# <u>The FMLP+</u>
# An Asymptotically Optimal Real-Time Locking Protocol for Suspension-Aware Analysis

ECRTS'14
July 9, 2014

Björn B. Brandenburg
bbb@mpi-sws.org

Max
Planck
Institute
for
Software Systems

# Suspension-Based Locking Protocols

## Semaphores in POSIX

```
pthread_mutex_lock(…)
// critical section
pthread_mutex_unlock(…)
```

semaphore:
a waiting task suspends,
**makes processor
available to other tasks**

# Suspension-Based Locking Protocols

## Semaphores in POSIX

```
pthread_mutex_lock(…)
// critical section
pthread_mutex_unlock(…)
```

**semaphore**:
a waiting task suspends,
**makes processor
available to other tasks**

**Priority Inversion Blocking**
➡ Locks cause **priority inversions**
   ≈ extra delay due to lock contention
➡ Short: **pi-blocking**

# Suspension-Based Locking Protocols

**Semaphores
in POSIX**

```
pthread_mutex_lock(…)
// critical section
pthread_mutex_unlock(…)
```

**semaphore**:
a waiting task suspends,
**makes processor
available to other tasks**

**Priority Inversion Blocking**

➡ Locks cause **priority inversions**
   ≈ extra delay due to lock contention
➡ Short: **pi-blocking**

**Blocking Analysis**

➡ For a **specific task set**, what is the
   **maximum duration of pi-blocking**
   incurred by each task?

# Suspension-Based Locking Protocols

## Semaphores in POSIX

```
pthread_mutex_lock(…)
// critical section
pthread_mutex_unlock(…)
```

**semaphore**:
a waiting task suspends,
**makes processor
available to other tasks**

**Priority Inversion Blocking**
➡ Locks cause **priority inversions**
  ≈ extra delay due to lock contention
➡ Short: **pi-blocking**

**Blocking Analysis**
➡ For a **specific task set**, what is the
  **maximum duration of pi-blocking**
  incurred by each task?

**Blocking Optimality**
➡ In general, what is the **maximum
  duration of pi-blocking** incurred by
  **any task in any task set**?

# Multiprocessor Real-Time Locking Optimality Classes

| Blocking Optimality<br>[— & Anderson, 2010] | suspension **oblivious** | suspension **aware** |
|---|---|---|
| How are suspensions analyzed? | CPU demand is **over-approximated** | CPU demand is modeled **accurately** |
| Advantage | simpler analysis | potentially less pessimistic |

[— & Anderson, 2010] *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.

# Asymptotically Optimal Locking Protocols

| JLFP<br>*job-level fixed-priority* | Suspension Oblivious<br>Any JLFP Scheduler | Suspension Aware<br>EDF w/ Implicit Deadlines | Suspension Aware<br>Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned**<br>(no migrations) | | | |
| **Global**<br>(jobs migrate freely) | | | |
| **Clustered**<br>(jobs migrate only among subset of processors) | | | |

[Block et al., 2007]     *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]     *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]     *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]     *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012] *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]     *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Asymptotically Optimal Locking Protocols

| JLFP _job-level fixed-priority_ | Suspension Oblivious Any JLFP Scheduler | Suspension Aware EDF w/ Implicit Deadlines | Suspension Aware Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** ✔ [— & Anderson, 2010] | | |
| **Global** (jobs migrate freely) | **G-OMLP** ✔ [— & Anderson, 2010] | | |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** ✔ [—, 2013] **C-OMLP** [— & Anderson, 2011] | | |

[Block et al., 2007]     *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]     *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]     *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]     *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012]     *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]     *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Asymptotically Optimal Locking Protocols

| JLFP *job-level fixed-priority* | Suspension Oblivious Any JLFP Scheduler | Suspension Aware EDF w/ Implicit Deadlines | Suspension Aware Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** ✔ [— & Anderson, 2010] | **SPFP** (asymptotical tightness) [— & Anderson, 2010] **P-FMLP⁺** (practical protocol) [—, 2011] | ✔ |
| **Global** (jobs migrate freely) | **G-OMLP** ✔ [— & Anderson, 2010] | | |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** ✔ [—, 2013] **C-OMLP** [— & Anderson, 2011] | | |

[Block et al., 2007]        *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]        *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]        *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]        *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012]        *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]        *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Asymptotically Optimal Locking Protocols

| JLFP<br>*job-level fixed-priority* | Suspension Oblivious<br>Any JLFP Scheduler | Suspension Aware<br>EDF w/ Implicit Deadlines | Suspension Aware<br>Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned**<br>(no migrations) | **P-OMLP** ✔<br>[— & Anderson, 2010] | **SPFP** (asymptotical tightness)<br>[— & Anderson, 2010]<br>**P-FMLP+** (practical protocol)<br>[—, 2011] ✔ | |
| **Global**<br>(jobs migrate freely) | **G-OMLP** ✔<br>[— & Anderson, 2010] | **FMLP** ⚠<br>[Block et al., 2007] | ❓ |
| **Clustered**<br>(jobs migrate only among subset of processors) | **OMIP** ✔<br>[—, 2013]<br>**C-OMLP**<br>[— & Anderson, 2011] | ❓ | ❓ |

[Block et al., 2007]   *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]   *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]   *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]   *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012]   *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]   *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

Support for **nested critical sections** added by **RNLP**.
[Ward & Anderson, 2012]

| JLFP *job-level fixed-priority* | Suspension Oblivious Any JLFP Scheduler | Suspension Aware EDF w/ Implicit Deadlines | Suspension Aware Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** ✔ [— & Anderson, 2010] | **SPFP** (asymptotical tightness) [— & Anderson, 2010] **P-FMLP+** (practical protocol) [—, 2011] | ✔ |
| **Global** (jobs migrate freely) | **G-OMLP** ✔ [— & Anderson, 2010] | **FMLP** ⚠ [Block et al., 2007] | ❓ |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** ✔ [—, 2013] **C-OMLP** [— & Anderson, 2011] | ❓ | ❓ |

[Block et al., 2007]       *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]     *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]     *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]                      *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012] *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]                      *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Asymptotically Optimal Locking Protocols

| JLFP<br>*job-level fixed-priority* | Suspension Oblivious<br>Any JLFP Scheduler | Suspension Aware<br>Any JLFP Scheduler |
|---|---|---|
| **Partitioned**<br>(no migrations) | **P-OMLP** ✔<br>[— & Anderson, 2010]<br>+ RNLP [Ward & Anderson, 2012] | **This Work**<br>The Generalized FMLP+<br>(FIFO Multiprocessor Locking Protocol)<br><br>✔<br><br>+ RNLP [Ward & Anderson, 2012] for nested critical sections |
| **Global**<br>(jobs migrate freely) | **G-OMLP** ✔<br>[— & Anderson, 2010]<br>+ RNLP [Ward & Anderson, 2012] | |
| **Clustered**<br>(jobs migrate only among subset of processors) | **OMIP**<br>[—, 2013] ✔<br>**C-OMLP**<br>[— & Anderson, 2011]<br>+ RNLP [Ward & Anderson, 2012] | |

[Block et al., 2007]    *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]    *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]    *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]    *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012]  *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]    *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Remainder of This Talk

