# Compiler Support for Ferroelectric Compute-in-Memory Solutions (and beyond)

Jeronimo Castrillon
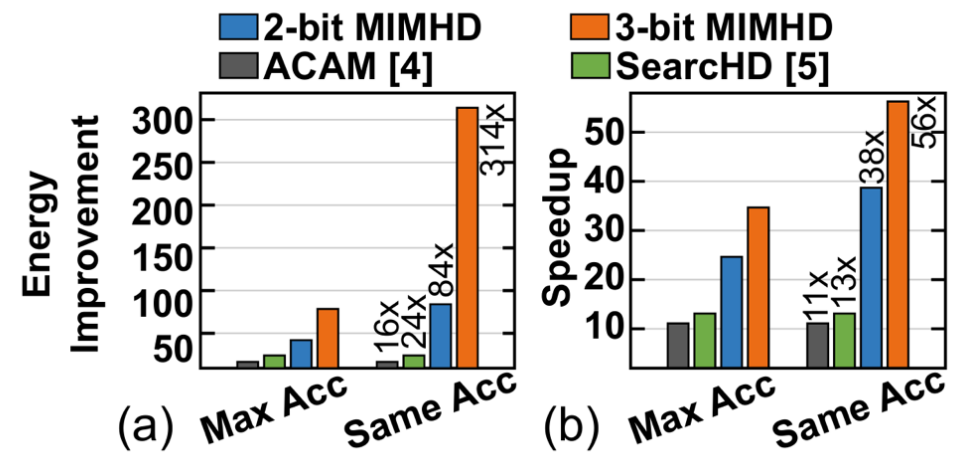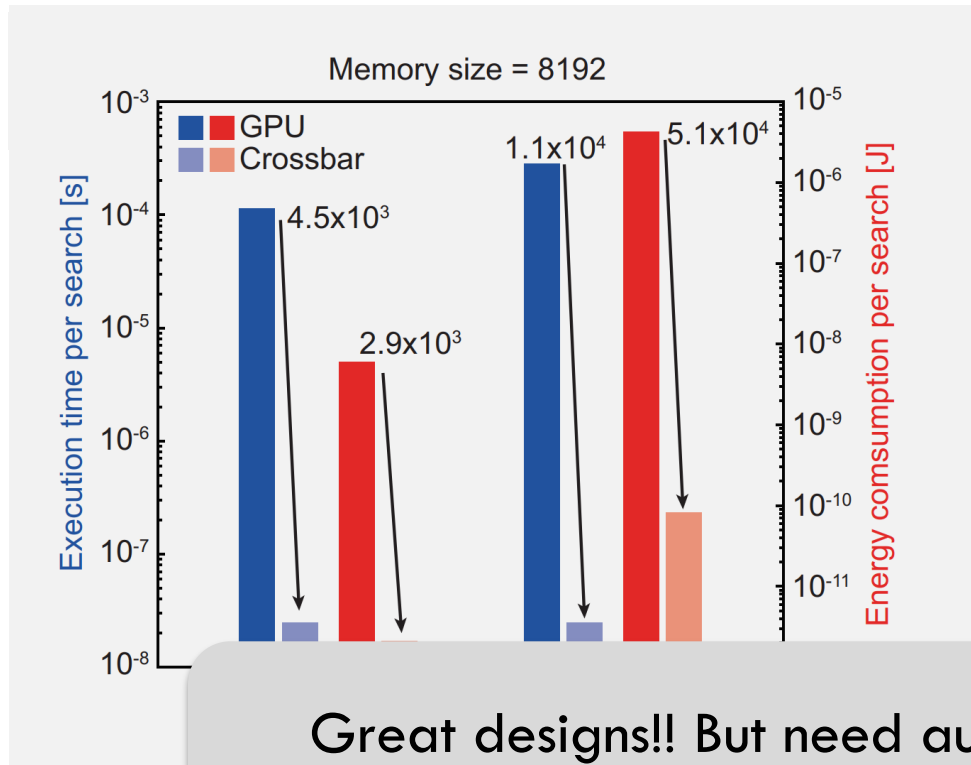
Chair for Compiler Construction (CCC), TU Dresden,
SCADS.AI Dresden/Leipzig & Center for Advancing Electronics (cfaed) Dresden

cfaed.tu-dresden.de

Memory size = 8192



GPU
Crossbar

Execution time per search [s]
Energy consumption per search [J]

$10^{-3}$
$10^{-4}$
$10^{-5}$
$10^{-6}$
$10^{-7}$
$10^{-8}$

$10^{-5}$
$10^{-6}$
$10^{-7}$
$10^{-8}$
$10^{-9}$
$10^{-10}$
$10^{-11}$

$4.5 \times 10^3$
$2.9 \times 10^3$
$1.1 \times 10^4$
$5.1 \times 10^4$



2-bit MIMHD
ACAM [4]
3-bit MIMHD
SearchHD [5]

Energy Improvement

300
250
200
150
100
50

314x
16x
24x
84x

(a) Max Acc    Same Acc

Speedup

50
40
30
20
10

56x
11x
13x
38x

(b) Max Acc    Same Acc

A. Kazemi, "Cross-Layer Design with Emerging Devices for Machine

Mao, Ruibin,
augmented

Great designs!! But need automation to **generalize** to other patterns and **optimize around kernels**

