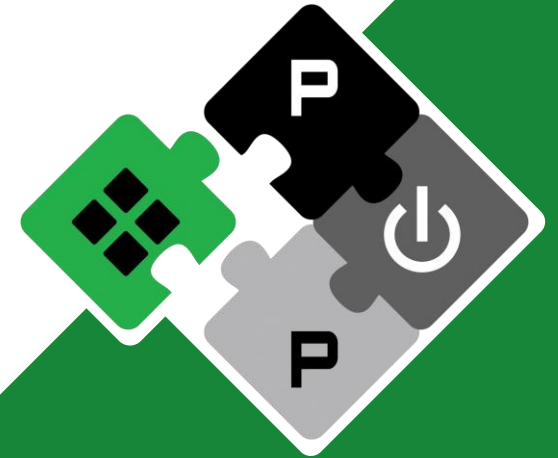# SoftHier Progress Update

**Chi Zhang**    chizhang@iis.ee.ethz.ch

**PULP Platform**
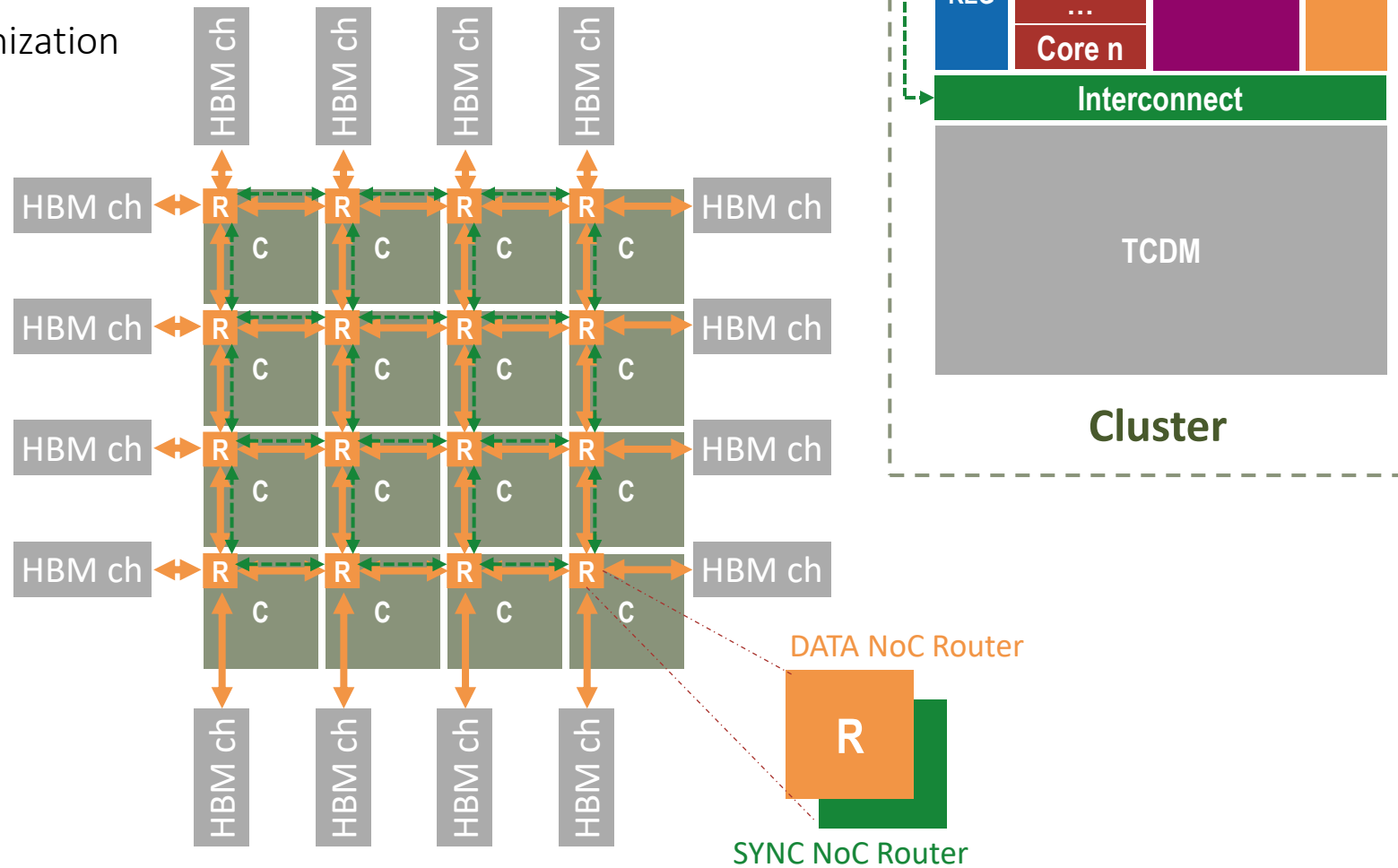Open Source Hardware, the way it should be!

@pulp_platform

pulp-platform.org

youtube.com/pulp_platform

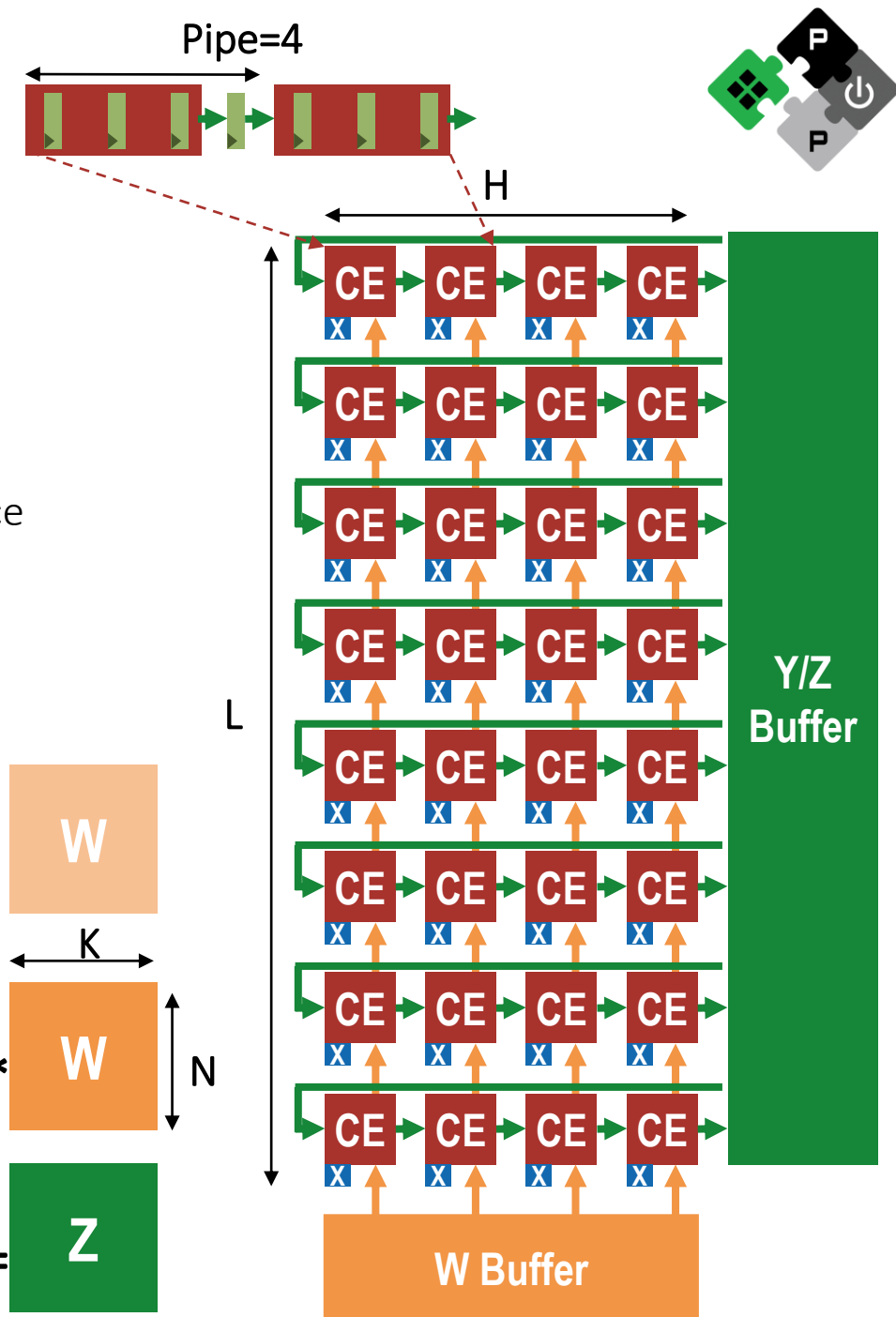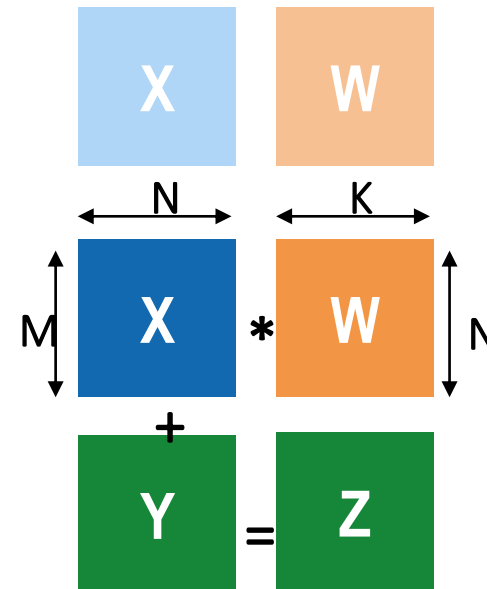# SoftHier: Parameterizable NoC-Based Scalable System

