



**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)

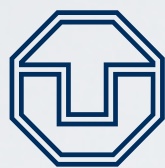
Wolfgang E. Nagel

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

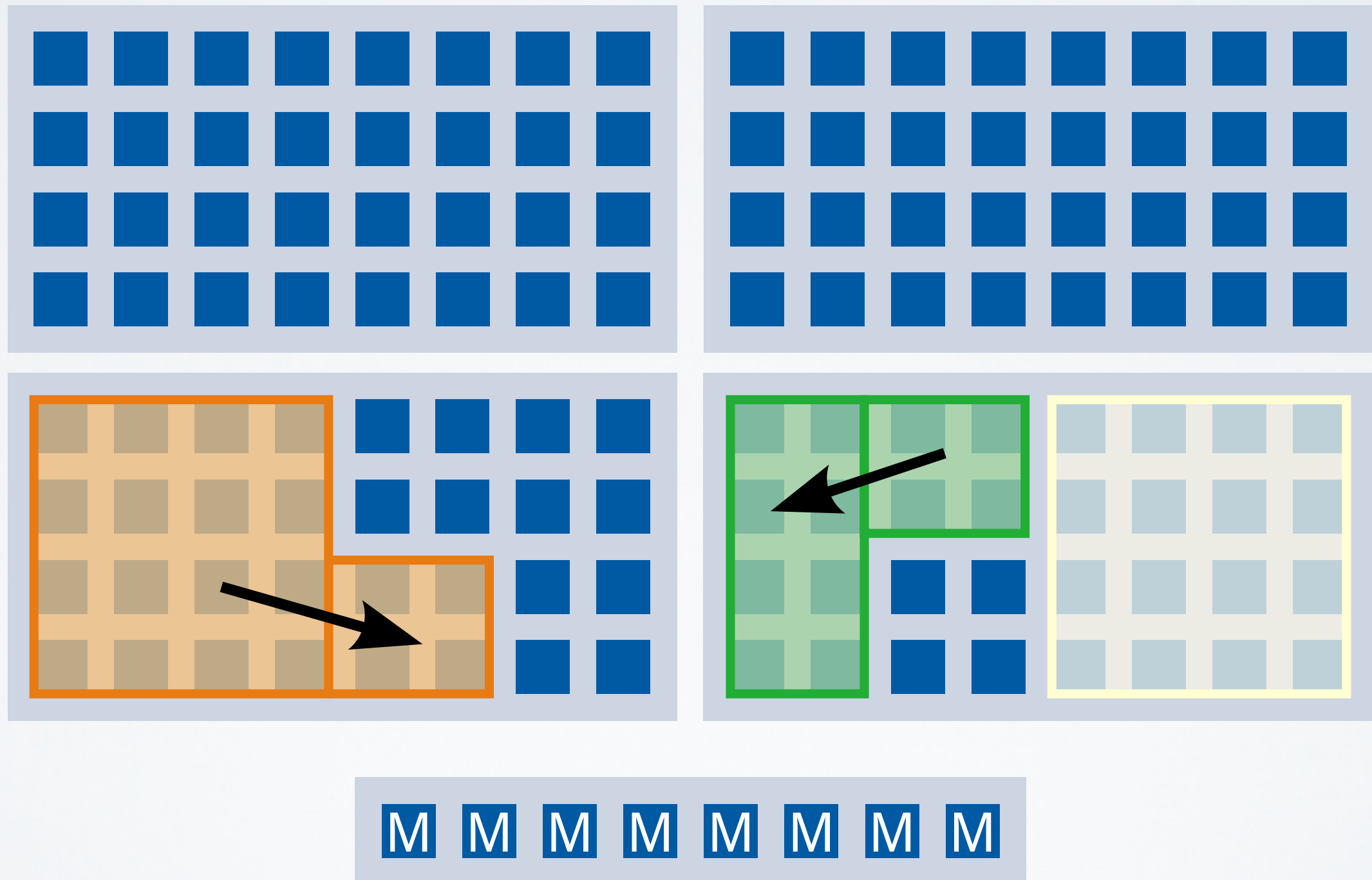
Alexander Reinefeld

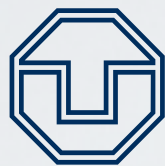
Konrad-Zuse-Zentrum für Informationstechnik Berlin (ZIB)