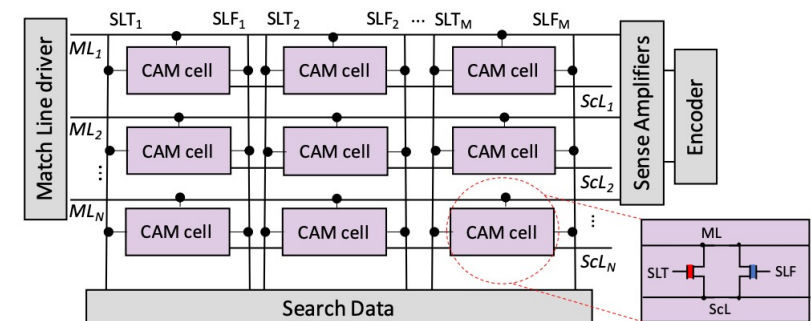CHAIR FOR COMPILER CONSTRUCTION

- ❑ Current practice
  - ❑ Low-level code: Hard to write, no formalization/generalization
  - ❑ Manual **app-by-app** optimization: Data mapping, synchronization, …

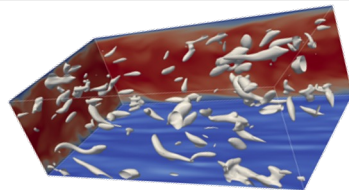- ❑ Towards high-level pythonic programming

```
1  i_flav = spectra.bnd_to_flav[i_strato][:,i_bnd]
2  i_eta  = j_eta[i_lay,i_flav,...]
3  i_eta  = np.stack((i_eta,i_eta+1), -1)[:,None,:,:]
4  a = n_prime_mix[None,i_lay,i_flav,None,:,None]
5  b = f_major[None,i_lay,i_flav,:,:,:]
6  c = spectra.k_major[gptS:gptE,i_p,i_T,i_eta]
7  result[gptS:gptE,:] = np.sum(a * b * c, (2,3,4))
```

© Prof. J. Castrillon. DATE Conference: W06. 2025

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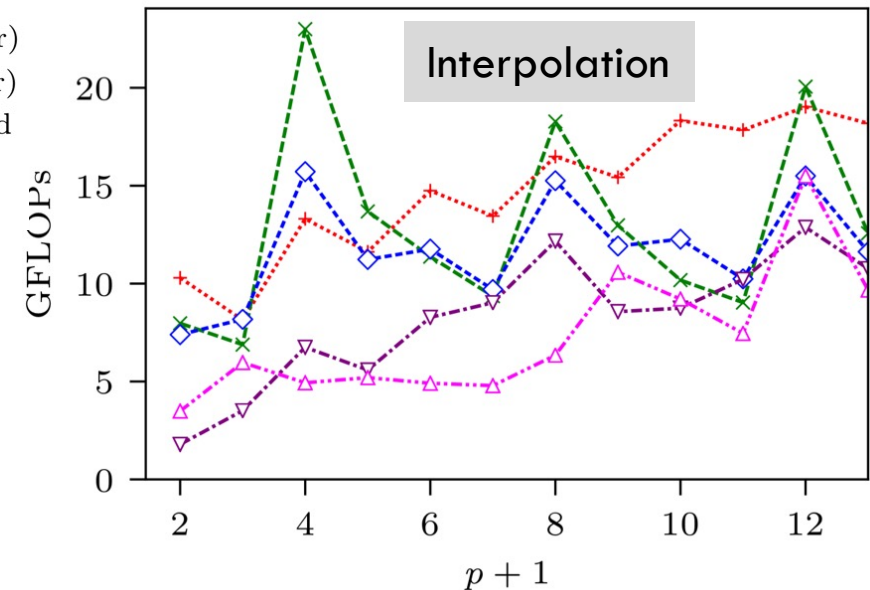
```
auto A = Matrix(m, n), B = Matrix(m, n),
     C = Matrix(m, n);
auto u = Tensor<3>(n, n, n);
auto v = (A*B*C)(u);
```

Legend:
- +----+ CFDlang(outer)
- ×--×- CFDlang(inner)
- ◇--◇ hand-optimized
- ▽--▽ DGEMM
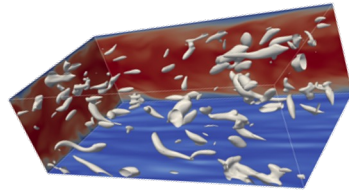- △--△ specialized



Interpolation

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, 79-92.

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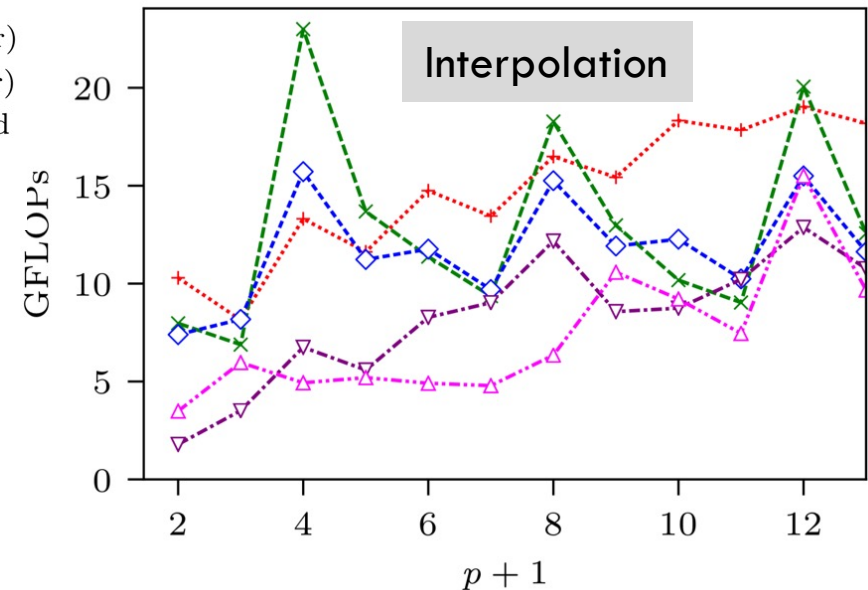
☐ CFDlang

