

Toward Performance & Portability & Productivity in Parallel Programming

**A Holistic Code *Generation*, *Optimization*, and *Execution* Approach
for Data-Parallel Computations Targeting Modern Parallel Architectures**

Thesis Defense Presentation
Ari Rasch
University of Münster, Germany
Advisor: Prof. Dr. Sergei Gorlatch
February 13, 2025

Introductory Remark

COMPUTER SCIENCE

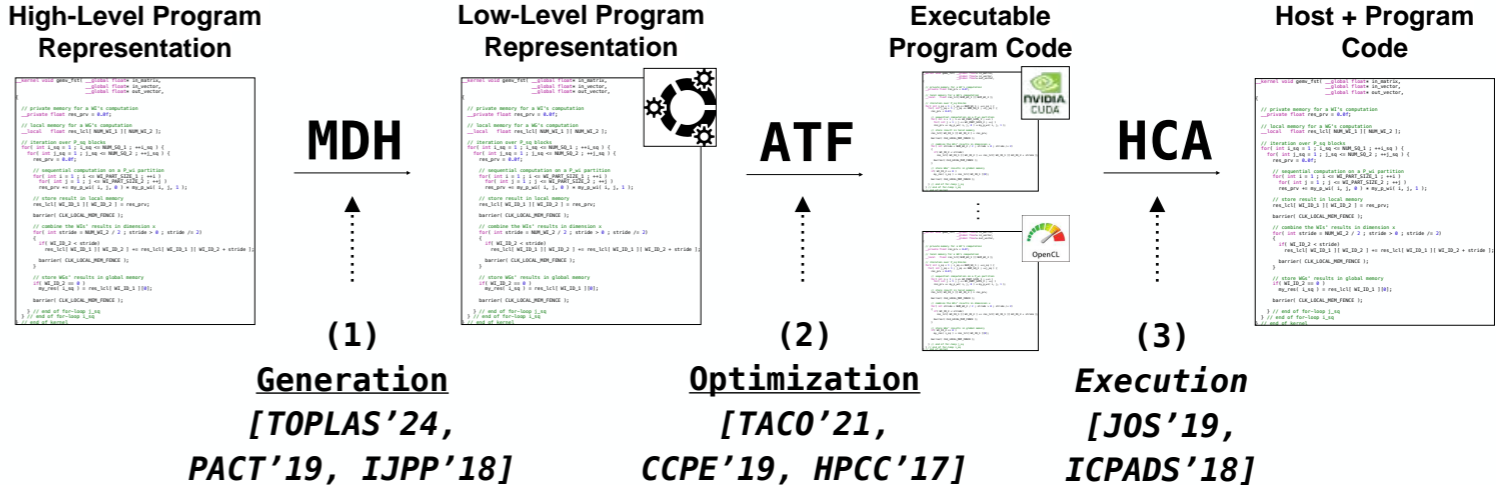
TOWARD PERFORMANCE & PORTABILITY & PRODUCTIVITY IN PARALLEL PROGRAMMING

A Holistic Code Generation, Optimization, and Execution Approach for Data-Parallel Computations Targeting Modern Parallel Architectures

Inaugural Dissertation
for the Award of a Doctoral Degree
Dr. rer. nat.
in the Field of Mathematics and Computer Science
from the Faculty of Mathematics und Natural Sciences
of the University of Münster, Germany

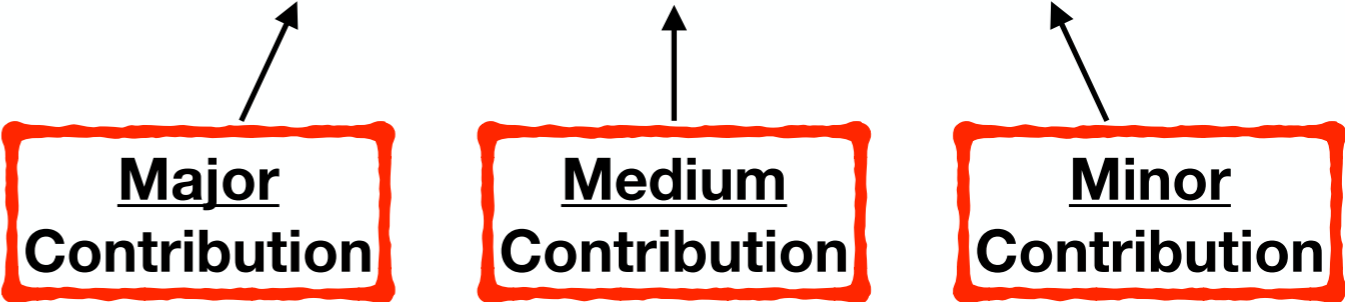
submitted by
ARI RASCH
born in Essen, Germany
– 2024 –

This thesis describes three major (sub-)projects:



The (sub-)projects complement each other to form a holistic approach to code

Generation & Optimization & Execution



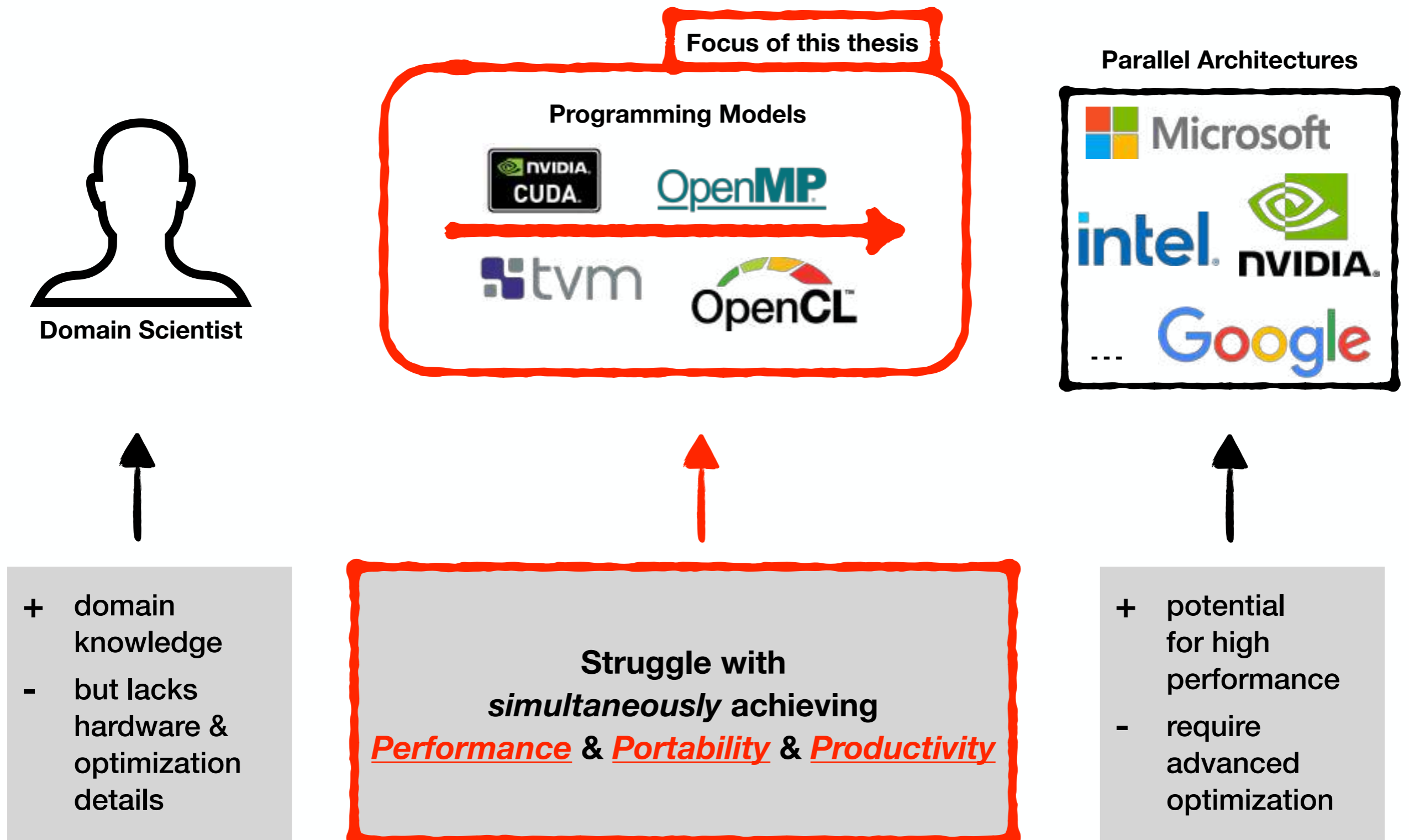
Part 1
(15 slides – 187 pages)

Part 2
(5 slides – 54 pages)

Part 3
(4 slides – 21 pages)

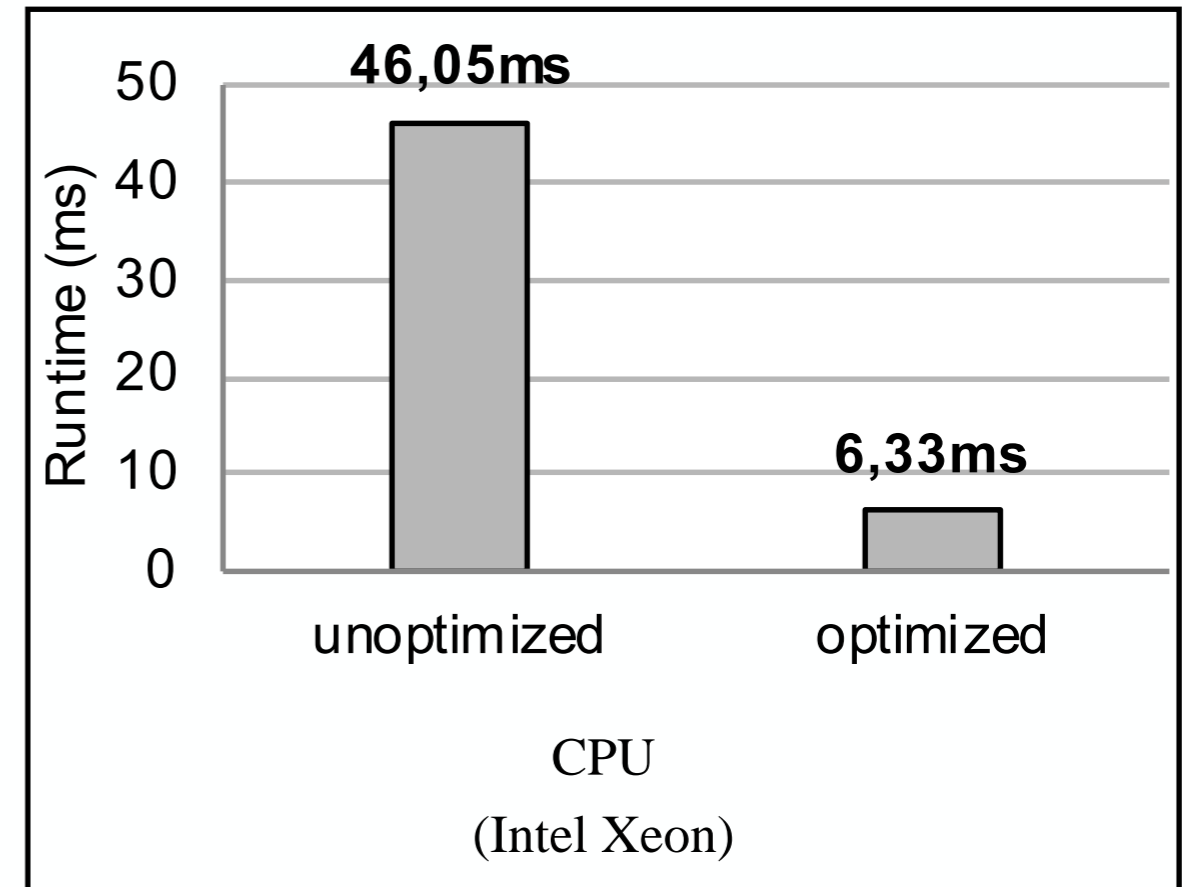
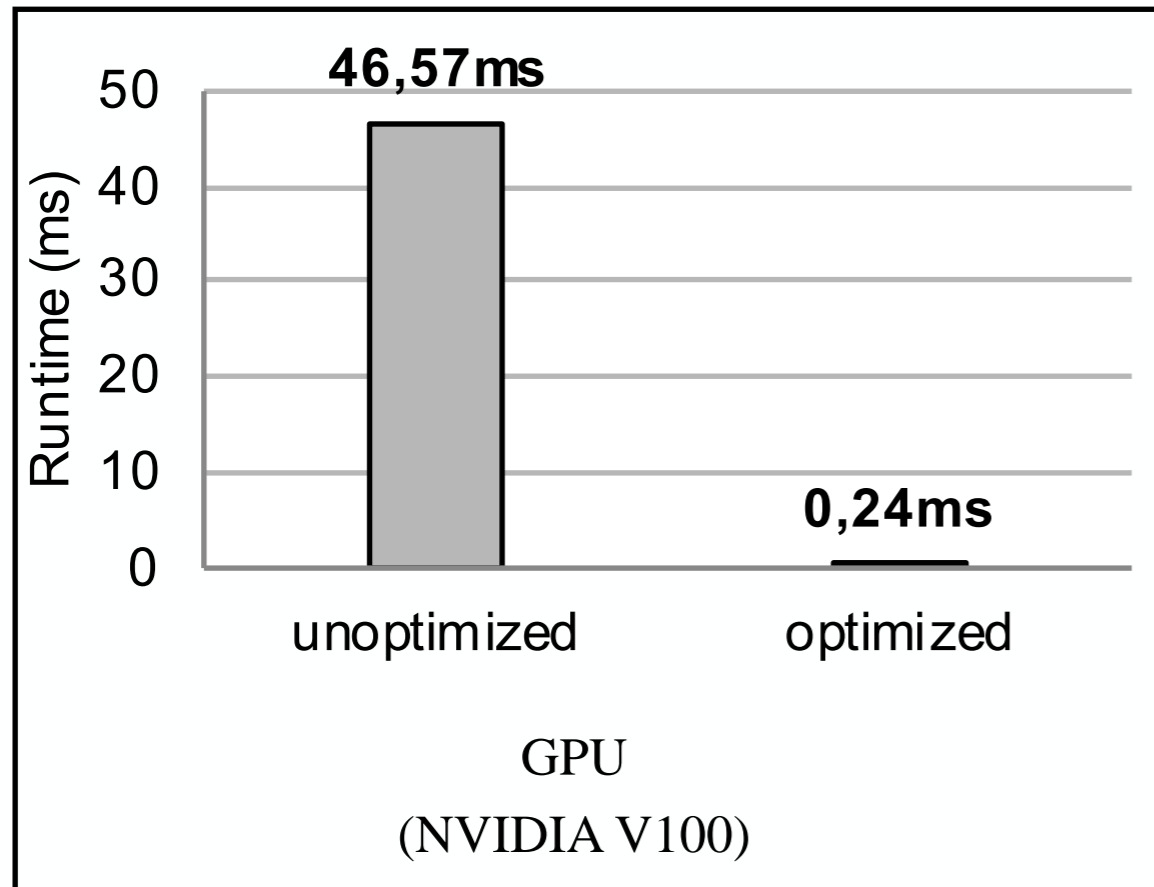
Parallel Programming in Today's World