CARSTEN WEINHOLD, TU DRESDEN

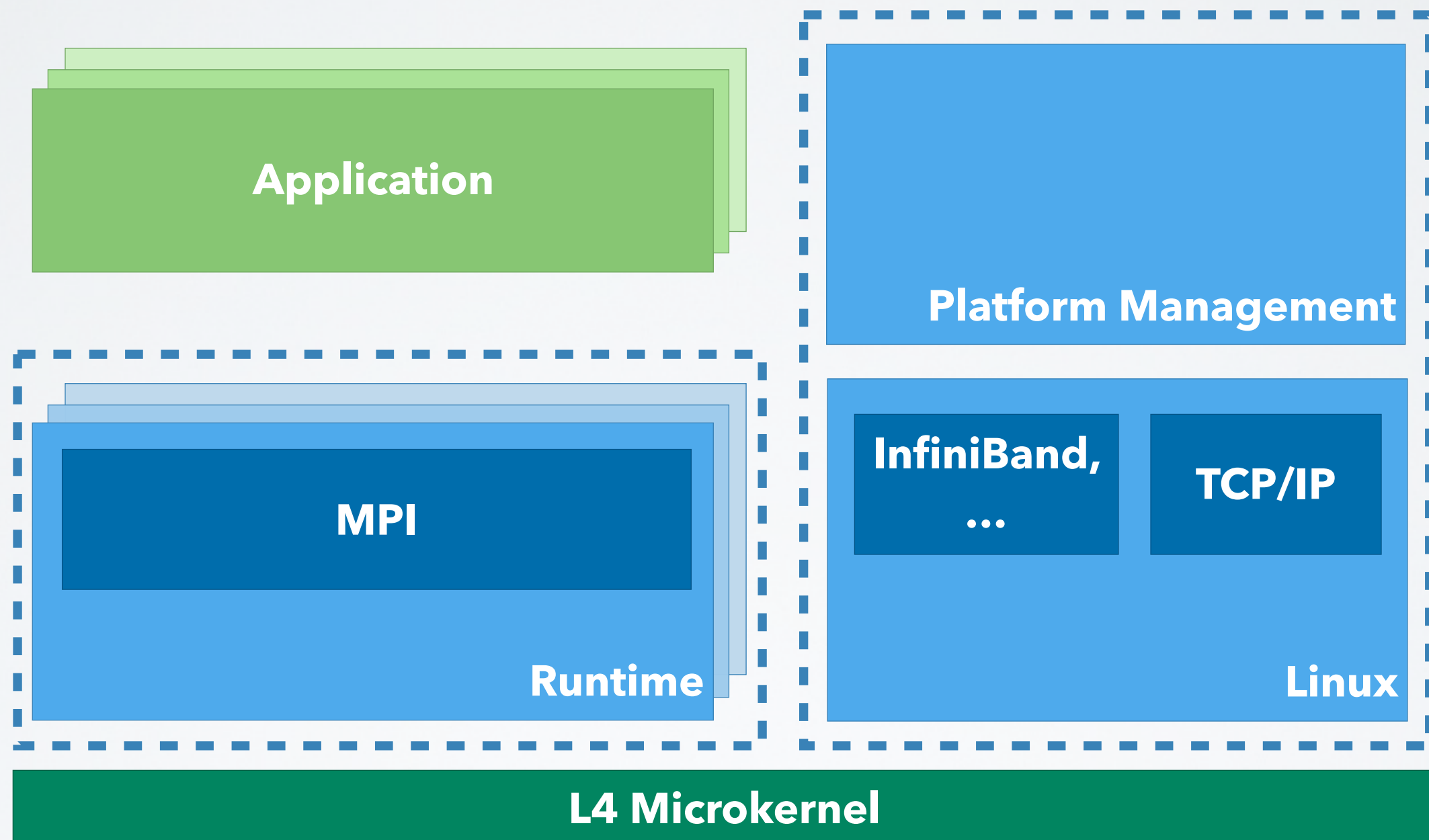


SYSTEM MODEL

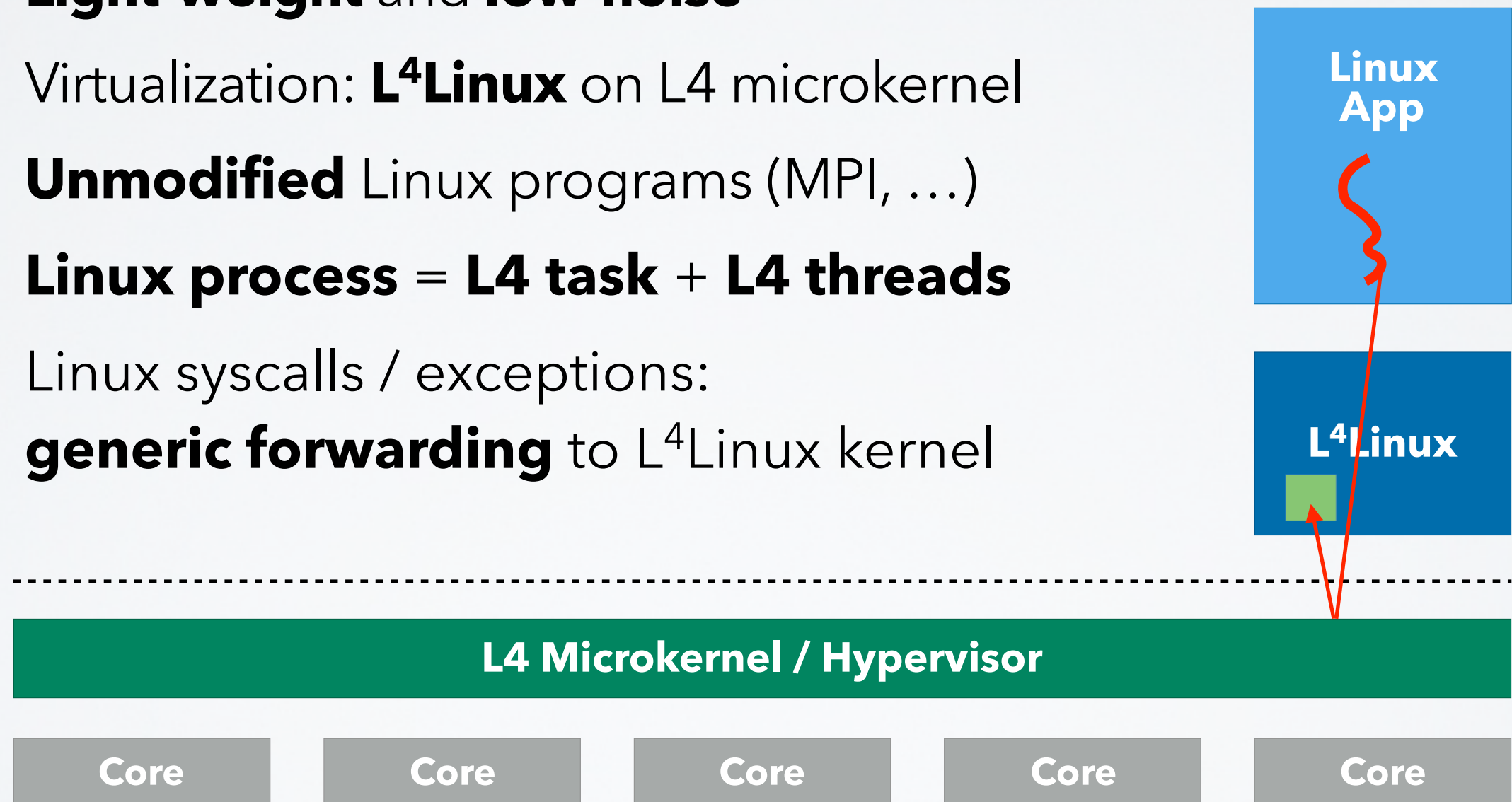




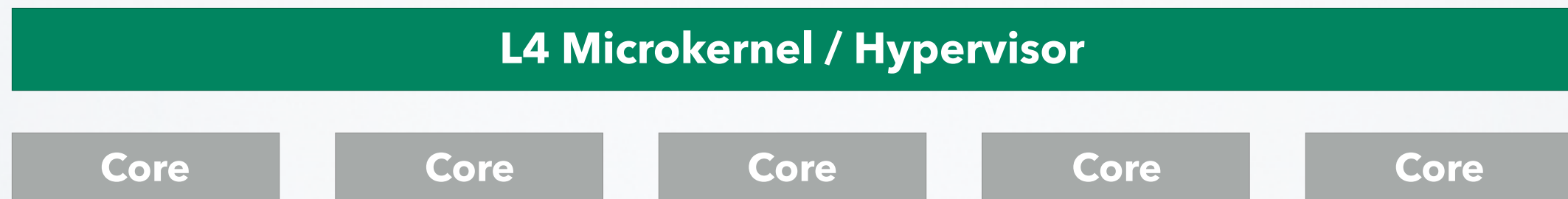
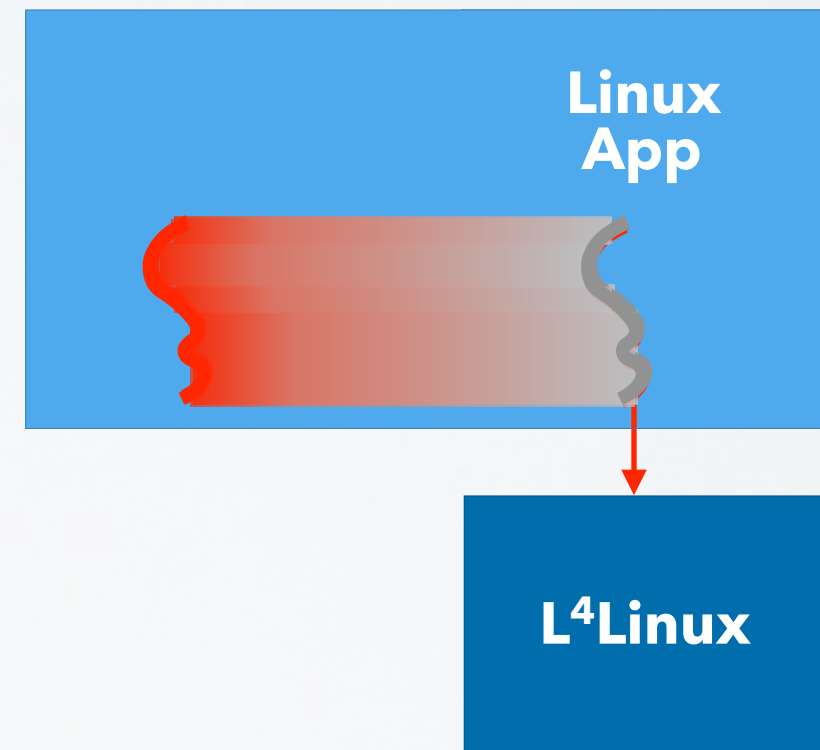
NODE ARCHITECTURE



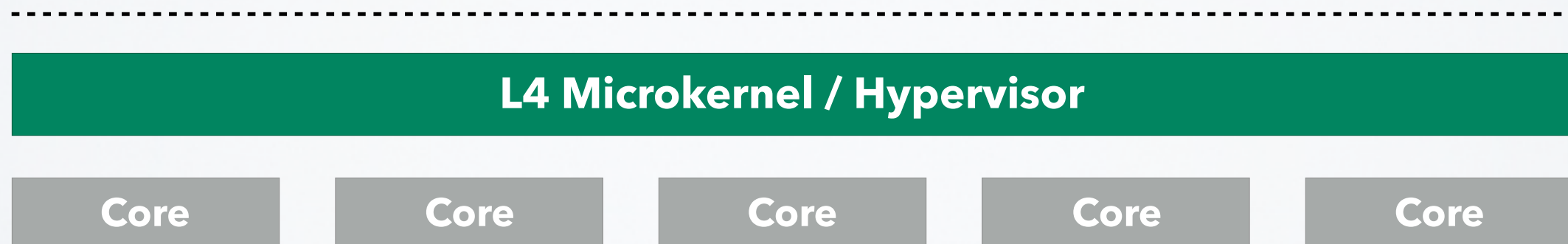
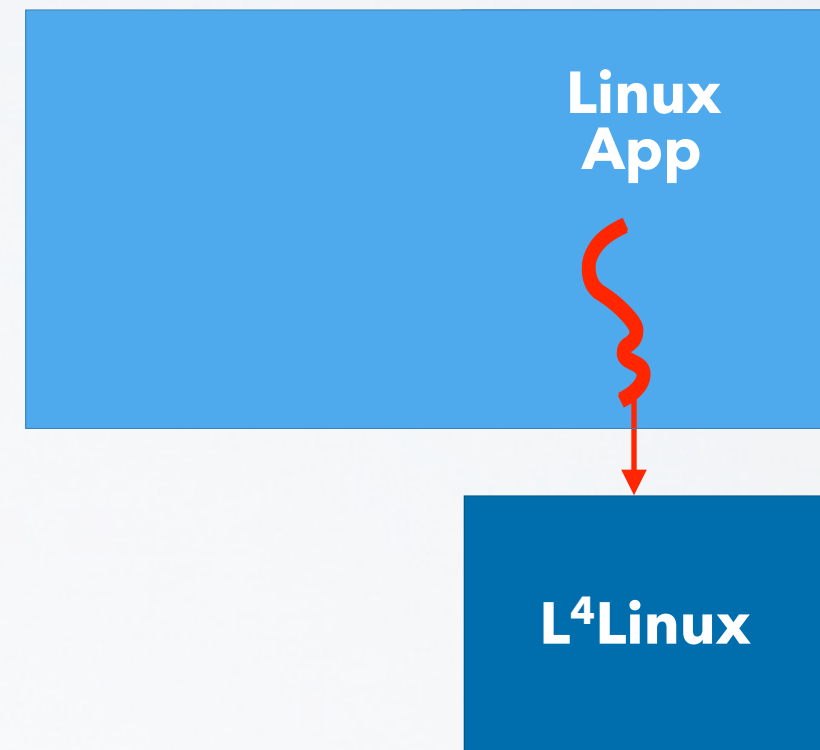
- **L4 microkernel** controls the node
- **Light-weight** and **low-noise**
- Virtualization: **L⁴Linux** on L4 microkernel
- **Unmodified** Linux programs (MPI, ...)
- **Linux process** = **L4 task** + **L4 threads**
- Linux syscalls / exceptions:
generic forwarding to L⁴Linux kernel



