# A Fully Preemptive
# Multiprocessor Semaphore Protocol for
# Latency-Sensitive Real-Time Applications

ECRTS'13
July 12, 2013

Björn B. Brandenburg
bbb@mpi-sws.org

Max
Planck
Institute
**for**
**Software Systems**

# A Rhetorical Question

*On  <u>uniprocessors</u>, why do we use*

*the* **priority inheritance protocol (PIP)**
*or the* **priority ceiling protocol (PCP)**

*instead of* **simple non-preemptive sections***?*

<u>AUTOSAR Non-Preemptive Critical Section:</u>

```
SuspendAllInterrupts(…);
// critical section
ResumeAllInterrupts(…);
```

# RT 101: Preemptive Synchronization Matters

# RT 101: Preemptive Synchronization Matters

**Deadline miss** due to **latency increase!**

uniprocessor, non-preemptive critical sections

release          deadline

**unrelated**, **latency-sensitive**
high-priority task

**less time-critical**
lower-priority tasks

long lower-priority
critical section (CS)

*time*

Long **non-preemptive** critical section.

# RT 101: Preemptive Synchronization Matters

## uniprocessor, with PIP

release     deadline

**unrelated**, **latency-sensitive**
high-priority task

**less time-critical**
lower-priority tasks

long lower-priority CS

*time*

# RT 101: Preemptive Synchronization Matters

**Latency-sensitive task**
**isolated** from **unrelated** critical section!

uniprocessor, with PIP

release    deadline

**unrelated**, **latency-sensitive**
high-priority task

**less time-critical**
lower-priority tasks

long lower-priority CS

*time*

Lower-priority critical section: **fully preemptive execution**.

# The Multiprocessor Case

*What if we host the same workload on a multiprocessor?*

partitioned multiprocessor scheduling

**unrelated**, **latency-sensitive**
high-priority task

**less time-critical**
lower-priority tasks
(*on same core*)

time

**No existing real-time <u>semaphore</u> protocol for <u>partitioned</u> or <u>clustered</u> scheduling isolates high-priority tasks from unrelated CSs.**

## partitioned multiprocessor scheduling

release          deadline

**unrelated**, **latency-sensitive**
high-priority task

**less time-critical**
lower-priority tasks
(*on same core*)

long lower-priority
critical section (CS)

*time*

MPCP, FMLP, FMLP+, OMLP, …

# This Paper

**Independence preservation** formalizes the idea that
*"tasks should never be delayed by <u>unrelated</u> critical sections."*

# This Paper

**Independence preservation** formalizes the idea that
*"tasks should never be delayed by <u>unrelated</u> critical sections."*

Independence preservation is
**impossible without (limited) job migrations.**

# This Paper

**Independence preservation** formalizes the idea that
*"tasks should never be delayed by <u>unrelated</u> critical sections."*

Independence preservation is
**impossible without (limited) job migrations.**

**First independence-preserving semaphore protocol**
for <u>clustered/partitioned</u> scheduling; the protocol also has
***asymptotically optimal blocking bounds***.

# Clustered JLFP Scheduling

*Job-Level Fixed-Priority Scheduling (**JLFP**)*

*$c$ … number of processors per cluster*

*$m$ … number of processors (total)*



| partitioned scheduling | clustered scheduling | global scheduling |
| --- | --- | --- |
| $c = 1$ | $1 \leq c \leq m$ | $c = m$ |

**This talk: Partitioned Fixed-Priority (P-FP) Scheduling**

*Job-Level Fixed-Priority Scheduling (JLFP)*

*$c$ … number of processors per cluster*

*$m$ … number of processors (total)*



partitioned scheduling

$c = 1$

clustered scheduling

$1 \leq c \leq m$

global scheduling

$c = m$

# Clustered JLFP Scheduling

*Job-Level Fixed-Priority Scheduling (**JLFP**)*

*__c__ … number of processors per cluster*

*__m__ … number of processors (total)*



partitioned scheduling     clustered scheduling     global scheduling

**Task model: implicit-deadline sporadic tasks**
(choice of deadline constraint irrelevant to results)

# Real-Time Semaphore Protocols

## **Binary Semaphores in POSIX**

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

*A blocked task **suspends** & yields the processor.*

# Real-Time Semaphore Protocols

## Binary Semaphores in POSIX

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

*A blocked task **suspends** & yields the processor.*

## Priority Inversion

*A job **should** be scheduled, but **is not**.*

*PI-Blocking: increase in worst-case response time due to priority inversions.*

# Real-Time Semaphore Protocols

**Binary Semaphores
in POSIX**

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

*A blocked task **suspends**
& yields the processor.*

Priority Inversion

*A job **should** be
scheduled, but **is not**.*

PI-Blocking: *increase in
worst-case response time
due to priority inversions.*

Goal: **bounded pi-blocking**.

*Bounded in terms of critical section lengths only!*

# Real-Time Semaphore Protocols

## Binary Semaphores in POSIX

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

*A blocked task **suspends** & yields the processor.*

## Priority Inversion

*A job **should** be scheduled, but **is not**.*

PI-Blocking: *increase in worst-case response time due to priority inversions.*

**Assumptions**
➡ Unnested critical sections.
➡ **Suspension-oblivious** schedulability analysis.

# Part 1

## Avoiding Delays due to Unrelated Critical Sections

# Independence Preservation

*(specific to s-oblivious analysis)*

*"Tasks should never be delayed by unrelated critical sections."*

# Independence Preservation

*(specific to s-oblivious analysis)*

Let $b_{i,q}$ denote the **maximum pi-blocking** incurred by task $T_i$ due to requests for resource $q$.

Let $N_{i,q}$ denote the maximum number of times that any job of $T_i$ **accesses** resource $q$.

Under an **<u>independence-preserving</u>** locking protocol,

if $N_{i,q} = 0$, then $b_{i,q} = 0$.

*"You only pay for what you use."*

# Independence Preservation

*(specific to s-oblivious analysis)*

Let $b_{i,q}$ denote the **maximum pi-blocking** incurred by task $T_i$ due to requests for resource $q$.

Let $N_{i,q}$ denote the maximum number of times that any job of $T_i$ **accesses** resource $q$.

Under an **independence-preserving** locking protocol,

if $N_{i,q} = 0$, then $b_{i,q} = 0.$

Isolation useful for:

**latency-sensitive** workloads (if no delay can be tolerated) or
if low-priority tasks contain **unknown** or **untrusted** critical sections.

# Real-Time Semaphore Protocols

| *real-time locking protocol* | **=** | *progress mechanism* | **+** | *queue structure* |
|---|---|---|---|---|

# Real-Time Semaphore Protocols

| real-time locking protocol | = | progress mechanism | + | queue structure |

**Ensure that a lock holder is scheduled** (while waiting tasks incur pi-blocking).

How to **order conflicting critical sections** (e.g., priority queue, FIFO queues).

# Real-Time Semaphore Protocols

| *real-time locking protocol* | = | *progress mechanism* | + | *queue structure* |

**partitioned** scheduling
**priority boosting**

**clustered** scheduling
**priority donation**

**global** scheduling
**priority inheritance**

**Priority boosting and Priority Donation:**

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

→ **effectively non-preemptive** → **not independence preserving**

| *real-time locking protocol* | = | *progress mechanism* | + | *queue structure* |

**partitioned** scheduling
**priority boosting**

**clustered** scheduling
**priority donation**

**global** scheduling
**priority inheritance**

# Real-Time Semaphore Protocols

| *real-time locking protocol* | = | *progress mechanism* | + | *queue structure* |

**partitioned** scheduling
**priority boosting**

**clustered** scheduling
**priority donation**

**global** scheduling
**priority inheritance**

Existing independence-preserving locking protocols:

**Global PIP**, **Global FMLP**, **Global OMLP**, …

# Observation

*Independence preservation + bounded priority inversion*
**requires *intra-cluster* job migrations**.



partitioned scheduling

clustered scheduling

# Example: Job Migration is Necessary

*three tasks, two cores, one resource, P-FP scheduling*

# Example: Job Migration is Necessary

*three tasks, two cores, one resource, P-FP scheduling*



$T_2$ starts executing **critical section**...

# Example: Job Migration is Necessary

*three tasks, two cores, one resource, P-FP scheduling*



Job of $T_1$ is released.

What to do with **in-progress** critical section?

# Case 1: **priority boosting** (=let $T_2$ continue).

*three tasks, two cores, one resource, P-FP scheduling*

**Case 1: priority boosting** (=let *T₂* continue).

*three tasks, two cores, one resource, P-FP scheduling*

**Case 1: priority boosting** (=let $T_2$ continue).

*three tasks, two cores, one resource, P-FP scheduling*



**Benefit:** $T_3$ incurs only **bounded pi-blocking**, meets deadline.

# Case 1: **priority boosting** (=let $T_2$ continue).

*three tasks, two cores, one resource, P-FP scheduling*



**Problem:** $T_1$ misses its deadline.

# Example: Job Migration is Necessary

*three tasks, two cores, one resource, P-FP scheduling*



Job of $T_1$ is released.

What to do with in-progress critical section?

Core 2

$T_3$

Core 1

$T_2$

$T_1$

*time*

# Case 2: independence preservation (= preempt $T_2$).

*three tasks, two cores, one resource, P-FP scheduling*

**Case 2: independence preservation** (= preempt *T₂*).

*three tasks, two cores, one resource, P-FP scheduling*



**Independence preservation:** *T₁* meets its deadline.

**Case 2: independence preservation** (= preempt $T_2$).

*three tasks, two cores, one resource, P-FP scheduling*

**Case 2: independence preservation** (= preempt $T_2$).

*three tasks, two cores, one resource, P-FP scheduling*



Core 2

$T_3$

*pi-blocked*

CS

**Problem:** $T_3$ incurs "**unbounded**" **pi-blocking**, misses deadline!

Core 1

$T_2$

$T_1$

*time*

# Partitioned Scheduling with Migrations?

# Partitioned Scheduling with Migrations?

Partitioned By Necessity



E.g., SoC with **heterogeneous cores** (ARM, PowerPC, x86, MIPS).

migrations **infeasible**
*for lack of technical capability*

# Partitioned Scheduling with Migrations?

Partitioned By Necessity



E.g., SoC with **heterogeneous cores** (ARM, PowerPC, x86, MIPS).

<u>migrations **infeasible**</u>
*for lack of technical capability*

➔ independence preservation **and** bounded priority inversion
**impossible to achieve!**

# Partitioned Scheduling with Migrations?

Partitioned By Necessity

Partitioned By Choice



migrations **infeasible**
*for lack of technical capability*

migrations **disallowed**
*but technically feasible*

**Occasional migrations** not desirable, but **possible**!

(Focus of this work.)

Partitioned By Necessity

Partitioned By Choice

$J_1$ $J_2$ $J_3$ $J_4$

Q1 Q2 Q3 Q4

Core 1 Core 2 Core 3 Core 4

Local Memory Local Memory Local Memory Local Memory

Shared Memory

$J_1$ $J_2$ $J_3$ $J_4$

Q1 Q2 Q3 Q4

Core 1 Core 2 Core 3 Core 4

L2 Cache L2 Cache

Main Memory

migrations **infeasible**
*for lack of technical capability*

migrations **disallowed**
*but technically feasible*

# Example: Job Migration is Necessary

*three tasks, two cores, one resource, P-FP scheduling*



Job of $T_1$ is released.

What to do with in-progress critical section?

**1) Ensure independence preservation** (= preempt *T₂*).

*three tasks, two cores, one resource, P-FP scheduling*



**Independence preservation:** *T₁* meets its deadline.

## 2) Ensure **bounded pi-blocking** (= schedule $T_2$).

*three tasks, two cores, one resource, P-FP scheduling*



## Easy fix: **migrate** $T_2$ when $T_3$ suspends.

# ...igration is Necessary

*three tasks, two cores, one resource, P-FP scheduling*

**Temporarily** move *T₃* to Core 2...



Core 2 — *T₃* — pi-blocked — CS

Core 1 — *T₂* — *T₁*

time

**Easy fix: migrate *T₂* when *T₃* suspends.**

# Example: Job Migration is Necessary

**Benefit:** $T_3$ incurs only **bounded pi-blocking**, meets deadline.



**Core 2**

$T_3$

*pi-blocked*

CS

**Core 1**

$T_2$

$T_1$

*time*

**Easy fix: migrate $T_2$ when $T_3$ suspends.**

# Theorem

*Under non-global scheduling (**c** ≠ **m**), it is **impossible** for a semaphore protocol to simultaneously*

***(i)*** *prevent **unbounded pi-blocking**,*

***(ii)*** *be **independence-preserving**, and*

***(iii)*** *avoid **inter-cluster job migrations**.*

*Pick any two…*

# Combinations of Properties

*Under non-global scheduling (**c ≠ m**), it is **impossible** for a semaphore protocol to simultaneously*

**(i)** *prevent* **unbounded pi-blocking***,*

**(ii)** *be* **independence-preserving***, and*

**(iii)** *avoid* **inter-cluster job migrations***.*

**(i) & (iii)**
➡ MPCP, Part. FMLP, FMLP+, OMLP, …

**(ii) & (iii)**
➡ Applying PIP to partitioned scheduling (**not sound**!)

**(i) & (ii)**
➡ **no such protocol known!**

# Part 2

## Independence Preservation
## +
## Asymptotically Optimal
## PI-Blocking

# High-Level Overview

| real-time locking protocol | = | progress mechanism | + | queue structure |

# High-Level Overview

| *real-time locking protocol* | = | *progress mechanism* | + | *queue structure* |

Must be
independence-preserving.

Must ensure
asymptotic optimality.

# High-Level Overview

| *real-time locking protocol* | = | *progress mechanism* | + | *queue structure* |

Must be
independence-preserving.

Must ensure
asymptotic optimality.

Adopt intuition from example:

when **lock holder is preempted**,
**migrate to blocked task's processor.**

# Migratory Priority Inheritance

**classic priority inheritance**

*inherit **priority** of blocked jobs*

# Migratory Priority Inheritance

**classic priority inheritance**

*inherit* **priority** *of blocked jobs*

**+**

**"cluster inheritance"**

*inherit* **eligibility to execute on assigned clusters** *from blocked jobs*

Jobs remain **fully preemptive** even in critical sections.

➜ enables **independence preservation**

**classic priority inheritance**

*inherit **priority** of blocked jobs*

**+**

**"cluster inheritance"**

*inherit **eligibility to execute
on assigned clusters**
from blocked jobs*

# High-Level Overview

| *real-time locking protocol* | = | *progress mechanism* ✓ | + | *queue structure* ? |

Must be independence-preserving.

Must ensure asymptotic optimality.

# High-Level Overview

| *real-time locking protocol* | = | *progress mechanism* ✅ | + | *queue structure* ❓ |

Must be
independence-preserving.

Must ensure
asymptotic optimality.

**Resolve (most) contention within clusters:
use a multi-level queue.**

# A 3-Level FIFO/FIFO/PRIO Queue

*one 3-level queue for each resource*

Cluster 1

Cluster 2

shared
resource

⋮

Cluster **K**

# A 3-Level FIFO/FIFO/PRIO Queue

*one 3-level queue for each resource*

Cluster 1

$PQ_1$ priority queue

shared resource

Cluster 2

$PQ_2$ priority queue

⋮

Cluster **K**

$PQ_K$ priority queue

# A 3-Level FIFO/FIFO/PRIO Queue

*one 3-level queue for each resource*

Cluster 1

priority queue

$PQ_1$

FIFO Queue $FQ_1$

shared resource

Cluster 2

priority queue

$PQ_2$

FIFO Queue $FQ_2$

Cluster $K$

priority queue

$PQ_K$

FIFO Queue $FQ_K$

# Queue

**Bounded length**: at most **c** jobs (in each cluster).

(**c** = *number of cores in cluster*)

Cluster 1

FIFO Queue **FQ₁**

**PQ₁** priority queue

🔒 shared resource

Cluster 2

FIFO Queue **FQ₂**

**PQ₂** priority queue

⋮

Cluster **K**

FIFO Queue **FQₖ**

**PQₖ** priority queue

# A 3-Le

**Priority queue** used only if more than **c** jobs contend.

(**c** = *number of cores in cluster*)

*one 3-level queue for each resource*

Cluster 1

priority
queue

*PQ₁*

FIFO Queue *FQ₁*

Cluster 2

shared
resource

priority
queue

*PQ₂*

FIFO Queue *FQ₂*

⋮

Cluster *K*

priority
queue

*PQ_K*

FIFO Queue *FQ_K*

# A 3-Level FIFO/FIFO/PRIO Queue

*one 3-level queue for each resource*

# The O($m$) Independence-Preserving Locking Protocol (OMIP)

| **The OMIP** | = | *migratory priority inheritance* | + | *3-level F/F/P queue* |
| :---: | :---: | :---: | :---: | :---: |

independence-preserving

**O($m$)** s-oblivious pi-blocking

# The O($m$) Independence-Preserving Locking Protocol (OMIP)

| **The OMIP** | = | *migratory priority inheritance* | + | *3-level F/F/P queue* |
|:---:|:---:|:---:|:---:|:---:|

independence-preserving

**O($m$)** s-oblivious pi-blocking

$\Omega(m)$ lower bound on s-oblivious pi-blocking (—  & Anderson, 2010)

→ The OMIP ensures **asymptotically optimal** s-oblivious pi-blocking.

# Part 3

Evaluation

# Prototype Implementation



Linux Testbed for Multiprocessor Scheduling in Real-Time Systems

## www.litmus-rt.org

**3-level queues**

➡ easy (reuse Linux wait queues)

➡ cheap compared to syscall

**Migratory priority inheritance**

➡ more tricky (need to avoid global locks)

➡ store bitmap of cores "offering" to
   schedule lock holder in each lock

# Response Times Experiment

*on an 8-core, 2-Ghz Xeon X7550 System*

Core 1                                    Core 8

$T_1$                                      $T_8$

$T_9$          ...                         $T_{16}$

$T_{17}$                                   $T_{24}$

$T_{25}$                                   $T_{32}$

**Setup**

➡ 4 tasks on each core (one independent & latency-sensitive)

➡ one shared resource

➡ max. critical section length: **~1ms**

# ...imes Experiment

*...Ghz Xeon X7550 System*

One **latency-sensitive**, independent task with **period = 1ms**.

Core 1                                                   Core 8

$T_1$                      ...                             $T_8$

$T_9$                                                      $T_{16}$

$T_{17}$                                                   $T_{24}$

$T_{25}$                                                   $T_{32}$

**Setup**
➡ 4 tasks on each core (one independent & latency-sensitive)
➡ one shared resource
➡ max. critical section length: **~1ms**

# Respo

*on an 8*

Core 1                    Core 8

**Three tasks** (on each core) with
**periods** **25 ms**, **100 ms**, and **1000 ms**.
Each job of these tasks locks the resource once.

$T_1$                                            $T_8$

...

$T_9$                                            $T_{16}$

$T_{17}$                                         $T_{24}$

$T_{25}$                                         $T_{32}$

**Setup**
➡ 4 tasks on each core (one independent & latency-sensitive)
➡ one shared resource
➡ max. critical section length: **~1ms**

# Response Times Experiment

*on an 8-core, 2-Ghz Xeon X7550 System*

Core 1

$T_1$

$T_9$

$T_{17}$

$T_{25}$

...

Core 8

$T_8$

$T_{16}$

$T_{24}$

$T_{32}$

**Three Configurations**

➡ **No locks** **(unsound!)**
‣ no blocking (baseline)

➡ **Clustered OMLP**
‣ priority donation

➡ **OMIP**
‣ migratory priority inheritance

**Experiment**

➡ Measured response times with `sched_trace`

➡ 30-minute traces

➡ more than **45 million jobs**

# Response Time CDF



*higher is better*

increasing response time

Fraction of jobs with response time at most **X**

P(response time ≤ X)

response time (in ms)

# Response Time CDF of 1-ms Tasks

# Response Time CDF of 1-ms Tasks



With **priority donation** (or priority boosting), **~20% of the jobs** of the 1ms-tasks **miss their deadline.**

Legend:
- NONE CDF
- OMIP CDF
- OMLP CDF

Y-axis: **P(response time ≤ X)**
X-axis: **response time (in ms)**

# Response Time CDF of 1-ms Tasks



Response time distribution **under OMIP equivalent** to case **without locks**.

(**OMIP** & **NONE** curves overlap)

Legend:
- NONE CDF
- OMIP CDF
- OMLP CDF

X-axis: **response time (in ms)**
Y-axis: **P(response time ≤ X)**

# Response Time CDF of 100-ms Tasks

# Response Time CDF of 100-ms Tasks

**NONE**

**OMLP**

**OMIP**



P(response time ≤ X) vs response time (in ms)

- NONE CDF
- OMIP CDF
- OMLP CDF

# Response Time CDF of 100-ms Tasks



**OMIP:** blocking **shifted** to **lower-priority** (= later-deadline) **jobs**.

Legend:
- NONE CDF
- OMIP CDF
- OMLP CDF

Y-axis: **P(response time ≤ X)**
X-axis: **response time (in ms)**

# Analytical Blocking/Latency Tradeoff

**Large-scale schedulability experiments**

➡ Varied #tasks, #cores, #resources, max. critical section lengths, etc.

➡ >150,000,000 task sets

➡ 678 schedulability plots, available in online appendix

# Analytical Blocking/Latency Tradeoff

**Large-scale schedulability experiments**

➡ Varied #tasks, #cores, #resources, max. critical section lengths, etc.

➡ >150,000,000 task sets

➡ 678 schedulability plots, available in online appendix

*In the presence of **latency-sensitive tasks**, the **OMIP** is generally the **only viable option**.*

# Analytical Blocking/Latency Tradeoff

**Large-scale schedulability experiments**

➡ Varied #tasks, #cores, #resources, max. critical section lengths, etc.

➡ >150,000,000 task sets

➡ 678 schedulability plots, available in online appendix

*In the presence of **latency-sensitive tasks**, the **OMIP** is generally the **only viable option**.*

***Without** latency-sensitive tasks, the **OMIP** does **not** offer substantial improvements.*

# Conclusion

# Summary

**Independence preservation** formalizes the idea that *"tasks should not be delayed by <u>unrelated</u> critical sections."*

Independence preservation is
**impossible without (limited) job migrations**.

The **OMIP** is the first independence-preserving <u>semaphore</u> protocol for clustered scheduling. It ensures **asymptotically optimal s-oblivious pi-blocking**.

Future Work

Nesting

Budget Overruns

Suspension-Aware Analysis

# Thanks!

**LITMUS^RT**

Linux Testbed for Multiprocessor Scheduling in Real-Time Systems

www.litmus-rt.org

**SchedCAT**

**Sched**ulability test
**C**ollection **A**nd **T**oolkit

www.mpi-sws.org/~bbb/
projects/schedcat

# Appendix

# Design Inspirations

**Migrate to Blocked Task's CPU**

➡ "**Local helping**" in TU Dresden's Fiasco/L4

  ‣ *Hohmuth & Peter (2001)*

➡ Multiprocessor **bandwidth inheritance** (MBWI)

  ‣ *Faggioli, Lipari, & Cucinotta (2010)*

**Queue Design**

➡ **Intra**-cluster queues adopted from **global OMLP**

  ‣ *— & Anderson (2010)*

➡ **Inter**-cluster queues similar to **clustered OMLP**

  ‣ *— & Anderson (2011)*

# What about overheads?



**Aren't job migrations expensive?**

➡ response time experiments **reflect all overheads in real system**

➡ latency-sensitive tasks do not migrate, only lower-priority tasks do

➡ only working set of critical section migrates (likely small), not entire task working set (likely much larger)

➡ the critical section would have been preempted anyway

# Blocking Analysis

priority queue

**PQ**

FIFO Queue **GQ**

FIFO Queue **FQ**

$m$ … number of processors (total)      $c$ … number of processors per cluster

# Blocking Analysis

priority
queue

**PQ**

FIFO Queue **FQ**

FIFO Queue **GQ**

at most
$K - 1 = m / c - 1$
queued jobs

at most
$c - 1$
queued jobs

*m … number of processors (total)*          *c … number of processors per cluster*

# Blocking Analysis

**PQ**

priority queue

FIFO Queue **GQ**

FIFO Queue **FQ**

at most
**$K - 1 = m / c - 1$**
queued jobs

at most
**$c - 1$**
queued jobs

at most
**$m / c - 1$**
blocking CS

at most
**$(c - 1) \cdot (m / c)$**
blocking CS

*$m$ … number of processors (total)*     *$c$ … number of processors per cluster*

At most

$$m / c - 1 + (c - 1) \cdot (m / c) = m - 1 = O(m)$$

blocking critical sections.

**PQ** priority queue

FIFO Queue **GQ**　　FIFO Queue **FQ**

at most
**K - 1 = m / c - 1**
queued jobs

at most
**c - 1**
queued jobs

at most
**m / c - 1**
blocking CS

at most
**(c - 1) · (m / c)**
blocking CS

**m** *… number of processors (total)*　　**c** *… number of processors per cluster*

At most

$$m \;/\; c \;-\; 1 \;+\; (c \;-\; 1) \cdot (m \;/\; c) \;=\; m \;-\; 1 \;=\; O(m)$$

blocking critical sections.

**PQ**   priority queue

FIFO Queue **GQ**   ←   FIFO Queue **FQ**

at most
$$K \;-\; 1 \;=\; m \;/\; c \;-\; 1$$
queued jobs

at most
$$c \;-\; 1$$
queued

Under **s-oblivious** analysis**:**
at most **O(m)** critical sections
cause **pi-blocking**.

at most
$$m \;/\; c \;-\; 1$$
blocking CS

at most
$$(c \;-\; 1) \cdot (m \;/\; c)$$
blocking CS

*m … number of processors (total)*          *c … number of processors per cluster*