

WESTFÄLISCHE
WILHELMS-UNIVERSITÄT
MÜNSTER

Multi-Dimensional Homomorphisms and Their Implementation in OpenCL

*An Algebraic Approach to Performance, Portability, and Productivity
for Data-Parallel Applications
on Multi- and Many-Core Architectures*

Ari Rasch, Richard Schulze, and Sergei Gorlatch

University of Münster, Germany

Motivation

Observation:

Applications

Linear Algebra (BLAS)

GEMM GEMV

 DOT

... 152 routines!

Tensor Contractions

...

stencils

...

Architectures

X

The central 'Architectures' section features a collection of logos for various hardware manufacturers. At the top left is the NVIDIA TESLA logo. To its right is the NVIDIA PASCAL TITAN X logo. Below these is the ARM logo in blue. To the right of ARM is the intel inside Xeon Phi logo. Below ARM is the AMD RADEON GRAPHICS logo. To the right of AMD is the intel inside XEON logo. At the bottom center is the IBM logo. A large red 'X' is positioned to the left of the ARM logo, and another large red 'X' is positioned to the right of the intel inside Xeon Phi logo.

Input Sizes

Machine Learning

$C = A \cdot B$

...

$C = A \cdot B$

The 'Machine Learning' section shows a diagram of matrix multiplication $C = A \cdot B$ with a 3D orange arrow pointing upwards and to the right, indicating increasing input sizes.

Numerical Computations

$C = A \cdot B$

...

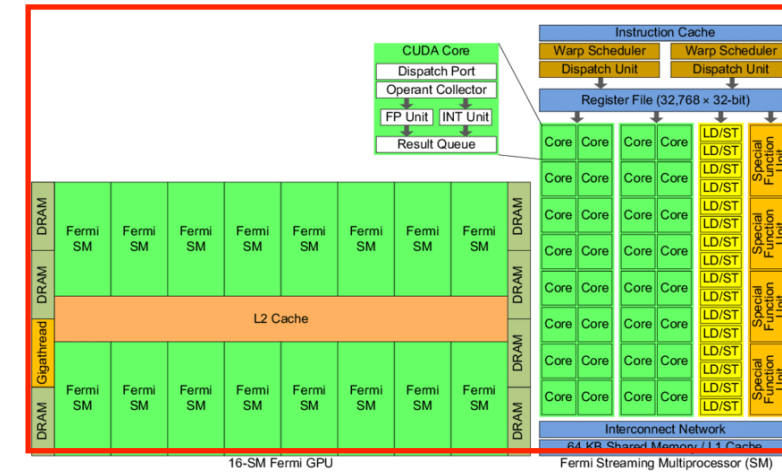
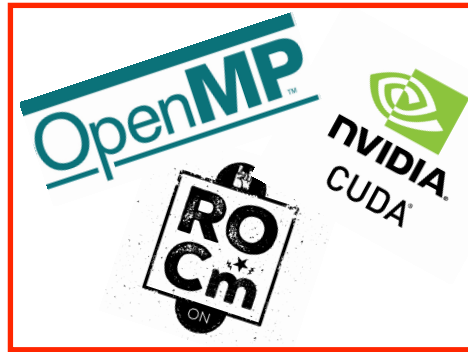
$C = A \cdot B$

The 'Numerical Computations' section shows a diagram of matrix multiplication $C = A \cdot B$ with a 3D orange arrow pointing upwards and to the right, indicating increasing input sizes.

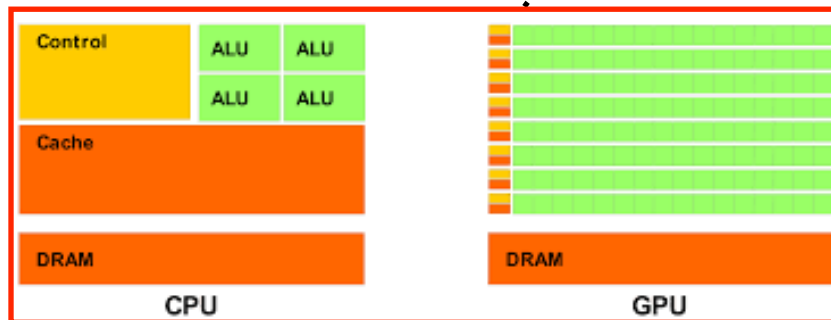
Combinatorial Explosion

Motivation

In a perfect world:



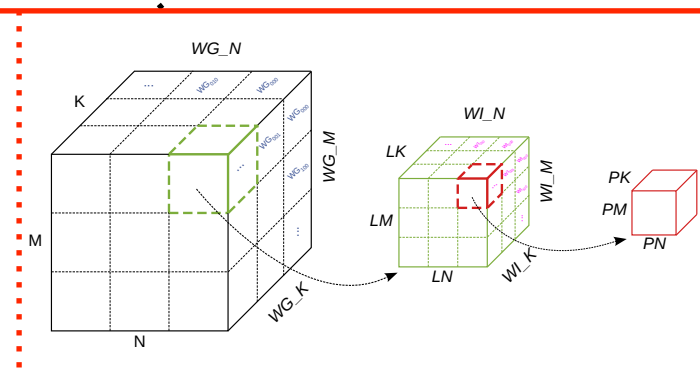
Write **one piece of code** for our application that provides **high performance**, is **performance-portable over different architectures** and **input sizes** and that is **easy to implement**.



```

gridDim.x = 4096
threadIdx.x  threadIdx.x  threadIdx.x  threadIdx.x
0 1 2 3 ... 255 0 1 2 3 ... 255 0 1 2 3 ... 255 ... 0 1 2 3 ... 255
blockIdx.x = 0  blockIdx.x = 1  blockIdx.x = 2  blockIdx.x = 4095

index = blockIdx.x * blockDim.x + threadIdx.x
index = (2) * (256) + (3) = 515
    
```



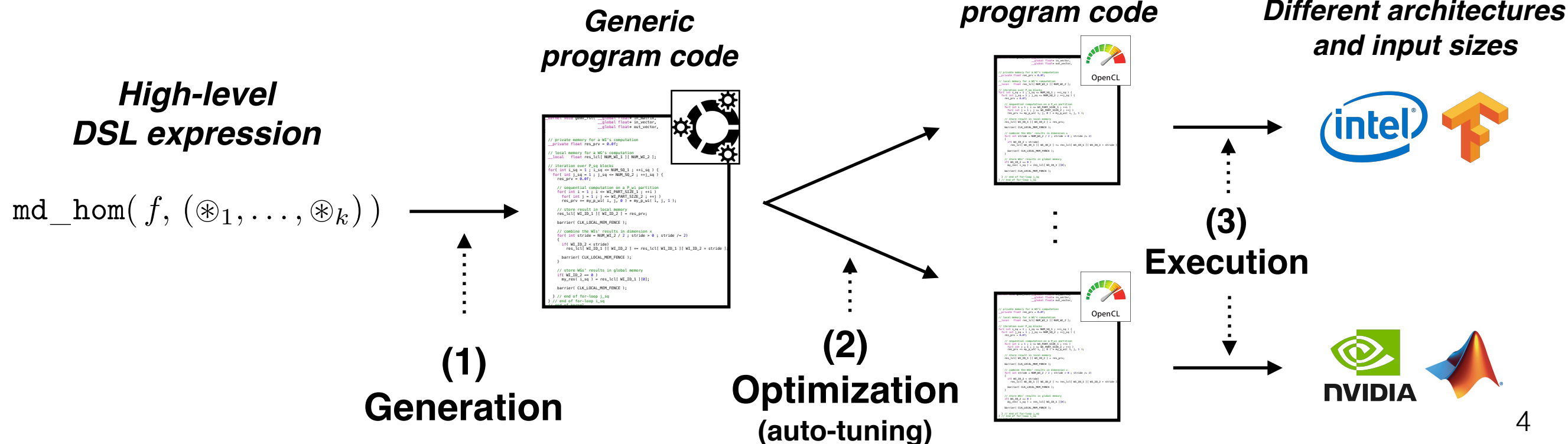
dimensions			$C = AB$	$C = A^T B^T$	small = 1
m	n	k	Shape	Shape	
large	large	large	$C = \begin{bmatrix} A & B \end{bmatrix}$	$C = \begin{bmatrix} A^T & B^T \end{bmatrix}$	gemm
large	large	small	$C = \begin{bmatrix} A \\ B \end{bmatrix}$	$C = \begin{bmatrix} A^T \\ B^T \end{bmatrix}$	ger
large	small	large	$C = \begin{bmatrix} A & B \end{bmatrix}$	$C = \begin{bmatrix} A^T \\ B^T \end{bmatrix}$	gemv
large	small	small	$C = \begin{bmatrix} A \\ B \end{bmatrix}$	$C = \begin{bmatrix} A^T \\ B^T \end{bmatrix}$	axpy ($\beta = 1$)
small	large	large	$C = \begin{bmatrix} A & B \end{bmatrix}$	$C = \begin{bmatrix} A^T \\ B^T \end{bmatrix}$	gemv
small	large	small	$C = \begin{bmatrix} A \\ B \end{bmatrix}$	$C = \begin{bmatrix} A^T \\ B^T \end{bmatrix}$	axpy ($\beta = 1$)
small	small	large	$C = \begin{bmatrix} A & B \end{bmatrix}$	$C = \begin{bmatrix} A^T \\ B^T \end{bmatrix}$	dot ($\alpha = \beta = 1$)
small	small	small	$C = \begin{bmatrix} A & B \end{bmatrix}$	$C = \begin{bmatrix} A^T \\ B^T \end{bmatrix}$	scalar mult.

Our Approach

We provide an approach to address all these challenges for our class of *Multi-Dimensional Homomorphisms*:

- **Multi-Dimensional Homomorphisms (MDHs)** are **formally-defined** class of functions that **cover important data-parallel applications**, e.g.: linear algebra (BLAS), stencils computations, ...
- We **enable conveniently expressing MDHs** by providing a **high-level DSL** for them.
- We provide a **DSL compiler** to **generate OpenCL code** so that we can target various parallel architectures (e.g., Intel CPU, NVIDIA GPU, ...).
- Our OpenCL code is **fully automatically optimizable** (auto-tunable) — for each combination of an **MDH, target architecture, and input size** — by being generated as targeted to the **OpenCL's abstract device models** (and not to a particular architecture, e.g., NVIDIA GPU) and as **parametrized in the models' performance-critical parameters**.