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

```
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 .
```



CFDlang(outer)
CFDlang(inner)
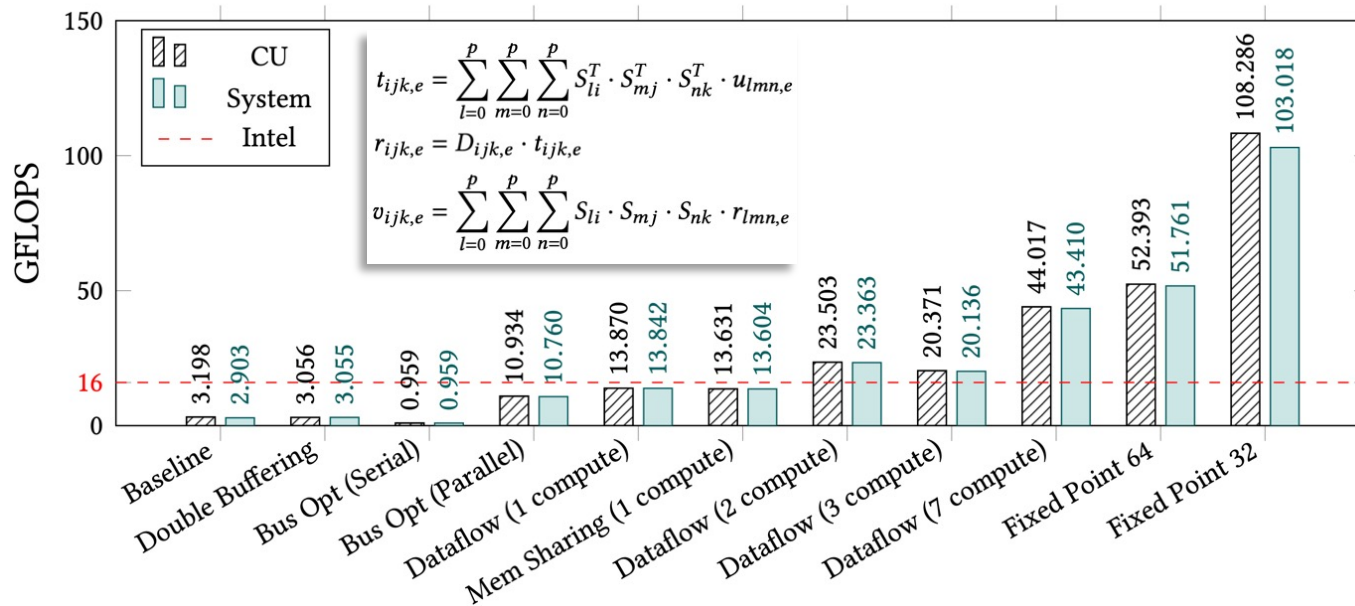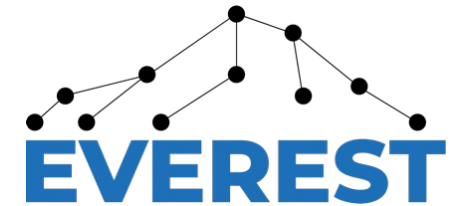hand-optimized
DGEMM
specialized

Interpolation

GFLOPs

$p + 1$

Plenty of other great DSL examples, e.g., Spiral, TACO, Halide, Lift, Firedrake, ML frameworks, … for CPUs, GPUs and TPUs.

```
auto v = (A*B*C)(u);
```

CHAIR FOR COMPILER CONSTRUCTION

# Own example for complex designs on HBM FPGA

- ❑ Transformations for a **17x speedup** (same precision)
- ❑ Variants with up to **24x better energy efficiency**



EVEREST

$$t_{ijk,e} = \sum_{l=0}^{p} \sum_{m=0}^{p} \sum_{n=0}^{p} S_{li}^{T} \cdot S_{mj}^{T} \cdot S_{nk}^{T} \cdot u_{lmn,e}$$

$$r_{ijk,e} = D_{ijk,e} \cdot t_{ijk,e}$$

$$v_{ijk,e} = \sum_{l=0}^{p} \sum_{m=0}^{p} \sum_{n=0}^{p} S_{li} \cdot S_{mj} \cdot S_{nk} \cdot r_{lmn,e}$$
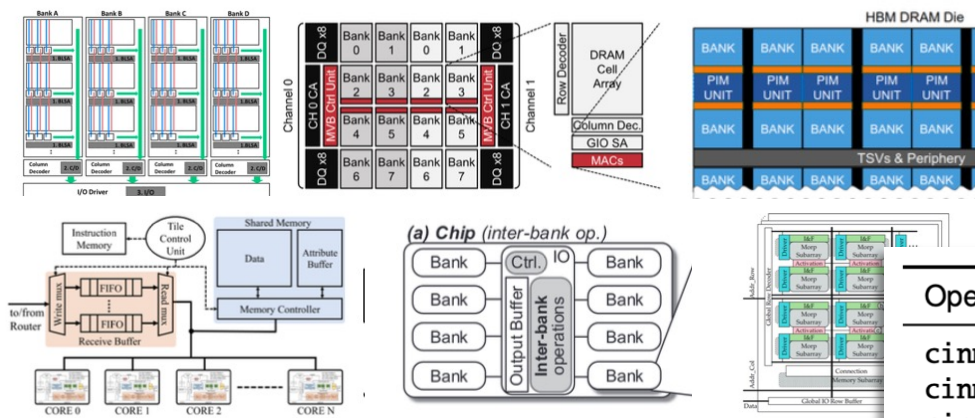
HBM DRAM | GPU