- Features

  - Two separate NoC bus

    - DATA NoC: wide link, transfer bulk data

    - SYNC NoC: narrow link, cluster synchronization

  - Infinite Instruction memory

    - Ignore I$ fetch overhead

- Fully Parameterizable

  - Configure file and push button
    - #Cores, RedMule config
    - L1 (TCDM) size & BW
    - #Clusters (#row, #col)
    - #HBM channels and placement
    - NoC link BW

- SW stack ready

# Design Space Exploration of RedMule



- Goal: We Want to Know

  - What is **optimal MatMul problem size** (M-N-K)to reach the best efficiency of RedMule **in a Cluster:** the **Optimal Efficiency Point**

    - When Matrix too small -> low RedMule utilization

    - When Matrix too large -> we introduce large and redundant TCDM space

    - We are seeking for **best MACs/SRAM in a CLuster**

  - What is the BW needed for RedMule at **Optimal Efficiency Point**

- Design Space Exploration Constraints

  - RedMule CE array constraint: L = H * Pipe, CE Pipe=4

- Key Metric

  - Efficiency Metric:

    - $\dfrac{Effective\ FLOP/Cycle}{TCDM\ Occupied\ Area} = \dfrac{RedMule\ Utilization * 2 * \#CEs}{Elem\ Size * 2 * (MN + NK + MK)}$
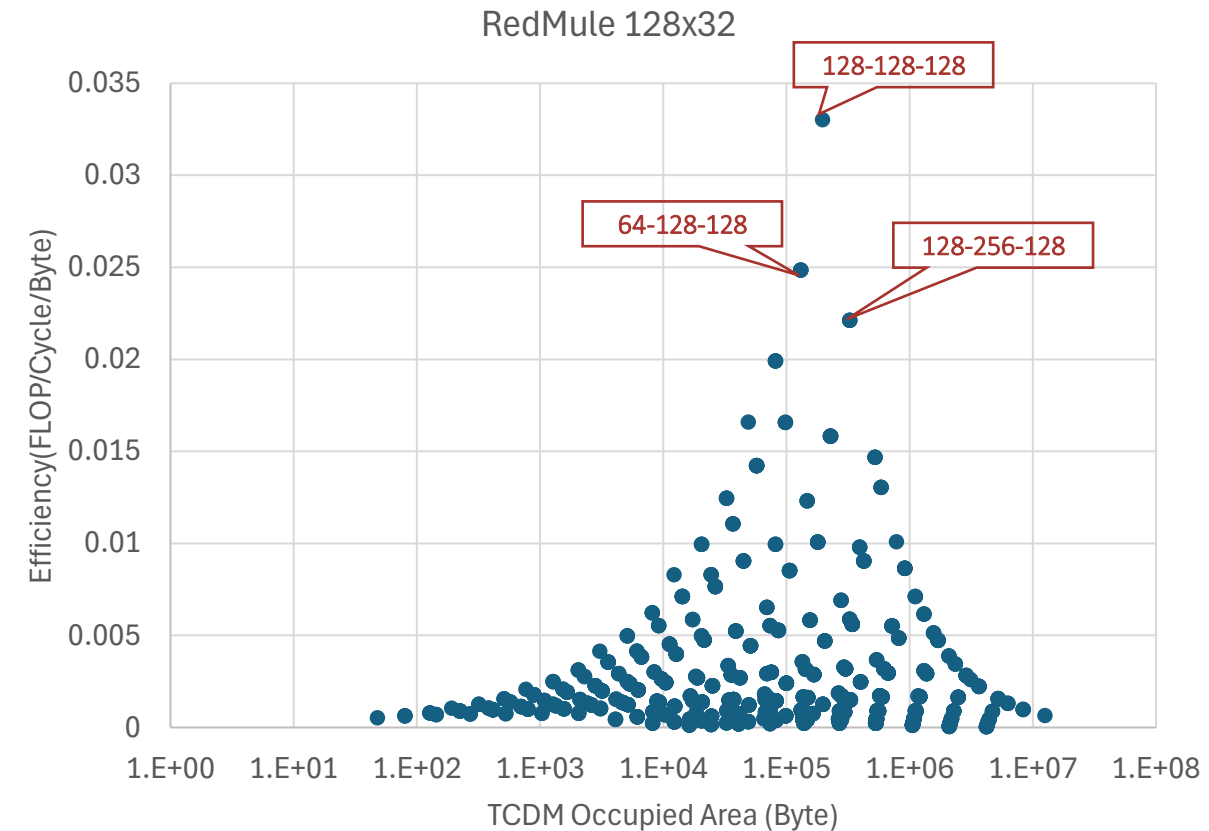
# Find Optimal GEMM Size(M-N-K) for RedMule (@ saturated BW)

