# Partitioned Fixed-Priority Scheduling of Parallel Tasks Without Preemptions

**Daniel Casini**[*], Alessandro Biondi[*], Geoffrey Nelissen[†], and Giorgio Buttazzo[*]

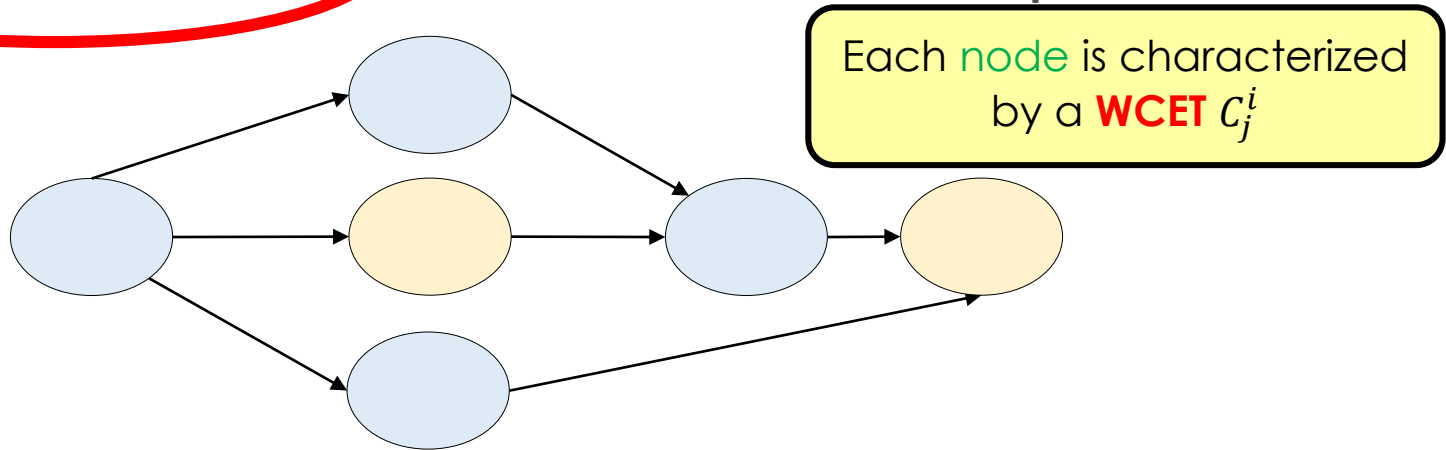[*] *ReTiS Lab, Scuola Superiore Sant'Anna, Pisa, Italy*

