# Discovering the In-Memory Kernels of 3D Dot-Product Engines

Muhammad Rashedul Haq Rashed[1], Sumit Kumar Jha[2], Rickard Ewetz[1]
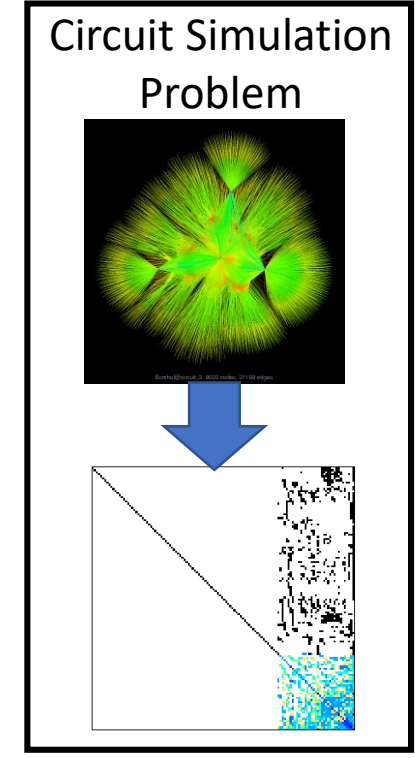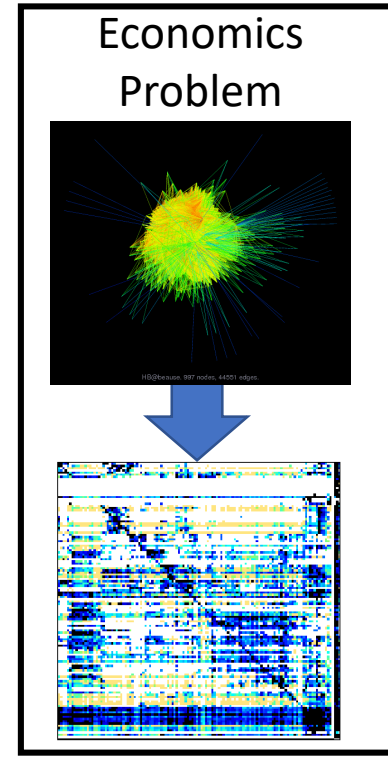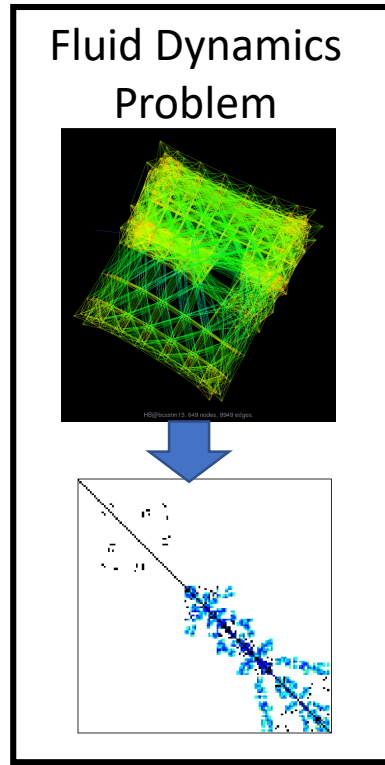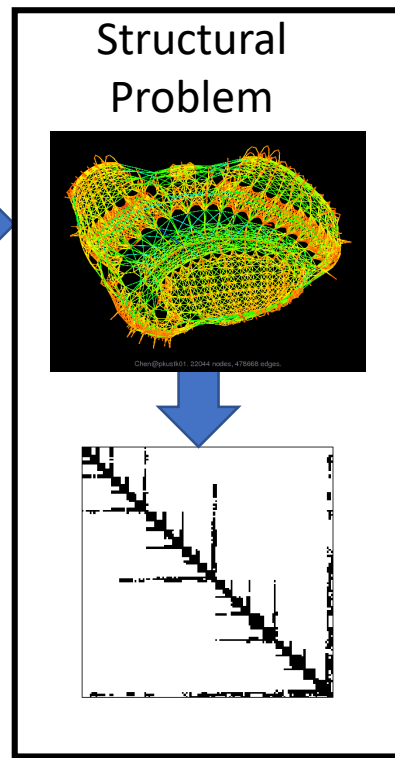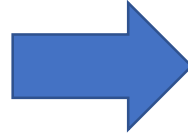
[1]University of Central Florida, [2]University of Texas at San Antonio

# Outline

- Preliminaries
- 3D DPE Level
  - Sharing of hardware
  - Hardware mapping
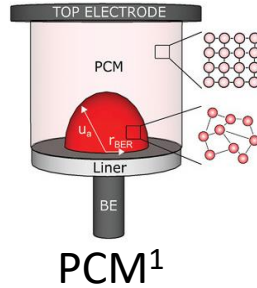- Application level
  - Methodology
- Evaluation
- Summary

# Why In-Memory Computing?

**Data-Intensive Systems**

| Structural Problem | Fluid Dynamics Problem | Economics Problem | Circuit Simulation Problem |
|---|---|---|---|

**CPU** ↔ **Expensive Data-Transfer** ↔ **Memory**

1. Davis, Timothy A., and Yifan Hu. "The University of Florida sparse matrix collection." *ACM Transactions on Mathematical Software (TOMS)* 38.1 (2011): 1-25.
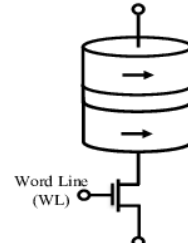
# NVM Technology and In-Memory Computing
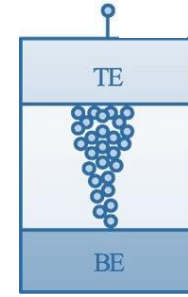
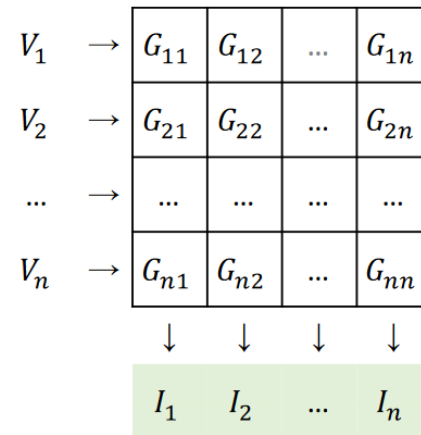## NVM Technologies



