

WESTFÄLISCHE
WILHELMS-UNIVERSITÄT
MÜNSTER



Array Programming via Multi-Dimensional Homomorphisms

Ari Rasch, Richard Schulze, Sergei Gorlatch

University of Münster, Germany

Existing Work

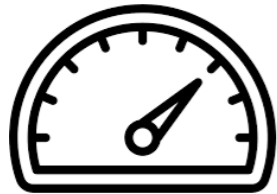
Array programming is at the heart of popular, existing approaches:

<i>Class</i>	<i>Popular Examples</i>	<i>Performance</i>	<i>Portability</i>	<i>Productivity</i>
Scheduling	TVM, Halide, Fireiron	✓	often require re-design/extension for new architectures	incorporate user into optimization process
Polyhedral	TC, PPCG, Pluto	struggle with reductions (e.g., dot in MatMul)	transformations chosen toward particular architectures and data characteristics	✓
Functional	Lift	✓	transformations designed toward particular architectures and data characteristics	often incorporate user into optimization process
Domain-Specific	cuBLAS, oneMKL	✓	hand-optimized toward particular architecture and data characteristics	(✓)
Higher-Level	<i>Futhark, Dex, ATL, Yang et al. [POPL'21], ...</i>	<i>We consider these approaches as greatly combinable with our approach 👍 (as frontends)</i>		

The existing approaches struggle with achieving simultaneously Performance & Portability & Productivity

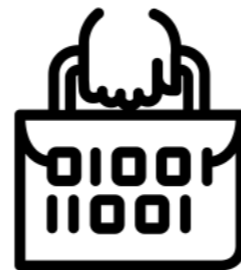
Existing Work

Why is **simultaneously** achieving **Performance & Portability & Productivity** challenging:



Performance

- hardware-specific code generation & optimization: exploiting core & memory hierarchy
- data-specific code generation & optimization: respecting data layouts and sizes
- ...



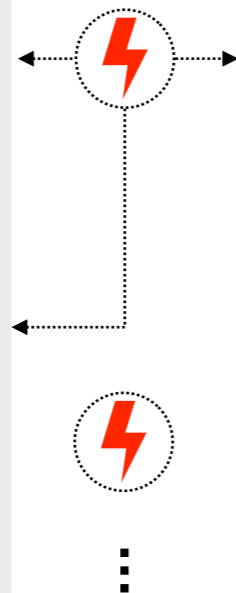
Portability

- flexible in hardware/data-specific optimizations (*performance portability*)
- easily extensible toward new target architectures (*functional portability*)
- ...



Productivity

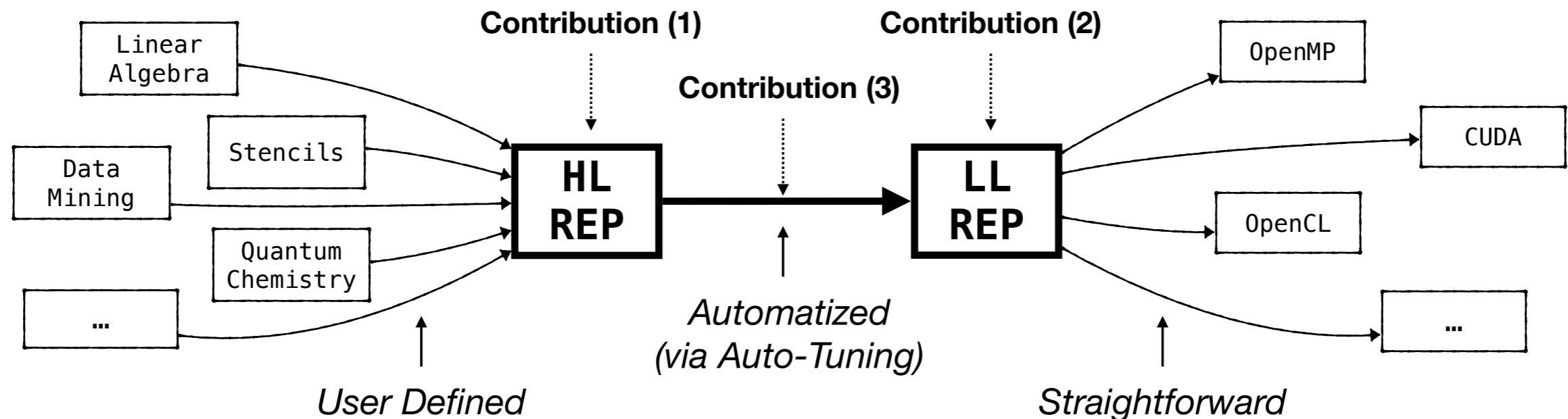
- freeing programmer from hardware and optimization details
- supporting wide range of target computations (*expressivity*)
- ...



***Achieving simultaneously
Performance & Portability & Productivity
poses challenging, often even contradicting, requirements***

Goal of this Work

Achieving **Performance** & **Portability** & **Productivity** in one approach, based on three major contributions:



1. **High-Level Functional Representation** that is formally defined, and that *expresses data-parallel computations* on a high-level of abstraction while still capturing all information relevant for generating high performing code
2. **Low-Level Functional Representation** that enables formally *expressing and reasoning about optimizations* for the memory and core hierarchies of state-of-the-art architectures, and that can be straightforwardly *transformed to executable program code* (e.g., in OpenMP, CUDA, OpenCL, ...)
3. **Lowering Process** that *fully automatically* lowers our *high-level representation* to an hardware/data-optimized instance of our *low-level representation*, in a formally sound and automatically optimizable (auto-tunable) manner

Agenda

1. **Contribution 1:** High-Level Representation
2. **Contribution 2:** Low-Level Representation
3. **Contribution 3:** Lowering: *High-Level Representation* → *Low-Level Representation*
4. Experimental Results (Performance & Portability & Productivity)
5. Conclusion
6. Future Work

High-Level Representation

Goals:

1. **Uniform:**

should be able to express any kind of data-parallel computation, without relying on computation-specific building blocks, extensions, etc.

