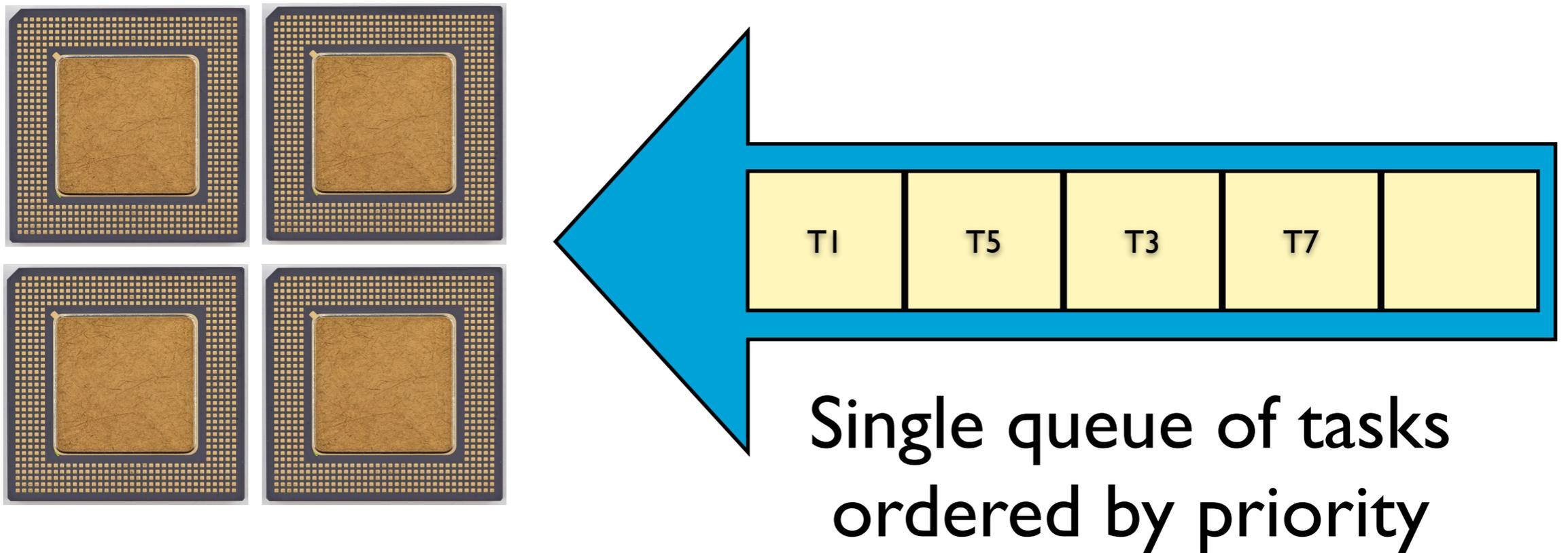# Scaling Global Scheduling with Message Passing

Max
Planck
Institute
for
Software Systems

**Felipe Cerqueira**
Manohar Vanga
Björn Brandenburg

# Global Scheduling



Single queue of tasks ordered by priority
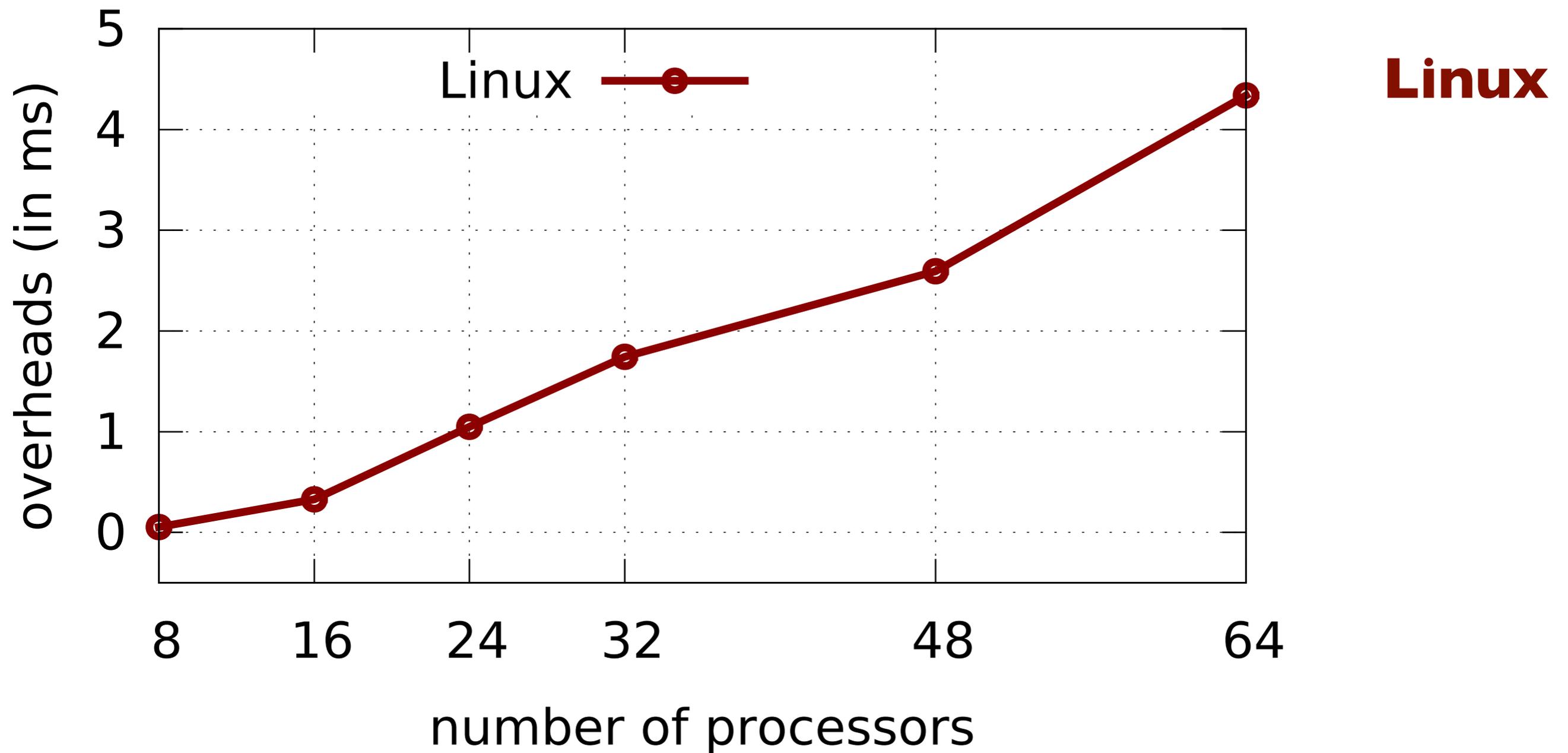
**Tasks can execute on any processor**