Our approach consists of three major steps:



Agenda

We discuss the three major steps of our approach:

- 1. Generation**
- 2. Optimization**
- 3. Execution**

Afterwards:

- 4. Experimental Results**
- 5. Current/Future Work**

Multi-Dimensional Homomorphisms

Our class of targeted applications is formally specified as:

Definition: [*Multi-Dimensional Homomorphisms [1]*]

Let T and T' be two arbitrary types. A function $h : T[N_1] \dots [N_d] \rightarrow T'$ on d -dimensional arrays is called a *Multi-Dimensional Homomorphism (MDH)* iff there exist *combine operators* $\otimes_1, \dots, \otimes_d : T' \times T' \rightarrow T'$, such that for each $k \in [1, d]$ and arbitrary, concatenated input MDA $a ++_k b$:

$$h(a ++_k b) = h(a) \otimes_k h(b)$$

Definition: [**md_hom**]

We write

$$\text{md_hom} (f, (\otimes_1, \dots, \otimes_d))$$

for the unique d -dimensional homomorphism with combine operators $\otimes_1, \dots, \otimes_d$ and action f on singleton arrays.

MDH Examples

Important functions are MDHs — we can express them conveniently using our `md_hom` pattern:

Linear Algebra

GEMM = `md_hom(*, (++, ++, +)) o view(A, B)(i, j, k)(A[i, k], B[k, j])`

Whats happening?

1. prepare the domain-specific input for `md_hom` as multi-dimensional array using function `view`
 - here: fuse matrices A and B to 3-dimensional array of pairs — array indices are mapped to the elements of A and B to multiply: $i, j, k \mapsto (A[i, k], B[k, j])$
2. apply multiplication (denoted as `*`) to each pair
3. combine results in dimension `k` by addition (`+`)
4. combine results in dimension `i` and `j` by concatenation (`++`)

MDH Examples

Important functions are MDHs — we can express them easily using our md_hom pattern:

Linear Algebra

```
GEMM = md_hom( *, (++, ++, +) ) o view( A,B )( i,j,k )( A[i,k], B[k,j] )
GEMV = md_hom( *, (++,      +) ) o view( A,B )( i,  k )( A[i,k], B[k]      )
DOT   = md_hom( *, (      , +) ) o view( A,B )(      k )( A[k]      , B[k]      )
```

- **GEMV** and **DOT** are expressed similarly as **GEMM**
- both are special cases of **GEMM**:
 - **GEMV**: **B** is **Kx1** matrix (i.e., no dimension **j**)
 - **DOT**: **A** is **1xK** matrix, **B** is **Kx1** matrix (i.e., no dimensions **i** and **j**)
- Note: **GEMV** and **DOT** are not required in our approach — we can use **GEMM** for them (provides the same high performance).

MDH Examples

Important functions are MDHs — we can express them easily using our md_hom pattern:

Linear Algebra

```
GEMM = md_hom( *, (++, ++, +) ) o view( A,B )( i,j,k )( A[i,k], B[k,j] )
GEMV = md_hom( *, (++,   +) ) o view( A,B )( i,  k )( A[i,k], B[k]   )
DOT   = md_hom( *, (    +) ) o view( A,B )(    k )( A[k]   , B[k]   )
```

Stencil Computations

```
Gaussian_2D = md_hom( G_func, (++,++) ) o view(...)
Jacobi_3D   = md_hom( J_func, (++,++,++) ) o view(...)
```

Data Mining

```
PRL = md_hom( weight, (++, max) ) o view(...)
```

Machine Learning

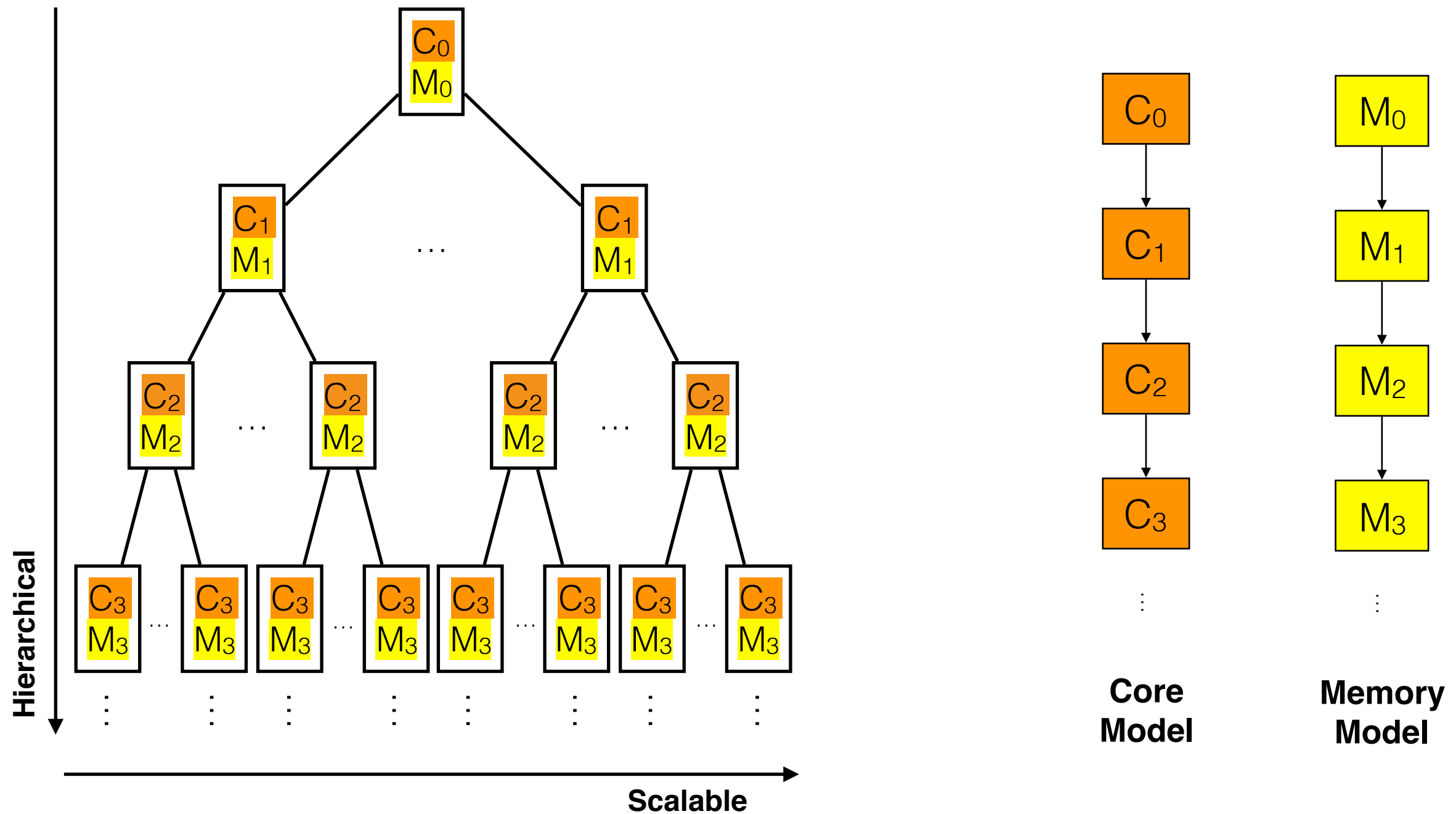
```
TC = md_hom( *, (++,...,++ , +,...,+) ) o view(...)
```

Further examples: MLP, SVM, ECC, ..., Mandelbrot, Parallel Reduction, ...

**Our DSL needs only the patterns
md_hom(...) and view(...)**

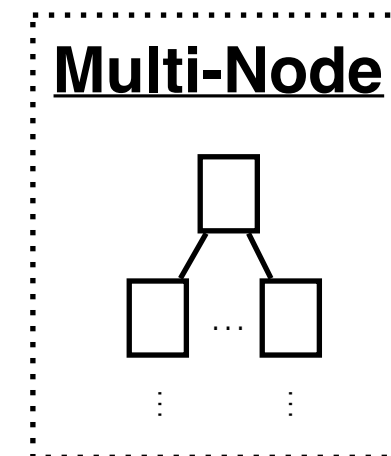
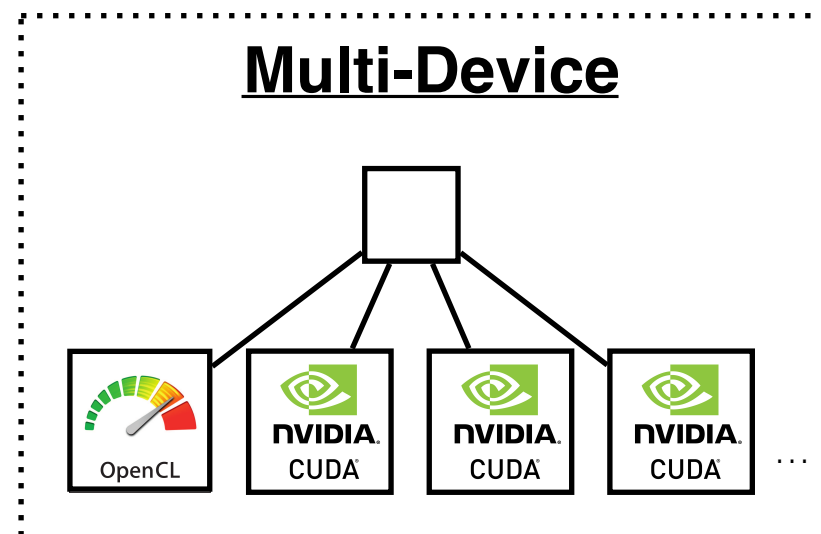
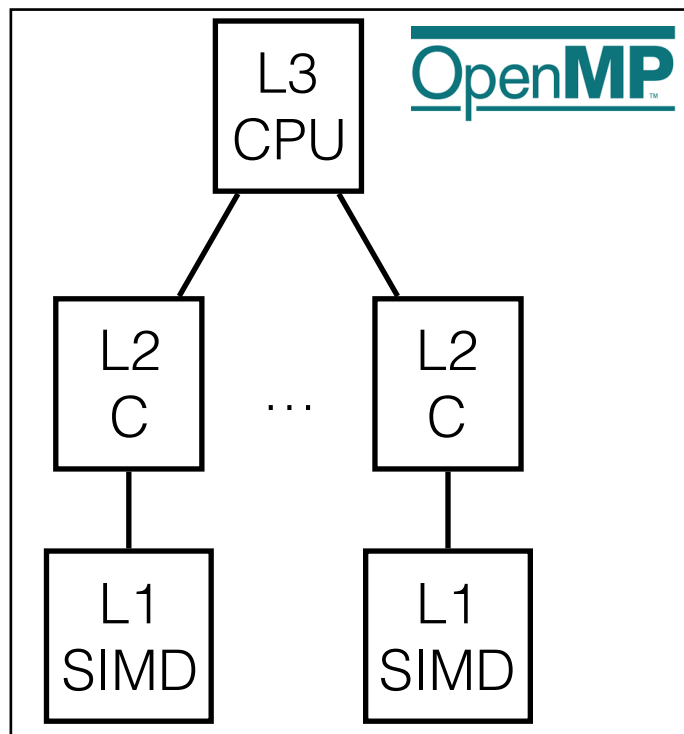
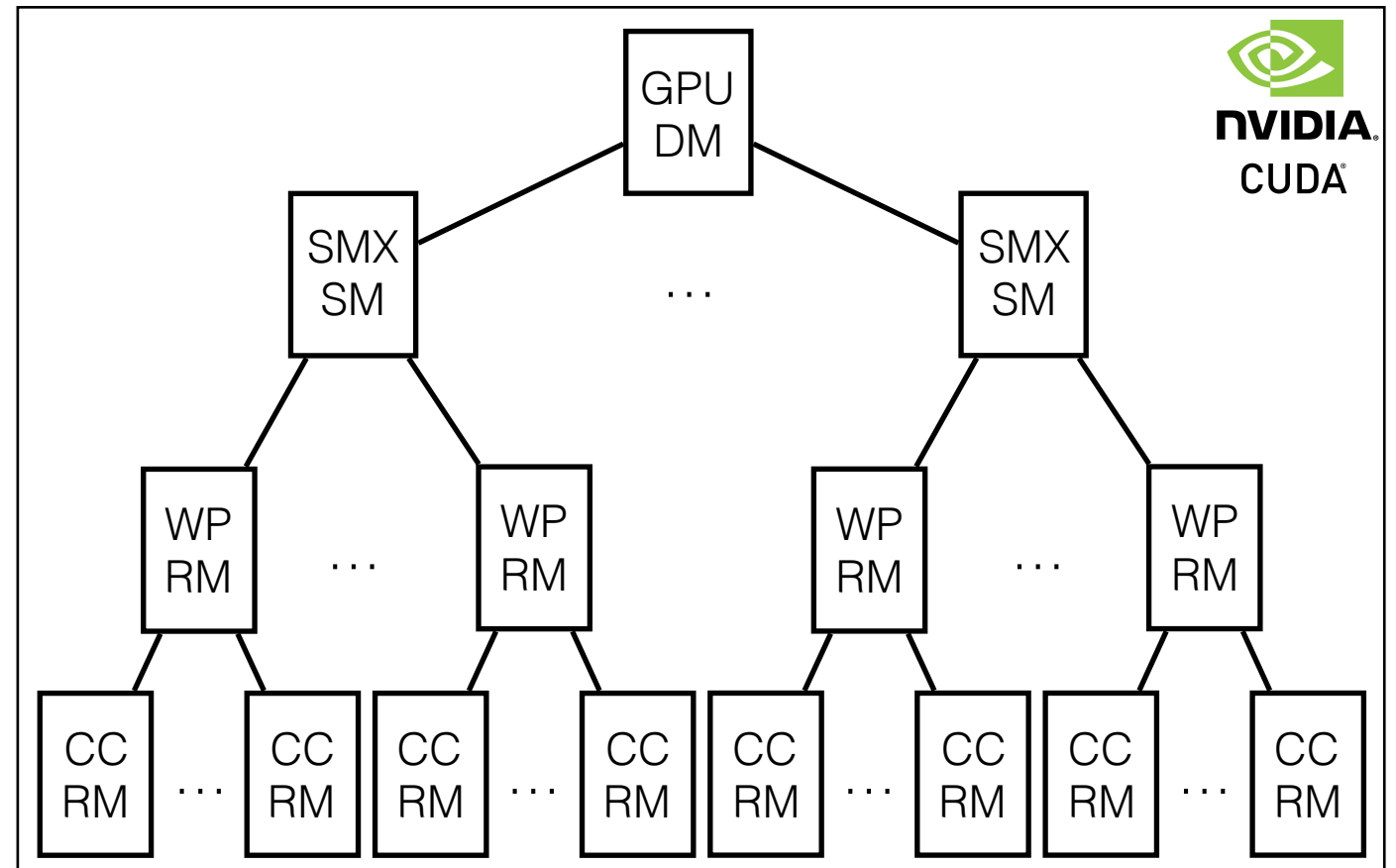
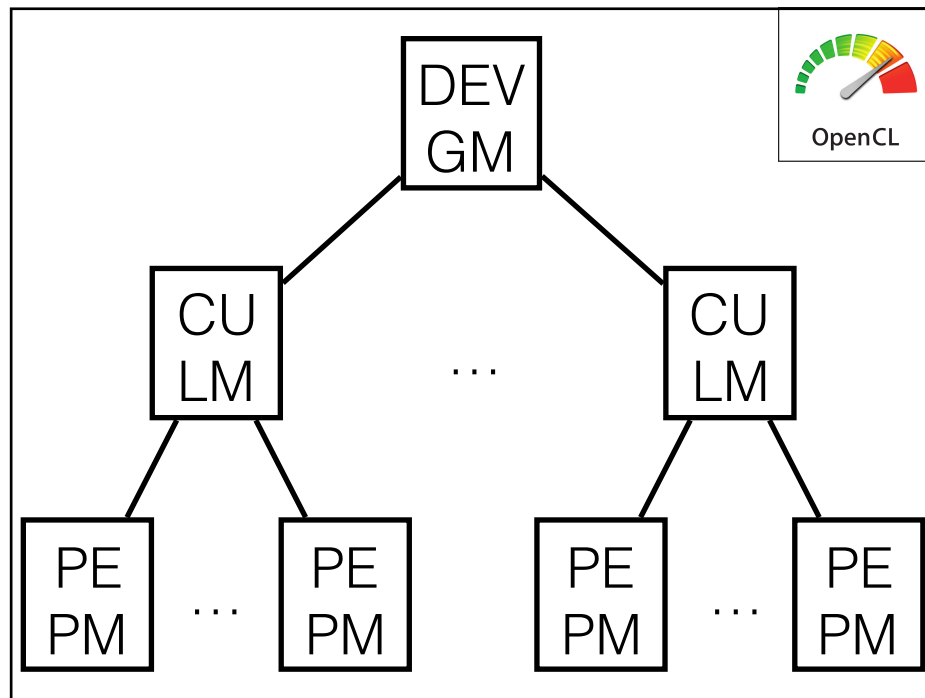
MDH – Target Machine Model

We **generate** code for MDHs targeting ***hierarchical, scalable*** machine models:



MDH – Target Machine Model

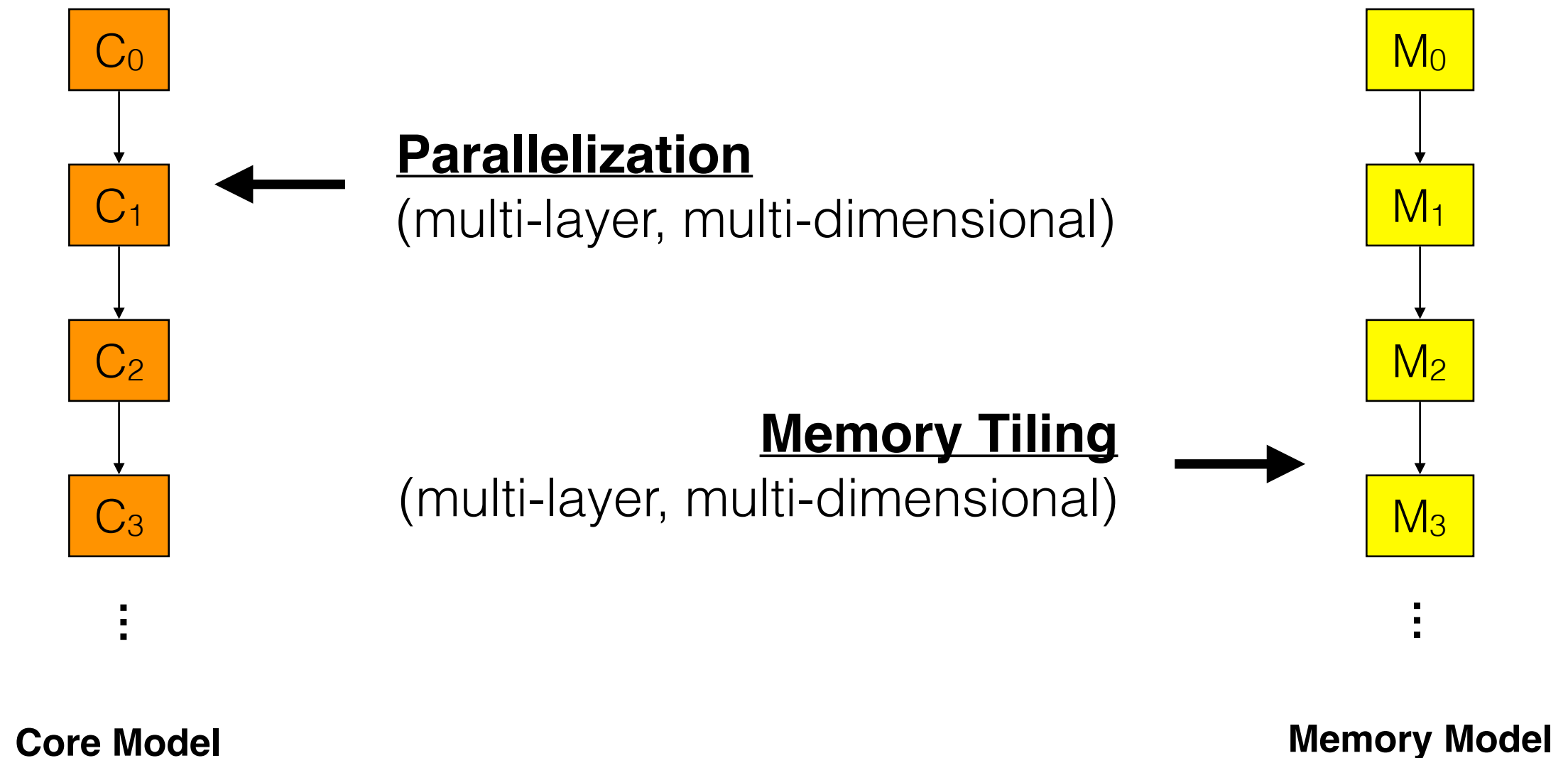
Examples of such *machine models*:



...

Code Generation for MDHs

Our **uniform md_hom representation** of MDHs enables **systematically generating code for such machine models** that can be **automatically optimized**:



In the following: Explain implementation at example of OpenCL.

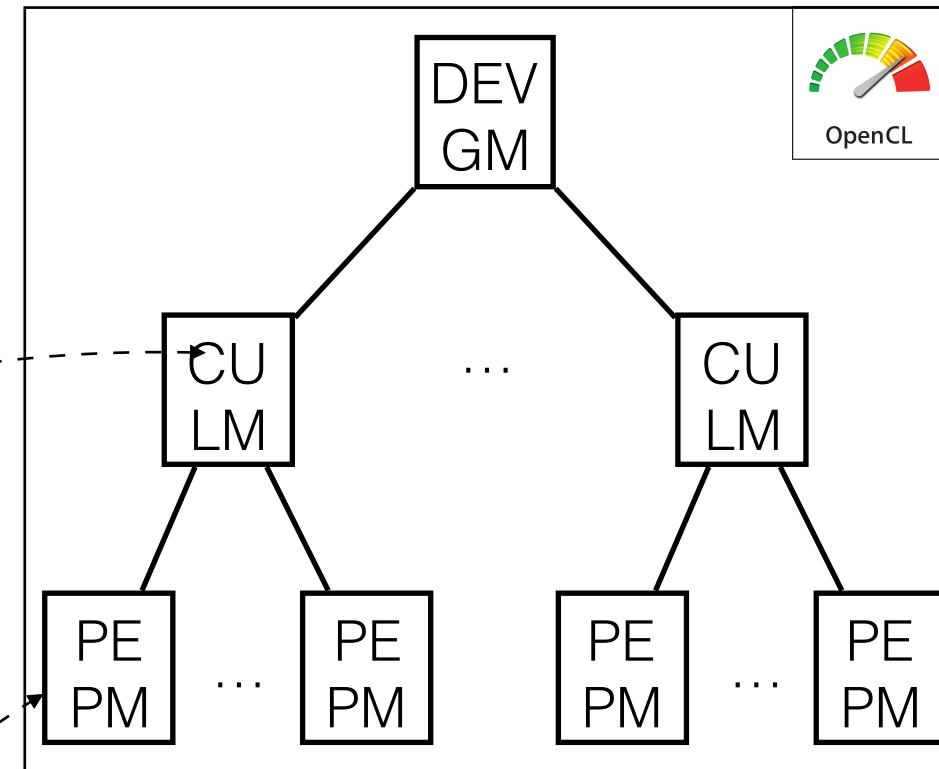
Optimization of MDHs

1. Parallelization (multi-layer, multi-dimensional):

```
#reduce  $\otimes_1$ 
for( i_1 = 1, ... , N_1 )
...
#reduce  $\otimes_d$ 
for( i_d = 1, ... , N_d )
{
  f( a[i_1]...[i_d] )
}
```

MDH
pseudocode for
 $md_hom(f, (\otimes_1, \dots, \otimes_d))$

```
#reduce  $\otimes_1$ 
parallel_for( i_1 = 1, ... , NUM_WG_1 )
...
#reduce  $\otimes_d$ 
parallel_for( i_d = 1, ... , NUM_WG_d )
```



```
#reduce  $\otimes_1$ 
parallel_for( ii_1 = 1, ... , NUM_WI_1 )
...
#reduce  $\otimes_d$ 
parallel_for( ii_d = 1, ... , NUM_WI_d )
```

- We parallelize for each of OpenCL's **two parallel layers**.
- We parallelize **on each layer** in **all d dimensions** of the MDH.
- ▶ We **auto-tune** the number of threads **on each layer** and **in each dimension**.

Optimization of MDHs

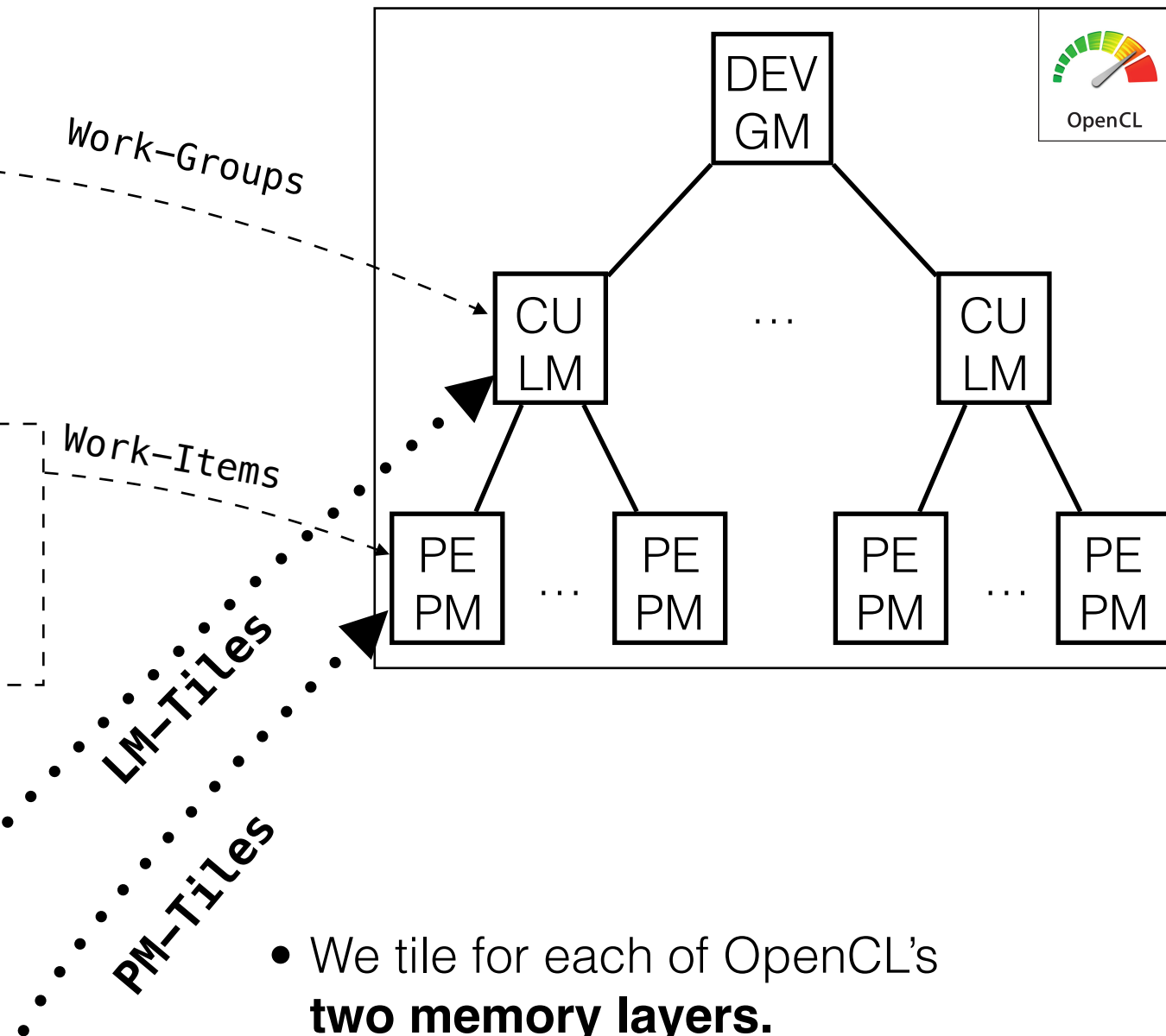
2. Memory Tiling (multi-layer, multi-dimensional):

```
#reduce ⊗1
parallel_for( i1 = 1, ... , NUM_WG_1 )
...
#reduce ⊗d
parallel_for( id = 1, ... , NUM_WG_d )
```

```
#reduce ⊗1
parallel_for( ii1 = 1, ... , NUM_WI_1 )
...
#reduce ⊗d
parallel_for( iid = 1, ... , NUM_WI_d )
```

```
for( j1 = 1, ... , NUM_LM_TL_1 )
...
for( jd = 1, ... , NUM_LM_TL_d )
```

```
for( j1 = 1, ... , NUM_PM_TL_1 )
...
for( jd = 1, ... , NUM_PM_TL_d )
{
    f( ... )
}
```



- We tile for each of OpenCL's **two memory layers**.
- We tile **on each layer** in **all dimensions** of the MDH.
- ▶ We **auto-tune** the sizes of tiles on each layer and in each dimension.

Automatic Performance Tuning

Our OpenCL implementation is generated as **generic in performance-critical parameters** (a.k.a. tuning parameters) of OpenCL's abstract device models:

	Name	Description
Platform Model	NUM_WG_i	number of Work-Groups
	NUM_WI_i	number of Work-Items
Memory Model	LM_TL_SIZE_i	local memory tile size
	PM_TL_SIZE_i	private memory tile size
	LM_CACHING	copying tiles to local memory on/off
	PM_CACHING	copying tiles to private memory on/off
	COMP_LM_TL	compute local tiles in: PM, LM, or GM
	COMP_PM_TL	compute private tiles in: PM, LM, or GM
Platform↔Memory Model	OCL_MDA_MAPPING	mapping OCL to MDA dimensions
	COMB_OP_ORDER	order of combine operators
	COMB_WG_RES_AFTER	when to combine WG results
	COMB_WI_RES_AFTER	when to combine WI results
	...	

All parameters are chosen independently of the target MDH, architecture, and input size!

Automatic Performance Tuning

We use our ***Auto-Tuning Framework (ATF)*** [1] to automatically choose optimized values of our performance-critical parameters.

	Domain-specific auto-tuning	OpenTuner	CLTune	ATF
Arbitrary Programming Language		✓		✓
Arbitrary Application Domain		✓	✓	✓
Arbitrary Tuning Objective	✓	✓		✓
Arbitrary Search Technique	✓	✓	✓	✓
Interdependent Parameters	✓		✓	✓
Large Parameter Ranges	✓	✓		✓
Directive-Based Auto-Tuning				✓
Automatic Cost Function Generation	✓		✓	✓

ATF combines major advantages over state-of-the-art auto-tuning approaches

Automatic Performance Tuning

ATF usage: We annotate our program code with easy-to-use *tuning directives*.

```
#atf::tp name      NUM_WG_1
        range      interval<int>( 1, N_1 )
```

```
#atf::tp name      NUM_WI_1
        range      interval<int>( 1, N_1 )
```

```
// ...
```

```
#atf::tp name      LM_SIZE_1
        range      interval<int>( 1, N_1 )
        constraint LM_SIZE_1 <= N_1
```

```
#atf::tp name      PM_SIZE_1
        range      interval<int>( 1, N_1 )
        constraint PM_SIZE_1 <= LM_SIZE_1
```

```
// ...
```

```
// OpenCL kernel code
```



ATF
automatically
determines
optimized parameter
values

(ATF is also available as C++/Python programming library)

Execution

We execute our *generated* and *optimized* OpenCL code using our own **dOCAL [1]** framework which:

1. provides **high-level abstractions** for simplifying implementing **OpenCL host code**, especially for multi-device systems (e.g., by automatically performing memory allocations and synchronization);
2. provides **asynchronous computation efficiency** (e.g., overlapping data transfers and/or kernel computations) by generating and maintaining a data-dependency graph transparently from the user;
3. enables conveniently executing OpenCL kernels on **remote nodes** (via *Boost.Asio*).

Execution

Illustration: using dOCAL for executing OpenCL kernel on GPU.

```
#include "docal.hpp"

int main()
{
    // 1. choose device
    auto device = docal::get_device( "NVIDIA", "Tesla" );

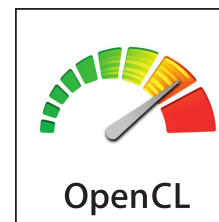
    // 2. declare kernel
    docal::kernel GEMM = docal::source( /* OpenCL Code */ );

    // 3. prepare kernels' inputs
    docal::buffer<float> A( N*N );
    docal::buffer<float> B( N*N );
    docal::buffer<float> C( N*N );

    std::generate( A.begin(), A.end(), std::rand );
    std::generate( B.begin(), B.end(), std::rand );

    // 4. start device computations
    device( GEMM
            ( nd_range( /* GS */, nd_range( /* LS */ ) )
            ( read( A ), read( B ), write( C ) ) );

    // 5. print result
    for( int i = 0 ; i < N*N ; ++i )
        std::cout << C[ i*N + j ];
}
```



Experimental Results

We compare our **automatically-generated and optimized code** using:

Applications

1. Linear Algebra Routines (GEMM, GEMV)
2. Stencil Computations (Gaussian Convolution 2D, Jacobi 3D)
3. Tensor Contractions

Competitors

- Performance-Portable approaches
- Domain-specific, hand-optimized approaches

Architectures

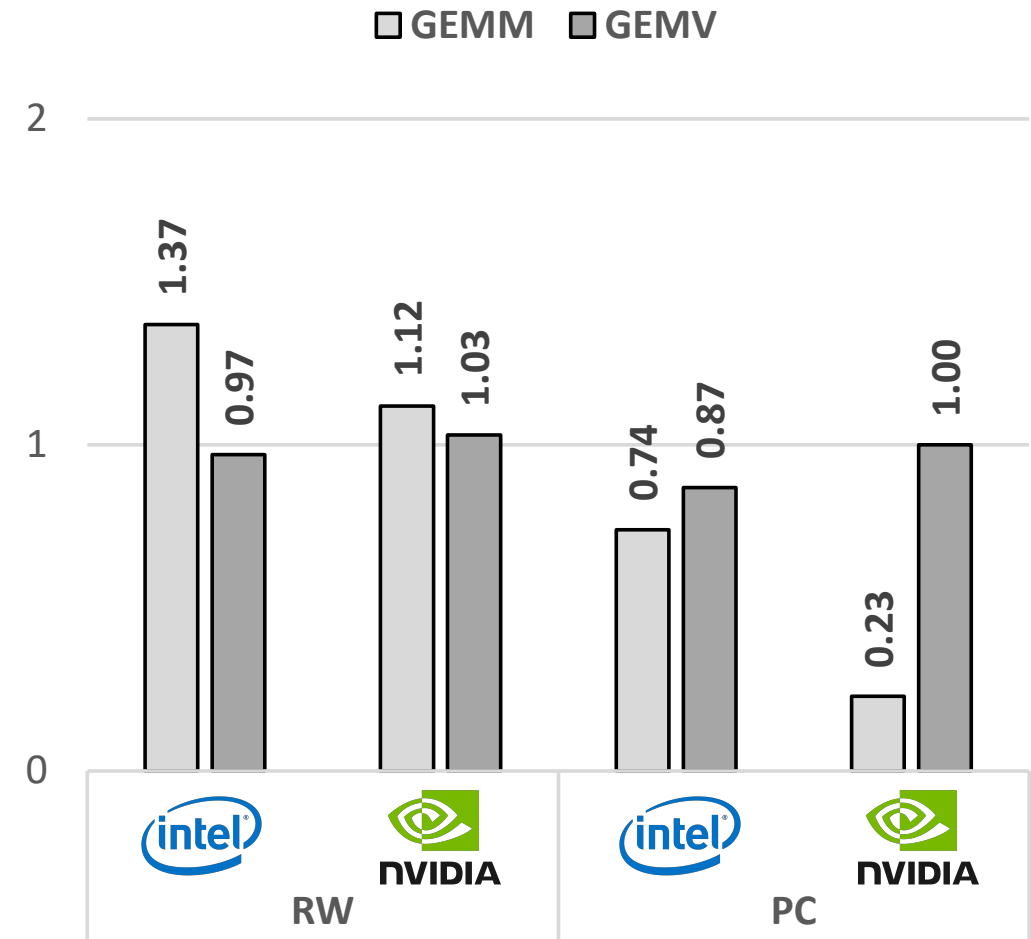
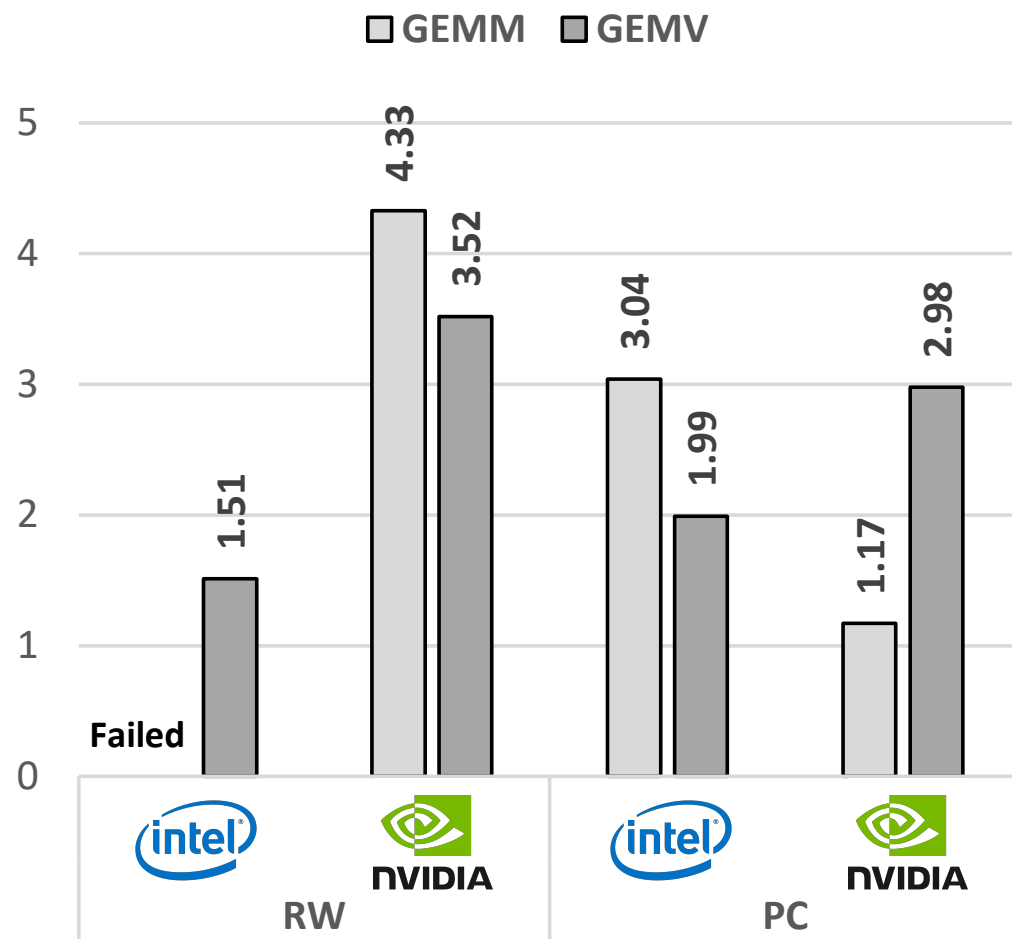
- Intel Xeon multi-core CPU (E5-2640)
- NVIDIA Tesla V100 GPU (SMX2-16GB)

Data Sets

- RW: Real-world sizes from Deep Learning
- PC: Sizes that are preferable for our competitors

Experimental Results

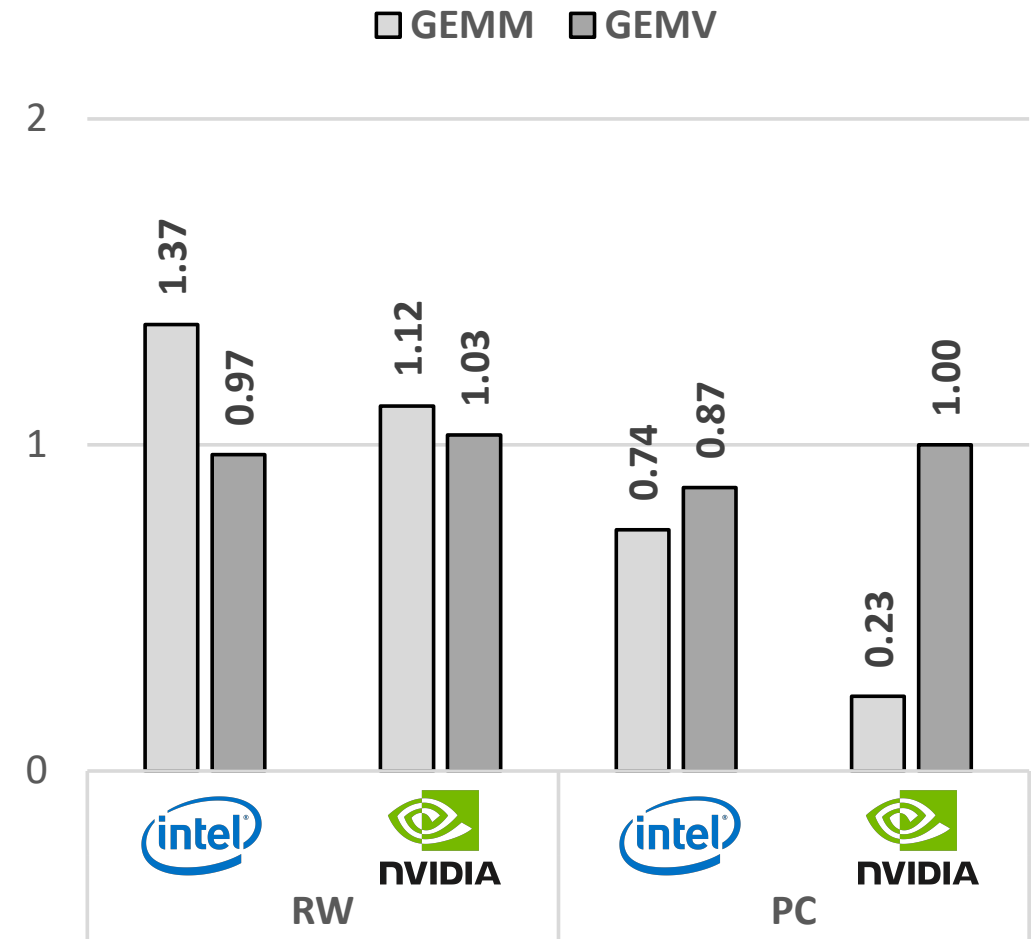
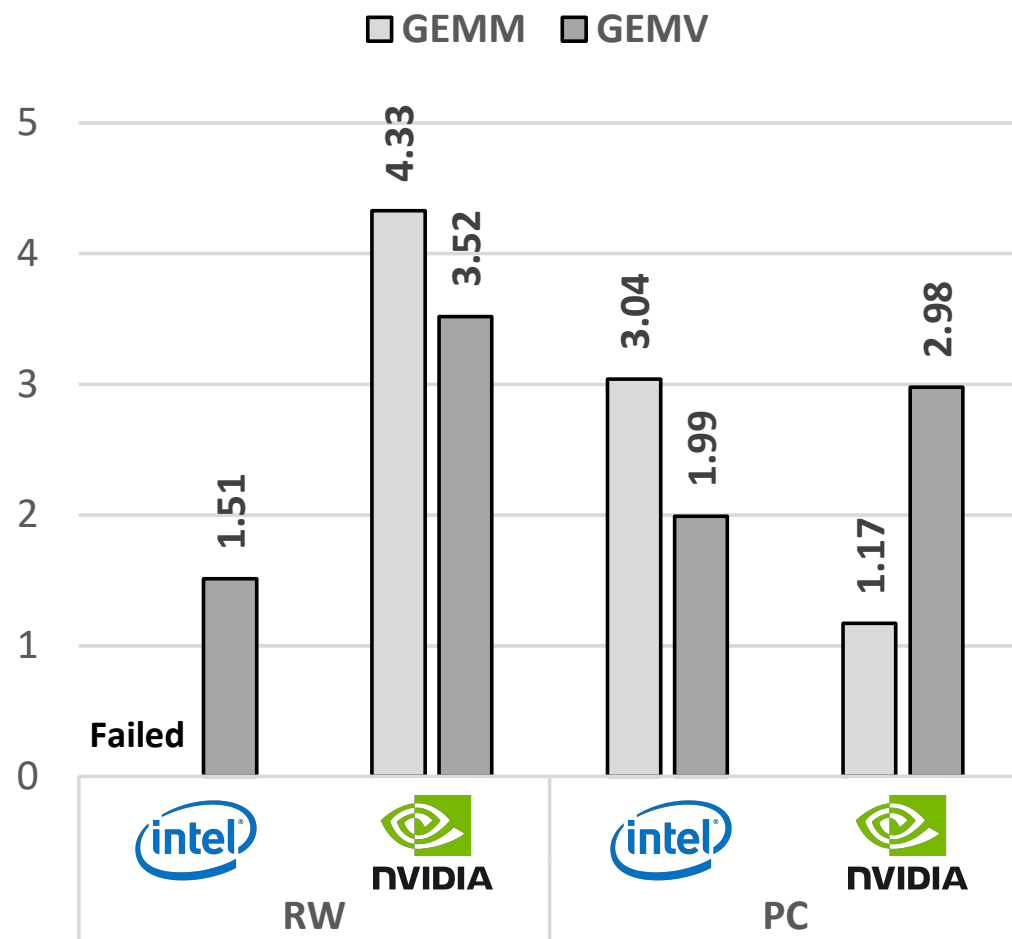
1. Linear Algebra:



- We demonstrate speedup of our approach for both **GEMM** and **GEMV** over
 - *Lift* (left) — currently one of the best-performing performance-portable approaches (based on OpenCL)
 - *Intel MKL & NVIDIA cuBLAS* (right) — hand-optimized, BLAS-specific approaches for Intel/NVIDIA hardware only
- Note: We use our **GEMM** implementation also for **GEMV** → demonstrates efficiency of our approach for very irregular sizes

Experimental Results

1. Linear Algebra:



Comparison to Lift [1,2]:

- Speedups up to **> 4x**.
- Our approach relies on generic optimizations + auto-tuning (rather than transformation rules).

Comparison to MKL/cuBLAS:

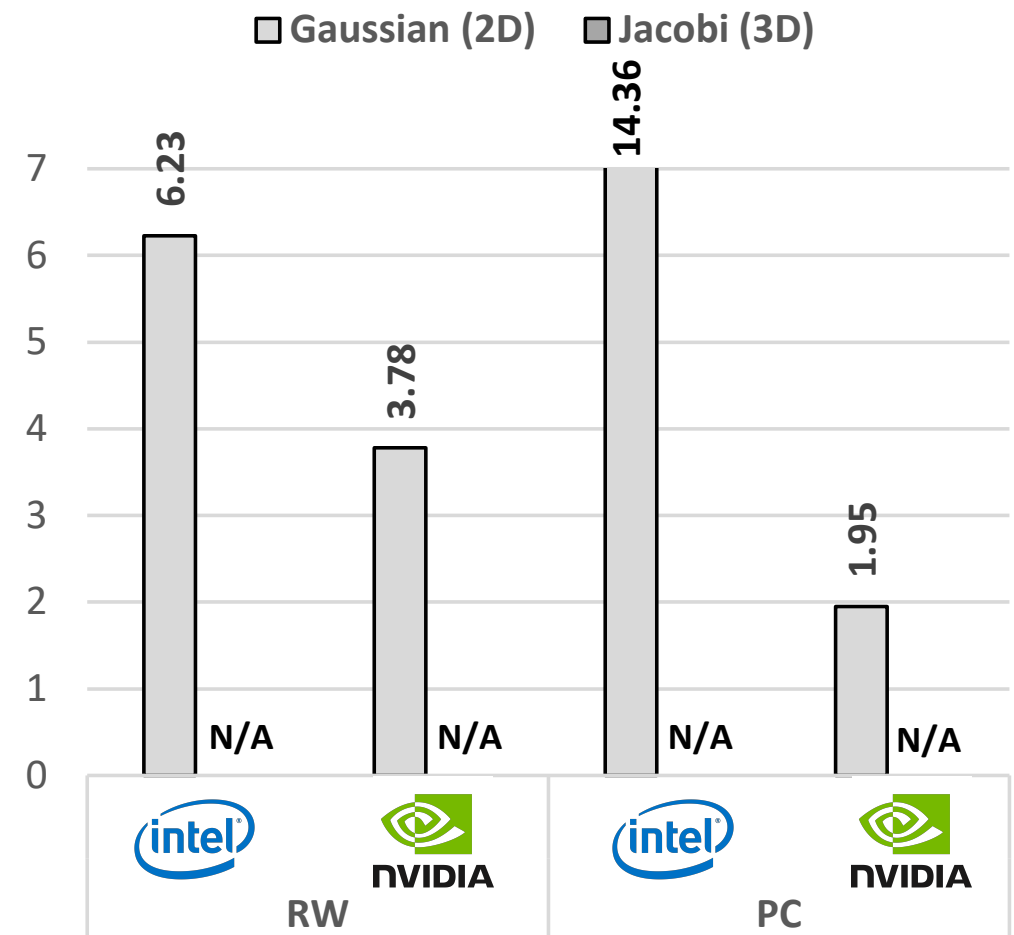
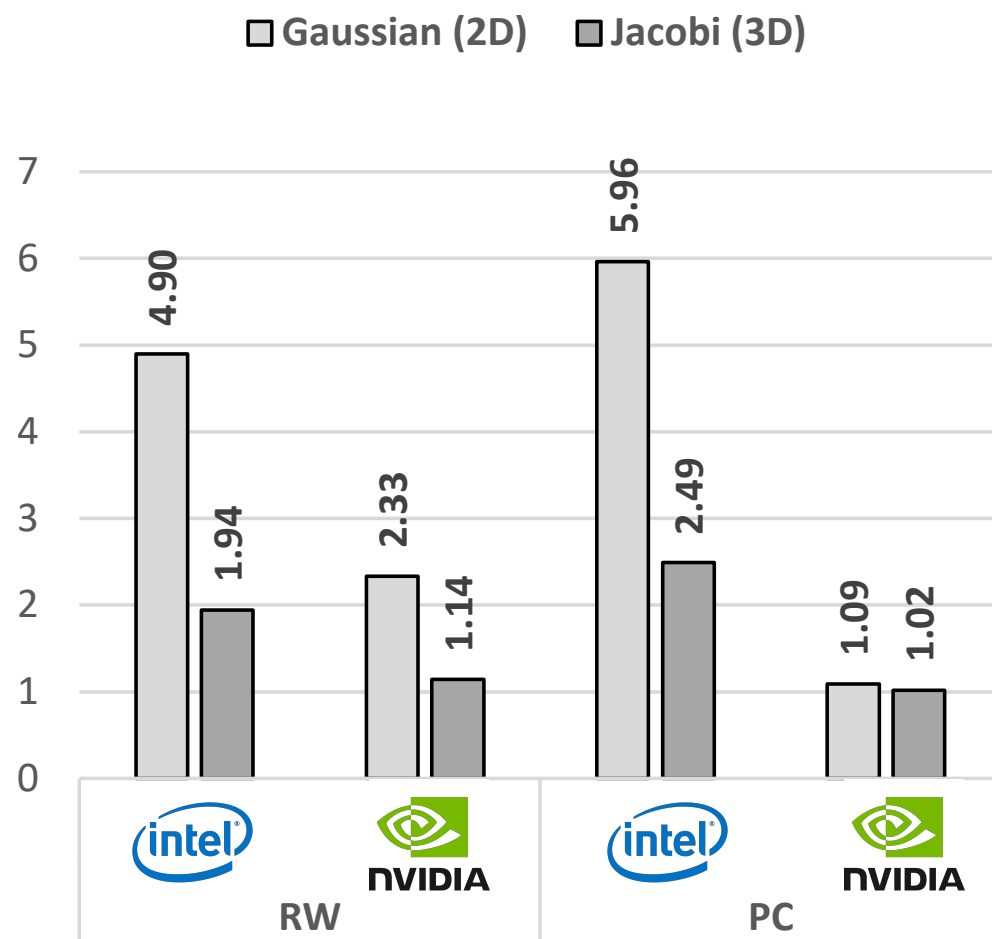
- Competitive performance for RW.
- Lower performance for PC because of assembly optimizations (e.g., CUDA tensor cores).

[1] Steuwer et. al, "Lift: A Functional Data-Parallel IR for High-Performance GPU Code Generation", CGO'17.

[2] Steuwer et. al, "Matrix Multiplication Beyond Auto-Tuning: Rewrite-Based GPU Code Generation.", CASES'16.

Experimental Results

2. Stencil Computations



Comparison to Lift [3]:

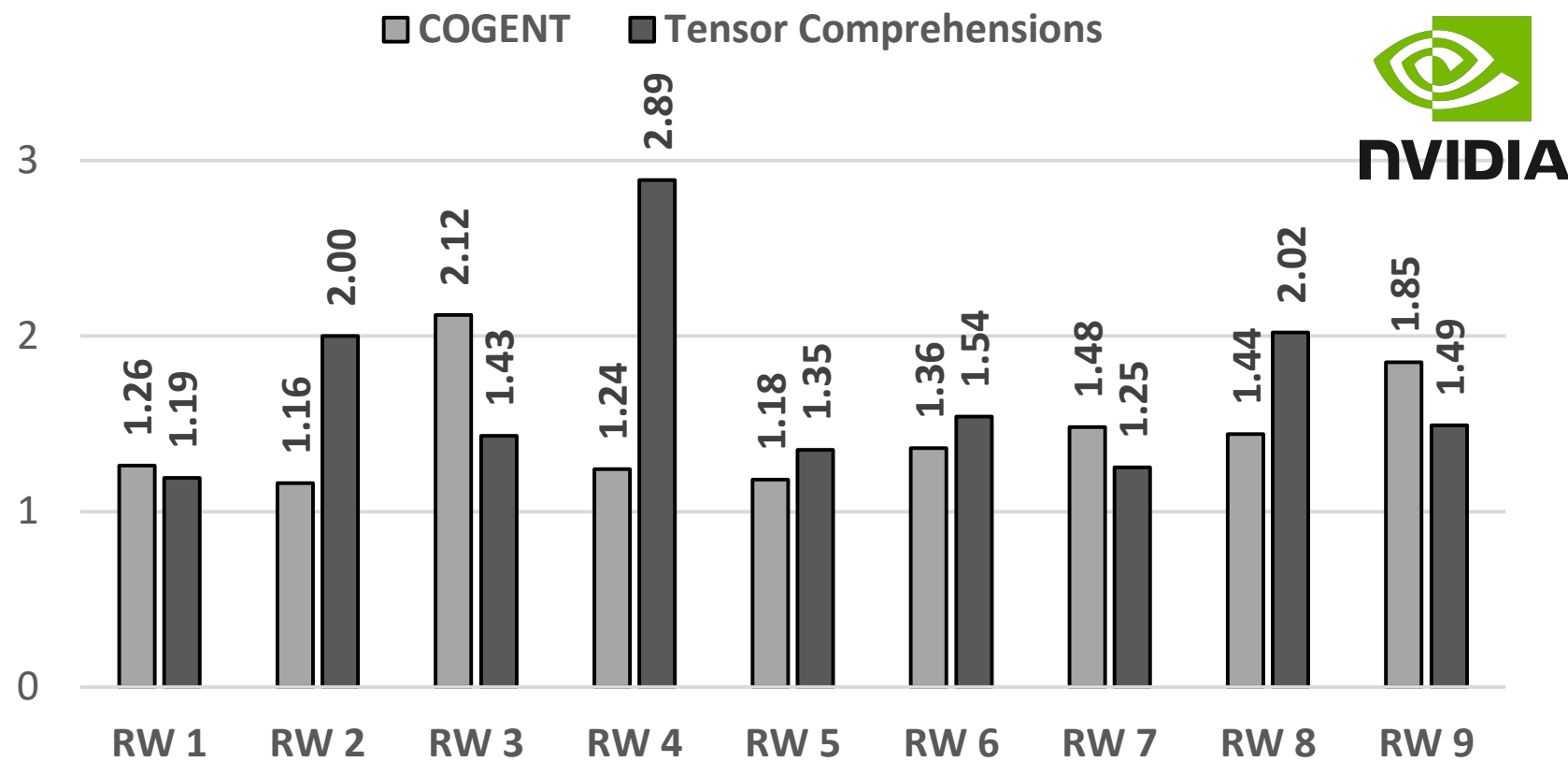
- Speedups up to **>5x**
- Our approach relies on generic optimizations + auto-tuning (rather than transformation rules)

Comparison to MKL-DNN/cuDNN:

- Better performance → MCC!
- MCC: speedup >2x over MKL-DNN
- MCC: speedup 0.5x over cuDNN

Experimental Results

3. Tensor Contractions:



Tensor Contractions on 9 real-world samples

- Speedups up to **>2x** over both competitors.
- We provide a more flexible implementations, e.g., tile sizes of also >1 (COGENT [1]).

Conclusion

We provide **performance, portability, and productivity** for MDH functions:

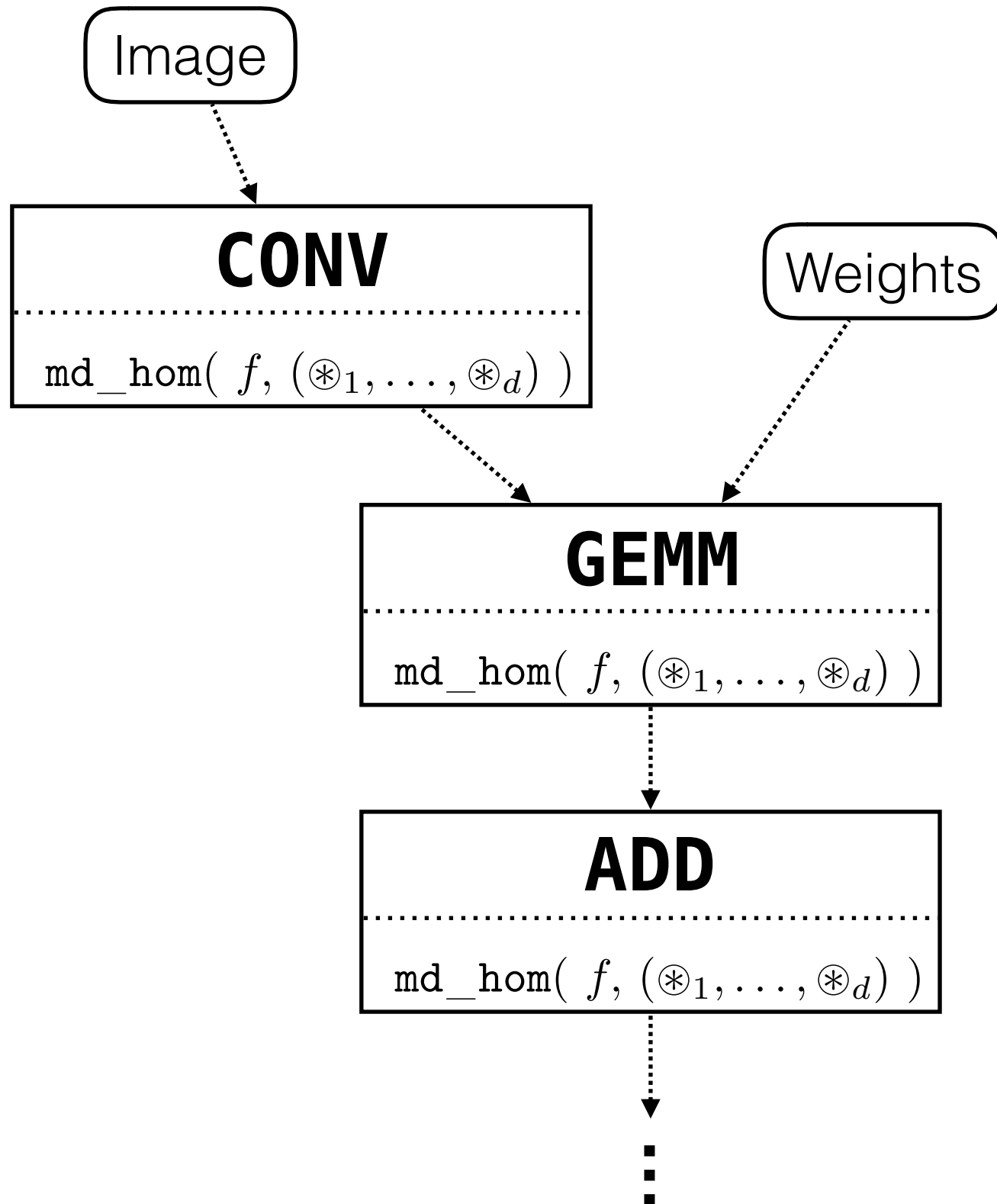
1. MDHs cover a **broad range of applications** (BLAS, Stencil, PRL, TC, ...).
2. Our implementation of MDHs provides **high performance** (speedups up to >5) on both CPU and GPU
3. MDHs are functionally and performance **portable** over architectures and input sizes.
4. MDHs can be **conveniently implemented** using our md_hom high-level abstraction.

Moreover:

- Our **Auto-Tuning Framework (ATF)** is a **general-purpose approach** that **supports auto-tuning of interdependent tuning parameters**.
- We provide our **dOCAL** framework for **conveniently executing OpenCL kernels on multi-device/multi-node systems**, and it automatically **provides asynchronous computation efficiency**.