2. **Minimalistic:**

should rely on less building blocks to keep language small and simple

3. **Structured:**

avoiding compositions and nestings of building blocks as much as possible, thereby further contributing to usability and simplicity of our language

```
MatVec<T∈TYPE | I,K∈ℕ> := out_view<T>( w:(i,k)↦(i) ) ◦  
                               md_hom<I,K>( *, ( #, + ) ) ◦  
                               inp_view<T,T>( M:(i,k)↦(i,k) , v:(i,k)↦(k) )
```

High-Level Representation of MatVec

High-Level Representation

Quick Reminder

Data-parallel computations: how are they characterized?

1. applying the *same function* f to each point in a multi-dimensional grid of data
2. combining the obtained results in the grid's different dimensions using *combine operators*

Matrix-Vector Multiplication:

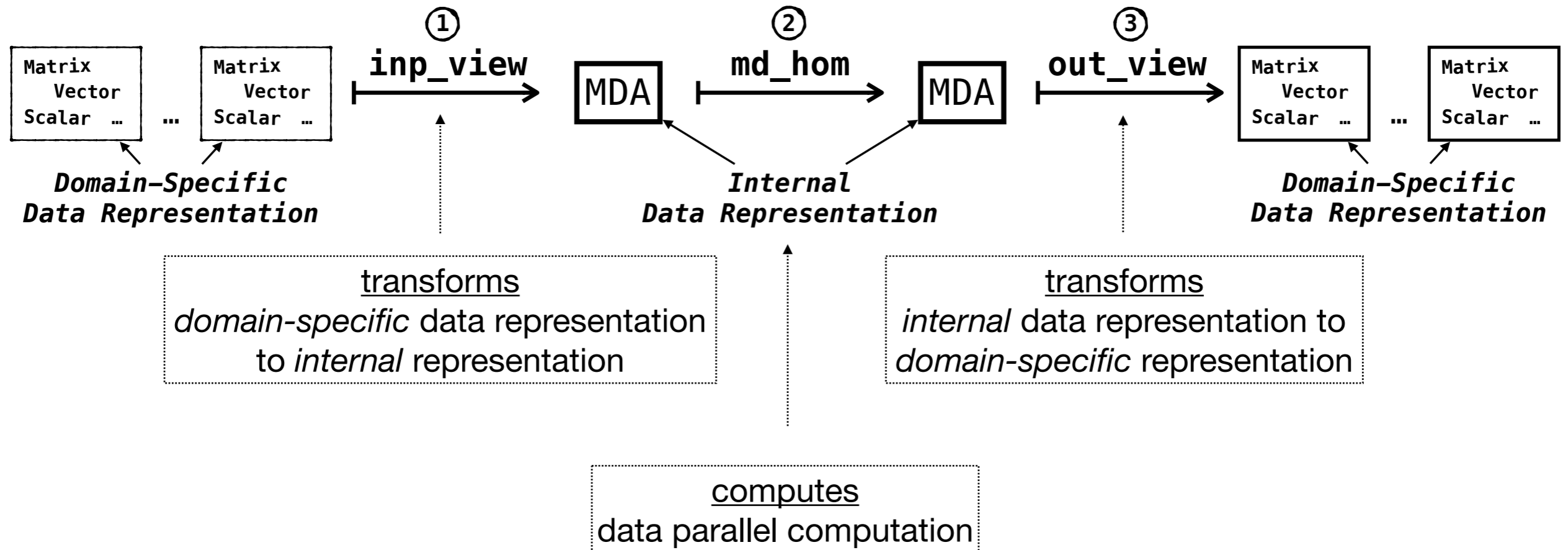
$$\begin{pmatrix} M_{1,1} & \dots & M_{1,K} \\ \vdots & \ddots & \vdots \\ M_{I,1} & \dots & M_{I,K} \end{pmatrix}, \begin{pmatrix} v_1 \\ \vdots \\ v_K \end{pmatrix} \xrightarrow{\text{MatVec}} \begin{array}{c} \xrightarrow{\otimes_2} \\ \left. \begin{array}{ccc} f(M_{1,1}, v_1) & \dots & f(M_{1,K}, v_K) \\ \vdots & \ddots & \vdots \\ f(M_{I,1}, v_1) & \dots & f(M_{I,K}, v_K) \end{array} \right| \\ \downarrow \otimes_1 \end{array} = \begin{pmatrix} M_{1,1} * v_1 + \dots + M_{1,K} * v_K \\ \vdots \\ M_{I,1} * v_1 + \dots + M_{I,K} * v_K \end{pmatrix} = \begin{pmatrix} w_1 \\ \vdots \\ w_I \end{pmatrix}$$

Jacobi 1D:

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

High-Level Representation

Overview:



Our high-level representation expresses data-parallel computations
— *agnostic from hardware and optimization details* —
using exactly three higher-order functions only

