

Programming abstractions for in and near-memory computing

Jeronimo Castrillon

Chair for Compiler Construction (CCC)

TU Dresden, Germany

In-Memory Architectures and Computing Applications Workshop (iMACAW'23)

Design Automation Conference (DAC'23). San Francisco, USA

July 9, 2023



Why new abstractions

$$v_{ijk,e} = \sum_{i'=0}^p \sum_{j'=0}^p \sum_{k'=0}^p A_{kk'} A_{jj'} A_{ii'} u_{i'j'k'} e$$

What we want

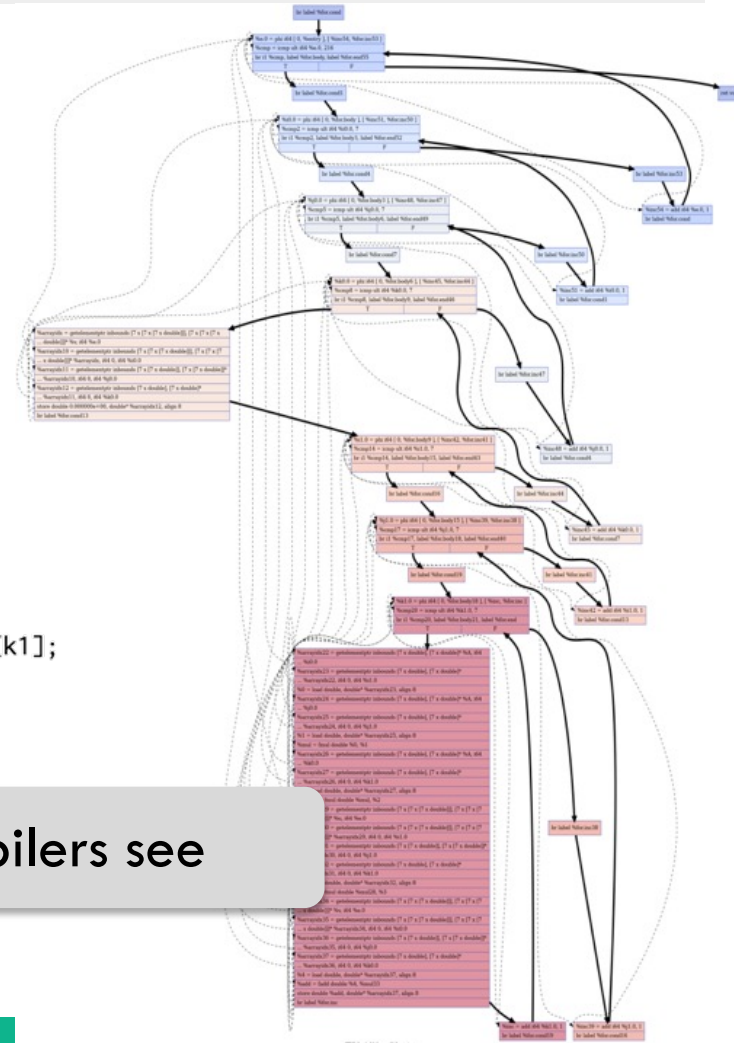
What we (naively) code

```

1 void cfd_kernel(
2   double A[restrict 7][7],
3   double u[restrict 216][7][7][7],
4   double v[restrict 216][7][7][7])
5 {
6   /* element loop: */
7   for(int e = 0; e < 216; e++) {
8     for(int i0 = 0; i0 < 7; i0++) {
9       for(int j0 = 0; j0 < 7; j0++) {
10        for(int k0 = 0; k0 < 7; k0++) {
11          v[e][i0][j0][k0] = 0.0;
12          for(int i1 = 0; i1 < 7; i1++) {
13            for(int j1 = 0; j1 < 7; j1++) {
14              for(int k1 = 0; k1 < 7; k1++) {
15                v[e][i0][j0][k0] += A[i0][i1]
16                                     * A[j0][j1]
17                                     * A[k0][k1]
18                                     * u[e][i1][j1][k1];
19              } } } } }
20          } /* end of element loop */
21        }

```

What compilers see



Semantic gap → performance gap

$$v_{ijk,e} = \sum_{i'=0}^p \sum_{j'=0}^p \sum_{k'=0}^p A_{kk'} A_{jj'} A_{ii'} u_{i'j'k'e}$$

What we want

What we (naively) code

```
1 void cfd_kernel(  
2     double A[restrict 7][7],  
3     double u[restrict 216][7][7][7],  
4     double v[restrict 216][7][7][7])  
5 {  
6     /* element loop: */  
7     for(int e = 0; e < 216; e++) {  
8         for(int i0 = 0; i0 < 7; i0++) {  
9             for(int j0 = 0; j0 < 7; j0++) {  
10                for(int k0 = 0; k0 < 7; k0++) {  
11                    v[e][i0][j0][k0] = 0.0;  
12                    for(int i1 = 0; i1 < 7; i1++) {  
13                        for(int j1 = 0; j1 < 7; j1++) {  
14                            for(int k1 = 0; k1 < 7; k1++) {  
15                                v[e][i0][j0][k0] += A[i0][i1]  
16                                    * A[j0][j1]  
17                                    * A[k0][k1]  
18                                    * u[e][i1][j1][k1];  
19                            } } } } } }  
20                } /* end of element loop */  
21            }  
22        }  
23    }
```

100X

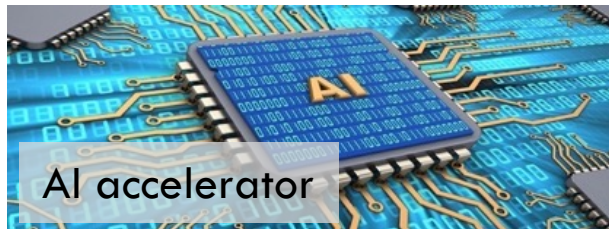
```
1 void cfd_kernel(  
2     double A[restrict 7][7],  
3     double u[restrict 216][7][7][7],  
4     double v[restrict 216][7][7][7])  
5 {  
6     /* element loop: */  
7     #pragma omp for  
8     for (int e = 0; e < 216; e++) {  
9         double t6[7][7][7];  
10        /* 1st contraction: */  
11        #pragma simd  
12        for (int i0 = 0; i0 < 7; i0++) {  
13            for (int i1 = 0; i1 < 7; i1++) {  
14                /* #pragma simd */  
15                for (int i2 = 0; i2 < 7; i2++) {  
16                    double t8 = 0.0;  
17                    for (int i3 = 0; i3 < 7; i3++)  
18                        t8 += A[i0][i3] * u[e][i1][i2][i3];  
19                    t6[i0][i1][i2] = t8;  
20                } } /* end of 1st contraction */  
21                double t7[7][7][7];  
22                /* 2nd contraction: */  
23                #pragma simd  
24                for (int i4 = 0; i4 < 7; i4++) {  
25                    for (int i5 = 0; i5 < 7; i5++) {  
26                        /* #pragma simd */  
27                        for (int i6 = 0; i6 < 7; i6++) {  
28                            double t9 = 0.0;  
29                            for (int i7 = 0; i7 < 7; i7++)  
30                                t9 += A[i4][i7] * t6[i5][i6][i7];  
31                            t7[i4][i5][i6] = t9;  
32                        } } } /* end of 2nd contraction */  
33                        /* 3rd contraction: */  
34                        #pragma simd  
35                        for (int i8 = 0; i8 < 7; i8++) {  
36                            for (int i9 = 0; i9 < 7; i9++) {  
37                                /* #pragma simd */  
38                                for (int i10 = 0; i10 < 7; i10++) {  
39                                    double t10 = 0.0;  
40                                    for (int i11 = 0; i11 < 7; i11++)  
41                                        t10 += A[i8][i11] * t7[i9][i10][i11];  
42                                    v[e][i8][i9][i10] = t10;  
43                                } } } /* end of third contraction */  
44                            } } } /* end of element loop */  
45                    } } }  
46                } } }  
47            } } }  
48        } } }  
49    }
```

What performance experts code

Semantic gap → performance gap

$$v_{ijk,e} = \sum_{i'=0}^p \sum_{j'=0}^p \sum_{k'=0}^p A_{kk'} A_{jj'} A_{ii'} u_{i'j'k'e}$$

What we want



AI accelerator

<https://www.hpcwire.com/2017/04/10/vidia-responds-google-tpu-benchmarking/>

Lee, Sukhan, et al. "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology: Industrial Product." ISCA 2021.

```

1 void cfd_kernel(
2   double A[restrict 7][7],
3   double u[restrict 216][7][7][7],
4   double v[restrict 216][7][7][7])
5 {
6   /* element loop: */
7   for(int e = 0; e < 216; e++) {
8     for(int i0 = 0; i0 < 7; i0++) {
9       for(int j0 = 0; j0 < 7; j0++) {
10        for(int k0 = 0; k0 < 7; k0++) {
11          v[e][i0][j0][k0] = 0.0;
12          for(int i1 = 0; i1 < 7; i1++)
13            for(int j1 = 0; j1 < 7; j1++)
14              for(int k1 = 0; k1 < 7; k1++) {
15                v[e][i0][j0][k0] += A[i0][i1]
16                                     * A[j0][j1]
17                                     * A[k0][k1];
18              }
19        }
20      }
21    }

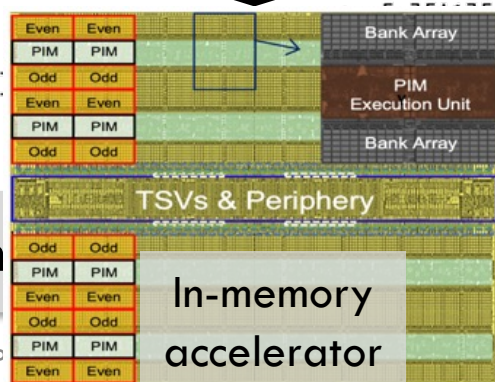
```

```

1 void cfd_kernel(
2   double A[restrict 7][7],
3   double u[restrict 216][7][7][7],
4   double v[restrict 216][7][7][7])
5 {
6   /* element loop: */
7   #pragma omp for
8   for (int e = 0; e < 216; e++) {
9     double t6[7][7][7];
10    /* 1st contraction: */
11    #pragma simd
12    for (int i0 = 0; i0 < 7; i0++) {
13      for (int i1 = 0; i1 < 7; i1++) {
14        /* #pragma simd */
15        for (int i2 = 0; i2 < 7; i2++) {
16          double t8 = 0.0;
17          for (int i3 = 0; i3 < 7; i3++)
18            t8 += A[i0][i3] * u[e][i1][i2][i3];
19          t6[i0][i1][i2] = t8;
20        } } /* end of 1st contraction */

```

100X



In-memory accelerator



HBM-FPGA

```

1 for (int i0 = 0; i0 < 7; i0++) {
2   double t10 = 0.0;
3   for (int i11 = 0; i11 < 7; i11++)
4     t10 += A[i8][i11] * t7[i9][i10][i11];
5 } } /* end of third contraction */
6 } /* end of element loop */

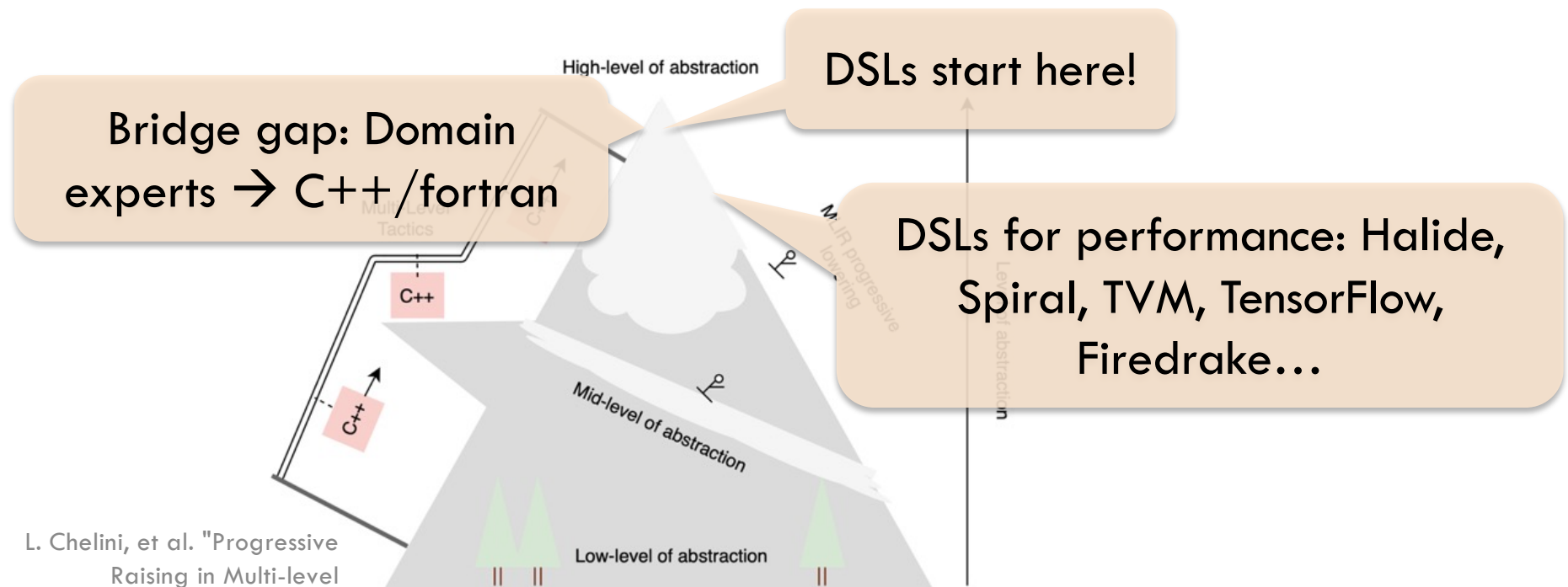
```

ports code

Wh

There is only so much we can do/reconstruct...