- GEMM Dimension

  - $M \quad \in [8, 16, 32, 64, 128, 256, 512]$

  - $N \quad \in [8, 16, 32, 64, 128, 256, 512]$

  - $K \quad \in [8, 16, 32, 64, 128, 256, 512]$

- RedMule Config

  - CE array  = 128x32

  - TCDM BW = 1024 Elem/Cycle

  - Element  = FP16

- Run GEMM on One Cluster with One RedMule
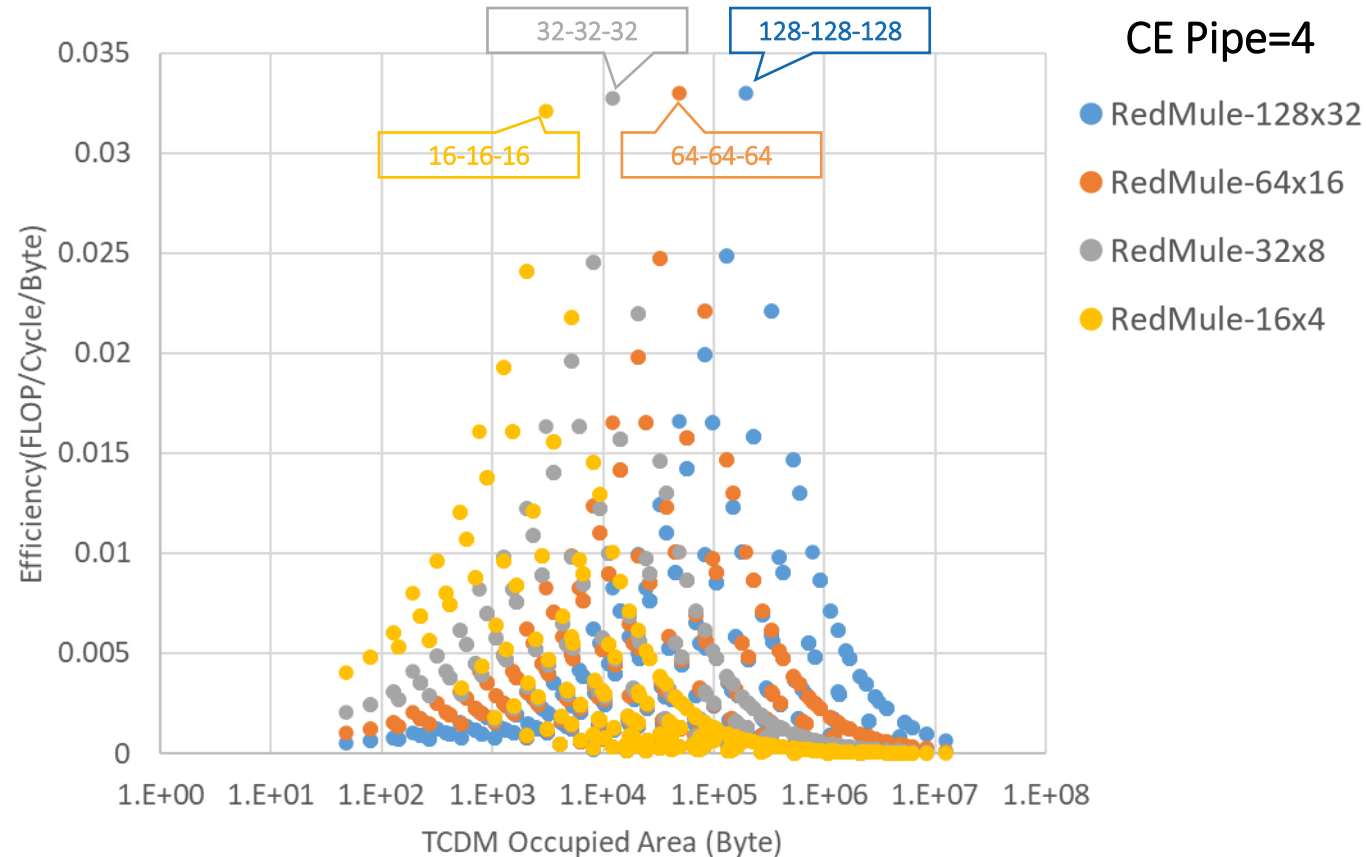
- Collect Data

  - Efficiency Metric:

    - $\dfrac{Effective\ FLOP/Cycle}{TCDM\ Occupied\ Area} = \dfrac{RedMule\ Utilization * 2 * \#CEs}{Elem\ Size * 2 * (MN+NK+MK)}$



RedMule 128x32

# Find Optimal GEMM Size (M-N-K) for RedMule (@ saturated BW)

- At RedMule Constraints: CE array constraint: L = H * Pipe



**CE Pipe=4**

- RedMule-128x32
- RedMule-64x16
- RedMule-32x8
- RedMule-16x4
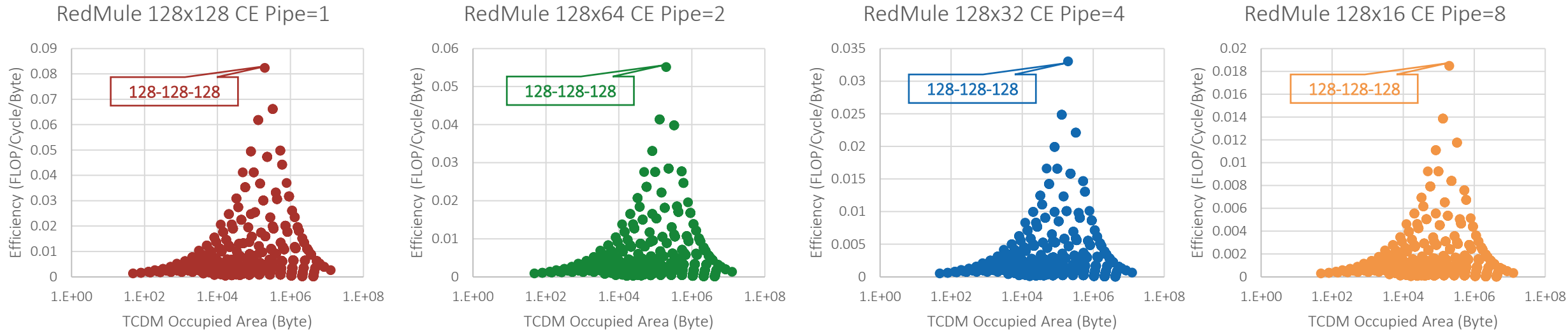
**To use RedMule at best efficiency (Optimal Efficiency Point)**

$$M = N = K = \sqrt{CE\ array\ *\ CE\ Pipeline} = L$$

# Find Optimal GEMM Size (M-N-K) for RedMule (@ saturated BW)

- At RedMule Constraints: CE array constraint: L = H * Pipe



## RedMule 128x128 CE Pipe=1
128-128-128

## RedMule 128x64 CE Pipe=2
128-128-128

## RedMule 128x32 CE Pipe=4
128-128-128

## RedMule 128x16 CE Pipe=8
128-128-128
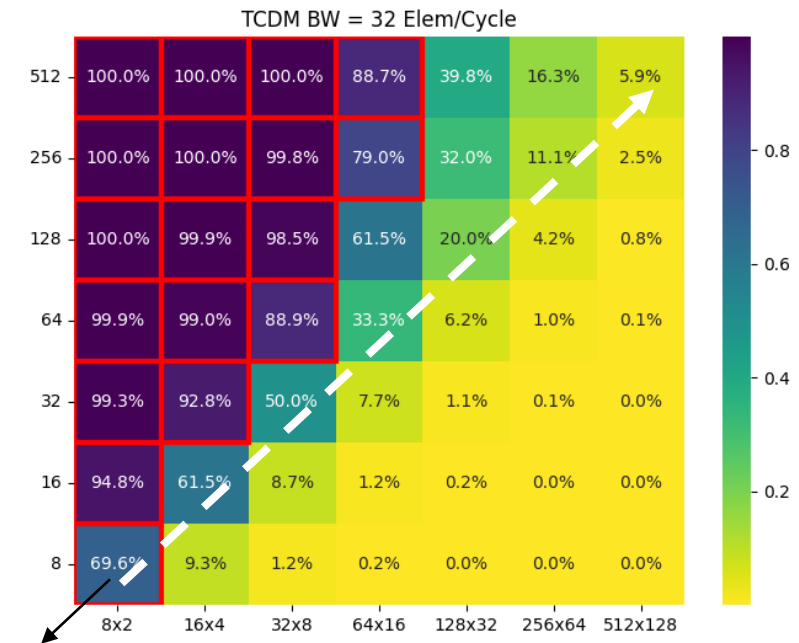
**To use RedMule at best efficiency (Optimal Efficiency Point)**

