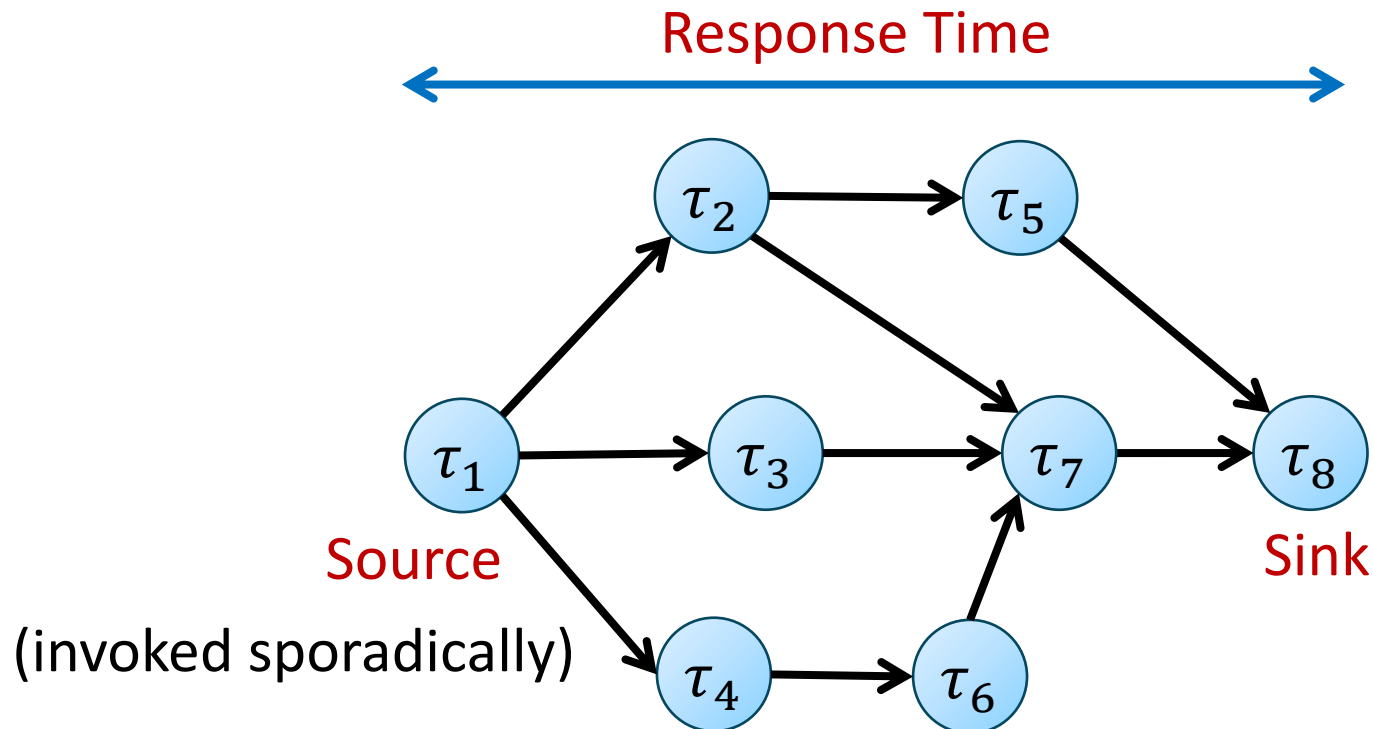# Scheduling Processing Graphs of Gang Tasks on Heterogeneous Platforms

Shareef Ahmed, Denver Massey and James H. Anderson
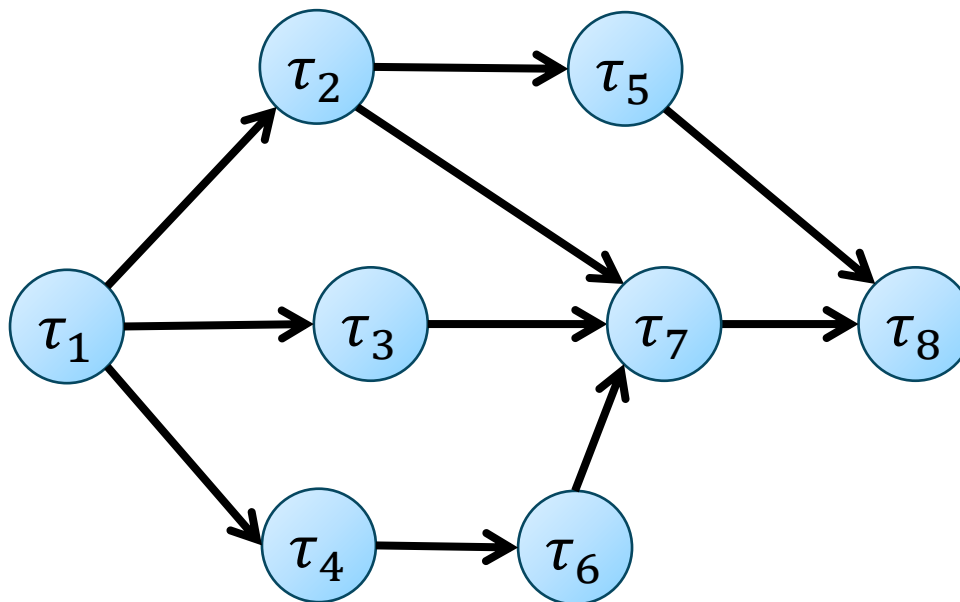
THE UNIVERSITY
of NORTH CAROLINA
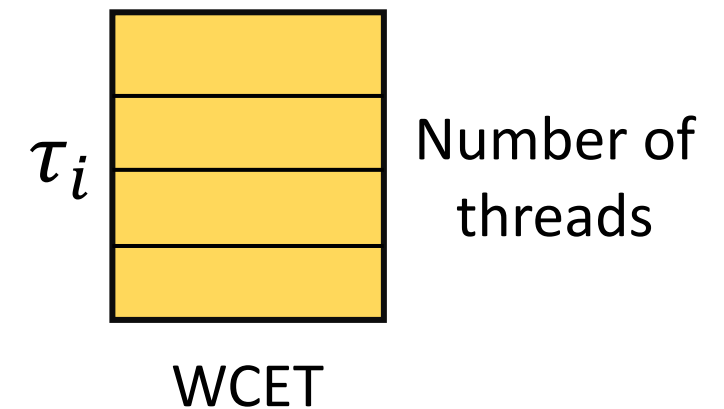at CHAPEL HILL

# Processing Graphs

Response Time

$\tau_2$ → $\tau_5$

$\tau_1$ → $\tau_3$ → $\tau_7$ → $\tau_8$

Source

(invoked sporadically)