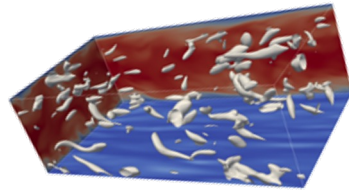
- Lots of progress: polyhedral compilers, trace-driven dynamic parallelization, patterns/ idiom extraction, ...



L. Chelini, et al. "Progressive Raising in Multi-level IR." CGO 2021

Sample DSLs

CFDlang

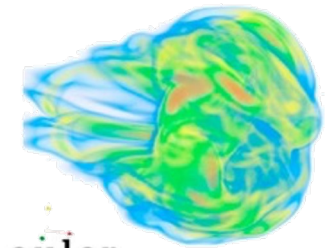


$$v_{ijk,e} = \sum_{i'=0}^p \sum_{j'=0}^p \sum_{k'=0}^p A_{kk'} A_{jj'} A_{ii'} u_{i'j'k'e}$$

```
source = ...
var input A      : matrix          &
var input u      : tensorIN        &
var input output v : tensorOUT     &
var input alpha  : []              &
var input beta   : []              &
v = alpha * (A # A # A # u .
  [[5 8] [3 7] [1 6]]) + beta * v
```

```
auto A = Matrix(m, n), B = Matrix(m, n),
      C = Matrix(m, n);
auto u = Tensor<3>(n, n, n);
auto v = (A*B*C)(u);
```

OpenPME



time loop

start: 0 stop: 1000

temporal method: explicit_euler

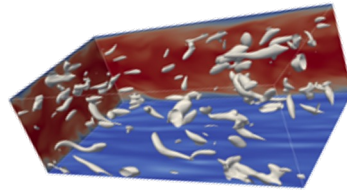
spatial method: DC-PSE

$$\frac{\partial u}{\partial t} = Du * \nabla^2 u - u * v^2 + F * (1 - u)$$

$$\frac{\partial v}{\partial t} = Dv * \nabla^2 v + u * v^2 - v * (F + k)$$

Sample DSLs

CFDlang

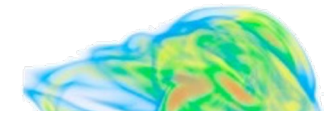


$$v_{ijk,e} = \sum_{i'=0}^p \sum_{j'=0}^p \sum_{k'=0}^p A_{kk'} A_{jj'} A_{ii'} u_{i'j'k'e}$$

```
source = ...
var input A : matrix &
var input u : tensorIN &
var input output v : tensorOUT &
var input alpha : [] &
var input beta : [] &
v = alpha * (A # A # A # u .
  [[5 8] [3 7] [1 6]]) + beta * v
```

```
auto A = Matrix(m, n), B = Matrix(m, n),
      C = Matrix(m, n);
auto u = Tensor<3>(n, n, n);
auto v = (A*B*C)(u);
```

OpenPME

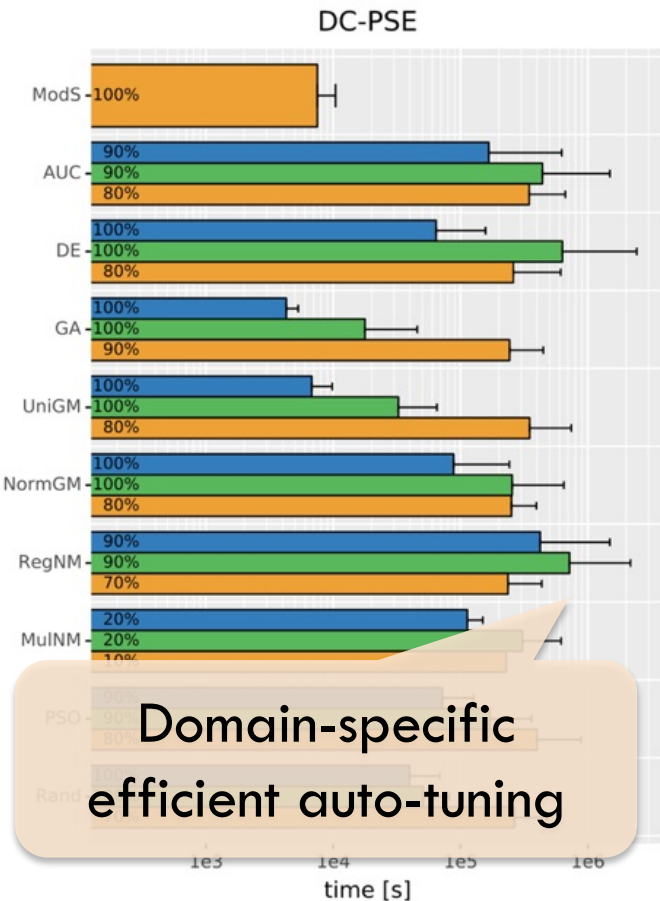
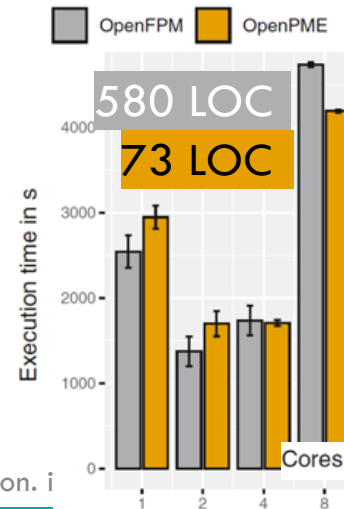


time loop

start: 0 stop: 1
temporal method: Runge-Kutta
spatial method: Finite Difference

$$\frac{\partial u}{\partial t} = Du * \nabla^2$$

$$\frac{\partial v}{\partial t} = Dv * \nabla^2$$



Closing the performance gap

- ❑ Not really optimization magic
 - ❑ Leverage expert knowledge
 - ❑ Algebraic identities

$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot (A_{jm} \cdot (A_{il} \cdot u_{lmn})))$$

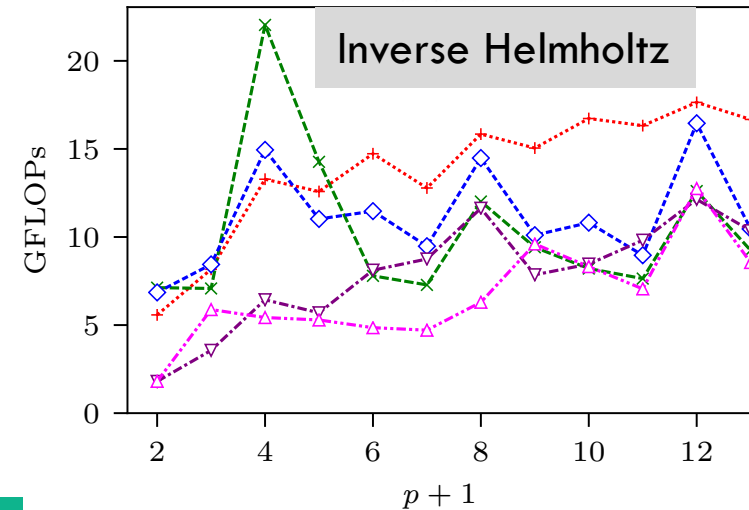
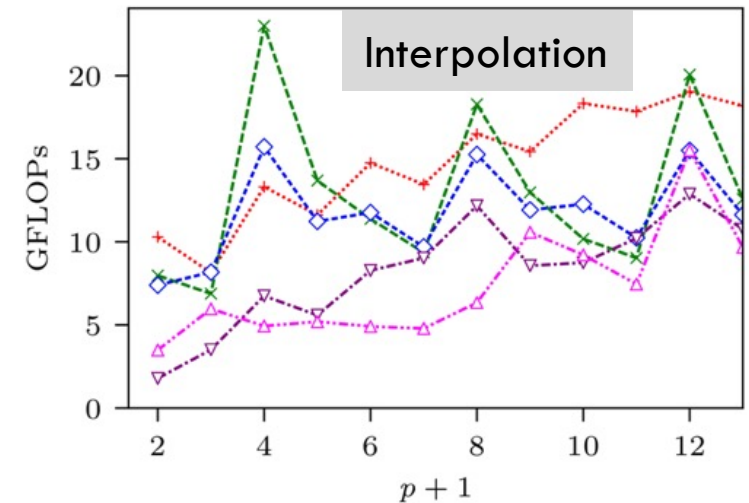
$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot A_{jm}) \cdot (A_{il} \cdot u_{lmn})$$

$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot ((A_{jm} \cdot A_{il}) \cdot u_{lmn}))$$

N. A. Rink, et al. "CFDlang: High-level code generation for high-order methods in fluid dynamics". RWDSL'18.

A. Susungi, et al., "Meta-programming for Cross-Domain Tensor Optimizations", GPCE'18 pp. 79-92.

- +...+ CFDlang(outer)
- x...x CFDlang(inner)
- ◇...◇ hand-optimized
- ▽...▽ DGEMM
- △...△ specialized



Closing the performance gap

- ❑ Not really optimization magic
 - ❑ Leverage expert knowledge
 - ❑ Algebraic identities

$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot (A_{jm} \cdot (A_{il} \cdot u_{lmn})))$$

$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot A_{jm}) \cdot (A_{il} \cdot u_{lmn})$$

$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot ((A_{jm} \cdot A_{il}) \cdot u_{lmn}))$$

N. A. Rink, et al. "CFDlang: High-level code generation for high-order methods in fluid dynamics". RWDSL'18.

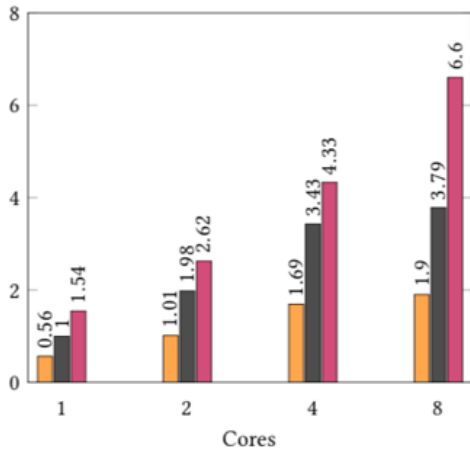
A. Susungi, et al., "Meta-programming for Cross-Domain Tensor Optimizations", GPCE'18 pp. 79-92.

Easy to generate,
hard to transform

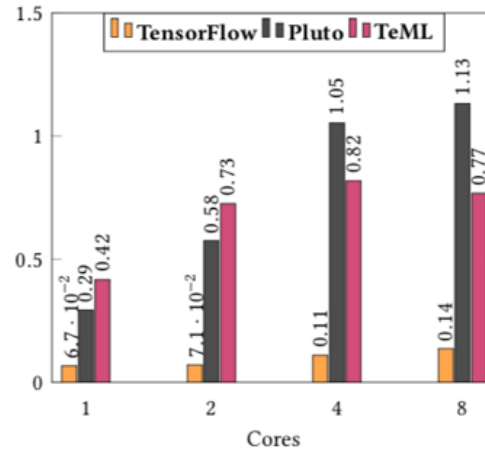
Actual code variants

Cross-domain optimizations

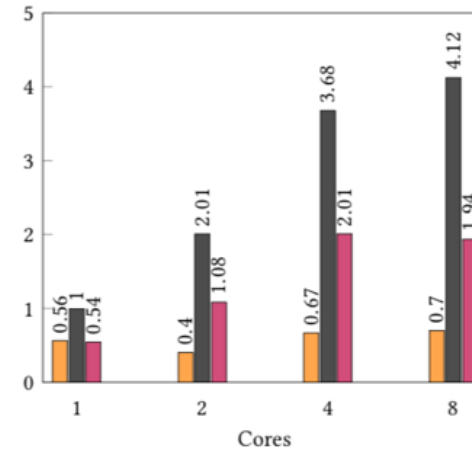
(a) mttkrp



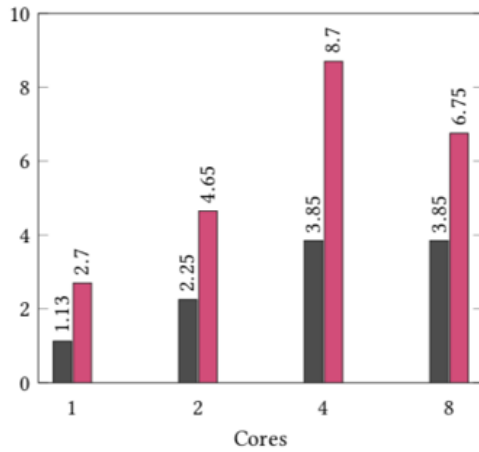
(b) bmm



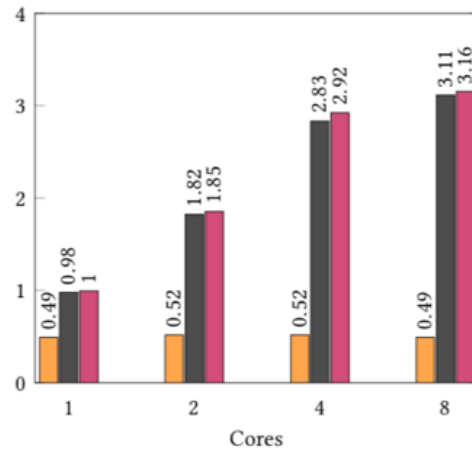
(c) sddmm



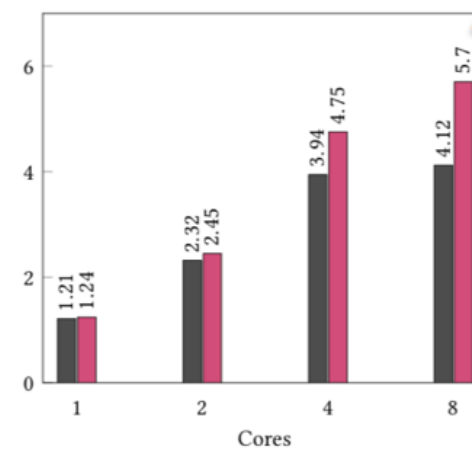
(d) gconv



(e) interp



(f) helm



Performance of Pluto could be reproduced

Higher abstraction → more optimization potential

A. Susungi, et al. "Meta-programming for cross-domain tensor optimizations", GPCE'18, 79-92

Tell: Tensor Intermediate language

- Formalized core common to multiple tensor languages
- Index-free notation and strong type system
- Provably** no out-of-bound accesses

```
A = placeholder((m,h), name='A')
B = placeholder(h, name='B')
C = compute((m, lambda i, j:
    sum(A[k, i] * B[k, j], axis=k)))
```

$$C_{ij} = \sum_{k=1}^h A_{ki} B_{kj}$$

$\llbracket \cdot \rrbracket : \text{Context} \rightarrow \text{Memory} \rightarrow (\text{list of Nat}) \rightarrow \mathbb{D}$

