



**TECHNISCHE
UNIVERSITÄT
DRESDEN**



The Hebrew University
of Jerusalem

FFMK: A FAST AND FAULT-TOLERANT MICROKERNEL- BASED SYSTEM FOR EXASCALE COMPUTING

Amnon Barak

Hebrew University Jerusalem (HUJI)

Hermann Härtig

TU Dresden, Operating Systems Group (TUDOS)

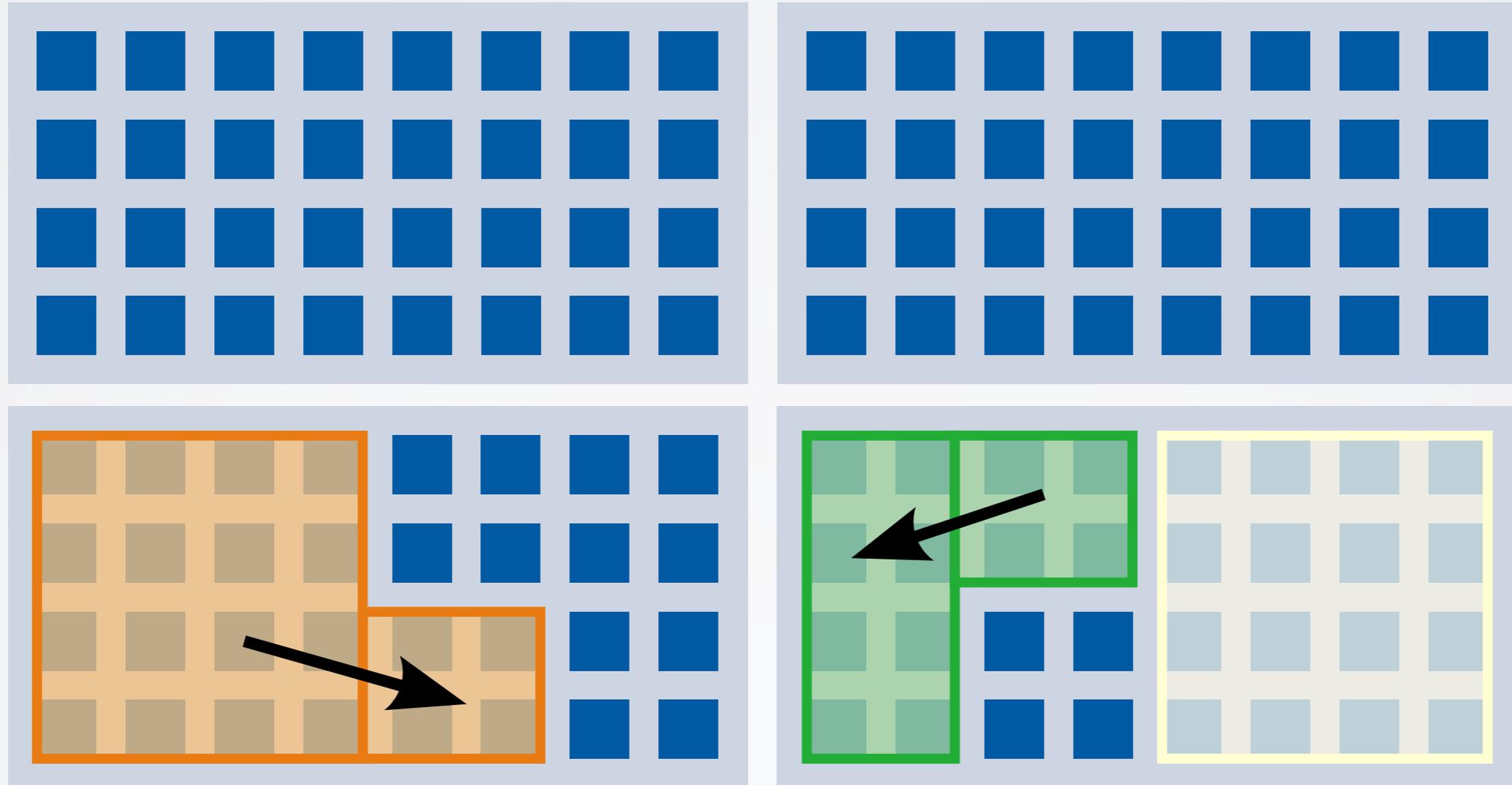
Wolfgang E. Nagel

TU Dresden, Center for Information Services and HPC (ZIH)

Alexander Reinefeld

Konrad-Zuse-Zentrum für Informationstechnik Berlin (ZIB)

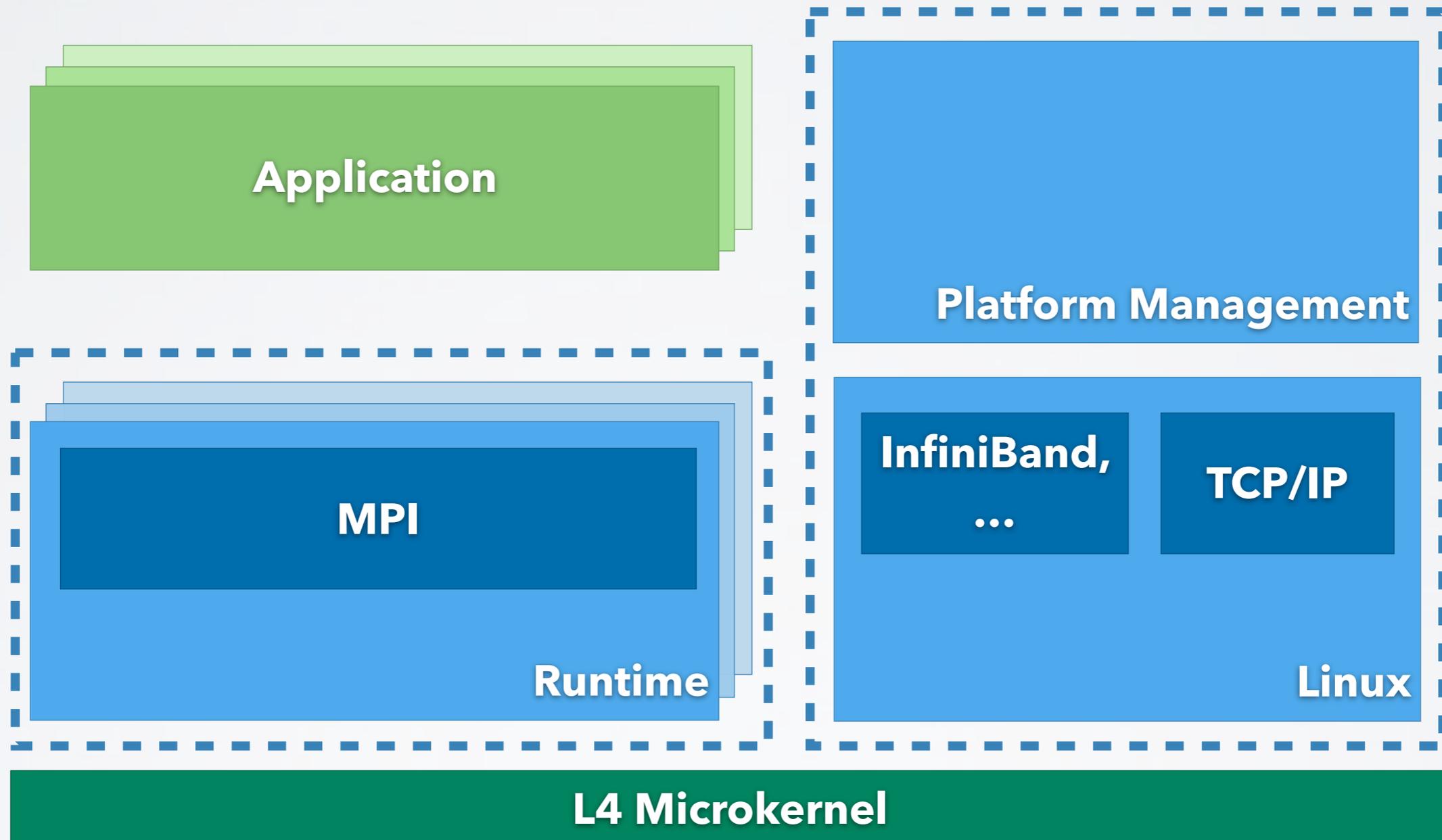
CARSTEN WEINHOLD, TU DRESDEN



M M M M M M M M

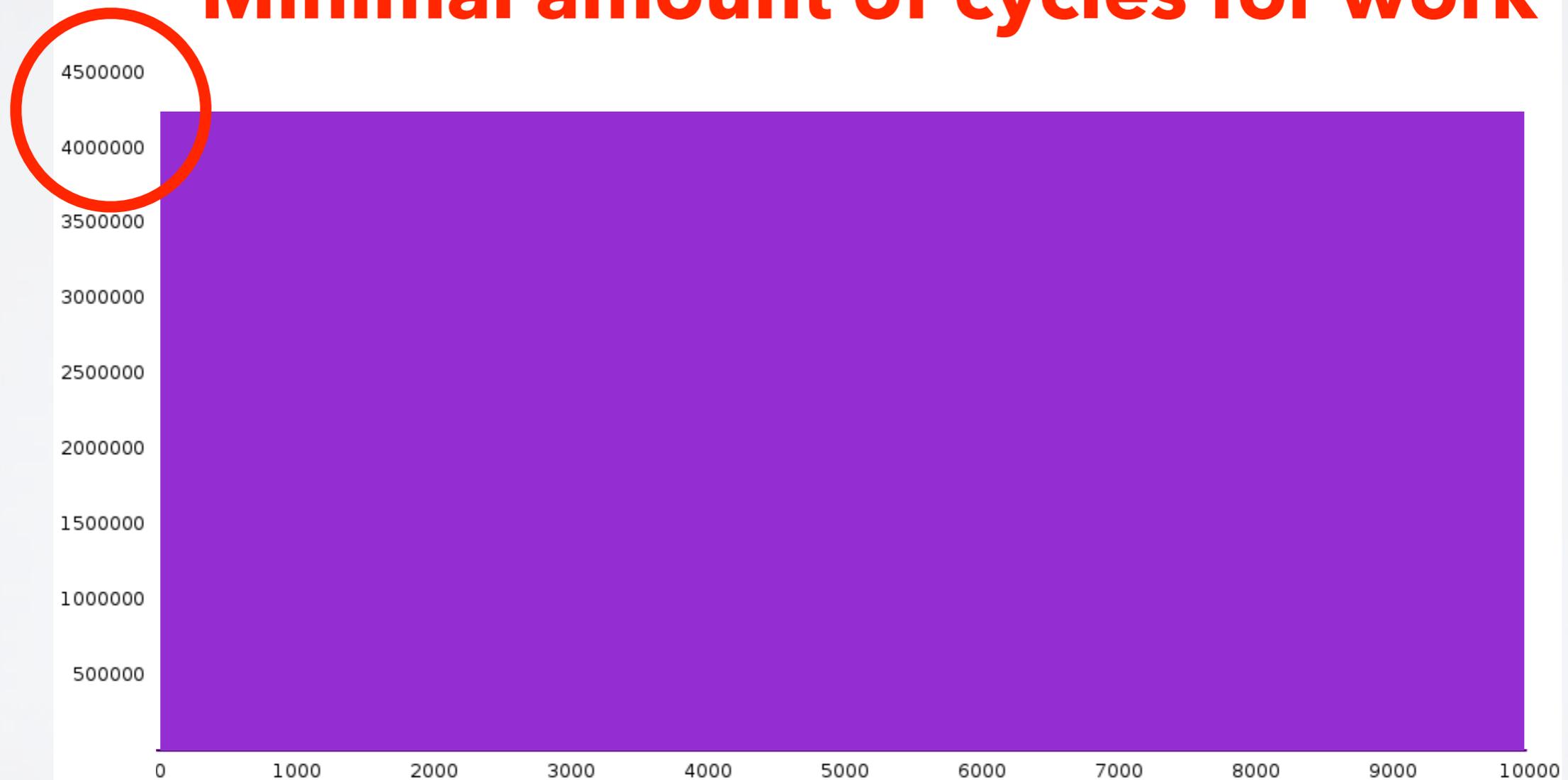


NODE ARCHITECTURE



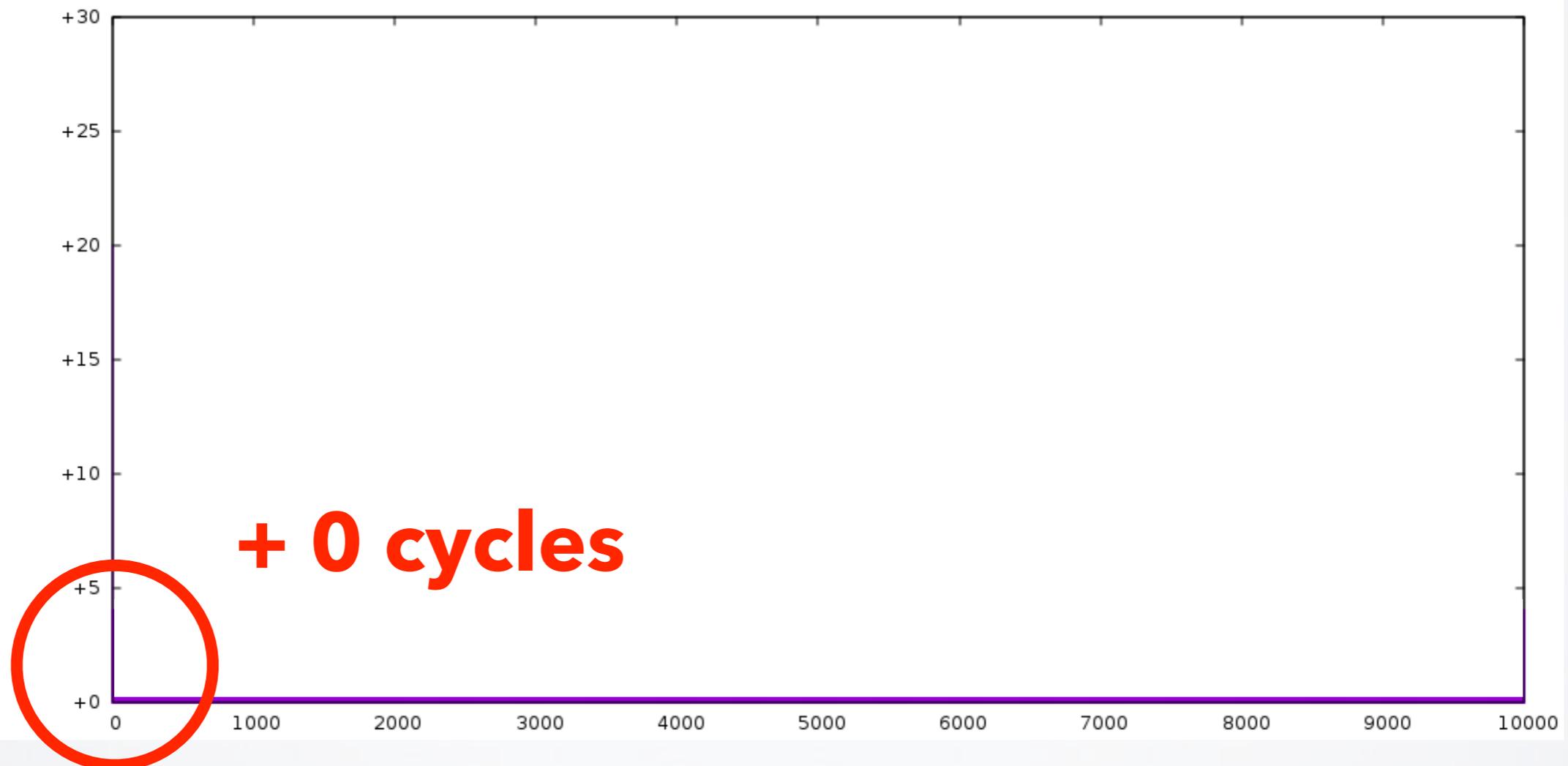


Minimal amount of cycles for work



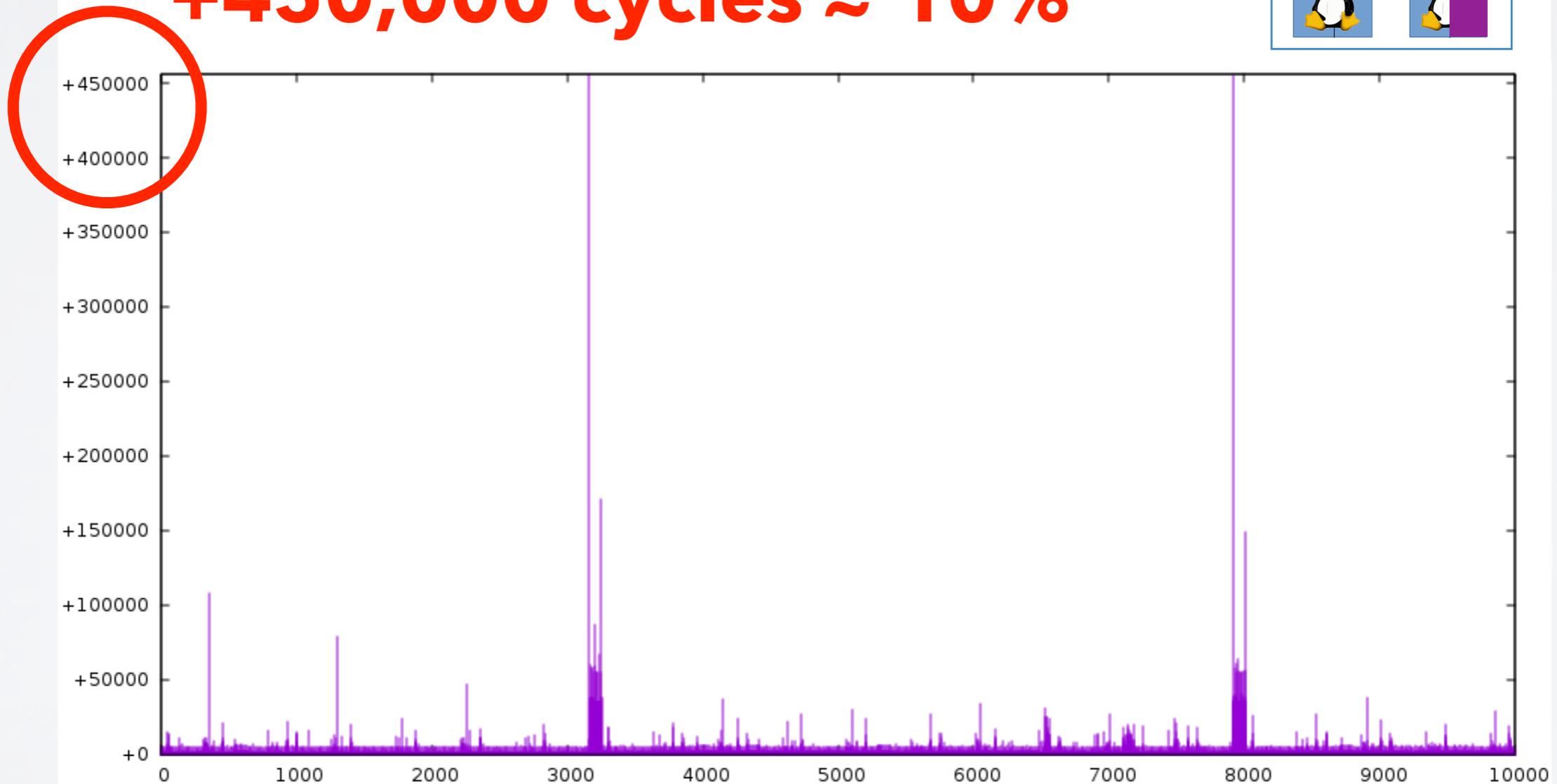
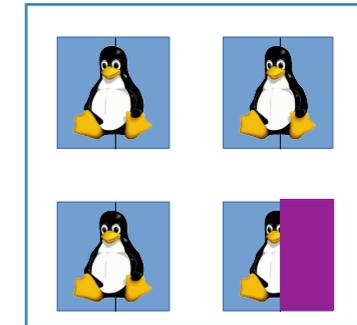