Current/Future Work

Graph-Level Optimizations (e.g. as required for deep learning):



Exploiting md_hom representation for fusing (better locality)

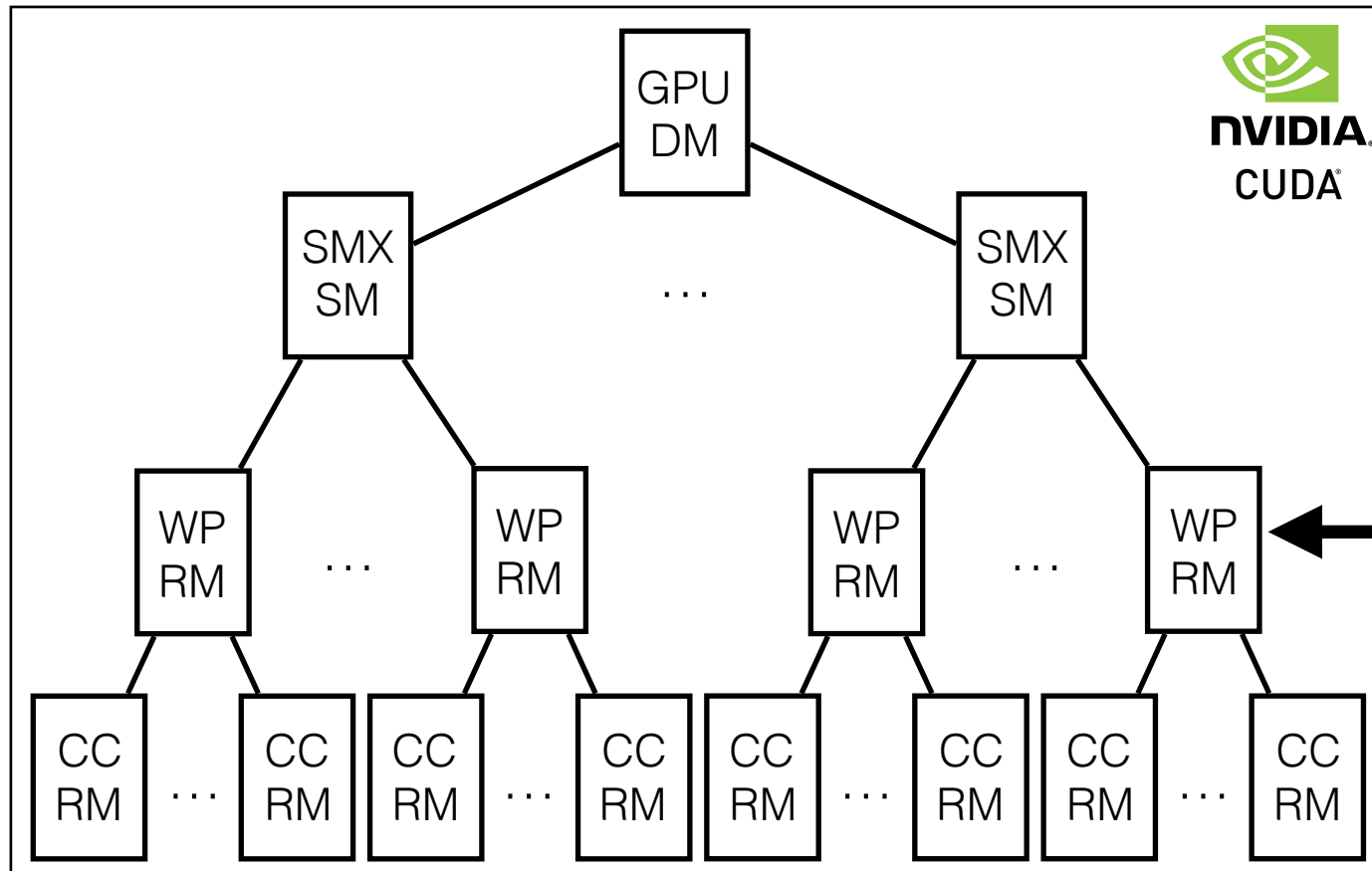
Fused Kernel

$\text{md_hom}(f, (\otimes_1, \dots, \otimes_d))$

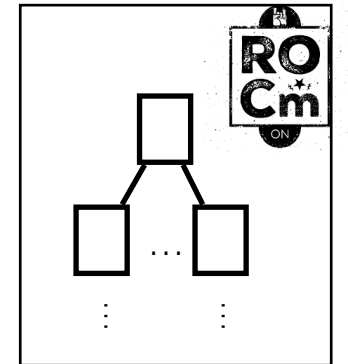
```
kernel void gemv_test(__global float* in_matrix,
                    __global float* in_vector,
                    __global float* out_vector,
                    // private memory for a WI's computation
                    __private float res_prv = 0.0f;
                    // local memory for a WG's computation
                    __local float res_lcl[ NUM_WI_1 ][ NUM_WI_2 ];
                    // iteration over P_sq blocks
                    for( int i_sq = 1; i_sq <= NUM_SO_1; ++i_sq ) {
                        for( int j_sq = 1; j_sq <= NUM_SO_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 WGs' results in global memory
                        if( WI_ID_2 == 0 )
                            my_res[ i_sq ] = res_lcl[ WI_ID_1 ][0];
                        barrier( CLK_LOCAL_MEM_FENCE );
                    } // end of for-loop j_sq
                } // end of for-loop i_sq
            } // end of kernel
```

Current/Future Work

Further backends:

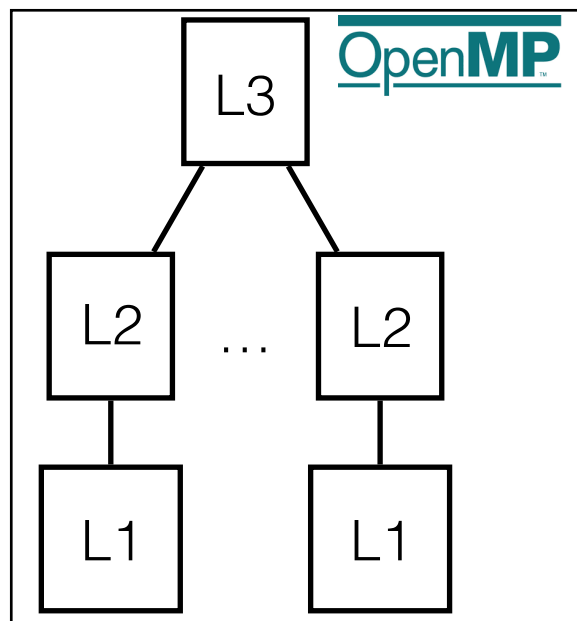


- Shuffle Operations
 - Tensor Cores
- comparable performance to cuBLAS

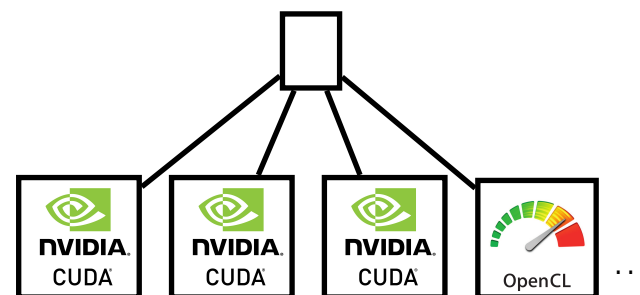


Assembler Backends

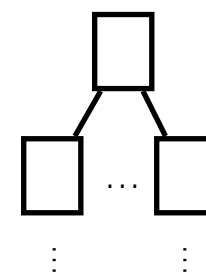
(e.g., NVIDIA PTX, Intel x86-64, ...)



Multi-Device



Multi-Node



Current/Future Work

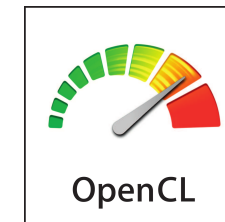
Parallelizing sequential C programs via MDH approach:

```
int main()
{
    // ...

    #pragma mdh ( , , +=:C[i][j] )
    for( int i = 0 ; i < M ; ++i )
        for( int j = 0 ; j < N ; ++j )
            for( int k = 0 ; k < K ; ++k )
            {
                C[ i ][ j ] += A[ i ][ k ] * B[ k ][ j ];
            }

    // ...
}
```

MDH
Compiler



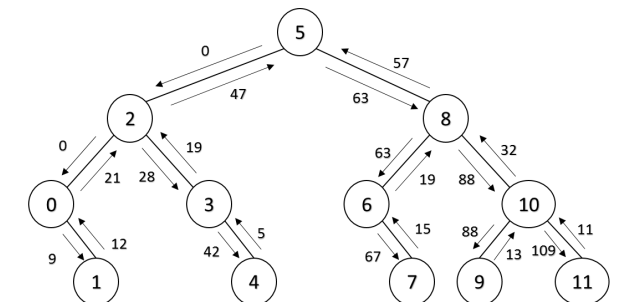
md_hom expression and view are automatically derived by MDH compiler:
md_hom(loop_body, (++ , ++ , +=:C[i][j]))
view(A,B)(i,j,k)=(A[i][k], B[k][j])

Annotating sequential C-Code with simple *MDH directives* (similarly as in *OpenMP/OpenACC*).

Current/Future Work

Demonstrating MDH approach's efficiency for more applications:

Benchmark	Typical Bottleneck of an Unoptimized Implementation	Optimizations Applied	Optimized Implementation Bottleneck	Potential Improvements
cutcp	Contention, Locality	Scatter-to-Gather, Binning, Regularization, Coarsening	Instruction Throughput	Minimizing Reads/Checks of Irrelevant Bin Data
mri-q	Poor Locality	Data Layout Transformation, Tiling, Coarsening	Instruction Throughput	
gridding	Contention, Load Imbalance	Scatter-to-Gather, Binning, Compaction, Regularization, Coarsening	Instruction Throughput	Minimizing Reads/Checks of Irrelevant Bin Data
sad	Locality	Tiling, Coarsening	Memory Bandwidth/Latency	Target Devices with Higher Register Capacities
stencil	Locality	Coarsening, Tiling	Bandwidth	
tpacf	Locality, Contention	Tiling, Privatization, Coarsening	Instruction Throughput	
lbn	Bandwidth	Data Layout Transformation	Bandwidth	
sgemm	Bandwidth	Coarsening, Tiling	Instruction Throughput	
spmv	Bandwidth	Data Layout Transformation	Bandwidth	
bfs	Contention, Load Imbalance	Privatization, Compaction, Regularization	Bandwidth	Avoiding Global Barriers / Better Kernels for Midsized Frontiers
histogram	Contention, Bandwidth	Privatization, Scatter-to-Gather	Bandwidth	Reducing Reads of Irrelevant Input (alleviated by cache)



Rank	0	1	2	3	4	5	6	7	8	9	10	11
n	9	12	7	14	5	16	4	15	6	13	8	11
S	9	21	28	42	47	63	67	82	88	101	109	120

Prefix sum

Parboil



Rodinia

TABLE I
RODINIA APPLICATIONS AND KERNELS (*DENOTES KERNEL).

Application / Kernel	Dwarf	Domain
K-means	Dense Linear Algebra	Data Mining
Needleman-Wunsch	Dynamic Programming	Bioinformatics
HotSpot*	Structured Grid	Physics Simulation
Back Propagation*	Unstructured Grid	Pattern Recognition
SRAD	Structured Grid	Image Processing
Leukocyte Tracking	Structured Grid	Medical Imaging
Breadth-First Search*	Graph Traversal	Graph Algorithms
Stream Cluster*	Dense Linear Algebra	Data Mining
Similarity Scores*	MapReduce	Web Mining

Questions?