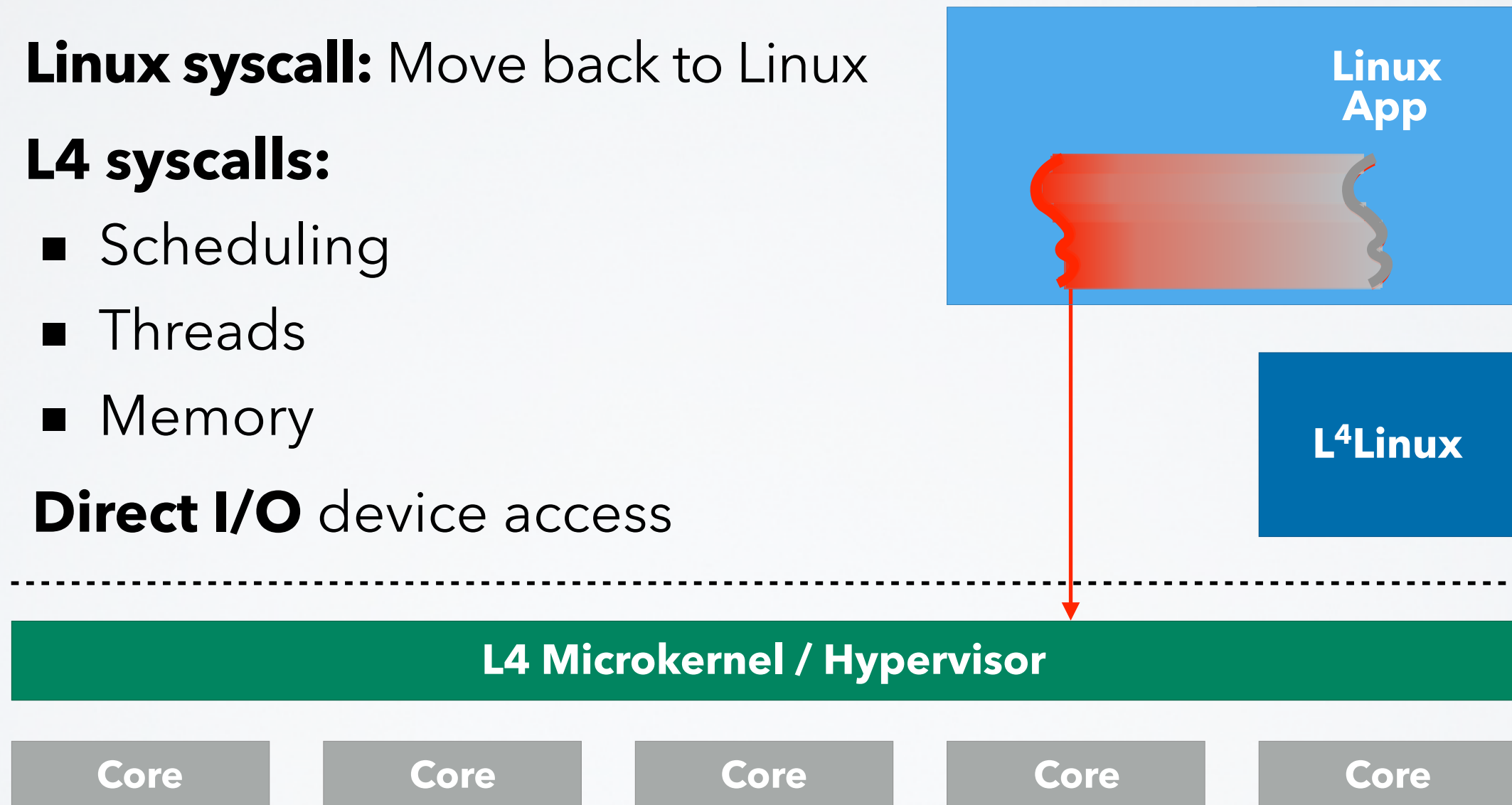
- **Decoupling:** move Linux thread to new L4 thread on its own core

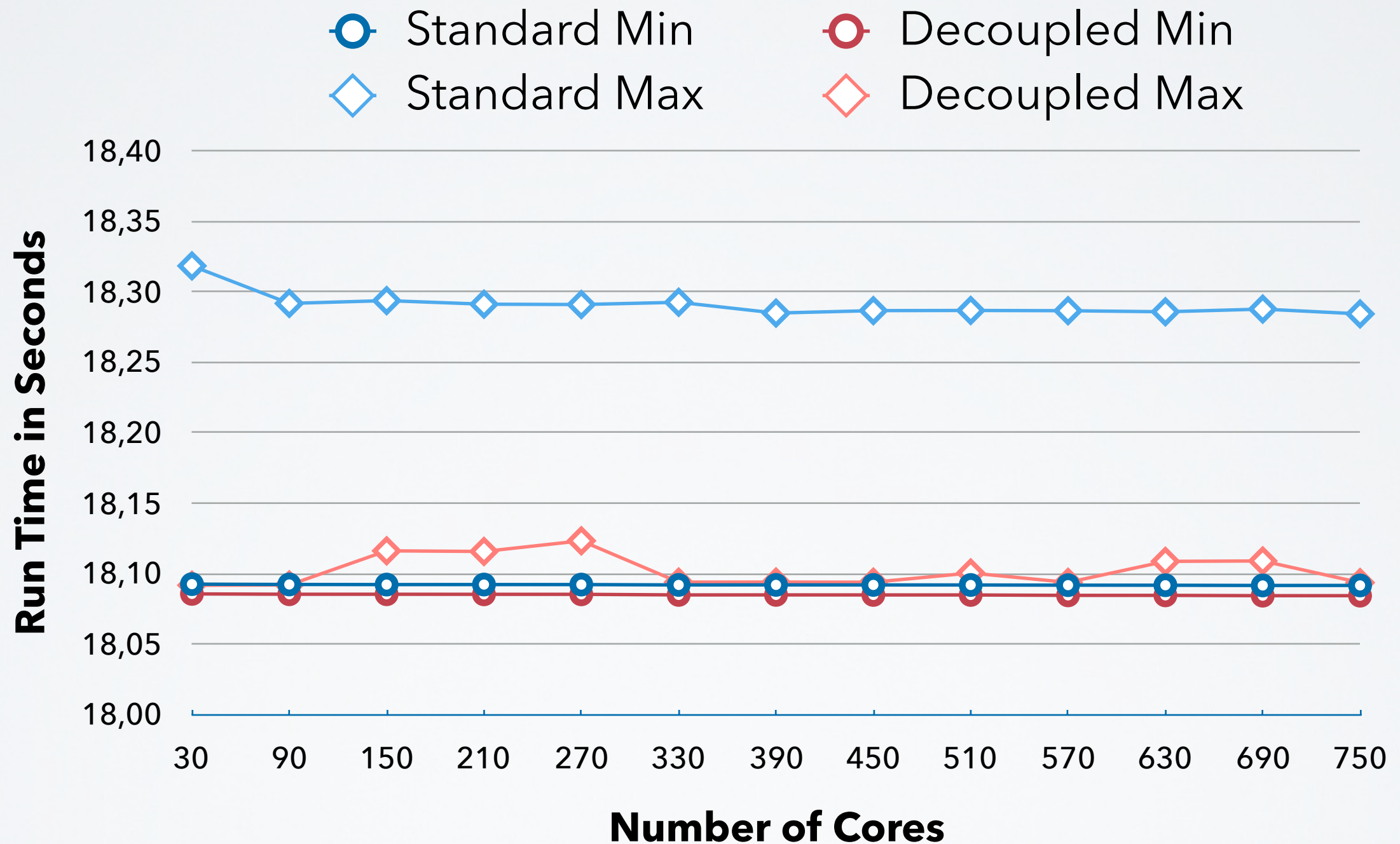


- **Decoupling:** move Linux thread to new L4 thread on its own core
- **Linux syscall:** Move back to Linux

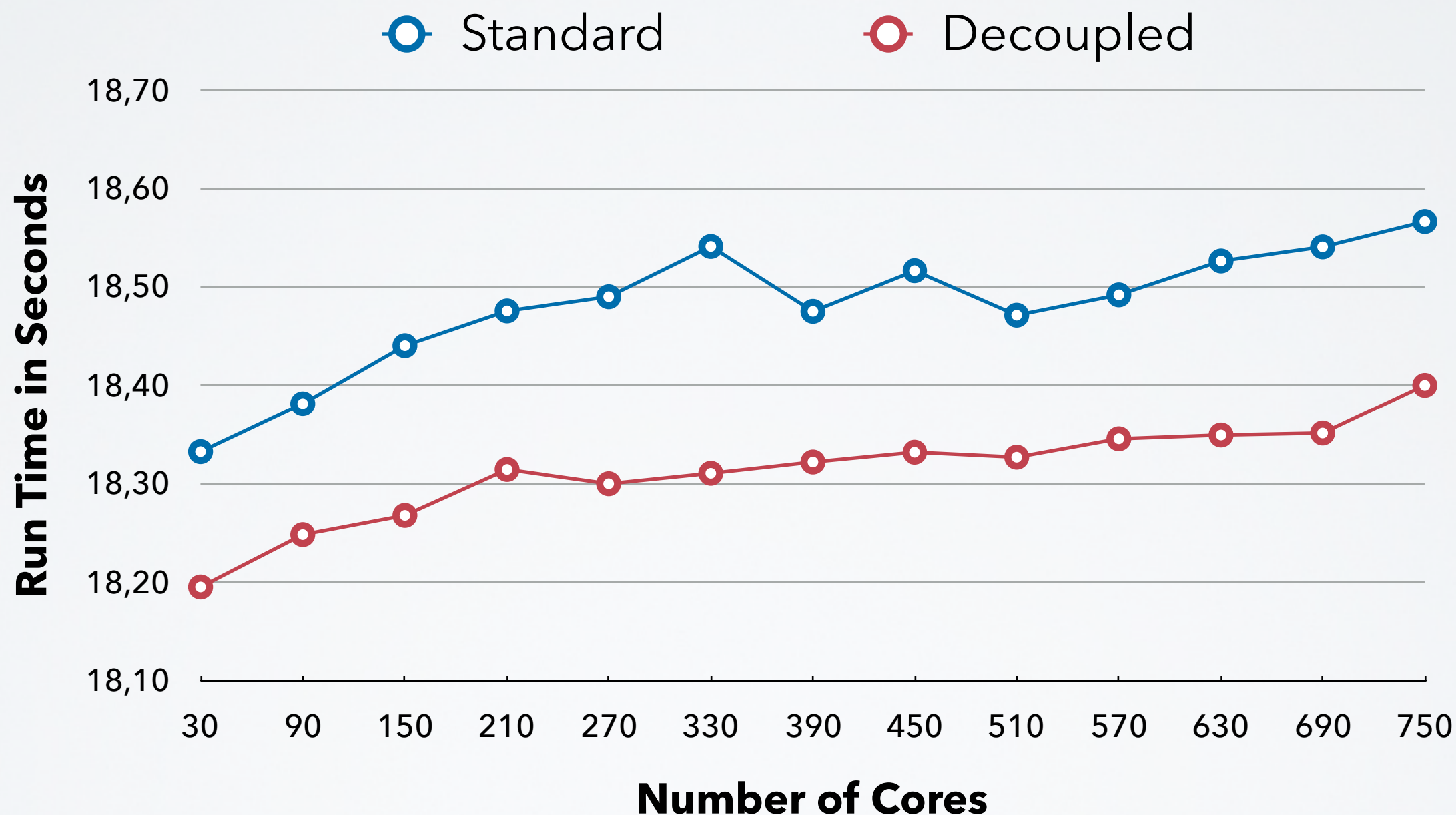


- **Decoupling:** move Linux thread to new L4 thread on its own core
- **Linux syscall:** Move back to Linux
- **L4 syscalls:**
 - Scheduling
 - Threads
 - Memory
- **Direct I/O** device access