$\tau_4$ → $\tau_6$

Sink

- Node = Task

- Edge = Precedence constraint

- Goal:
  - Response time $\leq$ Deadline

# Processing Graphs of Gang Nodes



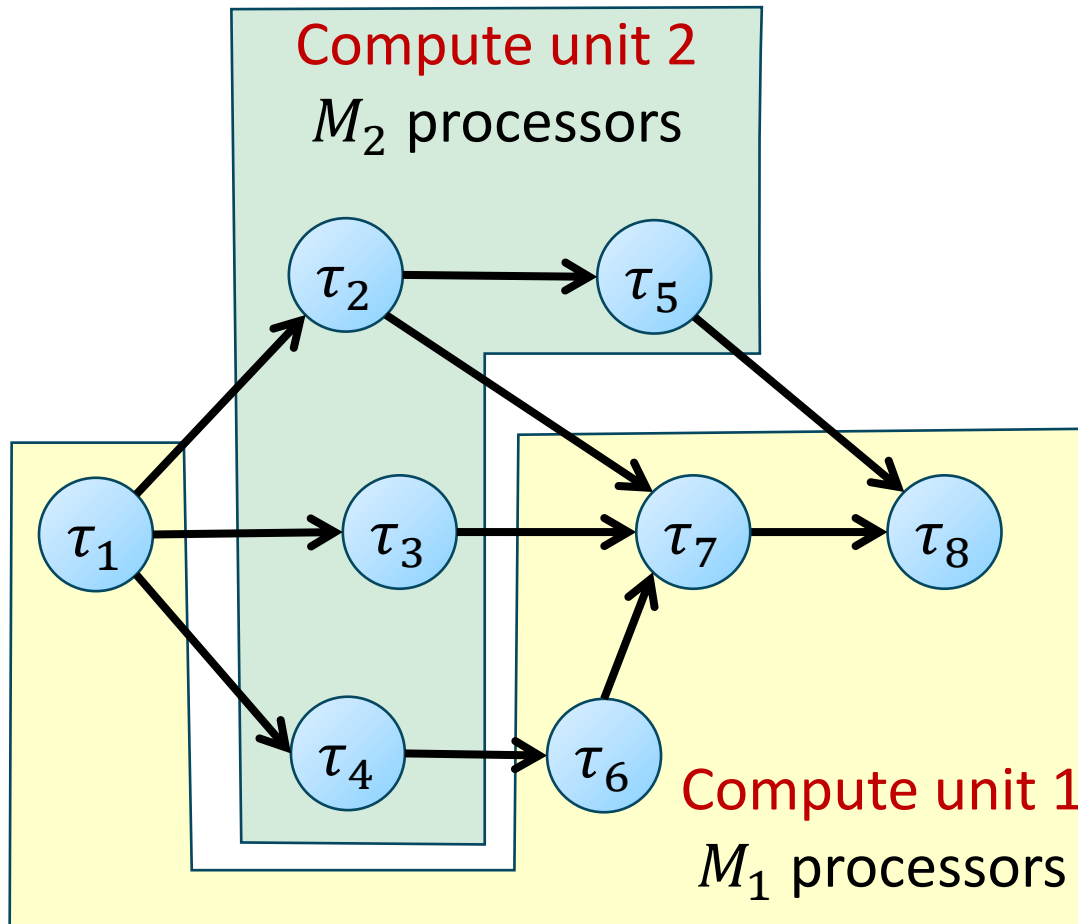- Nodes are <u>**rigid gang**</u> tasks

$\tau_i$ Number of threads

WCET

- Rigid = Same number of threads for all jobs of $\tau_i$

# Heterogenous Compute Platform

Compute unit 2
$M_2$ processors

Compute unit 1
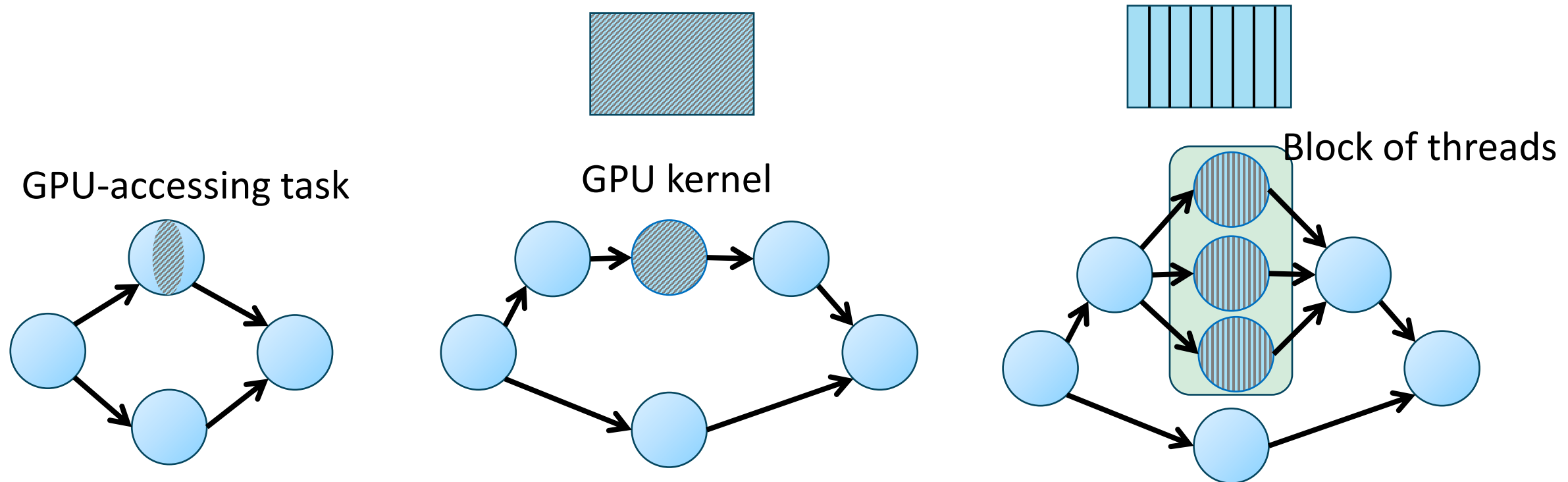$M_1$ processors

- Multiple compute units

- Each node assigned to a compute unit

- Each compute unit has multiple same-speed processors

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# A Use Case

Scheduling processing graphs on multicore+GPU



GPU-accessing task

GPU kernel

Block of threads

# Problem
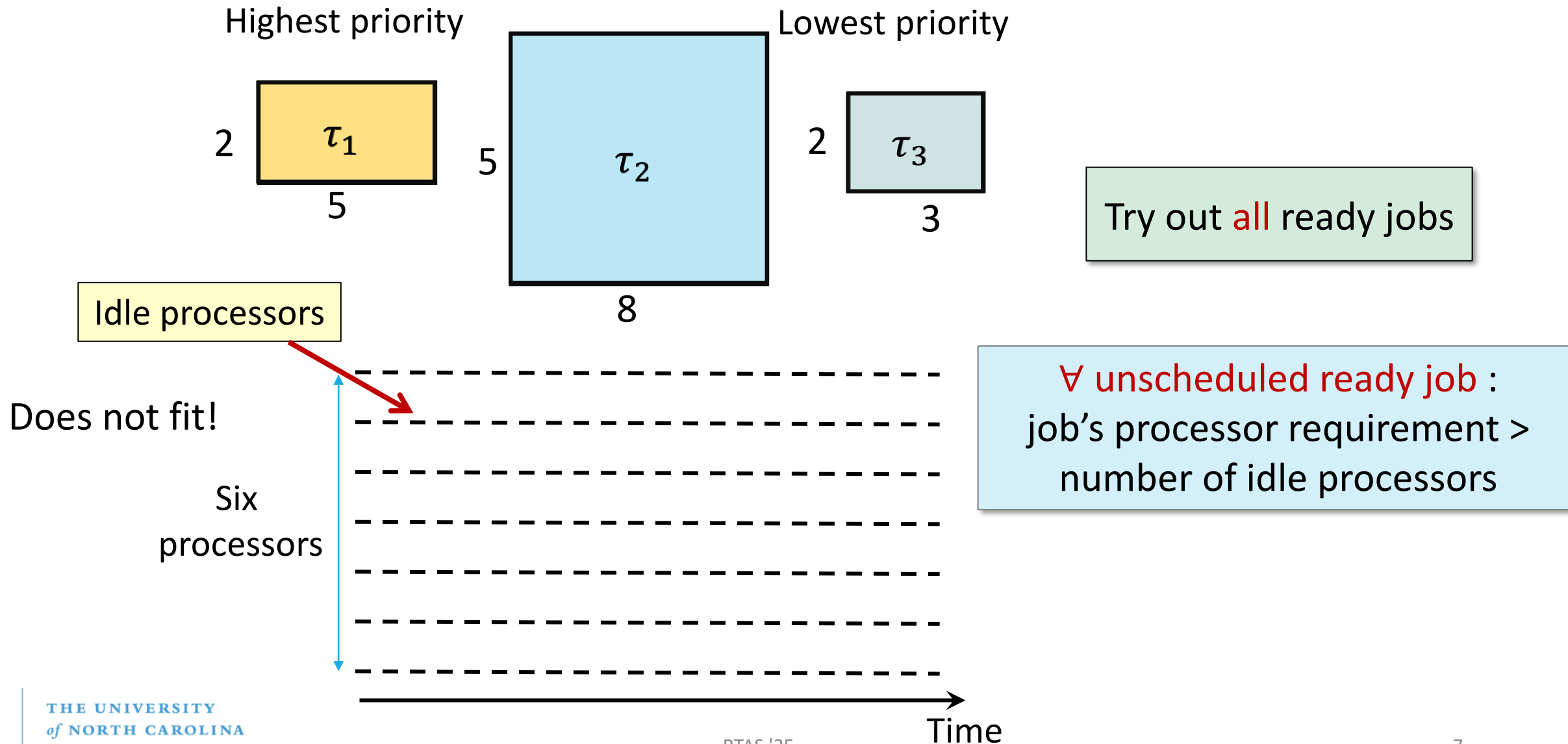
Determine response-time bound of DAGs formed by gang tasks