PCM[1]

STT-RAM[2]
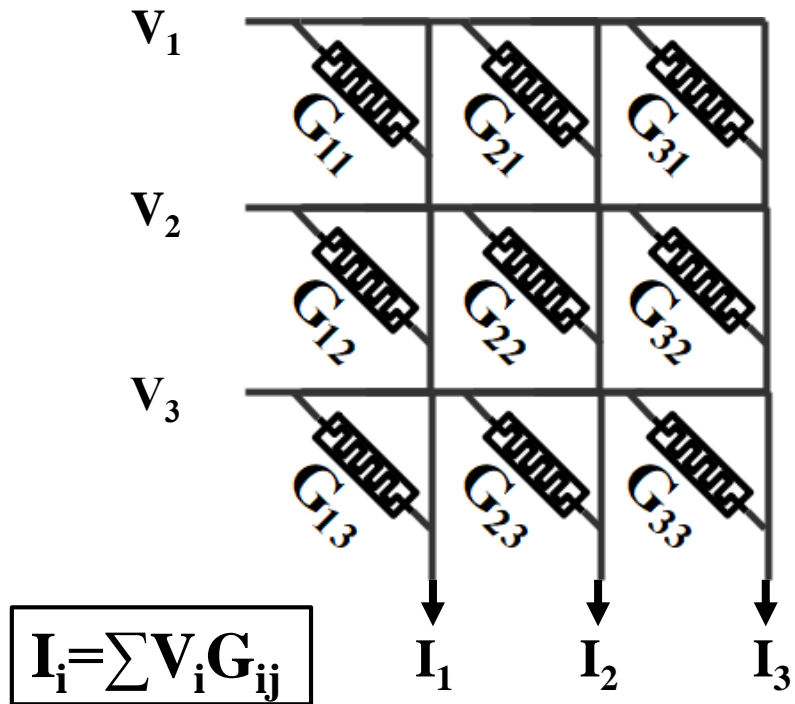
ReRAM

### Matrix Vector Multipication[3]



*Pro*

- In-situ computation
- Extremely energy efficient

1. Ghazi Sarwat, Syed, et al. "Projected Mushroom Type Phase-Change Memory." *Advanced Functional Materials* 31.49 (2021): 2106547.
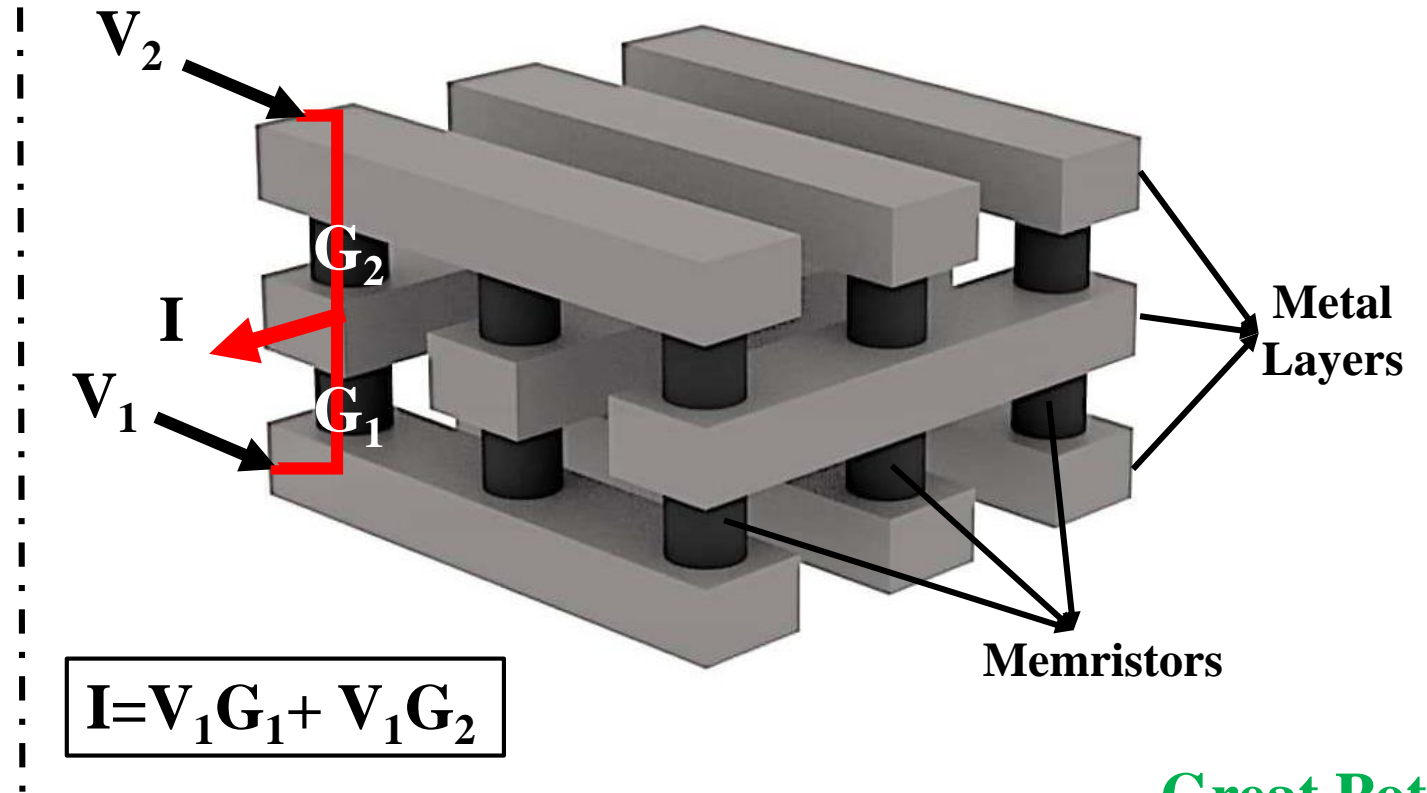2. Seyedzadeh, Seyed Mohammad & Maddah, Rakan & Jones, Alex & Melhem, Rami. (2016). Leveraging ECC to Mitigate Read Disturbance, False Reads and Write Faults in STT-RAM.
3. 1. Wang, Feng, Guangyu Sun, and Guojie Luo. "SSR: A Skeleton-based Synthesis Flow for Hybrid Processing-in-RRAM Modes." *2021 ICCAD*. IEEE, 2021.