**What is Suspension-Aware PI-Blocking?**
➡ Assumptions & quick review

**Finding a Suitable Progress Mechanism**
➡ How to deal with lock-holder preemptions

**Closing the Suspension-Aware Optimality Gap**
➡ New progress mechanism:
   **restricted segment boosting**
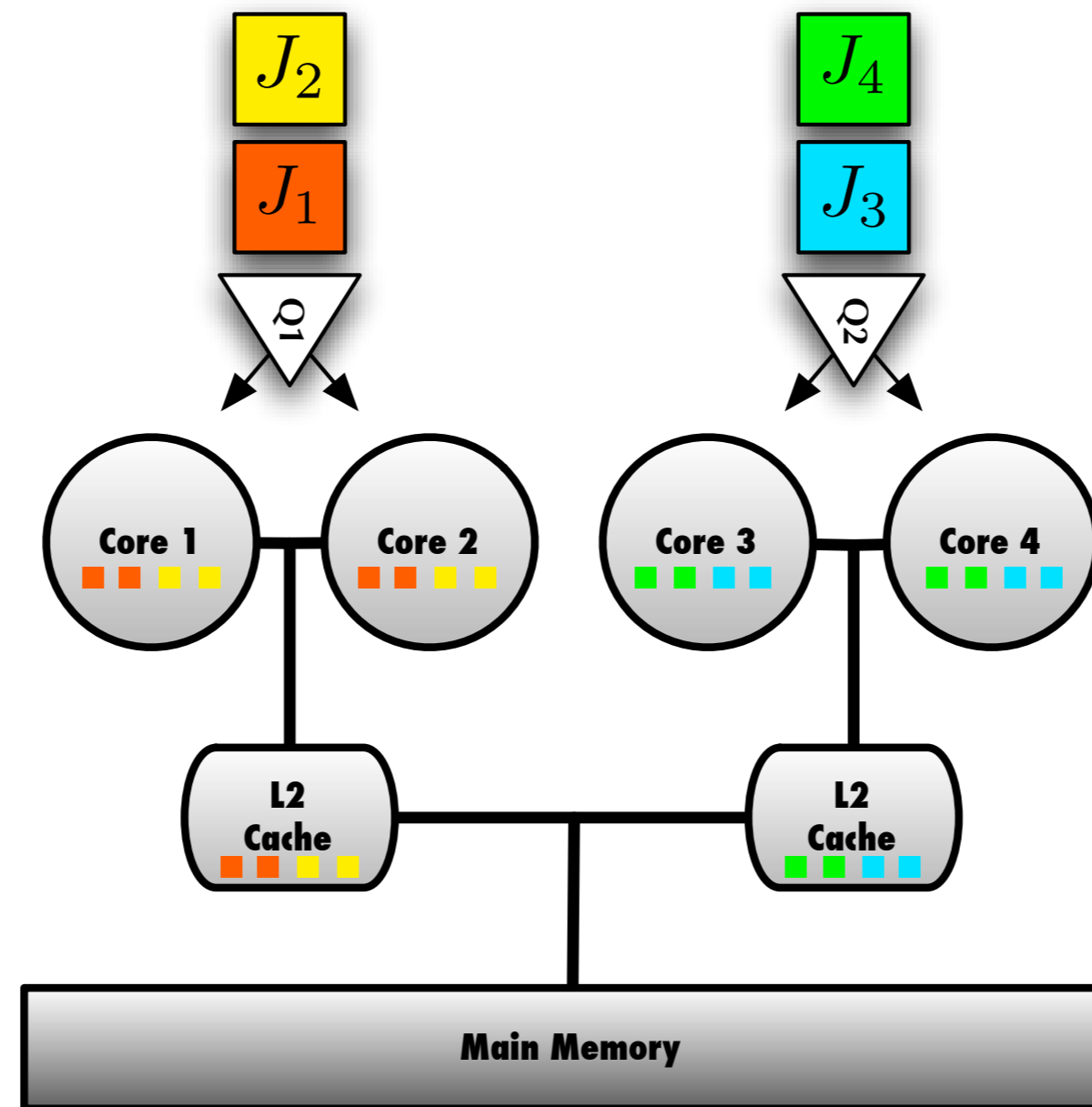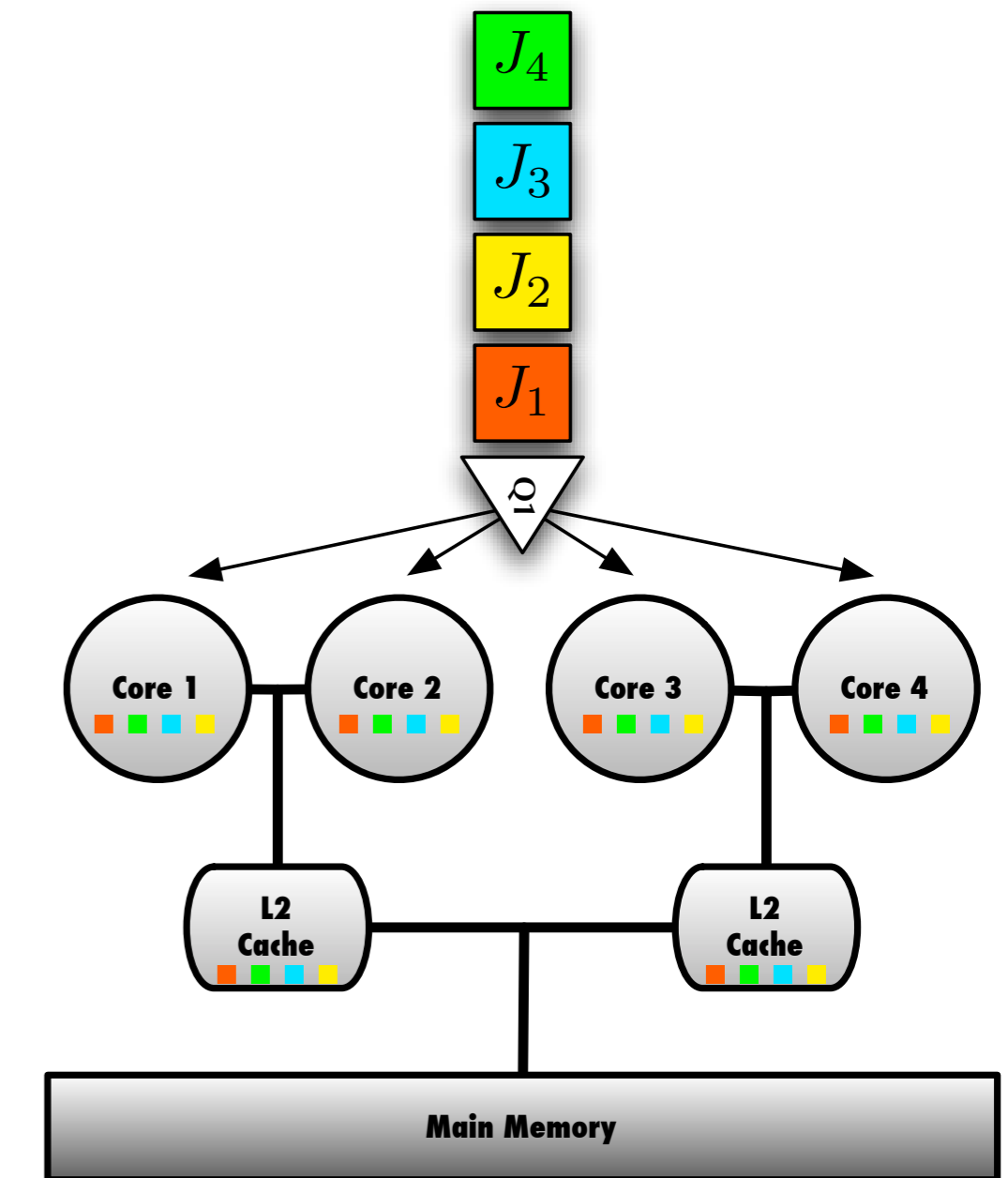➡ Achieving asymptotic optimality with the
   **generalized FMLP+**

# Assumptions
# & Review
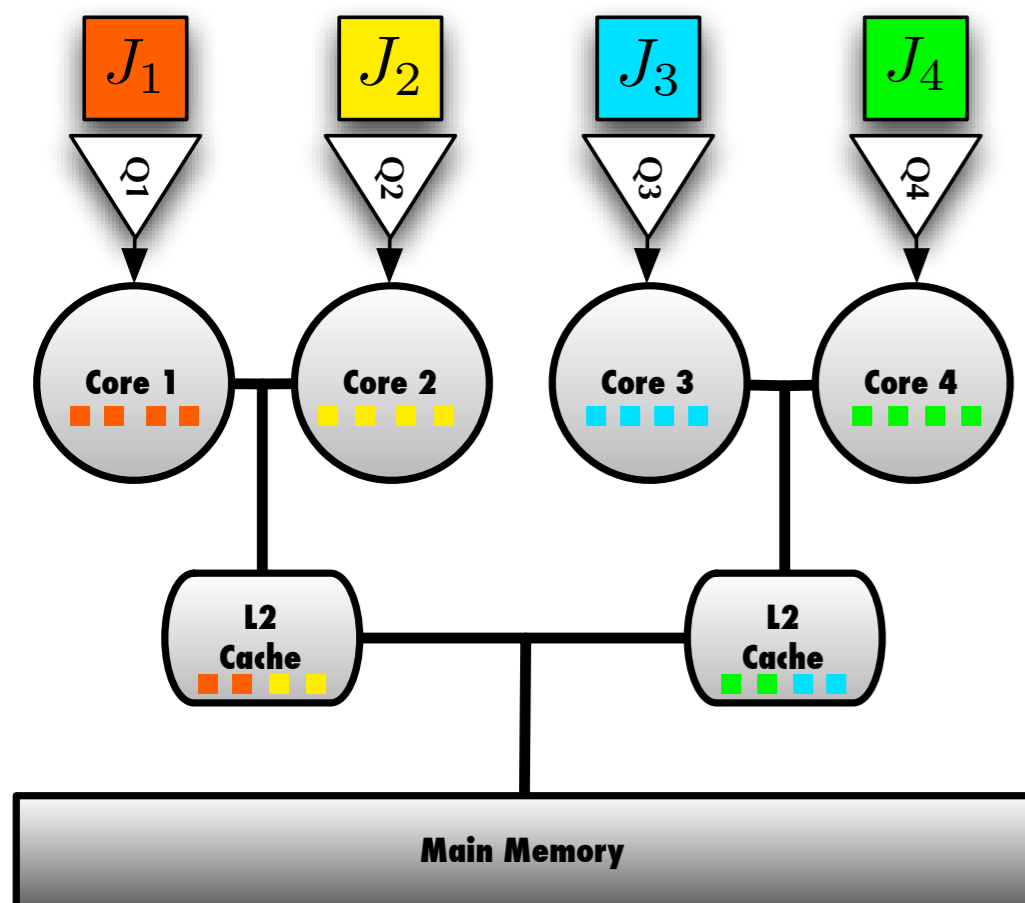# of Suspension-Aware
# PI-Blocking

# System Model



partitioned scheduling

**clustered scheduling**

global scheduling

**Clustered Scheduling**

➡ disjoint clusters of processors
  ‣ special cases: partitioned & global
➡ *job-level fixed-priority policy* (JLFP)
  ‣ e.g., EDF, <u>static task priorities</u>
➡ **cluster size may be non-uniform**

