$\rightarrow \mathfrak{b}_1^{\text{IB}} : \downarrow\text{-mem}(1,D)[\sigma_{\downarrow\text{-mem}}(1,D)(1), \dots, \sigma_{\downarrow\text{-mem}}(1,D)(D_1^{\text{IB}})]$	
1818	$\vdots$	$\vdots$	$\vdots$
1819	$\mathfrak{b}_{B^{\text{IB}}}^{\text{IB}} : \downarrow\text{-mem}(1,1)[\sigma_{\downarrow\text{-mem}}(1,1)(1), \dots, \sigma_{\downarrow\text{-mem}}(1,1)(D_{B^{\text{IB}}}^{\text{IB}})]$	$\mathfrak{b}_{B^{\text{IB}}}^{\text{IB}} : \downarrow\text{-mem}(1,D)[\sigma_{\downarrow\text{-mem}}(1,D)(1), \dots, \sigma_{\downarrow\text{-mem}}(1,D)(D_{B^{\text{IB}}}^{\text{IB}})]$	$\leftrightarrow \sigma_{\downarrow\text{-ord}}(1,D)$
1820	$\underbrace{\leftrightarrow \sigma_{\downarrow\text{-ord}}(1,1)}$		$\vdots$
1821	$\vdots$		$\vdots$
1822	$\vdots$		$\vdots$
1823	$\oplus_1^{\leftrightarrow_{\text{prt-ass}}(L,1)}$	$\dots$	$\oplus_D^{\leftrightarrow_{\text{prt-ass}}(L,D)}$
1824	$p_1^L \in \#PRT(L,1)$		$p_D^L \in \#PRT(L,D)$
1825	$\rightarrow \mathfrak{b}_1^{\text{IB}} : \downarrow\text{-mem}(L,1)[\sigma_{\downarrow\text{-mem}}(L,1)(1), \dots, \sigma_{\downarrow\text{-mem}}(L,1)(D_1^{\text{IB}})]$	$\rightarrow \mathfrak{b}_1^{\text{IB}} : \downarrow\text{-mem}(L,D)[\sigma_{\downarrow\text{-mem}}(L,D)(1), \dots, \sigma_{\downarrow\text{-mem}}(L,D)(D_1^{\text{IB}})]$	
1826	$\vdots$	$\vdots$	$\vdots$
1827	$\mathfrak{b}_{B^{\text{IB}}}^{\text{IB}} : \downarrow\text{-mem}(L,1)[\sigma_{\downarrow\text{-mem}}(L,1)(1), \dots, \sigma_{\downarrow\text{-mem}}(L,1)(D_{B^{\text{IB}}}^{\text{IB}})]$	$\mathfrak{b}_{B^{\text{IB}}}^{\text{IB}} : \downarrow\text{-mem}(L,D)[\sigma_{\downarrow\text{-mem}}(L,D)(1), \dots, \sigma_{\downarrow\text{-mem}}(L,D)(D_{B^{\text{IB}}}^{\text{IB}})]$	$\leftrightarrow \sigma_{\downarrow\text{-ord}}(L,D)$
1828	$\underbrace{\leftrightarrow \sigma_{\downarrow\text{-ord}}(L,1)}$		$\vdots$
1829	$\downarrow \mathfrak{a}_f^{<p_1^1, \dots, p_D^1   \dots   p_1^L, \dots, p_D^L >}$		
1830			
1831	(a) De-Composition Phase		
1832			
1833	$\downarrow \mathfrak{a}_f^{<p_1^1, \dots, p_D^1   \dots   p_1^L, \dots, p_D^L >} \xrightarrow{\tilde{f}} \uparrow \mathfrak{a}_f^{<p_1^1, \dots, p_D^1   \dots   p_1^L, \dots, p_D^L >}$		
1834	$\rightarrow < \sigma_{f\text{-ord}}(1,1), \dots, \sigma_{f\text{-ord}}(L,D) >$		
1835	$\rightarrow < \leftrightarrow_{f\text{-ass}}(1,1), \dots, \leftrightarrow_{f\text{-ord}}(L,D) >$		
1836	$\rightarrow \mathfrak{b}_1^{\text{IB}} : f^{\downarrow}\text{-mem}[\sigma_{f\downarrow\text{-mem}}(1), \dots, \sigma_{f\downarrow\text{-mem}}(D_1^{\text{IB}})], \dots, \mathfrak{b}_{B^{\text{IB}}}^{\text{IB}} : f^{\downarrow}\text{-mem}[\sigma_{f\downarrow\text{-mem}}(1), \dots, \sigma_{f\downarrow\text{-mem}}(D_{B^{\text{IB}}}^{\text{IB}})]$		
1837	$\rightarrow \mathfrak{b}_1^{\text{OB}} : f^{\uparrow}\text{-mem}[\sigma_{f\uparrow\text{-mem}}(1), \dots, \sigma_{f\uparrow\text{-mem}}(D_1^{\text{OB}})], \dots, \mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} : f^{\uparrow}\text{-mem}[\sigma_{f\uparrow\text{-mem}}(1), \dots, \sigma_{f\uparrow\text{-mem}}(D_{B^{\text{OB}}}^{\text{OB}})]$		
1838			
1839	(b) Scalar Phase		
1840			
1841	$\mathfrak{b}_1^{\text{OB}}, \dots, \mathfrak{b}_{B^{\text{OB}}}^{\text{OB}}$	$\xrightarrow{\text{out\_view}}$	$\uparrow \mathfrak{a} :=$
1842	$\oplus_1^{\leftrightarrow_{\text{prt-ass}}(1,1)}$	$\dots$	$\oplus_D^{\leftrightarrow_{\text{prt-ass}}(1,D)}$
1843	$p_1^1 \in \#PRT(1,1)$		$p_D^1 \in \#PRT(1,D)$
1844	$\rightarrow \mathfrak{b}_1^{\text{OB}} : \uparrow\text{-mem}(1,1)[\sigma_{\uparrow\text{-mem}}(1,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,1)(D_1^{\text{OB}})]$	$\rightarrow \mathfrak{b}_1^{\text{OB}} : \uparrow\text{-mem}(1,D)[\sigma_{\uparrow\text{-mem}}(1,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,D)(D_1^{\text{OB}})]$	
1845	$\vdots$	$\vdots$	$\vdots$
1846	$\mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(1,1)[\sigma_{\uparrow\text{-mem}}(1,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,1)(D_{B^{\text{OB}}}^{\text{OB}})]$	$\mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(1,D)[\sigma_{\uparrow\text{-mem}}(1,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(1,D)(D_{B^{\text{OB}}}^{\text{OB}})]$	$\leftrightarrow \sigma_{\uparrow\text{-ord}}(1,D)$
1847	$\underbrace{\leftrightarrow \sigma_{\uparrow\text{-ord}}(1,1)}$		$\vdots$
1848	$\vdots$		$\vdots$
1849	$\vdots$		$\vdots$
1850	$\oplus_1^{\leftrightarrow_{\text{prt-ass}}(L,1)}$	$\dots$	$\oplus_D^{\leftrightarrow_{\text{prt-ass}}(L,D)}$
1851	$p_1^L \in \#PRT(L,1)$		$p_D^L \in \#PRT(L,D)$
1852	$\rightarrow \mathfrak{b}_1^{\text{OB}} : \uparrow\text{-mem}(L,1)[\sigma_{\uparrow\text{-mem}}(L,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(L,1)(D_1^{\text{OB}})]$	$\rightarrow \mathfrak{b}_1^{\text{OB}} : \uparrow\text{-mem}(L,D)[\sigma_{\uparrow\text{-mem}}(L,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(L,D)(D_1^{\text{OB}})]$	
1853	$\vdots$	$\vdots$	$\vdots$
1854	$\mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(L,1)[\sigma_{\uparrow\text{-mem}}(L,1)(1), \dots, \sigma_{\uparrow\text{-mem}}(L,1)(D_{B^{\text{OB}}}^{\text{OB}})]$	$\mathfrak{b}_{B^{\text{OB}}}^{\text{OB}} : \uparrow\text{-mem}(L,D)[\sigma_{\uparrow\text{-mem}}(L,D)(1), \dots, \sigma_{\uparrow\text{-mem}}(L,D)(D_{B^{\text{OB}}}^{\text{OB}})]$	$\leftrightarrow \sigma_{\uparrow\text{-ord}}(L,D)$
1855	$\underbrace{\leftrightarrow \sigma_{\uparrow\text{-ord}}(L,1)}$		$\vdots$
1856	$\uparrow \mathfrak{a}_f^{<p_1^1, \dots, p_D^1   \dots   p_1^L, \dots, p_D^L >}$		
1857			
1858	(c) Re-Composition Phase		

Fig. 19. Generic low-level expression for data-parallel computations

No.	Name	Range	Description
1863	0	#PRT	MDH-LVL $\rightarrow \mathbb{N}$
1864			number of parts
1865	D1	$\sigma_{\downarrow\text{-ord}}$	MDH-LVL $\leftrightarrow$ MDH-LVL
1866	D2	$\leftrightarrow_{\downarrow\text{-ass}}$	MDH-LVL $\leftrightarrow$ ASM-LVL
1867	D3	$\downarrow\text{-mem}^{<\text{ib}>}$	MDH-LVL $\rightarrow$ MR
1868	D4	$\sigma_{\downarrow\text{-mem}}^{<\text{ib}>}$	memory regions of input BUFs (ib)
1869			memory layouts of input BUFs (ib)
1870	S1	$\sigma_f\text{-ord}$	MDH-LVL $\leftrightarrow$ MDH-LVL
1871	S2	$\leftrightarrow_f\text{-ass}$	MDH-LVL $\leftrightarrow$ ASM-LVL
1872	S3	$f\downarrow\text{-mem}^{<\text{ib}>}$	MR
1873	S4	$\sigma_{f\downarrow\text{-mem}}^{<\text{ib}>}$	memory region of input BUF (ib)
1874			memory layout of input BUF (ib)
1875	S5	$f\uparrow\text{-mem}^{<\text{ob}>}$	MR
1876	S6	$\sigma_{f\uparrow\text{-mem}}^{<\text{ob}>}$	memory region of output BUF (ob)
1877			memory layout of output BUF (ob)
1878	R1	$\sigma_{\uparrow\text{-ord}}$	MDH-LVL $\leftrightarrow$ MDH-LVL
1879	R2	$\leftrightarrow_{\uparrow\text{-ass}}$	MDH-LVL $\leftrightarrow$ ASM-LVL
1880	R3	$\uparrow\text{-mem}^{<\text{ob}>}$	MDH-LVL $\rightarrow$ MR
1881	R4	$\sigma_{\uparrow\text{-mem}}^{<\text{ob}>}$	memory regions of output BUFs (ob)
1882			memory layouts of output BUFs (ob)
1883			
1884			
1885			
1886			
1887			
1888			
1889			
1890			
1891			
1892			
1893			
1894			

Table 1. Tuning parameters of our low-level expressions

### 3.4 Generic Low-Level Expression

Figure 19 shows a generic expression in our low-level representation: it targets an arbitrary but fixed  $L$ -layered ASM instance, and it implements – on low level – the generic instance of our high-level expression in Figure 15. Inserting into the low-level expression a particular value for ASM’s numbers of layer  $L$ , as well as particular values for the generic parameters of the high-level expression in Figure 15 (dimensionality  $D$ , combine operators  $\otimes_1, \dots, \otimes_d$ , and input/output views) results in an instance of the expression in Figure 19 that remains generic in tuning parameters only; this auto-tunable instance will be the focus of our discussion in the remainder of this section.

In Section 4, we show that we fully automatically compute the auto-tunable low-level expression for a concrete ASM instance and high-level expression, and we automatically optimize this tunable

<sup>17</sup> For simplicity, we refrain from annotating identifier ASM-LVL with values  $L$  and  $D$  (e.g.,  $\text{ASM-LVL}^{<L,D>}$ ), because both values will usually be clear from the context.

1912 expression for a particular target device and characteristics of the input and output data using  
 1913 auto-tuning [Rasch et al. 2021]. Our final outcome is a concrete (non-generic) low-level expression  
 1914 (as in Figure 17) that is auto-tuned for the particular target device (represented via an ASM  
 1915 instance, e.g., ASM instance  $\text{ASM}_{\text{CUDA}}$  when targeting an NVIDIA Ampere GPU) and high-level  
 1916 MDH expression; from this auto-tuned low-level expression, we can straightforwardly generate  
 1917 executable program code, because all the major optimization decisions have already been made in  
 1918 the previous auto-tuning step. Our overall approach is illustrated in Figure 4.

1919 *Auto-Tunable Parameters.* Table 1 lists the tuning parameters of our auto-tunable low-level expres-  
 1920 sions. Different values of tuning parameters lead to semantically equal variants of the auto-tunable  
 1921 low-level expression (which we prove formally in Section 4), but the variants will be translated to  
 1922 differently optimized code variants.  
 1923

1924 In the following, we explain the 15 tuning parameters in Table 1. We give our explanations in a  
 1925 general, formal setting that is independent of a particular computation and programming model; the  
 1926 parameters are discussed afterwards for the concrete example computation *matrix multiplication* in  
 1927 models OpenMP, CUDA, and OpenCL.

1928 In Table 1, we refer to combine operators in Figure 17 using pairs  $(l, d)$  to which we refer to as  
 1929 *MDH Levels* (MDH-LVL) (terminology motivated in the Appendix, Section C.6). We use the pairs as  
 1930 enumeration for operators in the de-composition or re-composition phase, respectively.

1931 **Definition 16** (MDH Level). We define set MDH-LVL as:

$$\text{MDH-LVL} := \{ (l, d) \mid l \in [1, L]_{\mathbb{N}}, d \in [1, D]_{\mathbb{N}} \}^{18}$$

1934 For example, in the de-composition phase of Figure 17 (right part of the figure), pair  $(1, 1) \in$   
 1935 MDH-LVL points to the first combine operator, as the operator operates on the first layer  $l = 1$  and  
 1936 in the first dimension  $d = 1$  (discussed in Section 3.1). Analogously, pairs  $(1, 2), (2, 1) \in$  MDH-LVL  
 1937 point to the second and third operator, etc. An operator's enumeration can be easily deduced from  
 1938 its corresponding  $p$  variable: the variable's superscript indicates the operator's corresponding layer  
 1939  $l$  and the variable's subscript indicates its dimension  $d$ .  
 1940

1941 Our tuning parameters in Table 1 have constraints: 1) *algorithmic constraints* which have to  
 1942 be satisfied by all target programming models; 2) *model constraints* which are specific for par-  
 1943 ticular programming models only (CUDA-specific constraints, OpenCL-specific constraints, etc),  
 1944 e.g., that the results of CUDA's thread blocks can only be combined in designated memory re-  
 1945 gions [NVIDIA 2022f]. We discuss algorithmic constraints in the following, together with our  
 1946 tuning parameters; model constraints are discussed in our Appendix, Section C.1, for the interested  
 1947 reader.

1948 In the following, we present our 15 tuning parameters in Table 1. Dotted lines separate parameters  
 1949 for different phases: parameters D1-D4 customize the de-composition phase, parameters S1-S6 the  
 1950 scalar phase, and parameters R1-R4 the re-composition phase, correspondingly; the parameter  $\theta$   
 1951 impacts all three phases (separated by a straight line in the table).

1952 Note that our parameters do not aim to introduce novel optimization techniques, but to unify,  
 1953 generalize, and combine together well-proven optimizations, based on a formal foundation, toward  
 1954 an efficient, overall optimization process that applies to various combinations of data-parallel  
 1955 computations, architectures, and characteristics of input and output data (e.g., their size and  
 1956 memory layout).

1957 <sup>18</sup> The same as for identifier  $\text{ASM-LVL}$ , we refrain from annotating identifier  $\text{MDH-LVL}$  with values  $L$  and  $D$ . Note that  
 1958  $\text{MDH-LVL}$  and  $\text{ASM-LVL}$  (Definition 15) both refer to the same set of pairs, but we use identifier  $\text{MDH-LVL}$  when referring to  
 1959 MDH levels and identifier  $\text{ASM-LVL}$  when referring to ASM levels, correspondingly, for better clarity.

1961 *Parameter 0*: Parameter #PRT is a function that maps pairs in MDH-LVL to natural numbers; the  
 1962 parameter determines *how much* data are grouped together into parts in our low-level expression  
 1963 in Figure 19 (and consequently also in our generated code later) by setting the particular number  
 1964 of parts (a.k.a. *tiles*) used in our expression. For example, in Figure 17, we use #PRT(1, 1) := 2  
 1965 which causes combine operators  $+_1^{(\text{HM}, \times)}$  and  $\otimes_1^{(\text{HM}, \times)}$  to iterate over interval  $[0, 2]_{\mathbb{N}_0}$ , and we use  
 1966 #PRT(1, 2) := 4 to let operators  $+_2^{(\text{HM}, y)}$  and  $\otimes_2^{(\text{HM}, \times)}$  iterate over interval  $[0, 4]_{\mathbb{N}_0}$ , etc.  
 1967

1968 To ensure a full partitioning (so that we obtain singleton MDAs in the scalar phase to which scalar  
 1969 function  $f$  can be applied, as discussed above), we require the following algorithmic constraint for  
 the parameter:  
 1970

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

1972 Here,  $N_d$  is the input size in dimension  $d$  (Figure 15).  
 1973

1974 In our generated code, the number of parts directly translates to the number of *tiles* which are  
 1975 computed either sequentially (a.k.a. *cache blocking* [Lam et al. 1991]) or in parallel in our code,  
 1976 depending on the combine operators's tags (which are chosen via Parameters D2, S2, R2 as discussed  
 1977 soon). In our example from Figure 17, we process the parts belonging to combine operators tagged  
 1978 with HM and L1 sequentially, via for-loops, because HM and L1 correspond to ASM's memory layers  
 1979 (note that Parameter 0 only chooses the number of tiles; the parameter has no effect on explicitly  
 1980 copying data into fast memory resources, which is the purpose of Parameters D3, R3, S1, S2); COR  
 1981 parts are computed in parallel in our code, because COR corresponds to ASM's core layer, and thus,  
 1982 the number of COR parts determines the number of threads used in our generated code.  
 1983

1984 An optimized number of tiles is known to be essential for high performance [Bacon et al. 1994],  
 1985 e.g., due to its impact for locality-aware data accesses (number of sequentially computed tiles) and  
 1986 efficiently exploiting parallelism (number of tiles computed in parallel, which corresponds to the  
 1987 number of threads used in our generated code).  
 1988

1989 *Parameters D1, S1, R1*: These three parameters are permutations on MDH-LVL, determining *when*  
 1990 data is accessed and combined. The parameters specify the order (indicated by symbol  $\hookrightarrow$  in  
 1991 Figure 19) of combine operators in the de-composition and re-composition phases (parameters D1  
 1992 and R1), and the order of applying scalar function  $f$  to parts (parameter S1). Thereby, the parameters  
 1993 specify when parts are processed during computation.  
 1994

1995 In our generated code, combine operators are implemented as sequential/parallel loops such that  
 1996 the parameters enable optimizing loop orders (a.k.a. *loop permutation* [McKinley et al. 1996]). For  
 1997 combine operators assigned (via parameter R2) to ASM's core layer and thus computed in parallel,  
 1998 parameter R1 particularly determines when the computed results of threads are combined: if we  
 1999 used in the re-composition phase of Figure 17 combine operators tagged with (COR,  $\times$ ) and (COR,  $y$ )  
 2000 immediately after applying scalar function  $f$  (i.e., in steps ⑩ and ⑪, rather than steps ⑫ and ⑬), we would  
 2001 combine the computed intermediate results of threads multiple times, repeatedly  
 2002 after each individual computation step of threads, and using the two operators at the end of the  
 2003 re-composition phase (in steps ⑭ and ⑮) would combine the result of threads only once, at the  
 2004 end of the re-composition phase.  
 2005

2006 Note that each phase corresponds to an individual loop nest which we fuse together when  
 2007 parameters D1, S1, R1 (as well as parameters D2, S2, R2) coincide.  
 2008

2009 *Parameters D2, S2, R2*: These parameters (symbol  $\leftrightarrow$  in the table denotes bijection) map MDH  
 2010 levels to ASM levels, thereby determining *which* data is assigned to which ASM level (similar to  
 2011 the concept of bind in scheduling languages [Apache TVM Documentation 2022a]), by setting  
 2012

2010 the operators' tags. Consequently, the parameters determine which parts should be computed  
 2011 sequentially in our generated code and which parts in parallel. For example, in Figure 17, we use  
 2012  $\leftrightarrow_{\text{ass}} (2, 1) := (\text{COR}, x)$  and  $\leftrightarrow_{\text{ass}} (2, 2) := (\text{COR}, y)$ , thereby assigning the computation of MDA  
 2013 parts on layer 2 in both dimensions to ASM's COR layers in the de-composition phase, which causes  
 2014 processing the parts in parallel in our generated code. In CUDA-like approaches which have a  
 2015 notion of cores on multiple layers, the parameter particularly determines the thread layer to be  
 2016 used for parts computed in parallel (e.g., block or thread in CUDA).

2017 Using these parameters, we are able to flexibly set data access patterns in our generated code. In  
 2018 Figure 17, we assign parts on layer 2 to COR layers, which results in a so-called block access pattern  
 2019 of cores: we start  $8 * 16$  threads, according to the  $8 * 16$  core parts, and each thread processes a  
 2020 part of the input MDA representing a block of  $32 * 64$  MDA elements within the input data. If  
 2021 we had assigned in the figure the first computation layer to ASM's COR layer (in the figure, this  
 2022 layer is assigned to ASM's HM layer), we would start  $2 * 4$  threads and each thread would process  
 2023 MDA parts of size  $(8 * 32) * (16 * 64)$ ; assigning the last MDH layer to CORs would result in  
 2024  $(2 * 8 * 32) * (4 * 16 * 64)$  threads each processing MDA parts containing a single MDA element  
 2025 only (a.k.a. *strided access*).

2026 The parameters also enable expressing so-called *swizzle* access patterns [Phothilimthana et al.  
 2027 2019]. For example, in CUDA, processing consecutive data elements in data dimension 1 by threads  
 2028 that are consecutive in thread dimension 2 (a.k.a `threadIdx.y` dimension in CUDA) can achieve  
 2029 higher performance due to the hardware design of fast memory resources in NVIDIA GPUs.  
 2030 Such swizzle patterns can be easily expressed and auto-tuned in our approach; for example, by  
 2031 interchanging in Figure 17 tags  $(\text{COR}, x)$  and  $(\text{COR}, y)$ . For memory layers (such as HM and L1),  
 2032 the dimension tags  $x$  and  $y$  currently have no effect on our generated code, as the programming  
 2033 models we target at the moment (OpenMP, CUDA, and OpenCL) have no explicit notion of tiles.  
 2034 However, this might change in the future when targeting new kinds of programming models, e.g.,  
 2035 for upcoming architectures.

2036

2037

2038 *Parameters D3, R3 and S3, S5:* Parameters D3 and R3 set for each BUF the memory region to  
 2039 use, thereby determining *where* data is read from or written to, respectively. In the table, we use  
 2040  $ib \in [1, B^{\text{IB}}]_{\mathbb{N}}$  to refer to a particular input BUF and  $ob \in [1, B^{\text{OB}}]_{\mathbb{N}}$  to refer to an output BUF,  
 2041 correspondingly. Parameter D3 specifies the memory region to read from, and parameter R3 the  
 2042 regions to write to. The set  $MR := [1, \text{NUM\_MEM\_LYRs}]_{\mathbb{N}}$  denotes the ASM's memory regions.

2043 Similarly to parameters D3 and R3, parameters S3 and S5 set for scalar function  $f$  the memory  
 2044 regions for the input and output scalar.

2045 Exploiting fast memory resources of architectures is a fundamental optimization [Bondhugula  
 2046 2020; Hristea et al. 1997; Mei et al. 2014; Salvador Rohwedder et al. 2023], particularly due to the  
 2047 performance gap between processors' cores and their memory system [Oliveira et al. 2021; Wilkes  
 2048 2001].

2049

2050

2051 *Parameters D4, R4 and S4, S6:* These parameters set the memory layouts of BUFs, thereby de-  
 2052 termining *how* data is accessed in memory; for brevity in Table 1, we denote the set of all BUF  
 2053 permutations  $[1, D]_{\mathbb{N}} \leftrightarrow [1, D]_{\mathbb{N}}$  (Definition 13) as  $[1, \dots, D]_{\mathcal{S}}$  (symbol  $\mathcal{S}$  is taken from the notation  
 2054 of *symmetric groups* [Sagan 2001]). In the case of our matrix-vector multiplication example in  
 2055 Figure 17, we use a standard memory layout for all matrices, which we express via the parameters  
 2056 by setting them to the identity function, e.g.,  $\sigma_{\text{-mem}}^{<\text{MDH}>} (1, 1) := id$  (Parameter D4) for the matrix read  
 2057 by operator  $+_1^{(\text{HM}, x)}$ .

2058

2059 An optimized memory layout is important to access data in a locality-aware and thus efficient  
2060 manner.

2061

### 2062 3.5 Examples

2063 Figures 20-23 show how our low-level representation is used for expressing the (de/re)-compositions  
2064 of concrete, state-of-the-art implementations, using the popular example of matrix multiplication  
2065 (MatMul) on a real-world input size taken from the ResNet-50 [He et al. 2015] deep learning  
2066 neural network (training phase). To challenge our formalism: i) we express implementations gener-  
2067 ated and optimized according to notably different approaches: scheduling approach TVM using its  
2068 recent Ansor [Zheng et al. 2020a] optimization engine which is specifically designed and optimized  
2069 toward deep learning computations; polyhedral compilers PPCG and Pluto with auto-tuned tile  
2070 sizes; ii) we consider optimizations for two fundamentally different kinds of architectures: NVIDIA  
2071 Ampere GPU and Intel Skylake CPU. We consider our study as challenging for our formalism,  
2072 because it needs to express – in the same formal framework – the (de/re)-compositions of imple-  
2073 mentations generated and optimized according to notably different approaches (scheduling-based  
2074 and polyhedral-based) and for significantly different kinds of architectures (GPU and CPU). Experi-  
2075 mental results for TVM, PPCG, and Pluto (including the MatMul study used in this section) are  
2076 presented and discussed in Section 5, as the focus of this section is on analyzing and discussing the  
2077 expressivity of our low-level representation, rather than its experimental evaluation.

2078 In Figures 20-23, we list our low-level representation’s particular tuning parameter values for  
2079 expressing the TVM- and PPCG/Pluto-generated implementations, which concisely describe the  
2080 (de/re)-composition strategies used by TVM, PPCG and Pluto for MatMul on GPU or CPU for  
2081 the ResNet-50’s input size; inserting these tuning parameter values into our generic low-level  
2082 expression in Figure 19 results in the formal representation of the (de/re)-composition strategies  
2083 used by TVM, PPCG and Pluto.

2084 In the following, we describe the columns of the tables in Figures 20-23, each of which listing  
2085 particular values of tuning parameters in Table 1: column 0 in Figures 20-23 lists values of tuning  
2086 parameter 0 in Table 1, column D1 of tuning parameter D1, etc. As all four tables follow the  
2087 same structure, we focus on describing the particular example table in Figure 20 (example chosen  
2088 arbitrarily), which shows the (de/re)-composition used in TVM’s generated CUDA code for MatMul  
2089 on NVIDIA Ampere GPU using input matrices of sizes  $16 \times 2048$  and  $2048 \times 1000$  taken from  
2090 ResNet-50<sup>19</sup>. Note that for clarity, we use in the figures domain-specific aliases, instead of numerical  
2091 values, to better differentiate between different ASM layers and memory regions. For example,  
2092 we use in Figure 20 as aliases DEV := 1, SHR := 2, and REG := 3 to refer to CUDA’s three memory  
2093 layers (device memory layer DEV, shared memory layer SHR, and register memory layer REG), and  
2094 we use DM := 1, SM := 2, and RM := 3 to refer to CUDA’s memory regions device DM, shared SM, and  
2095 register RM; aliases BLK := 4 and THR := 5 refer to CUDA’s core layers which are programmed via  
2096 blocks and threads in CUDA. We differentiate between memory layers and memory regions for  
2097 the following reason: for example, using tuning parameter 0 in Table 1, we partition input data  
2098 hierarchically for each particular memory layer of the target architecture (sometime possibly into  
2099 one part only, which is equivalent to not partitioning). However, depending on the value of tuning  
2100 parameter D3, we do not necessarily copy the input’s parts always into the corresponding memory  
2101 regions (e.g., a part on SHR layer is not necessarily copied into shared memory SM), for example,  
2102 when the architecture provides automatically managed memory regions (as caches in CPUs) or  
2103

2104  
2105 <sup>19</sup> For the interested reader, TVM’s corresponding, Ansor-generated scheduling program is presented in our Appendix,  
2106 Section C.8.

2108 when only some parts of the input are accessed multiple times (e.g., the input vector in the case of  
 2109 matrix-vector multiplication, but not the input matrix), etc.

2110 *Column 0.* The column lists the particular number of parts (a.k.a. *tiles*) used in TVM’s multi-  
 2111 layered, multi-dimensional partitioning strategy for MatMul on the ResNet-50’s input matrices of  
 2112 sizes  $16 \times 2048$  and  $2048 \times 1000$ . We can observe from this column that the input MDA, which is  
 2113 initially of size  $(16, 1000, 2048)$  for the ResNet-50’s input matrices (follows from MatMul’s particular  
 2114 input view functions shown in Figure 16), is partitioned into  $(2 * 50 * 1)$ -many parts (indicated by  
 2115 the first three rows in column 1): 2 parts in the first dimension, 50 parts in the second dimension,  
 2116 and 1 part in the third dimension. Each of these parts is then further partitioned into  $(2 * 1 * 8)$ -many  
 2117 parts (rows 4,5,6), and these parts are again partitioned into  $(4 * 20 * 1)$ -many further parts (rows  
 2118 7,8,9), etc.

2119  
 2120 *Columns D1, S1, R1.* These 3 columns describe the order in which parts are processed in the  
 2121 different phases: de-composition (column D1), scalar phase (column S1), and re-composition (col-  
 2122 umn R1). For example, we can observe from column R1 that TVM’s CUDA code first combines parts  
 2123 on layer 1 in dimensions 1, 2, 3 (indicated via 1, 2, 3 in rows 1, 2, 3 of column R1); afterwards, parts  
 2124 on layer 3 in dimensions 1, 2, 3 are processed (indicated via 1, 2, 3 in rows 7, 8, 9 of column R1), etc.

2125 Note that TVM uses the same order in all three phases (i.e., columns D1, S1, R1 coincide). Most  
 2126 likely this is because in CUDA, iteration over memory tiles are programmed via for-loops such that  
 2127 columns D1, S1, R1 represent loop orders – using the same order of loops in columns D1, S1, R1  
 2128 consequently allows TVM to generate the loop nests of the three different phases as fused, in one  
 2129 single loop nest (rather than three individual nests), which usually achieves high performance in  
 2130 CUDA.

2131 Note further that the order of parts that are processed in parallel (columns D2, S2, R2, described  
 2132 in the next paragraph, determine if parts are processed in parallel or not) effects when results of  
 2133 blocks and threads are combined (a.k.a. *parallel reduction* [Harris et al. 2007]), e.g., early in the  
 2134 computation and thus often (but probably at the cost of low memory consumption) or late and thus  
 2135 only once (but requiring possibly more memory), etc.

2136  
 2137 *Columns D2, S2, R2.* The columns determine how computations are assigned to the target archi-  
 2138 tecture. In our example in Figure 20, the  $(2 * 50 * 1)$ -many parts in the first partitioning layer (rows  
 2139 1-3) are assigned to CUDA blocks (BLK); consequently, each such part is computed by an individual  
 2140 block (i.e., TVM uses a so-called *grid size* of 2, 50, 1 in its generated CUDA code for MatMul). The  
 2141  $(4 * 20 * 1)$ -many parts in the third partitioning layer (rows 7-9) are processed by CUDA threads  
 2142 (THR) (i.e., the CUDA *block size* in the TVM-generated code is 4, 20, 1). All other parts, e.g., those  
 2143 belonging to the  $(2 * 1 * 8)$ -many parts in the second partitioning layer (rows 4-6), are assigned to  
 2144 CUDA’s memory layers (denoted as DEV, SHR, REG in Figure 20) and thus processed sequentially,  
 2145 via for-loops.

2146  
 2147 *Columns D3, S3, S5, R3.* While column 0 sets the multi-layered, multi-dimensional partitioning  
 2148 strategy used in TVM’s CUDA code (according to the CUDA model’s multi-layered memory and  
 2149 core hierarchies, shown in Example 11), column 0 does not indicate how CUDA’s fast memory  
 2150 regions are exploited in the TVM-generated CUDA code for MatMul – column 0 only describes  
 2151 TVM’s partitioning of the input/output computations such that parts of the input and output data  
 2152 can potentially fit into fast memory resources.

2153 The actual assignment of parts to memory regions is set via columns D3 (memory regions to be  
 2154 used for input data), columns S3 and S5 (memory regions to be used for the scalar computations),  
 2155 and column R3 (memory regions to be used for storing the computed output data). For example,  
 2156 column D3 indicates that in TVM’s CUDA code for MatMul, parts of the *A* and *B* input matrices are

2157 stored in CUDA's fast shared memory SM (rows 10-15), and column R3 indicates that each thread  
 2158 computes its results within CUDA's faster register memory RM (rows 10-15).

2159 Our flexibility of separating the tiling strategy (Parameter 0) from the actual usage of memory  
 2160 regions (columns D3, S3, S5, R3) allows us, for example, to store parts belonging to one input buffer  
 2161 into fast memory resources (e.g., the input vector of matrix-vector multiplication, whose values  
 2162 are accessed multiple times), but not parts of other buffers (e.g., the input matrix of matrix-vector  
 2163 multiplication, whose values are accessed only once) or only subparts of buffers, etc.

2164 Note that in the case of Figure 23 which shows Pluto's (de/re)-composition for OpenMP code,  
 2165 the memory tags in columns D3, S3, S5, R3 have no effect on the generated code: OpenMP relies  
 2166 on its target device's implicit memory management mechanism (e.g., CPU caches), rather than  
 2167 exposing the memory hierarchy explicitly to the programmer. Consequently, the memory tags are  
 2168 ignored by our OpenMP code generator and only emphasize the implementer's intention, e.g., that  
 2169 in Figure 23, each of the  $(2 * 962 * 218)$ -many tiles (rows 7-9) are intended by the implementer  
 2170 to be processed in L2 memory (even though this decision is eventually made automatically in the  
 2171 CPU hardware's automatic cache engine).

2172 *Columns D4, S4, S6, R4.* These columns set the memory layout to be used for memory allocations  
 2173 in the CUDA code. TVM chooses in all cases CUDA's standard transpositions layout (indicated by  
 2174  $[1, 2]$ , called *row-major* layout, instead of  $[2, 1]$  which is known as *column-major*). Since the  
 2175 same layout and memory region is used on consecutive layers, the same memory allocation is  
 2176 re-used in the CUDA code. For example, only one memory region is allocated in shared memory  
 2177 for input buffer A; the region is accessed in the computations of all SHR, REG tiles as well as CC tiles  
 2178 in dimensions x and z. Similarly, only one region is allocated in register memory for computing  
 2179 the results of SHR, REG, and DEV tiles.

TVM's (de/re)-composition for MatMul in CUDA on NVIDIA Ampere GPU															
Part.	De-Comp. Phase				Scalar Phase						Re-Comp. Phase				
	0	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	C	C	
2188	2	1	BLK,y	DM	DM	[1,2]	1	BLK,y				1	BLK,y	DM	[1,2]
2189	50	2	BLK,x	DM	DM	[1,2]	2	BLK,x				2	BLK,x	DM	[1,2]
2190	1	3	BLK,z	DM	DM	[1,2]	3	BLK,z				3	BLK,z	DM	[1,2]
2191	2	14	DEV,x	SM	SM	[1,2]	14	DEV,x				14	DEV,x	RM	[1,2]
2192	1	15	DEV,y	SM	SM	[1,2]	15	DEV,y				15	DEV,y	RM	[1,2]
2193	8	7	DEV,z	DM	DM	[1,2]	7	DEV,z				7	DEV,z	RM	[1,2]
2194	4	4	THR,y	DM	DM	[1,2]	4	THR,y				4	THR,y	DM	[1,2]
2195	20	5	THR,x	DM	DM	[1,2]	5	THR,x	RM	RM	[1,2]	5	THR,x	DM	[1,2]
2196	1	6	THR,z	DM	DM	[1,2]	6	THR,z				6	THR,z	DM	[1,2]
2197	1	9	SHR,x	SM	SM	[1,2]	9	SHR,x				9	SHR,x	RM	[1,2]
2198	1	10	SHR,y	SM	SM	[1,2]	10	SHR,y				10	SHR,y	RM	[1,2]
2199	128	8	SHR,z	SM	SM	[1,2]	8	SHR,z				8	SHR,z	RM	[1,2]
2200	1	12	REG,x	SM	SM	[1,2]	12	REG,x				12	REG,x	RM	[1,2]
2201	1	13	REG,y	SM	SM	[1,2]	13	REG,y				13	REG,y	RM	[1,2]
2202	2	11	REG,z	SM	SM	[1,2]	11	REG,z				11	REG,z	RM	[1,2]

2203 Fig. 20. 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															
Part.	De-Comp. Phase				Scalar Phase						Re-Comp. Phase				
	0	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	
1	1	WG,1	GM	GM	[1,2]	1	WG,1					1	WG,1	GM	[1,2]
125	2	WG,0	GM	GM	[1,2]	2	WG,0					2	WG,0	GM	[1,2]
1	3	WG,2	GM	GM	[1,2]	3	WG,2					3	WG,2	GM	[1,2]
1	14	GLB,0	LM	LM	[1,2]	14	GLB,0					14	GLB,0	PM	[1,2]
1	15	GLB,1	LM	LM	[1,2]	15	GLB,1					15	GLB,1	PM	[1,2]
16	7	GLB,2	GM	GM	[1,2]	7	GLB,2					7	GLB,2	PM	[1,2]
16	4	WI,1	GM	GM	[1,2]	4	WI,1					4	WI,1	GM	[1,2]
1	5	WI,0	GM	GM	[1,2]	5	WI,0	PM	PM	[1,2]	PM	5	WI,0	GM	[1,2]
1	6	WI,2	GM	GM	[1,2]	6	WI,2					6	WI,2	GM	[1,2]
1	9	LCL,0	LM	LM	[1,2]	9	LCL,0					9	LCL,0	PM	[1,2]
1	10	LCL,1	LM	LM	[1,2]	10	LCL,1					10	LCL,1	PM	[1,2]
128	8	LCL,2	LM	LM	[1,2]	8	LCL,2					8	LCL,2	PM	[1,2]
1	12	PRV,0	LM	LM	[1,2]	12	PRV,0					12	PRV,0	PM	[1,2]
8	13	PRV,1	LM	LM	[1,2]	13	PRV,1					13	PRV,1	PM	[1,2]
1	11	PRV,2	LM	LM	[1,2]	11	PRV,2					11	PRV,2	PM	[1,2]

Fig. 21. 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				
	0	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	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. 22. 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				
	0	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	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. 23. Pluto's (de/re)-composition for MatMul in OpenMP on CPU expressed in our low-level representation

## 4 LOWERING: FROM HIGH LEVEL TO LOW LEVEL

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 17). We confirm this by parameterizing the high-level expression in Figure 15 – 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 19.