# 2D vs 3D DPE



$$I_i = \sum V_i G_{ij}$$

**2D DPE**

$$I = V_1 G_1 + V_1 G_2$$

Metal Layers

Memristors

**3D DPE**

Great Potential for hardware sharing

# Previous Work

## Table 1: Architectural comparison of dot-product engines

| Work in | Form of DPE | # Metal layers |
|---|---|---|
| [9, 14, 18] | 2D | 2 |
| [1, 4, 6, 13] | 3D | 2-3 |
| **This work** | 3D | **2-8** |

- Previous works focused on mitigation of parasitics on the algorithmic level

- DAC and ADC operations are time multiplexed

- Efficient  utilization of 3D ReRAM crossbars with more than three layers have not been explored

[1] G. C. Adam et al. 3d reram arrays and crossbars: Fabrication, characterization and applications. In *2017 Ieee-Nano*, pages 844–849. IEEE, 2017

[4] B. Chakrabarti et al. A multiplyadd engine with monolithically integrated 3d memristor crossbar/cmos hybridcircuit. *Scientific reports*, 7(1):1–10, 2017.

[6] F. Chen et al., A novel zero-free dataflow accelerator for generative adversarial networks in 3d reram. In *Proceedings of the 56th DAC 2019*, pages 1–6, 2019.
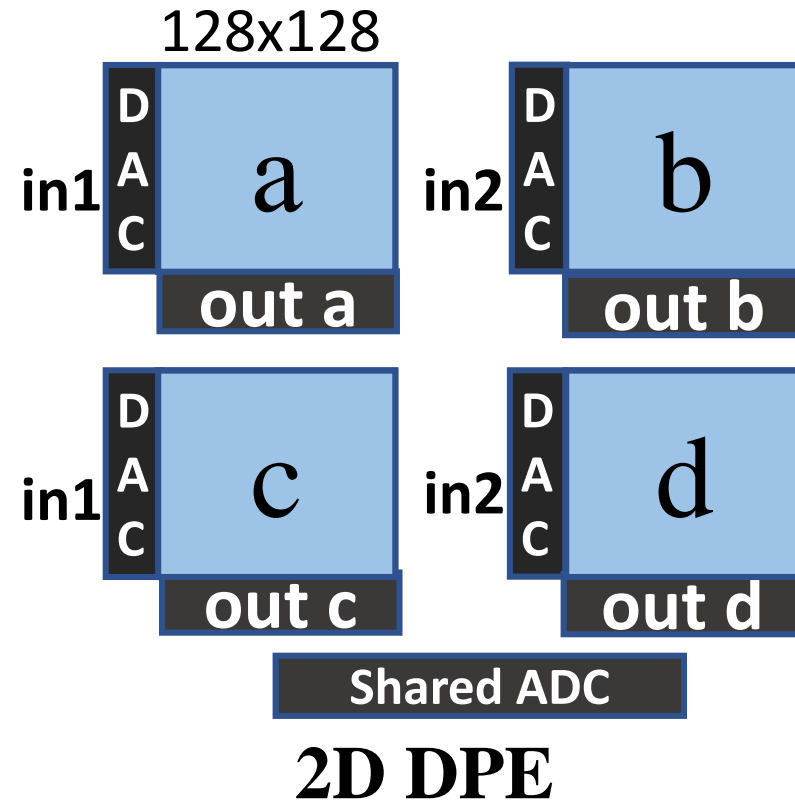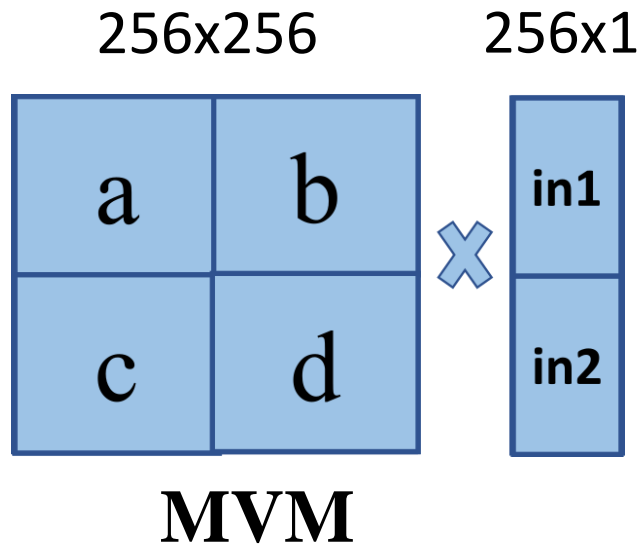
[9] B. Feinberg et. Al, Enabling scientific computing on memristive accelerators. In *2018 ACM/IEEE 45th ISCA*, pages 367–382. IEEE, 2018.

[13] M. A. Lastras-Montano et. Al, 3d-dpe: A 3d high-bandwidth dot-product engine for high-performance neuromorphic computing. In *2017 DATE*, pages 1257–1260. IEEE, 2017.
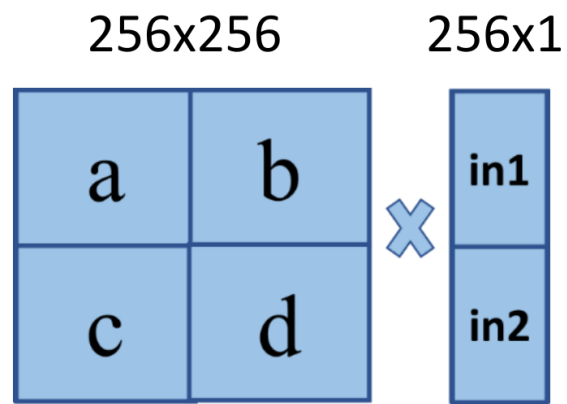
[14] C. Li et al. Analogue signal and image processing with large memristor crossbars. *Nature Electronics*, 1(1):52, 2018.