$\llbracket x \rrbracket \Gamma \mu \bar{i} = \mu x \bar{i}$

$\llbracket (e) \rrbracket \Gamma \mu \bar{i} = \llbracket e \rrbracket \Gamma \mu \bar{i}$

$\llbracket \text{add } e_0 e_1 \rrbracket \Gamma \mu \bar{i} = \llbracket e_0 \rrbracket \Gamma \mu \bar{i} + \llbracket e_1 \rrbracket \Gamma \mu \bar{i}$

$\llbracket \text{mul } e_0 e_1 \rrbracket \Gamma \mu \bar{i} = \begin{cases} \llbracket e_0 \rrbracket \Gamma \mu [] \cdot \llbracket e_1 \rrbracket \Gamma \mu \bar{i}, & \text{if } \text{type}_{\Gamma}(e_0) = [] \\ \llbracket e_0 \rrbracket \Gamma \mu \bar{i} \cdot \llbracket e_1 \rrbracket \Gamma \mu \bar{i}, & \text{otherwise} \end{cases}$

$\llbracket \text{prod } e_0 e_1 \rrbracket \Gamma \mu (\bar{i}_0 \# \bar{i}_1) = \llbracket e_0 \rrbracket \Gamma \mu \bar{i}_0 \cdot \llbracket e_1 \rrbracket \Gamma \mu \bar{i}_1,$
if $\text{rank}_{\Gamma}(e_0) = \text{length}(\bar{i}_0)$ and $\text{rank}_{\Gamma}(e_1) = \text{length}(\bar{i}_1)$

$\llbracket \text{red}_+ i e \rrbracket \Gamma \mu [j_1, \dots, j_{i-1}, j_i, \dots, j_k] = \sum_{m=1}^n \llbracket e \rrbracket \Gamma \mu [j_1, \dots, j_{i-1}, m, j_i, \dots, j_k],$ if $\text{type}_{\Gamma}(e) = [n_1, \dots, n_{i-1}, n, n_{i+1}, \dots, n_{k+1}]$

$\llbracket \text{transp } i_0 i_1 e \rrbracket \Gamma \mu [j_1, \dots, j_{i_0}, \dots, j_{i_1}, \dots, j_k] =$

$\llbracket e \rrbracket \Gamma \mu [j_1, \dots, j_{i_1}, \dots, j_{i_0}, \dots, j_k]$

$\llbracket \text{diag } i_0 i_1 e \rrbracket \Gamma \mu [j_1, \dots, j_{i_0-1}, j_{i_0}, j_{i_0+1}, \dots, j_{i_1-1}, j_{i_1}, \dots, j_k] =$

$\llbracket e \rrbracket \Gamma \mu [j_1, \dots, j_{i_0-1}, j_{i_0}, j_{i_0+1}, \dots, j_{i_1-1}, j_{i_1}, \dots, j_k]$

$\llbracket \text{exp } i n e \rrbracket \Gamma \mu [j_1, \dots, j_{i-1}, j_i, j_{i+1}, \dots, j_k] =$

$\llbracket e \rrbracket \Gamma \mu [j_1, \dots, j_{i-1}, j_{i+1}, \dots, j_k]$

$\llbracket \text{proj } i m e \rrbracket \Gamma \mu [j_1, \dots, j_{i-1}, j_i, \dots, j_k] =$

$\llbracket e \rrbracket \Gamma \mu [j_1, \dots, j_{i-1}, m, j_i, \dots, j_k]$

N.A. Rink, N. A. and J. Castrillon. "Tell: a type-safe imperative Tensor Intermediate Language", ARRAY'19, pp. 57-68

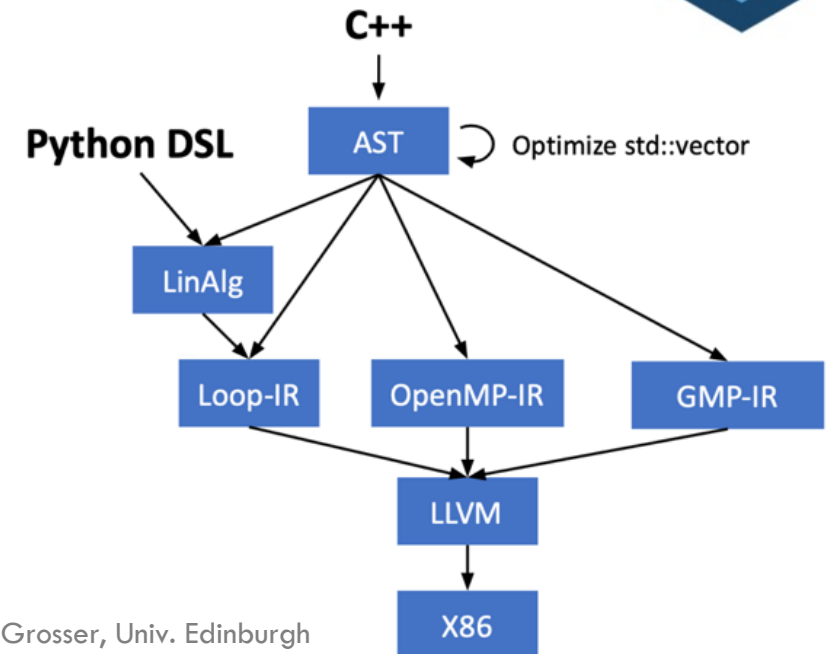
Intermediate abstractions in MLIR

- Started by Google ~2018, now in public domain

Lattner, Chris, et al. "Mlir: Scaling compiler infrastructure for domain specific computation." 2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO). IEEE, 2021.



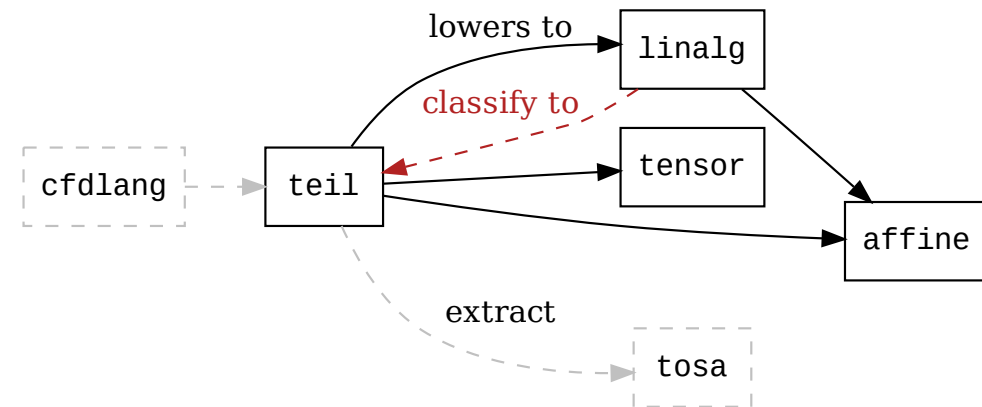
- Not an IR, but an extensible framework
 - to describe intermediate abstractions (called **dialects**),
 - to optimize representations between dialects (**transform, lower or raise**),
 - that builds on the success of LLVM to build community/infrastructure and reuse ("LLVM-quality" all the way)



Source: T. Grosser, Univ. Edinburgh

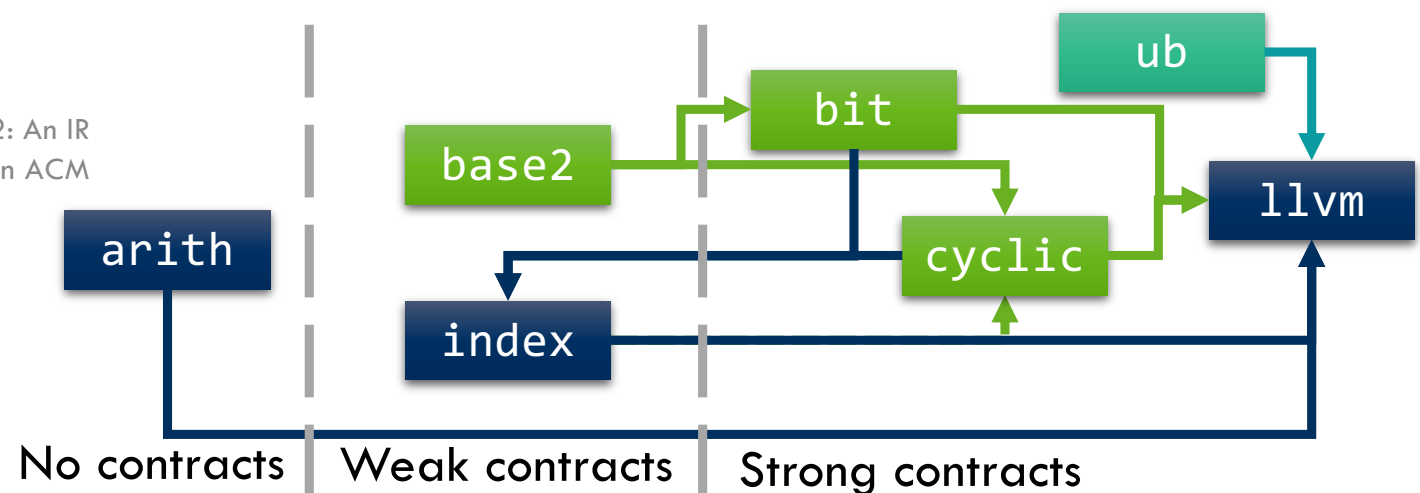
Tell in MLIR

- ❑ Primitive ops instead of index maps
 - ❑ Easier to express identities (big-O trfs)
 - ❑ Uses symbolic math, infinite precision



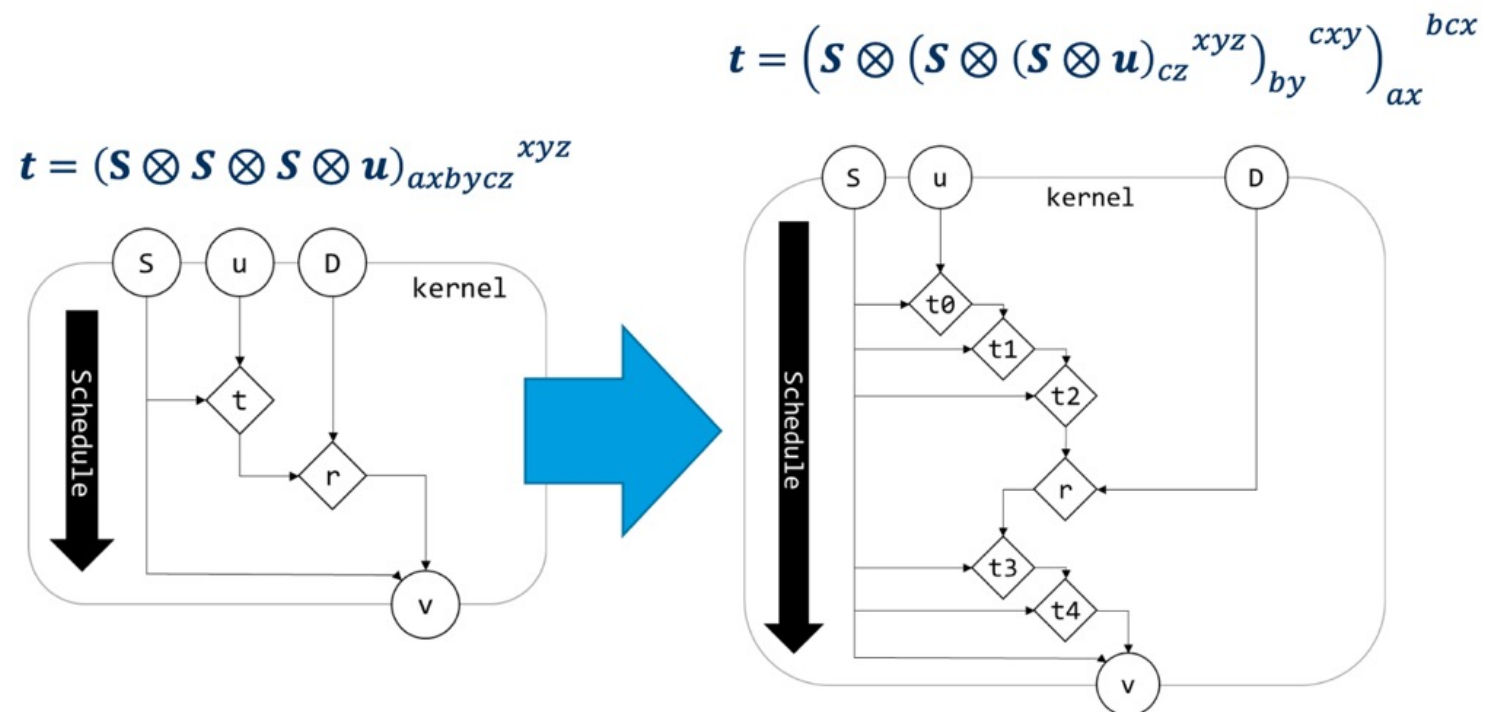
- ❑ Specialization path to custom hardware