PARAMETER 0. According to Definition 12, we use the  $L$ -layered,  $D$ -dimensional,  $P$ -partitioned low-level representation  ${}^l\alpha_f$  of  $\alpha$  for

$$P := \left( \underbrace{(\#PRT(1,1), \dots, \#PRT(1,D))}_{\text{Dimension 1}} , \dots , \underbrace{(\#PRT(L,1), \dots, \#PRT(L,D))}_{\text{Dimension } D} \right) \underbrace{\quad}_{\text{Layer 1}} \quad \underbrace{\quad}_{\text{Layer } L}$$

where  $\#PRT$  denotes the number of partitions (Parameter 0 in Table 1).

Applying the homomorphic property (Definition 4)  $L * D$  times, we get:

$$\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(a) =$$

$$\underbrace{p_1^1 \in \#PRT(1,1)}_{\text{Dimension 1}} \dots \underbrace{p_D^1 \in \#PRT(1,D)}_{\text{Dimension } D} \dots \underbrace{p_1^L \in \#PRT(L,1)}_{\text{Dimension 1}} \dots \underbrace{p_D^L \in \#PRT(L,D)}_{\text{Dimension } D}$$

$$\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))({}^l\alpha_f^{p_1^1, \dots, p_D^1} \mid \dots \mid p_1^L, \dots, p_D^L)$$

2304 where  $\downarrow a_f^{<p_1^1, \dots, p_D^1 | \dots | p_1^L, \dots, p_D^L>} \text{ is defined as:}$

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 1, discussed in Section 3.4), it holds

$$\text{md\_hom}(f, (\otimes_1, \dots, \otimes_D))(\downarrow_{\mathfrak{a}_f^{< p_1^1, \dots, p_D^1}}}^{\downarrow \mathfrak{a}_f^{< p_1^1, \dots, p_D^1} \mid \dots \mid p_1^L, \dots, p_D^L}) = \vec{f}(\downarrow_{\mathfrak{a}_f^{< p_1^1, \dots, p_D^1}}}^{\downarrow \mathfrak{a}_f^{< p_1^1, \dots, p_D^1} \mid \dots \mid p_1^L, \dots, p_D^L})$$

for  $\vec{f}$  defined as in Definition 4.

2318  
 2319     PARAMETER D1. Parameter D1 reorders concatenation operators  $\text{++}_1, \dots, \text{++}_D$  (Example 1) which  
 2320     are of types  $\text{++}_d \in \text{CO}^{< d \mid T \mid D \mid d >}$  for arbitrary but fixed  $T \in \text{TYPE}$  and  $D \in \mathbb{N}$ . We prove our assumption  
 2321     w.l.o.g. for the case  $D = 2$ ; the general case  $D \in \mathbb{N}$  follows analogously.

Let  $+_d \in \text{CO}^{<id|T|D|d_1>}$  and  $+_d \in \text{CO}^{<id|T|D|d_2>}$  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 1.

PARAMETERS D2, S2, R2. These parameters replaces combine operators (Definition 2) by low-level combine operators (Definition 15) which has no effect on semantics.  $\square$

PARAMETERS D3, S3, S5, R3. These parameters set the memory tags of low-level BUFs (Definition 13), which have no effect on semantics.  $\square$

PARAMETERS D4, S4, S6, R4. The parameters change the memory layout of low-level BUFs (Definition 13), which does not affect extensional equality.  $\square$

PARAMETERS S1. This parameter sets the order in which function  $f$  is applied to parts, which is trivially sound for any order.  $\square$

PARAMETERS R1. Similarly to parameter D1, parameter R1 reorders combine operators  $\otimes_1, \dots, \otimes_d$ .

PARAMETERS R1. Similarly to parameter D1, parameter R1 reorders combine operators  $\otimes_1, \dots, \otimes_d$ .

PARAMETERS R1. Similarly to parameter D1, parameter R1 reorders combine operators  $\otimes_1, \dots, \otimes_D$ . In contrast to parameter D1, the combine operators reordered by parameter R1 are not restricted to be concatenation.

We prove our assumption by exploiting the MDH property (Definition 3) together with the proof of parameter D1, as follows:

```

2345
2346
2347 = md_hom(...)( (a1 ⊗d1 a2) ⊗d2 (a3 ⊗d1 a4) )
2348 = md_hom(...)( (a1 ++d1 a2) ++d2 (a3 ++d1 a4) )
2349 = (a1 ⊗d1 a3) ⊗d2 (a2 ⊗d1 a4) ✓

```

## 2353 5 EXPERIMENTAL RESULTS

2354 All experiments described in this section can be reproduced using [TOPLAS Artifact 2022].

2355 We experimentally evaluate our approach by comparing it to popular representatives of four  
2356 important classes:

- 2357 (1) *scheduling approach*: TVM [Chen et al. 2018a] which generates GPU and CPU code from  
2358 programs expressed in TVM’s high-level program representation;
- 2359 (2) *polyhedral compilers*: PPCG [Verdoolaege et al. 2013] for GPUs<sup>20</sup> and Pluto [Bondhugula et al.  
2360 2008b] for CPUs, which automatically generate executable program code in CUDA (PPCG)  
2361 or OpenMP (Pluto) from straightforward, unoptimized C programs;
- 2362 (3) *functional approach*: Lift [Steuwer et al. 2015] which generates OpenCL code from a Lift-  
2363 specific, functional program representation;
- 2364 (4) *domain-specific libraries*: NVIDIA cuBLAS [NVIDIA 2022b] and NVIDIA cuDNN [NVIDIA  
2365 2022e], as well as Intel oneMKL [Intel 2022c] and Intel oneDNN [Intel 2022b], which offer the  
2366 user easy-to-use, domain-specific building blocks for programming. The libraries internally  
2367 rely on pre-implemented assembly code that is optimized by experts for their target applica-  
2368 tion domains: linear algebra (cuBLAS and oneMKL) or convolutions (cuDNN and oneDNN),  
2369 respectively. To make comparison against the libraries challenging for us, we compare to  
2370 all routines provided by the libraries. For example, the cuBLAS library offers three, seman-  
2371 tically equal but differently optimized routines for computing MatMul: cublasSgemm (the  
2372 default MatMul implementation in cuBLAS), cublasGemmEx which is part of the cuBLASEx  
2373 extension of cuBLAS [NVIDIA 2022c], and the most recent cublasLtMatmul which is part  
2374 of the cuBLASLt extension [NVIDIA 2022d]; each of these three routines may perform dif-  
2375 ferently on different problem sizes (NVIDIA usually recommends to naively test which one  
2376 performs best for the particular target problem). To make comparison further challenging  
2377 for us, we exhaustively test for each routine all of its so-called cublasGemmAlgo\_t variants,  
2378 and report the routine’s runtime for the best performing variant only. In the case of oneMKL,  
2379 we compare also to its *JIT engine* [Intel 2019] which is specifically designed and optimized  
2380 toward for small problem sizes. We also compare to library *EKR* [Hentschel et al. 2008] which  
2381 computes data mining example PRL (Figure 16) on CPUs – the library is implemented in  
2382 the Java programming language and parallelized via *Java Threads*, and the library is used in  
2383 practice by the *Epidemiological Cancer Registry* in North Rhine-Westphalia (Germany) which  
2384 is the currently largest cancer registry in Europe.

2385 We compare to the approaches experimentally in terms of:

- 2386 i) *performance*: via a runtime comparison of our generated code against code that is generated  
2387 according to the related approaches;
- 2388 ii) *portability*: based on the *Pennycook Metric* [Pennycook et al. 2019] which mathematically  
2389 defines portability<sup>21</sup> as:

$$2390 \Phi(a, p, H) = \begin{cases} \frac{|H|}{\sum_{i \in H} \frac{1}{e_i(a, p)}} & \text{if } i \text{ is supported, } \forall i \in H \\ 0 & \text{otherwise} \end{cases}$$

2395  
2396 <sup>20</sup> We cannot compare to polyhedral compiler TC [Vasilache et al. 2019] which is optimized toward deep learning com-  
2397 putations on GPUs, because TC is not under active development anymore and thus is not working for newer CUDA  
2398 architectures [Facebook Research 2022]. Rasch et al. [2019a] show that our approach – already in its proof-of-concept  
2399 version – achieves higher performance than TC for popular computations on real-world data sets.

2400 <sup>21</sup> Pennycook’s metric is actually called *Performance Portability (PP)*. Since performance portability particularly includes  
2401 functional portability, we refer to Pennycook’s PP also more generally as *Portability* only.

2402 In words: "for a given set of platforms  $H$ , the *performance portability*  $\Phi$  of an application  $a$   
 2403 solving problem  $p$  is defined as  $\Phi(a, p, H)$ , where  $e_i(a, p)$  is the performance efficiency (i.e.  
 2404 a ratio of observed performance relative to some proven, achievable level of performance)  
 2405 of application  $a$  solving problem  $p$  on platform  $i$ ; value  $\Phi(a, p, H)$  is 0, if any platform in  $H$   
 2406 is unsupported by  $a$  running  $p$ ." [Pennycook et al. 2019]. Consequently, Pennycook defines  
 2407 portability as a real value in the interval  $[0, 1]_{\mathbb{R}}$  such that a value close to 1 indicates *high*  
 2408 portability and a value close to 0 indicates *low* portability. Here, platforms  $H$  represents a  
 2409 set of devices (CPUs, GPUs, ...), an application  $a$  is in our context a framework (such as  
 2410 TVM, a polyhedral compiler, or our approach), problems  $p$  are our case studies, and  $e_i(a, p)$   
 2411 is computed as the runtime  $a_{p,i}^{\text{best}}$  of the application that achieves the best observed runtime  
 2412 for problem  $p$  on platform  $i$ , divided by the runtime of application  $a$  for problem  $p$  running  
 2413 on platform  $i$ .

- 2414 iii) *productivity*: by intuitively arguing that our approach achieves the same/lower/higher productivity  
 2415 as the related approaches, using the representative example computation *Matrix-Vector  
 2416 Multiplication* (MatVec) (Figure 6) – classical code metrics, like *Lines of Code (LOC)*, CO-  
 2417 COMO [Boehm et al. 1995], McCabe's cyclomatic complexity [McCabe 1976], and Halstead  
 2418 development effort [Halstead 1977] are not meaningful for comparing the short and concise  
 2419 programs in high-level languages as proposed by the related work as well as our approach.  
 2420

2421 In the following, after discussing our application case studies, experimental setup, auto-tuning  
 2422 system, and code generator, we compare our approach to each of the four above mentioned classes  
 2423 of approaches in Sections 5.1-5.4.

## 2424 Application Case Studies

2425 We use for experiments in this section popular example computations from Figure 16 that belong  
 2426 to different classes of computations:

- 2427 • Linear Algebra Subroutines (BLAS): *Matrix Multiplication* (MatMul) and *Matrix-Vector Multi-  
 2428 pllication* (MatVec);
- 2429 • Stencil Computations: *Jacobi Computation* (Jacobi3D) and *Gaussian Convolution* (Conv2D)  
 2430 which differ from linear algebra routines by accessing neighboring elements in their input  
 2431 data;
- 2432 • Quantum Chemistry: *Coupled Cluster* (CCSD(T)) computations which differ from linear  
 2433 algebra routines and stencil computations by accessing their high-dimensional input data in  
 2434 complex, transposed fashions;
- 2435 • Data Mining: *Probabilistic Record Linkage* (PRL) which differs from the previous computations  
 2436 by relying on a PRL-specific combine operator and scalar function, instead of straightforward  
 2437 additions or multiplications as previous computations;
- 2438 • Deep Learning: the most time-intensive computations within the popular neural networks  
 2439 ResNet-50 [He et al. 2015], VGG-16 [Simonyan and Zisserman 2014], and MobileNet [Howard  
 2440 et al. 2017], according to their TensorFlow implementations [TensorFlow 2022a,b,c]. Deep  
 2441 learning computations rely on advanced variants of linear algebra routines and stencil  
 2442 computations, e.g., MCC and MCC\_Capsule for computing convolution-like stencils, instead  
 2443 of the classical Conv2D variant of convolution (Figure 16) – the deep learning variants are  
 2444 considered as significantly more challenging to optimize than their classical variants [Barham  
 2445 and Isard 2019].

2446 We use for experiments this subset of computations from Figure 16 to make experimenting chal-  
 2447 lenging for us: state-of-practice approaches, such as TVM and vendor libraries from NVIDIA and  
 2448 Intel, are specifically designed and optimized toward these computations. To make experimenting  
 2449

2451 further challenging for us, we consider data sizes and characteristics either taken from real-world  
2452 computations (e.g., from the *TCCG* benchmark suite [Springer and Bientinesi 2016] for quantum  
2453 chemistry computations) or sizes that are preferable for our competitors, e.g., powers of two for  
2454 which many competitors are highly optimized, e.g., vendor libraries. For the deep learning case  
2455 studies, we use data characteristics (sizes, strides, padding strategy, image/filter formats, etc.)  
2456 taken from the particular implementations of the neural networks when computing the popular  
2457 *ImageNet* [Krizhevsky et al. 2012] data set (the particular characteristics are listed in our Appendix,  
2458 Section D.1, for the interested reader). For all experiments, we use single precision floating point  
2459 numbers (a.k.a. float or fp32), as such precision is the default in TensorFlow and many other  
2460 frameworks.

2461

## 2462 **Experimental Setup**

2463 We run our experiments on a cluster containing two different kinds of GPUs and CPUs:

2464

- 2465 • NVIDIA Ampere GPU A100-PCIE-40GB
- 2466 • NVIDIA Volta GPU V100-SXM2-16GB
- 2467 • Intel Xeon Broadwell CPU E5-2683 v4 @ 2.10GHz
- 2468 • Intel Xeon Skylake CPU Gold-6140 @ 2.30GHz

2469 We represent the two CUDA GPUs in our formalism using model  $\text{ASM}_{\text{CUDA+WRP}}$  (Example 11). We  
2470 rely on model  $\text{ASM}_{\text{CUDA+WRP}}$ , rather than the CUDA’s standard model  $\text{ASM}_{\text{CUDA}}$  (also in Example 11),  
2471 in order to exploit CUDA’s (implicit) warp level for a fair comparison to competitors: warp-level  
2472 optimizations are exploited by our competitors, e.g., for *shuffle operations* [NVIDIA 2018] which  
2473 combine the results of threads within a warp with high performance. To fairly compare to TVM and  
2474 PPCG, we avoid exploiting warps’ *tensor core intrinsics* [NVIDIA 2017] in all experiments, which  
2475 compute the multiplication of small matrices with high performance [Feng et al. 2022], because  
2476 these intrinsics are not used in the TVM- and PPCG-generated CUDA code. For the two CPUs,  
2477 we rely on model  $\text{ASM}_{\text{OpenCL}}$  (Example 11) for generating OpenCL code. The same as our approach,  
2478 TVM also generates OpenCL code for CPUs; Pluto relies on the OpenMP approach for CPUs.

2479

2480 For all experiments, we use the currently newest versions of frameworks, libraries, and compilers,  
2481 as follows. We compile our generated GPU code using library CUDA NVRTC [NVIDIA 2022h] from  
2482 CUDA Toolkit 11.4, and we use Intel’s OpenCL runtime version 18.1.0.0920 for compiling  
2483 CPU code. For both compilers, we do not set any flags so that they run in their default modes. For  
2484 the related approaches, we use the following versions of frameworks, libraries, and compilers:

2485

- 2486 • 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  
2487 CPU;
- 2488 • 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  
2489 GPUs, and we use flag --sizes followed by auto-tuned tile sizes – we rely on the *Auto-Tuning  
2490 Framework (ATF)* [Rasch et al. 2021] for choosing optimized tile size values (as we discuss in  
2491 the next subsection);
- 2492 • 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-  
2494 tuned tile sizes for Pluto; the Pluto-generated OpenMP code is compiled via Intel’s icx  
2495 compiler version 2022.0.0 using the Pluto-recommended optimization flags -O3 -fopenmp;
- 2496 • NVIDIA cuBLAS [NVIDIA 2022b] from CUDA Toolkit 11.4, using the NVIDIA-recommended  
2497 compiler flags -fast -O3 -DNDEBUG;

2498

2499

- 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${MKLROOT}/lib/intel64 -liomp5 -lpthread -lm -ldl`, as recommended for oneMKL by Intel's *Link Line Advisor* tool [Intel 2022a], as well as standard flags `-O3 -DNDEBUG`;
- Intel oneDNN [Intel 2022b] also compiled with Intel's icpx compiler version 2022.0.0, using flags `-I${DNNLROOT}/include -L${DNNLROOT}/lib -ldnnl`, according to oneDNN's documentation, as well as standard flags `-O3 -DNDEBUG`;
- 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]. Libraries NVIDIA cuBLAS and NVIDIA cuDNN are also based on the CUDA programming model; thus, we profile them also via CUDA events. To measure the runtimes of the EKR Java library, we use Java function `System.currentTimeMillis()`.

All measurements of CUDA and OpenCL programs contain the pure program runtime only (a.k.a. *kernel runtime*). The runtime of *host code*<sup>22</sup> is not included in the reported runtimes, as performance of host code is not relevant for this work and the same for all approaches.

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

## 2524 Auto-Tuning

The auto-tuning process of our approach relies completely on the generic *Auto Tuning Framework (ATF)* [Rasch et al. 2021]. The ATF framework has proven to be efficient for exploring large search spaces that are based on constrained tuning parameters (as our space introduced in Section 3.4). We use ATF, out of the box, exactly as described by Rasch et al. [2021]: 1) we straightforwardly represent in ATF our search space (Table 1) via *tuning parameters* which express the parameters in the table and their constraints; 2) we use ATF's pre-implemented cost functions for CUDA and OpenCL to measure the cost of our generated OpenCL and CUDA codes (in this work, we consider as cost program's runtime, rather than its energy consumption, etc); 3) we start the tuning process using ATF's default search technique (*AUC bandit* [Ansel et al. 2014]). ATF then fully automatically determines a well-performing tuning parameter configuration for the particular combination of a case study, architecture, and input/output characteristics (size, memory layout, etc).

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

To achieve the best possible performance results for TVM, PPCG, and Pluto, we auto-tune each of these frameworks individually for each particular combination of case study, architecture, and

<sup>22</sup> *Host code* is required in approaches CUDA and OpenCL for program execution – it compiles the CUDA and OpenCL programs, performs data transfers between host and device, etc.

2549 input/output characteristics, the same as for our approach. For example, we start for TVM one tuning  
2550 run when auto-tuning case study MatMul for architecture GPU on one input size, and another, new  
2551 tuning run for a new input size, etc.

2552 Hand-optimized libraries NVIDIA cuBLAS/cuDNN and Intel oneMKL/oneDNN rely on heuristics  
2553 provided by experts, rather than auto-tuning. By relying on heuristics, the libraries avoid the  
2554 time-intensive process of auto-tuning. However, auto-tuning is well amortized in many application  
2555 areas (e.g., deep learning), because the auto-tuned implementations are re-used in many program  
2556 runs. Moreover, auto-tuning avoids the complex and costly process of hand optimization by experts,  
2557 and it often achieves higher performance than hand optimizations, as confirmed later by our  
2558 experiments.

2559 For a fair comparison, we use for each tuning run uniformly the same tuning time of 12h. Even  
2560 though for many computations well-performing tuning results could often also be found in less  
2561 than 12h for our approach as well as for other frameworks, we use such generous tuning time  
2562 for all frameworks to avoid auto-tuning issues in our reported results – analyzing, improving,  
2563 and accelerating the auto-tuning process is beyond the scope of this work and intended for our  
2564 future work (as also outlined in Section 8). In particular, TVM’s Ansor optimizer was often able  
2565 to find well performing optimization decisions in 6h of tuning time or less. This is because Ansor  
2566 explores a small search space that is specifically designed and optimized toward deep learning  
2567 computations – Ansor’s space is a proper subset of our space, as our space aims at capturing general  
2568 optimizations that apply to a broad class of data-parallel computations. Thereby, Ansor usually  
2569 struggles with achieving high performance for computations not taken from the deep learning  
2570 area, as we confirm in our experiments later.

2571 To improve the auto-tuning efficiency for our implementations, we rely on a straightforward  
2572 cost model that shrinks our search space in Table 1 before starting our ATF-based auto-tuning pro-  
2573 cess: i) we always use the same values for Parameters D1, S1, R1 as well as for Parameters D2, S2, R2,  
2574 thereby generating the same loop structure for all three phases (de-composition, scalar, and re-  
2575 composition) such that the structures can be generated as a fused loop nest; ii) we restrict Pa-  
2576 rameters D2, S2, R2 to two values – one value that let threads process outer parts (a.k.a. *blocked*  
2577 *access* or *outer parallelism*, respectively) and one to let threads process inner parts (*strided access*  
2578 or *inner parallelism*); all other permutations are currently ignored for simplicity or because they  
2579 have no effect on the generated code (e.g., permutations of Parameters D2, S2, R2 that only differ in  
2580 dimension tags belonging to memory layers, as discussed in the previous sections); iii) we restrict  
2581 Parameters D3, S3, S5, R3 such that each parameter is invariant under different values of  $d$  of its  
2582 input pairs  $(l, d) \in \text{MDH-LVL}$ , i.e., we always copy full tiles in memory regions (and not a full tile  
2583 of one input buffer and a half tile of another input buffer, which sometimes might achieve higher  
2584 performance when memory is a limited resource). Our cost model is straightforward and might  
2585 filter out configurations of our search space that achieve potentially higher performance than we  
2586 report for our approach in Sections 5.1-5.4. We aim to substantially improve our naive cost model  
2587 in future work, based on *operational semantics* for our low-level representation, in order to improve  
2588 the auto-tuning quality and to reduce (or even avoid) tuning time.

## 2589 **Code Generator**

2591 We provide a straightforward, proof-of-concept code generator, implemented in C++. Our generator  
2592 takes as input the high-level representation of the target computation (Figure 16), in the form of a  
2593 straightforward text file (see Appendix, Section A.4), and it fully automatically generates auto-tuned  
2594 program code, based on the concepts and methodologies introduced and discussed in this paper  
2595 and the ATF auto-tuning framework. In our future work, we aim to integrate our code generation  
2596 approach into the *MLIR* compiler framework [Lattner et al. 2021], building on work-in-progress  
2597

2598 results [[Google SIG MLIR Open Design Meeting 2020](#)], thereby making our work better accessible  
 2599 for the community.

2600

### 2601 5.1 Scheduling Approaches

2602 *Performance.* Figures 24-29 report the performance of the TVM-generated code, which is in CUDA  
 2603 for GPUs and in OpenCL for CPUs. We observe that we usually achieve the high performance  
 2604 of TVM and often perform even better. For example, in Figure 28, we achieve a speedup  $> 2\times$   
 2605 over TVM on NVIDIA Ampere GPU for matrix multiplications as used in the inference phase of the  
 2606 ResNet-50 neural network – an actually favorable example for TVM which is designed and optimized  
 2607 toward deep learning computations executed on modern GPUs. Our performance advantage over  
 2608 TVM is because we parallelize and optimize more efficiently reduction-like computations – in the  
 2609 case of MatMul (Figure 16), its 3rd-dimension (a.k.a.  $k$ -dimension). The difficulties of TVM with  
 2610 reduction computations becomes particularly obvious when computing dot products (Dot) on GPUs  
 2611 (Figure 24): the Dot’s main computation part is a reduction computation (via point-wise addition,  
 2612 see Figure 16), thus requiring reduction-focussed optimization, in particular when targeting the  
 2613 highly-parallel architecture of GPUs: in the case of Dot (Figure 24), our MDH-generated CUDA  
 2614 code exploits parallelization over CUDA blocks, whereas the Ansor-generated TVM code exploits  
 2615 parallelization over threads within in a single block only, because TVM currently cannot use blocks for  
 2616 parallelizing reduction computations [[Apache TVM Community 2022a](#)]. Furthermore, while TVM’s  
 2617 Ansor rigidly parallelizes outer dimensions [[Zheng et al. 2020a](#)], our ATF-based tuning process has  
 2618 auto-tuned our tuning parameters D2, S2, R2 in Table 1 to exploit parallelism for inner dimensions,  
 2619 which achieves higher performance for this particular MatMul example used in ResNet-50. Also,  
 2620 for MatMul-like computations, Ansor always caches parts of the input in GPU’s shared memory,  
 2621 and it computes these cached parts always in register memory. In contrast, our caching strategy  
 2622 is auto-tunable (via parameters D3, S3 S5, R3 in Table 1), and ATF has determined to not cache  
 2623 the input matrices into fast memory resources for the MatMul example in ResNet-50. Surprisingly,  
 2624 Ansor does not exploit fast memory resources for Jacobi stencils (Figure 25), as required to achieve  
 2625 high performance for them: our MDH-generated and ATF-tuned CUDA kernel for Jacobi uses  
 2626 register memory for both inputs (image buffer and filter) when targeting NVIDIA Ampere GPU (small  
 2627 input size), thereby achieving a speedup over TVM+Ansor of  $1.93\times$  for Jacobi. Most likely, Ansor  
 2628 fails to foresee the potential of exploiting fast memory resources for Jacobi stencils, because the  
 2629 Jacobi’s index functions used for memory accesses (Figure 16) are injective. For the MatMul example  
 2630 of ResNet-50’s training phase (Figure 28), we achieve a speedup over TVM on NVIDIA Ampere GPU  
 2631 of  $1.26\times$ , because auto-tuning determined to store parts of input matrix  $A$  as transposed into fast  
 2632 memory (via parameter D4 in Table 1). Storing parts of the input/output data as transposed is not  
 2633 considered by Ansor as optimization, perhaps because such optimization must be expressed in  
 2634 TVM’s high-level language, rather than scheduling language [[Apache TVM Community 2022c](#)]. For  
 2635 MatVec on NVIDIA Ampere GPU (Figure 24), we achieve a speedup over TVM of  $1.22\times$  for the small  
 2636 input size, by exploiting a so-called *swizzle pattern* [[Phothilimthana et al. 2019](#)]: our ATF tuner  
 2637 has determined to assign threads that are consecutive in CUDA’s  $x$ -dimension to the second MDA  
 2638 dimension (via parameters D2, S2, R2 in Table 1), thereby accessing the input matrix in a GPU-  
 2639 efficient manner (a.k.a *coalesced global memory accesses* [[NVIDIA 2022f](#)]). In contrast, for MatVec  
 2640 computations, Ansor assigns threads with consecutive  $x$ -ids always to the first data dimension, in a  
 2641 non-tunable manner, causing lower performance.

2642 Our positive speedups over TVM on CPU are for the same reasons as discussed above for GPU. For  
 2643 example, we achieve a speedup of  $> 3\times$  over TVM on Intel Skylake CPU for MCC (Figure 29) as used  
 2644 in the training phase of the MobileNet neural network, because we exploit fast memory resources  
 2645 more efficiently than TVM: our auto-tuning process has determined to use register memory for the

2647 MCC's second input (the filter buffer  $F$ , see Table 16) and using no fast memory for the first input  
2648 (image buffer  $I$ ), whereas Ansor uses shared memory rigidly for both inputs of MCC. Moreover,  
2649 our auto-tuning process has determined to parallelize the inner dimensions of MCC, while Ansor  
2650 always parallelizes outer dimensions. We achieve the best speedup over TVM for MCC on an input size  
2651 taken from TVM's own tutorials [Apache TVM Documentation 2022b] (Figure 25), rather than from  
2652 neural networks (as in Figures 28 and 29). This is because TVM's MCC size includes large reduction  
2653 computations, which are not efficiently optimized by TVM (as discussed above).

2654 The TVM compiler achieves higher performance than our approach for some examples in Figures  
2655 24-29. However, in most cases, this has a technical reason only: TVM uses the NVCC compiler for  
2656 compiling CUDA code, while our proof-of-concept code generator relies on NVIDIA's NVRTC library  
2657 which surprisingly generates less efficient CUDA assembly than NVCC. In three cases, the higher  
2658 performance of TVM over our approach is because our ATF auto-tuning framework was not able  
2659 to find a better performing tuning configuration than TVM's Ansor optimization engine; the three  
2660 cases are: 1) MCC (capsule variant) from ResNet-50's training phase on Intel Skylake CPU, 2) MCC  
2661 from VGG-16's inference phase on NVIDIA Ampere GPU, and 3) MCC (capsule variant) from VGG-16's  
2662 training phase on NVIDIA Ampere GPU. However, when we manually set the Ansor-found tuning  
2663 configuration also for our approach (analogously as done in Section 3.5), instead of using the  
2664 ATF-found configuration, we achieve for these three cases exactly the same high performance as  
2665 TVM+Ansor, i.e., the well-performing configurations are contained in our MDH-based search space.  
2666 Most likely, Ansor was able to find this well-performing configuration automatically, because it  
2667 explores a smaller search space that is particularly designed for deep learning computations. To  
2668 avoid such tuning issues in our approach, we aim to substantially improve our auto-tuning process  
2669 in future work: we plan to introduce an analytical cost model that assists (or even replaces) our  
2670 auto-tuner, as we also outline in Section 8.

2671 Note that the TVM compiler crashes for our data mining example PRL, because TVM has difficulties  
2672 with computations relying on user-defined combine operators [Apache TVM Community 2022d].

2673 *Portability.* Figure 30 reports the portability of the TVM compiler over our GPU and CPU architectures.  
2674 The portability measurements are based on the Pennycook metric where a value close to 1 indicates  
2675 high portability and a value close to 0 low portability, correspondingly. We observe that except  
2676 for the example of transposed matrix multiplication  $\text{GEMM}^\top$ , we always achieve higher portability  
2677 than TVM. The higher portability of TVM for  $\text{GEMM}^\top$  is because TVM achieves for this example higher  
2678 performance than our approach on NVIDIA Volta GPU. However, the higher performance of TVM is  
2679 only due to the fact that TVM uses NVIDIA's NVCC for compiling CUDA code, while we currently  
2680 rely on NVIDIA's NVRTC library which surprisingly generates less efficient CUDA assembly, as  
2681 discussed above.

2682 *Productivity.* Listing 1 shows how matrix-vector multiplication (MatVec) is implemented in TVM's  
2683 high-level program representation which is embedded into the Python programming language.  
2684 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  
2685 declared, in the form of function parameters; the matrix and vector are named  $M$  and  $v$  and both are  
2686 assumed to contain elements of scalar type  $T = \text{float32}$  (floating point numbers). Line 5 defines  
2687 a so-called *reduction axis* in TVM in which all values are combined in line 8 via `te.sum` (addition).  
2688 The basic computation part of MatVec – multiplying matrix element  $M[i, k]$  with vector element  
2689  $v[k]$  – is also specified in line 8.

2690 While we see the MatVec implementations of TVM (Listing 1) and our approach (Figure 6) basically  
2691 on the same level of abstraction, we consider our approach in general as more expressive.  
2692 This is because our approach supports multiple reduction dimensions that may rely on different  
2693 combine operators, e.g., as required for expressing the MBBS example in Figure 16 with which TVM  
2695

```

2696 1 def MatVec(I, K):
2697 2     M = te.placeholder((I, K), name='M', dtype='float32')
2698 3     v = te.placeholder((K,), name='v', dtype='float32')
2699 4
2700 5     k = te.reduce_axis((0, K), name='k')
2701 6     w = te.compute(
2702 7         (I,),
2703 8         lambda i: te.sum(M[i, k] * v[k], axis=k)
2704 9     )
2705 10
2706 11     return [M, v, w]

```

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

2706  
2707

2708  
2709 is struggling – adding support for multiple, different reduction dimensions is considered in the  
2710 TVM community as a non-trivial extension of TVM [Apache TVM Community 2020, 2022b]. Also,  
2711 we consider our approach as slightly less error-prone: we automatically compute the expected  
2712 sizes of matrix  $M$  (as  $I \times K$ ) and vector  $v$  (as  $K$ ), based on the user-defined input size  $(I, K)$  in  
2713 line 1 and index functions  $(i, k) \mapsto (i, k)$  for the matrix and  $(i, k) \mapsto (k)$  for the vector in line 8  
2714 (see Definition 8). In contrast, TVM redundantly requests these matrix and vector sizes from the  
2715 user for computing the function specification of its generated MatVec code (once in lines 2 and 3 of  
2716 Listing 1, and again in lines 5 and 7), which lets TVM generate incorrect low-level code – without  
2717 issuing an error message – when the user sets sizes different from  $I \times K$  for the matrix and  $K$  for  
2718 the vector in lines 2 and 3 [Apache TVM Community 2022f].

## 2719 5.2 Polyhedral Compilers

2720 *Performance.* Figures 24–29 report for our application case studies the performance of the PPCG-  
2721 generated CUDA code for GPUs and of the OpenMP-annotated C code generated by polyhedral  
2722 compiler Pluto for CPUs. For a fair comparison, we report for both polyhedral compilers their  
2723 performance achieved for ATF-tuned tile sizes (denoted as PPCG+ATF/Pluto+ATF in Figures 24–29),  
2724 as well as the performance of the two compilers when relying on their internal heuristics instead  
2725 of auto-tuning (denoted as PPCG and Pluto in the figures). In some cases, PPCG’s heuristic crashed  
2726 with error too many resources requested for launch, because the heuristic seems to not take  
2727 into account device-specific constraints, e.g., limited availability of GPUs’ fast memory resources.

2728 We observe that our MDH-based approach achieves better performance than PPCG and Pluto in  
2729 all cases – sometimes by multiple orders of magnitude – in particular for deep learning computations  
2730 (Figures 28 and 29). This is caused by the rigid optimization goals of PPCG and Pluto, e.g., always  
2731 parallelizing outer dimensions, which causes severe performance losses. For example, we achieve a  
2732 speedup over PPCG of  $> 13\times$  on NVIDIA Ampere GPU and of  $> 60\times$  over Pluto on Intel Skylake CPU  
2733 for MCC as used in the inference phase of the real-world ResNet-50 neural network. Compared to  
2734 PPCG, our better performance for this MCC example is because PPCG has difficulties with efficiently  
2735 parallelizing computations relying on more than 3 dimension; most likely, this is because CUDA  
2736 offers per default 3 dimensions for parallelization (called x, y, z dimension in CUDA). However,  
2737 MCC relies on 7 parallelizable dimensions (as shown in Figure 16); exploiting the parallelization  
2738 opportunities of the 4 further dimensions (as done in our generated CUDA code) is essential to  
2739 achieve high performance for this MCC example from ResNet-50. Our performance advantage over  
2740 Pluto for the MCC example is because Pluto parallelizes the outer dimensions of MCC only; however,  
2741 the dimension has a size of only 1 for this real-world example, resulting in starting only 1 thread in  
2742 the Pluto-generated OpenMP code.

2744

2745 For dot products Dot (Figure 24), we can observe that PPCG fails to generate parallel CUDA code,  
 2746 because PPCG cannot parallelize and optimize computations which rely solely on combine operators  
 2747 different from concatenation, as we also discuss in Section 6.2. In Section 6.2, we also discuss that  
 2748 we do not consider the performance issues of PPCG and Pluto as weaknesses of the polyhedral  
 2749 approach in general, but of the particular polyhedral transformations used in PPCG and Pluto.

2750 Note that Pluto crashes for our data mining examples (Figure 27), with Error extracting  
 2751 polyhedra from source file, which is due to the complex scalar function of the example,  
 2752 which involves if-statements. Moreover, Intel’s icx compiler struggles with compiling the Pluto-  
 2753 generated OpenMP code for quantum chemistry computations (Figure 26): we aborted icx’s  
 2754 compilation process after 24h compilation time. The icx’s issues with the Pluto-generated code  
 2755 are most likely because of too aggressive loop unrolling of Pluto – the Pluto-generated OpenMP  
 2756 code has often a size > 50MB for our real-world quantum chemistry examples.

2757 *Portability.* Since PPCG and Pluto are each designed for particular architectures only, they achieve  
 2758 the lowest portability of 0 for all our studies, according to the Pennycook metric. To simplify for  
 2759 PPCG and Pluto the portability comparison with our approach, we compute the Pennycook metric  
 2760 additionally also for two restricted sets of devices: only GPUs to make comparison against our  
 2761 approach easier for PPCG, and only CPUs for Pluto.

2762 Figures 31-35 report the portability of PPCG when considering only GPUs, as well as the portability  
 2763 of Pluto for only CPUs. We observe that we achieve higher portability for all our studies, as we  
 2764 constantly achieve higher performance than the two polyhedral compilers for the studies.

2765 Note that even when restricting our set of devices to only GPUs for PPCG or only CPUs for Pluto,  
 2766 the two polyhedral compilers still achieve a portability of 0 for some examples, because they fail to  
 2767 generate code for them (as discussed above).

2768 *Productivity.* Listing 2 shows the input program of polyhedral compilers PPCG and Pluto for  
 2769 MatVec. Both take as input easy-to-implement, straightforward, sequential C code. We consider  
 2770 these two polyhedral compilers as more productive than our approach (as well as scheduling/func-  
 2771 tional approaches and polyhedral compilers relying on a DSL, rather than sequential programs, as  
 2772 TC [Vasilache et al. 2019]), because both compilers fully automatically generate optimized parallel  
 2773 code from unoptimized, sequential programs.

2774 Rasch et al. [2020a,b] show that our approach can achieve the same, high user productivity as  
 2775 polyhedral compilers, by using a polyhedral frontend for our approach: we can alternatively take  
 2776 as input the same sequential user programs as PPCG and Pluto, instead of programs implemented  
 2777 in our high-level program representation (as in Figure 6). The sequential input program is then  
 2778 transformed via polyhedral tool *pet* [Verdoolaege and Grosser 2012] to its polyhedral representation  
 2779 which is then automatically transformed to our high-level program representation, according to  
 2780 the methodology presented by Rasch et al. [2020a,b].

```
2783 1 for( int i = 0 ; i < M ; ++i )
2784 2   for( int k = 0 ; k < K ; ++k )
2785 3     w[i] += M[i][k] * v[k];
```

2786 Listing 2. PPCG/Pluto program expressing Matrix-Vector Multiplication (MatVec)

### 2788 5.3 Functional Approaches

2789 Our previous work [Rasch et al. 2019a] already shows that while functional approaches provide  
 2790 a solid formal foundation for computations, they typically have performance and portability

2794 issues; for this, our previous work used the state-of-the-art Lift [Steuwer et al. 2015] framework  
 2795 as running example (which, to the best of our knowledge, has so far not been improved toward  
 2796 higher performance and/or better portability). Therefore, we refrain from a further performance and  
 2797 portability evaluation of Lift and focus in the following on analyzing and discussing the productivity  
 2798 potentials of functional approaches, using again the state-of-the-art Lift approach as running  
 2799 example. We discuss the performance and portability issues of functional approaches from a general  
 2800 perspective thoroughly in Section 6.3.

2801  
 2802 *Performance/Portability.* Already experimentally evaluated in previous work [Rasch et al. 2019a]  
 2803 and discussed in general terms in Section 6.3.