Ideal: zero extra cycles

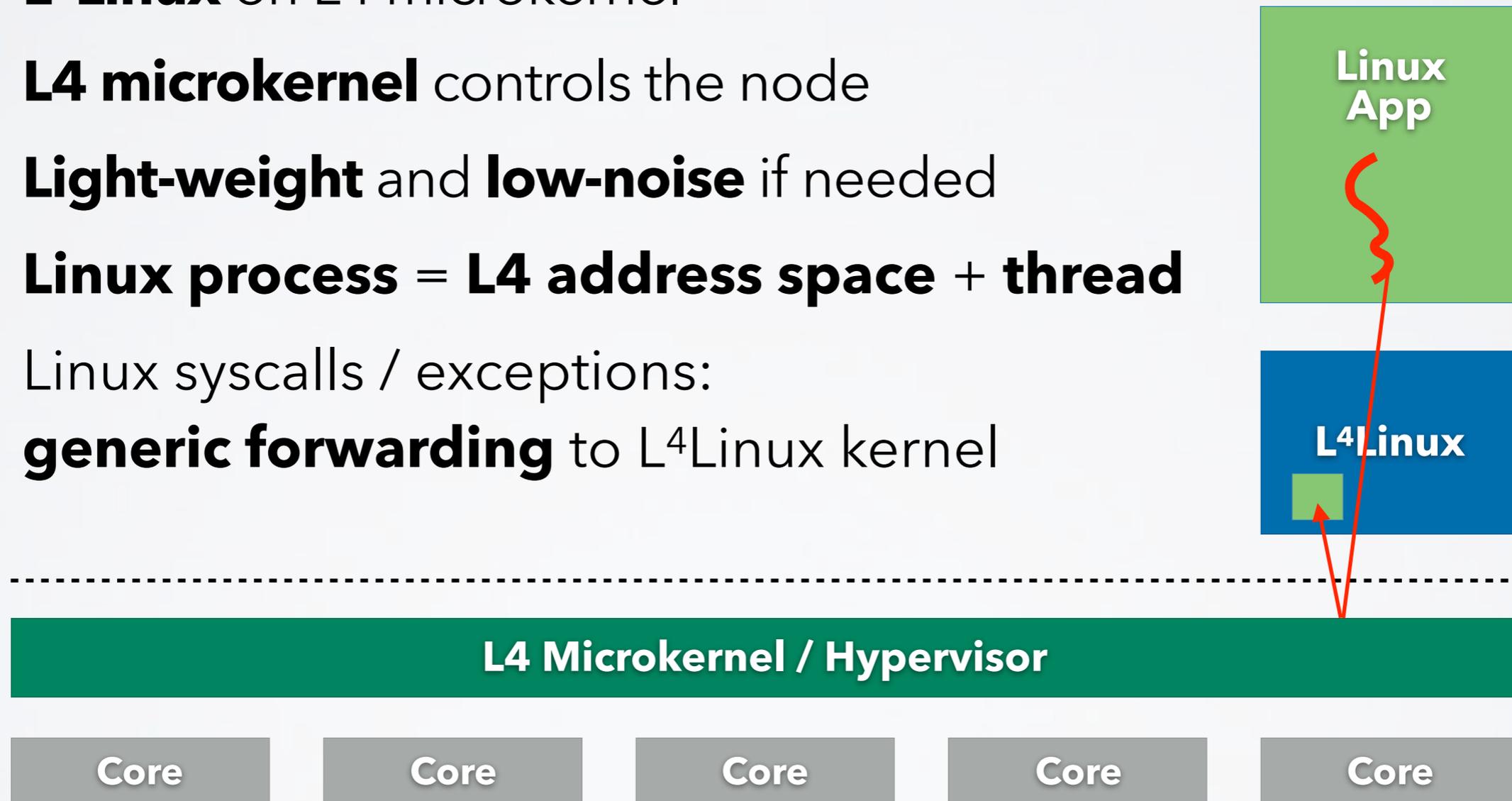


Real-World HPC Linux

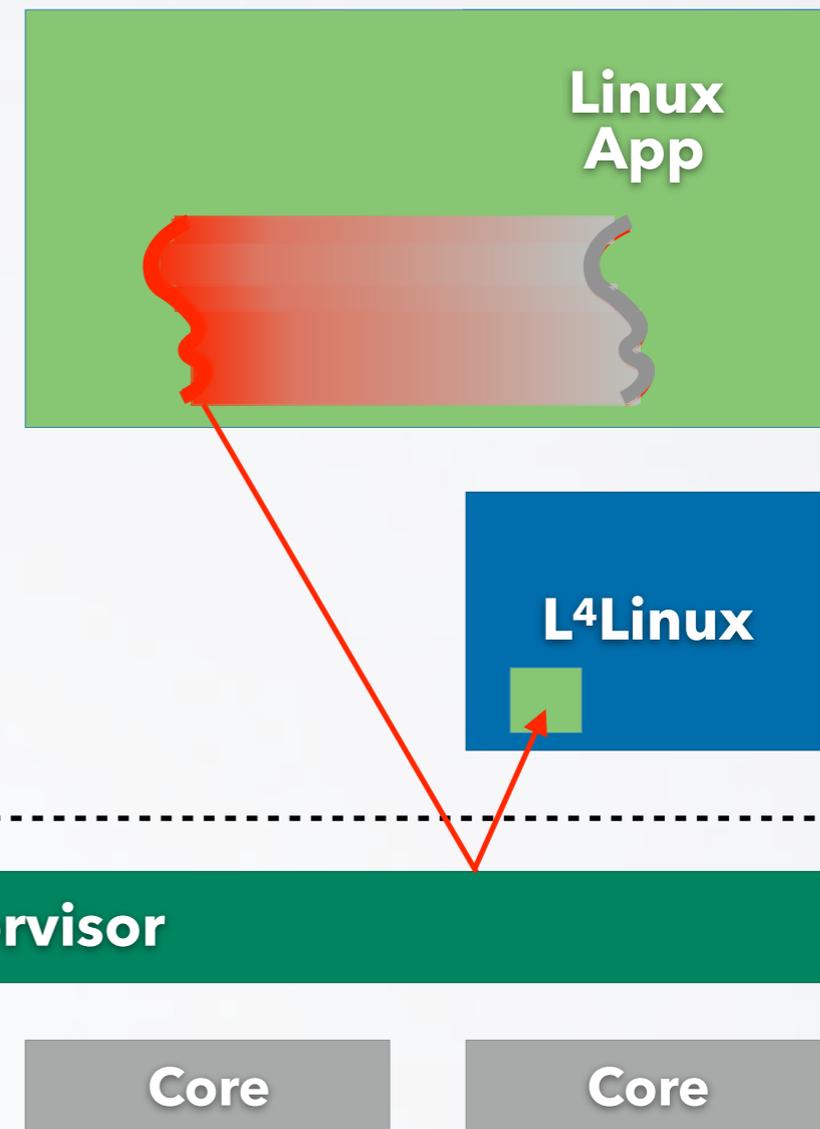
+450,000 cycles \approx 10%



- **Unmodified** Linux programs (MPI, ...)
- **L⁴Linux** on L4 microkernel
- **L4 microkernel** controls the node
- **Light-weight** and **low-noise** if needed
- **Linux process = L4 address space + thread**
- Linux syscalls / exceptions:
generic forwarding to L⁴Linux kernel



- **Decoupling:** move Linux thread to new L4 thread on its own core
- **Linux syscall:** Move back to Linux
- **Direct I/O** device access
- **L4 syscalls:**
 - Memory
 - Threads & Scheduling
 - Interrupts



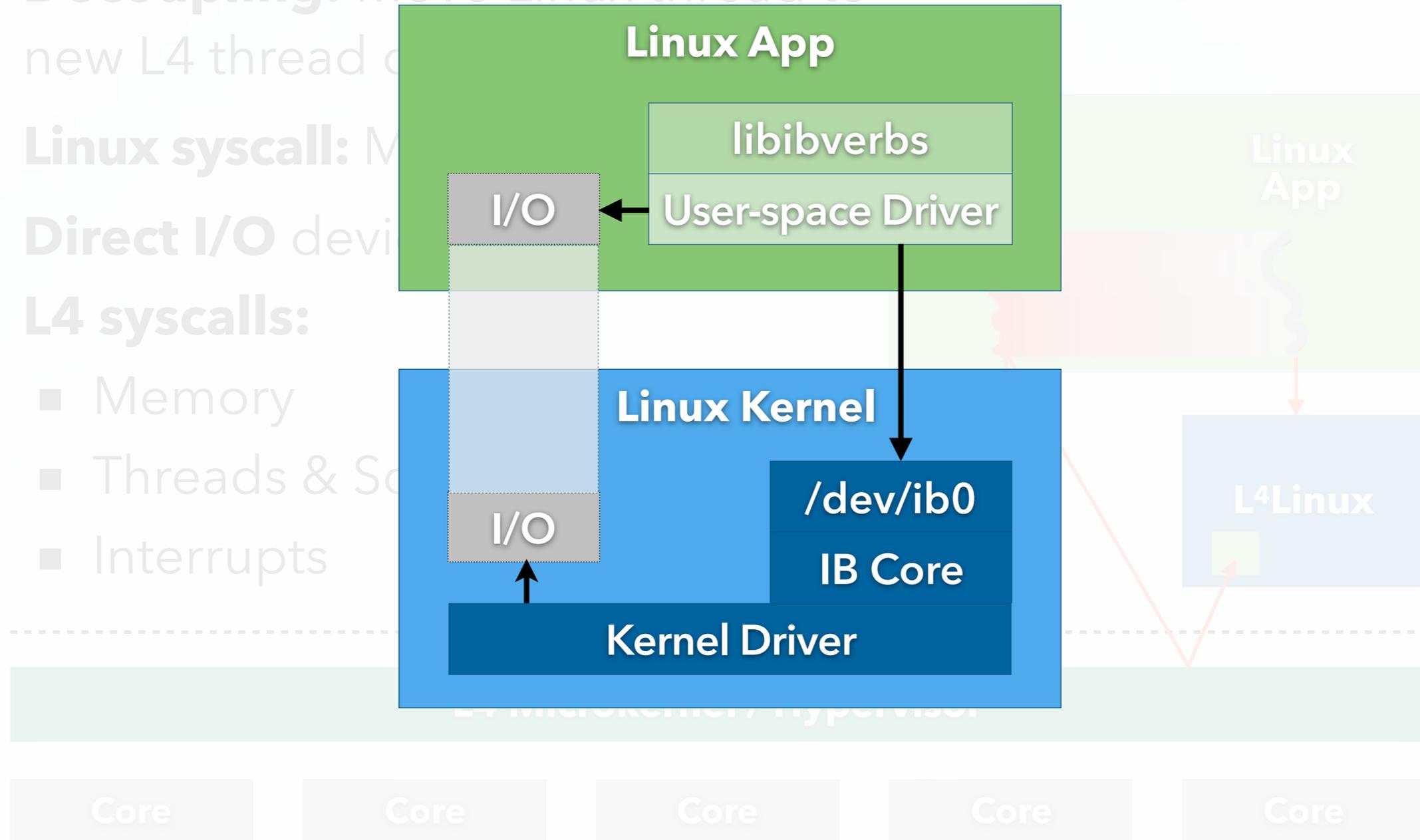
- **Decoupling:** move Linux thread to new L4 thread

- **Linux syscall:** M

- **Direct I/O** device

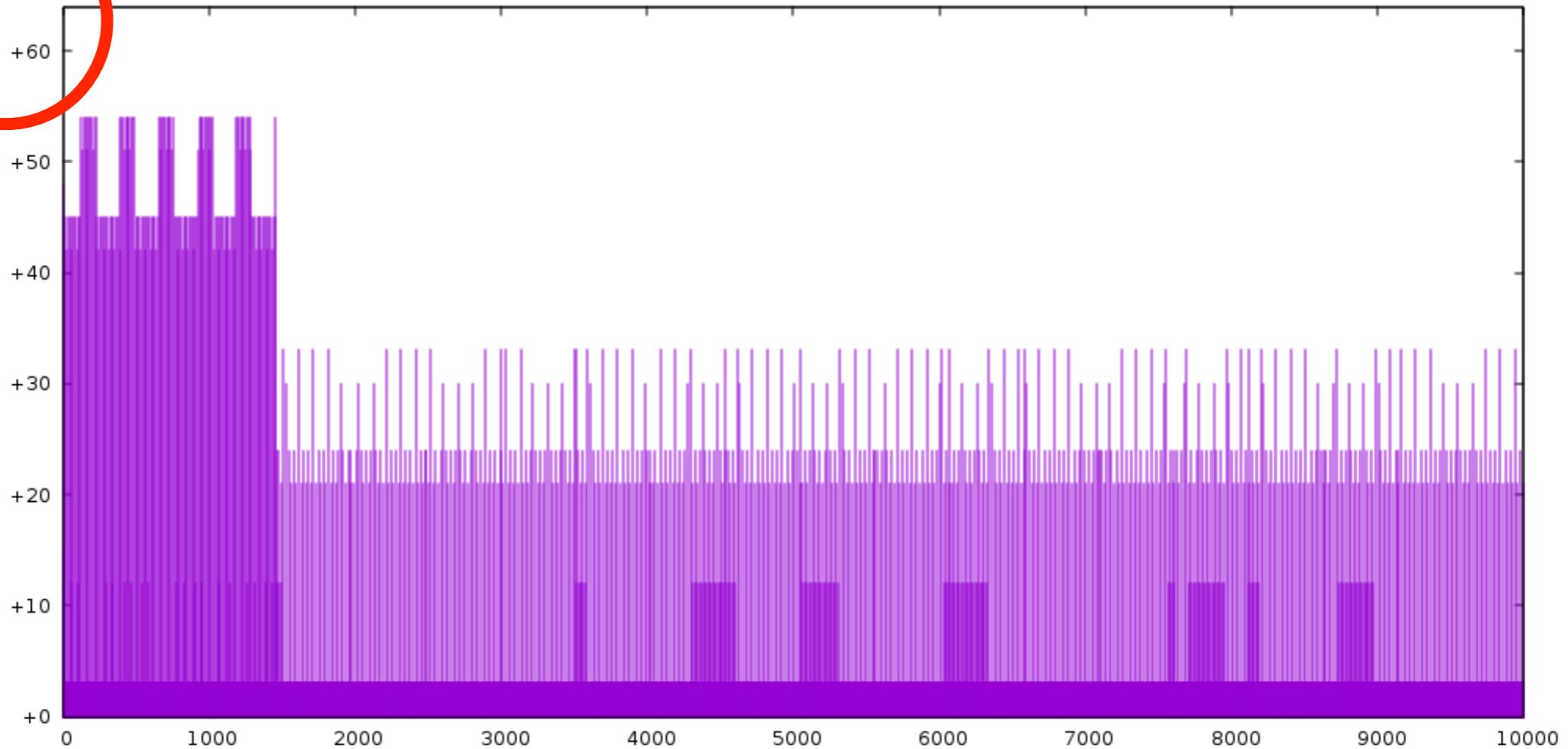
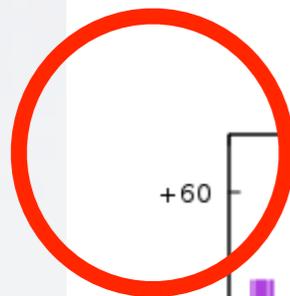
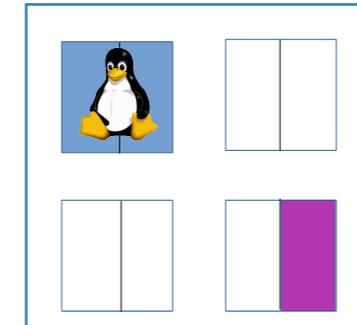
- **L4 syscalls:**

- Memory
- Threads & S
- Interrupts

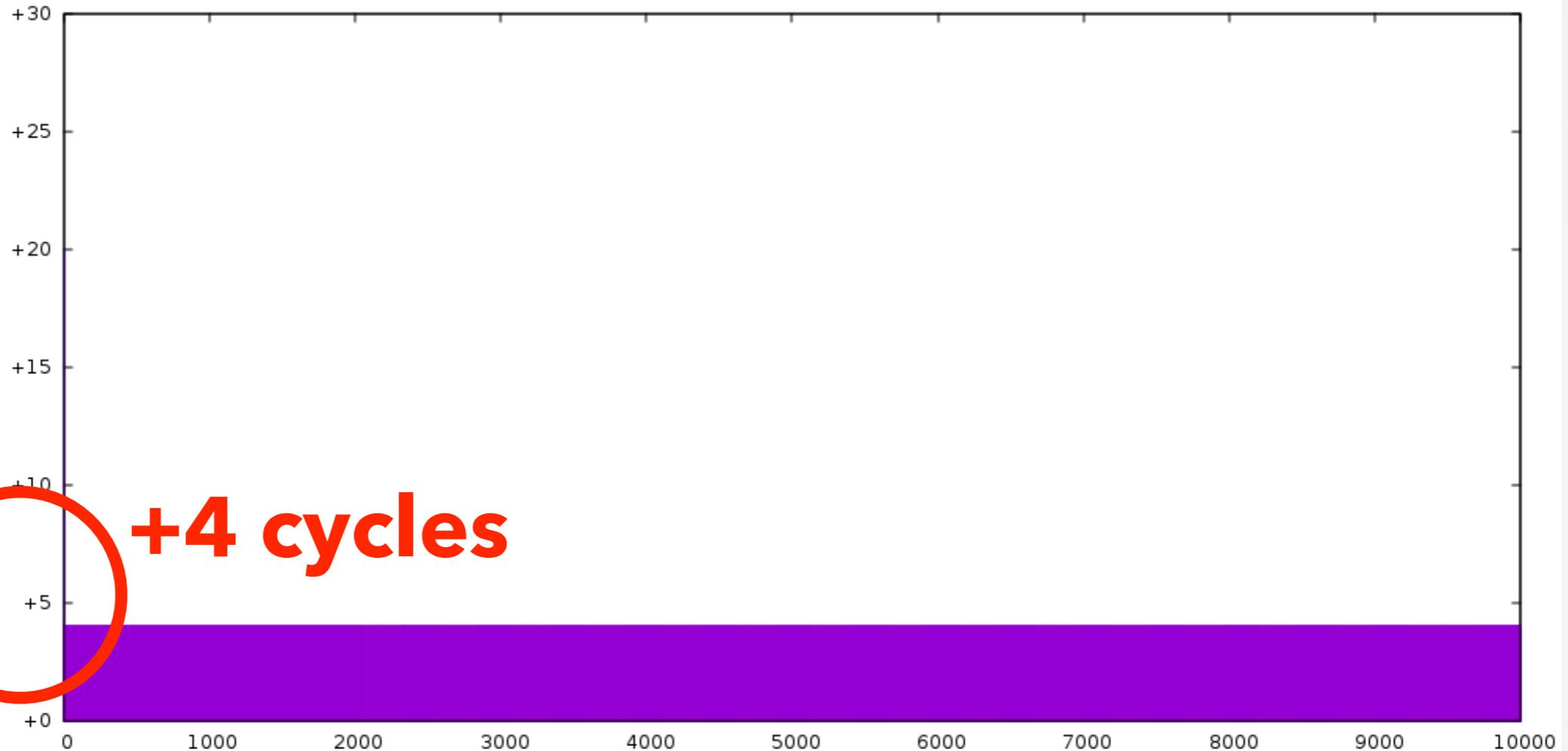
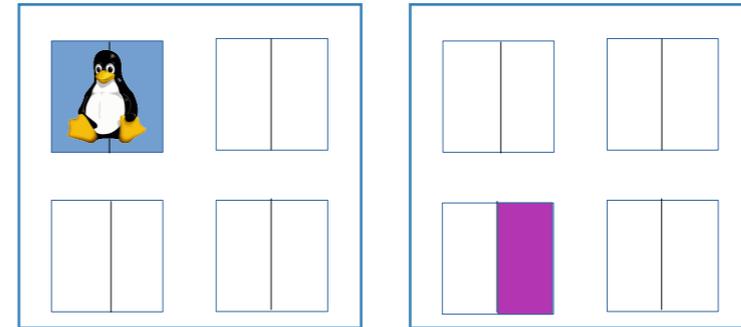