Assumption: Constrained deadline

Scheduling: Work-conserving, Semi work-conserving

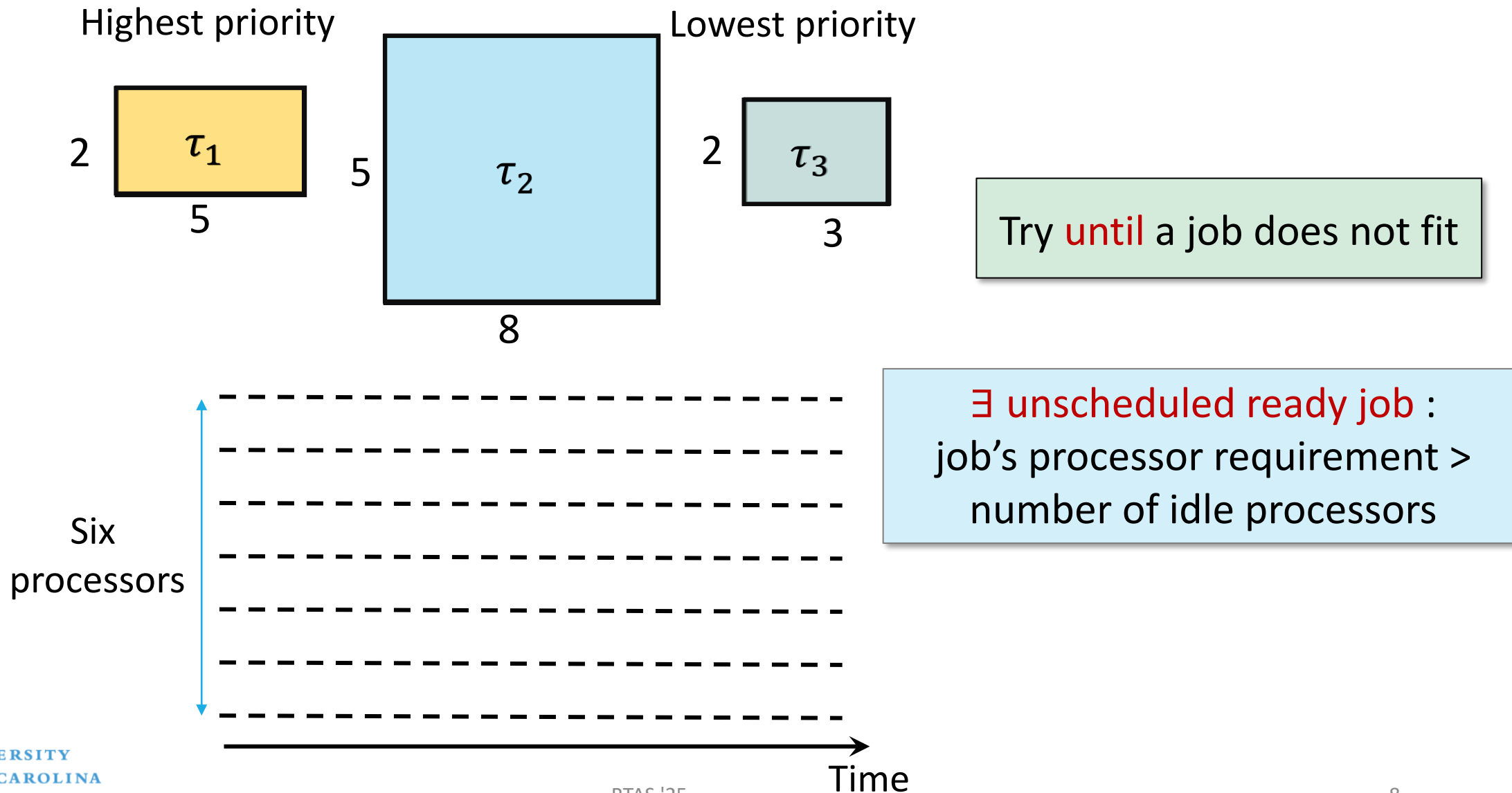Each DAG receives dedicated number of processors on each compute unit

# Work-Conserving Scheduling

Highest priority                          Lowest priority

$$\tau_1$$

2    $\tau_1$

5

$$5 \quad \tau_2$$

8

$$2 \quad \tau_3$$

3

Try out all ready jobs

Idle processors

Does not fit!

Six
processors

∀ unscheduled ready job :
job's processor requirement >
number of idle processors

Time

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Semi Work-Conserving Scheduling



Highest priority        Lowest priority

2  $\tau_1$

5

5  $\tau_2$

8

2  $\tau_3$

3

Try until a job does not fit

∃ unscheduled ready job :
job's processor requirement >
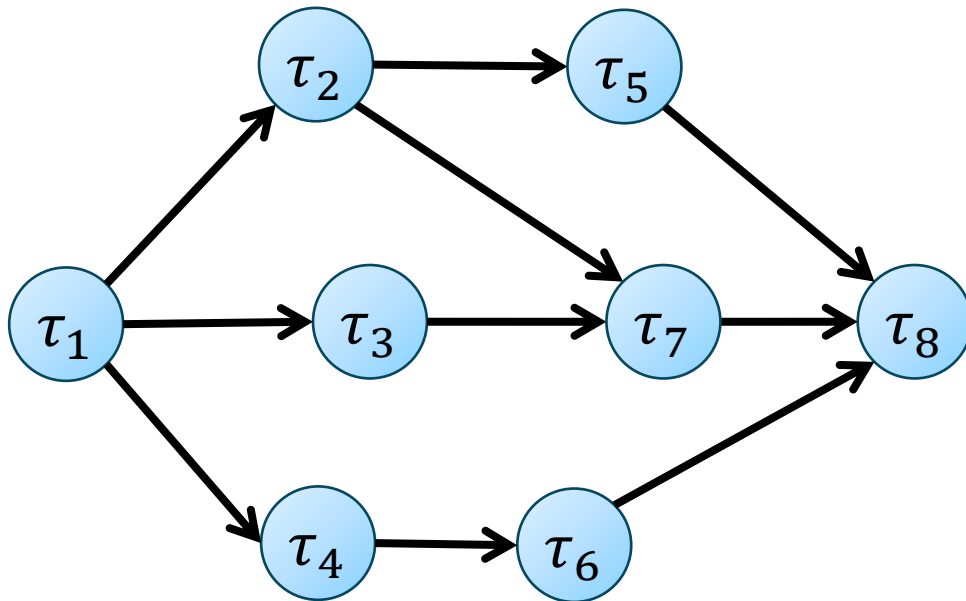number of idle processors

