

TeraPool: Boosting Wireless Communications by Pooling 1000s cores with PULP

Integrated Systems Laboratory (ETH Zürich)

Marco Bertuletti

Yichao Zhang

Alessandro Vanelli-Coralli

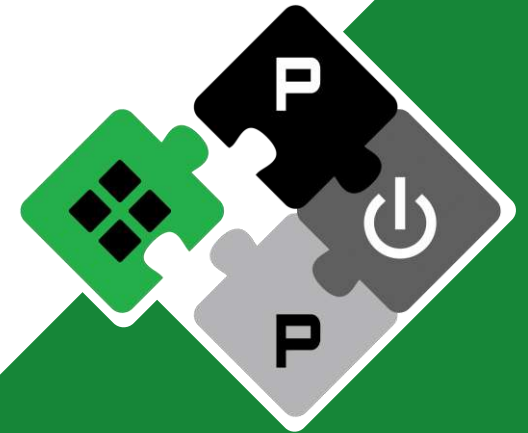
Luca Benini

mbertuletti@iis.ee.ethz.ch

yiczhang@iis.ee.ethz.ch

avanelli@iis.ee.ethz.ch

lbenini@iis.ee.ethz.ch



PULP Platform

Open Source Hardware, the way it should be!

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



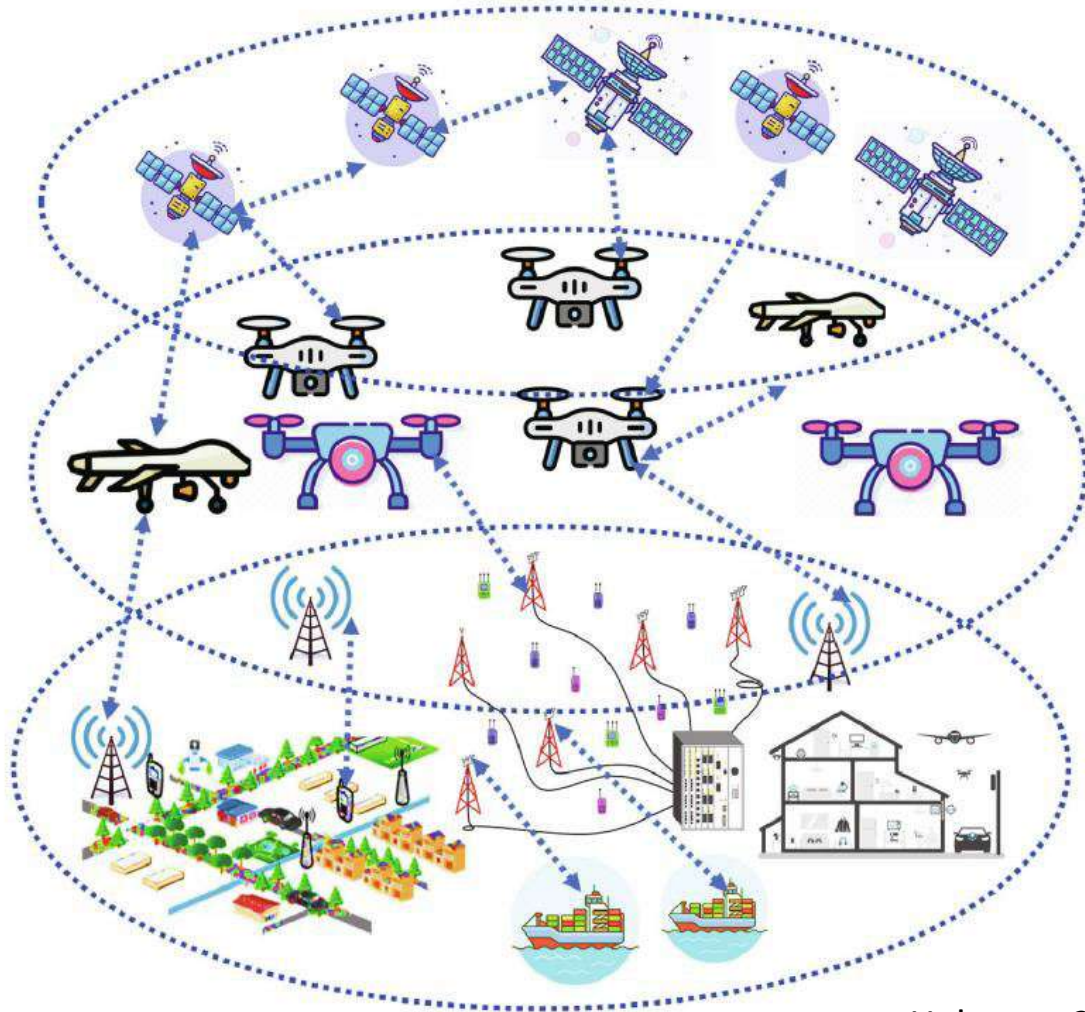
pulp-platform.org



youtube.com/pulp_platform



6G: A tight mesh of network tranceivers



- **Underwater Network**

- Autonomous Underwater Vehicles

- **Terrestrial Network**

- CRAN (CU + DU + RU)
- Vehicular, M2M, IoT

- **Aerial Network**

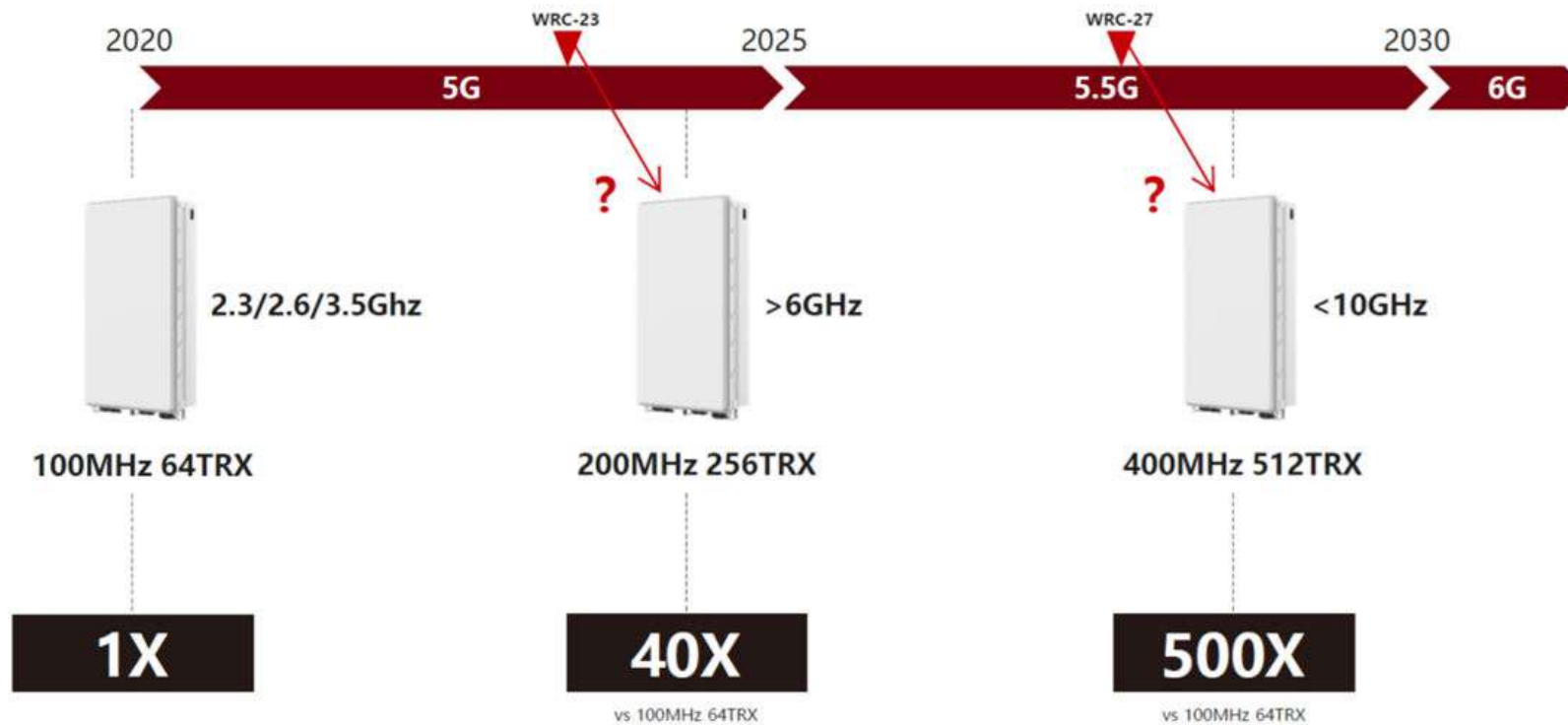
- Drones (Unmanned Aerial Vehicles)
- LEO, MEO, GEO Satellites

Hakeem, 2022

This huge complexity is a computing problem!



6G puts constraints on the computing load of the Base-stations, in terms of throughput and latency...



Khan, 2022

More BW, lower latency, more devices!

- < 10 GHz → **3 THz**
- 20 Gbps → **1 Tbps**
- 1 ms → **0.1/0.01 ms**
- 10^6 dev/km² → **10^7 dev/km²**

Rappaport, 2019

6G drives needs of today's hardware



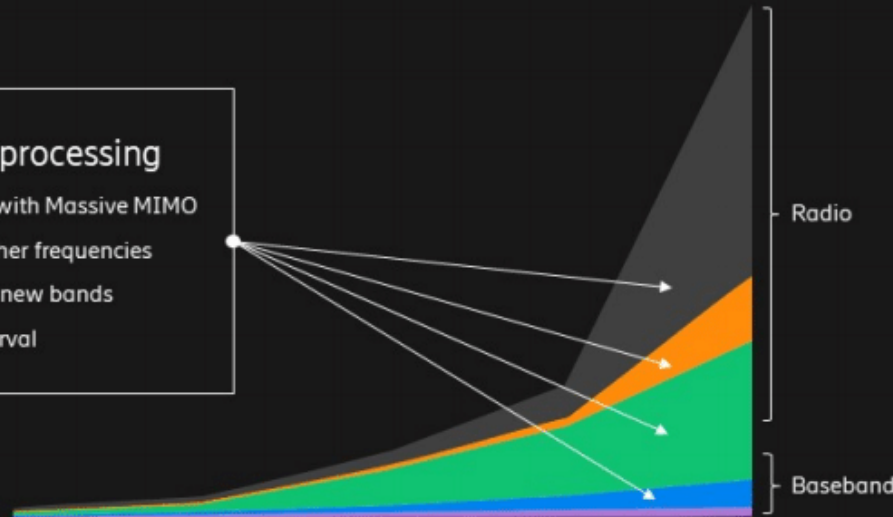
ISSCC,
Ekholm,
2023

Mobile Network Traffic is Driving Compute Needs at a Higher Rate Than Moore's Law

Trends driving higher processing

- # Antenna branches growing with Massive MIMO
- More carrier bandwidth at higher frequencies
- Wider spectrum allocations in new bands
- Shorter transmission time interval

Carrier bandwidth	20 MHz	20 MHz	100 MHz	100 MHz	100 MHz
Antenna branches	2T2R	4T4R	4T4R	8T8R	64T64R
Transmission time interval	1 ms	1 ms	0.5 ms	0.5 ms	0.5 ms



Ericsson Silicon

A range of purpose built ASICs in advanced technology augmented by partner innovation.