HBM+FPGA

C. Pilato, et al. "EVEREST: A design environment for extreme-scale big data analytics on heterogeneous platforms", DATE 2021
S. Soldavini, K. F. A. Friebel, et al. "Automatic Creation of High-Bandwidth Memory Architectures from Domain-Specific Languages: The Case of Computational Fluid Dynamics". In: ACM TRETS, 2023.

CHAIR FOR COMPILER CONSTRUCTION

# Colorful landscape

- ❑ Lots **in and near memory** systems!



### The Landscape of Compute-near-memory and Compute-in-memory: A Research and Commercial Overview

ASIF ALI KHAN, TU Dresden, Germany
JOÃO PAULO C. DE LIMA, TU Dresden and ScaDS.AI, Germany
HAMID FARZANEH, TU Dresden, Germany
JERONIMO CASTRILLON, TU Dresden and ScaDS.AI, Germany

In today's data-centric world, where data fuels numerous application domains, with machine learning at the

A. Khan, et al "The Landscape of Compute-near-memory and Compute-in-memory: A Research and Commercial Overview." arXiv:2401.1442 (2024)
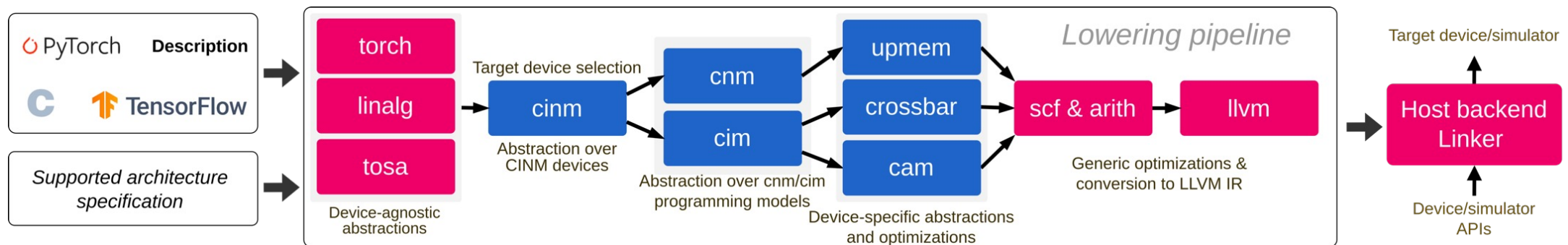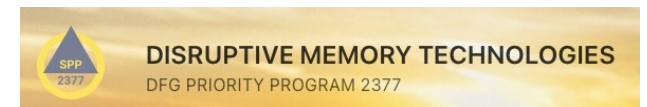
- ❑ Commonalities
  - ❑ Hierarchical HW
  - ❑ Common high-level operations

A. Khan et al, "CINM (Cinnamon): A Compilation Infrastructure for Heterogeneous Compute In-Memory and Compute Near-Memory Paradigms", ASPLOS'25