[18] A. Shafiee et. al., : A convolutional neural network accelerator with in-situ analog arithmetic in crossbars. *ACM SIGARCH Computer Architecture News*, 44(3):14–26, 2016.
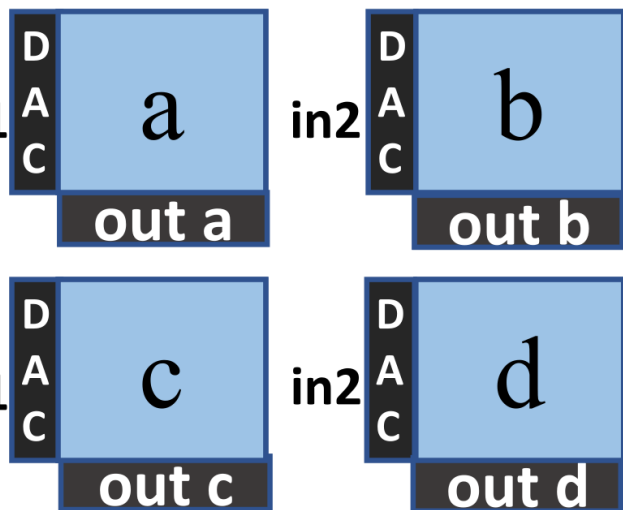
# MVM in In-memory DPE

256x256    256x1



**MVM**

128x128



in1  D A C  a    in2  D A C  b

out a    out b

in1  D A C  c    in2  D A C  d

out c    out d

Shared ADC

**2D DPE**

**Hardware constraint:**
**Maximum Crossbar Dimension**

256x256 · 256x1

a   b   × in1

c   d     in2

**MVM**

128x128

in1 DAC a · out a

in2 DAC b · out b

in1 DAC c · out c

in2 DAC d · out d

Shared ADC

**2D DPE**

in2

in1

Shared DAC 2

Shared DAC 1

L5

C4 · out c+

L4

C3 · out d

L3 · C2 · out a+

L2 · out b

C1

L1

Shared ADC

**3D DPE**

# Overhead Comparison

**Table 2: Overhead comparison of 2D vs 3D architecture.**

|  | 2D Architecture | 3D Architecture |
|---|:---:|:---:|
| Crossbar footprint on chip | 4X | 1X |
| # DACs | 4X | 2X |
| ADC timesteps | 4X | 2X |

# Problem 1: Hardware Mapping

**Matrix Segment**

**3D Hardware**
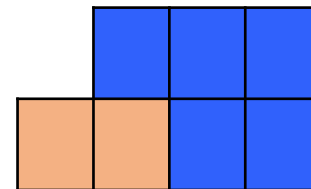
$L_8$
$L_7$
$L_6$
$L_5$
$L_4$
$L_3$
$L_2$
$L_1$

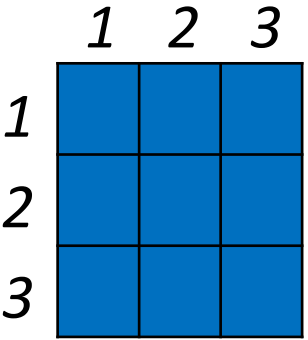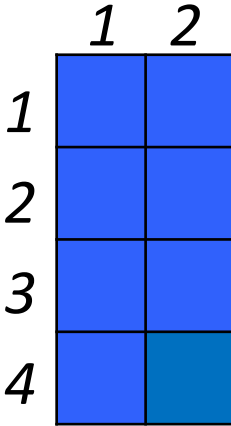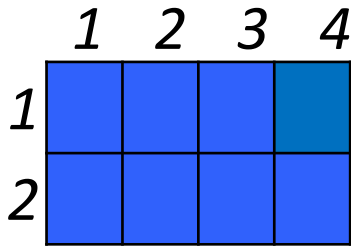Substrate

