

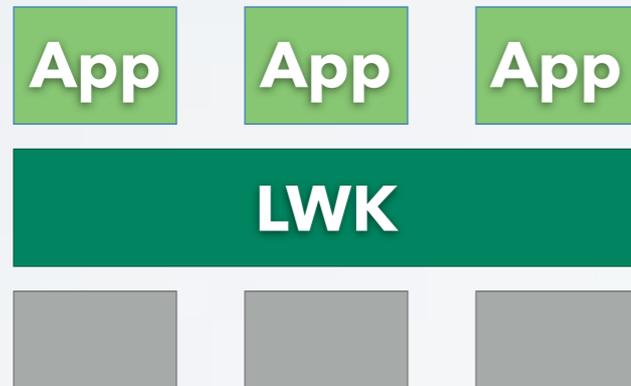


**TECHNISCHE
UNIVERSITÄT
DRESDEN**

TAILORING THE LINUX KERNEL AND MPI TO A FAST AND NOISE-FREE MICROKERNEL

CARSTEN WEINHOLD, TU DRESDEN

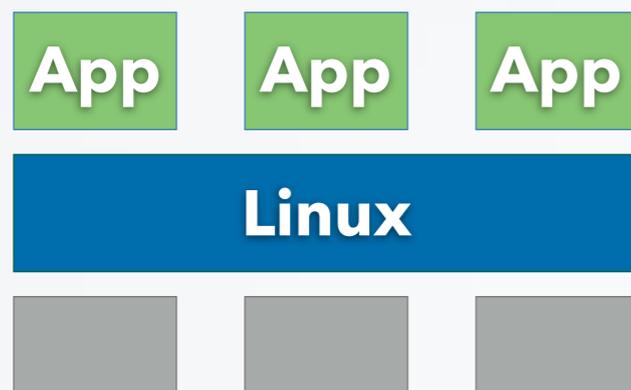
ADAM LACKORZYNSKI, JULIAN STECKLINA, MARKUS PARTHEYMÜLLER,
JAN BIERBAUM, MICHAEL JAHN, HERMANN HÄRTIG, STAFF@ZIH, STAFF@JSC



LWK

- ⊕ No Noise
- ⊖ Compatibility
- ⊖ Features

CNK

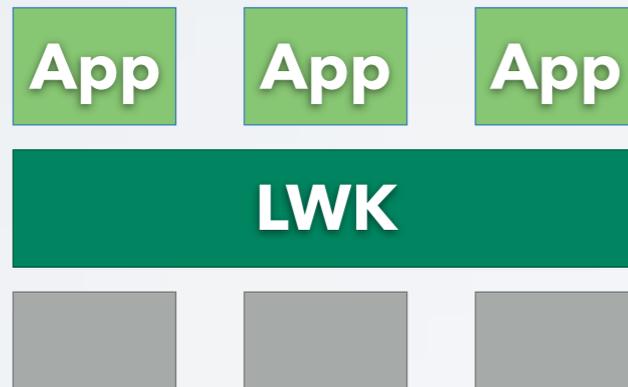


Linux (tweaked)

- ⊙ Low Noise
- ⊕ Compatibility
- ⊕ Features
- ⊖ Fast moving target

Cray

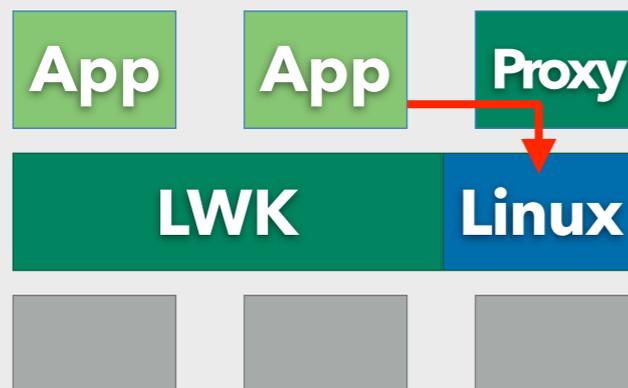
Argo



LWK

- ⊕ No Noise
- ⊖ Compatibility
- ⊖ Features

CNK



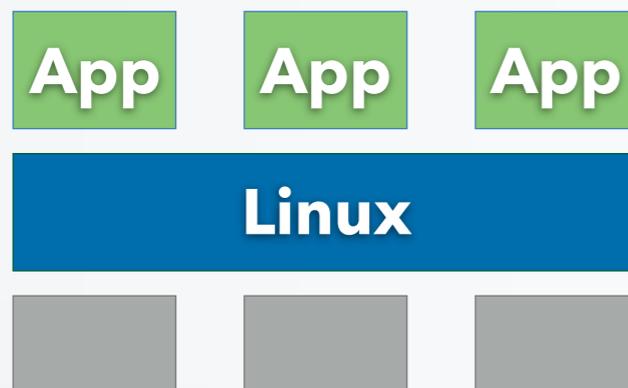
LWK + Linux

- ⊕ No Noise
- ⊕ Compatibility
- ⊕ Features
- ⊖ **Effort**

mOS

McKernel

Hobbes/Kitten



Linux (tweaked)