# System Model



partitioned scheduling

**clustered scheduling**

global scheduling

## Clustered Scheduling
➡ disjoint clusters of processors
 ‣ special cases: partitioned & global
➡ *job-level fixed-priority policy* (JLFP)
 ‣ e.g., EDF, <u>static task priorities</u>
➡ **cluster size may be non-uniform**

## Sporadic Tasks
➡ **arbitrary** deadlines
➡ shared resources
 ‣ in the paper: also nested CSs
 ‣ in the talk: <u>only unnested CSs</u>
➡ *locking-unrelated self-suspensions*

# Definition: S-Aware PI-Blocking

*A job **J** assigned to a cluster with **c** CPUs
incurs **s-aware pi-blocking** at a time **t** iff*

**(1)** **J** *is **not scheduled** at time **t**, and*

**(2)** *fewer than **c** **higher-priority jobs** are scheduled.*

[— & Anderson, 2010]

**Intuition**
➡ Locking-related delays are **not problematic** iff
**J** would **not** have been scheduled anyway…

[— & Anderson, 2010] *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.

# Maximum PI-Blocking

$b_i$ — bound on max. pi-blocking incurred by task $T_i$

max $\{b_i\}$ — **maximum pi-blocking** of any task in task set

# Maximum PI-Blocking

$b_i$ — bound on max. pi-blocking incurred by task $T_i$

max $\{b_i\}$ — **maximum pi-blocking** of any task in task set

---

*There exist task sets such that under s-aware analysis*

$$\max \{b_i\} = \Omega(n)$$

*under **any** suspension-based locking protocol.*

[— & Anderson, 2010]

---

(assuming constant critical section lengths)

[— & Anderson, 2010] *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.

# Maximum PI-Blocking

$b_i$ — bound on max. pi-blocking incurred by task $T_i$

max $\{b_i\}$ — **maximum pi-blocking** of any task in task set

*There exist task sets such that under s-aware analysis*

$$\max \{b_i\} = \Omega(n)$$

*under **any** suspension-based locking protocol.*

[— & Anderson, 2010]

(assuming constant critical section lengths)

→ O($n$) max. s-aware pi-blocking is asymptotically optimal.

[— & Anderson, 2010] *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.

# Objective

Define a locking protocol such that

$$max \{ b_i \} = O(n)$$

for **any task set** under **any clustered JLFP scheduler**.

**Need to define queue order**
➡ FIFO works

**Need to define progress mechanism**
➡ To deal with risk of lock-holder preemption
➡ Ensure timely completion of critical sections
➡ Classic example: priority inheritance

# Finding a Suitable Progress Mechanism

# Progress Mechanism Choices

| JLFP<br>*job-level fixed-priority* | Suspension Oblivious<br>Any JLFP Scheduler | Suspension Aware<br>EDF w/ Implicit Deadlines | Suspension Aware<br>Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned**<br>(no migrations) | **P-OMLP**<br>[— & Anderson, 2010] | **SPFP** (asymptotical tightness)<br>[— & Anderson, 2010]<br>**P-FMLP+** (practical protocol)<br>[—, 2011] | |
| **Global**<br>(jobs migrate freely) | **G-OMLP**<br>[— & Anderson, 2010] | **FMLP**<br>[Block et al., 2007] | **?** |
| **Clustered**<br>(jobs migrate only among subset of processors) | **OMIP**<br>[—, 2013]<br>**C-OMLP**<br>[— & Anderson, 2011] | **?** | **?** |

[Block et al., 2007]     *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]     *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]     *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]     *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[—, 2013]     *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Progress Mechanism Choices

(1) The classic choice: (variants of) **priority inheritance**.

| job-level fixed priority | Any JLFP Scheduler | Deadlines | Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** [— & Anderson, 2010] | **SPFP** (asymptotical tightness) [— & Anderson, 2010] **P-FMLP⁺** (practical protocol) [—, 2011] | |
| **Global** (jobs migrate freely) | **G-OMLP** [— & Anderson, 2010] | **FMLP** [Block et al., 2007] | ？ |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** [—, 2013] **C-OMLP** [— & Anderson, 2011] | ？ | ？ |

[Block et al., 2007]    *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]    *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]    *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]    *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[—, 2013]    *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Progress Mechanism Choices

(1) The classic choice: (variants of) **priority inheritance**.

| | job-level fixed priority Any JLFP Scheduler | Deadlines | Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** [— & Anderson, 2010] | **SPFP** (asymptotical tightness) [— & Anderson, 2010] **P-FMLP⁺** (practical protocol) [— 2011] | |
| **Global** (jobs migrate freely) | **G-OMLP** [— & Anderson, 2010] | **FMLP** [Block et al., 2007] | ? |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** [— 2013] **C-OMLP** [— & Anderson, 2011] | ? | ? |

[Block et al., 2007]   *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.

(2) The partitioned & clustered choice: (variants of) **priority boosting**.

# Sub-Optimality of Priority <u>Inheritance</u>

It is **impossible to construct**

an **asymptotically optimal** locking protocol
*(w.r.t. s-aware analysis)*

under global JLFP scheduling

based on **priority <u>inheritance</u>**.

[—, 2011]

(And hence also under clustered JLFP scheduling.)

[—, 2011]: *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, University of North Carolina at Chapel Hill, 2011.

# Sub-Optimality of Priority <u>Inheritance</u>

**<u>Priority Inheritance Example Schedule</u>**

4 tasks on 2 processors

**global fixed-priority scheduling**

Task priorities (high to low): $T_1 > T_2 > T_3 > T_4$

| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | ☐ | ◻ | ↑ | | ≡ | |
| | | | ↓ | deadline | | |
| Processor 2 | ▨ | ▨ | ⊤ | job completion | ⊢PI⊣ | priority inversion |

# Sub-Optimality of Priority <u>Inheritance</u>

# Sub-Optimality of Priority <u>Inheritance</u>

$T_4$ :  short period, <u>long deadline</u> ($d_4 > p_4$), shares lock with $T_2$



| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | | | ↑ | | | |
| | | | ↓ | deadline | | |
| Processor 2 | | | ⊤ | job completion | PI | priority inversion |

# Sub-Optimality of Priority <u>Inheritance</u>



$T_4$ :  short period, <u>long deadline</u> ($d_4 > p_4$), shares lock with $T_2$

$T_2$ & $T_1$:  short period, <u>implicit deadline</u> ($d_{\{1,2\}} = p_{\{1,2\}}$)

| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | □ | ◻ | ↑ | | ≡ | |
| Processor 2 | ▨ | ▨ | ↓ | deadline | PI | priority inversion |
| | | | ⊤ | job completion | | |

# Sub-Optimality of Priority <u>Inheritance</u>

**$T_4$ :** short period, <u>long deadline</u> ($d_4 > p_4$), shares lock with **$T_2$**

**$T_3$ :** long period, <u>implicit deadline</u>, "**victim task**"

**$T_2$ & $T_1$:** short period, <u>implicit deadline</u> ($d_{\{1,2\}} = p_{\{1,2\}}$)

$T_4$

$T_3$

$T_2$

$T_1$

0          5          10          15          20      time

| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | ☐ | ▨ | ↑ | | ≡ | |
| Processor 2 | ▦ | ▨ | ↓ | deadline | | |
| | | | ⊤ | job completion | ⊢ PI ⊣ | priority inversion |

# Independent Job Incurs Priority Inversion



Legend:

Processor 1: scheduled (white box), critical section (white box with diagonal line)

Processor 2: scheduled (gray box), critical section (gray box with diagonal line)

↑ job release
↓ deadline
⊤ job completion

≡ job suspended
PI priority inversion

# Independent Job Incurs Priority Inversion

$T_4$ **acquires lock first...**



...but is **preempted** by arrival of higher-priority jobs.

scheduled     critical section          ↑  job release          ≡  job suspended

Processor 1   [ ]           [ ╱ ]       ↓  deadline

Processor 2   [ ]           [ ╱ ]       ⊤  job completion       ⊢PI⊣  priority inversion