Parallel programming is hard:



Challenges: Performance & Portability & Productivity

The Performance challenge:

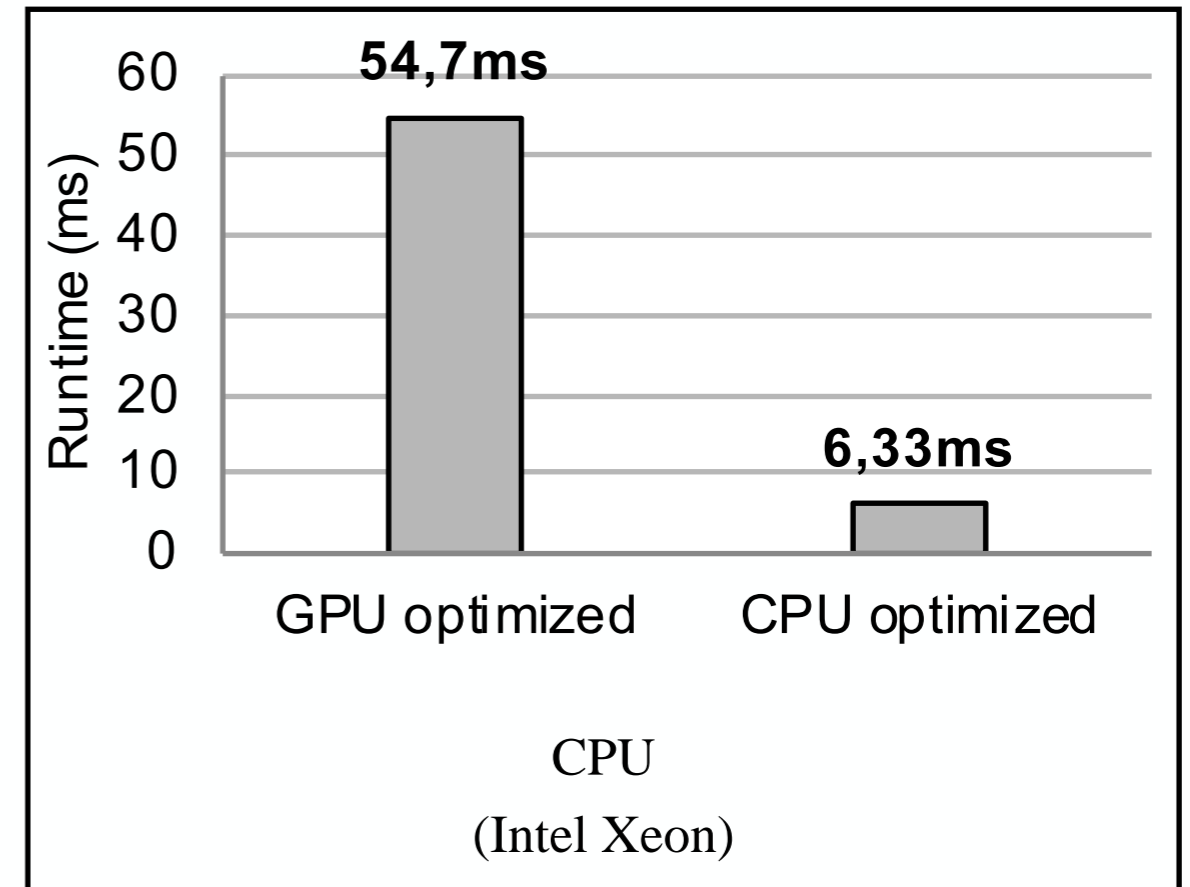
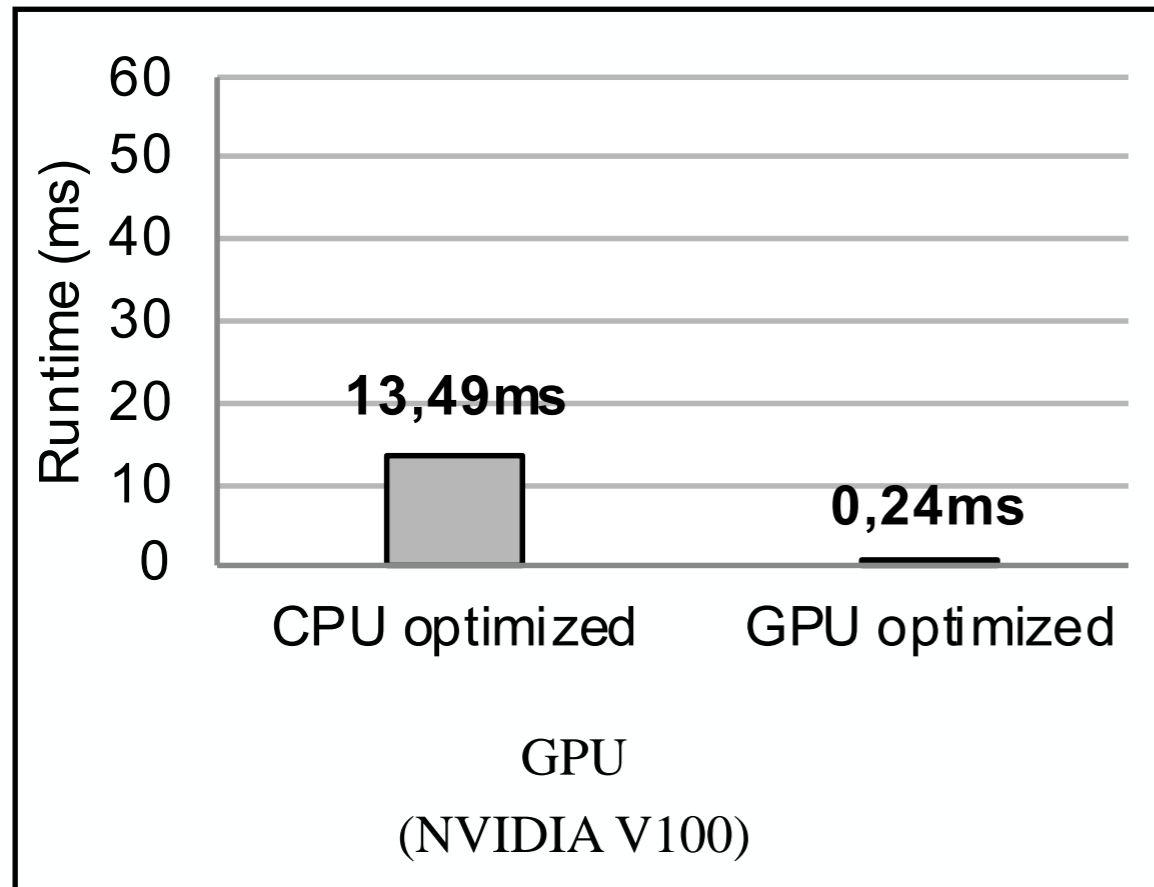


Runtime (lower is better) of unoptimized vs optimized matrix multiplication on GPU (left) and CPU (right).

High *Performance* requires *complex* optimizations

Challenges: Performance & Portability & Productivity

The *Portability* challenge:



Runtime (lower is better) of GPU/CPU-optimized matrix multiplication on GPU (left) and CPU (right).

High *Portability* requires *architecture(/data)-specific* optimizations

Challenges: Performance & Portability & Productivity

The Productivity challenge:

```
1  __kernel void MatMul( __global const float A[M][K] ,
2                        __global const float B[K][N] ,
3                        __global float C[M][N] )
4  {
5      int i = get_global_id(0);
6      int j = get_global_id(1);
7
8      for( int k=0 ; k<K ; ++k )
9          C[i][j] += A[i][k] * B[k][j];
10 }
```

Naive OpenCL implementation of matrix multiplication

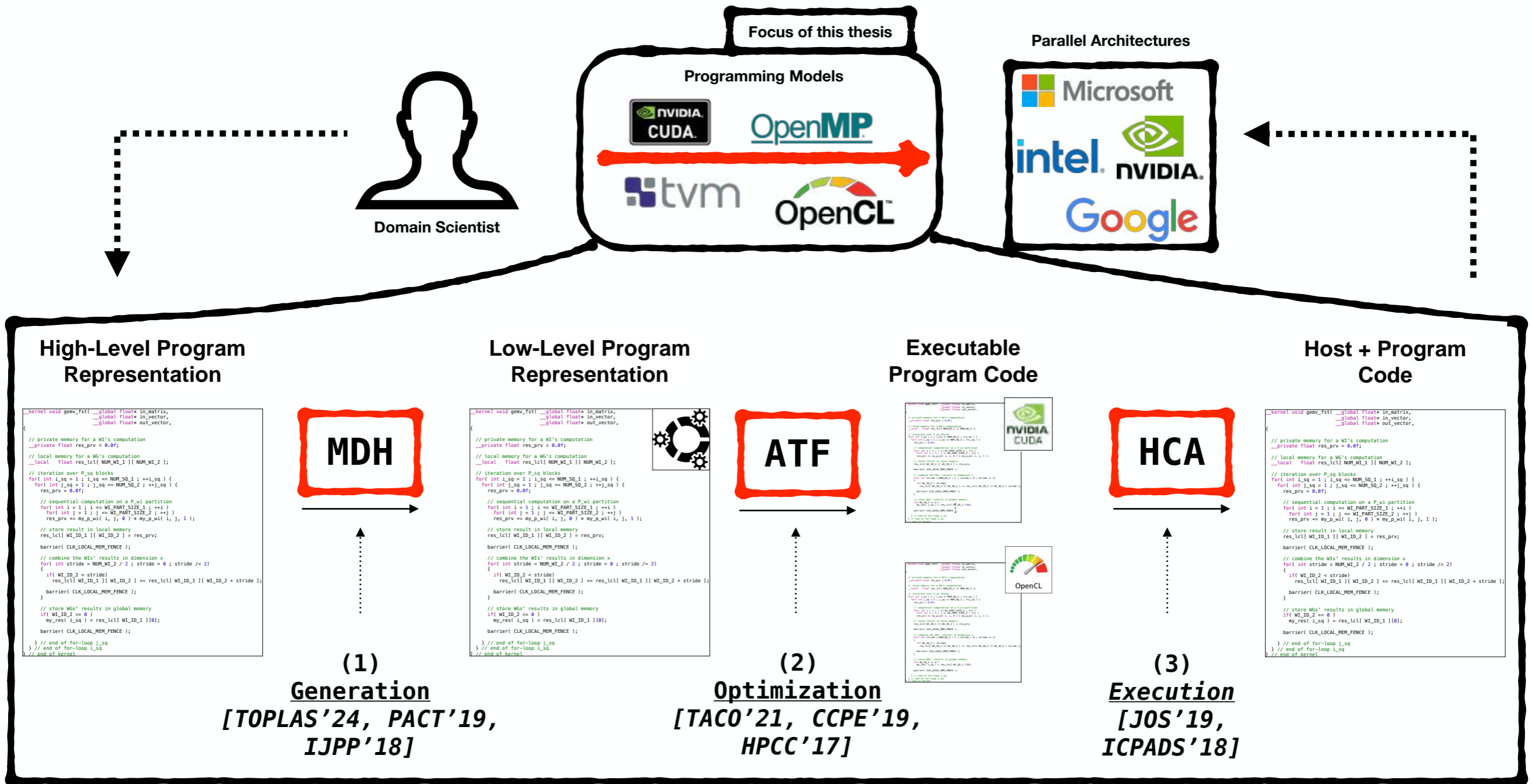
```
1  __kernel void MatMul( /* ... */ )
2  {
3      const size_t i_wg_l_1 = get_group_id(2);
4      // ... 5 lines skipped
5
6      __private TYPE_TS res_p[/*...*/][/*...*/];
7      {
8          // ... 7 lines skipped
9          for (size_t p_iteration_l_1 = 0; p_iteration_l_1 < (2);
10              ++p_iteration_l_1) {
11              for (size_t p_iteration_l_2 = 0; p_iteration_l_2 < (1)
12                  ; ++p_iteration_l_2) {
13                  size_t p_iteration_r_1 = 0;
14                  res_p[p_step_l_1][((p_iteration_l_1) * 1 + 0)][(0)][
15                      p_step_l_2][((p_iteration_l_2) * 1 + 0)] = f(
16                      a[(((l_step_l_1 * (32 / 1) + ((p_step_l_1 *
17                          (2) + ((p_iteration_l_1) * 1 + 0)) / 1) * 1 +
18                          i_wi_l_1 * 1 + (((p_iteration_l_1) * 1 +
19                          0)) % 1))) / 1) * (64 * 1) + i_wg_l_1 * 1 +
20                      (((p_step_l_1 * (2) + ((p_iteration_l_1)
21                          * 1 + 0)) / 1) * 1 + i_wi_l_1 * 1 + (((
22                          p_iteration_l_1) * 1 + 0)) % 1))) % 1))) *
23                      1024 + (((l_step_r_1 * (2 / 1) + ((
24                          p_step_r_1 * (1) + ((p_iteration_r_1) * 1 +
25                          0)) / 1) * 1 + i_wi_r_1 * 1 + (((
26                          p_iteration_r_1) * 1 + 0)) % 1))) / 1) * (2
27                          * 1) + i_wg_r_1 * 1 + (((p_step_r_1 * (1) +
28                          ((p_iteration_r_1) * 1 + 0)) / 1) * 1 +
29                          i_wi_r_1 * 1 + (((p_iteration_r_1) * 1 + 0)
30                          ) % 1))) % 1))],
31
32                  // ... 107 lines skipped
33              }
34          }
35 }
```