# Shareability Opportunity

- Two adjacent crossbars in a 3D stack can share a set of DACs if the mapped matrix segments are in the same column i.e., to be multiplied with the same input vector
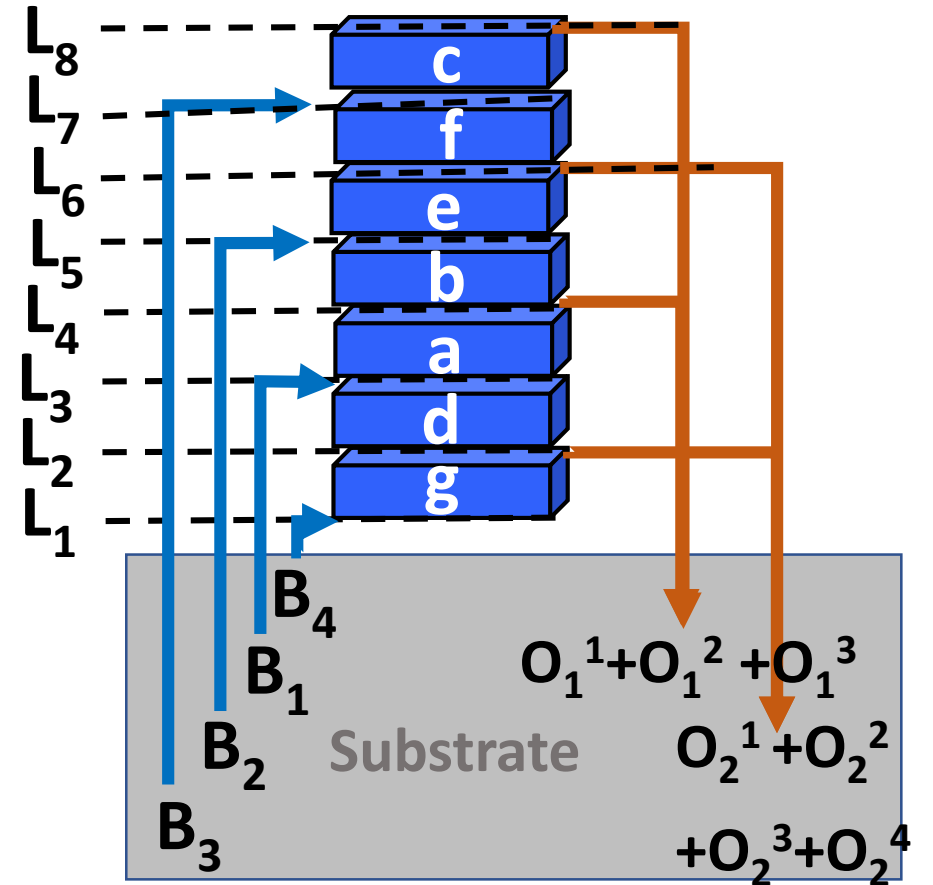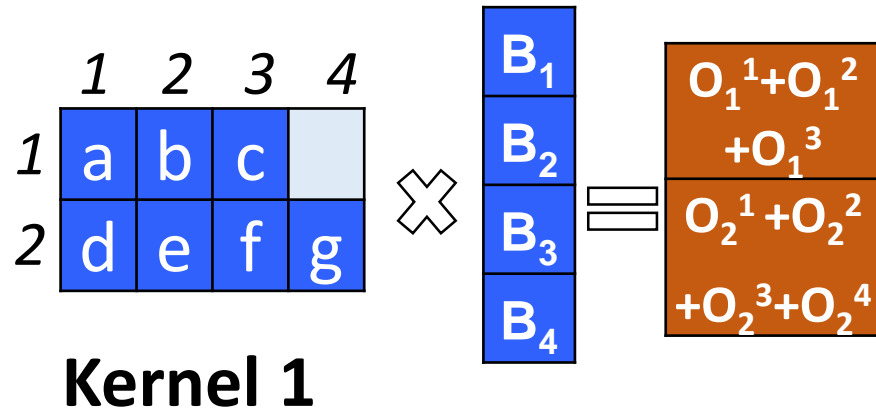