$$M = N = K = \sqrt{CE\ array\ *\ CE\ Pipeline} = L$$
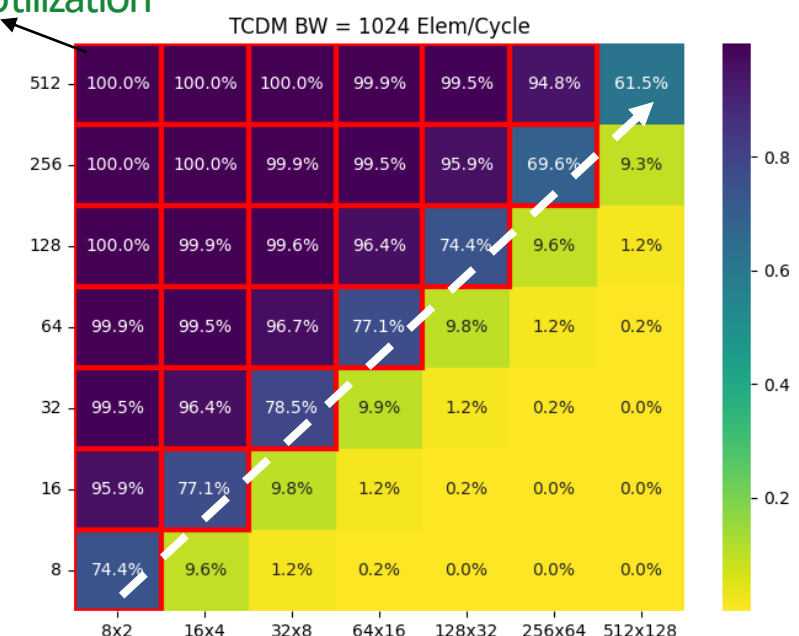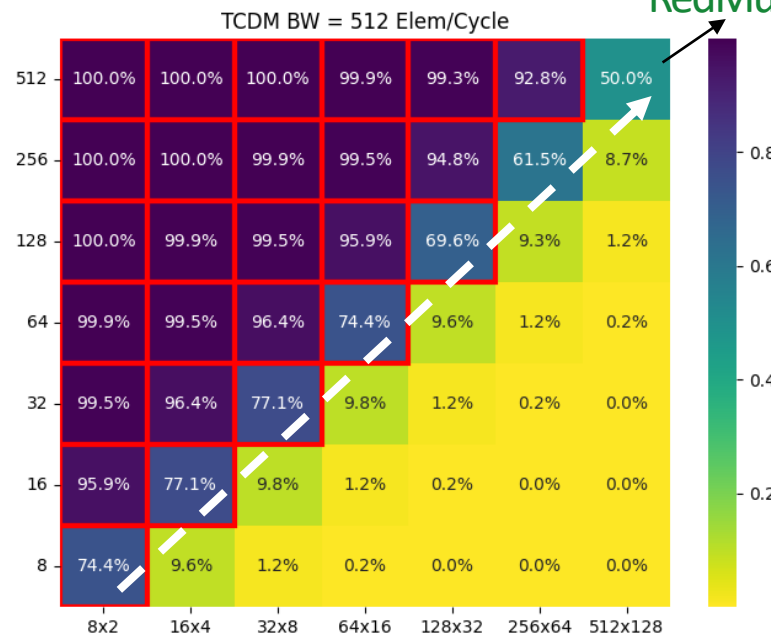
# TCDM Bandwidth Requirement

- RedMule Utilization Heatmap
  - GEMM Size vs RedMule CE array
  - Vary TCDM BW
  - We're satisfied at > 69% uti

**To feed RedMule efficiently TCDM Bandwidth (elem per cycle):**
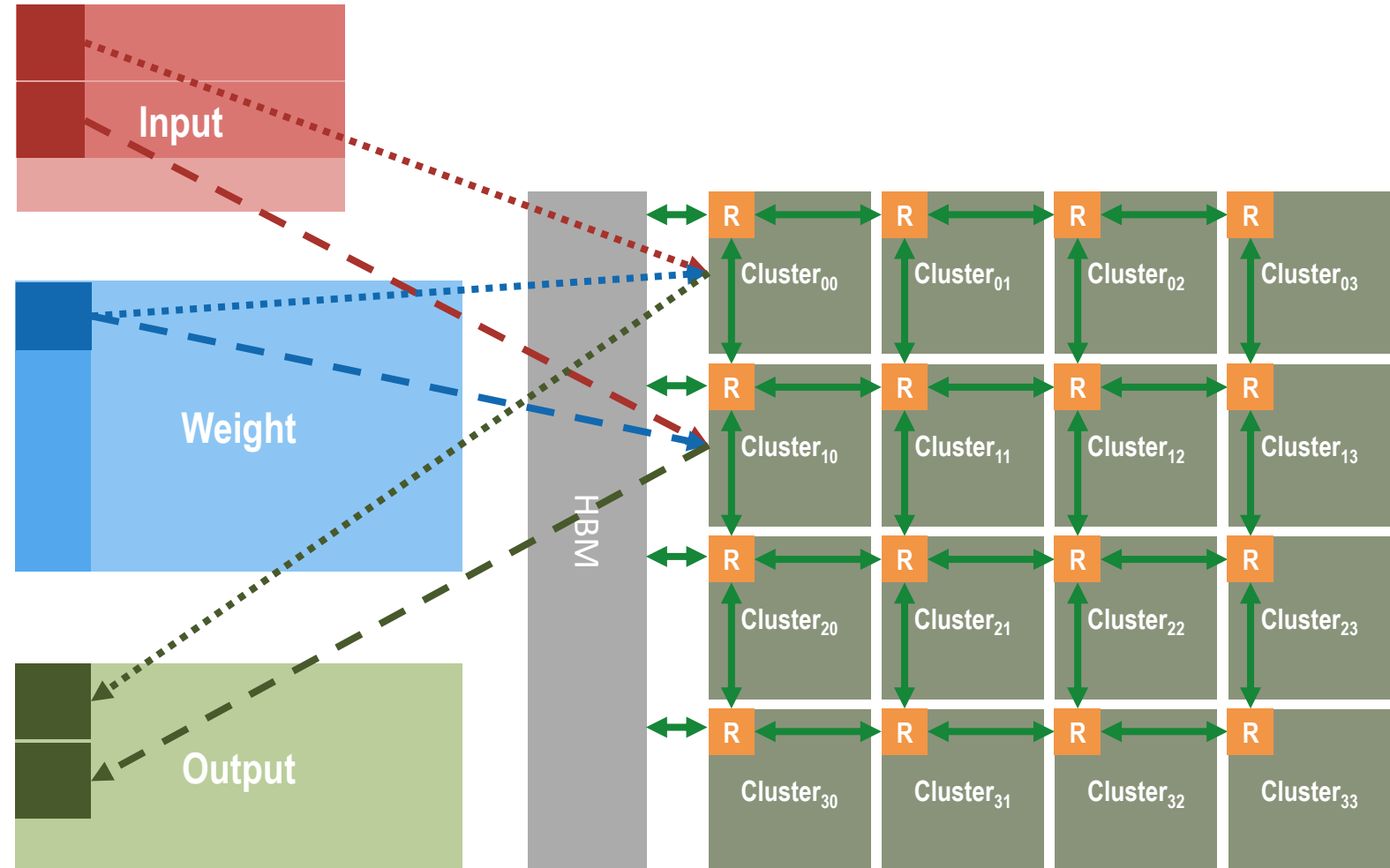
$$BW \geq 4\sqrt{CE\ array\ *\ CE\ Pipeline} = 4L$$