# Independent Job Incurs Priority Inversion

# Independent Job Incurs Priority Inversion

**Priority inheritance: $T_4$ inherits from $T_2$**

**→ $T_3$ is preempted, incurs s-aware pi-blocking**



| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | ☐ | ◹ | ↑ | | ≡ | |
| Processor 2 | ▨ | ◺ | ↓ | deadline | PI | priority inversion |
| | | | ⊤ | job completion | | |

# Independent Job Incurs Priority Inversion

**How often** can this scenario **repeat**?



| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | | | ↑ | job release | ≡ | |
| | | | ↓ | deadline | | |
| Processor 2 | | | ⊤ | job completion | PI | priority inversion |

# "Victim Task" Accumulates PI-Blocking



$\Omega(\phi)$ *pi-blocking is possible, where* $\phi$ *= {max response time} / {min period}*

# "Victim Task" Accumulates PI-Blocking

$\Omega(\phi)$ *pi-blocking is possible, where* $\phi$ = *{max response time} / {min period}*



**Bounded only by**

**the number of jobs released by *T1*, *T₂*, and *T₄* while *T₃* is pending.**

*How many jobs?* → $\phi$ = *{max response time} / {min period}*

# Sub-Optimality of Priority <u>Inheritance</u>



$\Omega(\phi)$ *pi-blocking* is possible, where $\phi$ = *{max response time} / {min period}*

$\phi$ is *not* bounded by the **number of tasks *n*.**

→ **not asymptotically optimal.**

scheduled

Processor 1

Processor 2

# What about Priority Boosting?

(1) The classic choice: (variants of) **priority inheritance.**

| *job-level fixed priority* | Any JLFP Scheduler | Deadlines | Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** [— & Anderson, 2010] | **SPFP** (asymptotical tightness) [— & Anderson, 2010] **P-FMLP+** (practical protocol) [— 2011] | |
| **Global** (jobs migrate freely) | **G-OMLP** [— & Anderson, 2010] | **FMLP** [Block et al., 2007] | ? |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** [— 2013] **C-OMLP** [— & Anderson, 2011] | ? | ? |

[Block et al., 2007] *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.

(2) The partitioned & clustered choice: (variants of) **priority boosting.**

# What about Priority Boosting?



(1) The classic choice: (variants of) priority inheritance.

| | Any JLFP-Scheduler | Deadlines | Any JLFP-Scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** [— & Anderson, 2010] | **SPFP** (asymptotical tightness) [— & Anderson, 2010] **P-FMLP⁺** (practical protocol) [— 2011] | |
| **Global** (jobs migrate freely) | **G-OMLP** [— & Anderson, 2010] | **FMLP** [Block et al., 2007] | ? |
| **Clustered** (jobs migrate only among subset of processors) | **C-OMLP** [— & Anderson, 2011] | ? | ? |

[Block et al., 2007] *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.

(2) The partitioned & clustered choice: (variants of) **priority boosting.**

# Sub-Optimality of <u>Unrestricted</u> Priority Boosting

*lock-holding jobs **always** have higher
effective priority than non-lock-holding jobs*

It is ***impossible to construct***

an ***asymptotically optimal*** locking protocol
*(w.r.t. s-aware analysis)*

under global JLFP scheduling

based on **<u>unrestricted</u> priority boosting**.

[−, 2011]

(And hence also under clustered JLFP scheduling.)

[−, 2011]: *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, University of North Carolina at Chapel Hill, 2011.

# Sub-Optimality of <u>Unrestricted</u> Priority Boosting

## **Unrestricted Priority Boosting Example Schedule**

**[same task set & arrival sequence as before]**

4 tasks on 2 processors

**global fixed-priority scheduling**

Task priorities (high to low): $T_1 > T_2 > T_3 > T_4$

| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | ☐ | ⧄ | ↑ | job release | ☰ | job suspended |
| Processor 2 | ▨ | ◪ | ↓ | deadline | | priority inversion |
| | | | ⊤ | job completion | ⊢PI⊣ | priority inversion |

# Lock-Holding Jobs Cannot Be Preempted



**scheduled**   **critical section**   ↑ job release   ≡ job suspended

Processor 1   □   ▨   ↓ deadline

Processor 2   ▨(gray)   ▨(gray)   ⊤ job completion   ⊢PI⊣ priority inversion

# Lock-Holding Jobs Cannot Be Preempted



**Priority Boosting**

$T_4$ cannot be preempted
by newly released, **non-resource-holding job...**

| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | □ | ◸ | ↓ | deadline | | priority inversion |
| Processor 2 | ▨ | ◩ | ⊤ | job completion | PI | |

# PI-Blocking Shifted, not Avoided

**Priority Boosting**

$T_4$ **cannot be preempted**
by newly released, **non-resource-holding job...**



Processor 1 — scheduled (white box), critical section (white box with diagonal line)
Processor 2 — scheduled (grey box), critical section (grey box with diagonal line)

↑ job release
↓ deadline
⊤ job completion
≡ job suspended
PI — priority inversion

# PI-Blocking Shifted, not Avoided



**Priority Boosting**

$T_4$ **cannot be preempted**
by newly released, **non-resource-holding job…**

… but $T_3$ is still **preempted** by $T_1$.

→ **s-aware pi-blocking**

Legend:

scheduled · critical section · job release · job suspended

Processor 1 / Processor 2 · deadline · job completion · priority inversion

# "Victim Task" Accumulates PI-Blocking

$\Omega(\phi)$ *pi-blocking is possible, where* $\phi$ = *{max response time} / {min period}*



**Legend:**

|  | scheduled | critical section | | job release / deadline / job completion | | job suspended / priority inversion |
|---|---|---|---|---|---|---|
| Processor 1 | □ (white) | ▧ (white) | ↑ | job release | ☰ | job suspended |
| Processor 2 | ▪ (grey) | ▨ (grey) | ↓ | deadline | PI | priority inversion |
|  |  |  | ⊤ | job completion |  |  |

# Sub-Optimality of <u>Unrestricted</u> Priority Boosting

$\Omega(\phi)$ *pi-blocking is possible, where* $\phi$ = *{max response time} / {min period}*



Also repeats  $\phi$ = *{max response time} / {min period}*  times…

→ **not asymptotically optimal.**

**Remark:** examples use **single resource** shared by **only two tasks**

→ **queue order irrelevant** (FIFO-ordered, priority-ordered, etc.)
→ **cannot simplify problem** with coarser-grained locking

# We need something new...

**(1) The classic choice: (variants of) priority inheritance.**

| job-level-fixed-priority Any JLFP-scheduler | | Deadlines | Any JLFP-scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** [— & Anderson, 2010] | **SPFP** (asymptotical tightness) [— & Anderson, 2010] **P-FMLP+** (practical protocol) [— 2011] | |
| **Global** (jobs migrate freely) | **G-OMLP** [— & Anderson, 2010] | **FMLP** [Block et al., 2007] | ? |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** [— 2013] **C-OMLP** [— & Anderson, 2011] | ? | ? |

[Block et al., 2007]    *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[
[
[
[

**(2) The partitioned & clustered choice: (variants of) priority boosting.**

# We need something new...



| | Any JLFP Scheduler | Deadlines | Any JLFP Scheduler |
|---|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** [— & Anderson, 2010] | **SPFP** (asymptotical tightness) [— & Anderson, 2010] **P-FMLP⁺** (practical protocol) [— 2011] | |
| **Global** (jobs migrate freely) | **G-OMLP** [— & Anderson, 2010] | **FMLP** [Block et al., 2007] | ? |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** [— 2013] **C-OMLP** [— & Anderson, 2011] | ? | ? |

(1) The classic choice: (variants of) priority inheritance.

(2) The partitioned & clustered choice: (variants of) priority boosting.

[Block et al., 2007] *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.

# Observation: *O(n)* PI-Blocking Possible



| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | □ | ▨ | ↑ | | | |
| | | | ↓ | deadline | | priority inversion |
| Processor 2 | ▨ | ▨ | ⊤ | job completion | PI | |

# Observation: O(*n*) PI-Blocking Possible

*This schedule:* **unrestricted** *priority boosting.*



| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | ☐ | ◿ | ↑ | | ≡ | |
| | | | ↓ | deadline | | |
| Processor 2 | ▨ | ◿ | ⊤ | job completion | PI | priority inversion |

# Observation: O(*n*) PI-Blocking Possible

This schedule: **unrestricted** priority boosting.



PI-blocking **cannot be avoided**, but it can be **shifted** to **new jobs**.

→ **prevent accumulation** of pi-blocking in individual jobs.

# Idea: <u>Protect</u> Existing Independent Jobs



*Alternative possible schedule.*

Legend:

| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | ☐ | ◹ | ↑ | deadline | ≡ | |
| Processor 2 | ▩ | ◸ | ↓ | | PI | priority inversion |
| | | | ⊤ | job completion | | |

# Idea: <u>Protect</u> Existing Independent Jobs

*Alternative possible schedule.*



**$T_3$ is not preempted while $T_4$ holds a lock!**

**Both $T_1$ and $T_2$ incur pi-blocking** instead.

# $O(n)$ PI-Blocking Per Job



**Alternative possible schedule.**

**$T_3$ "protected": incurs no pi-blocking in this example.**

*Alternative possible schedule.*



**Each job incurs only a limited amount of pi-blocking.**

→ **asymptotically optimal.**

This is actually a **Generalized FMLP⁺** schedule.