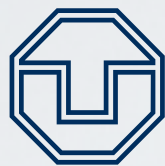




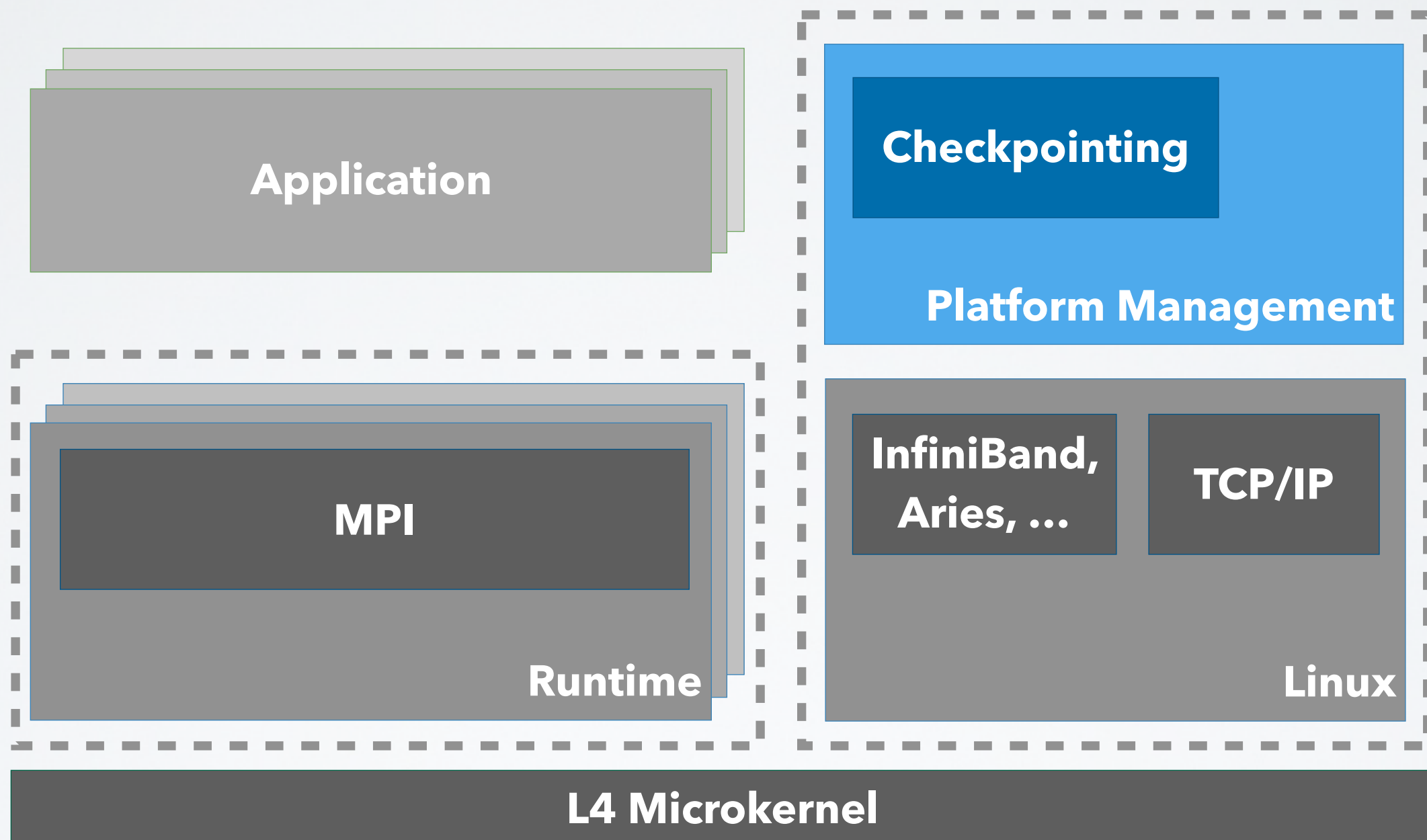
Adam Lackorzynski, Carsten Weinhold, Hermann Härtig, "Decoupled: Low-Effort Noise-Free Execution on Commodity Systems", ROSS 2016, June 2016, Kyoto, Japan

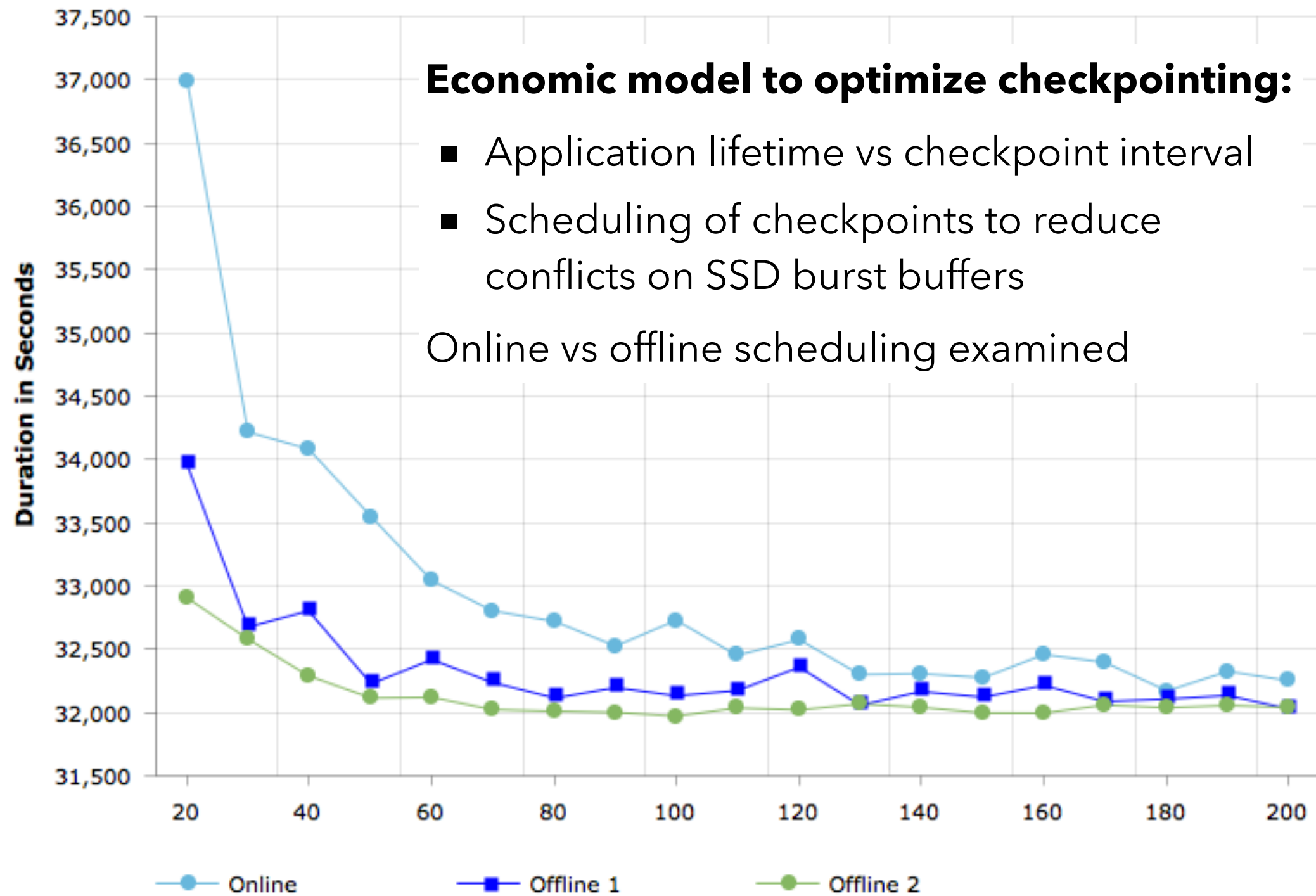


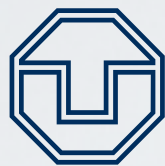
Adam Lackorzynski, Carsten Weinhold, Hermann Härtig, "Decoupled: Low-Effort Noise-Free Execution on Commodity Systems", ROSS 2016, June 2016, Kyoto, Japan