Optimized OpenCL implementation of matrix multiplication

High *Productivity* requires *automatic optimization*

Contributions of this Thesis

This thesis introduces a novel, holistic approach to **Generating & Optimizing & Executing** code:



The ultimate goal of **MDH+ATF+HCA** is to simultaneously achieve **Performance & Portability & Productivity**

Outline

This talk(/thesis) is structured into three main parts:

High-Level Program Representation

```
__kernel void gemv_fst( __global float* in_matrix,
                     __global float* in_vector,
                     __global float* out_vector,
                     __private float res_prv = 0.0f;
// local memory for a WI's computation
__local float res_lcl[ NUM_WI_1 ][ NUM_WI_2 ];
// iteration over P sq blocks
for( int l_sq = 1; l_sq <= NUM_SQ_1; ++l_sq ) {
  for( int j_sq = 1; j_sq <= NUM_SQ_2; ++j_sq ) {
    res_prv = 0.0f;
// sequential computation on a P-wi partition
for( int i = 1; i <= WI_PART_SIZE_1; ++i )
  for( int j = 1; j <= WI_PART_SIZE_2; ++j )
    res_prv += my_p_wi( i, j, 0 ) * my_p_wi( i, j, 1 );
// store result in local memory
res_lcl[ WI_ID_1 ][ WI_ID_2 ] = res_prv;
barrier( CLK_LOCAL_MEM_FENCE );
// combine the WIs' results in dimension x
for( int stride = NUM_WI_2 / 2; stride > 0; stride /= 2 )
  if( WI_ID_2 < stride )
    res_lcl[ WI_ID_1 ][ WI_ID_2 ] += res_lcl[ WI_ID_1 ][ WI_ID_2 + stride ];
barrier( CLK_LOCAL_MEM_FENCE );
// store WIs' results in global memory
if( WI_ID_2 == 0 )
  my_res( l_sq ) = res_lcl[ WI_ID_1 ][0];
barrier( CLK_LOCAL_MEM_FENCE );
} // end of for-loop j_sq
} // end of for-loop l_sq
// end of kernel
```

MDH

(1)

Generation

[*TOPLAS'24, PACT'19, IJPP'18*]

Low-Level Program Representation

```
__kernel void gemv_fst( __global float* in_matrix,
                     __global float* in_vector,
                     __global float* out_vector,
                     __private float res_prv = 0.0f;
// local memory for a WI's computation
__local float res_lcl[ NUM_WI_1 ][ NUM_WI_2 ];
// iteration over P sq blocks
for( int l_sq = 1; l_sq <= NUM_SQ_1; ++l_sq ) {
  for( int j_sq = 1; j_sq <= NUM_SQ_2; ++j_sq ) {
    res_prv = 0.0f;
// sequential computation on a P-wi partition
for( int i = 1; i <= WI_PART_SIZE_1; ++i )
  for( int j = 1; j <= WI_PART_SIZE_2; ++j )
    res_prv += my_p_wi( i, j, 0 ) * my_p_wi( i, j, 1 );
// store result in local memory
res_lcl[ WI_ID_1 ][ WI_ID_2 ] = res_prv;
barrier( CLK_LOCAL_MEM_FENCE );
// combine the WIs' results in dimension x
for( int stride = NUM_WI_2 / 2; stride > 0; stride /= 2 )
  if( WI_ID_2 < stride )
    res_lcl[ WI_ID_1 ][ WI_ID_2 ] += res_lcl[ WI_ID_1 ][ WI_ID_2 + stride ];
barrier( CLK_LOCAL_MEM_FENCE );
// store WIs' results in global memory
if( WI_ID_2 == 0 )
  my_res( l_sq ) = res_lcl[ WI_ID_1 ][0];
barrier( CLK_LOCAL_MEM_FENCE );
} // end of for-loop j_sq
} // end of for-loop l_sq
// end of kernel
```



ATF

(2)

Optimization

[*TACO'21, CCPE'19, HPCC'17*]

Executable Program Code

```
__kernel void gemv_fst( __global float* in_matrix,
                     __global float* in_vector,
                     __global float* out_vector,
                     __private float res_prv = 0.0f;
// local memory for a WI's computation
__local float res_lcl[ NUM_WI_1 ][ NUM_WI_2 ];
// iteration over P sq blocks
for( int l_sq = 1; l_sq <= NUM_SQ_1; ++l_sq ) {
  for( int j_sq = 1; j_sq <= NUM_SQ_2; ++j_sq ) {
    res_prv = 0.0f;
// sequential computation on a P-wi partition
for( int i = 1; i <= WI_PART_SIZE_1; ++i )
  for( int j = 1; j <= WI_PART_SIZE_2; ++j )
    res_prv += my_p_wi( i, j, 0 ) * my_p_wi( i, j, 1 );
// store result in local memory
res_lcl[ WI_ID_1 ][ WI_ID_2 ] = res_prv;
barrier( CLK_LOCAL_MEM_FENCE );
// combine the WIs' results in dimension x
for( int stride = NUM_WI_2 / 2; stride > 0; stride /= 2 )
  if( WI_ID_2 < stride )
    res_lcl[ WI_ID_1 ][ WI_ID_2 ] += res_lcl[ WI_ID_1 ][ WI_ID_2 + stride ];
barrier( CLK_LOCAL_MEM_FENCE );
// store WIs' results in global memory
if( WI_ID_2 == 0 )
  my_res( l_sq ) = res_lcl[ WI_ID_1 ][0];
barrier( CLK_LOCAL_MEM_FENCE );
} // end of for-loop j_sq
} // end of for-loop l_sq
// end of kernel
```



HCA

(3)

Execution

[*JOS'19, ICPADS'18*]

Host + Program Code

```
__kernel void gemv_fst( __global float* in_matrix,
                     __global float* in_vector,
                     __global float* out_vector,
                     __private float res_prv = 0.0f;
// local memory for a WI's computation
__local float res_lcl[ NUM_WI_1 ][ NUM_WI_2 ];
// iteration over P sq blocks
for( int l_sq = 1; l_sq <= NUM_SQ_1; ++l_sq ) {
  for( int j_sq = 1; j_sq <= NUM_SQ_2; ++j_sq ) {
    res_prv = 0.0f;
// sequential computation on a P-wi partition
for( int i = 1; i <= WI_PART_SIZE_1; ++i )
  for( int j = 1; j <= WI_PART_SIZE_2; ++j )
    res_prv += my_p_wi( i, j, 0 ) * my_p_wi( i, j, 1 );
// store result in local memory
res_lcl[ WI_ID_1 ][ WI_ID_2 ] = res_prv;
barrier( CLK_LOCAL_MEM_FENCE );
// combine the WIs' results in dimension x
for( int stride = NUM_WI_2 / 2; stride > 0; stride /= 2 )
  if( WI_ID_2 < stride )
    res_lcl[ WI_ID_1 ][ WI_ID_2 ] += res_lcl[ WI_ID_1 ][ WI_ID_2 + stride ];
barrier( CLK_LOCAL_MEM_FENCE );
// store WIs' results in global memory
if( WI_ID_2 == 0 )
  my_res( l_sq ) = res_lcl[ WI_ID_1 ][0];
barrier( CLK_LOCAL_MEM_FENCE );
} // end of for-loop j_sq
} // end of for-loop l_sq
// end of kernel
```

1. Part: How to generate automatically optimizable (auto-tunable) code?

2. Part: How to optimize (auto-tune) code?

3. Part: How to execute code on (distr.) multi-dev. systems?

→ *Interface Kinds* for MDH+ATF+HCA are outlined at the end of talk

Code Generation via MDH



Overview Getting Started Code Examples Publications Citations Contact

MDH

Multi-Dimensional Homomorphisms (MDH)
An Algebraic Approach Toward Performance & Portability & Productivity
for Data-Parallel Computations

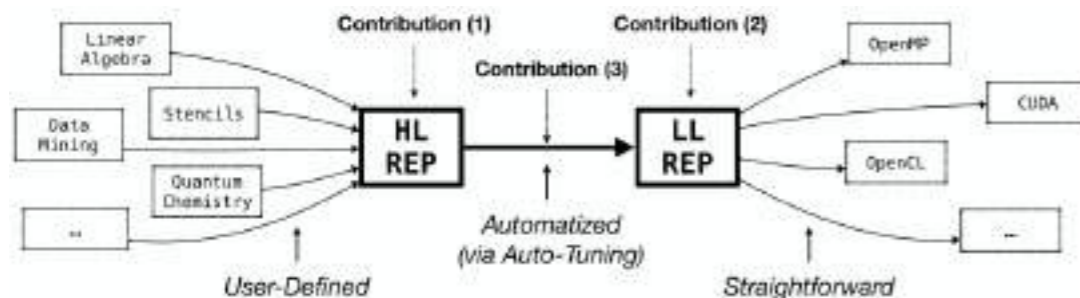
Overview

The approach of **Multi-Dimensional Homomorphisms (MDH)** is an algebraic formalism for systematically reasoning about *de-composition* and *re-composition* strategies of data-parallel computations (such as linear algebra routines and stencil computations) for the memory and core hierarchies of state-of-the-art parallel architectures (GPUs, multi-core CPU, multi-device and multi-node systems, etc).

The MDH approach (formally) introduces:

1. **High-Level Program Representation (Contribution 1)** that enables the user conveniently implementing data-parallel computations, agnostic from hardware and optimization details;
2. **Low-Level Program Representation (Contribution 2)** that expresses device- and data-optimized de- and re-composition strategies of computations;
3. **Lowering Process (Contribution 3)** that fully automatically lowers a data-parallel computation expressed in its high-level program representation to an optimized instance in its low-level representation, based on concepts from automatic performance optimization (a.k.a. *auto-tuning*), using the **Auto-Tuning Framework (ATF)**.

The MDH's low-level representation is designed such that **Code Generation** from it (e.g., in **OpenMP** for CPUs, **CUDA** for NVIDIA GPUs, or **OpenCL** for multiple kinds of architectures) becomes straightforward.



Our **Experiments** report encouraging results on GPUs and CPUs for MDH as compared to state-of-practice approaches, including NVIDIA **cuBLAS/cuDNN** and Intel **oneMKL/oneDNN** which are hand-optimized libraries provided by vendors.

ACM TOPLAS 2024

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

ARI RASCH, University of Muenster, Germany

Data-parallel computations, such as linear algebra routines and stencil computations, constitute one of the most relevant classes in parallel computing, e.g., due to their importance for deep learning. Efficiently de-composing such computations for the memory and core hierarchies of modern architectures and re-composing the computed intermediate results back to the final result—we say *(de/re)-composition* for short—is key to achieve high performance for these computations on, e.g., GPU and CPU. Current high-level approaches to generating data-parallel code are often restricted to a particular subclass of data-parallel computations and architectures (e.g., only linear algebra routines on only GPU or only stencil computations), and/or 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.