# For Large GEMM: Enable Cluster-to-Cluster Comm or Not?

- No Cluster-to-Cluster Comm.
  - Each clusters take care of different output tiles (384x384).
  - iDMA transfers input matrix tile + weight matrix tile from HBM
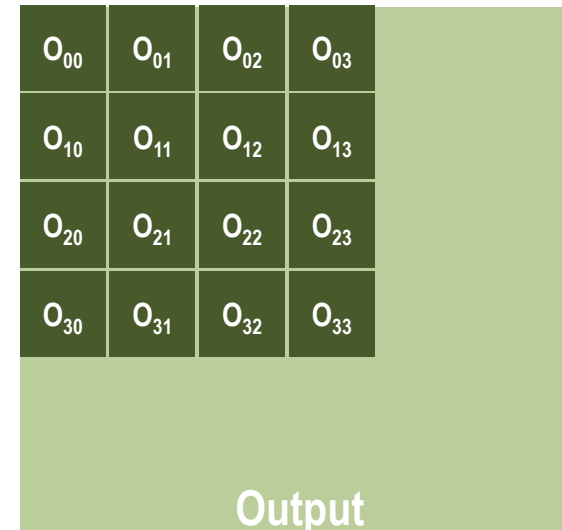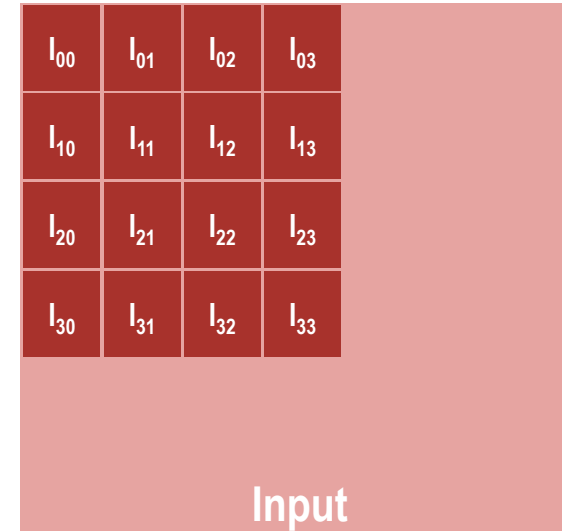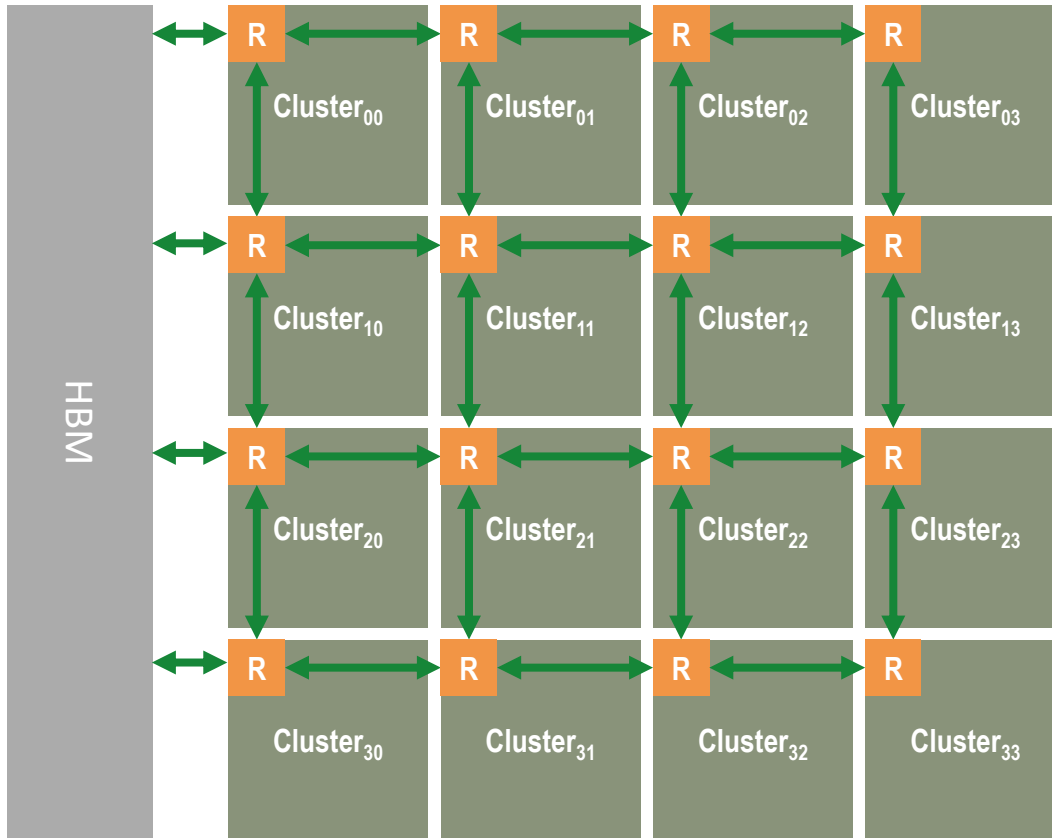  - No inter-cluster tile reuses
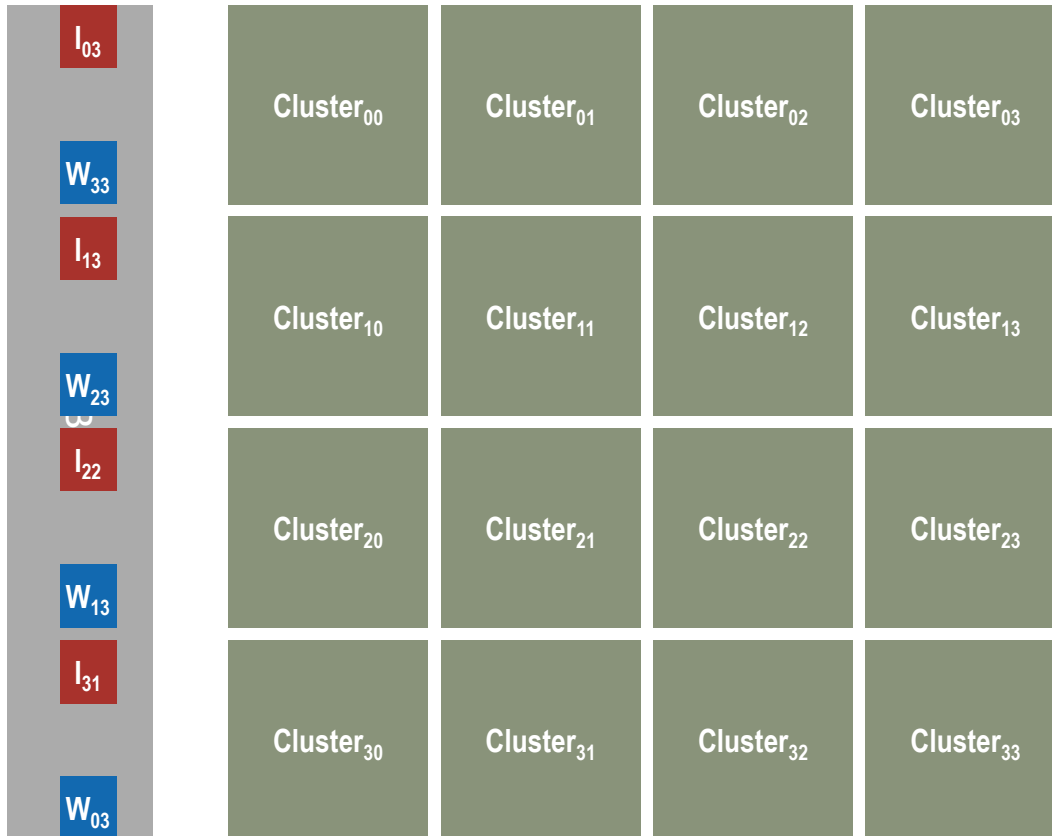
# For Large GEMM: Enable Cluster-to-Cluster Comm or Not?

- Leverage Cluster-to-Cluster Comm.
  - Reuse tiles, reduce HBM accesses, save BW limitation
  - Need "smart" tile mapping and scheduling