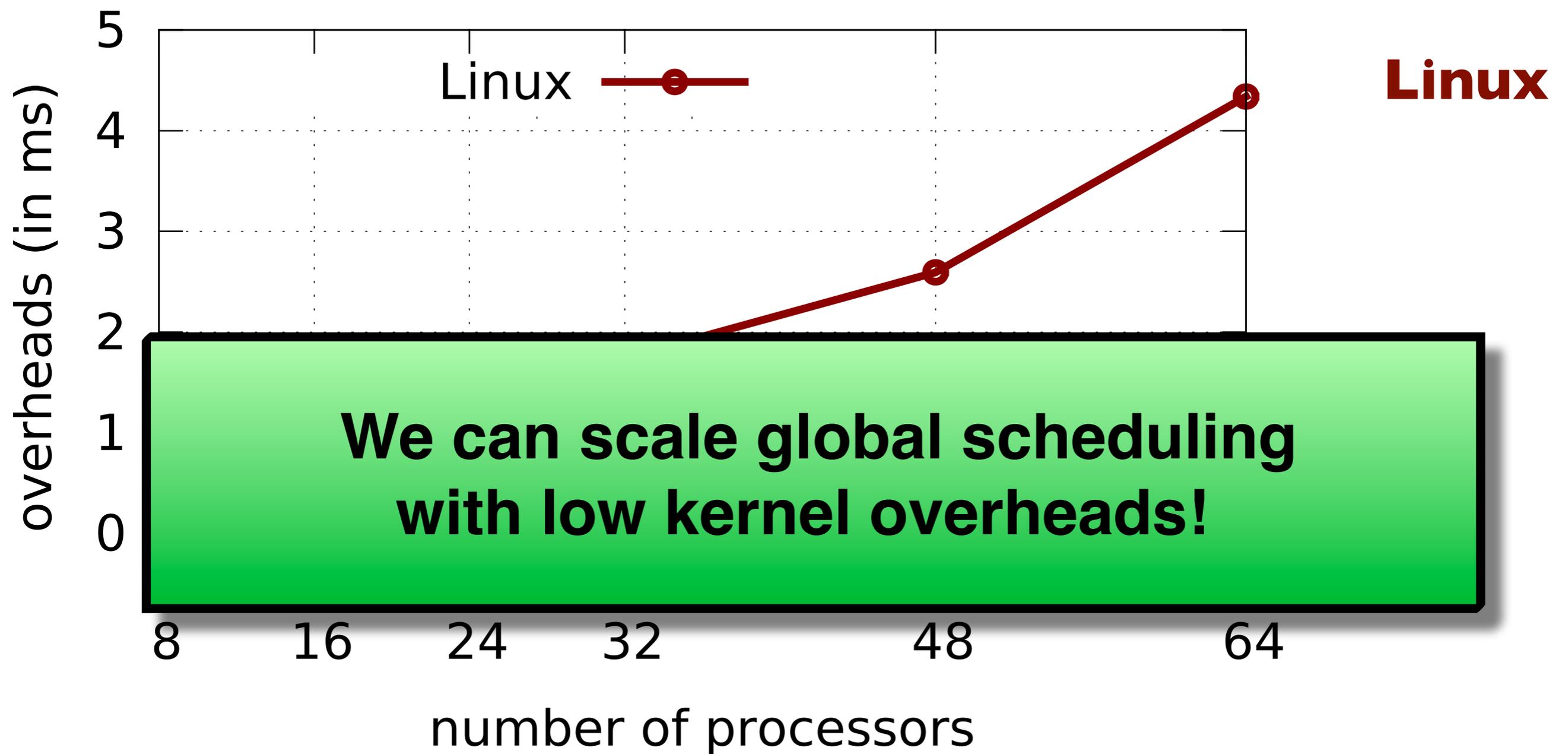# Global Scheduling

In **theory,** desirable analytical properties
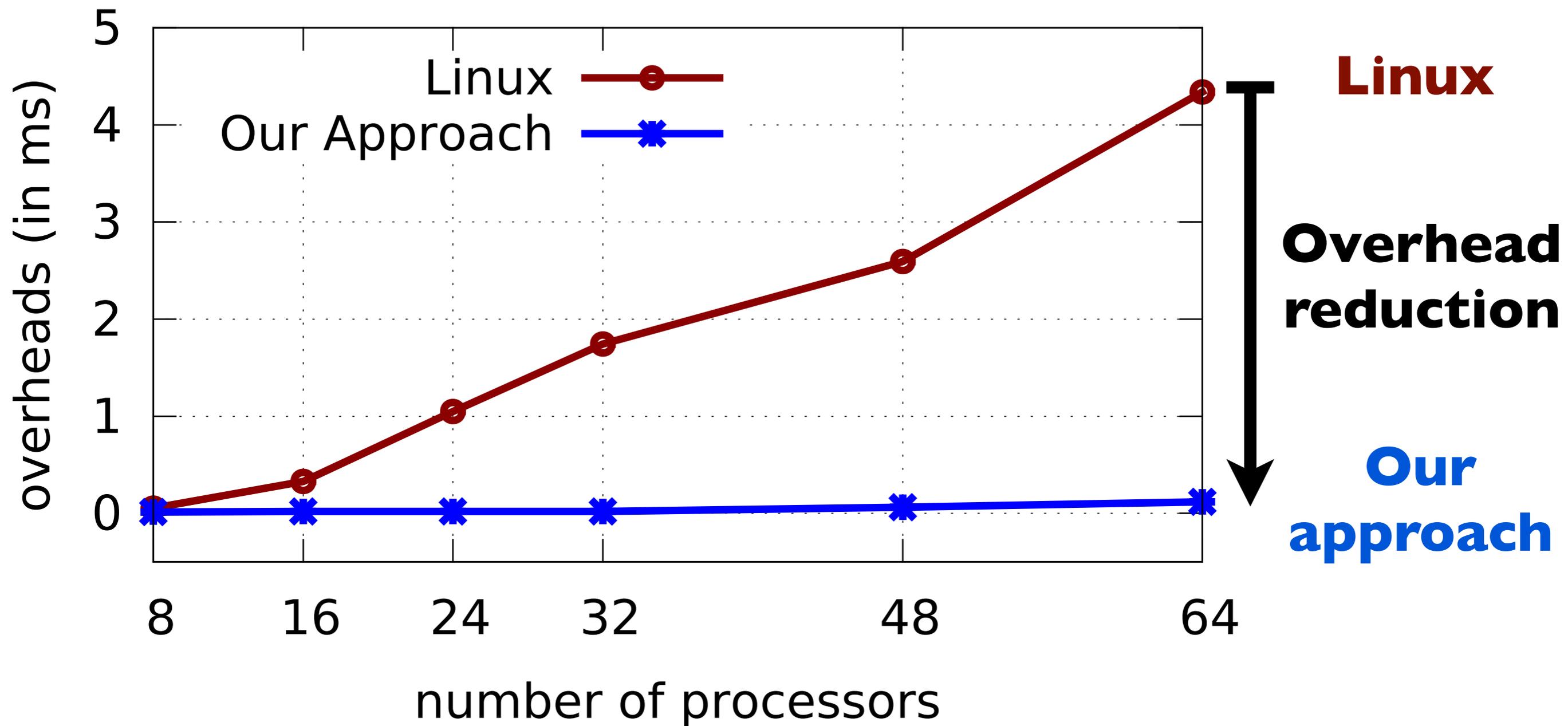
In **practice**,

   not scalable due to high overheads

# Making Global Scheduling Practical

# Making Global Scheduling Practical



**Linux**

We can scale global scheduling
with low kernel overheads!

# Making Global Scheduling Practical

# This Talk
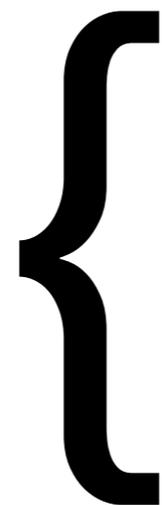
1) Why global scheduling?

2) Current implementations

3) Root causes of overhead

4) How to scale global scheduling?

5) Evaluation

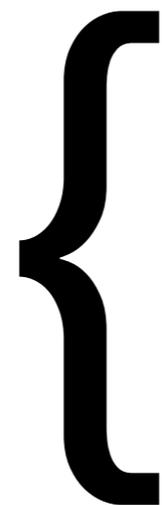# This Talk

# Why Global Scheduling?

Reasons {
- Optimal schedulers
- Work-conserving
- Soft-real-time
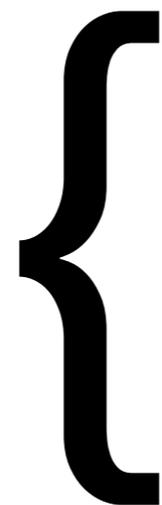- and more...

# Why Global Scheduling?

Reasons {
**Optimal schedulers**
Work-conserving
Soft-real-time
and more...

**Optimal** real-time schedulers are global
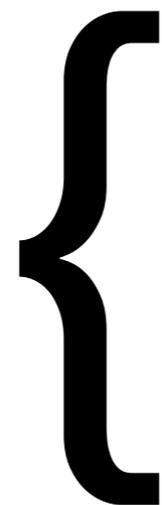
# Why Global Scheduling?

Reasons {
Optimal schedulers
**Work-conserving**
Soft-real-time
and more...

Good for **open** and **dynamic systems**
Resilient to overloads
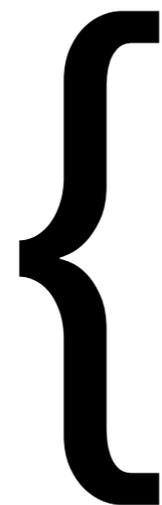
# Why Global Scheduling?

Reasons {
- Optimal schedulers
- Work-conserving
- **Soft-real-time**
- and more...

Some global schedulers guarantee **bounded tardiness** without utilization loss

# Why Global Scheduling?

Reasons {
Optimal schedulers
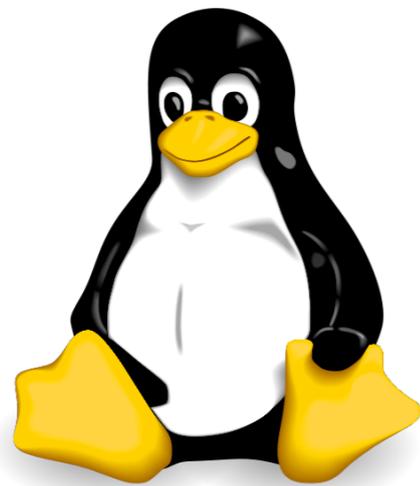Work-conserving
Soft-real-time
**and more...**

Supports **priority inheritance**
Useful in **race-to-idle** energy conservation

# Why Global Scheduling?

Reasons {
Optimal schedulers
Work-conserving
Soft-real-time
and more...
}

**Properties not fully guaranteed by Partitioned and Clustered Scheduling!**

# Global Schedulers
# in Practice



Default scheduler for Linux, QNX and VXWorks.

# This Talk

# This Talk

1) Why global scheduling?

2) **Current implementations**

3) Root causes of overhead
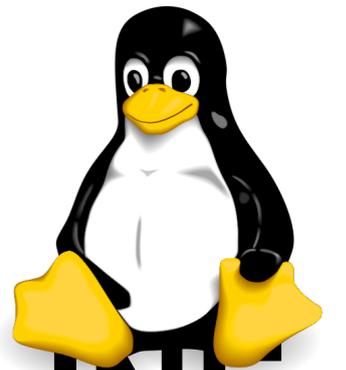
4) How to scale global scheduling?

5) Evaluation

**G-EDF as a representative of global scheduling**

# Comparing Two Extremes

**GSN-EDF**

**SCHED_DEADLINE**

# Comparing Two Extremes



**GSN-EDF**

Globally shared state, single lock

Distributed state, multiple locks
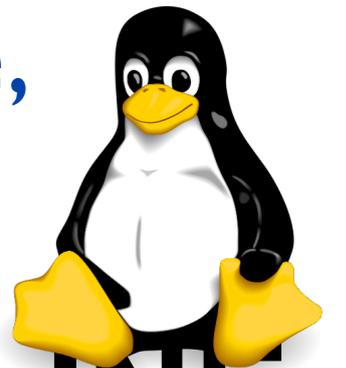
**SCHED_DEADLINE**

# Comparing Two Extremes

### LITMUS^RT
Linux Testbed for Multiprocessor Scheduling in Real-Time Systems

## GSN-EDF

Globally shared state,
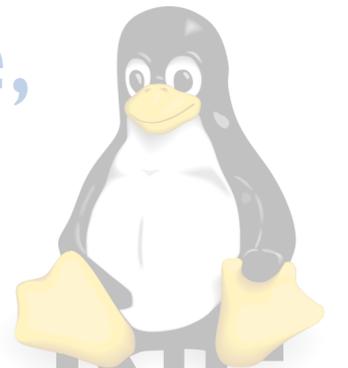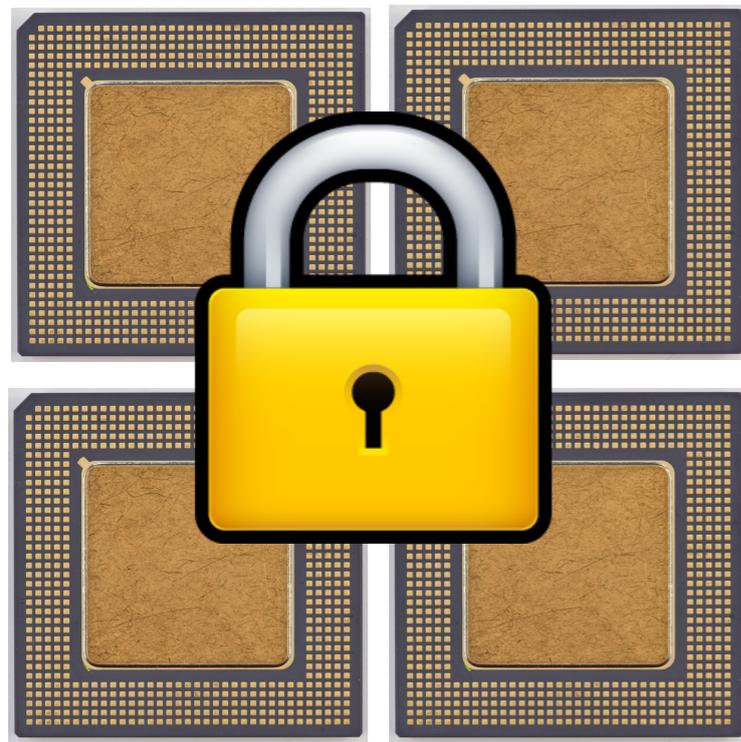single lock

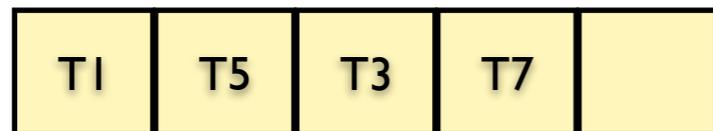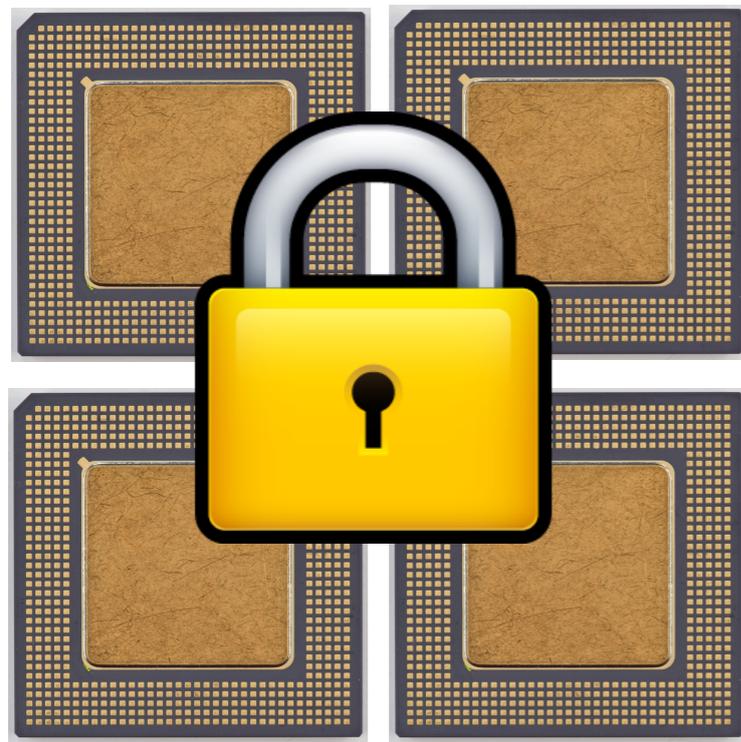Distributed state,
multiple locks