We formally introduce a systematic *(de/re)-composition* approach, based on the algebraic formalism of Multi-Dimensional Homomorphisms (MDHs). Our approach is designed as general enough to be applicable to a wide range of data-parallel computations and for various kinds of target parallel architectures. To efficiently target the deep and complex memory and core hierarchies of contemporary architectures, we exploit our introduced *(de/re)-composition* approach for a correct-by-construction, parametrized cache blocking, and parallelization strategy. We show that our approach is powerful enough to express, in the same formalism, the *(de/re)-composition* strategies of different classes of state-of-the-art approaches (scheduling-based, polyhedral, etc.), and we demonstrate that the parameters of our strategies enable systematically generating code that can be fully automatically optimized (auto-tuned) for the particular target architecture and characteristics of the input and output data (e.g., their sizes and memory layouts). Particularly, our experiments confirm that via auto-tuning, we achieve higher performance than state-of-the-art approaches, including hand-optimized solutions provided by vendors (such as NVIDIA cuBLAS/cuDNN and Intel oneMKL/oneDNN), on real-world datasets and for a variety of data-parallel computations, including linear algebra routines, stencil and quantum chemistry computations, data mining algorithms, and computations that recently gained high attention due to their relevance for deep learning.

CCS Concepts: • **Computing methodologies** → **Parallel computing methodologies**; *Machine learning*; • **Theory of computation** → **Program semantics**; • **Software and its engineering** → **Compilers**;

Additional Key Words and Phrases: Code generation, data parallelism, auto-tuning, GPU, CPU, OpenMP, CUDA, OpenCL, linear algebra, stencils computation, quantum chemistry, data mining, deep learning

A full version of this article is provided by Rasch [2024], which presents our novel concepts with all of their formal details. In contrast to the full version, this article relies on a simplified formal foundation for better illustration and easier understanding. We often refer the interested reader to Rasch [2024] for formal details that should not be required for understanding the basic ideas and concepts of our approach.

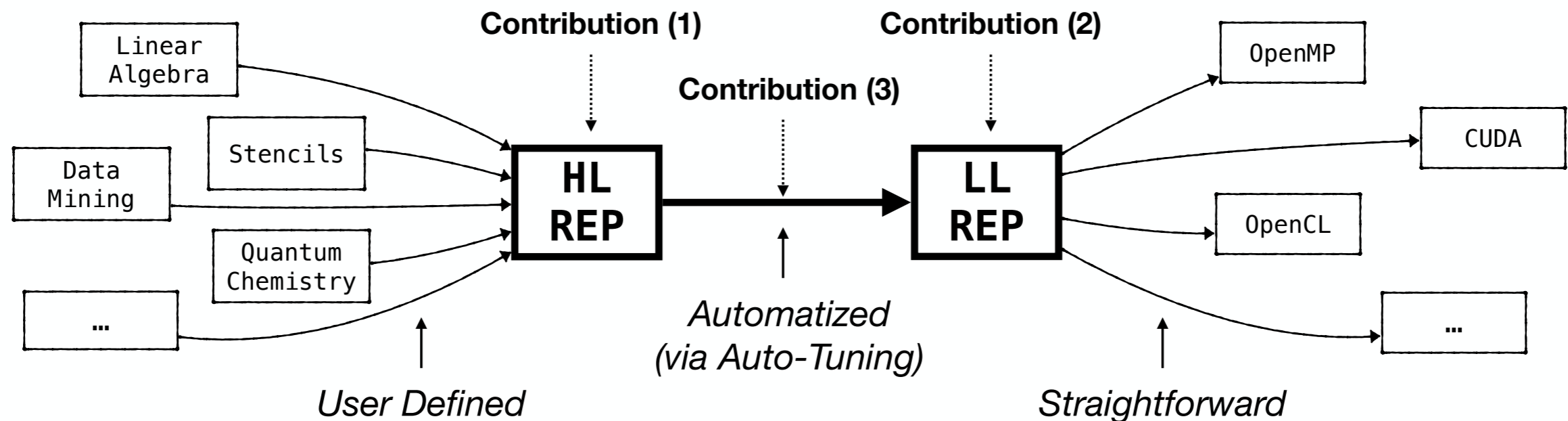
This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)—project PPP-DL (470527619).

Author's Contact Information: Ari Rasch (Corresponding author), University of Muenster, Muenster, Germany; e-mail: a.rasch@uni-muenster.de.

<https://mdh-lang.org>

Goal of MDH

MDH is a (formal) framework for *expressing & optimizing* data-parallel computations:



1. **Contribution 1 (HL-REP):** defines *data parallelism*, based on *common algebraic properties of computations* & introduces *higher-order functions* for expressing these computations, agnostic from hardware and optimization details while still capturing high-level information relevant for generating high-performing code
2. **Contribution 2 (LL-REP):** allows *expressing and reasoning about optimizations* for the memory and core hierarchies of contemporary parallel architectures & generalizes these optimizations to apply to arbitrary combinations of data-parallel computations and architectures
3. **Contribution 3 (→):** introduces a *structured optimization process* — for arbitrary combinations of data-parallel computations and parallel architectures — that allows *fully automatic optimization* (auto-tuning)

MDH: High-Level Representation

Goals:

1. **Uniform:**

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

2. **Minimalistic:**

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

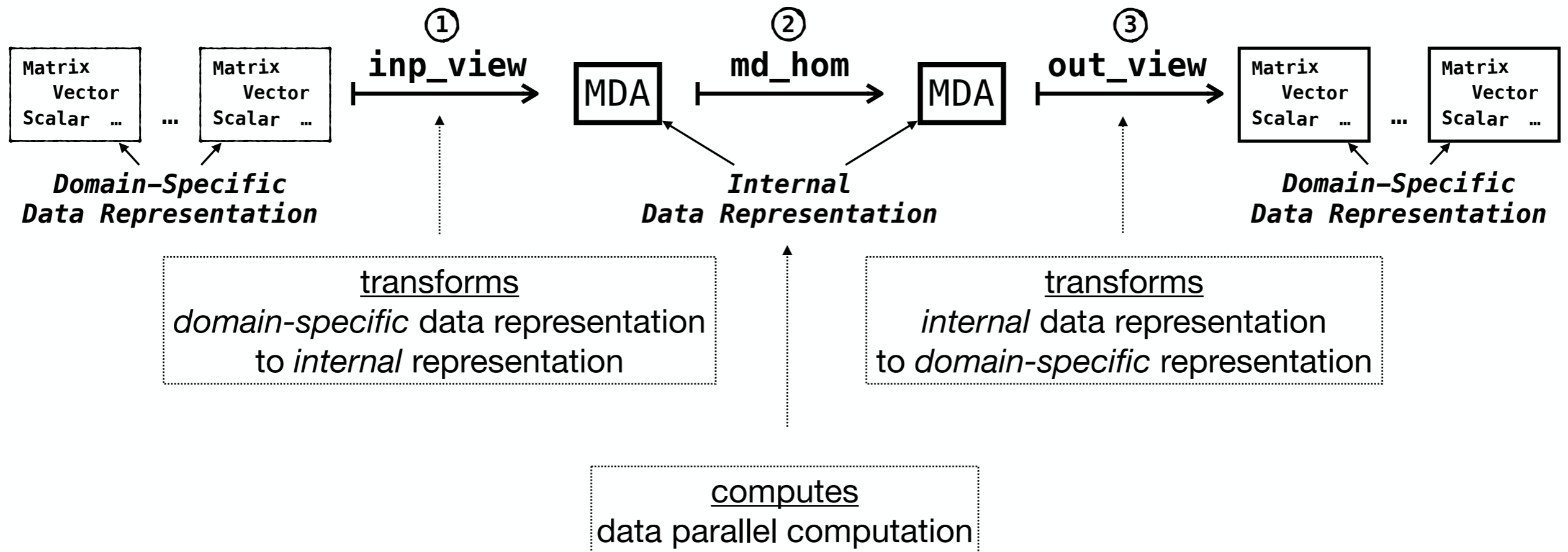
3. **Structured:**

avoiding compositions and nestings of building blocks as much as possible, thereby further contributing to the 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) )
```

***MDH High-Level Representation of MatVec
(discussed later)***

MDH: High-Level Representation



Our high-level representation defines data-parallel computations as Multi-Dimensional Homomorphisms (MDH) and it expresses data-parallel computations using exactly three straightforwardly composed higher-order functions only

MDH: High-Level Representation

Example: MatVec expressed in MDH

$$\text{MatVec}^{\langle T \in \text{TYPE} \mid I, K \in \mathbb{N} \rangle} := \text{out_view} \langle T \rangle (w : (i, k) \mapsto (i)) \circ \text{md_hom} \langle I, K \rangle (*, (\#, +)) \circ \text{inp_view} \langle T, T \rangle (M : (i, k) \mapsto (i, k) , v : (i, k) \mapsto (k))$$

MDH

MDH 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