Decoupled Linux thread

+60 cycles



Decoupled Linux thread

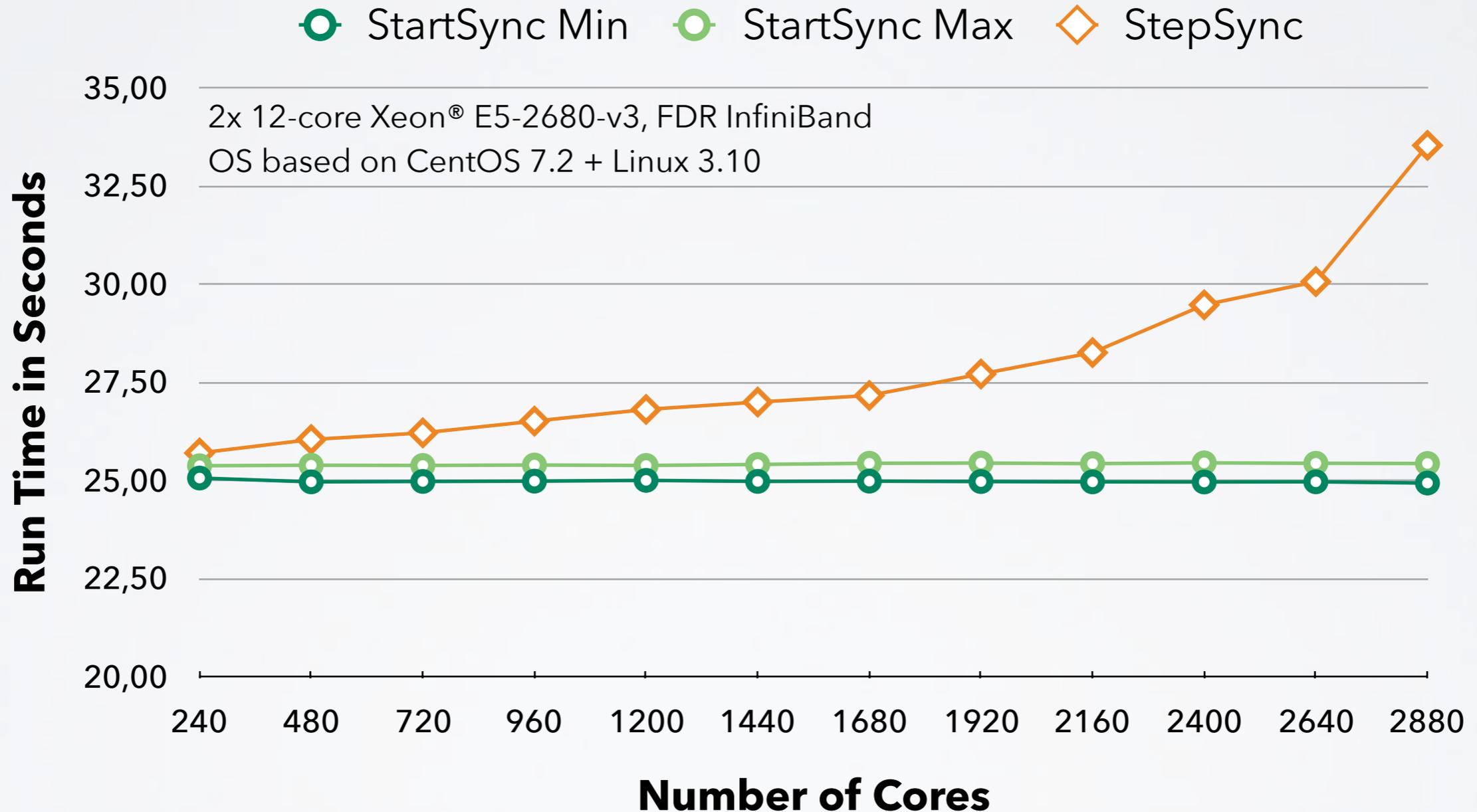




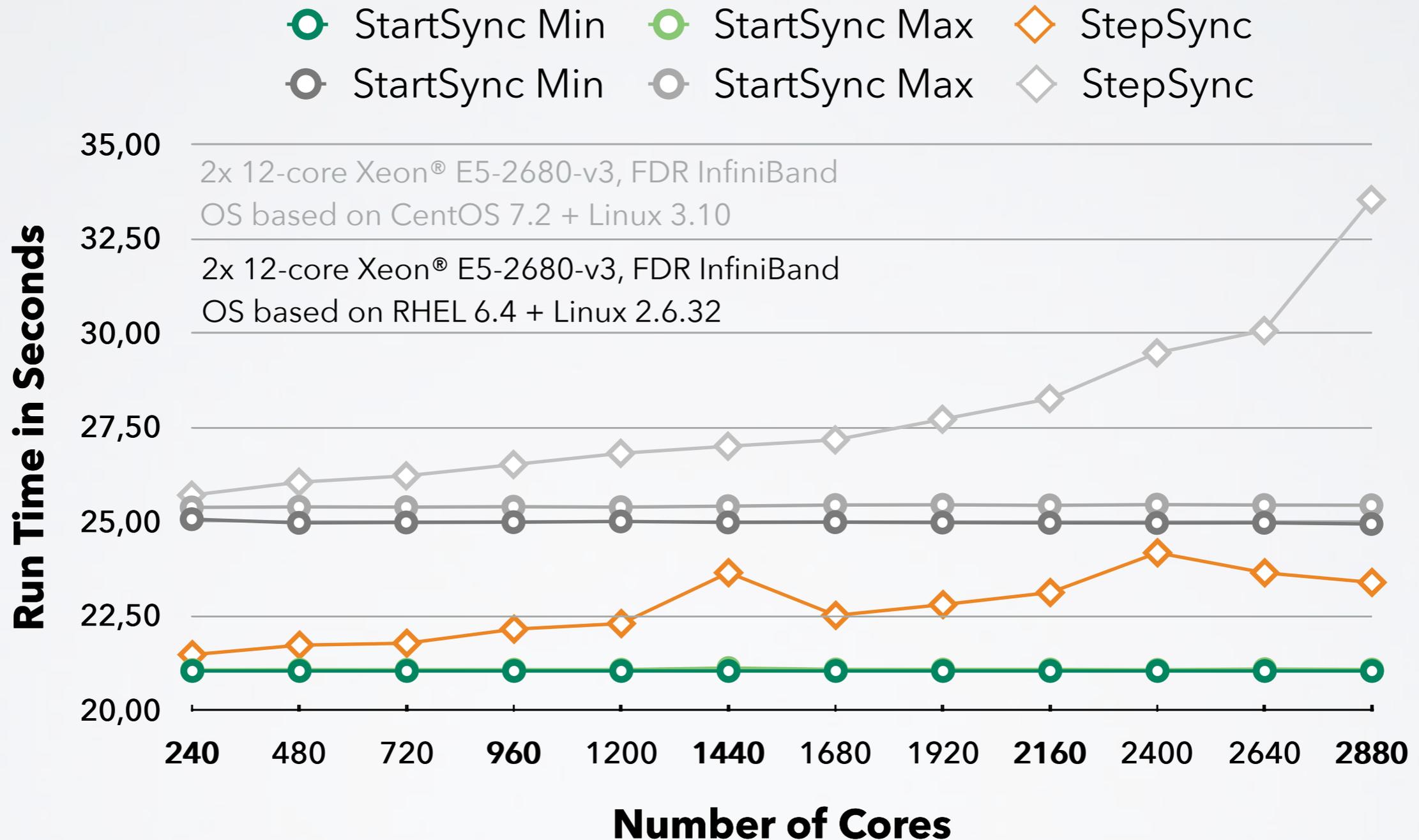
Behavior: **embarrassingly parallel**



Behavior: **bulk-synchronous**

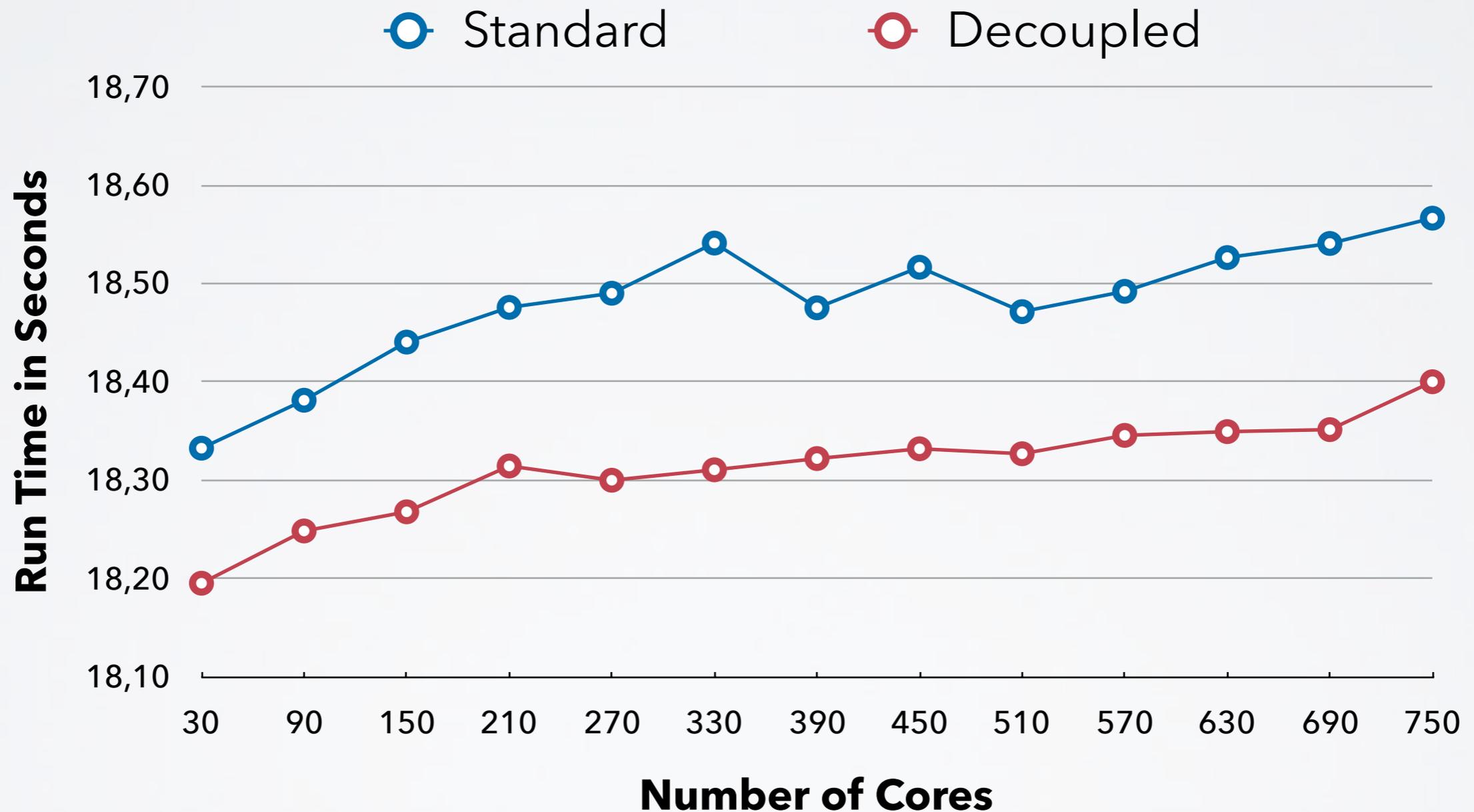


Adam Lackorzynski, Carsten Weinhold, Hermann Härtig, "Decoupled: Low-Effort Noise-Free Execution on Commodity Systems", ROSS 2016, June 2016, Kyoto, Japan

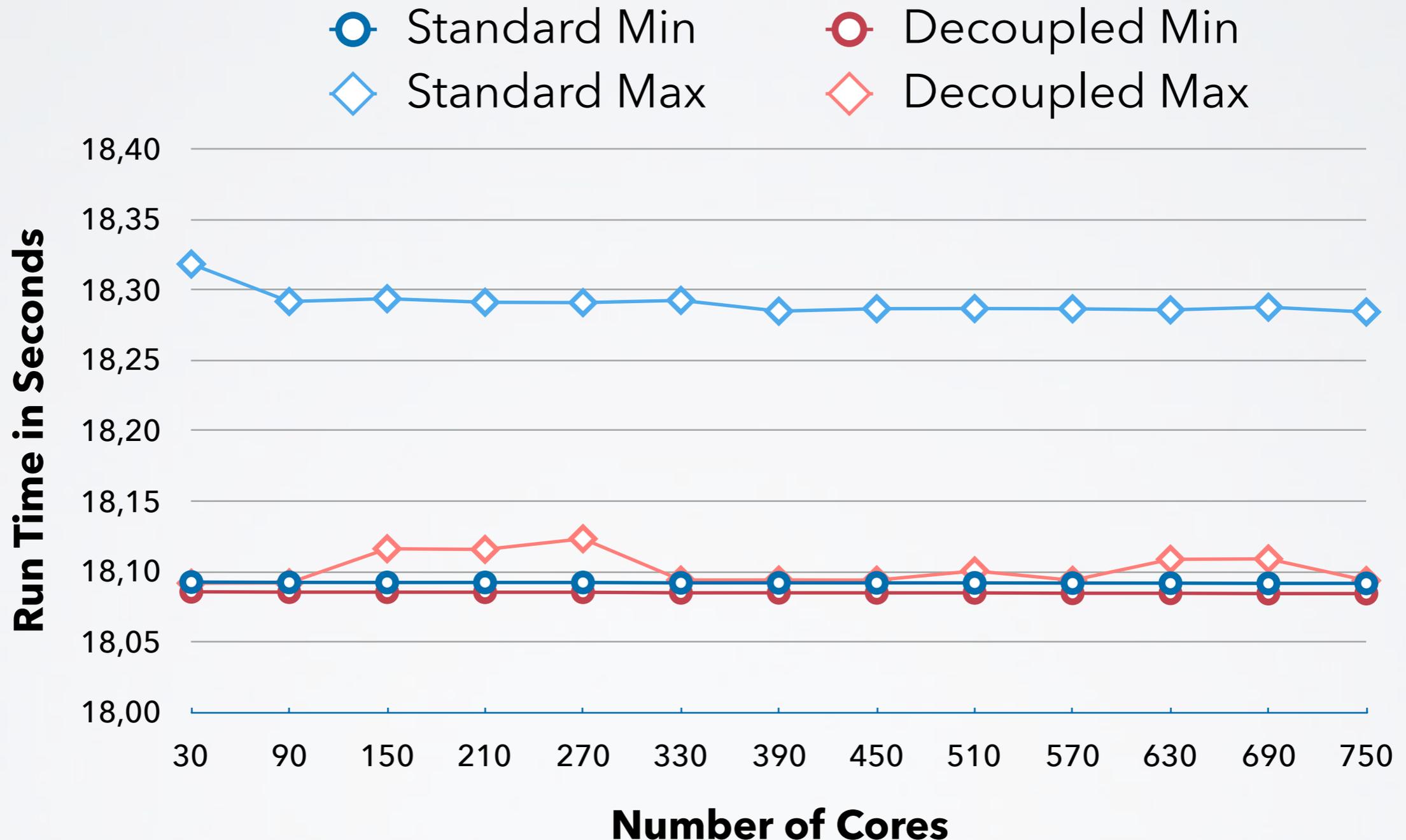


Adam Lackorzynski, Carsten Weinhold, Hermann Härtig, "Decoupled: Low-Effort Noise-Free Execution on Commodity Systems", ROSS 2016, June 2016, Kyoto, Japan

- **Bare-metal access** to Taurus:
 - Little time
 - Fewer cores
 - Different type of nodes
- Vendor OS: **Linux 2.6.32** or **3.10** ...
- Decoupled threads: **L⁴Linux 4.4**
- Custom Linux distribution



Adam Lackorzynski, Carsten Weinhold, Hermann Härtig, "Decoupled: Low-Effort Noise-Free Execution on Commodity Systems", ROSS 2016, June 2016, Kyoto, Japan



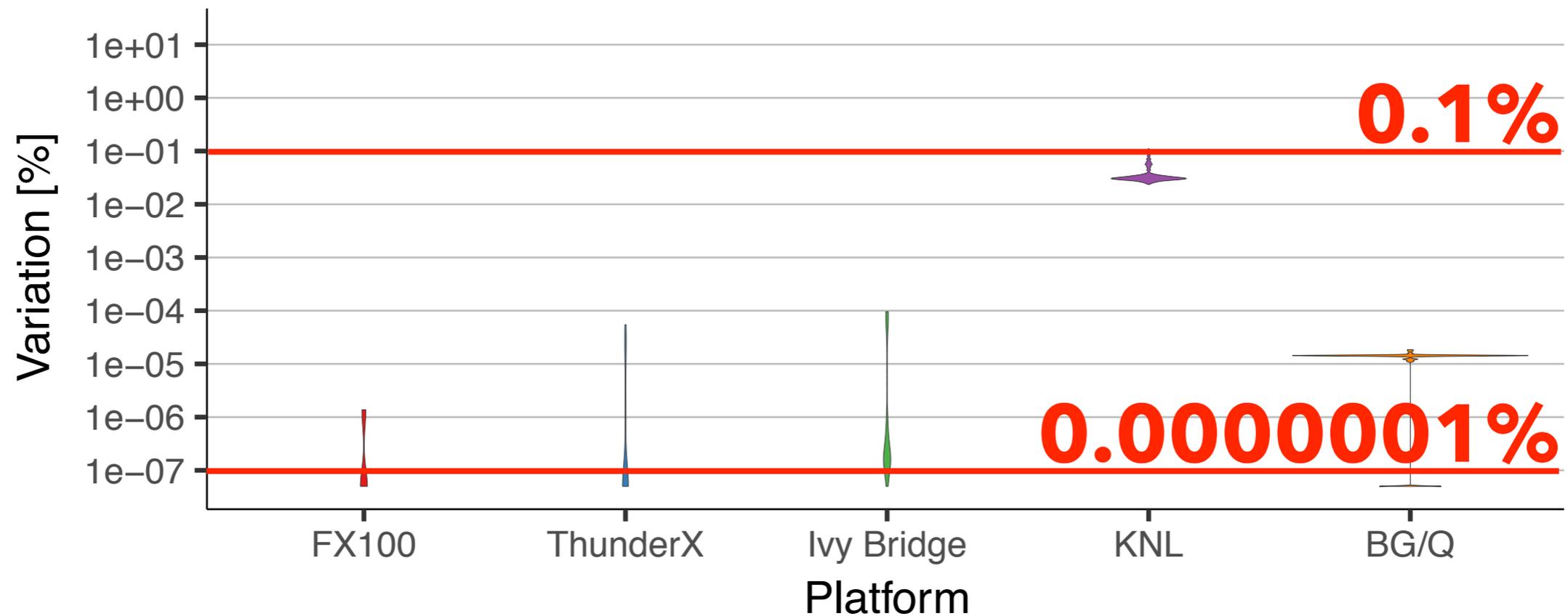
Adam Lackorzynski, Carsten Weinhold, Hermann Härtig, "Decoupled: Low-Effort Noise-Free Execution on Commodity Systems", ROSS 2016, June 2016, Kyoto, Japan

- PhD student: internship at RIKEN, Japan
- Comparative study:
 - Hardware performance variation
 - 5 different CPU architectures
 - Light-weight kernel (McKernel)

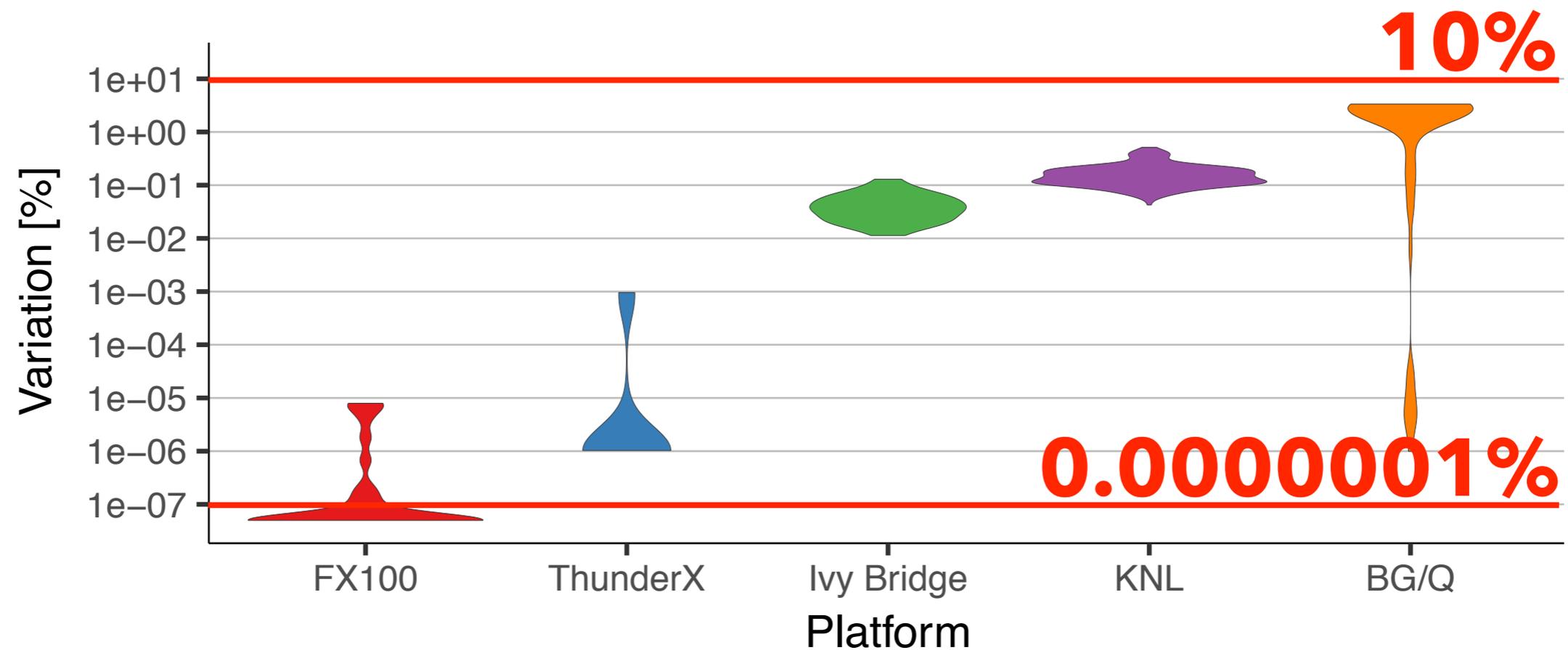
Hannes Weisbach, Brian Kocoloski, Hermann Härtig, Yutaka Ishikawa, Balazs Gerofi, „Hardware Performance Variation: A Comparative Study using Lightweight Kernels“, ISC'18, Frankfurt, Germany, June 2018



FWQ BENCHMARK



Hannes Weisbach, Brian Kocoloski, Hermann Härtig, Yutaka Ishikawa, Balazs Gerofi, „Hardware Performance Variation: A Comparative Study using Lightweight Kernels“, ISC'18, Frankfurt, Germany, June 2018



Hannes Weisbach, Brian Kocoloski, Hermann Härtig, Yutaka Ishikawa, Balazs Gerofi, „Hardware Performance Variation: A Comparative Study using Lightweight Kernels“, ISC'18, Frankfurt, Germany, June 2018