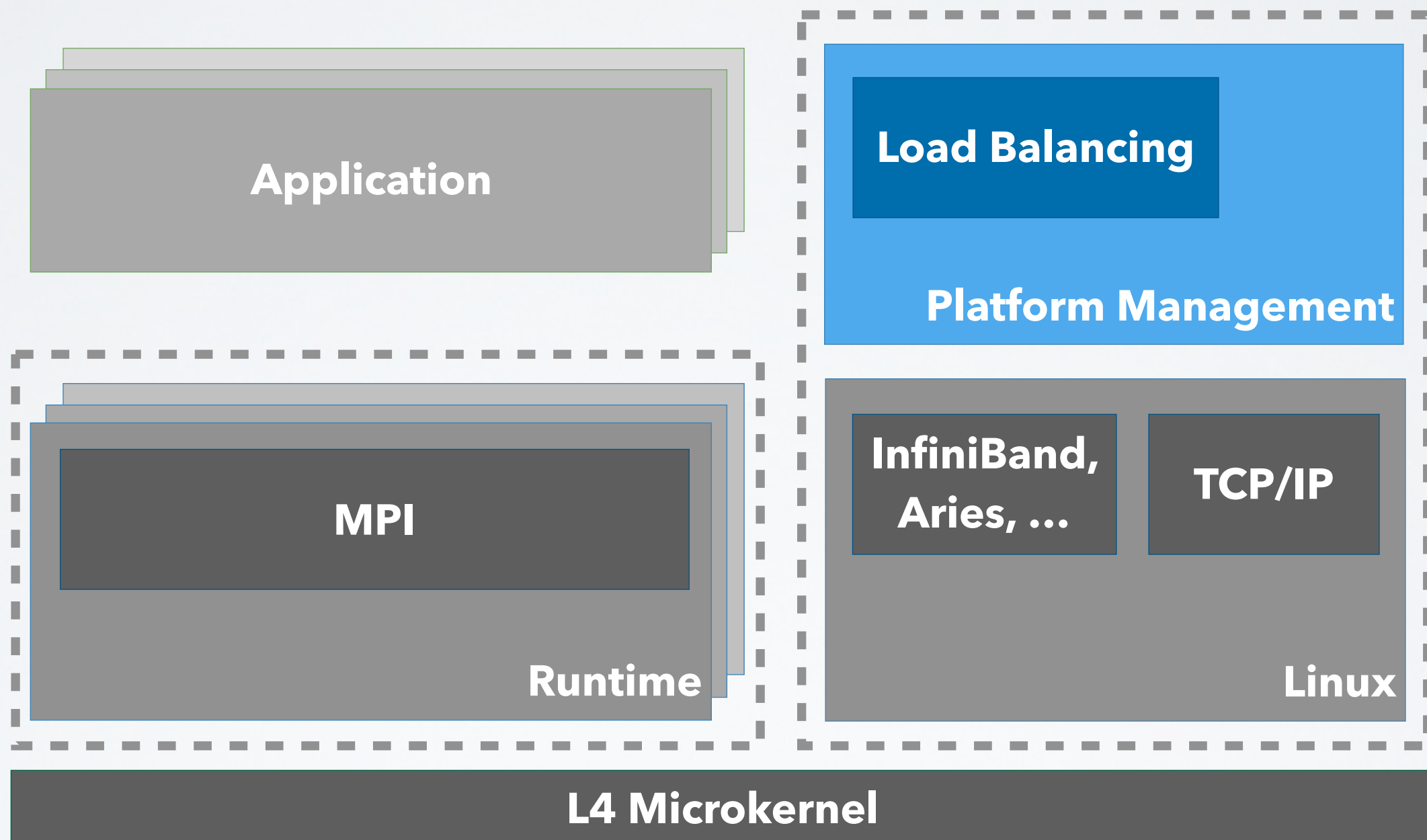
NODE ARCHITECTURE



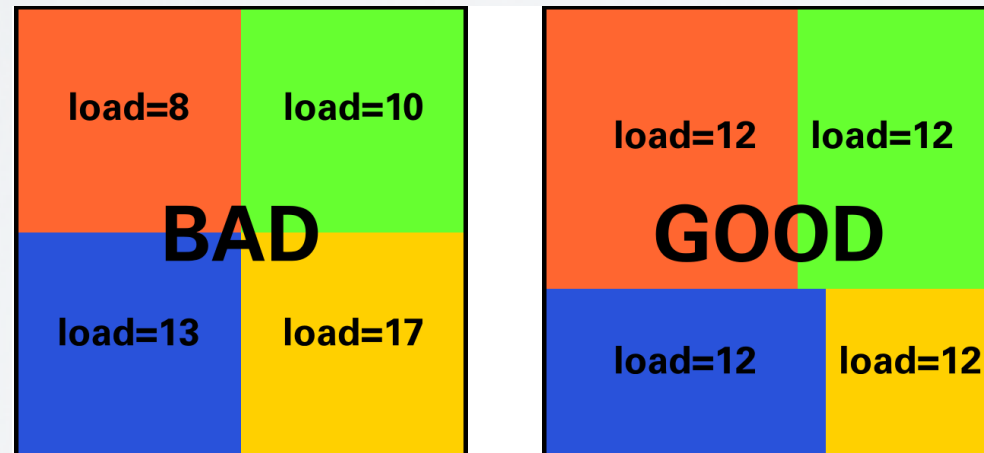




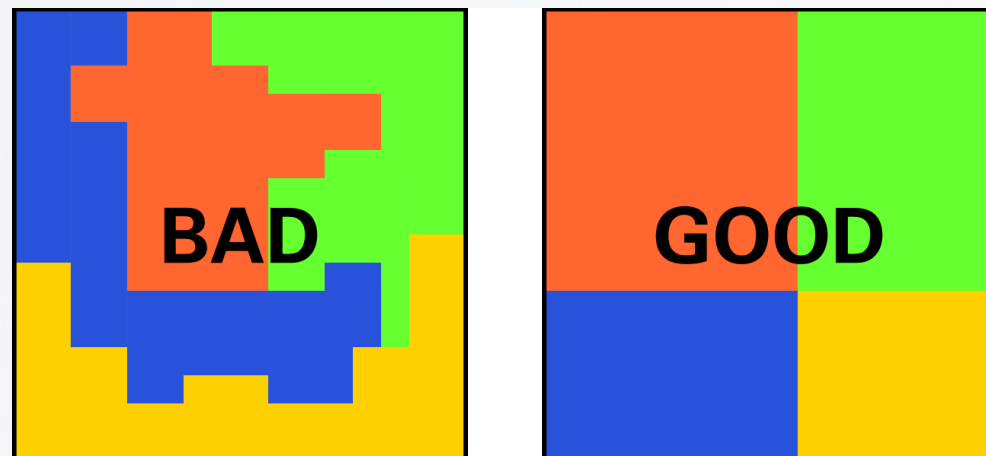
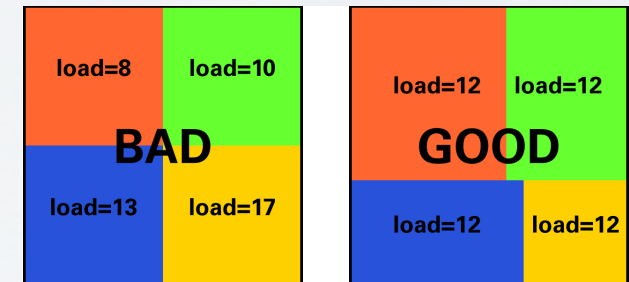
NODE ARCHITECTURE



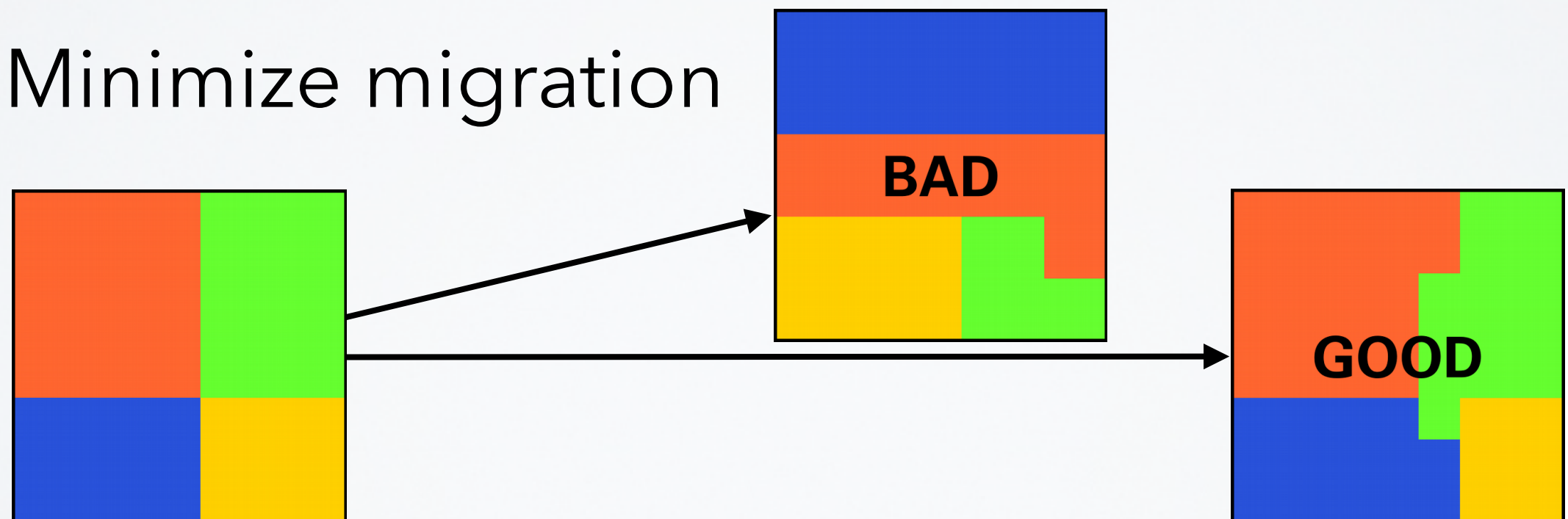
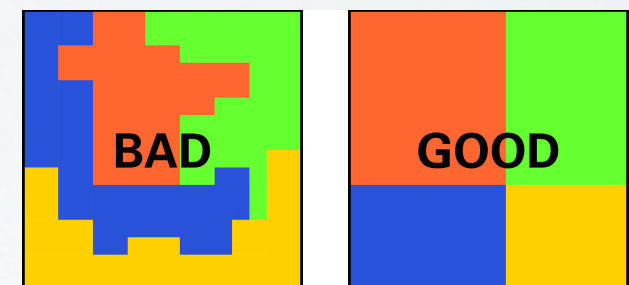
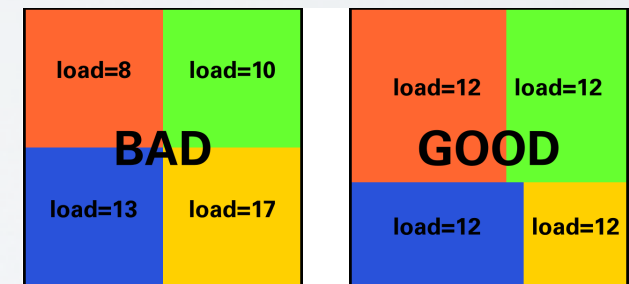
■ Balance workload