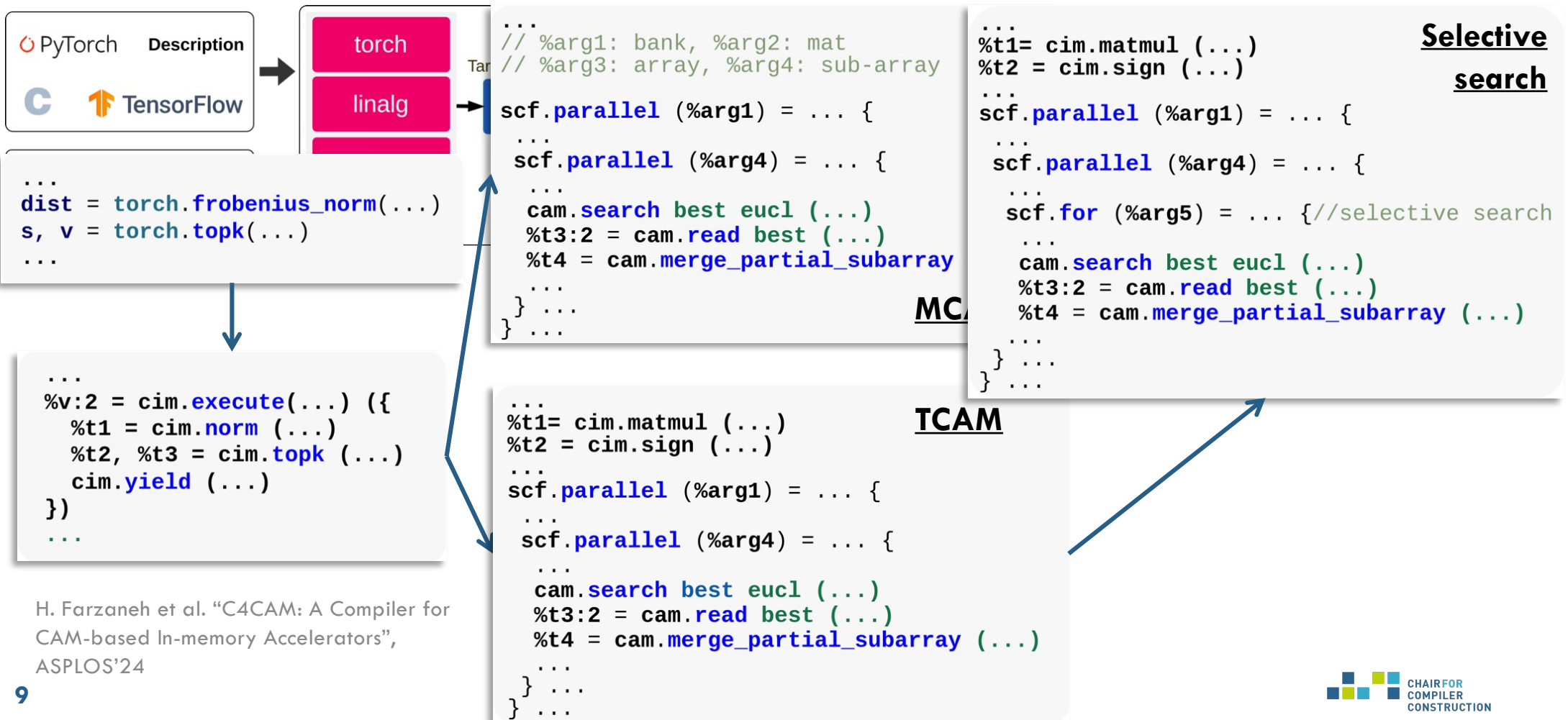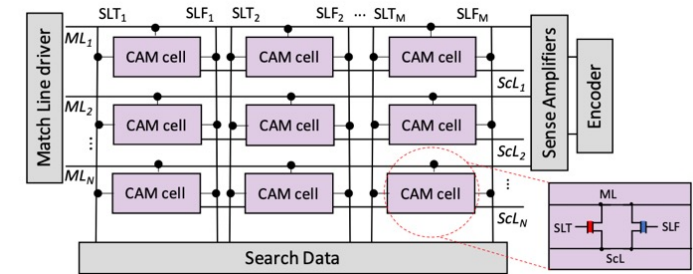| Operation | Type |
|---|---|
| cinm.{add,sub,mul,div,min,max}(%lhs, %rhs) | $T \times T \to T$ |
| cinm.{and,or,xor,not}(%lhs, %rhs) | $T \times T \to T$ |
| cinm.gemv(%lhs, %rhs) | $S^{m \times n} \times S^n \to S^m$ |
| cinm.gemm(%lhs, %rhs) | $S^{m \times k} \times S^{k \times n} \to S^{m \times n}$ |
| cinm.transpose(%in, %perms) | $S^n \times \mathbb{N}^n \to S'$ |
| cinm.{histogram,majority}(%in) | $S^n \to S^k$ |
| cinm.topk(%in, %k) | $S^n \times \mathbb{N} \to S^k \times \mathbb{N}^k$ |
| cinm.simSearch #E, #k (%in1, %in2) | $E \times \mathbb{N}^k \times S^n \times S^n \times \mathbb{N} \to S^k$ |
| cinm.mergePartial #op #dir (%lhs, %rhs) | $E \times D \times T \times T \to T$ |
| cinm.popCount(%in) | $T \to \mathbb{N}$ |
| cinm.reduce #op (%in) | $E \times S^n \to S$ |
| cinm.scan #op (%in) | $E \times S^n \to S^n$ |

# CINM: Generalized MLIR infrastructure

- ❑ From linear algebra abstractions (common to ML frameworks and beyond)
- ❑ Intermediate languages for **in and near memory computing**
- ❑ **Pattern recognition**, target-specific **models and optimizations**
- ❑ Targets: memristive crossbars, CAMs, logic in memory, UPMEM, …



A. Khan et al, "CINM (Cinnamon): A Compilation Infrastructure for Heterogeneous Compute In-Memory and Compute Near-Memory Paradigms", ASPLOS'25

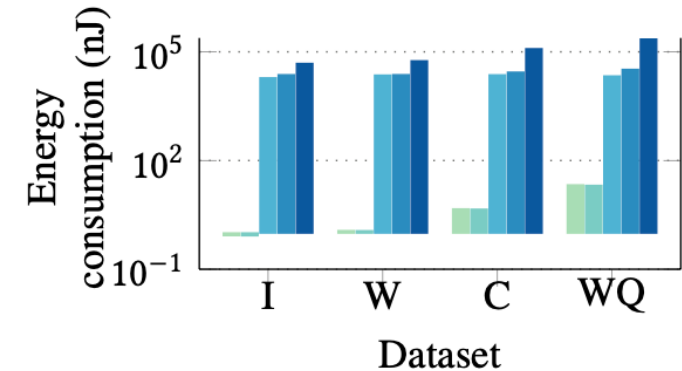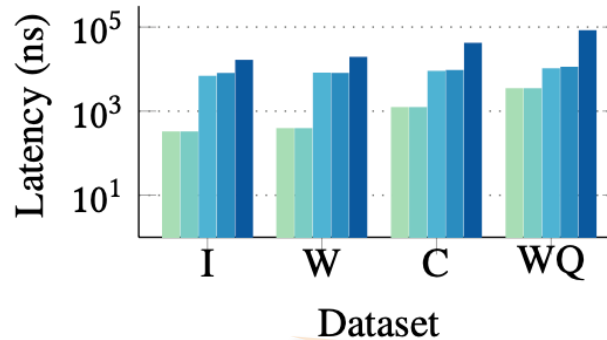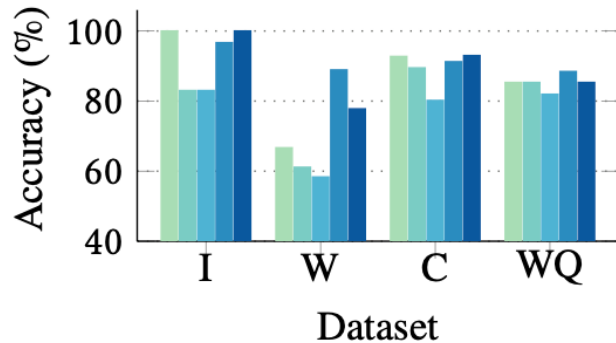# Lowering for different CAM-based accelerators