Six
processors

Time

# Why Semi Work-Conserving Scheduling?

Scheduling in NVIDIA GPU is semi work-conserving when all GPU work is submitted from the same address space (and some more constraints)
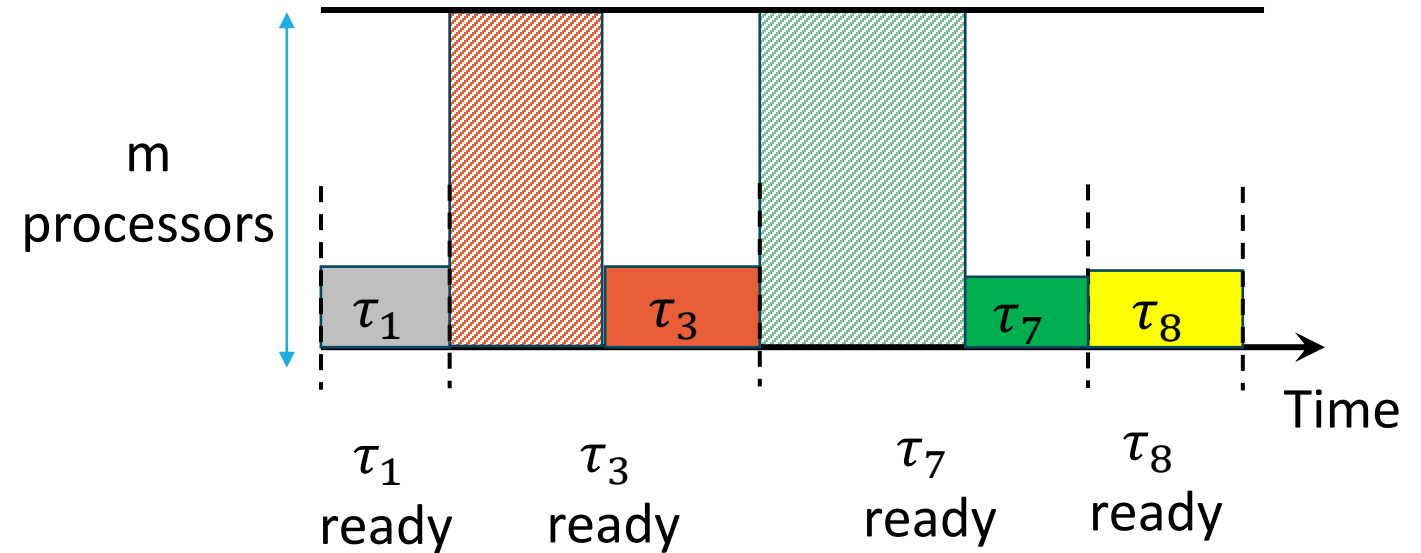
1. Amert et al., RTSS 2017
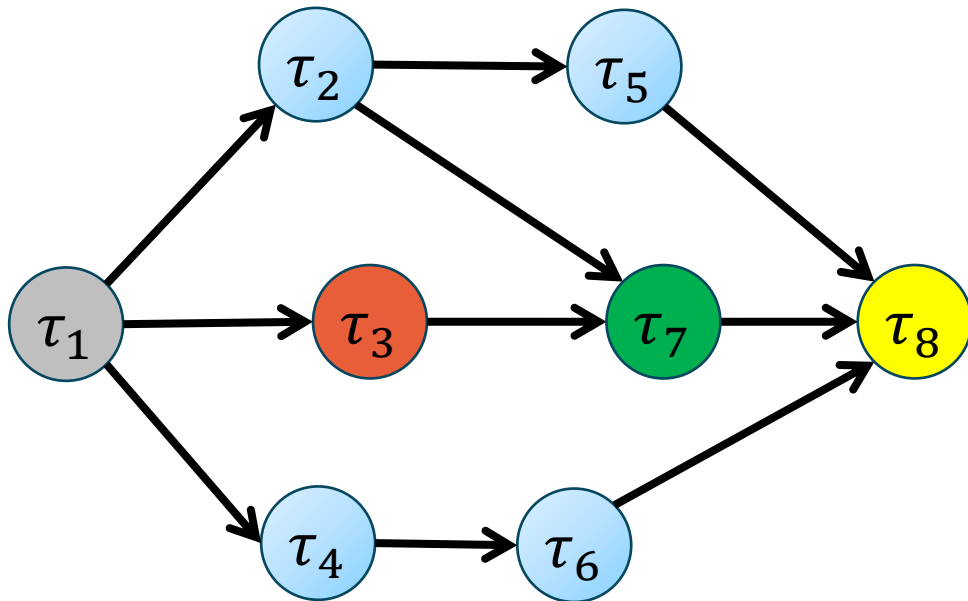2. Bakita and Anderson, RTAS 2024

# Response-Time Bound



m processors

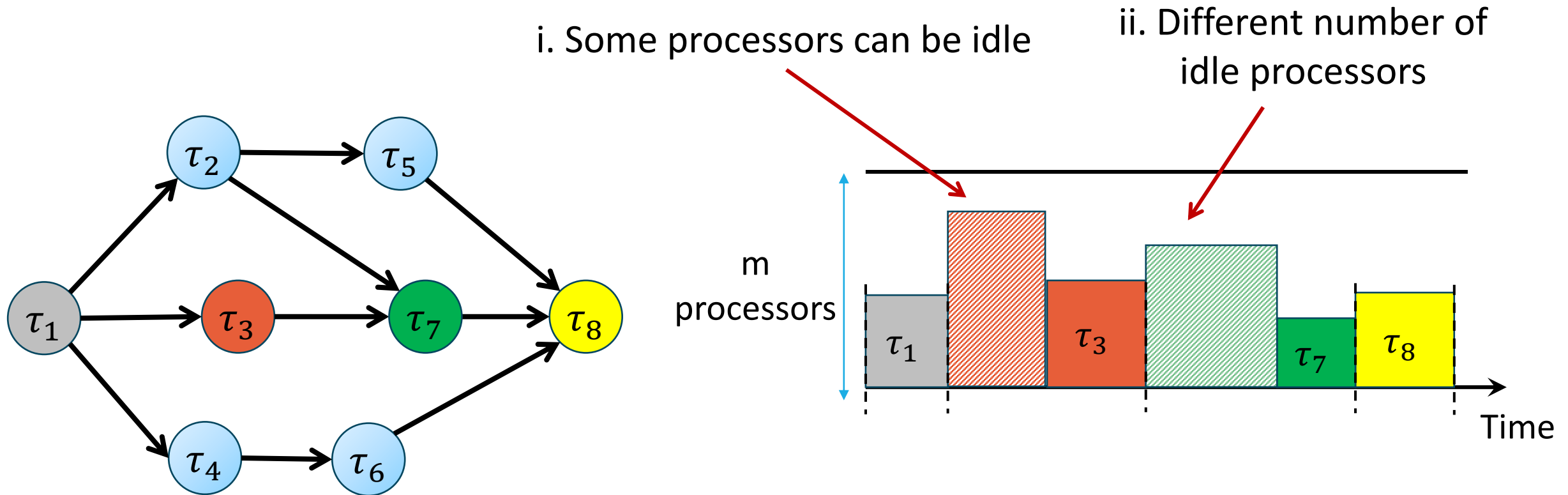$\tau_1$ ready          $\tau_3$ ready          $\tau_7$ ready          $\tau_8$ ready