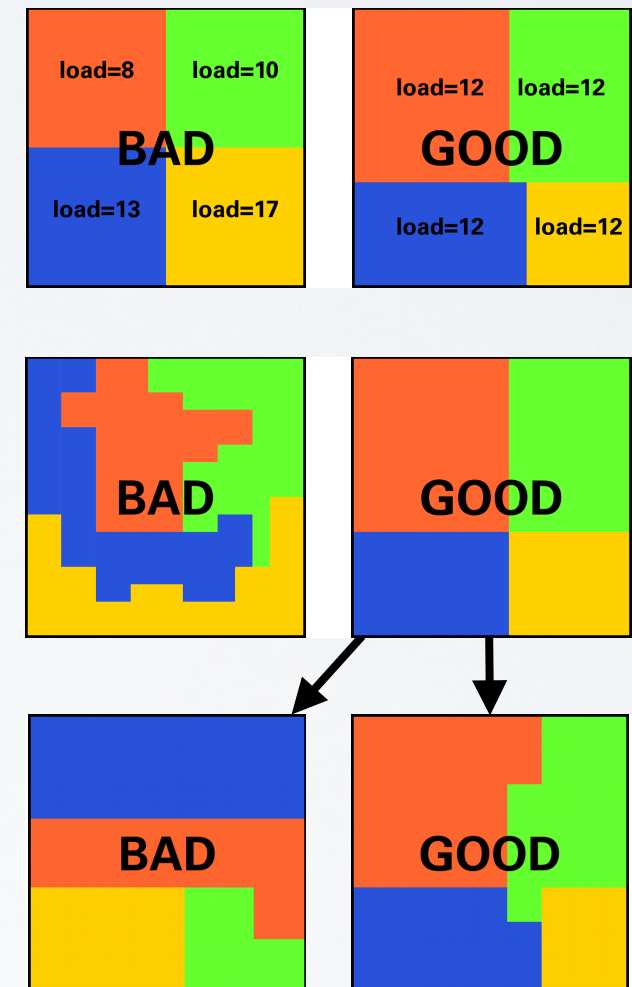
- Balance workload
- Minimize communication between partitions

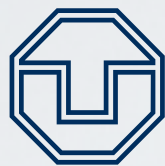


- Balance workload
- Minimize communication between partitions
- Minimize migration

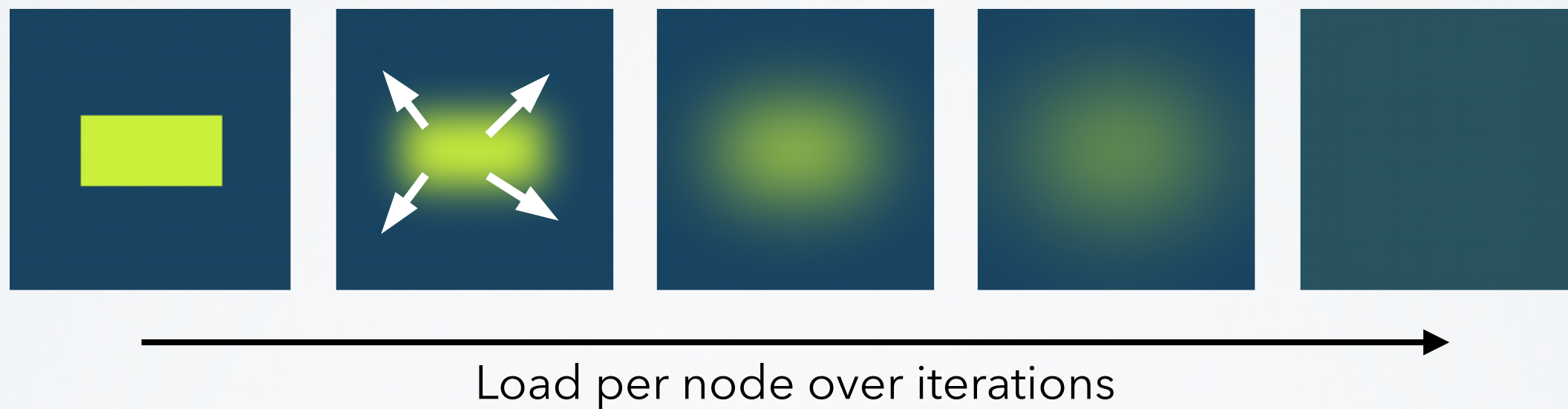


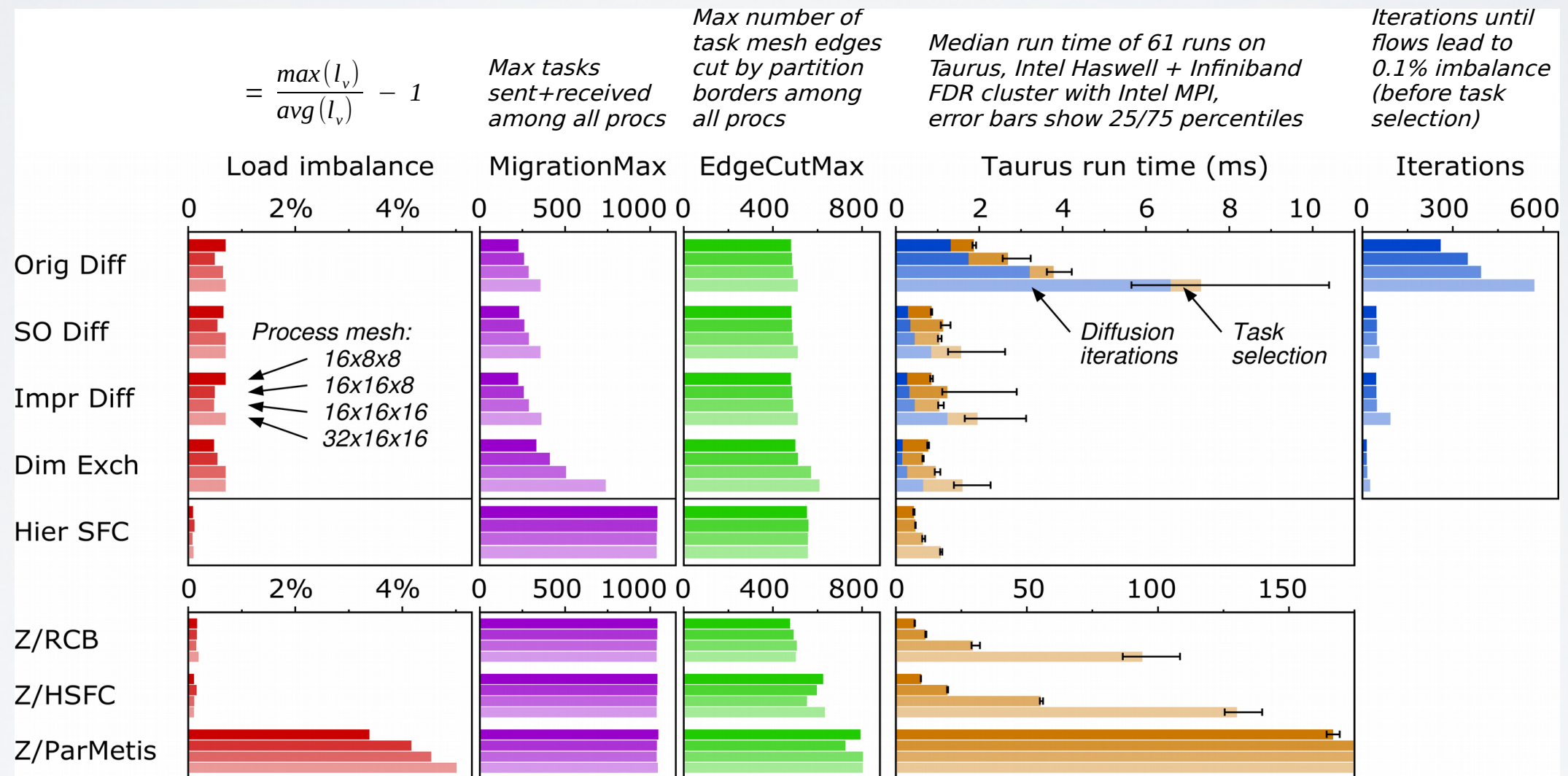
- Balance workload
- Minimize communication between partitions
- Minimize migration
- Compute new partitions fast