- ⊖ Low Noise
- ⊕ Compatibility
- ⊕ Features
- ⊖ Fast moving target

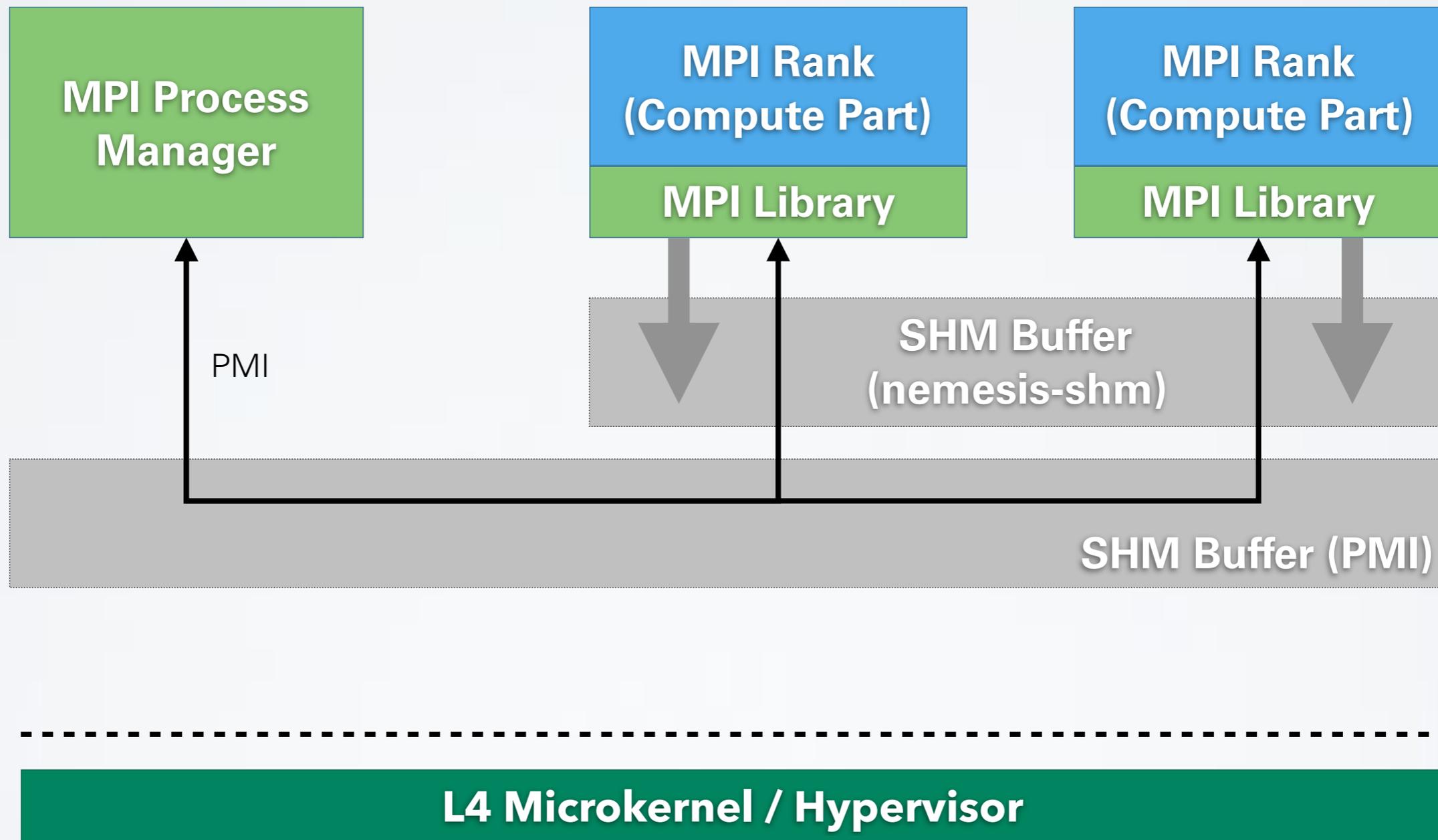
Cray

Argo

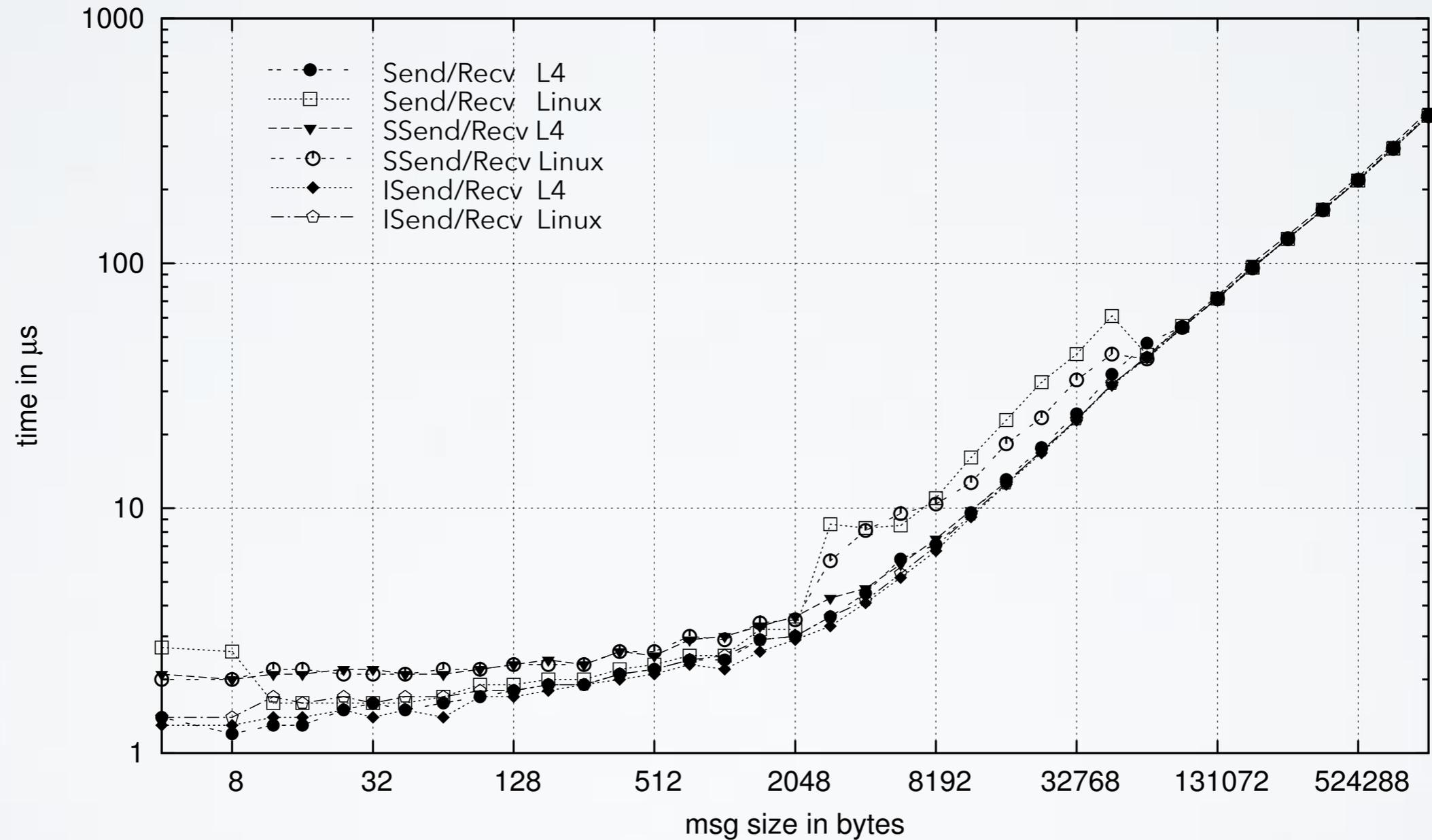


VERSION 0

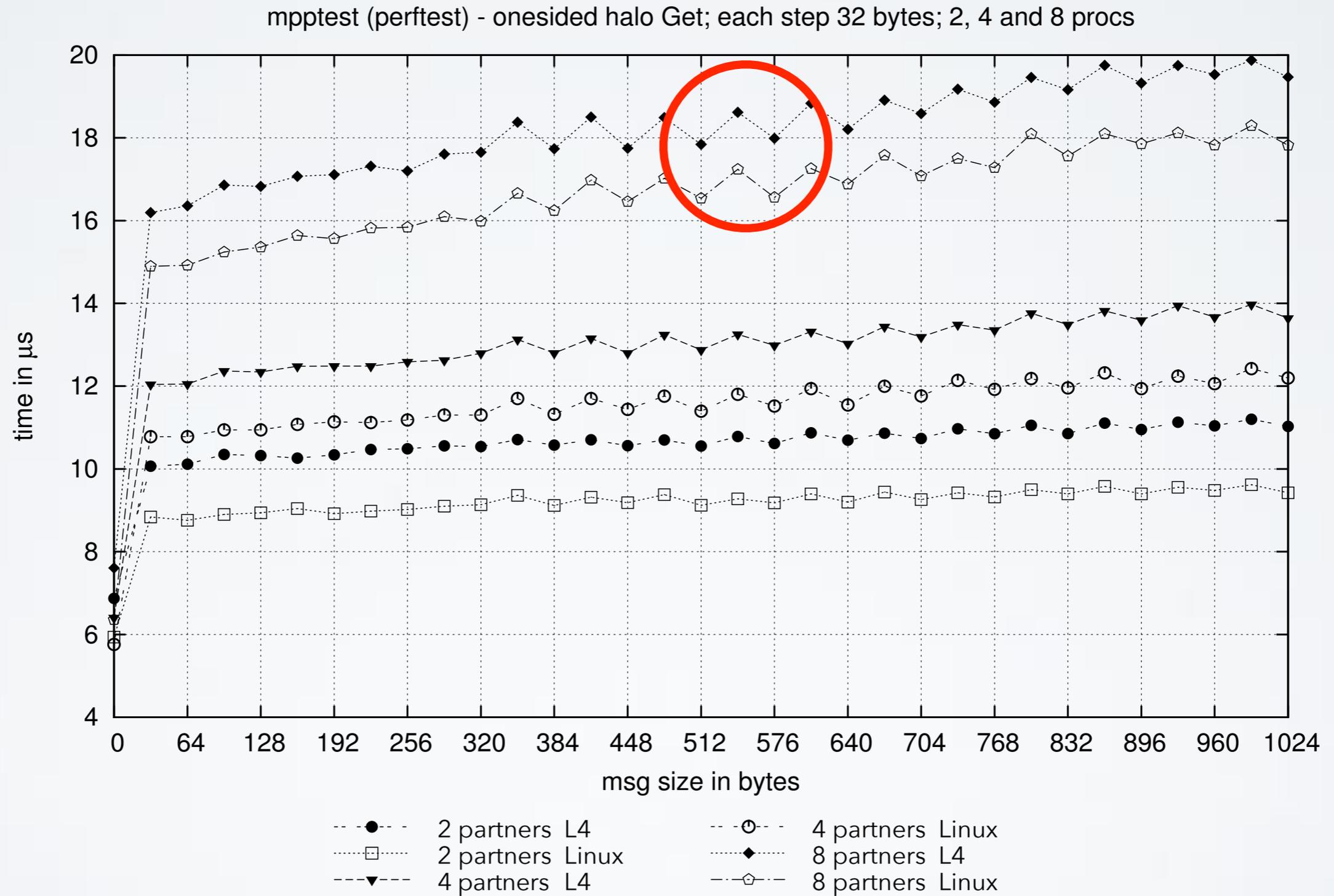
Pure L4



skampi - PingPong tests pt. I; 2 (of 10) procs



MPICH ON L4: PERFORMANCE



```
MPICH-L4| 0>--- p1MD version "0815" ---
MPICH-L4| MPI_Wtime
MPICH-L4| MPI_Wtime
MPICH-L4| 0> [compiled: "2013-08-12" ]
MPICH-L4| 0>
MPICH-L4| 0><i> Simulation started :: 01.01.1970 00:00:03.993 (UTC +0)
MPICH-L4| 0>
MPICH-L4| 0>
MPICH-L4| 0>Command: rom/p1md rom/md_10k_nochkpt.nml
MPICH-L4| 0>
MPICH-L4| 0>
MPICH-L4| 0><i> 4 processes

                <more output here>

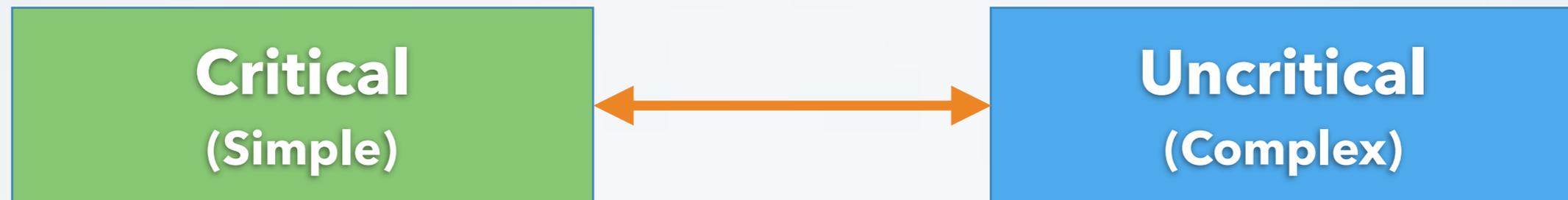
MPICH-L4| 0><cpu> wall-clock timer function: MPI_Wtime
MPICH-L4| 0> Work count sum: 184102828
MPICH-L4| 1>
MPICH-L4| 2>
MPICH-L4| 3><cpu> MD step performance: 2.936 steps/s
MPICH-L4| 0>
MPICH-L4| 1><cpu> MD step performance: 2.938 steps/s
MPICH-L4| 2><cpu> MD step performance: 2.939 steps/s
MPICH-L4| 3><cpu> MD interaction performance: 0.270 Mega interactions/s
MPICH-L4| 0><cpu> MD step performance: 2.945 steps/s
MPICH-L4| 1><cpu> MD interaction performance: 0.271 Mega interactions/s
MPICH-L4| 2><cpu> MD interaction performance: 0.271 Mega interactions/s
MPICH-L4| 0><cpu> MD interaction performance: 0.271 Mega interactions/s
MPICH-L4| 0>
MPICH-L4| 0>
MPICH-L4| 0><i> Simulation finished :: 01.01.1970 00:12:58.468 (UTC +0)
MPICH-L4| 0>
```

- It worked:
 - p1MD, MPICH test, perftest, ...
 - HPL, CP2K, ...
- Performance comparable to Linux^(*)
- Major limitation: no network backend



VERSION 1

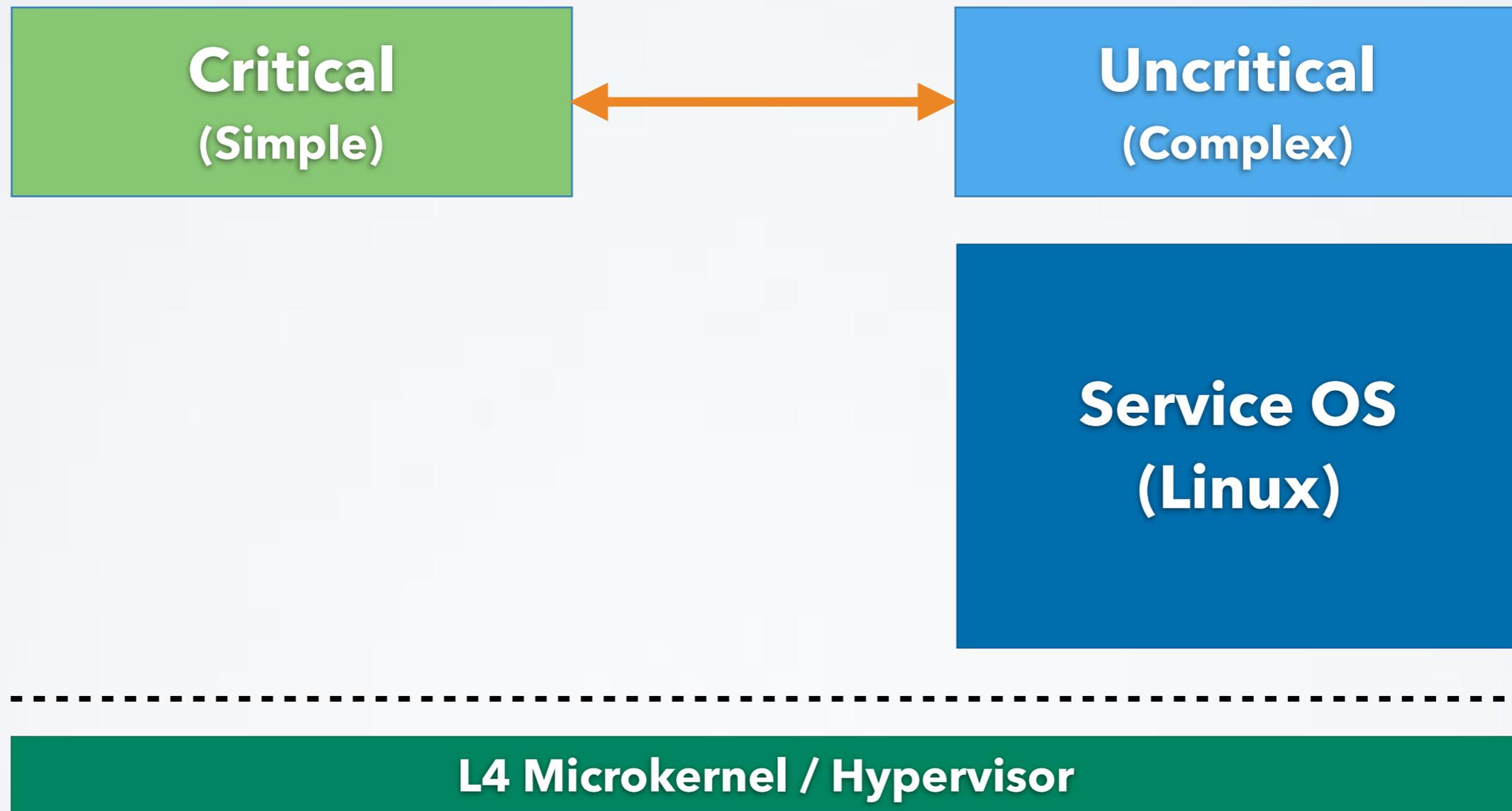
Hybrid L4 + Linux

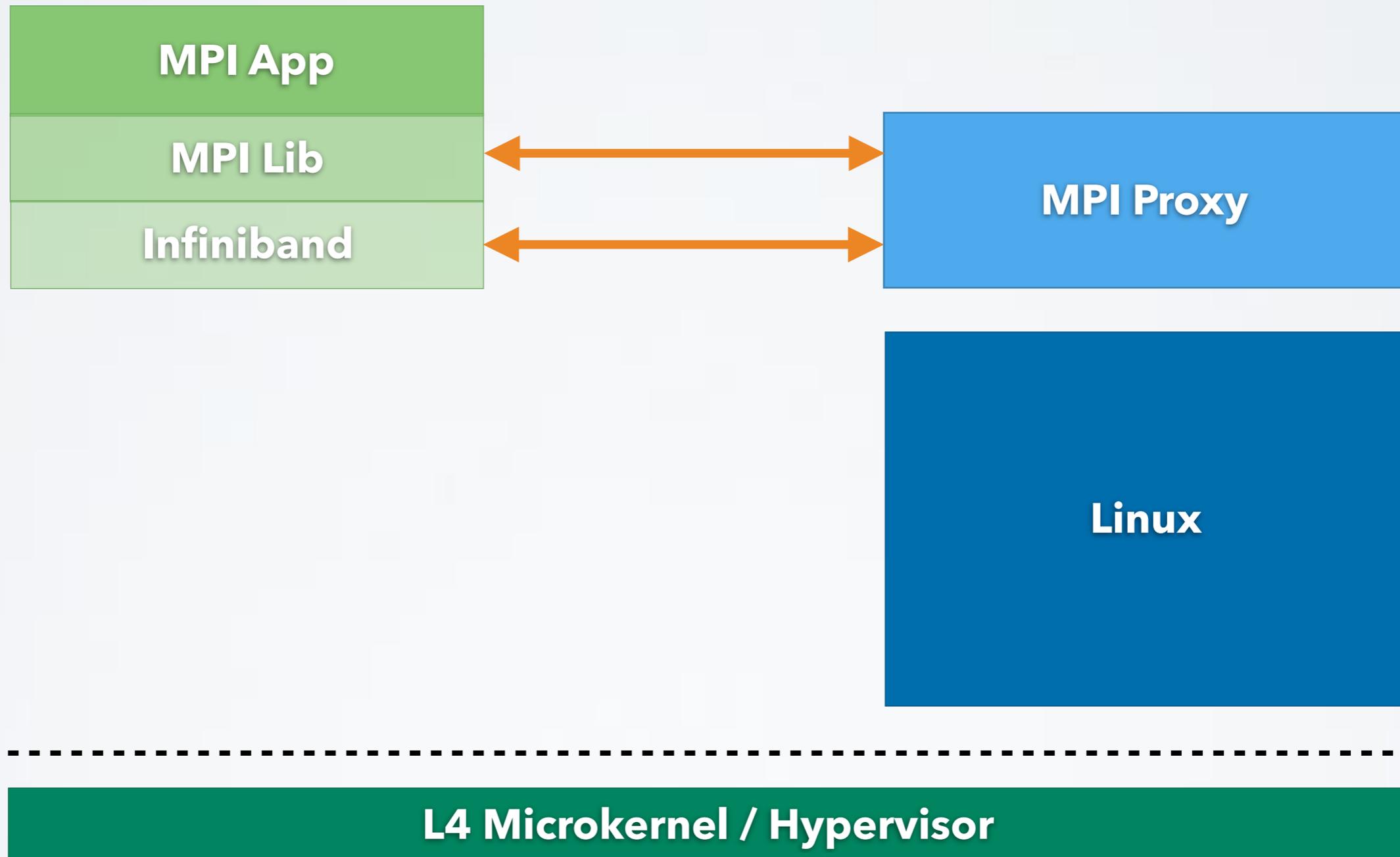


Real-time

Security

Resilience





Fast path (direct access to MMIO regs)

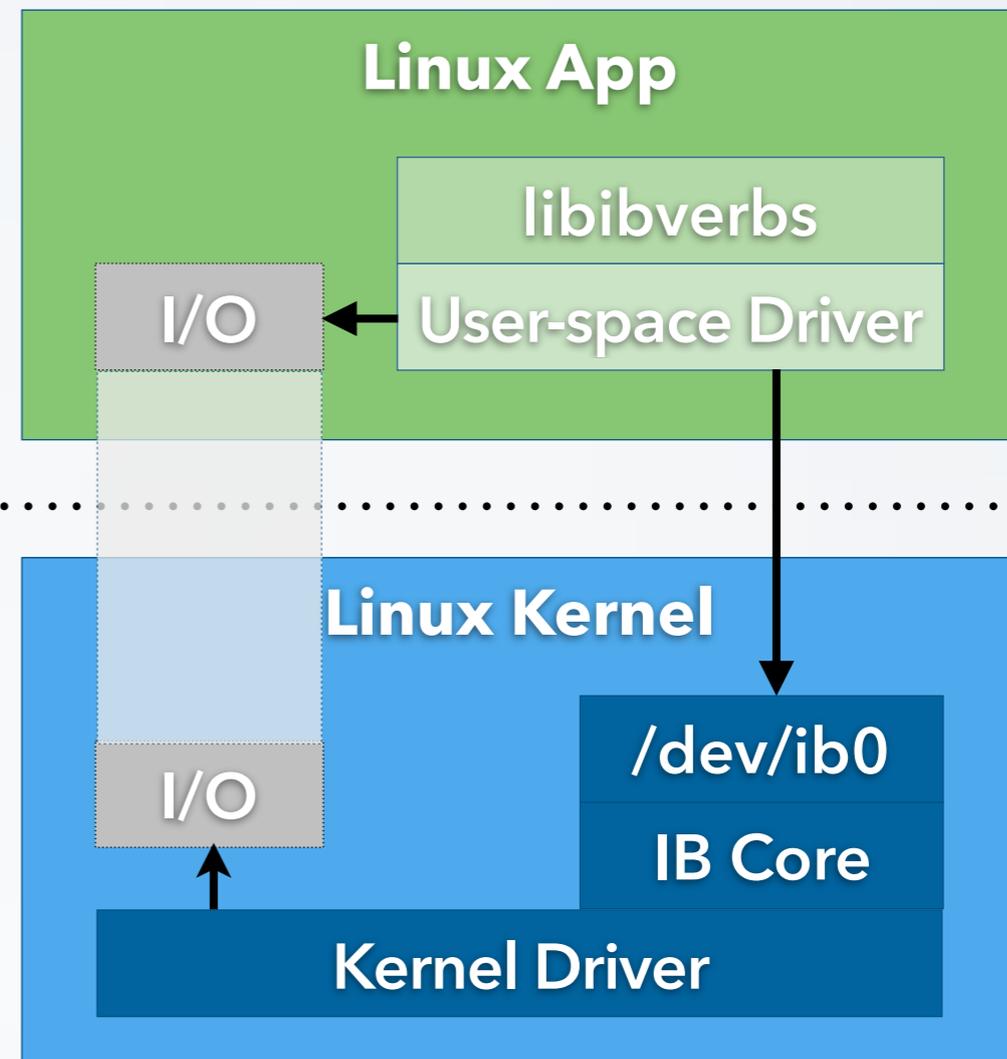
- Send / recv
- RDMA read / write
- Event polling

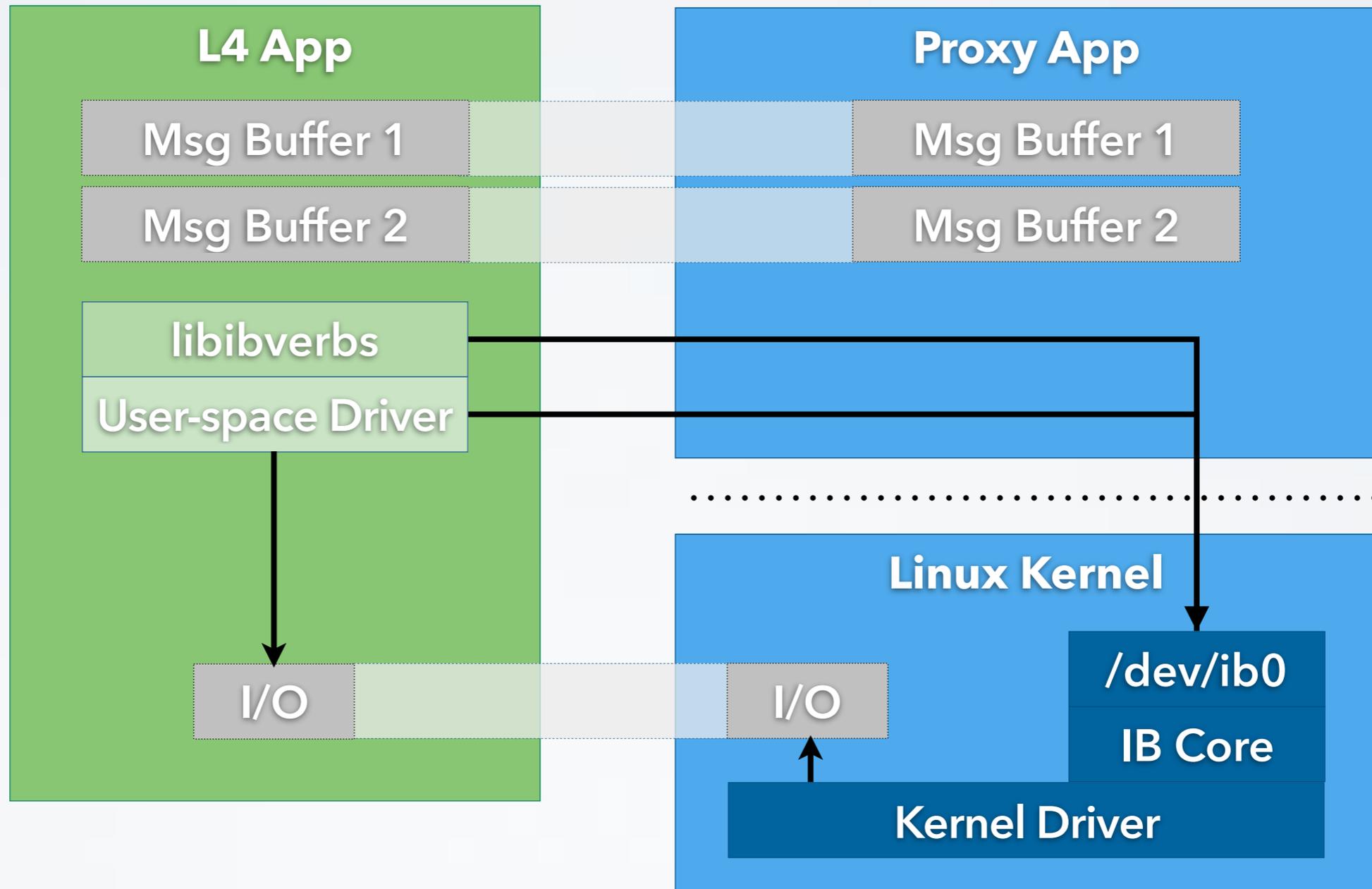
User Space

Kernel Space

Slow path (syscall required)

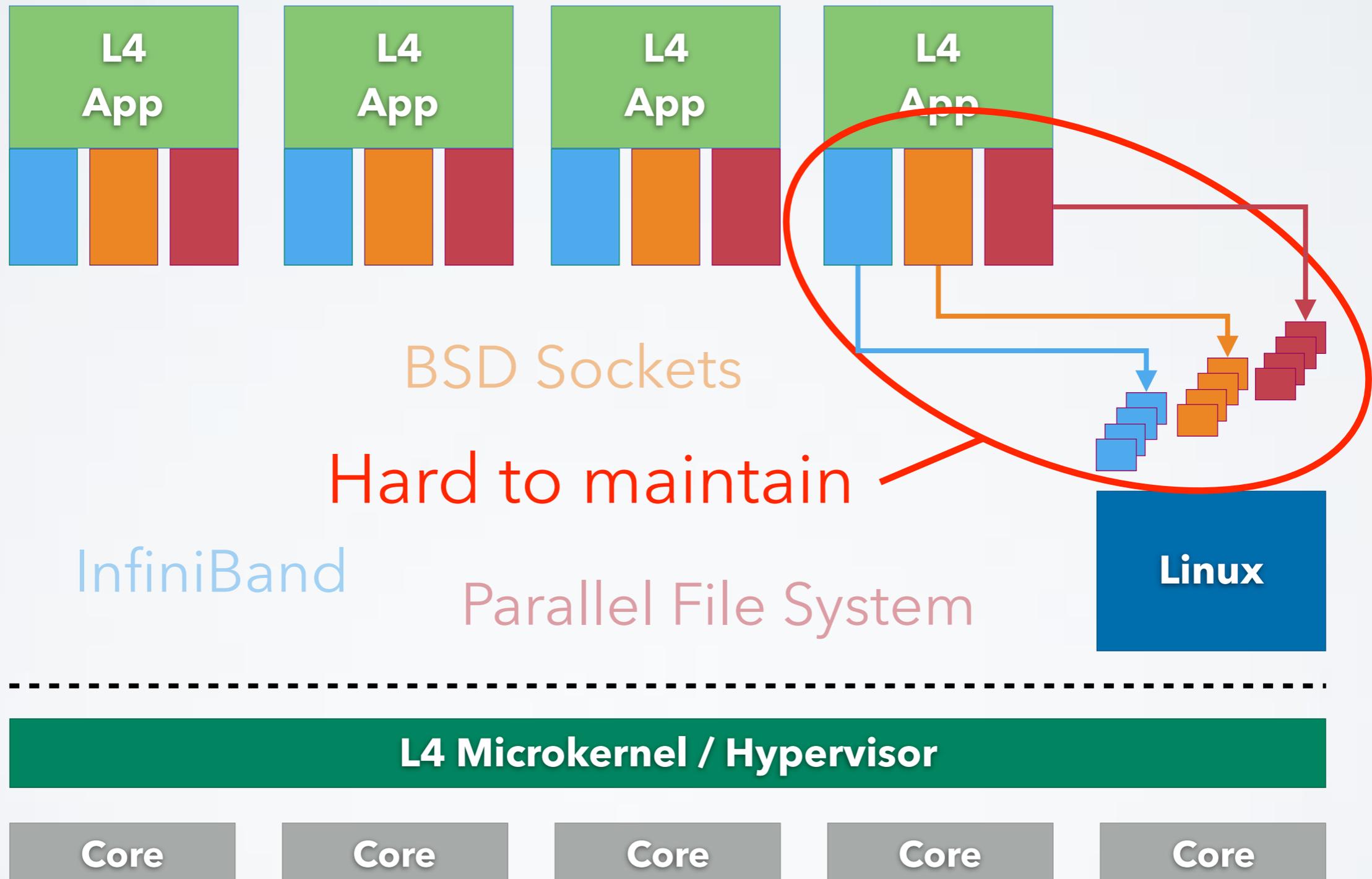
- Connection setup / teardown
- Buffer registration
- Event blocking