PyTorch | Description

C | TensorFlow

torch

linalg

```
...
dist = torch.frobenius_norm(...)
s, v = torch.topk(...)
...
```

```
...
%v:2 = cim.execute(...) ({
  %t1 = cim.norm (...)
  %t2, %t3 = cim.topk (...)
  cim.yield (...)
})
...
```
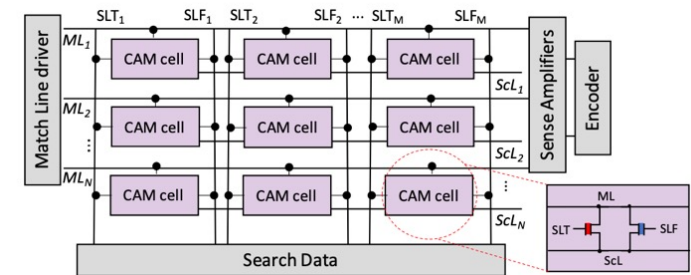
H. Farzaneh et al. "C4CAM: A Compiler for CAM-based In-memory Accelerators", ASPLOS'24

```
...
// %arg1: bank, %arg2: mat
// %arg3: array, %arg4: sub-array

scf.parallel (%arg1) = ... {
 ...
 scf.parallel (%arg4) = ... {
  ...
  cam.search best eucl (...)
  %t3:2 = cam.read best (...)
  %t4 = cam.merge_partial_subarray
  ...
 } ...
} ...
```
**MCA**

```
...
%t1= cim.matmul (...)
%t2 = cim.sign (...)
...
scf.parallel (%arg1) = ... {
 ...
 scf.parallel (%arg4) = ... {
  ...
  cam.search best eucl (...)
  %t3:2 = cam.read best (...)
  %t4 = cam.merge_partial_subarray (...)
  ...
 } ...
} ...
```
**TCAM**

**Selective search**

```
...
%t1= cim.matmul (...)
%t2 = cim.sign (...)
...
scf.parallel (%arg1) = ... {
 ...
 scf.parallel (%arg4) = ... {
  ...
  scf.for (%arg5) = ... {//selective search
   ...
   cam.search best eucl (...)
   %t3:2 = cam.read best (...)
   %t4 = cam.merge_partial_subarray (...)
   ...
  } ...
 } ...
} ...
```

- Pattern matching in high-level TorchScript code
- Automatic flow **matches manual designs**



H. Farzaneh, et al. "C4CAM: A Compiler for CAM-based In-memory Accelerators", ASPLOS, 2024

KNN results (128x128 CAM): **14x faster and ~$10^4$ less energy** compared to GPU