K. F. A. Friebel, J. Bi, J. Castrillon, "BASE2: An IR for Binary Numeral Types" (to appear), In ACM HEART 2023



Domain-specific optimization

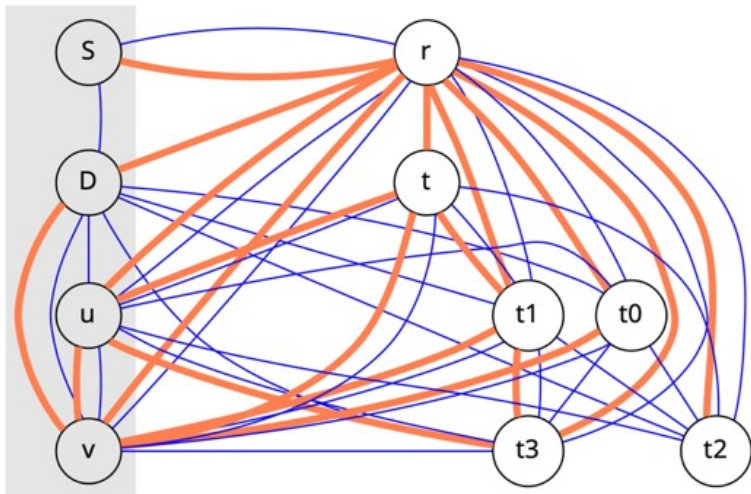
- ❑ Encode algebraic transformations (Interpolation as example)
- ❑ Direct feedback to expert via DSL export



High-level buffer re-use

- Generate host code and accelerator code (for HLS)
- Generate liveness info (buffering, memory subsystem gen)

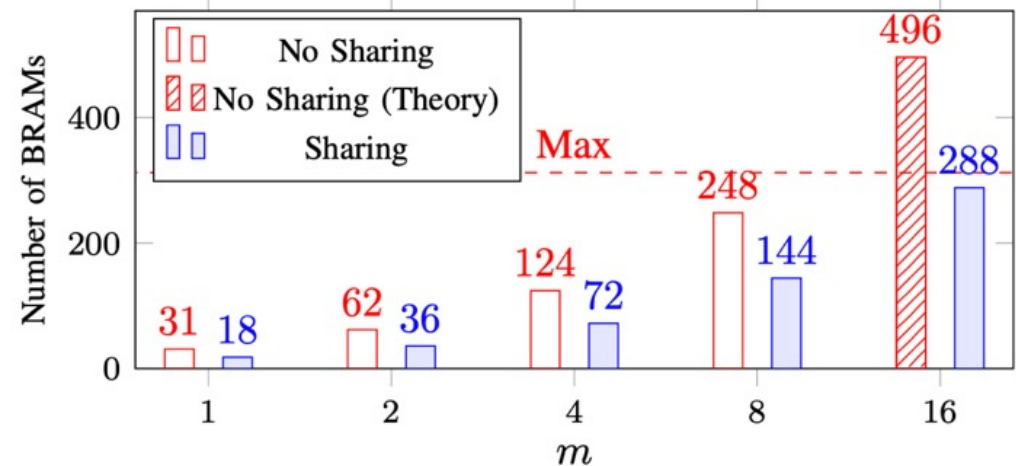
```
t = S # S # S # u . [[1 6] [3 7] [5 8]]
r = D * t
v = S # S # S # r . [[0 6] [2 7] [4 8]]
```



memory-interface and address-space compatibilities



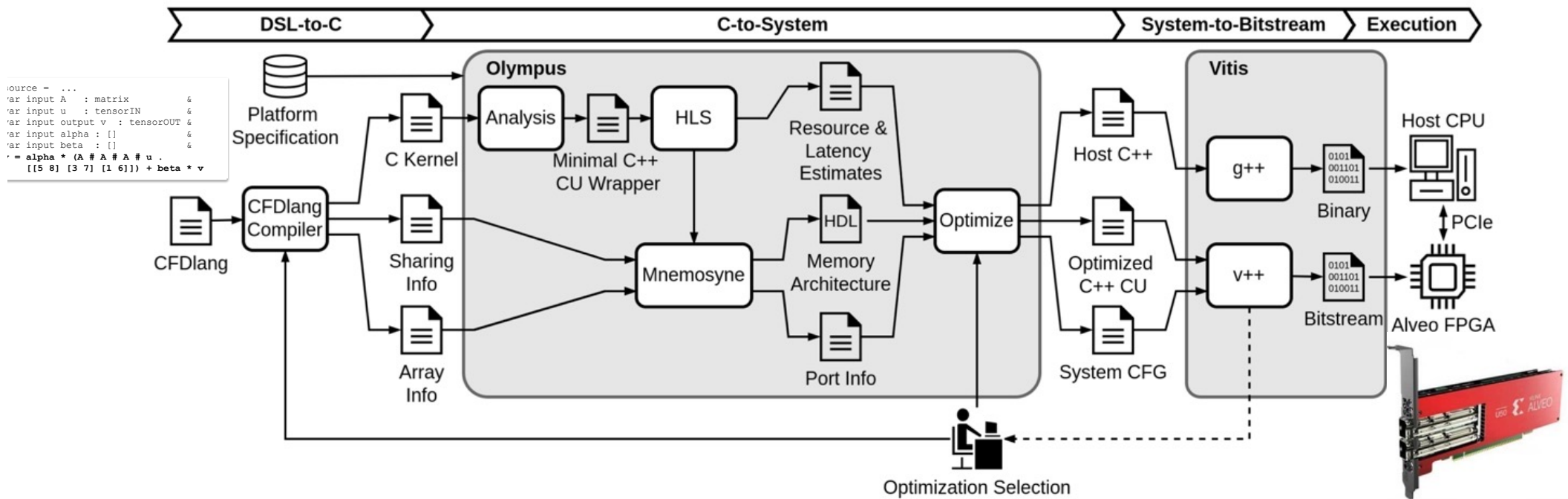
<https://everest-h2020.eu>



K. F. A. Friebel, et al., "From Domain-Specific Languages to Memory-Optimized Accelerators for Fluid Dynamics", Proceedings of the FPGA for HPC Workshop, held in conjunction with IEEE Cluster 2021, Sep 2021

Putting it all together

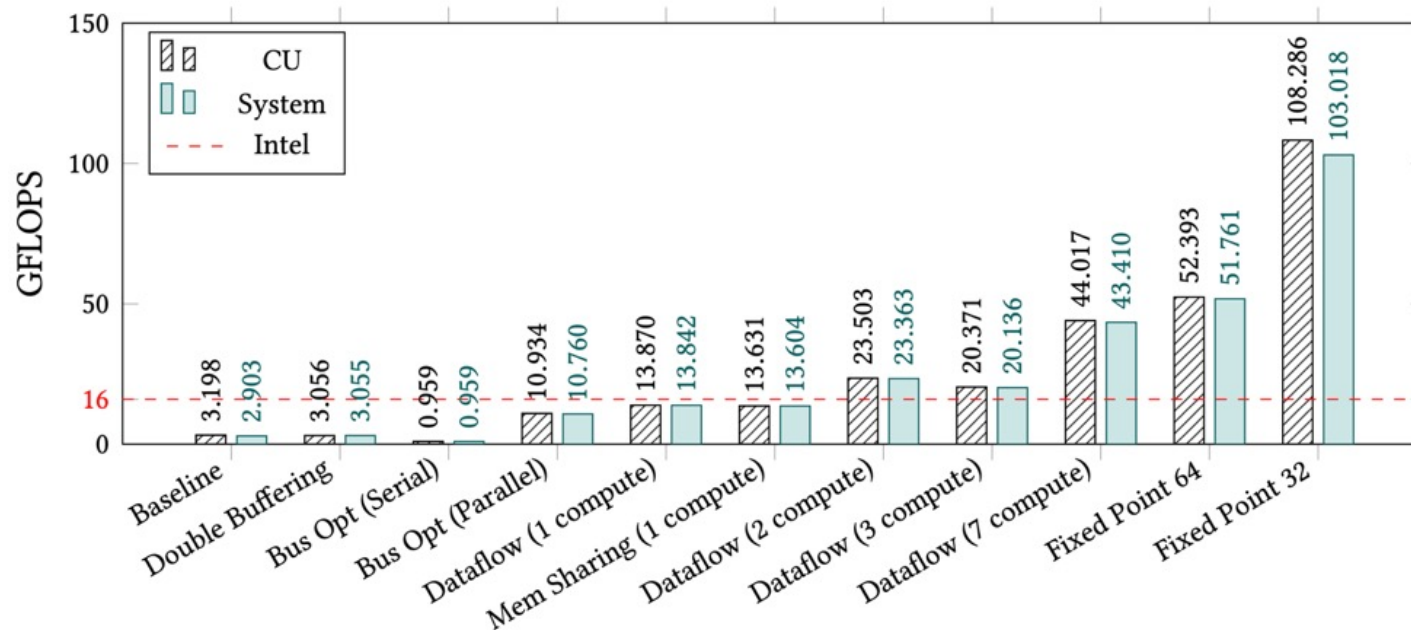
Complex compilation/design flow from DSL to system-level architecture



S. Soldavini, K. F. A. Friebel, M. Tibaldi, G. Hempel, J. Castrillon, and C. Pilato. "Automatic Creation of High-Bandwidth Memory Architectures from Domain-Specific Languages: The Case of Computational Fluid Dynamics". In: ACM TRET, Sept. 2022.

FPGA code generation: HBM FPGA

- ❑ H2020 EU Project: Convergence HPC, Big Data and ML
- ❑ Transformations for a **17x speedup** (same precision)



<https://everest-h2020.eu>



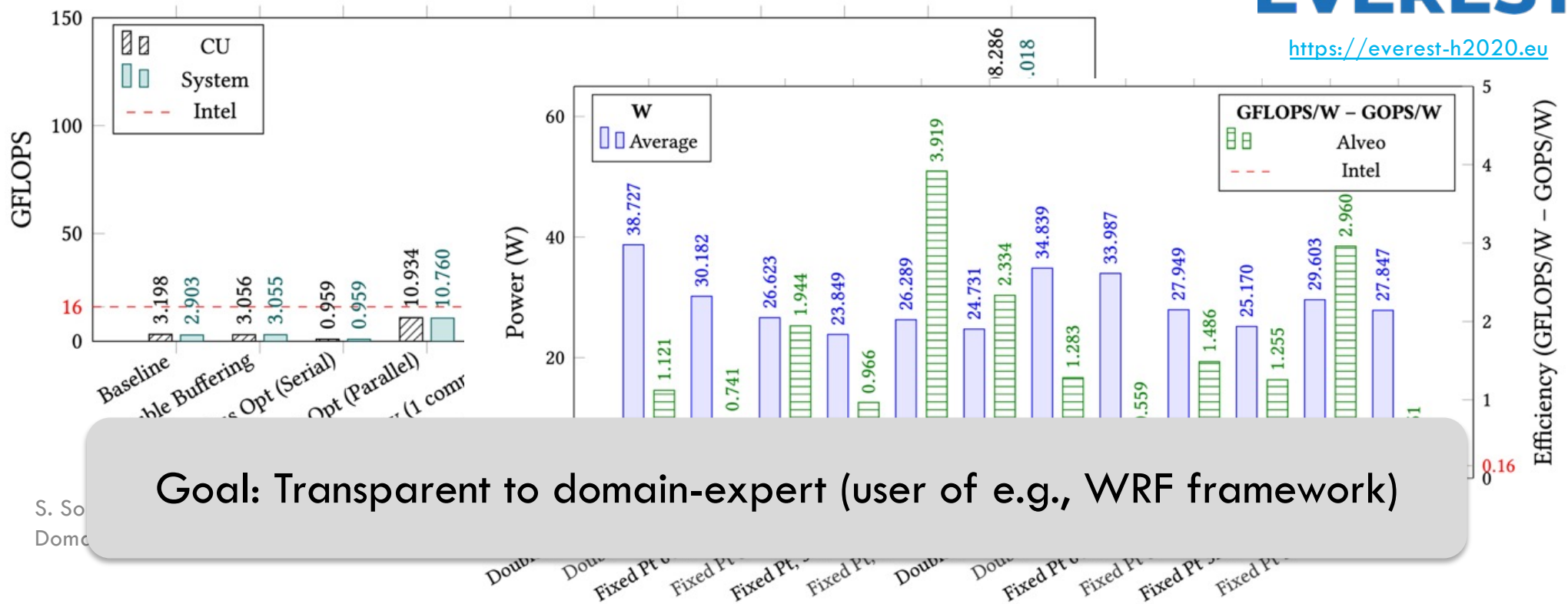
S. Soldavini, K. F. A. Friebel, M. Tibaldi, G. Hempel, J. Castrillon, and C. Pilato. “Automatic Creation of High-Bandwidth Memory Architectures from Domain-Specific Languages: The Case of Computational Fluid Dynamics”. In: ACM TRET, Sept. 2022.