2804  
 2805 *Productivity.* Listing 3 shows the implementation of MatVec in the Lift approach. In line 1, type  
 2806 parameters  $n$  and  $m$  are declared via Lift building block `nFun` – type parameters are limited to  
 2807 natural numbers in the Lift formalism. Line 2 declares a function in Lift that takes as input a  
 2808 matrix of size  $m \times n$  and a vector of size  $n$ , correspondingly; the matrix and vector are both assumed  
 2809 to consist of `float` numbers (floating point numbers). The computation of MatVec is specified in  
 2810 lines 3 and 4: the `map` function of Lift in line 3 iterates over all rows of the matrix, each row is  
 2811 pair-wise combined with the input vector via Lift pattern `zip`, multiplication `*` is applied to each  
 2812 pair in the combined vector via Lift's `map` pattern in line 4, and the obtained products are finally  
 2813 combined via addition `+` using Lift's `reduce` pattern.

```
2814
2815 1  nFun(n => nFun(m =>
2816 2    fun(matrix: [[float]]n)m => fun(xs: [float]n =>
2817 3      matrix :>> map(fun(row =>
2818 4        zip(xs, row) :>> map(*) :>> reduce(+, 0)
2819 5        )))) ) )
```

2820 Listing 3. Lift program expressing Matrix-Vector Multiplication (MatVec)

2821  
 2822  
 2823 Already for expressing MatVec, we can observe that Lift relies on a vast set of small, functional  
 2824 building blocks (five building blocks for MatVec: `nFun`, `fun`, `map`, `zip`, and `reduce`), and the blocks  
 2825 have to be composed and nested in complex ways for expressing computations. Consequently,  
 2826 we consider programming in Lift-like approaches as complex and their productivity for the user  
 2827 as limited. Moreover, Lift-like approaches often need fundamental extension for targeting new  
 2828 kinds of computations, e.g., so-called *macro-rules* which had to be added to Lift for efficiently  
 2829 targeting matrix multiplications [Remmehg et al. 2016] and primitives `slide` and `pad` together with  
 2830 optimization *overlapped tiling* for expressing stencil computations [Hagedorn et al. 2018]. This  
 2831 need for extensions limits the expressiveness of the Lift language and thus hinders productivity.

2832 In contrast to Lift, our approach relies on exactly three higher-order functions (Figure 5)  
 2833 to express various kinds of data-parallel computations (Figure 16): 1) `inp_view` (Definition 8)  
 2834 which prepares the input data; our `inp_view` function is designed as general enough to sub-  
 2835 sume – in a structured manner – the subset of all Lift patterns intended to change the view on  
 2836 input data, including patterns `zip`, `pad`, and `slide`; 2) `md_hom` (Definition 3) expresses the actual  
 2837 computation and subsumes all Lift patterns performing actual computations (`fun`, `map`, `reduce`,  
 2838 `...`); 3) `out_view` (Definition 10) expresses the view on output data and is designed to work similarly  
 2839 as function `inp_view` (Lemma 2). Our three functions are always composed straightforwardly in  
 2840 the same, fixed order (Figure 5), and they do not rely on complex function nesting for expressing  
 2841 computations.

2843 Note that even though our language is designed as minimalistic, it should cover the expressivity  
2844 of the Lift language<sup>23</sup> and beyond: for example, we are currently not aware of any Lift program  
2845 being able to express the prefix-sum examples in Figure 16. For the above reasons, we consider  
2846 programming in our high-level language as more productive for the user than programming in  
2847 Lift-like functional languages. Furthermore, as discussed in Section 5.2, our approach can take as  
2848 input also straightforward, sequential program code, which further contributes to the productivity  
2849 of our approach.

2850

2851

2852

2853

2854

2855

2856

2857

2858

2859

2860

2861

2862

2863

2864

2865

2866

2867

2868

2869

2870

2871

2872

2873

2874

2875

2876

2877

2878

2879

2880

2881

2882

2883

2884

2885

2886

2887

2888 <sup>23</sup> This work is focussed on dense computations. Lift supports sparse computations [Pizzuti et al. 2020] which we consider  
2889 as future work (as also outlined in Section 8). We consider Lift’s approach, based on their so-called *position dependent*  
2890 *arrays*, as a great inspiration for our future goal.

2891

2892	2893	2894	2895	2896	2897	2898	2899	2900	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		1.00		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		-				

2901	2902	2903	2904	2905	2906	2907	2908	2909	2910	2911	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		-						

2912	2913	2914	2915	2916	2917	2918	2919	2920	2921	2922	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		-						

2923	2924	2925	2926	2927	2928	2929	2930	2931	2932	2933	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		-						

2935 Fig. 24. Speedup (higher is better) of our approach for linear algebra routines on GPUs and CPUs over:  
2936 i) scheduling approach TVM, ii) polyhedral compilers PPCG (GPU) and Pluto (CPU), as well as iii) hand-  
2937 optimized libraries provided by vendors. Dash symbol "-" means this framework does not support this  
2938 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. 25. 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									
	Quantum Chemistry	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
2990	TVM+Ansor	1.15	1.07	1.25	1.00	1.36	1.05	1.00	1.15
2991	PPCG	10585.85	10579.40	9819.81	11211.57	10181.14	10482.81	11693.21	10585.85
2992	PPCG+ATF	11.19	15.60	14.06	11.45	11.81	12.06	11.72	11.19

NVIDIA Volta GPU									
	Quantum Chemistry	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
3000	TVM+Ansor	1.09	0.93	1.04	1.03	1.01	1.11	1.01	1.09
3001	PPCG	6466.22	6019.64	6300.31	6468.40	6608.80	5256.49	6602.22	6466.22
3002	PPCG+ATF	8.28	9.61	9.38	7.21	6.60	5.14	7.77	8.28

Intel Skylake CPU									
	Quantum Chemistry	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
3010	TVM+Ansor	1.60	1.50	2.06	1.70	1.20	2.12	1.56	1.60
3011	Pluto	147.45	151.55	206.60	162.58	157.43	145.17	321.66	147.45
3012	Pluto+ATF	1.89	2.01	1.89	1.80	1.82	1.92	1.84	1.89

Intel Broadwell CPU									
	Quantum Chemistry	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
3019	TVM+Ansor	1.06	1.28	1.16	1.15	1.29	1.13	2.07	1.06
3020	Pluto	-	-	-	-	-	-	-	-
3021	Pluto+ATF	-	-	-	-	-	-	-	-

3025 Fig. 26. Speedup (higher is better) of our approach for quantum chemistry computations Coupled Cluster  
 3026 (CCSD(T)) on GPUs and CPUs over: i) scheduling approach TVM, and ii) polyhedral compilers PPCG (GPU)  
 3027 and Pluto (CPU). Dash symbol "-" means this framework does not support this combination of architecture,  
 3028 computation, and data characteristic.

	Quantum Chemistry	abcdef-gdab-efgc	abcdef-gdac-efgb	abcdef-gdbc-efga	abcdef-geab-dfgc	abcdef-geac-dfgb	abcdef-gebc-dfga	abcdef-gfab-degc	abcdef-gfbc-dega
3029	TVM+Ansor	1.15	1.07	1.25	1.00	1.36	1.05	1.00	1.15
3030	PPCG	10585.85	10579.40	9819.81	11211.57	10181.14	10482.81	11693.21	10585.85
3031	PPCG+ATF	11.19	15.60	14.06	11.45	11.81	12.06	11.72	11.19
3032	Pluto	-	-	-	-	-	-	-	-
3033	Pluto+ATF	-	-	-	-	-	-	-	-
3034									
3035									
3036									
3037									
3038									

Data Mining		NVIDIA Ampere GPU					
		$2^{15}$	$2^{16}$	$2^{17}$	$2^{18}$	$2^{19}$	$2^{20}$
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					
		$2^{15}$	$2^{16}$	$2^{17}$	$2^{18}$	$2^{19}$	$2^{20}$
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					
		$2^{15}$	$2^{16}$	$2^{17}$	$2^{18}$	$2^{19}$	$2^{20}$
TVM+Ansor	-	-	-	-	-	-	-
Pluto	-	-	-	-	-	-	-
Pluto+ATF	-	-	-	-	-	-	-
EKR	6.18	5.39	9.62	19.87	26.42	24.78	

Data Mining		Intel Broadwell CPU					
		$2^{15}$	$2^{16}$	$2^{17}$	$2^{18}$	$2^{19}$	$2^{20}$
TVM+Ansor	-	-	-	-	-	-	-
Pluto	-	-	-	-	-	-	-
Pluto+ATF	-	-	-	-	-	-	-
EKR	8.01	9.17	23.58	66.90	119.33	167.19	

Fig. 27. Speedup (higher is better) of our approach for data mining algorithm Probabilistic Record Linkage (PRL) on GPUs and CPUs over: i) scheduling approach TVM, and ii) polyhedral compilers PPCG (GPU) and Pluto (CPU), as well as the iii) hand-implemented Java CPU implementation used by EKR – the largest cancer registry in Europa. Dash symbol “-” means this framework does not support this combination of architecture, computation, and data characteristic.

NVIDIA Ampere GPU											
3088	3089	3090	ResNet-50				VGG-16			MobileNet	
			Training		Inference		Training		Inference	Training	Inference
			MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MCC	
TVM+Ansor	1.00	1.26	1.05	2.22	0.93	1.42	0.88	1.14	0.94	1.00	
PPCG	3456.16	8.26	—	7.89	1661.14	7.06	5.77	5.08	2254.67	7.55	
PPCG+ATF	3.28	2.58	13.76	5.44	4.26	3.92	9.46	3.73	3.31	10.71	
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	—	—	
NVIDIA Volta GPU											
3101	3102	3103	ResNet-50				VGG-16			MobileNet	
			Training		Inference		Training		Inference	Training	Inference
			MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MCC	
TVM+Ansor	0.75	1.21	0.72	1.79	1.00	1.11	1.06	1.00	1.00	1.00	1.00
PPCG	1976.38	5.88	—	5.64	994.16	3.41	8.21	2.51	1411.92	7.26	
PPCG+ATF	3.43	3.54	3.42	4.93	3.85	3.15	8.13	2.05	3.49	3.56	
cuDNN	1.21	—	1.29	—	2.80	—	3.50	—	2.32	3.14	
cuBLAS	—	1.33	—	1.14	—	1.09	—	1.04	—	—	
cuBLASEx	—	1.21	—	1.07	—	1.04	—	1.03	—	—	
cuBLASLt	—	1.00	—	1.07	—	1.04	—	1.02	—	—	
NVIDIA Ampere GPU											
3114	3115	3116	ResNet-50				VGG-16			MobileNet	
			Training	Inference	Training	Inference	Training	Inference	Training	Inference	
			MCC_Capsule								
TVM+Ansor	0.96	1.00	0.79	—	1.02	—	0.88	—	0.99	—	
PPCG	4642.24	—	1013.55	—	—	—	4017.74	—	—	—	
PPCG+ATF	25.98	85.33	4.41	—	13.64	—	8.89	—	22.12	—	
cuDNN	—	—	—	—	—	—	—	—	—	—	
NVIDIA Volta GPU											
3123	3124	3125	ResNet-50				VGG-16			MobileNet	
			Training	Inference	Training	Inference	Training	Inference	Training	Inference	
			MCC_Capsule								
TVM+Ansor	0.95	1.01	1.05	—	0.97	—	1.04	—	0.87	—	
PPCG	2935.40	—	945.16	—	—	—	2885.90	—	—	—	
PPCG+ATF	19.24	19.68	8.28	—	12.29	—	8.84	—	6.41	—	
cuDNN	—	—	—	—	—	—	—	—	—	—	

3132 Fig. 28. Speedup (higher is better) of our approach for the most time-intensive computations used in deep  
 3133 learning neural networks ResNet-50, VGG-16, and MobileNet on GPUs over: i) scheduling approach TVM,  
 3134 ii) polyhedral compilers PPCG (GPU), as well as iii) hand-optimized libraries provided by vendors. Dash  
 3135 symbol “—” means this framework does not support this combination of architecture, computation, and data  
 3136 characteristic.

3137 3138 3139 3140 3141 3142 3143 3144 3145 3146 3147		Intel Skylake CPU									
		ResNet-50				VGG-16				MobileNet	
		Training		Inference		Training		Inference		Training	Inference
MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MCC
TVM+Ansor	1.53	1.05	1.14	1.20	1.97	1.14	2.38	1.27	3.01	1.40	
Pluto	355.81	49.57	364.43	13.93	130.80	93.21	186.25	36.30	152.14	75.37	
Pluto+ATF	13.08	19.70	170.69	6.57	3.11	6.29	53.61	8.29	3.50	25.41	
oneDNN	0.39	—	5.07	—	1.22	—	9.01	—	1.05	4.20	
oneMKL	—	0.44	—	1.09	—	0.88	—	0.53	—	—	
oneMKL(JIT)	—	6.43	—	8.33	—	27.09	—	9.78	—	—	

3148 3149 3150 3151 3152 3153 3154 3155 3156 3157 3158 3159 3160		Intel Broadwell CPU									
		ResNet-50				VGG-16				MobileNet	
		Training		Inference		Training		Inference		Training	Inference
MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MatMul	MCC	MCC
TVM+Ansor	1.53	1.60	1.29	1.53	1.32	1.00	1.27	1.02	2.42	1.92	
Pluto	4349.20	40.41	137.21	15.96	1865.07	53.57	113.40	24.10	2255.00	53.85	
Pluto+ATF	6.43	8.93	61.60	6.91	5.07	4.38	42.63	4.45	6.43	29.18	
oneDNN	1.30	—	1.81	—	2.94	—	2.85	—	1.83	4.47	
oneMKL	—	1.45	—	1.36	—	1.35	—	0.50	—	—	
oneMKL(JIT)	—	19.78	—	9.77	—	50.58	—	10.70	—	—	

3161 3162 3163 3164 3165 3166 3167 3168 3169 3170 3171		Intel Skylake CPU									
		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	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule
TVM+Ansor	0.94	1.14	3.50	1.18	2.94	1.59					
Pluto	209.36	265.77	—	166.45	160.49	159.34					
Pluto+ATF	14.33	265.77	3.33	60.66	4.40	57.21					
oneDNN	—	—	—	—	—	—	—	—	—	—	—

3172 3173 3174 3175 3176 3177 3178 3179 3180		Intel Broadwell CPU									
		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	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule	MCC_Capsule
TVM+Ansor	2.61	1.30	3.55	1.00	1.32	2.24					
Pluto	—	—	—	—	—	—	—	—	—	—	—
Pluto+ATF	4418.82	56.17	75.77	2173.72	202.34	158.52					
oneDNN	—	—	—	—	—	—	—	—	—	—	—

3181 Fig. 29. Speedup (higher is better) of our approach for the most time-intensive computations used in deep  
3182 learning neural networks ResNet-50, VGG-16, and MobileNet on CPUs over: i) scheduling approach TVM,  
3183 ii) polyhedral compilers Pluto (CPU), as well as iii) hand-optimized libraries provided by vendors. Dash  
3184 symbol “—” means this framework does not support this combination of architecture, computation, and data  
3185 characteristic.

Pennycook Metric								
Linear Algebra	Dot		MatVec		MatMul		MatMult	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
	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		
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	abcdefg-gfbc-dega
	MDH+ATF	1.00	0.98	1.00	1.00	1.00	1.00	1.00
	TVM+Ansor	0.82	0.82	0.73	0.82	0.82	0.74	0.71
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>		
	MDH+ATF	1.00	1.00	1.00	1.00	1.00		
	TVM+Ansor	0.00	0.00	0.00	0.00	0.00		
Pennycook Metric								
Deep Learning	ResNet-50				VGG-16			MobileNet
	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	0.54	0.75	0.96	0.97	0.96		
Pennycook Metric								
Deep Learning (Capsule)	ResNet-50				VGG-16			MobileNet
	Training		Inference		Training		Inference	
	MCC_Capsule	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	0.97	0.96	0.96
TVM+Ansor	0.71	0.90	0.44	0.95	0.63	0.63	0.69	0.69

Fig. 30. 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.

3235 3236 3237 3238 3239 3240 3241 3242 3243 3244 3245	3246 3247 3248 3249 3250 3251 3252 3253 3254 3255	Pennycook Metric (GPUs only)							
		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
MDH+ATF	MDH+ATF	1.00	1.00	1.00	1.00	1.00	0.45	0.89	1.00
TVM+Ansor	TVM+Ansor	0.01	0.01	0.71	0.88	1.00	0.42	1.00	0.92
PPCG	PPCG	0.00	0.00	0.25	0.43	0.56	0.15	0.30	0.01
PPCG+ATF	PPCG+ATF	0.00	0.00	0.30	0.40	0.91	0.21	0.71	0.34
cuBLAS	cuBLAS	0.93	0.91	0.89	0.96	0.50	0.42	0.52	0.60
cuBLASEx	cuBLASEx	0.00	0.00	0.00	0.00	0.67	0.98	0.60	0.00
cuBLASL <sub>t</sub>	cuBLASL <sub>t</sub>	0.00	0.00	0.00	0.00	0.83	0.48	0.60	0.00
Pennycook Metric (CPUs only)									
MDH+ATF	MDH+ATF	0.78	0.48	0.48	0.74	0.79	0.67	1.00	0.90
TVM+Ansor	TVM+Ansor	0.15	0.06	0.44	0.32	0.71	0.60	0.94	0.86
Pluto	Pluto	0.15	0.07	0.18	0.24	0.20	0.08	0.16	0.07
Pluto+ATF	Pluto+ATF	0.15	0.07	0.23	0.37	0.31	0.18	0.24	0.55
oneMKL	oneMKL	0.99	1.00	1.00	0.41	0.17	1.00	0.37	1.00
oneMKL(JIT)	oneMKL(JIT)	0.00	0.00	0.00	0.00	0.98	0.00	0.83	0.00

3256 Fig. 31. Portability (higher is better), according to Pennycook metric, for linear algebra routines computed on  
 3257 only GPUs or CPUs, respectively. The restriction simplifies for frameworks with limited architectural support  
 3258 (such as polyhedral compilers and vendor libraries) the portability comparisons against our approach.

3260 3261 3262 3263 3264 3265 3266 3267 3268 3269	3270 3271 3272 3273 3274 3275 3276 3277	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	MDH+ATF	1.00	1.00	1.00	1.00	1.00
TVM+Ansor	TVM+Ansor	0.50	0.51	1.00	0.46	0.64
PPCG	PPCG	0.18	0.10	0.66	0.49	0.00
PPCG+ATF	PPCG+ATF	0.95	0.99	0.82	0.74	0.11
cuDNN	cuDNN	0.00	0.00	0.42	0.23	0.29
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	MDH+ATF	1.00	1.00	1.00	1.00	1.00
TVM+Ansor	TVM+Ansor	0.44	0.58	0.42	0.31	0.30
Pluto	Pluto	0.35	0.46	0.43	0.56	0.01
Pluto+ATF	Pluto+ATF	0.65	0.83	0.52	0.86	0.01
oneDNN	oneDNN	0.10	0.11	0.15	0.10	0.17

3278 Fig. 32. Portability (higher is better), according to Pennycook metric, for stencil computations computed on  
 3279 only GPUs or CPUs, respectively. The restriction simplifies for frameworks with limited architectural support  
 3280 (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
MDH+ATF	1.00	0.96	1.00	1.00	1.00	1.00	1.00	1.00
TVM+Ansor	0.90	0.96	0.87	0.99	0.84	0.93	0.99	1.00
PPCG	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
PPCG+ATF	0.10	0.08	0.09	0.11	0.11	0.12	0.10	0.15

Pennycook Metric (CPUs 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
MDH+ATF	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
TVM+Ansor	0.75	0.72	0.62	0.70	0.80	0.62	0.55	0.72
Pluto	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Pluto+ATF	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Fig. 33. Portability (higher is better), according to Pennycook metric, for quantum chemistry computation Coupled Cluster (CCSD(T)) 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)						
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>
MDH+ATF	1.00	1.00	1.00	1.00	1.00	1.00
TVM+Ansor	0.00	0.00	0.00	0.00	0.00	0.00
PPCG	0.77	0.91	0.90	0.80	0.69	0.59
PPCG+ATF	0.75	0.77	0.67	0.59	0.51	0.43

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>	2 <sup>20</sup>
MDH+ATF	1.00	1.00	1.00	1.00	1.00	1.00
TVM+Ansor	0.00	0.00	0.00	0.00	0.00	0.00
Pluto	0.00	0.00	0.00	0.00	0.00	0.00
Pluto+ATF	0.00	0.00	0.00	0.00	0.00	0.00
EKR	0.14	0.14	0.06	0.02	0.01	0.01

Fig. 34. Portability (higher is better), according to Pennycook metric, for data mining algorithm Probabilistic Record Linkage (PRL) 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.

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	
3345											
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	
3357											
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	
3367											
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. 35. 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.

3382 **5.4 Domain-Specific Approaches**

3383 *Performance.* Figures 24-29 report for completeness also performance results achieved by domain-  
 3384 specific approaches. Since domain-specific approaches are specifically designed and optimized  
 3385 toward particular applications domains and often also architectures (e.g., only linear algebra routines  
 3386 on only GPU), we consider comparing to them as challenging for us: our approach is designed and  
 3387 optimized toward data-parallel computations in general, from arbitrary application domains (the  
 3388 same as also TVM, polyhedral compilers, and many functional approaches), and our approach is  
 3389 also flexible in the target parallel architecture.

3390 We observe in Figures 24-29 that the domain-specific libraries NVIDIA cuBLAS/cuDNN (for lin-  
 3391 ear algebra routines and convolutions on GPUs) and Intel oneMKL/oneDNN (for linear algebra  
 3392 routines and convolutions on CPUs) sometimes perform better and sometimes worse than our  
 3393 approach. The better performance of libraries over our approach is most likely<sup>24</sup> because the  
 3394 libraries internally rely on assembly-level optimizations, while we currently focus on the higher  
 3395 CUDA/OpenCL abstraction level which offers less optimization opportunities [Goto and Geijn  
 3396 2008; Lai and Seznec 2013]. The cuBLASEx extension of cuBLAS achieves in one case – MatMul  
 3397 on NVIDIA Volta GPU for square  $1024 \times 1024$  input matrices – significantly higher performance  
 3398 than our approach. The high performance is achieved by cuBLASEx when using its algorithm  
 3399 variant CUBLAS\_GEMM\_ALG01\_TENSOR\_OP which casts the float-typed inputs implicitly to the  
 3400 half precision type (a.k.a. half or fp16), allowing cuBLASEx exploiting the GPU’s tensor core  
 3401 extension [NVIDIA 2017]. Thereby, cuBLASEx achieves significantly higher performance than our  
 3402 approach, because tensor cores compute small matrix multiplication immediately in hardware; how-  
 3403 ever, at the cost of a significant precision loss: the half scalar type achieves only half the accuracy  
 3404 achieved by scalar type float. When using cuBLASEx’s default algorithm CUBLAS\_GEMM\_DEFAULT  
 3405 (rather than algorithm CUBLAS\_GEMM\_ALG01\_TENSOR\_OP), which retains the float type and thus  
 3406 meets the accuracy expected from the computation, we achieve a speedup of  $1.11 \times$  over cuBLASEx.  
 3407 For the interested reader, we report in our Appendix, Section D.2, the runtime of cuBLASEx for all  
 3408 its algorithm variants, including reports for the accuracy achieved by the different variants.

3409 The reason for the better performance of our approach over NVIDIA and Intel libraries is most  
 3410 likely because our approach allows generating code that is also optimized (auto-tuned) for data  
 3411 characteristics, which is important for high performance [Tillet and Cox 2017]. In contrast, the  
 3412 vendor libraries usually rely for each computation on pre-implemented implementations each  
 3413 optimized toward only average high performance for a range of data characteristics (size, memory  
 3414 layout, etc). By relying on these fixed, pre-implemented implementations, the libraries avoid the  
 3415 auto-tuning overhead. However, auto-tuning is often amortized, particularly for deep learning  
 3416 computations – the main target of libraries NVIDIA cuDNN and Intel oneDNN – because the auto-  
 3417 tuned implementations are re-used in many program runs. Moreover, we achieve better performance  
 3418 for convolutions (Figure 25), because the libraries re-use optimizations for these computations  
 3419 originally intended for linear algebra routines [Li et al. 2016], while our optimization space (Table 1)  
 3420 is designed as general and not oriented toward linear algebra.

3421 Compared to the EKR library (Figure 27), we achieve higher performance, because the EKR’s Java  
 3422 implementation inefficiently handles memory: the library is implemented using Java’s ArrayList  
 3423 data structure which is convenient to use for the Java programmer, but inefficient in terms of  
 3424 performance, because the structure internally performs costly memory reallocations.

3425 *Portability.* Similarly to polyhedral compilers PPCG and Pluto, the domain-specific approaches  
 3426 work for particular architectures only. The domain-specific approaches are also restricted to a

3429 <sup>24</sup> Since the Intel and NVIDIA libraries are not open source, we cannot explain their performance behavior with certainty.

3431 narrow set of studies, e.g., only linear algebra routines as NVIDIA cuBLAS and Intel oneMKL or  
 3432 only data mining example PRL as EKR. Consequently, the approaches achieve for these unsupported  
 3433 studies the lowest portability of only 0 in Figure 30 as well as Figures 31-35. For their target  
 3434 studies, domain-specific approaches can achieve high portability. This is because the approaches are  
 3435 specifically designed and optimized toward these studies, e.g., via application-specific assembly-level  
 3436 optimizations which are currently beyond the scope of our work.

3437  
 3438 *Productivity.* Listing 4 shows the implementation of MatVec in domain-specific approach NVIDIA  
 3439 cuBLAS; the implementation of MatVec in other domain-specific approaches, e.g., Intel oneMKL,  
 3440 is analogous to the implementation in Listing 4. We consider domain-specific approaches as  
 3441 most productive for their target domain: in the case of MatVec, the user simply calls the high-level  
 3442 function `cublasSgemv` and passes to it the input matrices (omitted via ellipsis in the listing) together  
 3443 with some meta information (memory layout of matrices, etc); cuBLAS then automatically starts  
 3444 the GPU computation for MatVec.

3445 Besides the fact that domain-specific approaches typically target only particular target architec-  
 3446 tures, a further fundamental productivity issue of domain-specific approaches is that they can only  
 3447 be used for a narrow class of computations, e.g., only linear algebra routines as NVIDIA cuBLAS and  
 3448 Intel oneMKL. Moreover, in the case of domain-specific libraries from NVIDIA and Intel, it is often  
 3449 up to the user to manually choose among different, semantically equal but differently performing  
 3450 implementations for high performance. For example, the cuBLAS library offers three different  
 3451 routines for computing matrix multiplications – routines `cublasSgemm` (part of standard cuBLAS),  
 3452 `cublasGemmEx` (part of the cuBLASEx extension of cuBLAS), and routine `cublasLtMatmul` (part of  
 3453 the cuBLASLt extension) – and the routines often also offer different, so-called *algorithms* (e.g., 42  
 3454 algorithm variants in the case cuBLASEx) which impact the internal optimization process. When  
 3455 striving for the highest performance potentials of libraries, the user is in charge of naively testing  
 3456 each possible combination of routine and algorithm variant (as we have done in Figures 24-29  
 3457 to make experimenting challenging for us). In addition, the user must be aware that different  
 3458 combinations of routines and algorithms can produce results of reduced accuracy (as discussed  
 3459 above), which can be critical for accuracy-sensitive use cases.

3460

3461 1 `cublasSgemv( /* ... */ );`

3462  
 3463 Listing 4. cuBLAS program expressing Matrix-Vector Multiplication (MatVec)

3464

3465

3466

## 3467 6 RELATED WORK

3468 Three major classes of approaches currently focus on code generation and optimization for data-  
 3469 parallel computations: 1) scheduling, 2) polyhedral, and 3) functional. In the following, we compare  
 3470 in Sections 6.1-6.3 our approach to each of these three classes – in terms of *performance*, *portability*,  
 3471 and *productivity*. In contrast to Section 5, which has compared our approach against these classes  
 3472 experimentally, this section is focussed on discussions in a more general, non-experimental context.  
 3473 Afterwards, we outline domain-specific approaches in Section 6.4, which are specifically designed  
 3474 and optimized toward their target application domains. In Section 6.5, we briefly outline approaches  
 3475 focussing on optimizations at the algorithmic level of abstraction which we consider as higher-level  
 3476 optimizations than proposed by our approach and as greatly combinable with our work. Finally,  
 3477 we discuss in Section 6.6 the differences between our approach introduced in this paper and the  
 3478 already existing work on MDHs.

3479

## 3480 6.1 Scheduling Approaches

3481 Popular examples of scheduling approaches include TVM [Chen et al. 2018a], Halide [Ragan-Kelley  
 3482 et al. 2013], Elevate [Hagedorn et al. 2020b], DaCe [Ben-Nun et al. 2019], Tiramisu [Baghdadi et al.  
 3483 2019], ChiLL [Chen et al. 2008; Khan et al. 2013], Clay [Bagnères et al. 2016], UTF [Kelly and Pugh  
 3484 1998], URUK [Girbal et al. 2006], Fireiron [Hagedorn et al. 2020a], DISTAL [Yadav et al. 2022], and  
 3485 LoopStack [Wasti et al. 2022]. While scheduling approaches usually achieve high performance, they  
 3486 often have difficulties with achieving portability and productivity, as we discuss in the following.<sup>25</sup>  
 3487

3488 *Performance.* Scheduling approaches usually achieve high performance. For this, the approaches  
 3489 incorporate human expert knowledge into their optimization process which is based on two  
 3490 major steps: 1) a human expert implements an optimization program (a.k.a *schedule*) in a so-called  
 3491 *scheduling language* – the program specifies the basic optimizations to perform, such as tiling  
 3492 and parallelization; 2) an auto-tuning system (or, alternatively, a human hardware expert) chooses  
 3493 performance-critical parameter values of the optimizations implemented in the schedule, e.g.,  
 3494 particular values of tile sizes and concrete numbers of threads.

3495 Our experiments in Section 5 show that compared to scheduling approach TVM (using its recent  
 3496 Ansor optimizer [Zheng et al. 2020a] for schedule generation), our approach achieves competitive  
 3497 and sometimes even better performance, e.g., speedups up to 2.22× on GPU and 3.55× on CPU over  
 3498 TVM+Ansor for computations taken from TVM’s favorable application domain (deep learning).  
 3499 Section 5 discusses that our better performance is due to the design and structure of our general  
 3500 optimization space (Table 1) which can be efficiently explored fully automatically using state-of-  
 3501 the-art auto tuning techniques [Rasch et al. 2021]. We focus on TVM in our experiments (rather  
 3502 than, e.g. Halide) to make experimenting challenging for us: TVM+Ansor has proved to achieve  
 3503 higher performance on GPUs and CPUs than popular state-of-practice approaches [Zheng et al.  
 3504 2020a], including Halide, pyTorch [Paszke et al. 2019], and the recent FlexTensor optimizer [Zheng  
 3505 et al. 2020b].

3506 Recent approach TensorIR [Feng et al. 2022] is a compiler for deep learning computations that  
 3507 achieves higher performance than TVM on NVIDIA GPUs. However, this performance gain over  
 3508 TVM is mainly achieved by exploiting the domain-specific *tensor core* [NVIDIA 2017] extensions  
 3509 of NVIDIA GPUs, which compute in hardware the multiplications of small, low-precision 4 × 4  
 3510 matrices. For this, TensorIR introduces the concept of *blocks* which represent sub-computations,  
 3511 e.g., for computing matrix multiplication on 4 × 4 sub-matrices. These blocks are then mapped by  
 3512 TensorIR to domain-specific hardware extensions, such as NVIDIA’s tensor cores, leading to high  
 3513 performance.

3514 While domain-specific hardware extensions are not targeted by this paper, we can naturally ex-  
 3515 ploit them in our approach, similar to TensorIR, as we plan for our future work: the sub-computations  
 3516 targeted by the current hardware extensions, like matrix multiplication on 4 × 4 matrices, can be  
 3517 straightforwardly expressed in our approach (Figure 16). Thus, we can match these sub-expressions  
 3518 in our low-level representation and map them to hardware extensions in our generated code. For  
 3519 this, instead of relying on a full partitioning in our low-level representation (as in Figure 17) such  
 3520 that we can apply scalar function  $f$  to the fully de-composed data (consisting of a single data  
 3521 element only in the case of a full partitioning), we plan to rely on a coarser-grained partitioning  
 3522 schema, e.g., down to only 4 × 4 matrices (rather than 1 × 1 matrices, as in the case of a full parti-  
 3523 tioning). This allows us replacing scalar function  $f$  (which in the case of matrix multiplication is

3524 <sup>25</sup> Rasch et al. [2023] introduce (optionally) a scheduling language for MDH to incorporate expert knowledge into MDH’s  
 3525 optimization process, e.g., to achieve 1) better optimization, as an auto-tuning system might not always make the same  
 3526 high-quality optimization decisions as a human expert, or 2) faster auto-tuning, as some (or even all) optimization decisions  
 3527 might be made by the expert user and thus are not left to the costly auto-tuner.

3529 a simple scalar multiplication  $*$ ) with the operation supported by the hardware extension, such  
3530 as matrix multiplication on  $4 \times 4$  matrices. We expect for our future work to achieve the same  
3531 advantages over TensorIR as over TVM, because apart from supporting domain-specific hardware  
3532 extensions, TensorIR is very similar to TVM.

3533  
3534 *Portability.* While scheduling approaches achieve high performance, they tend to struggle with  
3535 achieving portability. This is because even though the approaches often provide different, pre-  
3536 implemented backends (e.g., a CUDA backend to target NVIDIA GPUs and an OpenCL backend for  
3537 CPUs), they do not propose any structured methodology about how new backends can be added, e.g.,  
3538 for potentially upcoming architectures, with potentially deeper memory and core hierarchies than  
3539 GPUs and CPUs. This might be particularly critical (or requiring significant development effort)  
3540 for the application area of deep learning which is the main target of many scheduling approaches,  
3541 e.g., TVM and TensorIR, and for which new architectures are arising continuously [Hennessy and  
3542 Patterson 2019].

3543 In contrast, we introduce in this paper a formally precise recipe for correct-by-construction code  
3544 generation in different backends (including OpenMP, CUDA, and OpenCL), generically in the target  
3545 architecture: we introduce an architecture-agnostic low-level representation (Section 3) as target  
3546 for our high-level programs (Section 2), and we describe formally how our high-level programs are  
3547 automatically lowered to our low-level representation (Section 4), based on the architecture-agnostic  
3548 optimization space in Table 1. Our Appendix, Section E, outlines how executable, imperative-style  
3549 program code is straightforwardly generated from low-level expressions, which we plan to discuss  
3550 and illustrate in detail in our future work.

3551  
3552 *Productivity.* Scheduling approaches rely on a two-step optimization process, as discussed  
3553 above: implementing a schedule (first step) and choosing optimized values of performance-critical  
3554 parameters within that schedule (second step). While the second step often can be easily automa-  
3555 tized, e.g., via auto-tuning [Chen et al. 2018b], the first step – implementing a schedule – usually  
3556 has to be conducted manually by the user for high performance, which requires expert knowledge  
3557 and thus hinders productivity. The lack of formal foundation of many scheduling approaches  
3558 further complicates implementing schedules for the user, as implementation becomes error prone  
3559 and hardly predictable. For example, Fireiron’s schedules can achieve high performance, close to  
3560 GPUs’ peak, but schedules in Fireiron can easily generate incorrect low-level code: Fireiron cannot  
3561 guarantee that optimizations expressed in its scheduling language are semantics preserving, e.g.,  
3562 based on a formal foundation as done in this work, making programming Fireiron’s schedules error  
3563 prone and complex for the user. Similarly, TVM is sometimes unable to detect user errors in both  
3564 its high-level language (as discussed in Section 5.1) as well as scheduling language [Apache TVM  
3565 Community 2022e]. Safety in parallel programming is an ongoing major demand, in particular  
3566 from industry [Khronos 2022a].

3567 Auto schedulers, such as Halide’s optimization engine [Mullapudi et al. 2016] and TVM’s recent  
3568 Ansor [Zheng et al. 2020a], aim at automatically generating well-performing, correct schedules  
3569 for the user. However, a major flaw of the current auto schedulers is that even though they work  
3570 well for some computations (e.g., from deep learning, as TVM’s Ansor), they may perform worse  
3571 for others. For example, our approach achieves a speedup over TVM+Ansor of  $> 100\times$  already  
3572 for straightforward dot products. This is because Ansor does not exploit multiple thread blocks  
3573 and uses only a small number of threads for reduction computations, which is usually beneficial  
3574 for reductions as computed within convolutions and matrix multiplications used in deep learning  
3575 applications (because parallelization can be better exploited for outer loops of these computations),  
3576 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 (Section 4). 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 (Figure 19). 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 16. Also, TVM has difficulties with user-defined combine operators [Apache TVM Community 2022d] and thus crashes for example *Probabilistic Record Linkage* in Figure 16. 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], PPCG [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.,

3627 transformations toward fine-grained parallelization) might not achieve their full performance  
3628 potential when compiled with a particular polyhedral compiler.

3629 In contrast to the currently existing polyhedral compilers, we have designed our optimization  
3630 process as generic in goals: for example, our space is designed such that the degree of parallelization  
3631 (coarse, fine, ...) is completely auto-tunable for the particular combination of target architecture  
3632 and computation to optimize. We consider it as an interesting future work to investigate the strength  
3633 and weaknesses of the polyhedral model for expressing our generic optimization space.

3634 We see the second reason for potential performance issues in polyhedral compilers in their  
3635 difficulties with reduction-like computations. This is mainly caused by the fact that the polyhedral  
3636 model captures less semantic information than the high-level program representation introduced  
3637 in Section 2 of this paper: combine operators which are used to combine the intermediate results of  
3638 computations (e.g., operator + from Example 2 for combining the intermediate results of the dot  
3639 products within matrix multiplication) are not explicitly represented in the polyhedral model; the  
3640 polyhedral model is rather focussed on modeling memory accesses and their relative order only.  
3641 Most likely, these semantic information are missing in the polyhedral model, because polyhedral  
3642 approaches were originally intended to fully automatically optimize sequential code (such as Pluto  
3643 and PPCG) – extracting combine operators automatically from sequential code is challenging and  
3644 often even impossible (Rice’s theorem).

3645 In contrast, our proposed high-level representation explicitly captures combine operators (Figure  
3646 16), by requesting these operators explicitly from the user. This is important, because the  
3647 operators are often required for generating code that fully utilizes the highly parallel hardware of  
3648 state-of-the-art parallel architectures (GPUs, etc), as discussed in Section 5. Similarly to our  
3649 approach, polyhedral compiler TC also requests combine operators explicitly from the user. However,  
3650 TC is restricted to operators + (addition), \* (multiplication), min (minimum), and max (maximum)  
3651 only, thereby TC is not able to express important examples in Figure 16, e.g., PRL which is popular  
3652 in data mining. Moreover, TC outsources the computation of its combine operators to the NVIDIA  
3653 CUB library [NVIDIA 2022a]; most likely as a workaround, because TC relies on the polyhedral  
3654 model which is not designed to capture and exploit semantic information about combine opera-  
3655 tors for optimization. Thereby, TC is dependent on external approaches for computing combine  
3656 operators, which might not always be available (e.g., for upcoming architectures).

3657 Workarounds have been proposed by the polyhedral community to target reduction-like computa-  
3658 tions [Doerfert et al. 2015; Reddy et al. 2016]. However, these approaches are limited to a subset  
3659 of computations, e.g., by not supporting user-defined scalar types [Doerfert et al. 2015] (as required  
3660 for our PRL example in Figure 16), or by being limited to GPUs only [Reddy et al. 2016]. Comparing  
3661 the semantic information captured in the polyhedral model vs our MDH-based representation have  
3662 been the focus of discussions between polyhedral experts and MDH developers [Google SIG MLIR  
3663 Open Design Meeting 2020].