- Retargetable compiler for CAM exploration: Sizes and features
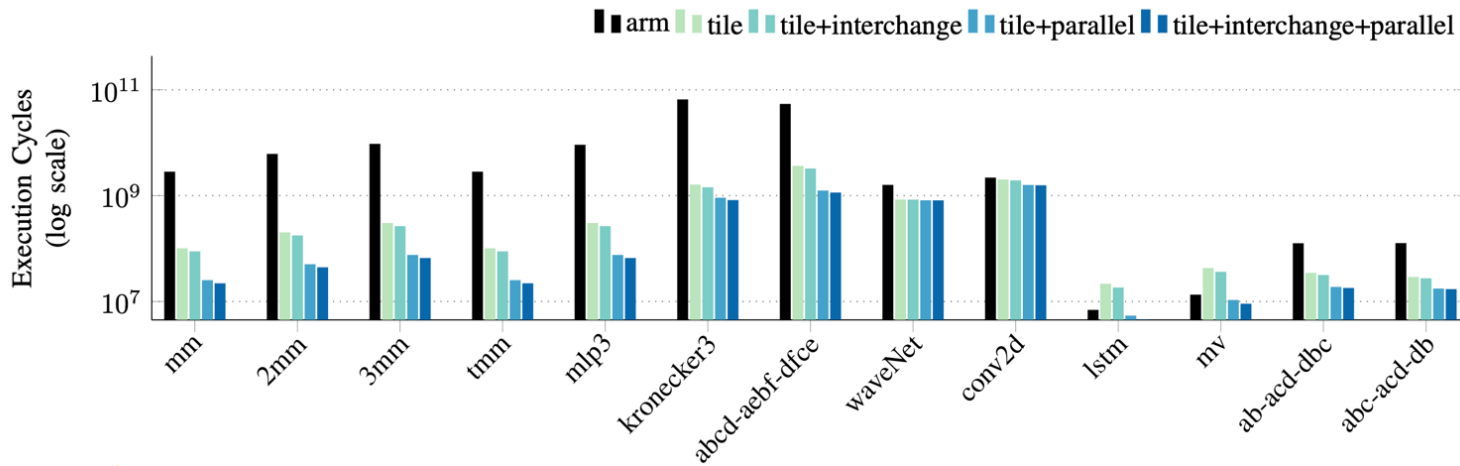- Compiler flags for optimization target
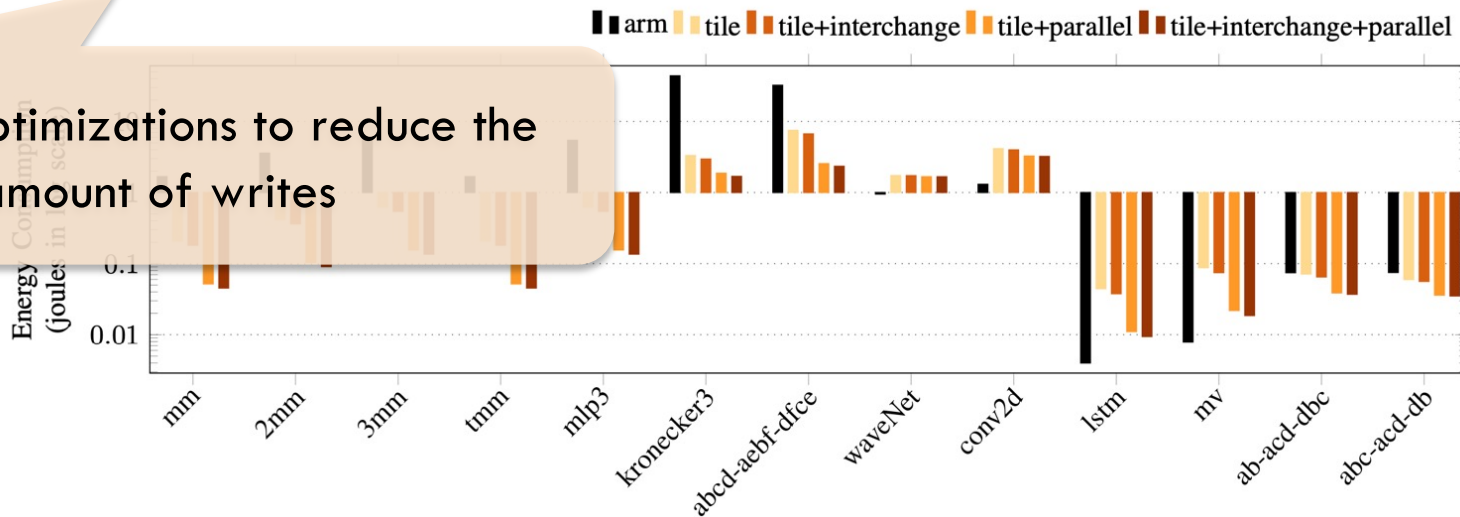


cam-base ■ cam-density ■ cam-power ■ cam-density+power



H. Farzaneh, et al. "C4CAM: A Compiler for CAM-based In-memory Accelerators", ASPLOS, 2024

Different flags expose trade-offs w/o manual re-coding.

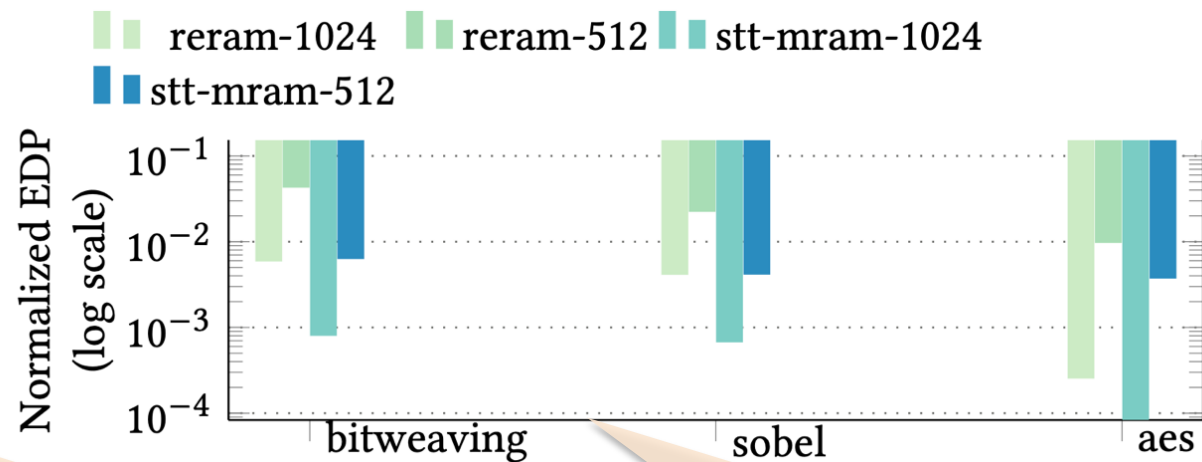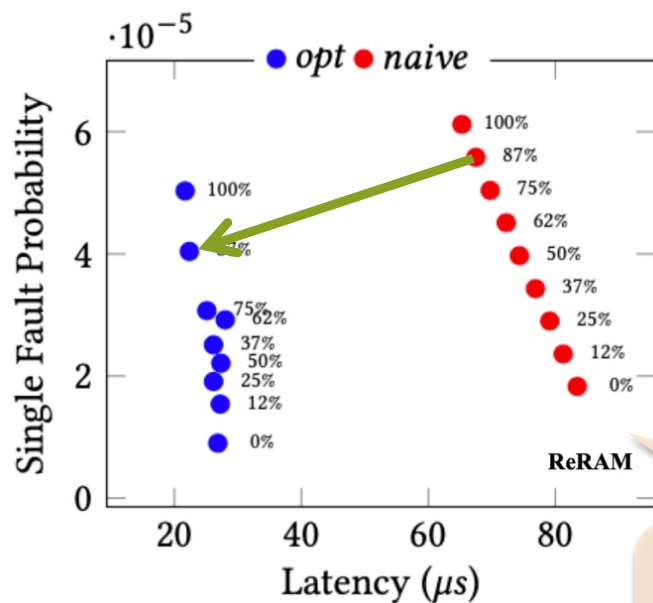# Optimization results: Crossbars beyond matmult