# Example Solution for GEMM: Step1 Distribute Tiles to Clusters

- Load All (16 Input tiles, 16 weight tiles) from HBM
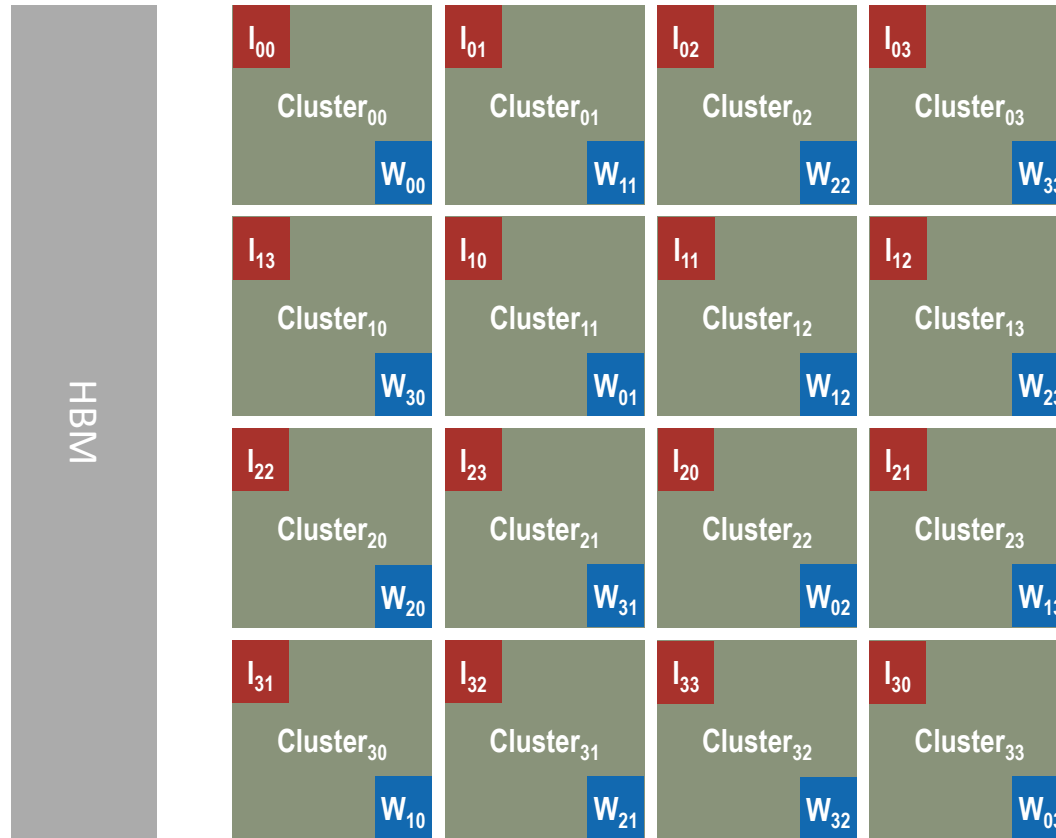- Cluster contain different tiles from each other

$$O_{00} = \boxed{I_{00} \cdot W_{00}} + I_{01} \cdot W_{10} + I_{02} \cdot W_{20} + I_{03} \cdot W_{30}$$

$$O_{01} = I_{00} \cdot W_{01} + \boxed{I_{01} \cdot W_{11}} + I_{02} \cdot W_{21} + I_{03} \cdot W_{31}$$

$$O_{02} = I_{00} \cdot W_{02} + I_{01} \cdot W_{12} + \boxed{I_{02} \cdot W_{22}} + I_{03} \cdot W_{32}$$

$$O_{03} = I_{00} \cdot W_{03} + I_{01} \cdot W_{13} + I_{02} \cdot W_{23} + \boxed{I_{03} \cdot W_{33}}$$

$$O_{10} = I_{10} \cdot W_{00} + I_{11} \cdot W_{10} + I_{12} \cdot W_{20} + \boxed{I_{13} \cdot W_{30}}$$

$$O_{11} = \boxed{I_{10} \cdot W_{01}} + I_{11} \cdot W_{11} + I_{12} \cdot W_{21} + I_{13} \cdot W_{31}$$

$$O_{12} = I_{10} \cdot W_{02} + \boxed{I_{11} \cdot W_{12}} + I_{12} \cdot W_{22} + I_{13} \cdot W_{32}$$

$$O_{13} = I_{10} \cdot W_{03} + I_{11} \cdot W_{13} + \boxed{I_{12} \cdot W_{23}} + I_{13} \cdot W_{33}$$

$$O_{20} = I_{20} \cdot W_{00} + I_{21} \cdot W_{10} + \boxed{I_{22} \cdot W_{20}} + I_{23} \cdot W_{30}$$

$$O_{21} = I_{20} \cdot W_{01} + I_{21} \cdot W_{11} + I_{22} \cdot W_{21} + \boxed{I_{23} \cdot W_{31}}$$

$$O_{22} = \boxed{I_{20} \cdot W_{02}} + I_{21} \cdot W_{12} + I_{22} \cdot W_{22} + I_{23} \cdot W_{32}$$

$$O_{23} = I_{20} \cdot W_{03} + \boxed{I_{21} \cdot W_{13}} + I_{22} \cdot W_{23} + I_{23} \cdot W_{33}$$

$$O_{30} = I_{30} \cdot W_{00} + \boxed{I_{31} \cdot W_{10}} + I_{32} \cdot W_{20} + I_{33} \cdot W_{30}$$

$$O_{31} = I_{30} \cdot W_{01} + I_{31} \cdot W_{11} + \boxed{I_{32} \cdot W_{21}} + I_{33} \cdot W_{31}$$

$$O_{32} = I_{30} \cdot W_{02} + I_{31} \cdot W_{12} + I_{32} \cdot W_{22} + \boxed{I_{33} \cdot W_{32}}$$

$$O_{33} = \boxed{I_{30} \cdot W_{03}} + I_{31} \cdot W_{13} + I_{32} \cdot W_{23} + I_{33} \cdot W_{33}$$

Column (HBM): $I_{03}$, $W_{33}$, $I_{13}$, $W_{23}$, $I_{22}$, $W_{13}$, $I_{31}$, $W_{03}$

| Cluster$_{00}$ | Cluster$_{01}$ | Cluster$_{02}$ | Cluster$_{03}$ |
|---|---|---|---|
| Cluster$_{10}$ | Cluster$_{11}$ | Cluster$_{12}$ | Cluster$_{13}$ |
| Cluster$_{20}$ | Cluster$_{21}$ | Cluster$_{22}$ | Cluster$_{23}$ |
| Cluster$_{30}$ | Cluster$_{31}$ | Cluster$_{32}$ | Cluster$_{33}$ |