3664 *Portability.* The polyhedral approach, in its general form, is a framework offering transformation  
3665 rules (affine transformations), and each individual polyhedral compiler implements a set of such  
3666 transformations which are then instantiated (e.g., with particular tile sizes) and applied when  
3667 compiling a particular application. However, individual polyhedral compilers (e.g., PPCG and Pluto)  
3668 apply a fixed set of affine transformations, thereby directly generating a schedule that is optimized  
3669 for a particular target architecture only, e.g., only GPU (as PPCG) or only CPU (as Pluto), and  
3670 it remains open which affine transformations have to be used and how for other architectures,  
3671 e.g., upcoming accelerators for deep learning computations [Hennessy and Patterson 2019] with  
3672 potentially more complex memory and core hierarchies than GPUs and CPUs. Moreover, while  
3673 we introduce an explicit low-level representation (Section 3), the polyhedral approach does not  
3674

3676 introduce representations on different abstraction levels: the model relies on one representation  
 3677 that is transformed via affine transformations. Apart from the ability of our low-level representation  
 3678 to handle combine operators (which we consider as complex and important), we see the advantages  
 3679 of our explicit low-level representation in, for example, explicitly representing memory regions,  
 3680 which allows formally defining important correctness constraints, e.g., that GPU architectures  
 3681 allow combining the results of threads in designated memory regions only – shared and device  
 3682 memory, but not registers (outlined in Section C.1 of our Appendix). Furthermore, our low-level  
 3683 representation also allows straightforwardly generating executable code from it (focus of Section E  
 3684 in our Appendix, and planned to be discussed thoroughly in future work). In contrast, code  
 3685 generation from the polyhedral model has proven challenging [Bastoul et al. 2022; Grosser et al.  
 3686 2015; Vasilache et al. 2022].

3687 *Productivity.* Most polyhedral compilers achieve high user productivity, by fully automatically  
 3688 parallelizing and optimizing straightforward sequential code (as Pluto and PPCG). Our approach  
 3689 currently relies on a DSL (Domain-Specific Language) for expressing computations, as discussed  
 3690 in Section 2; thus, our approach can be considered as less productive than many polyhedral  
 3691 compilers. However, Rasch et al. [2020a,b] confirm that DSL programs in our approach can be  
 3692 automatically generated from sequential code (optionally annotated with simple, OpenMP-like  
 3693 directives for expressing combine operators, enabling advanced optimizations), by using polyhedral  
 3694 tool pet [Verdoolaege and Grosser 2012] as a frontend for our approach. Thereby, we are able to  
 3695 achieve the same, high user productivity as polyhedral compilers. We consider this direction –  
 3696 combining the polyhedral model with our approach – as promising, as it enables benefitting from  
 3697 the advantages of both directions: optimizing sequential programs and making them parallelizable  
 3698 using polyhedral techniques (like *loop skewing*, as also outlined in Section 6.5), and mapping the  
 3699 optimized and parallelizable code eventually to parallel architectures based on the concepts and  
 3700 methodologies introduced in this paper.

### 3702 6.3 Functional Approaches

3703 Functional approaches map data-parallel computations that are expressed via small, formally  
 3704 defined building blocks (a.k.a. patterns [Gorlatch and Cole 2011]), such as *map* and *reduce*, to the  
 3705 memory and core hierarchies of parallel architectures, based on a strong formal foundation. Notable  
 3706 functional approaches include Accelerate [Chakravarty et al. 2011], Obsidian [Svensson et al. 2011],  
 3707 so-called *skeleton libraries* [Aldinucci et al. 2017; Enmyren and Kessler 2010; Ernstsson et al. 2018;  
 3708 Steuwer et al. 2011], and the modern Lift approach [Steuwer et al. 2015] (recently also known as  
 3709 RISE [Steuwer et al. 2022]).

3710 In the following, as functional approaches usually follow the same basic concepts and methodologies,  
 3711 we focus on comparing to Lift, because Lift is more recent than, e.g., Accelerate and  
 3712 Obsidian.

3713 *Performance.* Functional approaches tend to struggle with achieving their full performance  
 3714 potential, often caused by the design of their optimization spaces. For example, analogously to our  
 3715 approach, functional approach Lift relies on an internal low-level representation [Steuwer et al.  
 3716 2017] that is used as target for Lift’s high-level programs. However, Lift’s transformation process,  
 3717 from high level to low level, turned out to be challenging: Lift’s lowering process relies on an  
 3718 infinitely large optimization space – identifying a well-performing configuration within that space  
 3719 is too complex to be done automatically, in general, due to the space’s large and complex structure.  
 3720 As a workaround, Lift currently uses approach Elevate [Hagedorn et al. 2020b] to incorporate  
 3721 user knowledge into the optimization process; however, at the cost of productivity, as manually  
 3722 expressing optimization is challenging, particularly for non-expert users.

3725 In contrast, our optimization process is designed as auto-tunable (Table 1), thereby achieving fully  
3726 automatically high performance, as confirmed in our experiments (Section 5), without involving the  
3727 user for optimization decisions. In particular, our previous work already showed that our approach –  
3728 even in its original version [Rasch and Gorlatch 2016; Rasch et al. 2019a] – can significantly  
3729 outperform Lift on GPU and CPU [Rasch et al. 2019a]. Our performance advantage over Lift is  
3730 mainly caused by the design of our optimization process: relying on formally defined tuning  
3731 parameters (Table 1) – rather than on formal transformation rules that span a too large and complex  
3732 search space, as in Lift – thereby contributing to a simpler, fully auto-tunable optimization process.  
3733

3734 *Portability.* The current functional approaches usually are designed and optimized toward code  
3735 generation in a particular programming model only. For example, Lift inherently relies on the  
3736 OpenCL programming model, because OpenCL works for multiple kinds of architectures: NVIDIA  
3737 GPU, Intel CPU, etc. However, we see two major disadvantages in addressing the portability issue via  
3738 OpenCL only: 1) GPU-specific optimizations (such as *shuffle operations* [NVIDIA 2018]) are available  
3739 only in the CUDA programming model, but not OpenCL; 2) the set of OpenCL-compatible devices is  
3740 broad but still limited; in particular, in the *new golden age for computer architectures* [Hennessy and  
3741 Patterson 2019], upcoming architectures are arising continuously and may not support the OpenCL  
3742 standard. We consider targeting new programming models as challenging for Lift, as its formal  
3743 low-level representation is inherently designed for OpenCL [Steuwer et al. 2017]; targeting further  
3744 programming models with Lift would require the design and implementation of new low-level  
3745 representations, which we do not consider as straightforward.

3746 To allow easily targeting new programming models with our approach, we have designed our  
3747 formalism as generic in the target model: our low-level representation (Figure 19) and optimization  
3748 space (Table 1) are designed and optimized toward an abstract system model (Definition 11) which  
3749 is capable of representing the device models of important programming approaches, including  
3750 OpenMP, CUDA, and OpenCL (Example 11). Furthermore, we have designed our high- and low-  
3751 level representations as minimalistic (Figures 15 and 19), e.g., by relying on three higher-order  
3752 functions only for expressing programs at the high abstraction level, which simplifies and reduces  
3753 the development effort for implementing code generators for programming models.

3754 In addition, we believe that compared to our approach, the following basic design decisions  
3755 of Lift (and similar functional approaches) complicate the process of code generation for them  
3756 and increase the development effort for implementing code generators: 1) relying on a vast set  
3757 of small patterns for expressing computations, rather than aiming at a minimalistic design as  
3758 we do (also discussed in Section 5.3); 2) relying on complex function nestings and compositions  
3759 for expressing computations, rather than avoiding nesting and relying on a fixed composition  
3760 structure of functions, as in our approach (Figure 5); 3) requiring new patterns for targeting new  
3761 classes of data-parallel computations (such as patterns *slide* and *pad* for stencils [Hagedorn et al.  
3762 2018]), which have to be non-trivially integrated into Lift’s type and optimization system (often via  
3763 extensions of the systems [Hagedorn et al. 2018; Remmelt et al. 2016]), instead of relying on a fixed  
3764 set of expressive patterns (Figure 15) and generalized optimizations (Table 1) that work for various  
3765 kinds of data-parallel computations (Figure 16); 4) expressing high-level and low-level concepts in  
3766 the same language, instead of separating high-level and low-level concepts for a more structured  
3767 and thus simpler code generation process (Figure 4). We consider these four design decisions as  
3768 disadvantageous for code generation, because they require from a code generator handling various  
3769 kinds of patterns (decision 1), and the patterns need to be translated to significantly different code  
3770 variants, depending on their nesting level and composition order (decision 2). Moreover, each  
3771 extension of patterns (decision 3) might affect code generation also for the already supported  
3772 patterns, because the existing patterns need to be combined with the new ones via composition and  
3773

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*  $\rightarrow$  *low-level language*  $\rightarrow$  *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 16).

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). Work-in-progress results 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 [Rasch et al. 2020a,b].

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>26</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>26</sup> Previous work [Rasch et al. 2021] shows that MDH (already in its original, proof-of-concept implementation) achieves higher performance than *ATLAS*.

3823 oneDNN [Intel 2022b] work well for convolution computations on either NVIDIA GPU (cuDNN) or  
3824 Intel CPU (oneDNN), respectively.

3825 In the following, we discuss domain-specific approaches in terms of *performance*, *portability*, and  
3826 *productivity*.

3828 *Performance.* Domain-specific approaches, such as cuBLAS and cuDNN, usually achieve high  
3829 performance. This is because the approaches are hand-optimized by performance experts – on  
3830 the assembly level – to exploit the full performance potential of their target architecture. In our  
3831 experiments, we show that our approach often achieves competitive and sometimes even better  
3832 performance than domain-specific approaches provided by NVIDIA and Intel, which is mainly due  
3833 to their portability issues over different data characteristics, as we discuss in the next paragraph.  
3834

3835 *Portability.* Domain-specific approaches usually struggle with achieving portability over different  
3836 architectures. This is because the approaches are often implemented in architecture-specific  
3837 assembly code to achieve high performance. The domain-specific approaches often also struggle  
3838 with achieving portability over different characteristics of their input and output data: they usually  
3839 rely on a set of pre-implemented implementations that are each designed and optimized toward  
3840 average high performance over a range of input characteristic (e.g., their sizes). In contrast, our  
3841 approach (as well as many scheduling and polyhedral approaches) allow automatically optimizing  
3842 (auto-tuning) kernels for particular data characteristics, which is important for performance [Tillet  
3843 and Cox 2017], thereby often achieving higher performance than domain-specific approaches  
3844 for advanced data characteristics (small, uneven, irregularly shaped, . . .), e.g., as used in deep  
3845 learning. The costly time for auto-tuning is well amortized in many application areas, because the  
3846 auto-tuned implementations are re-used in many program runs. Moreover, auto-tuning avoids the  
3847 time-intensive and costly process of hand optimization by human experts.  
3848

3849 *Productivity.* Domain-specific approaches usually achieve highest productivity for their target  
3850 domain (e.g., linear algebra), by providing easy to use high-level abstractions. However, the ap-  
3851 proaches suffer from significant expressivity issues, because – per design – they are inherently  
3852 restricted to their target application domain only. Also, the approaches are often inherently bound  
3853 to only particular architectures, e.g., only GPU (as NVIDIA cuBLAS and cuDNN) or only CPU  
3854 (as Intel oneMKL and oneDNN). Domain-specific vendor libraries, such as NVIDIA cuBLAS and  
3855 Intel oneMKL, also tend to offer the user differently performing variants of computations; the  
3856 variants have to be naively tested by the user when striving for the full performance potentials of  
3857 approaches (as discussed in Section 5.4), which is cumbersome for the user.  
3858

## 3859 **6.5 Higher-Level Approaches**

3860 There is a broad range of existing work that is focused on higher-level optimizations than proposed  
3861 by this work. We consider such higher-level approaches as greatly combinable with our approach.  
3862 For example, the polyhedral approach is capable of expressing algorithmic-level optimizations, like  
3863 *loop skewing* [Wolf and Lam 1991], to make programs parallelizable; such optimizations are beyond  
3864 the scope of this work, but they can be combined with our approach as demonstrated by Rasch et al.  
3865 [2020a,b]. Similarly, we consider approaches introduced by Farzan and Nicolet [2019]; Frigo et al.  
3866 [1999]; Gunnels et al. [2001]; Yang et al. [2021], which also focus on algorithmic-level optimizations,  
3867 as greatly combinable with our approach: algorithmically optimizing user code according to the  
3868 approaches’ techniques, and using our methodologies to eventually map the optimized code to  
3869 low-level code for parallel architectures.  
3870

3872 Futhark [Henriksen et al. 2017], Dex [Paszke et al. 2021], and ATL [Liu et al. 2022] are further  
 3873 approaches focussed on high-level program transformations, like advanced *flattening* mechanisms  
 3874 [Henriksen et al. 2019], thereby optimizing programs at the algorithmic level of abstraction.  
 3875 We consider using our work as backend for these approaches as promising: the three approaches  
 3876 often struggle with mapping their algorithmically optimized program variants eventually to the  
 3877 multi-layered memory and core hierarchies of state-of-the-art parallel architectures, which is  
 3878 exactly the focus of this work.

## 3879 6.6 Existing Work on MDH

3880 Our work is inspired by the algebraic approach of Multi-Dimensional Homomorphisms (MDHs)  
 3881 which is introduced in the work-in-progress paper [Rasch and Gorlatch 2016]. The MDH approach,  
 3882 as presented in the previous work, relies on a semi-formal foundation and focuses on code generation  
 3883 for the OpenCL programming model only [Rasch et al. 2019a]. This work makes major contributions  
 3884 over the existing work on MDHs and its OpenCL code generation approach.

3885 We introduce a full formalization of MDH’s high-level program representation. In our new  
 3886 formalism, we rely on expressive typing: for example, we encode MDHs’ data sizes into our type  
 3887 system, e.g., by introducing both *index sets* for MDAs (Definition 1) and *index set functions* for  
 3888 combine operators (Definition 2), and we respect and maintain these sets and functions thoroughly  
 3889 during MDH computations (Definition 3). Our expressive typing significantly contributes to correct  
 3890 and simplified code generation, as all relevant type and data size information are contained in  
 3891 our formal, low-level program representation (Figure 19) from which we eventually generate exe-  
 3892 cutable program code (Section 3). In contrast, the existing MDH work considers multi-dimensional  
 3893 arrays (MDAs) of arbitrary sizes and dimensionalities to be all of the same, straightforward type,  
 3894 which has greatly simplified the design of the proof-of-concept MDH formalism introduced by Rasch  
 3895 and Gorlatch [2016] (in particular, the definition and usage of combine operators), but at the cost  
 3896 of significantly harder and error-prone code generation: all the missing, type-relevant information  
 3897 need to be elaborated by the implementer of the code generator in the existing MDH work, e.g.,  
 3898 allocation sizes of fast memory resources used for caching input data or for storing computed  
 3899 intermediate results. Furthermore, while the original MDH work [Rasch and Gorlatch 2016] is  
 3900 focused on introducing higher-order function `md_hom` only, this work in particular also introduces  
 3901 higher order functions `inp_view` and `out_view` (Section 2.3) which express input and output  
 3902 views in a formally structured and concise manner, and which are central building blocks for  
 3903 expressing computations (Figure 16). Also, by introducing and exploiting the index set concept  
 3904 for MDAs, we have improved the definition of the concatenation operator `++` (Example 1) toward  
 3905 commutativity, which is required for important optimizations. e.g., loop permutations (expressed  
 3906 via Parameters  $D1, S1, R1$  in Table 1).

3907 A further substantial improvement is the introduction of our low-level representation (Section 3).  
 3908 It relies on a novel combination of tuning parameters (Table 1) that enhance, generalize, and extend  
 3909 the existing, proof-of-concept MDH parameters which capture a subset of OpenCL-orientated  
 3910 features only [Rasch et al. 2019a]. Moreover, while the existing MDH work introduces formally only  
 3911 parameters for flexibly choosing numbers of threads [Rasch and Gorlatch 2016] (which corresponds  
 3912 to a very limited variant of our tuning parameter  $\theta$  in Table 1, because our parameter  $\theta$  also chooses  
 3913 numbers of memory tiles and is not restricted to OpenCL); the other OpenCL parameters are  
 3914 introduced and discussed by Rasch et al. [2019a] only informally, from a technical perspective. With  
 3915 our novel parameter set, we are able to target various kinds of programming models (e.g., also CUDA,  
 3916 as in Section 5) and also to express important optimizations that are beyond the existing work on  
 3917 MDH, e.g., optimizing the memory access pattern of MDH computations: for example, we achieve  
 3918 speedups  $> 2\times$  over existing MDH for the deep learning computations discussed in Section 5.  
 3919

3921 Our new tuning parameters are expressive enough to represent state-of-the-art, data-parallel  
3922 implementations, e.g., as generated by scheduling and polyhedral approaches (Figures 20-23), and  
3923 our experiments in Section 5 confirm that auto-tuning our parameters enables performance beyond  
3924 the state of the art, including hand-optimized solutions provided by vendors, which is not possible  
3925 when using the existing MDH approach. The expressivity of our parameters in particular also  
3926 enables comparing significantly differently optimized implementations (e.g., scheduling-optimized  
3927 vs. polyhedral-optimized, as in Section 3.5), based on the values of formally specified tuning  
3928 parameters, which we consider as promising for structured performance analysis in future work.  
3929 Moreover, our new low-level representation targets architectures that may have arbitrarily deep  
3930 memory and core hierarchies, by having optimized our representation toward an *Abstract System*  
3931 *Model* (Definition 11). In contrast, the existing MDH work is focused on OpenCL-compatible  
3932 architectures only.

3933 Our experimental evaluation extends previous MDH experiments by comparing also to the popu-  
3934 lar state-of-practice approach TVM which is attracting increasing attention from both academia [[Apache](#)  
3935 [Software Foundation 2021](#)] and industry [[OctoML 2022](#)]. Also, we compare to the popular poly-  
3936 hedral compilers PPCG and Pluto, as well as the currently newest versions of hand-optimized  
3937 high-performance libraries provided by vendors. Furthermore, we have included a real-world  
3938 case study in our experiments, considering the most time-intensive computations within the three  
3939 popular deep learning neural networks ResNet-50, VGG-16, and MobileNet; the study also includes  
3940 Capsule-style convolution computations, which are considered as challenging to optimize [[Barham](#)  
3941 and [Isard 2019](#)]. Moreover, Table 16 analyzes MDH’s expressivity using new examples: it shows  
3942 that MDH – based on the new contributions of this work (e.g., view functions) – is capable of  
3943 expressing computations bMatMul, MCC\_Capsule, Histo, scan, and MBBS, which have not been  
3944 expressed via MDH in previous work. Our experiments confirm that we achieve high performance  
3945 for bMatMul and MCC\_Capsule on GPUs and CPUs, and our future work aims to thoroughly analyze  
3946 our approach for computations Histo, scan, and MBBS in terms of performance, portability, and  
3947 productivity.

3948  
3949 

## 7 CONCLUSION

  
3950 We introduce a formal (de/re)-composition approach for data-parallel computations targeting  
3951 state-of-the-art parallel architectures. Our approach aims to combine three major advantages  
3952 over related approaches – performance, portability, and productivity – by introducing formal  
3953 program representations on both: 1) high level, for conveniently expressing – in one uniform  
3954 formalism – various kinds of data-parallel computations (including linear algebra routines, stencil  
3955 computations, data mining algorithms, and quantum chemistry computations), agnostic from  
3956 hardware and optimization details, while still capturing all information relevant for generating high-  
3957 performance program code; 2) low level, which allows uniformly reasoning – in the same formalism –  
3958 about optimized (de/re)-compositions of data-parallel computations targeting different kinds of  
3959 parallel architectures (GPUs, CPUs, etc). We lower our high-level representation to our low-level  
3960 representation, in a formally sound manner, by introducing a generic search space that is based on  
3961 performance-critical parameters. The parameters of our lowering process enable fully automatically  
3962 optimizing (auto-tuning) our low-level representations for a particular target architecture and  
3963 characteristics of the input and output data, and our low-level representation is designed such  
3964 that it can be straightforwardly transformed to executable program code in imperative-style  
3965 programming approaches (including OpenMP, CUDA, and OpenCL). Our experiments confirm  
3966 that due to the design and structure of our generic search space in combination with auto-tuning,  
3967 our approach achieves higher performance on GPUs and CPUs than popular state-of-practice  
3968 approaches, including hand-optimized libraries provided by vendors.

## 3970 8 FUTURE WORK

3971 We consider this work as a promising starting point for future directions. A major future goal is  
 3972 to extend our approach toward expressing and optimizing simultaneously multiple data-parallel  
 3973 computations (e.g., matrix multiplication followed by convolution), rather than optimizing computa-  
 3974 tions individually and thus independently from each other only (e.g., only matrix multiplication or  
 3975 only convolution). Such extension enables optimizations, such as *kernel fusion*, which is important  
 3976 for the overall application performance and considered as challenging [Fukuhara and Takimoto  
 3977 2022; Li et al. 2022; Wahib and Maruyama 2014]. We see this work as a promising foundation for  
 3978 our future goal, because it enables expressing and reasoning about different computations in the  
 3979 same formal framework. Targeting computations on sparse input/output data formats, inspired  
 3980 by Ben-Nun et al. [2017]; Hall [2020]; Kjolstad et al. [2017]; Pizzuti et al. [2020], is a further major  
 3981 goal, which requires extending our approach toward irregularly-shaped input and output data,  
 3982 similarly as done by Pizzuti et al. [2020]. Regarding our optimization process, we aim to introduce  
 3983 an analytical cost model for computations expressed in our formalism – based on operational  
 3984 semantics – thereby accelerating (or even avoiding) the auto-tuning overhead, similarly as done  
 3985 by Li et al. [2021]; Muller and Hoffmann [2021]. Moreover, we aim to incorporate machine-learning  
 3986 based methods into our optimization process [Leather et al. 2014], instead of relying on empirical  
 3987 auto-tuning methods only. To make our work better accessible for the community, we aim to  
 3988 implement our approach into *MLIR* [Lattner et al. 2021] which offers a reusable compiler infras-  
 3989 tructure. The contributions of this work give a precise, formal recipe of how to implement our  
 3990 introduced methods into approaches like *MLIR*. Moreover, relying on the *MLIR* framework will  
 3991 contribute to a structured code generation process in assembly-level programming models, like  
 3992 LLVM [Lattner and Adve 2004] and NVIDIA PTX [NVIDIA 2022i]. We consider targeting assembly  
 3993 languages as important for our future work: assembly code offers further, low-level optimization  
 3994 opportunities [Goto and Geijn 2008; Lai and Seznec 2013], thereby enabling our approach to poten-  
 3995 tially achieve higher performance than presented in this work for our MDH-generated CUDA and  
 3996 OpenCL code. Also, we aim to extend our approach toward distributed multi-device systems that  
 3997 are heterogeneous, inspired by dynamic load balancing approaches [Chen et al. 2010] and advanced  
 3998 data distributions techniques [Yadav et al. 2022]. Targeting domain-specific hardware extensions,  
 3999 such as *NVIDIA Tensor Cores* [NVIDIA 2017], is also an important goal for our future work, as such  
 4000 extensions allow significantly accelerating computations for the target of the extensions (e.g., deep  
 4001 learning [Markidis et al. 2018]).

## 4003 REFERENCES

- 4004 Marco Aldinucci, Marco Danelutto, Peter Kilpatrick, and Massimo Torquati. 2017. Fastflow: high-level and efficient streaming  
 4005 on multi-core. *Programming multi-core and many-core computing systems, parallel and distributed computing* (2017).
- 4006 Jason Ansel, Shoaib Kamil, Kalyan Veeramachaneni, Jonathan Ragan-Kelley, Jeffrey Bosboom, Una-May O'Reilly, and  
 4007 Saman Amarasinghe. 2014. OpenTuner: An Extensible Framework for Program Autotuning. In *Proceedings of the 23rd  
 4008 International Conference on Parallel Architectures and Compilation (PACT '14)*. Association for Computing Machinery,  
 4009 New York, NY, USA, 303–316. <https://doi.org/10.1145/2628071.2628092>
- 4010 Apache. 2022. TVM: Open Deep Learning Compiler Stack. <https://github.com/apache/tvm>.
- 4011 Apache Software Foundation. 2021. TVM and Open Source ML Acceleration Conference. <https://www.tvmcon.org>.
- 4012 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>.
- 4013 Apache TVM Community. 2022a. Bind reduce axis to blocks. <https://discuss.tvm.apache.org/t/bind-reduce-axis-to-blocks-2907>.
- 4014 Apache TVM Community. 2022b. Expressing nested reduce operations. <https://discuss.tvm.apache.org/t/expressing-nested-reduce-operations/8784>.
- 4015 Apache TVM Community. 2022c. Implementing Array Packing via cache\_read. <https://discuss.tvm.apache.org/t/implementing-array-packing-via-cache-read/13360>.

- 4019 Apache TVM Community. 2022d. Invalid comm\_reducer. <https://discuss.tvm.apache.org/t/invalid-comm-reducer/12788>.  
4020 Apache TVM Community. 2022e. Undetected parallelization issue. <https://discuss.tvm.apache.org/t/undetected-parallelization-issue/13224>.  
4021 Apache TVM Community. 2022f. Undetected type issue. <https://discuss.tvm.apache.org/t/undetected-type-issue/13223>.  
4022 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>.  
4023 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).  
4024 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>  
4025 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>  
4026 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>  
4027 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>  
4028 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>  
4029 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>  
4030 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)*.  
4031 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>  
4032 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.  
4033 Guy E. Blelloch. 1990. *Prefix Sums and Their Applications*. Technical Report CMU-CS-90-190. School of Computer Science, Carnegie Mellon University.  
4034 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>  
4035 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)  
4036 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.  
4037 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.  
4038 C++ reference. 2022. Date and time utilities. <https://en.cppreference.com/w/cpp/chrono>.  
4039 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.  
4040 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>  
4041 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.  
4042  
4043  
4044  
4045  
4046  
4047  
4048  
4049  
4050  
4051  
4052  
4053  
4054  
4055  
4056  
4057  
4058  
4059  
4060  
4061  
4062  
4063  
4064  
4065  
4066  
4067

- 4068 Long Chen, Oreste Villa, Sriram Krishnamoorthy, and Guang R. Gao. 2010. Dynamic load balancing on single- and  
 4069 multi-GPU systems. In *2010 IEEE International Symposium on Parallel Distributed Processing (IPDPS)*. 1–12. <https://doi.org/10.1109/IPDPS.2010.5470413>
- 4070 Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang,  
 4071 Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. 2018a. TVM: An Automated End-to-End Optimizing  
 4072 Compiler for Deep Learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*.  
 4073 USENIX Association, Carlsbad, CA, 578–594. <https://www.usenix.org/conference/osdi18/presentation/chen>
- 4074 Tianqi Chen, Lianmin Zheng, Eddie Yan, Ziheng Jiang, Thierry Moreau, Luis Ceze, Carlos Guestrin, and Arvind Kr-  
 4075 ishnamurthy. 2018b. Learning to Optimize Tensor Programs. In *Advances in Neural Information Processing Systems*,  
 4076 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc.  
 4077 <https://proceedings.neurips.cc/paper/2018/file/8b5700012be65c9da25f49408d959ca0-Paper.pdf>
- 4078 Peter Christen. 2012. *Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*.  
 Springer Publishing Company, Incorporated.
- 4079 Krzysztof Ciesielski. 1997. *Set theory for the working mathematician*. Number 39. Cambridge University Press.
- 4080 Basile Clément and Albert Cohen. 2022. End-to-End Translation Validation for the Halide Language. In *OOPSLA 2022 -  
 4081 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>
- 4082 MURRAY I. COLE. 1995. PARALLEL PROGRAMMING WITH LIST HOMOMORPHISMS. *Parallel Processing Letters* 05, 02  
 4083 (1995), 191–203. <https://doi.org/10.1142/S0129626495000175> arXiv:<https://doi.org/10.1142/S0129626495000175>
- 4084 Haskell B. Curry. 1980. Some Philosophical Aspects of Combinatory Logic. In *The Kleene Symposium*, Jon Barwise, H. Jerome  
 4085 Keisler, and Kenneth Kunen (Eds.). Studies in Logic and the Foundations of Mathematics, Vol. 101. Elsevier, 85–101.  
 4086 [https://doi.org/10.1016/S0049-237X\(08\)71254-0](https://doi.org/10.1016/S0049-237X(08)71254-0)
- 4087 Johannes Doerfert, Kevin Streit, Sebastian Hack, and Zino Benaissa. 2015. Polly’s Polyhedral Scheduling in the Presence of  
 4088 Reductions. *CoRR* abs/1505.07716 (2015). arXiv:1505.07716 <http://arxiv.org/abs/1505.07716>
- 4089 Vincent Dumoulin and Francesco Visin. 2018. A guide to convolution arithmetic for deep learning. arXiv:stat.ML/1603.07285
- 4090 Johan Enmyren and Christoph W. Kessler. 2010. SkePU: A Multi-Backend Skeleton Programming Library for Multi-GPU  
 4091 Systems. In *Proceedings of the Fourth International Workshop on High-Level Parallel Programming and Applications (HLPP  
 ’10)*. Association for Computing Machinery, New York, NY, USA, 5–14. <https://doi.org/10.1145/1863482.1863487>
- 4092 August Ernstsson, Lu Li, and Christoph Kessler. 2018. SkePU 2: Flexible and Type-Safe Skeleton Programming for Hetero-  
 4093 geneous Parallel Systems. *International Journal of Parallel Programming* 46, 1 (2018), 62–80. <https://doi.org/10.1007/s10766-017-0490-5>
- 4094 Facebook Research. 2022. Tensor Comprehensions. <https://github.com/facebookresearch/TensorComprehensions>.
- 4095 Azadeh Farzan and Victor Nicolet. 2019. Modular Divide-and-Conquer Parallelization of Nested Loops. In *Proceedings of  
 4096 the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2019)*. Association for  
 4097 Computing Machinery, New York, NY, USA, 610–624. <https://doi.org/10.1145/3314221.3314612>
- 4098 Siyuan Feng, Bohan Hou, Hongyi Jin, Wuwei Lin, Junru Shao, Ruihang Lai, Zihao Ye, Lianmin Zheng, Cody Hao Yu,  
 4099 Yong Yu, and Tianqi Chen. 2022. TensorIR: An Abstraction for Automatic Tensorized Program Optimization. <https://doi.org/10.48550/ARXIV.2207.04296>
- 4100 M. Frigo and S.G. Johnson. 1998. FFTW: an adaptive software architecture for the FFT. In *Proceedings of the 1998 IEEE  
 4101 International Conference on Acoustics, Speech and Signal Processing, ICASSP ’98 (Cat. No.98CH36181)*, Vol. 3. 1381–1384  
 4102 vol.3. <https://doi.org/10.1109/ICASSP.1998.681704>
- 4103 M. Frigo, C.E. Leiserson, H. Prokop, and S. Ramachandran. 1999. Cache-oblivious algorithms. In *40th Annual Symposium on  
 4104 Foundations of Computer Science (Cat. No.99CB37039)*. 285–297. <https://doi.org/10.1109/SFCS.1999.814600>
- 4105 Junji Fukuhara and Munehiro Takimoto. 2022. Automated Kernel Fusion for GPU Based on Code Motion. In *Proceedings of  
 4106 the 23rd ACM SIGPLAN/SIGBED International Conference on Languages, Compilers, and Tools for Embedded Systems (LCTES  
 4107 2022)*. Association for Computing Machinery, New York, NY, USA, 151–161. <https://doi.org/10.1145/3519941.3535078>
- 4108 Sylvain Girbal, Nicolas Vasilache, Cédric Bastoul, Albert Cohen, David Parello, Marc Sigler, and Olivier Temam. 2006.  
 4109 Semi-Automatic Composition of Loop Transformations for Deep Parallelism and Memory Hierarchies. *International  
 4110 Journal of Parallel Programming* 34, 3 (2006), 261–317. <https://doi.org/10.1007/s10766-006-0012-3>
- 4111 Horacio González-Vélez and Mario Leyton. 2010. A survey of algorithmic skeleton frameworks: high-level structured parallel  
 4112 programming enablers. *Software: Practice and Experience* 40, 12 (2010), 1135–1160. <https://doi.org/10.1002/spe.1026>  
 arXiv:<https://onlinelibrary.wiley.com/doi/10.1002/spe.1026>
- 4113 Google SIG MLIR Open Design Meeting. 2020. Using MLIR for Multi-Dimensional Homomorphisms. <https://drive.google.com/file/d/1bS4vapyzf7705wWj7t3WzwWkcbxZyF6/view>
- 4114 Sergei Gorlatch. 1999. Extracting and implementing list homomorphisms in parallel program development. *Science of  
 4115 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)

- 4117 Sergei Gorlatch and Murray Cole. 2011. Parallel skeletons. In *Encyclopedia of parallel computing*. Springer-Verlag GmbH, 1417–1422.
- 4118 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>
- 4121 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>
- 4122 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>
- 4123 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>
- 4126 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>
- 4129 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>
- 4132 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>
- 4135 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>
- 4137 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>
- 4139 Maurice H Halstead. 1977. *Elements of Software Science (Operating and programming systems series)*. Elsevier Science Inc.
- 4140 Mark Harris et al. 2007. Optimizing Parallel Reduction in CUDA. *NVIDIA Developer Technology* (2007).
- 4141 Haskell Wiki. 2013. Parameter Order. [https://wiki.haskell.org/Parameter\\_order](https://wiki.haskell.org/Parameter_order)
- 4142 Haskell.org. 2022. Haskell: An advanced, purely functional programming language. <https://www.haskell.org>.
- 4143 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>
- 4144 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>
- 4146 Troels Henriksen, Sune Hellfritsch, 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>
- 4148 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>
- 4151 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>
- 4152 K Hentschel et al. 2008. Das Krebsregister-Manual der Gesellschaft der epidemiologischen Krebsregister in Deutschland e.V. Zuckschwerdt Verlag.
- 4156 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>
- 4157 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>
- 4161 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>
- 4163 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)*.