- The dot product output if two adjacent crossbars in a 3D stack can combined if the mapped matrix segments are in the same row i.e., contributing to the same output vector
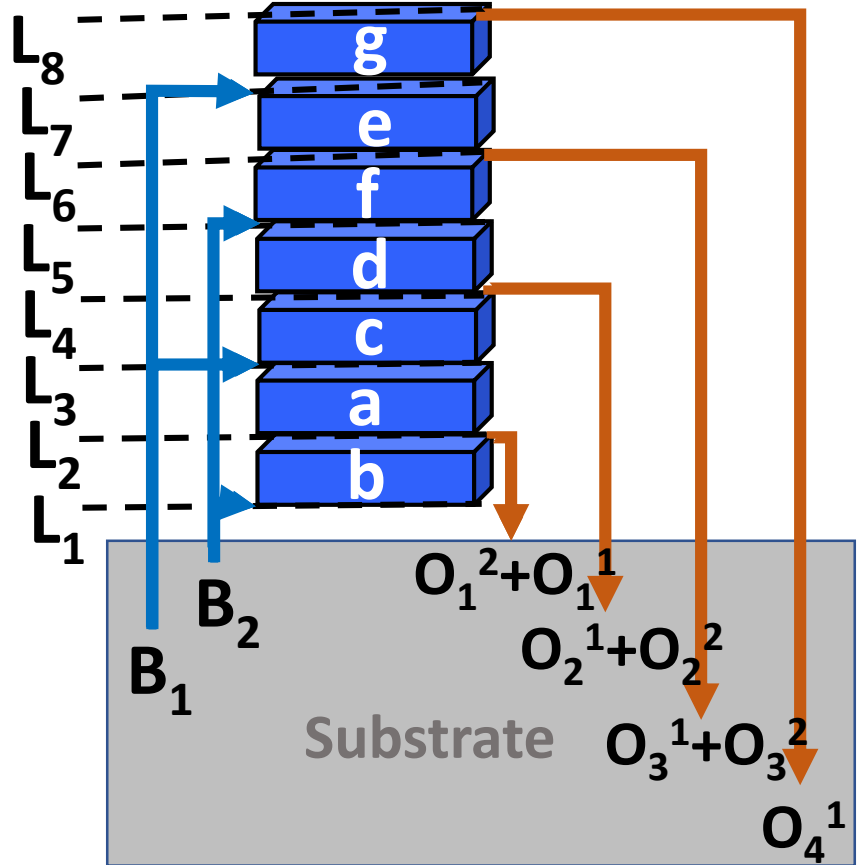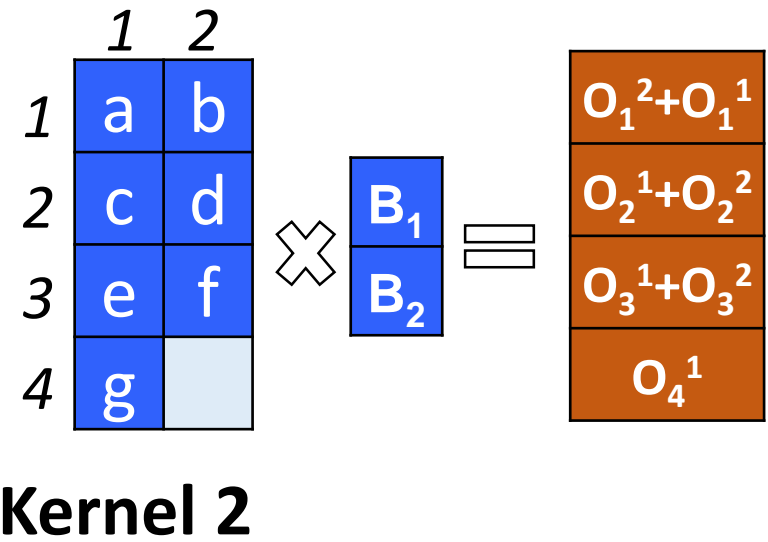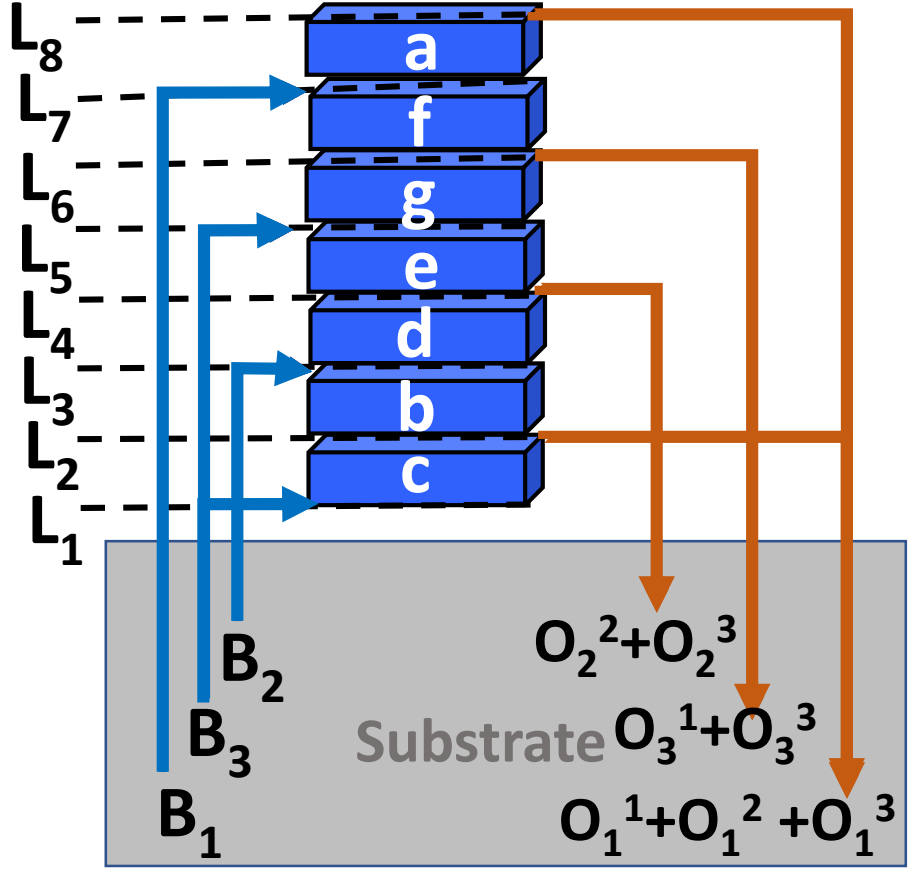
# Problem 2: Discovering In-memory 3D DPE Kernel