# Example Solution for GEMM: Step1 Distribute Tiles to Clusters
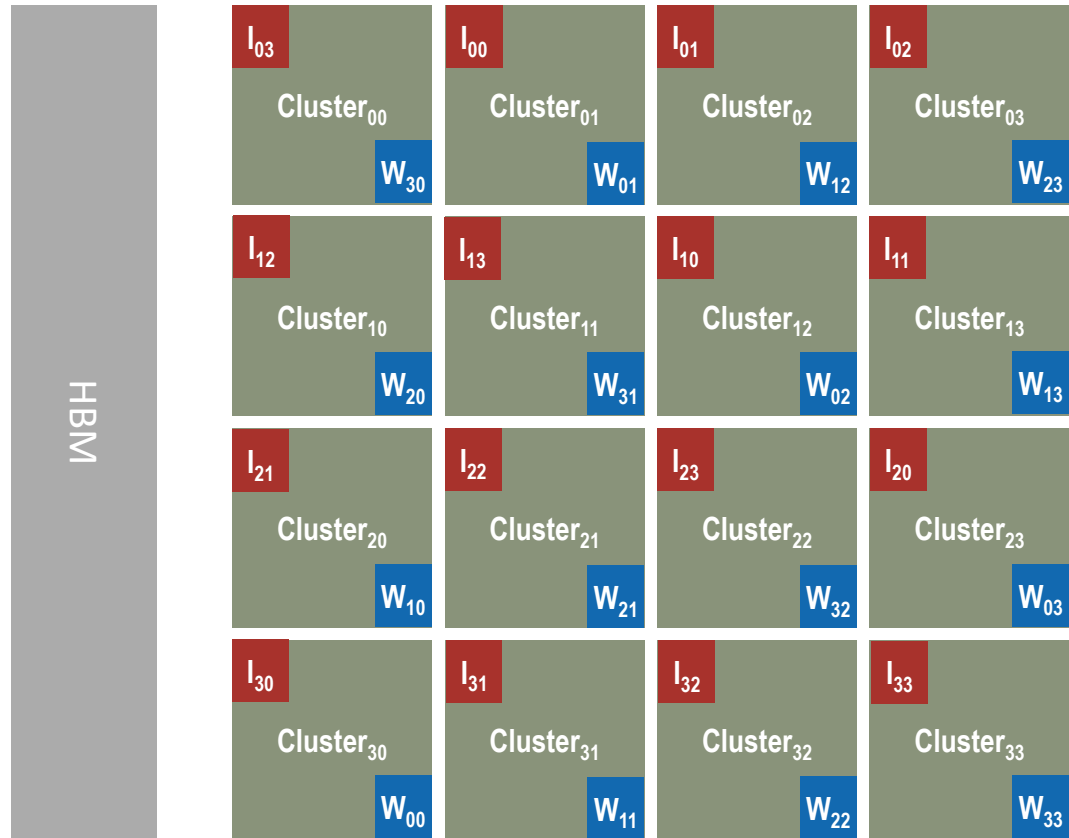
- Load All (16 Input tiles, 16 weight tiles) from HBM
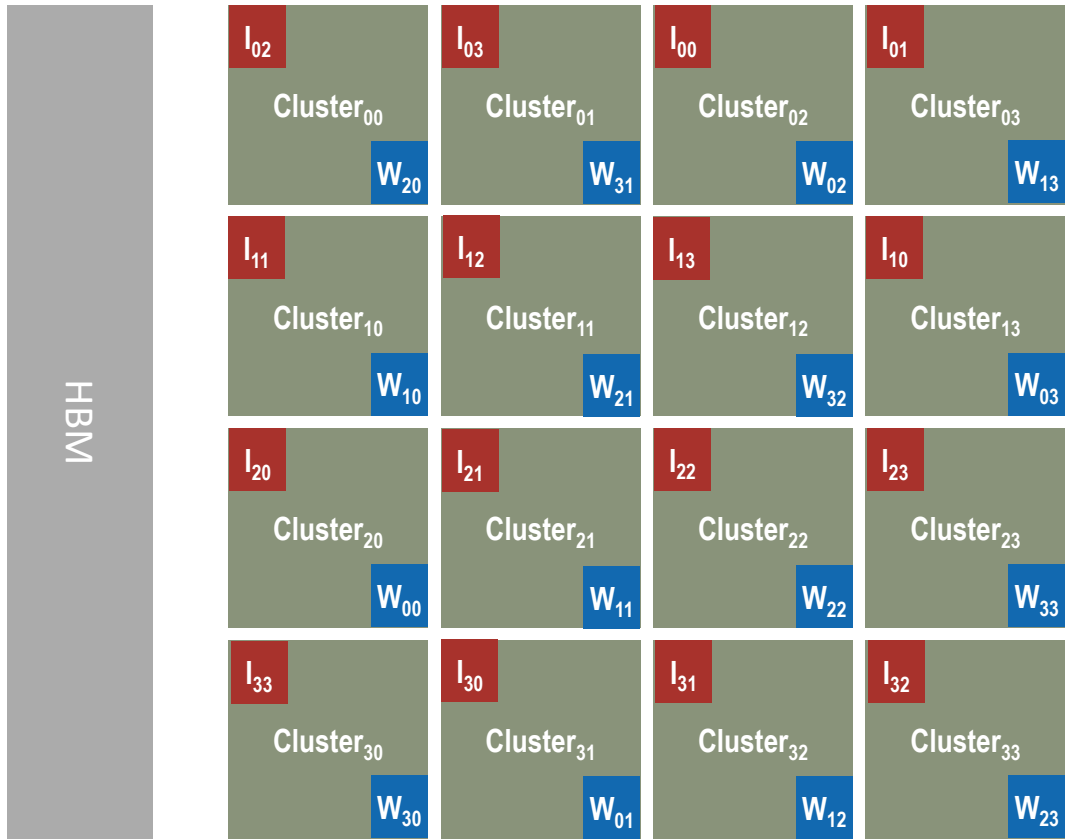- Cluster contain different tiles from each other



$$O_{00} = \boxed{I_{00} \cdot W_{00}} + I_{01} \cdot W_{10} + I_{02} \cdot W_{20} + I_{03} \cdot W_{30}$$

$$O_{01} = I_{00} \cdot W_{01} + \boxed{I_{01} \cdot W_{11}} + I_{02} \cdot W_{21} + I_{03} \cdot W_{31}$$

$$O_{02} = I_{00} \cdot W_{02} + I_{01} \cdot W_{12} + \boxed{I_{02} \cdot W_{22}} + I_{03} \cdot W_{32}$$

$$O_{03} = I_{00} \cdot W_{03} + I_{01} \cdot W_{13} + I_{02} \cdot W_{23} + \boxed{I_{03} \cdot W_{33}}$$

$$O_{10} = I_{10} \cdot W_{00} + I_{11} \cdot W_{10} + I_{12} \cdot W_{20} + \boxed{I_{13} \cdot W_{30}}$$

$$O_{11} = \boxed{I_{10} \cdot W_{01}} + I_{11} \cdot W_{11} + I_{12} \cdot W_{21} + I_{13} \cdot W_{31}$$

$$O_{12} = I_{10} \cdot W_{02} + \boxed{I_{11} \cdot W_{12}} + I_{12} \cdot W_{22} + I_{13} \cdot W_{32}$$

$$O_{13} = I_{10} \cdot W_{03} + I_{11} \cdot W_{13} + \boxed{I_{12} \cdot W_{23}} + I_{13} \cdot W_{33}$$

$$O_{20} = I_{20} \cdot W_{00} + I_{21} \cdot W_{10} + \boxed{I_{22} \cdot W_{20}} + I_{23} \cdot W_{30}$$

$$O_{21} = I_{20} \cdot W_{01} + I_{21} \cdot W_{11} + I_{22} \cdot W_{21} + \boxed{I_{23} \cdot W_{31}}$$

$$O_{22} = \boxed{I_{20} \cdot W_{02}} + I_{21} \cdot W_{12} + I_{22} \cdot W_{22} + I_{23} \cdot W_{32}$$

$$O_{23} = I_{20} \cdot W_{03} + \boxed{I_{21} \cdot W_{13}} + I_{22} \cdot W_{23} + I_{23} \cdot W_{33}$$

$$O_{30} = I_{30} \cdot W_{00} + \boxed{I_{31} \cdot W_{10}} + I_{32} \cdot W_{20} + I_{33} \cdot W_{30}$$

$$O_{31} = I_{30} \cdot W_{01} + I_{31} \cdot W_{11} + \boxed{I_{32} \cdot W_{21}} + I_{33} \cdot W_{31}$$

$$O_{32} = I_{30} \cdot W_{02} + I_{31} \cdot W_{12} + I_{32} \cdot W_{22} + \boxed{I_{33} \cdot W_{32}}$$

$$O_{33} = \boxed{I_{30} \cdot W_{03}} + I_{31} \cdot W_{13} + I_{32} \cdot W_{23} + I_{33} \cdot W_{33}$$

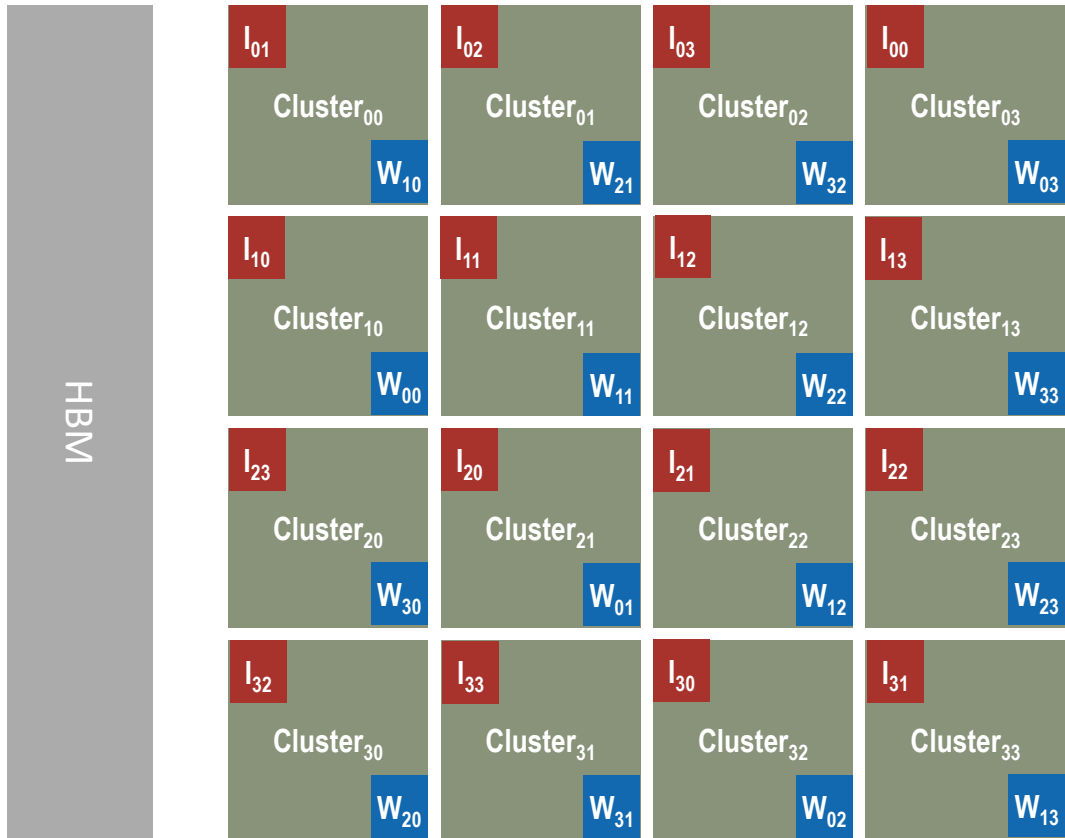# Example Solution for GEMM: Step2 Inter-Cluster Tile Exchanging