SCHED_DEADLINE

# GSN-EDF

**G**lobal-EDF with support for
**S**uspension-based protocols and O(1)
**N**on-preemptable sections

Link-based scheduler
(Block et al., 07)

| T1 | T5 | T3 | T7 | |
|----|----|----|----|--|

LITMUS^RT
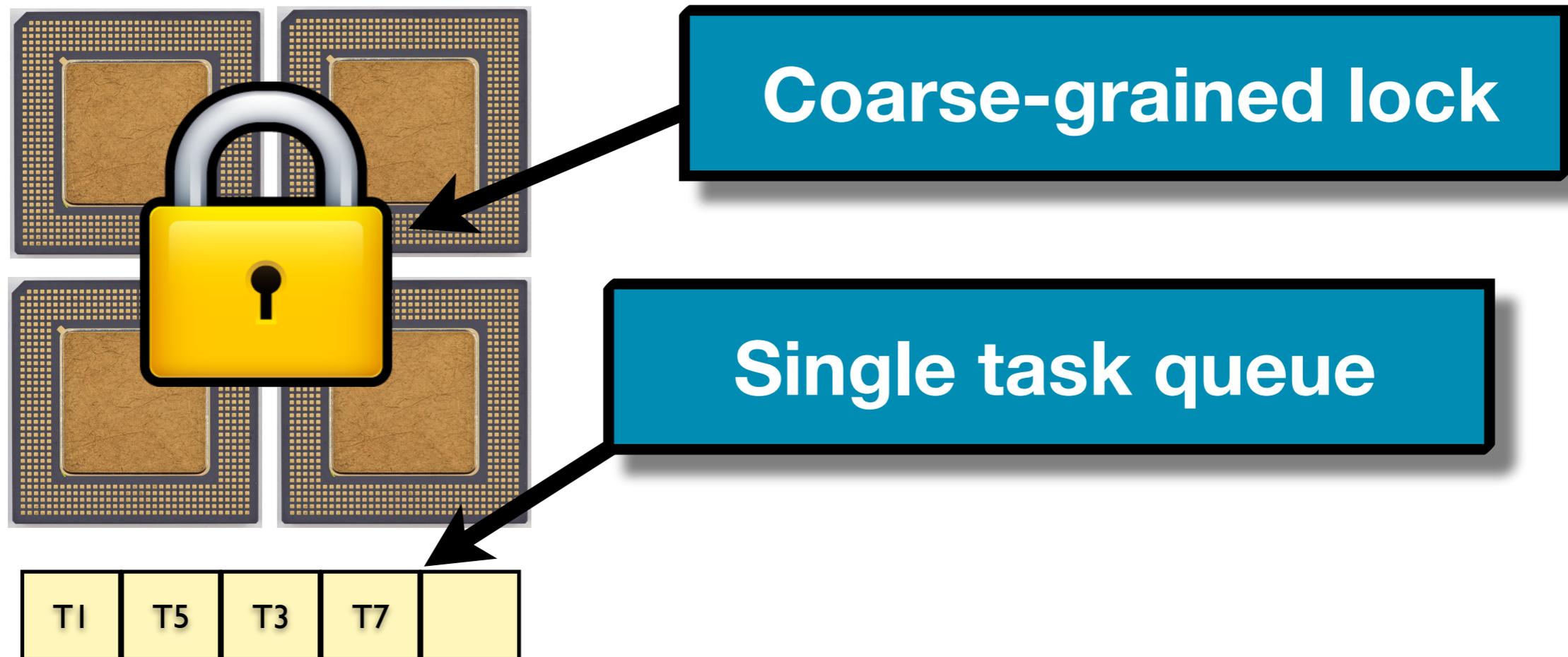Linux Testbed for Multiprocessor Scheduling in Real-Time Systems

# GSN-EDF



**G**lobal-EDF with support for
**S**uspension-based protocols and O(1)
**N**on-preemptable sections

Link-based scheduler
(Block et al., 07)
➡ allows simplified locking

# GSN-EDF



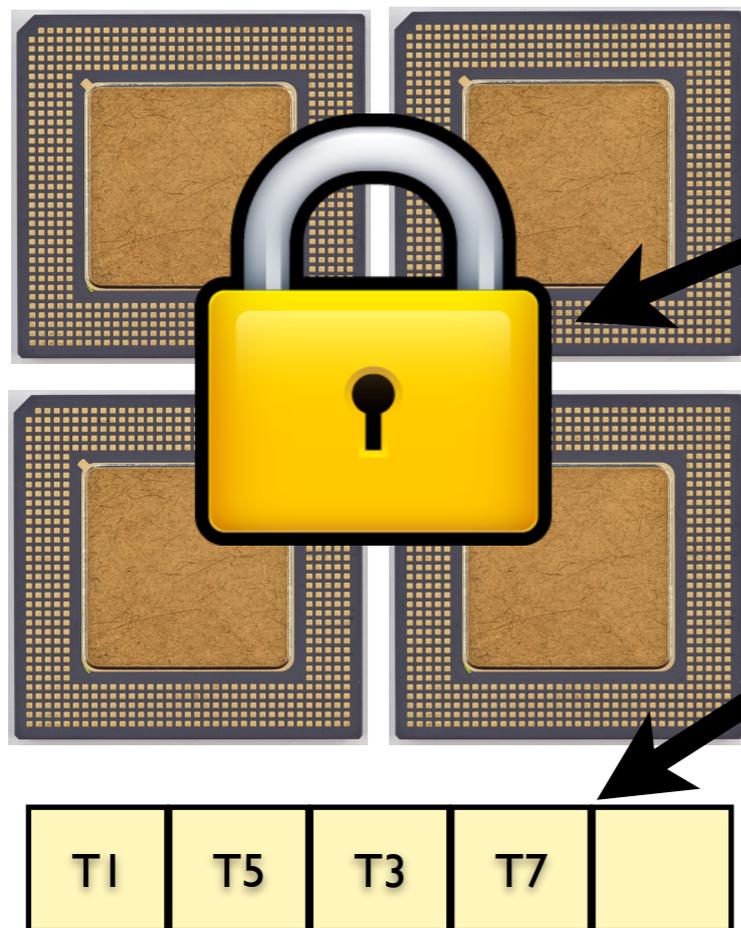**Coarse-grained lock**

**Single task queue**

| T1 | T5 | T3 | T7 | |

# GSN-EDF

**Coarse-grained lock**

**Single task queue**

| T1 | T5 | T3 | T7 | |
|----|----|----|----|--|

Simple
implementation!

LITMUS^RT

Linux Testbed for Multiprocessor Scheduling in Real-Time Systems

# GSN-EDF

**Coarse-grained lock**

**Single task queue**

| T1 | T5 | T3 | T7 | |

Simple implementation!
How does it scale?

LITMUS^RT

Linux Testbed for Multiprocessor Scheduling in Real-Time Systems
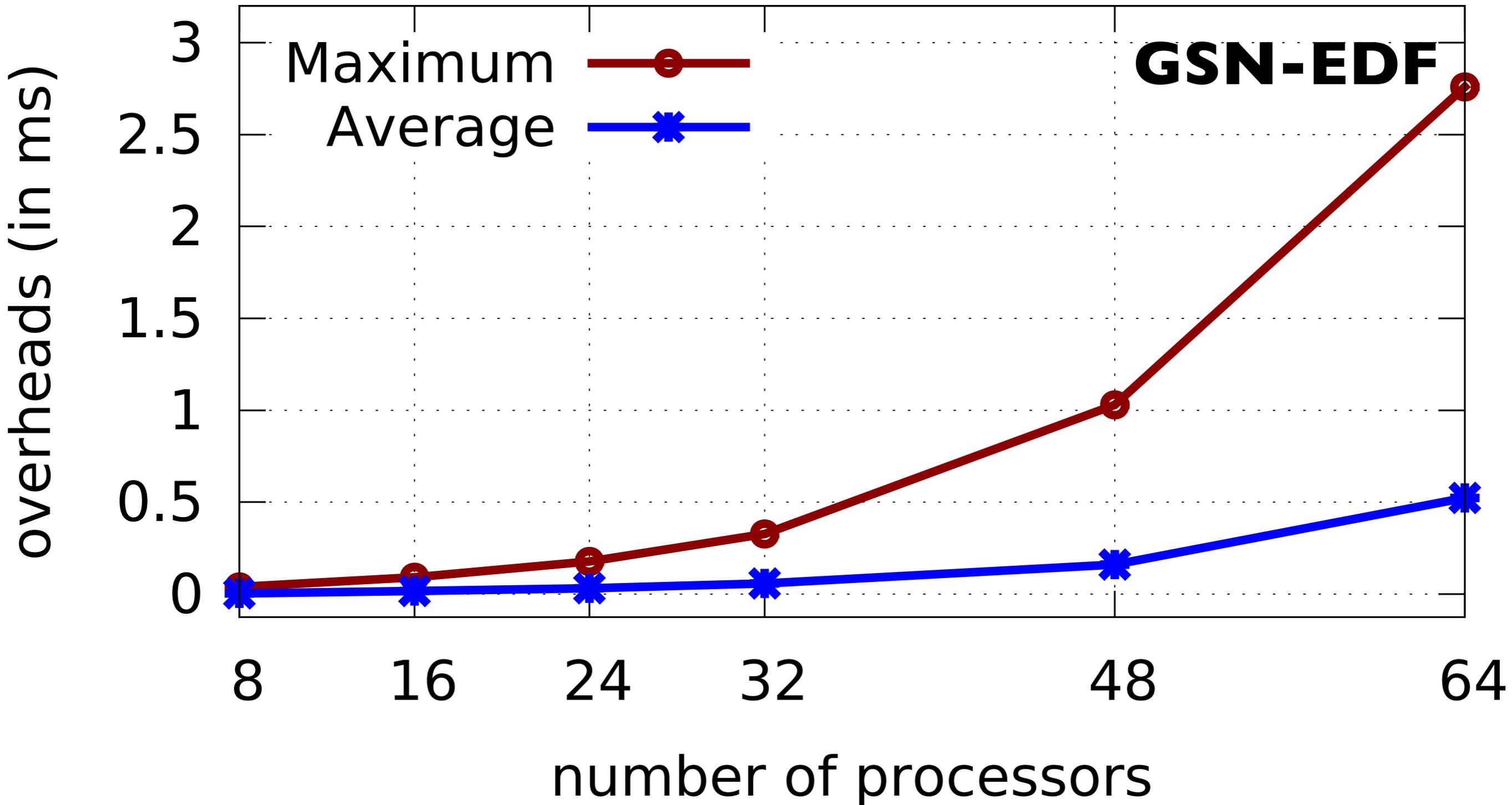
# Experimental Setup