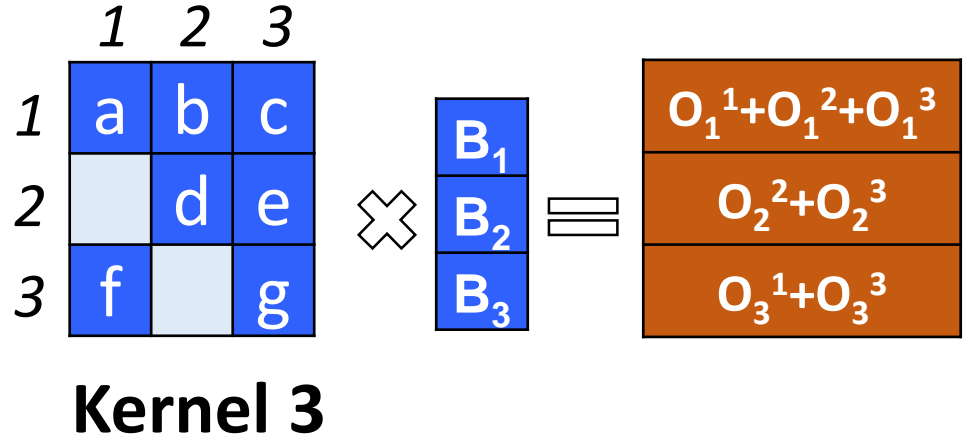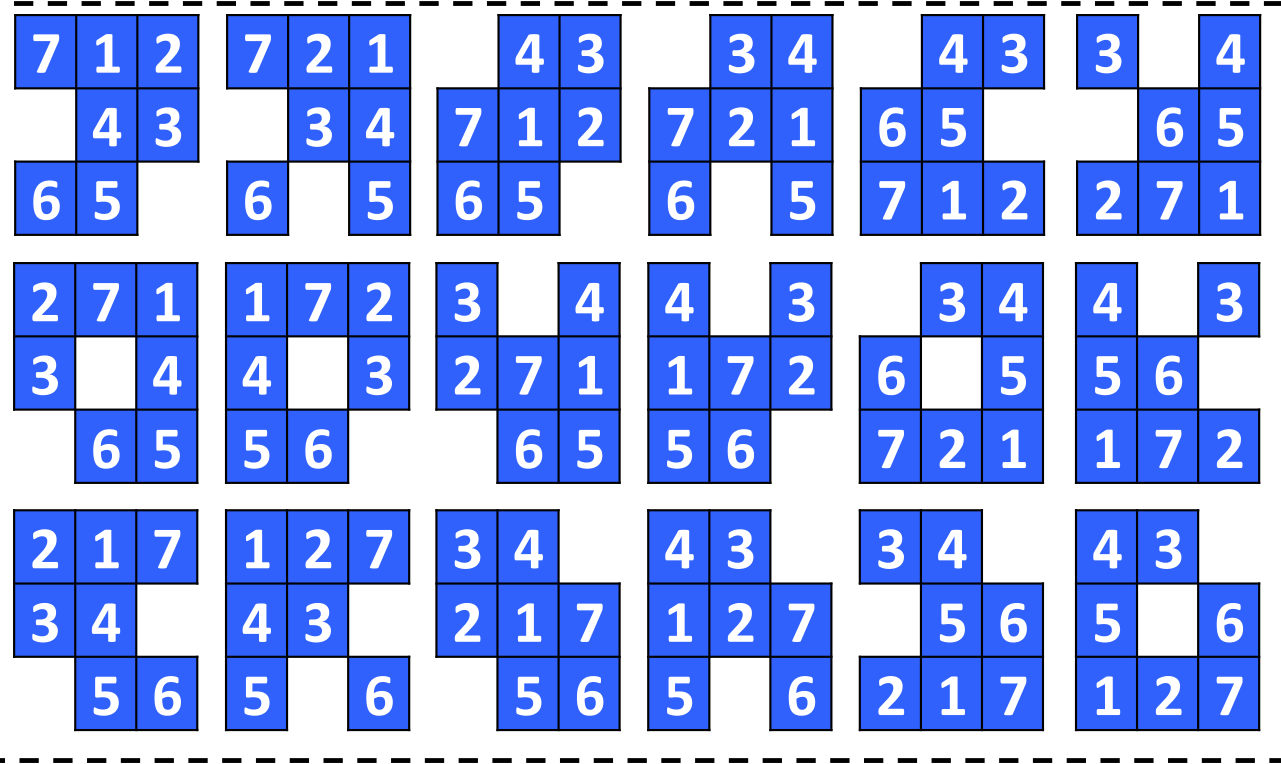# Discovering In-memory 3D DPE Kernel
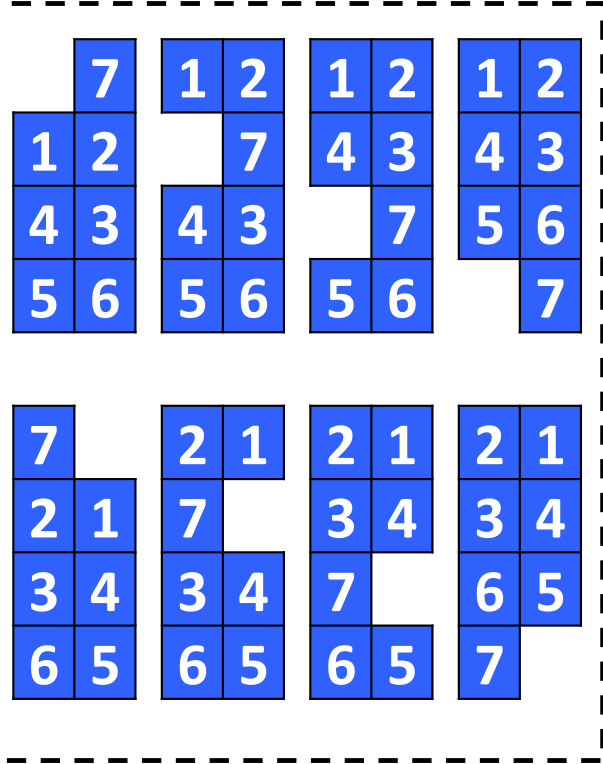