*Key question: how to specify that "$T_3$ must be protected"?*

# Closing the S-Aware Asymptotic Optimality Gap

# Key Problem: Preemptions due to Later-Started Critical Sections

**On Uniprocessors**

➡ a job is **blocked only** by critical sections that **are already in progress** when the job is released / resumed.

# Key Problem: Preemptions due to Later-Started Critical Sections

**On Uniprocessors**

➡ a job is **blocked only** by critical sections that **are already in progress** when the job is released / resumed.

**On Multiprocessors**

➡ Priority boosting example…

## Priority Boosting / Priority Inheritance Examples

*T3* blocked due to **φ** requests issued **after** *T3* started executing.

(**→ root cause: parallel scheduling of lower-priority jobs**)

## On Uniprocessors

➡ a job is **blocked only** by critical sections that **are already in progress** when the job is released / resumed.

## On Multiprocessors

➡ Priority boosting example…

# What if this is disallowed…?

*strawman rule: jobs cannot be preempted due to later-issued requests*



cannot preempt

may not preempt

may preempt (one per task)

*critical sections*

$CS_b$

$CS_x$

$CS_a$

*analyzed job*

$J_i$

*released*

*time*

➡ ***O(n)*** preemptions

**This is the desired effect, but the simple rule fails in corner cases.**

*strawman rule: jobs cannot be preempted due to later-issued requests*

cannot preempt

may <u>not</u> preempt

may preempt (one per task)

*critical sections*

$CS_b$

$CS_x$

$CS_a$

*analyzed job*

$J_i$

*released*

*time*

➡ **O(n)** preemptions

# Independent & Request Segments

a job at runtime:

| independent segment | request segment | independent segment | … |

*holding a lock or suspended*     *ready & does not require a lock*

**Note: exact segments known only at runtime**
➡ potentially complex, non-linear **control flow** determines which resources are required and in which order
➡ approach *not* limited to linear, branch-free tasks

## **Key Concept:** Segment Start Time

*Simply the start time of a job's current segment.*

a job at runtime:

independent segment | request segment | independent segment | ...

*holding a lock or suspended*    *ready & does not require a lock*

**Note: exact segments known only at runtime**
➡ potentially complex, non-linear **control flow** determines which resources are required and in which order
➡ approach ***not*** limited to linear, branch-free tasks

# A Lock Holder's Co-Boosting Set

Key idea underlying the Generalized FMLP+

*If a job is **priority-boosted**,*
*then **certain other jobs** must also be **co-boosted**.*

# A Lock Holder's Co-Boosting Set

*If a job $J_b$ holds a lock at **time t**, then its **co-boosting set** is defined as:*

$$\left\{ J_y \middle| \begin{array}{l} J_y \text{ executes an } \underline{\textbf{\textit{independent segment}}} \text{ at } \textbf{\textit{time t}} \text{ and} \\[4pt] J_y \text{ has } \underline{\textbf{\textit{higher priority}}} \text{ than } J_b \text{ and} \\[4pt] J_y\text{'s current segment } \underline{\textbf{\textit{started before}}} J_b\text{'s segment.} \end{array} \right\}$$

(Note: in this talk, I'll use "task" and "job" interchangeably.)

## Intuition

➡ The set of jobs **at risk** of **accumulating pi-blocking** due to $J_b$.

# Example: { *T₃* } is *T₄*'s Co-Boosting Set



*Generalized FMLP⁺ schedule.*

**Legend:**

scheduled | critical section | ↑ job release | ≡ job suspended
Processor 1 (white) | (white diagonal) | ↓ deadline
Processor 2 (grey) | (grey diagonal) | ⊤ job completion | ⊢PI⊣ priority inversion

# Example: { $T_3$ } is $T_4$'s Co-Boosting Set



Generalized FMLP⁺ schedule.

$T_3$ executes an **_independent segment_** at **_time t_** and

$T_3$ has **_higher priority_** than $T_4$ and

$T_3$'s current segment **_started before_** $T_4$'s segment.

$T_1$ and $T_2$ execute ***independent segments*** at ***time t*** and

$T_1$ and $T_2$ have ***higher priority*** than $T_4$ but

$T_1$ and $T_2$'s current segments ***did NOT start before*** $T_4$'s segment.

**scheduled** **critical section** | **job release**
Processor 1 — job suspended
Processor 2 | **deadline** | **priority inversion**
job completion

# Restricted Segment Boosting

*In a cluster with **c** CPUs, at any point in time **t**, schedule the following jobs:*

# Restricted Segment Boosting

*In a cluster with **c** CPUs, at any point in time **t**, schedule the following jobs:*

## A <span style="color:green">Single</span> Boosted Job $J_b$

The **lock-holding ready job** (if any) with the **earliest segment start time**.

*(any ties broken arbitrarily but consistently)*

# Restricted Segment Boosting

*In a cluster with **c** CPUs, at any point in time **t**, schedule the following jobs:*

### A **Single** Boosted Job $J_b$

The **lock-holding ready job** (if any) with the **earliest segment start time**.

### Up to $c - 1$ jobs from $J_b$'s Co-Boosting Set

Select the (up to) $c - 1$ jobs with the **earliest segment start times**.

*(any ties broken arbitrarily but consistently)*

# Restricted Segment Boosting

*In a cluster with **c** CPUs, at any point in time **t**, schedule the following jobs:*

## A **Single** Boosted Job $J_b$

The **lock-holding ready job** (if any) with the **earliest segment start time**.

## Up to **c - 1** jobs from $J_b$'s **Co-Boosting Set**

Select the (up to) **c - 1** jobs with the **earliest segment start times**.

## If less than **c** jobs scheduled so far: **any other ready jobs**

Select the **highest-priority ready jobs** not yet scheduled (may hold locks).

*(any ties broken arbitrarily but consistently)*

# Restricted Segment Boosting at <u>Time 1</u>



Generalized FMLP⁺ schedule.

**(1)** The **lock-holding ready job** (if any) with the **earliest segment start time**.

*Generalized FMLP+ schedule.*

Legend:

Processor 1 — scheduled / critical section

Processor 2 — scheduled / critical section

↑ job release
↓ deadline
⊤ job completion

≡ job suspended
PI priority inversion

**(1)** The **lock-holding ready job** (if any) with the **earliest segment start time**.

**(2)** Up to $c - 1 = 1$ jobs from $T_4$'s **co-boosting set** = {$T_3$}.

Legend:
- Processor 1 / Processor 2
- scheduled
- critical section
- job release
- deadline
- job completion
- job suspended
- PI priority inversion

**(1)** The **lock-holding ready job** (if any) with the **earliest segment start time**.

**(2) Up to** $c - 1 = 1$ **jobs from** $T_4$'s **co-boosting set =** $\{T_3\}$.

**(3)** If less than $c = 2$ jobs scheduled so far: any other ready jobs.

**At time 1**: **no more CPUs available** after steps 1 & 2.

sched

Processor 1

Processor 2

job completion     priority inversion

# The Generalized <u>FIFO Multiprocessor Locking Protocol</u> (FMLP+)

*Restricted Segment Boosting* + *Per-Resource FIFO Queues*

# The Generalized <u>FIFO Multiprocessor Locking Protocol</u> (FMLP+)

*Restricted Segment Boosting* + *Per-Resource FIFO Queues*

## S-Aware PI-Blocking per Segment

➡ …during <u>request segment</u>: *O($n$)*.

‣ <u>Proof</u>: rather straightforward ➡ see paper.

➡ …during <u>independent segment</u>: *O($n$)*.

‣ <u>Proof</u>: rather involved ➡ see paper.

# The Generalized <u>FIFO Multiprocessor Locking Protocol</u> (FMLP⁺)

*Restricted Segment Boosting* + *Per-Resource FIFO Queues*

**S-Aware PI-Blocking per Segment**

➡ …during <u>request segment</u>: **O($n$)**.
  ‣ <u>Proof</u>: rather straightforward ➞ see paper.

➡ …during <u>independent segment</u>: **O($n$)**.
  ‣ <u>Proof</u>: rather involved ➞ see paper.

**Number of segments**

➡ Constant **#requests** and **#self-suspensions** per job
  ➞ **constant number of segments**.

# The Generalized <u>FIFO Multiprocessor Locking Protocol</u> (FMLP⁺)

> *Restricted Segment Boosting*  +  *Per-Resource FIFO Queues*

## S-Aware PI-Blocking per Segment

➡ …during <u>request segment</u>: **O(*n*)**.
  ‣ <u>Proof</u>: rather straightforward ➞ see paper.

➡ …during <u>independent segment</u>: **O(*n*)**.
  ‣ <u>Proof</u>: rather involved ➞ see paper.

