

NextRAN-AI – 05/09/2025

Integrated Systems Laboratory (ETH Zürich)

Marco Bertuletti

mbertuletti@iis.ee.ethz.ch

Yichao Zhang

yiczhang@iis.ee.ethz.ch

Mahdi Abdollahpour

mahdi.abdollahpour@unibo.it

Alessandro Vanelli-Coralli

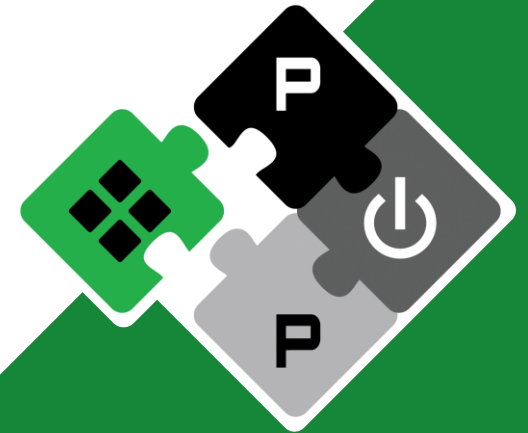
avanelli@iis.ee.ethz.ch

Luca Benini

lbenini@iis.ee.ethz.ch

PULP Platform

Open Source Hardware, the way it should be!



pulp-platform.org

[@pulp_platform](https://twitter.com/pulp_platform)

[company/pulp-platform](https://www.linkedin.com/company/pulp-platform)

[youtube.com/pulp_platform](https://www.youtube.com/pulp_platform)



Roadmap



- **TensorPool:**

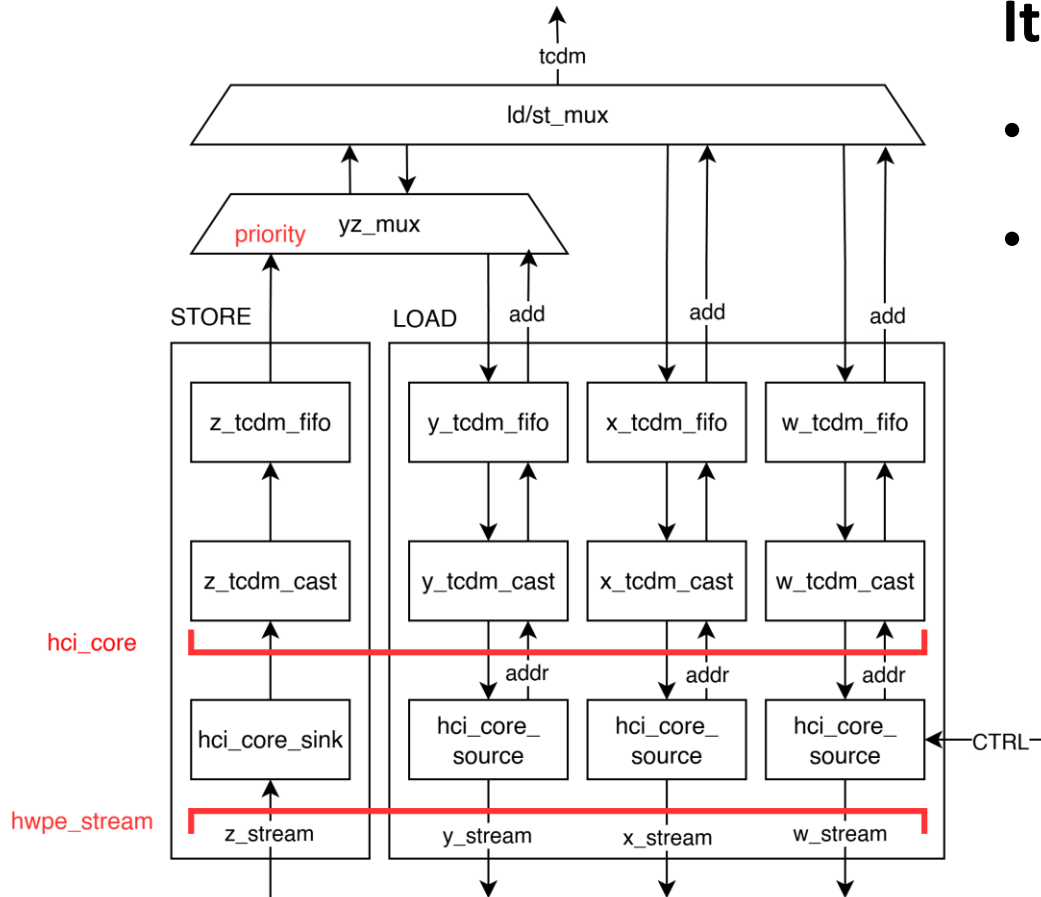
- Optimize RedMule to L1 connection
- TeraPool physical design in 7nm
- TeraPool + RedMule (TensorPool) physical design in 7nm
- PPA on model microkernels and operators (combined RedMule&Cores)