Assumption 1: One compute unit
Assumption 2: Sequential node (one thread per task)

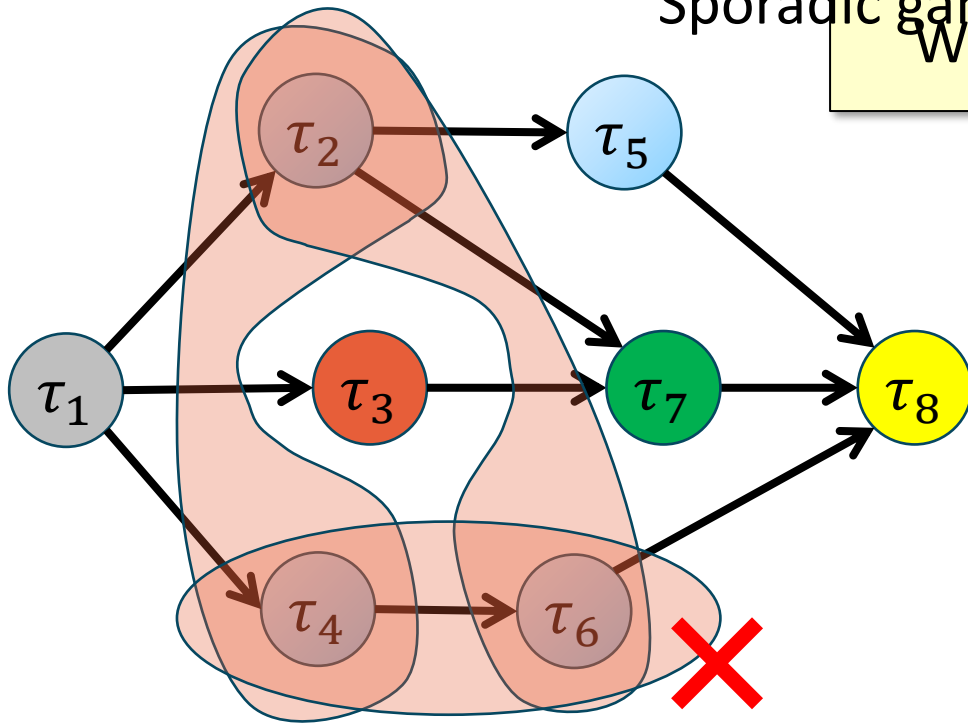Graham, Siam J. of Appl. Math., 1969

# Response-Time Bound



Interfering workload / m

m processors

$\tau_1$

$\tau_3$

$\tau_7$

$\tau_8$

Time

Longest path

Graham, Siam J. of Appl. Math., 1969

Assumption 1: One compute unit

Assumption 2: Sequential node (one thread per task)

# Response-Time Bound

i. Some processors can be idle

ii. Different number of idle processors



Assumption 1: One compute unit

Assumption 2: ~~Sequential~~ node (~~one thread~~ per task)
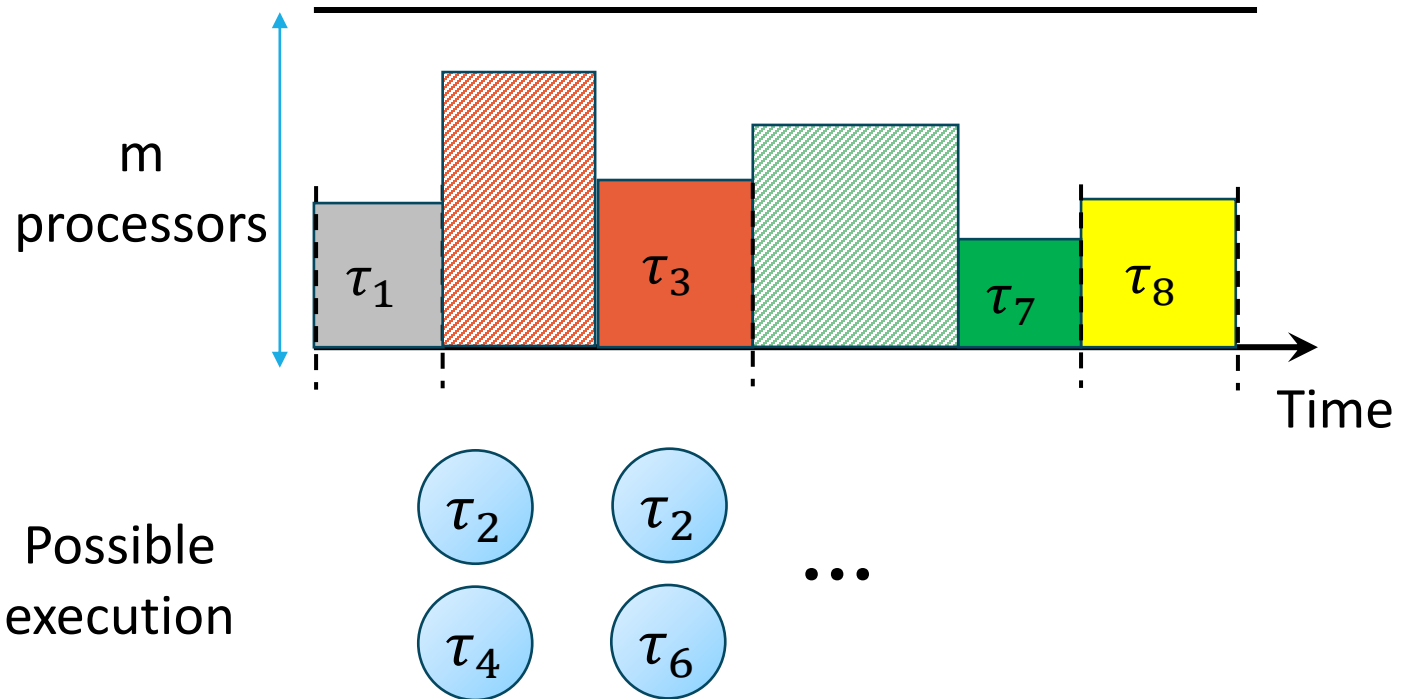              Gang              Multiple threads

# Response-Time Bound

<u>Step 1</u>: Determine the minimum number of busy processors when $\tau_i$ is ready but unscheduled