## Number of segments

➡ Constant **#requests** and **#self-suspensions** per job
  ➞ **constant number of segments**.

## Overall Max. S-Aware PI-Blocking

➡ **O(*n*)** under clustered JLFP scheduling.

# Multiprocessor Real-Time Locking Optimality Results

| JLFP<br>*job-level fixed-priority* | Suspension Oblivious<br>Any JLFP Scheduler | Suspension Aware<br>Any JLFP Scheduler |
|---|---|---|
| **Partitioned**<br>(no migrations) | **P-OMLP** ✅<br>[— & Anderson, 2010]<br>+ RNLP [Ward & Anderson, 2012] | **The Generalized FMLP+**<br>**(restricted segment boosting)**<br><br><br>✅<br><br>+ RNLP [Ward & Anderson, 2012] for nested critical sections |
| **Global**<br>(jobs migrate freely) | **G-OMLP** ✅<br>[— & Anderson, 2010]<br>+ RNLP [Ward & Anderson, 2012] | |
| **Clustered**<br>(jobs migrate only among subset of processors) | **OMIP** ✅<br>[—, 2013]<br>**C-OMLP**<br>[— & Anderson, 2011]<br>+ RNLP [Ward & Anderson, 2012] | |

[Block et al., 2007]     *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]     *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]     *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]                *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012]  *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]                *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

**The Generalized FMLP+ based on Restricted Segment Boosting closes the s-aware asymptotic optimality gap.**

*See paper & online appendix for **large-scale empirical evaluation**.*
*(Summary: **the FMLP+ works well** if the schedulability analysis is **accurate enough**.)*

| JLFP *job-level fixed-priority* | Suspension Oblivious Any JLFP Scheduler | Suspension Aware Any JLFP Scheduler |
|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** ✔ [— & Anderson, 2010] + RNLP [Ward & Anderson, 2012] | **The Generalized FMLP+ (restricted segment boosting)** ✔ |
| **Global** (jobs migrate freely) | **G-OMLP** ✔ [— & Anderson, 2010] + RNLP [Ward & Anderson, 2012] | |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** [—, 2013] ✔ **C-OMLP** [— & Anderson, 2011] + RNLP [Ward & Anderson, 2012] | + RNLP [Ward & Anderson, 2012] for nested critical sections |

[Block et al., 2007]         *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]      *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]      *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]                         *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012]  *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]                         *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Conclusion

# Summary

## The Generalized FMLP+

➡ priority boosting & inheritance **unsuitable**

➡ based instead on **restricted segment boosting**

‣ Key idea: **co-boosting of independent jobs**

## In the Paper

➡ Empirical evaluation.

➡ How to integrate with **locking-unrelated self-suspensions**…

‣ … also **within critical sections**.

➡ How to integrate with Ward & Anderson's RNLP [2012] for asymptotically optimal pi-blocking given **nested critical sections**.

## Online Appendix

➡ **Fine-grained blocking analysis** based on linear programming framework [—, 2013].

➡ Complete evaluation results (5760 graphs).

[Ward & Anderson, 2012]  *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]      *Improved Analysis and Evaluation of Real-Time Semaphore Protocols for P-FP Scheduling*, RTAS 2013.

# Future Work & Open Questions

**Apply this technique to reader-writer locks?**
➡ Lower bounds on s-aware pi-blocking?

**Apply this technique to k-exclusion locks?**
➡ GPUs & other co-processors
➡ Lower bounds on s-aware pi-blocking?

**Overheads of restricted segment boosting?**
➡ Tracking **segment start times** is simple and cheap.
➡ But… additional preemptions?

# Multiprocessor Real-Time Locking Optimality Results

| JLFP *job-level fixed-priority* | Suspension Oblivious Any JLFP Scheduler | Suspension Aware Any JLFP Scheduler |
|---|---|---|
| **Partitioned** (no migrations) | **P-OMLP** ✔ [— & Anderson, 2010] + RNLP [Ward & Anderson, 2012] | The Generalized FMLP+ (restricted segment boosting) ✔ |
| **Global** (jobs migrate freely) | **G-OMLP** ✔ [— & Anderson, 2010] + RNLP [Ward & Anderson, 2012] | |
| **Clustered** (jobs migrate only among subset of processors) | **OMIP** [—, 2013] ✔ **C-OMLP** [— & Anderson, 2011] + RNLP [Ward & Anderson, 2012] | + RNLP [Ward & Anderson, 2012] for nested critical sections |

[Block et al., 2007]      *A Flexible Real-Time Locking Protocol for Multiprocessors*, RTCSA 2007.
[— & Anderson, 2010]      *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.
[— & Anderson, 2011]      *Real-Time Resource-Sharing under Clustered Scheduling: Mutex, Reader-Writer, and k-Exclusion Locks*, EMSOFT 2011.
[—, 2011]      *Scheduling and Locking in Multiprocessor Real-Time Operating Systems*, PhD thesis, UNC, 2011.
[Ward & Anderson, 2012]      *Supporting Nested Locking in Multiprocessor Real-Time Systems*, ECRTS 2012.
[—, 2013]      *A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications*, ECRTS 2013.

# Appendix

# Multiprocessor Real-Time Locking Optimality Classes

| Blocking Optimality [— & Anderson, 2010] | suspension **oblivious** | suspension **aware** |
|---|---|---|
| How are suspensions analyzed? | CPU demand is **over-approximated** | CPU demand is modeled **accurately** |
| Lower bound on <u>maximum priority inversion blocking</u> $max_i\{b_i\}$ | $\Omega(m)$ <br> $m$ = #CPUs | $\Omega(n)$ <br> $n$ = #tasks |

[— & Anderson, 2010] *Optimality Results for Multiprocessor Real-Time Locking*, RTSS 2010.

# S-Aware vs. S-Oblivious Analysis

# S-Aware vs. S-Oblivious Analysis

## Suspension-oblivious (*s-oblivious*) Analysis



*task set* → **s-oblivious blocking analysis** → *modified task set with inflated WCETs* → **s-oblivious schedulability test** → *schedulable* / *not schedulable*
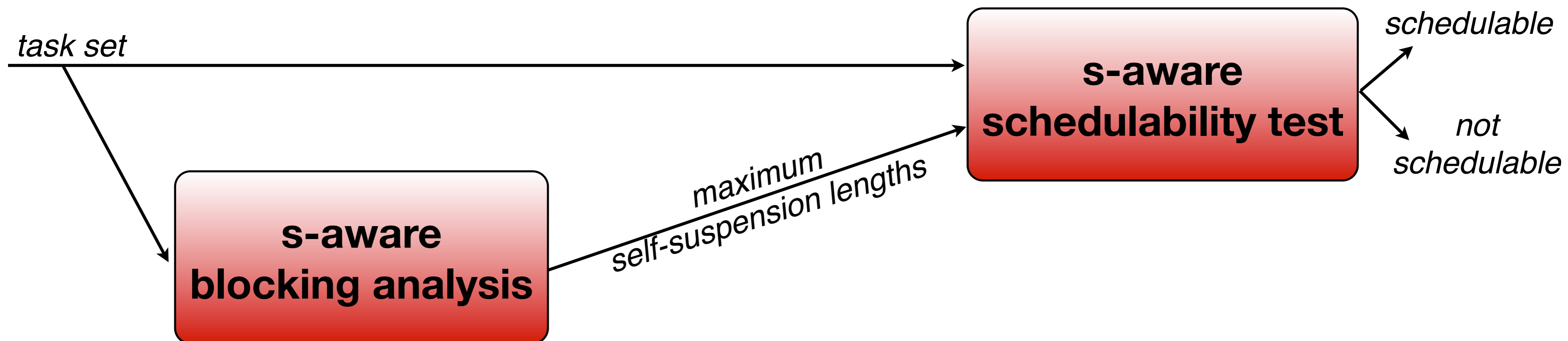
**Motivation**: *reuse existing schedulability analyses that assume **independent**, **always ready** tasks.*

# S-Aware vs. S-Oblivious Analysis

## Suspension-oblivious (*s-oblivious*)  Analysis

*task set* → **s-oblivious blocking analysis** → *modified task set with inflated WCETs* → **s-oblivious schedulability test** → *schedulable* / *not schedulable*