- 4166      Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/509593.509638>
- 4167      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>.
- 4169      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>.
- 4171      Intel. 2022b. oneDNN. [https://oneapi-src.github.io/oneDNN/group\\_dnnl\\_api.html](https://oneapi-src.github.io/oneDNN/group_dnnl_api.html).
- 4172      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>.
- 4173      Wayne Kelly and William Pugh. 1998. *A framework for unifying reordering transformations*. Technical Report. Technical Report UMIACS-TR-92-126.1.
- 4174      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>
- 4175      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>.
- 4176      Khronos. 2022b. OpenCL: Open Standard For Parallel Programming of Heterogeneous Systems. <https://www.khronos.org/opencl/>.
- 4177      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>
- 4178      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>
- 4179      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)
- 4180      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>
- 4181      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>
- 4182      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>
- 4183      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>
- 4184      Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, 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>
- 4185      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>
- 4186      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>
- 4187      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>

- 4215 Xiaqing Li, Guangyan Zhang, H. Howie Huang, Zhufan Wang, and Weimin Zheng. 2016. Performance Analysis of  
4216 GPU-Based Convolutional Neural Networks. In *2016 45th International Conference on Parallel Processing (ICPP)*. 67–76.  
4217 <https://doi.org/10.1109/ICPP.2016.15>
- 4218 Amanda Liu, Gilbert Louis Bernstein, Adam Chlipala, and Jonathan Ragan-Kelley. 2022. Verified Tensor-Program Op-  
4219 timization via High-Level Scheduling Rewrites. *Proc. ACM Program. Lang.* 6, POPL, Article 55 (jan 2022), 28 pages.  
4220 <https://doi.org/10.1145/3498717>
- 4221 Stefano Markidis, Steven Wei Der Chien, Erwin Laure, Ivy Bo Peng, and Jeffrey S. Vetter. 2018. NVIDIA Tensor Core  
4222 Programmability, Performance & Precision. In *2018 IEEE International Parallel and Distributed Processing Symposium  
4223 Workshops (IPDPSW)*. 522–531. <https://doi.org/10.1109/IPDPSW.2018.00091>
- 4224 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>
- 4225 Kathryn S. McKinley, Steve Carr, and Chau-Wen Tseng. 1996. Improving Data Locality with Loop Transformations. *ACM  
4226 Trans. Program. Lang. Syst.* 18, 4 (jul 1996), 424–453. <https://doi.org/10.1145/233561.233564>
- 4227 Xinxin Mei, Kaiyong Zhao, Chengjian Liu, and Xiaowen Chu. 2014. Benchmarking the Memory Hierarchy of Modern  
4228 GPUs. In *Network and Parallel Computing*, Ching-Hsien Hsu, Xuanhua Shi, and Valentina Salapura (Eds.). Springer Berlin  
Heidelberg, Berlin, Heidelberg, 144–156.
- 4229 Michael Kruse. 2022. Polyhedral Parallel Code Generation. <https://github.com/Meinersbur/ppcg>, commit = 8a74e46, date =  
4230 19.11.2020.
- 4231 Ravi Teja Mullapudi, Andrew Adams, Dillon Sharlet, Jonathan Ragan-Kelley, and Kayvon Fatahalian. 2016. Automatically  
4232 Scheduling Halide Image Processing Pipelines. *ACM Trans. Graph.* 35, 4, Article 83 (jul 2016), 11 pages. <https://doi.org/10.1145/2897824.2925952>
- 4233 Stefan K. Muller and Jan Hoffmann. 2021. Modeling and Analyzing Evaluation Cost of CUDA Kernels. *Proc. ACM Program.  
4234 Lang.* 5, POPL, Article 25 (Jan. 2021), 31 pages. <https://doi.org/10.1145/3434306>
- 4235 NVIDIA. 2017. Programming Tensor Cores in CUDA 9. [https://developer.nvidia.com/blog/](https://developer.nvidia.com/blog/programming-tensor-cores-cuda-9/)
- 4236 NVIDIA. 2018. Warp-level Primitives. <https://developer.nvidia.com/blog/using-cuda-warp-level-primitives/>
- 4237 NVIDIA. 2022a. CUB. <https://docs.nvidia.com/cuda/cub/>.
- 4238 NVIDIA. 2022b. cuBLAS. <https://developer.nvidia.com/cublas>.
- 4239 NVIDIA. 2022c. cuBLAS – BLAS-like Extension. <https://docs.nvidia.com/cuda/cublas/index.html#blas-like-extension>
- 4240 NVIDIA. 2022d. cuBLAS – Using the cuBLASLT API. <https://docs.nvidia.com/cuda/cublas/index.html#using-the-cublaslt-api>
- 4241 NVIDIA. 2022e. CUDA Deep Neural Network library. <https://developer.nvidia.com/cudnn>
- 4242 NVIDIA. 2022f. CUDA Programming Guide. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/>.
- 4243 NVIDIA. 2022g. CUDA Toolkit Documentation. <https://docs.nvidia.com/cuda/>.
- 4244 NVIDIA. 2022h. NVRTC. <https://docs.nvidia.com/cuda/nvrtc>.
- 4245 NVIDIA. 2022i. Parallel Thread Execution ISA. <https://docs.nvidia.com/cuda/parallel-thread-execution>.
- 4246 OctoML. 2022. Accelerated Machine Learning Deployment. <https://octoml.ai>.
- 4247 Geraldo F. Oliveira, Juan Gómez-Luna, Lois Orosa, Saugata Ghose, Nandita Vijaykumar, Ivan Fernandez, Mohammad  
4248 Sadrosadati, and Onur Mutlu. 2021. DAMOV: A New Methodology and Benchmark Suite for Evaluating Data Movement  
4249 Bottlenecks. *IEEE Access* 9 (2021), 134457–134502. <https://doi.org/10.1109/ACCESS.2021.3110993>
- 4250 OpenMP. 2022. The OpenMP API Specification for Parallel Programming. <https://www.openmp.org>.
- 4251 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin,  
4252 Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison,  
Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An  
4253 Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems*,  
H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc.  
<https://proceedings.neurips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf>
- 4254 Adam Paszke, Daniel D. Johnson, David Duvenaud, Dimitrios Vytiniotis, Alexey Radul, Matthew J. Johnson, Jonathan Ragan-  
4255 Kelley, and Dougal Maclaurin. 2021. Getting to the Point: Index Sets and Parallelism-Preserving Autodiff for Pointful  
4256 Array Programming. *Proc. ACM Program. Lang.* 5, ICFP, Article 88 (aug 2021), 29 pages. <https://doi.org/10.1145/3473593>
- 4257 S.J. Pennycook, J.D. Sewall, and V.W. Lee. 2019. Implications of a metric for performance portability. *Future Generation  
4258 Computer Systems* 92 (2019), 947–958. <https://doi.org/10.1016/j.future.2017.08.007>
- 4259 Phitchaya Mangpo Phothilimthana, Archibald Samuel Elliott, An Wang, Abhinav Jangda, Bastian Hagedorn, Henrik Barthels,  
4260 Samuel J. Kaufman, Vinod Grover, Emina Torlak, and Rastislav Bodik. 2019. Swizzle Inventor: Data Movement Synthesis  
4261 for GPU Kernels. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming  
4262 Languages and Operating Systems (ASPLOS '19)*. Association for Computing Machinery, New York, NY, USA, 65–78.  
<https://doi.org/10.1145/3297858.3304059>

- 4264 Federico Pizzuti, Michel Steuwer, and Christophe Dubach. 2020. Generating Fast Sparse Matrix Vector Multiplication from a  
 4265 High Level Generic Functional IR. In *Proceedings of the 29th International Conference on Compiler Construction (CC 2020)*.  
 4266 Association for Computing Machinery, New York, NY, USA, 85–95. <https://doi.org/10.1145/3377555.3377896>
- 4267 Victor Podlozhnyuk. 2007. Image Convolution with CUDA. *NVIDIA Corporation White Paper* (2007).
- 4268 M. Puschel, J.M.F. Moura, J.R. Johnson, D. Padua, M.M. Veloso, B.W. Singer, Jianxin Xiong, F. Franchetti, A. Gacic, Y.  
 4269 Voronenko, K. Chen, R.W. Johnson, and N. Rizzolo. 2005. SPIRAL: Code Generation for DSP Transforms. *Proc. IEEE* 93, 2  
 4270 (2005), 232–275. <https://doi.org/10.1109/JPROC.2004.840306>
- 4271 Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Amarasinghe. 2013.  
 4272 Halide: A Language and Compiler for Optimizing Parallelism, Locality, and Recomputation in Image Processing Pipelines.  
 4273 In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '13)*.  
 4274 Association for Computing Machinery, New York, NY, USA, 519–530. <https://doi.org/10.1145/2491956.2462176>
- 4275 Ari Rasch and Sergei Gorlatch. 2016. Multi-Dimensional Homomorphisms and Their Implementation in OpenCL. In  
 4276 *International Workshop on High-Level Parallel Programming and Applications (HLPP)*. 101–119.
- 4277 Ari Rasch, Richard Schulze, and Sergei Gorlatch. 2019a. Generating Portable High-Performance Code via Multi-Dimensional  
 4278 Homomorphisms. In *2019 28th International Conference on Parallel Architectures and Compilation Techniques (PACT)*.  
 4279 354–369. <https://doi.org/10.1109/PACT.2019.00035>
- 4280 Ari Rasch, Richard Schulze, and Sergei Gorlatch. 2020a. `md_poly`: A Performance-Portable Polyhedral Compiler based on  
 4281 Multi-Dimensional Homomorphisms. In *Proceedings of the International Workshop on Polyhedral Compilation Techniques  
 4282 (IMPACT'20)*. 1–4.
- 4283 Ari Rasch, Richard Schulze, and Sergei Gorlatch. 2020b. `md_poly`: A Performance-Portable Polyhedral Compiler based on  
 4284 Multi-Dimensional Homomorphisms. In *ACM SRC Grand Finals Candidates, 2019 – 2020*. 1–5.
- 4285 Ari Rasch, Richard Schulze, Waldemar Gorus, Jan Hiller, Sebastian Bartholomäus, and Sergei Gorlatch. 2019b. High-  
 4286 Performance Probabilistic Record Linkage via Multi-Dimensional Homomorphisms. In *Proceedings of the 34th ACM/SI-  
 4287 GAPP Symposium on Applied Computing (SAC '19)*. Association for Computing Machinery, New York, NY, USA, 526–533.  
 4288 <https://doi.org/10.1145/3297280.3297330>
- 4289 Ari Rasch, Richard Schulze, Denys Shabalin, Anne Elster, Sergei Gorlatch, and Mary Hall. 2023. (De/Re)-Compositions  
 4290 Expressed Systematically via MDH-Based Schedules. In *Proceedings of the 32nd ACM SIGPLAN International Conference  
 4291 on Compiler Construction (CC 2023)*. Association for Computing Machinery, New York, NY, USA, 61–72. <https://doi.org/10.1145/3578360.3580269>
- 4292 Ari Rasch, Richard Schulze, Michel Steuwer, and Sergei Gorlatch. 2021. Efficient Auto-Tuning of Parallel Programs with  
 4293 Interdependent Tuning Parameters via Auto-Tuning Framework (ATF). *ACM Trans. Archit. Code Optim.* 18, 1, Article 1  
 4294 (Jan. 2021), 26 pages. <https://doi.org/10.1145/3427093>
- 4295 Chandan Reddy, Michael Kruse, and Albert Cohen. 2016. Reduction Drawing: Language Constructs and Polyhedral  
 4296 Compilation for Reductions on GPU. In *Proceedings of the 2016 International Conference on Parallel Architectures and  
 4297 Compilation (PACT '16)*. Association for Computing Machinery, New York, NY, USA, 87–97. <https://doi.org/10.1145/2967938.2967950>
- 4298 Toomas Remmelg, Thibaut Lutz, Michel Steuwer, and Christophe Dubach. 2016. Performance Portable GPU Code Generation  
 4299 for Matrix Multiplication. In *Proceedings of the 9th Annual Workshop on General Purpose Processing Using Graphics  
 4300 Processing Unit (GPGPU '16)*. Association for Computing Machinery, New York, NY, USA, 22–31. <https://doi.org/10.1145/2884045.2884046>
- 4301 Bertrand Russell. 2020. *The principles of mathematics*. Routledge.
- 4302 Bruce Sagan. 2001. *The symmetric group: representations, combinatorial algorithms, and symmetric functions*. Vol. 203. Springer  
 4303 Science & Business Media.
- 4304 Caio Salvador Rohwedder, Nathan Henderson, João P. L. De Carvalho, Yufei Chen, and José Nelson Amaral. 2023. To Pack  
 4305 or Not to Pack: A Generalized Packing Analysis and Transformation. In *Proceedings of the 21st ACM/IEEE International  
 4306 Symposium on Code Generation and Optimization (CGO 2023)*. Association for Computing Machinery, New York, NY,  
 4307 USA, 14–27. <https://doi.org/10.1145/3579990.3580024>
- 4308 Richard Schulze. 2023. MDH Python Implementation.
- 4309 Karen Simonyan and Andrew Zisserman. 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition.  
 4310 <https://doi.org/10.48550/ARXIV.1409.1556>
- 4311 Paul Springer and Paolo Bientinesi. 2016. Design of a high-performance GEMM-like Tensor-Tensor Multiplication. *CoRR*  
 4312 (2016). arXiv:quant-ph/1607.00145 <http://arxiv.org/abs/1607.00145>
- 4313 Michel Steuwer, Christian Fensch, Sam Lindley, and Christophe Dubach. 2015. Generating Performance Portable Code  
 4314 Using Rewrite Rules: From High-Level Functional Expressions to High-Performance OpenCL Code. In *Proceedings of  
 4315 the 20th ACM SIGPLAN International Conference on Functional Programming (ICFP 2015)*. Association for Computing  
 4316 Machinery, New York, NY, USA, 205–217. <https://doi.org/10.1145/2784731.2784754>

- 4313 Michel Steuwer, Philipp Kegel, and Sergei Gorlatch. 2011. SkelCL - A Portable Skeleton Library for High-Level GPU  
4314 Programming. In *2011 IEEE International Symposium on Parallel and Distributed Processing Workshops and Phd Forum*.  
4315 1176–1182. <https://doi.org/10.1109/IPDPS.2011.269>
- 4316 Michel Steuwer, Thomas Koehler, Bastian Köpcke, and Federico Pizzuti. 2022. RISE & Shine: Language-Oriented Compiler  
4317 Design. *CoRR* abs/2201.03611 (2022). arXiv:2201.03611 <https://arxiv.org/abs/2201.03611>
- 4318 Michel Steuwer, Toomas Remmelg, and Christophe Dubach. 2017. LIFT: A functional data-parallel IR for high-performance  
4319 GPU code generation. In *2017 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. 74–85.  
4320 <https://doi.org/10.1109/CGO.2017.7863730>
- 4321 StreamHPC. 2016. Comparing Syntax for CUDA, OpenCL and HiP. <https://streamhpc.com/blog/2016-04-05/comparing-syntax-cuda-opencl-hip/>.
- 4322 Yifan Sun, Nicolas Bohm Agostini, Shi Dong, and David R. Kaeli. 2019. Summarizing CPU and GPU Design Trends with  
4323 Product Data. *CoRR* abs/1911.11313 (2019). arXiv:1911.11313 <http://arxiv.org/abs/1911.11313>
- 4324 Joel Svensson, Mary Sheeran, and Koen Claessen. 2011. Obsidian: A Domain Specific Embedded Language for Parallel  
4325 Programming of Graphics Processors. In *Implementation and Application of Functional Languages*, Sven-Bodo Scholz and  
4326 Olaf Chitil (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 156–173.
- 4327 Walid Taha and Tim Sheard. 1997. Multi-Stage Programming with Explicit Annotations. In *Proceedings of the 1997 ACM  
4328 SIGPLAN Symposium on Partial Evaluation and Semantics-Based Program Manipulation (PEPM '97)*. Association for  
4329 Computing Machinery, New York, NY, USA, 203–217. <https://doi.org/10.1145/258993.259019>
- 4330 TensorFlow. 2022a. MobileNet v1 models for Keras. <https://github.com/keras-team/keras/blob/master/keras/applications/mobilenet.py>.
- 4331 TensorFlow. 2022b. ResNet models for Keras. <https://github.com/keras-team/keras/blob/master/keras/applications/resnet.py>.
- 4332 TensorFlow. 2022c. VGG16 model for Keras. <https://github.com/keras-team/keras/blob/master/keras/applications/vgg16.py>.
- 4333 Philippe Tillet and David Cox. 2017. Input-Aware Auto-Tuning of Compute-Bound HPC Kernels. In *Proceedings of the  
4334 International Conference for High Performance Computing, Networking, Storage and Analysis (SC '17)*. Association for  
4335 Computing Machinery, New York, NY, USA, Article 43, 12 pages. <https://doi.org/10.1145/3126908.3126939>
- 4336 TIOBE. 2023. The Software Quality Company. <https://www.tiobe.com/tiobe-index/>.
- 4337 TOPLAS Artifact. 2022. [https://gitlab.com/mdh-project/toplas22\\_artifact](https://gitlab.com/mdh-project/toplas22_artifact).
- 4338 Uday Bondhugula. 2022. Pluto: An automatic polyhedral parallelizer and locality optimizer. <https://github.com/bondhugula/pluto>, commit = 12e075a, date = 31.10.2021.
- 4339 Nicolas Vasilache, Oleksandr Zinenko, Aart J. C. Bik, Mahesh Ravishankar, Thomas Raoux, Alexander Belyaev, Matthias  
4340 Springer, Tobias Gysi, Diego Caballero, Stephan Herhut, Stella Laurenzo, and Albert Cohen. 2022. Composable  
4341 and Modular Code Generation in MLIR: A Structured and Retargetable Approach to Tensor Compiler Construction.  
arXiv:cs.PL/2202.03293
- 4342 Nicolas Vasilache, Oleksandr Zinenko, Theodoros Theodoridis, Priya Goyal, Zachary Devito, William S. Moses, Sven  
4343 Verdoolaege, Andrew Adams, and Albert Cohen. 2019. The Next 700 Accelerated Layers: From Mathematical Expressions  
4344 of Network Computation Graphs to Accelerated GPU Kernels, Automatically. *ACM Trans. Archit. Code Optim.* 16, 4,  
4345 Article 38 (Oct. 2019), 26 pages. <https://doi.org/10.1145/3355606>
- 4346 Sven Verdoolaege, Juan Carlos Juega, Albert Cohen, José Ignacio Gómez, Christian Tenllado, and Francky Catthoor. 2013.  
4347 Polyhedral Parallel Code Generation for CUDA. *ACM Trans. Archit. Code Optim.* 9, 4, Article 54 (Jan. 2013), 23 pages.  
4348 <https://doi.org/10.1145/2400682.2400713>
- 4349 Sven Verdoolaege and Tobias Grosser. 2012. Polyhedral Extraction Tool. In *International Workshop on Polyhedral Compilation  
4350 Techniques (IMPACT'12), Paris, France*, Vol. 141.
- 4351 J. von Neumann. 1925. Eine Axiomatisierung der Mengenlehre. 1925, 154 (1925), 219–240. <https://doi.org/doi:10.1515/crll.1925.154.219>
- 4352 Mohamed Wahib and Naoya Maruyama. 2014. Scalable Kernel Fusion for Memory-Bound GPU Applications. In *SC '14:  
4353 Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 191–202.  
4354 <https://doi.org/10.1109/SC.2014.21>
- 4355 Bram Wasti, José Pablo Cambronero, Benoit Steiner, Hugh Leather, and Aleksandar Zlateski. 2022. LoopStack: a Lightweight  
4356 Tensor Algebra Compiler Stack. <https://doi.org/10.48550/ARXIV.2205.00618>
- 4357 R.C. Whaley and J.J. Dongarra. 1998. Automatically Tuned Linear Algebra Software. In *SC '98: Proceedings of the 1998  
4358 ACM/IEEE Conference on Supercomputing*. 38–38. <https://doi.org/10.1109/SC.1998.10004>
- 4359 Maurice V Wilkes. 2001. The memory gap and the future of high performance memories. *ACM SIGARCH Computer  
4360 Architecture News* 29, 1 (2001), 2–7.
- 4361 M.E. Wolf and M.S. Lam. 1991. A loop transformation theory and an algorithm to maximize parallelism. *IEEE Transactions  
4362 on Parallel and Distributed Systems* 2, 4 (1991), 452–471. <https://doi.org/10.1109/71.97902>

- 4362 Hongwei Xi and Frank Pfenning. 1999. Dependent Types in Practical Programming. In *Proceedings of the 26th ACM*  
4363 *SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '99)*. Association for Computing Machinery,  
4364 New York, NY, USA, 214–227. <https://doi.org/10.1145/292540.292560>
- 4365 Rohan Yadav, Alex Aiken, and Fredrik Kjolstad. 2022. DISTAL: The Distributed Tensor Algebra Compiler. In *Proceedings*  
4366 *of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI 2022)*.  
4367 Association for Computing Machinery, New York, NY, USA, 286–300. <https://doi.org/10.1145/3519939.3523437>
- 4368 Cambridge Yang, Eric Atkinson, and Michael Carbin. 2021. Simplifying Dependent Reductions in the Polyhedral Model.  
4369 *Proc. ACM Program. Lang.* 5, POPL, Article 20 (Jan. 2021), 33 pages. <https://doi.org/10.1145/3434301>
- 4370 Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo,  
4371 Koushik Sen, Joseph E. Gonzalez, and Ion Stoica. 2020a. AnsoR: Generating High-Performance Tensor Programs for Deep  
4372 Learning. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. USENIX Association,  
4373 863–879. <https://www.usenix.org/conference/osdi20/presentation/zheng>
- 4374 Size Zheng, Yun Liang, Shuo Wang, Renze Chen, and Kaiwen Sheng. 2020b. *FlexTensor: An Automatic Schedule Exploration*  
4375 *and Optimization Framework for Tensor Computation on Heterogeneous System*. Association for Computing Machinery,  
4376 New York, NY, USA, 859–873. <https://doi.org/10.1145/3373376.3378508>
- 4377
- 4378
- 4379
- 4380
- 4381
- 4382
- 4383
- 4384
- 4385
- 4386
- 4387
- 4388
- 4389
- 4390
- 4391
- 4392
- 4393
- 4394
- 4395
- 4396
- 4397
- 4398
- 4399
- 4400
- 4401
- 4402
- 4403
- 4404
- 4405
- 4406
- 4407
- 4408
- 4409
- 4410

4411 **APPENDIX**

4412 Our Appendix provides details for the interested reader, which should not be required for under-  
 4413 standing the basic ideas presented in this paper.

4414

4415 **A MATHEMATICAL FOUNDATION**

4416 We rely on a set theoretical foundation, based on ZFC set theory [Ciesielski 1997]. We avoid class  
 4417 theory, such as NBG [von Neumann 1925], by assuming, for example, that our universe of types  
 4418 contains all relevant representatives (int, float, struct, etc), but is not the "class of all types".  
 4419 Thereby, we avoid fundamental issues [Russell 2020] which are not relevant for this work.

4420

4421 **A.1 Family**

4422 **Definition 17** (Family). Let  $I$  and  $A$  be two sets. A *family*  $F$  from  $I$  to  $A$  is any set

4423

$$4424 F := \{ (i, a) \mid i \in I \wedge a \in A \}$$

4425

such that the following two properties are satisfied:

4426

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

4427 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  
 4428 refer to  $F$  also as *ordered family*.

4429

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

4430

We write:

4431

- $F_i$  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  $(\dots (F_{i_1, \dots, i_n})_{i_n \in I_n} \dots)_{i_1 \in I_1}$ .

4432

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

4433

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

4434

For nested families, each index set  $I_k$  may depend on the earlier-defined values  $i_1, \dots, i_{k-1}$  (not  
 4435 explicitly stated above for brevity).

4436

**Definition 18** (Tuple). We identify  $n$ -tuples as families that have index set  $[1, n]_{\mathbb{N}}$ .

4437

**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  
 4438  $F_2 = b$ .

4439

## A.2 Scalar Types

4440

We denote by

4441

$$4442 \text{TYPE} := \{ \text{int}, \text{int8}, \text{int16}, \dots, \text{float}, \text{double}, \dots, \text{struct}, \dots \}$$

4443

4444 our set of *scalar types*, where *int8* and *int16* represent 8-bit/16-bit integer numbers, *float* and  
 4445 *double* are the types of single/double precision floating point numbers (IEEE 754 standards),  
 4446 *structs* contain a fixed set of other scalar types, etc. For simplicity, we interpret integer types  
 4447 (*int*, *int8*, *int16*, ...) uniquely as integers  $\mathbb{Z}$ , floating point number types (*float* and *double*) as  
 4448 rationale numbers  $\mathbb{Q}$ , etc.

4449

4460 For high flexibility, we avoid fixing  $\text{TYPE}$  to a particular set of scalar types, i.e., we assume that  
 4461  $\text{TYPE}$  contains all practice-relevant types. This is legal, because our formalism makes no assumptions  
 4462 on the number and kinds of scalar types.

4463 We consider operations on scalar types (addition, multiplication, etc) to be: 1) *atomic*: we do  
 4464 not aim at parallelizing or otherwise optimizing operations on scalar values in this work; 2) *size  
 4465 preserving*: we assume that all values of a scalar type have the same arbitrary but fixed size.

4466 Note that we can potentially also define, for example, the set of arbitrarily sized matrices  
 4467  $\{T^{m \times n} \mid m, n \in \mathbb{N}, T \in \text{TYPE}\}$  as scalar type in our approach. However, this would prevent any kind  
 4468 of formal reasoning about performance and type correctness of matrix-related operations (e.g.,  
 4469 matrix multiplication), like parallelization (due to our atomic assumption above) or type correctness  
 4470 (e.g., assuring in matrix multiplication that number of columns of the first input matrix coincides  
 4471 with and number of rows of the second matrix: due to our size preservation assumption above, we  
 4472 would not be able to distinguish matrices based on their sizes).

### 4473 A.3 Functions

4474 **Definition 19** (Function). Let  $A \in \text{TYPE}$  and  $B \in \text{TYPE}$  be two scalar types.

4475 A (*total*) *function*  $f$  is a tuple of the form

$$4476 f \in \{ (\underbrace{(A, B)}_{\text{function type}}, \underbrace{G_f}_{\text{function graph}}) \mid G_f \subseteq \{ (a, b) \mid a \in A \wedge b \in B \} \}$$

4481 that satisfies the following two properties:

- 4482 • *left-total*:  $\forall a \in A : \exists b \in B : (a, b) \in G_f$ ;
- 4483 • *right-unique*:  $(a, b) \in G_f \wedge (a, b') \in G_f \Rightarrow b = b'$ .

4484 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  
 4485  $A \rightarrow B$ , and we write  $f : A \rightarrow B$  to state that  $f$  has function type  $A \rightarrow B$ .

4486 We refer to:

- 4488 •  $A$  as the *domain* of  $f$
- 4489 •  $B$  as the *co-domain* (or *range*) of  $f$
- 4490 •  $(A, B)$  as the *type* of  $f$
- 4491 •  $G_f$  as the *graph* of  $f$

4492 If  $f$  does not satisfy the left total property, we say  $f$  is *partial*, and we denote  $f$ 's type as  $f : A \rightarrow_p B$   
 4493 (where  $\rightarrow$  is replaced by  $\rightarrow_p$ ).

4494 We allow functions to have so-called *dependent types* [Xi and Pfenning 1999] for expressive typing.  
 4495 For example, dependent types enable encoding the sizes of families into the type system, which  
 4496 contributes to better error checking. We refer to dependently typed functions as *meta-functions*, as  
 4497 outlined in the following.

4499 **Definition 20** (Meta-Function). We refer to any family of functions

$$4500 (f^{<i>} : A^{<i>} \rightarrow B^{<i>})^{<i \in I>}$$

4502 as *meta-function*. In the context of meta-functions, we refer to index  $i \in I$  also as *meta-parameter*,  
 4503 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  
 4504 meta-parameter  $i \in I$ ), and to  $A^{<i>} \rightarrow B^{<i>}$  for concrete  $i$  as *meta-function*  $f$ 's ordinary function  
 4505 type.

4506 In the following, we often write:

- 4507 •  $f^{<i \in I>} : A^{<i>} \rightarrow B^{<i>}$  instead of  $(f^{<i>} : A^{<i>} \rightarrow B^{<i>})^{<i \in I>}$ ;

- 4509 •  $f^{<i>} : A' \rightarrow B^{<1>} \text{ (or } f^{<i>} : A^{<i>} \rightarrow B') \text{ iff } A^{<i>} = A' \text{ (or } B^{<i>} = B') \text{ for all } i \in I.$

4510  
4511 We use *multi-stage meta-functions* as a concept analogous to *multi staging* [Taha and Sheard  
4512 1997] in programming and similar to *currying* in mathematics.

4513 **Definition 21** (Multi-Stage Meta-Function). A *multi-stage meta-function* is a nested family of  
4514 functions:

4515

$$4516 \overbrace{f^{<i_1 \in I_1^{<1>} | \dots | i_S \in I_S^{<i_1, \dots, i_{s-1}>}}}^{\text{stage 1}} : \underbrace{A^{<i_1, \dots, i_S>} \rightarrow B^{<i_1, \dots, i_S>}}_{\text{function instance}}$$

4517

4520 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  
4521 meta-parameters of the previous stages  $i_1, \dots, i_{s-1}$ . We refer to such meta-functions also as *S-stage*  
4522 *meta-functions*, and we denote their type also as

4524

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

4525 and access to them as

4527

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

4528 where different stages are separated by vertical bars.

4530 We allow partially applying parameters (meta and ordinary) of meta-functions.

4532 **Definition 22** (Partial Meta-Function Application). Let

4533

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

4535 be a meta-function (meta-parameters of meta-types  $I_1, \dots, I_S$  omitted for brevity).

- 4536 • The *partial application* of meta-function  $f$  on stage  $s$  to meta-parameter  $\hat{i}_s$  is the meta-function

4538

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

4540 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:

4542

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

4544 We write for  $f'$ 's type also

4545

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

4547 where  $f'$  is replaced by  $f$  and  $i_s \in I_s$  is replaced by the concrete value  $\hat{i}_s$ .

- 4548 • The *partial application* of meta-function  $f$  to ordinary parameter  $x$  is the meta-function

4550

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

4553 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  $i_1 \in \hat{I}_1, \dots, i_S \in \hat{I}_S$ .  
4554 The function is defined as:

4556

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

We allow *generalizing* meta-parameters. For example, by generalizing meta-parameters for input sizes, we can use functions on arbitrarily sized inputs (a.k.a. *dynamic size* in programming). In our generated code, meta-parameters are available at compile time such that concrete meta-parameter values can be exploited for generating well-performing code (e.g., for static loop boundaries). Consequently, generalization increases expressivity of the generated code, e.g., by being able to process differently sized inputs without requiring re-compilation for unseen input sizes, but usually at the cost of performance.

**Definition 23** (Generalized Meta-Parameters). Let

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

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

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

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

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

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

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

which is defined as:

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

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

We write for  $f$ 's type also

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

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

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

We use *postponed meta-parameters* to change the order of meta-parameters of already defined meta-functions. For example, we use postponed meta-parameters in Definition 8 to compute the values of meta-parameters based on the particular meta-parameter values of later stages.

**Definition 24** (Postponed Meta-Parameters). Let

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

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

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

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

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

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

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

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

4607 and defined as:

$$4608 \quad f^{<i_1 | \dots | i_{s-1} | i_{s+1} | \dots | i_s> <i_s>} (a) = f^{<i_1 | \dots | i_{s-1} | i_s | i_{s+1} | \dots | i_s>} (a)$$

4609

4610 We write for  $f'$ 's type also

$$4611 \quad f^{<i_1 \in I_1^{<\dots>} | \dots | i_{s-1} \in I_{s-1}^{<\dots>} | \rightarrow | i_{s+1} \in I_{s+1}^{<\dots>} | \dots | i_s \in I_s^{<\dots>} > <i_s \in I_s^{<\dots>} >} : A^{<i_1, \dots, i_s, \dots, i_s>} \rightarrow B^{<i_1, \dots, i_s, \dots, i_s>} \\ 4612$$

4613 where  $f'$  is replaced by  $f$  and  $i_s$  by symbol " $\rightarrow$ ". For access to  $f'$ , we write

$$4614 \quad f^{<i_1 | \dots | i_{s-1} | \rightarrow | i_{s+1} | \dots | i_s> <i_s>} (x) \\ 4615$$

4616 When using a binary function for combining a family of elements, we often use the following  
4617 notation.

4618 **Notation 5** (Iterative Function Application). Let  $\otimes : T \times T \rightarrow T$  be an arbitrary associative and  
4619 commutative binary function on scalar type  $T \in \text{TYPE}$ . Let further  $x$  be an arbitrary family that has  
4620 index set  $I := \{i_1, \dots, i_N\}$  and image set  $\{x_i\}_{i \in I} \subseteq T$ .

4621 We write  $\otimes_{i \in I} x_i$  instead of  $x_{i_1} \otimes \dots \otimes x_{i_N}$  (infix notation).

#### 4622 A.4 MatVec Expressed in MDH DSL

4623 Our MatVec example from Figure 6 is expressed in MDH's high-level *Domain-Specific Language (DSL)*,  
4624 used as input by our proof-of-concept MDH compiler, as follows:

```
4628 1 MatVec<T in TYPE | I,K in IN> := 
4629 2     out_view<T>( w:(i,k)->(i) ) o
4630 3     md_hom<I,K>( *, (++,+) ) o
4631 4     inp_view<T,T>( M:(i,k)->(i,k) ,
4632 5             v:(i,k)->(k) )
```

4633 Listing 5. MatVec expressed in MDH's DSL

4634

4635 Our compiler takes an expression as in Listing 5 as input, and it generates auto-tunable code from  
4636 it, according to the methods presented in this paper (particularly in Section E). Our compile times  
4637 currently vary for our examples in Figure 16 from < 1 Minute for low-dimensional examples (e.g.,  
4638 linear algebra routines) to 5 Minutes for high-dimensional examples (such as quantum chemistry  
4639 computations). We aim to significantly reduce our compile times and to describe our code generation  
4640 process in our future work, as code generation is not the focus of this paper.

4641

4642 In addition, we also offer an implementation of our MDH approach and its high-level DSL  
4643 embedded in the *Python* programming language [Schulze 2023], as Python is becoming increasingly  
4644 popular in both academia and industry [TIOBE 2023].

4645

## 4646 B ADDENDUM SECTION 2

### 4647 B.1 Design Decisions: Combine Operators

4648 We list some design decisions for combine operators (Definition 2).

4649

#### 4650 Note 1.

4651

- 4652 • We deliberately restrict index set function  $\Rightarrow_{\text{MDA}}^{\text{MDA}}$  to compute the index set in the particular  
4653 dimension  $d$  only, and not of all  $D$  Dimensions (i.e., the function's output is in MDA-IDX-SETS $^d$   
4654 and not MDA-IDX-SETS $^D$ ), because this enables applying combine operator  $\otimes$  iteratively:

$$4655 \quad (\dots ( (a_1 \otimes^{<(P,Q)>} a_2) \otimes^{<(P \cup Q, R)>} a_3) \otimes^{<(P \cup Q \cup R, \dots)>} \dots$$

4656

4656 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  $d$ .  
 4657 This is because the index set of the output MDA changes only in dimension  $d$ , to the new  
 4658 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  
 4659 input for a new application of  $\otimes$ .

- 4660 • It is a design decision that a combine operator's index set function  $\Rightarrow_{MDA}^{MDA}$  takes as input the  
 4661 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.
- 4662 • For better convenience, we could potentially define the meta-type of combine operators (Definition 2) 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 21). 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  $b \in T[I_1, \Rightarrow_{MDA}^{MDA}(Q)]$ .

4663 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 1), prefix-sum (Example 15), and point-wise combination (Example 2) 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.

## 4682 B.2 Generalized Notion of MDHs

4683 The MDH Definition 3 can be generalized to have an arbitrary algebraic structure as input.

4684 **Definition 25** (Multi-Dimensional Homomorphism). Let

$$4685 \mathcal{A}^\downarrow := (T^{INP}[\Rightarrow_{MDA}^{MDA\downarrow}(*), \dots, \Rightarrow_{MDA}^{MDA\downarrow}(*)], (\downarrow \otimes_d)_{d \in [1, D]_{\mathbb{N}}})$$

4686 and

$$4687 \mathcal{A}^\uparrow := (T^{OUT}[\Rightarrow_{MDA}^{MDA\uparrow}(*), \dots, \Rightarrow_{MDA}^{MDA\uparrow}(*)], (\uparrow \otimes_d)_{d \in [1, D]_{\mathbb{N}}})$$

4688 be two algebraic structures, where

$$4689 (\downarrow \otimes_d \in CO^{<\Rightarrow_{MDA}^{MDA\downarrow} | T^{INP} | D | d >})_{d \in [1, D]_{\mathbb{N}}}$$

4690 and

$$4691 (\uparrow \otimes_d \in CO^{<\Rightarrow_{MDA}^{MDA\uparrow} | T^{OUT} | D | d >})_{d \in [1, D]_{\mathbb{N}}}$$

4692 are tuples of combine operators, for  $D \in \mathbb{N}$ ,  $T^{INP}, T^{OUT} \in \text{TYPE}$ ,  $\Rightarrow_{MDA}^{MDA\downarrow}, \Rightarrow_{MDA}^{MDA\uparrow} : MDA\text{-IDX-SETS} \rightarrow MDA\text{-IDX-SETS}$ , and the two structures' carrier sets

$$4693 T^{INP}[\Rightarrow_{MDA}^{MDA\downarrow}(*), \dots, \Rightarrow_{MDA}^{MDA\downarrow}(*)]$$

4694 and

$$4695 T^{OUT}[\Rightarrow_{MDA}^{MDA\uparrow}(*), \dots, \Rightarrow_{MDA}^{MDA\uparrow}(*)]$$

4705 denote the set of MDAs that are in the function domain of combine operators (the star symbol is  
 4706 used for indicating the function range of index functions).

4707 A *Multi-Dimensional Homomorphism* (MDH) from the algebraic structure  $\mathcal{A}^\downarrow$  to the structure  $\mathcal{A}^\uparrow$   
 4708 is any function

$$4709 \quad h^{<I_1, \dots, I_D \in \text{MDA-IDX-SETS}>} : T^{\text{INP}}[\overset{1}{\Rightarrow}_{\text{MDA}}(I_1), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(I_D)] \rightarrow T^{\text{OUT}}[\overset{1}{\Rightarrow}_{\text{MDA}}(I_1), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(I_D)]$$

4711 that satisfies the *homomorphic property*:

$$4713 \quad h(\alpha_1 \downarrow \otimes_d \alpha_2) = h(\alpha_1) \uparrow \otimes_d h(\alpha_2)$$

4714 The MDH Definition 3 is a special case of our generalized MDH Definition 25 above, for  $\downarrow \otimes_d =$   
 4715  $+^{<T^{\text{INP}} | D | d>} \text{ (Example 1).}$

4716 Higher-order function `md_hom` (originally introduced in Definition 4) is defined for the generalized  
 4717 MDH Definition 25 as follows.

4718 **Definition 26** (Higher-Order Function `md_hom`). The higher-order function `md_hom` is of type

$$4719 \quad \text{md\_hom}^{<T^{\text{INP}}, T^{\text{OUT}} \in \text{TYPE} \mid D \in \mathbb{N} \mid (\overset{d}{\Rightarrow}_{\text{MDA}} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS})_{d \in [1, D]_{\mathbb{N}}},}$$

$$4720 \quad (\overset{d}{\Rightarrow}_{\text{MDA}} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS})_{d \in [1, D]_{\mathbb{N}}} > :$$

$$4721 \quad \underbrace{(\text{CO}^{<\overset{1}{\Rightarrow}_{\text{MDA}} \mid T^{\text{OUT}} | D | 1>} \times \dots \times \text{CO}^{<\overset{D}{\Rightarrow}_{\text{MDA}} \mid T^{\text{OUT}} | D | D>})}_{\downarrow \otimes_1, \dots, \downarrow \otimes_D} \times$$

$$4722 \quad \underbrace{\text{SF}^{<T^{\text{INP}}, T^{\text{OUT}}>}}_f \times$$

$$4723 \quad \underbrace{(\text{CO}^{<\overset{1}{\Rightarrow}_{\text{MDA}} \mid T^{\text{OUT}} | D | 1>} \times \dots \times \text{CO}^{<\overset{D}{\Rightarrow}_{\text{MDA}} \mid T^{\text{OUT}} | D | D>})}_{\uparrow \otimes_1, \dots, \uparrow \otimes_D}$$

$$4724 \quad \rightarrow_p \quad \underbrace{\text{MDH}^{<T^{\text{INP}}, T^{\text{OUT}} \mid D \mid (\overset{d}{\Rightarrow}_{\text{MDA}})_{d \in [1, D]_{\mathbb{N}}}, (\overset{d}{\Rightarrow}_{\text{MDA}})_{d \in [1, D]_{\mathbb{N}}}>}}_{\text{md\_hom}(\downarrow \otimes_1, \dots, \downarrow \otimes_D, f, \uparrow \otimes_1, \dots, \uparrow \otimes_D)}$$

4725 The function takes as input a scalar function  $f$  and two tuples of  $D$ -many combine operators  
 4726  $(\downarrow \otimes_1, \dots, \downarrow \otimes_D)$  and  $(\uparrow \otimes_1, \dots, \uparrow \otimes_D)$ , and it yields a function  $\text{md\_hom}(\downarrow \otimes_1, \dots, \downarrow \otimes_D, f, \uparrow \otimes_1, \dots, \uparrow \otimes_D)$   
 4727 which is defined as:

$$4728 \quad \downarrow \alpha \in T^{\text{INP}}[\overset{1}{\Rightarrow}_{\text{MDA}}(I_1), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(I_D)]$$

4729  $=:$

$$4730 \quad \downarrow \otimes_1 \dots \downarrow \otimes_D \downarrow \alpha^{<i_1, \dots, i_d>} \in T^{\text{INP}}[\overset{1}{\Rightarrow}_{\text{MDA}}(\{i_1\}), \dots, \overset{D}{\Rightarrow}_{\text{MDA}}(\{i_D\})]$$

4731  $\mapsto$

$$4732 \quad +_1 \dots +_D \downarrow \alpha_f^{<i_1, \dots, i_d>} \in T^{\text{INP}}[\{i_1\}, \dots, \{i_D\}]$$

4733  $\overset{\vec{f}}{\mapsto}$

$$\begin{aligned}
 4754 \quad & \uparrow_{i_1 \in I_1} \dots \uparrow_{i_D \in I_D} \uparrow_{\mathfrak{a}} f^{<i_1, \dots, i_d>} \in T^{\text{OUT}}[\{i_1\}, \dots, \{i_D\}] \\
 4755 \quad & \mapsto \\
 4756 \quad & \uparrow_{i_1 \in I_1} \dots \uparrow_{i_D \in I_D} \uparrow_{\mathfrak{a}} f^{<i_1, \dots, i_d>} \in T^{\text{OUT}}[\stackrel{1}{\Rightarrow} \underset{\text{MDA}}{\uparrow}(\{i_1\}), \dots, \stackrel{D}{\Rightarrow} \underset{\text{MDA}}{\uparrow}(\{i_D\})] \\
 4757 \quad & =: \\
 4758 \quad & \uparrow_{\mathfrak{a}} \in T^{\text{OUT}}[\stackrel{1}{\Rightarrow} \underset{\text{MDA}}{\uparrow}(I_1), \dots, \stackrel{D}{\Rightarrow} \underset{\text{MDA}}{\uparrow}(I_D)]
 \end{aligned}$$

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

$$\vec{f}^{< i_1, \dots, i_D \in \mathbb{N} >} : T^{\text{INP}} \lceil \{i_1\}, \dots, \{i_D\} \rceil \rightarrow T^{\text{OUT}} \lceil \{i_1\}, \dots, \{i_D\} \rceil$$

and defined as

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

We refer to the first application of  $\mapsto$  as *de-composition*, to the application of  $\overset{f}{\mapsto}$  as *scalar function application*, and to the second application of  $\mapsto$  as *re-composition*.

