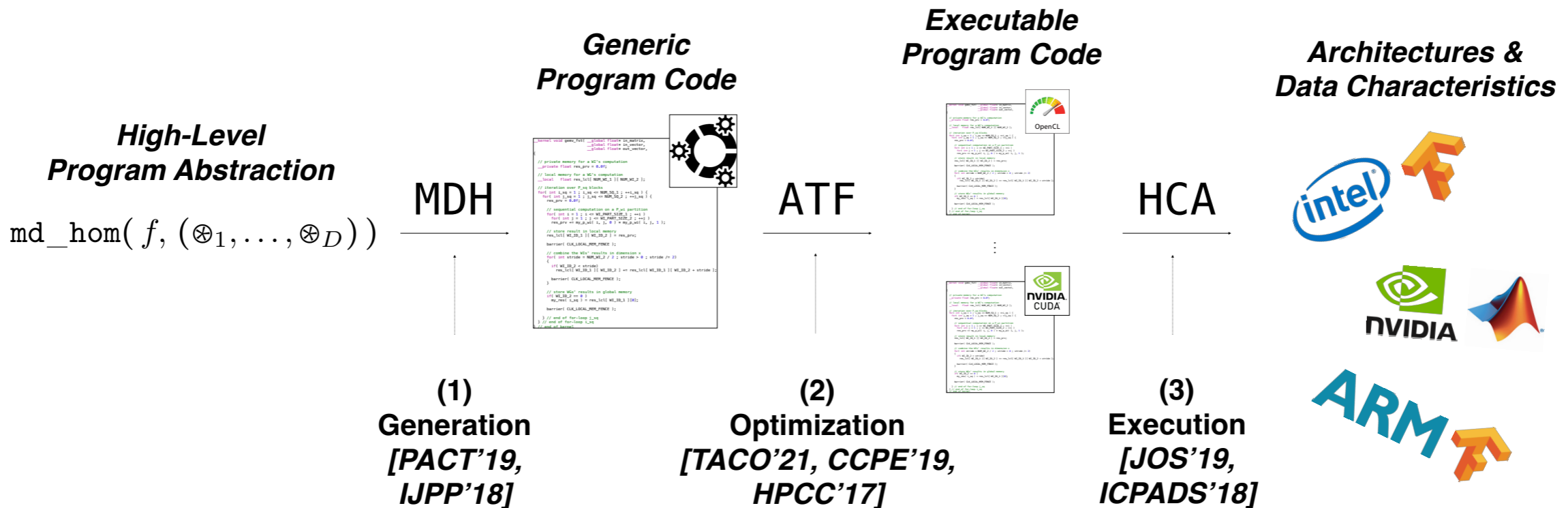# (De/Re)-Compositions Expressed Systematically via MDH-Based Schedules

Ari Rasch, Richard Schulze, Denys Shabalin,
Anne Elster, Sergei Gorlatch, Mary Hall
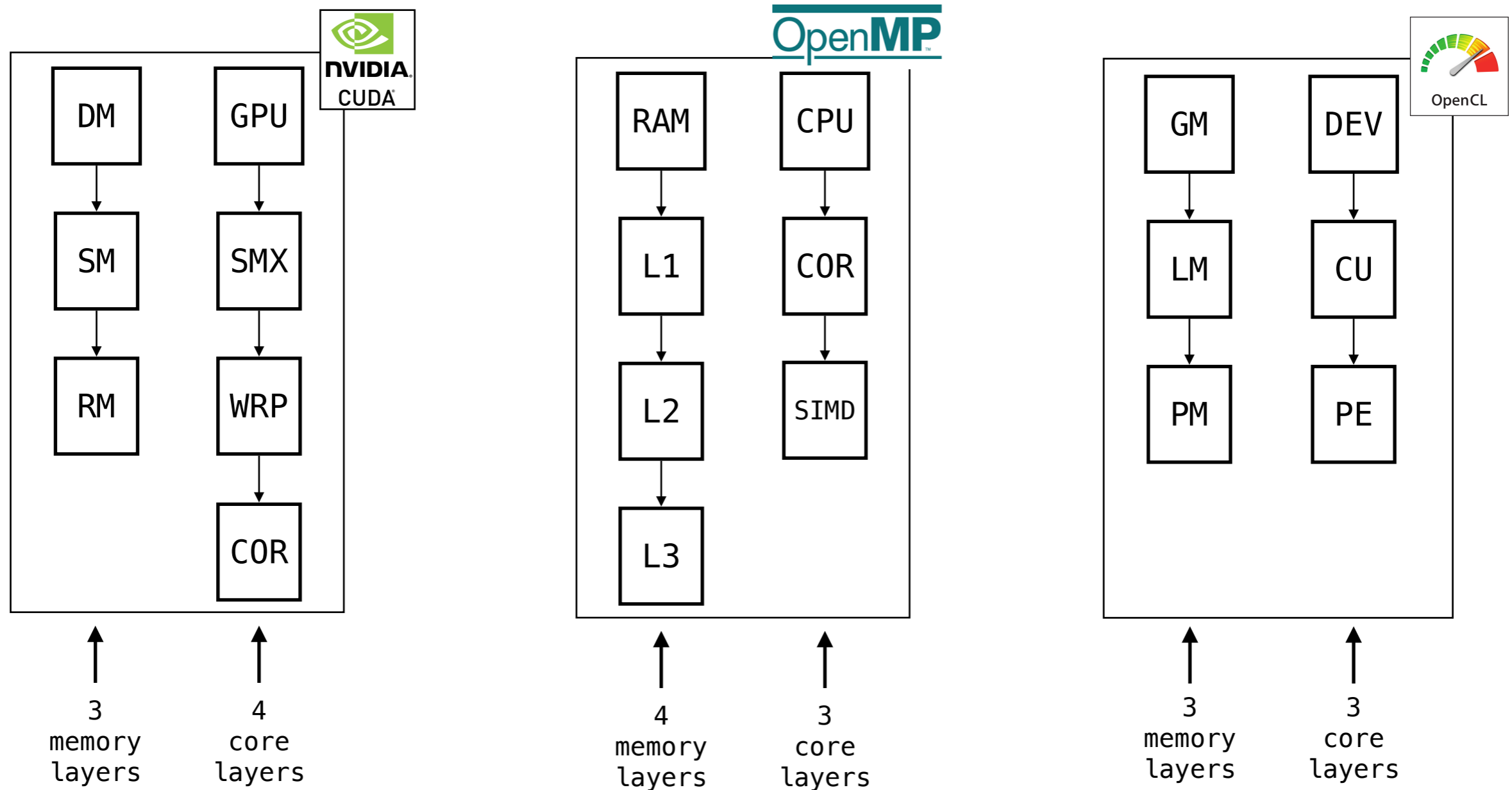
# The MDH(+ATF+HCA) Approaches