FPGA code generation: HBM FPGA

- ❑ H2020 EU Project: Convergence HPC, Big Data and ML
- ❑ Variants with up to **24x better energy efficiency**



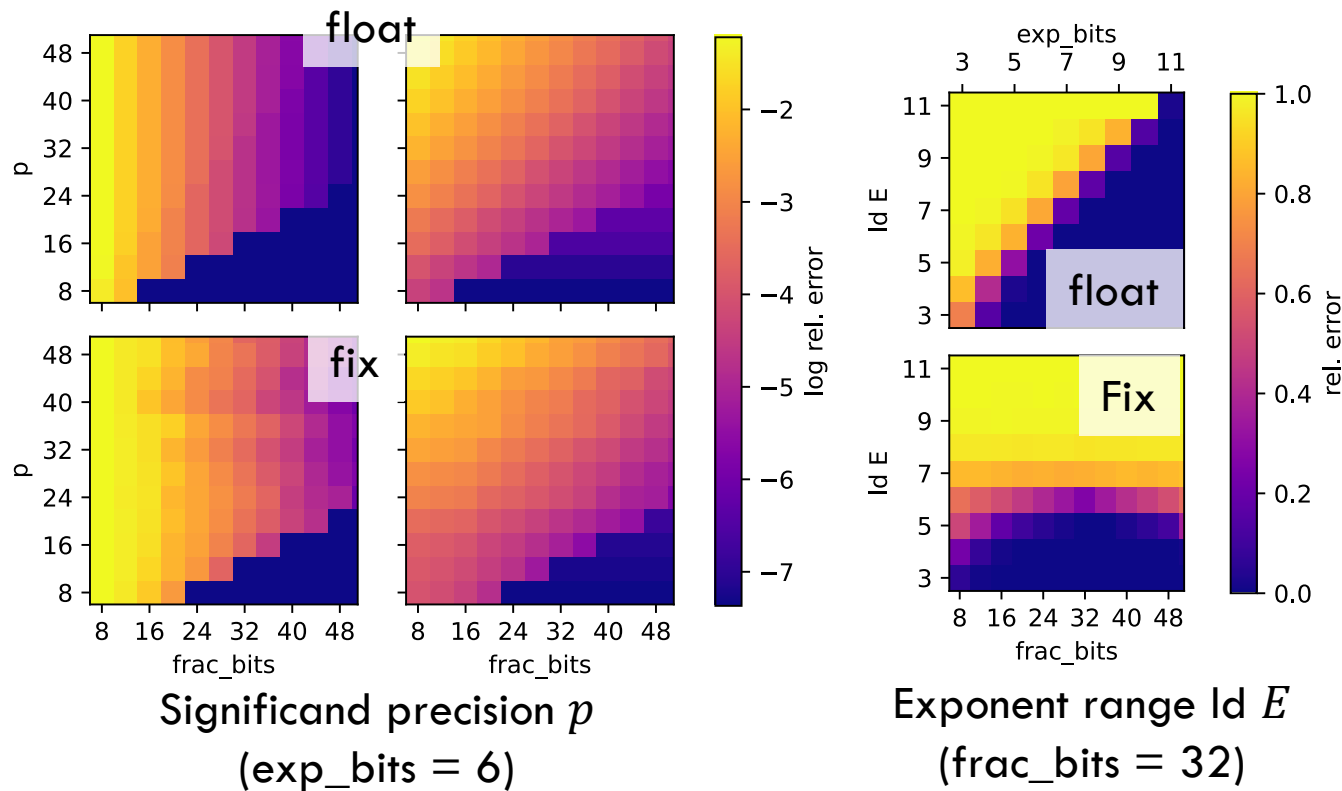
<https://everest-h2020.eu>



Goal: Transparent to domain-expert (user of e.g., WRF framework)

Base2: Custom precision analysis

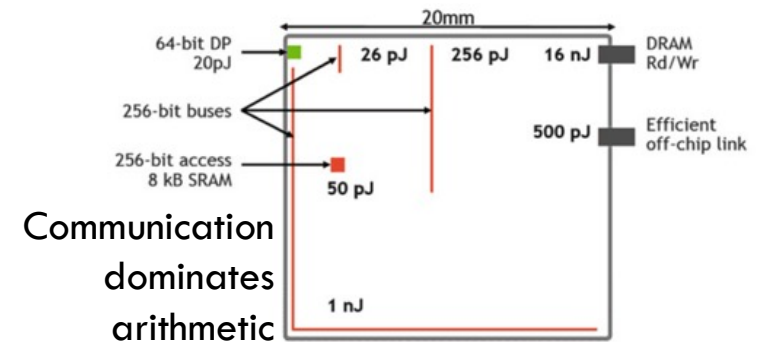
□ Interpolation
$$v_{ijk,e} = \sum_{i'=0}^p \sum_{j'=0}^p \sum_{k'=0}^p A_{kk'} A_{jj'} A_{ii'} u_{i'j'k'} e$$



K. F. A. Friebel, J. Bi, J. Castrillon, "BASE2: An IR for Binary Numeral Types" (to appear), In ACM HEART 2023

Emerging data-centric architectures

- ❑ Compute (almost) in-place, avoid data movement, transformations to match primitives
- ❑ Novel architectures for near-memory and in-memory computing



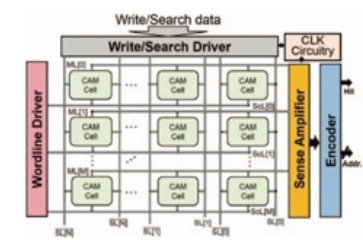
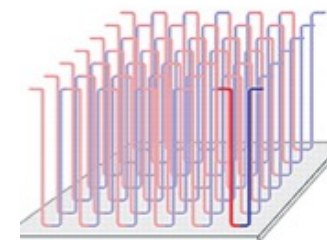
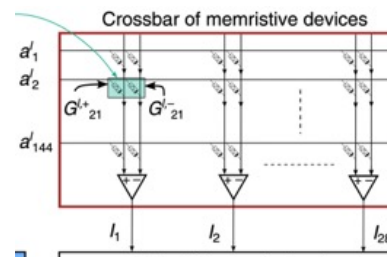
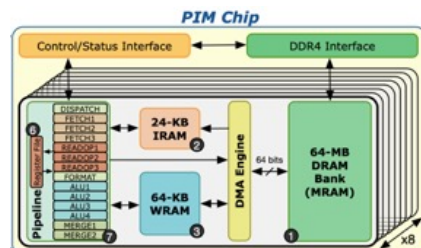
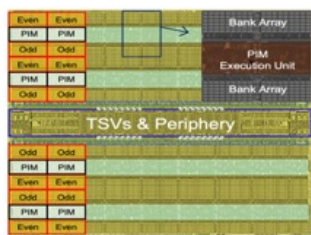
Source: Dally, NVIDIA

Samsung, Lee, Sukhan, et al. ISCA 2021

UPMEM by Gómez-Luna, Juan, et al. arXiv:2105.03814 (2021)

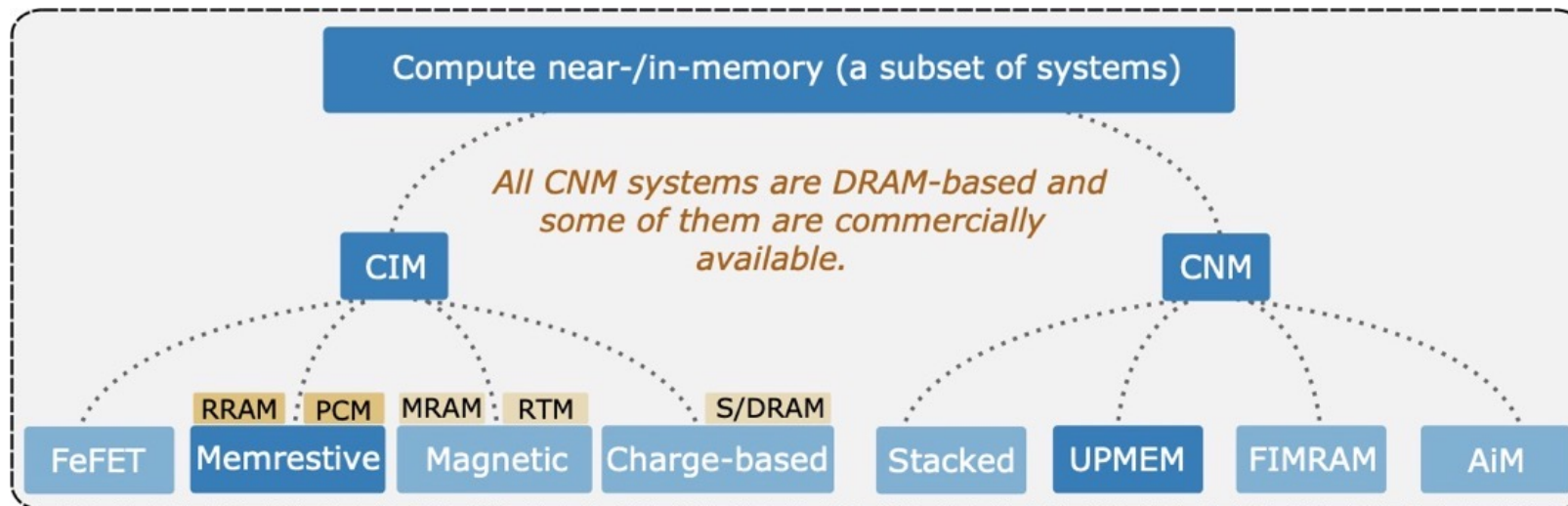
In-PCM Computing: Joshi, Vinay, et al. Nature Communications 11.1 (2020): 1-13.

CAM accelerators: Hu, Sharon, et al. 2021 IEDM



Compilation for heterogeneous CIM/CNM systems

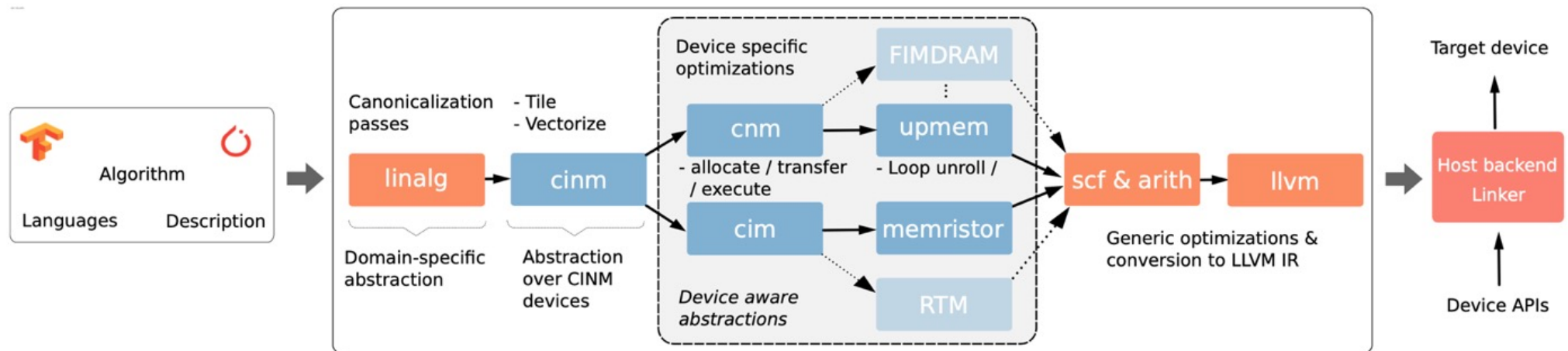
- A partial landscape/taxonomy of the CIM and CNM systems



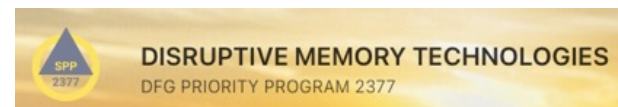
A. Khan et al, "CINM (Cinnamon): A Compilation Infrastructure for Heterogeneous Compute In-Memory and Compute Near-Memory Paradigms", arXiv, Jan 2023

Towards a generalized MLIR infrastructure

- ❑ Entry: linear algebra abstraction (common to ML frameworks and beyond)
- ❑ Intermediate languages for **in and near memory computing**
- ❑ Target-specific models and optimizations