High-Level Representation

Example: MatVec expressed in MDH's High-Level Representation¹

```
MatVec<T∈TYPE | I,K∈ℕ> := out_view<T>( w:(i,k)↦(i) ) ◦  
                             md_hom<I,K>( *, (⊕,+) ) ◦  
                             inp_view<T,T>( M:(i,k)↦(i,k) , v:(i,k)↦(k) )
```

High-Level Representation of MatVec

What is happening here:

- `inp_view` captures the accesses to input data
- `md_hom` expresses the data-parallel computation
- `out_view` captures the accesses to output data

¹We can generate such MDH expressions also automatically from straightforward (annotated) C code [IMPACT'19]

High-Level Representation

md_hom	f	⊗ ₁	⊗ ₂	⊗ ₃	⊗ ₄	Views	inp_view		out_view
							A	B	C
Dot	*	+	/	/	/	Dot	(k) ↦ (k)	(k) ↦ (k)	(k) ↦ ()
MatVec	*	++	+	/	/	MatVec	(i,k) ↦ (i,k)	(i,k) ↦ (k)	(i,k) ↦ (i)
MatMul	*	++	++	+	/	MatMul	(i,j,k) ↦ (i,k)	(i,j,k) ↦ (k,j)	(i,j,k) ↦ (i,j)
MatMul ^T	*	++	++	+	/	MatMul ^T	(i,j,k) ↦ (k,i)	(i,j,k) ↦ (j,k)	(i,j,k) ↦ (j,i)
bMatMul	*	++	++	++	+	bMatMul	(b,i,j,k) ↦ (b,i,k)	(b,i,j,k) ↦ (b,k,j)	(b,i,j,k) ↦ (b,i,j)

1) Linear Algebra Routines

md_hom	f	⊗ ₁	⊗ ₂	⊗ ₃	⊗ ₄	⊗ ₅	⊗ ₆	⊗ ₇	⊗ ₈	⊗ ₉	⊗ ₁₀
Conv2D	*	++	++	+	+	/	/	/	/	/	/
MCC	*	++	++	++	++	+	+	+	/	/	/
MCC.Capsule	*	++	++	++	++	+	+	+	++	++	+

md_hom	f	⊗ ₁	⊗ ₂	Views	inp_view	out_view
					A	Out
MBBS	id	++ _{prefix-sum(+)}	+	MBBS	(i,j) ↦ (i,j)	(i) ↦ (i)

8) Maximum Bottom Box Sum

Views	inp_view		out_view
	I	F	O
Conv2D	(p,q,r,s) ↦ (p+r,q+s)	(p,q,r,s) ↦ (r,s)	(p,q,r,s) ↦ (p,q)
MCC	(n,p,...) ↦ (n,p+r,q+s,c)	(n,p,...) ↦ (k,r,s,c)	(n,p,...) ↦ (n,p,q,k)
MCC.Capsule	(n,p,...) ↦ (n,p+r,q+s,c,cm,ck)	(n,p,...) ↦ (k,r,s,c,ck,cn)	(n,p,...) ↦ (n,p,q,k,cm,cn)

2) Convolution Stencils

md_hom	f	⊗ ₁	⊗ ₂	Views	inp_view	out_view
					I	O
Jacobi1D	J _{1D}	++	/	Jacobi1D	(i) ↦ (i+0) , (i) ↦ (i+1) , (i) ↦ (i+2)	(i) ↦ (i)
Jacobi2D	J _{2D}	++	++	Jacobi2D	(i,j) ↦ (i,j+1) , (i,j) ↦ (i+1,j) , ...	(i,j) ↦ (i,j)

3) Jacobi Stencils

md_hom	f	⊗ ₁	⊗ ₂	Views	inp_view	out_view	
					N	E	M
PRL	wght	++	max _{PRL}	PRL	(i,j) ↦ (i)	(i,j) ↦ (j)	(i,j) ↦ (i)

4) Probabilistic Record Linkage

md_hom	f	⊗ ₁	Views	inp_view	out_view	
				I	O ₁	O ₂
map(f)	f	++	map(f)	(i) ↦ (i)	(i) ↦ (i)	/
reduce(⊕)	id	⊕	reduce(⊕)	(i) ↦ (i)	(i) ↦ ()	/
reduce(⊕, ⊗)	(x) ↦ (x,x)	(⊕, ⊗)	reduce(⊕, ⊗)	(i) ↦ (i)	(i) ↦ ()	(i) ↦ ()

6) Map/Reduce Patterns

md_hom	f	⊗ ₁	Views	inp_view	out_view
				I	O
scan(⊕)	id	++ _{prefix-sum(⊕)}	scan(⊕)	(i) ↦ (i)	(i) ↦ (i)

7) Scan Pattern

md_hom	f	⊗ ₁	⊗ ₂	Views	inp_view	out_view	
					Bins	Elms	Out
Histo	f _{Histo}	++	+	Histo	(b,e) ↦ (b)	(b,e) ↦ (e)	(b,e) ↦ (b)

5) Histogram

Our high-level representation is capable of expressing the various kinds of data-parallel computations which often differ in major characteristics

Summary: High-Level Representation

- **Uniform:** for all the different kinds of data-parallel computations, by exploiting the common algebraic properties of data-parallel computations (rather than relying on domain-specific building blocks)
- **Minimalistic:** by relying on exactly three higher-order functions only
- **Structured:**
 - `inp_view` prepares the *input data*
 - `md_hom` computes *data-parallel computation*
 - `out_view` prepares the *output data*
- **Full formal foundation**

Definition 5 (Buffer). Let $T \in \text{TYPE}$ be an arbitrary scalar type, $D \in \mathbb{N}_0$ a natural number⁹, and $N := (N_1, \dots, N_D) \in \mathbb{N}^D$ a sequence of natural numbers.
A *Buffer (BUF)* b that has *dimensionality* D , *size* N , and *scalar type* T is a function with the following signature:

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

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 as set of *MDA index sets*. 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 sequence of D -many MDA index sets, and $N := (N_1, \dots, N_D) := (|I_1|, \dots, |I_D|)$ the sequence of index sets' sizes.