- Intel Xeon X7550 @2.0GHz, with 64 cores

- Linux 3.10 with patches

  ➡ LITMUS^RT 2013.1 and SCHED_DEADLINE v8

- Lightweight build — disabled most drivers and debugging options
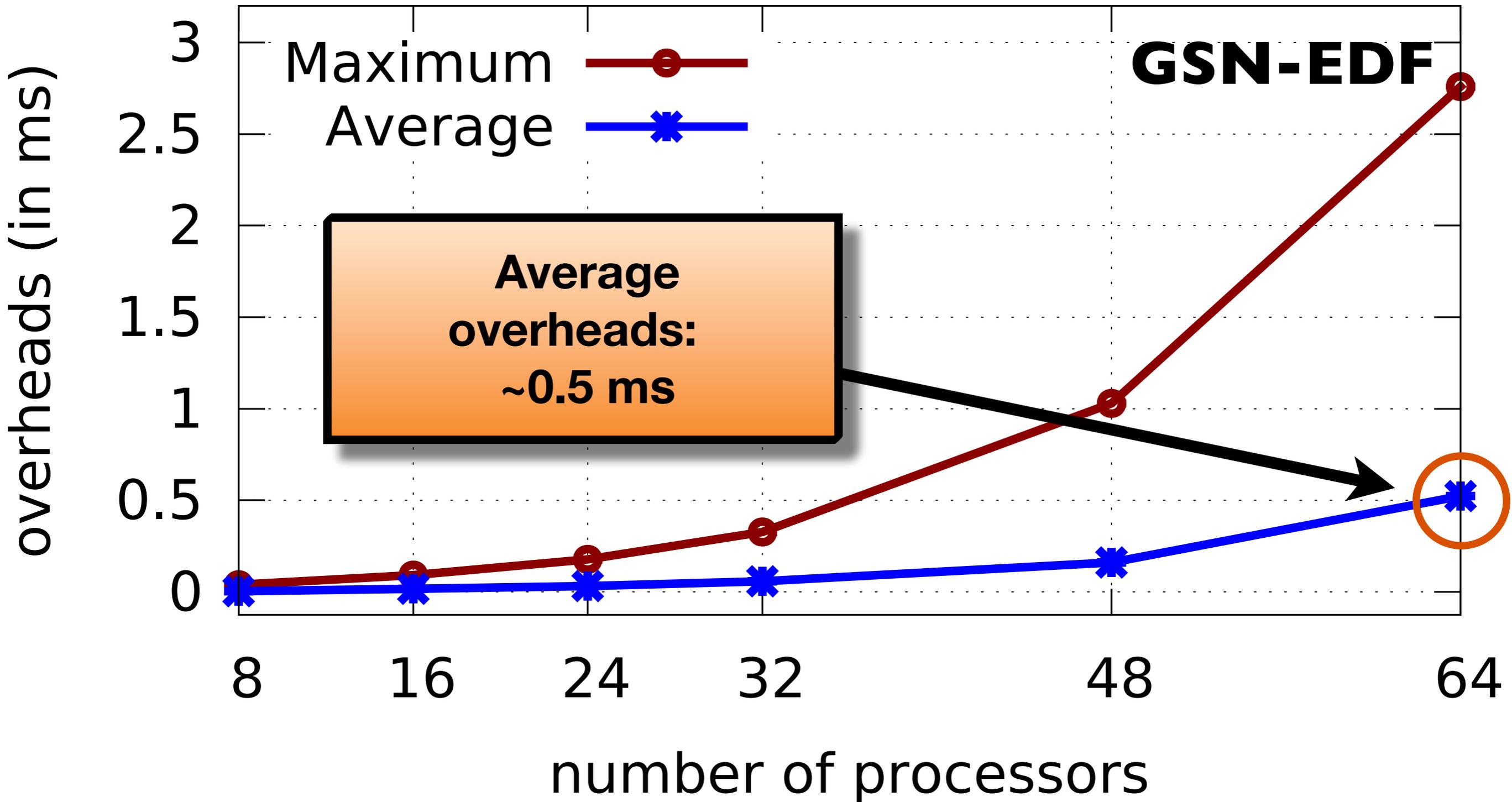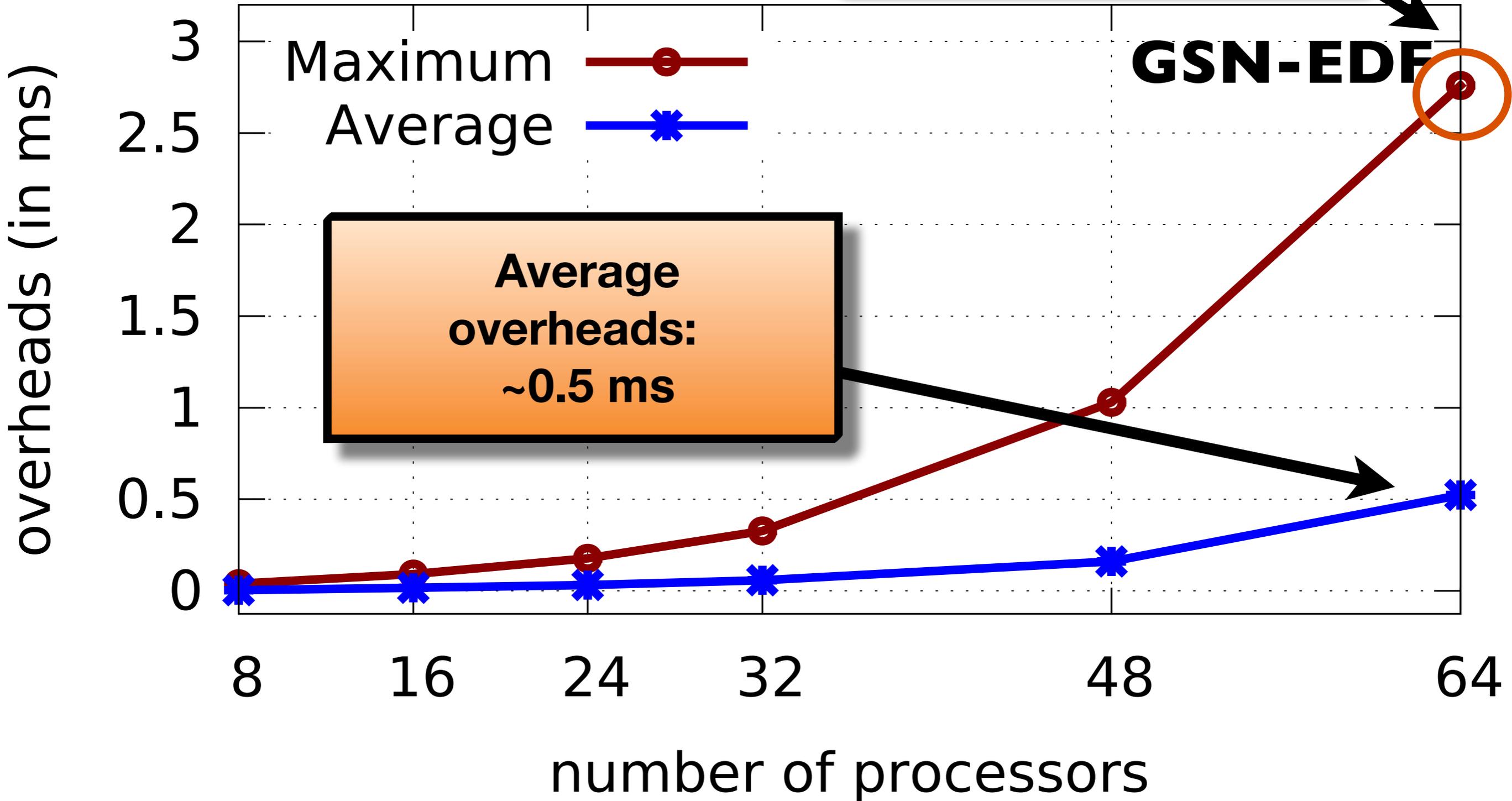
# Overheads Under GSN-EDF



overheads (in ms)

3
2.5
2
1.5
1
0.5
0

Scheduling Overheads
**Higher is worse**

8    16    24    32         48

number of processors

# Global Lock Does Not Scale!



**GSN-EDF**

Legend: Maximum, Average

y-axis: overheads (in ms), from 0 to 3

x-axis: number of processors — 8, 16, 24, 32, 48, 64

# Global Lock Does Not Scale!

# Global Lock Do[____]e!



For 64 CPUs, maximum overheads of ~3 ms

Average overheads: ~0.5 ms

GSN-EDF

- Maximum
- Average

overheads (in ms)

number of processors

# Comparing Two Extremes

LITMUS<sup>RT</sup>

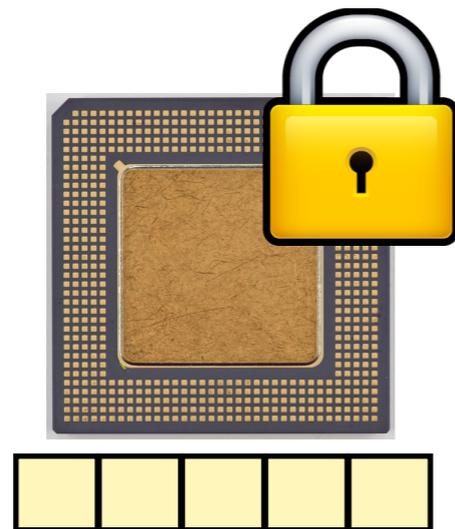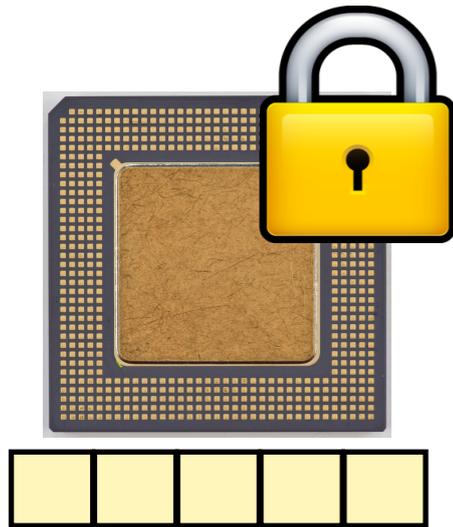Linux Testbed for Multiprocessor Scheduling in Real-Time Systems
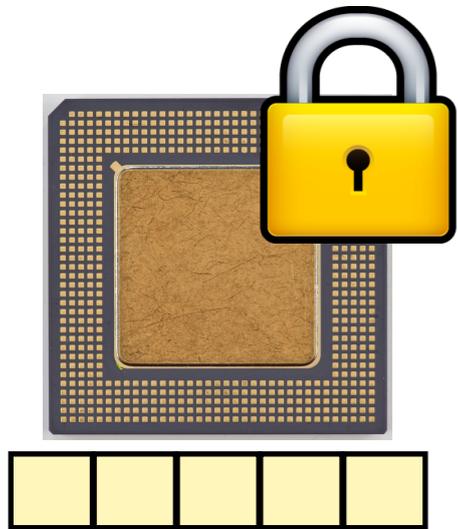
**GSN-EDF**

Globally share state, single lock
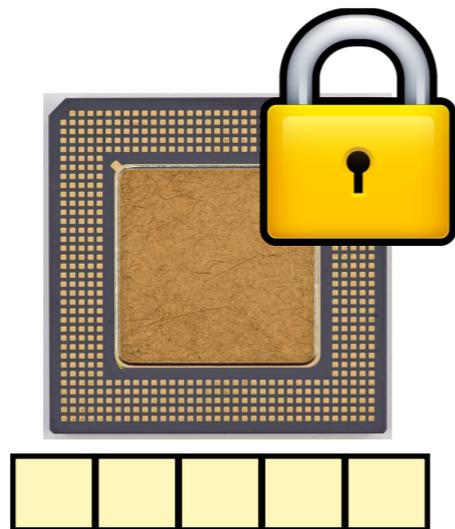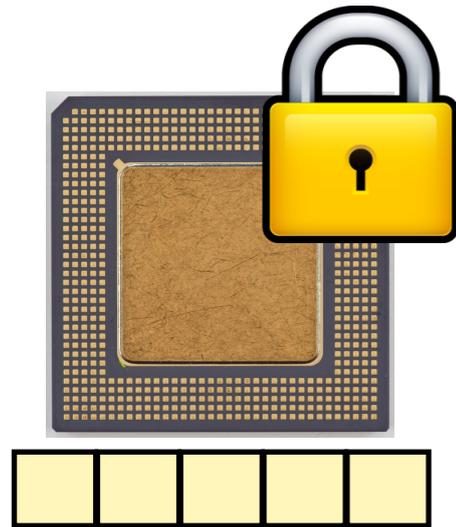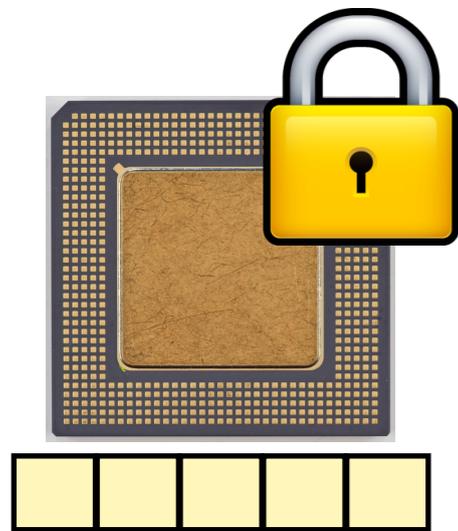
Distributed state, multiple locks