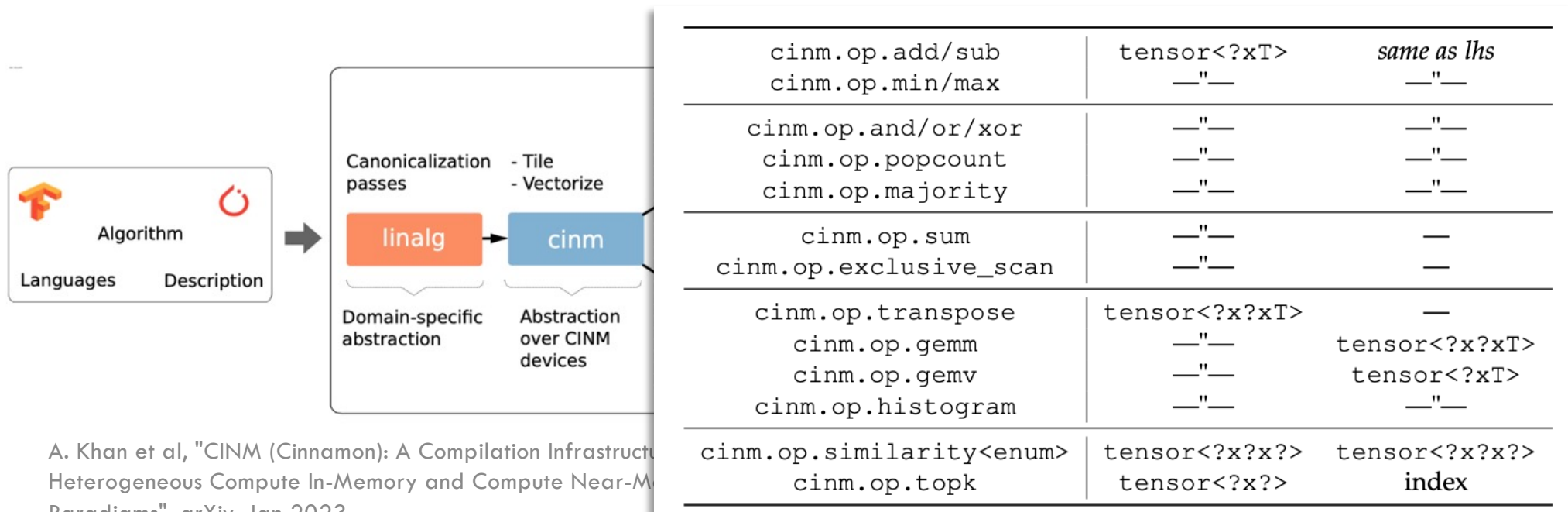


A. Khan et al, "CINM (Cinnamon): A Compilation Infrastructure for Heterogeneous Compute In-Memory and Compute Near-Memory Paradigms", arXiv, Jan 2023



Towards a generalized MLIR infrastructure

- ❑ Entry: linear algebra abstraction (common to ML frameworks and beyond)
- ❑ Intermediate languages for **in and near memory computing**
- ❑ Target-specific models and optimizations



A. Khan et al, "CINM (Cinnamon): A Compilation Infrastructure for Heterogeneous Compute In-Memory and Compute Near-Memory Paradigms", arXiv, Jan 2023

UPMEM example: Matmult

```
def mm(int32(64, 64) A, int32(64, 64) B) -> (int32(64, 64) C) {
    C(i,j) += A(i,k) * B(k,j)
        where i in 0:64, k in 0:64, j in 0:64
}
```

```
uint32_t mram_base_addr_A = (uint32_t) (DPU_MRAM_HEAP_POINTER );
uint32_t mram_base_addr_B = (uint32_t) (DPU_MRAM_HEAP_POINTER + ROWS * COLS *
↪ sizeof(T));
uint32_t mram_base_addr_C = (uint32_t) (DPU_MRAM_HEAP_POINTER + 2 * ROWS * COLS
↪ * sizeof(T));
for(int i = (tasklet_id * point_per_tasklet) ; i < (
↪ (tasklet_id+1)*point_per_tasklet ) ; i++) {
    if( new_row != row ){
        ...
        mram_read((__mram_ptr void const*) (mram_base_addr_A + mram_offset_A),
↪ cache_A, COLS * sizeof(T));
    }
    mram_read((__mram_ptr void const*) (mram_base_addr_B + mram_offset_B),
↪ cache_B, COLS * sizeof(T));
    dot_product(cache_C, cache_A, cache_B, number_of_dot_products);
    ...
}
...
mram_write( cache_C, (__mram_ptr void *) (mram_base_addr_C + mram_offset_C),
↪ point_per_tasklet * sizeof(T));
}
```


UPMEM example: Matmult

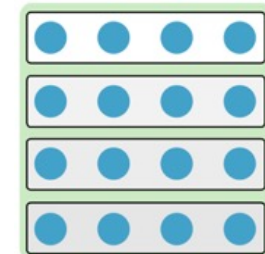
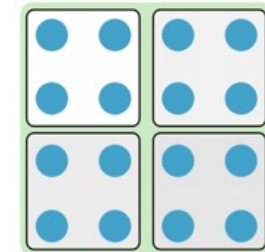
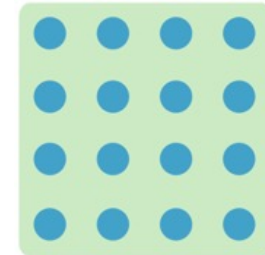
```
def mm(int32(64, 64) A, int32(64, 64) B) -> (int32(64, 64) C) {
  C(i,j) += A(i,k) * B(k,j)
    where i in 0:64, k in 0:64, j in 0:64
}
```



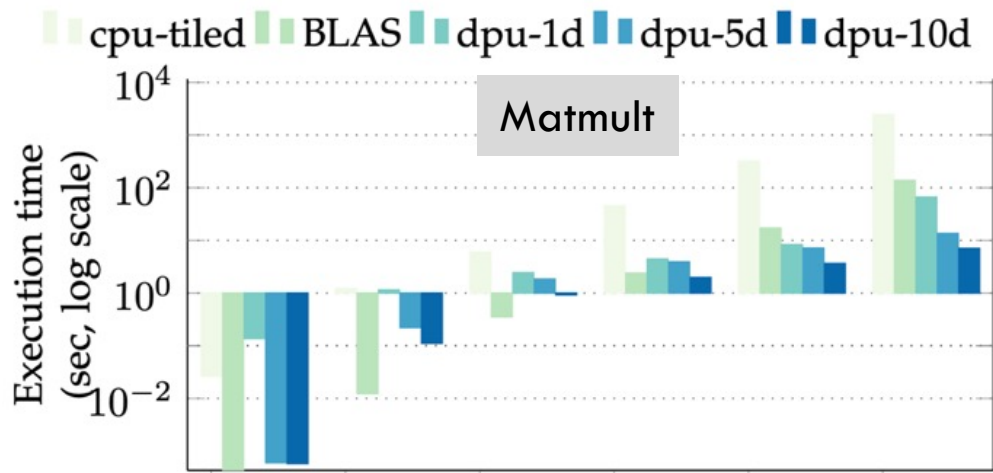
```
%D = linalg.matmul ins(%A, %B : tensor<64x64xi32>, tensor<64x64xi32>)
  outs(%C: tensor<64x64xi32>)
return %D : tensor<64x64xi32>
```



```
%C = scf.for %i = %cst0_i to %cst64_i step %cst16_i
<...>
  %B_tile = tensor.extract_slice %B[%o2, %o1][16, 16][1, 1]:
    tensor<64x64xi32> to tensor<16x16xi32>
  %A_t_b = bufferization.to_memref %A_tile: memref<16x16xi32>
  %B_t_b = bufferization.to_memref %B_tile: memref<16x16xi32>
  %A_dev = cnm.load matrix %A_t_b[%c0, %c0] {leadDimension = 16: index}
    : memref<16x16xi32> -> !cnm.matrix<16x16xi32>
  %B_dev = cnm.load matrix %B_t_b[%c0, %c0] {leadDimension = 16: index}
    : memref<16x16xi32> -> !cnm.matrix<16x16xi32>
  %C_part = cnm.op.gemm %A_dev, %B_dev: tensor<16x16xi32>, tensor<16x16xi32>
  %out_tile = arith.addf %in_tile, %C_part: tensor<16x16xi32>
  scf.yield %out_tile : tensor<16x16xi32>
}
2 %out_result = tensor.insert_slice %C_tile, %in_result[%o0, %o1][16, 16][1, 1]:
<...>
```

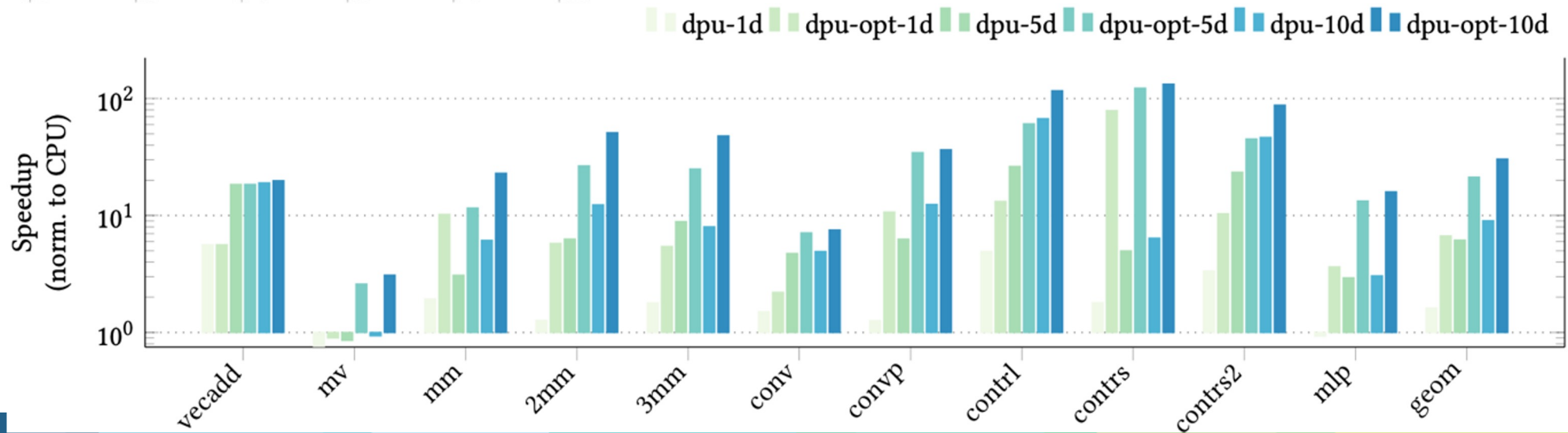


UPMEM example: Results



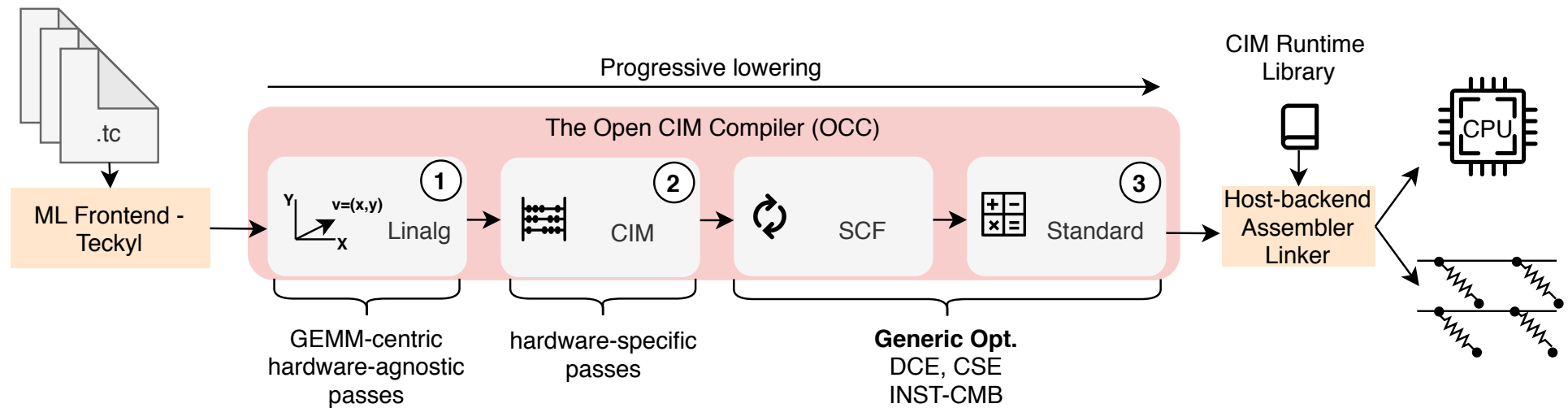
1-DIMM 128 DPUs
 5-DIMMs 640 DPUs
 10-DIMMs 1280-DPUs

6.1×, (1 DIMM)
 21.3× (5 DIMM) and
 30.4× (10 DIMM) wrt
 host CPU



CIM Pipeline: Cross-bar example

- ❑ MLIR dialect for general tensor expressions (Tensor Comprehensions)
 - ❑ Reuse GEMM transformations from linalg
 - ❑ Lower to CIM dialect (co-existing with SCF and Standard)
 - ❑ Lower CIM dialect to runtime APIs



A. Siemieniuk, L. Chelini, A. A. Khan, J. Castrillon, A. Drebes, H. Corporaal, T. Grosser, M. Kong, "OCC: An Automated End-to-End Machine Learning Optimizing Compiler for Computing-In-Memory", In IEEE TCAD, 2021