- Digital front-end processing
- Beamforming processing
- Layer 1 processing
- Layer 2 processing
- Packet processing function
- Radio control function

© 2023 IEEE
International Solid-State Circuits Conference

Can we meet the 6G tight requirements with PULP?



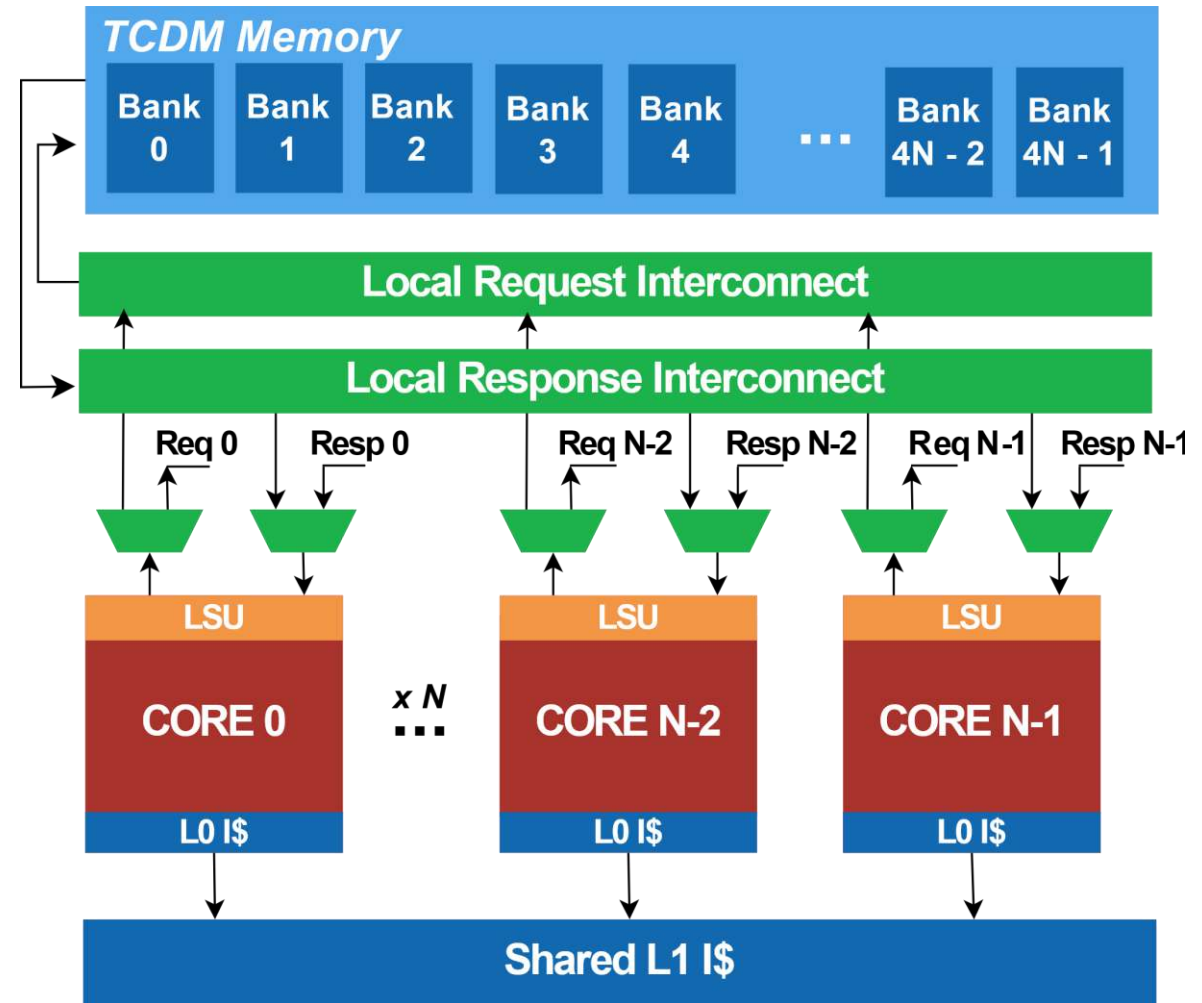
Parallel architectures

- Shared-Memory programming
- Low-Latency Tightly-Coupled Data Memory (TCDM)

ISA specialization

- General-Purpose vs. ASICs
- Open ISA (RISC-V) for architecture specialization

Huge amount of data flowing in 5G/6G workloads!

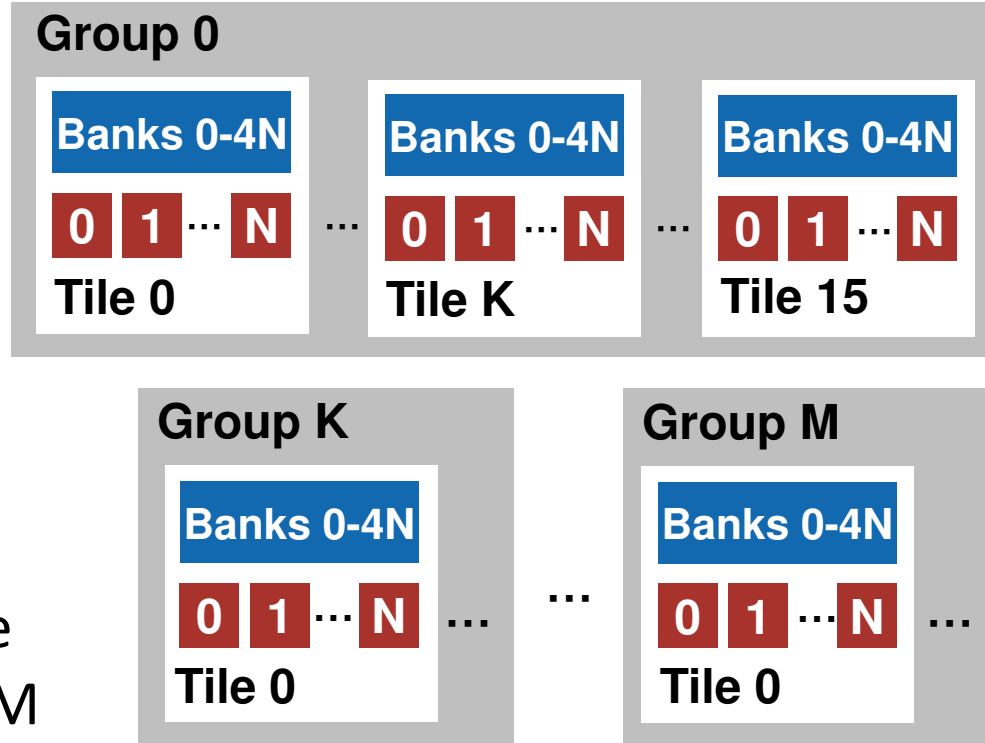
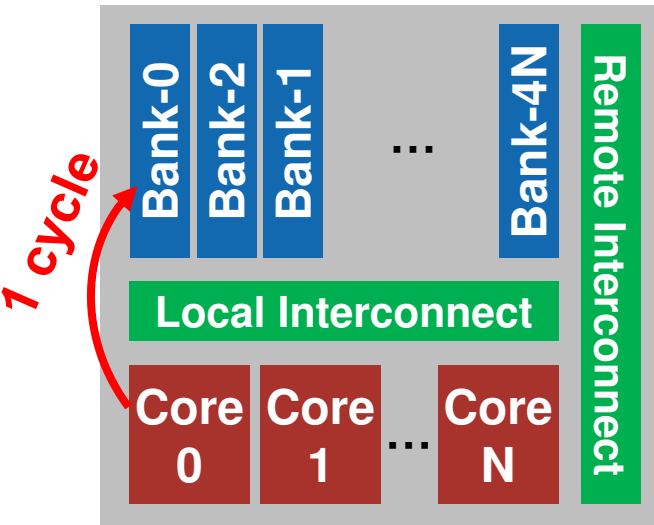


Use more cores and more memory!



MemPool → 256 cores

TeraPool → 1024 cores



Any core can access any bank, some interconnection resources are shared → NUMA

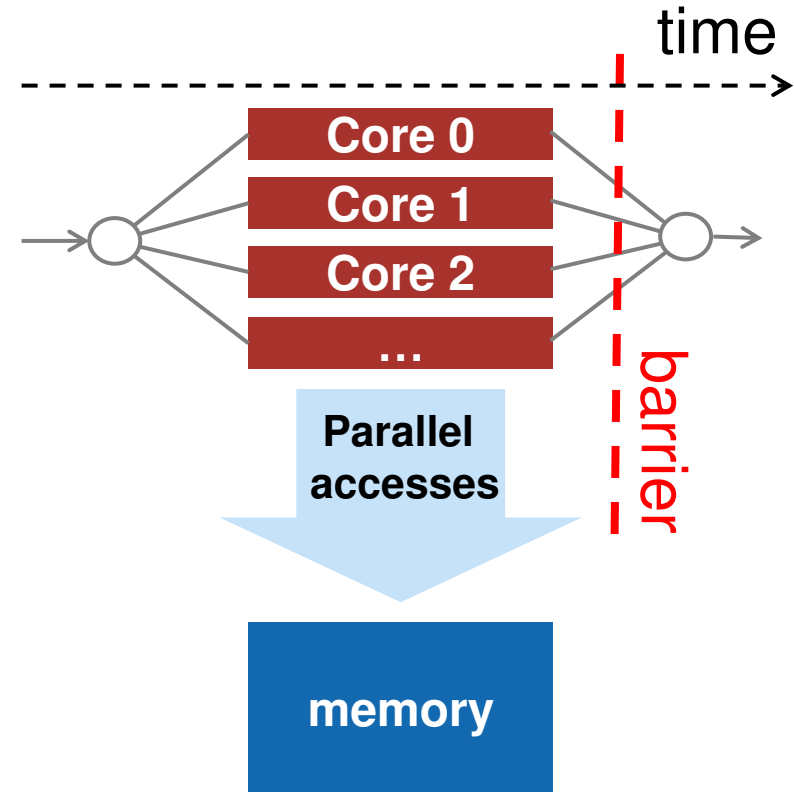
Snitch Cores (rv32ima) are grouped in Tiles with TCDM