```
void MatVec( T[] M, T[] v, T[] w )
{
    for( int i=0 ; i < I ; ++i )
        for( int k=0 ; k < K ; ++k )
            w[i] += M[i][k] * v[k];
}
```

MatVec in C++



¹We can generate such MDH expressions also automatically from straightforward (annotated) code in Python or C

MDH: High-Level Representation

md_hom	f	\otimes_1	\otimes_2	\otimes_3	\otimes_4
Dot	*	+			
MatVec	*	++	+		
MatMul	*	++	++	+	
MatMul ^T	*	++	++	+	
bMatMul	*	++	++	++	+

Views	inp_view		out_view
	A	B	C
Dot	$(k) \mapsto (k)$	$(k) \mapsto (k)$	$(k) \mapsto ()$
MatVec	$(i,k) \mapsto (i,k)$	$(i,k) \mapsto (k)$	$(i,k) \mapsto (i)$
MatMul	$(i,j,k) \mapsto (i,k)$	$(i,j,k) \mapsto (k,j)$	$(i,j,k) \mapsto (i,j)$
MatMul ^T	$(i,j,k) \mapsto (k,i)$	$(i,j,k) \mapsto (j,k)$	$(i,j,k) \mapsto (j,i)$
bMatMul	$(b,i,j,k) \mapsto (b,i,k)$	$(b,i,j,k) \mapsto (b,k,j)$	$(b,i,j,k) \mapsto (b,i,j)$

md_hom	f	\otimes_1	\otimes_2	\otimes_3	\otimes_4	\otimes_5	\otimes_6	\otimes_7	\otimes_8	\otimes_9	\otimes_{10}
Conv2D	*	++	++	+	+						
MCC	*	++	++	++	++	+	+	+			
MCC_Capsule	*	++	++	++	++	+	+	+	++	++	+

1) Linear Algebra Routines

Views	inp_view		out_view
	I	F	O
Conv2D	$(p,q,r,s) \mapsto (p+r,q+s)$	$(p,q,r,s) \mapsto (r,s)$	$(p,q,r,s) \mapsto (p,q)$
MCC	$(n,p,\dots) \mapsto (n,p+r,q+s,c)$	$(n,p,\dots) \mapsto (k,r,s,c)$	$(n,p,\dots) \mapsto (n,p,q,k)$
MCC_Capsule	$(n,p,\dots) \mapsto (n,p+r,q+s,c,mi,mk)$	$(n,p,\dots) \mapsto (k,r,s,c,mk,mj)$	$(n,p,\dots) \mapsto (n,p,q,k,mi,mj)$

2) Convolution Stencils

md_hom	f	\otimes_1	...	\otimes_6	\otimes_7
CCSD(T)	*	++	...	++	+

Views	inp_view		out_view
	A	B	C
I1	$(a,\dots,g) \mapsto (g,d,a,b)$	$(a,\dots,g) \mapsto (e,f,g,c)$	$(a,\dots,g) \mapsto (a,\dots,f)$
I2	$(a,\dots,g) \mapsto (g,d,a,c)$	$(a,\dots,g) \mapsto (e,f,g,b)$	$(a,\dots,g) \mapsto (a,\dots,f)$

3) Quantum Chemistry

md_hom	f	\otimes_1	\otimes_2
PRL	wght	++	maxPRL

Views	inp_view		out_view
	N	E	M
PRL	$(i,j) \mapsto (i)$	$(i,j) \mapsto (j)$	$(i,j) \mapsto (i)$

5) Probabilistic Record Linkage

md_hom	f	\otimes_1
map(f)	f	++
reduce(\oplus)	id	\oplus
reduce(\oplus, \otimes)	$(x) \mapsto (x,x)$	(\oplus, \otimes)

Views	inp_view	out_view	
	I	O ₁	O ₂
map(f)	$(i) \mapsto (i)$	$(i) \mapsto (i)$	
reduce(\oplus)	$(i) \mapsto (i)$	$(i) \mapsto ()$	
reduce(\oplus, \otimes)	$(i) \mapsto (i)$	$(i) \mapsto ()$	$(i) \mapsto ()$

7) Map/Reduce Patterns

md_hom	f	\otimes_1	\otimes_2
Histo	f _{Histo}	+	++
GenHisto	f	\oplus	++

Views	inp_view		out_view
	Elms	Bins	Out
Histo	$(e,b) \mapsto (e)$	$(e,b) \mapsto (b)$	$(e,b) \mapsto (b)$
GenHisto	$(e,b) \mapsto (e)$	$(e,b) \mapsto (b)$	$(e,b) \mapsto (b)$

6) Histogram

md_hom	f	\otimes_1	\otimes_2
scan(\oplus)	id	++ _{prefix-sum(\oplus)}	
MBBS	id	++ _{prefix-sum(+)}	+

Views	inp_view	out_view
	A	Out
scan(\oplus)	$(i) \mapsto (i)$	$(i) \mapsto (i)$
MBBS	$(i,j) \mapsto (i,j)$	$(i,j) \mapsto (i)$

8) Prefix Sum Computations

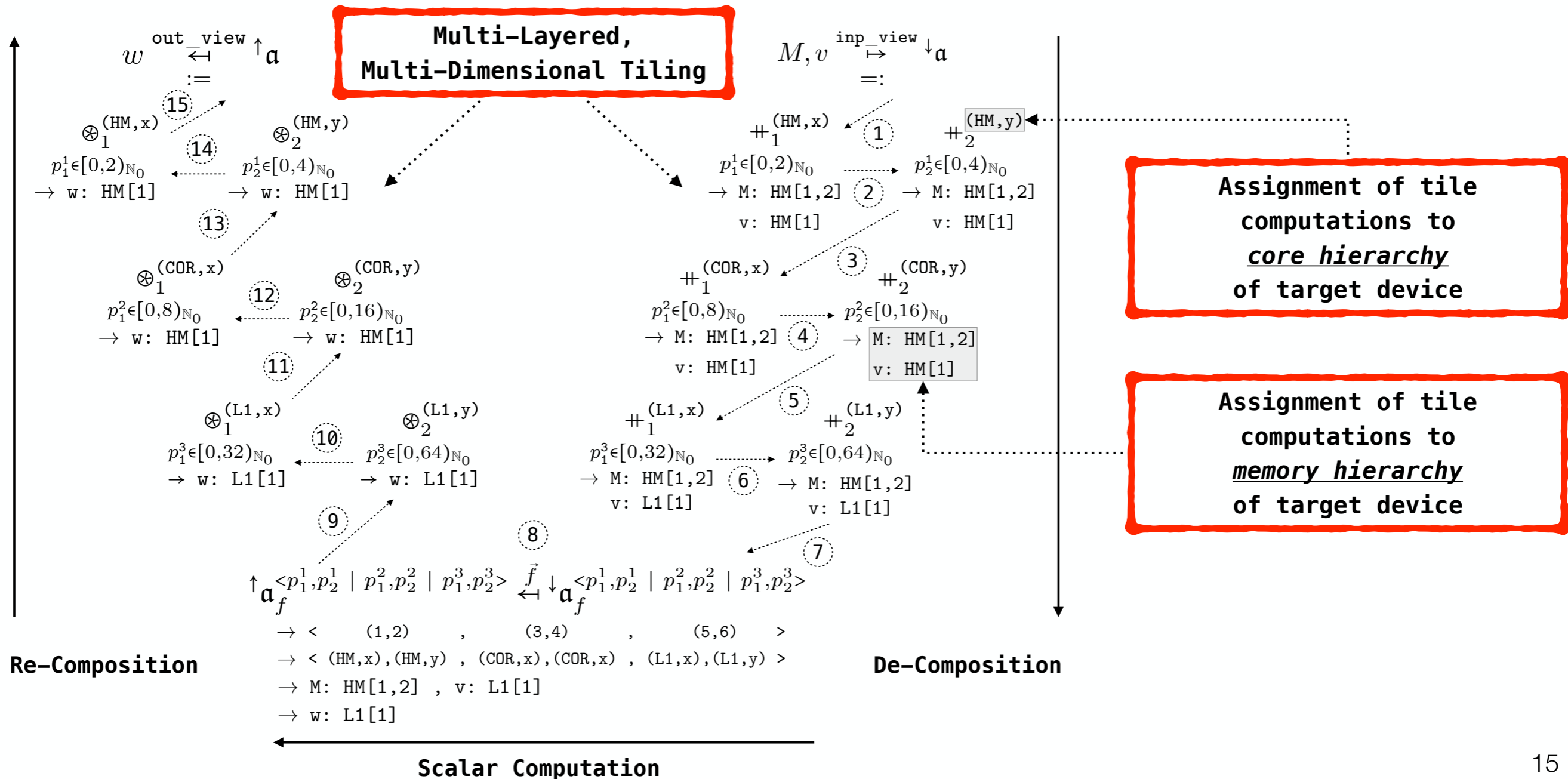
The MDH high-level representation is capable of expressing various kinds of data-parallel computations (with significantly different characteristics)

MDH: Low-Level Representation

Basic Idea

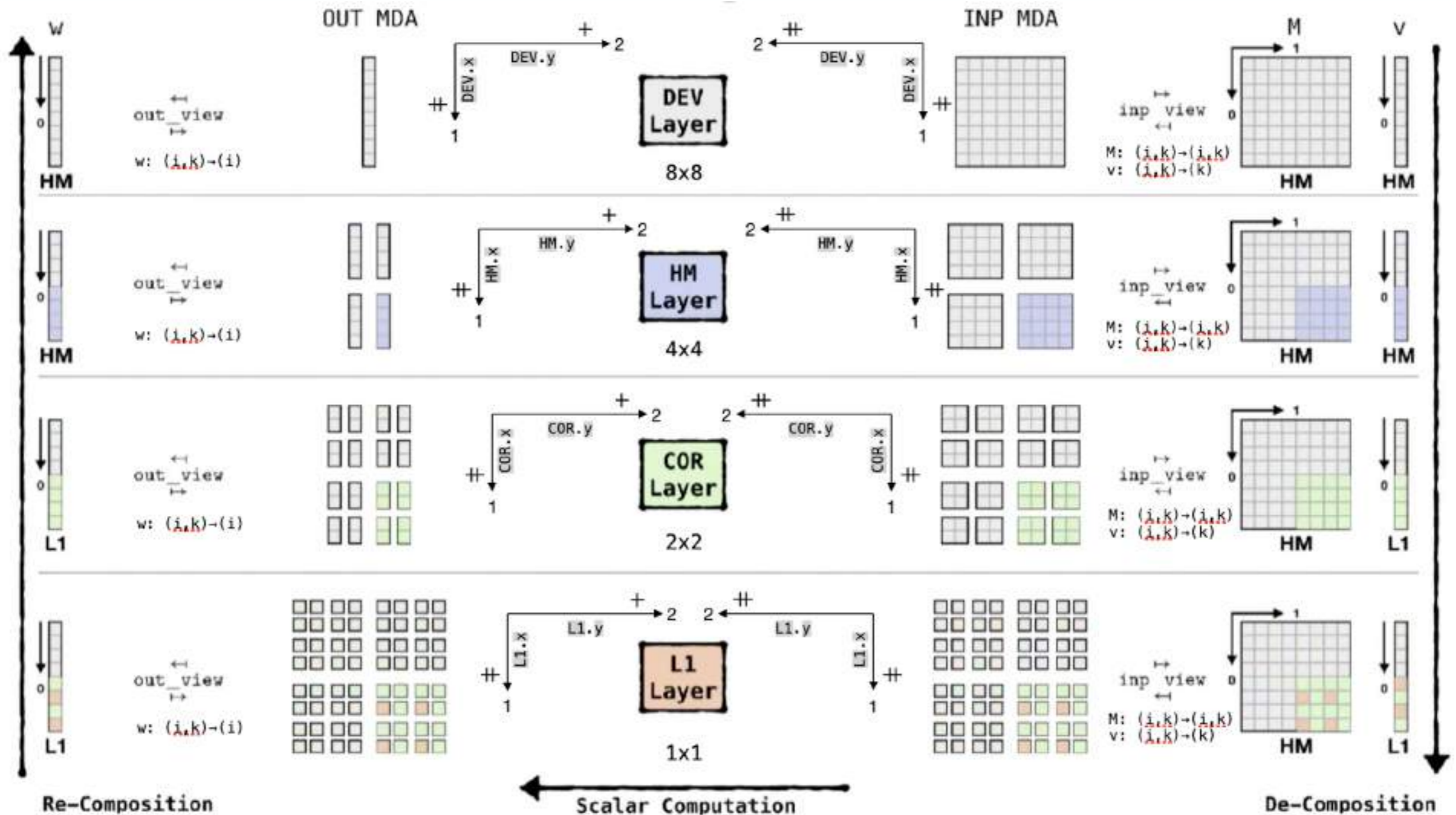
Goals:

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



MDH: Low-Level Representation

Excursion: visualizing low-level instances



The (structured) design of our Low-Level Representation allows uniformly visualizing optimizations

MDH: Lowering: *High Level* → *Low-Level*

Based on (formally defined) performance-critical parameters, for a *structured* optimization process:

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

*exploiting core hierarchy
(parallelization)*

*exploiting memory hierarchy
(data movements)*

**Our parameters
unify & generalize & combine
well-proven optimizations
(e.g., tiling, data movements,
and parallelization)**

We use our **Auto-Tuning Framework (ATF)** to automatically determine optimized values of parameters

MDH: Experimental Results

MDH is experimentally evaluated in terms of ***Performance & Portability & Productivity***:

Competitors:

1. Scheduling Approach:
 - Apache TVM [1] (GPU & CPU)
2. Polyhedral Compilers:
 - PPCG [2] (GPU)
 - Pluto [3] (CPU)
3. Functional Approach:
 - Lift [4] (GPU & CPU)
4. Domain-Specific Libraries:
 - NVIDIA cuBLAS & cuDNN (GPU)
 - Intel oneMKL & oneDNN (CPU)

[1] Chen et al., “TVM: An Automated End-to-End Optimizing Compiler for Deep Learning”, OSDI’18

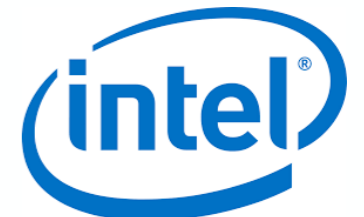
[2] Verdoolaege et al., “Polyhedral Parallel Code Generation for CUDA”, TACO’13

[3] Bondhugula et al., “PLuTo: A Practical and Fully Automatic Polyhedral Program Optimization System”, PLDI’08

[4] Steuwer et al., “Generating Performance Portable Code using Rewrite Rules”, ICFP’15

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)



MDH: Experimental Results

Performance Evaluation: (via runtime comparison)

Highlights only

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	-	-

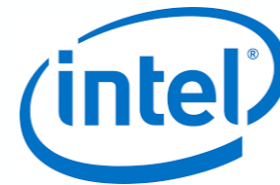


MDH speedup over

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

NVCC vs NVRTC

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	-	-



MDH speedup over

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

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

MDH: 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

The other related approaches achieve lowest portability — of “0.00” — only, because they are designed for particular architectures and/or application classes only

MDH: 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

```
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)

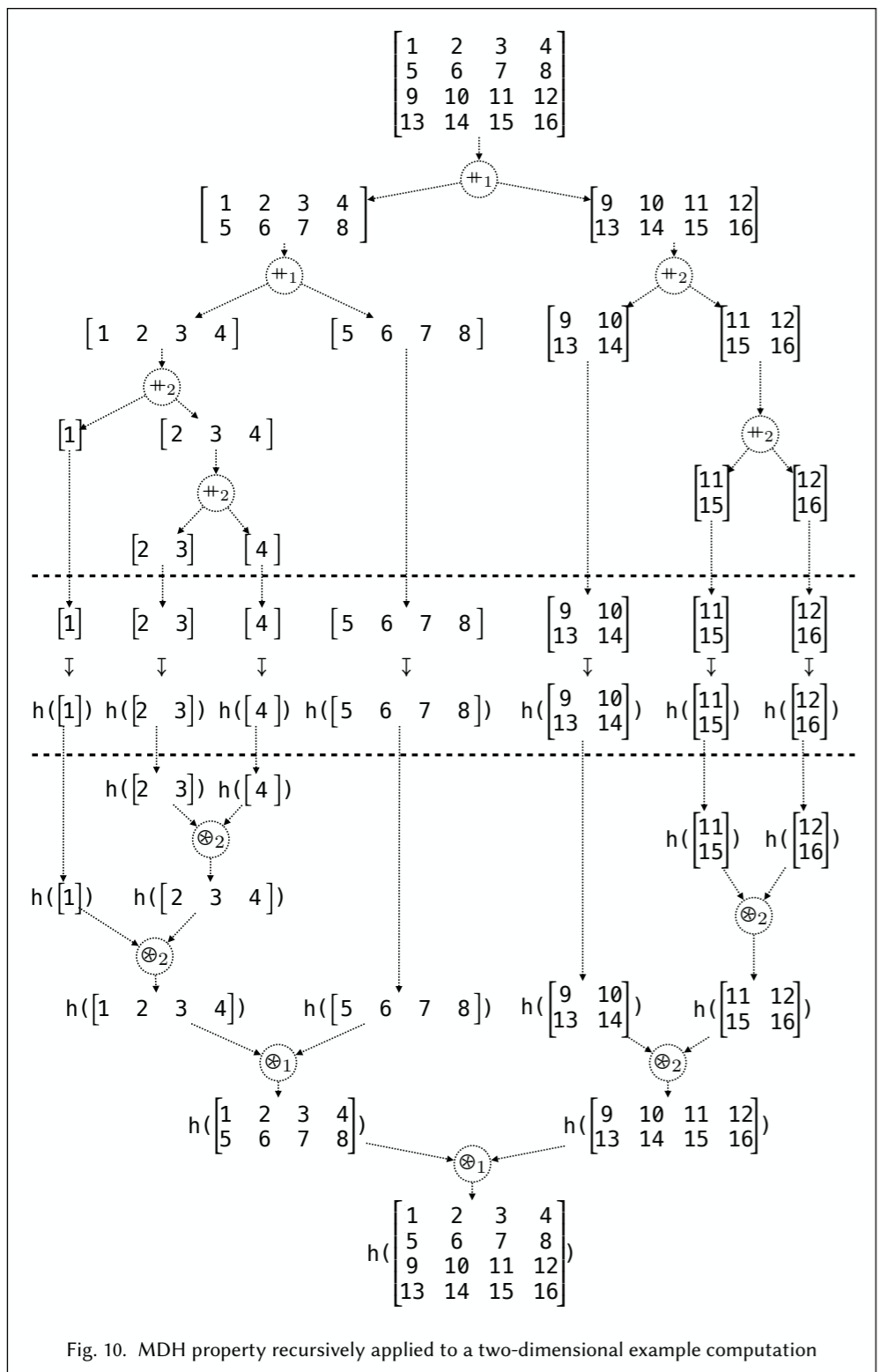
¹MDH can also take (annotated) C code as input [IMPACT'20]

MDH: Summary

- MDH combines three key goals — ***Performance & Portability & Productivity*** — as compared to related approaches
- For this, MDH **formally introduces 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 for the memory and core hierarchies of contemporary parallel architectures (GPUs, CPUs, etc)
- MDH **lowers** instances in its high-level representation to device- and data-optimized instances in its low-level representation, in a **formally sound** manner, by introducing a generic search space that is based on **performance-critical parameters & auto-tuning**
- Our **experiments** confirm that MDH often achieves higher ***Performance & Portability & Productivity*** than popular state-of-practice approaches, including hand-optimized libraries provided by vendors

MDH: Summary

That was a very quick, informal dive!