IMBALANCED WORKLOADS

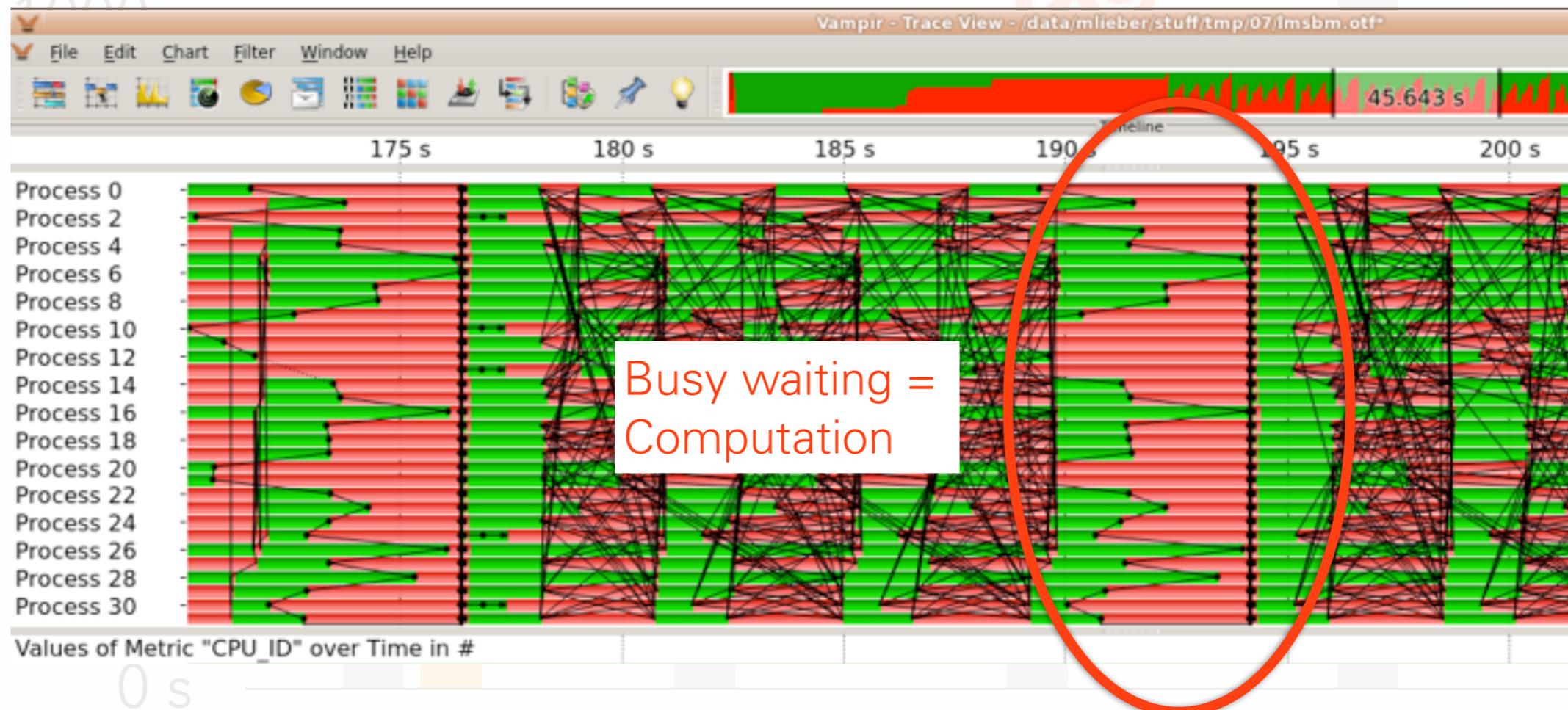


OVERDECOMPOSITION

■ Polling

■ Blocking

Application: COSMO-SPECS+FD4



Busy waiting =
Computation

of MPI Processes per Node

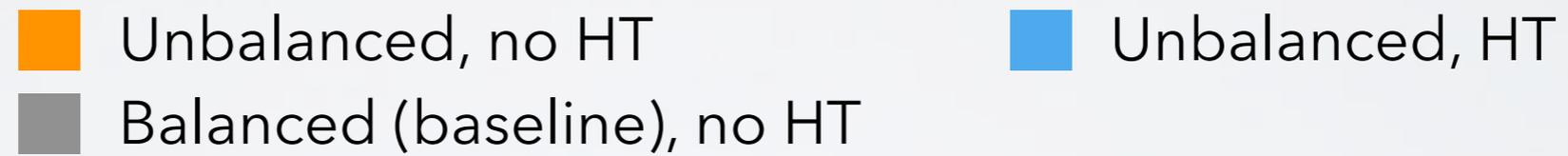
of Nodes

4 / 4

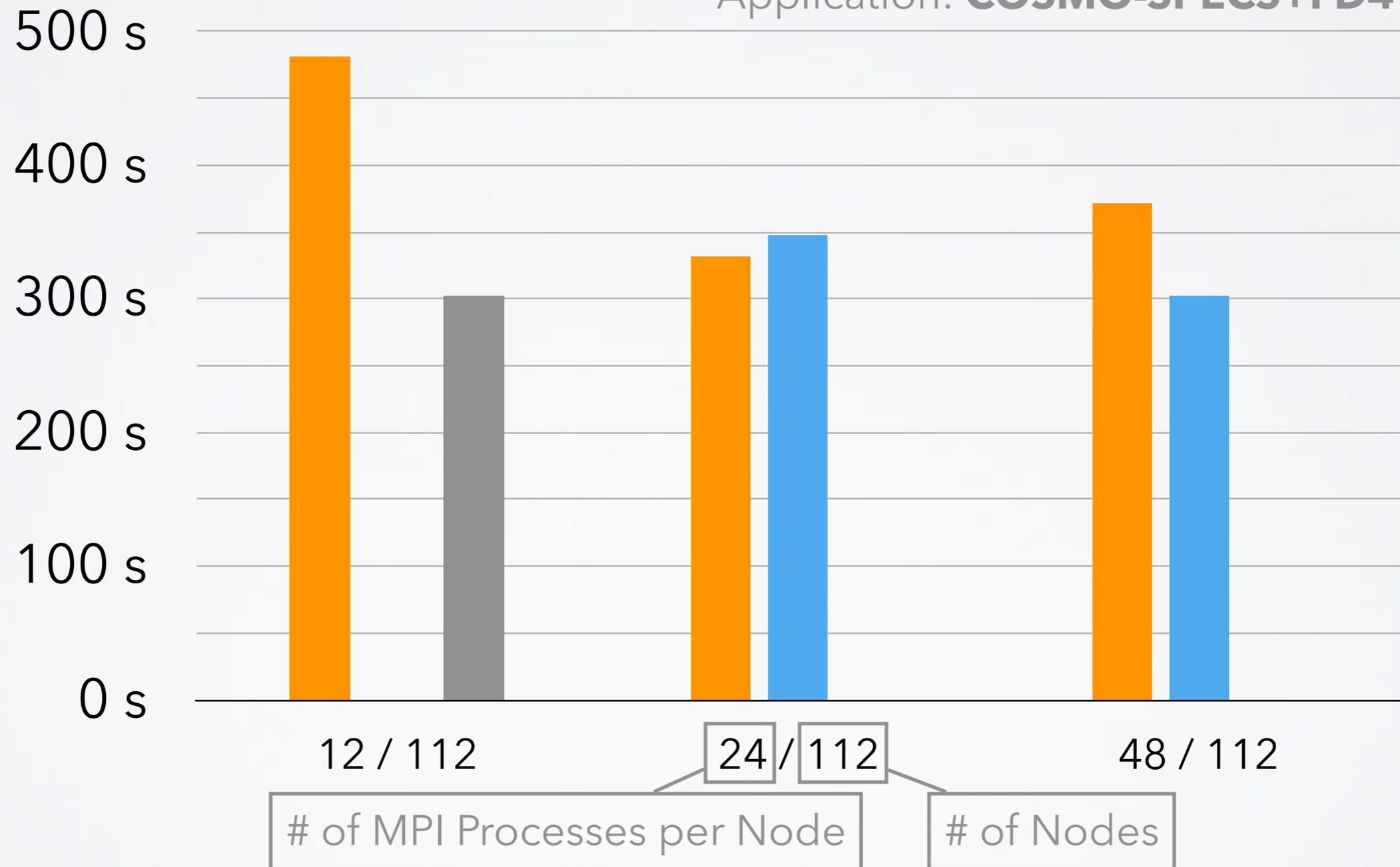
8 / 4

16 / 4

32 / 4

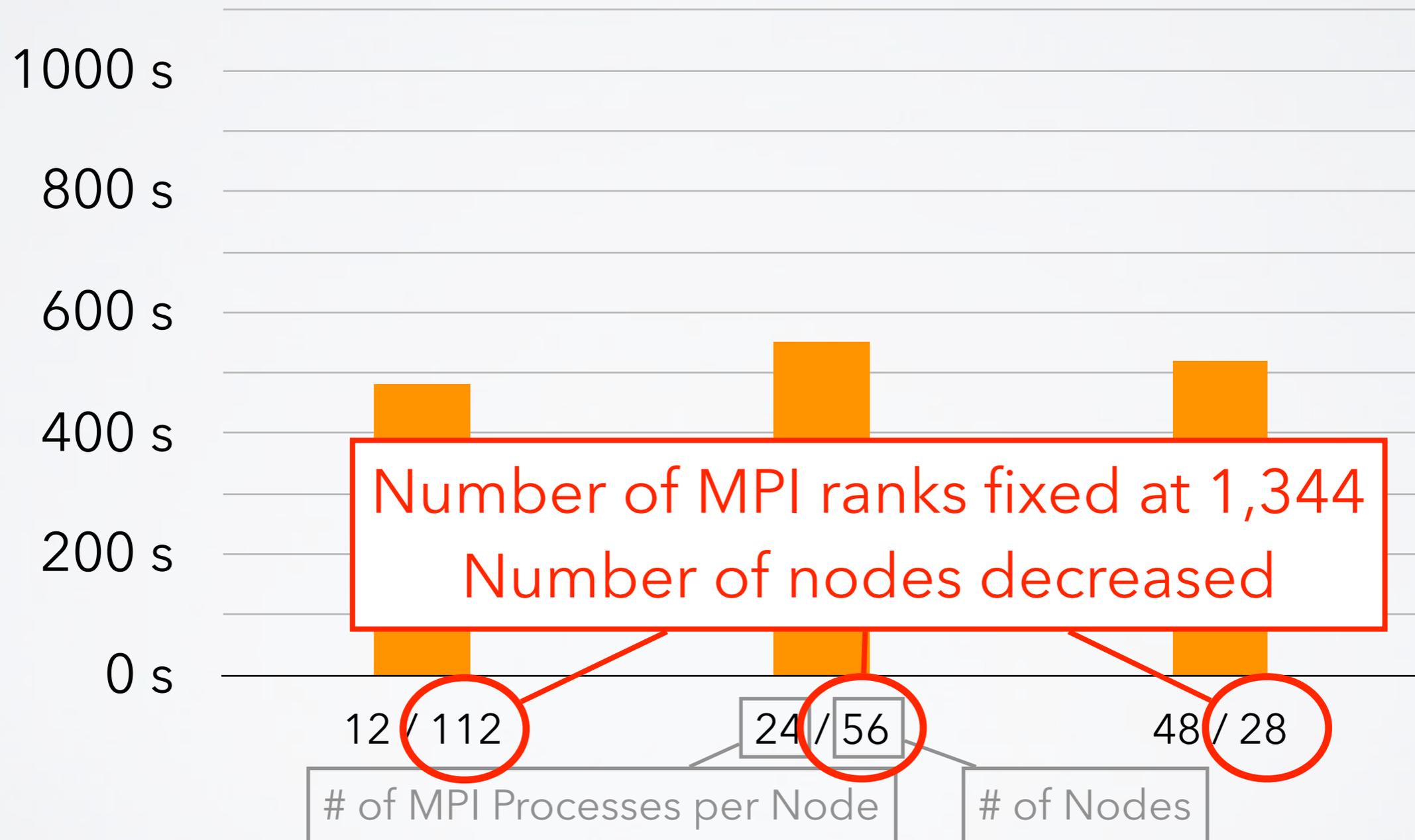


Application: **COSMO-SPECS+FD4**



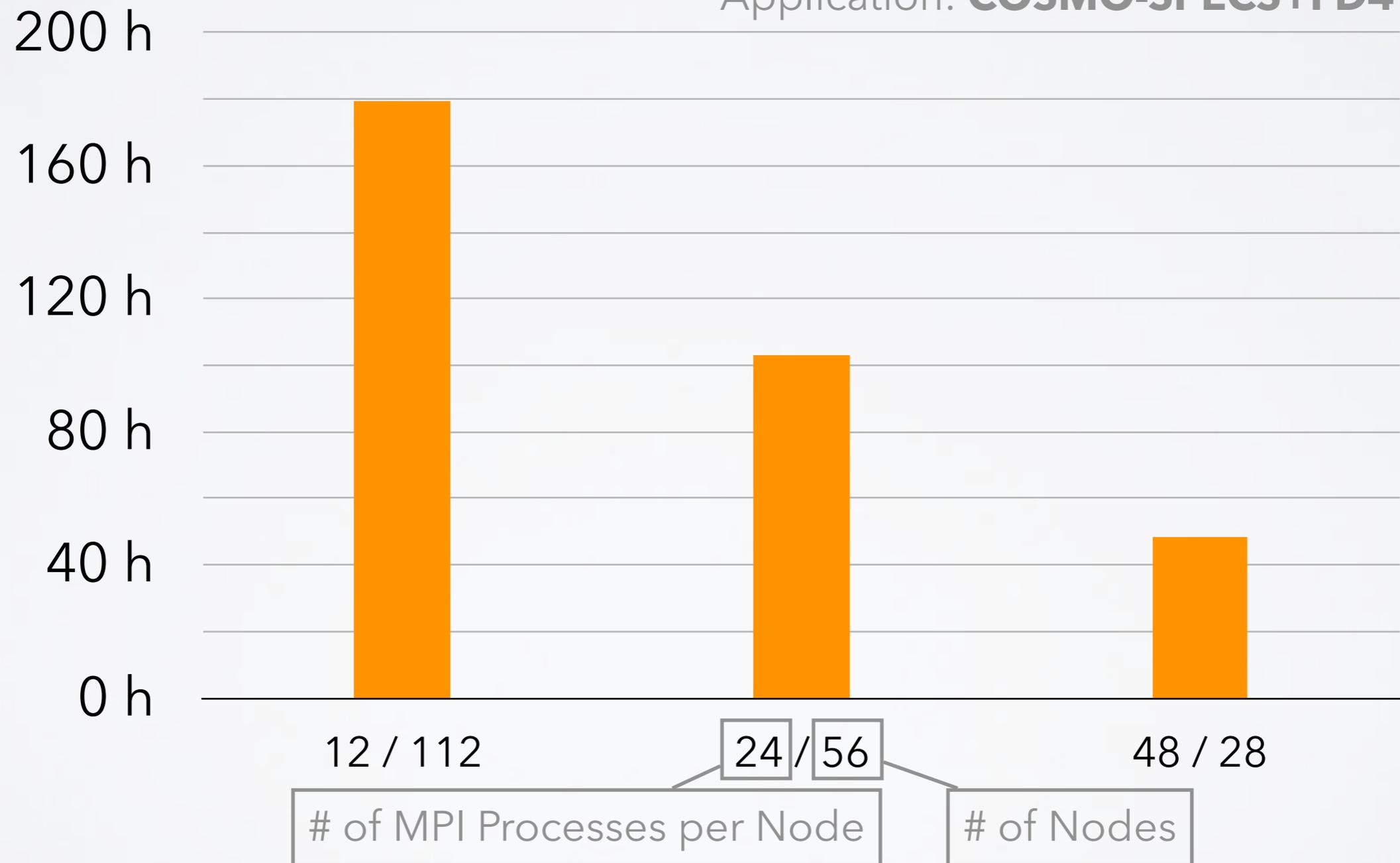
Unbalanced, no HT

Application: **COSMO-SPECS+FD4**



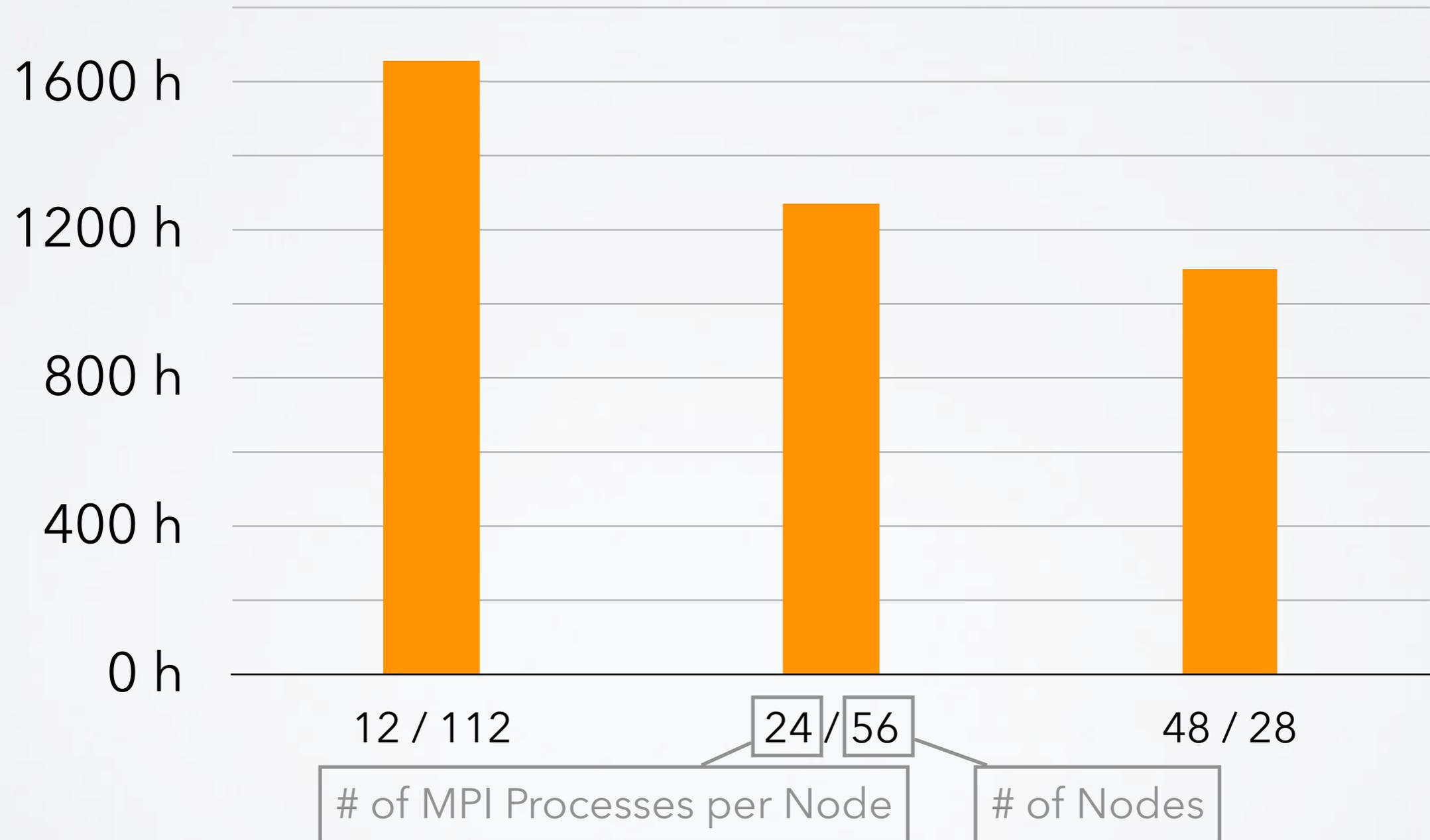
Unbalanced, no HT

Application: **COSMO-SPECS+FD4**



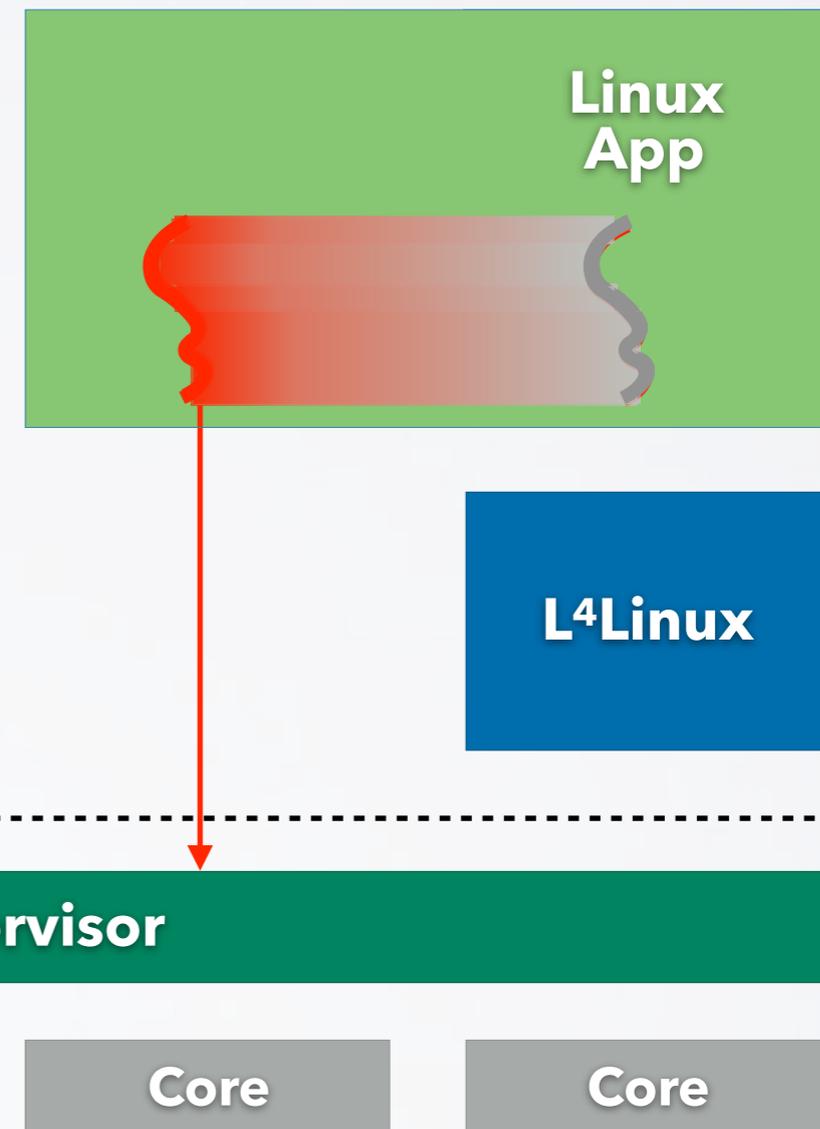
Unbalanced, no HT

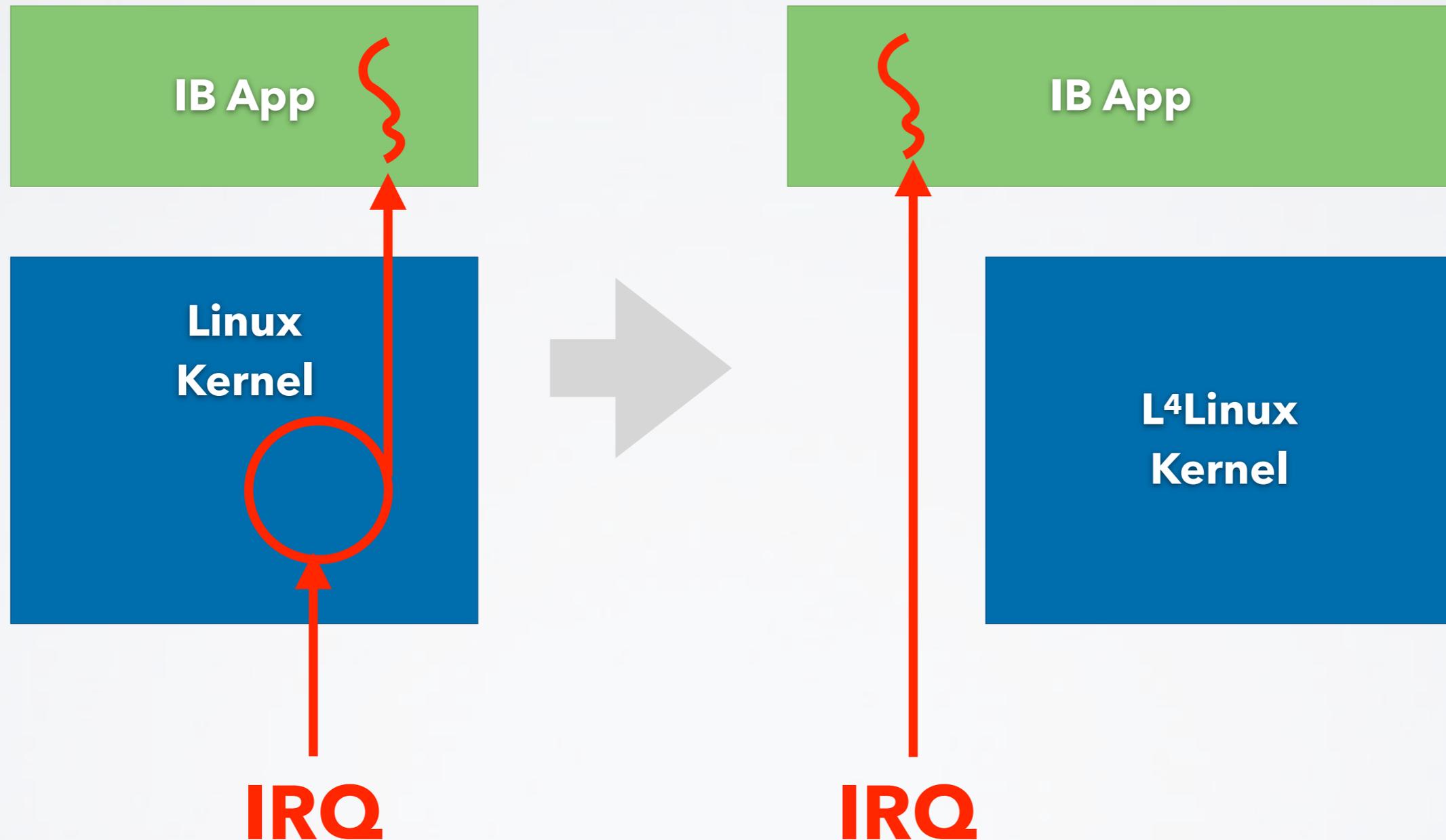
Application: **CP2K**

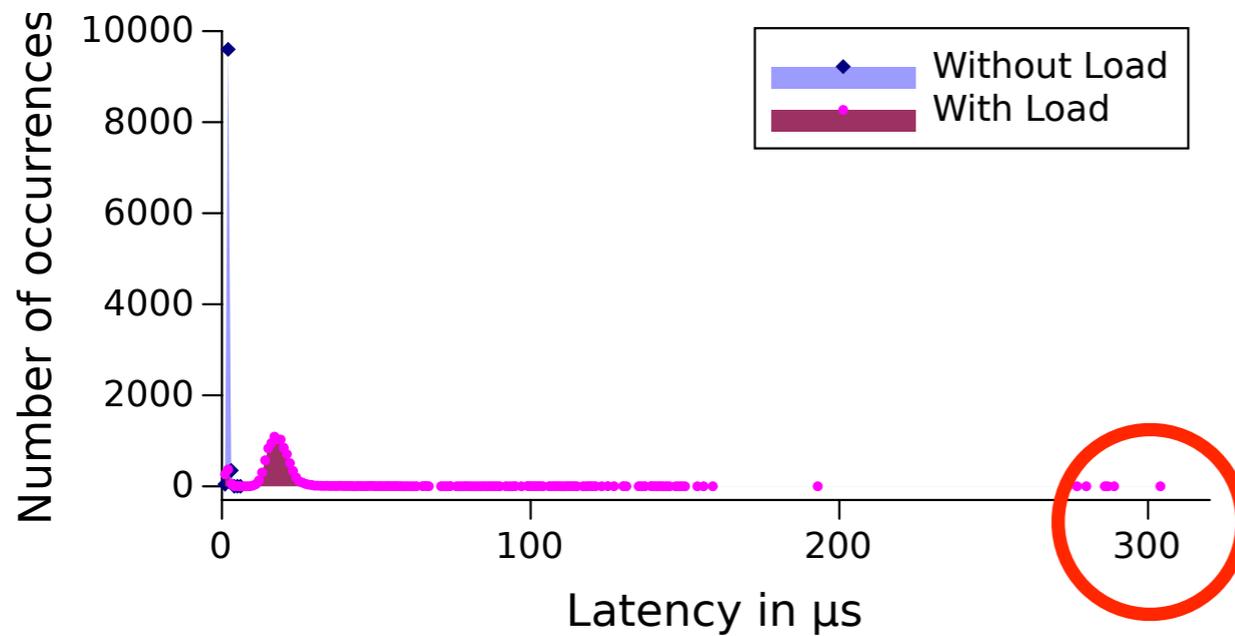


WIP: DECOUPLED INTERRUPTS

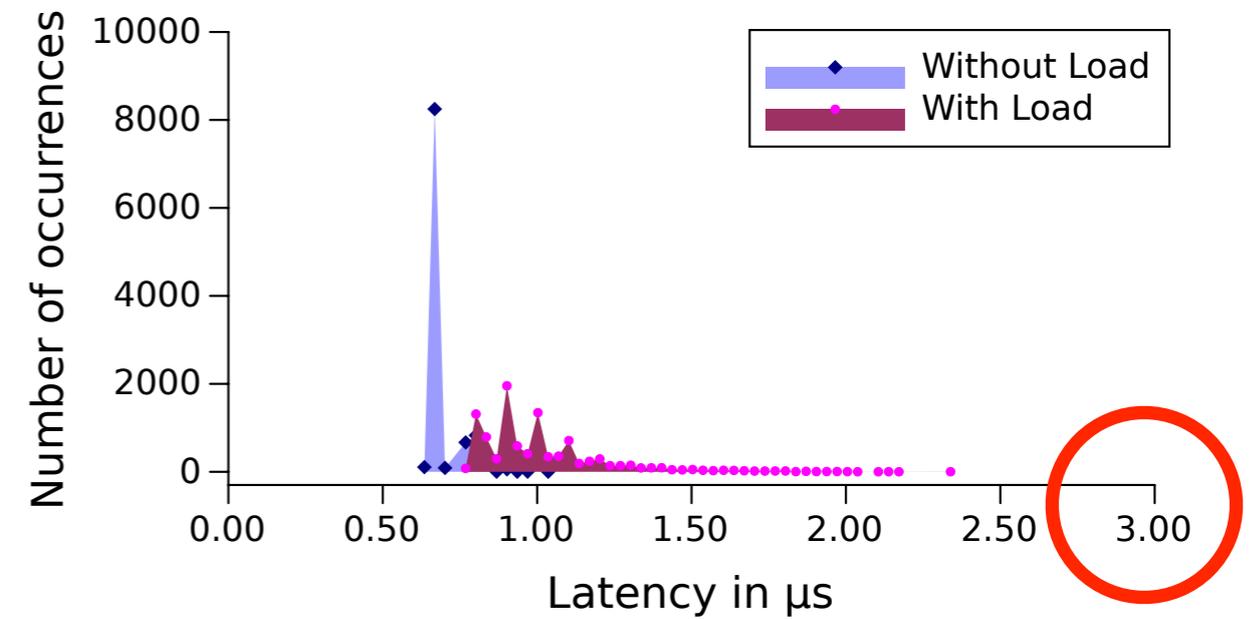
- **Decoupling:** move Linux thread to new L4 thread on its own core
- **Linux syscall:** Move back to Linux
- **Direct I/O** device access
- **L4 syscalls:**
 - Memory
 - Threads & Scheduling
 - Interrupts







Linux



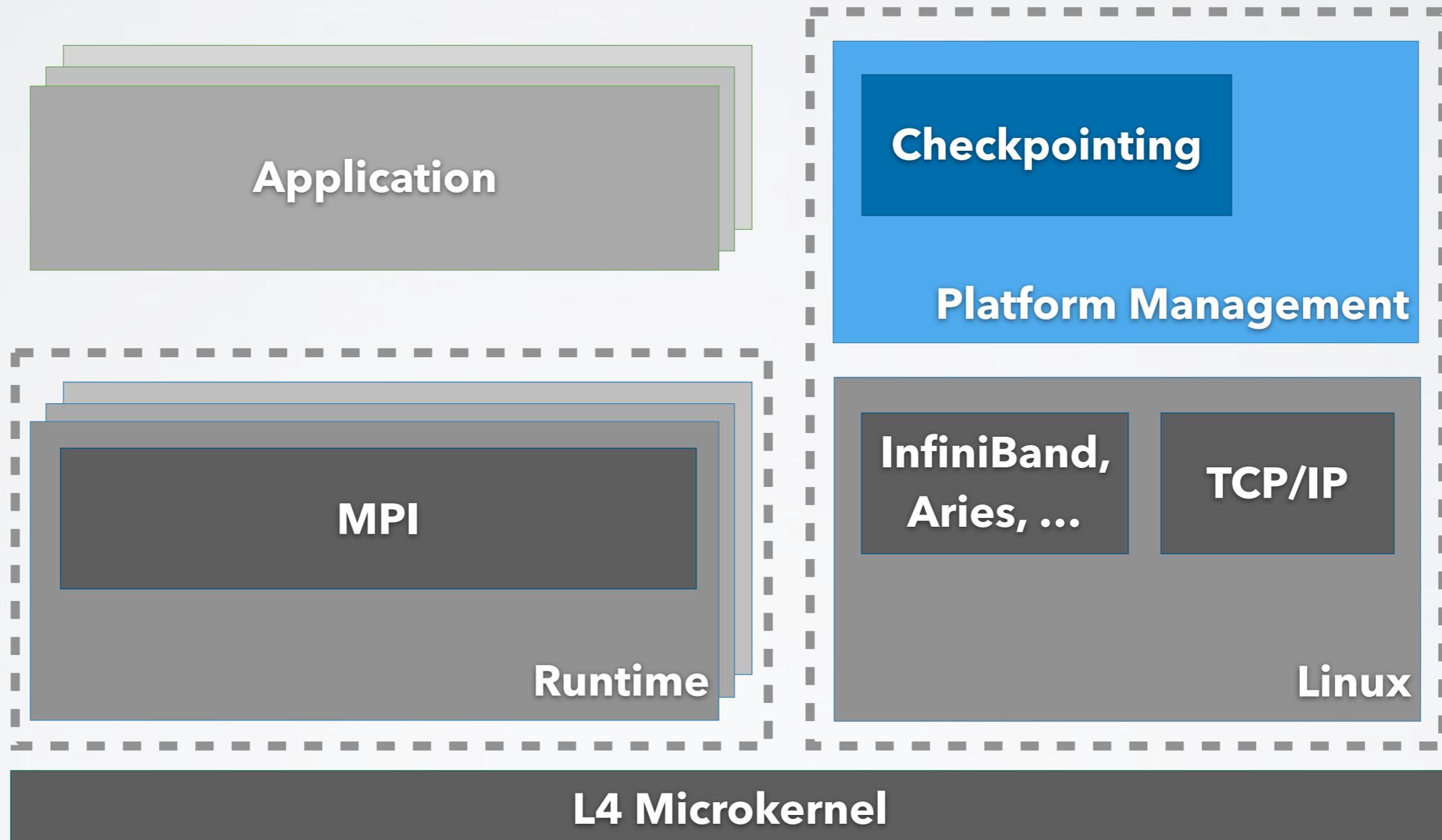
L4

Work in progress: User-space handling of InfiniBand HCA interrupts

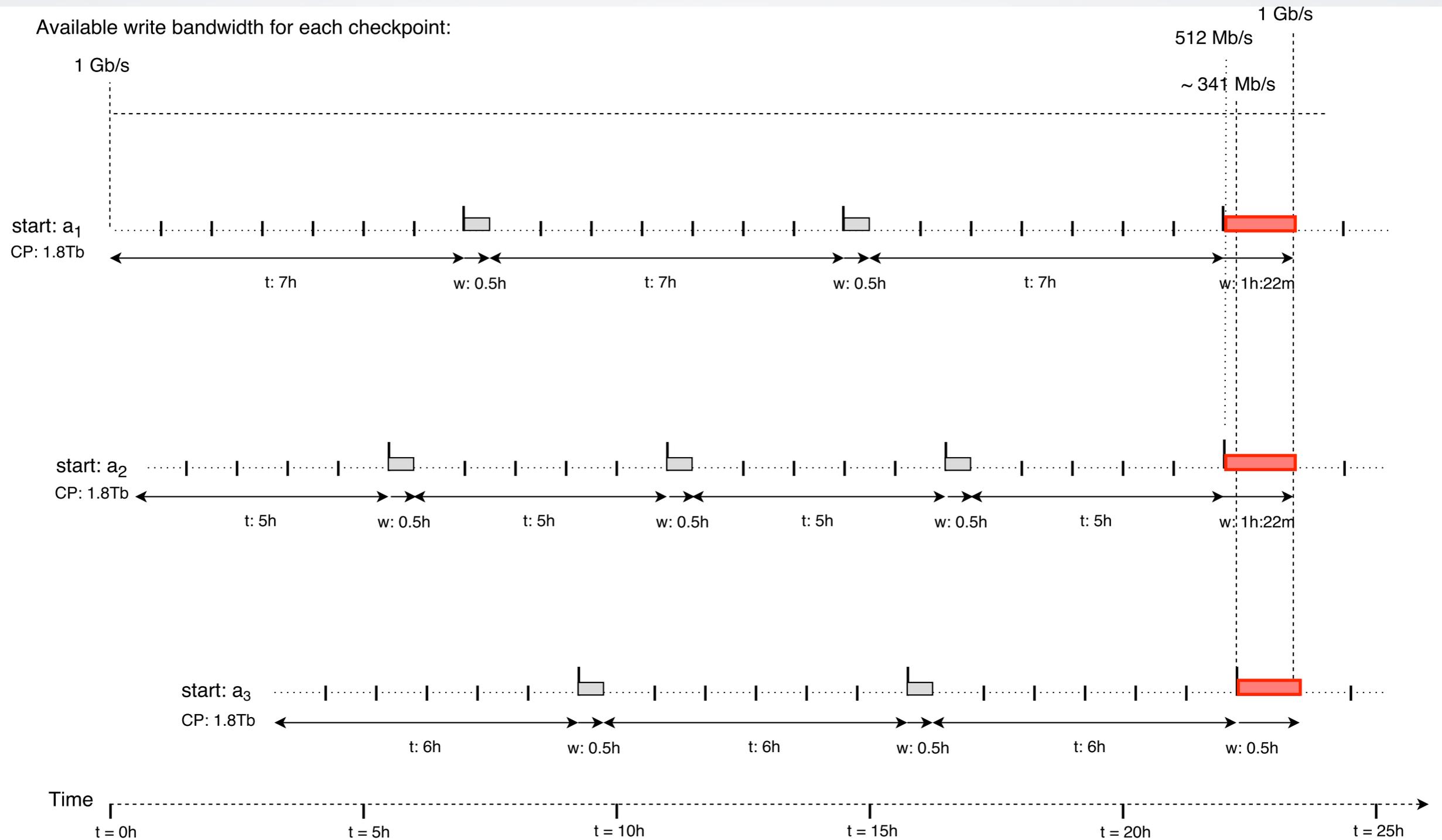
Adam Lackorzynski, Carsten Weinhold, Hermann Härtig, „Predictable Low-Latency Interrupt Response with General-Purpose Systems“, OSPERT 2017, Dubrovnik, Croatia, June 2017



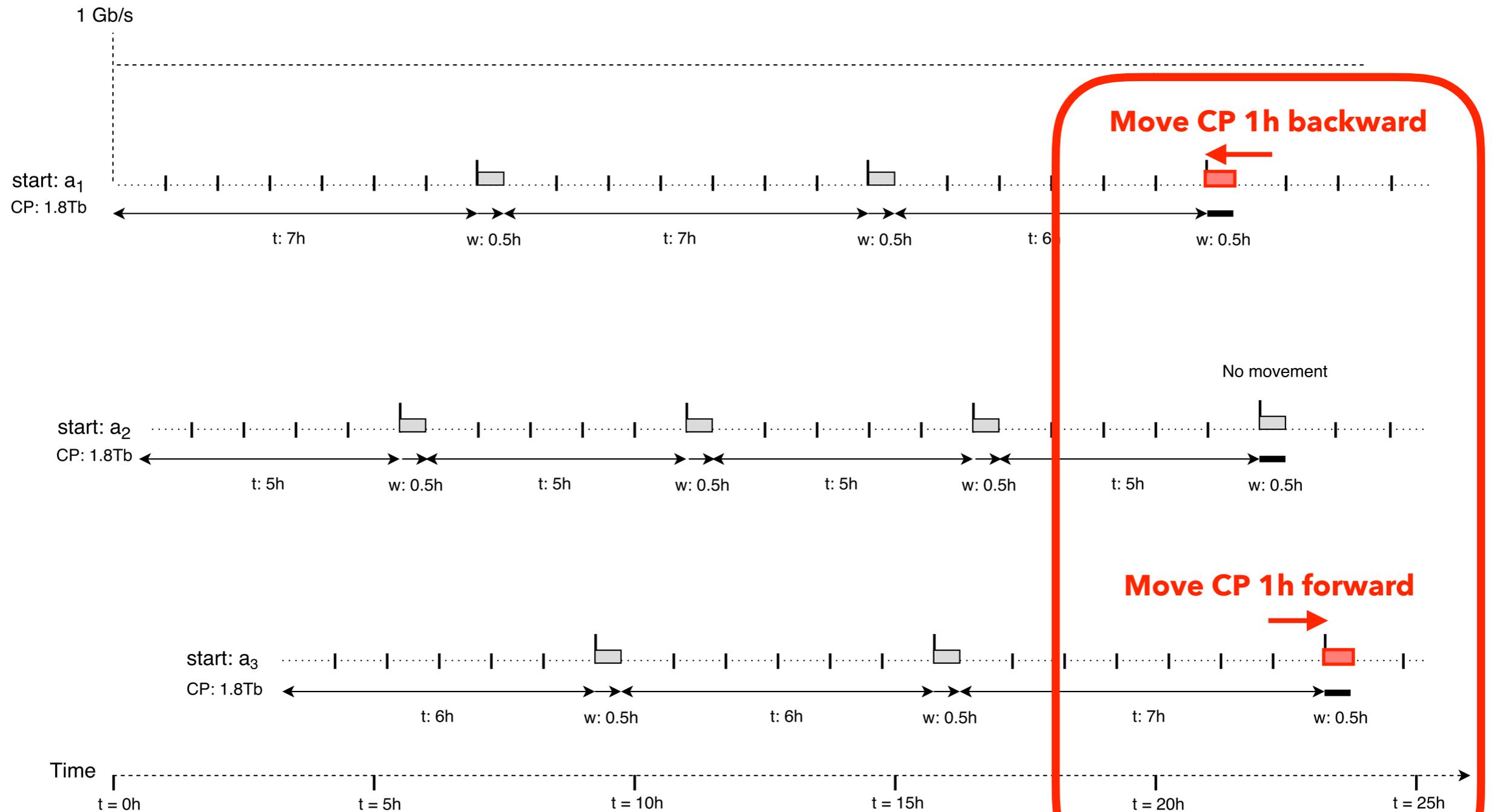
NODE ARCHITECTURE



Available write bandwidth for each checkpoint:

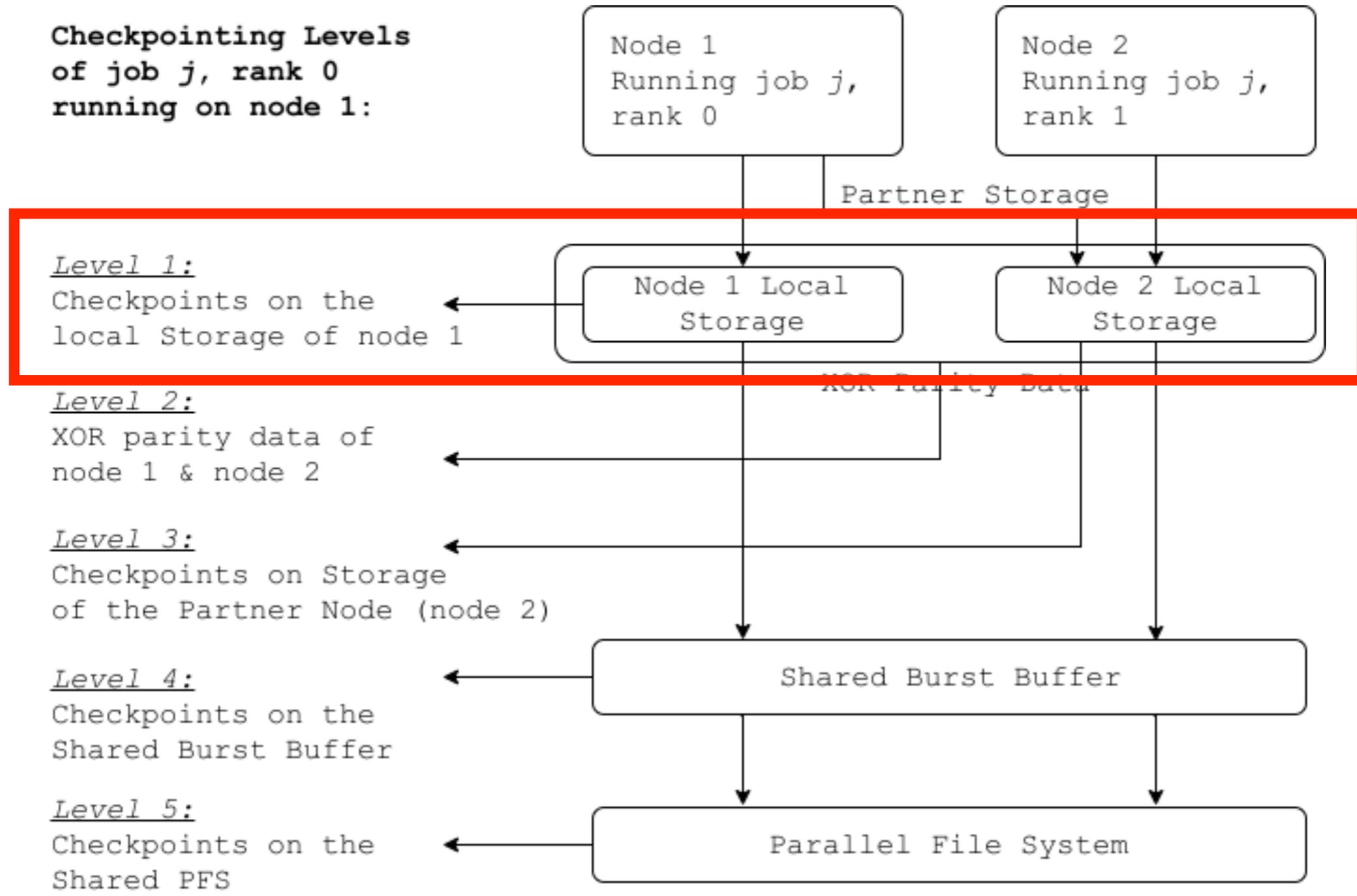


Available write bandwidth for each checkpoint:



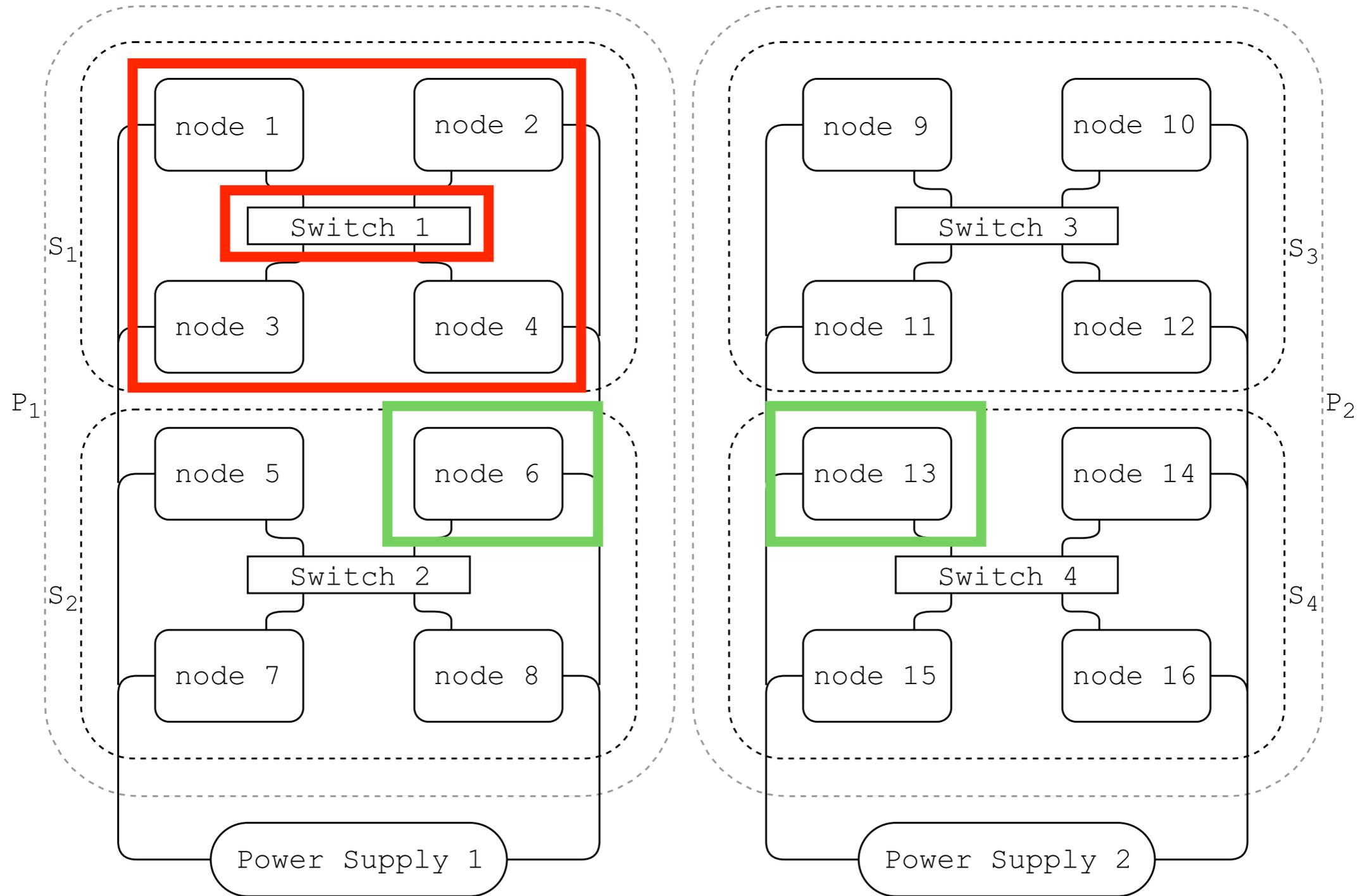


**Checkpointing Levels
of job j , rank 0
running on node 1:**



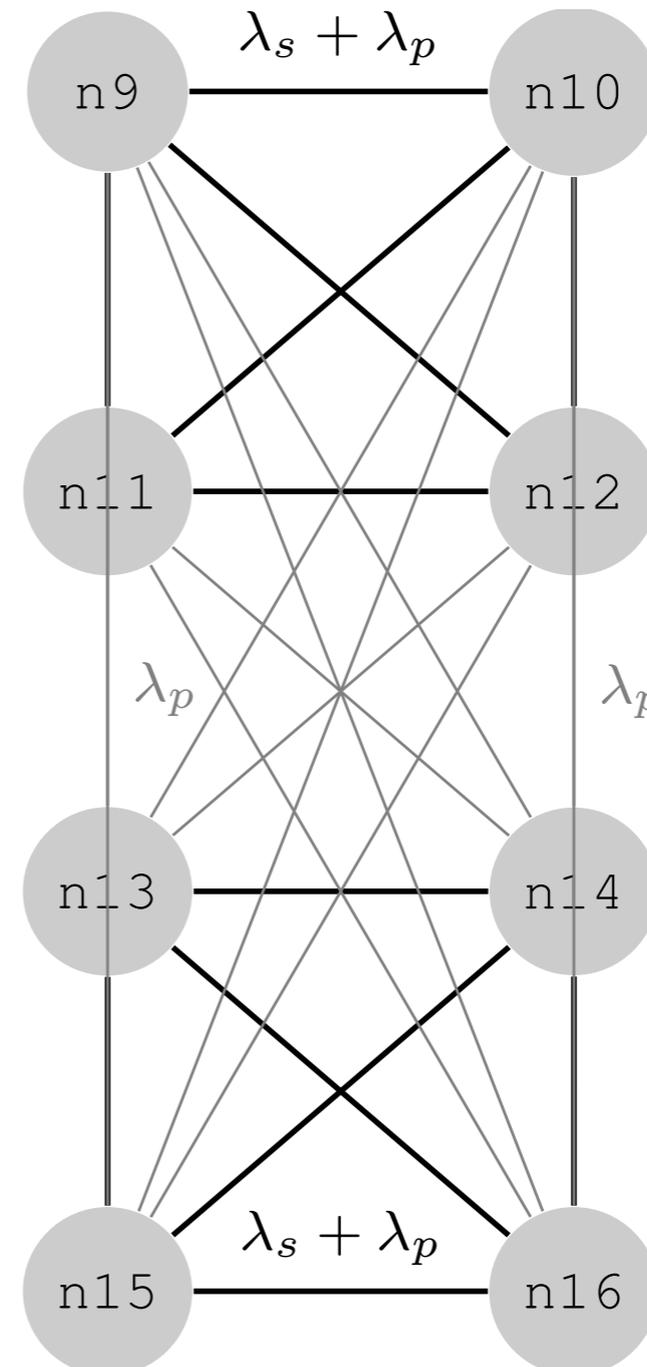
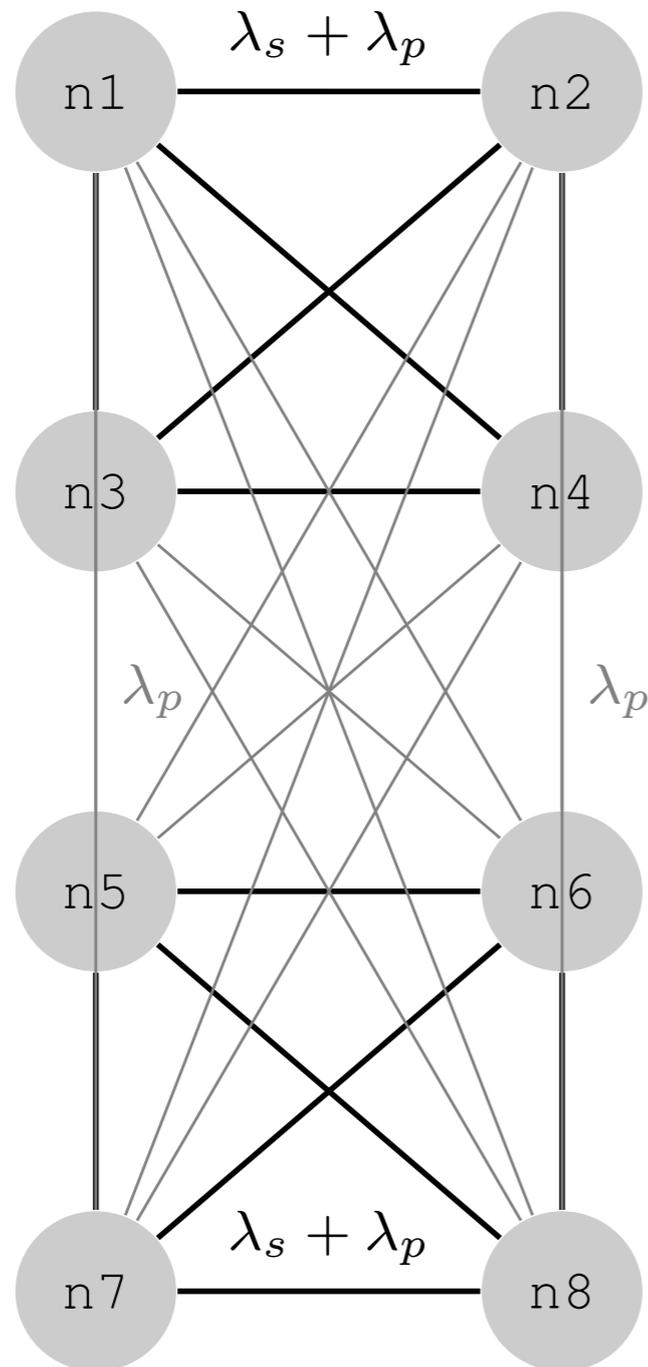


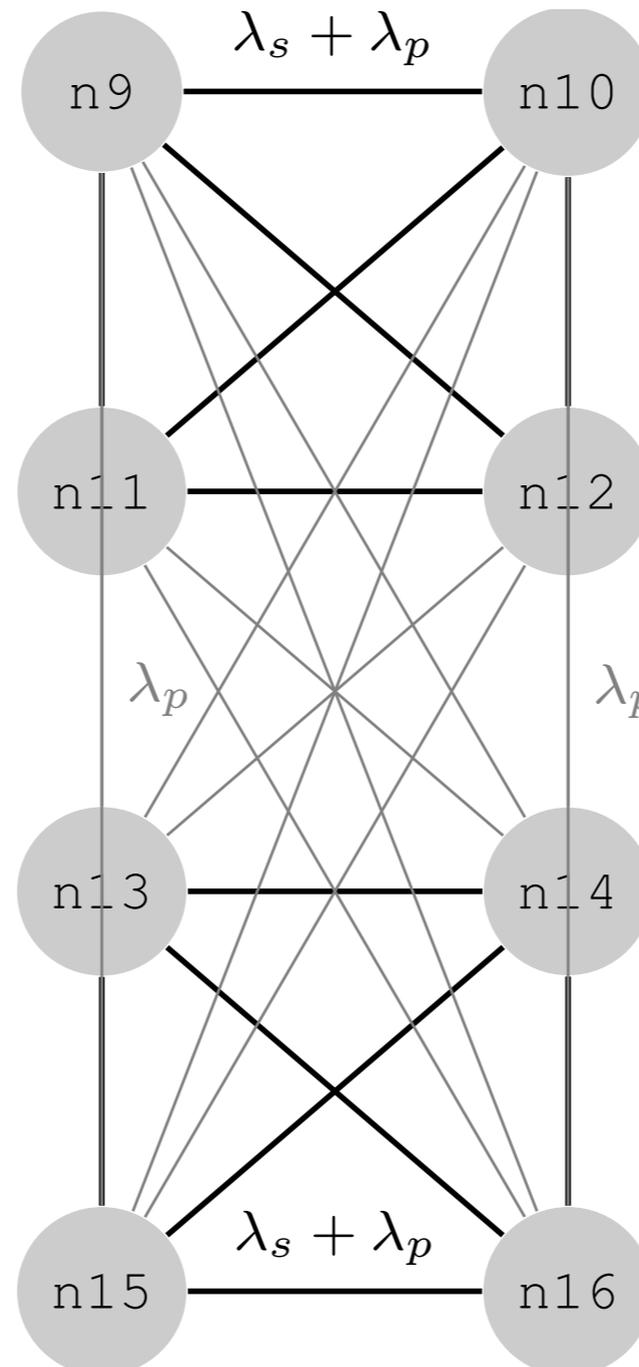
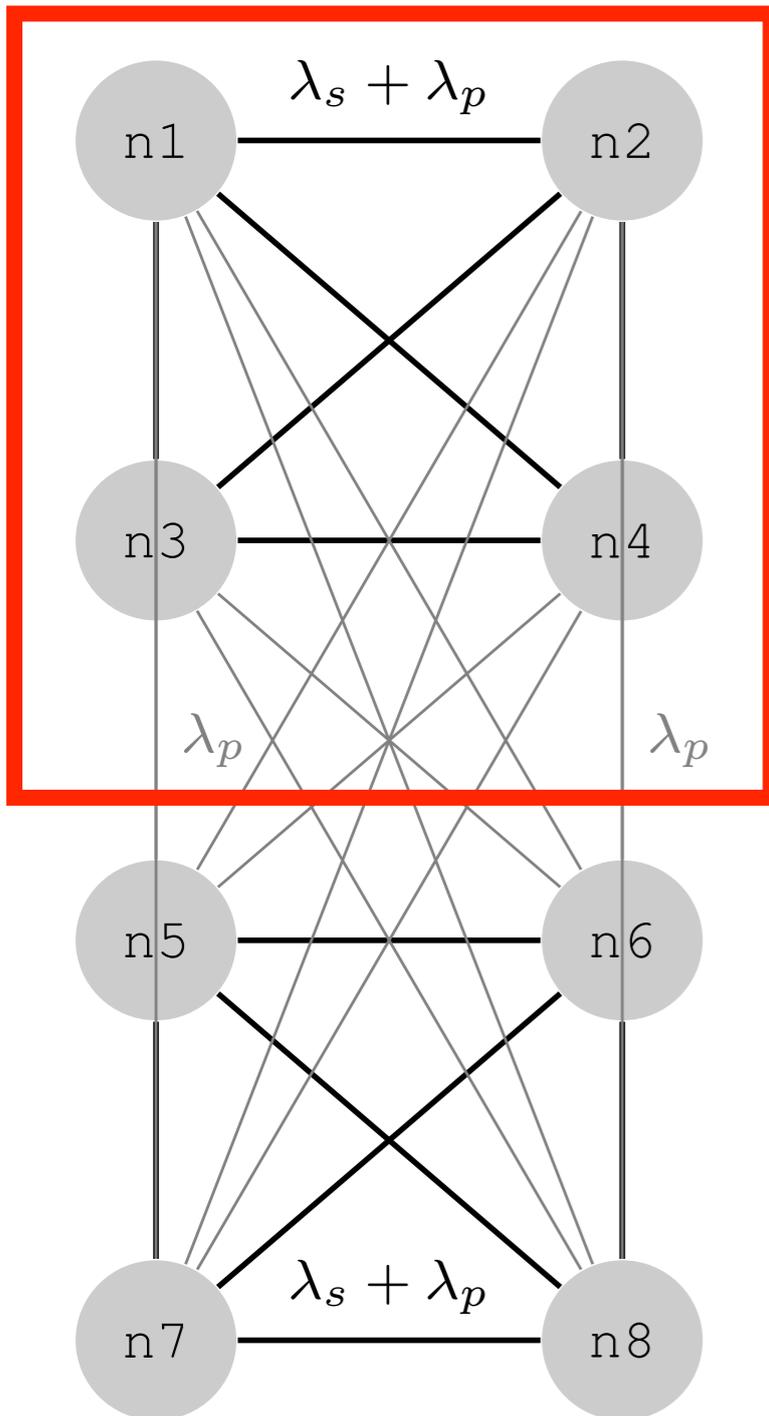
CORRELATED FAILURE





CORRELATED FAILURE





Graph problem:

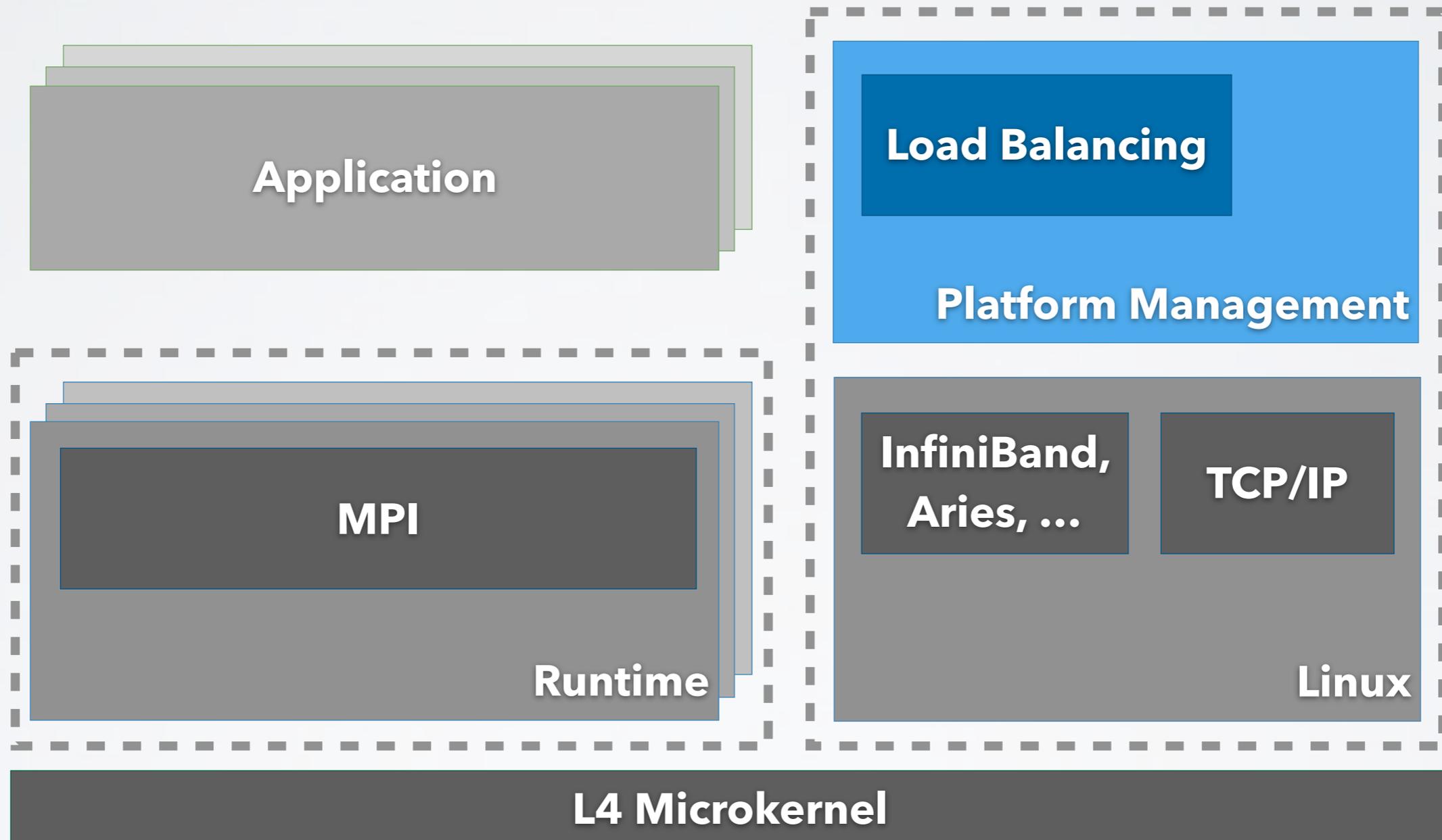
- Find disjoint independent sets
- Find dominating subgraphs („least correlated nodes“)

Optimization problem:

- least correlated nodes for checkpoint distribution
- Consider: job run time, C/R cost, MTTI
- Minimize run time

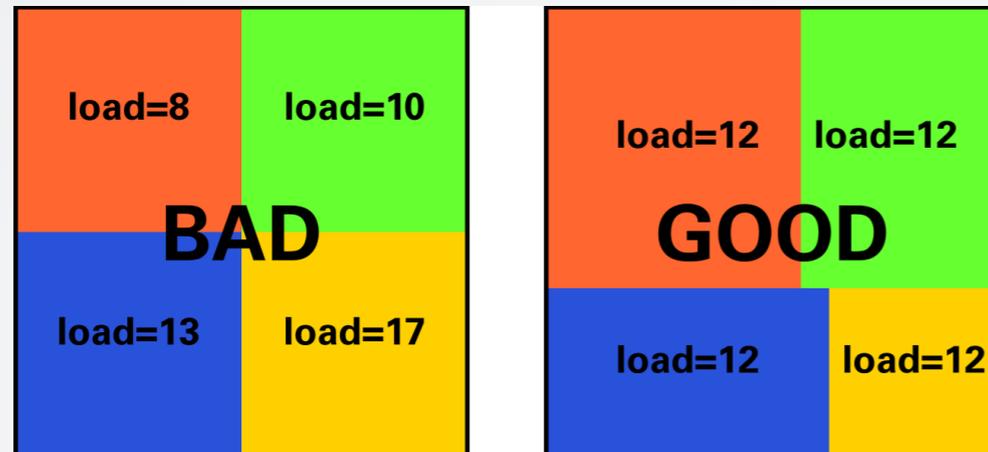


NODE ARCHITECTURE

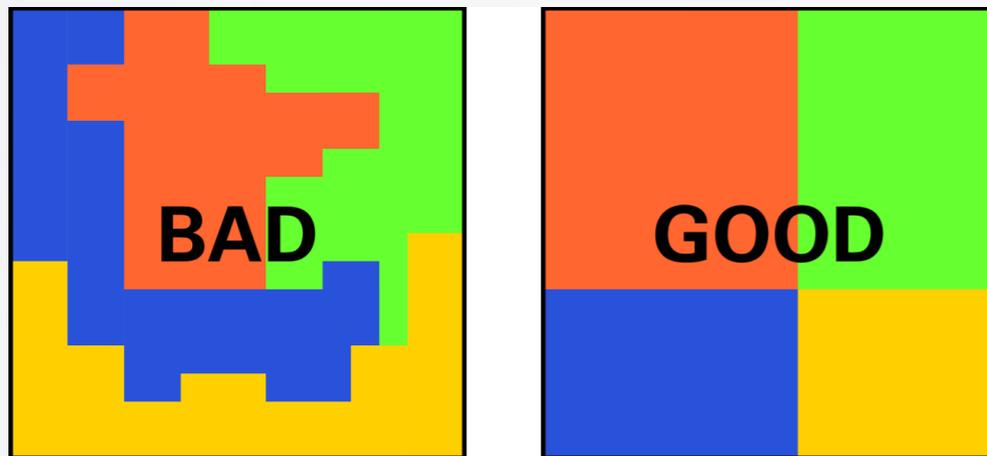
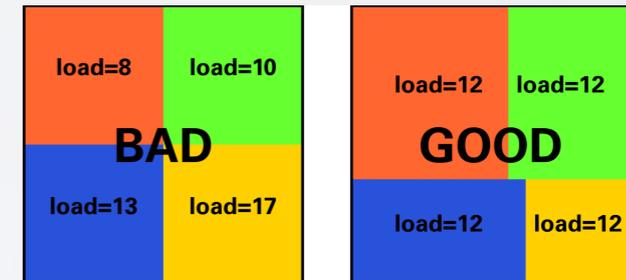




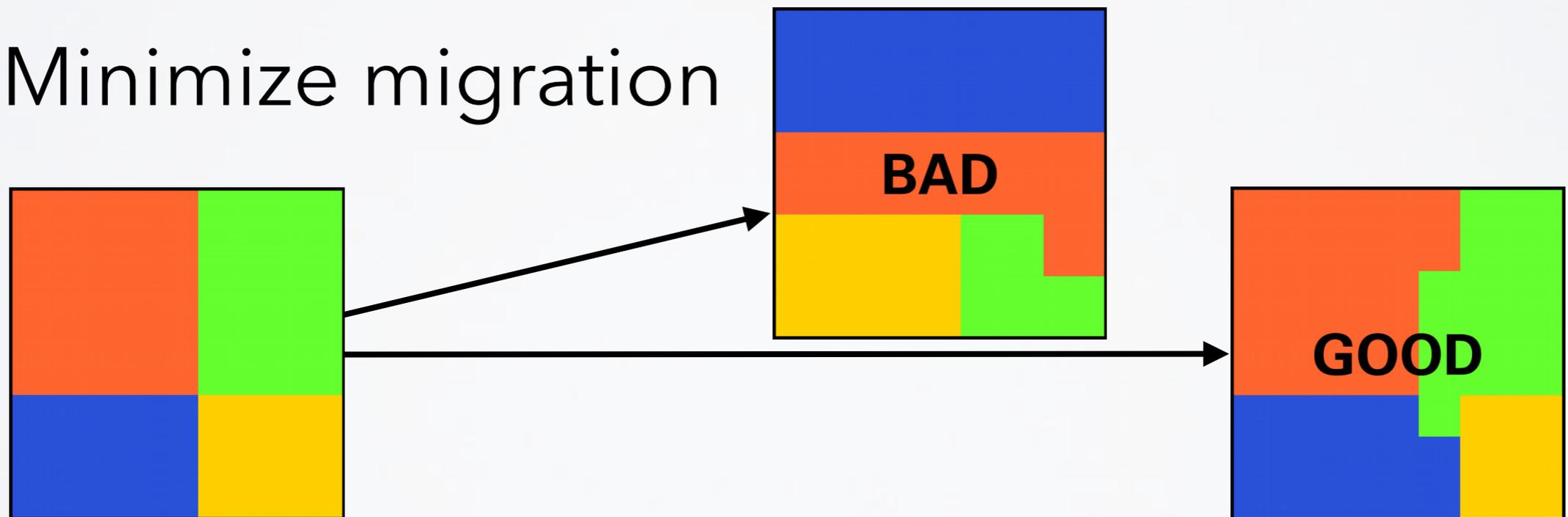
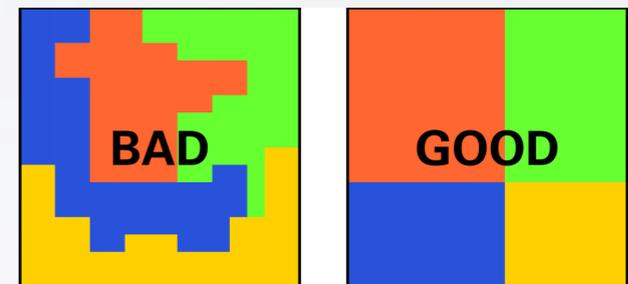
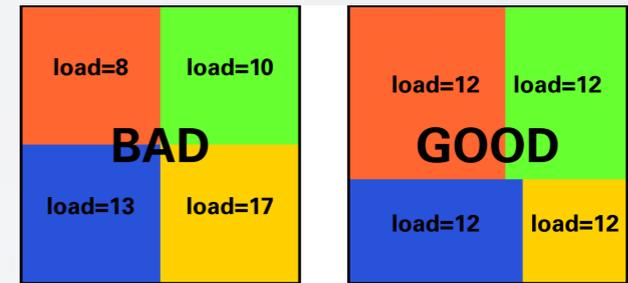
■ Balance workload



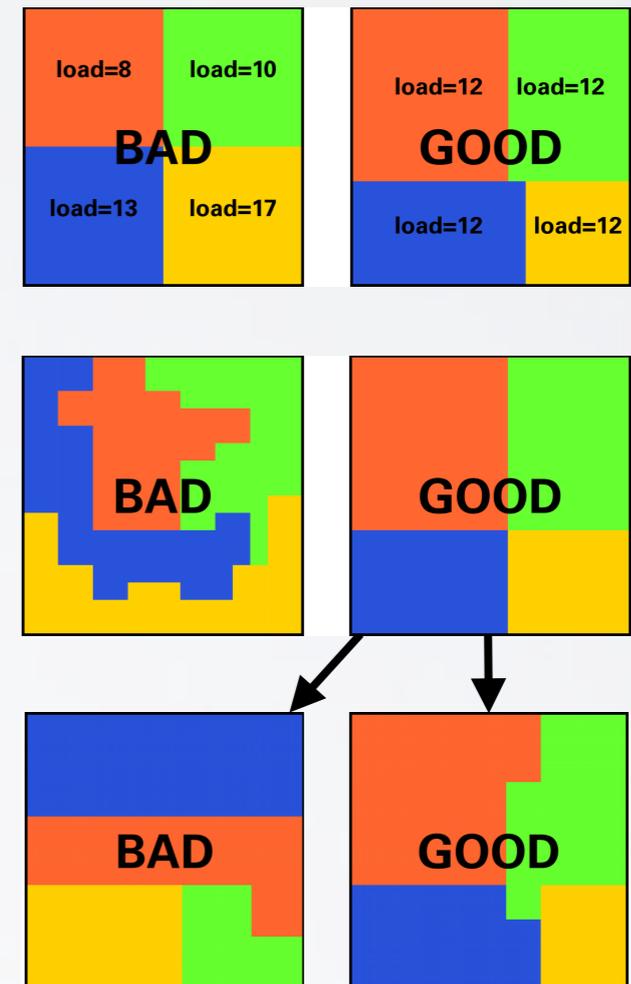
- Balance workload
- Minimize communication between partitions

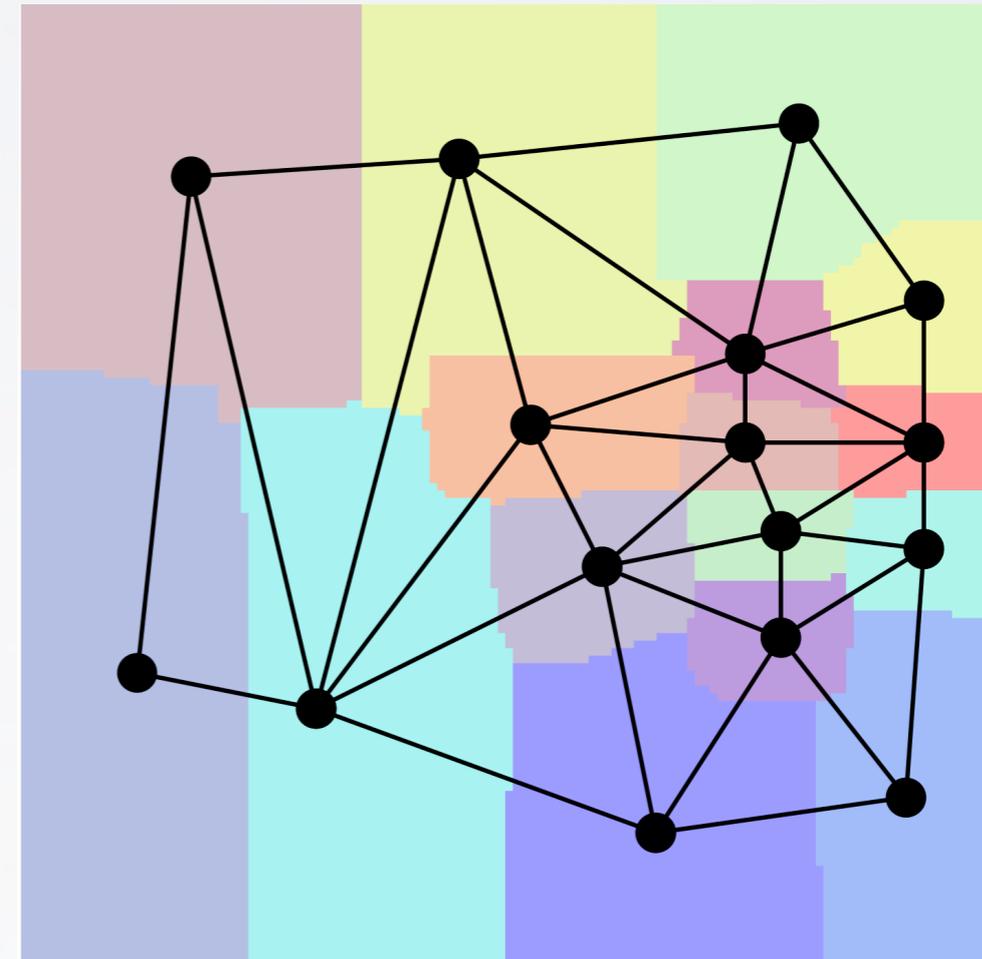


- Balance workload
- Minimize communication between partitions
- Minimize migration



- Balance workload
- Minimize communication between partitions
- Minimize migration
- Compute new partitions fast





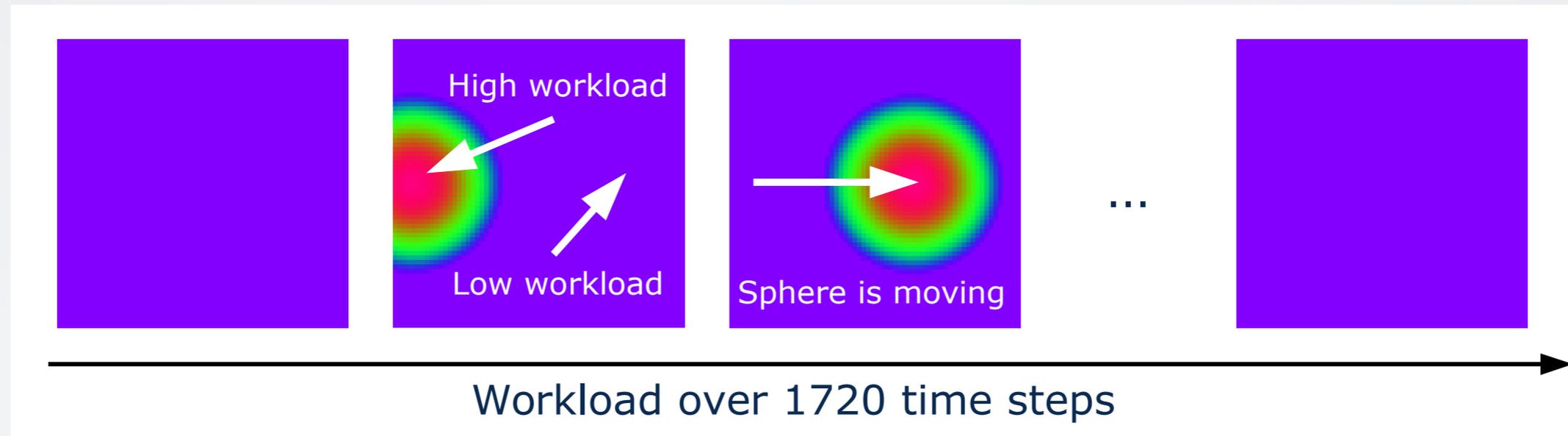
Diffusion graph topology from application topology

Diffusion coefficient weighted by interface length:

- Tasks migrated between neighbor partitions
- Better partition shape

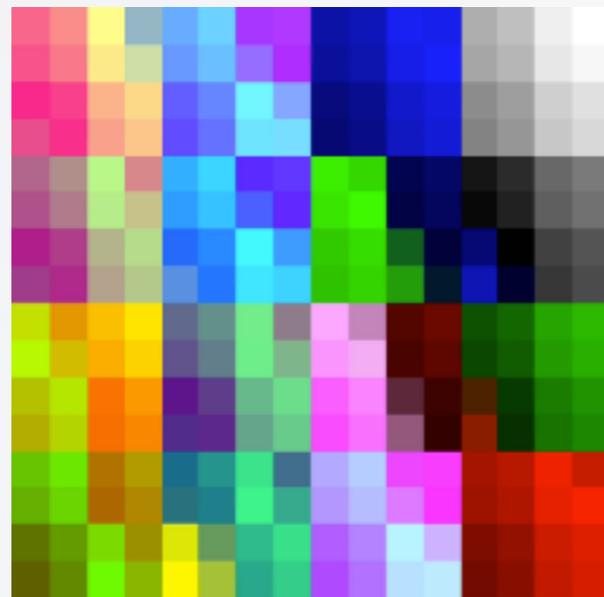
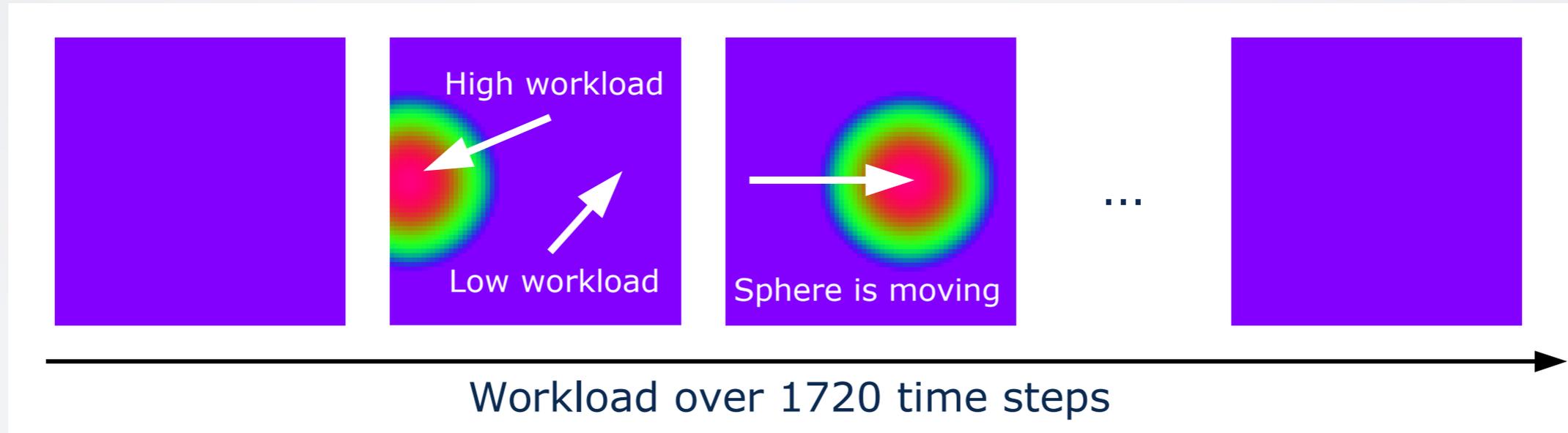


DIFFUSION EXAMPLE





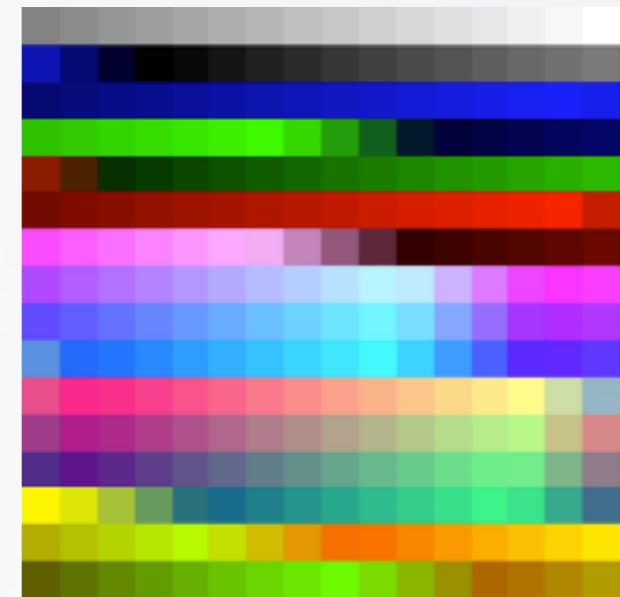
DIFFUSION EXAMPLE



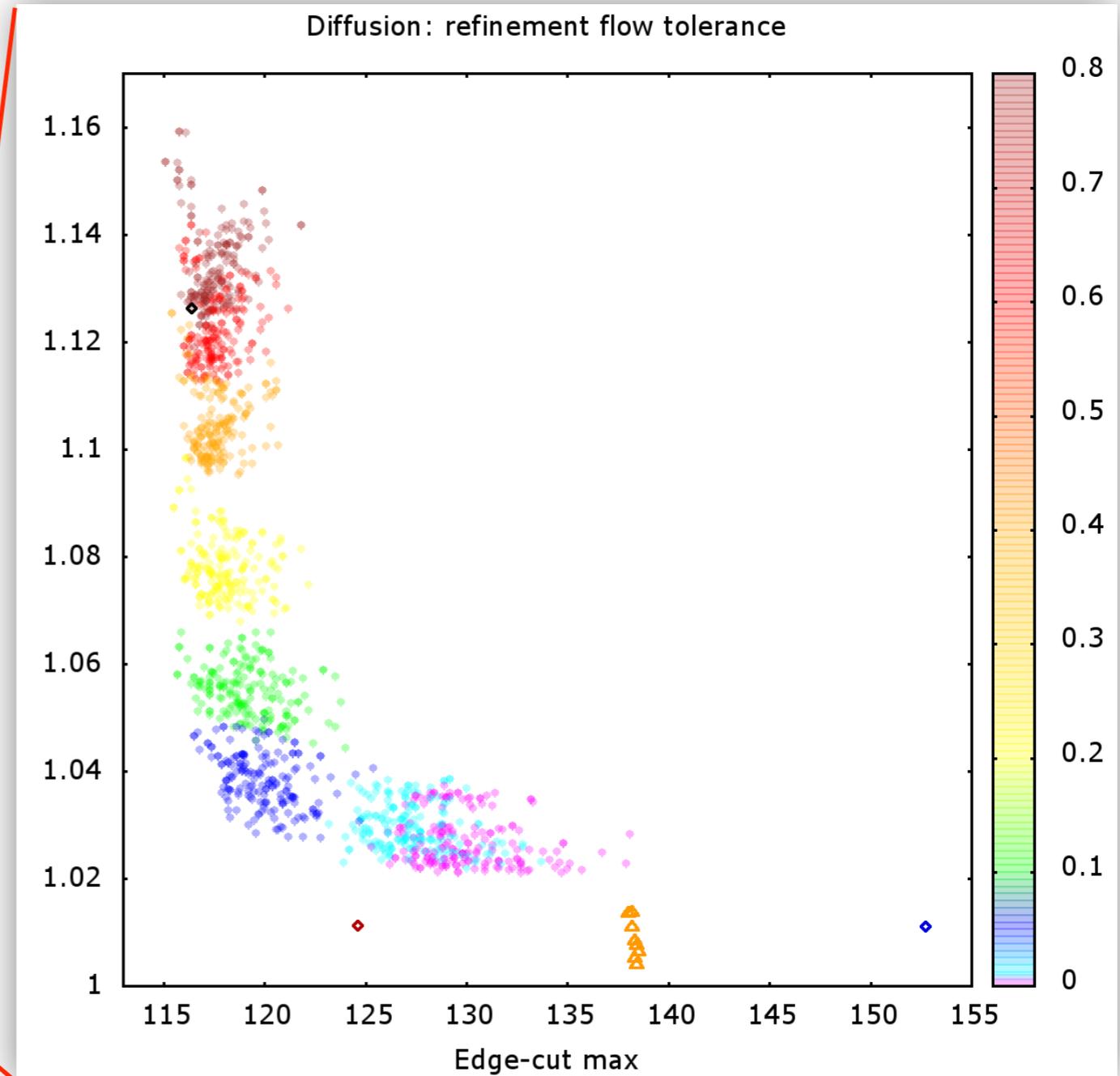
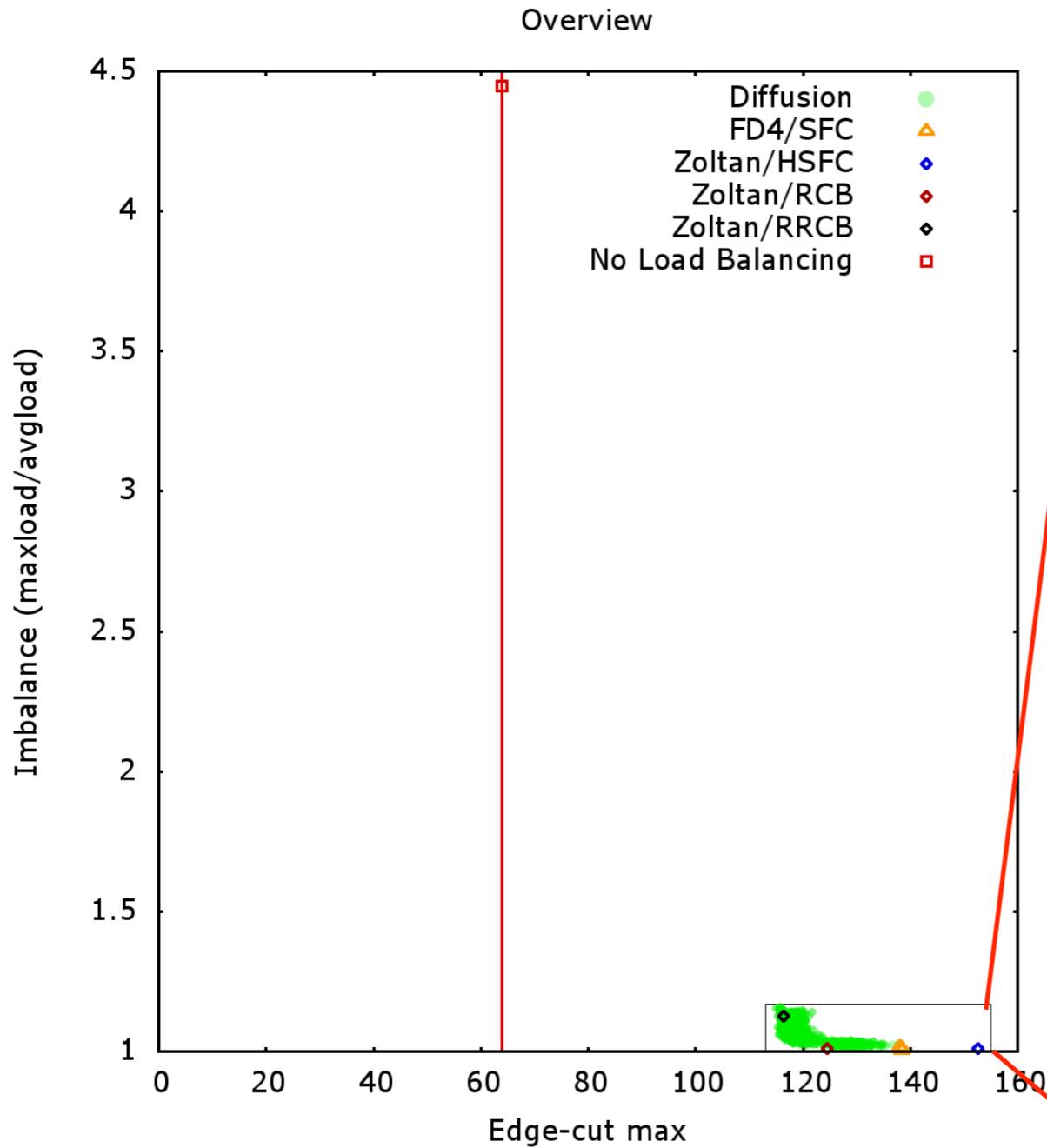
Zoltan



Space-filling
Curves



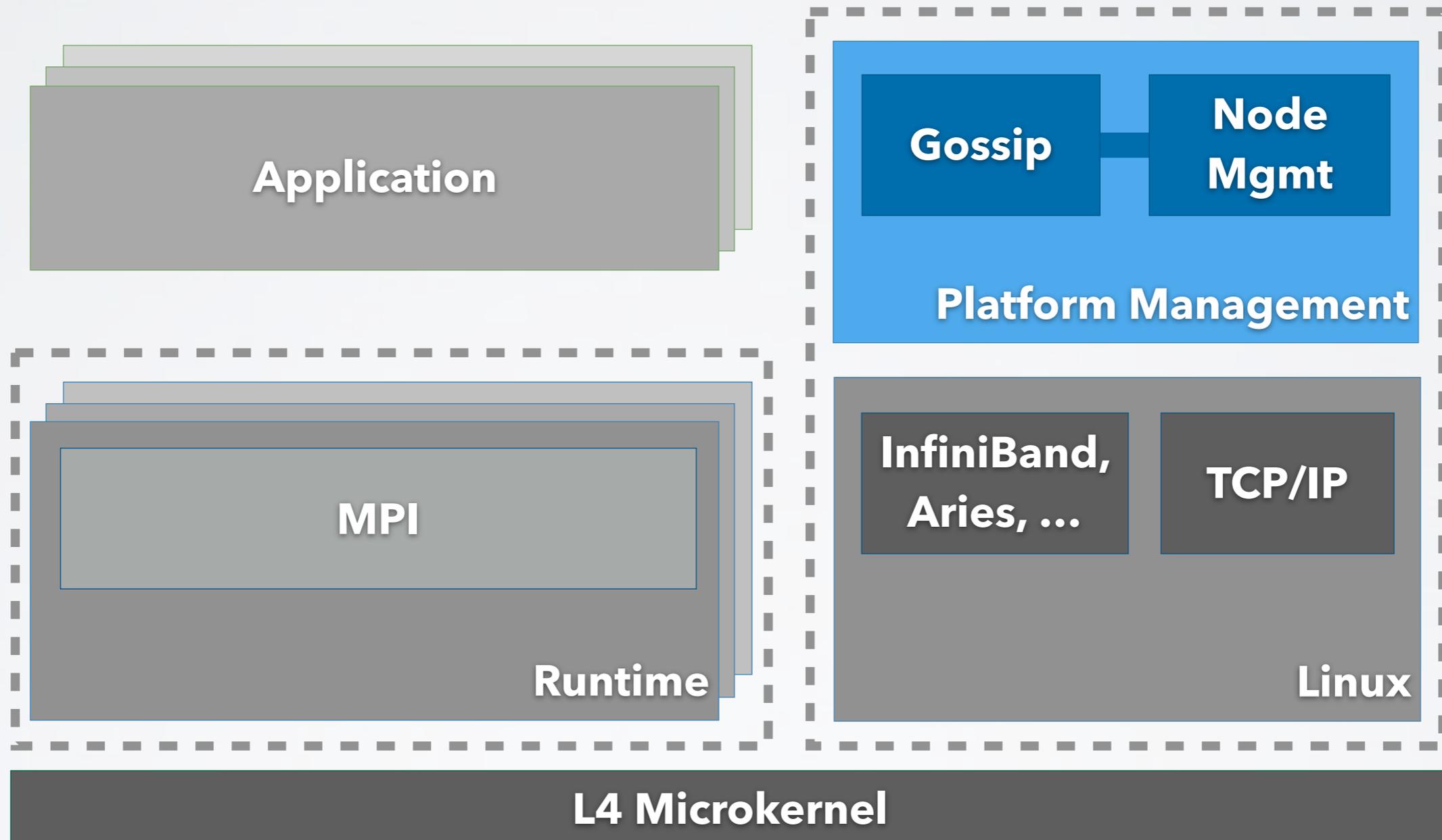
Diffusion

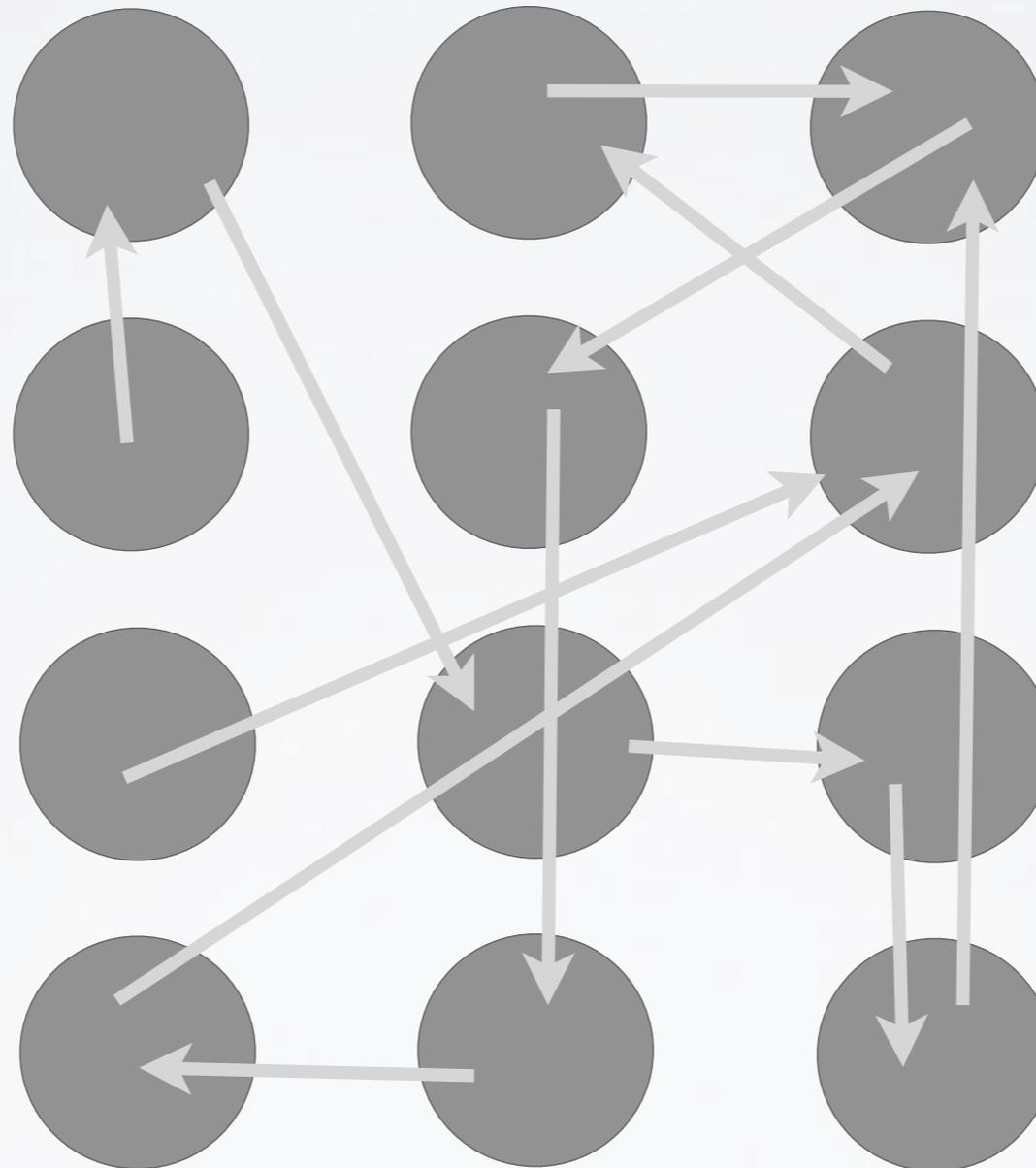


- **Best method to reduce:**
 - Migrations (less data movement)
 - Edge cut (less communication)
- **Load balance** good, but not superior
- **Flexible:** uses communication graph specific to application