*How do you program
1000s cores?*

We target fork-join parallelism

Fork-join programming model

- Serial execution forks to parallel execution
- Cores access memory concurrently
- Cores are synchronized and parallel execution joins to serial

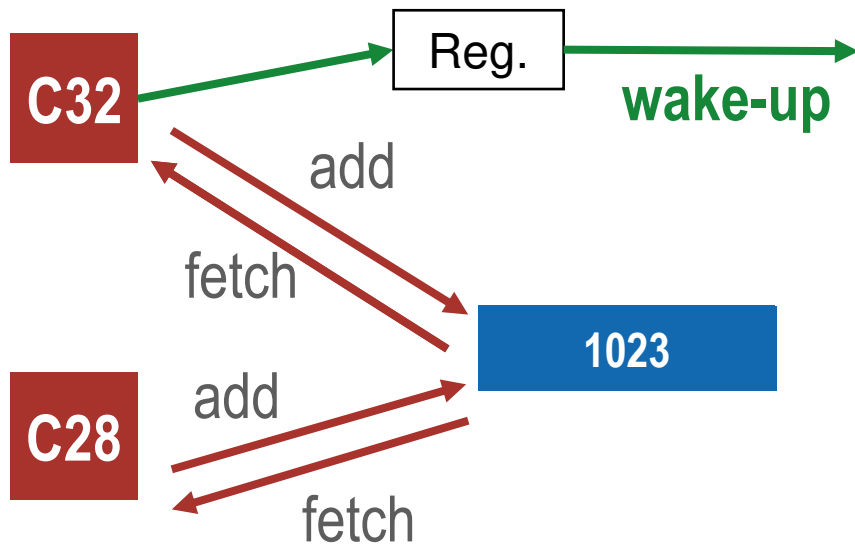


Synchronization: a log-tree approach

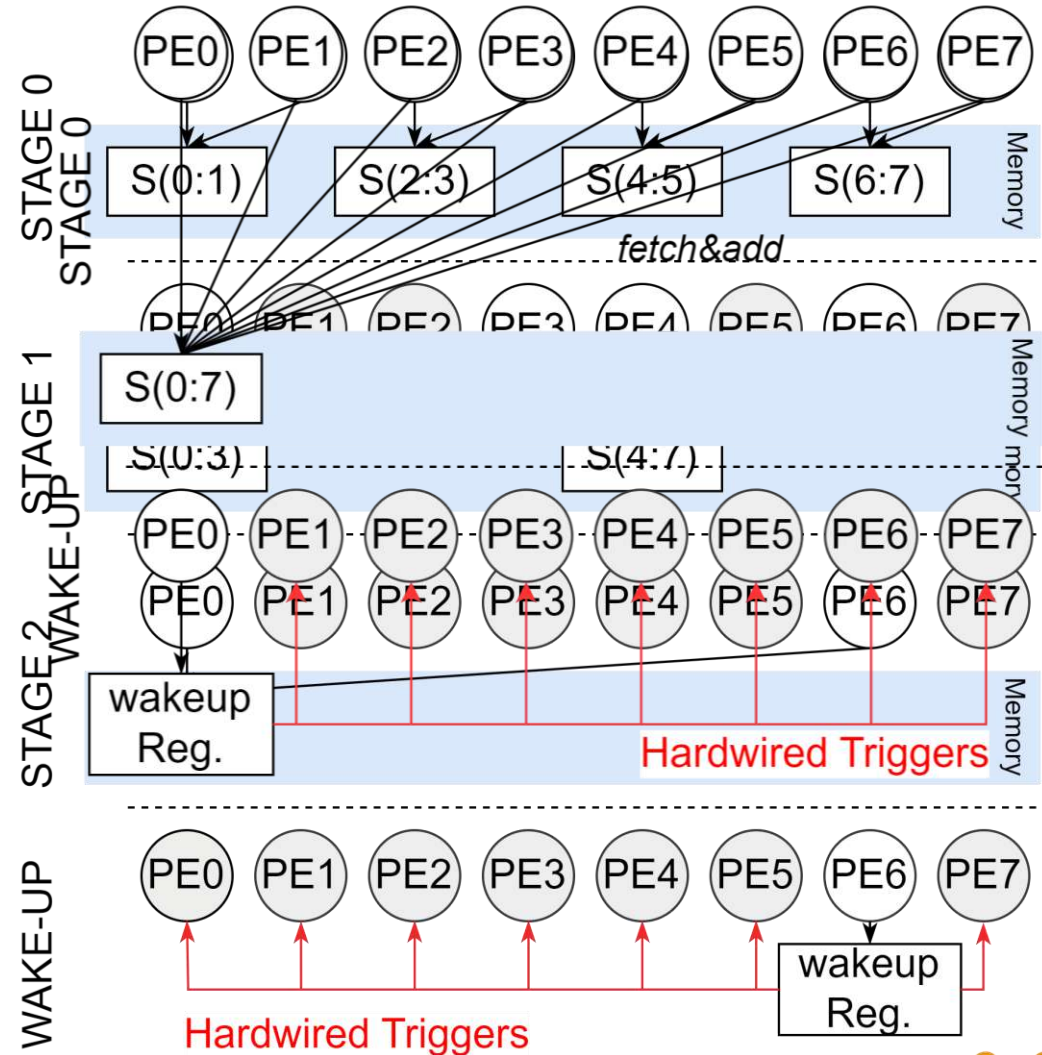


Synchronization barriers

- Arrival = atomic writes to a synch variable
- Hardwired **wake-up triggers** for departure



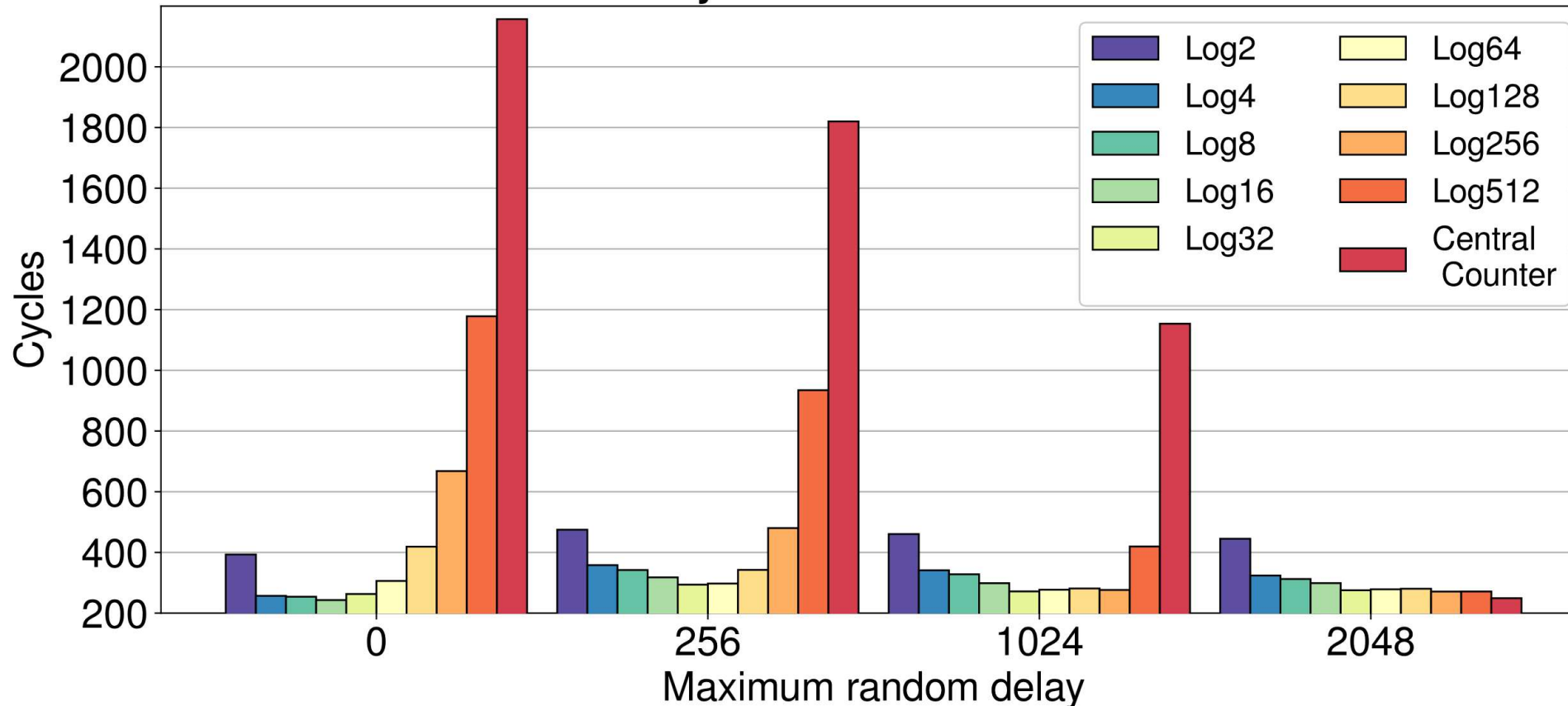
ISSUE: cores arriving all together will contend for the same memory resource!



Overhead depends on cores' arrival time

- Small delays → small radices are better,
- Large delays → central-counter wins

Cycles for barrier call



Choose the barrier depending on the kernel



AXPY

2 ¹⁰	0.86	0.76	0.63	0.46
2 ⁹	0.76	0.63	0.47	0.31
2 ⁸	0.63	0.48	0.33	0.20
2 ⁷	0.50	0.35	0.22	0.12
2 ⁶	0.40	0.27	0.16	0.09
2 ⁵	0.37	0.24	0.14	0.08
2 ⁴	0.35	0.23	0.13	0.07
2 ³	0.36	0.23	0.13	0.08
2 ²	0.37	0.24	0.14	0.08
2	0.49	0.34	0.21	0.12
	65536	131072	262144	524288

Radix

LEN

DCT

2 ¹⁰	0.93	0.58	0.27	0.16
2 ⁹	0.88	0.43	0.17	0.09
2 ⁸	0.81	0.29	0.10	0.06
2 ⁷	0.72	0.20	0.06	0.03
2 ⁶	0.64	0.15	0.04	0.02
2 ⁵	0.60	0.13	0.04	0.02
2 ⁴	0.58	0.12	0.04	0.02
2 ³	0.60	0.13	0.04	0.03
2 ²	0.61	0.13	0.05	0.03
2	0.70	0.19	0.07	0.04
	2x4096	4x4096	8x4096	16x4096

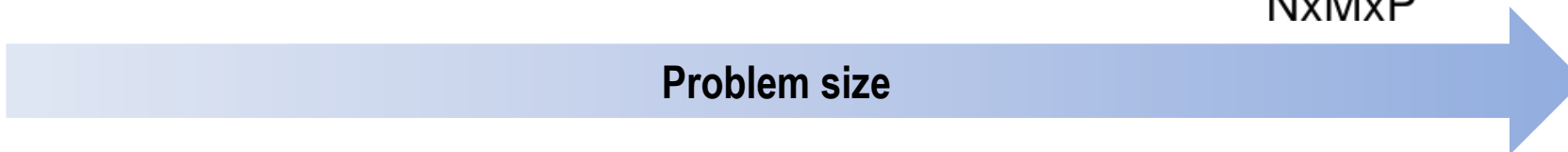
LEN

MATMUL

2 ¹⁰	0.64	0.35	0.11
2 ⁹	0.48	0.20	0.07
2 ⁸	0.34	0.13	0.05
2 ⁷	0.24	0.12	0.05
2 ⁶	0.18	0.11	0.05
2 ⁵	0.16	0.11	0.05
2 ⁴	0.15	0.13	0.06
2 ³	0.16	0.14	0.06
2 ²	0.17	0.14	0.06
2	0.24	0.15	0.07
	128x32x128	128x128x128	256x128x256

NxMxP

Different kernels have different cores' arrival times: choose barrier accordingly!