Definition 3 (Multi-Dimensional Homomorphism). Let $T^{INP}, T^{OUT} \in \text{TYPE}$ be two arbitrary scalar types, $D \in \mathbb{N}$ a natural number, and $\Rightarrow_{MDA}^1, \dots, \Rightarrow_{MDA}^D : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$ functions on MDA index sets. Let further $\#_d := \#^{<T^{INP} | D | d>} \in \text{CO}^{<id | T^{INP} | D | d>}$ denote concatenation (Definition 1) in dimension $d \in [1, D]_{\mathbb{N}}$ on D -dimensional MDAs that have scalar type T^{INP} . A function

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

is a *Multi-Dimensional Homomorphism (MDH)* that has *input scalar type* T^{INP} , *output scalar type* T^{OUT} , *dimensionality* D , and *index set functions* $\Rightarrow_{MDA}^1, \dots, \Rightarrow_{MDA}^D$, iff for each $d \in [1, D]_{\mathbb{N}}$, there exists a combine operator $\otimes_d \in \text{CO}^{<\Rightarrow_{MDA}^d | T^{OUT} | D | d>}$ (Definition 2), such that for any concatenated input MDA $a_1 \#_d a_2$ in dimension d , the *homomorphic property* is satisfied:

$$h(a_1 \#_d a_2) = h(a_1) \otimes_d h(a_2)$$

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

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 $\Rightarrow_{MDA} : \text{MDA-IDX-SETS} \rightarrow \text{MDA-IDX-SETS}$ 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)

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

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

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

... (~140 pages)

All concepts are fully formally defined in the thesis

Code Optimization via ATF



Overview Getting Started Code Examples Publications Citations Contact



Auto-Tuning Framework (ATF)

Efficient Auto-Tuning of Parallel Programs with Constrained Tuning Parameters

Overview

The **Auto-Tuning Framework (ATF)** is a general-purpose auto-tuning approach: given a program that is implemented as generic in performance-critical program parameters (a.k.a. *tuning parameters*), such as sizes of tiles and numbers of threads, ATF fully automatically determines a hardware- and data-optimized configuration of such parameters.

Key Feature of ATF

A key feature of ATF is its support for **Tuning Parameter Constraints**. Parameter constraints allow auto-tuning programs whose tuning parameters have so-called *interdependencies* among them, e.g., the value of one tuning parameter has to evenly divide the value of another tuning parameter.

ATF's support for parameter constraints is important: modern parallel programs target novel parallel architectures, and such architectures typically have deep memory and core hierarchies thus requiring constraints on tuning parameters, e.g., the value of a tile size tuning parameter on an upper memory layer has to be a multiple of a tile size value on a lower memory layer.

For such parameters, ATF introduces novel concepts for **Generating & Storing & Exploring** the search spaces of constrained tuning parameters, thereby contributing to a substantially more efficient overall auto-tuning process for such parameters, as confirmed in our **Experiments**.

Generality of ATF

For wide applicability, ATF is designed as generic in:

1. The target program's **Programming Language**, e.g., C/C++, CUDA, OpenMP, or OpenCL. ATF offers *pre-implemented* cost functions for conveniently auto-tuning C/C++ programs, as well as CUDA and OpenCL kernels which require host code for their execution which is automatically generated and executed by ATF's pre-implemented CUDA and OpenCL cost functions. ATF also offers a pre-implemented generic cost function that can be used for conveniently auto-tuning programs in any other programming language different from C/C++, CUDA, and OpenCL.

<https://atf-tuner.org>

ACM TACO 2021

Efficient Auto-Tuning of Parallel Programs with Interdependent Tuning Parameters via Auto-Tuning Framework (ATF)

ARI RASCH and RICHARD SCHULZE, University of Muenster, Germany
MICHEL STEUWER, University of Edinburgh, United Kingdom
SERGEI GORLATCH, University of Muenster, Germany

Auto-tuning is a popular approach to program optimization: it automatically finds good configurations of a program's so-called tuning parameters whose values are crucial for achieving high performance for a particular parallel architecture and characteristics of input/output data. We present three new contributions of the Auto-Tuning Framework (ATF), which enable a key advantage in *general-purpose auto-tuning*: efficiently optimizing programs whose tuning parameters have *interdependencies* among them. We make the following contributions to the three main phases of general-purpose auto-tuning: (1) ATF *generates* the search space of interdependent tuning parameters with high performance by efficiently exploiting parameter constraints; (2) ATF *stores* such search spaces efficiently in memory, based on a novel chain-of-trees search space structure; (3) ATF *explores* these search spaces faster, by employing a multi-dimensional search strategy on its chain-of-trees search space representation. Our experiments demonstrate that, compared to the state-of-the-art, general-purpose auto-tuning frameworks, ATF substantially improves generating, storing, and exploring the search space of interdependent tuning parameters, thereby enabling an efficient overall auto-tuning process for important applications from popular domains, including stencil computations, linear algebra routines, quantum chemistry computations, and data mining algorithms.

CCS Concepts: • **General and reference** → **Performance**; • **Computer systems organization** → **Parallel architectures**; • **Software and its engineering** → **Parallel programming languages**;

Additional Key Words and Phrases: Auto-tuning, parallel programs, interdependent tuning parameters

ACM Reference format:

Ari Rasch, Richard Schulze, Michel Steuwer, and Sergei Gorlatch. 2021. Efficient Auto-Tuning of Parallel Programs with Interdependent Tuning Parameters via Auto-Tuning Framework (ATF). *ACM Trans. Archit. Code Optim.* 18, 1, Article 1 (January 2021), 26 pages.
<https://doi.org/10.1145/3427093>

This is a new paper, not an extension of a conference paper.

Authors' addresses: A. Rasch, R. Schulze, and S. Gorlatch, University of Muenster, Muenster, Germany; emails: {a.rasch, r.schulze, gorlatch}@uni-muenster.de; M. Steuwer, University of Edinburgh, Edinburgh, United Kingdom; email: michel.steuwer@glasgow.ac.uk.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2021 Association for Computing Machinery.

1544-3566/2021/01-ART1

Goal of ATF

Advantage of Auto-Tuning Framework (ATF) over state-of-the-art general-purpose AT approaches:

ATF finds values of performance-critical parameters with

interdependencies among them

via optimized processes to

generating & storing & exploring

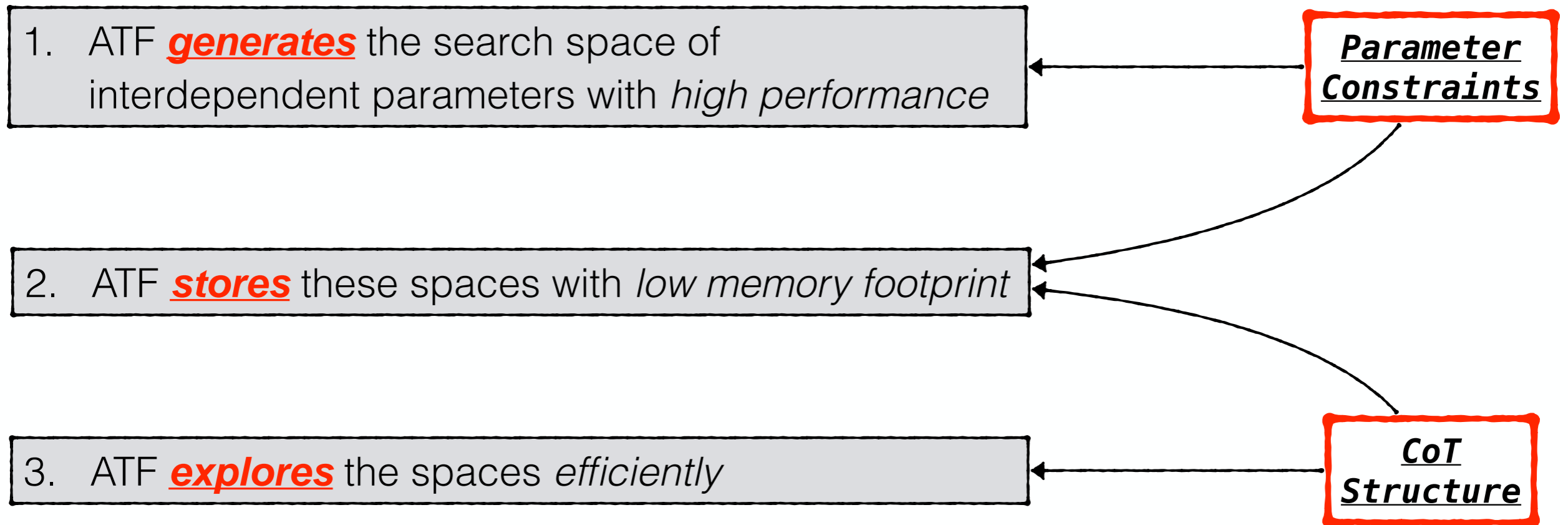
the spaces of interdependent parameters

→ We illustrate ATF by comparing it to MIT's *OpenTuner* [PACT'14] & *CLTune* [MCSoc'15] which is the foundation of many related approaches (e.g., KernelTuner & KTT).

Note: BaCO [ASPLOS'23] & KTT [FGCS'20] recently adopted the ATF techniques to efficiently handle interdependencies among tuning parameters.

ATF: Contributions

The three major contributions of ATF:



ATF introduces novel processes to ***generating*** & ***storing*** & ***exploring*** the spaces of interdependent tuning parameters, based on a novel ***constraint design*** and ***search space structure***

ATF: Internals in a Nutshell

The three major contributions of *Auto-Tuning Framework (ATF)*:

Based on
Traditional Constraints

```

for ( v1 : r1 )
  ⋮
  for ( vk : rk )
    if( sc(v1, ..., vk) )
      add_config( v1, ..., vk );
    
```

```

parallel_for ( G : {G1, ..., Gn} )
{
  parallel_for ( v1G : r1G )
    if( pc1G(v1G) )
      ⋮
      parallel_for ( vtgG : rtgG )
        if( pctgG(vtgG) )
          for ( vtg+1G : rtg+1G )
            if( pc(vtg+1G) )
              ⋮
              for( vkG : rkG )
                if( pc(vkG) )
                  add_config( v1G, ..., vkG );
            }
        }
    }
}
    
```

Based on
ATF's Parameter Constraints

Fast Space Generation:

- Break loop nest
- Parallelize loop nest
- Skip loops iterations

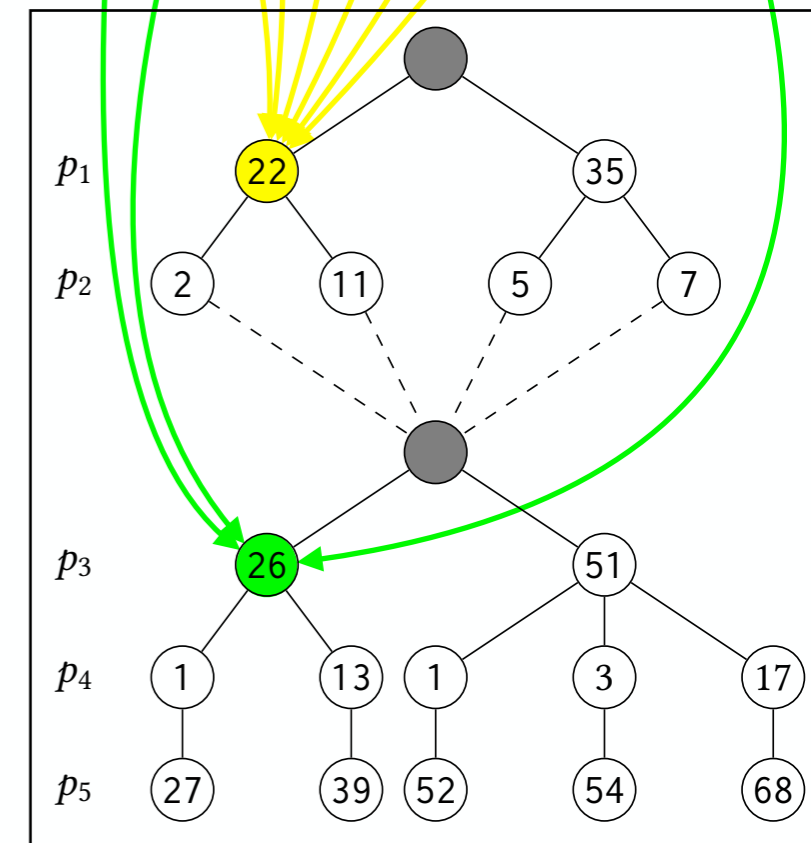
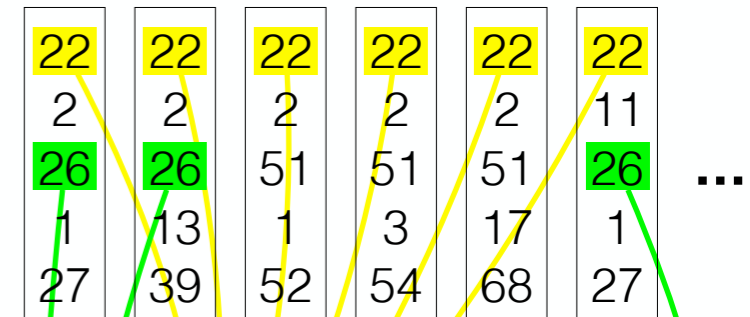
Memory-Friendly Storing:

- Avoid Redundancies
- interdependent parameters
→ "tree"
- independent parameters
→ "chain"

Efficient Exploration:

- CoT maps to Coordinate Space (CS)
- CS efficient structure for search techniques

Straightforward Search Space



ATF's CoT Search Space

ATF: Summary

ATF's efficiently handles tuning parameters with *interdependencies* among them:

ATF introduces novel concepts to
Generating & Storing & Exploring
the search spaces of interdependent parameters, based on its *novel*
constraint design and ***search space representation***

Further ATF features (not presented on slides for brevity):

- ATF's has a DSL-based **user interface** that is **arguably simpler** to use and more expressive than existing auto-tuners (including: OpenTuner & CLTune)
- ATF offers different kinds of **search techniques** and **abort conditions** (extensible)
- ATF offers a **DSL-based** user interface (*offline tuning*), as well as **GPL-based** interfaces (*online auto-tuning*):



pyATF

github.com/atf-tuner/pyATF



cppATF

github.com/atf-tuner/cppATF

... (future work)

Code Execution via HCA



[Overview](#) [Contact](#)

HCA

Host Code Abstraction (HCA)

*A High-Level Abstraction for Host Code Programming
Designed for Distributed, Heterogeneous Systems*

Overview

The **Host Code Abstraction (HCA)** is a high-level programming abstraction that simplifies implementing and optimizing so-called host code which is required in modern parallel programming approaches (e.g., **CUDA** and **OpenCL**) to execute code on the devices of distributed, heterogeneous systems.

More details will follow soon!

Contact

<https://hca-project.org>

**Just a brief outline of HCA
(for brevity)**

Journal of Supercomputing 2019

The Journal of Supercomputing (2020) 76:5117–5138
<https://doi.org/10.1007/s11227-019-02829-2>



dOCAL: high-level distributed programming with OpenCL and CUDA

Ari Rasch¹ · Julian Bigge¹ · Martin Wrodczyk¹ · Richard Schulze¹ · Sergei Gorlatch¹

Published online: 30 March 2019
© Springer Science+Business Media, LLC, part of Springer Nature 2019

Abstract

In the state-of-the-art parallel programming approaches OpenCL and CUDA, so-called host code is required for program's execution. Efficiently implementing host code is often a cumbersome task, especially when executing OpenCL and CUDA programs on systems with multiple nodes, each comprising different devices, e.g., multi-core CPU and graphics processing units; the programmer is responsible for explicitly managing node's and device's memory, synchronizing computations with data transfers between devices of potentially different nodes and for optimizing data transfers between devices' memories and nodes' main memories, e.g., by using pinned main memory for accelerating data transfers and overlapping the transfers with computations. We develop distributed OpenCL/CUDA abstraction layer (dOCAL)—a novel high-level C++ library that simplifies the development of host code. dOCAL combines major advantages over the state-of-the-art high-level approaches: (1) it simplifies implementing both OpenCL and CUDA host code by providing a simple-to-use, high-level abstraction API; (2) it supports executing arbitrary OpenCL and CUDA programs; (3) it allows conveniently targeting the devices of different nodes by automatically managing node-to-node communications; (4) it simplifies implementing data transfer optimizations by providing different, specially allocated memory regions, e.g., pinned main memory for overlapping data transfers with computations; (5) it optimizes memory management by automatically avoiding unnecessary data transfers; (6) it enables interoperability between OpenCL and CUDA host code for systems with devices from different vendors. Our experiments show that dOCAL significantly simplifies the development of host code for heterogeneous and distributed systems, with a low runtime overhead.

Goal of HCA

HCA simplifies implementing host code (e.g., OpenCL & CUDA) for the programmer:

- HCA offers an **easy-to-use high-level API** for host code programming that frees the user from boilerplate low-level commands (e.g., for memory allocations and data transfers)
- HCA allows **conveniently targeting multiple devices** and of potentially **different nodes** (by automatically managing device-to-device communications)
- HCA simplifies implementing **host code optimizations** (e.g., using pinned main memory for overlapping data transfers with computations)
- HCA **optimizes memory management** by automatically avoiding unnecessary data transfers (by internally generating & maintaining a data-dependency graph)
- ...



**Related approaches usually offer only a subset of these
HCA advantages**

HCA: Overview

CUDA vs. HCA — for executing parallel reduction CUDA kernel:

```
int main(int argc, char **argv)
{
    // initialization
    int i, j, gpuBase, GPU_N;
    cudaGetDeviceCount(&GPU_N);
    //...

    /* ... prepare input data ... */

    // Allocate device and host memory
    for (i = 0; i < GPU_N; i++) {
        cudaSetDevice(i);
        cudaStreamCreate(&plan[i].stream);
        cudaMalloc((void **)&plan[i].d_Data, plan[i].dataN *
sizeof(float));
        cudaMalloc((void **)&plan[i].d_Sum, ACCUM_N * sizeof(float));
        cudaMallocHost((void **)&plan[i].h_Sum_from_device, ACCUM_N *
sizeof(float));
        cudaMallocHost((void **)&plan[i].h_Data, plan[i].dataN *
sizeof(float));
        for (j = 0; j < plan[i].dataN; j++)
            plan[i].h_Data[j] = (float)rand()/(float)RAND_MAX;
    }

    // Perform data transfers and start device computations
    for (i = 0; i < GPU_N; i++) {
        cudaSetDevice(i);
        cudaMemcpyAsync(plan[i].d_Data, plan[i].h_Data, plan[i].dataN *
sizeof(float), cudaMemcpyHostToDevice, plan[i].stream);
        reduceKernel<<<BLOCK_N, THREAD_N, 0,
plan[i].stream>>>(plan[i].d_Sum, plan[i].d_Data, plan[i].dataN);
        cudaMemcpyAsync(plan[i].h_Sum_from_device, plan[i].d_Sum,
ACCUM_N *sizeof(float), cudaMemcpyDeviceToHost, plan[i].stream);
    }

    // combine GPUs' results
    for (i = 0; i < GPU_N; i++) {
        float sum;
        cudaSetDevice(i);
        cudaStreamSynchronize(plan[i].stream);
        sum = 0;
        for (j = 0; j < ACCUM_N; j++)
            sum += plan[i].h_Sum_from_device[j];
        *(plan[i].h_Sum) = (float)sum;
        cudaFreeHost(plan[i].h_Sum_from_device);
        cudaFree(plan[i].d_Sum);
        cudaFree(plan[i].d_Data);
        cudaStreamDestroy(plan[i].stream);
    }

    /* ... Compare GPU and CPU results ... */
}
```



Excerpt of NVIDIA SDK CUDA host code for executing parallel reduction CUDA kernel

```
#include "hca.hpp"

int main()
{
    int N = /* arbitrary chunk size */;

    // 1. choose devices
    auto devices = hca::get_all_devices<CUDA>();

    // 2. declare kernel
    hca::kernel reduction = cuda::source( /* kernel */ );

    const int GS = 32, BS = 256;

    // 3. prepare kernels' inputs
    hca::buffer<float> in ( N * devices.size() );
    hca::buffer<float> out( GS*BS * devices.size() );

    std::generate(in.begin(), in.end(), std::rand);

    // 4. start device computations
    for( auto& dev : devices )
        dev( reduction
            ( dim3( GS ), dim3( BS )
            ( write(out.begin()+dev.id()* GS*BS, GS*BS ),
              read (in.begin() +dev.id()* N , N
                  N
                  );
            );

    auto res = std::accumulate( out.begin(), out.end(), std::plus<float>() );

    std::cout << res << std::endl;
}
```

HCA

HCA Code (Equivalent to CUDA code left)

HCA relies on a high user-level of abstraction instead of complex, error-prone low-level functions

HCA: Summary

Goal of HCA:

HCA is a programming library for **simplifying host code programming** by abstracting from low-level details

Further HCA Features:

- enables interoperability between OpenCL and CUDA host code
- supports auto-tuning
- compatible with existing libraries (e.g., NVIDIA cuBLAS/cuDNN)
- simplifies profiling
- ...



MDH+ATF+HCA: User Interfaces

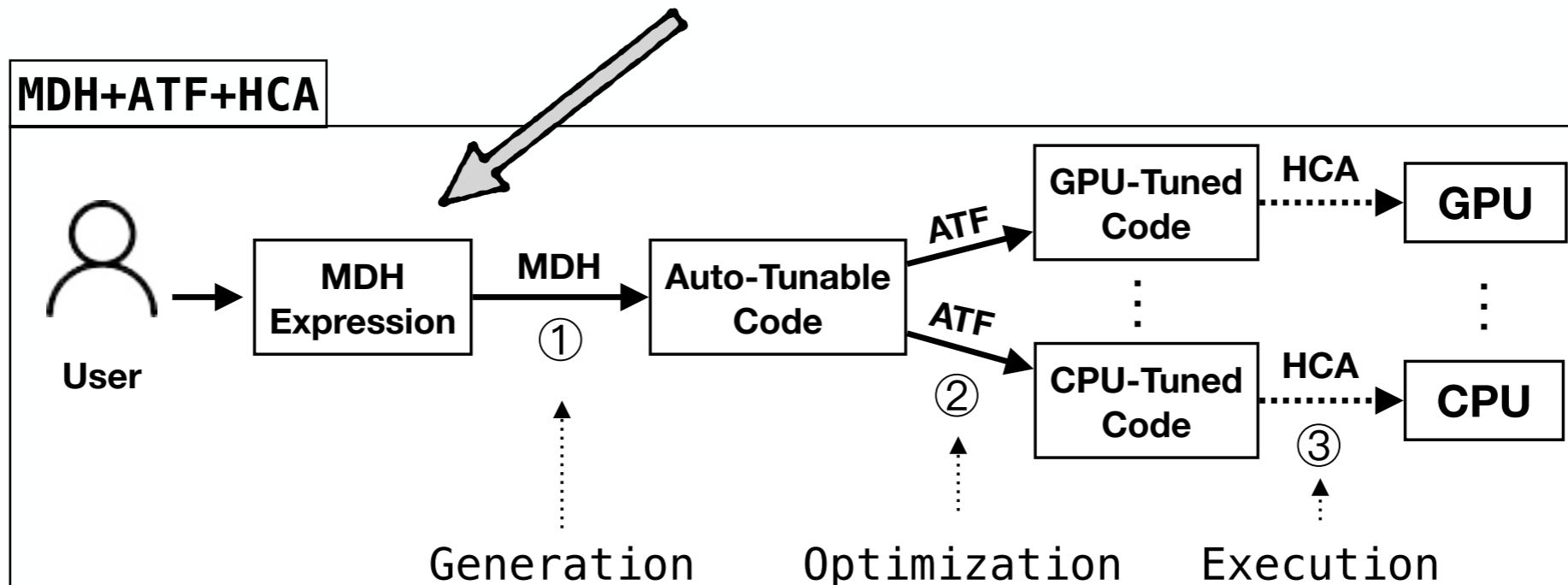
A Domain-Specific Language (DSL):

```
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) )
```

MDH's Formal Representation of MatVec

MDH's DSL Representation of MatVec

```
MatVec<T in TYPE | I,K in IN> := out_view<T>( w:(i,k)->(i) ) ◦
  md_hom<I,K>( *, (++,+) ) ◦
  inp_view<T,T>( M:(i,k)->(i,k),
  v:(i,k)->(k) )
```