Approaches to code *generation* (MDH) & *optimization* (ATF) & *execution* (HCA):

(1) MDH (*Multi-Dimensional Homomorphisms*): How to generate automatically optimizable (auto-tunable) code?

(2) ATF (*Auto-Tuning Framework*): How to optimize (auto-tune) code?

(3) HCA (*Host Code Abstraction*): How to execute code on (distr.) multi-dev. systems?

# Observation

State-of-the-art architectures rely on deep *memory & core hierarchies*:



**Optimizations** are required for **both hierarchy kinds**
— *memory* and *core* —
to achieve the full performance potential of architectures

# Observation

- Modern high performance compilers include: TVM, Halide, …

- These compilers efficiently target modern architectures, by allowing expert users to explicitly express *code optimizations* in form of so-called *scheduling programs*

- Flaw:

   ***The existing scheduling languages usually rely on a vast set of low-level commands, and the commands have to be combined in complex ways to achieve high performance***

We show that this design decision of the existing approaches makes them expressive, but complicates:

1. achieving *high performance*

2. guaranteeing *safety*

3. offering *auto-tuning*

4. *enabling applicability*

5. allowing *visualization*

```
1   # exploiting fast memory resources for "C":
2   matmul_local, = s.cache_write([matmul], "local"
        )
3   matmul_1, matmul_2, matmul_3 = tuple(
        matmul_local.op.axis) + tuple(matmul_local.
        op.reduce_axis)
4   SHR_1, REG_1 = s[matmul_local].split(matmul_1,
        factor=1)
5   # 9 further split commands
6   s[matmul_local].reorder(BLK_1, BLK_2, DEV_1,
        DEV_2, THR_1, THR_2, DEV_3, SHR_3, SHR_1,
        SHR_2, REG_3, REG_1, REG_2)
7
8   # ... (loop unrolling)
9
10  # tiling:
11  matmul_1, matmul_2, matmul_3 = tuple(matmul.op.
        axis) + tuple(matmul.op.reduce_axis)
12  THR_1, SHR_REG_1 = s[matmul].split(matmul_1,
        factor=1)
13  # 5 further split commands
14  s[matmul].reorder(BLK_1, BLK_2, DEV_1, DEV_2,
        THR_1, THR_2, SHR_REG_1, SHR_REG_2)
15  s[matmul_local].compute_at(s[matmul], THR_2)
16
17  # block/thread assignments:
18  BLK_fused = s[matmul].fuse(BLK_1, BLK_2)
19  s[matmul].bind(BLK_fused, te.thread_axis("
        blockIdx.x"))
20  # ... (similar to lines 18 and 19)
21
22  # exploiting fast memory resources for "A":
23  A_shared = s.cache_read(A, "shared", [
        matmul_local])
24  A_shared_ax0, A_shared_ax1 = tuple(A_shared.op.
        axis)
25  A_shared_ax0_ax1_fused = s[A_shared].fuse(
        A_shared_ax0, A_shared_ax1)
26  A_shared_ax0_ax1_fused_o,
        A_shared_ax0_ax1_fused_i = s[A_shared].
        split(A_shared_ax0_ax1_fused, factor=1)
27  s[A_shared].vectorize(A_shared_ax0_ax1_fused_i)
28  # ...
29  s[A_shared].compute_at(s[matmul_local], DEV_3)
30
31  # exploiting fast memory resources for "B":
32  # ... (analogous to lines 23 − 29)
```

**Listing 3.** TVM+Ansor schedule (shortened for brevity) for Matrix Multiplication as used in ResNet-50 network on NVIDIA Ampere GPU