- **System Performance:**

- Simulation speed is impaired by large design size → from RTL simulation to higher abstraction level, e.g. GVSoC
- TensorPool GVSoC model developement
- Data-Movement and end-end performance

-
- a)** High-level block diagram of RedMulE. The architecture consists of a TCDM (288b) connected to a REDMULE CAST (288b), which feeds into a STREAMER (256b). The STREAMER outputs to X FIFO (256b), W FIFO (256b), and Y FIFO (256b). X FIFO feeds into a SCHEDULER, which feeds into an X-BUFFER (256b). W FIFO feeds into a W-BUFFER (256b). Y FIFO feeds into a Y/Z-BUFFER (256b). The X-BUFFER, W-BUFFER, and Y/Z-BUFFER all feed into the DATAPATH. The DATAPATH outputs to the Y/Z-BUFFER. The Y/Z-BUFFER outputs to a Z FIFO (256b). The Z FIFO feeds into a PERIPH INTERCO CONTROLLER. The PERIPH INTERCO CONTROLLER feeds into the DATAPATH.
- b)** Detailed datapath showing a 3x4 grid of CE blocks (CE 0,0 to CE 11,3). Each CE block has inputs X, W, and Y, and outputs Z. The inputs X, W, and Y are connected to the CE blocks via multiplexers. The outputs Z are connected to the CE blocks via multiplexers. The CE blocks are arranged in rows (ROW 0, ROW 1, ROW 11) and columns (COLUMN 0, COLUMN 1, COLUMN 2, COLUMN 3).
- c)** RedMulE Parameters
- | Parameter | Description |
|-----------|---|
| L | # of rows of CEs |
| H | # of columns of CEs
(# of CEs per row) |
| P | # of pipeline stages per CE |