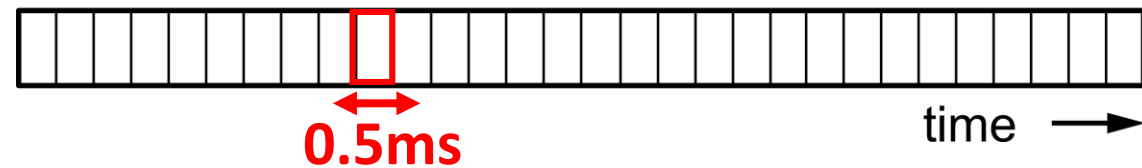
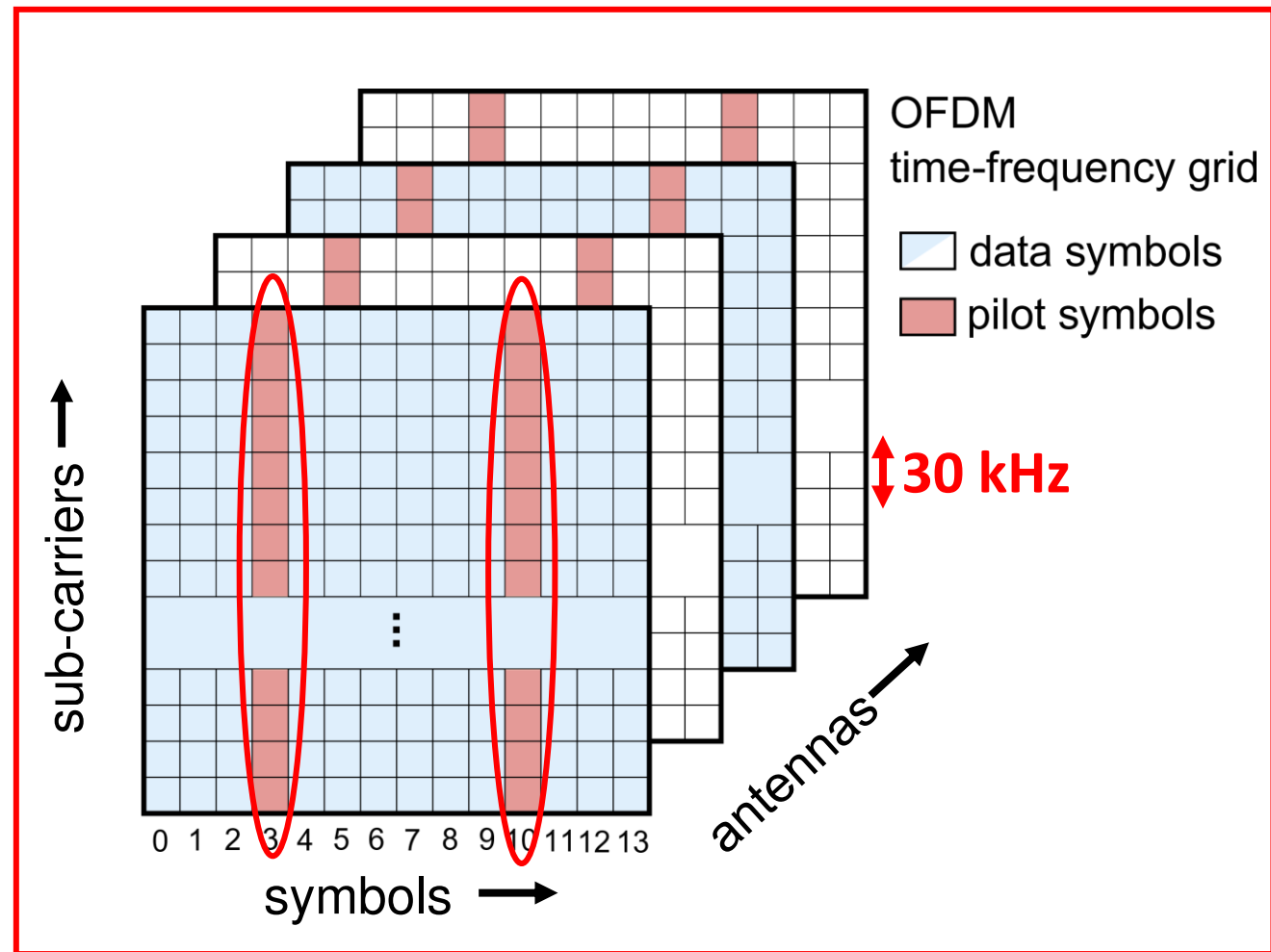


PUSCH processing

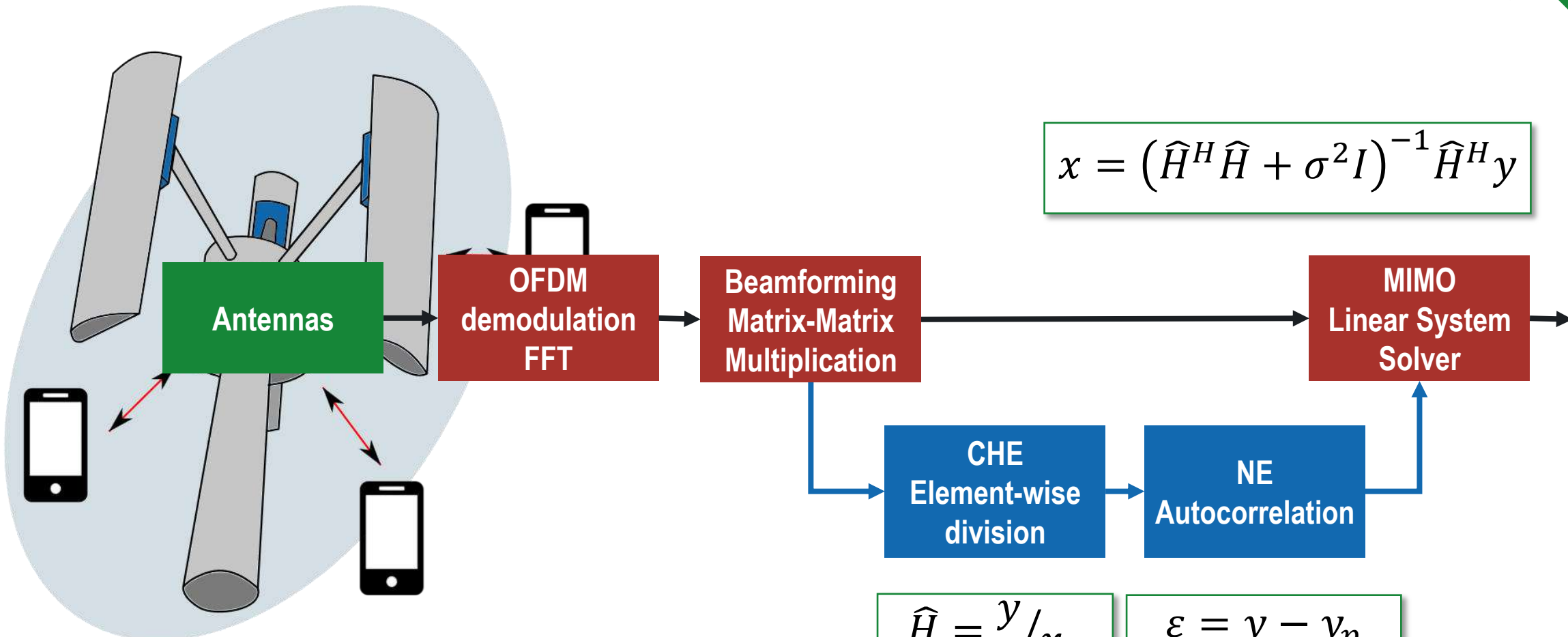
We receive **frequency-multiplexed transmissions** = symbols

- Orthogonal subcarriers
- From multiple antennas
- 14 symbols in Transmission Time-Interval (0.5ms)

(Pilot symbols, are known at the RX + TX, and allow the reconstruction of the channel)



PUSCH processing



$$x = (\hat{H}^H \hat{H} + \sigma^2 I)^{-1} \hat{H}^H y$$

$$\hat{H} = y / x_p$$
$$y_p = \hat{H} x_p$$

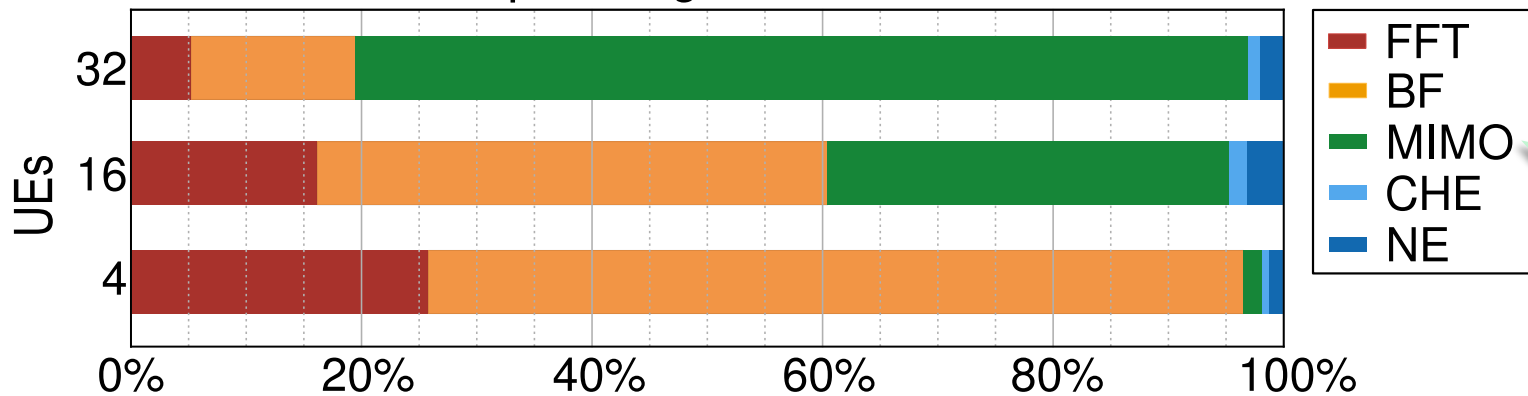
$$\varepsilon = y - y_p$$
$$\sigma^2 I = R_{\varepsilon\varepsilon}$$