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

$$a : I_1 \times \dots \times I_D \rightarrow T$$

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

Definition 5 (Buffer). Let $T \in \text{TYPE}$ be an arbitrary scalar type, $D \in \mathbb{N}_0$ a natural number⁹, and $N := (N_1, \dots, N_D) \in \mathbb{N}^D$ a sequence of natural numbers.

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

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

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

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

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

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

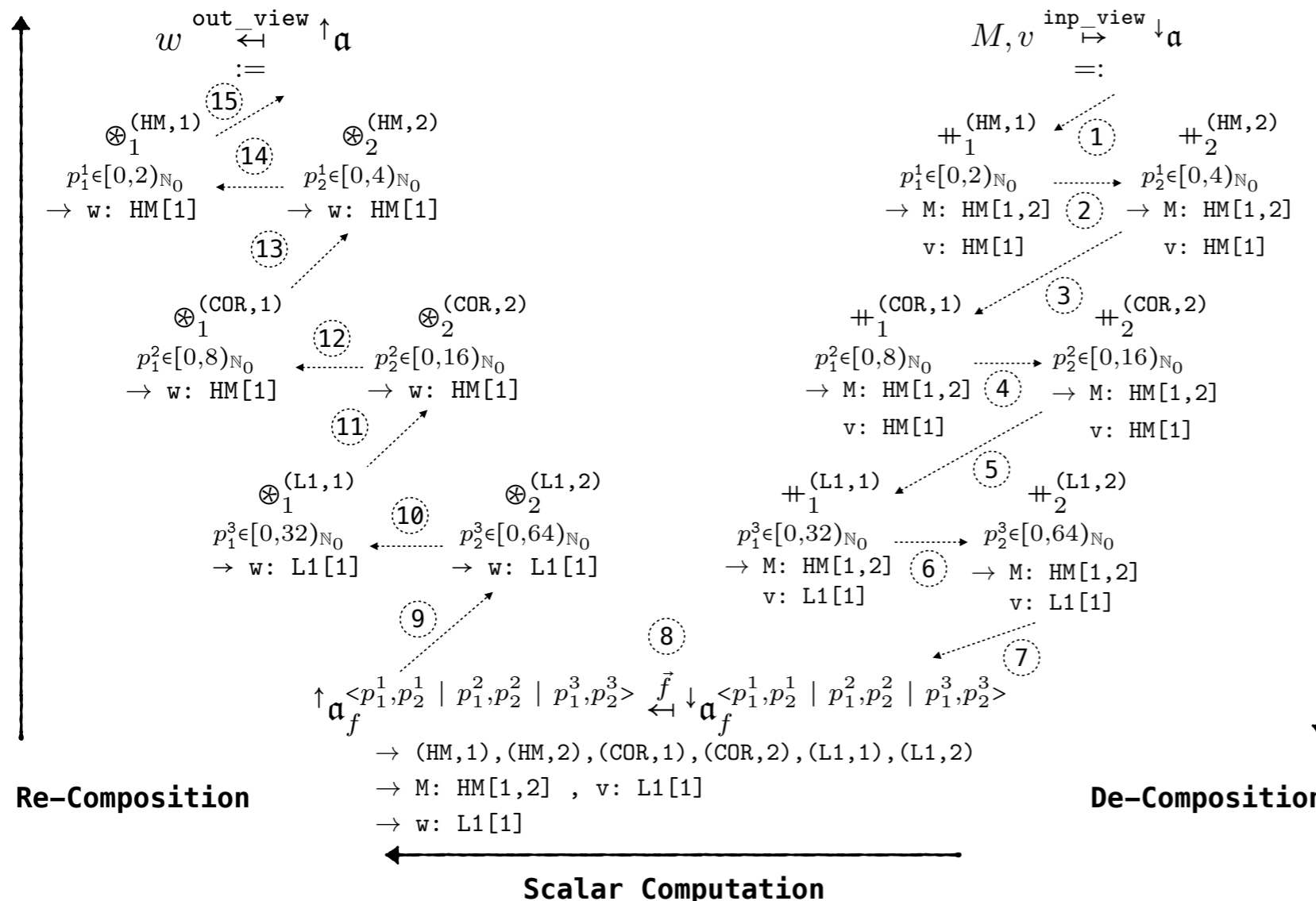
$$T[I_1, \dots, \underbrace{\Rightarrow_{\text{MDA}}^{\text{MDA}}(P), \dots, I_D}_{\uparrow d}] \times T[I_1, \dots, \underbrace{\Rightarrow_{\text{MDA}}^{\text{MDA}}(Q), \dots, I_D}_{\uparrow d}] \rightarrow T[I_1, \dots, \underbrace{\Rightarrow_{\text{MDA}}^{\text{MDA}}(P \cup Q), \dots, I_D}_{\uparrow d}]$$