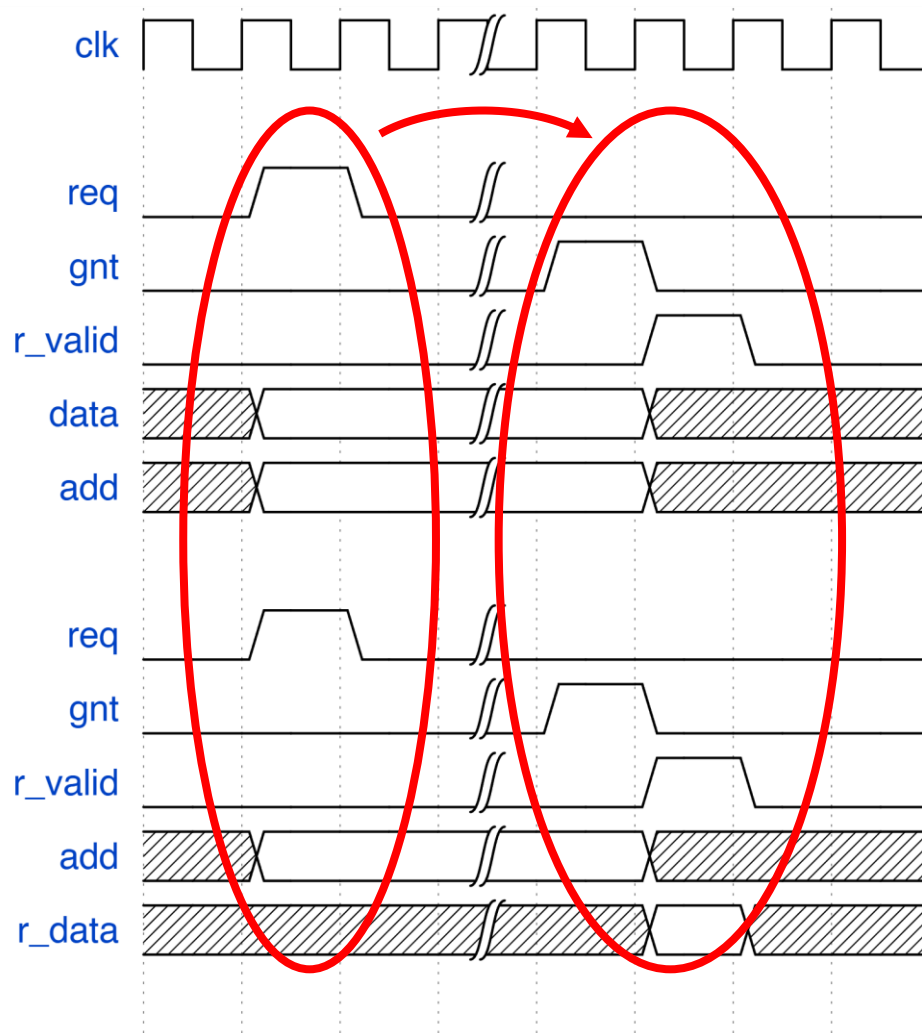
3



It is RedMulE's memory interface:

- Control generates addresses (streams)
- Streamer:
 - Generates wide $H \cdot (P+1) \cdot 16\text{bit} + 32\text{bit}$ TCDM request, based on the addresses from control
 - Applies conversions (e.g. 16/8bit operands)
 - Multiplexes x, w, y, z streams

RedMulE Streamer: Protocol



Requests are based on a req/gnt/r_valid protocol:

- req is asserted on valid requests
- A new request can be issued when the slave grants the transaction
- The response must be valid one cycle after the transaction is granted

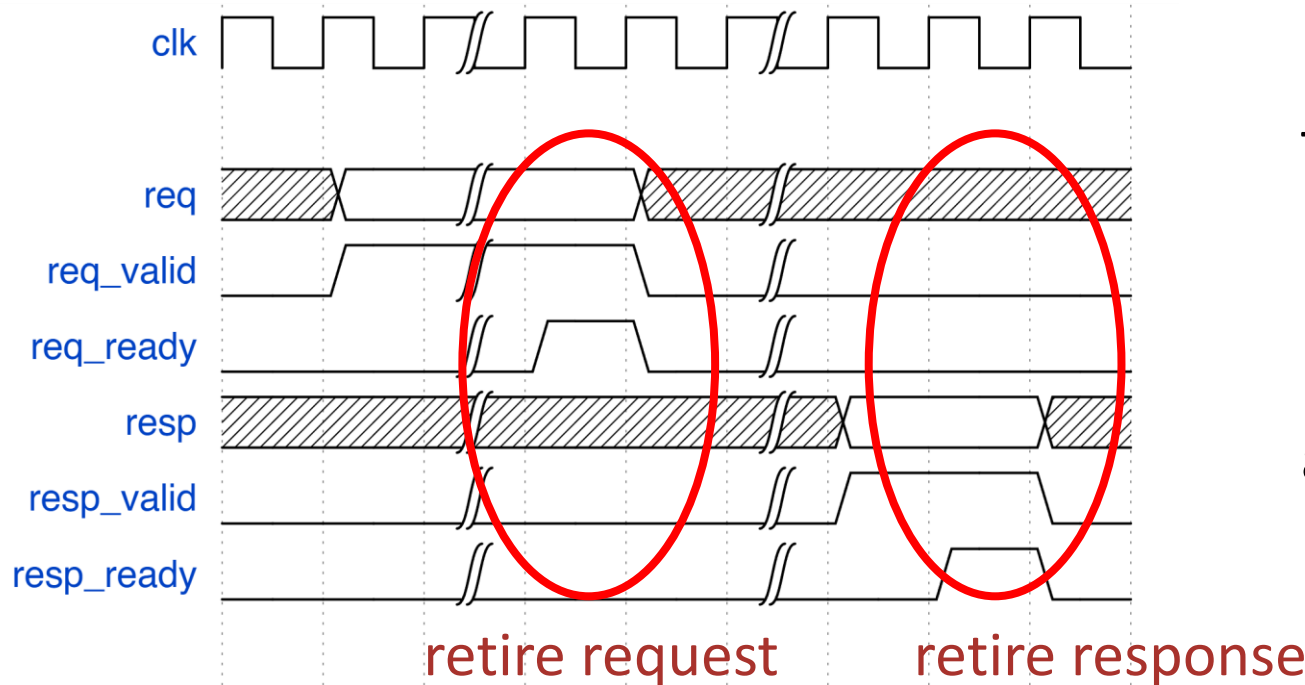
When the memory takes multiple cycles to handshake no new request can be posted on the interconnect!