[†] *CISTER, ISEP, Polytechnic Institute of Porto, Portugal*

Retis
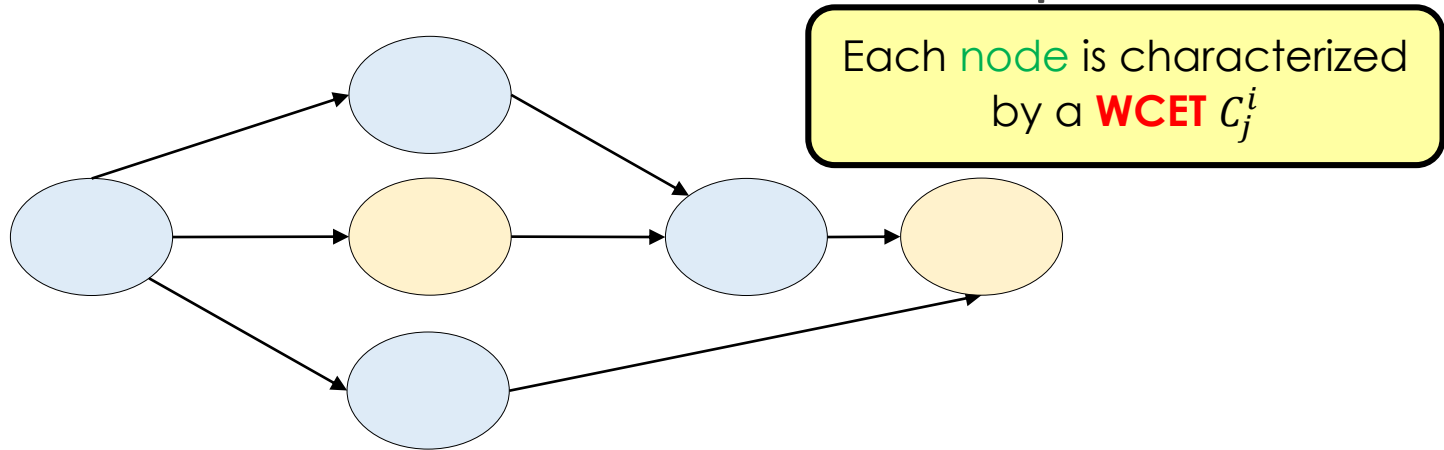Real-Time Systems Laboratory

CISTER
Research Centre in
Real-Time & Embedded
Computing Systems

Partitioned Fixed-Priority Scheduling of Parallel Tasks Without Preemptions

Each node is characterized by a **WCET** $C_j^i$

- Each task is represented by a **Direct Acyclic Graph**, and is characterized by

i. a minimum inter-arrival time $T_i$

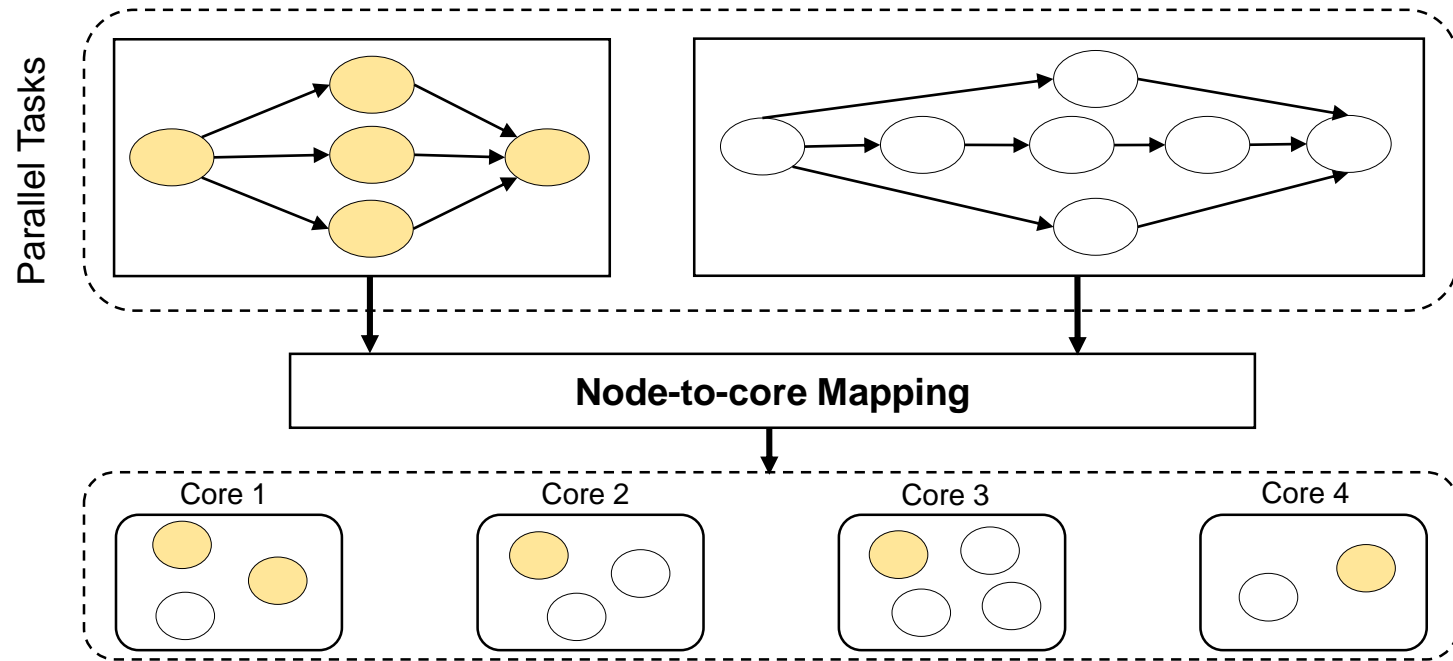ii. a constrained deadline $D_i \leq T_i$

iii. a **fixed priority** $\pi_i$

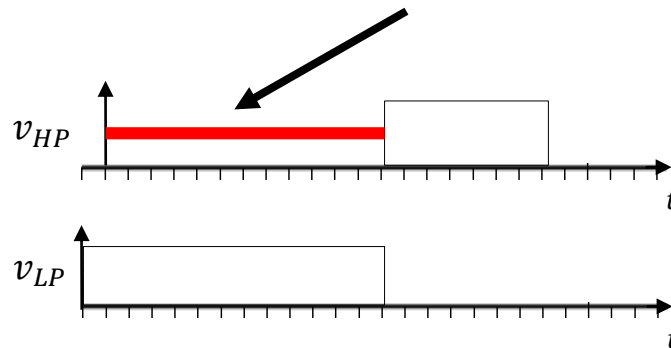# Partitioned Fixed-Priority Scheduling of Parallel Tasks Without Preemptions