Sporadic gang task: Dong and Liu, RTSS 2017
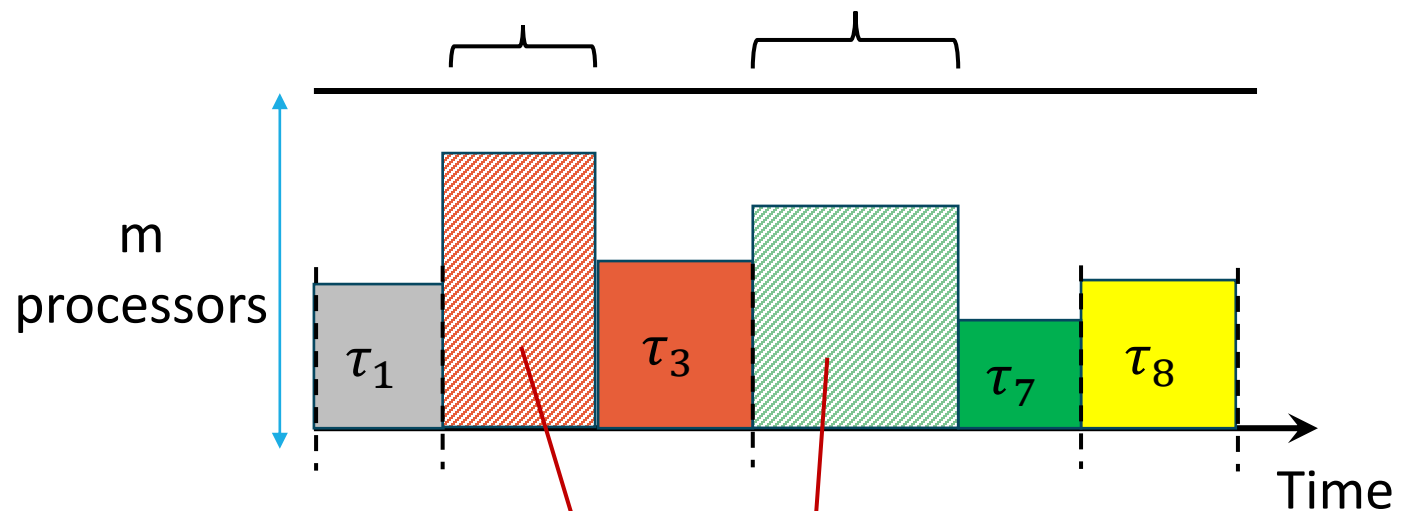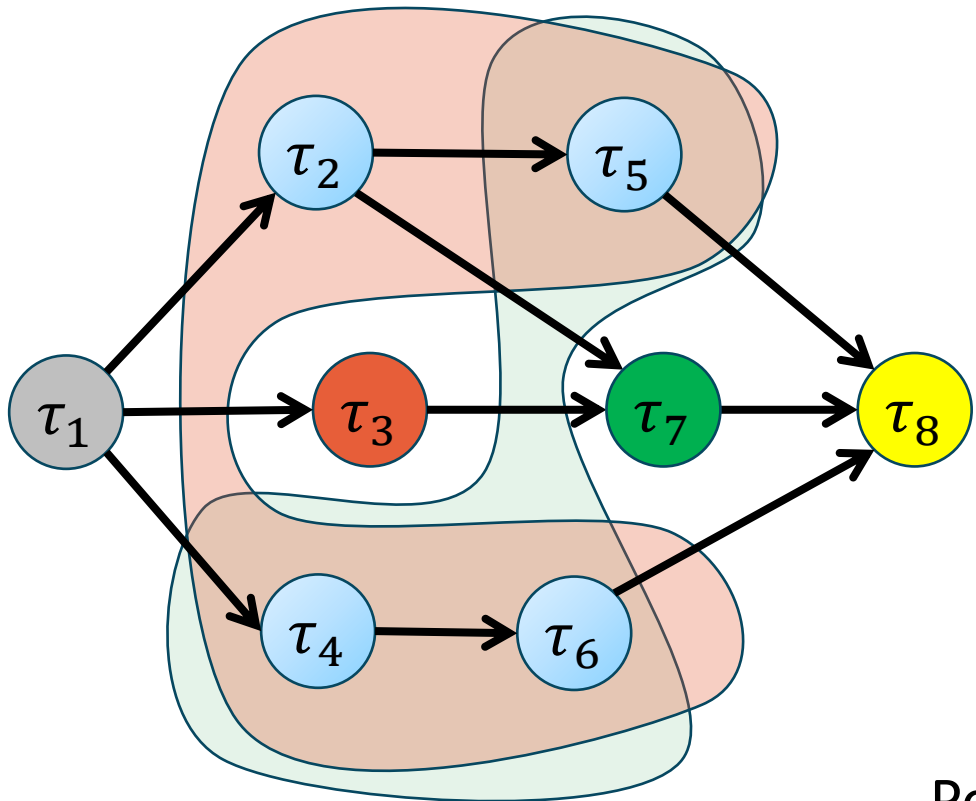
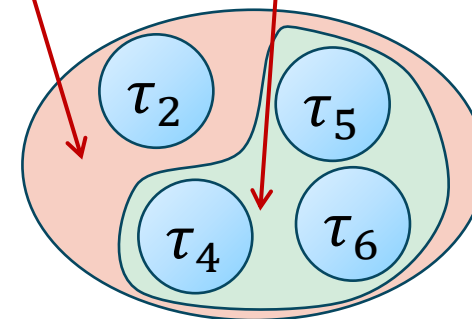We give a dynamic programming algorithm



Assumption 1: One compute unit

Possible execution

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Response-Time Bound

Step 2: Upper bound total interference time



m processors

Time

Possible interfering tasks

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL
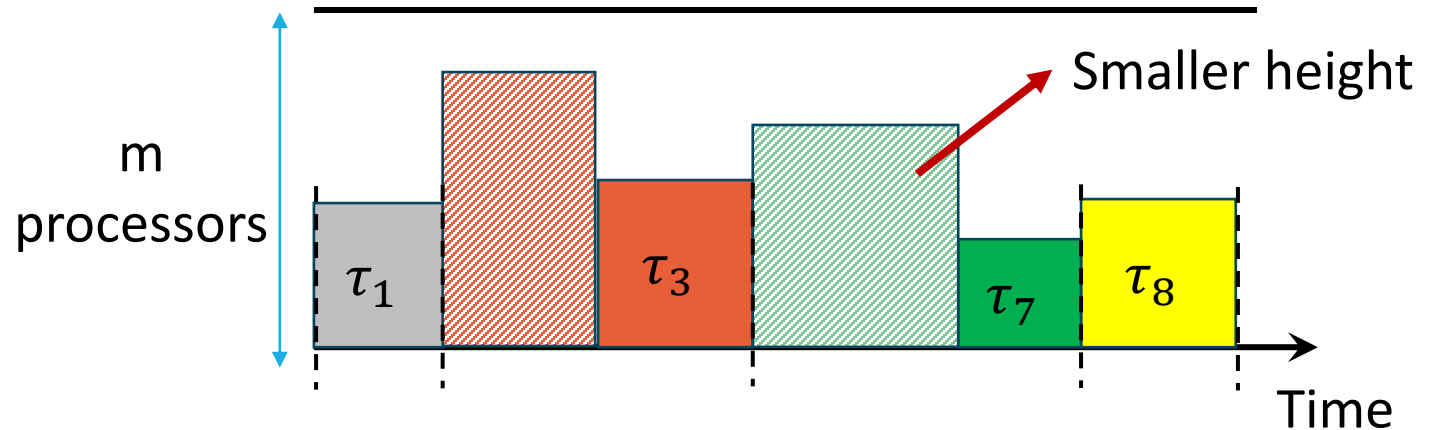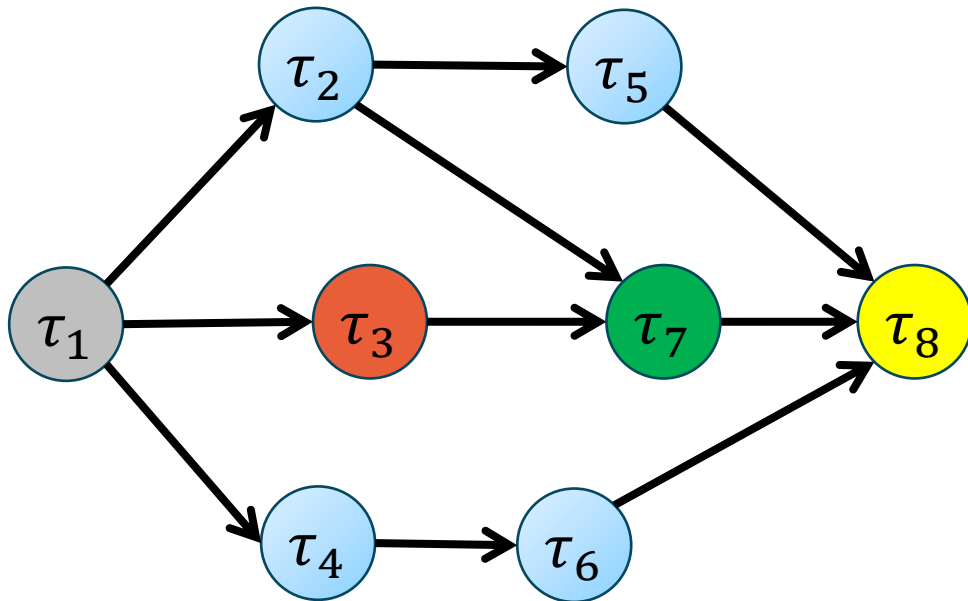
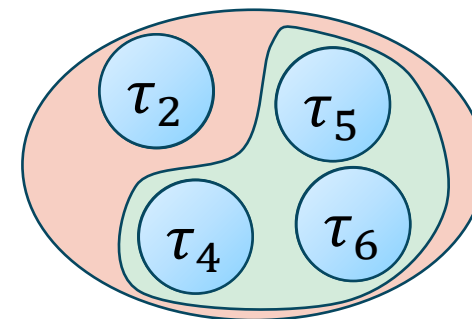# Response-Time Bound