# Experiment Setup: Enable Cluster-to-Cluster Comm or Not?

- RedMule
  - Each Cluster has One RedMule
  - RedMule CE array 128x32 (8TFLOPs @1GHz)
  - TCDM BW = 1024 GB/s

- HBM
  - Place HBM on left side
  - Each HBM CTRL mange 4 HBM2E channel
    - 256Byte address interleaving
    - Each HBM CTRL provide Max.BW = 205 GB/s

- NoC
  - Mesh 4x4
  - NoC link width = 2048 bits
    - Link BW = 256GB/s

# Experiment Result: Enable Cluster-to-Cluster Comm or Not?

- Benchmark GEMM
  - 16384 x16384 x16384
  - GEMM Elem Size = FP16
  - RedMule 128x32

- Varying Cluster TCDM Size
  - Tiling strategy: Max possible tile fit in TCDM
    - Tile dimension $M = N = K = \sqrt{\frac{TCDM\ Area}{5*Elem\_Size}}$

- Results
  - No inter-cluster comm
    - HBM BW limited
    - Need 4x more TCDM size to saturate compute power
  - Enable inter-cluster comm
    - Reduce HBM traffic, better area/power efficiency



Legend:
- HBM BW Uti--Enable Cluster Comm
- HBM BW Uti--No Cluster Comm
- Compute Uti--Enable Cluster Comm
- Compute Uti--No Cluster Comm

| TCDM Size of Cluster: | 10KB | 40KB | 160KB | 640KB | 2.56MB | 10.24MB |
|---|---|---|---|---|---|---|
| Max Tile Dimension: | 32 | 64 | 128 | 256 | 512 | 1024 |

# Enable Cluster-to-Cluster Comm is Scalable

- Benchmark GEMM
  - 16384 x16384 x16384
  - GEMM Elem Size = FP16
  - RedMule 128x32

- Scale-out SoftHier System
  - 4x4    Clusters + 16 HBM2E channels
  - 8x8    Clusters + 32 HBM2E channels
  - 16x16 Clusters + 64 HBM2E channels
    - 2048 TFLOPS @FP16
    - 3.2 TB/s HBM BW



Legend:
- 4x4  Clusters----Enable Cluster Comm
- 4x4  Clusters----No Cluster Comm
- 8x8  Clusters----Enable Cluster Comm
- 8x8  Clusters----No Cluster Comm
- 16x16 Clusters--Enable Cluster Comm
- 16x16 Clusters--No Cluster Comm

| TCDM Size in Cluster: | 10KB | 40KB | 160KB | 640KB | 2.56MB | 10.24MB |
|---|---|---|---|---|---|---|
| Max Tile Dimension: | 32 | 64 | 128 | 256 | 512 | 1024 |

Y-axis: Compute Utilization (0% to 100%)

# Further Discussion ...

- New Questions Comes Out
  - @ RedMule optimal efficiency point
    - RedMule in cluster reaches 70% comp uti
    - But end-to-end the system shows 44% uti
    - Why? How can we optimize this?
  - Tile mapping and inter-cluster scheduling scheme
    - Is there any more schemes for large GEMM
    - How can we also leverage inter-cluster comm for MHA?
  - Inter-cluster comm vs multi-broadcasting
    - Which one is better?



Legend:
- 4x4  Clusters----Enable Cluster Comm
- 4x4  Clusters----No Cluster Comm
- 8x8  Clusters----Enable Cluster Comm
- 8x8  Clusters----No Cluster Comm
- 16x16 Clusters--Enable Cluster Comm
- 16x16 Clusters--No Cluster Comm

RedMule 128x32 Optimal Efficiency Point

Y-axis: Compute Utilization (0% to 100%)

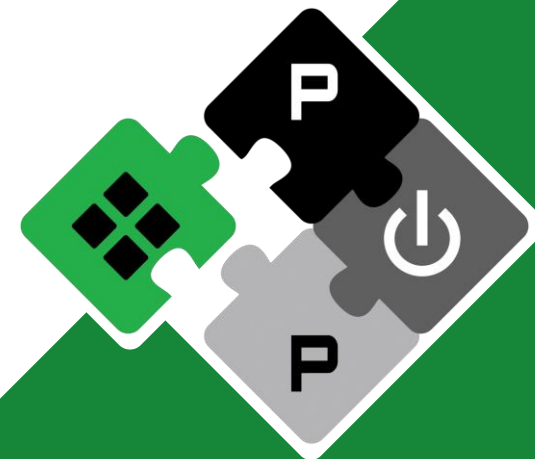| TCDM Size in Cluster: | 10KB | 40KB | 160KB | 640KB | 2.56MB | 10.24MB |
|---|---|---|---|---|---|---|
| Max Tile Dimension: | 32 | 64 | 128 | 256 | 512 | 1024 |

# SoftHier Progress Update

Thomas Benz   tbenz@iis.ee.ethz.ch

**PULP Platform**
Open Source Hardware, the way it should be!

# iDMA

- Various contributions and bugfixes merged in iDMA
  - Detailed tracer will soon be enabled in Snitch (port to GVSoC pending)
  - Release early next week

- Transposition engine implemented by student overhauled
  - Support for packed SIMD types is still ongoing effort
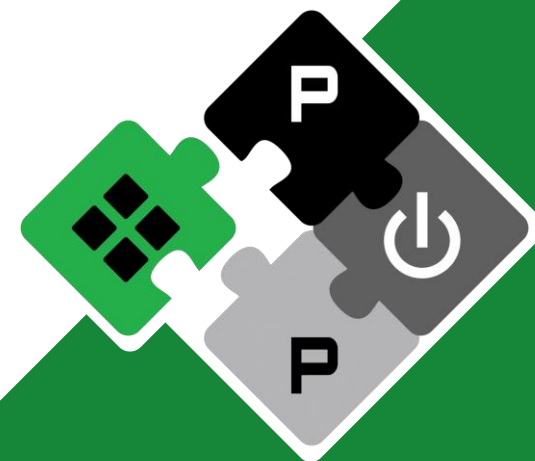  - Update in Snitch is ongoing

# SoftHier Progress Update

Luca Colagrande        colluca@iis.ee.ethz.ch

**PULP Platform**
Open Source Hardware, the way it should be!

@pulp_platform

pulp-platform.org

youtube.com/pulp_platform

# Snitch cluster

- A lot of maintenance work:
  - Merged PRs: #165, #163, #161
  - WIP PRs: #115, #71, #158

- Implemented MHA and MLP layers for Snitch/Occamy [Link]
  - Work on streamlining data generation functions and scripts, to reuse base layer functions (e.g. GEMM, Layernorm) in composite layers (e.g. MHA, MLP)
  - Set up a proper Python package infrastructure to cross-reference these functions

- WIP on full encoder block

**ETH** *zürich*    ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA