**TECHNISCHE UNIVERSITÄT DRESDEN**

ZIB

The Hebrew University of Jerusalem

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

Amnon Barak             Hebrew University Jerusalem (HUJI)
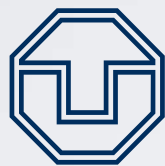Hermann Härtig          TU Dresden, Operating Systems Group (TUDOS)
Wolfgang E. Nagel       TU Dresden, Center for Information Services and HPC (ZIH)
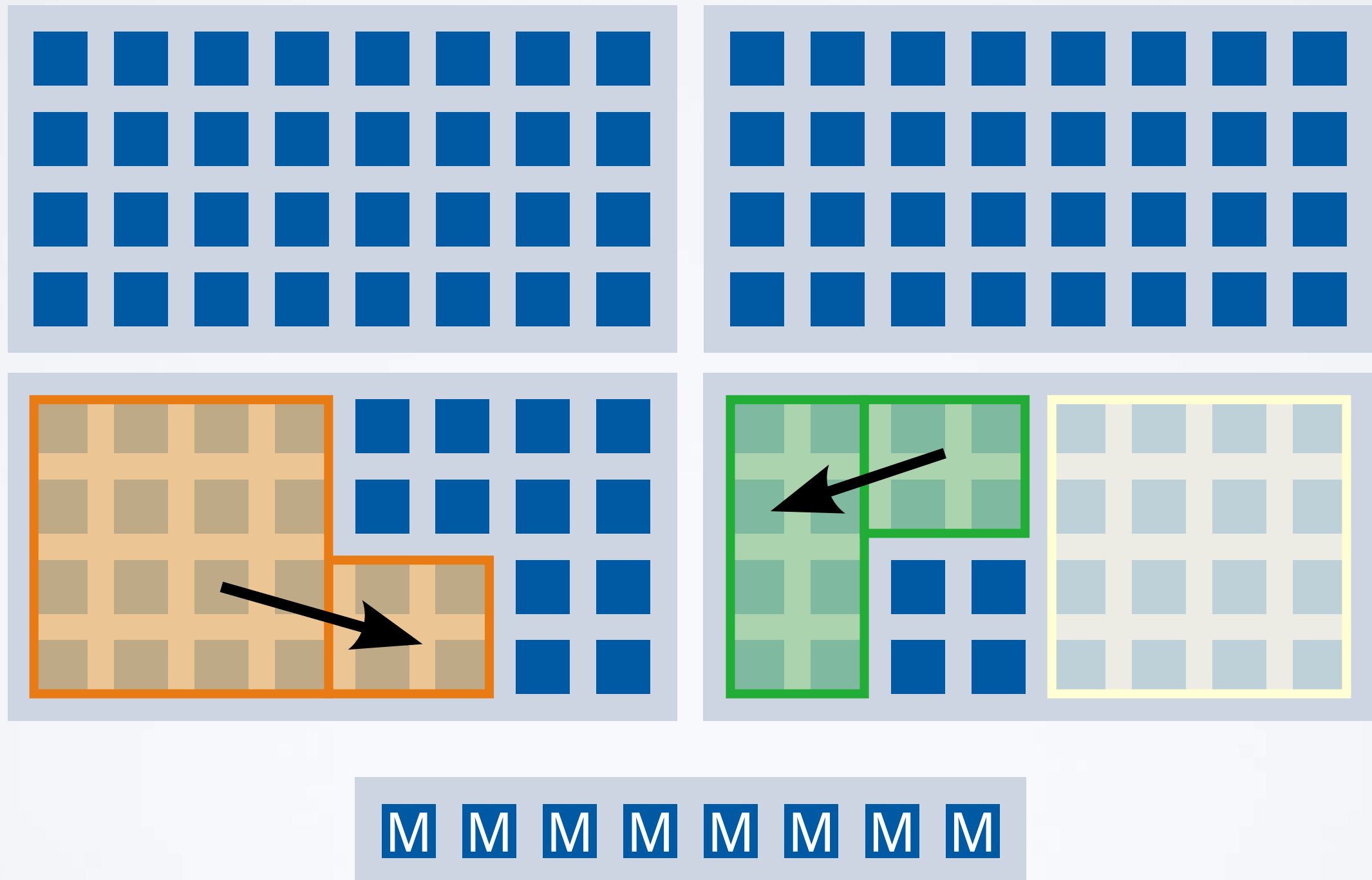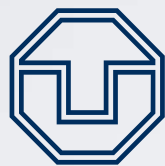Alexander Reinefeld     Konrad-Zuse-Zentrum für Informationstechnik Berlin (ZIB)
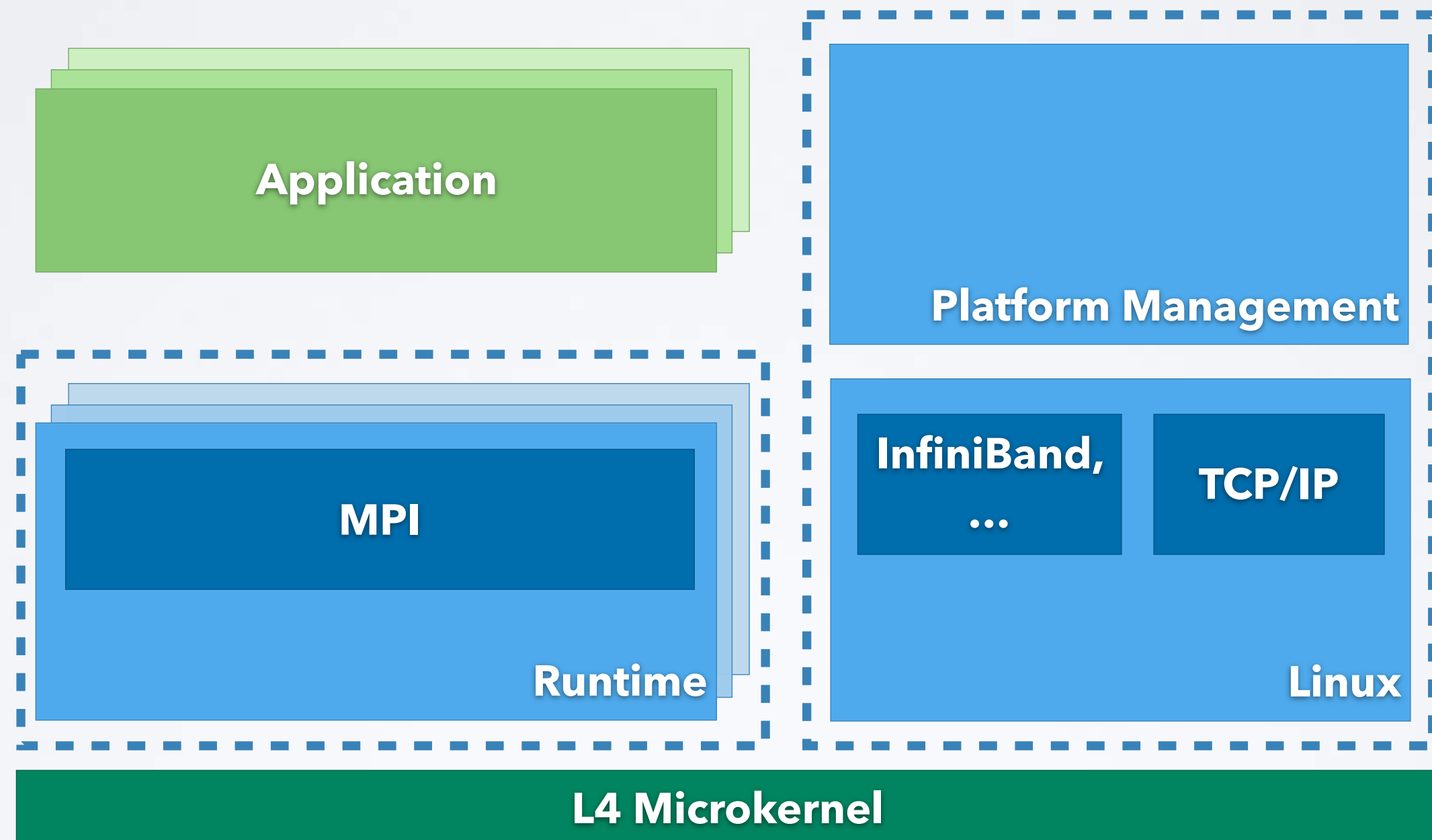
**CARSTEN WEINHOLD**, TU DRESDEN

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIB

The Hebrew University
of Jerusalem



Application

Platform Management

MPI

Runtime

InfiniBand, ...

TCP/IP

Linux

L4 Microkernel
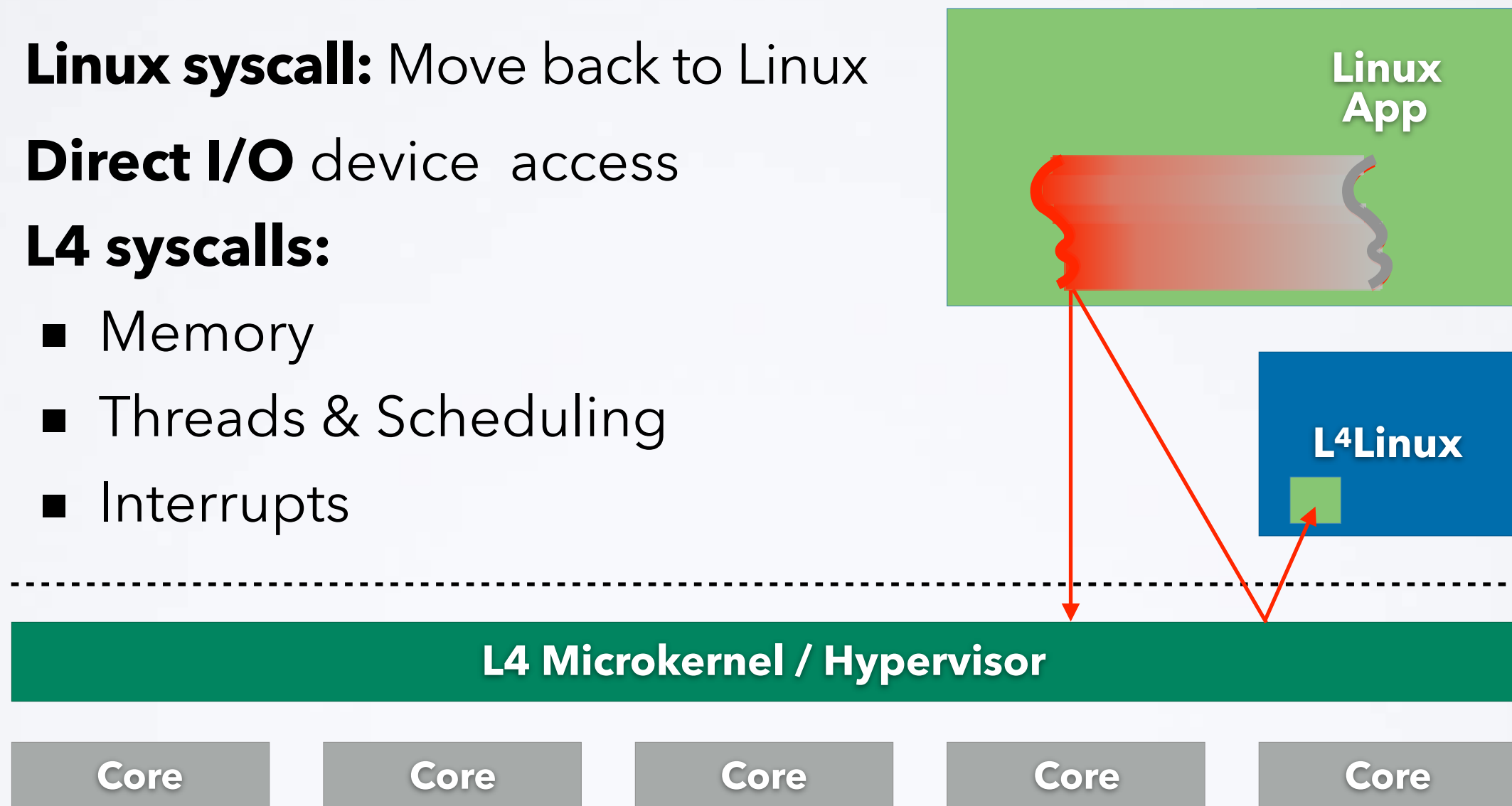
- **Decoupling:** move Linux thread to new L4 thread on its own core

- **Linux syscall:** Move back to Linux

- **Direct I/O** device access

- **L4 syscalls:**
  - Memory
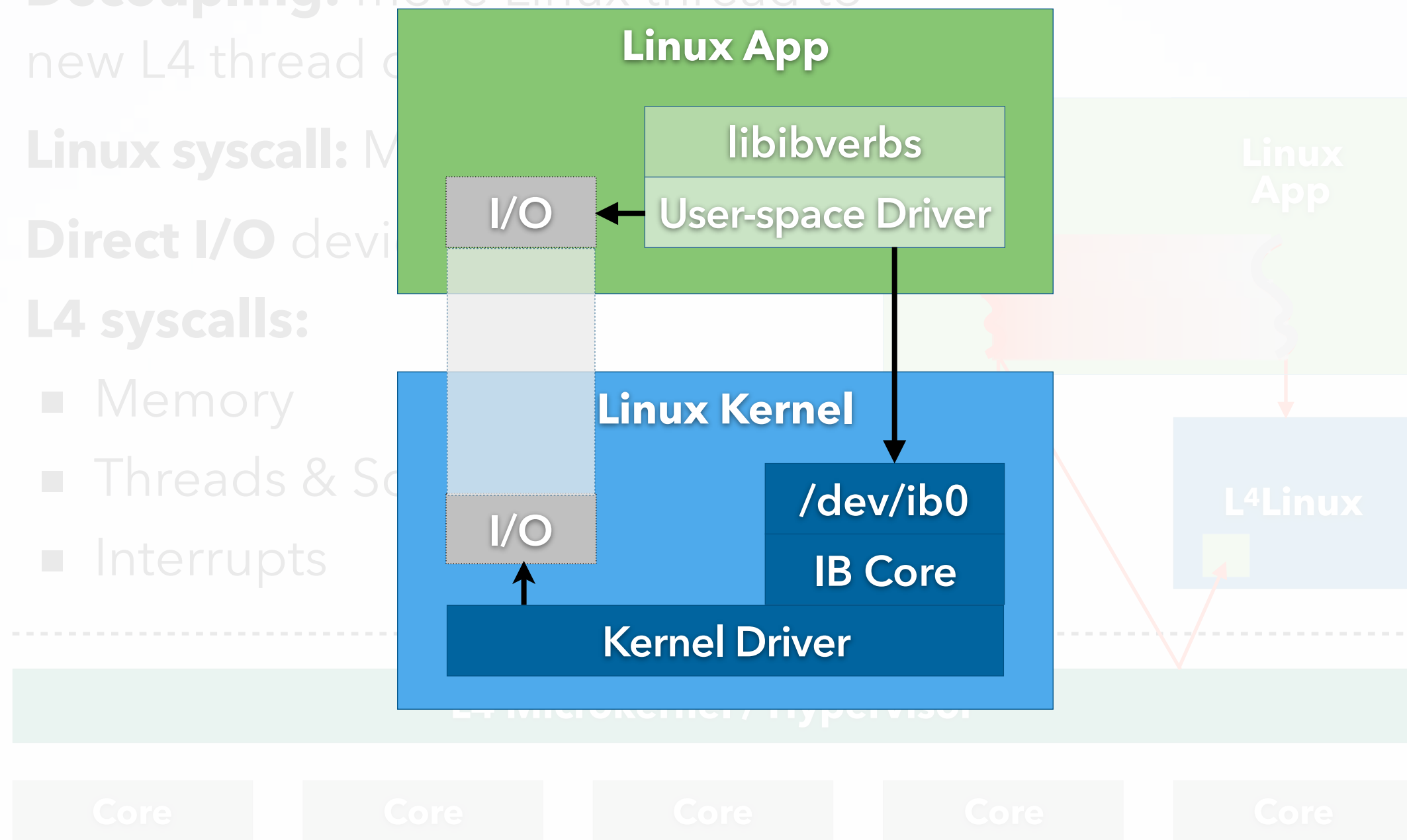  - Threads & Scheduling
  - Interrupts

**Linux App**

**L⁴Linux**

**L4 Microkernel / Hypervisor**

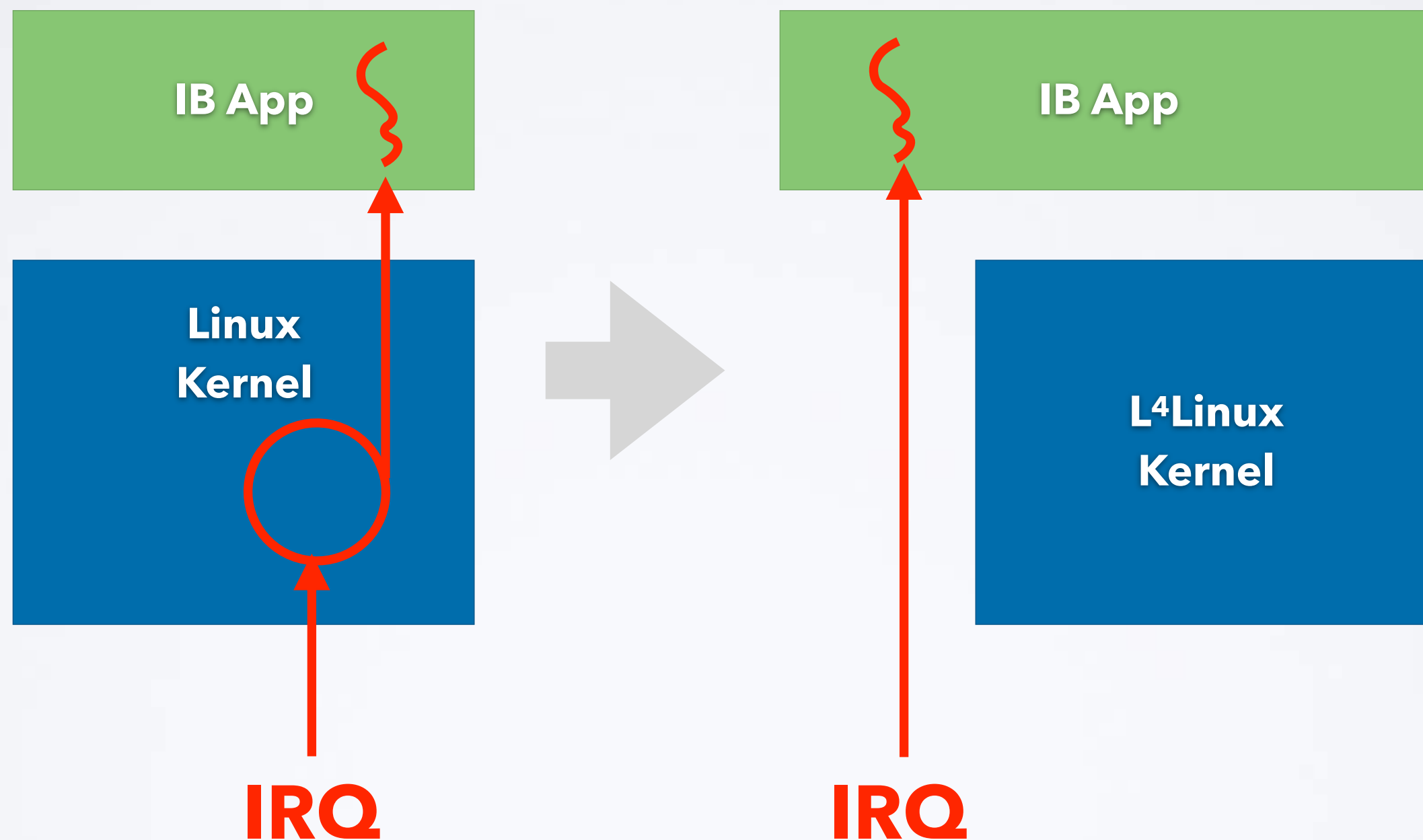| Core | Core | Core | Core | Core |
|------|------|------|------|------|

- **Decoupling:** move Linux thread to new L4 thread o
- **Linux syscall:** M
- **Direct I/O** devi
- **L4 syscalls:**
  - Memory
  - Threads & Sc
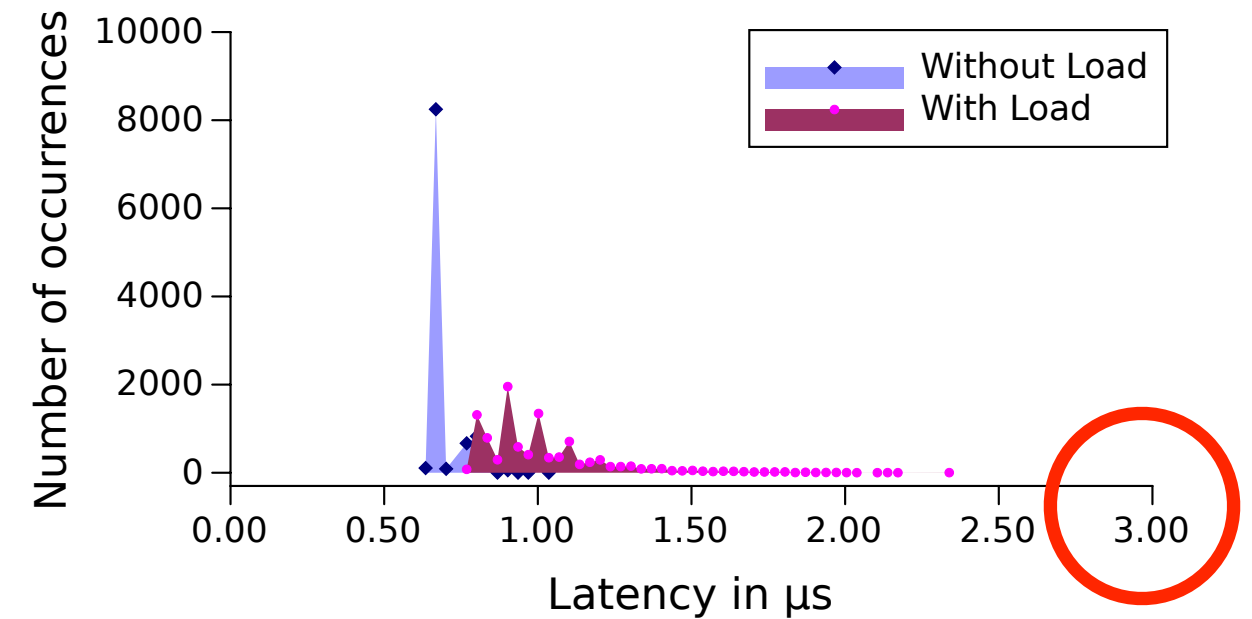  - Interrupts

**Linux App**

libibverbs

User-space Driver

I/O

**Linux Kernel**

/dev/ib0

IB Core

Kernel Driver

I/O

Linux App

L⁴Linux

Core    Core    Core    Core    Core

**IB App**

**Linux Kernel**

**IRQ**

**IB App**

**L⁴Linux Kernel**

**IRQ**

Linux

L4

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

- PhD student: internship at RIKEN, Japan

- Comparative study:
  - Hardware performance variation
  - 5 different CPU architectures
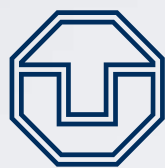  - Light-weight kernel (McKernel)

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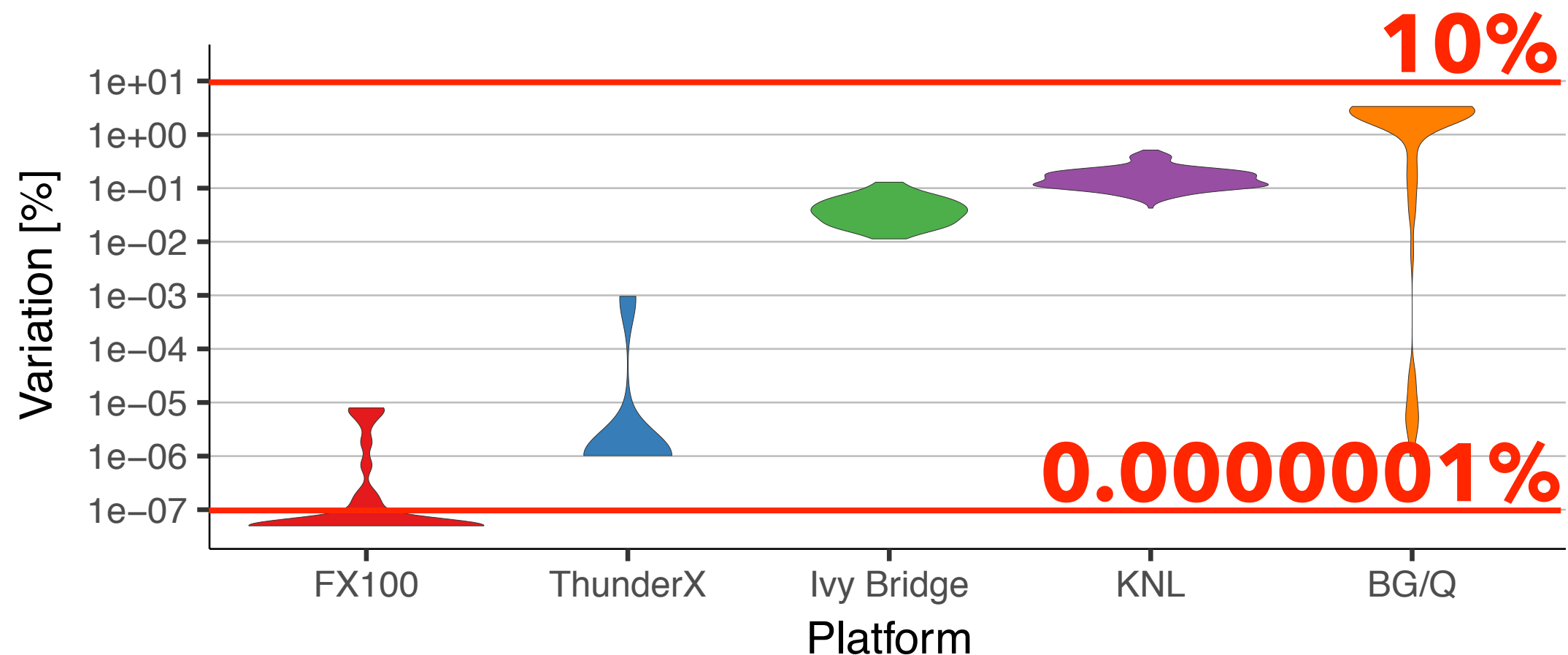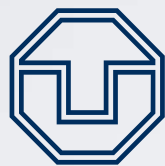
Hannes Weisbach, Brian Kocoloski, Hermann Härtig, Yutaka Ishikawa, Balazs Gerofi, „Hardware Performance Variation: A Comparative Study using Lightweight Kernels", ISC'18, Frankfurt, Germany, June 2018
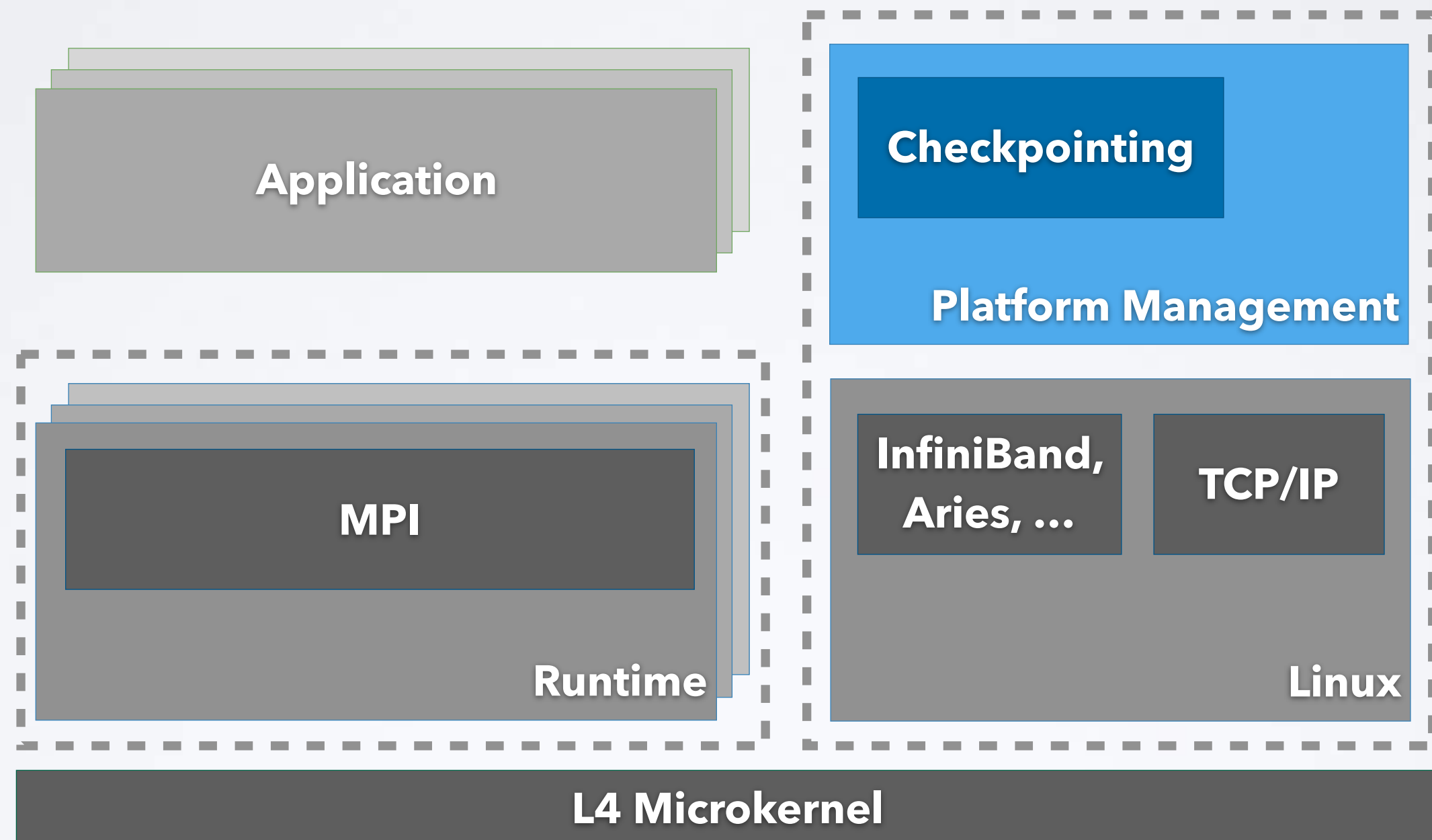
Hannes Weisbach, Brian Kocoloski, Hermann Härtig, Yutaka Ishikawa, Balazs Gerofi, „Hardware Performance Variation: A Comparative Study using Lightweight Kernels", ISC'18, Frankfurt, Germany, June 2018
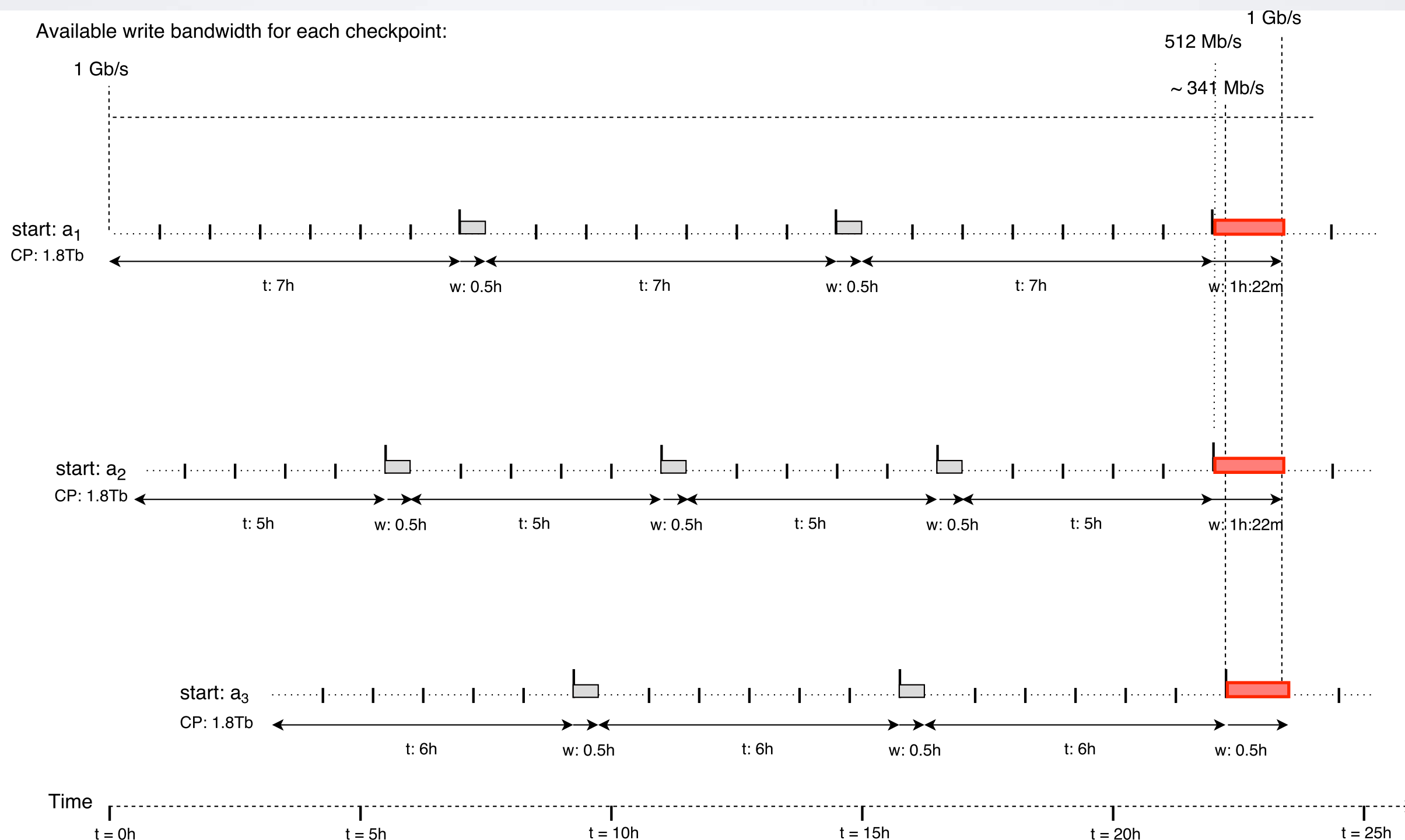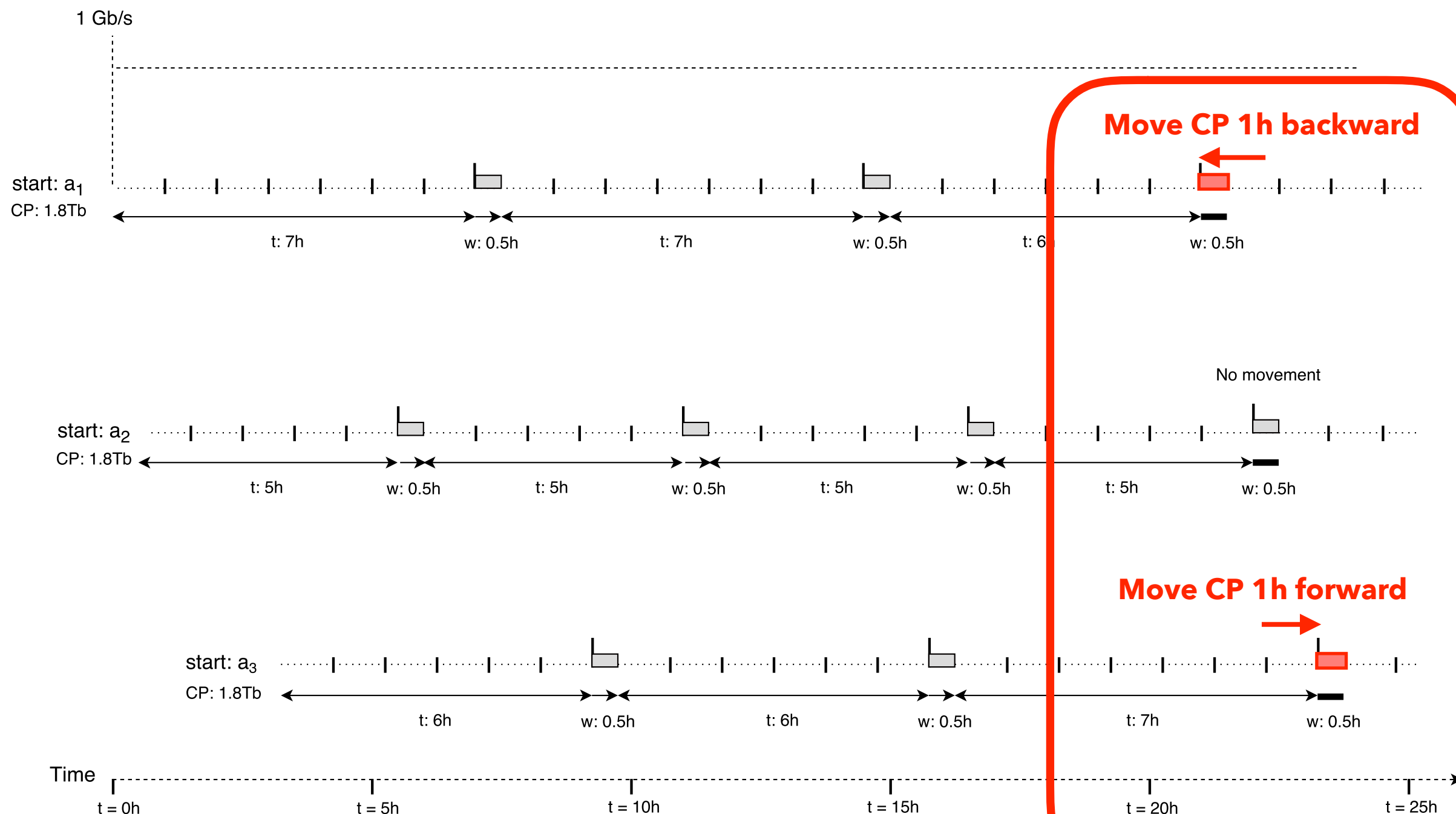
**TECHNISCHE UNIVERSITÄT DRESDEN**

ZIB

The Hebrew University of Jerusalem

Application

Checkpointing

**Platform Management**

MPI

**Runtime**

InfiniBand, Aries, ...

TCP/IP

**Linux**

**L4 Microkernel**

Available write bandwidth for each checkpoint:

1 Gb/s

**Move CP 1h backward**

start: $a_1$
CP: 1.8Tb

t: 7h          w: 0.5h          t: 7h          w: 0.5h          t: 6          w: 0.5h

No movement

start: $a_2$
CP: 1.8Tb

t: 5h          w: 0.5h          t: 5h          w: 0.5h          t: 5h          w: 0.5h          t: 5h          w: 0.5h

**Move CP 1h forward**

start: $a_3$
CP: 1.8Tb

t: 6h          w: 0.5h          t: 6h          w: 0.5h          t: 7h          w: 0.5h

Time

t = 0h          t = 5h          t = 10h          t = 15h          t = 20h          t = 25h

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

*Level 1:*
Checkpoints on the local Storage of node 1

*Level 2:*
XOR parity data of node 1 & node 2
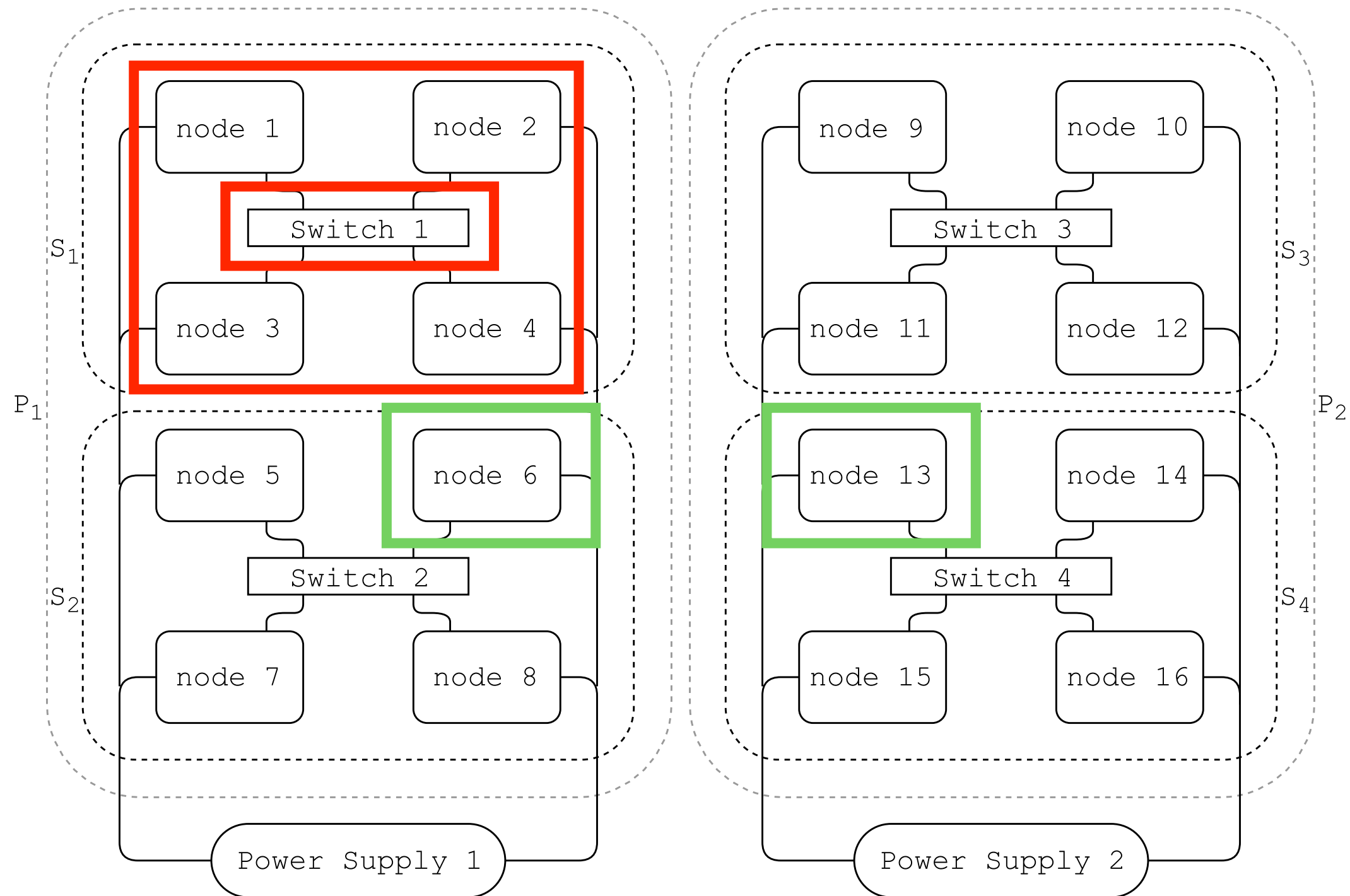
*Level 3:*
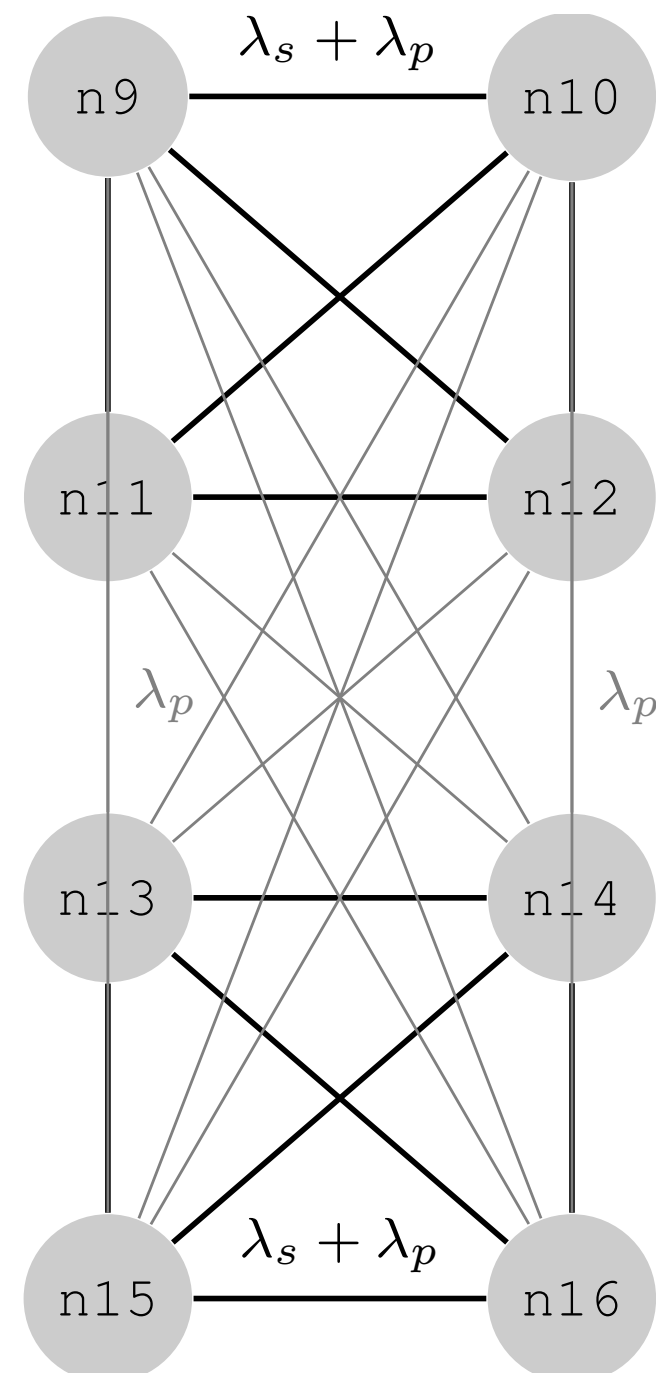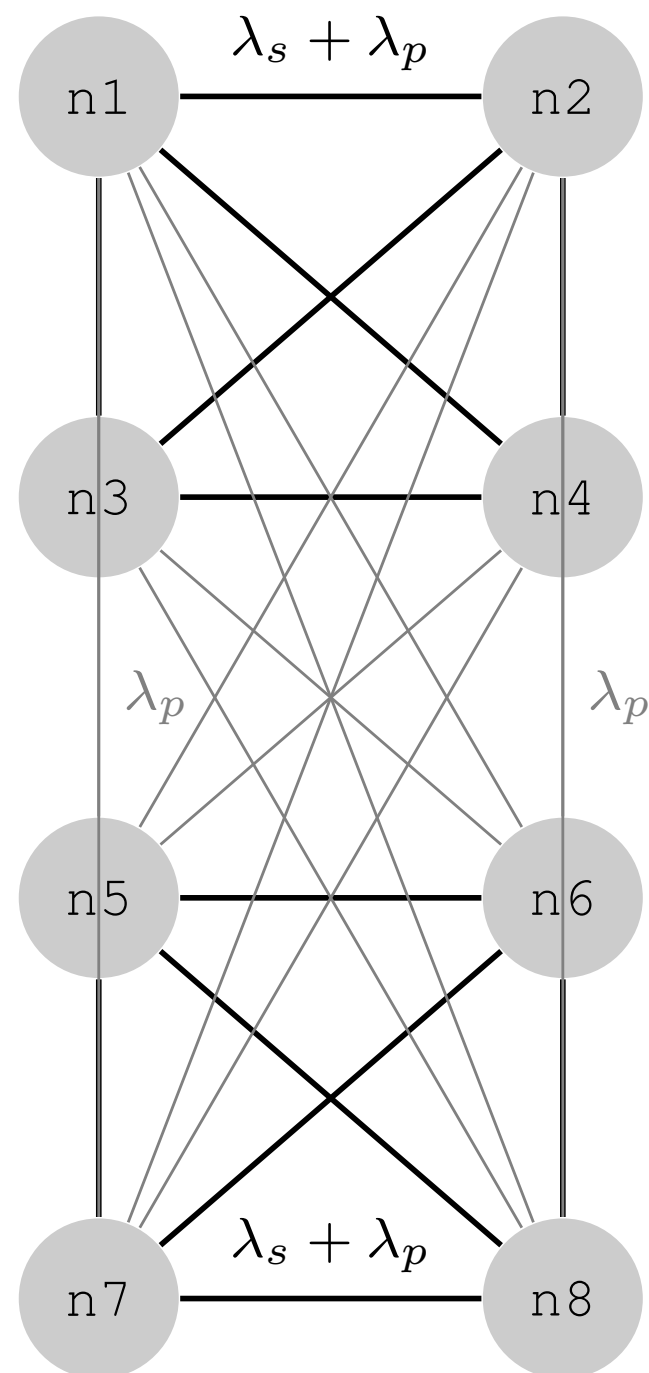Checkpoints on Storage of the Partner Node (node 2)
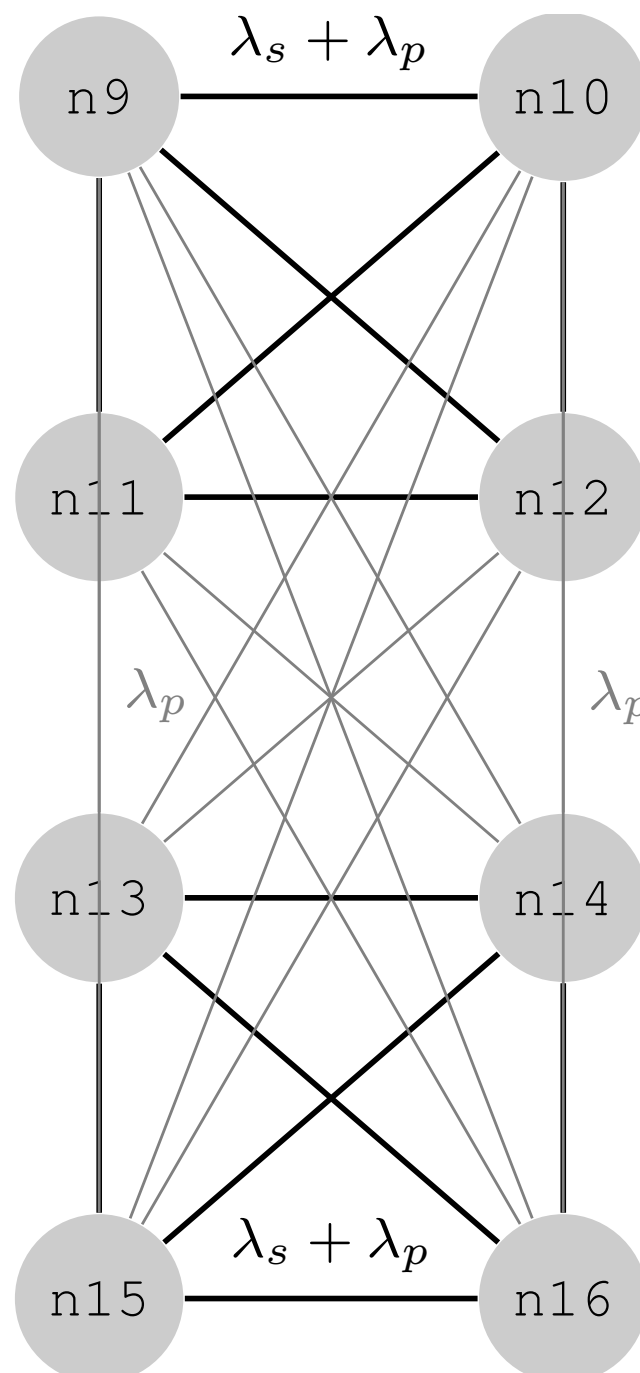
*Level 4:*
Checkpoints on the Shared Burst Buffer

*Level 5:*
Checkpoints on the Shared PFS
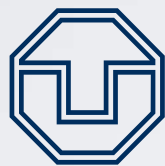
Node 1 Running job $j$, rank 0

Node 2 Running job $j$, rank 1

Partner Storage

Node 1 Local Storage

Node 2 Local Storage

XOR Parity Data

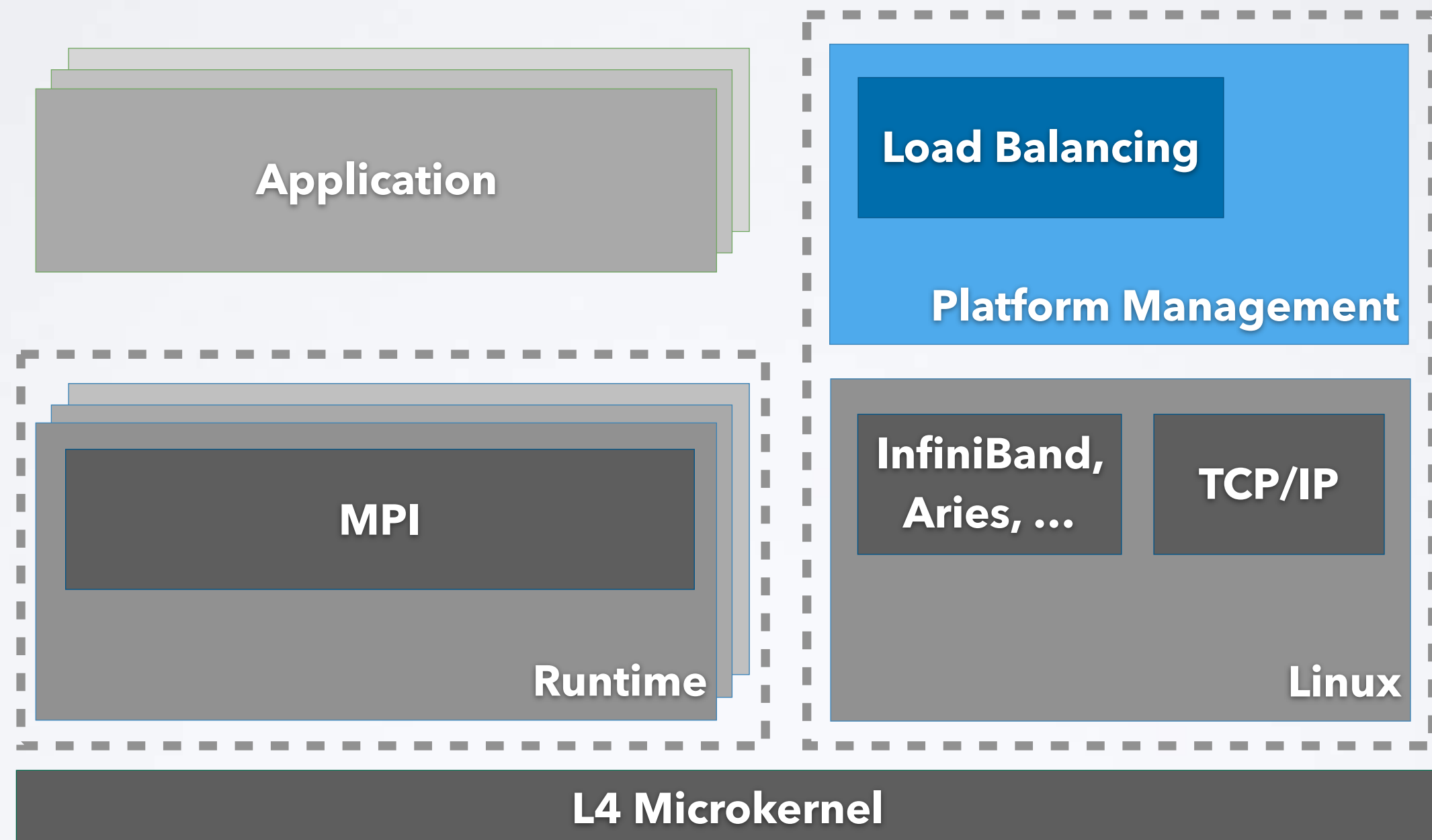Shared Burst Buffer

Parallel File System

**Graph problem:**

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

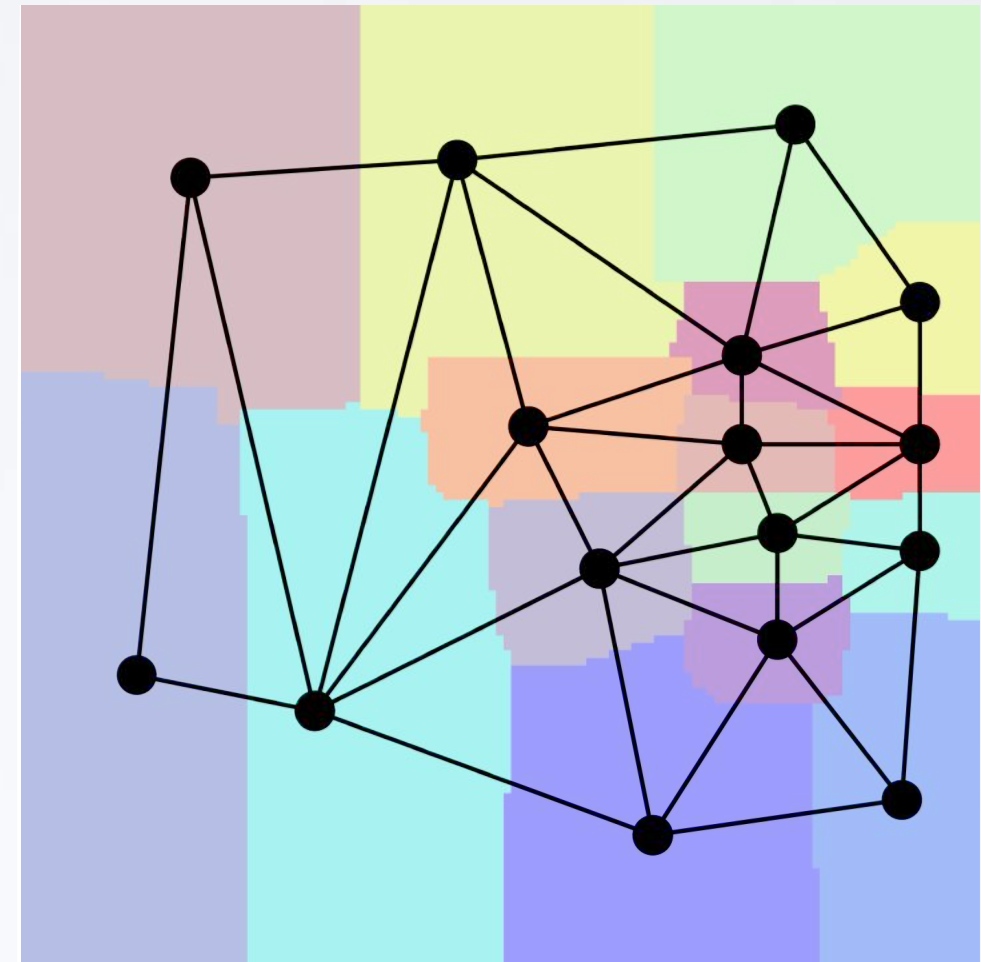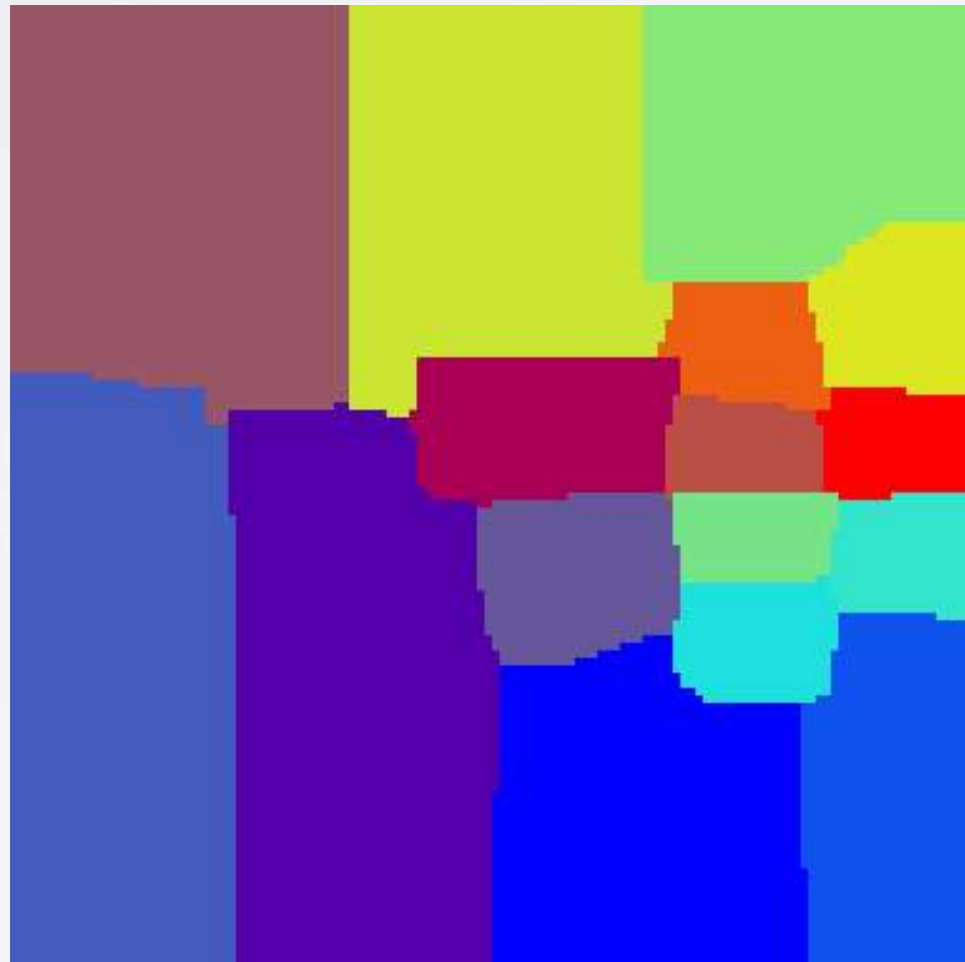**Optimization problem:**

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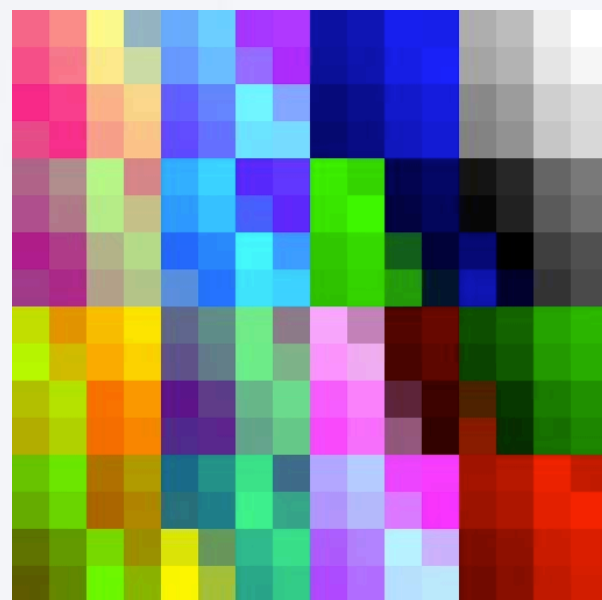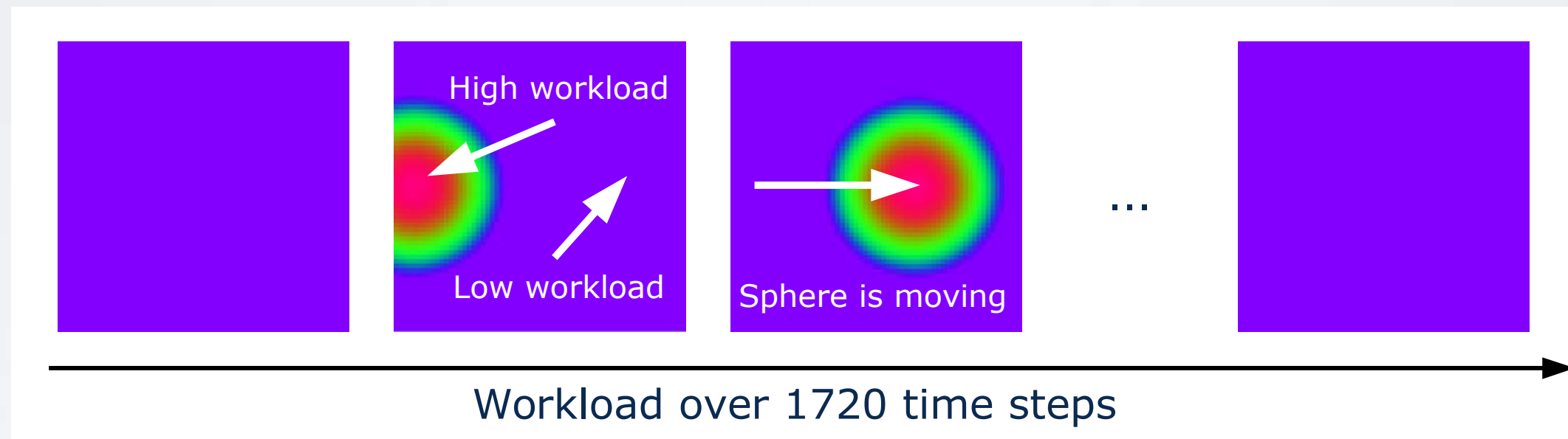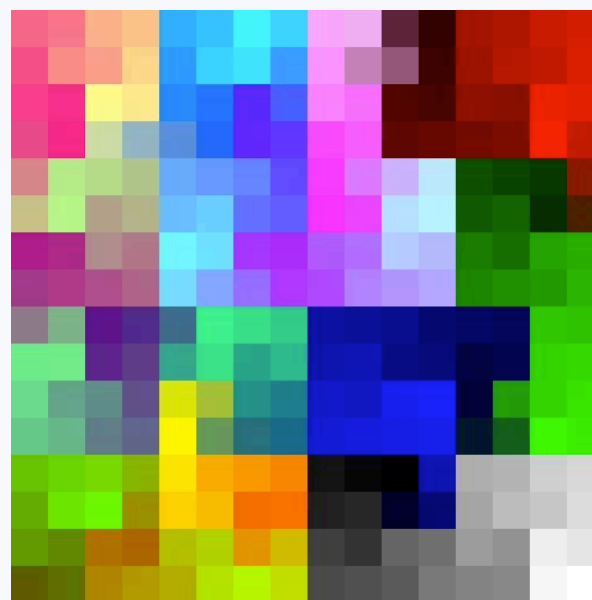- Consider: job run time, C/R cost, MTTI
- Minimize run time

Diffusion graph topology from application topology
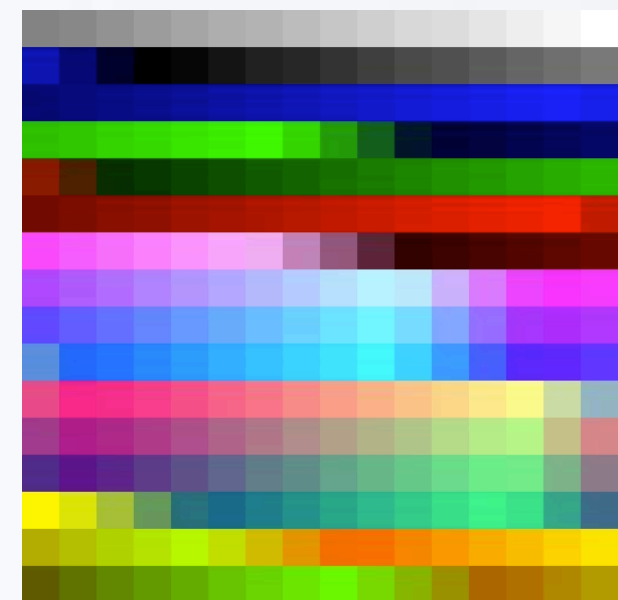
Diffusion coefficient weighted by interface length:
- Tasks migrated between neighbor partitions
- Better partition shape

TECHNISCHE UNIVERSITÄT DRESDEN

The Hebrew University of Jerusalem



High workload

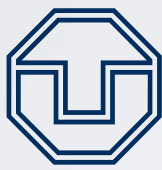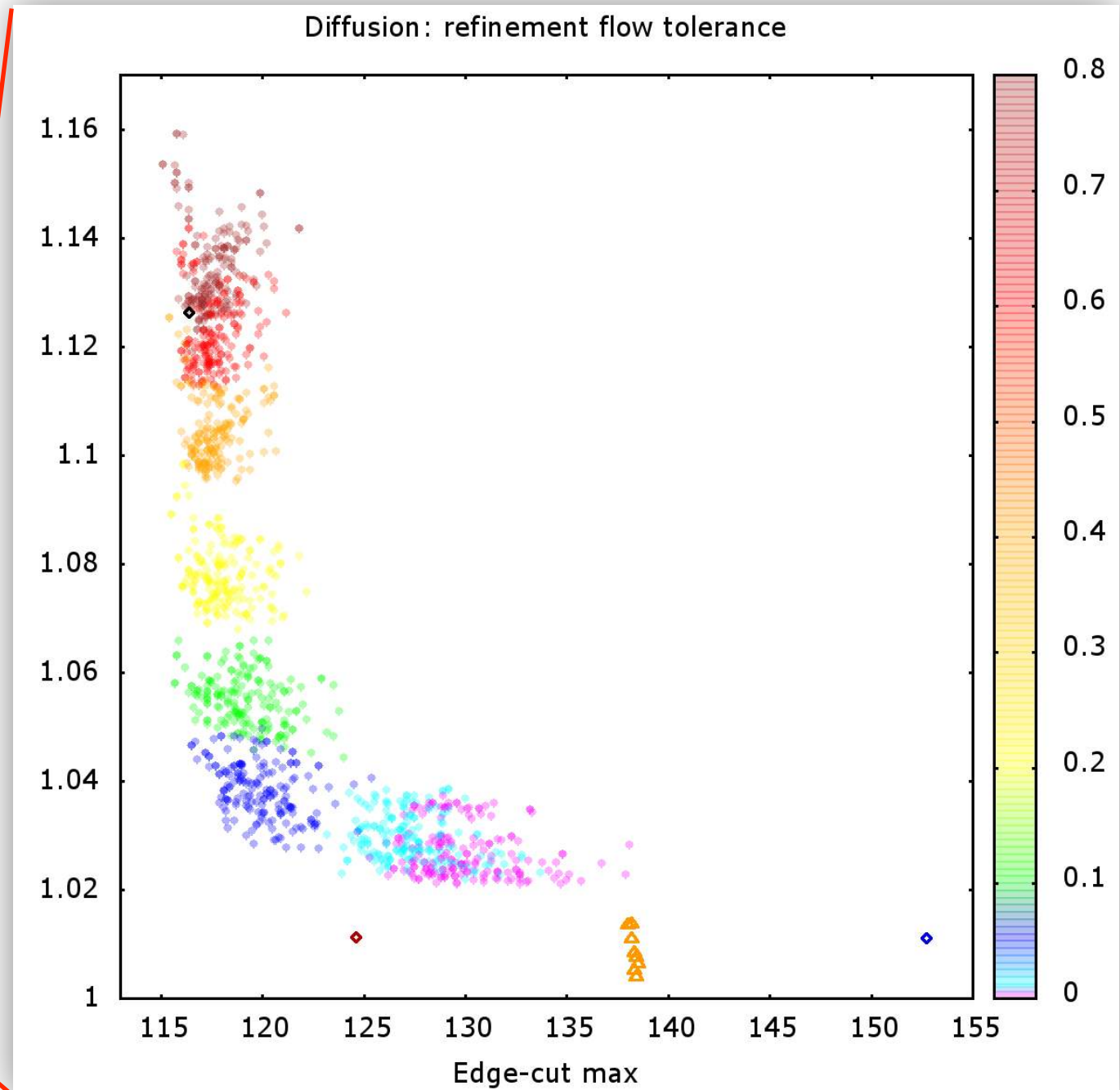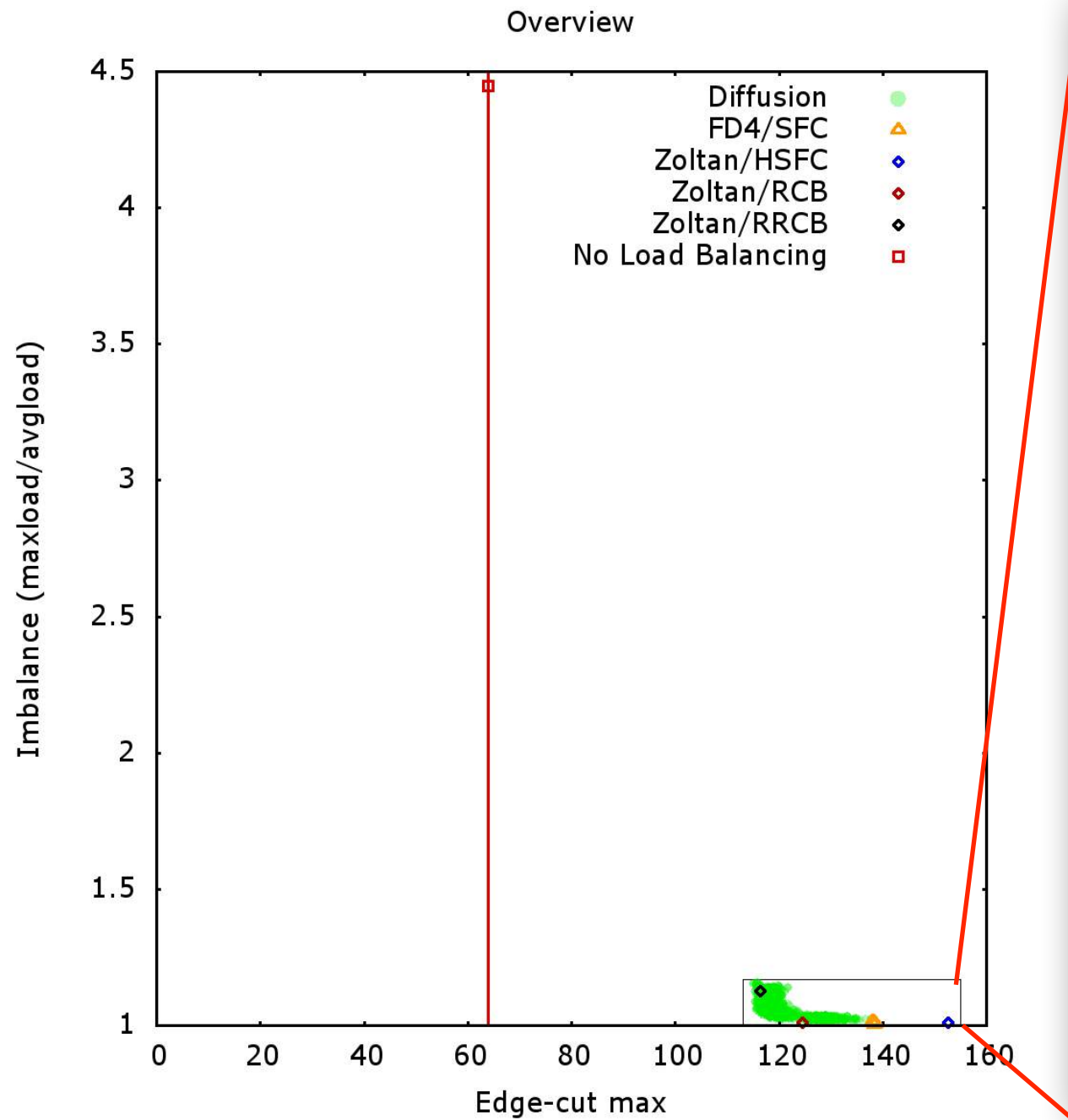Low workload

Sphere is moving

...

Workload over 1720 time steps



Zoltan



Space-filling Curves
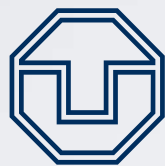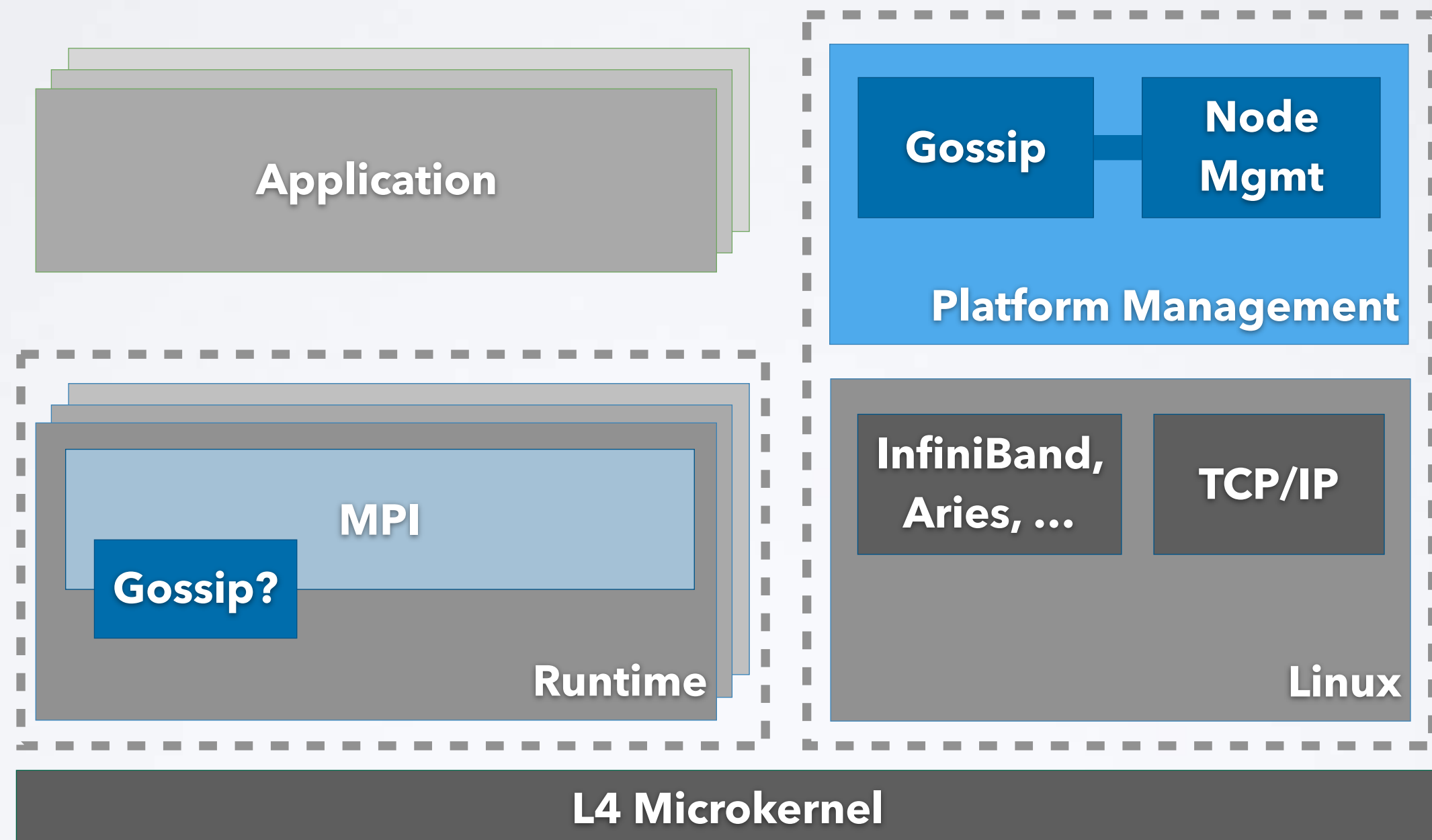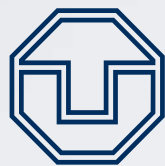


Diffusion

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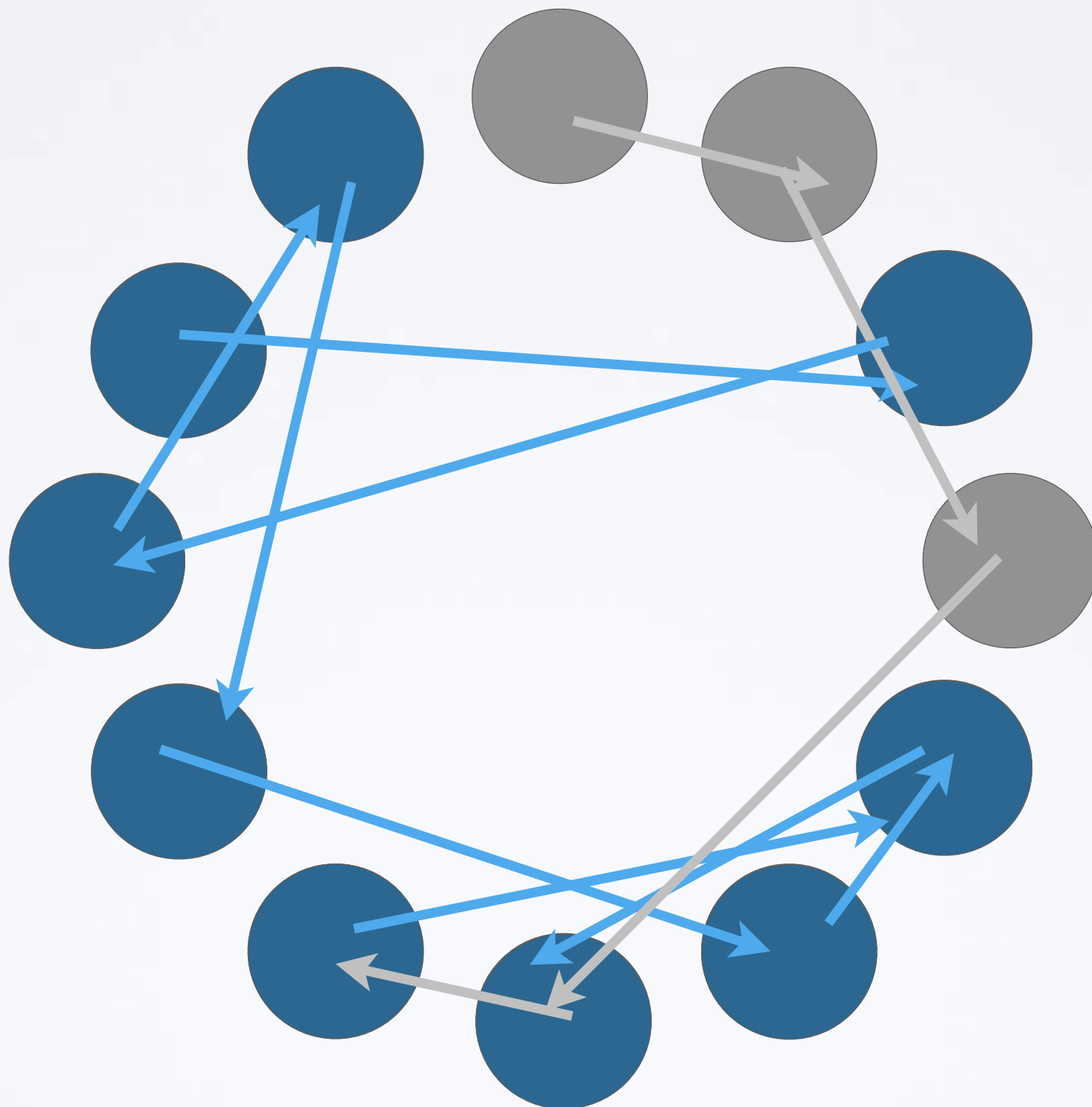- **Flexible:** uses communication graph specific to application

TECHNISCHE UNIVERSITÄT DRESDEN
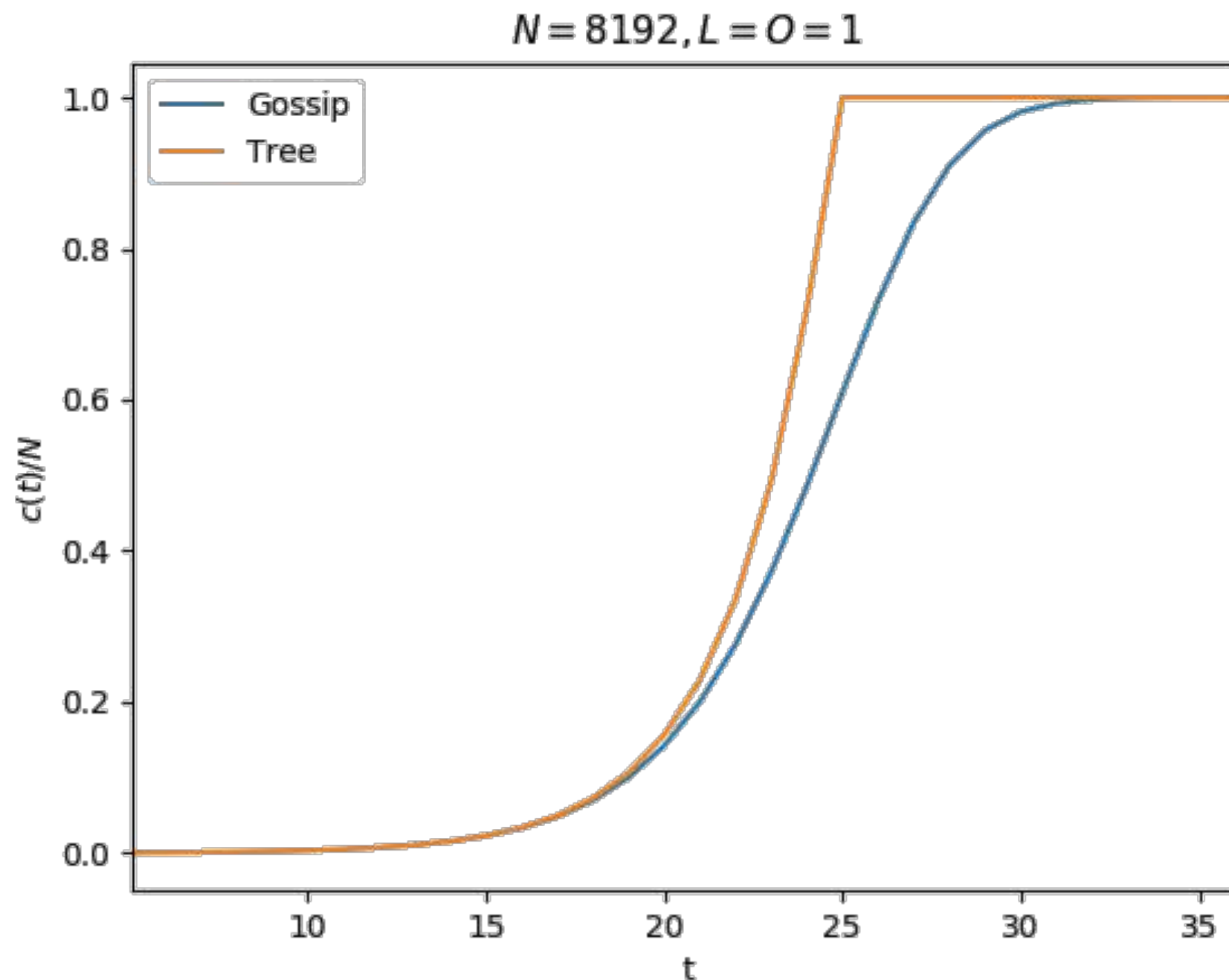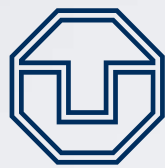
ZIB

The Hebrew University of Jerusalem

FFMK: Building an Exascale Operating System
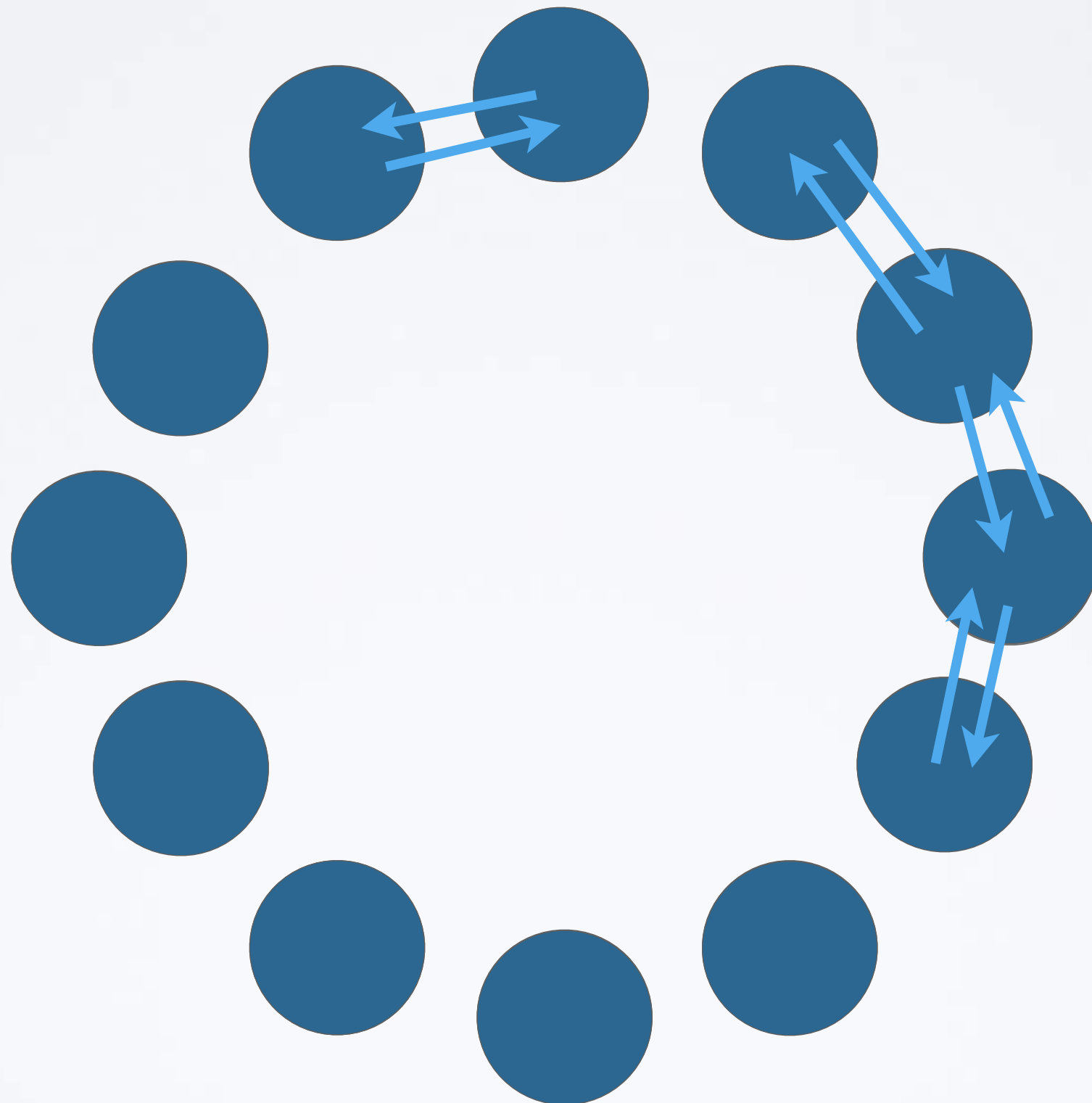
- Fault-tolerant **broadcast:** published[*]

- 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 interrupts:** faster wakeup

- **Checkpointing:** Global optimization

- **Diffusion:** Promising

- **Corrected Gossip & Trees:** fault-tolerant collective operations (maybe for MPI)

- **Integrated:** gossip + decision making

- **WIP:** integrate monitoring + migration