DIFFUSION



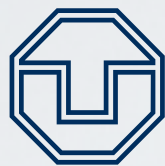


Diffusion leads to smallest migration

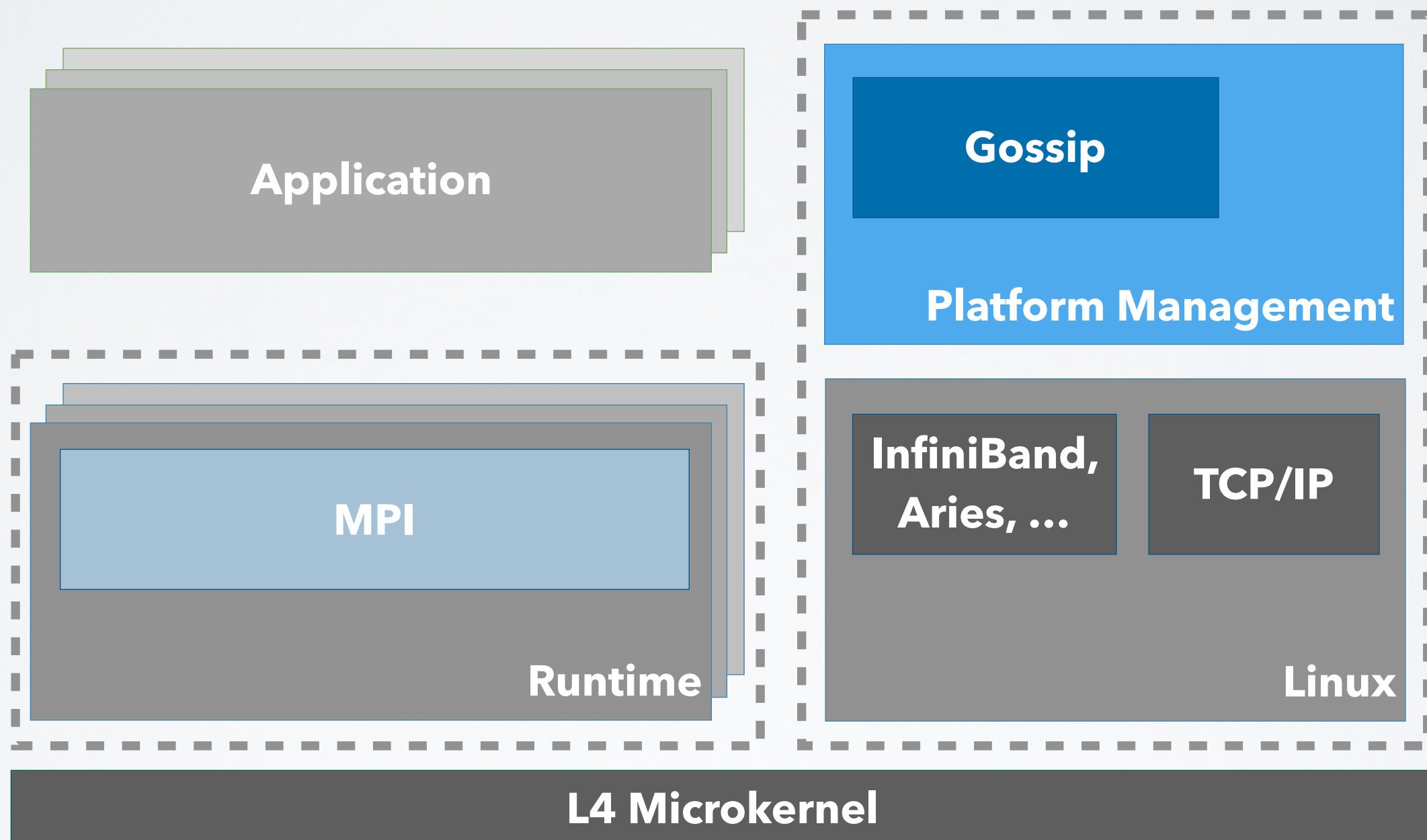
Diffusion achieves very good edge cut

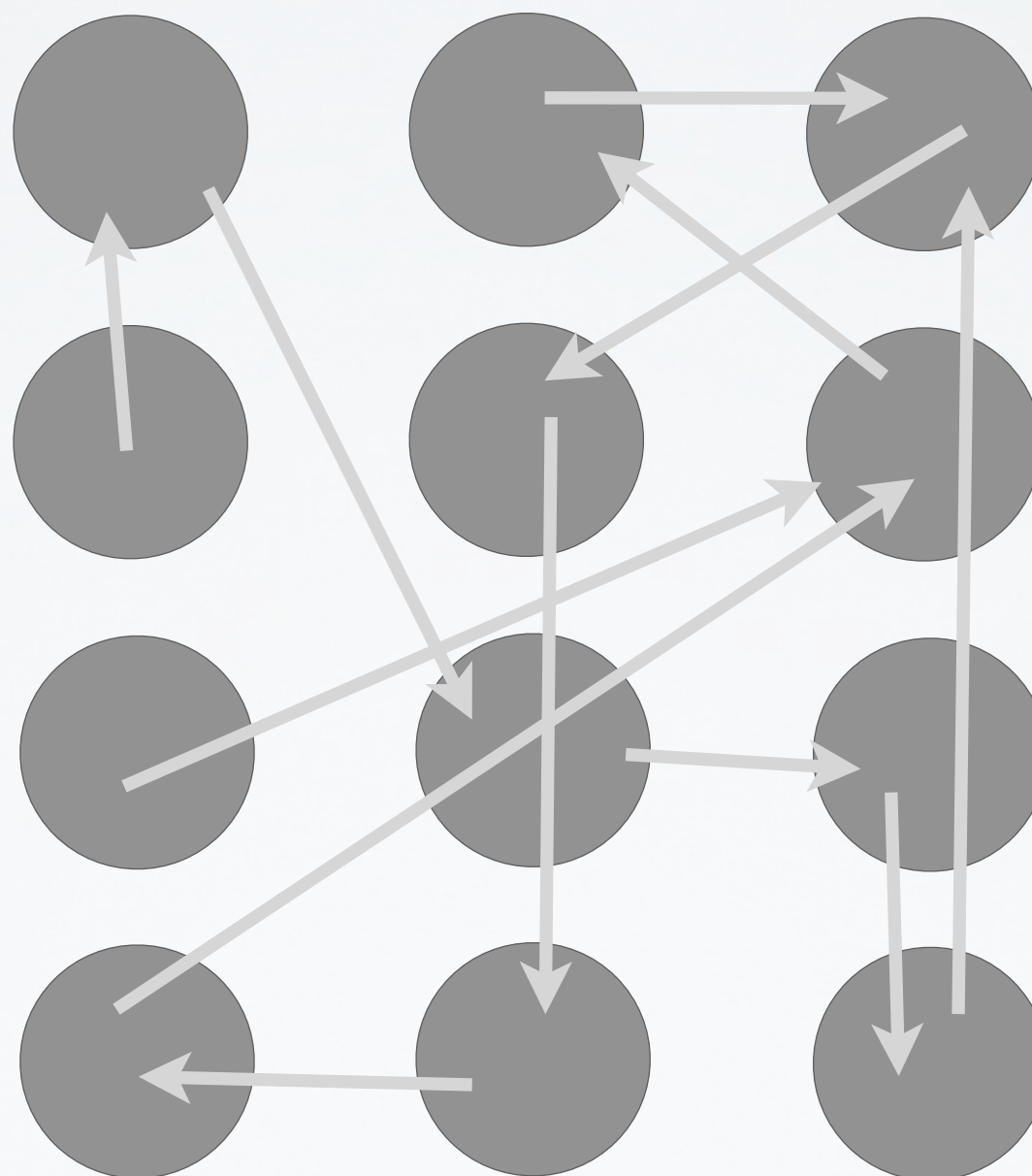
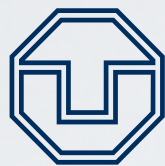
Diffusion run time ~2 ms for 8192 processes, Zoltan much slower

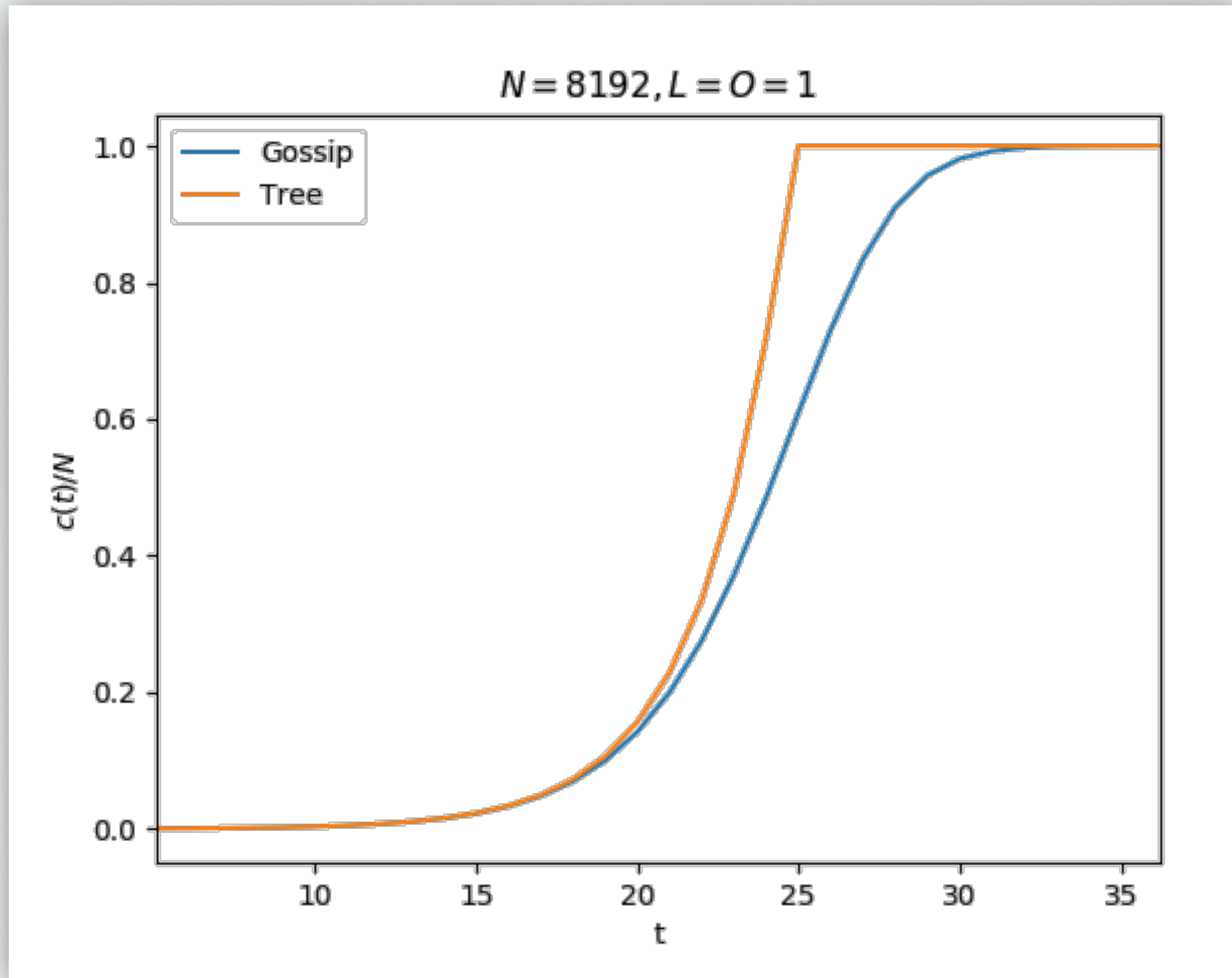
Matthias Lieber, Kerstin Gößner, Wolfgang E. Nagel, "The Potential of Diffusive Load Balancing at Large Scale", EuroMPI 2016, June 2016, Edinburgh, United Kingdom