L4 Microkernel / Hypervisor

VERSION 1: SPLIT MVAPICH2

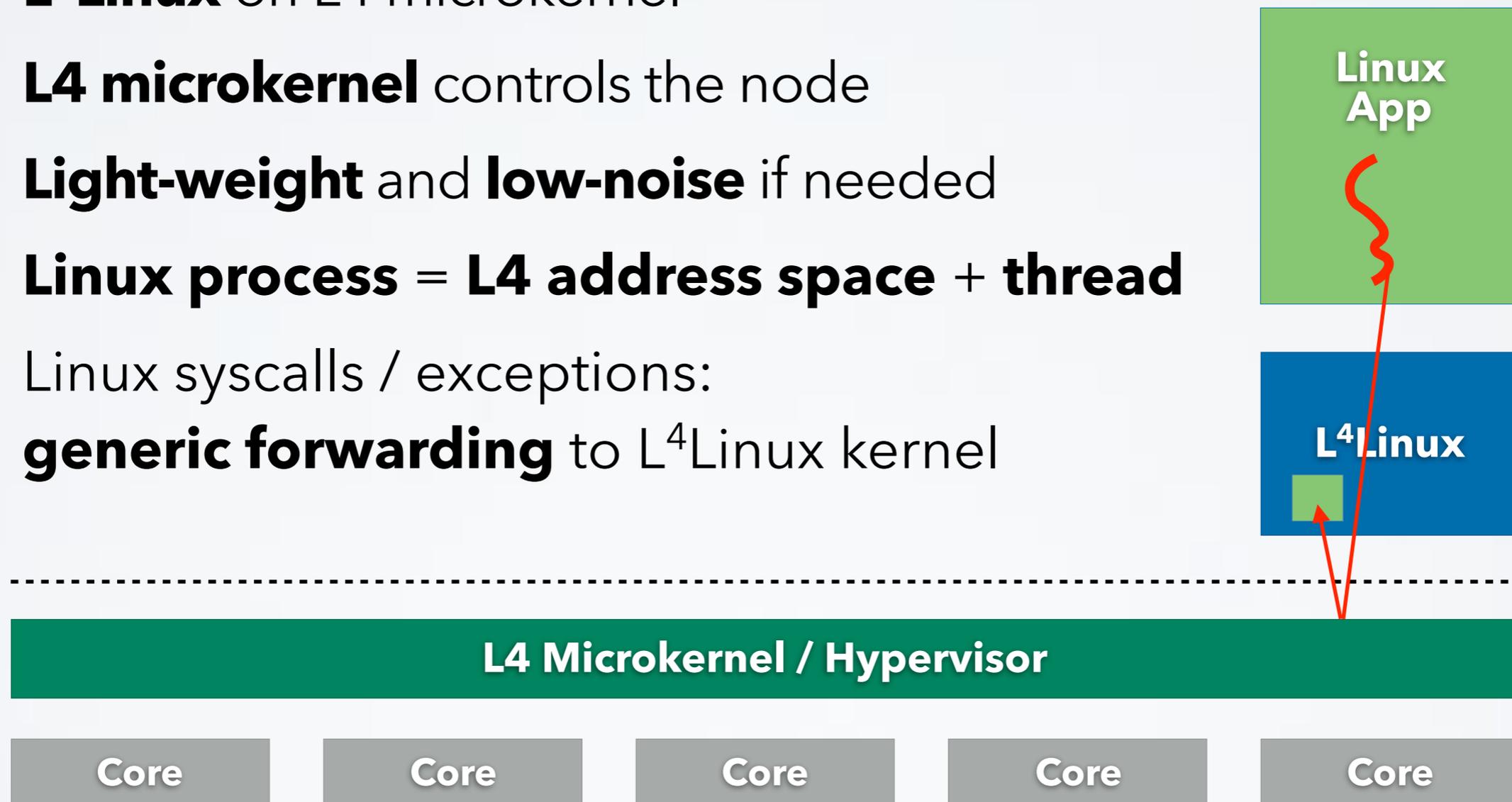




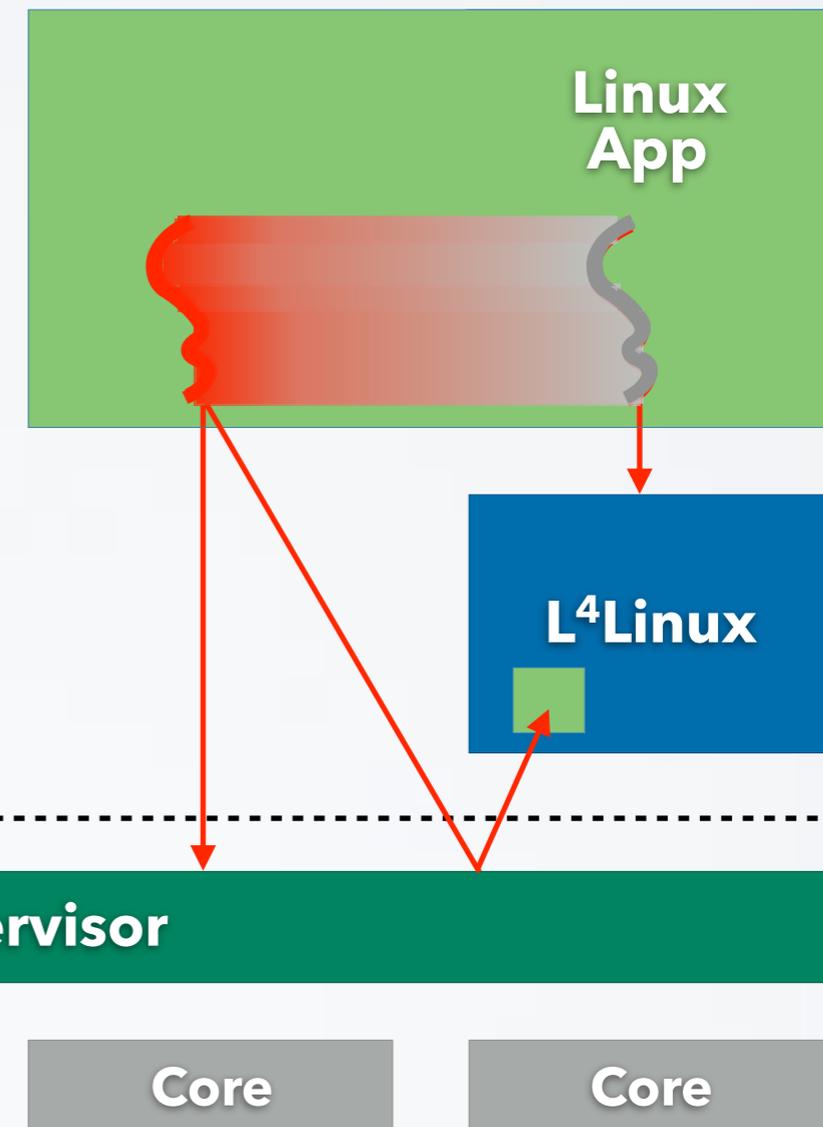
VERSION 2

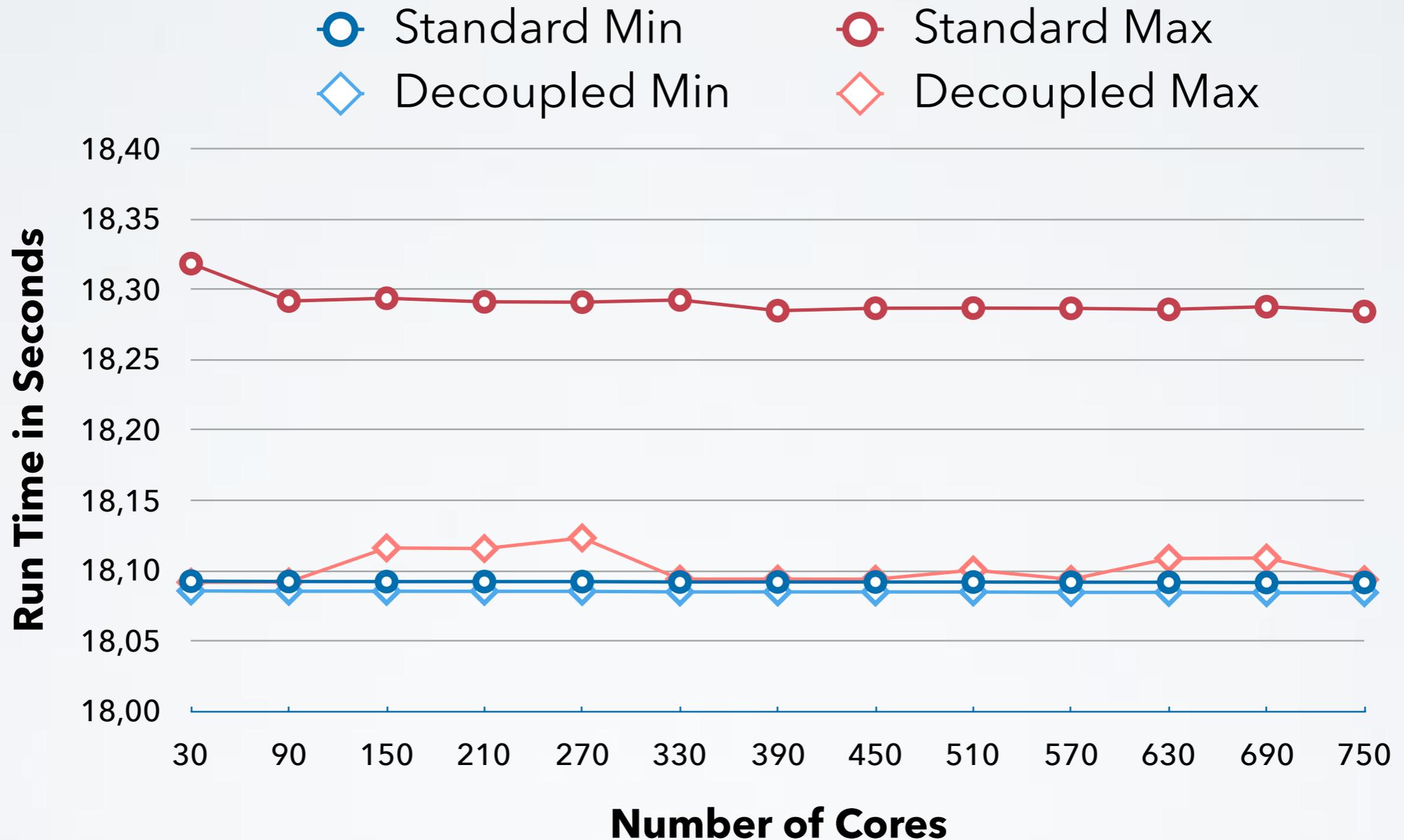
Decoupled Threads

- **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

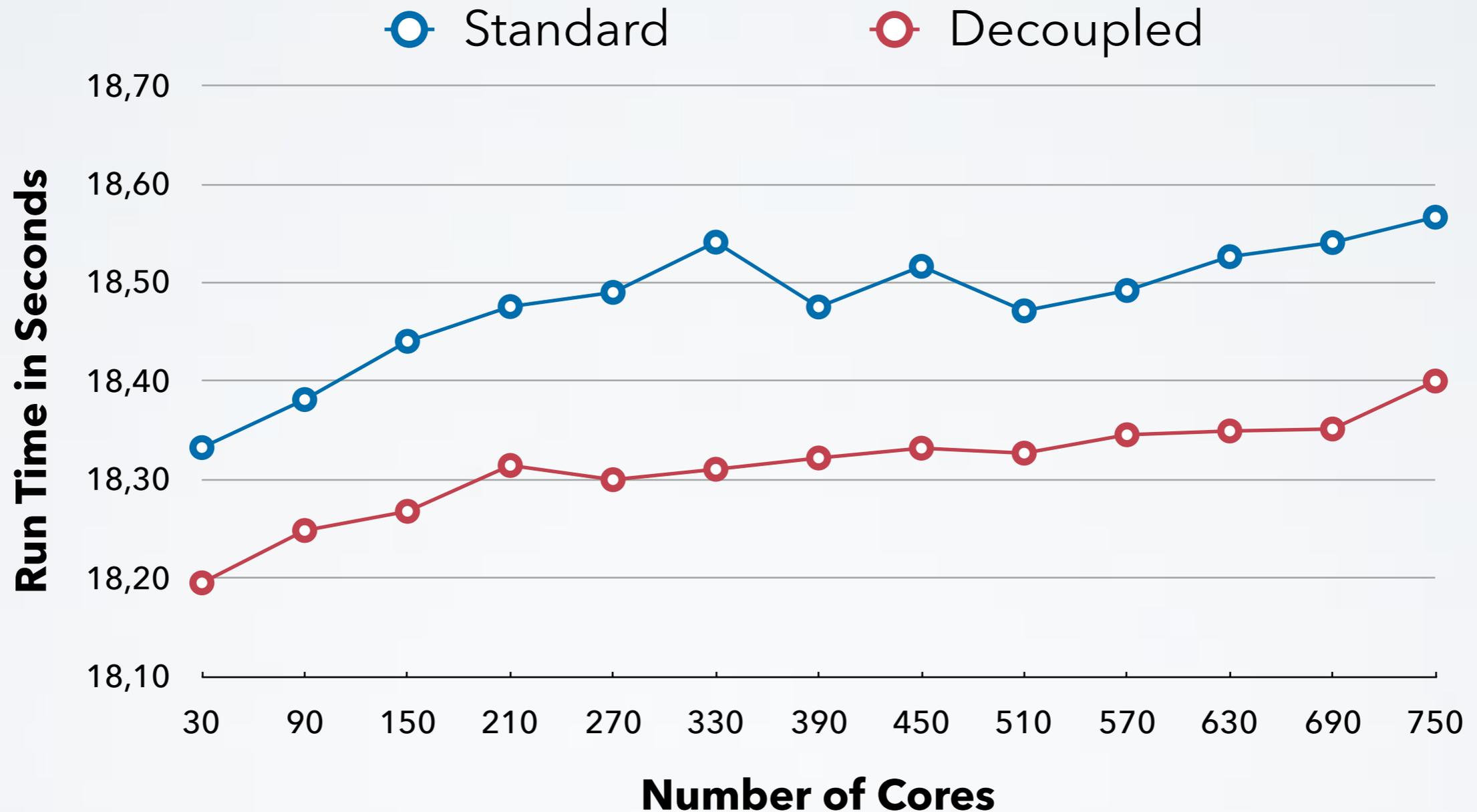


- **Decoupling:** move Linux thread to new L4 thread on its own core
- **Linux syscall:** Move back to Linux
- **L4 syscalls:**
 - Scheduling
 - Threads
 - Memory
- **Direct I/O** device access