Lowering examples (somewhat beyond matmul)

```
def contr(int16(K,L,M) A, int16(L,K,N) B)
  -> (int16(M,N) C)
{
  C(m,n) += A(k,l,m) * B(l,k,n)
}
```

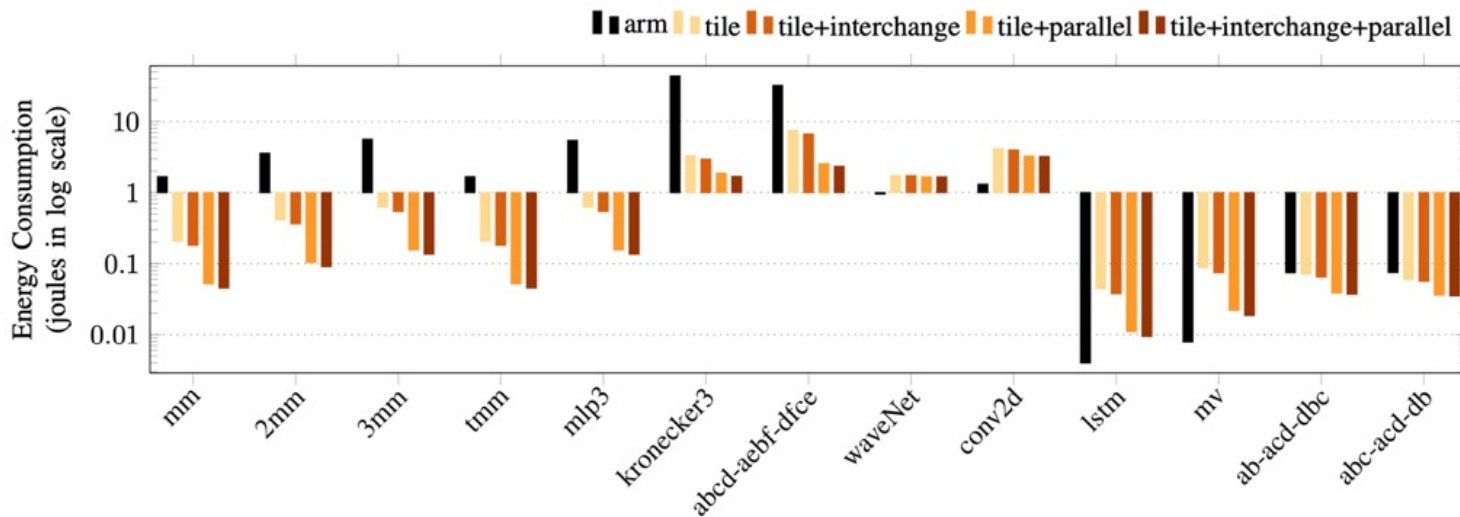
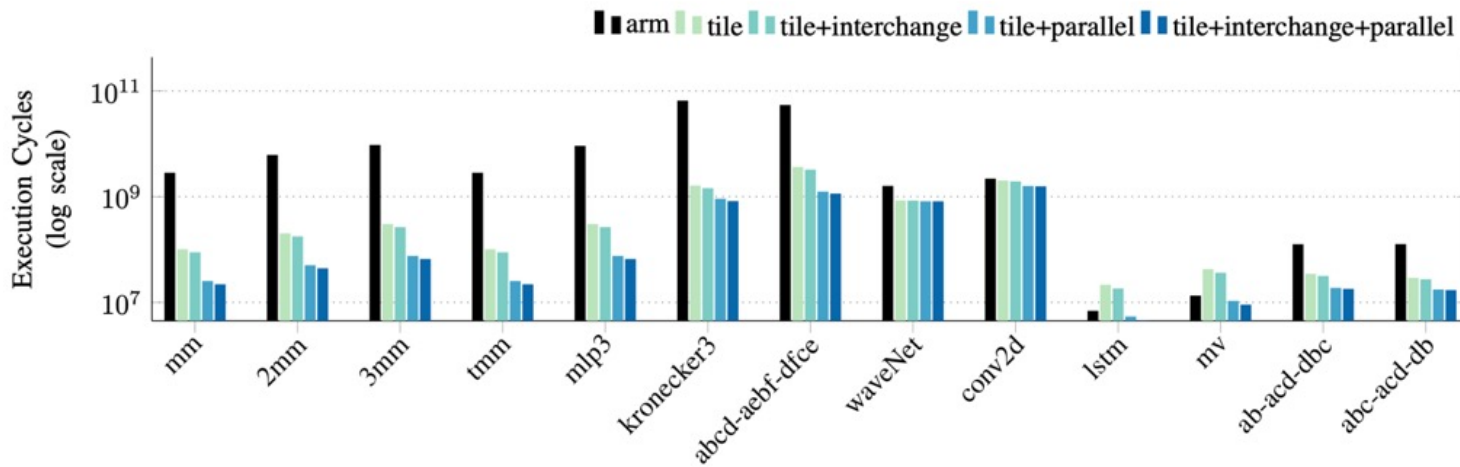
⇓ lowers to

```
%0 = linalg.transpose(%A, {2, 0, 1})
%1 = linalg.transpose(%B, {1, 0, 2})
%2 = linalg.reshape(%0, {0, {1, 2}})
%3 = linalg.reshape(%1, {{0, 1}, 2})
// eligible for offloading to CIM
linalg.matmul(%2, %3, %C)
```

```
// loop interchanged GEMM
scf.for %k = %c0 to %numTiles step %c1 {
  scf.for %j = %c0 to %tiledCols step %c1 {
    %tileB = cim.copyTile(%B, %k, %j)
    cim.write(%id, %tileB)
    scf.for %i = %c0 to %tiledRows step %c1 {
      %tileC = cim.copyTile(%C, %i, %j)
      ...
      cim.storeTile(%tileC, %C, %i, %j)
    }
  }
}
```

```
linalg.matmul(%A, %B, %C)
  ⇓ lowers to
// tiled GEMM in the CIM dialect
%c0 = constant 0 : i32
%c1 = constant 1 : i32
%id = constant 0 : i32 // tile id
scf.for %i = %c0 to %tiledRows step %c1 {
  scf.for %j = %c0 to %tiledCols step %c1 {
    %tileC = cim.copyTile(%C, %i, %j)
    %tempTile = cim.allocDuplicate(%tileC)
    scf.for %k = %c0 to %numTiles step %c1 {
      %tileA = cim.copyTile(%A, %i, %k)
      %tileB = cim.copyTile(%B, %k, %j)
      cim.write(%id, %tileB)
      cim.matmul(%id, %tileA, %tempTile)
      cim.barrier(%id)
      // tileC += tempTile
      cim.accumulate(%tileC, %tempTile)
    }
    cim.storeTile(%tileC, %C, %i, %j)
  }
}
```

Optimization results: Crossbars beyond matmult

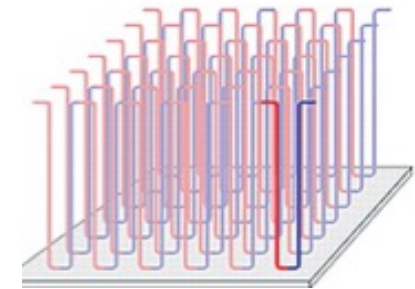
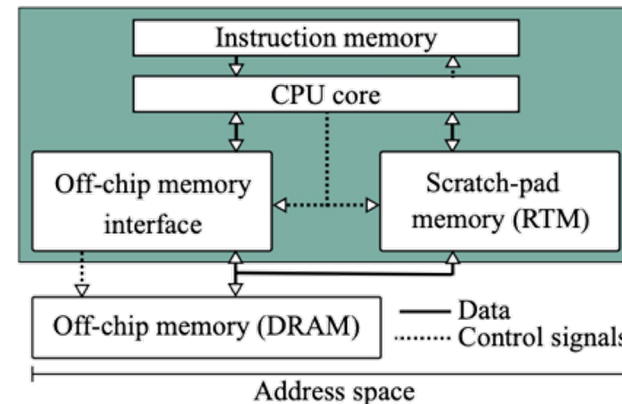


A. Siemieniuk, L. Learning Optim

Machine

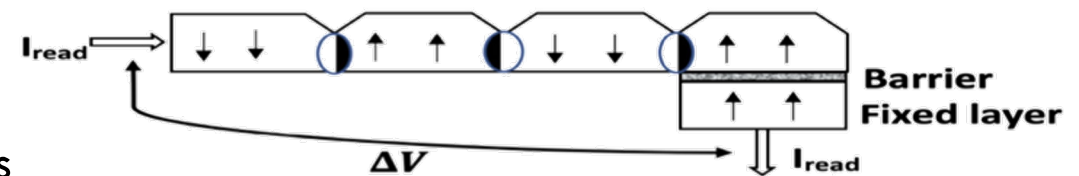
Racetrack memories

- ❑ Racetrack memories (RTM)
 - ❑ Extreme density
 - ❑ Sequential bit access per cell
 - ❑ Sequentiality on top of locality
 - ❑ Ex.: Placement for tensor contraction



- ❑ Different approaches for in-RTM computing proposed

- ❑ Example: Transverse reads
- ❑ Interesting data-allocation problems

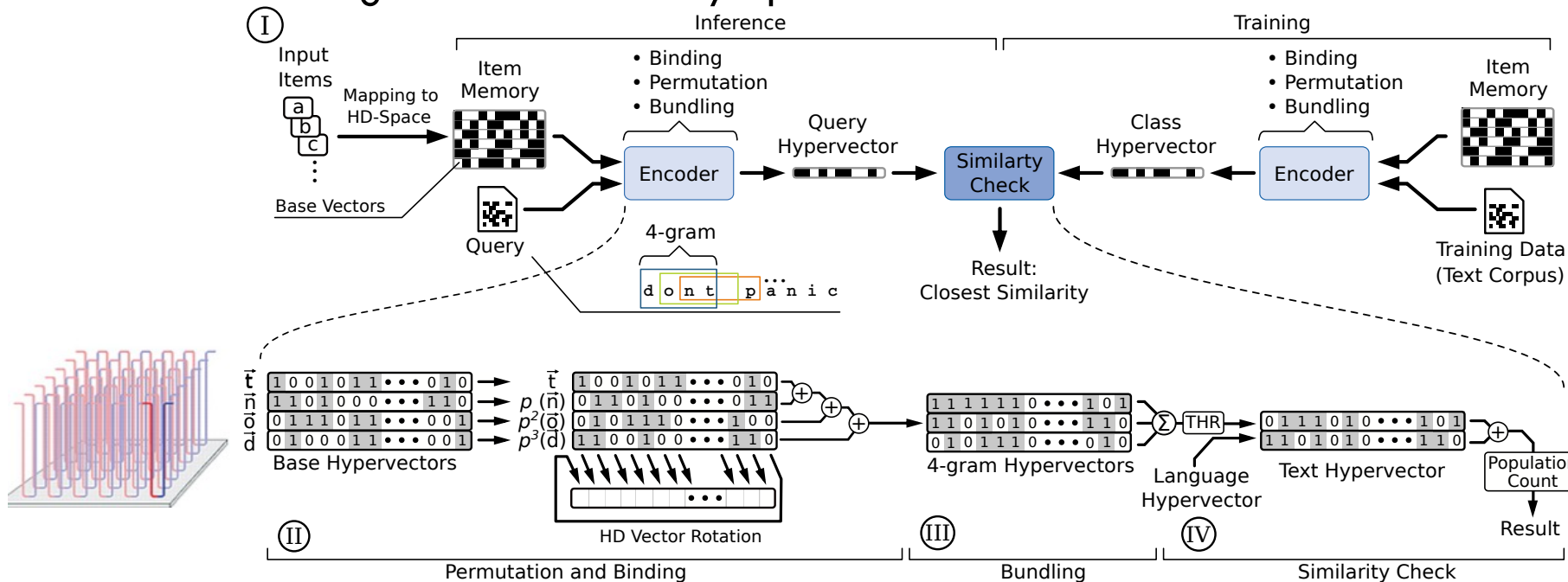


K. Roxy, IEEE T Nano 2020

Example: Hyper dimensional computing (HDC)

❑ HDC: Embed data in 10 k-dimensions – Von-Neuman Bottleneck!