4

# Contribution of this Work

We introduce a *__new scheduling language__* for
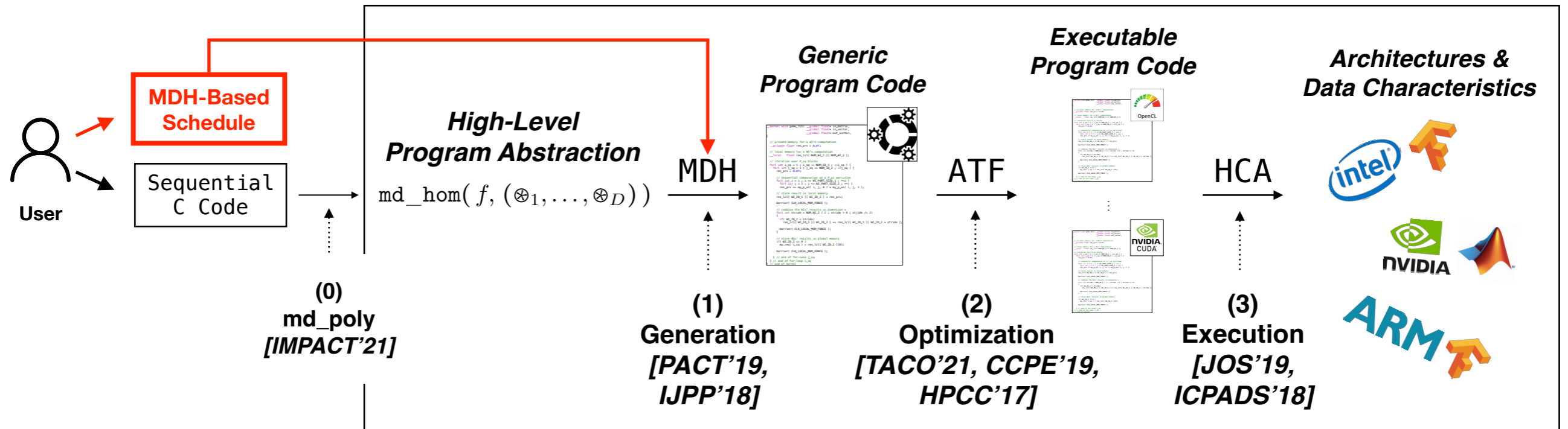**expressing code optimizations
in a systematic way**

Our systematic language design enables:

1. **Performance:** we can achieve (and often even outperform) the state-of-the-art TVM+Ansor compiler

2. **Safety:** we offer strong error checking, backed by MDH formalism

3. **Auto-Tuning:** any particular optimization decisions can be optionally left for auto-tuning (schedules can be recommended)

4. **Applicability:** our language is used analogously for multiple kinds of programming models (CUDA, OpenMP, OpenCL, …)

5. **Visualization:** our schedules can be visualized and also be generated from visual inputs

```
1   // initialization
2   0: (de/re)-comp( 16,1000,2048 )
3                   ( A:DM[1,2],B:DM[1,2] ;
4                     C:DM[1,2]           )
5                   ( GPU.y,GPU.x,GPU.z )
6
7   // parallelization over CUDA Blocks
8   1: (de/re)-comp( 8,20,^ )
9                   ( ^,^ ; ^ )
10                  ( BLK.y,BLK.x,BLK.z )
11
12  // tiling 1
13  6: (de/re)-comp( 4,^,^ )
14                  ( ^,^ ; ^ )
15                  ( FOR.1,FOR.2,FOR.3 )
16
17  // parallelization over CUDA Threads &
18  // utilization of CUDA Register Memory
19  2: (de/re)-comp( 1,1,^ )
20                  ( ^,^ ; C:RM[1,2] )
21                  ( THR.y,THR.x,THR.z )
22
23  // utilization of CUDA Shared Memory
24  3: (de/re)-comp( ^,^,256 )
25                  ( A:SM[1,2],B:SM[1,2] ; ^ )
26                  ( FOR.2,FOR.3,FOR.1 )
27
28  // tiling 2
29  4: (de/re)-comp( ^,^,2 )
30                  ( ^,^ ; ^ )
31                  ( ^,^,^ )
32
33  // tiling 3
34  5: (de/re)-comp( ^,^,1 )
35                  ( ^,^ ; ^ )
36                  ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

# Overview



**In this work:**

We extend the existing MDH+ATF+HCA pipeline,

by allowing expert users to explicitly express some/all optimizations

via **MDH-Based Schedules**

Advantages over existing MDH+ATF+HCA:

1. Better Optimization: an auto-tuning system might not always make the same high-quality optimization decisions as an expert user

2. Faster Auto-Tuning: as some (or even all) optimization decisions are made by the expert user and thus are not left to the auto-tuning system

6

# MDH-Based Schedules

- We allow expert users to express optimizations via *MDH-Based Schedules*

- Our language consists of exactly one primitive which has the following basic structure:

  ```
  (de-)comp( /* sub-problem size */ )
           ( /* memory hierarchy assignments */ )
           ( /* core    hierarchy assignments */ )
  ```

- We illustrate our primitive using the example of Matrix Multiplication on 16x2048 & 2048x1000 matrices (taken from ResNet-50):

```
 1   // initialization
 2   0: (de/re)-comp( 16,1000,2048 )
 3                  ( A:DM[1,2],B:DM[1,2] ;
 4                    C:DM[1,2]           )
 5                  ( GPU.y,GPU.x,GPU.z )
 6
 7   // parallelization over CUDA Blocks
 8   1: (de/re)-comp( 8,20,^ )
 9                  ( ^,^ ; ^ )
10                  ( BLK.y,BLK.x,BLK.z )
11
12   // tiling 1
13   6: (de/re)-comp( 4,^,^ )
14                  ( ^,^ ; ^ )
15                  ( FOR.1,FOR.2,FOR.3 )
16
17   // parallelization over CUDA Threads &
18   // utilization of CUDA Register Memory
19   2: (de/re)-comp( 1,1,^ )
20                  ( ^,^ ; C:RM[1,2] )
21                  ( THR.y,THR.x,THR.z )
22
23   // utilization of CUDA Shared Memory
24   3: (de/re)-comp( ^,^,256 )
25                  ( A:SM[1,2],B:SM[1,2] ; ^ )
26                  ( FOR.2,FOR.3,FOR.1 )
27
28   // tiling 2
29   4: (de/re)-comp( ^,^,2 )
30                  ( ^,^ ; ^ )
31                  ( ^,^,^ )
32
33   // tiling 3
34   5: (de/re)-comp( ^,^,1 )
35                  ( ^,^ ; ^ )
36                  ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

# MDH-Based Schedules

- We allow expert users to express optimizations via *MDH-Based Schedules*

- Our language consists of exactly one primitive which has the following basic structure:

```
(de-)comp( /* sub-problem size */ )
         ( /* memory hierarchy assignments */ )
         ( /* core   hierarchy assignments */ )
```