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

Low-Level Representation

Goals:

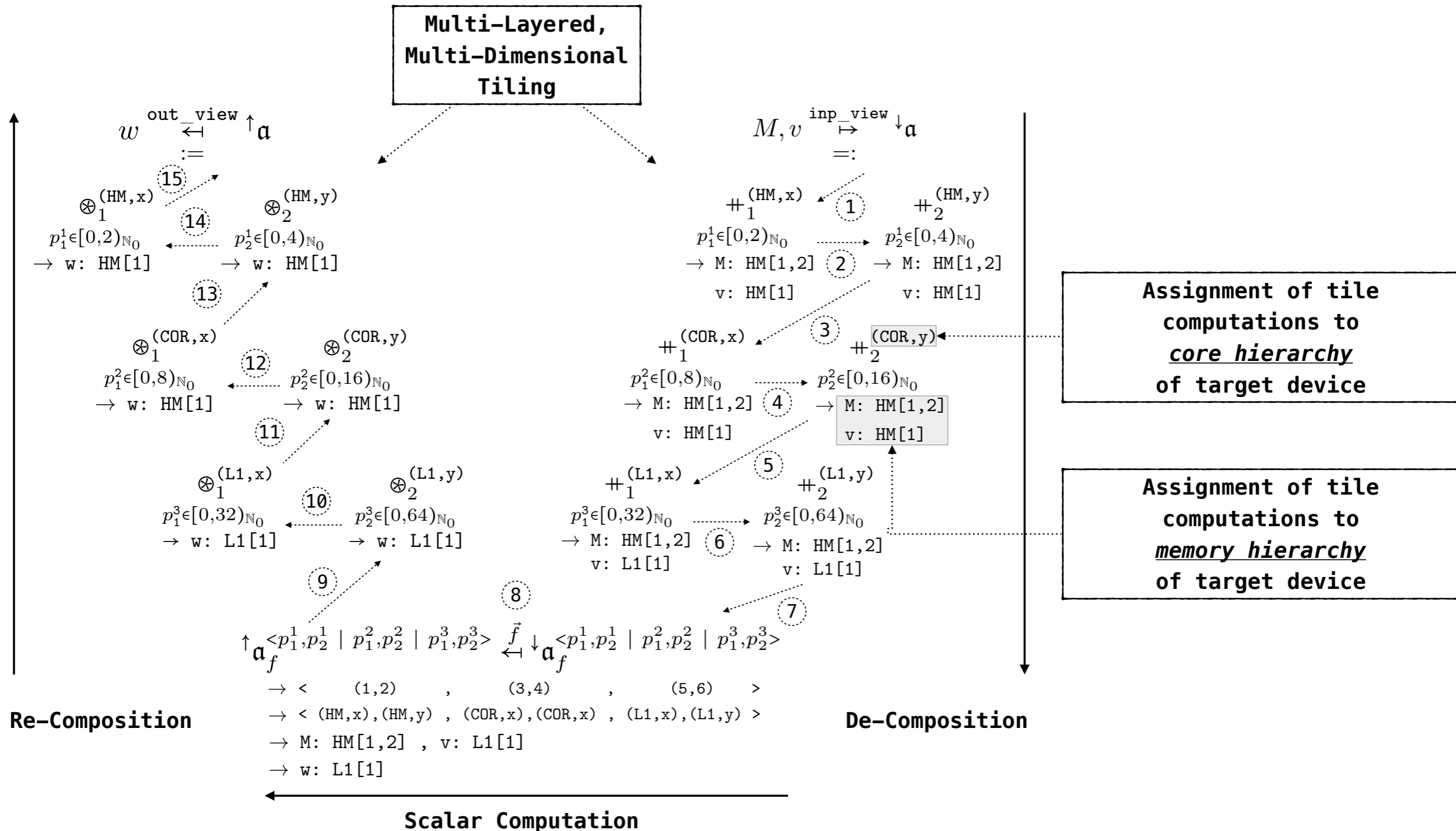
1. Expressing a hardware- & data-optimized *de-composition* and *re-composition* of data-parallel computations
2. Being straightforwardly transformable to executable program code (e.g., in OpenMP, CUDA, and OpenCL) — optimization decisions are explicitly expressed in low-level representation !



*Low-Level Representation of MatVec
(particular, straightforward instance,
for an artificial architecture)*

Low-Level Representation

How do we express *(De/Re)-Compositions* in our low-level representation:



Low-Level Representation

Code generation:

```

1 // 0.1.2. combine operators
2
3 // pre-implemented combine operators
4
5 // inverse concatenation
6  $\forall d \in \mathbb{N}$ :
7  $cc\_inv \langle\langle d \rangle\rangle \langle I_1, \dots, I_{d-1}, I_{d+1}, \dots, I_D \in \text{MDA-IDX-SETS}, (P, Q) \in \text{MDA-IDX-SETS} \times \text{MDA-IDX-SETS} \rangle ($ 
8    $T^{INP}[I_1, \dots, I_{d-1}, id(P \cup Q), I_{d+1}, \dots, I_D] \text{ lhs } ,$ 
9    $T^{INP}[I_1, \dots, I_{d-1}, id(Q), I_{d+1}, \dots, I_D] \text{ rhs } )$ 
10 {
11   int  $i_{-1} \in I_1$ 
12    $\vdots$ 
13   int  $i_{-d-1} \in I_{d-1}$ 
14   int  $i_{-d+1} \in I_{d+1}$ 
15    $\vdots$ 
16   int  $i_{-D} \in I_D$ 
17   {
18     int  $i_{-d} \in P$ 
19      $res[i_{-1}, \dots, i_{-d}, \dots, i_{-D}] =: lhs[i_{-1}, \dots, i_{-d}, \dots, i_{-D}];$ 
20     int  $i_{-d} \in Q$ 
21      $res[i_{-1}, \dots, i_{-d}, \dots, i_{-D}] =: rhs[i_{-1}, \dots, i_{-d}, \dots, i_{-D}];$ 
22   }
23 }