For  $\text{md\_hom}((\otimes_1^\downarrow, \dots, \otimes_d^\downarrow), f, (\otimes_1^\uparrow, \dots, \otimes_d^\uparrow))$ , we require per definition the homomorphic property (Definition 25), i.e., for each  $d \in [1, D]_{\mathbb{N}}$ , it must hold:

$$\mathrm{md\_hom}(\mathbf{(\downarrow_{\otimes_1}, \dots, \downarrow_{\otimes_D})}, f, \mathbf{(\uparrow_{\otimes_1}, \dots, \uparrow_{\otimes_D})})(\mathbf{q_1 \downarrow_{\otimes_d} q_2}) \equiv$$

$$\mathrm{md\_hom}((\downarrow_{\otimes_1}, \dots, \downarrow_{\otimes_D}), f, (\uparrow_{\otimes_1}, \dots, \uparrow_{\otimes_D}))(\mathfrak{q}_1)$$

$$\text{md\_hom}((\downarrow_{\oplus} \quad \quad \downarrow_{\oplus-}) \quad f \quad (\uparrow_{\oplus} \quad \quad \uparrow_{\oplus-})) (g)$$

### B.3 Simple MDH Examples

**Function Mapping.** Function map $^{ 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  $I := (I_1, \dots, I_D) \in \text{MDA-IDX-SETS}^D$ .$

The function is of type

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

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

and it is computed as:

$$\mathfrak{a} \quad \text{map}^{<T^{\text{INP}}, T^{\text{OUT}} | D | (I_1, \dots, I_D)}(f) \mapsto \begin{smallmatrix} +1 & \dots & +D \end{smallmatrix} \vec{f}_{\text{map}}(\mathfrak{a}|_{\{i_1\} \times \dots \times \{i_D\}})$$

Here,  $+_d := +^{< T^{\text{OUT}} | D | d >}$  denotes concatenation (Example 1) in dimension  $d$ , MDA  $\mathfrak{a}|_{\{i_1\} \times \dots \times \{i_D\}}$  is the restriction of  $\mathfrak{a}$  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

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

and defined as

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

4803

4804 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  
 4805 combine operators are concatenation  $\text{++}_{d \in \text{CO}^{< id | T^{\text{OUT}} | D | d >}}$  in all dimensions  $d \in [1, D]_{\mathbb{N}}$ .  
 4806

4807 We have chosen map function's order of stages –  $T^{\text{INP}}, T^{\text{OUT}} \in \text{TYPE}$  (stage 1),  $D \in \mathbb{N}$  (stage 2), and  
 4808  $(I_1, \dots, I_D) \in \text{MDA-IDX-SETS}^D$  (stage 3) – according to the recommendations in [Haskell Wiki \[2013\]](#),  
 4809 i.e., earlier stages (such as the scalar types  $T^{\text{INP}}, T^{\text{OUT}}$ ) are expected to change less frequently than  
 4810 later stages (e.g., the MDAs' index sets  $I_1, \dots, I_D$ ).

4811 *Reduction.* Function  $\text{red}^{< T | D | (I_1, \dots, I_D) >} (\oplus)$  combines all elements within an MDA that has scalar  
 4812 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  
 4813 associative and commutative binary function  $\oplus : T \times T \rightarrow T$ .

4814 The function is of type

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

4815 and it is computed as:

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

4816 Here,  $\overrightarrow{\bullet}_d (\oplus) := \overrightarrow{\bullet}^{< T | D | d >} (\oplus)$  denotes point-wise combination (Example 2) in dimension  $d$ , MDA  
 4817  $\text{a}|_{\{i_1\} \times \dots \times \{i_D\}}$  is defined as above, and  $\vec{f}_{\text{red}}$  is the function of type

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

4818 that is defined as

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

4819 It is easy to see that  $\text{red}^{< T | D >} (\oplus)$  is an MDH of type  $\text{MDH}^{< T, T | 0_f, \dots, 0_f >}$  whose combine operators  
 4820 are point-wise addition  $\overrightarrow{\bullet}_d (\oplus) \in \text{CO}^{< id | T | D | d >}$  in all dimensions  $d \in [1, D]_{\mathbb{N}}$ . The same as for  
 4821 function  $\text{map}$ , function  $\text{red}$ 's order of meta-parameter stages are chosen according to [\[Haskell Wiki 2013\]](#).  
 4822

#### 4823 B.4 Design Decisions: `md_hom`

4824 We list some design decisions for higher-order function `md_hom` (Definition 4).

4825 **Note 2.** For some MDHs (such as Mandelbrot), the scalar function  $f$  (Definition 4) is dependent  
 4826 on the position in the input MDA, i.e., it takes as arguments, in addition to  $\text{a}[i_1, \dots, i_D]$ , also the  
 4827 indices  $i_1, \dots, i_D$ . Such MDHs can be easily expressed via `md_hom` after a straightforward type  
 4828 adjustment: type  $\text{SF}^{< T^{\text{INP}}, T^{\text{OUT}} >}$  has to be defined as the set of functions  $f : T^{\text{INP}} \times \text{MDA-IDX-SETS}^D \rightarrow T^{\text{OUT}}$   
 4829 (rather than of functions  $f : T^{\text{INP}} \rightarrow T^{\text{OUT}}$ , as in Definition 4).

4830 Since we do not aim at forcing scalar functions to always take MDA indices as input arguments –  
 4831 for expressing most computations, this is not required (Figure 16) and only causes additional  
 4832 complexity – we assume in the following two different definitions of pattern `md_hom`: one variant  
 4833 exactly as in Definition 4, and one variant with the adjusted type for scalar functions and that  
 4834 passes automatically indices  $i_1, \dots, i_D$  to  $f$ . The two variants can be easily differentiated, via an  
 4835 additional, boolean meta-parameter `USE_MDA_INDICES`: first variant iff `USE_MDA_INDICES = false`  
 4836 and second variant iff `USE_MDA_INDICES = true`.

4852 For simplicity, we focus in this paper on the first variant (as in Definition 4), because, in practice,  
 4853 it is the more common variant, and because all insights presented in this work apply to both  
 4854 variants.

## 4855 B.5 Proof `md_hom` Lemma

4856 We prove Lemma 1.

4857 PROOF. Let  $\alpha_1 \in T[I_1^{\alpha_1}, \dots, I_D^{\alpha_2}]$  and  $\alpha_2 \in T[I_1^{\alpha_2}, \dots, I_D^{\alpha_2}]$  be two arbitrary MDAs for which  $I_d^{\alpha_1} \cap I_d^{\alpha_2} = \emptyset$ ,  
 4858 for  $d \in [1, D]_{\mathbb{N}}$  arbitrary but fixed, i.e., the two MDAs are concatenable in dimension  $d$ .

4859 According to Definition 4, we have to show that

$$4860 \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\alpha_1 \mathbin{+}_d \alpha_2) =$$

$$4861 \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\alpha_1) \mathbin{\oplus}_d \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\alpha_2)$$

4862 For this, we first show for arbitrary  $k \in [1, D]_{\mathbb{N}}$  that

$$4863 \dots \underset{i_k \in I_k}{\oplus}_k \underset{i_{k+1} \in I_{k+1}}{\oplus}_{k+1} \dots x|_{\dots, \{i_k\}, \{i_{k+1}\}, \dots} = \dots \underset{i_{k+1} \in I_{k+1}}{\oplus}_{k+1} \underset{i_k \in I_k}{\oplus}_k \dots x|_{\dots, \{i_k\}, \{i_{k+1}\}, \dots}$$

4864 from which follows

$$4865 \underset{i_1 \in I_1}{\oplus}_1 \dots \underset{i_D \in I_D}{\oplus}_D x|_{\{i_1\}, \dots, \{i_D\}} = \underset{i_\sigma(1) \in I_\sigma(1)}{\oplus}_\sigma(1) \dots \underset{i_\sigma(D) \in I_\sigma(D)}{\oplus}_\sigma(D) x|_{\{i_1\}, \dots, \{i_D\}}$$

4866 for any permutation  $\sigma : \{1, \dots, D\} \leftrightarrow \{1, \dots, D\}$ . Afterwards, in our assumption above, we can  
 4867 assume w.l.o.g. that  $d = 1$ .

4868 Case 1:  $[\oplus_k = \oplus_{k+1}]$  Follows immediately from the commutativity of  $\mathbin{+}$  or  $\vec{\bullet}(\oplus)$  for com-  
 4869 mutative  $\oplus$ , respectively. ✓

4870 Case 2:  $[\oplus_k \neq \oplus_{k+1}]$  Trivial, as it is either  $\oplus_k = \mathbin{+}$  or  $\oplus_{k+1} = \mathbin{+}$ , and

$$4871 (\underset{i_d \in I_d}{\mathbin{+}_d} x|_{\dots, \{i_d\}, \dots})[i_1, \dots, i_D] = (x|_{\dots, \{i_d\}, \dots})[i_1, \dots, i_D]$$

4872 according to the definition of MDA concatenation  $\mathbin{+}$  (Example 1). ✓

4873 Let now be  $d = 1$  (see assumption above), it holds:

$$4874 \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\alpha_1 \mathbin{+}_1 \alpha_2) \\ 4875 = \underset{i_1 \in I_1}{\oplus}_1 \dots \underset{i_D \in I_D}{\oplus}_D \vec{f}((\alpha_1 \mathbin{+}_1 \alpha_2)|_{\{i_1\} \times \dots \times \{i_D\}}) \\ 4876 = \underset{i_1 \in I_1^{\alpha_1}}{\oplus}_1 \dots \underset{i_D \in I_D^{\alpha_2}}{\oplus}_D \vec{f}(\alpha_1|_{\{i_1\} \times \dots \times \{i_D\}}) \underset{i_1 \in I_1^{\alpha_2}}{\oplus}_1 \dots \underset{i_D \in I_D^{\alpha_2}}{\oplus}_D \vec{f}(\alpha_2|_{\{i_1\} \times \dots \times \{i_D\}}) \\ 4877 = \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\alpha_1) \underset{i_1 \in I_1^{\alpha_2}}{\oplus}_1 \text{md\_hom}(f, (\oplus_1, \dots, \oplus_D))(\alpha_2) \quad \checkmark$$

4878 □

## 4879 B.6 Examples of Index Functions

4880 We present examples of index functions (Definition 6).

4881 **Example 13** (Matrix-Vector Multiplication). The index functions we use for expressing Matrix-  
 4882 Vector Multiplication (MatVec) are:

4883

- Input Matrix:

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

for  $\Rightarrow_{\text{BUF}}^{\text{MDA}}(I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_1^{\text{MDA}})]_{N_0}, [0, \max(I_2^{\text{MDA}})]_{N_0}$

- Input Vector:

`iidx(i, k) := (k) ∈ MDA-IDX-to-BUF-IDX<D=2, D1=1 | =>MDABUF`

for  $\Rightarrow_{\text{BUF}}^{\text{MDA}}(I_1^{\text{MDA}}, I_2^{\text{MDA}}) := [0, \max(I_2^{\text{MDA}})]_{N_0}$

- Output Vector:

$\text{idx}(i, k) := (i) \in \text{MDA-TRX-to-BUF-TRX}^{< D=2, D_1=1 } \Rightarrow \text{MDA-BUF}$

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

**Example 14** (Jacobi 1D). The index functions we use for expressing Jacobi 1D (Jacobi 1D) are:

- Input Buffer, 1. Access:

$\text{idx}(i) := (i + 0) \in \text{MDA-TDX-to-BUF}^{< D=1, D_1=1 | \Rightarrow \text{MDA-BUF}^{>}}$

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

- Input Buffer, 2. Access:

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

for  $\Rightarrow_{\text{PHE}}^{\text{MDA}}(J_1^{\text{MDA}}) \equiv [0, \max(J_1^{\text{MDA}}) + 1]_{\mathbb{N}_0}$

- Input Buffer 3 Access:

$\Rightarrow_{\text{MDA}} (i+2) \in \text{MDA}_{\text{TDY}} \Rightarrow_{\text{PUE}} \text{TDY} \leq D=1, D_1=1 \Rightarrow_{\text{PUE}} \text{MDA}$

for  $\rightarrow^{\text{MDA}}(I^{\text{MDA}}) \leftarrow [0, \max(I^{\text{MDA}}) + 2]$

- Output Buffer:

$\text{is}_{\text{RUE}}(i) : (i) \in \text{MDA} \wedge \text{ID}_{\text{RUE}} \neq \text{RUE} \wedge \text{ID}_{\text{RUE}} < D=1, D_1=1 \Rightarrow \text{RUE} \geq \text{MDA}$

for  $\rightarrow^{\text{MDA}}(I^{\text{MDA}}) \leftarrow [0, \max(I^{\text{MDA}})]$

## B.7 Representation of Scalar Values

Scalar values can be considered as 0-dimensional BUFs (Definition 5). Consequently, in Definition 5, 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

4950  $b_{1D}[0] := s$ , or a 2-dimensional BUF  $b_{2D}[0, 0] := s$ , or a 3-dimensional BUF  $b_{3D}[0, 0, 0] := s$ , etc.  
4951 However, this results in an ambiguous representation of scalar values, which we aim to avoid by  
4952 considering scalars as 0-dimensional BUFs, as described above.

## 4954 B.8 Runtime Complexity of Histograms

Our implementation of Histograms (Subfigure 5 in Figure 16) 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  $B$  (which has a step complexity of  $\mathcal{O}(1)$ ) and take into account its reduction dimension  $B$  only, which has a step complexity of  $\log(B)$  (parallel reduction [Harris et al. 2007]). In contrast, related approaches [Henriksen et al. 2020] are often work efficient, by having a work complexity of  $\mathcal{O}(B)$ ; however, their high work efficiency is at the cost of their step complexity which is also  $\mathcal{O}(B)$ , rather than  $\mathcal{O}(\log(B))$  as for our implementation in Subfigure 5, thereby being asymptotically less efficient for parallel machines consisting of many cores. Our future work will show that the work-efficient Histogram implementation introduced in Henriksen et al. [2020] can also be expressed in our approach, by using for scalar function  $f$  an optimized micro kernel for Histogram computation, similarly as done in the related work.

## B.9 Combine Operator of Prefix-Sum Computations

We define *prefix-sum* which is the combine operator of compute pattern `scan` and example MBBS in Section 2.5.

**Example 15** (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

4977  $\oplus_{\text{prefix-sum}} < T \in \text{TYPE} \mid D \in \mathbb{N} \mid d \in [1, D]_{\mathbb{N}} \mid (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} \rangle :$   
 4978  $\underbrace{T \times T \rightarrow T}_{\oplus} \rightarrow \underbrace{T[I_1, \dots, \underbrace{id(P), \dots, I_D}_{d}] \times T[I_1, \dots, \underbrace{id(Q), \dots, I_D}_{d}] \rightarrow T[I_1, \dots, \underbrace{id(P \sqcup Q), \dots, I_D}_{d}]}_{\text{prefix-sum (according to } \oplus\text{)}}$   
 4979  
 4980  
 4981  
 4982  
 4983  
 4984

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  $\max(P) < \max(Q)$  for commutativity):

```

4988
4989      $\otimes_{\text{prefix-sum}}^{<T|D|d|(I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D), (P, Q)>} (\oplus)(\mathfrak{a}_1, \mathfrak{a}_2)[i_1, \dots, i_d, \dots, i_D]$ 
4990
4991     := 
$$\begin{cases} \mathfrak{a}_1[i_1, \dots, i_d, \dots, i_D] & , i_d \in P \\ \mathfrak{a}_1[i_1, \dots, \max(P), \dots, i_D] \oplus \mathfrak{a}_2[i_1, \dots, i_d, \dots, i_D] & , i_d \notin Q \end{cases}$$

4992

```

Function  $\oplus^{<T|D|d>}_{\text{prefix-sum}}(\oplus)$  (meaning:  $\oplus_{\text{prefix-sum}}$  is partially applied to ordinary function parameter  $\oplus$ ; formal details provided in the Appendix, Definition 22) is a combine operator of type  $\mathbb{C}^{<id|T|D|d>}$  for any binary operator  $\oplus: T \times T \rightarrow T$

4999		
5000		
5001		
5002	No.	Constraint
5003		
5004		
5005	0	$\prod_{d \in [1, D]_{\mathbb{N}}} \#PRT(CC, d) \leq 1024$ (Number of CCs limited)
5006		
5007		
5008	R3	$\#PRT(BLK, d) > 1 \wedge \otimes_{\substack{\uparrow-MDH \\ d}} \neq +_{\substack{\uparrow-MDH \\ d}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\text{BLK}, d) \in \{\text{DM}\}$ (SMXs combine in DM)
5009		
5010		
5011		
5012		

Table 2. CUDA model constraints on tuning parameters

5017	No.	Constraint	
5018			
5019			
5020	0	$\prod_{d \in [1, D]_{\mathbb{N}}} \#PRT(CC, \overset{\bullet \leftarrow}{d}) \leq 1024$	(Number of CCs limited)
5021			
5022			
5023	R3	$\#PRT(BLK, \overset{\uparrow \leftarrow}{d}) > 1 \wedge \otimes_{\overset{\uparrow \leftarrow}{d}} \neq +_{\overset{\uparrow \leftarrow}{d}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\overset{\uparrow \leftarrow}{BLK}, \overset{\uparrow \leftarrow}{d}) \in \{\text{DM}\}$	(SMXs combine in DM)
5024			
5025			
5026		$\#PRT(WRP, \overset{\uparrow \leftarrow}{d}) > 1 \wedge \otimes_{\overset{\uparrow \leftarrow}{d}} \neq +_{\overset{\uparrow \leftarrow}{d}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\overset{\uparrow \leftarrow}{WRP}, \overset{\uparrow \leftarrow}{d}) \in \{\text{DM, SM}\}$	(WRPs combine in DM/SM)
5027			
5028			

Table 3. CUDA+WRP model constraints on tuning parameters

5033	No.	Constraint	
5034			
5035			
5036	0	$\prod_{d \in [1, D]_N} \#PRT(\overset{\bullet}{WI}, \overset{\leftarrow}{d}) \leq C_{DEV}$	(Number of PEs limited)
5037			
5038			
5039	R3	$\#PRT(\overset{\uparrow}{WG}, \overset{\leftarrow}{d}) > 1 \wedge \otimes_{\overset{\leftarrow}{d}} \neq +_{\overset{\leftarrow}{d}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\overset{\uparrow}{WG}, \overset{\leftarrow}{d}) \in \{\text{GM}\}$	(CUs combine in GM)
5040			
5041			
5042		$\#PRT(\overset{\uparrow}{WI}, \overset{\leftarrow}{d}) > 1 \wedge \otimes_{\overset{\leftarrow}{d}} \neq +_{\overset{\leftarrow}{d}} \Rightarrow \uparrow\text{-mem}^{<\text{ob}>}(\overset{\uparrow}{WI}, \overset{\leftarrow}{d}) \in \{\text{GM, LM}\}$	(PEs combine in GM/LM)
5043			
5044			

Table 4. OpenCL model constraints on tuning parameters

5048 **C ADDENDUM SECTION 3**

5049 **C.1 Constraints of Programming Models**

5050 Constraints of programming models can be expressed in our formalism; we demonstrate this  
 5051 using the example models CUDA and OpenCL. For this, we add to the general, model-unspecific  
 5052 constraints (described in Section 3.4) the new, model-specific constraints listed in Tables 2 or 3 for  
 5053 CUDA or the constraints in Table 4 for OpenCL, respectively. For brevity, we use in the following:  
 5054

$$5055 \quad (l_{\text{MDH}}, d_{\text{MDH}}) := \leftrightarrow_{\bullet\text{-ass}}^{-1} (l_{\text{ASM}}, d_{\text{ASM}}), \bullet \in \{\downarrow, f, \uparrow\}$$

5056 In Tables 2 and 3 for CUDA, the constraint No. 1 (which constrains tuning parameter No. 1 in  
 5057 Table 1) limits the number of cuda cores (CC) to 1024, according to the CUDA specification [NVIDIA  
 5058 2022g]. The constraints on tuning parameter 14 specify that the results of SMX can be combined in  
 5059 device memory (DM) only, and that of CCs/WRPs in only device memory (DM) or shared memory (SM).  
 5060 Note that in the case of Table 3, CCs are not constrained in parameter 14, as CCs within a WRP have  
 5061 access to all CUDA memory regions: DM, SM, as well as RM (via warp shuffles [NVIDIA 2018]).  
 5062

5063 In Table 4 for OpenCL, the constraints are similar to the CUDA's constraints in Tables 2: they limit  
 5064 the number of PEs to  $C_{\text{DEV}}$  (which is a device-specific constant in OpenCL), and they specify the valid  
 5065 memory regions for combining the results of cores, according to the OpenCL specification [Khronos  
 5066 2022b].

5067 Note that the tables present some important example constraints only and are not complete: for  
 5068 example, CUDA and OpenCL devices are also constrained regarding their memory sizes (shared/private  
 5069 memory), which is not considered in the tables.

5070 **C.2 Inverse Concatenation**

5071 **Definition 27** (Inverse Concatenation). The inverse of operator *concatenation* (Example 1) is  
 5072 function  $\text{++}^{-1}$  which is of type

$$5073 \quad \text{++}^{-1} < T \in \text{TYPE} \mid D \in \mathbb{N} \mid d \in [1, D]_{\mathbb{N}} \mid (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D) \in \text{MDA-IDX-SETS}^{D-1}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} > :$$

$$5074 \quad T[I_1, \dots, \underbrace{id(P \cup Q)}, \dots, I_D] \rightarrow T[I_1, \dots, \underbrace{id(P)}, \dots, I_D] \times T[I_1, \dots, \underbrace{id(Q)}, \dots, I_D]$$

5075  $\uparrow$

5076  $\uparrow$

5077  $\uparrow$

5078 where  $\text{id} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$  is the identity function on MDA index sets. The  
 5079 function is computed as:

$$5080 \quad \text{++}^{-1} < T \mid D \mid d \mid (I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D), (P, Q) > (\alpha) := (\alpha_1, \alpha_2)$$

5081 for

$$5082 \quad \alpha_1[i_1, \dots, i_d, \dots, i_D] := \alpha[i_1, \dots, i_d, \dots, i_D], i_d \in P$$

5083 and

$$5084 \quad \alpha_2[i_1, \dots, i_d, \dots, i_D] := \alpha[i_1, \dots, i_d, \dots, i_D], i_d \in Q$$

5085 i.e.,  $\alpha_1$  and  $\alpha_2$  behave exactly as MDA  $\alpha$  on their restricted index sets  $P$  or  $Q$ , respectively.

5086 We often write for  $(\alpha_1, \alpha_2) := \text{++}^{-1} < \dots > (\alpha)$  (meta-parameters omitted via ellipsis) also  $\alpha =: 5087 \alpha_1 \text{++} < \dots > \alpha_2$ . Our notation is justified by the fact that the inverse of MDA  $\alpha$  is uniquely determined  
 5088 as the two MDAs  $\alpha_1$  and  $\alpha_2$  which are equal to  $\alpha$  when concatenating them (follows immediately  
 5089 from the definition of inverse functions).

5090

5097 **C.3 Example 17 in Verbose Math Notation**

5098 Figures 36-38 show our low-level representation from Example 17 in verbose math notation. The  
 5099 symbols  $\blacksquare_1, \dots, \blacksquare_f$  used in the figures are a textual abbreviation for:  
 5100

5101	$\blacksquare_1$	$::=$	$*, *$	$ $	$*, *$	$ $	$*, *$
5102	$\blacksquare_1^1$	$::=$	$*, *$	$ $	$*, *$	$ $	$*, *$
5103	$\blacksquare_2^1$	$::=$	$p_1^1, *$	$ $	$*, *$	$ $	$*, *$
5104	$\blacksquare_1^2$	$::=$	$p_1^1, p_2^1$	$ $	$*, *$	$ $	$*, *$
5105	$\blacksquare_2^2$	$::=$	$p_1^1, p_2^1$	$ $	$p_1^2, *$	$ $	$*, *$
5106	$\blacksquare_1^3$	$::=$	$p_1^1, p_2^1$	$ $	$p_1^2, p_2^2$	$ $	$*, *$
5107	$\blacksquare_2^3$	$::=$	$p_1^1, p_2^1$	$ $	$p_1^2, p_2^2$	$ $	$p_1^3, *$
5108	$\blacksquare_f$	$::=$	$p_1^1, p_2^1$	$ $	$p_1^2, p_2^2$	$ $	$p_1^3, p_2^3$

5112 where symbol  $*$  indicates generalization in meta-parameters (Definition 23).

5113 In Example 17, the arrow annotation of combine operators is formally an abbreviation. For  
 5114 example, operator  $\text{++}_2^{(\text{COR}, y)}$  in Figure 17 is annotated with  $\rightarrow \text{ M: HM}[1, 2], v: \text{HM}[1]$  which  
 5115 abbreviates

$$\dots \downarrow a_2^{2 < p_1^1, p_2^1 | p_1^2, p_2^2 := * | p_1^3 := *, p_2^3 := * >} =: \text{++}_2^{(\text{COR}, y)} \dots$$

$p_2^2 \in [0, 16]_{\mathbb{N}_0}$

5116 Here,  $\downarrow a_2^2$  represents the low-level MDA (Definition 12) that is already partitioned for layer 1 in  
 5117 dimensions 1 and 2, and for layer 2 in dimension 1 (because in Figure 17, operators  $\text{++}_1^{(\text{HM}, x)}, \text{++}_2^{(\text{HM}, y)}$ ,  
 5118  $\text{++}_1^{(\text{COR}, x)}$  appear before operator  $\text{++}_2^{(\text{COR}, y)}$ ), but not yet for layer 2 in dimension 2 as well as for layer 3  
 5119 in both dimensions (indicated by symbol  $*$  which is described formally in Definition 23 of our  
 5120 Appendix). In our generated code (discussed in Section E of our Appendix), we store low-level MDAs,  
 5121 like  $\downarrow a_2^2$ , using their domain-specific data representation, as the domain-specific representation is  
 5122 usually more efficient: in the case of MatVec, we physically store matrix  $M$  and vector  $v$  for the  
 5123 input MDA, and vector  $w$  for the output MDA. For example, low-level MDA

$$\downarrow a_2^{2 < p_1^1, p_2^1 | p_1^2, p_2^2 := * | p_1^3 := *, p_2^3 := * >}$$

5124 can be transformed via view functions (Definitions 8 and 10) to *low-level BUFS* (Definition 13)

$$M_2^{2 < \text{HM} | id > < p_1^1, p_2^1 | p_1^2, p_2^2 := * | p_1^3 := *, p_2^3 := * >}$$

5125 and

$$v_2^{2 < \text{HM} | id > < p_1^1, p_2^1 | p_1^2, p_2^2 := * | p_1^3 := *, p_2^3 := * >}$$

5126 and back (Lemma 2). Similarly as for data structures in low-level programming models (e.g., C  
 5127 arrays as in OpenMP, CUDA, and OpenCL), low-level BUFS are defined to have an explicit notion  
 5128 of memory regions and memory layouts.

5129 In Figure 36, we de-compose the input MDA  $\downarrow a$ , step by step, for the MDH levels  $(1, 1), \dots, (3, 2)$ :

$$\downarrow a =: \downarrow a_1^{< \blacksquare_1 >} \rightarrow \downarrow a_1^{1 < \blacksquare_1^1 >} \rightarrow \downarrow a_2^{1 < \blacksquare_2^1 >} \rightarrow \downarrow a_2^{2 < \blacksquare_1^2 >} \rightarrow \downarrow a_2^{2 < \blacksquare_2^2 >} \rightarrow \downarrow a_1^{3 < \blacksquare_1^3 >} \rightarrow \downarrow a_2^{3 < \blacksquare_2^3 >} \rightarrow \downarrow a_f^{< \blacksquare_f >}$$

5130 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 12).  
 5131 We use as partitioning schema  $\bar{P}$  (Definition 12)

$$P := ((P_1^1, P_2^1), (P_1^2, P_2^2), (P_1^3, P_2^3)) = ((2, 4), (8, 16), (32, 64))$$

5146 and we use the index sets  $I_d$  from Definition 28 (which define a uniform index set partitioning) as  
 5147 follows:

$$5148 \quad \left( \begin{array}{c} \stackrel{d}{\Rightarrow}_{MDA} (I_d^{<P_1^1, P_2^1 | P_1^2, P_2^2 | P_1^3, P_2^3>} )^{<(P_1^1, P_2^1) \in P_1^1 \times P_2^1 | (P_1^2, P_2^2) \in P_1^2 \times P_2^2 | (P_1^3, P_2^3) \in P_1^3 \times P_2^3>} \\ \oplus_d \end{array} \right)_{d \in [1, D]_{\mathbb{N}}}$$

5150 Here,  $\stackrel{d}{\Rightarrow}_{MDA}$  denotes the index set function of combine operator concatenation (Example 1), which is  
 5151 the identity function and explicitly stated for the sake of completeness only. Note that in Figure 36,  
 5152 we access low-level MDAs  $\downarrow \alpha_d^l$  as generalized in some partition sizes, via  $*$  (Definition 23), according  
 5153 to the definitions of the  $\blacksquare_d^l$ .  
 5154

5155 Each MDA  $\downarrow \alpha$  can be transformed to its domain-specific data representation matrix  $\downarrow M$  and  
 5156 vector  $\downarrow v$  and vice versa, using the view functions, as discussed above.

5157 Figure 37 shows our scalar phase, which is formally trivial.

5158 In Figure 38, we re-compose the computed data  $\uparrow \alpha_f^{<\blacksquare_f>}$ , step by step, to the final result  $\uparrow \alpha$ :

$$5159 \quad \uparrow \alpha_f^{<\blacksquare_f>} \rightarrow \uparrow \alpha_1^{<\blacksquare_1^1>} \rightarrow \uparrow \alpha_2^{<\blacksquare_2^2>} \rightarrow \uparrow \alpha_2^{<\blacksquare_2^2>} \rightarrow \uparrow \alpha_1^{<\blacksquare_1^3>} \rightarrow \uparrow \alpha_2^{<\blacksquare_2^3>} \rightarrow \uparrow \alpha_{\perp}^{<\blacksquare_{\perp}>} =: \uparrow \alpha$$

5160 Analogously to the de-composition phase, each output MDA  $(\uparrow \alpha_d^l)_{l \in [1, 3]_{\mathbb{N}}, d \in [1, 2]_{\mathbb{N}}}$  and  $\uparrow \alpha_f$  is a  
 5161 low-level MDA representation, for  $P$  as defined above and index sets  
 5162

$$5163 \quad \left( \begin{array}{c} \stackrel{d}{\Rightarrow}_{MDA} (I_d^{<P_1^1, P_2^1 | P_1^2, P_2^2 | P_1^3, P_2^3>} )^{<(P_1^1, P_2^1) \in P_1^1 \times P_2^1 | (P_1^2, P_2^2) \in P_1^2 \times P_2^2 | (P_1^3, P_2^3) \in P_1^3 \times P_2^3>} \\ \otimes_d \end{array} \right)_{d \in [1, D]_{\mathbb{N}}}$$

5164 where  $\stackrel{d}{\Rightarrow}_{MDA}$  are the index set functions of the combine operators (Definition 2) used in the re-  
 5165 composition phase. The same as in the de-composition phase, we access the output low-level MDAs  
 5166 as generalized in some partition sizes, according to our definitions of the  $\blacksquare_d^l$ , and we identify each  
 5167 MDA with its domain-specific data representation (the output vector  $w$ ).  
 5168  
 5169  
 5170

5171

5172

5173

5174

5175

5176

5177

5178

5179

5180

5181

5182

5183

5184

5185

5186

5187

5188

5189

5190

5191

5192

5193

5194

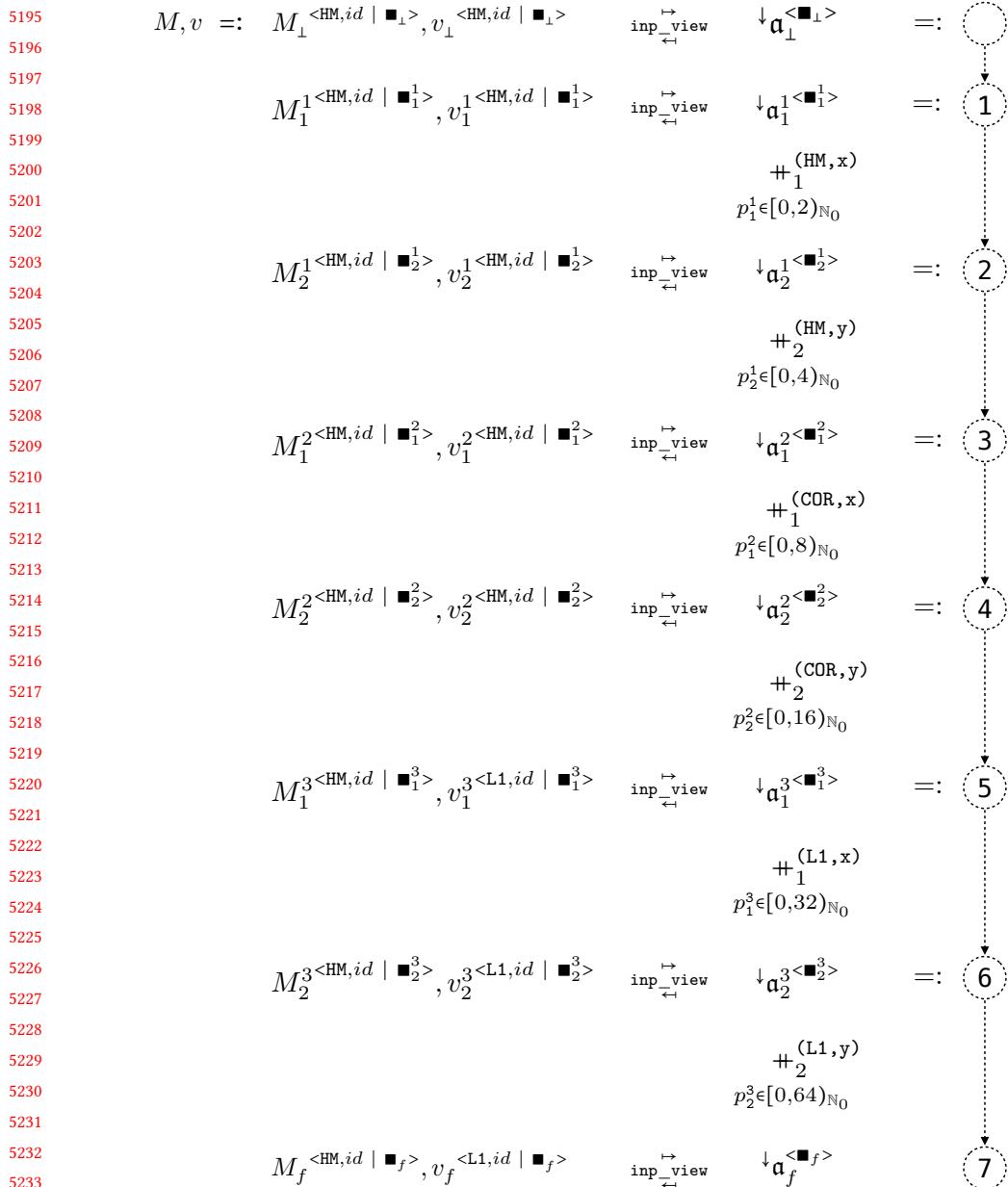


Fig. 36. De-composition phase of Example 17 in verbose math notation.

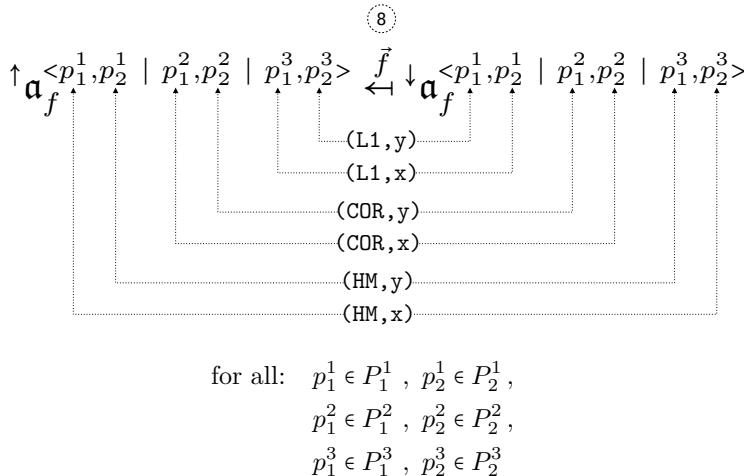
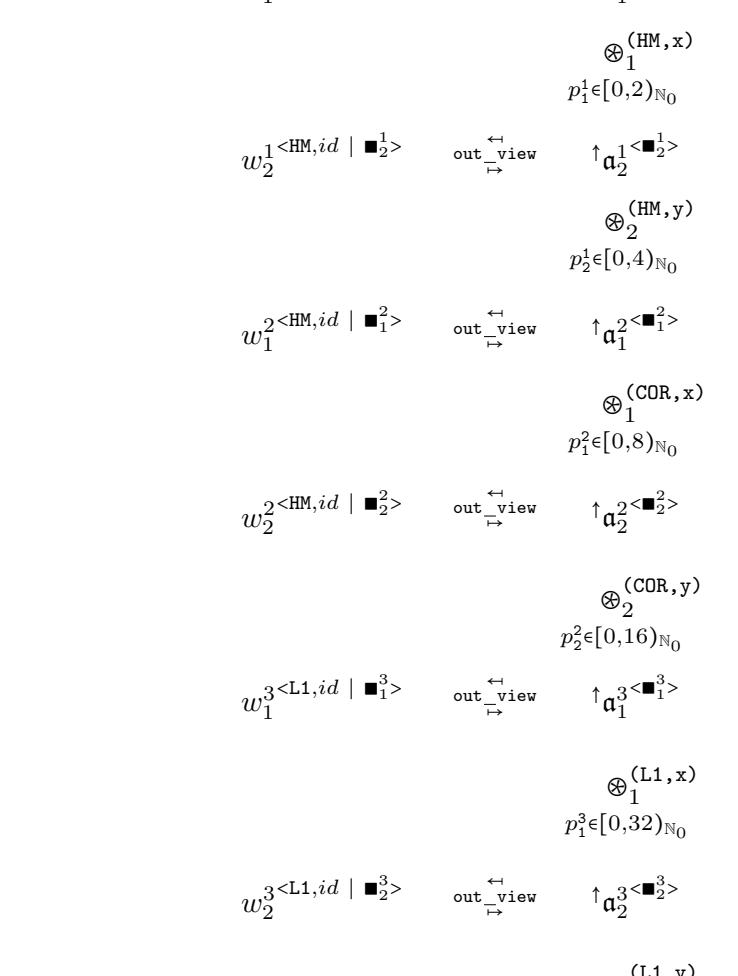


Fig. 37. Scalar phase of Example 17 in verbose math notation.