PUSCH processing: Computational complexity



- A computational complexity analysis shows that most of the MACs are in the FFT, the BF and the MIMO stages
- We therefore focus on the optimization of these steps

MACs per stage in PUSCH chain

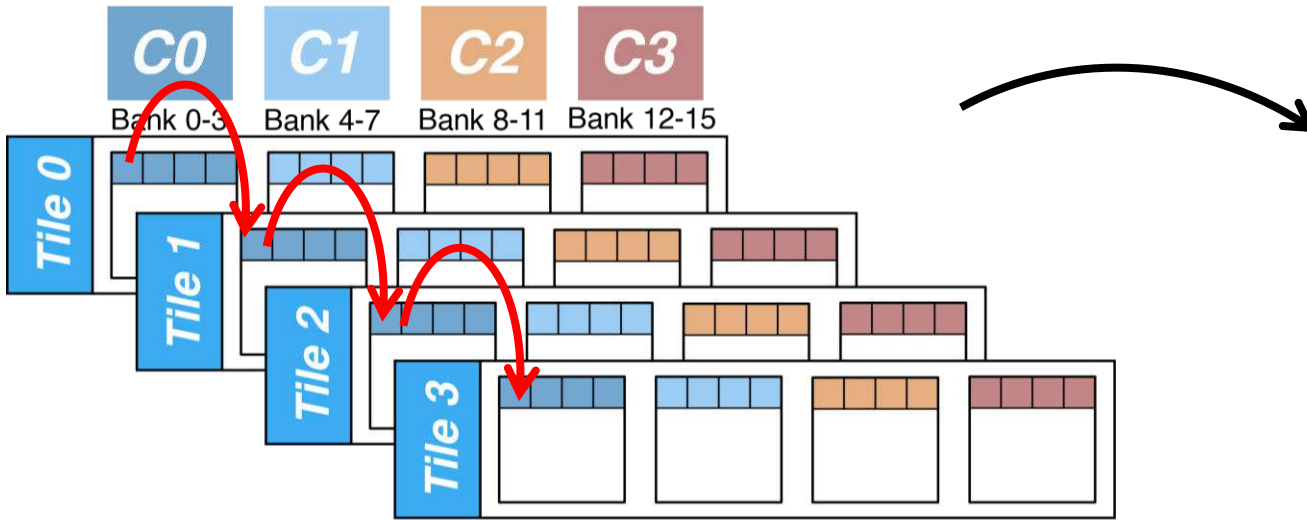


Impact of MIMO stage depends on the number of UEs transmitting on the same sub-carrier.

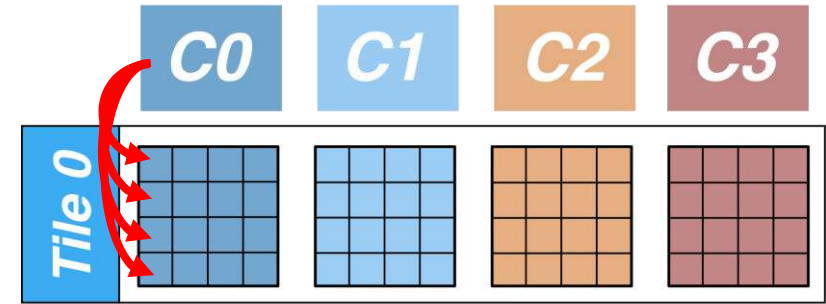
Low access latency programming



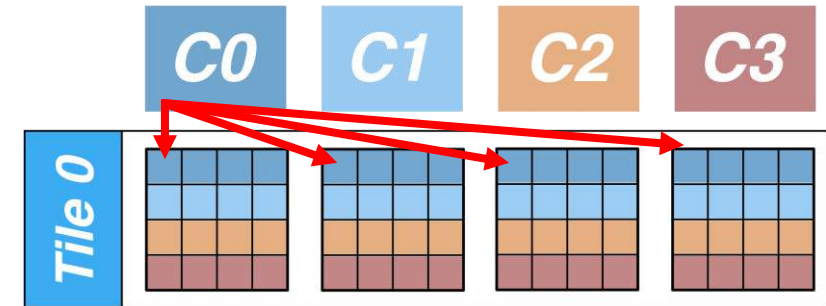
The kernels compute and store locally → reduce load latency of cores that are using data in the next phase of the computation



Load access pattern



Store access pattern

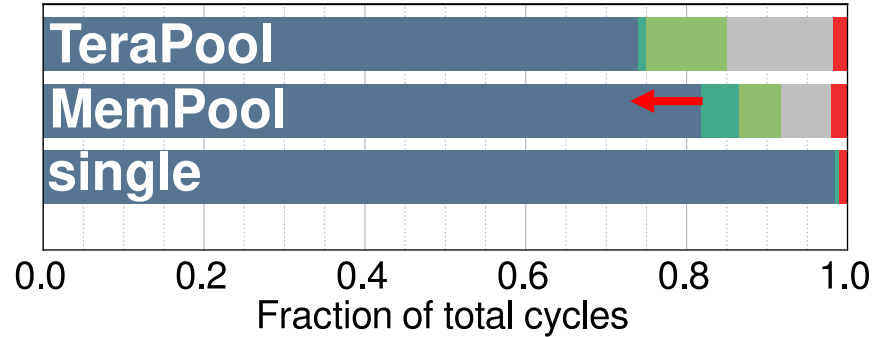


High IPC is obtained on all benchmarks



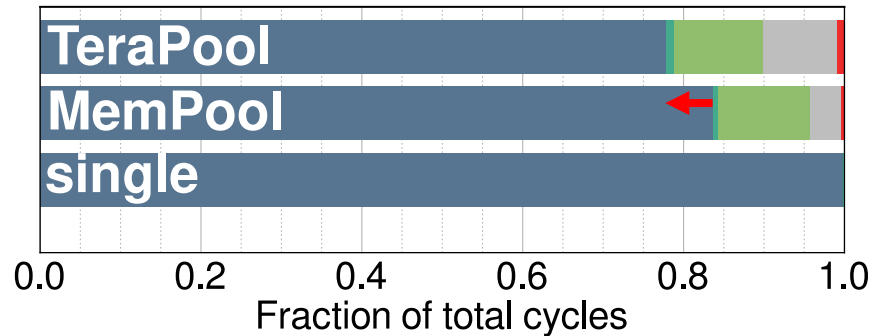
FFT

4096-points
(16 independent
FFTs run between
barriers)



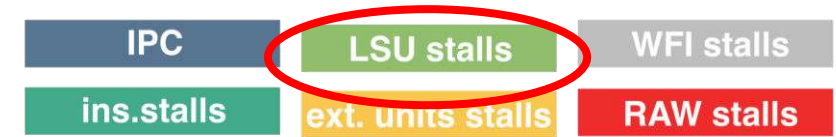
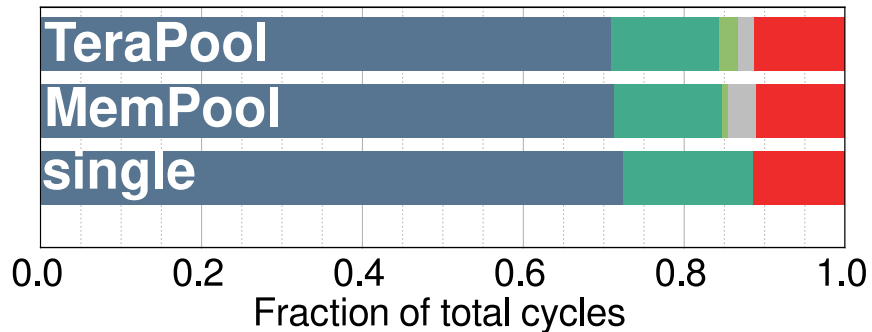
MMM

(Input 1 4096x64
Input 2 64x32)



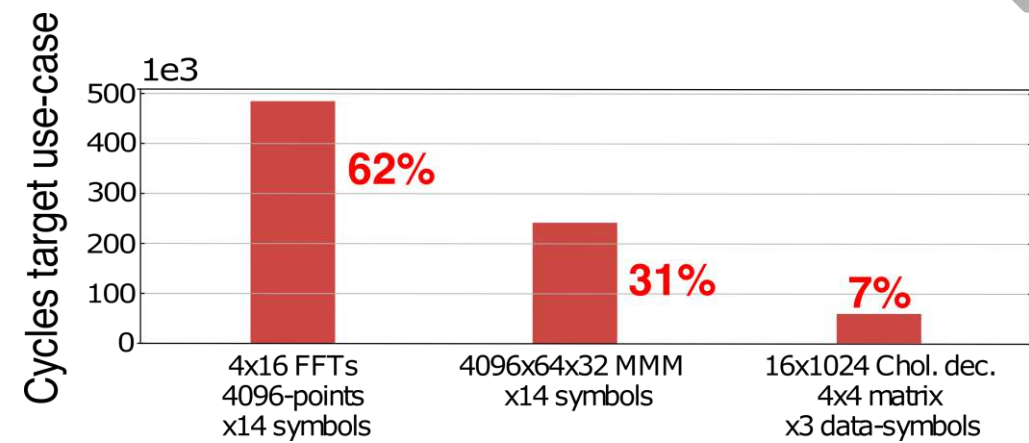
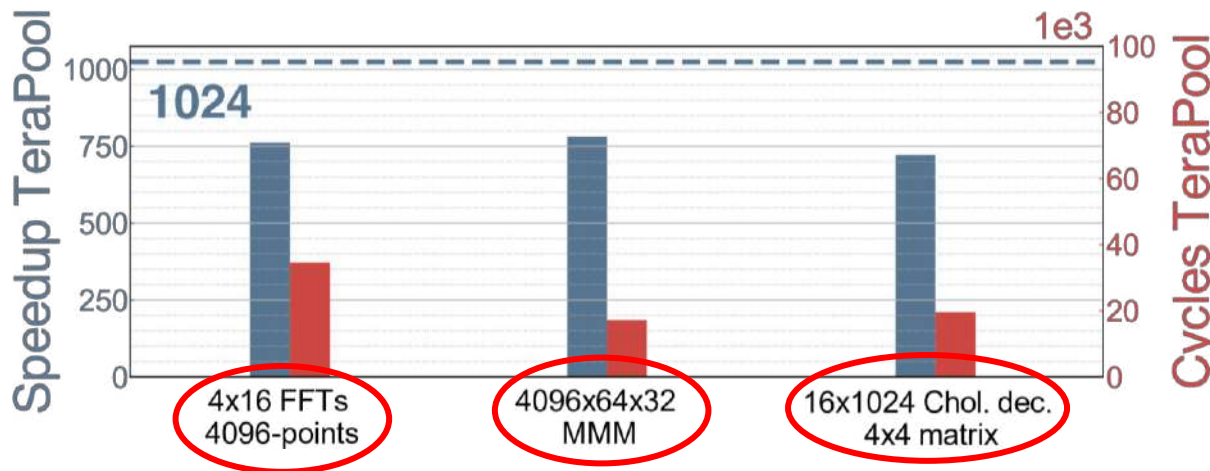
Cholesky

4x4 matrix
(16 independent
dec. Run between
barriers)



- TeraPool scales well compared to MemPool (overhead = synchronization)
- **LSU stalls** are reduced to **less than 10%** of the total execution time

Quasi-ideal speed-up and low latency



4096 subcarriers, 64 antennas, 32 beams, 4 UEs on the same subcarrier

Speedups → **762, 781, 722**

Runtime → **0.785ms @1GHz**

Further improvement from architecture specialization!