Enabling outstanding requests



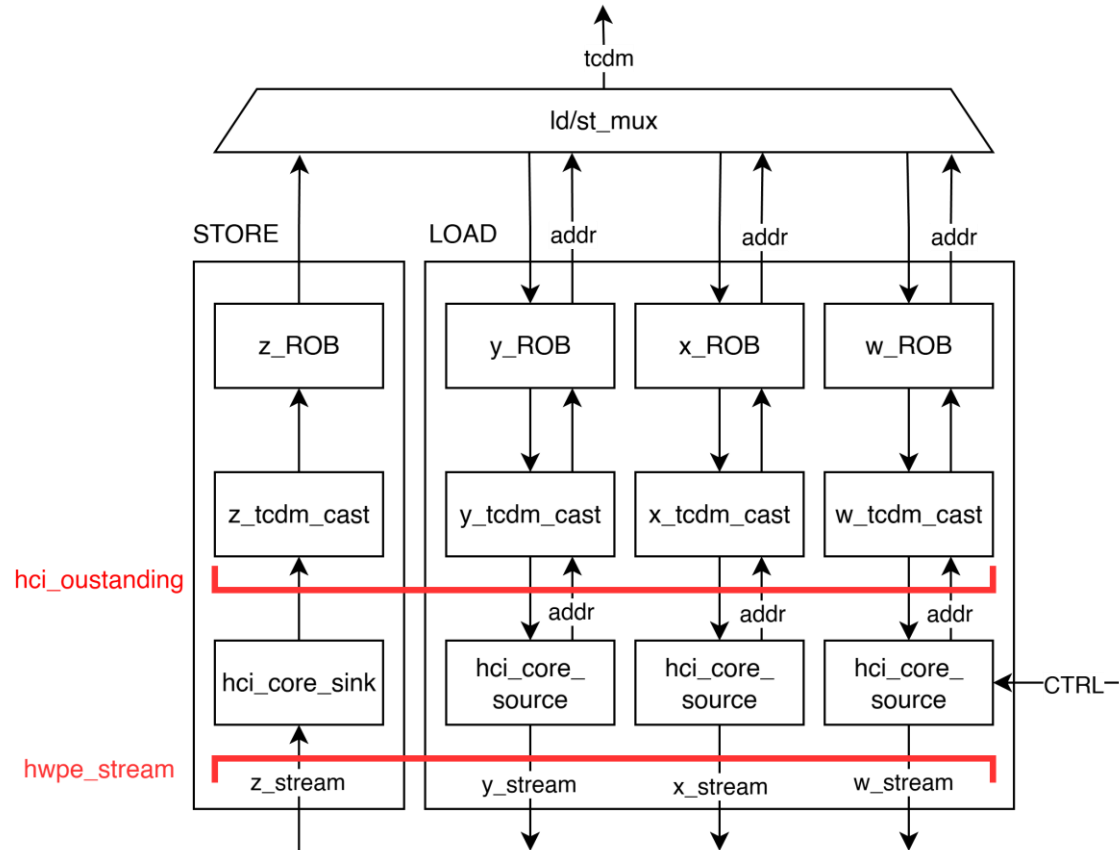
A req_valid/ready, resp_valid/ready protocol decouples req and resp