Any path can be a critical path

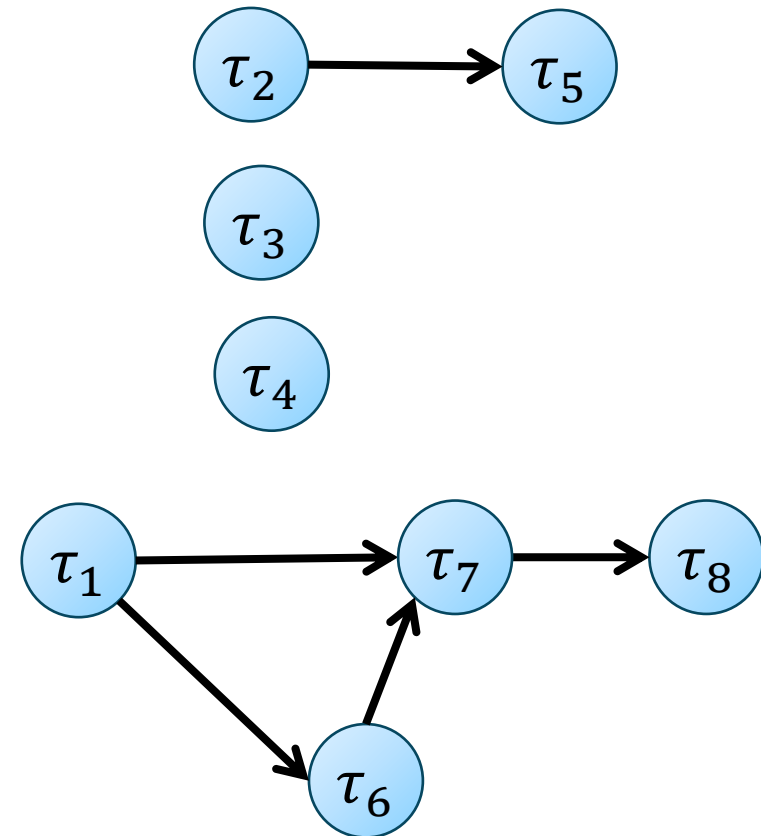Determine a set of nodes (not necessarily on a path) that upper bounds interference time



m processors

Smaller height

Time

Possible interfering tasks

# Response-Time Bound



Compute unit 2
$M_2$ processors

Compute unit 1
$M_1$ processors

Assumption 1: ~~One~~ compute ~~unit~~
Multiple        units

Scheduling in different compute units can be different

# Evaluation

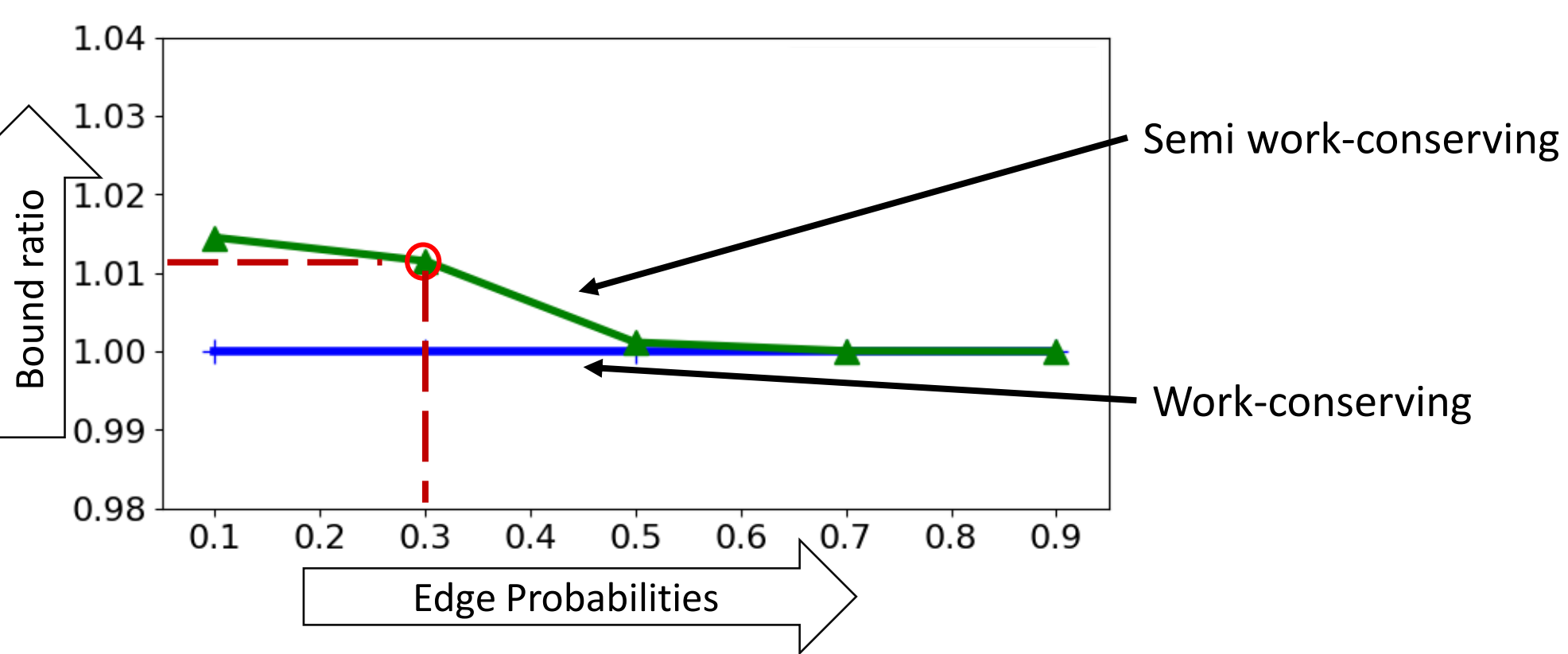$$\frac{\text{Bound under X}}{\text{Bound under work-conserving}}$$