- We illustrate our primitive using the example of Matrix Multiplication on 16x2048 & 2048x1000 matrices (taken from ResNet-50):

**Initialization (optional):**

— the initial iteration space has a size of **16,1000,2048**

— *CUDA's Device Memory* (**DM**) is used for A & B input matrices and C output matrix

— computation is performed by a **GPU**

```
 1  // initialization
 2  0: (de/re)-comp( 16,1000,2048 )
 3               ( A:DM[1,2],B:DM[1,2] ;
 4                 C:DM[1,2]           )
 5               ( GPU.y,GPU.x,GPU.z )
 6
 7  // parallelization over CUDA Blocks
 8  1: (de/re)-comp( 8,20,^ )
 9               ( ^,^ ; ^ )
10               ( BLK.y,BLK.x,BLK.z )
11
12  // tiling 1
13  6: (de/re)-comp( 4,^,^ )
14               ( ^,^ ; ^ )
15               ( FOR.1,FOR.2,FOR.3 )
16
17  // parallelization over CUDA Threads &
18  // utilization of CUDA Register Memory
19  2: (de/re)-comp( 1,1,^ )
20               ( ^,^ ; C:RM[1,2] )
21               ( THR.y,THR.x,THR.z )
22
23  // utilization of CUDA Shared Memory
24  3: (de/re)-comp( ^,^,256 )
25               ( A:SM[1,2],B:SM[1,2] ; ^ )
26               ( FOR.2,FOR.3,FOR.1 )
27
28  // tiling 2
29  4: (de/re)-comp( ^,^,2 )
30               ( ^,^ ; ^ )
31               ( ^,^,^ )
32
33  // tiling 3
34  5: (de/re)-comp( ^,^,1 )
35               ( ^,^ ; ^ )
36               ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

# MDH-Based Schedules

- We allow expert users to express optimizations via *MDH-Based Schedules*

- Our language consists of exactly one primitive which has the following basic structure:

```
(de-)comp( /* sub-problem size */ )
        ( /* memory hierarchy assignments */ )
        ( /* core   hierarchy assignments */ )
```

- We illustrate our primitive using the example of Matrix Multiplication on 16x2048 & 2048x1000 matrices (taken from ResNet-50):

**Block Parallelization:**

— iteration space is split into tiles of size **8,20,2048**

— no memory optimizations

— each tile is computed by a *CUDA Block* (**BLK**)

```
 1  // initialization
 2  0: (de/re)-comp( 16,1000,2048 )
 3                ( A:DM[1,2],B:DM[1,2] ;
 4                  C:DM[1,2]           )
 5                ( GPU.y,GPU.x,GPU.z )
 6
 7  // parallelization over CUDA Blocks
 8  1: (de/re)-comp( 8,20,^ )
 9                ( ^,^ ; ^ )
10                ( BLK.y,BLK.x,BLK.z )
11
12  // tiling 1
13  6: (de/re)-comp( 4,^,^ )
14                ( ^,^ ; ^ )
15                ( FOR.1,FOR.2,FOR.3 )
16
17  // parallelization over CUDA Threads &
18  // utilization of CUDA Register Memory
19  2: (de/re)-comp( 1,1,^ )
20                ( ^,^ ; C:RM[1,2] )
21                ( THR.y,THR.x,THR.z )
22
23  // utilization of CUDA Shared Memory
24  3: (de/re)-comp( ^,^,256 )
25                ( A:SM[1,2],B:SM[1,2] ; ^ )
26                ( FOR.2,FOR.3,FOR.1 )
27
28  // tiling 2
29  4: (de/re)-comp( ^,^,2 )
30                ( ^,^ ; ^ )
31                ( ^,^,^ )
32
33  // tiling 3
34  5: (de/re)-comp( ^,^,1 )
35                ( ^,^ ; ^ )
36                ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

9

# MDH-Based Schedules

- We allow expert users to express optimizations via *MDH-Based Schedules*

- Our language consists of exactly one primitive which has the following basic structure:

```
(de-)comp( /* sub-problem size */ )
         ( /* memory hierarchy assignments */ )
         ( /* core   hierarchy assignments */ )
```

- We illustrate our primitive using the example of Matrix Multiplication on 16x2048 & 2048x1000 matrices (taken from ResNet-50):

**Classical Tiling:**

— iteration space is split into tiles of size **4,20,2048**

— no memory optimizations

— no parallelization

```
1    // initialization
2    0: (de/re)-comp( 16,1000,2048 )
3                  ( A:DM[1,2],B:DM[1,2] ;
4                    C:DM[1,2]          )
5                  ( GPU.y,GPU.x,GPU.z )
6
7    // parallelization over CUDA Blocks
8    1: (de/re)-comp( 8,20,^ )
9                  ( ^,^ ; ^ )
10                 ( BLK.y,BLK.x,BLK.z )
11
12   // tiling 1
13   6: (de/re)-comp( 4,^,^ )
14                 ( ^,^ ; ^ )
15                 ( FOR.1,FOR.2,FOR.3 )
16
17   // parallelization over CUDA Threads &
18   // utilization of CUDA Register Memory
19   2: (de/re)-comp( 1,1,^ )
20                 ( ^,^ ; C:RM[1,2] )
21                 ( THR.y,THR.x,THR.z )
22
23   // utilization of CUDA Shared Memory
24   3: (de/re)-comp( ^,^,256 )
25                 ( A:SM[1,2],B:SM[1,2] ; ^ )
26                 ( FOR.2,FOR.3,FOR.1 )
27
28   // tiling 2
29   4: (de/re)-comp( ^,^,2 )
30                 ( ^,^ ; ^ )
31                 ( ^,^,^ )
32
33   // tiling 3
34   5: (de/re)-comp( ^,^,1 )
35                 ( ^,^ ; ^ )
36                 ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

10

# MDH-Based Schedules

- We allow expert users to express optimizations via *MDH-Based Schedules*

- Our language consists of exactly one primitive which has the following basic structure:

```
(de-)comp( /* sub-problem size */ )
         ( /* memory hierarchy assignments */ )
         ( /* core    hierarchy assignments */ )
```

- We illustrate our primitive using the example of Matrix Multiplication on 16x2048 & 2048x1000 matrices (taken from ResNet-50):

**Thread Parallelization & Register Memory Utilization:**

– iteration space is split into tiles of size **1,1,2048**

– *CUDA Register Memory* (**RM**) is used for computed intermediate results of C output matrix

– each tile is computed by a *CUDA Thread* (**THR**)

```
 1  // initialization
 2  0: (de/re)-comp( 16,1000,2048 )
 3               ( A:DM[1,2],B:DM[1,2] ;
 4                 C:DM[1,2]           )
 5               ( GPU.y,GPU.x,GPU.z )
 6
 7  // parallelization over CUDA Blocks
 8  1: (de/re)-comp( 8,20,^ )
 9               ( ^,^ ; ^ )
10               ( BLK.y,BLK.x,BLK.z )
11
12  // tiling 1
13  6: (de/re)-comp( 4,^,^ )
14               ( ^,^ ; ^ )
15               ( FOR.1,FOR.2,FOR.3 )
16
17  // parallelization over CUDA Threads &
18  // utilization of CUDA Register Memory
19  2: (de/re)-comp( 1,1,^ )
20               ( ^,^ ; C:RM[1,2] )
21               ( THR.y,THR.x,THR.z )
22
23  // utilization of CUDA Shared Memory
24  3: (de/re)-comp( ^,^,256 )
25               ( A:SM[1,2],B:SM[1,2] ; ^ )
26               ( FOR.2,FOR.3,FOR.1 )
27
28  // tiling 2
29  4: (de/re)-comp( ^,^,2 )
30               ( ^,^ ; ^ )
31               ( ^,^,^ )
32
33  // tiling 3
34  5: (de/re)-comp( ^,^,1 )
35               ( ^,^ ; ^ )
36               ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

11

# MDH-Based Schedules

- We allow expert users to express optimizations via *MDH-Based Schedules*

- Our language consists of exactly one primitive which has the following basic structure:

```
(de-)comp( /* sub-problem size */ )
         ( /* memory hierarchy assignments */ )
         ( /* core   hierarchy assignments */ )
```

- We illustrate our primitive using the example of Matrix Multiplication on 16x2048 & 2048x1000 matrices (taken from ResNet-50):

**Shared Memory Utilization:**

— iteration space is split into tiles of size **1,1,256**

— *CUDA Shared Memory* (**SM**) is used for A & B input matrices

— no parallelization

```
 1   // initialization
 2   0: (de/re)-comp( 16,1000,2048 )
 3                  ( A:DM[1,2],B:DM[1,2] ;
 4                    C:DM[1,2]           )
 5                  ( GPU.y,GPU.x,GPU.z )
 6
 7   // parallelization over CUDA Blocks
 8   1: (de/re)-comp( 8,20,^ )
 9                  ( ^,^ ; ^ )
10                  ( BLK.y,BLK.x,BLK.z )
11
12   // tiling 1
13   6: (de/re)-comp( 4,^,^ )
14                  ( ^,^ ; ^ )
15                  ( FOR.1,FOR.2,FOR.3 )
16
17   // parallelization over CUDA Threads &
18   // utilization of CUDA Register Memory
19   2: (de/re)-comp( 1,1,^ )
20                  ( ^,^ ; C:RM[1,2] )
21                  ( THR.y,THR.x,THR.z )
22
23   // utilization of CUDA Shared Memory
24   3: (de/re)-comp( ^,^,256 )
25                  ( A:SM[1,2],B:SM[1,2] ; ^ )
26                  ( FOR.2,FOR.3,FOR.1 )
27
28   // tiling 2
29   4: (de/re)-comp( ^,^,2 )
30                  ( ^,^ ; ^ )
31                  ( ^,^,^ )
32
33   // tiling 3
34   5: (de/re)-comp( ^,^,1 )
35                  ( ^,^ ; ^ )
36                  ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

# MDH-Based Schedules

- We allow expert users to express optimizations via *MDH-Based Schedules*

- Our language consists of exactly one primitive which has the following basic structure:

```
(de-)comp( /* sub-problem size */ )
         ( /* memory hierarchy assignments */ )
         ( /* core   hierarchy assignments */ )
```

- We illustrate our primitive using the example of Matrix Multiplication on 16x2048 & 2048x1000 matrices (taken from ResNet-50):

**Classical Tiling:**

– iteration space is split in tiles of size **1,1,2**

– no memory optimizations

– no parallelization

```
1   // initialization
2   0: (de/re)-comp( 16,1000,2048 )
3                 ( A:DM[1,2],B:DM[1,2] ;
4                   C:DM[1,2]           )
5                 ( GPU.y,GPU.x,GPU.z )
6
7   // parallelization over CUDA Blocks
8   1: (de/re)-comp( 8,20,^ )
9                 ( ^,^ ; ^ )
10                ( BLK.y,BLK.x,BLK.z )
11
12  // tiling 1
13  6: (de/re)-comp( 4,^,^ )
14                ( ^,^ ; ^ )
15                ( FOR.1,FOR.2,FOR.3 )
16
17  // parallelization over CUDA Threads &
18  // utilization of CUDA Register Memory
19  2: (de/re)-comp( 1,1,^ )
20                ( ^,^ ; C:RM[1,2] )
21                ( THR.y,THR.x,THR.z )
22
23  // utilization of CUDA Shared Memory
24  3: (de/re)-comp( ^,^,256 )
25                ( A:SM[1,2],B:SM[1,2] ; ^ )
26                ( FOR.2,FOR.3,FOR.1 )
27
28  // tiling 2
29  4: (de/re)-comp( ^,^,2 )
30                ( ^,^ ; ^ )
31                ( ^,^,^ )
32
33  // tiling 3
34  5: (de/re)-comp( ^,^,1 )
35                ( ^,^ ; ^ )
36                ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

# MDH-Based Schedules

- We allow expert users to express optimizations via *MDH-Based Schedules*

- Our language consists of exactly one primitive which has the following basic structure:

```
(de-)comp( /* sub-problem size */ )
        ( /* memory hierarchy assignments */ )
        ( /* core    hierarchy assignments */ )
```

- We illustrate our primitive using the example of Matrix Multiplication on 16x2048 & 2048x1000 matrices (taken from ResNet-50):

**Classical Tiling:**

– iteration space is split in tiles of size **1,1,1**

– no memory optimizations

– no parallelization

```
1  // initialization
2  0: (de/re)-comp( 16,1000,2048 )
3                  ( A:DM[1,2],B:DM[1,2] ;
4                    C:DM[1,2]           )
5                  ( GPU.y,GPU.x,GPU.z )
6
7  // parallelization over CUDA Blocks
8  1: (de/re)-comp( 8,20,^ )
9                  ( ^,^ ; ^ )
10                 ( BLK.y,BLK.x,BLK.z )
11
12 // tiling 1
13 6: (de/re)-comp( 4,^,^ )
14                 ( ^,^ ; ^ )
15                 ( FOR.1,FOR.2,FOR.3 )
16
17 // parallelization over CUDA Threads &
18 // utilization of CUDA Register Memory
19 2: (de/re)-comp( 1,1,^ )
20                 ( ^,^ ; C:RM[1,2] )
21                 ( THR.y,THR.x,THR.z )
22
23 // utilization of CUDA Shared Memory
24 3: (de/re)-comp( ^,^,256 )
25                 ( A:SM[1,2],B:SM[1,2] ; ^ )
26                 ( FOR.2,FOR.3,FOR.1 )
27
28 // tiling 2
29 4: (de/re)-comp( ^,^,2 )
30                 ( ^,^ ; ^ )
31                 ( ^,^,^ )
32
33 // tiling 3
34 5: (de/re)-comp( ^,^,1 )
35                 ( ^,^ ; ^ )
36                 ( ^,^,^ )
```

**Listing 5.** MDH-based schedule for optimizing matrix multiplication on NVIDIA A100 GPU according to the optimization decisions of TVM+Ansor in Listing 3

# Language Features

We show that our systematic language design enables:

1. **Performance**

2. **Safety**

3. **Auto-Tuning**

4. **Applicability**

5. **Visualization**

# Language Features

1. **Performance:** "*Deep Learning*" case study (TVM's favorable application class!)

## NVIDIA Ampere GPU

| Deep Learning | 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 | 1.00 | 1.42 | 1.00 | 1.14 | 1.00 | 1.00 |
| NVIDIA cuDNN | 0.92 | – | 1.85 | – | 1.22 | – | 1.94 | – | 1.81 | 2.14 |
| NVIDIA cuBLAS | – | 1.58 | – | 2.67 | – | 0.93 | – | 1.04 | – | – |

## NVIDIA Volta GPU

| Deep Learning | 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 | 1.21 | 1.00 | 1.79 | 1.00 | 1.11 | 1.06 | 1.00 | 1.00 | 1.00 |
| NVIDIA cuDNN | 1.21 | – | 1.29 | – | 2.80 | – | 3.50 | – | 2.32 | 3.14 |
| NVIDIA cuBLAS | – | 1.33 | – | 1.14 | – | 1.09 | – | 1.04 | – | – |

## Intel Skylake CPU

| Deep Learning | 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 |
| Intel oneDNN | 0.39 | – | 5.07 | – | 1.22 | – | 9.01 | – | 1.05 | 4.20 |
| Intel oneMKL | – | 0.44 | – | 1.09 | – | 0.88 | – | 0.53 | – | – |
| TVM+Ansor (DeVM) | 1.20 | 0.67 | 0.90 | 0.26 | 1.42 | 0.76 | 0.66 | 0.76 | 0.56 | 0.36 |

## Intel Broadwell CPU

| Deep Learning | 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 | 1.60 | 1.29 | 1.53 | 1.32 | 1.00 | 1.27 | 1.02 | 2.42 | 1.92 |
| Intel oneDNN | 1.30 | – | 1.81 | – | 2.94 | – | 2.85 | – | 1.83 | 4.47 |
| Intel oneMKL | – | 1.45 | – | 1.36 | – | 1.35 | – | 0.50 | – | – |
| TVM+Ansor (LLVM) | 1.45 | 1.35 | 1.06 | 0.72 | 1.63 | 0.85 | 0.98 | 0.79 | 1.14 | 0.52 |

*Speedup (higher is better) of our approach over TVM+Ansor*

**We achieve competitive and often higher performance than TVM+Ansor**

- TVM [OSDI'18] is a state-of-the-art compiler based on schedules

- We use for TVM the schedules generated by Ansor [OSDI'20]

- We report performance for approach using the schedules automatically recommended by our system

- Better performance of our approach is because our language:
  - ‣ i) supports *more optimization* (e.g., data layout changes);
  - ‣ ii) has *more potential for auto-tuning* (discussed on next slides)

# Language Features

2. **Safety:** We **formally guarantee correctness** of our scheduling programs, by checking the formal constraints defined by the MDH formalism

### CUDA

- Tile Size on lower layer <= Tile Size on upper layer

- BLKs combine in `{DM}`
- WRPs combine in `{SM,DM}`
- THRs combine in `{RM,SM,DM}`

- Number of THRs limited

- as well as:
  - `BLK/THR.{x,y,z}` can be used only once
  - `(de/re)-comp` order must be permutation
  - …

### OpenCL

- Tile Size on lower layer <= Tile Size on upper layer

- WGs combine in `{GM}`
- WIs combine in `{SM,GM}`

- Number of `WIs` limited

- as well as:
  - `WG/WI.{1,2,3,…}` can be used only once
  - `(de/re)-comp` order must be permutation
  - …

### OpenMP

…

**In related approaches, e.g., Fireiron [PACT'20], it is possible to implement schedules from which incorrect low-level code is generated, without issued error messages**

# Language Features

3.  **Auto-Tuning:** our language is designed such that optimizations can be left for auto-tuning (via symbol "?")

```
(de-)comp( 32,32,32 )
        ( A:SM[1,2] , B:SM[1,2] ;
          C:RM[1,2]              )
        ( BLK.x , BLK.y , BLK.z )
```

**No
tuning**

```
(de-)comp( 32,?,? )
        ( A:SM[1,2] , B:SM[1,2] ;
          C:RM[1,2]              )
        ( BLK.x , BLK.y , BLK.z )
```

**tile size
tuning**

```
(de-)comp( 32,?,? )
        ( A:?[?,?] , B:SM[1,2] ;
          C:?[1,2]              )
        ( BLK.x , BLK.y , BLK.z )
```
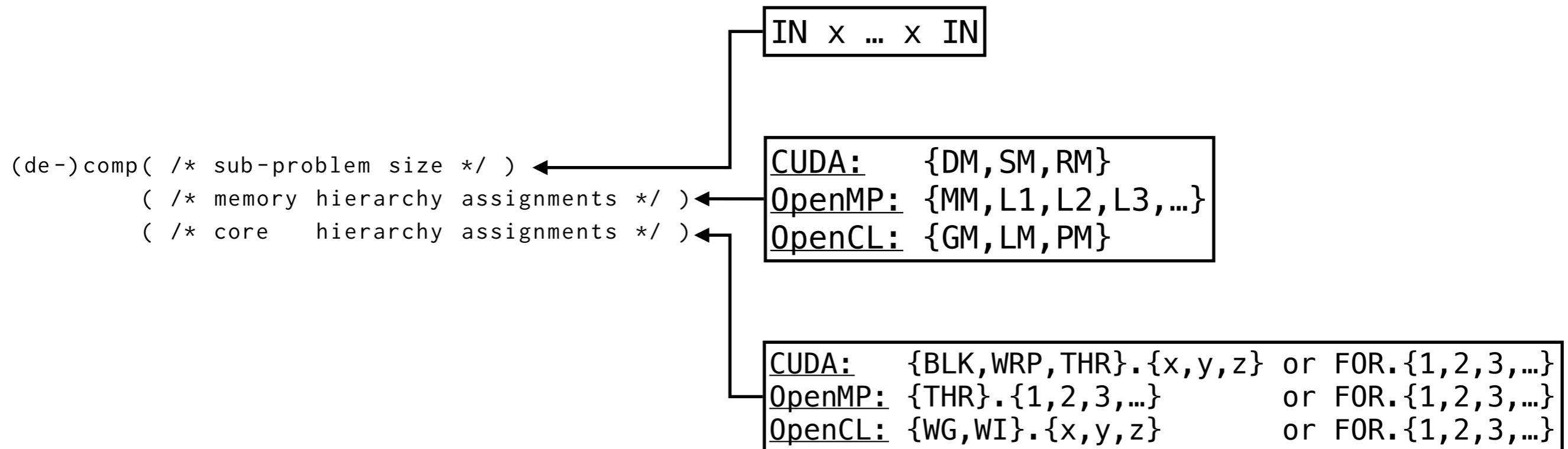
**tile size & memory
tuning**

```
(de-)comp( 32,?,? )
        ( A:?[?,?] , B:SM[1,2] ;
          C:?[1,2]              )
        ( ?.? , BLK.y , ?.? )
```

**tile size & memory & parallelization
tuning**

- **Our language is designed such that <u>any(!)</u> optimization decision can be left for auto-tuning.**

- **In contrast, the language design of other approaches (such as TVM) support auto-tuning for <u>some</u> optimizations (e.g, choosing tile size values), <u>but not for others</u> (e.g., binding parallelization to inner/outer tiles, using fast memory regions or not, etc).**

# Language Features

4. **Applicability:** our language is used similarly for different kinds of programming models

```
                                       ┌─────────────┐
                                       │ IN x … x IN │
                                       └─────────────┘

                                       ┌──────────────────────────┐
(de-)comp( /* sub-problem size */ ) ◄──│ CUDA:   {DM,SM,RM}       │
         ( /* memory hierarchy assignments */ ) ◄─│ OpenMP: {MM,L1,L2,L3,…} │
         ( /* core    hierarchy assignments */ ) ◄─│ OpenCL: {GM,LM,PM}       │
                                       └──────────────────────────┘

                                       ┌─────────────────────────────────────────────┐
                                       │ CUDA:   {BLK,WRP,THR}.{x,y,z} or FOR.{1,2,3,…} │
                                       │ OpenMP: {THR}.{1,2,3,…}        or FOR.{1,2,3,…} │
                                       │ OpenCL: {WG,WI}.{x,y,z}        or FOR.{1,2,3,…} │
                                       └─────────────────────────────────────────────┘
```

**Our system can be used/extended for C-based programming models targeting arbitrarily deep memory & core hierarchies**

5. **Visualization:** our schedules can be visualized & also be generated from visual inputs



0. (De/Re)-Composition
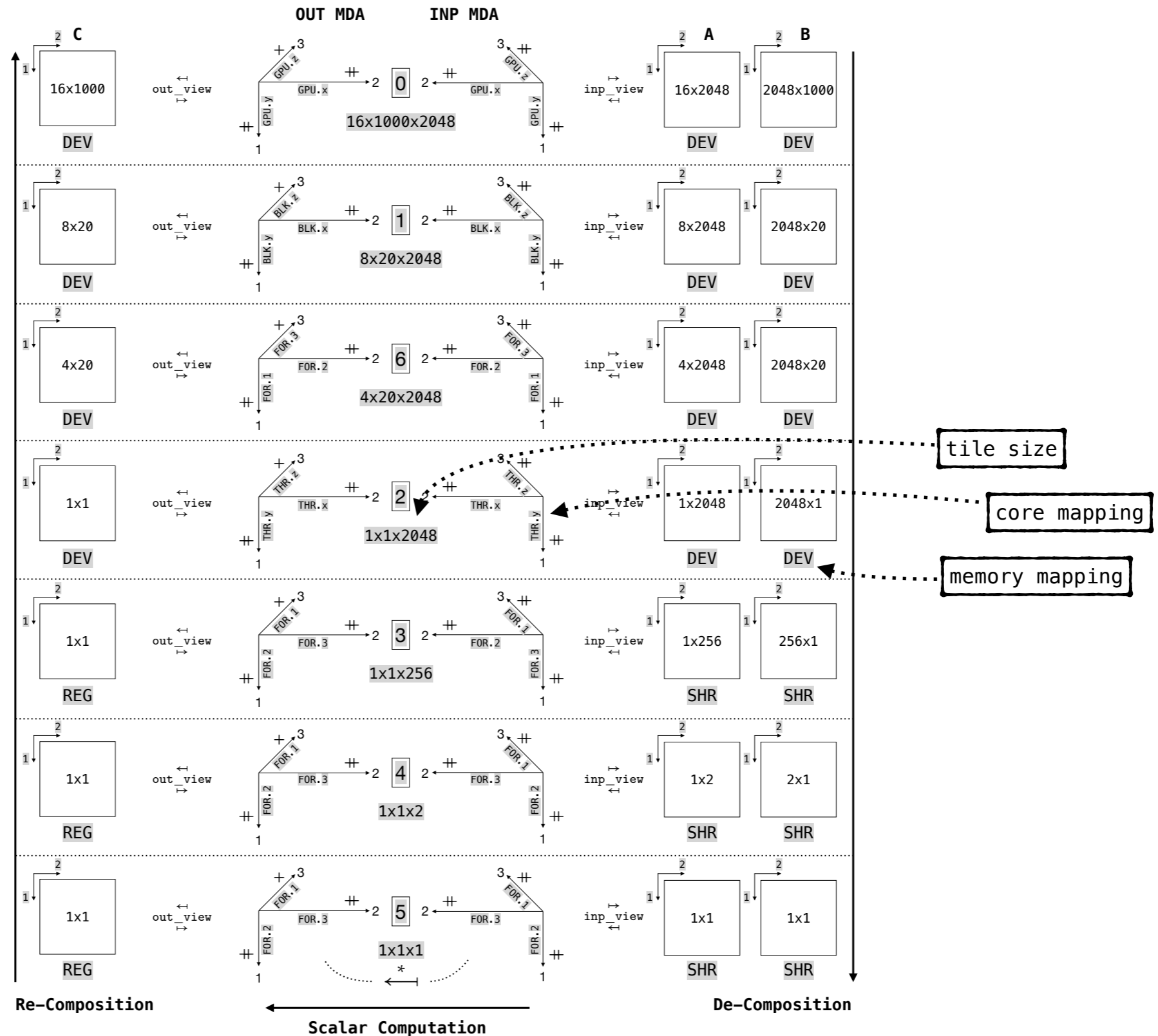
1. (De/Re)-Composition

2. (De/Re)-Composition

3. (De/Re)-Composition

4. (De/Re)-Composition

5. (De/Re)-Composition

6. (De/Re)-Composition

# Related Work

- Popular scheduling approaches include: `TVM [OSDI '18], Halide [PLDI '13], Elevate [ICFP '20], DaCe [SC '19], Tiramisu [CGO '19], CUDA-CHiLL [TACO '13], Fireiron [PACT '20], Distal [PLDI '22], and LoopStack [arXiv '22]`

- All these approaches have in common that their scheduling languages rely on fine-grained low-level primitives which are expressive but complex to use, often even for experts

- Our language design allows combining the following advantages over the related work:

  1. **Performance:** competitive to TVM and often higher

  2. **Safety:** backed by MDH formalism

  3. **Auto-Tuning:** any optimization decision can be optionally left for auto-tuning

  4. **Applicability:** our language is used similarly for multiple kinds of programming models

  5. **Visualization:** our schedules can be visualized and also be generated from visual inputs

# Conclusion & Future Work

## Conclusion:

- We introduce a new scheduling language, based on the formalism of Multi-Dimensional Homomorphisms (MDH)

- The goal of our language design is to express (de/re)-compositions of computations in a systematic, structured way to simplify the complex and error-prone optimization process for performance experts

- Our language design enables: 1) *Performance*, 2) *Safety*, 3) Auto-Tuning, 4) *Applicability*, and 5) *Visualization*

## Future Work:

- Computations consisting of multiple loop nests (currently limited to individual nests)

- Targeting domain-specific hardware extensions, e.g., *NVIDIA Tensor Cores*

- Targeting further models, e.g., LLVM to benefit from assembly-level optimizations