- A new valid req can be posted as soon as req_valid = resp_ready = 1'b1
- Valid responses can be collected later, when resp_valid = resp_ready = 1'b1



This however implies that multiple transactions could be in-flight, therefore we need to send transaction IDs and reorder responses

Enabling outstanding requests

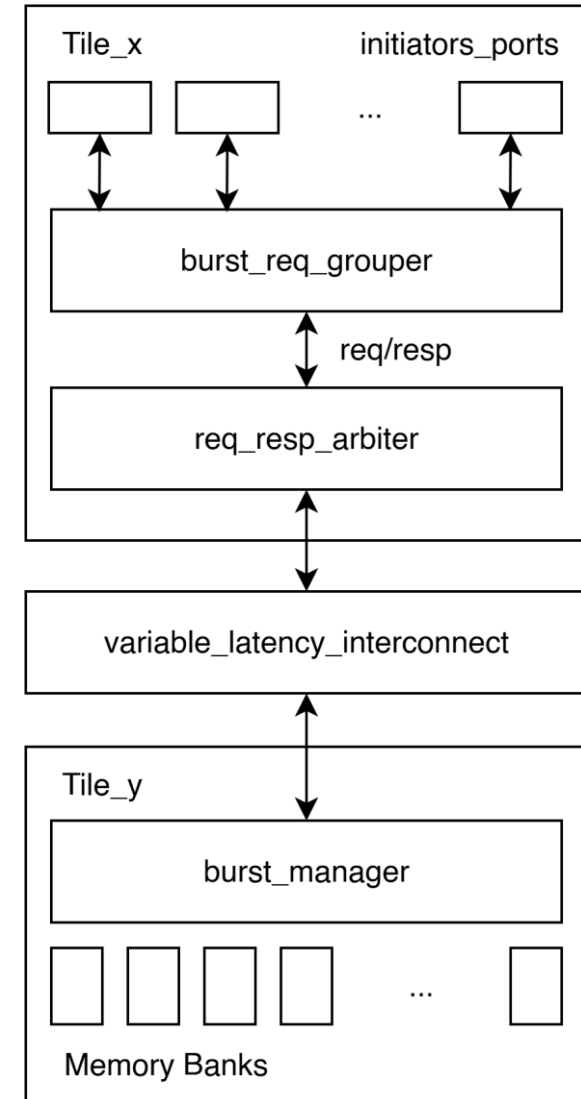


Progress:

- Replaced the req/gnt/r_valid protocol with a req_valid/ready resp_valid/ready protocol
- Added ROB with parametrizable number of transactions on each stream