<https://arxiv.org/pdf/2210.09196.pdf>

People will ask questions...

Is this physically feasible?

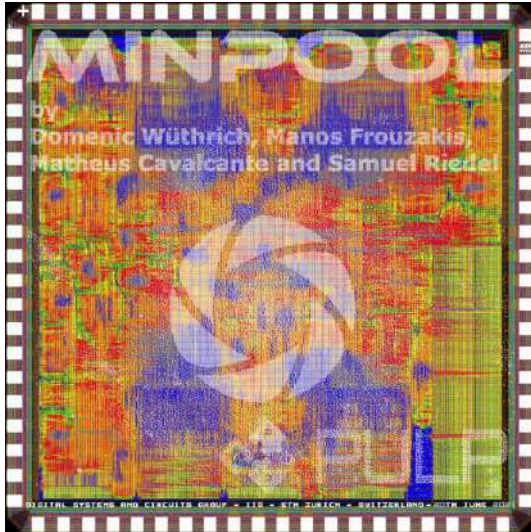
Does it place & route?

Was TeraPool ever taped-out?

Let's get Physical!

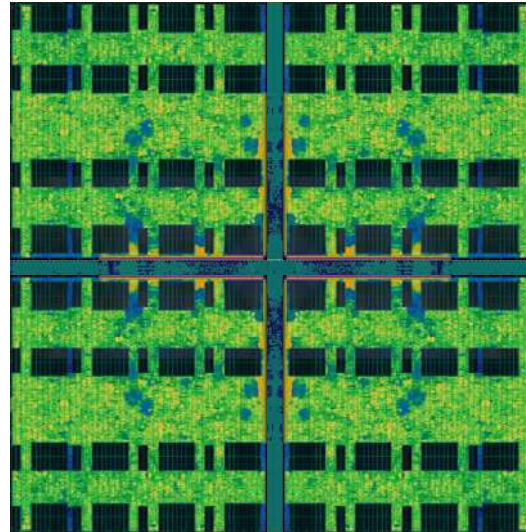


Hierarchical low-latency interconnect + Many Latency-Tolerant Cores (Snitch)



MinPool: first tape-out

- **16** cores, 64 KiB, **3** cycles
- TSMC 65



MemPool: main driver

- **256** cores, 1 MiB, **5** cycles
- GF 22FDX
 - 500 MHz (WC)
- **MemPool-3D**



TeraPool:

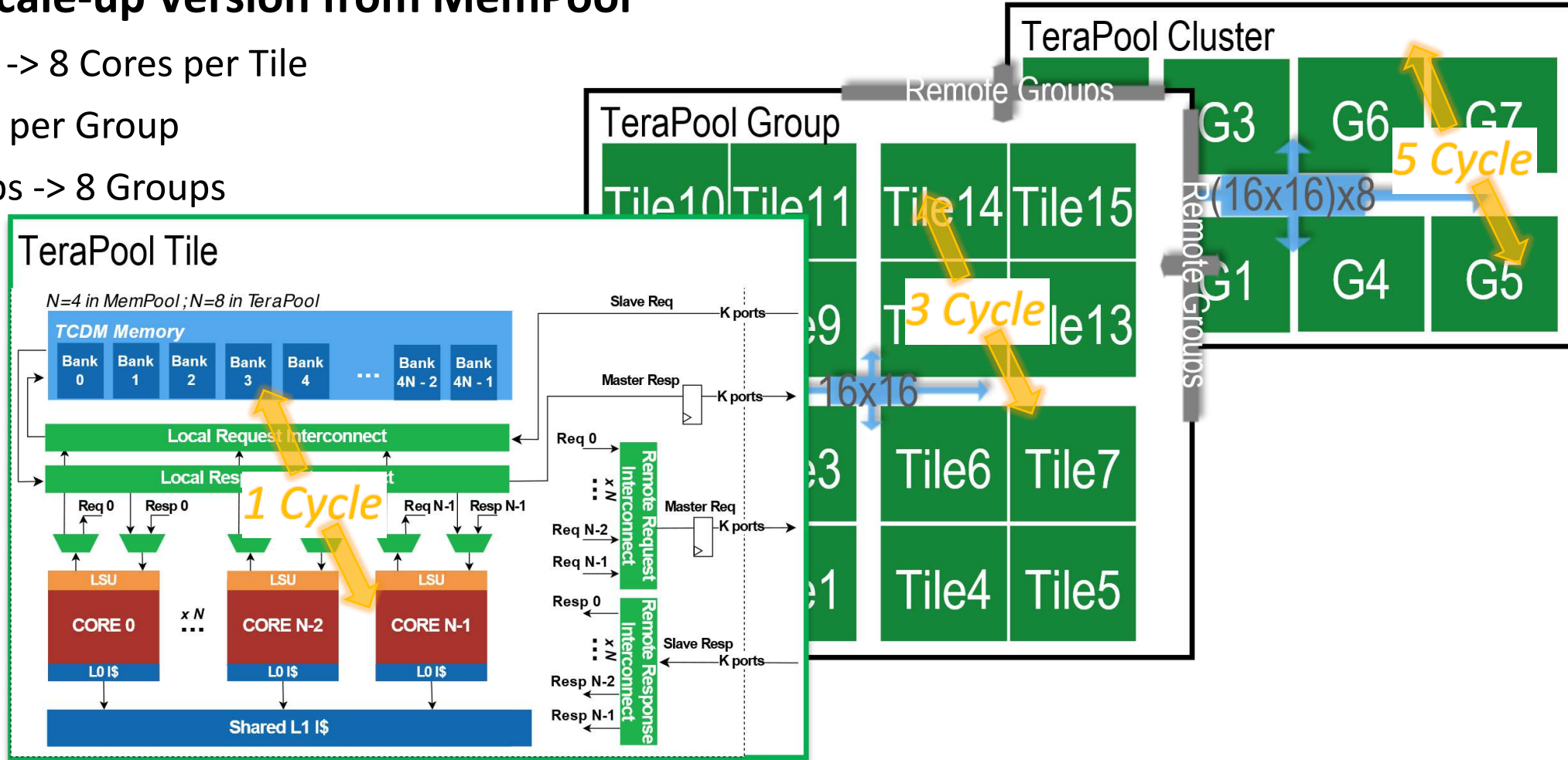
- **1024** cores, 4 MiB, ?cycles
- GF 12 LP+
- **How connecting?**
- **How go Physically?**