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



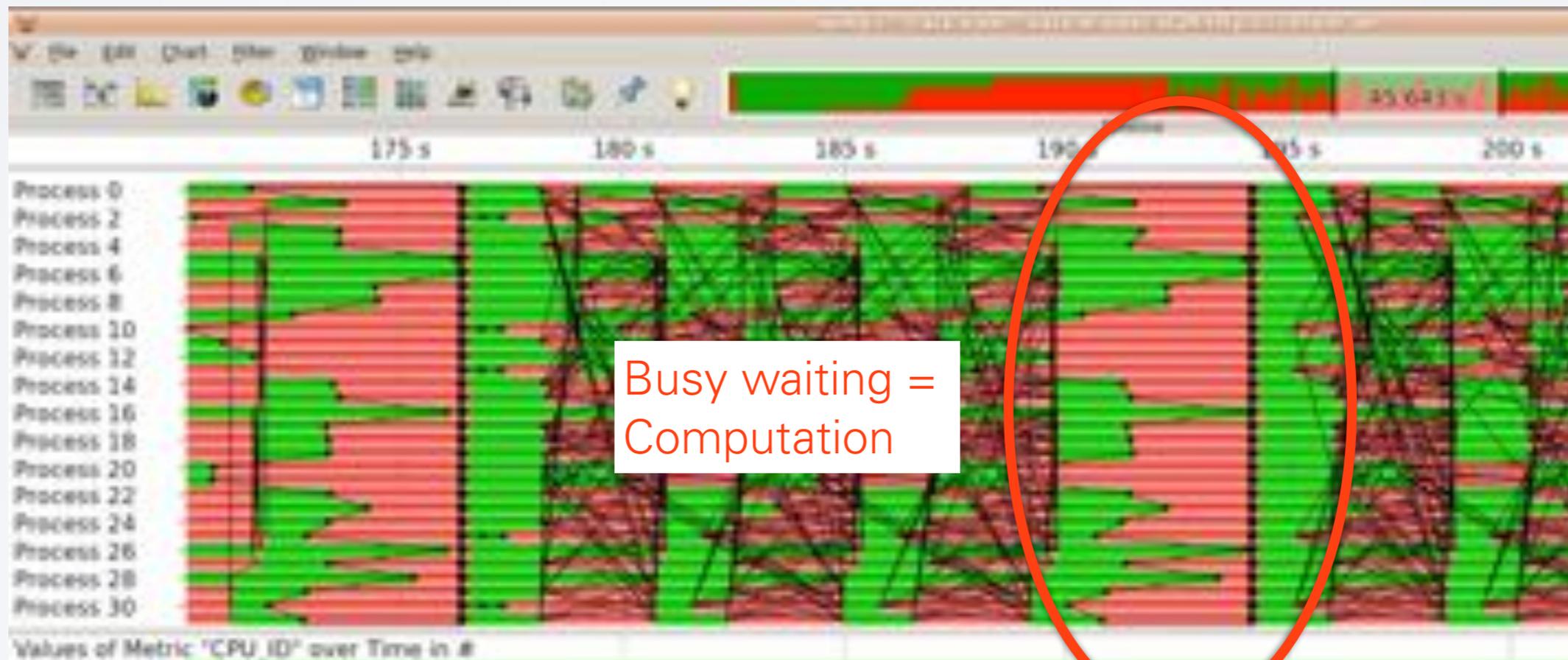
Adam Lackorzynski, Carsten Weinhold, Hermann Härtig, "Decoupled: Low-Effort Noise-Free Execution on Commodity Systems", ROSS 2016, June 2016, Kyoto, Japan