Integration in MemPool: Recap on bursts

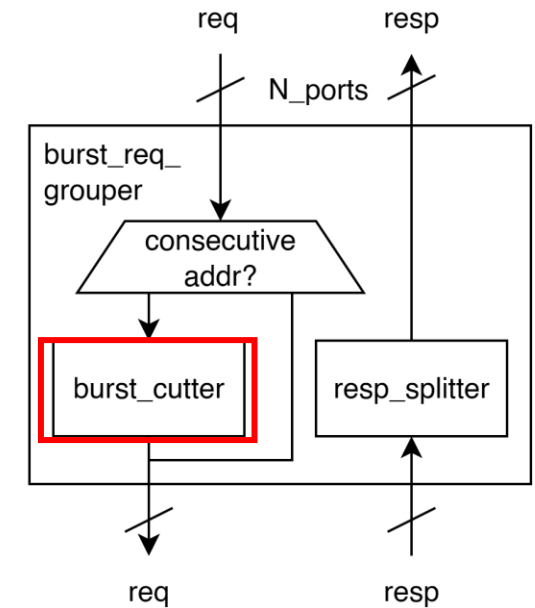
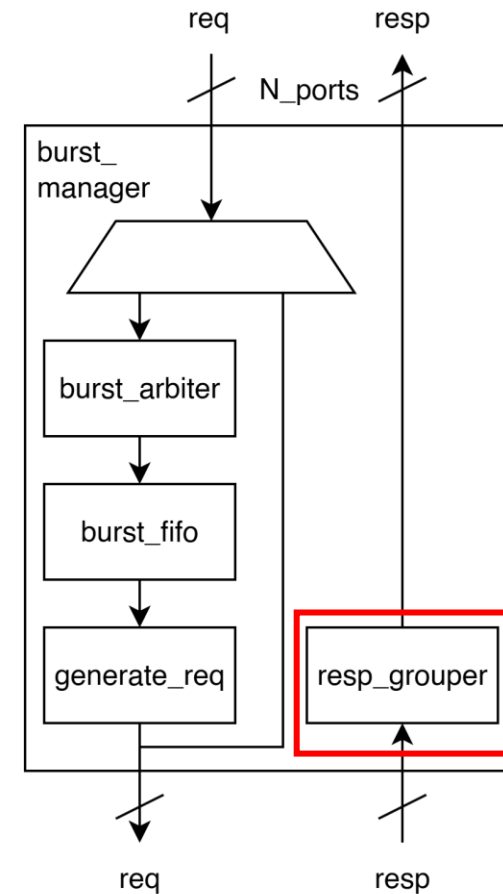
- System with distributed shared-memory and hierarchical core-memory interconnect
- Arbitration of wide parallel TCDM requests in shared interconnect resources → performance penalty
- **Burst the transaction:**
 - Group a request to consecutive addresses in a burst (only first port is valid and is arbitrated)
 - Propagate request in the shared-memory hierarchical interconnect
 - Identify burst request and redistribute to memory banks with consecutive addresses in the destination Tile
 - Send back wide-response



Integration in MemPool: Recap on bursts

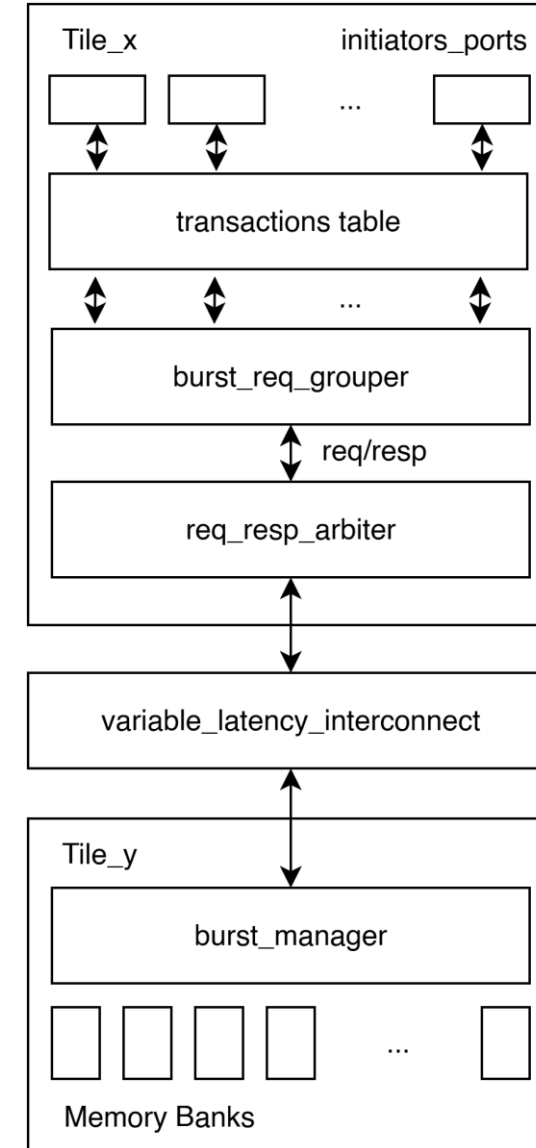


- We can reduce arbitration also in the response path by grouping responses on the same response transactions
(Done by the **resp_grouper**)
- When requests cross the boundary of a Tile two burst must be sent
(Done by the **burst_cutter**)