Work-conserving vs. Semi work-conserving scheduling

≥ 1 means higher bound



Semi work-conserving

Work-conserving

Bound ratio

Edge Probabilities

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Evaluation

GPU as a shared resource vs. scheduling platform

GPU kernel

Block of threads

With locking

Without locking

CPU-only DAG response-time bound

DAG of gang tasks response-time bound

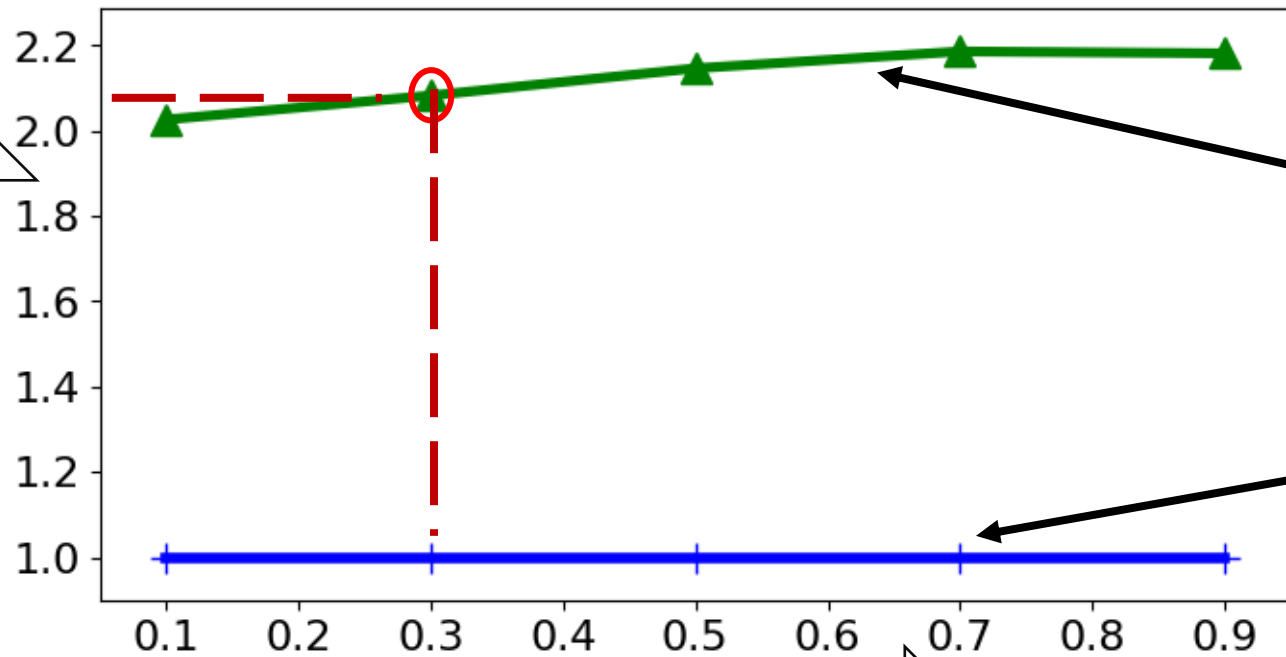# Evaluation

Bound under X
Bound under default scheduler

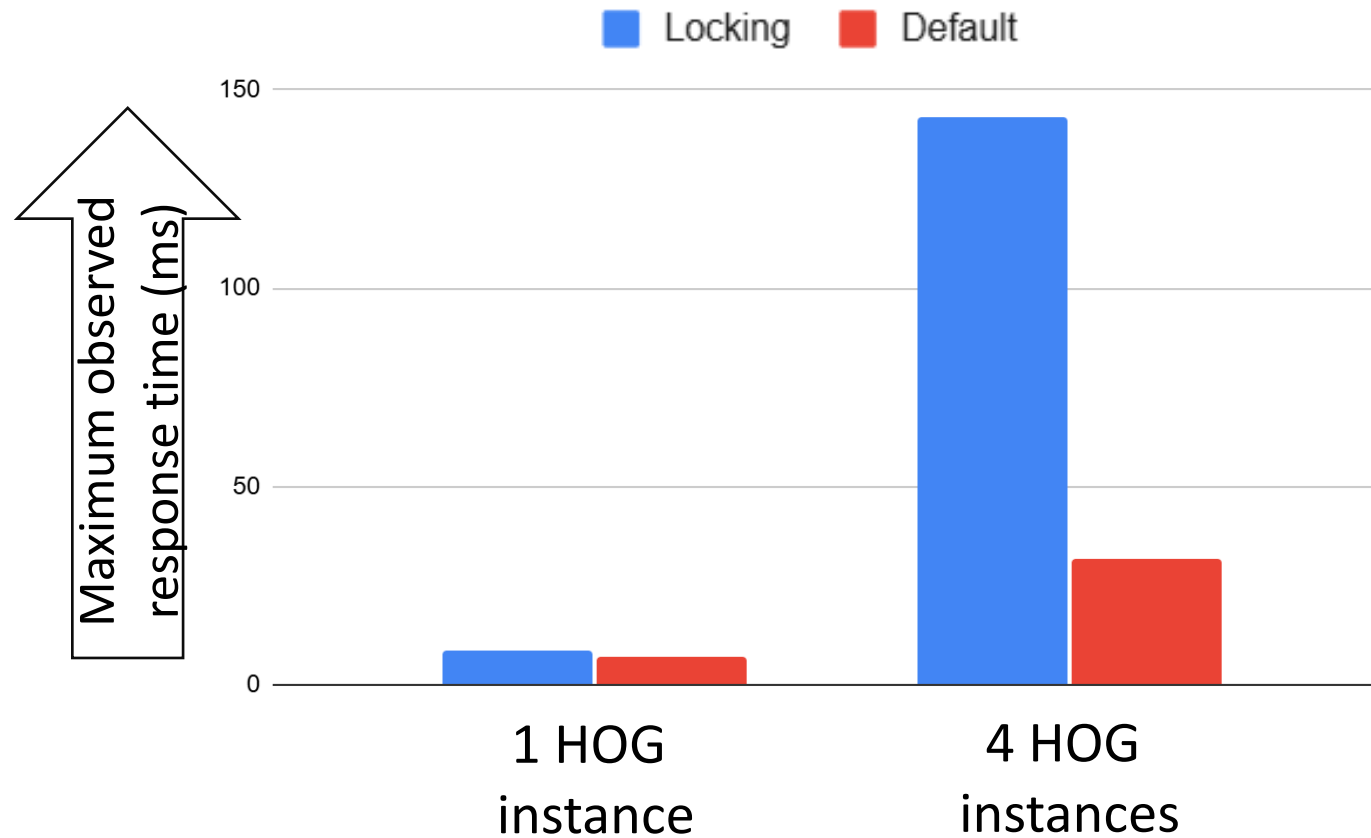Locking-based GPU access vs. Default GPU scheduling



≥ 1 means higher bound

Locking-based
[He et al., RTSS 2022]

Default

Bound ratio

Edge Probabilities

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Evaluation



Histogram of Oriented Gradients

GPU partitioning using `libsmctrl`
[Bakita and Anderson, RTAS 2023]

# Conclusion & Thank You!