## Suspension-aware (*s-aware*)  Analysis

*task set* → **s-aware schedulability test** → *schedulable* / *not schedulable*

**s-aware blocking analysis** → *maximum self-suspension lengths* → **s-aware schedulability test**
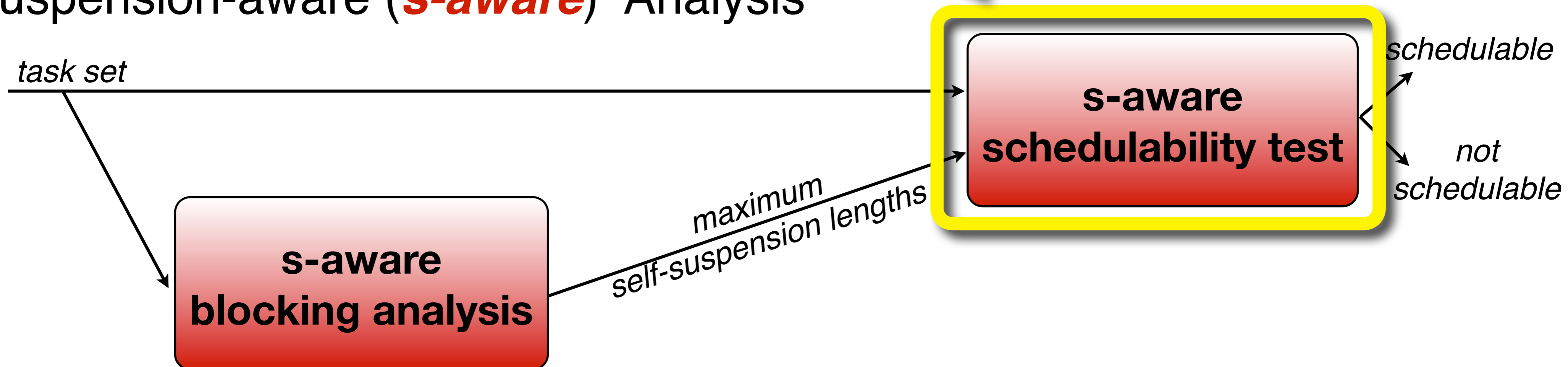
**Requires availability of a schedulability test**
that accounts (reasonably accurately) for **self-suspensions.**

(…which can be tricky to derive)
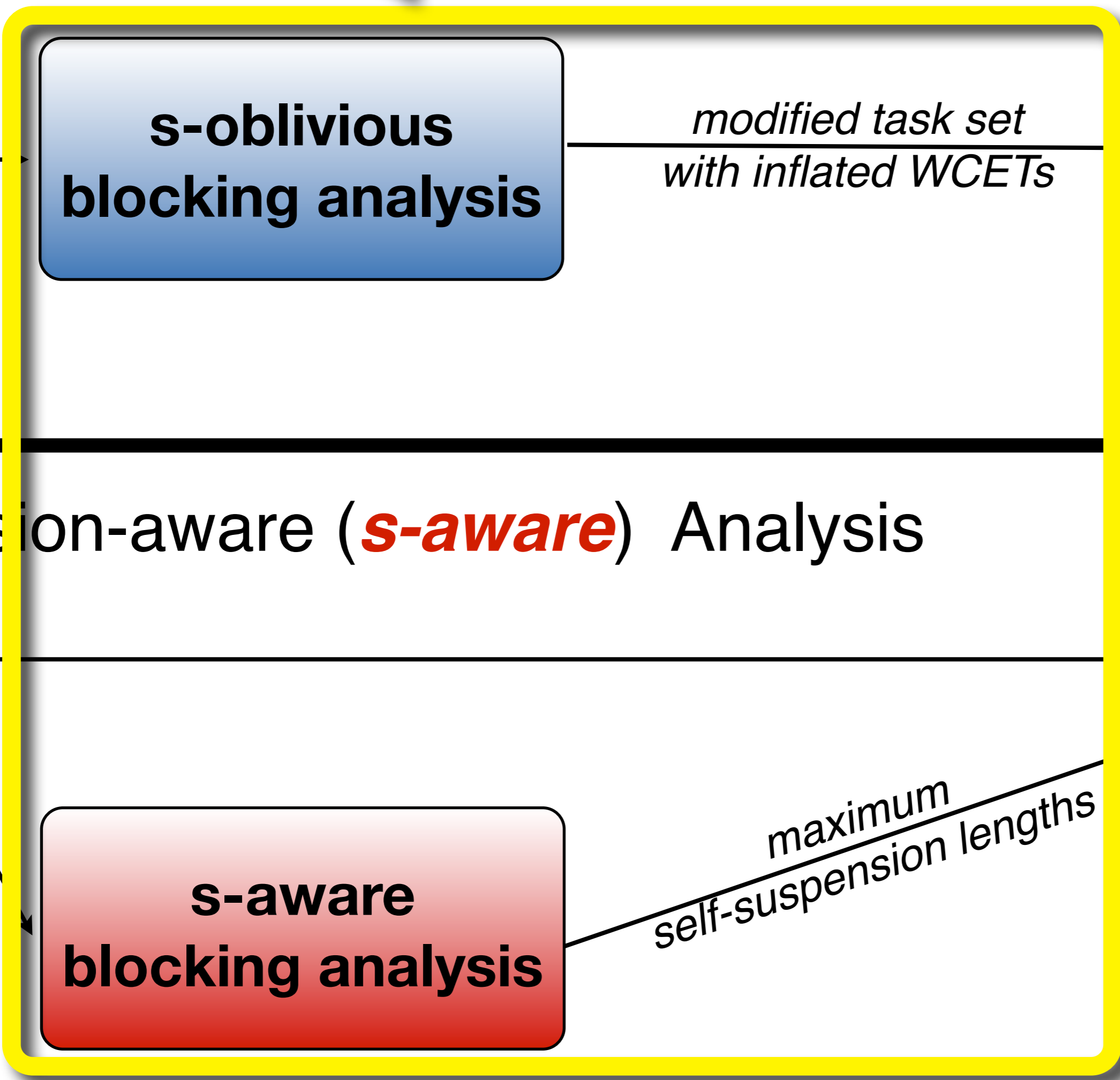
Suspension-oblivious (*s-oblivious*) Analysis

task set → **s-oblivious blocking analysis** → *modified task set with inflated WCETs* → **s-oblivious schedulability test** → *schedulable* / *not schedulable*

Suspension-aware (*s-aware*) Analysis

task set → **s-aware schedulability test** → *schedulable* / *not schedulable*

**s-aware blocking analysis** → *maximum self-suspension lengths* → **s-aware schedulability test**
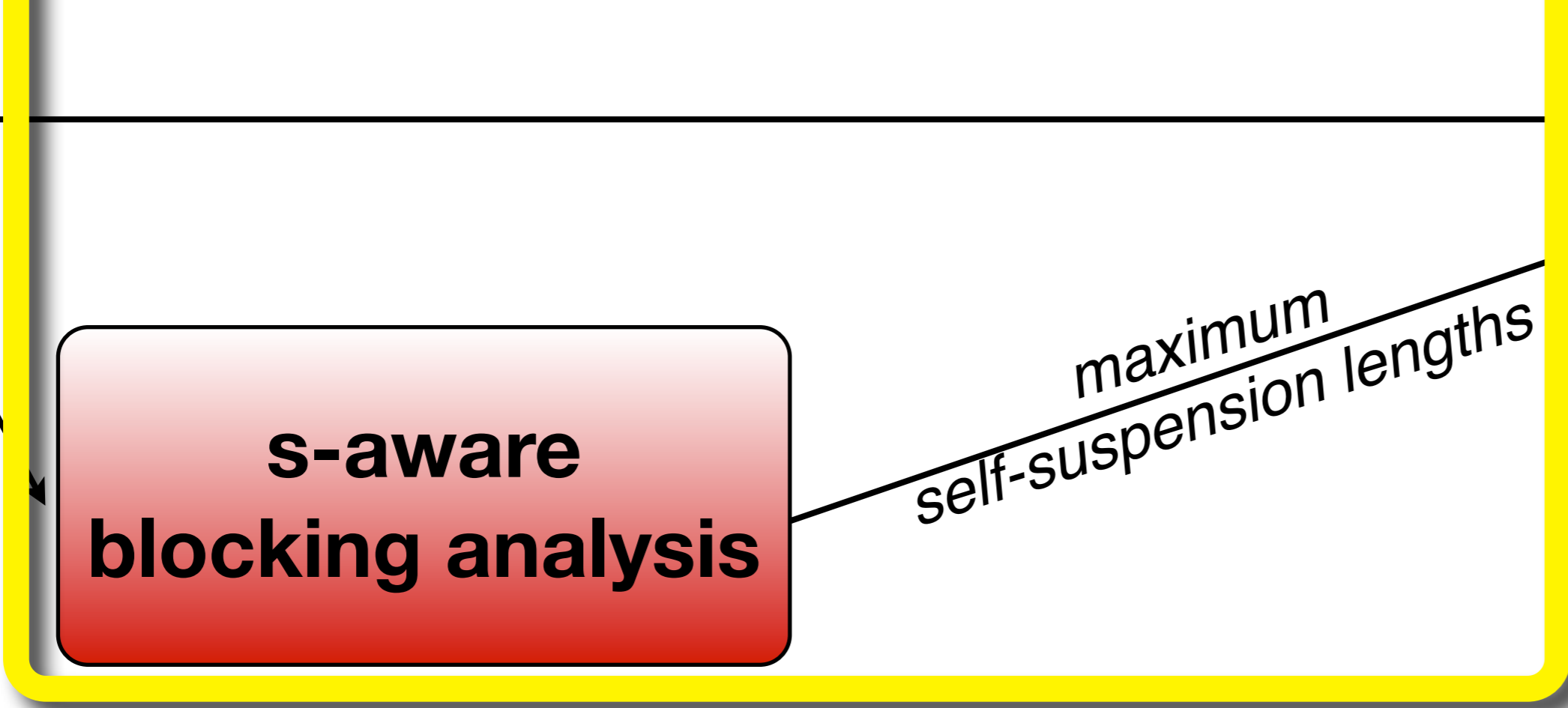
Different notions of **"processor demand"**
→ different definitions of **"priority inversion"**.