Each node is characterized by a **WCET** $C_j^i$

- Each task is represented by a **Direct Acyclic Graph**, and is characterized by

    i.  a minimum inter-arrival time $T_i$

    ii. a constrained deadline $D_i \le T_i$

    iii. a **fixed priority** $\pi_i$

# Partitioned Fixed-Priority Scheduling of Parallel Tasks Without Preemptions



- Each node is statically assigned to a core

- Nodes of the same task can be allocated to different cores

Partitioned Fixed-Priority Scheduling of Parallel Tasks ~~Without Preemptions~~

**Without Preemptions**

**Non-preemptive blocking**

$v_{HP}$

$t$

$v_{LP}$

$t$

As soon a **node** starts executing, it cannot be preempted

# Why non-preemptive scheduling?

Predictable management of local memories

e.g., nodes can pre-load data from scratchpads before start executing

Memory Feasibility Analysis of Parallel Tasks
Running on Scratchpad-Based Architectures

*Daniel Casini, Alessandro Biondi, Geoffrey Nelissen and Giorgio Buttazzo*

*This morning @ RTSS*

Predictable management of local memories
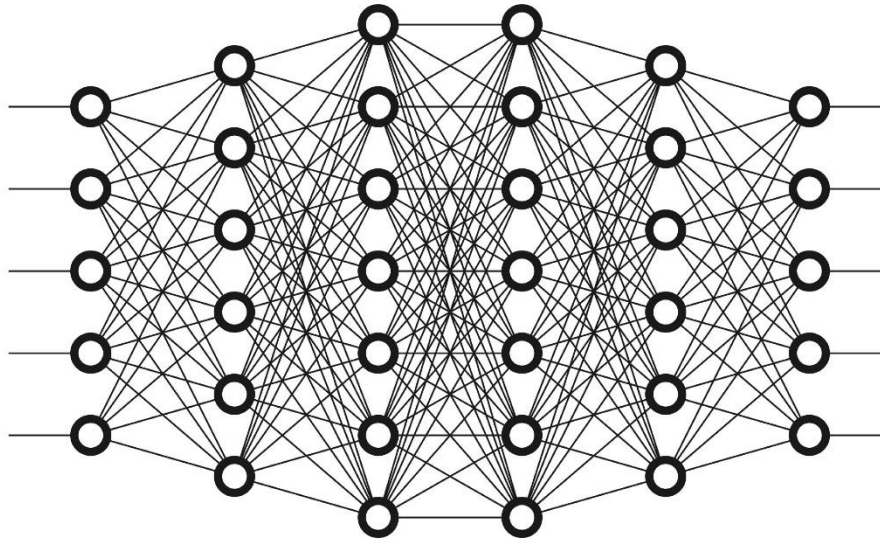
Reduces context-switch overhead

Simplifies WCET Analysis

Use of HW accelerators and GPUs

Can be a good choice for executing deep neural networks

- We profiled a deep neural network executed by Tensorflow a 8-core Intel i7 machine @ 3.5GHz