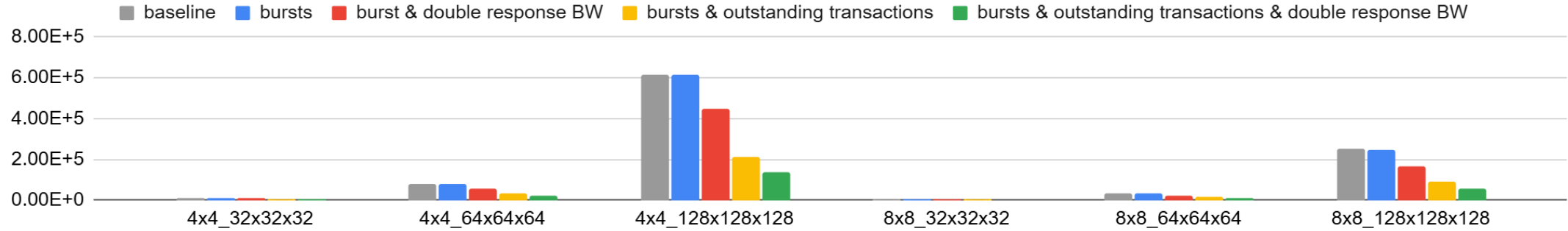
Integration in MemPool: transactions table

- Narrow responses come back out-of-order
- With outstanding transactions a transaction table reorders the narrow memory requests

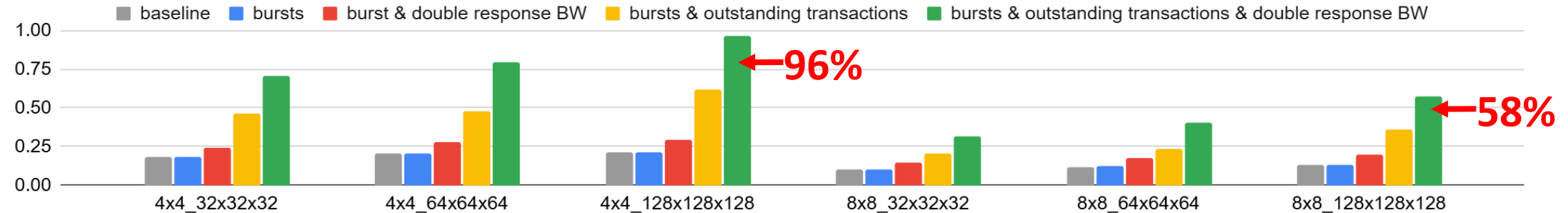


Results (single RedMulE)

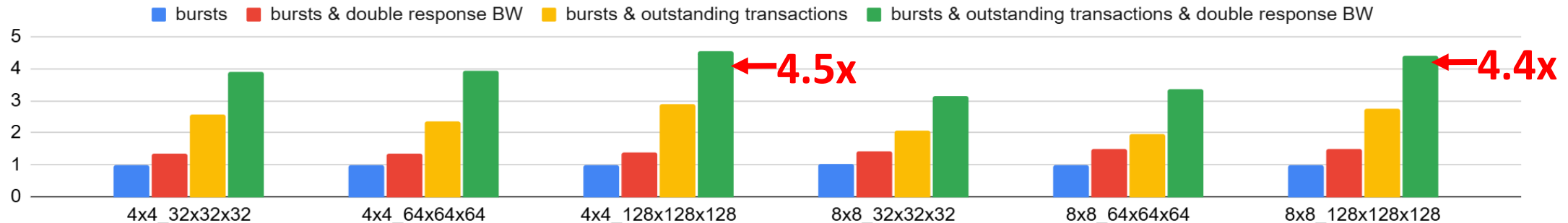
Runtime (cycles)