5293  $w := w_{\perp}^{<\text{HM}, id \mid \blacksquare_{\perp}>} \quad \text{out}_{\overleftarrow{\rightarrow}}^{\text{view}} \quad \uparrow \mathfrak{a}_{\perp}^{<\blacksquare_{\perp}>} :=$  

5294

5295

5296  $w_1^{1 <\text{HM}, id \mid \blacksquare_1^1>} \quad \text{out}_{\overleftarrow{\rightarrow}}^{\text{view}} \quad \uparrow \mathfrak{a}_1^1 <\blacksquare_1^1> :=$  15

5297

5298

5299  $\oplus_1^{(\text{HM}, x)}$

5300  $p_1^1 \in [0, 2)_{\mathbb{N}_0}$

5301  $w_2^1 <\text{HM}, id \mid \blacksquare_2^1> \quad \text{out}_{\overleftarrow{\rightarrow}}^{\text{view}} \quad \uparrow \mathfrak{a}_2^1 <\blacksquare_2^1> :=$  14

5302

5303

5304  $\oplus_2^{(\text{HM}, y)}$

5305  $p_2^1 \in [0, 4)_{\mathbb{N}_0}$

5306

5307  $w_1^2 <\text{HM}, id \mid \blacksquare_1^2> \quad \text{out}_{\overleftarrow{\rightarrow}}^{\text{view}} \quad \uparrow \mathfrak{a}_1^2 <\blacksquare_1^2> :=$  13

5308

5309

5310  $\oplus_1^{(\text{COR}, x)}$

5311  $p_1^2 \in [0, 8)_{\mathbb{N}_0}$

5312

5313  $w_2^2 <\text{HM}, id \mid \blacksquare_2^2> \quad \text{out}_{\overleftarrow{\rightarrow}}^{\text{view}} \quad \uparrow \mathfrak{a}_2^2 <\blacksquare_2^2> :=$  12

5314

5315

5316  $\oplus_2^{(\text{COR}, y)}$

5317  $p_2^2 \in [0, 16)_{\mathbb{N}_0}$

5318  $w_1^3 <\text{L1}, id \mid \blacksquare_1^3> \quad \text{out}_{\overleftarrow{\rightarrow}}^{\text{view}} \quad \uparrow \mathfrak{a}_1^3 <\blacksquare_1^3> :=$  11

5319

5320

5321  $\oplus_1^{(\text{L1}, x)}$

5322  $p_1^3 \in [0, 32)_{\mathbb{N}_0}$

5323

5324

5325  $w_2^3 <\text{L1}, id \mid \blacksquare_2^3> \quad \text{out}_{\overleftarrow{\rightarrow}}^{\text{view}} \quad \uparrow \mathfrak{a}_2^3 <\blacksquare_2^3> :=$  10

5326

5327

5328  $\oplus_2^{(\text{L1}, y)}$

5329  $p_2^3 \in [0, 64)_{\mathbb{N}_0}$

5330

5331  $w_f <\text{L1}, id \mid \blacksquare_f> \quad \text{out}_{\overleftarrow{\rightarrow}}^{\text{view}} \quad \uparrow \mathfrak{a}_f <\blacksquare_f> :=$  9

5332

Fig. 38. Re-composition phase of Example 17 in verbose math notation.

5342 **C.4 Multi-Dimensional ASM Arrangements**

5343 We demonstrate how we arrange memory regions and cores of ASM-represented systems (Section  
 5344 3.2) in multiple dimensions using the example of CUDA.

5345

5346 *Cores (COR)*: In CUDA, SMX cores are programmed via so-called *CUDA Blocks*, and CUDA's  
 5347 CC cores are programmed via *CUDA Threads*. CUDA has native support for arranging its blocks  
 5348 and threads in up to three dimensions which are called x, y, and z in CUDA [NVIDIA 2022f].  
 5349 Consequently, even though the original CUDA specification [NVIDIA 2022g] introduces SMX and  
 5350 CC without having an order, the CUDA programmer benefits from imagining SMX and CC as three-  
 5351 dimensionally arranged.

5352 Additional dimensions can be explicitly programmed in CUDA. For example, to add a fourth  
 5353 dimension to CUDA, we can embed the additional dimension in the CUDA's z dimension, thereby  
 5354 splitting CUDA dimension z in the explicitly programmed dimensions z\_1 (third dimension) and  
 5355 z\_2 (fourth dimension), as follows:

5356 
$$z_1 := z \% Z_1 \text{ and } z_2 := z / Z_1$$

5357

5358 Here,  $Z_1$  represents the number of threads in the additional dimension, and symbol  $\%$  the modulo  
 5359 operator.

5360

5361 *Memory (MEM)*: In CUDA, memory is managed via *C arrays* which may be multi-dimensional: to  
 5362 arrange  $(\text{DIM}_1 \times \dots \times \text{DIM}_D)$ -many memory regions, each of size N, we use a CUDA array of the  
 5363 following type (pseudocode):

5364 
$$\text{array}[\text{DIM}_1] \dots [\text{DIM}_D][\text{N}]$$

5365

5366 Note that CUDA implicitly arranges its *shared* and *private* memory allocations in multiple  
 5367 dimensions, depending on the number of blocks and threads: a shared memory array of type  
 5368 *shared\_array*[ $\text{DIM}_1$ ]...[ $\text{DIM}_D$ ][ $\text{N}$ ] is internally managed in CUDA as *shared\_array*[  
 5369  $\text{blockIdx.x}$ ][ $\text{blockIdx.y}$ ][ $\text{blockIdx.z}$ ][ $\text{DIM}_1$ ]...[ $\text{DIM}_D$ ][ $\text{N}$ ], i.e., each CUDA  
 5370 block has its own shared memory region. Analogously, a private memory array *private\_array*[  
 5371  $\text{DIM}_1$ ]...[ $\text{DIM}_D$ ][ $\text{N}$ ] is managed in CUDA as *private\_array*[ $\text{blockIdx.x}$ ][ $\text{blockIdx.y}$ ]  
 5372 ][ $\text{blockIdx.z}$ ][ $\text{threadIdx.x}$ ][ $\text{threadIdx.y}$ ][ $\text{threadIdx.z}$ ][ $\text{DIM}_1$ ]...[ $\text{DIM}_D$ ][ $\text{N}$ ], correspondingly. Our arrangement methodology continues the CUDA's approach by explicitly  
 5373 programming the additional arrangement dimensions  $\text{DIM}_1, \dots, \text{DIM}_D$ .

5374

5375 Figure 39 illustrates our multi-dimensional core and memory arrangement using the example of  
 5376 CUDA for  $D = 2$  (two-dimensional arrangement).

5377

5378 **C.5 ASM Levels**

5379 ASM levels are pairs  $(l_{\text{ASM}}, d_{\text{ASM}})$  consisting of an ASM layer  $l_{\text{ASM}} \in \mathbb{N}$  and ASM dimension  $d_{\text{ASM}} \in \mathbb{N}$ .

5380 Figure 40 illustrates ASM levels using the example of CUDA's thread hierarchy. The figure shows  
 5381 that thread hierarchies can be considered as a tree in which each level is uniquely determined by a  
 5382 particular combination of a layer (block or thread in the case of CUDA) and dimension (x, y, or  
 5383 z). In the figure, we use *lvl* as an abbreviation for *level*, 1 for *layer*, and d for *dimension*.

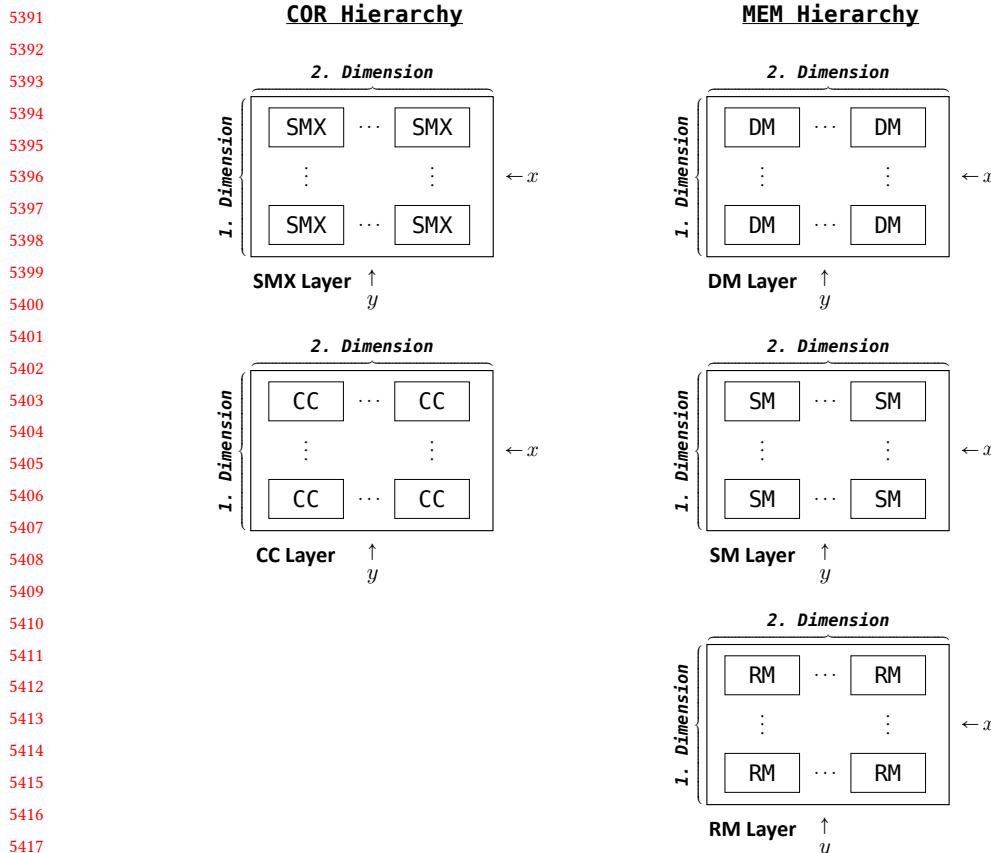
5384 For ASM layers and dimensions, we usually use their domain-specific identifiers, e.g., BLK/CC  
 5385 and x/y/z as aliases for numerical values of layers and dimensions.

5386

5387 **C.6 MDH Levels**

5388 MDH levels are pairs  $(l_{\text{MDH}}, d_{\text{MDH}})$  consisting of an MDH layer  $l_{\text{MDH}} \in \mathbb{N}$  and MDH dimension  $d_{\text{MDH}} \in \mathbb{N}$ .

5389



5419 Fig. 39. Multi-dimensional ASM arrangement illustrated using CUDA for the case  $D = 2$  (two dimensions)  
5420

5421 Figure 41 illustrates MDH levels using as example the de-composition phase in Figure 17. The  
5422 levels  $(l_{MDH}, d_{MDH})$  can be derived from the super- and subscripts of combine operators' variables  
5423  $p_{d_{MDH}}^{l_{MDH}}$ .  
5424

5425  
5426  
5427  
5428  
5429  
5430  
5431  
5432  
5433  
5434  
5435  
5436  
5437  
5438  
5439

5440

5441     d = x {

5442         blockIdx.x = 0

5443         blockIdx.x = 1 ... ← lvl = (BLK,x)

5444     l = BLK { d = y {

5445         blockIdx.y = 0

5446         blockIdx.y = 1 ... ← lvl = (BLK,y)

5447         blockIdx.y = 0

5448         blockIdx.y = 1 ... ← lvl = (BLK,z)

5449         blockIdx.z = 0

5450         blockIdx.z = 1 ... ← lvl = (BLK,z)

5451         blockIdx.z = 0

5452         blockIdx.z = 1 ... ← lvl = (BLK,z)

5453         blockIdx.z = 0

5454         blockIdx.z = 1 ... ← lvl = (BLK,z)

5449     d = x {

5450         threadIdx.x = 0

5451         threadIdx.x = 1 ... ← lvl = (THR,x)

5452     l = THR { d = y {

5453         threadIdx.y = 0

5454         threadIdx.y = 1 ... ← lvl = (THR,x)

5455         threadIdx.y = 0

5456         threadIdx.y = 1 ... ← lvl = (THR,x)

5457         threadIdx.y = 0

5458         threadIdx.y = 1 ... ← lvl = (THR,x)

5459         threadIdx.z = 0

5460         threadIdx.z = 1 ... ← lvl = (THR,x)

5461         threadIdx.z = 0

5462         threadIdx.z = 1 ... ← lvl = (THR,x)

5463         threadIdx.z = 0

5464         threadIdx.z = 1 ... ← lvl = (THR,x)

Fig. 40. ASM levels illustrated using CUDA's thread hierarchy

$\text{#1}^{(\text{HM}, 1)}$	$\text{#2}^{(\text{HM}, 2)}$	$\text{#1}^{(\text{COR}, 1)}$	$\text{#2}^{(\text{COR}, 2)}$	$\text{#1}^{(\text{L1}, 1)}$	$\text{#2}^{(\text{L1}, 2)}$	
$p_1^1 \in [0,2)_{\mathbb{N}_0}$	$p_2^1 \in [0,4)_{\mathbb{N}_0}$	$p_1^2 \in [0,8)_{\mathbb{N}_0}$	$p_2^2 \in [0,16)_{\mathbb{N}_0}$	$p_1^3 \in [0,32)_{\mathbb{N}_0}$	$p_2^3 \in [0,64)_{\mathbb{N}_0}$	
$\rightarrow \text{M: HM}[1, 2]$	$\rightarrow \text{M: HM}[1, 2]$	$\rightarrow \text{M: HM}[1, 2]$				
$\underbrace{\text{v: HM}[1]}$	$\underbrace{\text{v: HM}[1]}$	$\underbrace{\text{v: HM}[1]}$	$\underbrace{\text{v: HM}[1]}$	$\underbrace{\text{v: L1}[1]}$	$\underbrace{\text{v: L1}[1]}$	
<b>lvl:</b>	<b>(1,1)</b>	<b>(1,2)</b>	<b>(2,1)</b>	<b>(2,2)</b>	<b>(3,1)</b>	<b>(3,2)</b>

Fig. 41. MDH levels illustrated using as example the de-composition phase in Figure 17

## C.7 MDA Partitioning

We demonstrate how we partition MDAs into equally sized parts (a.k.a. *uniform partitioning*).

**Definition 28** (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 12) 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

5489 the Part's Size). The partitioned MDA  $\alpha_{\text{prt}}$  is defined as:

$$5490 \quad \alpha =: \underbrace{\begin{array}{c} +_1 \dots +_D \\ p_1^1 \in P_1^1 \quad p_D^1 \in P_D^1 \end{array}}_{\text{Layer 1}} \dots \underbrace{\begin{array}{c} +_1 \dots +_D \\ p_1^L \in P_1^L \quad p_D^L \in P_D^L \end{array}}_{\text{Layer } L} \alpha_{\text{prt}}^{< p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L >}$$

5495 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  
5496 original MDA  $\alpha$ .

## 5497 C.8 TVM Schedule for MatMul

5499 Listing 6 shows TVM's Ansor-generated schedule program for MatMul on input matrices of sizes  
5500  $16 \times 2048$  and  $2048 \times 1000$  taken from ResNet-50's training phase, discussed in Section 3.5. Code  
5501 formatting, like names of variables and comments, have been shortened and adapted in the listing  
5502 for brevity.

5503  
5504  
5505  
5506  
5507  
5508  
5509  
5510  
5511  
5512  
5513  
5514  
5515  
5516  
5517  
5518  
5519  
5520  
5521  
5522  
5523  
5524  
5525  
5526  
5527  
5528  
5529  
5530  
5531  
5532  
5533  
5534  
5535  
5536  
5537

```

5538 1 # exploiting fast memory resources for computed results
5539 2 matmul_local, = s.cache_write([matmul], "local")
5540 3 matmul_1, matmul_2, matmul_3 = tuple(matmul_local.op.axis) + tuple(matmul_local
5541 4 .op.reduce_axis)
5542 5 SHR_1, REG_1 = s[matmul_local].split(matmul_1, factor=1)
5543 6 THR_1, SHR_1 = s[matmul_local].split(SHR_1, factor=1)
5544 7 DEV_1, THR_1 = s[matmul_local].split(THR_1, factor=4)
5545 8 BLK_1, DEV_1 = s[matmul_local].split(DEV_1, factor=2)
5546 9 SHR_2, REG_2 = s[matmul_local].split(matmul_2, factor=1)
5547 10 THR_2, SHR_2 = s[matmul_local].split(SHR_2, factor=1)
5548 11 DEV_2, THR_2 = s[matmul_local].split(THR_2, factor=20)
5549 12 BLK_2, DEV_2 = s[matmul_local].split(DEV_2, factor=1)
5550 13 SHR_3, REG_3 = s[matmul_local].split(matmul_3, factor=2)
5551 14 DEV_3, SHR_3 = s[matmul_local].split(SHR_3, factor=128)
5552 15 s[matmul_local].reorder(BLK_1, BLK_2, DEV_1, DEV_2, THR_1, THR_2, DEV_3, SHR_3,
5553 16 SHR_1, SHR_2, REG_3, REG_1, REG_2)
5554 17 # low-level optimizations:
5555 18 s[matmul_local].pragma(BLK_1, "auto_unroll_max_step", 512)
5556 19 s[matmul_local].pragma(BLK_1, "unroll_explicit", True)
5557 20 # tiling
5558 21 matmul_1, matmul_2, matmul_3 = tuple(matmul.op.axis) + tuple(matmul.op.
5559 22 .reduce_axis)
5560 23 THR_1, SHR_REG_1 = s[matmul].split(matmul_1, factor=1)
5561 24 DEV_1, THR_1 = s[matmul].split(THR_1, factor=4)
5562 25 BLK_1, DEV_1 = s[matmul].split(DEV_1, factor=2)
5563 26 THR_2, SHR_REG_2 = s[matmul].split(matmul_2, factor=1)
5564 27 DEV_2, THR_2 = s[matmul].split(THR_2, factor=20)
5565 28 BLK_2, DEV_2 = s[matmul].split(DEV_2, factor=1)
5566 29 s[matmul].reorder(BLK_1, BLK_2, DEV_1, DEV_2, THR_1, THR_2, SHR_REG_1,
5567 30 SHR_REG_2)
5568 31 s[matmul_local].compute_at(s[matmul], THR_2)
5569 32 # block/thread assignments:
5570 33 BLK_fused = s[matmul].fuse(BLK_1, BLK_2)
5571 34 s[matmul].bind(BLK_fused, te.thread_axis("blockIdx.x"))
5572 35 DEV_fused = s[matmul].fuse(DEV_1, DEV_2)
5573 36 s[matmul].bind(DEV_fused, te.thread_axis("vthread"))
5574 37 THR_fused = s[matmul].fuse(THR_1, THR_2)
5575 38 s[matmul].bind(THR_fused, te.thread_axis("threadIdx.x"))
5576 39 # exploiting fast memory resources for first input matrix:
5577 40 A_shared = s.cache_read(A, "shared", [matmul_local])
5578 41 A_shared_ax0, A_shared_ax1 = tuple(A_shared.op.axis)
5579 42 A_shared_ax0_ax1_fused = s[A_shared].fuse(A_shared_ax0, A_shared_ax1)
5580 43 A_shared_ax0_ax1_fused_o, A_shared_ax0_ax1_fused_i = s[A_shared].split(
5581 44 A_shared_ax0_ax1_fused, factor=1)
5582 45 s[A_shared].vectorize(A_shared_ax0_ax1_fused_i)
5583 46 A_shared_ax0_ax1_fused_o_o, A_shared_ax0_ax1_fused_o_i = s[A_shared].split(
5584 47 A_shared_ax0_ax1_fused_o, factor=80)
5585 48 s[A_shared].bind(A_shared_ax0_ax1_fused_o_i, te.thread_axis("threadIdx.x"))
5586 49 s[A_shared].compute_at(s[matmul_local], DEV_3)
5587 50 # exploiting fast memory resources for second input matrix:
5588 51 # ... (analogous to lines 40 - 47)

```

5589 Listing 6. TVM schedule for Matrix Multiplication on NVIDIA Ampere GPU (variable names shortened for brevity)

5590 Proc. ACM Program. Lang., Vol. 1, No. CONF, Article 1. Publication date: January 2018.

5587 **D ADDENDUM SECTION 5**

5588 **D.1 Data Characteristics used in Deep Neural Networks**

5589 Figure 42 shows the data characteristics used for the deep learning experiments in Figures 28 and 29  
 5590 of Section 5. We use real-world characteristics taken from the neural networks ResNet-50, VGG-16,  
 5591 and MobileNet. For each network, we consider computations MCC and MatMul (Table 16), because  
 5592 these are the networks' most time-intensive building blocks. Each computation is called in each  
 5593 network on different data characteristics – we use for each combination of network and computation  
 5594 the two most time-intensive characteristics. Note that the MobileNet network does not use MatMul  
 5595 in its implementation.

5596 The capsule variants MCC\_Capsule in Figures 28 and 29 of Section 5 have the same characteristics  
 5597 as those listed for MCCs in Figure 42 – the only difference is that MCC\_Capsule, in addition to the  
 5598 dimensions N, H, W, K, R, S, C, uses three additional dimensions MI, MJ, MK, each with a fixed size of 4.  
 5599 This is because MCC\_Capsule operates on  $4 \times 4$  matrices, rather than scalars as MCC does.

Network	Phase	N	H	W	K	R	S	C	Stride H	Stride W	Padding	P	Q	Image Format	Filter Format	Output Format
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			VALID	112	112	NHWC	KRSC	NPQK
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			VALID	224	224	NHWC	KRSC	NPQK
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			VALID	112	112	NHWC	KRSC	NPQK

5601 5602 (a) Data characteristics used for MCC experiments

Network	Phase	M	N	K	Transposition
ResNet-50	Training	16	1000	2048	NN
	Inference	1	1000	2048	NN
VGG-16	Training	16	4096	25088	NN
	Inference	1	4096	25088	NN

5603 5604 (b) Data characteristics used for MatMul experiments

5605 5606 Fig. 42. Data characteristics used for experiments in Section 5

5617 **D.2 Runtime and Accuracy of cuBLASEx**

5618 Listing 7 shows the runtime of cuBLASEx for its different *algorithm* variants. For demonstration, we  
 5619 use the example of matrix multiplication MatMul on NVIDIA Volta GPU for square input matrices  
 5620 of sizes  $1024 \times 1024$ . For each algorithm variant, we list both: 1) the runtime achieved by cuBLASEx  
 5621 (in nanoseconds ns), as well as 2) the maximum absolute deviation ( $\delta_{\max}$  values) compared  
 5622 to a straightforward, sequential CPU computation. For example, the  $\delta_{\max}$  value of algorithm  
 5623 CUBLAS\_GEMM\_DEFAULT is  $3.14713e-05$ , i.e., at least one value  $c_{i,j}^{\text{GPU}}$  in the GPU-computed output  
 5624 matrix deviates by  $3.14713e-05$  from its corresponding, sequentially computed value  $c_{i,j}^{\text{seq}}$  such that  
 5625  $|c_{i,j}^{\text{GPU}}| = |c_{i,j}^{\text{seq}}| + 3.14713e-05$  (bar symbols  $|\dots|$  denote absolute value). All other GPU-computed  
 5626 values  $c_{i',j'}^{\text{GPU}}$  deviate from their sequentially computed CPU-variant by  $3.14713e-05$  or less.

5627 5628 Note that cuBLASEx offers 42 algorithm variants, but not all of them are supported for all potential  
 5629 characteristics of the input and output data (size, memory layout, ...). For our MatMul example,  
 5630 the list of unsupported variants includes: CUBLAS\_GEMM\_ALG01, CUBLAS\_GEMM\_ALG012, etc.

```

5636
5637 CUBLAS_GEMM_DEFAULT: 188416ns (delta_max: 3.14713e-05)
5638 CUBLAS_GEMM_ALG02: 190464ns (delta_max: 6.86646e-05)
5639 CUBLAS_GEMM_ALG03: 186368ns (delta_max: 6.86646e-05)
5640 CUBLAS_GEMM_ALG04: 185344ns (delta_max: 6.86646e-05)
5641 CUBLAS_GEMM_ALG05: 181248ns (delta_max: 6.86646e-05)
5642 CUBLAS_GEMM_ALG06: 181248ns (delta_max: 6.86646e-05)
5643 CUBLAS_GEMM_ALG07: 178176ns (delta_max: 4.1008e-05)
5644 CUBLAS_GEMM_ALG08: 189440ns (delta_max: 4.1008e-05)
5645 CUBLAS_GEMM_ALG09: 171008ns (delta_max: 4.1008e-05)
5646 CUBLAS_GEMM_ALG010: 188416ns (delta_max: 4.1008e-05)
5647 CUBLAS_GEMM_ALG011: 191488ns (delta_max: 4.1008e-05)
5648 CUBLAS_GEMM_ALG018: 185344ns (delta_max: 2.67029e-05)
5649 CUBLAS_GEMM_ALG019: 172032ns (delta_max: 2.67029e-05)
5650 CUBLAS_GEMM_ALG020: 192512ns (delta_max: 2.67029e-05)
5651 CUBLAS_GEMM_ALG021: 201728ns (delta_max: 1.90735e-05)
5652 CUBLAS_GEMM_ALG022: 177152ns (delta_max: 1.90735e-05)
5653 CUBLAS_GEMM_ALG023: 194560ns (delta_max: 1.90735e-05)
5654 CUBLAS_GEMM_DEFAULT_TENSOR_OP: 184320ns (delta_max: 3.14713e-05)
5655 CUBLAS_GEMM_ALG00_TENSOR_OP: 62464ns (delta_max: 0.0131454)
5656 CUBLAS_GEMM_ALG01_TENSOR_OP: 52224ns (delta_max: 0.0131454)
5657 CUBLAS_GEMM_ALG02_TENSOR_OP: 190464ns (delta_max: 3.14713e-05)
5658 CUBLAS_GEMM_ALG03_TENSOR_OP: 189440ns (delta_max: 3.14713e-05)
5659 CUBLAS_GEMM_ALG04_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5660 CUBLAS_GEMM_ALG05_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5661 CUBLAS_GEMM_ALG06_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5662 CUBLAS_GEMM_ALG07_TENSOR_OP: 189440ns (delta_max: 3.14713e-05)
5663 CUBLAS_GEMM_ALG08_TENSOR_OP: 183296ns (delta_max: 3.14713e-05)
5664 CUBLAS_GEMM_ALG09_TENSOR_OP: 189440ns (delta_max: 3.14713e-05)
5665

```

Listing 7. Runtime of cuBLASEx for its different *algorithm* variants on NVIDIA Volta GPU when computing MatMul on square  $1024 \times 1024$  input matrices

```

5666
5667
5668
5669
5670
5671
5672
5673
5674
5675
5676
5677
5678
5679
5680
5681
5682
5683
5684

```

5685 **E CODE GENERATION**

5686 This section outlines how imperative-style pseudocode is generated from our low-level program  
 5687 representation in Section 3. Code optimizations that are below the abstraction level of our low-  
 5688 level representation (like loop fusion and loop unrolling) are considered in this work as *code-level*  
 5689 *optimizations* and briefly outlined in Section E.3. We aim to discuss and illustrate our code generation  
 5690 approach in detail in future work.

5691 In the following, we highlight tuning parameters gray in our pseudocode which are substituted  
 5692 by concrete, optimized values in our executable program code. Static parameters, such as scalar  
 5693 types and the number of input/output buffers, are denoted in math font and also substituted by  
 5694 concrete values in our executable code. We list meta-parameters in angle brackets  $\langle \dots \rangle$ , and  
 5695 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  
 5696 denoting in our pseudocode index function  $\text{idx}_{1,1}^{\text{OUT}}$  used in Figure 15.  
 5697

5698 **Overall Structure**

5699 Listing 8 shows the overall structure of our generated code. We implement a particular expression  
 5700 in our low-level representation (Figure 19) as a compute kernel that is structured in the following  
 5701 phases: 0) preparation (Section E.0), 1) de-composition phase (Section E.1), 2) scalar phase  
 5702 (Section E.2), 3) re-composition phase (Section E.3).

```
5704 1 kernel mdh(
5705 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  $\langle\langle \perp \rangle\rangle \langle\langle B^{\text{IB}} \rangle\rangle \langle\langle \star, \dots, \star \rangle\rangle,$ 
5706 3    $T_1^{\text{OB}}$  trans_ll_OB  $\langle\langle \perp \rangle\rangle \langle\langle 1 \rangle\rangle \langle\langle \star, \dots, \star \rangle\rangle, \dots, T_B^{\text{OB}}$  trans_ll_OB  $\langle\langle \perp \rangle\rangle \langle\langle B^{\text{OB}} \rangle\rangle \langle\langle \star, \dots, \star \rangle\rangle$ 
5707 4 )
5708 5
5709 6   // 0. preparation
5710 7   ...
5711 8   // 2. de-composition phase
5712 9   ...
5713 10  // 3. scalar phase
5714 11  ...
5715 12  // 4. re-composition phase
5716 13  ...
5717 14 }
```

5716 Listing 8. Overall structure of our generated code  
 5717  
 5718

5719 **E.0 Preparation**

5720 Listing 9 shows the preparation phase. It prepares in five sub-phases the basic building blocks used  
 5721 in our low-level representation: 1) md\_hom (Section E.0.1), 2) inp\_view (Section E.0.2), 3) out\_view  
 5722 (Section E.0.3), 4) BUFS (Section E.0.4), 5) MDAs (Section E.0.5).

```
5723 1 // 0. preparation
5724 2   // 0.1. md_hom
5725 3   ...
5726 4   // 0.2. inp_view
5727 5   ...
5728 6   // 0.3. out_view
5729 7   ...
5730 8   // 0.4. BUFS
5731 9   ...
5732 10  // 0.5. MDAs
5733 11  ...
```

## Listing 9. Preparation Phase

## E.0.1 md\_hom.

Listing 10 shows the user-defined scalar function and low-level combine operators (Definition 15) which are both provided by the user via higher-order function `md_hom` (Definition 4).

Listing 11 shows how we pre-implement for the user the two combine operators *concatenation* (Example 1) and *point-wise combination* (Example 2).

Listing 12 shows how we pre-implement the *inverse of concatenation* (Definition C.2), which we will use in the de-composition phase (via Definition 28).

```

5744 1 // 0.1. md_hom
5745 2
5746 3     // 0.1.1. scalar function
5747 4     f(  $T^{INP}$  inp ) ->  $T^{OUT}$  out
5748 5     {
5749 6         // ... (user defined)
5750 7     }
5751 8
5752 9     // 0.1.2. combine operators
5753 10     $\forall d \in [1, D]_{\mathbb{N}}$ :
5754 11         $co<<d>><I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in MDA\text{-IDX-SETS}, (P, Q) \in MDA\text{-IDX-SETS} \times MDA\text{-IDX-SETS}>(<$ 
5755 12             $T^{OUT}[I_1, \dots, I_{d-1}, \xrightarrow{d}^{MDA}(P), I_{d+1}, \dots, I_D] \text{ lhs} ,$ 
5756 13             $T^{OUT}[I_1, \dots, I_{d-1}, \xrightarrow{d}^{MDA}(Q), I_{d+1}, \dots, I_D] \text{ rhs} ) \rightarrow T^{OUT}[I_1, \dots, I_{d-1}, \xrightarrow{d}^{MDA}(P \cup Q), I_{d+1}, \dots, I_D] \text{ res}$ 
5757 14        {
5758 15            // ... (user defined)
5759 16        }
```

## Listing 10. Scalar Function &amp; Combine Operators

```

5760 1 // 0.1.2. combine operators
5761 2
5762 3     // pre-implemented combine operators
5763 4
5764 5     // concatenation
5765 6      $\forall d \in \mathbb{N}$ :
5766 7         $cc<<d>><I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in MDA\text{-IDX-SETS}, (P, Q) \in MDA\text{-IDX-SETS} \times MDA\text{-IDX-SETS}>(<$ 
5767 8             $T^{OUT}[I_1, \dots, I_{d-1}, id(P), I_{d+1}, \dots, I_D] \text{ lhs} ,$ 
5768 9             $T^{OUT}[I_1, \dots, I_{d-1}, id(Q), I_{d+1}, \dots, I_D] \text{ rhs} ) \rightarrow T^{OUT}[I_1, \dots, I_{d-1}, id(P \cup Q), I_{d+1}, \dots, I_D] \text{ res}$ 
5769 10       {
5770 11           int i_1  $\in I_1$ 
5771 12           ..
5772 13           int i_{d-1}  $\in I_{d-1}$ 
5773 14           int i_{d+1}  $\in I_{d+1}$ 
5774 15           ..
5775 16           int i_D  $\in I_D$ 
5776 17           {
5777 18               int i_d  $\in P$ 
5778 19               res[ i_1, \dots, i_d, \dots, i_D ] := lhs[ i_1, \dots, i_d, \dots, i_D ];
5779 20               int i_d  $\in Q$ 
5780 21               res[ i_1, \dots, i_d, \dots, i_D ] := rhs[ i_1, \dots, i_d, \dots, i_D ];
5781 22           }
5782 23       }
5783 24
5784 25     // point-wise combination
```

```

5783 26   $\forall d \in \mathbb{N}:$ 
5784 27   $\text{pw} << d >> < I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in \text{MDA-IDX-SETS}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} > ($ 
5785 28   $\oplus : T^{\text{OUT}} \times T^{\text{OUT}} \rightarrow T^{\text{OUT}}) ( T^{\text{OUT}}[I_1, \dots, I_{d-1}, 0_f(P), I_{d+1}, \dots, I_D] \text{ lhs} ,$ 
5786 29   $T^{\text{OUT}}[I_1, \dots, I_{d-1}, 0_f(Q), I_{d+1}, \dots, I_D] \text{ rhs } )$ 
5787 30   $\rightarrow T^{\text{OUT}}[I_1, \dots, I_{d-1}, 0_f(P \cup Q), I_{d+1}, \dots, I_D] \text{ res}$ 
5788 31  {
5789 32  int i_1  $\in I_1$ 
5790 33   $\dots$ 
5791 34  int i_{d-1}  $\in I_{d-1}$ 
5792 35  int i_{d+1}  $\in I_{d+1}$ 
5793 36   $\dots$ 
5794 37  int i_D  $\in I_D$ 
5795 38  {
5796 39   $\text{res}[ i_1, \dots, i_{d-1}, 0, i_{d+1}, \dots, i_D ]$ 
5797 40   $:= \text{lhs}[ i_1, \dots, i_{d-1}, 0, i_{d+1}, \dots, i_D ]$ 
5798 41   $\text{atomic}(\oplus) \text{ rhs}[ i_1, \dots, i_{d-1}, 0, i_{d+1}, \dots, i_D ];$ 
5799 42  }
5800 43  }

```

Listing 11. Pre-Implemented Combine Operators

```

5800
5801 1 // 0.1.2. combine operators
5802 2
5803 3 // pre-implemented combine operators
5804 4
5805 5 // inverse concatenation
5806 6  $\forall d \in \mathbb{N}:$ 
5807 7  $\text{cc\_inv} << d >> < I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in \text{MDA-IDX-SETS}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} > ($ 
5808 8   $T^{\text{INP}}[I_1, \dots, I_{d-1}, id(P \cup Q), I_{d+1}, \dots, I_D] \text{ res } ) \rightarrow ( T^{\text{INP}}[I_1, \dots, I_{d-1}, id(P), I_{d+1}, \dots, I_D] \text{ lhs} ,$ 
5809 9   $T^{\text{INP}}[I_1, \dots, I_{d-1}, id(Q), I_{d+1}, \dots, I_D] \text{ rhs } )$ 
5810 10 {
5811 11 int i_1  $\in I_1$ 
5812 12  $\dots$ 
5813 13 int i_{d-1}  $\in I_{d-1}$ 
5814 14 int i_{d+1}  $\in I_{d+1}$ 
5815 15  $\dots$ 
5816 16 int i_D  $\in I_D$ 
5817 17 {
5818 18   int i_d  $\in P$ 
5819 19    $\text{res}[ i_1, \dots, i_d, \dots, i_D ] =: \text{lhs}[ i_1, \dots, i_d, \dots, i_D ];$ 
5820 20   int i_d  $\in Q$ 
5821 21    $\text{res}[ i_1, \dots, i_d, \dots, i_D ] =: \text{rhs}[ i_1, \dots, i_d, \dots, i_D ];$ 
5822 22   }
5823 23 }

```

Listing 12. Pre-Implemented Combine Operators

5820  
5821  
5822 *E.0.2 inp\_view.*  
5823 Listing 13 shows the user-defined index functions provided by the user via higher-order function  
5824 *inp\_view* (Definition 8).

```

5826 1 // 0.2. inp_view
5827 2
5828 3 // index functions
5829 4  $\forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}, a \in [1, A_b^{\text{IB}}]_{\mathbb{N}}: \forall d \in [1, D_b^{\text{IB}}]_{\mathbb{N}}:$ 
5830 5 static
5831 6  $\text{idx} << \text{INP} >> << b, a >> << d >> ( \text{int} i_{\text{MDA\_1}}, \dots, i_{\text{MDA\_D}} ) \rightarrow \text{int} i_{\text{BUF\_d}}$ 
5832

```

```

5832 7  {
5833 8  // ... (user defined)
5834 9  }

```

Listing 13. Index Functions (input)

5835

5836

5837 *E.0.3 out\_view.*  
 5838 Listing 14 shows the user-defined index functions provided by the user via higher-order function  
 5839 *out\_view* (Definition 10).

```

5840
5841 1 // 0.3. out_view
5842 2
5843 3 // index functions
5844 4  $\forall b \in [1, B^{OB}]_{\mathbb{N}}, a \in [1, A_b^{OB}]_{\mathbb{N}}: \forall d \in [1, D_b^{OB}]_{\mathbb{N}}:$ 
5845 5 static
5846 6  $\text{idx} << \text{OUT} >> << b, a >> << d >> ( \text{int } i_{\text{MDA\_1}}, \dots, i_{\text{MDA\_D}} ) \rightarrow \text{int } i_{\text{BUF\_d}}$ 
5847 7
5848 8 // ... (user defined)
5849 9

```

Listing 14. Index Functions (output)

5849

5850

5851 *E.0.4 BUFs.*

5852 Listing 15 shows our implementation of low-level BUFs (Definition 13). We compute BUFs' sizes  
 5853 using the ranges of their index functions (Definitions 8 and 10). Moreover, we partially evaluate  
 5854 BUFs' meta-parameters *MEM* (memory region) and  $\sigma$  (memory layout) immediately, as the same  
 5855 values are re-used for them during program runtime.