MDH+ATF+HCA: User Interfaces

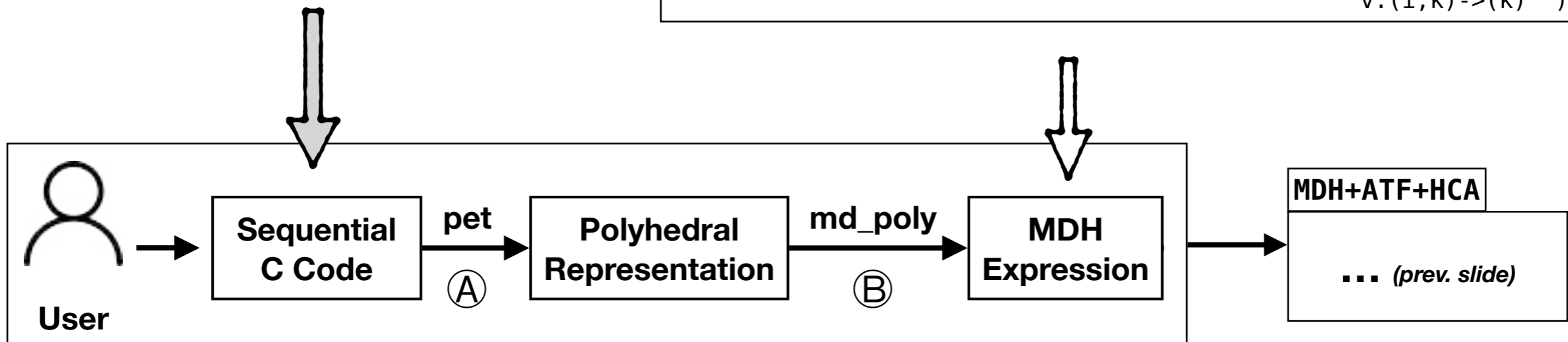
Automatic Parallelization:

C Implementation of MatVec

```
1 for( int i = 0; i < I; ++i )
2   for( int k = 0; k < K; ++k )
3   {
4     w[i] += M[i][k] * v[k];
5   }
```

MDH's DSL Representation of MatVec

```
MatVec<T in TYPE | I,K in IN> := out_view<T>( w:(i,k)->(i) ) o
                                md_hom<I,K>( *, (++,+) ) o
                                inp_view<T,T>( M:(i,k)->(i,k),
                                                v:(i,k)->(k) )
```



MDH pragma enables advanced optimizations:

```
#pragma mdh ( w[i] : ++ , + )
```

Conclusion

The **MDH+ATF+HCA** approach is a novel approach toward **Performance & Portability & Productivity** for data-parallel computations targeting modern parallel architectures:

COMPUTER SCIENCE

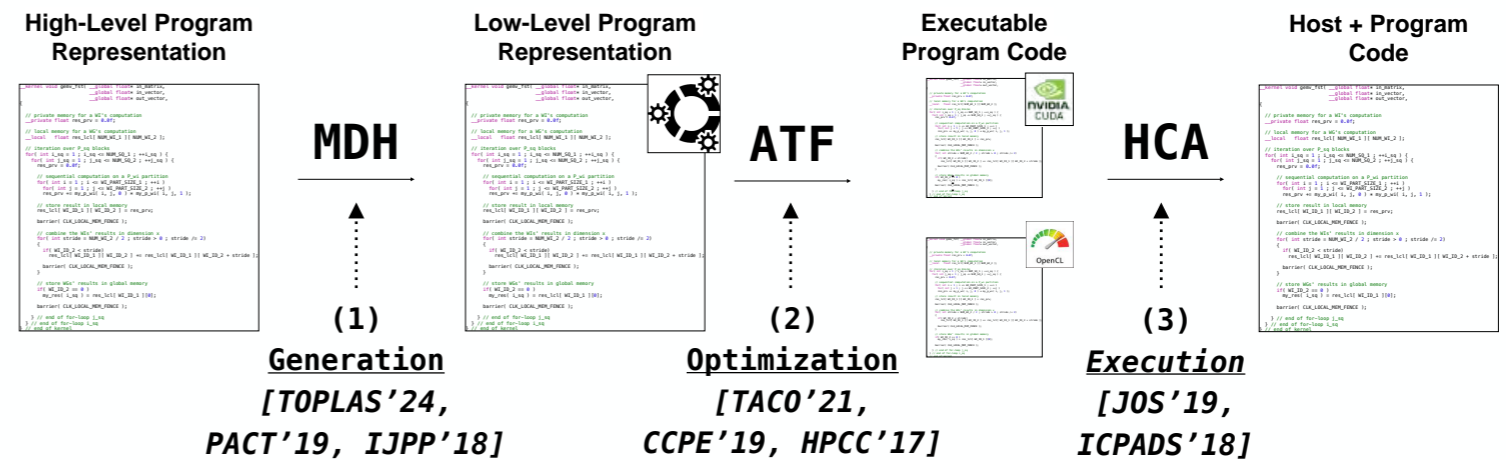
**TOWARD
PERFORMANCE & PORTABILITY & PRODUCTIVITY
IN PARALLEL PROGRAMMING**

A Holistic Code *Generation, Optimization, and Execution* Approach
for Data-Parallel Computations Targeting Modern Parallel Architectures

Inaugural Dissertation
for the Award of a Doctoral Degree
Dr. rer. nat.
in the Field of Mathematics and Computer Science
from the Faculty of Mathematics und Natural Sciences
of the University of Münster, Germany

submitted by
ARI RASCH
born in Essen, Germany
- 2024 -

The MDH+ATF+HCA approach:

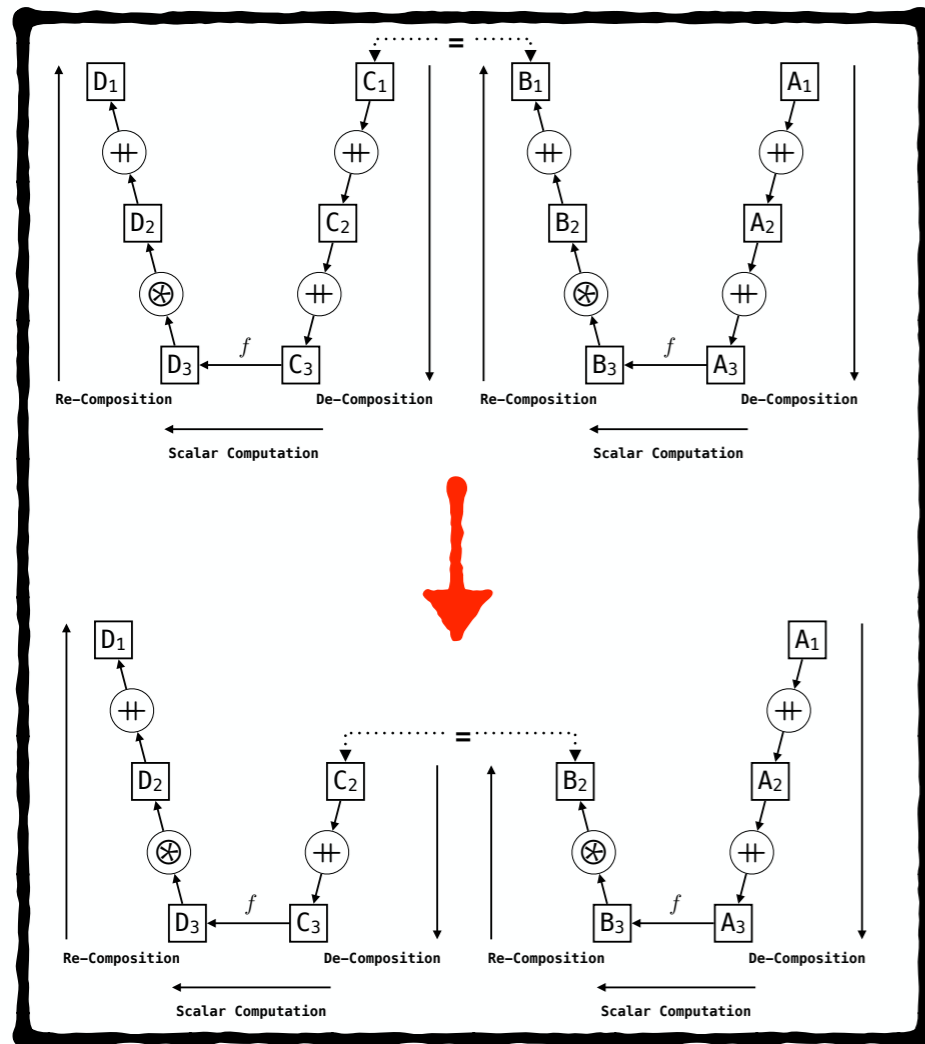


- The three sub-projects — MDH & ATF & HCA — complement each other to a holistic code *Generation & Optimization & Execution* approach
- Our holistic approach takes as input easy-to-use DSL programs or sequential (annotated) C code¹
- There are many (promising) future directions for MDH & ATF & HCA (one part of thesis dedicated to FW)

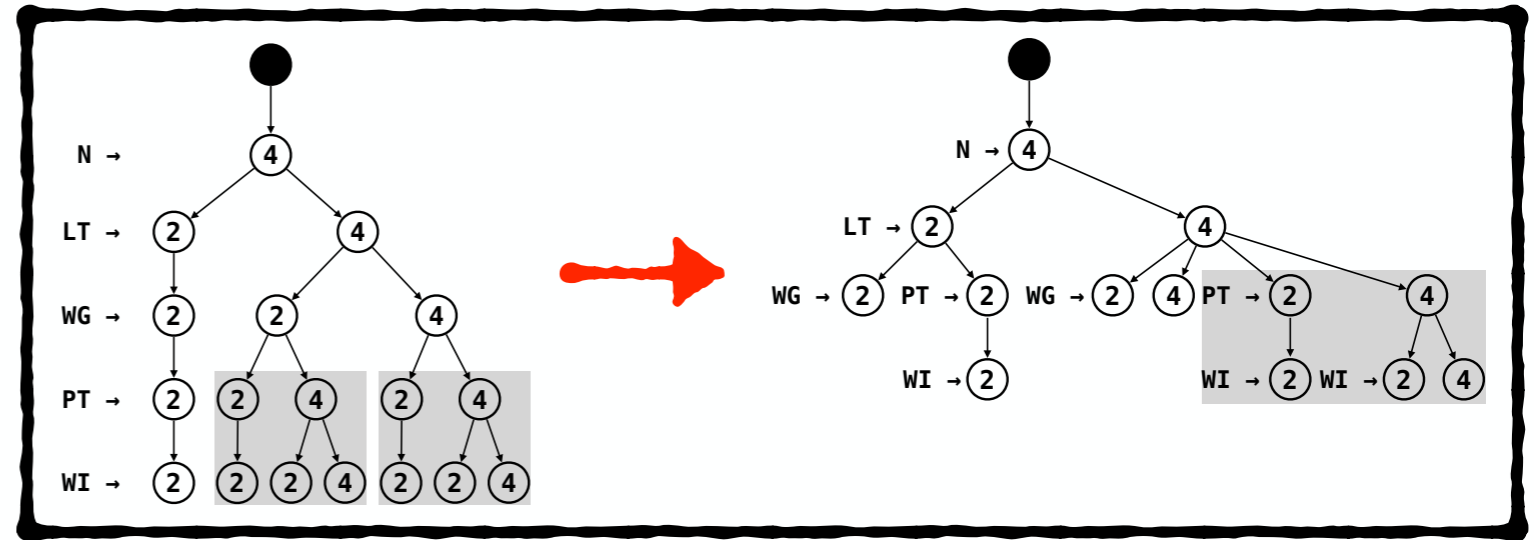
¹ A DSL-equivalent *MLIR* frontend will be introduced soon.

Stay Tuned !

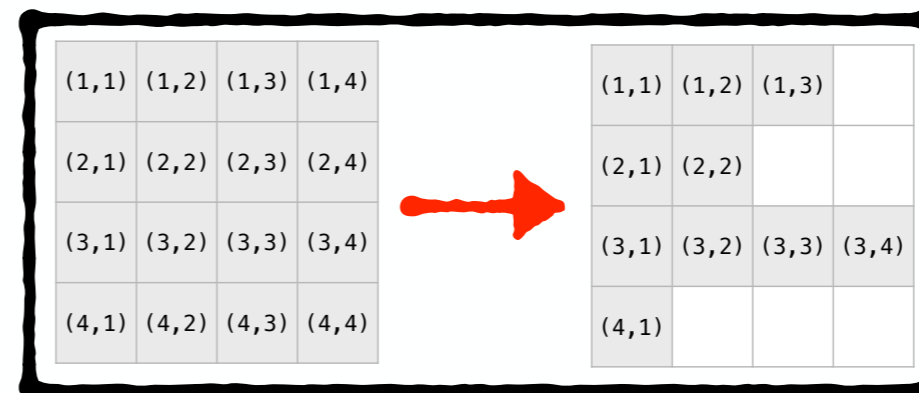
The approaches MDH & ATF & HCA are starting points for many future directions:



Fusion Optimization (MDH)



CoT Optimization (ATF)



Irregular Data Accesses (MDH)

... (many more)

Encouraging WIP results — already awarded & funded:



Questions