Kernel 1

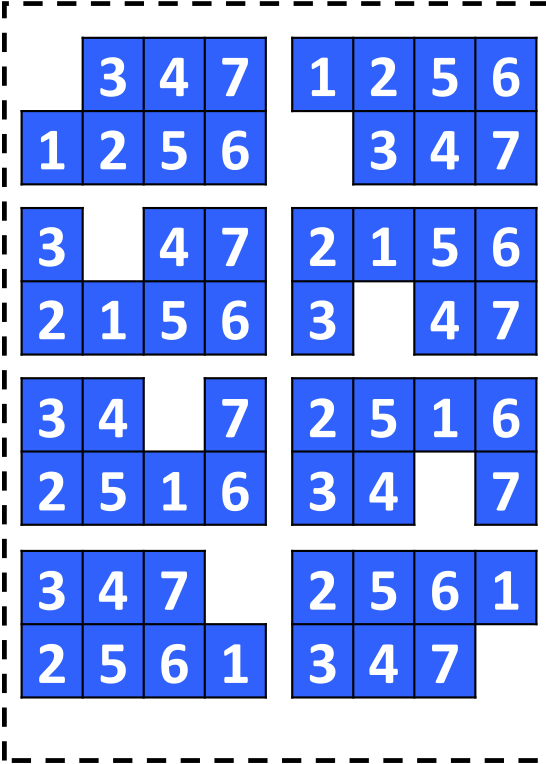# Discovering In-memory 3D DPE Kernel

# Discovering In-memory 3D DPE Kernel

# Library of In-memory 3D DPE Kernel
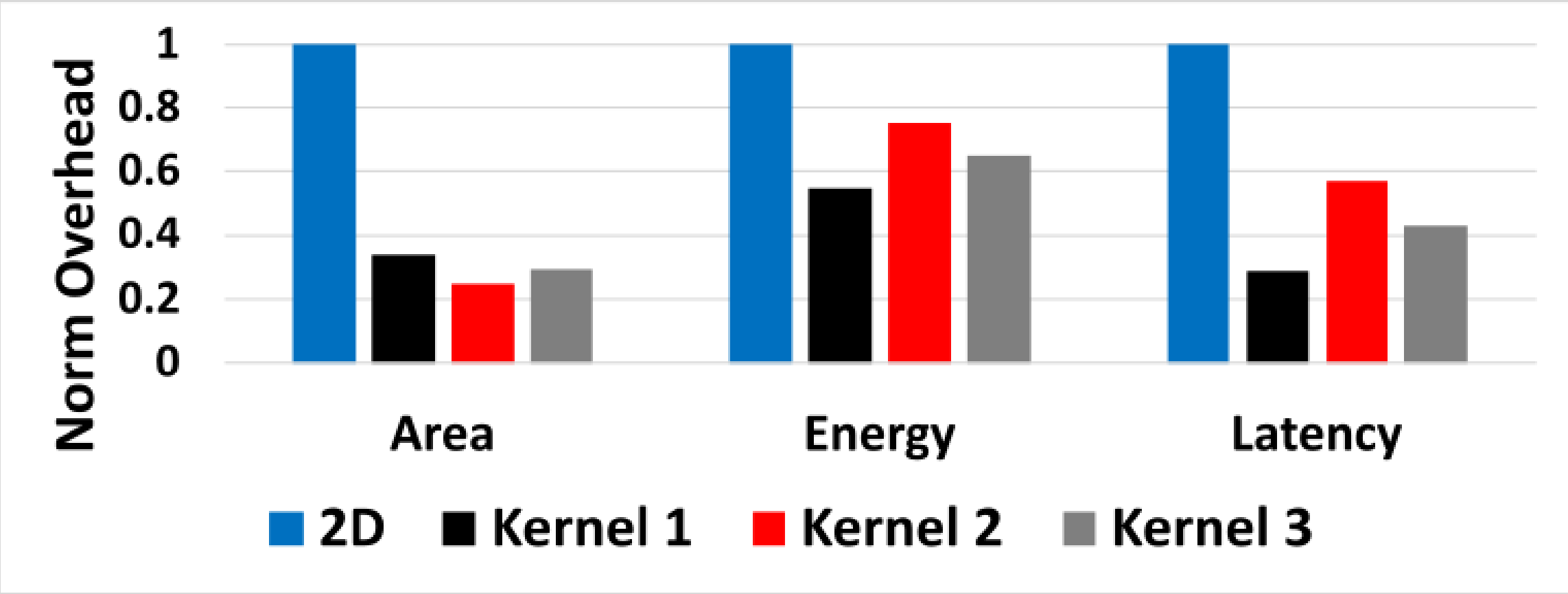


**Kernel 1**     **Kernel 2**                    **Kernel 3**
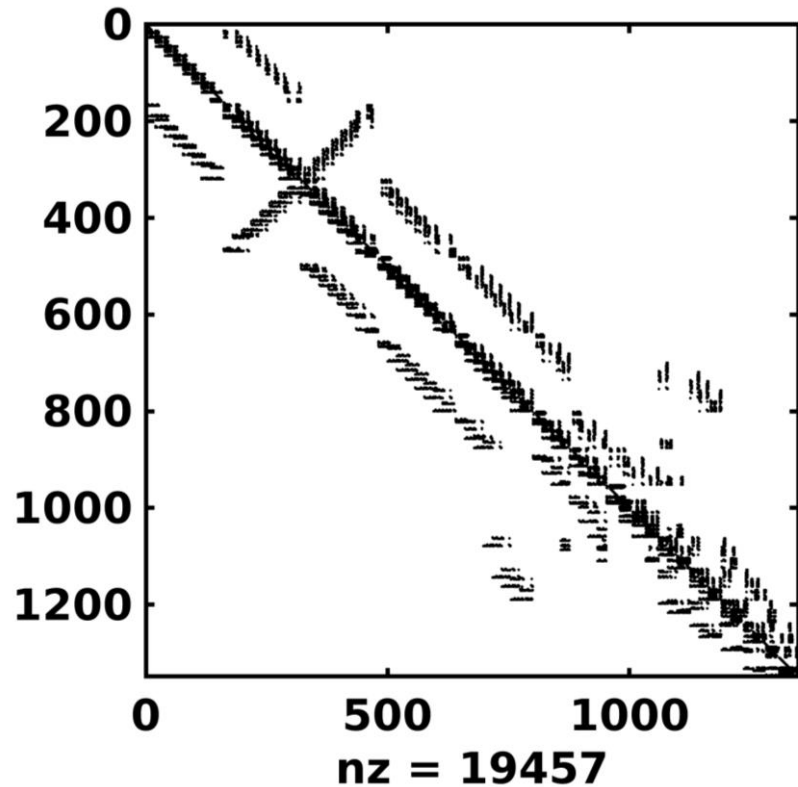
# Application level



Problem:
- Matrix dimension might exceed crossbar dimension

Solution:
- Partitioning Matrix into segments

# Greedy Kernel Cover of Sparse System



Target:
- Covering non-zero blocks using the kernel library

# Evaluation

**Area-power cost of architectural components**

| Component | Parameter | Specs | Area | Power |
|---|---|---|---|---|
| Crossbar | Size | $128 \times 128$ | $25\ \mu m^2$ | 0.3 mW |
| DAC | Resolution | 1 bit | $0.17\ \mu m^2$ | 0.004 mW |
| ADC | Resolution | 8 bits | $0.0012\ mm^2$ | 2 mW |
| Sample+Hold | # Unit | 1 | $0.04\ \mu m^2$ | 10 nW |
| Shift+Add | # Unit | 1 | $60\ \mu m^2$ | 0.05 mW |
| **2D Crossbar** | # Crossbar | 1 | **$51.25\ \mu m^2$** | **0.80 mW** |
| | # DACs | 128 | | |
| | # Sample+Hold | 128 | | |
| **3D Crossbar** | # Crossbar | 7 | **$122.28\ \mu m^2$** | **4.10 mW** |
| | # DACs | $4 \times 128$ | | |
| | # Sample+Hold | $2 \times 128$ | | |
| eDRAM Buffer | Size | 128 KB | $0.17\ mm^2$ | 41.4 mW |
| IR | Size | 4 KB | $4200\ \mu m^2$ | 2.48 mW |
| OR | Size | 512 B | $1500\ \mu m^2$ | 0.46 mW |
| Bus | Bandwidth | 128-bits | $15.7\ mm^2$ | 13 mW |

## Table 4: Overview of benchmarks from the SuiteSparse Matrix Collection.

| Applications | Systems | Matrix Dimensions | #Non-zeros |
| --- | --- | --- | --- |
| bcsstk34 | Structural Problem | $588 \times 588$ | 21418 |
| eris1176 | Power Network Problem | $1176 \times 1176$ | 18552 |
| coater1 | Computational Fluid Dynamics | $1348 \times 1348$ | 19457 |
| cegb2919 | Structural Problem | $2919 \times 2919$ | 321543 |
| mycielskian12 | Undirected Graph | $3071 \times 3071$ | 407200 |
| raefsky1 | Computational Fluid Dynamics | $3242 \times 3242$ | 293409 |
| crystk01 | Materials Problem | $4875 \times 4875$ | 315891 |
| fxm3_6 | Optimization Problem | $5026 \times 5026$ | 94026 |
| Na5 | Theoretical/Quantum Chemistry | $5832 \times 5832$ | 305630 |
| EX5 | Combinatorial Problem | $6545 \times 6545$ | 295680 |
| fp | Electromagnetics Problem | $7548 \times 7548$ | 834222 |
| ex40 | Computational Fluid Dynamics | $7740 \times 7740$ | 456188 |
| benzene | Theoretical/Quantum Chemistry | $8219 \times 8219$ | 242669 |
| bcsstk33 | Structural Problem | $8738 \times 8738$ | 591904 |
| graham1 | Computational Fluid Dynamics | $9035 \times 9035$ | 335472 |
| net25 | Optimization Problem | $9520 \times 9520$ | 401200 |
| bundle1 | Computer Graphics/Vision | $10581 \times 10581$ | 770811 |
| Si10H16 | Theoretical/Quantum Chemistry | $17077 \times 17077$ | 875923 |
| Goodwin_040 | Computational Fluid Dynamics | $17922 \times 17922$ | 561677 |
| pkustk06 | Structural Problem | $43164 \times 43164$ | 2571768 |

**Improvements:**
- **Area by 2.02X**
- **Latency by 2.45X**
- **Energy by 2.37X**

# Summary

# Thank You

Contact us:

rashed09@knights.ucf.edu

sumit.jha@utsa.edu

rickard.ewetz@ucf.edu