s Analysis

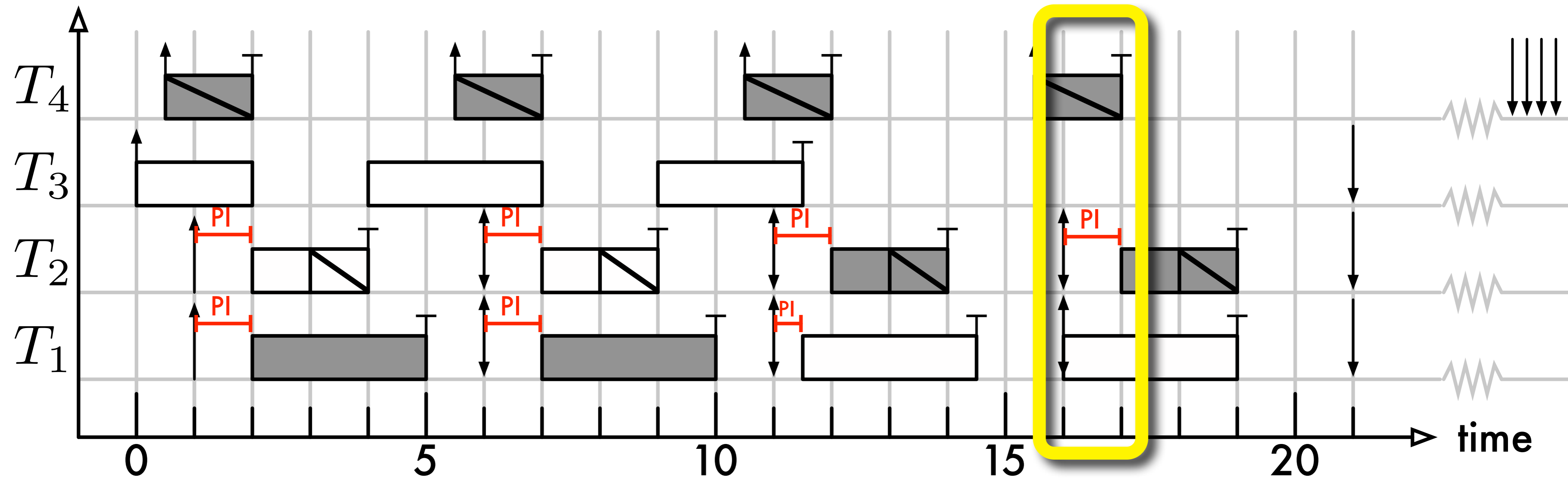Suspension-oblivious (*s-oblivious*) Analysis

task set → **s-oblivious blocking analysis** — *modified task set with inflated WCETs* → **s-oblivious schedulability test** → *schedulable* / *not schedulable*

Suspension-aware (*s-aware*) Analysis

task set → **s-aware schedulability test** → *schedulable* / *not schedulable*

**s-aware blocking analysis** — *maximum self-suspension lengths* →

# Restricted Segment Boosting at <u>Time 16</u>



Generalized FMLP+ schedule.

**(1)** The **lock-holding ready job** (if any) with the **earliest segment start time**.

Generalized FMLP+ schedule.



Legend:

| | scheduled | critical section | | |
|---|---|---|---|---|
| Processor 1 | □ (white) | ◺ (white) | ↑ job release | ≡ job suspended |
| Processor 2 | ■ (grey) | ◺ (grey) | ↓ deadline / ⊤ job completion | ⊢PI⊣ priority inversion |

**(1)** The **lock-holding ready job** (if any) with the **earliest segment start time**.

**(2)** Up to $c - 1 = 1$ jobs from $T_4$'s **co-boosting set** $= \varnothing$.

$T_4$

$T_3$

$T_2$

$T_1$

PI   PI   PI   PI

PI   PI   PI   PI

time

0   5   10   15   20

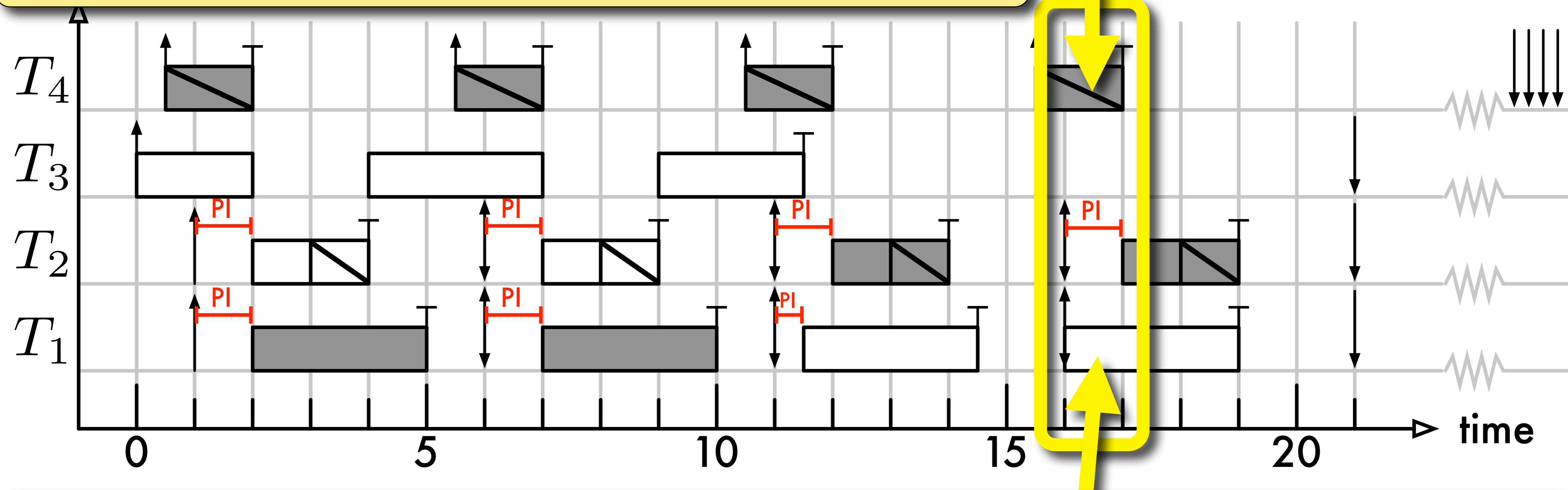| | scheduled | critical section | | job release | | job suspended |
|---|---|---|---|---|---|---|
| Processor 1 | ☐ | ⧄ | ↑ | | | |
| | | | ↓ | deadline | | priority inversion |
| Processor 2 | ▨ | ▨ | ⊤ | job completion | PI | |

**(1)** The **lock-holding ready job** (if any) with the **earliest segment start time**.

**(2)** Up to $c - 1 = 1$ jobs from $T_4$'s **co-boosting set** $= \varnothing$.
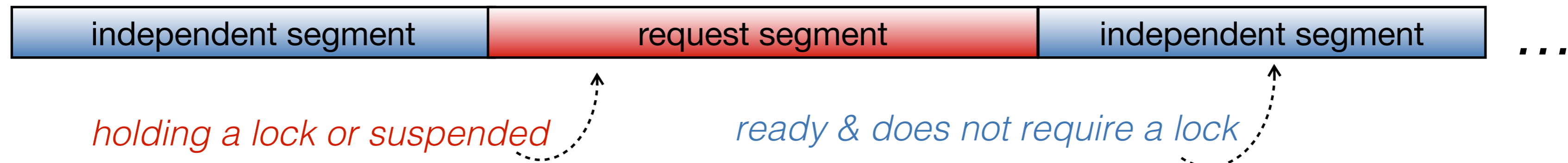


**(3)** If less than $c = 2$ jobs scheduled so far: **any other ready jobs.**

**At time 16**: **one CPU available** after steps 1 & 2
→ schedule **highest-priority task** $T_1$.

# Definition: Job Segments

<u>a job at runtime:</u>

| independent segment | request segment | independent segment | … |

*holding a lock or suspended*    *ready & does not require a lock*

**Independent segment**
➡ starts when a job is **released** or **resumed**,
   or when it **unlocks a resource**
➡ ends when job **completes**, **suspends**, or **requests a lock**

**Request segment**
➡ starts when a job **requests a lock**
➡ ends when it **unlocks the resource**