RedMulE PEs Utilization



SpeedUP



Next Steps



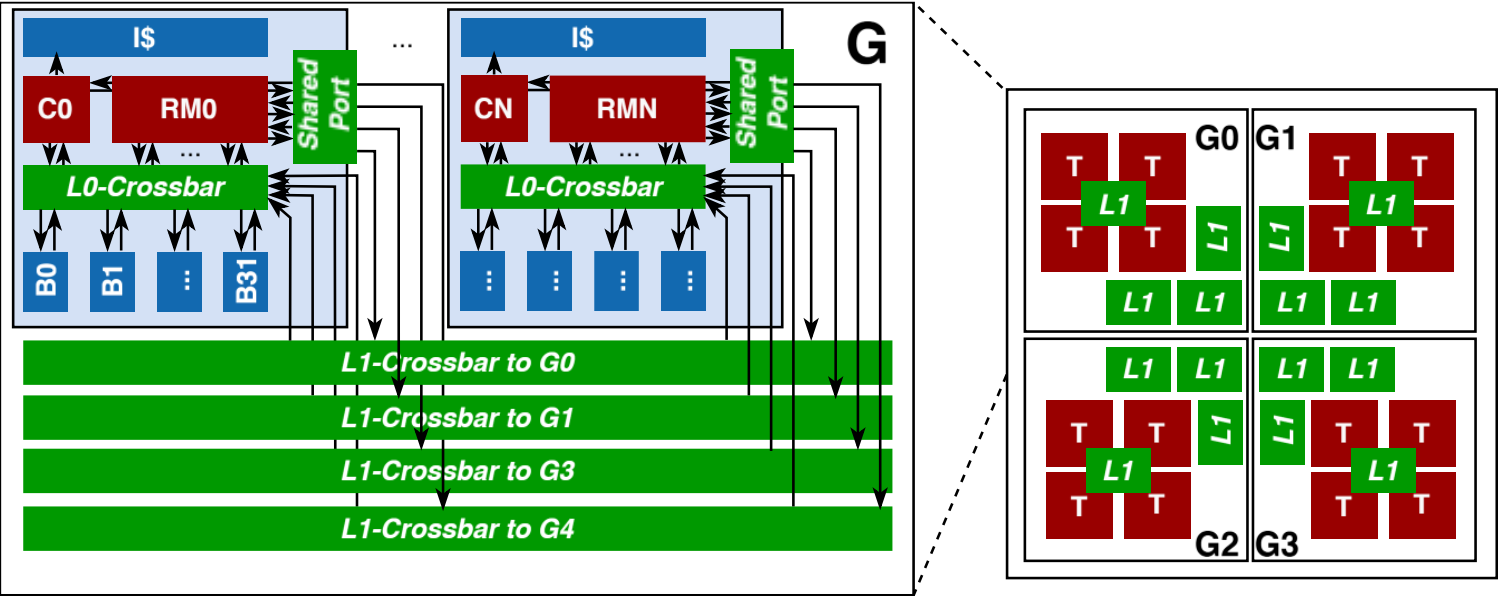
- **TensorPool:**

- Optimize RedMule to L1 connection → modify TeraPool interconnect for TensorPool-5 config.
- TeraPool physical design in 7nm
- TeraPool + RedMule (TensorPool) physical design in 7nm
- PPA on model microkernels and operators (combined RedMule&Cores)

- **System Performance:**

- Simulation speed is impaired by large design size → from RTL simulation to higher abstraction level, e.g. GVSoC
- TensorPool GVSoC model development
- Data-Movement and end-end performance

Recap TensorPool-5



Best option is highlighted and increases the BW of RedMule

FLOPS	NCores	NRM	NBank /Tile	NRM Tiles	LV1RM Ports	LV2RM Ports	NTiles	NG	NSGs	LEVEL 0		LEVEL 1		LEVEL 2	
										INxOUT	N	INxOUT	N	INxOUT	N
4096	0	16	512	1	0	0	0	0	0	512x512	1	-	-	-	-
4096	0	16	32	8	1	0	0	1	0	64x64	8	8x8	1	-	-
4096	0	16	32	16	1	0	0	1	0	32x32	16	16x16	1	-	-
4096	0	16	32	16	2	0	0	1	0	32x32	16	32x32	1	-	-
4096	0	16	32	16	1	0	0	4	0	32x32	16	4x4	16	-	-
4096	0	16	32	16	8	0	0	4	0	32x32	16	32x32	16	-	-