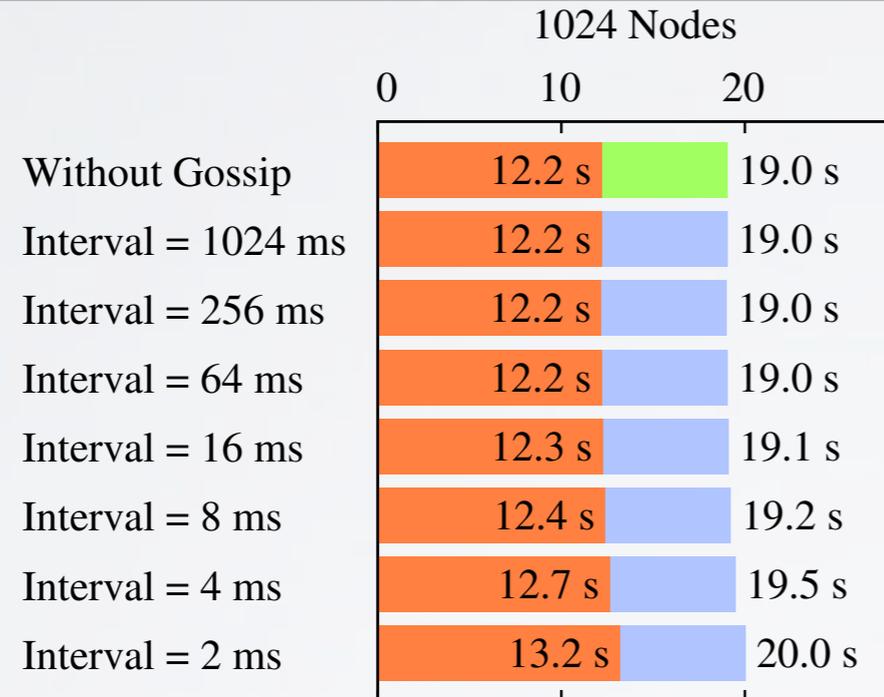
NODE ARCHITECTURE







GOSSIP SCALABILITY



MPI-FFT running on BG/Q "JUQUEEN"

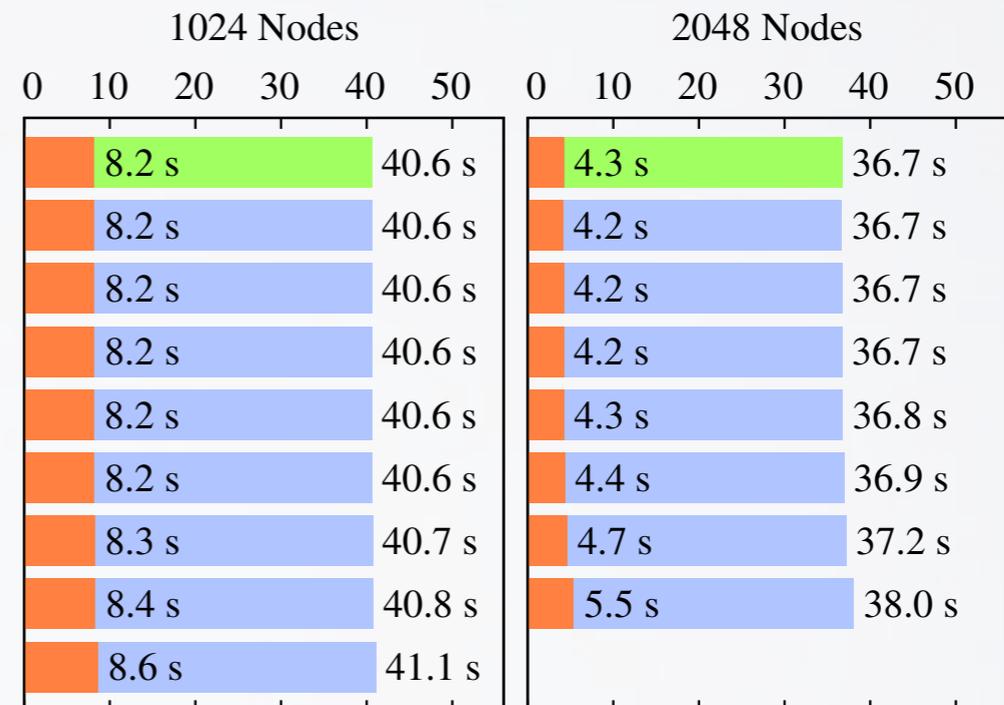
Low overhead:

No noticeable overhead at gossip interval of 64-256 ms

Without Gossip
Interval = 1024 ms
Interval = 256 ms
Interval = 64 ms
Interval = 16 ms
Interval = 8 ms
Interval = 4 ms
Interval = 2 ms
Interval = 1 ms

Quality of information:

Average age at nodes in the order of 2-3 s with gossip interval of 256 ms



COSMO-SPECS+FD4 on BG/Q "JUQUEEN"

E. Levy, A. Barak, A. Shiloh, M. Lieber, C. Weinhold, and H. Härtig, „Overhead of a Decentralized Gossip Algorithm on the Performance of HPC Applications“, ROSS 2014

Number of failed nodes per colony	Circulating <i>local windows</i> of size				
	16	32	64	128	256
0	11.74	9.67	8.66	8.20	8.07
1	11.71	9.72	8.67	8.21	8.07
2	11.75	9.68	8.70	8.21	8.08
4	11.81	9.73	8.70	8.23	8.11
8	11.83	9.79	8.72	8.28	8.17
16	11.95	9.90	8.79	8.34	8.20
32	12.12	10.05	8.96	8.48	8.36
Standard deviation	0.49	0.42	0.37	0.36	0.36
Increase rate	3.2%	3.9%	3.5%	3.4%	3.6%

Average age at master (1024 nodes per colony)

Gossip is fault tolerant:

Only slight increase in average age when substantial number of nodes fail (up to 32 of 1024 in each colony)

A. Barak, Z. Drezner, E. Levy, M. Lieber, and A. Shiloh, „Resilient gossip algorithms for collecting online management information in exascale clusters“, *Concurrency and Computation: Practice and Experience*, 2015



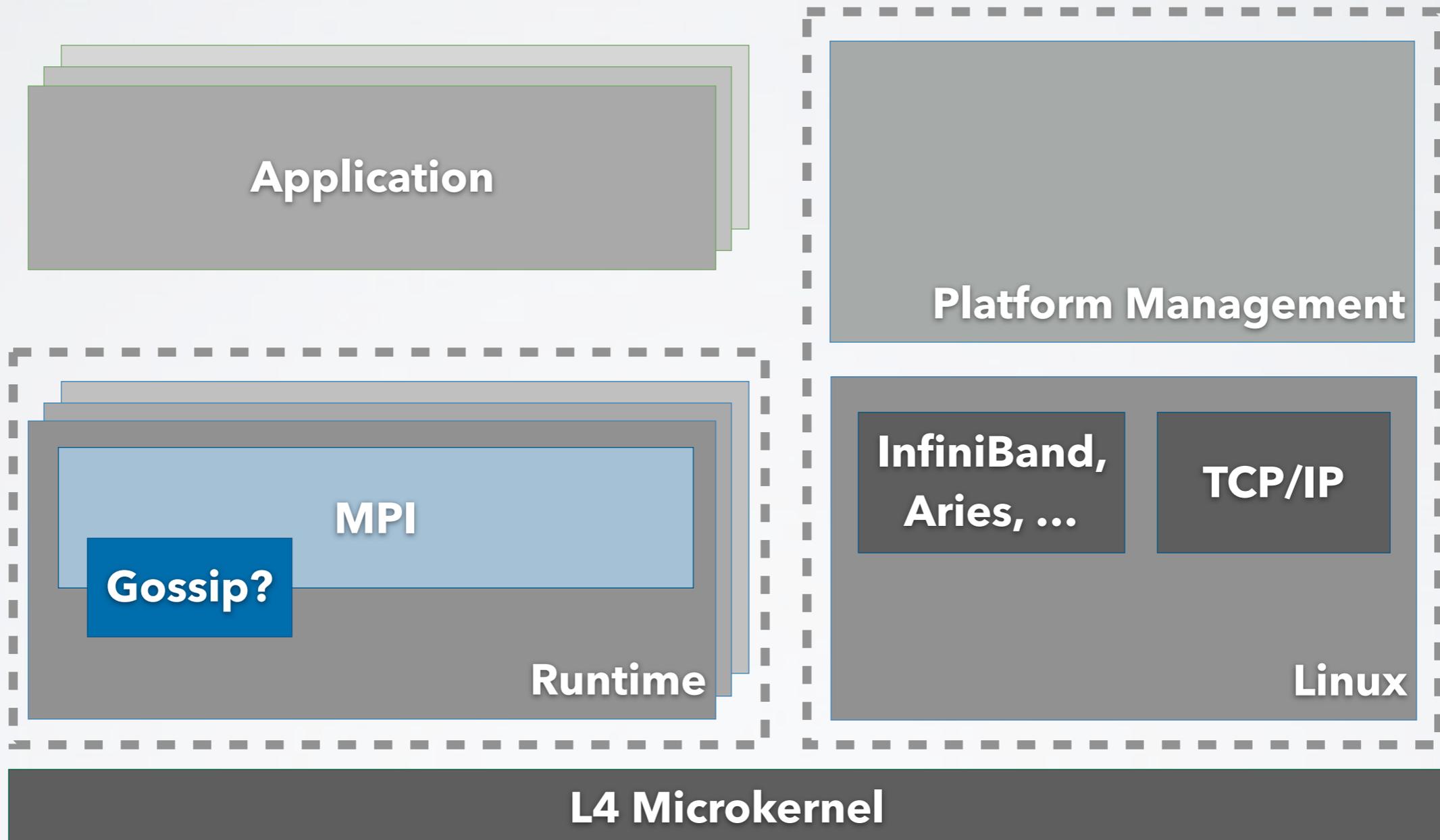
- Spread load+health info among nodes
- Analytic model ~ simulation ~ emulation
- **Negligible overhead** (64–256 ms intervals)
- **Good quality of information** (2–3 s old)
- **Fault tolerant** (simulated for up to 32 of 1024 nodes failing)

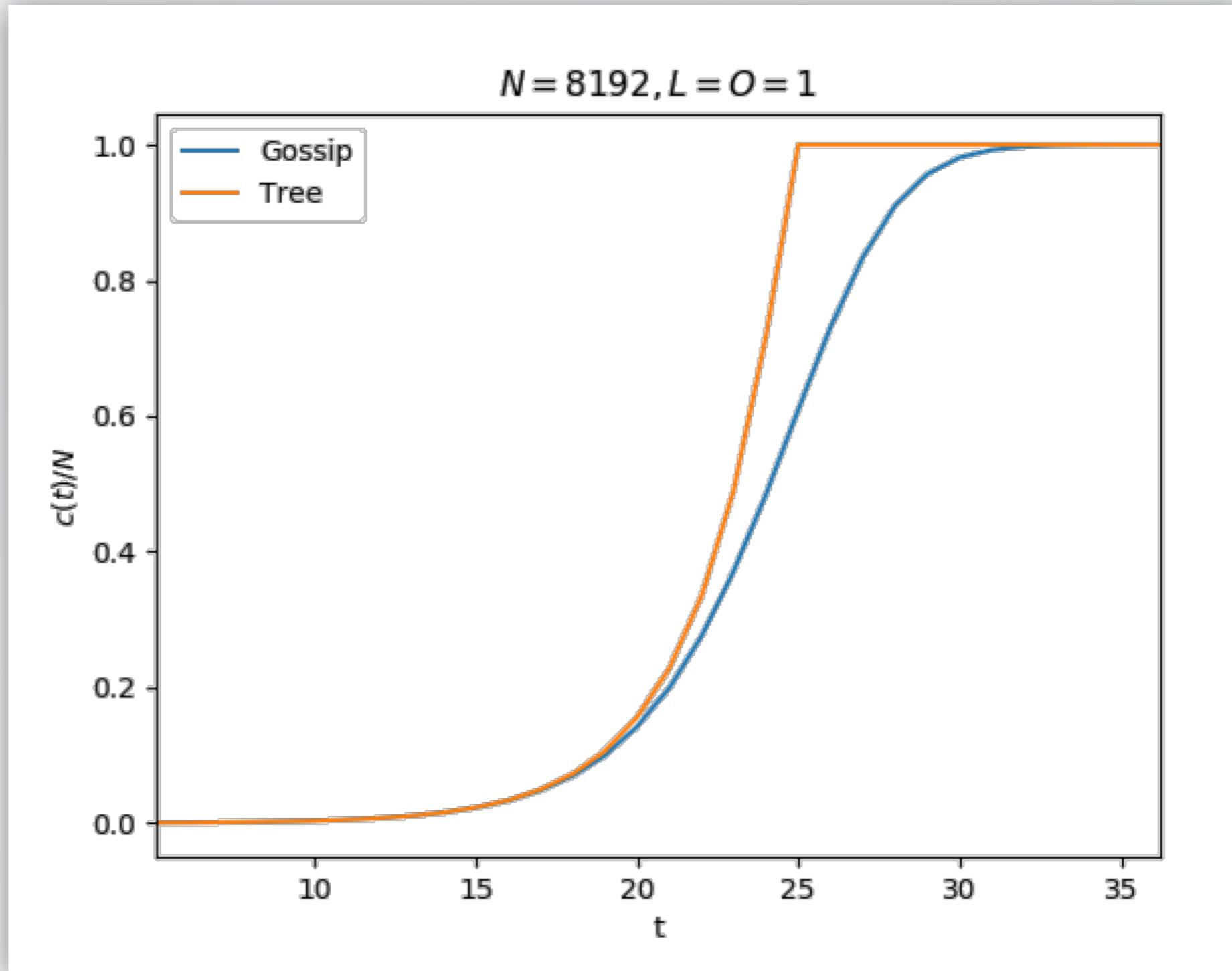
E. Levy, A. Barak, A. Shiloh, M. Lieber, C. Weinhold, and H. Härtig, „Overhead of a Decentralized Gossip Algorithm on the Performance of HPC Applications“, ROSS 2014

A. Barak, Z. Drezner, E. Levy, M. Lieber, and A. Shiloh, „Resilient Gossip Algorithms for Collecting Online Management information in Exascale Clusters“, Concurrency and Computation: Practice and Experience, 2015

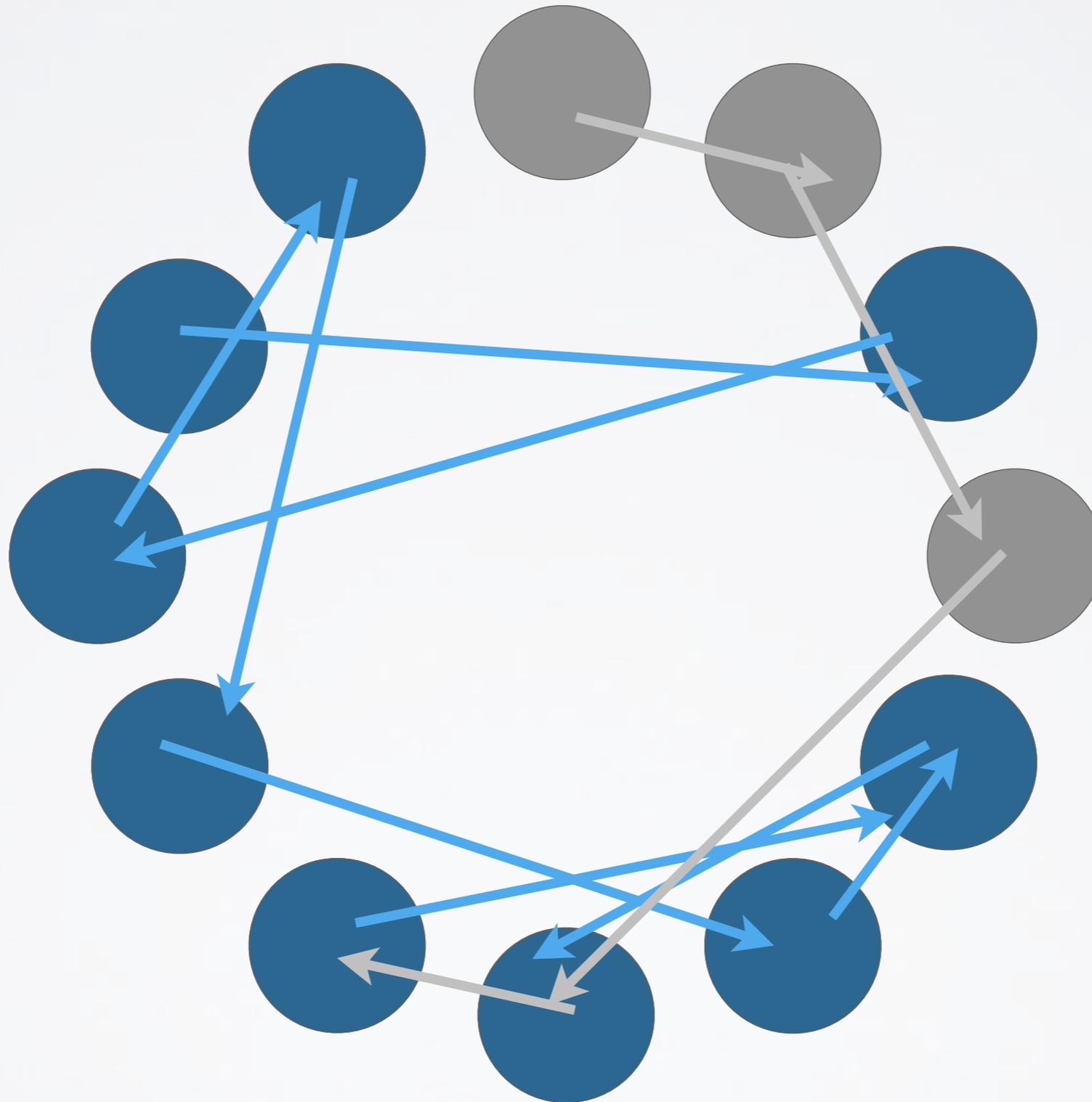


NODE ARCHITECTURE

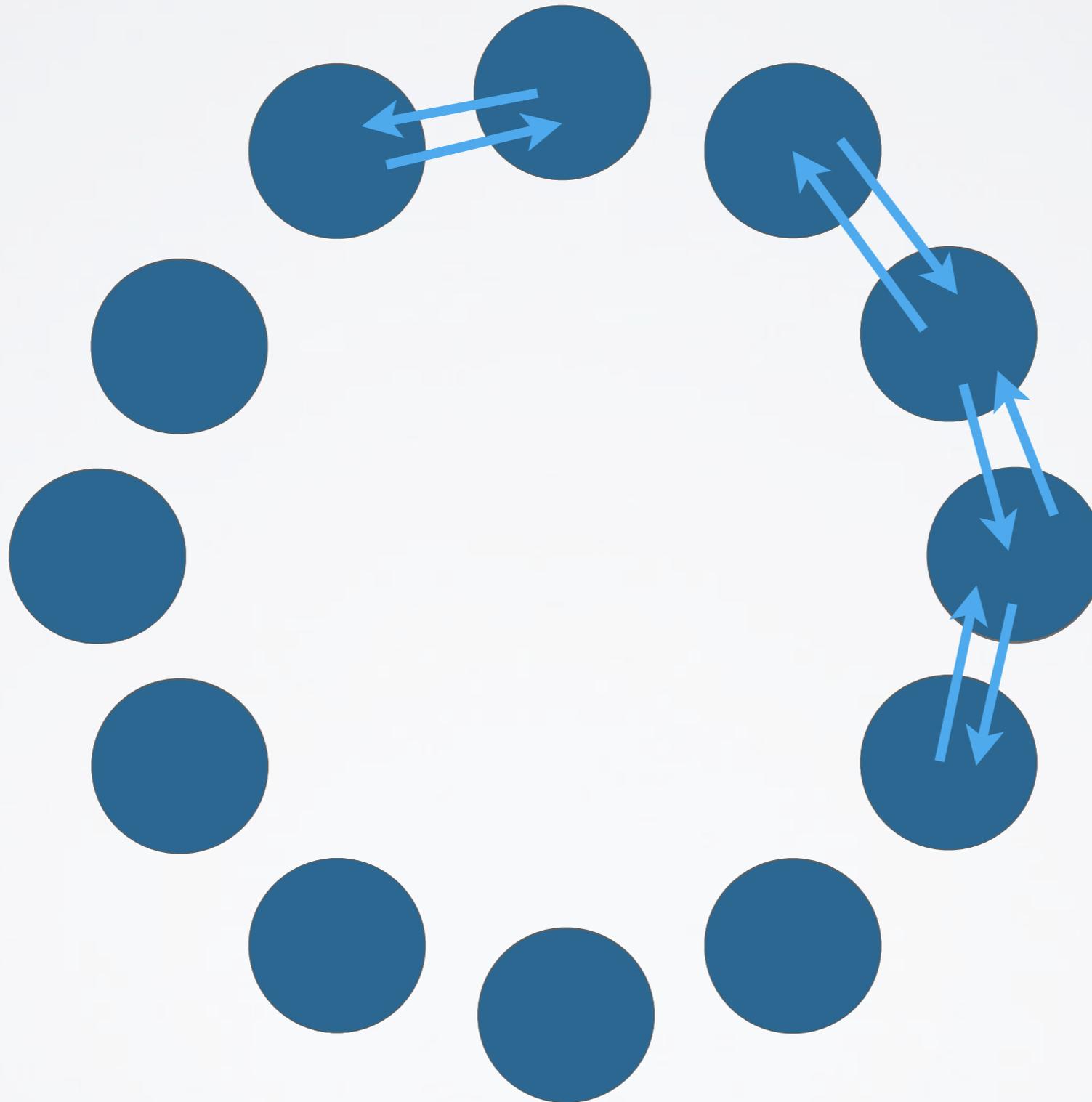




STEP 1: GOSSIP



STEP 2: CORRECTION



- Fault-tolerant **broadcast**: published^[*]
- Fault-tolerant **Reduce + Allreduce**, collectives with builtin **fault-detection**
 - Formal analysis, measurements show: log-scalable, sturdy in most cases
- Resiliency for **tree-based collectives**:
 - Succeed / complete with failing nodes
 - Latency comparable to non-ft algorithms

[*] Torsten Hoefler, Amnon Barak, Amnon Shiloh and Zvi Drezner, "Corrected Gossip Algorithms for Fast Reliable Broadcast on Unreliable Systems", IPDPS'17, Orlando, FL, USA

- **Decoupled execution:** low noise + latency
- **Checkpointing:** Coordinated + optimized
- **Diffusion:** Promising
- **Corrected Gossip & Trees:** fault-tolerant collective operations (maybe for MPI)
- **Integrated:** gossip + decision making
- **WIP:** integrate monitoring + migration