Scale Up From MemPool? - TeraPool₁₋₃₋₅ Design

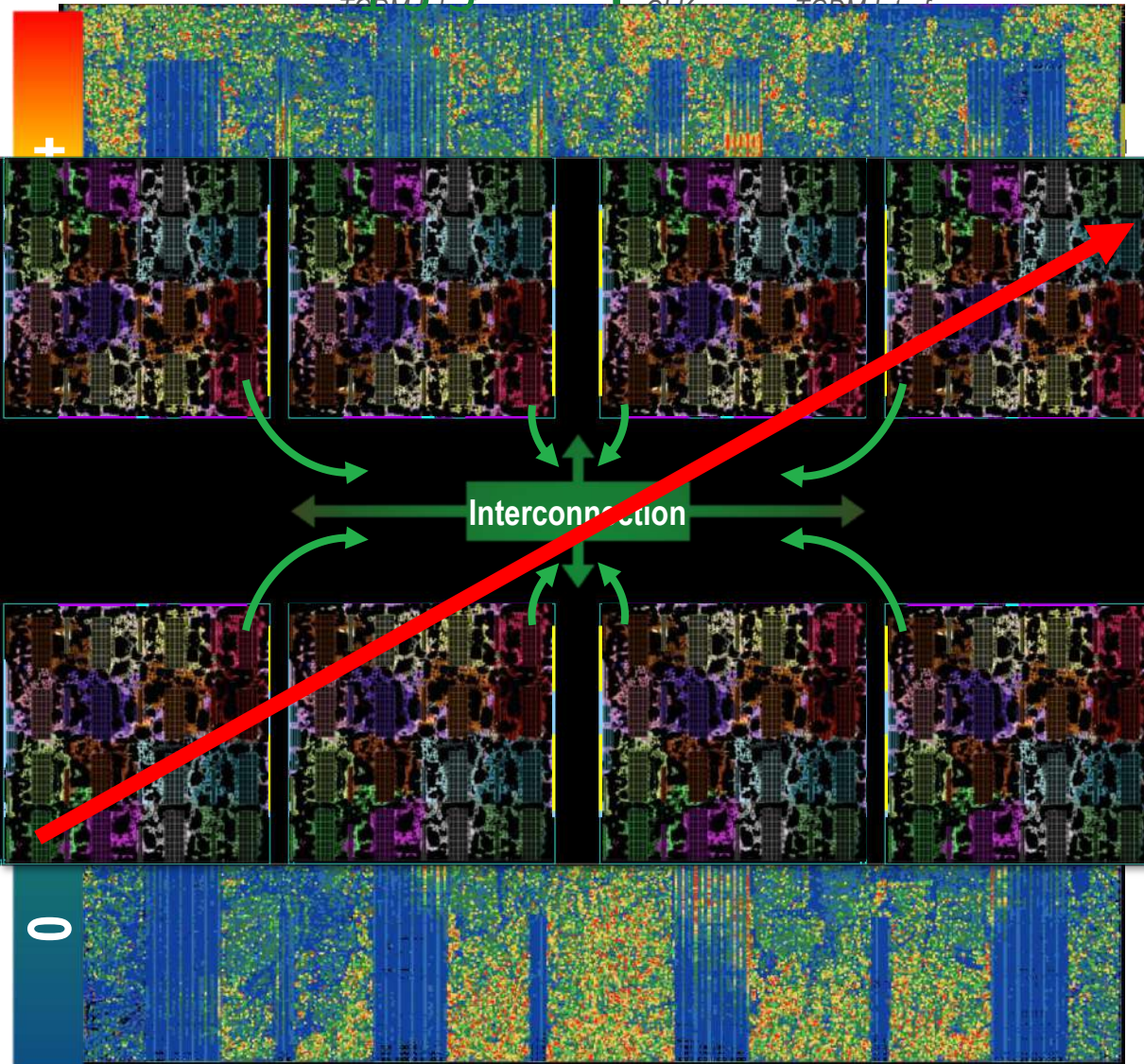


- **Direct Scale-up Version from MemPool**

- 4 Cores -> 8 Cores per Tile
- 16 Tiles per Group
- 4 Groups -> 8 Groups



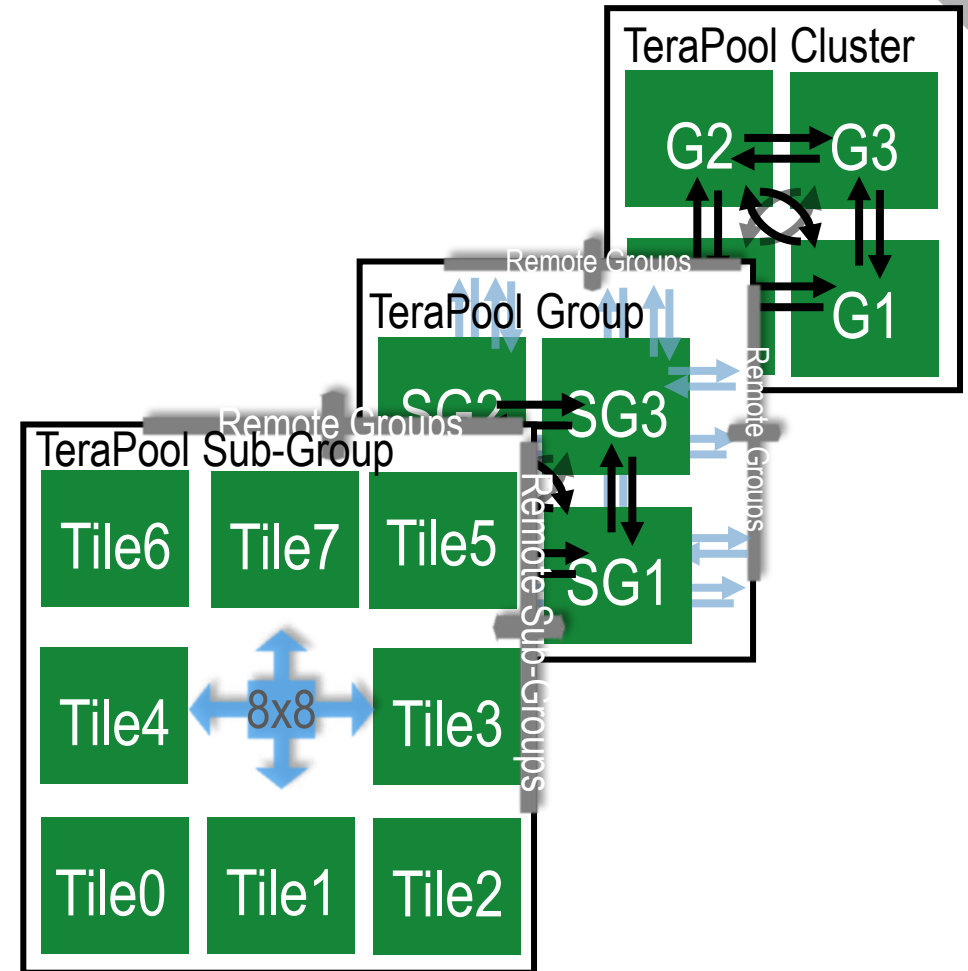
TeraPool_{1,3,5}: Implementation Challenges



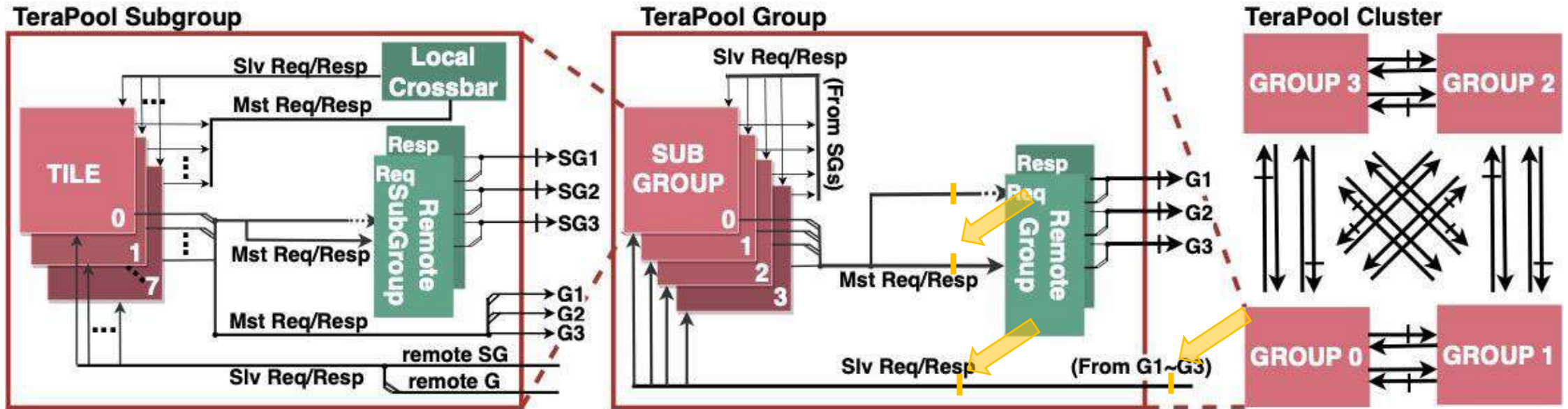
- **1.6ns+** critical path
 - The timing could not close at 500MHz, TT
- Design runtime is unmanageable
 - 1.5 Weeks for Group Level Design
- **Not routable**
 - ~33K DRCs
- Diagonal groups timing;
- Large routing channel required;

Brainstorming: what do we need?

- Short diagonal timing paths
 - Let's Keep 4 Groups per Cluster
- 256 Cores per Group, that's too large
 - New hierarchy: **Sub-Group** level
- Iteration runtime control
 - 8 Tiles per Subgroup.
- Frequency? or Latency?
 - Flexibility to add spill registers



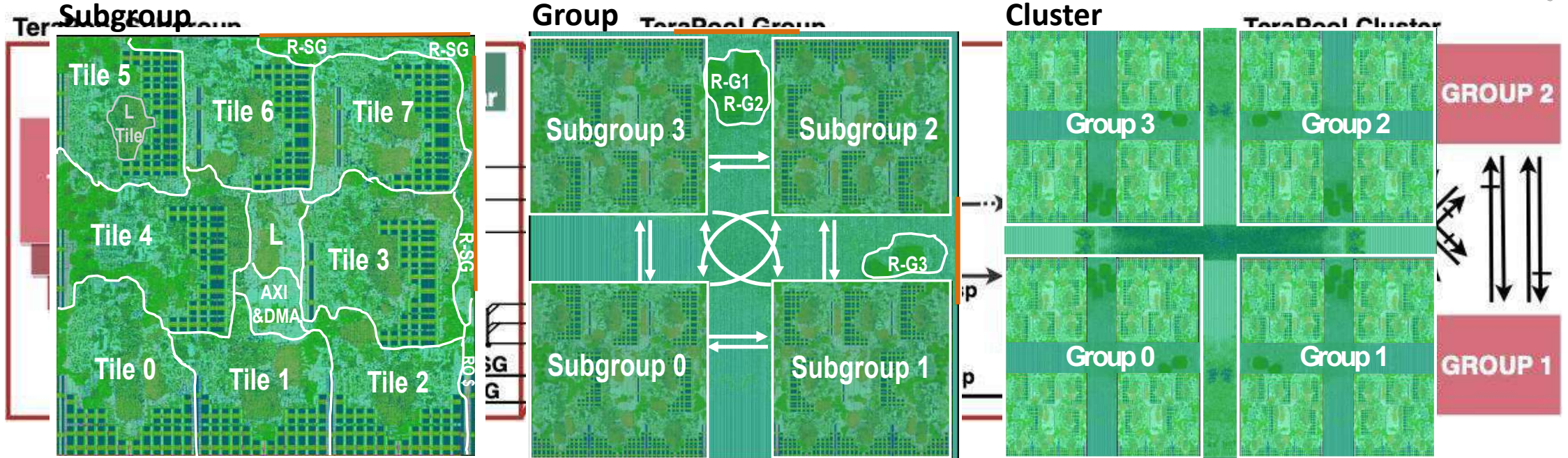
TeraPool_{1-3-5-X}: Hierarchical View



- Hierarchical Architecture View:
 - 4 Groups per Cluster
 - 4 Sub-Groups per Group
 - 8 Tiles per Sub-Group
 - 8 Cores per Tile

- Flexible spill registers adding:
 - Break long-distance remote accessing paths
 - Different latency/frequency targets
 - Hardware-configurable X= 7/9/11

TeraPool_{1-3-5-X}: Physically FEASIBLE!