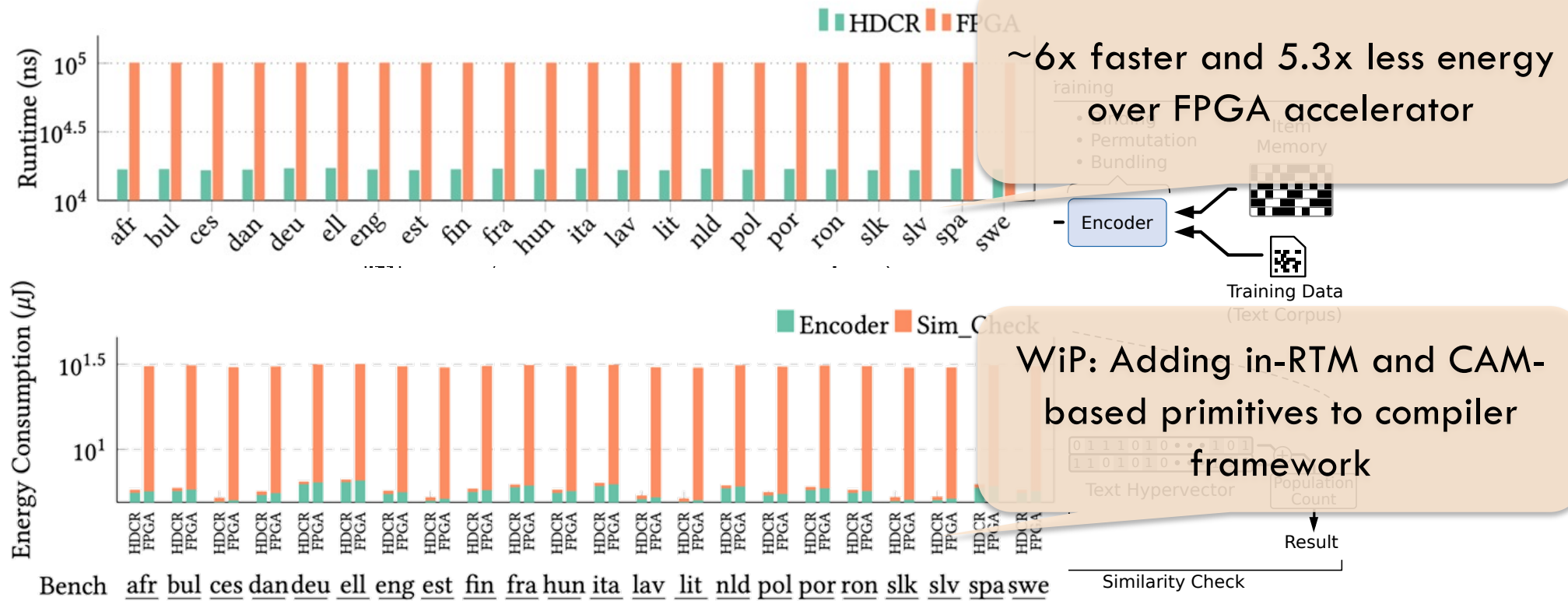
❑ Leverage bulk-wise binary operations



Khan, A. A., Ollivier, S., Longofono, S., Hempel, G., Castrillon, J., & Jones, A. K. (2022). Brain-inspired Cognition in Next Generation Racetrack Memories. In ACM TECS 2022

© Prof. J. Castrillon. iMACAW. 2023

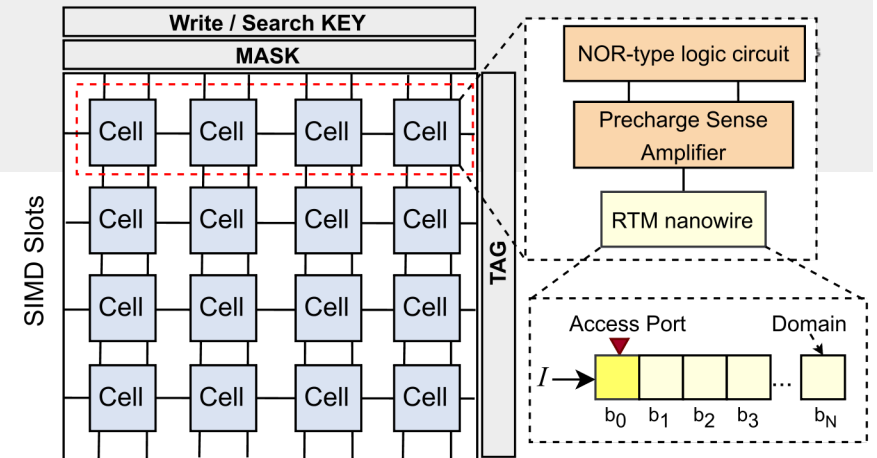
Example: Hyper dimensional computing (HDC)



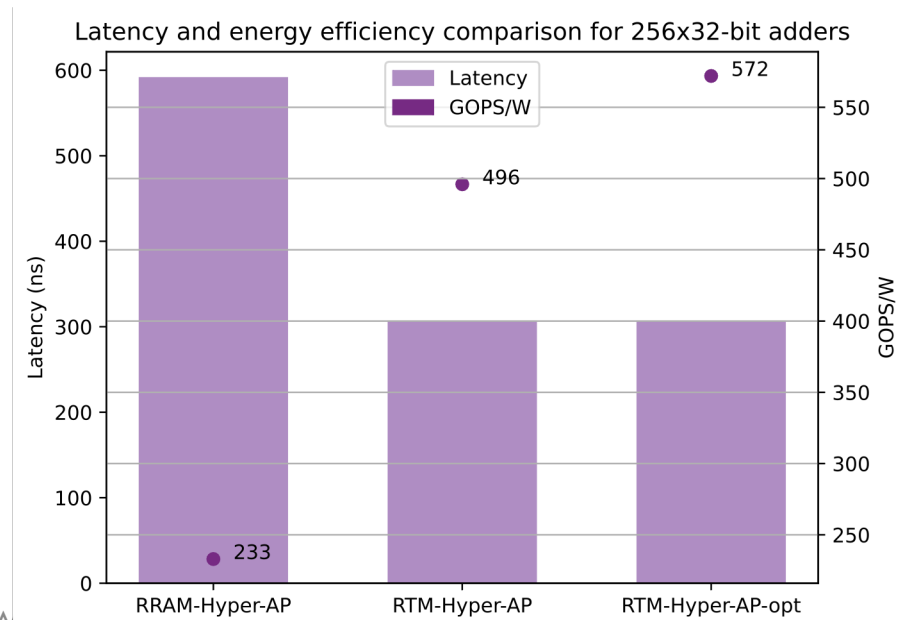
Khan, A. A., Ollivier, S., Longofono, S., Hempel, G., Castrillon, J., & Jones, A. K. (2022). "Brain-inspired Cognition in Next Generation Racetrack Memories" In ACM TECS 2022

Associative Processors with RTMs

- ❑ Search-write computational paradigm with CAM-based memories
- ❑ Use RTM nanowires for multi-bit serial storage in a cell
- ❑ Improved performance and energy efficiency!



J. P. C. de Lima, A. A. Khan, H. Farzaneh, J. Castrillon, "Efficient Associative Processing with RTM-TCAMs" 1st in-Memory Architectures and Computing Applications Workshop (iMACAW), 2pp, Jul 2023




Summary

- ❑ Challenging & exciting computing landscape!
- ❑ Abstractions important to target emerging and domain-specific systems
- ❑ Beyond infrastructure: Need execution models of non Von Neumann
- ❑ Beyond energy efficiency: understanding of full-life-cycle sustainability

$$v_{ijk,e} = \sum_{i'=0}^p \sum_{j'=0}^p \sum_{k'=0}^p A_{kk'} A_{jj'} A_{ii'} u_{i'j'k'e}$$

What we want



AI accelerator

<https://www.hpcwire.com/2017/04/10/nvidia-responds-google-tpu-benchmarking/>

Lee, Sukhan, et al. "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology: Industrial Product." ISCA 2021.

```

1 void cfd_kernel(
2   double A[restrict][7][7],
3   double u[restrict][216][7][7][7],
4   double v[restrict][216][7][7][7])
5 {
6   /* element loop: */
7   for(int e = 0; e < 216; e++) {
8     for(int i0 = 0; i0 < 7; i0++) {
9       for(int j0 = 0; j0 < 7; j0++) {
10        for(int k0 = 0; k0 < 7; k0++) {
11          v[e][i0][j0][k0] = 0.0;
12          for(int i1 = 0; i1 < 7; i1++) {
13            for(int j1 = 0; j1 < 7; j1++) {
14              for(int k1 = 0; k1 < 7; k1++) {
15                v[e][i0][j0][k0] += A[i0][i1]
16                  * A[j0][j1]
17                  * A[k0][k1];
18              }
19            }
20          }
21        }
22      }
23    }
24  }

```

100X

???

???


???

???

```

11 for(int e = 0; e < 216; e++) {
12   double t6[7][7];
13   /* 1st contraction: */
14   #pragma simd
15   for(int i0 = 0; i0 < 7; i0++) {
16     for(int i1 = 0; i1 < 7; i1++) {
17       /* #pragma simd */
18       for(int i2 = 0; i2 < 7; i2++) {
19         double t8 = 0.0;
20         for(int i3 = 0; i3 < 7; i3++)
21           t8 += A[i0][i3] * u[e][i1][i2][i3];
22         t6[i0][i1][i2] = t8;
23       } } /* end of 1st contraction */
24     } }

```



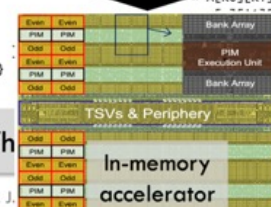
HBM-FPGA

```

11 for(int e = 0; e < 216; e++) {
12   double t10 = 0.0;
13   for(int i10 = 0; i10 < 7; i10++)
14     t10 += A[i8][i10] * t7[i9][i10][i11];
15   v[e][i8][i9][i10] = t10;
16 } } /* end of third contraction */

```

Why



In-memory accelerator

Why

starts code

© Prof. J. Castrillon. iMACAW. 2023

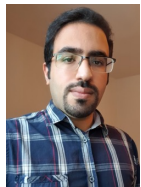
Thanks! & Acknowledgements



Hasna
Bouraoui



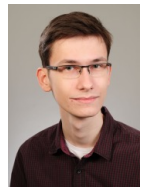
João P.
de Lima



Hamid
Farzaneh



Clément
Fournier



Karl
Friebe



Dr. Asif
Khan



Robert
Khasanov



Alexander
Brauckmann



Nesrine
Khouzami



Dr. Steffen
Köhler



Christian
Menard



Julian
Robledo



Lars
Schütze



Felix
Wittwer



Dr. Fazal
Hameed

..., and previous members of the group (Andres Goens, Norman Rink, Sven Karol, Sebastian Ertel), and collaborators (J. Fröhlich, I. Sbalzarini, T. Grosser, C. Pilato, S. Parkin)

Funded by

DFG Deutsche
Forschungsgemeinschaft

HetCIM (502388442), CO4RTM (450944241)

ScaDS.AI
DRESDEN LEIPZIG (901IS18026A)



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 957269

References

- [**RWDSL'18**] N. A. Rink, et al. "CFDlang: High-level code generation for high-order methods in fluid dynamics". 3rd International Workshop on Real World Domain Specific Languages (RWDSL 2018), Feb 2018.
- [**GPCE'18**] A. Susungi, et al. "Meta-programming for cross-domain tensor optimizations". Proceedings of 17th ACM SIGPLAN International Conference on Generative Programming: Concepts and Experiences (GPCE'18), Nov. 2018.
- [**ACM TOMS'18**] Karol, S., Nett, T., Castrillon, J., Sbalzarini, I. F. "A Domain-Specific Language and Editor for Parallel Particle Methods", ACM Transactions on Mathematical Software (TOMS), ACM, 2018, 44, 34:1-34:32.
- [**JCS'21**] N. Khouzami, F. Michel, P. Incardona, J. Castrillon, and I. F. Sbalzarini. "Model-based Autotuning of Discretization Methods in Numerical Simulations of Partial Differential Equations". In: Journal of Computational Science 57, 2021.
- [**ICCS'21**] N. Khouzami, L. Schütze, P. Incardona, L. Kraaz, T. Subic, J. Castrillon, and I. F. Sbalzarini. "The OpenPME Problem Solving Environment for Numerical Simulations". In: International Conference on Computational Science (ICCS'21).
- [**Array'19**] N.A. Rink, N. A. and J. Castrillon. "Tel: a type-safe imperative Tensor Intermediate Language", ACM SIGPLAN International Workshop on Libraries, Languages, and Compilers for Array Programming (ARRAY), ACM, 2019, pp. 57-68
- [**Cluster'21**] K. F. A. Friebel, S. Soldavini, G. Hempel, C. Pilato, J. Castrillon, "From Domain-Specific Languages to Memory-Optimized Accelerators for Fluid Dynamics", Proceedings of the FPGA for HPC Workshop, held in conjunction with IEEE Cluster 2021, Sep 2021
- [**TRETS'22**] S. Soldavini, K. F. A. Friebel, M. Tibaldi, G. Hempel, J. Castrillon, and C. Pilato. "Automatic Creation of High-Bandwidth Memory Architectures from Domain-Specific Languages: The Case of Computational Fluid Dynamics". In: ACM TRETS, Sept. 2022
- [**Heart'23**] K. F. A. Friebel, J. Bi, and J. Castrillon. "BASE2: An IR for Binary Numeral Types". In: 13th International Symposium on Highly-Efficient Accelerators and Reconfigurable Technologies (HEART 2023). June 2023.
- [**Arxiv'23**] A. A. Khan, H. Farzaneh, K. F. A. Friebel, L. Chelini, and J. Castrillon. "CINM (Cinnamon): A Compilation Infrastructure for Heterogeneous Compute In- Memory and Compute Near-Memory Paradigms". In Arxiv, Jan 2023.
- [**TCAD'21**] A. Siemieniuk, L. Chelini, A. A. Khan, J. Castrillon, A. Drebes, H. Corporaal, T. Grosser, M. Kong, "OCC: An Automated End-to-End Machine Learning Optimizing Compiler for Computing-In-Memory", In IEEE TCAD, 2021.
- [**TECS'22**] Khan, A. A., Ollivier, S., Longofono, S., Hempel, G., Castrillon, J., & Jones, A. K. (2022). Brain-inspired Cognition in Next Generation Racetrack Memories. In ACM TECS 2022
- [**iMACAW'23**] J. P. C. de Lima, A. A. Khan, H. Farzaneh, J. Castrillon, "Efficient Associative Processing with RTM-TCAMs" 1st in-Memory Architectures and Computing Applications Workshop (iMACAW), 2pp, Jul 2023