```

Listing 11. Pre-Implemented Combine Operators

```

1 // 3. re-composition phase
2
3 // 3.1. main
4 int  $p_{-}\sigma_{\uparrow\text{-ord}}(1,1) \in \langle \leftrightarrow_{\uparrow\text{-ass}}(1,1) \rangle \#PRT(\sigma_{\uparrow\text{-ord}}(1,1))$ 
5 {
6   int  $p_{-}\sigma_{\uparrow\text{-ord}}(1,2) \in \langle \leftrightarrow_{\uparrow\text{-ass}}(1,2) \rangle \#PRT(\sigma_{\uparrow\text{-ord}}(1,2))$ 
7   {
8      $\vdots$ 
9     int  $p_{-}\sigma_{\uparrow\text{-ord}}(L,D) \in \langle \leftrightarrow_{\uparrow\text{-ass}}(L,D) \rangle \#PRT(\sigma_{\uparrow\text{-ord}}(L,D))$ 
10    {
11       $ll\_out\_mda \langle\langle \sigma_{\uparrow\text{-ord}}(L,D) \rangle\rangle :=_{co \langle \sigma_{\uparrow\text{-ord}}(L,D) \rangle} out\_mda \langle\langle f \rangle\rangle;$ 
12    }
13     $\vdots$ 
14     $ll\_out\_mda \langle\langle \sigma_{\uparrow\text{-ord}}(1,2) \rangle\rangle :=_{co \langle \sigma_{\uparrow\text{-ord}}(1,2) \rangle} out\_mda \langle\langle \sigma_{\uparrow\text{-ord}}(1,3) \rangle\rangle;$ 
15  }
16   $ll\_out\_mda \langle\langle \sigma_{\uparrow\text{-ord}}(1,1) \rangle\rangle :=_{co \langle \sigma_{\uparrow\text{-ord}}(1,1) \rangle} out\_mda \langle\langle \sigma_{\uparrow\text{-ord}}(1,2) \rangle\rangle;$ 
17 }
18
19 // 3.2. finalization
20  $ll\_out\_mda \langle\langle 1 \rangle\rangle := ll\_out\_mda \langle\langle \sigma_{\uparrow\text{-ord}}(1,1) \rangle\rangle$ 

```

Listing 18. Re-Composition Phase

```

1 // 2. scalar phase
2 int  $p_{-}\sigma_{f\text{-ord}}(1,1) \in \langle \leftrightarrow_{f\text{-ass}}(1,1) \rangle \#PRT(\sigma_{f\text{-ord}}(1,1))$ 
3  $\vdots$ 
4 int  $p_{-}\sigma_{f\text{-ord}}(L,D) \in \langle \leftrightarrow_{f\text{-ass}}(L,D) \rangle \#PRT(\sigma_{f\text{-ord}}(L,D))$ 
5 {
6   (
7      $ll\_out\_mda \langle\langle f \rangle\rangle \langle\langle$ 
8        $p_{-}(1,1), \dots, p_{-}(1,D),$ 
9        $\dots$ 
10       $p_{-}(L,1), \dots, p_{-}(L,D) \rangle\rangle \langle\langle b, a \rangle\rangle ($ 
11         $\frac{1}{\otimes_{MDA}} \langle I \langle\langle 1 \rangle\rangle \langle\langle p_{-}(1,1), \dots, p_{-}(L,1) \rangle\rangle (\emptyset) \rangle ,$ 
12         $\vdots$ 
13         $\frac{D}{\otimes_{MDA}} \langle I \langle\langle D \rangle\rangle \langle\langle p_{-}(1,D), \dots, p_{-}(L,D) \rangle\rangle (\emptyset) \rangle )$ 
14       $\rangle_{b \in [1, B^{OB}]_{\mathbb{N}}, a \in [1, A_b^{IB}]_{\mathbb{N}}} := f( ( ll\_inp\_mda \langle\langle f \rangle\rangle \langle\langle$ 
15         $p_{-}(1,1), \dots, p_{-}(1,D),$ 
16         $\dots$ 
17         $p_{-}(L,1), \dots, p_{-}(L,D) \rangle\rangle \langle\langle b, a \rangle\rangle ($ 
18           $\frac{d}{\#_{MDA}} \langle I \langle\langle 1 \rangle\rangle \langle\langle p_{-}(1,1), \dots, p_{-}(L,1) \rangle\rangle (\emptyset) \rangle ,$ 
19           $\vdots$ 
20           $\frac{d}{\#_{MDA}} \langle I \langle\langle D \rangle\rangle \langle\langle p_{-}(1,D), \dots, p_{-}(L,D) \rangle\rangle (\emptyset) \rangle )$ 
21         $\rangle_{b \in [1, B^{IB}]_{\mathbb{N}}, a \in [1, A_b^{IB}]_{\mathbb{N}}}$ 
22      )
23 }