- Methodology:

- GlobalFoundries' 12P+ FinFET
- Synopsys' FusionCompiler 2022.03
- Synopsys' PrimeTime 2022.03
- WC: SS/0.72V/125C ;TT: TT/0.80V/25C

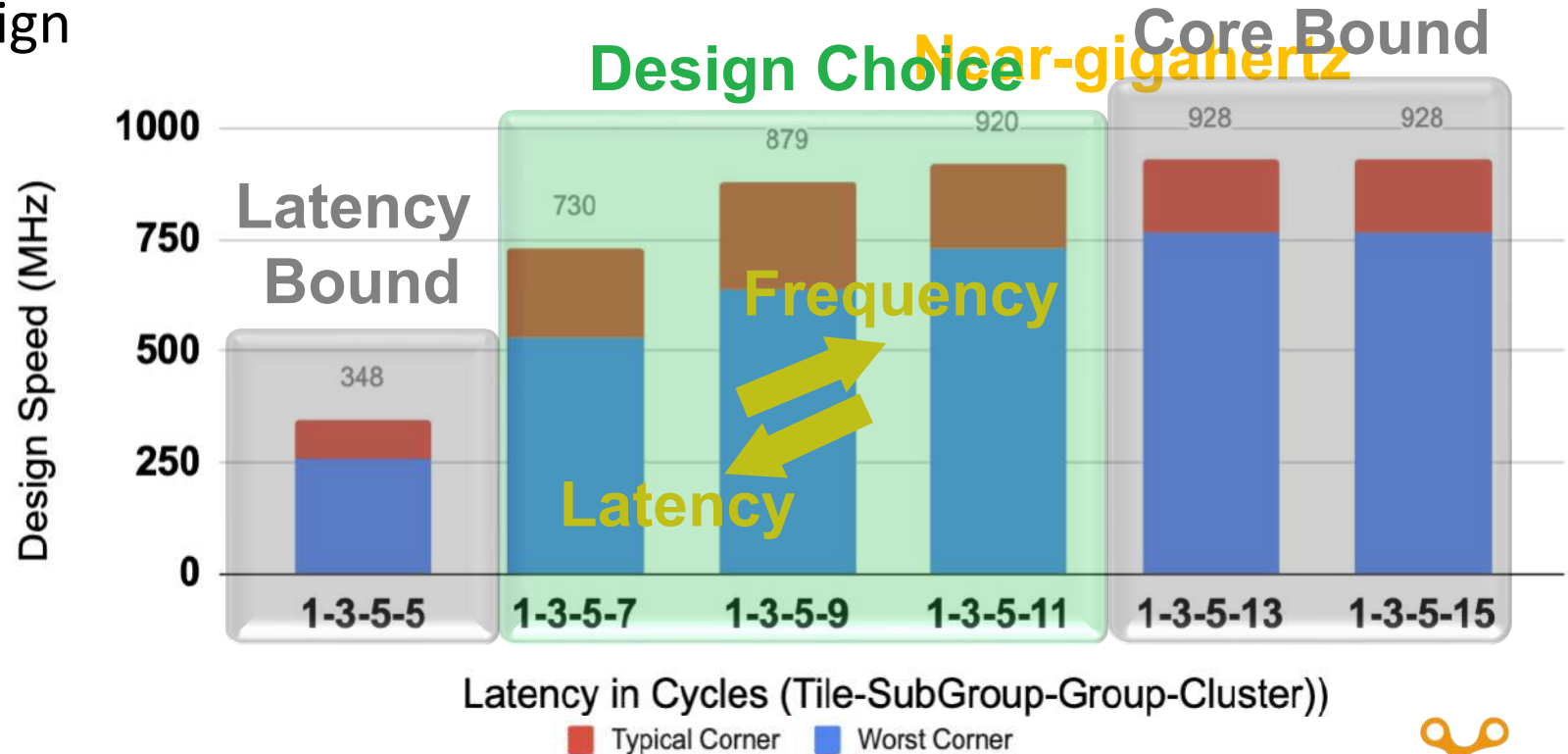
- Implement Area:

- Subgroup: 1.52 x 1.52 mm² (58% utilization)
- Group: 3.8 x 3.8 mm²
- Cluster: 8.3 x 8.3mm²

Latency vs. Operating Frequency



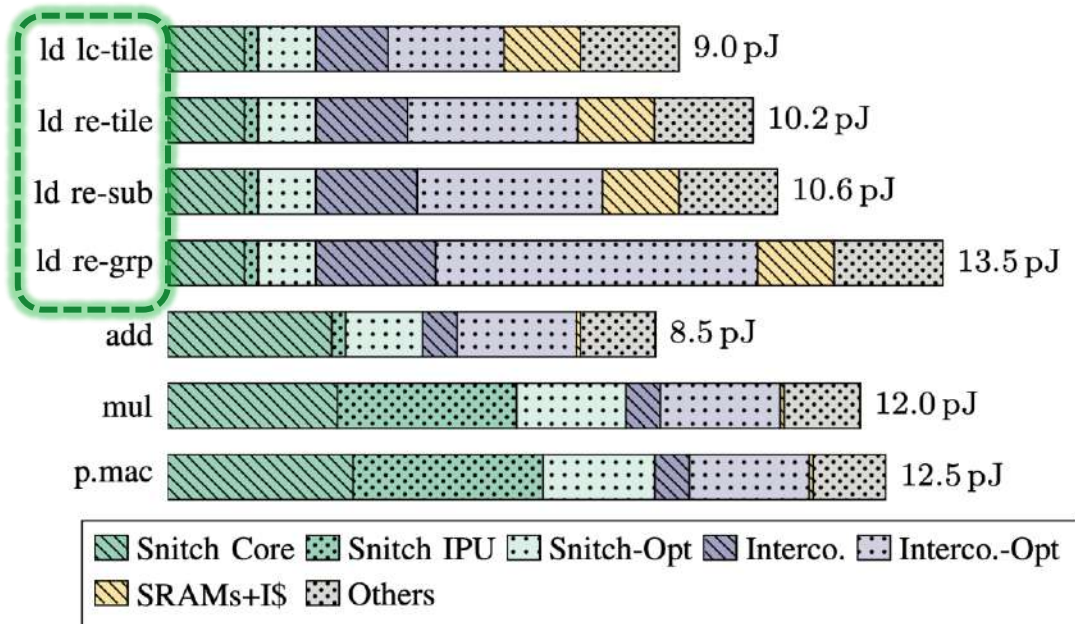
- Registers spill critical paths
- Latency-Frequency trade-off
- Near-GigaHertz Design



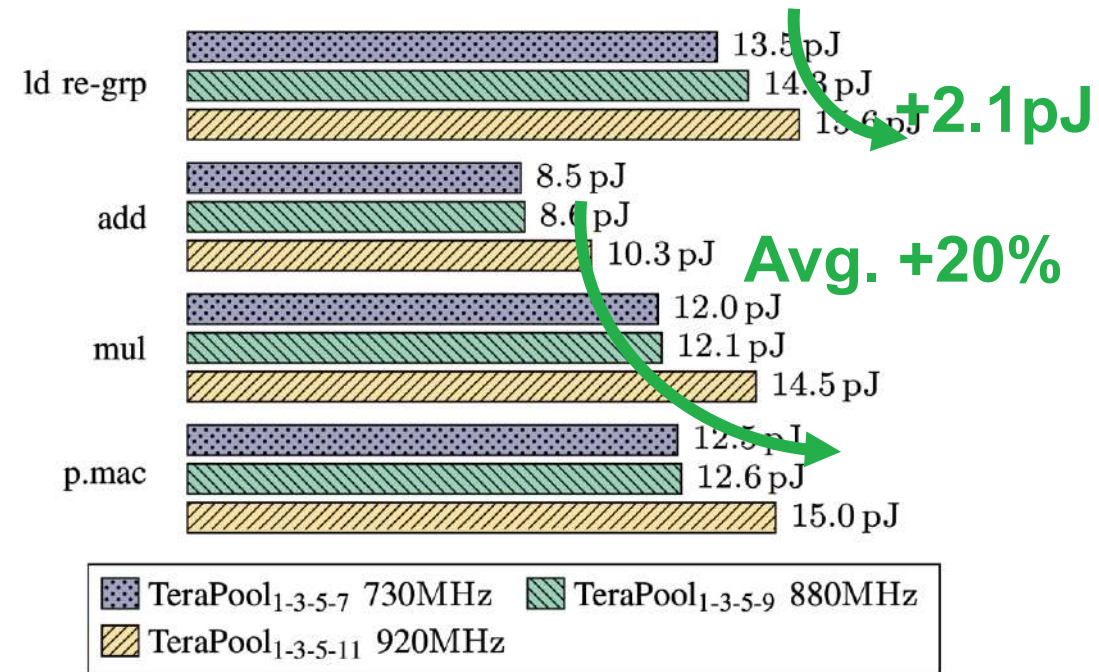
Energy Breakdown



- Energy-Efficiency:
 - Only 9-13.5pJ for whole L1 access
 - Even p.mac, only 12.5pJ



- With target frequency increases:
 - More design optimization cells add on
 - With larger driveability, SLVT → Power suffers



Beat MemPool, go go go



TeraPool is NOT MemPoolx4

- Larger Scale Cluster
 - Reduce system overhead
- Less data movement
 - Memory-Cluster
 - Cluster-Cluster
- Less workload distribution
 - Data allocation
 - Data Splitting

Core Number	256	1024		
Tech Node	22nm	12nm		
Die Area (mm ²)	12.9mm ²	68.9mm ²		
Hier-Latency (cycles)	1-3-5	1-3-5-7	1-3-5-9	1-3-5-11
Avg-Latency (zero load)(cycles)	4.7	6.4	7.9	9.3
Throughput (req/core/cycle)	0.33	0.23	0.24	0.25
Frequency (Typ. Condition)	587MHz	730MHz	880MHz	920MHz
Peak Performance (TOPS)	0.3	1.50	1.80	1.89
Area Efficiency (GPOS/MM ²)	23.3	21.8	26.1	27.4
Benchmark	Example: MatMul Kernel ¹			
Kernel Performace (TOPS)	0.17	0.60	0.66	0.75
Power Consumption (W)	2.25	4.16	4.87	6.78
Energy Efficiency (W)	136	144	136	111

More Efficiency

More Performance

What's next?

Much more will come...

More Configurations:

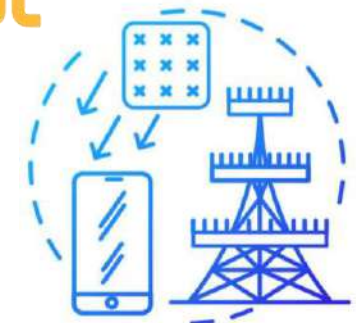
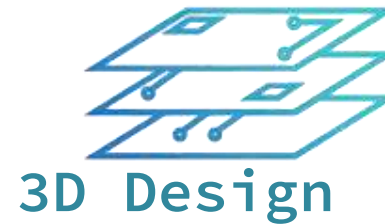
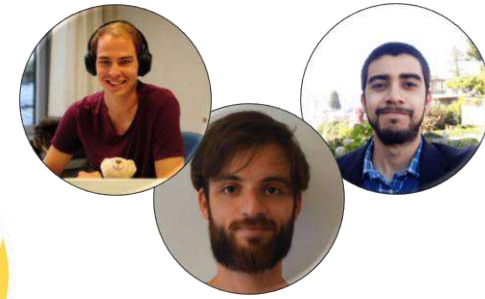
- 4 Cores/Tile -> more throughput
- 16 Cores/Tile -> less average latency

More Recipes:

- FPU implementation: 8bits, 16bits, even 32bits
- Spatz-TeraPool: Vector unit powered, more ILP

More Clusters:

- Multi-Cluster + Memory hierarchy
- **3D Implementation**
- **Full MIMO MMSE Receiving Chain**



github.com/pulp-platform/mempool

Marco Bertuletti
Yichao Zhang

mbertuletti@iis.ee.ethz.ch
yiczhang@iis.ee.ethz.ch



Institut für Integrierte Systeme – ETH Zürich
Gloriastrasse 35
Zürich, Switzerland

DEI – Università di Bologna
Viale del Risorgimento 2
Bologna, Italy