

DESIGN, AUTOMATION AND TEST IN EUROPE

THE EUROPEAN EVENT FOR ELECTRONIC SYSTEM DESIGN & TEST

31 MARCH – 2 APRIL 2025 LYON, FRANCE

CENTRE DE CONGRÈS DE LYON



TCDM Burst Access: Breaking the Bandwidth Barrier in Shared-L1 RVV Clusters Beyond 1000 FPUs

Diyou Shen^{1*}, Yichao Zhang^{1*}, Marco Bertuletti^{*}, Luca Benini^{*§} ^{*}ETH Zurich, Switzerland [§]University of Bologna, Italy ¹ These two authors contributes equally to the paper

ETH zürich

CORENEXT



Many-Core Cluster is Rocking

Emerging fields need more computation

- Attention-Based AI Models
- B5G Communication





[1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, & Illia Polosukhin. (2023). Attention Is All You Need.
[2] Marco Bertuletti, Yichao Zhang, Alessandro Vanelli-Coralli, & Luca Benini. (2022). Efficient Parallelization of 5G-PUSCH on a Scalable RISC-V Many-core Processor.

Diyou Shen / ETH Zurich



Many-Core Cluster is Rocking

Emerging fields need more computation

- Attention-Based AI Models
- B5G Communication
- How to keep efficient?
 - Reduce number of instructions
 - Reduce data duplication & memory transfers



More Cores

More Data Movements

Many-Core Cluster is Rocking

Emerging fields need more computation

- Attention-Based AI Models
- B5G Communication



- How to keep efficient?
 - Reduce number of instructions
 - Reduce data duplication & memory transfers



Many-Core Vector Shared-L1 Cluster

Hierarchical design is needed to route over hundreds cores

Next Hierarchy

Building Block: Tile



Diyou Shen / ETH Zurich

DAT





• Limited routing resources for connecting hundreds of cores

4/4/2025





- Limited routing resources for connecting hundreds of cores
- Hierarchical interco. designs
 - Physically feasible
 - Maintain low latency
 - But, limit the BW due to arbitration





- Limited routing resources for connecting hundreds of cores
- Hierarchical interco. designs
 - Physically feasible
 - Maintain low latency
 - But, limit the BW due to arbitration
 - More pronounced in Vector PEs due to unit-stride visiting patterns





- Limited routing resources for connecting hundreds of cores
- Hierarchical interco. designs
 - Physically feasible
 - Maintain low latency
 - But, limit the BW due to arbitration
 - More pronounced in Vector PEs due to unit-stride visiting patterns





- Limited routing resources for connecting hundreds of cores
- Hierarchical interco. designs
 - Physically feasible
 - Maintain low latency
 - But, limit the BW due to arbitration
 - More pronounced in Vector PEs due to unit-stride visiting patterns



- Fits well with Vector PE's access pattern
- No extra resource needed



C2

A[2]

C3

A[3]

Congestion

Pack the requests into a short burst

• Fits well with Vector PE's access pattern

C1

No extra resource needed

Req

CO

Baseline

C5

A[5]

 $\mathbf{C4}$

A[4]

C6

A[6]

C7

A[7]

C8

C9

C10

C11

- Fits well with Vector PE's access pattern
- No extra resource needed



- Fits well with Vector PE's access pattern
- No extra resource needed



- Fits well with Vector PE's access pattern
- No extra resource needed



- Fits well with Vector PE's access pattern
- No extra resource needed



- Fits well with Vector PE's access pattern
- No extra resource needed



Our Solution: TCDM Burst Access





4/4/2025

Diyou Shen / ETH Zurich

DATE

Group multiple data in a single transfer

- Have higher data bandwidth on response channel
- Minimize the cost: only widen the response data field
- Configurable: 2x data width (GF2), 4x data width (GF4), ...

Group multiple data in a single transfer

- Have higher data bandwidth on response channel
- Minimize the cost: only widen the response data field
- Configurable: 2x data width (GF2), 4x data width (GF4), ...





DAI



DA



DA



Diyou Shen / ETH Zurich

DAT

Fully Routable and Highly Efficient!

Implement on MemPool₆₄-Spatz₄ Cluster in GF12 tech.

3.26x Effective BW

Up to **2.76x** Performance Improvement



770 MHz @WW Less than 7.7% Area

Up to **1.9x** Energy Efficiency

We design TCDM Burst Access



• A physically feasible, high performance, low energy interconnection extension for large shared-L1 clusters.

ETH zürich CO

Come to our poster and discover more!

- Tested on different cluster sizes.
- Roofline models to evaluate the designs.
- Fully open-sourced!

github.com/pulp-platform/mempool

Diyou Shen <u>dishen@iis.ee.ethz.ch</u> ETZ, Gloriastrasse 35, 8092 Zürich @pulp_platform





27