5856 The BUFs in lines 30 and 45 as well as in lines 69 and 84 represent the BUFs' transposed function  
 5857 representation (Definition 13), and the BUFs in lines 23, 37, and 52 as well as in lines 62, 76, and 91  
 5858 are the transposed BUFs' ordinary low-level BUF representation.

```

5859 1 // 0.4. BUFs
5860 2
5861 3 // 0.4.1. compute BUF sizes
5862 4  $\forall IO \in \{IB, OB\}: \forall b \in [1, B^{IO}]_{\mathbb{N}} \forall d \in [1, D_b^{IO}]_{\mathbb{N}}:$ 
5863 5 static  $N << IO >> << b >> << d >> ( \text{mda\_idx\_set } I_1, \dots, I_D ) \rightarrow \text{int } N_{b\_d}$ 
5864 6
5865 7  $N_{b\_d} := 0;$ 
5866 8
5867 9  $i_1 \in I_1$ 
5868 10  $\dots$ 
5869 11  $i_D \in I_D$ 
5870 12  $\{$ 
5871 13  $\forall a \in [1, A_b^{IB}]_{\mathbb{N}}:$ 
5872 14  $N_{b\_d} :=_{\max} 1 + \text{idx} << IO >> << b, a >> << d >> ( i_1, \dots, i_D );$ 
5873 15  $\}$ 
5874 16  $\}$ 
5875 17
5876 18 // 0.4.2. input BUFs
5877 19
5878 20 // initial BUFs
5879 21  $\forall b \in [1, B^{IB}]_{\mathbb{N}}:$ 
5880 22 static  $l1_{IB} << \perp >> << b >> < \nabla_1^{(\perp)} \in \#PRT(1, 1), \dots, \nabla_D^{(\perp)} \in \#PRT(L, D) >> ( \text{int } i_1, \dots,$ 
5881 23  $\text{int } i_{D_b^{IB}} ) \rightarrow T_b^{IB} a$ 
5882 24  $\{$ 

```

5880

```

5881 25     a := trans_ll_IB<<_1>><<b>><<  $\nabla_1^{(1)}, \dots, \nabla_D^{(1)} >[ i_1, \dots, i_{D_b^{\text{IB}}} ];$ 
5882 26   }
5883 27
5884 28   // de-composition BUFS
5885 29    $\forall (l, d) \in \text{MDH-LVL}: \forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}:$ 
5886 30   auto trans_ll_IB<<l, d>><<b>><<  $\nabla_1^{(l, d)} \in \#PRT(1, 1), \dots, \nabla_D^{(l, d)} \in \#PRT(L, D) >$ 
5887 31   :=  $\downarrow \text{-mem}^{<b>} (l, d) T_b^{\text{IB}} [ N << \text{INP} >> <<b>> << \sigma_{\downarrow \text{-mem}}^{<b>} (l, d)(1) >> ( \frac{d_{\text{MDA}}}{\# \text{MDA}} (N_d) )_{d \in [1, D]_{\mathbb{N}} } ) ,$ 
5888 32   :
5889 33    $N << \text{INP} >> <<b>> << \sigma_{\downarrow \text{-mem}}^{<b>} (l, d)(D_b^{\text{IB}}) >> ( \frac{d_{\text{MDA}}}{\# \text{MDA}} (N_d) )_{d \in [1, D]_{\mathbb{N}} } ) ];$ 
5890 34
5891 35    $\forall (l, d) \in \text{MDH-LVL}: \forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}:$ 
5892 36   static ll_IB<<l, d>><<b>><<  $\nabla_1^{(l, d)} \in \#PRT(1, 1), \dots, \nabla_D^{(l, d)} \in \#PRT(L, D) >> ( \text{int } i_1, \dots,$ 
5893 37    $\text{int } i_{D_b^{\text{IB}}} ) \rightarrow T_b^{\text{IB}} a$ 
5894 38   {
5895 39   a := trans_ll_IB<<l, d>><<b>><<  $\nabla_1^{(l, d)}, \dots, \nabla_D^{(l, d)} >[ i_1 \sigma_{\downarrow \text{-mem}}^{<b>} (l, d)(1) , \dots ,$ 
5896 40    $i_{D_b^{\text{IB}}} \sigma_{\downarrow \text{-mem}}^{<b>} (l, d)(D_b^{\text{IB}}) ];$ 
5897 41   }
5898 42
5899 43   // scalar BUFS
5900 44    $\forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}:$ 
5901 45   auto trans_ll_IB<f>><<b>><<  $\nabla_1^{(f)} \in \#PRT(1, 1), \dots, \nabla_D^{(f)} \in \#PRT(L, D) >$ 
5902 46   :=  $\text{f} \downarrow \text{-mem}^{<b>} T_b^{\text{IB}} [ N << \text{INP} >> <<b>> << \sigma_{\text{f} \downarrow \text{-mem}}^{<b>} (1) >> ( \frac{d_{\text{MDA}}}{\# \text{MDA}} (N_d) )_{d \in [1, D]_{\mathbb{N}} } ) ,$ 
5903 47   :
5904 48    $N << \text{INP} >> <<b>> << \sigma_{\text{f} \downarrow \text{-mem}}^{<b>} (D_b^{\text{IB}}) >> ( \frac{d_{\text{MDA}}}{\# \text{MDA}} (N_d) )_{d \in [1, D]_{\mathbb{N}} } ) ];$ 
5905 49
5906 50    $\forall b \in [1, B^{\text{IB}}]_{\mathbb{N}}:$ 
5907 51   static ll_IB<<f>><<b>><<  $\nabla_1^{(f)} \in \#PRT(1, 1), \dots, \nabla_D^{(f)} \in \#PRT(L, D) >> ( \text{int } i_1, \dots,$ 
5908 52    $\text{int } i_{D_b^{\text{IB}}} ) \rightarrow T_b^{\text{IB}} a$ 
5909 53   {
5910 54   a := trans_ll_IB<<f>><<b>><<  $\nabla_1^{(f)}, \dots, \nabla_D^{(f)} >[ i_1 \sigma_{\text{f} \downarrow \text{-mem}}^{<b>} (1) , \dots , i_{D_b^{\text{IB}}} \sigma_{\text{f} \downarrow \text{-mem}}^{<b>} (D_b^{\text{IB}}) ];$ 
5911 55   }
5912 56
5913 57   // 0.4.3. output BUFS
5914 58
5915 59   // initial BUFS
5916 60    $\forall b \in [1, B^{\text{OB}}]_{\mathbb{N}}:$ 
5917 61   static ll_0B<<_1>><<b>><<  $\blacktriangle_1^{(\perp)} \in \#PRT(1, 1), \dots, \blacktriangle_D^{(\perp)} \in \#PRT(L, D) >> ( \text{int } i_1, \dots,$ 
5918 62    $\text{int } i_{D_b^{\text{OB}}} ) \rightarrow T_b^{\text{OB}} a$ 
5919 63   {
5920 64   a := trans_ll_0B<<_1>><<b>><<  $\blacktriangle_1^{(\perp)}, \dots, \blacktriangle_D^{(\perp)} >[ i_1 , \dots , i_{D_b^{\text{OB}}} ];$ 
5921 65   }
5922 66
5923 67   // re-composition BUFS
5924 68    $\forall (l, d) \in \text{MDH-LVL}: \forall b \in [1, B^{\text{OB}}]_{\mathbb{N}}:$ 
5925 69   auto trans_ll_0B<<l, d>><<b>><<  $\blacktriangle_1^{(l, d)} \in \#PRT(1, 1), \dots, \blacktriangle_D^{(l, d)} \in \#PRT(L, D) >$ 
5926
5927
5928
5929

```

```

5930 70   :=  $\uparrow \text{-mem}^{<\text{b}>}(l, d)$   $T_b^{\text{OB}}[ \text{N} << \text{OUT} >> << \sigma_{\uparrow\text{-mem}}^{<\text{b}>}(l, d)(1) >> ( (\xrightarrow{\text{d}_{\text{MDA}}}{\text{MDA}}(N_d))_{d \in [1, D]_{\mathbb{N}}} )$  ,
5931 71   :
5932 72    $\text{N} << \text{OUT} >> << \sigma_{\uparrow\text{-mem}}^{<\text{b}>}(l, d)(D_b^{\text{OB}}) >> ( (\xrightarrow{\text{d}_{\text{MDA}}}{\text{MDA}}(N_d))_{d \in [1, D]_{\mathbb{N}}} ) ];$ 
5933 73
5934 74    $\forall (l, d) \in \text{MDH-LVL} : \forall b \in [1, B^{\text{OB}}]_{\mathbb{N}} :$ 
5935 75   static  $\text{ll\_OB} << l, d >> << b >> < \mathbf{\Delta}_1^{(l, d)} \in \# \text{PRT}(1, 1), \dots, \mathbf{\Delta}_D^{(l, d)} \in \# \text{PRT}(L, D) >> ( \text{int } i_1, \dots,$ 
5936 76    $\text{int } i_{-D_b^{\text{OB}}} ) \rightarrow T_b^{\text{OB}} \text{ a}$ 
5937 77   {
5938 78    $\text{a} := \text{trans\_ll\_OB} << l, d >> << b >> < \mathbf{\Delta}_1^{(l, d)}, \dots, \mathbf{\Delta}_D^{(l, d)} >> [ i_{-} \sigma_{\uparrow\text{-mem}}^{<\text{b}>}(l, d)(1) , \dots ,$ 
5939 79    $i_{-} \sigma_{\uparrow\text{-mem}}^{<\text{b}>}(l, d)(D_b^{\text{OB}}) ];$ 
5940 80   }
5941 81
5942 82   // scalar BUFS
5943 83    $\forall b \in [1, B^{\text{OB}}]_{\mathbb{N}} :$ 
5944 84   auto  $\text{trans\_ll\_OB} << f >> << b >> < \mathbf{\Delta}_1^{(f)} \in \# \text{PRT}(1, 1), \dots, \mathbf{\Delta}_D^{(f)} \in \# \text{PRT}(L, D) >$ 
5945 85    $:= \mathbf{f} \uparrow\text{-mem}^{<\text{b}>} T_b^{\text{OB}}[ \text{N} << \text{OUT} >> << b >> < \sigma_{f \uparrow\text{-mem}}^{<\text{b}>}(1) >> ( (\xrightarrow{\text{d}_{\text{MDA}}}{\text{MDA}}(N_d))_{d \in [1, D]_{\mathbb{N}}} )$  ,
5946 86   :
5947 87    $\text{N} << \text{OUT} >> << b >> < \sigma_{f \uparrow\text{-mem}}^{<\text{b}>}(D_b^{\text{OB}}) >> ( (\xrightarrow{\text{d}_{\text{MDA}}}{\text{MDA}}(N_d))_{d \in [1, D]_{\mathbb{N}}} ) ];$ 
5948 88
5949 89    $\forall b \in [1, B^{\text{OB}}]_{\mathbb{N}} :$ 
5950 90   static  $\text{ll\_OB} << f >> << b >> < \mathbf{\Delta}_1^{(f)} \in \# \text{PRT}(1, 1), \dots, \mathbf{\Delta}_D^{(f)} \in \# \text{PRT}(L, D) >> ( \text{int } i_1, \dots,$ 
5951 91    $\text{int } i_{-D_b^{\text{OB}}} ) \rightarrow T_b^{\text{OB}} \text{ a}$ 
5952 92   {
5953 93    $\text{a} := \text{trans\_ll\_OB} << f >> << b >> < \mathbf{\Delta}_1^{(f)}, \dots, \mathbf{\Delta}_D^{(f)} >> [ i_{-} \sigma_{f \uparrow\text{-mem}}^{<\text{b}>}(1) , \dots , i_{-} \sigma_{f \uparrow\text{-mem}}^{<\text{b}>}(D_b^{\text{OB}}) ];$ 
5954 94   }

```

Listing 15. Low-Level BUFS

5960 where  $\mathbf{\Delta}_d^{\bullet}$  and  $\mathbf{\Delta}_d^{(\bullet)}$ , for  $\bullet \in \{\perp\} \cup \text{MDH-LVL} \cup \{f\}$ , are textually replaced by:

5962  
5963 
$$\mathbf{\Delta}_d^{(\bullet)}_1 = \begin{cases} p_d^1 & : \sigma_{\downarrow\text{-ord}}(l, d) < \bullet \\ * & : \text{else} \end{cases}$$

5964  
5965  
5966  
5967 
$$\mathbf{\Delta}_d^{(\bullet)} = \begin{cases} p_d^1 & : \sigma_{\uparrow\text{-ord}}(l, d) < \bullet \\ * & : \text{else} \end{cases}$$

5968 (symbol \* is taken from Definition 23) where < is defined according to the lexicographical order on  
5969  $\text{MDH-LVL} = [1, L]_{\mathbb{N}} \times [1, D]_{\mathbb{N}}$ , and:

5970  $\forall (l, d) \in \text{MDH-LVL} : \perp < (l, d) < f$

5971 Functions

5972 
$$\xrightarrow{\text{d}_{\text{MDA}}}{\text{MDA}}, \dots, \xrightarrow{D_{\text{MDA}}}{\text{MDA}}$$

5973

5979 are the index set functions *id* of combine operator concatenation `++` (Example 1), and functions

$$\xrightarrow[\otimes]{1} \text{MDA}, \dots, \xrightarrow[\otimes]{D} \text{MDA}$$

are the index set functions of combine operators  $\otimes_1, \dots, \otimes_D$ .

Note that we use generous BUFS sizes (lines 31-33, 46-48, 70-72, 85-87), as imperative-style programming models usually struggle with non-contiguous index ranges. We discuss optimizations targeting BUF sizes in Section E.3.

Note further that we do not need to initialize output buffers with neutral elements of combine operators in lines 64, 79, and 93 of Listing 15, as the buffers are initialized implicitly in the re-composition phase (Section E.3).

### E.0.5 MDAs.

Listing 16 shows our implementation of low-level MDAs (Definitions 12 and 28).

Note that for a particular choice of meta-parameters, low-level BUFs (Definition 13) are ordinary BUFs (Definition 5), as required by the types of functions `inp_view` and `out_view` (Definitions 8 and 10).

```

5995 1 // 0.5. MDAs
5996 2
5997 3 // 0.5.1. partitioned index sets
5998 4  $\forall d \in [1, D]_{\mathbb{N}} :$ 
5999 5 static  $I << d >> < p_d^1 \in [\#PRT(1, d)], \dots, p_d^L \in [\#PRT(L, d)] > ( \text{int } j' ) \rightarrow \text{int } i_j$ 
6000 6 {
6001 7      $i_j := ( p_d^1 * ( N_d / [\#PRT(1, d)] ) ) +$ 
6002 8      $\vdots$ 
6003 9      $p_d^L * ( N_d / [\#PRT(1, d)] * \dots * [\#PRT(L, d)] ) + j' );$ 
6004 10 }
6005 11
6006 12 // 0.5.2. input MDAs
6007 13  $\forall \bullet \in \{ \perp \} \cup \text{MDH-LVL} \cup \{ f \} :$ 
6008 14 static  $ll\_inp\_mda << \bullet >> < \overset{(\bullet)}{\nabla}_1^1 \in [\#PRT(1, 1)], \dots, \overset{(\bullet)}{\nabla}_D^L \in [\#PRT(L, D)] > ( \text{int } i_1, \dots,$ 
6009 15      $\text{int } i_D ) \rightarrow T_b^{IB} a$ 
6010 16 {
6011 17      $\forall b \in [1, B^{IB}]_{\mathbb{N}}, a \in [1, A_b^{IB}]_{\mathbb{N}} :$ 
6012 18      $a := ll\_IB << \bullet >> << b >> < \overset{(\bullet)}{\nabla}_1^1, \dots, \overset{(\bullet)}{\nabla}_D^L > ( \text{idx} << \text{INP} >> << b, a >> << 1 >> ( i_1, \dots, i_D ) ,$ 
6013 19      $\vdots$ 
6014 20      $\text{idx} << \text{INP} >> << b, a >> << D_b^{IB} >> ( i_1, \dots, i_D ) );$ 
6015 21 }
6016 22
6017 23 // 0.5.3. output MDAs
6018 24  $\forall \bullet \in \{ \perp \} \cup \text{MDH-LVL} \cup \{ f \} :$ 
6019 25 static  $ll\_out\_mda << \bullet >> < \overset{(\bullet)}{\Delta}_1^1 \in [\#PRT(1, 1)], \dots, \overset{(\bullet)}{\Delta}_D^L \in [\#PRT(L, D)] > ( \text{int } i_1, \dots,$ 
6020 26      $\text{int } i_D ) \rightarrow T_b^{OB} a$ 
6021 27 {
6022 28      $\forall b \in [1, B^{OB}]_{\mathbb{N}}, a \in [1, A_b^{OB}]_{\mathbb{N}} :$ 
6023 29      $a := ll\_OB << \bullet >> << b >> < \overset{(\bullet)}{\Delta}_1^1, \dots, \overset{(\bullet)}{\Delta}_D^L > ( \text{idx} << \text{OUT} >> << b, a >> << 1 >> ( i_1, \dots, i_D ) ,$ 
6024 30      $\vdots$ 
6025 31      $\text{idx} << \text{OUT} >> << b, a >> << D_b^{OB} >> ( i_1, \dots, i_D ) );$ 
6026 32 }

```

Listing 16. Low-Level MDAs

6028 For computing the partitioned index sets (lines 3-10), we exploit the following proposition.

6029 **Proposition 1.** Let  $\alpha \in T[N_1, \dots, N_D]$  be an arbitrary MDA that operates on contiguous index sets  
 6030  $[1, N_d]_{\mathbb{N}}, d \in [1, D]_{\mathbb{N}}$ . Let further be

$$6032 \quad \alpha^{<(p_1^1, \dots, p_d^1) \in P_1^1 \times \dots \times P_D^1 \mid \dots \mid (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>} : I_1^{<p_1^1 \dots p_1^L>} \times \dots \times I_D^{<p_D^1 \dots p_D^L>} \rightarrow T$$

6033 an arbitrary  $L$ -layered,  $D$ -dimensional,  $P$ -partitioning of MDA  $\alpha$ .

6034 It holds that  $j$ -th element within an MDA's part is accessed via index  $j$ :

$$6036 \quad I_d^{<p_d^1 \dots p_d^L>} = \left\{ \underbrace{\sum_{l \in [1, L]_{\mathbb{N}}} p_d^l * \frac{N_d}{\prod_{l' \in [1, l]_{\mathbb{N}}} P_d^{l'}} + 0,}_{\text{OS}} \quad \underbrace{\sum_{l \in [1, L]_{\mathbb{N}}} p_d^l * \frac{N_d}{\prod_{l' \in [1, l]_{\mathbb{N}}} P_d^{l'}} + 1, \dots}_{\text{OS}} \right\}$$

6040 PROOF. Since MDA  $\alpha$ 's index sets are contiguous ranges of natural numbers, it holds the  $i_j$  – the  
 6041 index to access the  $j$ -th element within an MDA's part (Definition 28) – is equal to  $j$  itself.  $\square$

## 6043 E.1 De-Composition Phase

6044 Listing 17 shows our implementation of the de-composition phase (Figure 19).

```
6045 1 // 1. de-composition phase
6046 2
6047 3 // 1.1. initialization
6048 4 ll_inp_mda<<1>> =: ll_inp_mda<<σ↓-ord(1,1)>>
6049 5
6050 6 // 1.2. main
6051 7 int p_ σ↓-ord(1,1) ∈ <↔-ass (1,1) > #PRT ( σ↓-ord(1,1) )
6052 8 {
6053 9     ll_inp_mda<<σ↓-ord(1,1)>> =: cc<σ↓-ord(1,1)> inp_mda<<σ↓-ord(1,2)>>;
6054 10    int p_ σ↓-ord(1,2) ∈ <↔-ass (1,2) > #PRT ( σ↓-ord(1,2) )
6055 11    {
6056 12        ll_inp_mda<<σ↓-ord(1,2)>> =: cc<σ↓-ord(1,2)> inp_mda<<σ↓-ord(1,3)>>;
6057 13        ..
6058 14        int p_ σ↓-ord(L,D) ∈ <↔-ass (L,D) > #PRT ( σ↓-ord(L,D) )
6059 15        {
6060 16            ll_inp_mda<<σ↓-ord(L,D)>> =: cc<σ↓-ord(L,D)> inp_mda<<f>>;
6061 17        }
6062 18        ..
6063 19    }
6064 20 }
```

6068 Listing 17. De-Composition Phase

6069 where

$$6070 \quad ll\_inp\_mda<<l, d>> =:_{cc<l, d>} ll\_inp\_mda<<l', d'>>$$

6072 abbreviates

$$6073 \quad ll\_inp\_mda<<l', d'>><\nabla_1^{(l', d')}, \dots, \nabla_D^{(l', d')}>, ll\_inp\_mda<<l, d>><\nabla_1^{(l, d)}, \dots, \nabla_D^{(l, d)}>>
 6074 \quad := cc\_inv<<d>><\frac{1}{MDA} \times MDA(I <<1>><\blacksquare_1^{(l, d)}, \dots, \blacksquare_1^{(l, d)}>>(\theta))>, // I_1
 6075 \quad : = cc\_inv<<d>><\frac{1}{MDA} \times MDA(I <<1>><\blacksquare_1^{(l, d)}, \dots, \blacksquare_1^{(l, d)}>>(\theta))>, // I_1$$

```

6077      :
6078       $\stackrel{D}{\Rightarrow}_{\#}^{\text{MDA}} ( \text{I} << D >> < \underset{1}{\blacksquare}_D, \dots, \underset{L}{\blacksquare}_D > (0) ), \quad // \text{I}_D$ 
6079
6080
6081       $\stackrel{d}{\Rightarrow}_{\#}^{\text{MDA}} ( \text{I} << d >> < \underset{1}{\boxplus}_d, \dots, \underset{L}{\boxplus}_d > (0) ), \quad // \text{P}$ 
6082       $\stackrel{d}{\Rightarrow}_{\#}^{\text{MDA}} ( \text{I} << d >> < \underset{1}{\boxtimes}_d, \dots, \underset{L}{\boxtimes}_d > (0) ) \quad // \text{Q}$ 
6083
6084      > ( 11\_inp\_mda << l, d >> < \underset{1}{\blacktriangledown}, \dots, \underset{L}{\blacktriangledown} > )
6085

```

6086 Here, functions  $\stackrel{1}{\#}^{\text{MDA}}, \dots, \stackrel{D}{\#}^{\text{MDA}}$  are the index set functions  $id$  of combine operator concatenation  $\#_1, \dots, \#_D$  (Example 1), and  $\underset{\bullet}{\blacksquare}_d, \underset{\bullet}{\boxplus}_d, \underset{\bullet}{\boxtimes}_d$ , for  $\bullet \in \text{MDH-LVL}$ , are textually replaced by:

$$\underset{\bullet}{\blacksquare}_d := \begin{cases} p_-(l, d) & : \sigma_{\downarrow\text{-ord}}(l, d) < \bullet \\ [p_-(l, d), \#PRT(l, d)]_{\mathbb{N}_0} & : (l, d) = \bullet \\ [0, \#PRT(l, d)]_{\mathbb{N}_0} & : \sigma_{\downarrow\text{-ord}}(l, d) > \bullet \end{cases}$$

$$\underset{\bullet}{\boxplus}_d := \begin{cases} p_-(l, d) & : \sigma_{\downarrow\text{-ord}}(l, d) < \bullet \\ p_-(l, d) & : (l, d) = \bullet \\ [0, \#PRT(l, d)]_{\mathbb{N}_0} & : \sigma_{\downarrow\text{-ord}}(l, d) > \bullet \end{cases}$$

$$\underset{\bullet}{\boxtimes}_d := \begin{cases} p_-(l, d) & : \sigma_{\downarrow\text{-ord}}(l, d) < \bullet \\ (p_-(l, d), \#PRT(l, d)]_{\mathbb{N}_0} & : (l, d) = \bullet \\ [0, \#PRT(l, d)]_{\mathbb{N}_0} & : \sigma_{\downarrow\text{-ord}}(l, d) > \bullet \end{cases}$$

6109 where  $<$  is defined as lexicographical order, according to Section E.0.4.

6110 Note that we re-use  $\text{inp\_mda} << l, d >>$  for the intermediate results given by different iterations of  
6111 variable  $p_-(l, d)$ . Correctness is ensured, as it holds:

$$A \subseteq B \Rightarrow \stackrel{d}{\#}^{\text{MDA}}(A) \subseteq \stackrel{d}{\#}^{\text{MDA}}(B)$$

6113 MDA  $\text{inp\_mda} << l, d >>$  has the following type when used for the intermediate result in a particular  
6114 iteration of  $p_-(l, d)$ :

$$\stackrel{1}{\#}^{\text{MDA}}(I_1^{< \underset{1}{\blacksquare}_1, \dots, \underset{D}{\blacksquare}_1 | \dots | \underset{1}{\blacksquare}_1, \dots, \underset{D}{\blacksquare}_D >}) \times \dots \times \stackrel{D}{\#}^{\text{MDA}}(I_D^{< \underset{1}{\blacksquare}_1, \dots, \underset{D}{\blacksquare}_1 | \dots | \underset{1}{\blacksquare}_1, \dots, \underset{D}{\blacksquare}_D >}) \rightarrow T^{\text{INP}}$$

6119 Here, for a set  $P \subseteq [0, \#PRT(l, d)]_{\mathbb{N}_0}$ , index set  $I_d^{< \dots | \dots P \dots | \dots >}$  denotes  $\bigcup_{p_d' \in P} I_d^{< \dots | \dots p_d' \dots | \dots >}$ .

## E.2 Scalar Phase

6124 Listing 18 shows our implementation of the scalar phase (Figure 19).

```

6126 1 // 2. scalar phase
6127 2 | int p_  $\sigma_{f\text{-ord}}(1, 1)$   $\in \leftrightarrow_{f\text{-ass}}(1, 1)$  #PRT(  $\sigma_{f\text{-ord}}(1, 1)$  )
6128 3 | `.
6129 4 | int p_  $\sigma_{f\text{-ord}}(L, D)$   $\in \leftrightarrow_{f\text{-ass}}(L, D)$  #PRT(  $\sigma_{f\text{-ord}}(L, D)$  )
6130 5 |
6131 6 {
6132 7 | (
6133 8 |     ll_out_mda<<f>><<
6134 9 |         p_(1, 1) , . . . , p_(1, D) ,
6135 10 |         . . .
6136 11 |         p_(L, 1) , . . . , p_(L, D)>><<b, a>>(
6137 12 |              $\stackrel{1}{\Rightarrow}_{MDA}$ ( I<<1>><p_(1, 1) , . . . , p_(L, 1)>(0) ) ,
6138 13 |             . . .
6139 14 |              $\stackrel{D}{\Rightarrow}_{MDA}$ ( I<<D>><p_(1, D) , . . . , p_(L, D)>(0) ) )
6140 15 |         )b \in [1, B^{OB}]_{\mathbb{N}}, a \in [1, A_b^{OB}]_{\mathbb{N}} := f( ( ll_inp_mda<<f>><< p_(1, 1) , . . . , p_(1, D) ,
6141 16 |             . . .
6142 17 |             p_(L, 1) , . . . , p_(L, D)>><<b, a>>(
6143 18 |                  $\stackrel{d}{\Rightarrow}_{MDA}$ ( I<<1>><p_(1, 1) , . . . , p_(L, 1)>(0) ) ,
6144 19 |                 . . .
6145 20 |                  $\stackrel{d}{\Rightarrow}_{MDA}$ ( I<<D>><p_(1, D) , . . . , p_(L, D)>(0) ) )
6146 21 |         )b \in [1, B^{IB}]_{\mathbb{N}}, a \in [1, A_b^{IB}]_{\mathbb{N}} )
6147
6148
6149 }

```

Listing 18. Scalar Phase

```

6150
6151
6152
6153
6154
6155
6156
6157
6158
6159
6160
6161
6162
6163
6164
6165
6166
6167
6168
6169
6170
6171
6172
6173
6174

```

### 6175 E.3 Re-Composition Phase

6176 Listing 19 shows our implementation of the re-composition phase (Figure 19).

6177

```

6178 1 // 3. re-composition phase
6179 2
6180 3 // 3.1. main
6181 4 int p_  $\sigma_{\uparrow\text{-ord}}(1, 1)$   $\in \text{leftrightarrow-ass}(1, 1)$  > #PRT (  $\sigma_{\uparrow\text{-ord}}(1, 1)$  )
6182 5 {
6183 6     int p_  $\sigma_{\uparrow\text{-ord}}(1, 2)$   $\in \text{leftrightarrow-ass}(1, 2)$  > #PRT (  $\sigma_{\uparrow\text{-ord}}(1, 2)$  )
6184 7     {
6185 8         ..
6186 9         int p_  $\sigma_{\uparrow\text{-ord}}(L, D)$   $\in \text{leftrightarrow-ass}(L, D)$  > #PRT (  $\sigma_{\uparrow\text{-ord}}(L, D)$  )
6187 10     {
6188 11         ll_out_mda <<  $\sigma_{\uparrow\text{-ord}}(L, D)$  >> :=  $_{\text{co}< \sigma_{\uparrow\text{-ord}}(L, D) >}$  out_mda << f >>;
6189 12     }
6190 13     ..
6191 14     ll_out_mda <<  $\sigma_{\uparrow\text{-ord}}(1, 2)$  >> :=  $_{\text{co}< \sigma_{\uparrow\text{-ord}}(1, 2) >}$  out_mda <<  $\sigma_{\uparrow\text{-ord}}(1, 3)$  >>;
6192 15     }
6193 16     ll_out_mda <<  $\sigma_{\uparrow\text{-ord}}(1, 1)$  >> :=  $_{\text{co}< \sigma_{\uparrow\text{-ord}}(1, 1) >}$  out_mda <<  $\sigma_{\uparrow\text{-ord}}(1, 2)$  >>;
6194 17 }
6195 18
6196 19 // 3.2. finalization
6197 20 ll_out_mda << 1 >> := ll_out_mda <<  $\sigma_{\uparrow\text{-ord}}(1, 1)$  >>
6200

```

6201 Listing 19. Re-Composition Phase

6202

6203 where

6204

6205  $\text{ll\_out\_mda} << l, d >> :=_{\text{co}< l, d >} \text{ll\_out\_mda} << l', d' >>$

6206

6207 abbreviates

6208

```

6209  $\text{ll\_out\_mda} << l, d >> < \overset{(l, d)}{\blacktriangle}_1, \dots, \overset{(l, d)}{\blacktriangle}_D >$ 
6210  $:= \text{co} << d >> < \underset{\otimes}{\overset{1}{\text{MDA}}} ( \text{I} << 1 >> < \overset{(l, d)}{\blacksquare}_1, \dots, \overset{(l, d)}{\blacksquare}_1 > (0) ), // I_1$ 
6211  $\vdots \qquad \qquad \qquad // \dots, I_{d-1}, I_{d+1}, \dots$ 
6212  $\underset{\otimes}{\overset{D}{\text{MDA}}} ( \text{I} << D >> < \overset{(l, d)}{\blacksquare}_D, \dots, \overset{(l, d)}{\blacksquare}_D > (0) ), // I_D$ 
6213
6214
6215  $\underset{\otimes}{\overset{d}{\text{MDA}}} ( \text{I} << d >> < \overset{(l, d)}{\boxplus}_d, \dots, \overset{(l, d)}{\boxplus}_d > (0) ), // P$ 
6216
6217  $\underset{\otimes}{\overset{d}{\text{MDA}}} ( \text{I} << d >> < \overset{(l, d)}{\boxtimes}_d, \dots, \overset{(l, d)}{\boxtimes}_d > (0) ) // Q$ 
6218
6219  $> ( \text{ll\_out\_mda} << l, d >> < \overset{(l, d)}{\blacktriangle}_1, \dots, \overset{(l, d)}{\blacktriangle}_D > ,$ 
6220  $\text{ll\_out\_mda} << l', d' >> < \overset{(l', d')}{\blacktriangle}_1, \dots, \overset{(l', d')}{\blacktriangle}_D > )$ 
6221
6222
6223

```

6224 Here, functions  $\xrightarrow[\otimes]{1}^{\text{MDA}}, \dots, \xrightarrow[\otimes]{D}^{\text{MDA}}$  are the index set function of combine operators  $\otimes_1, \dots, \otimes_D$  (Defi-  
 6225  
 6226 nition 2), and  $\bullet_{\text{d}}^{\text{I}}, \bullet_{\text{d}}^{\text{L}}, \bullet_{\text{d}}^{\text{R}}$ , for  $\bullet \in \text{MDH-LVL}$ , are textually replaced by:  
 6227

$$\bullet_{\text{d}}^{\text{I}} := \begin{cases} p_{-}(\text{l}, \text{d}) & : \sigma_{\uparrow\text{-ord}}(\text{l}, \text{d}) < \bullet \\ [0, p_{-}(\text{l}, \text{d})]_{\mathbb{N}_0} & : (\text{l}, \text{d}) = \bullet \\ [0, \#PRT(\text{l}, \text{d})]_{\mathbb{N}_0} & : \sigma_{\uparrow\text{-ord}}(\text{l}, \text{d}) > \bullet \end{cases}$$

$$\bullet_{\text{d}}^{\text{L}} := \begin{cases} p_{-}(\text{l}, \text{d}) & : \sigma_{\uparrow\text{-ord}}(\text{l}, \text{d}) < \bullet \\ [0, p_{-}(\text{l}, \text{d})]_{\mathbb{N}_0} & : (\text{l}, \text{d}) = \bullet \\ [0, \#PRT(\text{l}, \text{d})]_{\mathbb{N}_0} & : \sigma_{\uparrow\text{-ord}}(\text{l}, \text{d}) > \bullet \end{cases}$$

$$\bullet_{\text{d}}^{\text{R}} := \begin{cases} p_{-}(\text{l}, \text{d}) & : \sigma_{\uparrow\text{-ord}}(\text{l}, \text{d}) < \bullet \\ p_{-}(\text{l}, \text{d}) & : (\text{l}, \text{d}) = \bullet \\ [0, \#PRT(\text{l}, \text{d})]_{\mathbb{N}_0} & : \sigma_{\uparrow\text{-ord}}(\text{l}, \text{d}) > \bullet \end{cases}$$

6247 where  $<$  is defined as lexicographical order, according to Section E.0.4.  
 6248

6249 Note that we assume for index set functions  $\xrightarrow[\otimes]{d}^{\text{MDA}}$  that  
 6250

$$A \subseteq B \Rightarrow \xrightarrow[\otimes]{d}^{\text{MDA}}(A) \subseteq \xrightarrow[\otimes]{d}^{\text{MDA}}(B)$$

6253 (which holds for all kinds of index set functions used in this paper, e.g., in Examples 1 and 2) so that  
 6254 we can re-use `out_mda<<l, d>>` for the intermediate results given by different iterations of `p_(l, d)`.  
 6255 MDA `inp_mda<<l, d>>` has the following type when used for the intermediate result in a particular  
 6256 iteration of variable `p_(l, d)`:

$$6257 \quad T^{\text{INP}} \left[ \xrightarrow[\otimes]{1}^{\text{MDA}}(I_1^{<\bullet_1^{\text{I}, \text{d}}, \dots, \bullet_D^{\text{I}, \text{d}} | \dots | \bullet_1^{\text{L}, \text{d}}, \dots, \bullet_D^{\text{L}, \text{d}} >}), \dots, \xrightarrow[\otimes]{D}^{\text{MDA}}(I_D^{<\bullet_1^{\text{I}, \text{d}}, \dots, \bullet_D^{\text{I}, \text{d}} | \dots | \bullet_1^{\text{L}, \text{d}}, \dots, \bullet_D^{\text{L}, \text{d}} >}) \right]$$

6260 Note that in line 11 of Listing 19, we implicitly override the uninitialized value in `out_mda<<l, d>>`  
 6261 (not explicitly stated in the listing for brevity), thereby avoiding initializing output buffers with  
 6262 neutral elements of combine operators.  
 6263

## CODE-LEVEL OPTIMIZATIONS

6265 We consider optimizations that are below the abstraction level of our low-level representation  
 6266 (Section 3) as *code-level optimizations*. For some of these code-level optimizations, like loop fusion,  
 6267 we do not want to rely on the underlying compiler (e.g., the OpenMP/CUDA/OpenCL compiler): we  
 6268 exactly know the structure of our code presented in Section E, thus being able to implement  
 6269 code-level optimizations without requiring complex compiler analyses for optimizations.  
 6270

6271 Since code-level optimizations are not the focus of this work, we give a brief, basic overview of  
 6272 our these optimizations which our systems performs for the underlying compiler (OpenMP, CUDA,  
 6273

6273 OpenCL, etc). We will thoroughly discuss our code-level optimizations and how we apply them  
6274 to our code in future work. Further code-level optimizations, like loop unrolling, are currently  
6275 mostly left to the underlying compiler, e.g., the OpenMP, CUDA, or OpenCL compiler. In our future  
6276 work, we aim to incorporate code-level optimizations, such as loop unrolling, into our auto-tunable  
6277 optimization process.

6278 *Loop Fusion.* In Listings 17, 18, 19, the lines containing symbol " $\epsilon$ " are mapped to for-loops. These  
6279 loops can often be fused; for example, when parameters 3, 6, 11 in Table 1 coincide (as in Figure 17).  
6280 Besides reducing the overhead caused by loop control structures, loop fusion in particular allows  
6281 significantly reducing the memory footprint: we can re-use the same memory region for each BUF  
6282 partition (Definition 13), rather than allocating memory for all these partition.

6283 *Buffer Elimination.* In Listing 15, we allocate BUFs for each combination of a layer and dimension.  
6284 However, when memory regions and memory layouts of BUFs coincide, we can avoid a new BUF  
6285 allocation, by re-using the BUF of the upper level, thereby again reducing memory footprint.

6286 *Buffer Size Reduction.* We can reduce the sizes of BUFs when specific classes of index functions  
6287 are used for views (Definitions 8 and 10). For example, in the case of *Dot Product* (DOT) (Figure 16),  
6288 when accessing its input in a strided fashion – via index function  $k \mapsto (2 * k)$ , instead of function  
6289  $k \mapsto (k)$  (as in Figure 16) – we would have to allocate BUFs (Listing 15, lines 30 and 45) of size  
6290  $2 * K$  for an input size of  $K \in \mathbb{N}$ ; in these BUFs, each second value (accessed via indices  $2 * k + 1$ )  
6291 would be undefined. We avoid this waste of memory, by using index function  $k \mapsto (k)$  instead  
6292 of  $k \mapsto (2 * k)$  for allocated BUFs (Listing 16, lines 18-20 and 29-31, case " $\bullet \neq \perp$ "), which avoids  
6293 index functions' leading factors and potential constant additions. Thereby, we reduce the memory  
6294 footprint from  $2 * K$  to  $K$ . Furthermore, according to our partitioning strategy (Listing 16, line 5,  
6295 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  
6296 offset  $OS \in \mathbb{N}$ . We avoid such offset by using  $k \mapsto (2 * (k - OS))$ , thereby further reducing the  
6297 memory footprint.

6298 *Memory Operation Minimization.* In our code, we access BUFs uniformly via MDAs (Listing 16),  
6299 which may cause unnecessary memory operations. For example, in the de-composition phase  
6300 (Listing 17) of, e.g., matrix multiplication (MatMul) (Figure 16), we iterate over all dimensions of the  
6301 input (i.e., the  $i, j, k$  dimensions) for de-composition (Listing 12). However, the  $A$  input matrix of  
6302 MatMul is accessed via MDA indices  $i$  and  $k$  only (Figure 16). We avoid these unnecessary memory  
6303 operations ( $J$ -many in the case of an input MDA of size  $J$  in dimension  $j$ ) by using index 0 only  
6304 in dimension  $j$  for the de-composition of the  $A$  matrix. Analogously, we use index 0 only in the  $i$   
6305 dimension for the de-composition of MatMul's  $B$  matrix which is accessed via MDA indices  $k$  and  $j$   
6306 only. Moreover, we exploit all available parallelism for memory copy operations. For example, for  
6307 MatMul, we use also the threads intended for the  $j$  dimension when de-composing the  $A$  matrix,  
6308 and we use the threads in the  $i$  dimension for the  $B$  matrix. For this, we flatten the thread ids over  
6309 all dimensions  $i, j, k$  and re-structure them only in dimensions  $i, k$  (for the  $A$  matrix) or  $k, j$  (for the  
6310  $B$  matrix).

6311 *Constant Substitution.* We use constants whenever possible. For example, in CUDA, variable  
6312 `threadIdx.x` returns the thread id in dimension  $x$ . However, in our code, we use constant 0  
6313 instead of `threadIdx.x` when only one thread is started in dimension  $x$ , enabling the CUDA  
6314 compiler to significantly simplify arithmetic expressions.

6315

6316

6317

6318

6319

6320

6321