**SCHED_DEADLINE**

# SCHED_DEADLINE

Design inherited
from Linux scheduler

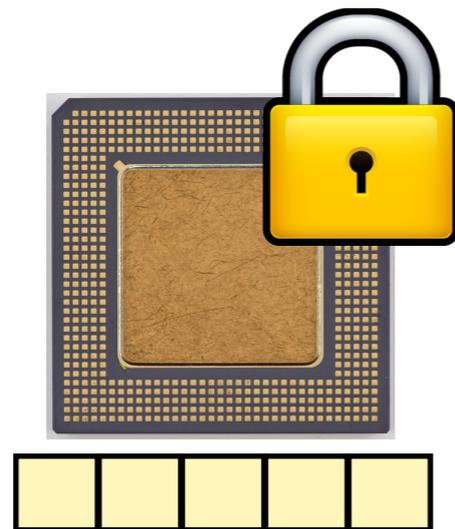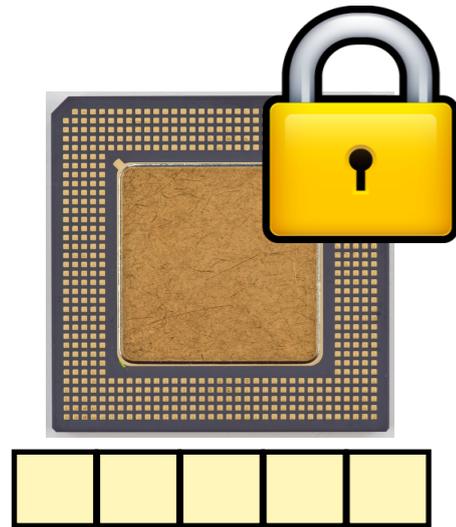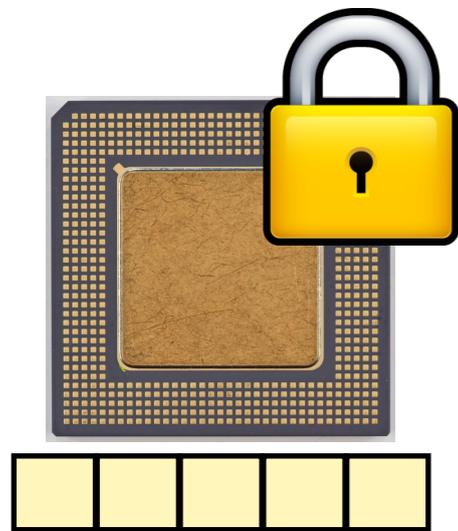# SCHED_DEADLINE

Design inherited
from Linux scheduler

**Per-CPU locks**

**Per-CPU task queues**

# SCHED_DEADLINE

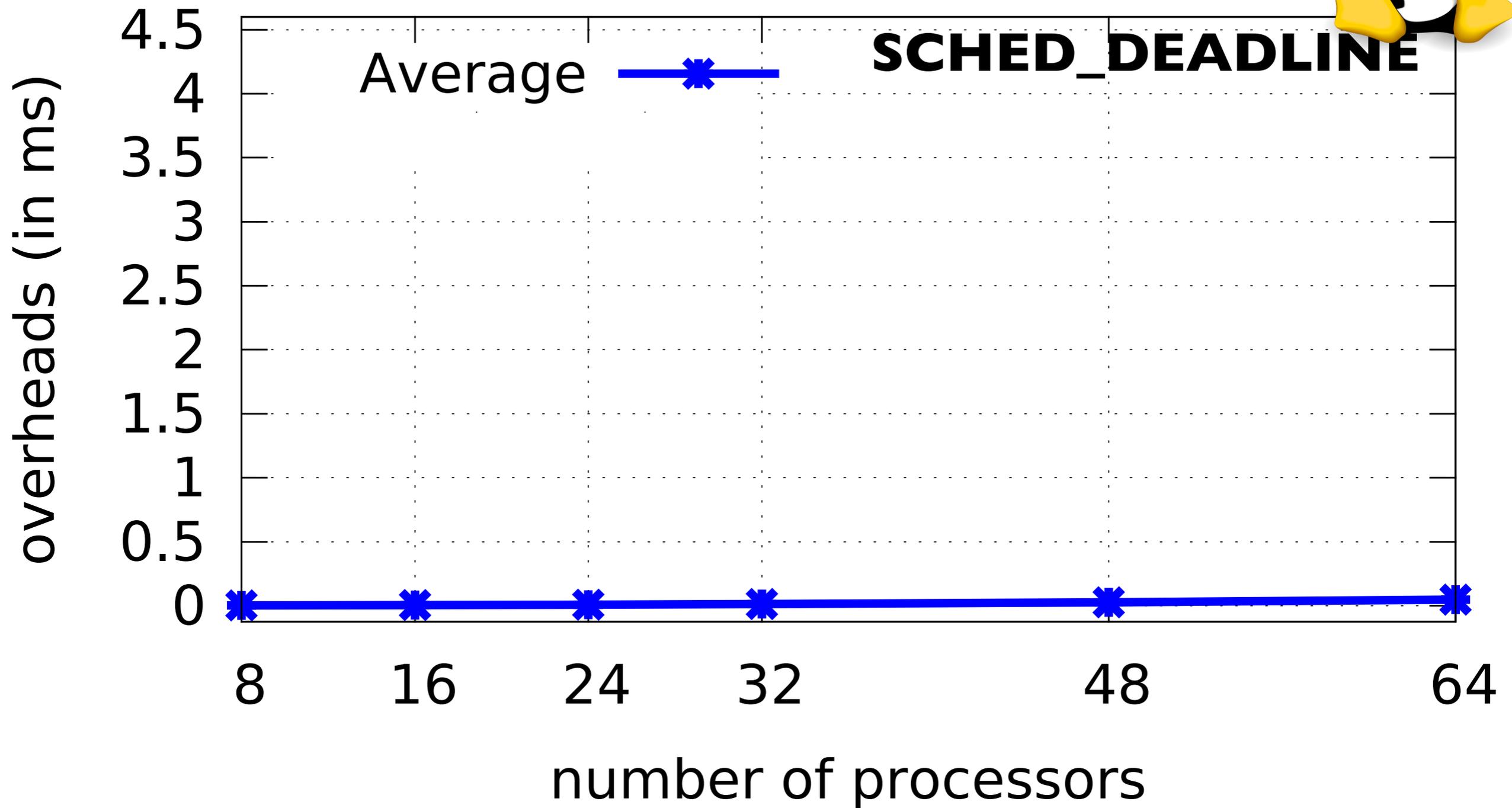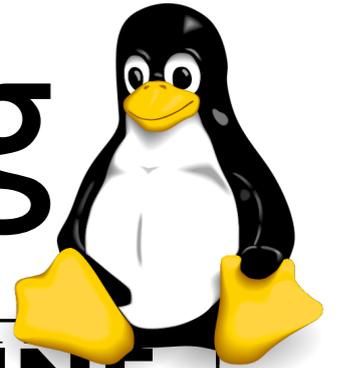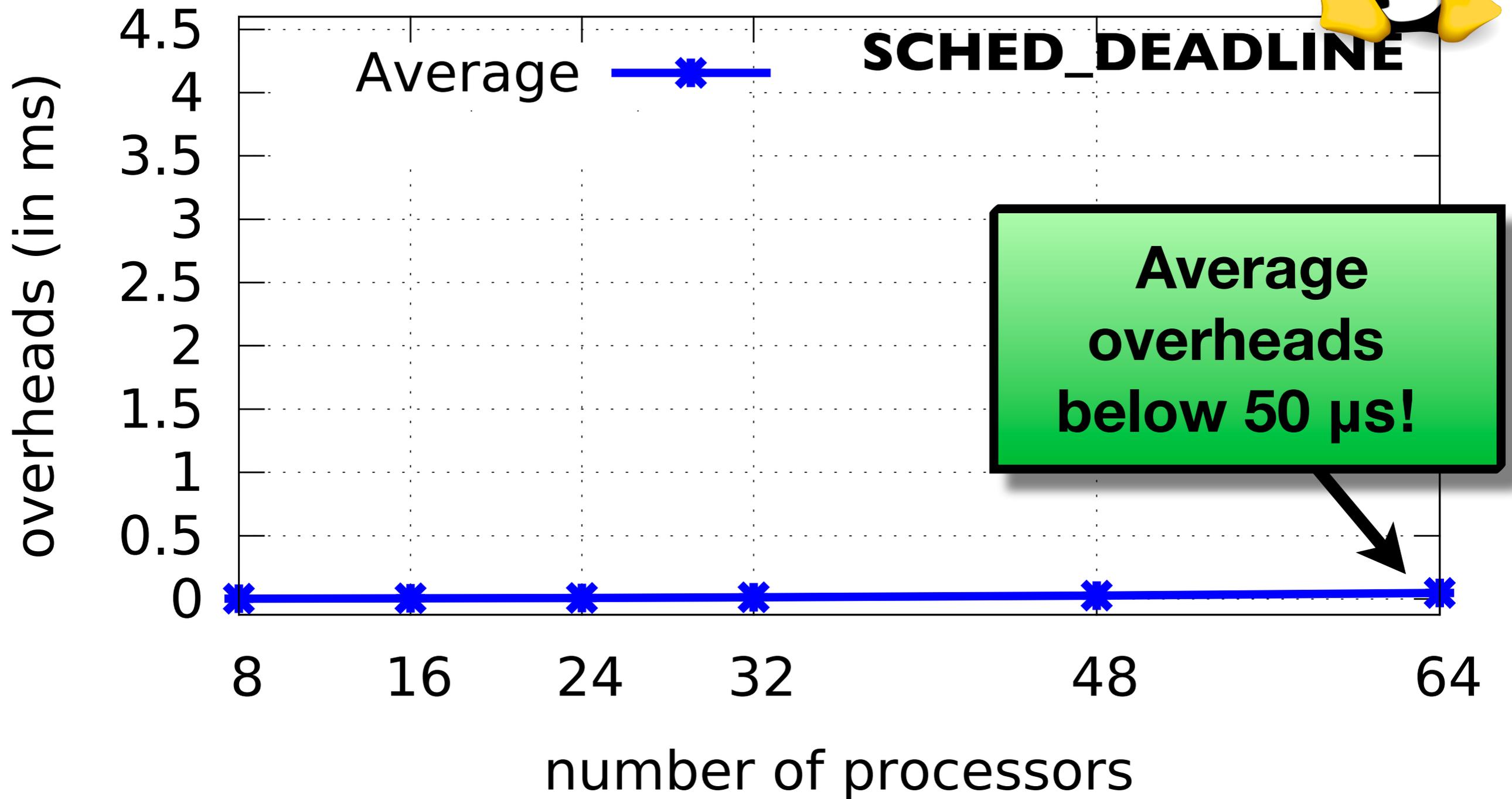Design inherited
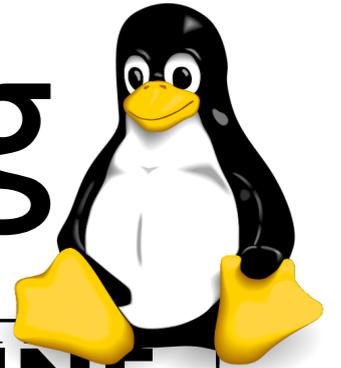from Linux scheduler

**Per-CPU locks**

**Per-CPU task queues**

Intuition:
Fine-grained locking
decreases contention

# Benefit of Fine-Grained Locking

overheads (in ms)

**SCHED_DEADLINE**

Average ✱——✱

4.5
4
3.5
3
2.5
2
1.5
1
0.5
0

8    16    24    32         48              64

number of processors

# Benefit of Fine-Grained Locking



**overheads (in ms)** vs **number of processors**

SCHED_DEADLINE

— Average

Average overheads below 50 µs!