```

Listing 17. Scalar Phase

Summary: Low-Level Representation

- Expresses **(de/re)-compositions** of data-parallel computations (BLAS, stencil, ...) for the **memory** and **core hierarchies** of state-of-the-art parallel architectures (GPU, CPU, ...)
- **Straightforwardly transformable** to imperative-style, executable **program code**
- Full **formal** foundation

Definition 12 (Low-Level MDA). Let be $L \in \mathbb{N}$ (representing an ASM's number of layers) and $D \in \mathbb{N}$ (representing an MDH's number of dimensions). Let further be $P = ((P_1^1, \dots, P_D^1), \dots, (P_1^L, \dots, P_D^L)) \in \mathbb{N}^{L \times D}$ an arbitrary sequence of L -many D -tuples of positive natural numbers, $T \in \text{TYPE}$ a scalar type, and $I := ((I_d^{<p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L>} \in \text{MDA-IDX-SETS})_{d \in [1, D]_{\mathbb{N}}})^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 \mid \dots \mid (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}$ an arbitrary collection of D -many MDA index sets (Definition 1) for each particular choice of indices $(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1$, ..., $(p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L$ ¹⁴ (illustrated in Figure 17).

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

$$\alpha_{ll}^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 \mid \dots \mid (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>} : \underbrace{[0, N_1^{<p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L>}}_{\text{Partitioning: Layer 1}} \times \dots \times \underbrace{[0, N_D^{<p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L>}}_{\text{Partitioning: Layer L}} \rightarrow T$$

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

$$(\text{NUM_MEM_LYRS}, \text{NUM_COR_LYRS}) \in \mathbb{N} \times \mathbb{N}$$

for which $\text{NUM_MEM_LYRS} + \text{NUM_COR_LYRS} = L$.

...

Definition 13 (Low-Level BUF). Let be $L \in \mathbb{N}$ (representing an ASM's number of layers) and $D \in \mathbb{N}$ (representing an MDH's number of dimensions). Let further $P = ((P_1^1, \dots, P_D^1), \dots, (P_1^L, \dots, P_D^L)) \in \mathbb{N}^{L \times D}$ be an arbitrary sequence of L -many D -tuples of positive natural numbers, $T \in \text{TYPE}$ a scalar type, and $N := ((N_d^{<p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L>} \in \mathbb{N})_{d \in [1, D]_{\mathbb{N}}})^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 \mid \dots \mid (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>}$ be a BUF's size (Definition 5) for each particular choice of p_1^1, \dots, p_D^L .
An L -layered, D -dimensional, P -partitioned *low-level BUF* that has scalar type T and size N is any function b_{ll} of type (\leftrightarrow denotes bijection):

$$b_{ll}^{<(p_1^1, \dots, p_D^1) \in P_1^1 \times \dots \times P_D^1 \mid \dots \mid (p_1^L, \dots, p_D^L) \in P_1^L \times \dots \times P_D^L>} : \underbrace{[0, \text{MEM} \in [1, \text{NUM_MEM_LYRS}]_{\mathbb{N}}]}_{\text{Memory Region}} \times \underbrace{[0, D]}_{\text{Memory Layout}} \times \underbrace{[0, N_1^{<p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L>}}_{\text{Partitioning: Layer 1}} \times \dots \times \underbrace{[0, N_D^{<p_1^1, \dots, p_D^1 \mid \dots \mid p_1^L, \dots, p_D^L>}}_{\text{Partitioning: Layer L}} \rightarrow T$$

Definition 14 (Low-Level Combine Operator). Let be $L \in \mathbb{N}$ (representing an ASM's number of layers) and $D \in \mathbb{N}$ (representing an MDH's number of dimensions). Let further be $\otimes \in \text{CO}^{<M_{\text{MDA}} \mid T \mid D \mid d>}$ an arbitrary D -dimensional combine operator (Definition 2).

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

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

has the same type and semantics as \otimes :

$$\otimes^{<l_{\text{ASM}}, d_{\text{ASM}}>} \in \text{CO}^{<M_{\text{MDA}} \mid T \mid D \mid d>} , \otimes^{<l_{\text{ASM}}, d_{\text{ASM}}>}(a, b) := \otimes(a, b)$$

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

Lowering: High Level → Low-Level

Our Lowering Process is formally proved, and it is generic in performance-critical parameters:

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

Table 1. Tuning parameters of our low-level expressions

tile sizes

exploiting core hierarchy
(parallelization)

exploiting memory hierarchy
(data movements)

...

...

...

...

We use our Auto-Tuning Framework (ATF) [1] to fully automatically determine optimized values of parameters

Experimental Results

We experimentally evaluate our MDH approach in terms of *Performance & Portability & Productivity*:

Competitors:

1. Scheduling Approach:

- Apache TVM [2] (GPU & CPU)

2. Polyhedral Compilers:

- PPCG [3] (GPU)
- Pluto [4] (CPU)

3. Functional Approach:

- Lift [5] (GPU & CPU)

4. Domain-Specific Libraries:

- NVIDIA cuBLAS & cuDNN (GPU)
- Intel oneMKL & oneDNN (CPU)

Case Studies:

1. Linear Algebra Routines:

- Matrix Multiplication (MatMul)
- Matrix-Vector Multiplication (MatVec)

2. Stencil Computations:

- Jacobi Computation (Jacobi1D)
- Gaussian Convolution (Conv2D)

3. Quantum Chemistry:

- Coupled Cluster (CCSD(T))