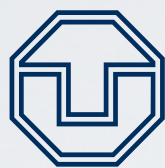


NODE ARCHITECTURE

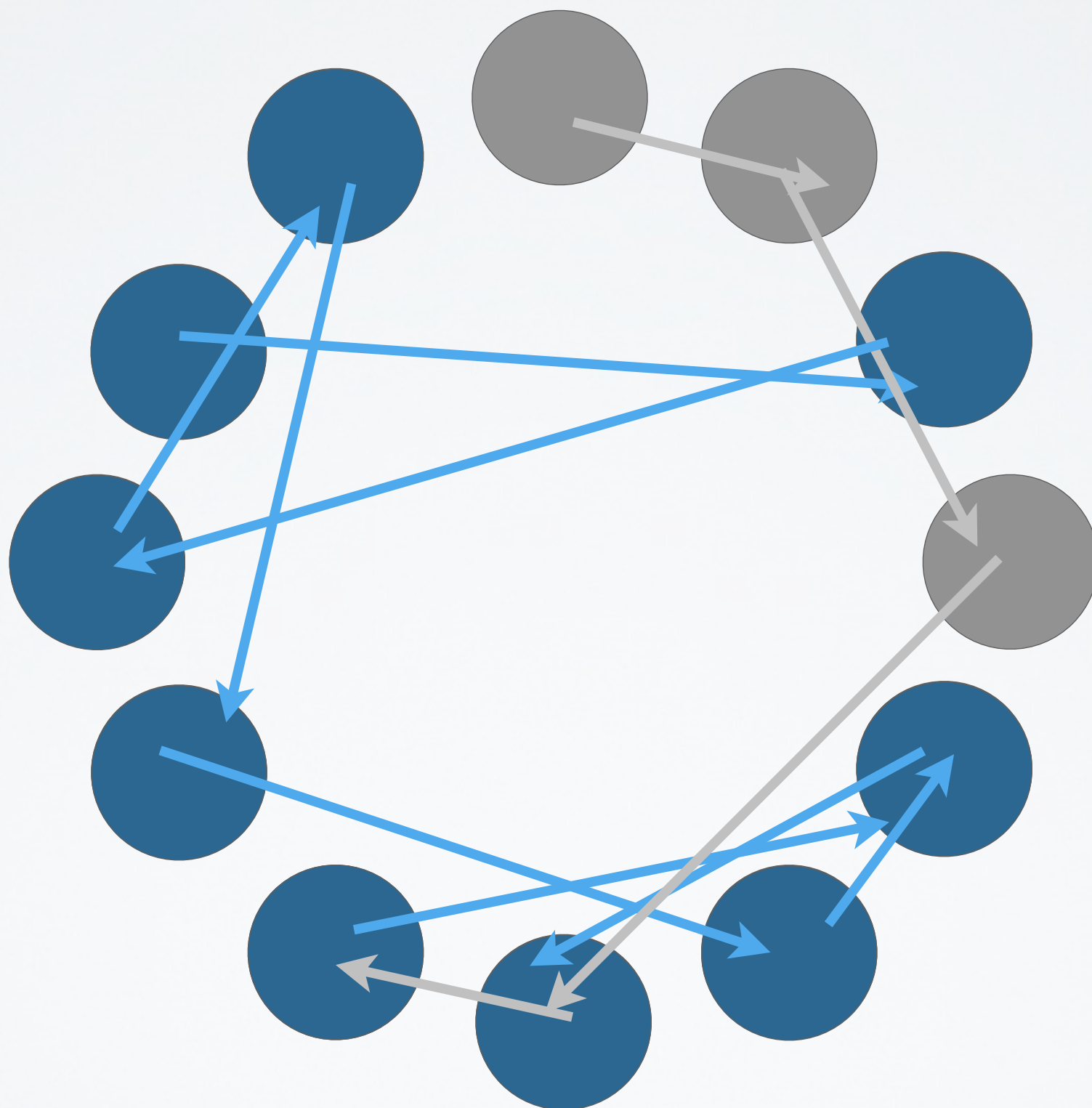






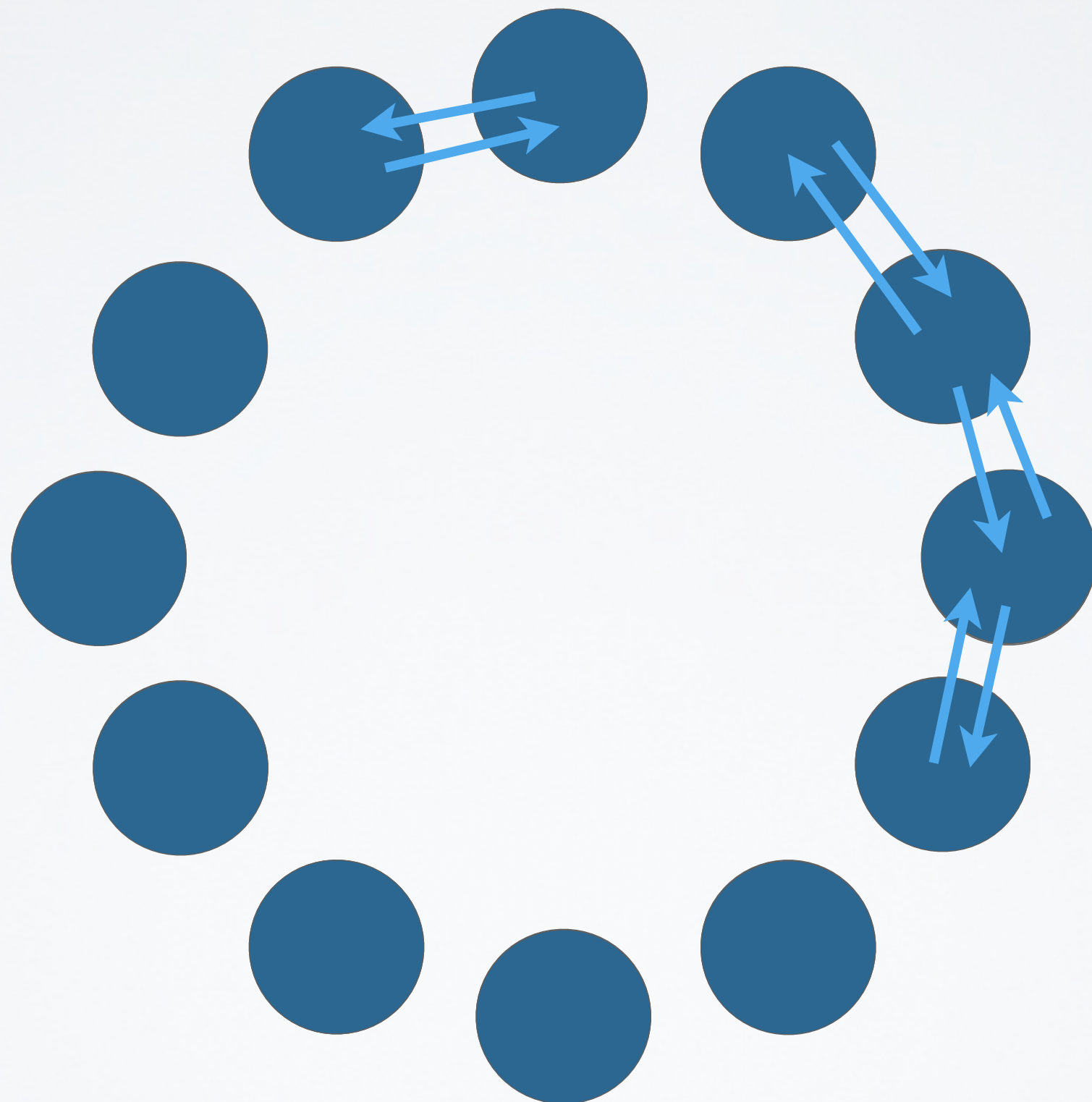


STEP 1: GOSSIP





STEP 2: CORRECTION



- **Two-stage algorithm:** gossip + correction
- **Main advantage:** scalability and resilience
(continues to work in presence of failures)
- **Works for:** fault-tolerant broadcast
- **Next step:** extend to operations that
include barrier semantics
- **Future:** use in MPI?

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

- **Decoupled threads:** reduced noise
- **Checkpointing:** Economic model
- **Diffusion:** may be efficient alternative
- **Corrected Gossip:** fault-tolerant broadcast
- **Work in progress:** integrate monitoring + gossip + decision making + migration