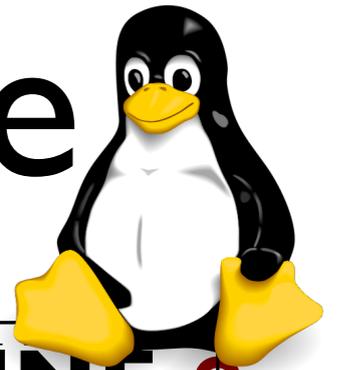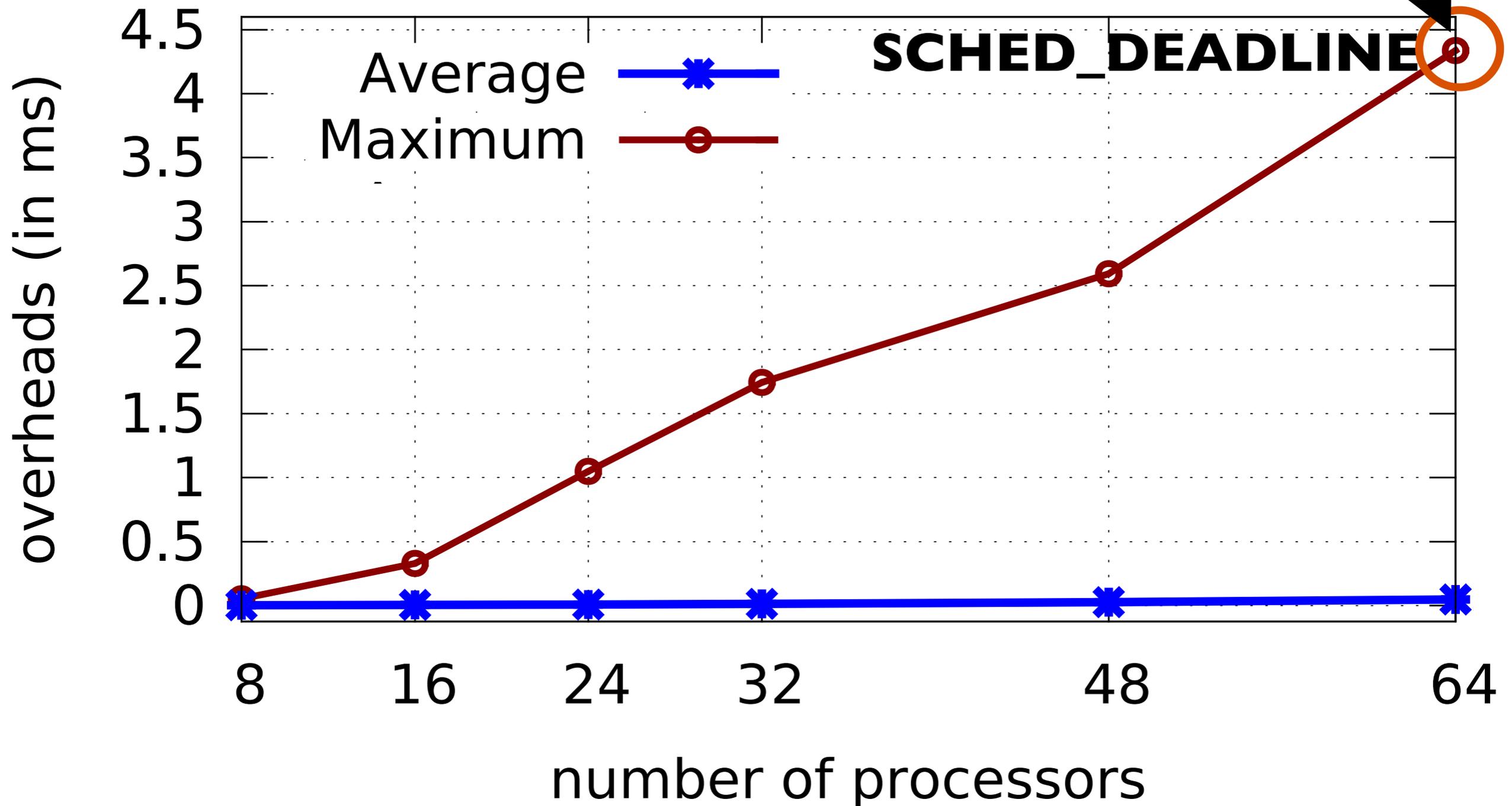# Fine-Grained Locking
# Fails in the Worst Case

# Fine-Grained Locking
# Fails in the Worst Case
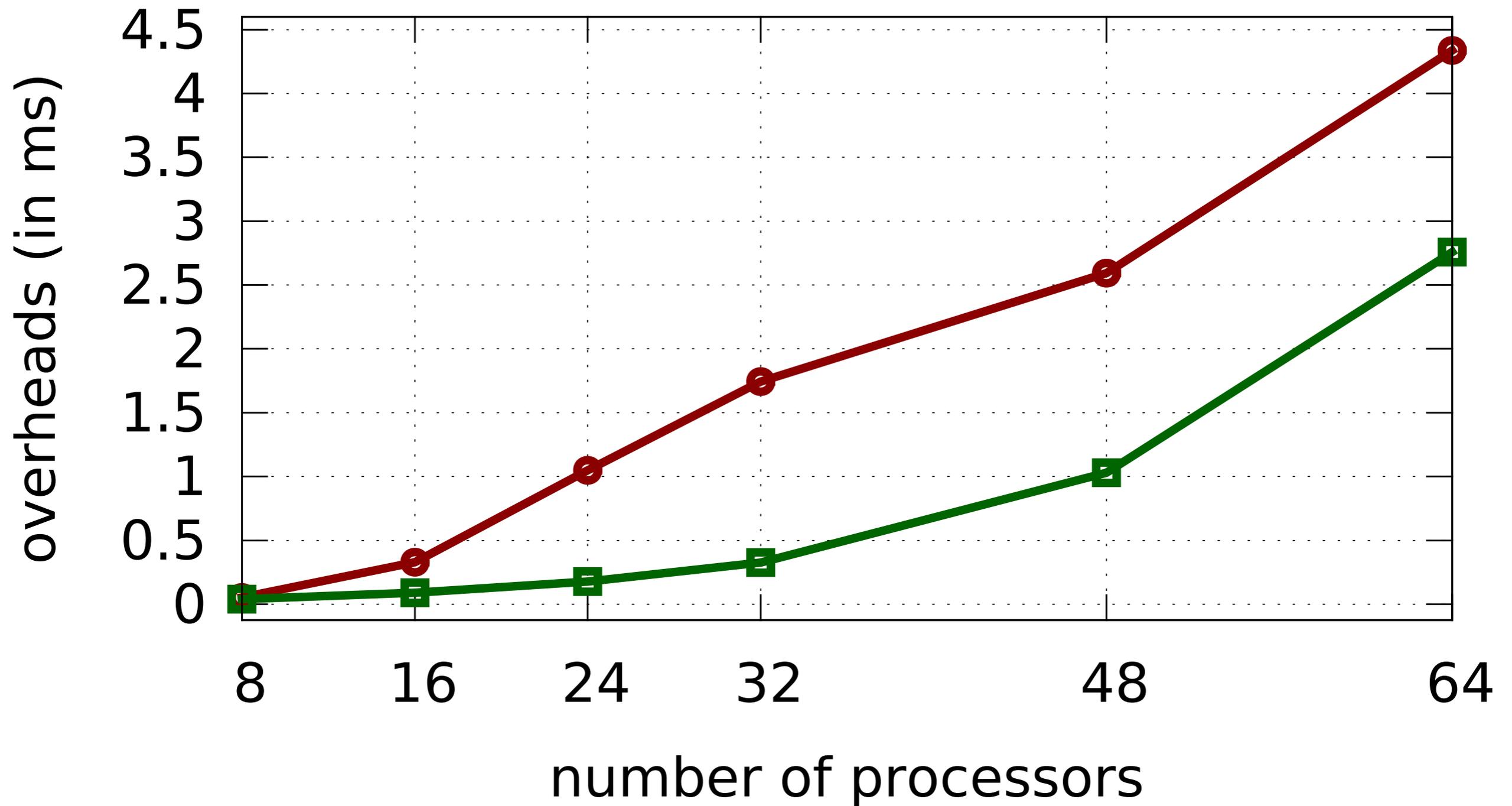
# Fine-Grained Fails in the Wo[rst]

**Very high overheads in the worst case!**

# Fine-Grained vs. Coarse-Grained Locks

overheads (in ms) vs. number of processors

SCHED_DEADLINE (Max)

GSN-EDF (Max)

# Fine-Grained vs. Coarse-Grained

**Both approaches do not scale in the worst case!**



**SCHED_DEADLINE (Max)**

**GSN-EDF (Max)**

overheads (in ms)

number of processors

# This Talk

1) Why global scheduling?

2) Current implementations

**3) Root causes of overhead**

4) How to scale global scheduling?

5) Evaluation

# Fine-grained Locking:
## Average Case

# Fine-grained Locking:
## Average Case

# Fine-grained Locking:
## Average Case

# Fine-grained Locking:
## Average Case



Task

Lock

# Fine-grained Locking:
## Average Case



Lock

Lock

**Low contention!**

# Fine-grained Locking:
# Worst Case

Locking **every** processor:
   O(m) iterations

# Fine-grained Locking:
# Worst Case

Locking **every** processor: O(m) iterations

O(m) processors **already** waiting for this lock

wait queue

Lock

Task

# Fine-grained Locking:
# Worst Case

Locking **every** processor:
O(m) iterations

O(m) processors **already**
waiting for this lock

# Fine-grained Locking:
# Worst Case

**Locking every processor:**
**O(m) iterations**

**O(m) processors already**
**waiting for this lock**



Task

Lock

wait queue

# Fine-grained Locking:
# Worst Case

Locking **every** processor:
O(m) iterations

O(m) processors **already**
waiting for this lock



Task

Lock

wait queue

# Fine-grained Locking:
# Worst Case

Locking **every** processor:
O(m) iterations

O(m) processors **already**
waiting for this lock

O(m) iterations x O(m) blocking
= **quadratic blocking times**

Lock

wait queue

# Peak Contention

**Observation #1:**
**Peak Contention** is more important
than synchronization granularity
with respect to **worst-case blocking.**

# Cache-Line Bouncing

**Cache-line ownership jumps from core to core**

**Scheduler state shared among all cores**



**GSN-EDF**

**SCHED_DEADLINE**

# Cache-Line Bouncing

**Observation #2:**
State sharing results in overheads due to **cache-line bouncing, even if it's distributed across cores**.

# Root Causes of Overhead

**Peak Contention**

**Cache-Line Bouncing**

# This Talk

1) Why global scheduling?