4. Data Mining:

- Probabilistic Record Linkage (PRL)

5. Deep Learning:

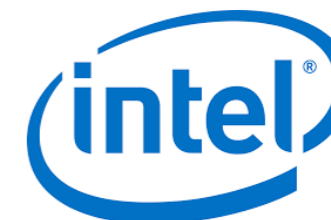
- Multi-Channel Convolution (MCC)
- Capsule-Style Convolution (MCC_Capsule)

[2] Chen et al., "TVM: An Automated End-to-End Optimizing Compiler for Deep Learning", OSDI'18

[3] Verdoolaege et al., "Polyhedral Parallel Code Generation for CUDA", TACO'13

[4] Bondhugula et al., "PLuTo: A Practical and Fully Automatic Polyhedral Program Optimization System", PLDI'08

[5] Steuwer et al., "Generating Performance Portable Code using Rewrite Rules", ICFP'15



Experimental Results

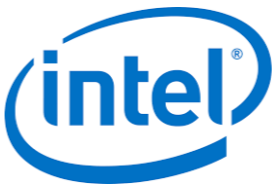
Highlights only

Performance Evaluation: (via runtime comparison)



MDH speedup over

- TVM: 0.88x – 2.22x
- PPCG: 2.58x – 13.76x
- (cuBLAS/cuDNN: 0.91x – 2.67x)



MDH speedup over

- TVM: 1.05 – 3.01x
- Pluto: 6.29x – 364.43x
- (oneMKL/oneDNN: 0.39x – 9.01x)

Deep Learning	NVIDIA Ampere GPU									
	ResNet-50				VGG-16				MobileNet	
	Training		Inference		Training		Inference		Training	Inference
	MCC	MatMul	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	-	-

Deep Learning	Intel Skylake CPU									
	ResNet-50				VGG-16				MobileNet	
	Training		Inference		Training		Inference		Training	Inference
	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	-	-

Case Study “Deep Learning” for which most competitors are highly optimized (most challenging for us!)

Significantly higher speedups for other case studies, e.g., >170x over TVM on GPU already for straightforward dot products

Experimental Results

Highlights only

Portability Evaluation: (via Pennycook Metric [6])

Deep Learning	Pennycook Metric									
	ResNet-50				VGG-16				MobileNet	
	Training		Inference		Training		Inference		Training	Inference
	MCC	MatMul	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	0.98	1.00
TVM+Ansor	0.53	0.62	0.89	0.59	0.76	0.81	0.70	0.61	0.54	0.75

Our other competitors achieve lowest portability — of “0.00” — only, because they are limited to particular architectures and/or application classes

Experimental Results

Highlights only

Productivity Evaluation: (via intuitive argumentation)

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

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

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

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

```
1 def MatVec(I, K):
2   M = te.placeholder((I, K), name='M', dtype='float32')
3   v = te.placeholder((K,), name='v', dtype='float32')
4
5   k = te.reduce_axis((0, K), name='k')
6   w = te.compute(
7     (I,),
8     lambda i: te.sum(M[i, k] * v[k], axis=k)
9   )
10  return [M, v, w]
```

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

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

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

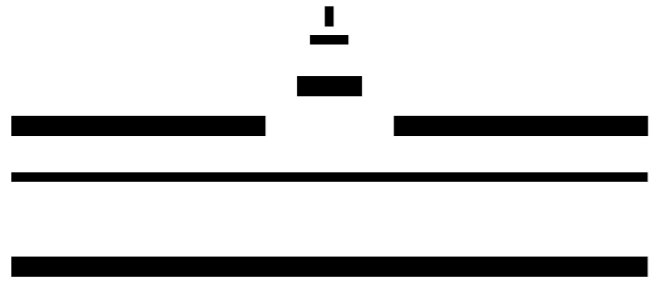
Conclusion

- Our approach **combines** together three major goals as compared to related work: ***Performance & Portability & Productivity***
- For this, we **formally introduce program representations** on both:
 - **high level**, for conveniently expressing, in one uniform formalism, the various kinds of data-parallel computations, agnostic from hardware and optimization details, while still capturing all information relevant for generating high-performance program code;
 - **low level**, which allows uniformly reasoning, in the same formalism, about optimized (de/re)-compositions of data-parallel computations targeting different kinds of architectures (GPUs, CPUs, etc).
- We **lower** our high-level representation to our low-level representation, in a **formally sound** manner, by introducing a generic search space that is based on **performance-critical parameters** and **auto-tuning**
- Our **experiments** confirm that our MDH approach often achieves higher ***Performance & Portability & Productivity*** than popular state-of-practice approaches, including hand-optimized libraries provided by vendors

Future Work

- Expressing and optimizing simultaneously **multiple data-parallel computations**, rather than optimizing computations individually and thus independently from each other only (a.k.a. *fusing optimization*).
- Supporting computations on **sparse data formats**
- Introduce an **analytical cost model** to accelerate (or even avoid) the auto-tuning overhead (possibly guided by machine learning techniques)
- Implement our approach into **MLIR** to make our work better accessible for the community
- Targeting **assembly languages** to benefit from assembly-level optimizations
- Extending our approach toward distributed multi-device systems that may be **heterogeneous**
- ...

This work provides the (formal) foundation for all our future goals!



WESTFÄLISCHE
WILHELMS-UNIVERSITÄT
MÜNSTER

Questions?

Grateful for any feedback



Ari Rasch
a.rasch@wwu.de



Richard Schulze
r.schulze@wwu.d