Studied optimizations to reduce the amount of writes

- Massively parallel multi-operand bit-wise operations in-memory
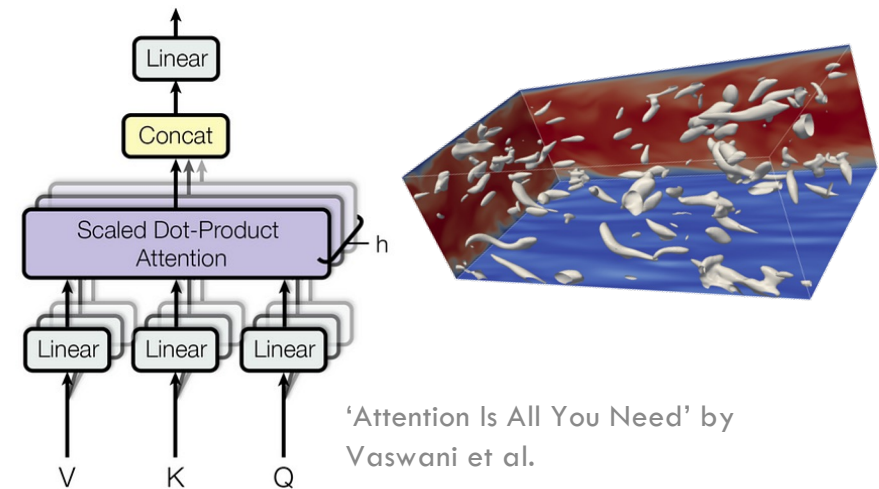- Complex mapping of operands, operations and temporaries to columns

H. Farzaneh, et al. ""SHERLOCK: Scheduling Efficient and Reliable Bulk Bitwise Operations in NVMs", DAC 2024

Optimized mapping: **Less latency (3x), better reliability (~1.4x)**

Orders of magnitude better EDP vs CPU baseline

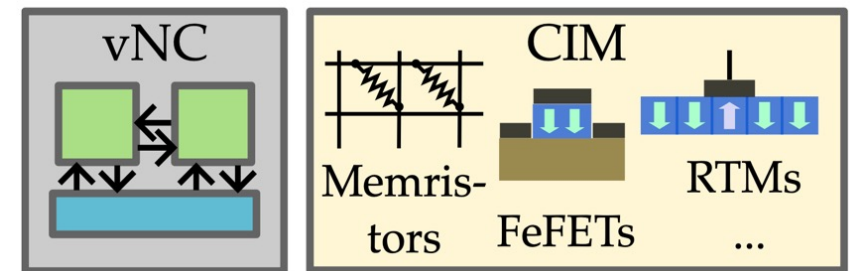© Prof. J. Castrillon. DATE Conference: W06. 2025

# Summary: Scratching the surface

❑ **Abstraction and compilation infrastructure**

    ❑ Capture "HW semantics" for execution

    ❑ Single flows: Crossbars, FeFET CAMs, Logic

    ❑ Leverage domain-specific abstractions

    ❑ Automated practices from manual designs

'Attention Is All You Need' by Vaswani et al.

❑ **Upcoming and challenges**

    ❑ End-to-end mapping on heterogeneous CIM

    ❑ Models (time, energy, endurance, resilience)

    ❑ Truly cross layer down to device parameters

    ❑ Runtime reconfigurability of mem arrays

# Thanks! & Acknowledgements

# References

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

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

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

**[DATE'21]** C. Pilato, et al. "EVEREST: A design environment for extreme-scale big data analytics on heterogeneous platforms", DATE 2021

**[TRETS'23]** 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, March 2023.

**[ASPLOS'25]** A. Khan et al, "CINM (Cinnamon): A Compilation Infrastructure for Heterogeneous Compute In-Memory and Compute Near-Memory Paradigms", arXiv, Aug 2023

**[ArXiv'24]** A. Khan, et al "The Landscape of Compute-near-memory and Compute-in-memory: A Research and Commercial Overview." arXiv:2401.1442 (2024)

**[TCAD'21]** A. Siemieniuk, et al. "OCC: An Automated End-to-End Machine Learning Optimizing Compiler for Computing-In-Memory", IEEE TCAD, 2021

**[ASPLOS'24]** H. Farzaneh, et al. "C4CAM: A Compiler for CAM-based In-memory Accelerators", ASPLOS 2024

**[DAC'24]** H. Farzaneh, et al. "SHERLOCK: Scheduling Efficient and Reliable Bulk Bitwise Operations in NVMs", DAC'24