2) Current implementations

3) Root causes of overhead

**4) How to scale global scheduling?**

5) Evaluation

# Candidate Solutions

Lock-free algorithms

# Candidate Solutions

Lock-free algorithms

multiple CAS in the same location, unpredictable fail-retry operations

# Candidate Solutions

Lock-free algorithms

multiple CAS in the same location, unpredictable fail-retry operations

Wait-free queue of events

# Candidate Solutions

Lock-free algorithms

multiple CAS in the same location, unpredictable fail-retry operations

Wait-free queue of events

complex garbage collection and serialization, didn't reduce cache-line bouncing

# Candidate Solutions

Lock-free algorithms

multiple CAS in the same location, unpredictable fail-retry operations

Wait-free queue of events

complex garbage collection and serialization, didn't reduce cache-line bouncing

All-to-all broadcast of events

# Candidate Solutions

**Lock-free algorithms**

multiple CAS in the same location, unpredictable fail-retry operations
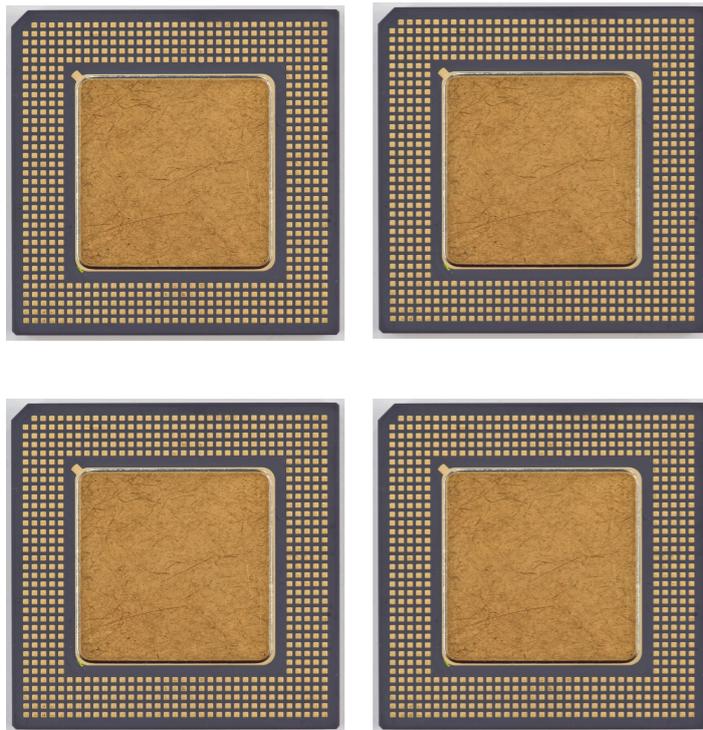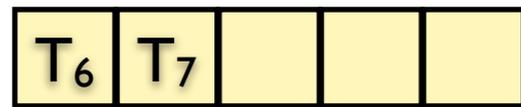
**Wait-free queue of events**

complex garbage collection and serialization, didn't reduce cache-line bouncing
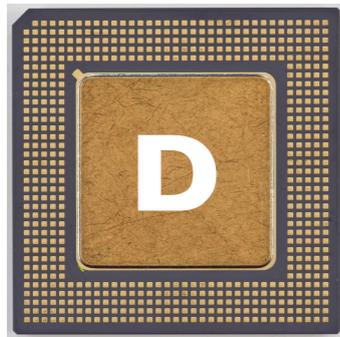
**All-to-all broadcast of events**

message ordering, consensus
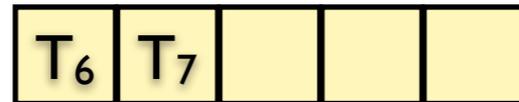
# Reducing Cache-Line Bouncing

Scheduler State

| T₆ | T₇ | | | |

# Reducing Cache-Line Bouncing

**D** Scheduler State

| $T_6$ | $T_7$ | | | |
|---|---|---|---|---|

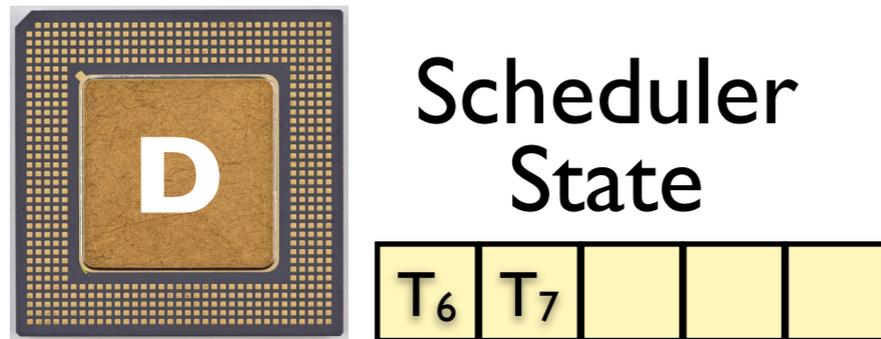## Dedicated Scheduler Processor
- Stores the full scheduler state
- Dedicated interrupt handling

# Reducing Cache-Line Bouncing

**D** Scheduler State | T₆ | T₇ | | | |

**Dedicated Scheduler Processor**
- Stores the full scheduler state
- Dedicated interrupt handling

**Client Processors**

- Only know which task they should schedule (local state)

**C** T₃   **C** T₂   **C** T₅

Local states

# Reducing Cache-Line Bouncing
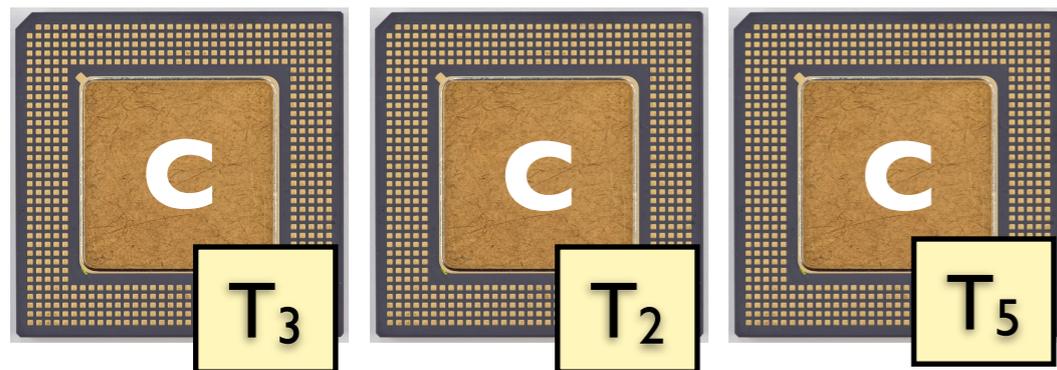


**Scheduler State**

T₆ | T₇ | | |

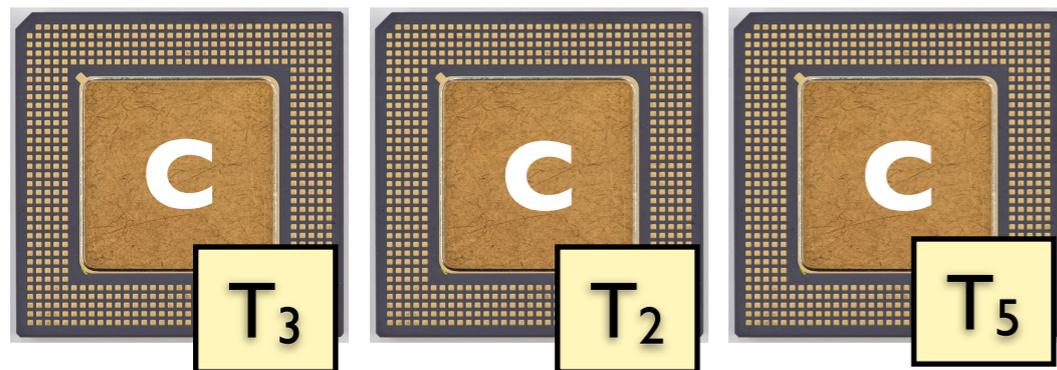**Dedicated Scheduler Processor**
- Stores the full scheduler state
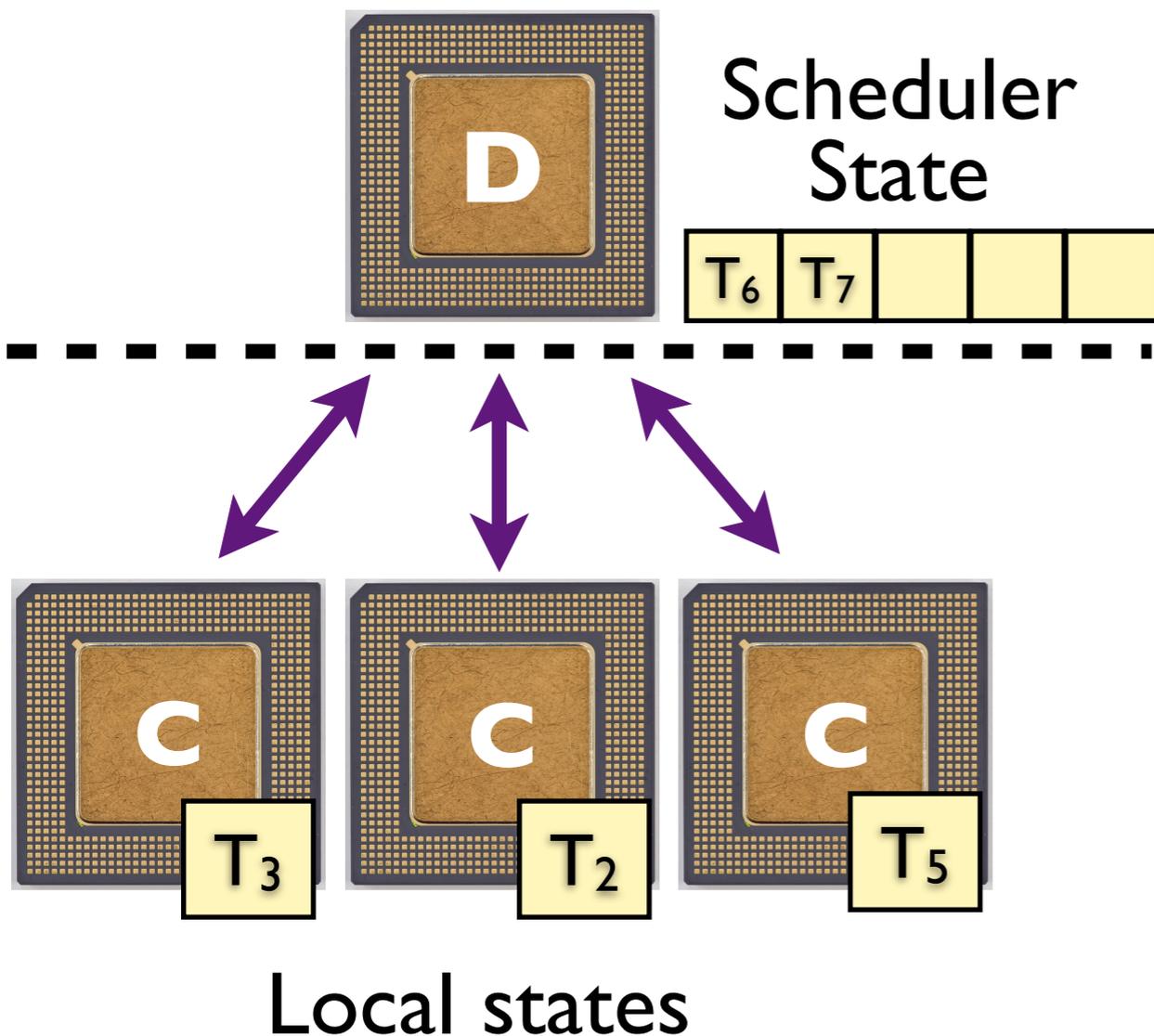- Dedicated interrupt handling

**Client Processors**
- Only know which task they should schedule (local state)

Local states

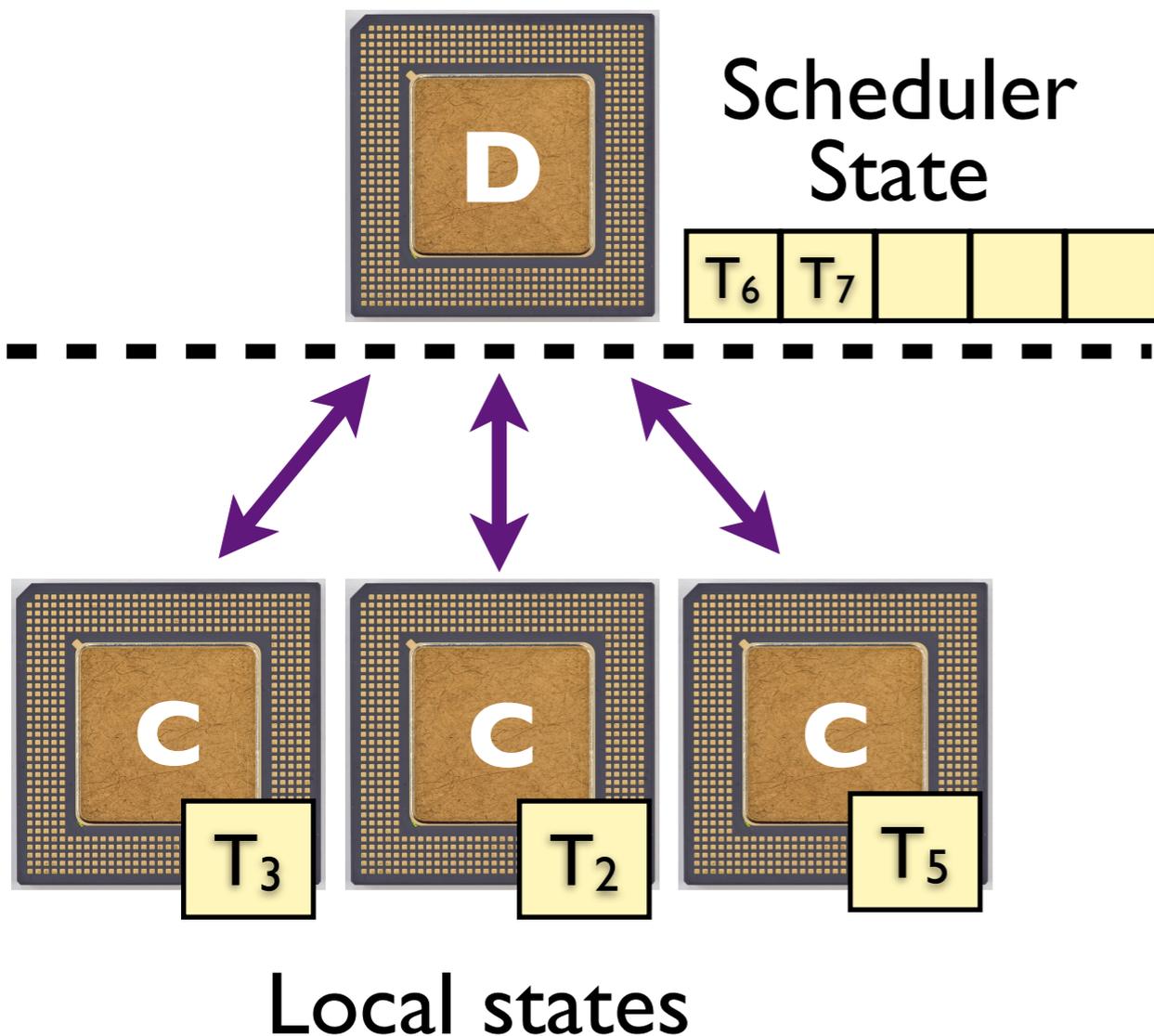**Centralized state reduces sharing**

# Communication with low
# Peak Contention



Scheduler State

**Centralized coordination**
- No interaction among clients
- Low-cost communication via message passing

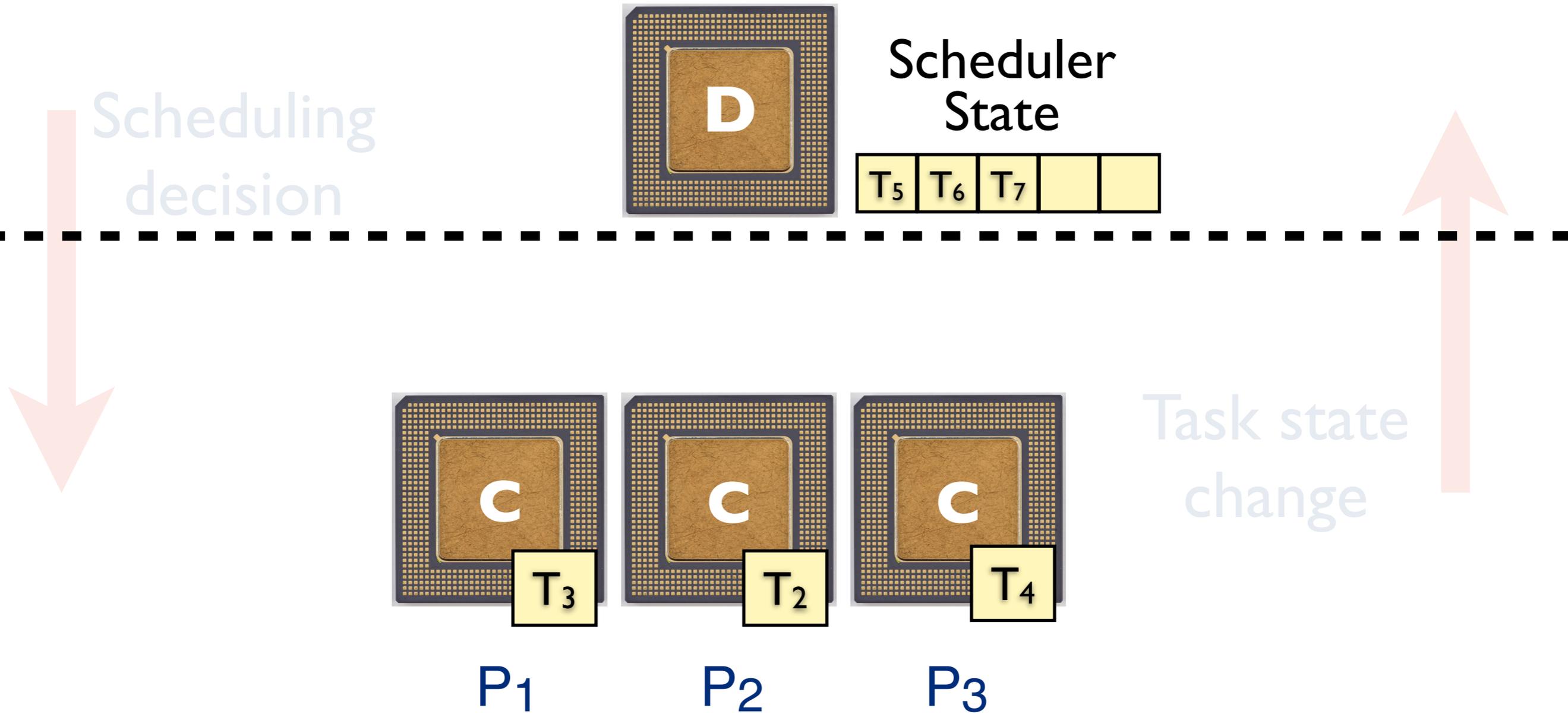Local states

# Communication with low
# Peak Contention

**D** (processor)

Scheduler
State

| $T_6$ | $T_7$ | | | |

**Centralized coordination**
- No interaction among clients
- Low-cost communication via message passing

**C** — $T_3$
**C** — $T_2$
**C** — $T_5$

Local states

**Contention limited to at most two processors**

# Message Passing

# Message Passing

# Message Passing

# Message Passing

# Implementing Messages Efficiently

- Message passing via per-cpu-socket mailboxes

- Shared-memory buffer with wait-free writes

Source code at
www.litmus-rt.org

# G-EDF-MP



Scheduler State

| T₆ | T₇ |  |  |  |

D

C — T₃

C — T₂

C — T₅

Global-EDF via **M**essage **P**assing

**Centralized Scheduling with Message passing**

# This Talk

1) Why global scheduling?

2) Current implementations

3) Root causes of overhead

4) How to scale global scheduling?

5) **Evaluation**

# Low Scheduling Overheads



**Maximum**

Legend: GSN-EDF, G-EDF-MP, SD

y-axis: overheads (in ms)
x-axis: number of processors

# Low Scheduling Overheads

## Maximum

# Low Scheduling Overheads

# Low Scheduling Overheads



Average

# Low Scheduling Overheads



**Average**

Legend:
- GSN-EDF
- G-EDF-MP
- SD

**G-EDF-MP incurs low average overheads!**

y-axis: overheads (in ms), 0 to 0.5

x-axis: number of processors, 8 16 24 32 48 64

# Low Scheduling Overheads

**Average**



But G-EDF-MP incurs additional message-passing overheads!

# Two Sources of Overhead

# Two Sources of Overhead

# Message-Passing



**Message passing overheads are significant!**

- Client Latency (Max)
- Callback Overhead (Max)
- Client Latency (Avg)
- Callback Overhead (Avg)

overheads (in ms)

number of processors

# Message-Passing



Message passing overheads are significant!

**What's the overall impact on schedulability?**

# Overhead-Aware Analysis

**Hard-real-time** $\longrightarrow$ **Max. overheads Schedulability test**

**Soft-real-time** $\longrightarrow$ **Avg. overheads Bounded Tardiness**

# Overhead-Aware Schedulability

schedulability ratio

1

0.8

0.6

0.4

0.2

0

Schedulability Results
for 64 CPUs
**Higher is better**

4   8   12   16   20   24   28   32   36   40   44

task set utilization

# Hard-Real-Time Schedulability

schedulability ratio

task set utilization

No overheads
G-EDF-MP
GSN-EDF
SD

**SCHED_DEADLINE does not implement dedicated interrupt handling, yielding a pessimistic analysis**

**Higher schedulability, even with additional message-passing delays**

schedulability ratio

No overheads
G-EDF-MP
GSN-EDF
SD

task set utilization

# Soft-Real-Time Schedulability



schedulability ratio vs. task set utilization

Legend:
- No overheads (red, solid)
- G-EDF-MP (black, solid)
- GSN-EDF (green, dash-dot)
- SD (blue, dashed)

# Schedulability

**SCHED_DEADLINE works well in the average case, but cannot be shown to do so analytically**

Legend:
- No overheads — (red solid line)
- G-EDF-MP — (black solid line)
- GSN-EDF — (green dash-dot line)
- SD — (blue dashed line)

Y-axis: schedulability ratio (0, 0.2, 0.4, 0.6, 0.8, 1)

X-axis: task set utilization (8, 16, 24, 32, 40, 48, 56, 64)

**SCHED_DEADLINE works well in the average case, but cannot be shown to do so analytically**

vs

**G-EDF-MP also performs well in the average case**

Legend:
- No overheads
- G-EDF-MP
- GSN-EDF
- SD

Y-axis: schedulability ratio (0, 0.2, 0.4, 0.6, 0.8, 1)

X-axis: task set utilization (8, 16, 24, 32, 40, 48, 56, 64)

# Global-EDF with Low Overheads

**Pair-wise coordination + Message passing** → **Scalable G-EDF implementation** up to 64 CPUs

# Limitations

Dedicated scheduling processor is still a **scalability bottleneck** *at **extreme core counts**.*

➜ G-EDF-MP scales *much further* than prior approaches.

G-EDF-MP is *inappropriate* for workloads that do not tolerate **excessive migration overheads.**

➜ Migrations are inherent to global scheduling policies, *irrespective of implementation.*

This approach can be applied to global scheduling in general, not just G-EDF.

# Conclusion

# Conclusion

Fine-grained locking is not enough. Scalability of worst-case overheads requires avoiding **peak contention** and **cache-line bouncing.**

# Conclusion

Fine-grained locking is not enough. Scalability of worst-case overheads requires avoiding **peak contention** and **cache-line bouncing.**

To reduce overheads, we used a **centralized scheduler** and **message passing.**

# Conclusion

Fine-grained locking is not enough. Scalability of worst-case overheads requires avoiding **peak contention** and **cache-line bouncing.**

To reduce overheads, we used a **centralized scheduler** and **message passing.**

G-EDF-MP's design can be applied to other global schedulers and **extends the range of processor counts** that can be practically supported.

# Thanks!



## www.litmus-rt.org

New release 2014.1 is now available!