VERSION 2.5

Decoupled Interrupts

WIP



OVERSUBSCRIPTION

■ Polling

■ Blocking

Application: **COSMO-SPECS+FD4**

1000 s
800 s
600 s
400 s
200 s
0 s



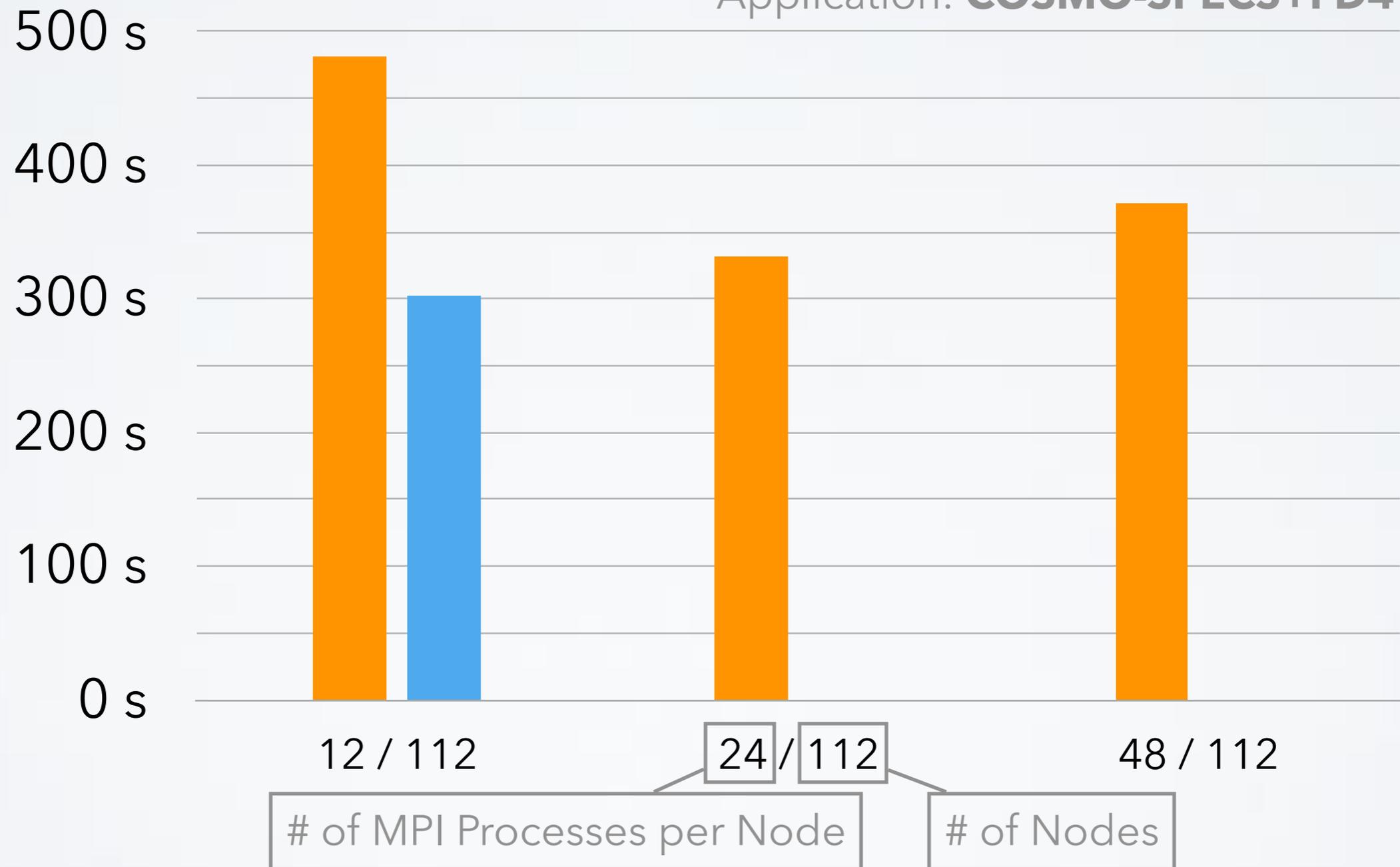
of MPI Processes per Node

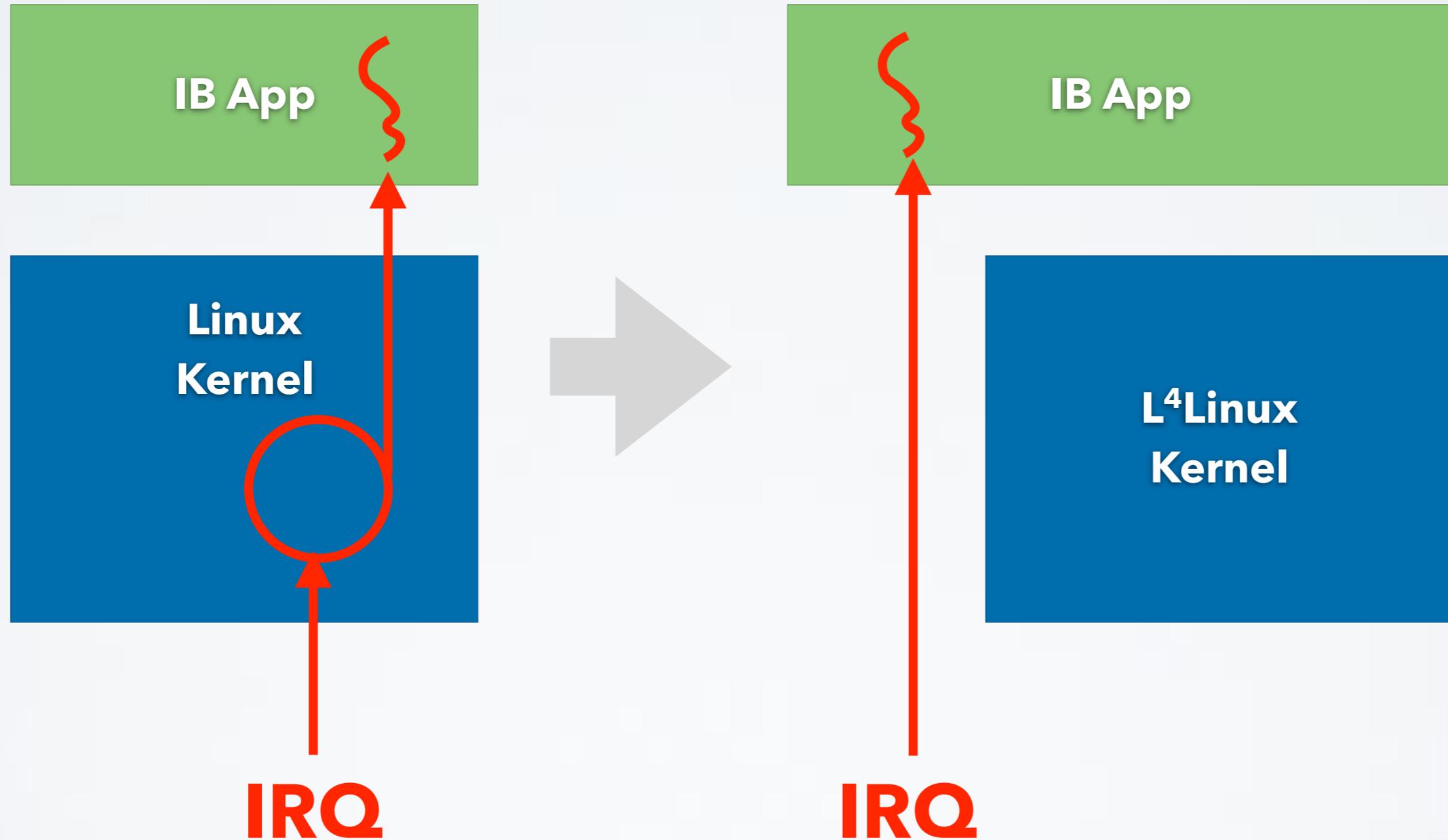
of Nodes

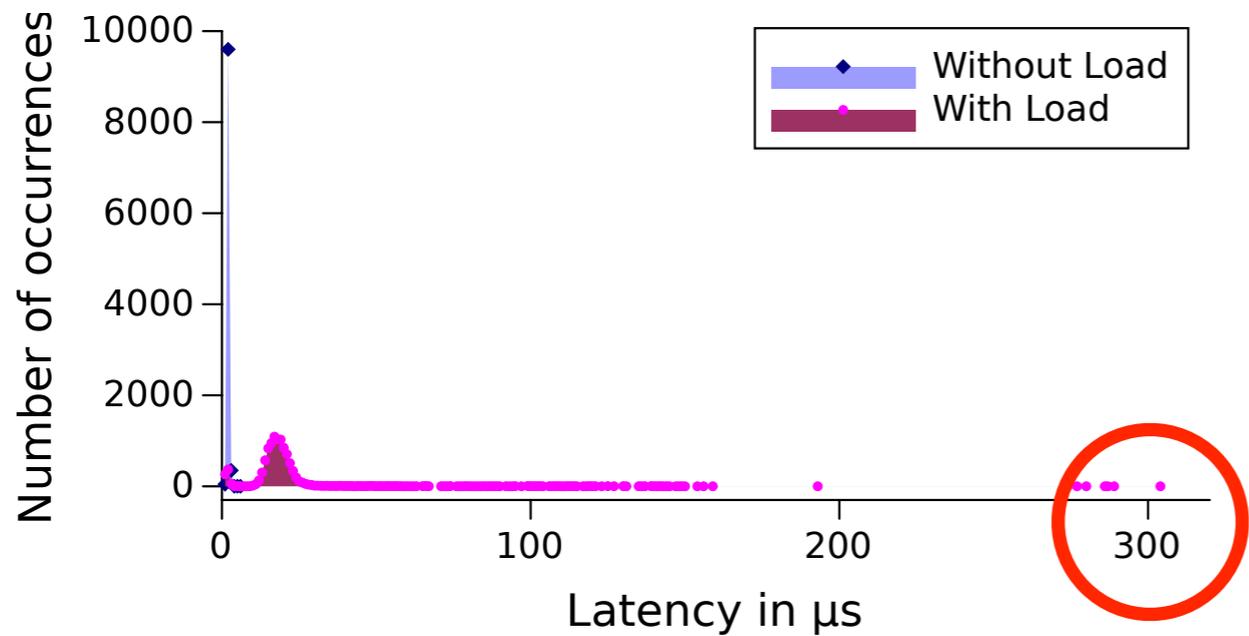
OVERSUBSCRIPTION

Unbalanced, no HT Balanced (baseline), no HT

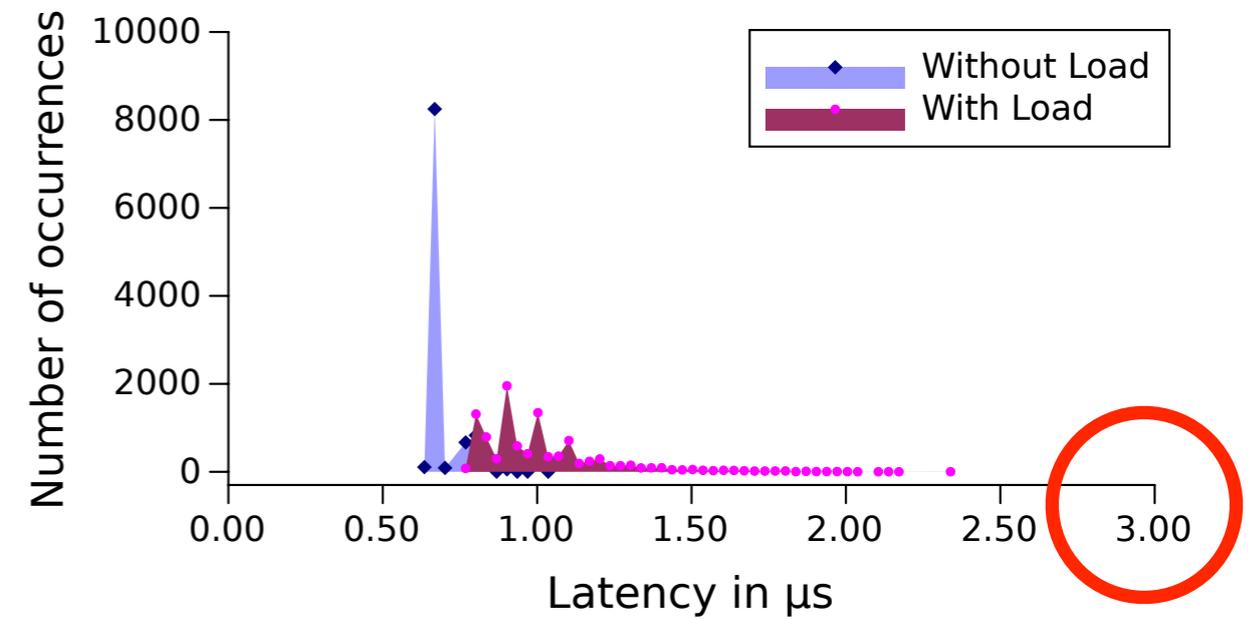
Application: **COSMO-SPECS+FD4**







Linux



L4

Wake from interrupt on L4/Nova: 900 cycles, 0.3 μs

(best case, on Intel Core i7-4770 CPU @ 3.40GHz)

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

- L4 Microkernel + L⁴Linux + Decoupling
- Binary compatible + `sysfs` interface
- No modification to MVAPICH2
- Lessons learned:
 - Maximize reuse + minimize critical path
 - Methodology: Start from Linux, not from L4

Source: <https://l4re.org>

Project: <https://ffmk.tudos.org>