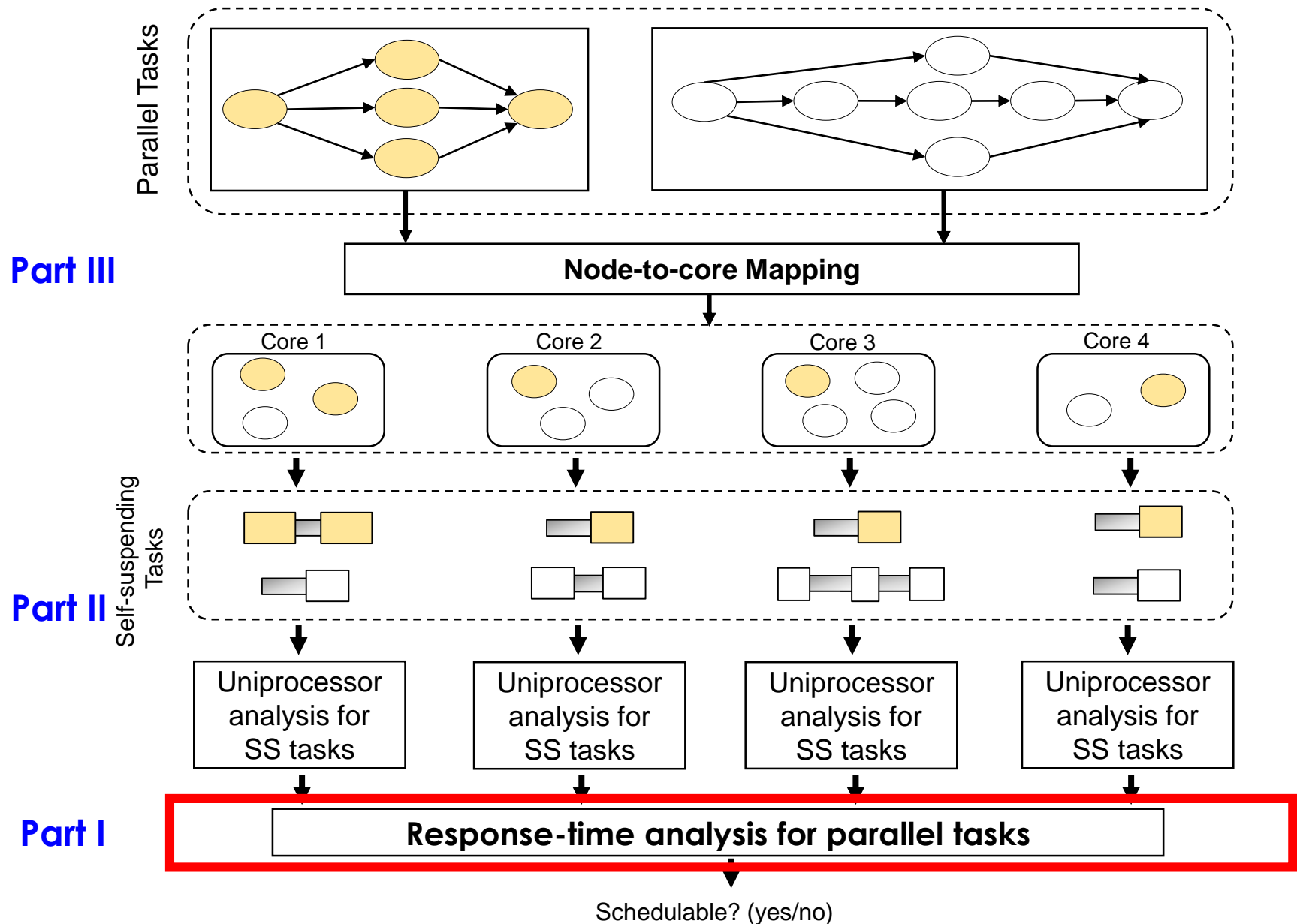**More than 34000 nodes** where only about **1.2%** of them have execution times larger than 100 microseconds

Parallel Tasks

**Part III**

Node-to-core Mapping

Core 1     Core 2     Core 3     Core 4

Self-suspending Tasks

**Part II**

| Uniprocessor analysis for SS tasks | Uniprocessor analysis for SS tasks | Uniprocessor analysis for SS tasks | Uniprocessor analysis for SS tasks |

**Part I**

**Response-time analysis for parallel tasks**

Schedulable? (yes/no)

# Overview of the analysis framework



**Parallel Tasks**

**Part III**

**Node-to-core Mapping**

Core 1     Core 2     Core 3     Core 4

**Self-suspending Tasks**

**Part II**

| Uniprocessor analysis for SS tasks | Uniprocessor analysis for SS tasks | Uniprocessor analysis for SS tasks | Uniprocessor analysis for SS tasks |

**Part I**     **Response-time analysis for parallel tasks**

Schedulable? (yes/no)

# Part I:
## Response-time analysis for parallel tasks

- Each core 'perceives' the execution of a parallel task as an interleaved sequence of execution and suspension regions



Suspension regions correspond to
execution regions on a **different core**

# Intuition

- Paths can be mapped to a self-suspending tasks

# Intuition

- Paths can be mapped to a self-suspending tasks

Path



$p_1$

$p_2$

execution region

suspension region

Segmented SS-tasks

$C_1^{SS} = C_1$      $C_2^{SS} = C_5$      $C_3^{SS} = C_7$

The length of each execution region directly maps to the WCET of a node in the graph

# Intuition

Segmented SS-tasks

$$S_1^{SS} = R_3 \qquad S_2^{SS} = R_6$$

The length of each suspension region depends on the
**response time** of nodes allocated to **different cores**

⚠️ Complex **inter-core dependencies** can arise

$p_1$: $v_1$ | $v_2$ | $v_4$ | | $v_5$ | | $v_7$

$p_2$: | $v_3$ | | $v_6$

- Recursive algorithm to unfold response-time dependencies:



$p_1$    $p_2$    $p_3$    $p_4$

$v_1 \rightarrow v_2 \rightarrow v_3 \rightarrow v_4 \rightarrow v_5 \rightarrow v_6 \rightarrow v_7$

$S_1^{SS1} = R^{SS2}$

$R^{SS1}$

$S_1^{SS2} = R^{SS3}$

$R^{SS2}$

$S_1^{SS3} = R^{SS4}$

$R^{SS3}$

$R^{SS4}$

- We extended this approach to work under non-preemptive scheduling

> Need for a fine-grained analysis for non-preemptive self-suspending tasks

**Our next contribution**

# Part II:

Analysis for non-preemptive self-suspending tasks

**Two different approaches:**

| 1 |

### Holistic analysis

- Computes the RT of a whole self-suspending task

$$C_i = \sum_{\text{all segments}} C_{i,j} \qquad S_i = \sum_{\text{all segments}} S_{i,j}$$

per-task response time bound

- Analytically dominates state-of-the-art analysis (Dong et al. 2018)

| 2 |

### Segment-based analysis

- Computes the RT of individual segments

per-segment response time bounds

**Two different approaches:**

1 **Holistic analysis**

- Computes the RT of a whole self-suspending task

$$c_i = \sum c$$

- A ... -art analysis (Dong et al. 2018)

2 **Segment-based analysis**

- Computes the RT of individual segments

per-segment response time bounds

*Hybrid model: Take the minimum of the two bounds*

**Interference** due to <u>higher-priority</u> segments

↓

**Response-time analysis**

↑

Non-preemptive **blocking** due to <u>lower-priority</u> segments

- Interference from higher-priority tasks is accounted by means of the following worst-case scenario*:



- The response time bound can be initially approximated to the task's deadline and iteratively refined

  - **Holistic** and **segmented** analyses are combined during the iterative refinement

*Jian-Jia Chen et al., "Many suspensions, many problems: a review of self-suspending tasks in real-time systems", Real-time System Journal.

# Fine-grained accounting of blocking

- With a multiset approach

Task release (arrow)  |  Segment release (dashed arrow)

$\Delta$

$\tau_{HP}$

$\tau_{LP}$

$t$

$t$

Two segmented SS-tasks contending for the same CPU

multiset

Contains the WCET of all the lower-priority segments that may block the task under analysis in a window of length $\Delta$

- With a multiset approach

| | |
|---|---|
| ↑ Task release | ⇡ Segment release |

$\Delta$

$\tau_{HP}$

$\tau_{LP}$

Two segmented SS-tasks contending for the same CPU

multiset

Contains the WCET of all the lower-priority segments that may block the task under analysis in a window of length $\Delta$

# Now we have our analysis!

**MISSION COMPLETE**

# Now we have our analysis!

**MISSION COMPLETE**

# Now we have our analysis!

$R'$

$R''$

$R'''$

$$R = \max(R', R'', R''')$$
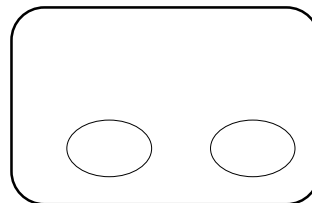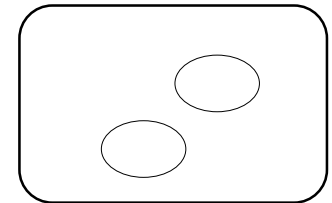
The RT of a parallel task can be derived from the maximum RT of all its paths

# Partitioning (meta-)algorithm

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph
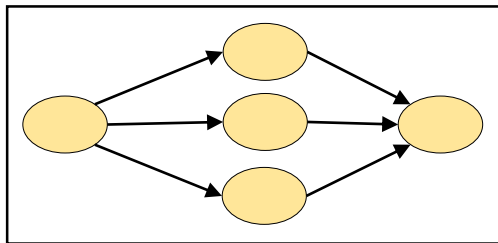
Inputs:

1. Strategy for ordering tasks

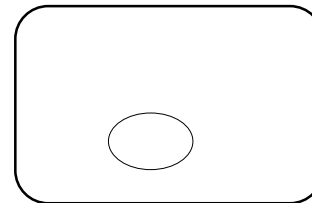2. Strategy for ordering cores

Output:

1. Node partitioning

Example:

# Partitioning (meta-)algorithm

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:
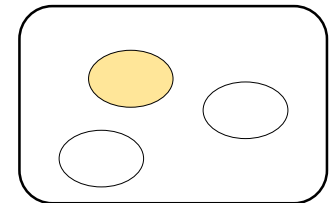
1. Node partitioning

Example:

Task under analysis

Task under analysis (during partitioning)

Core 1
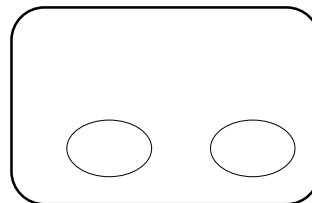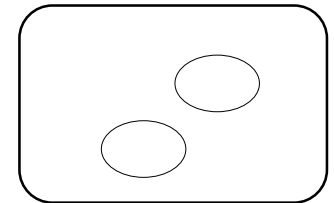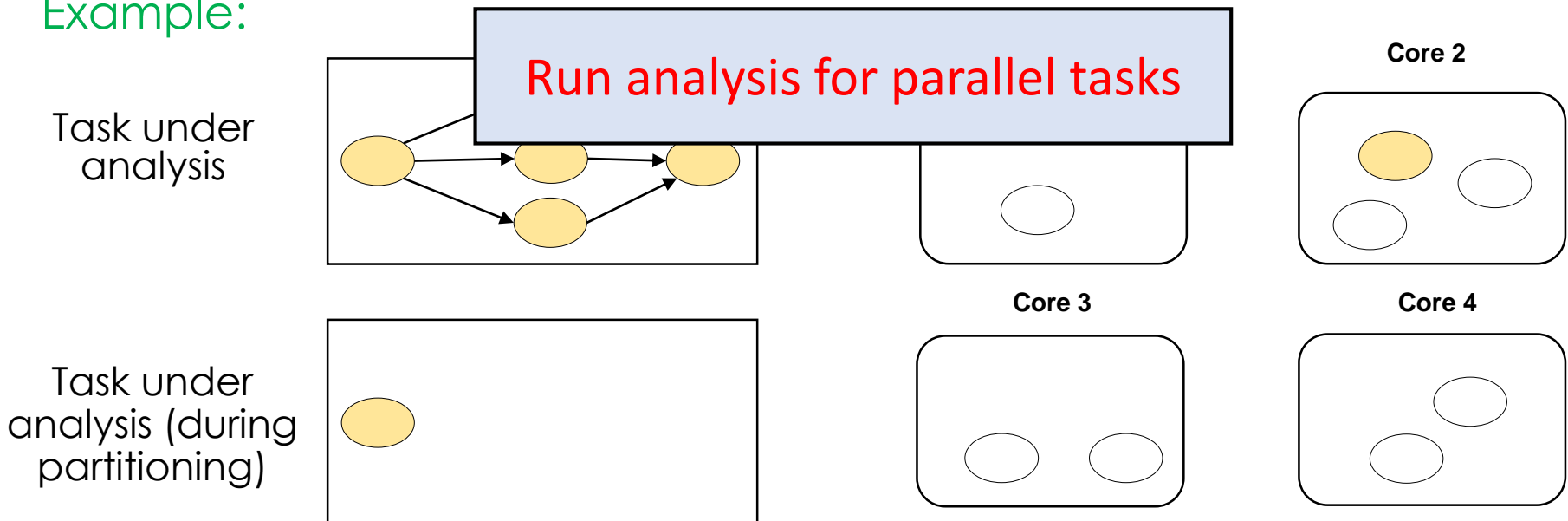
Core 2

Core 3

Core 4

# Partitioning (meta-)algorithm

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph
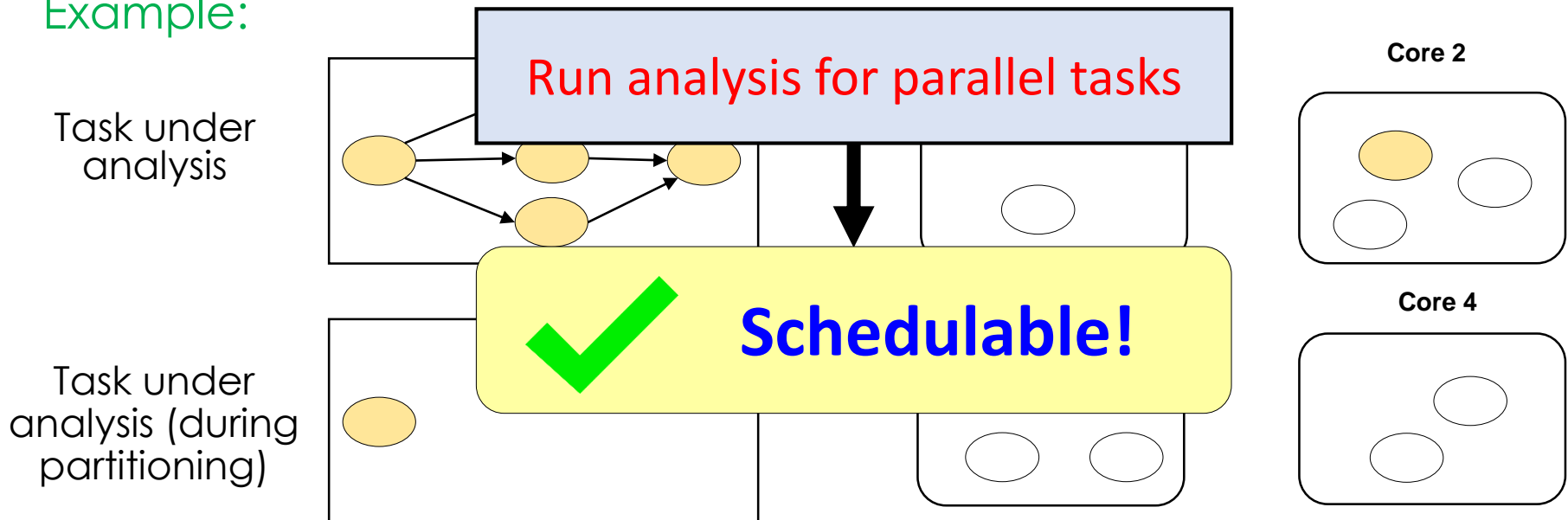
Inputs:

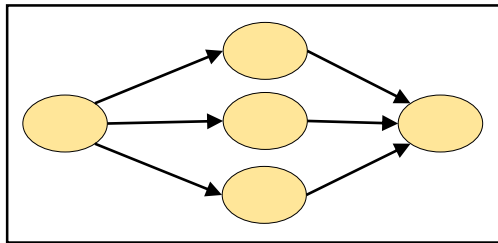1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

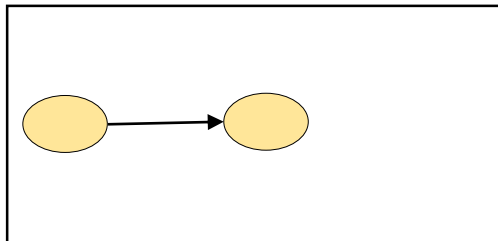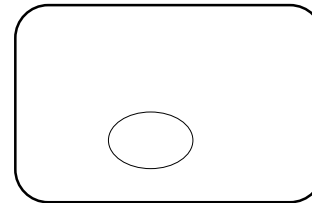1. Node partitioning

Example:



Task under analysis

Task under analysis (during partitioning)

Core 1

Core 2

Core 3

Core 4

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

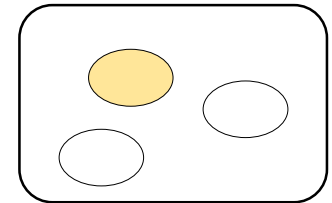Output:

1. Node partitioning

Example:

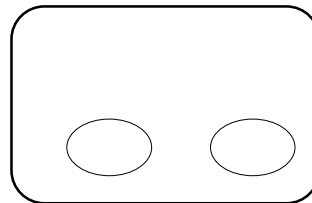Task under analysis

Task under analysis (during partitioning)

Run analysis for parallel tasks

Core 2

Core 3

Core 4

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:

Task under analysis

Task under analysis (during partitioning)

Run analysis for parallel tasks

Core 2

Core 4

✔ **Schedulable!**

# Partitioning (meta-)algorithm

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:

Task under analysis
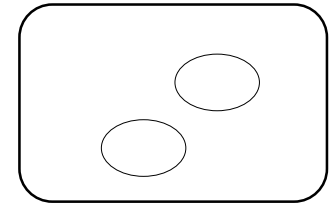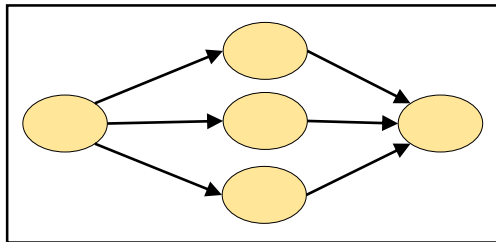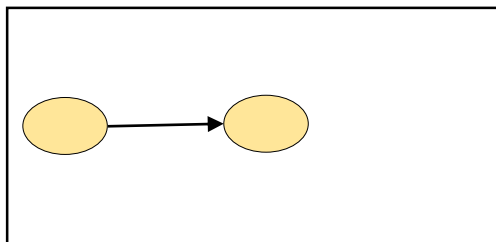
Task under analysis (during partitioning)

**Core 1**

**Core 2**

**Core 3**

**Core 4**

# Partitioning (meta-)algorithm

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:
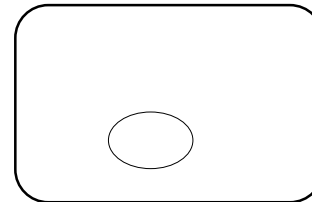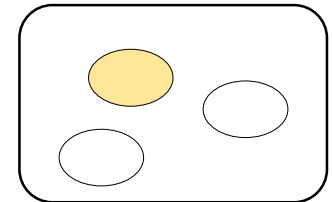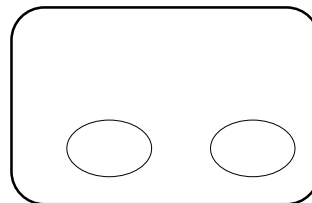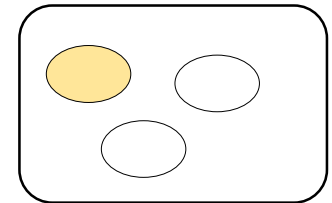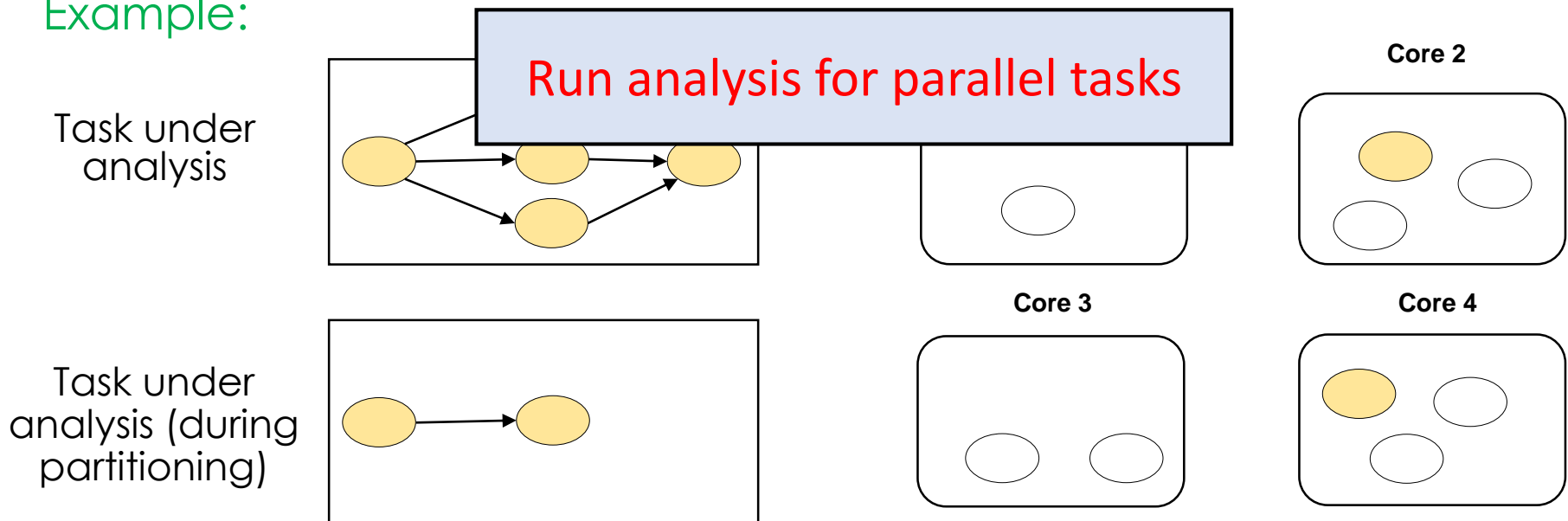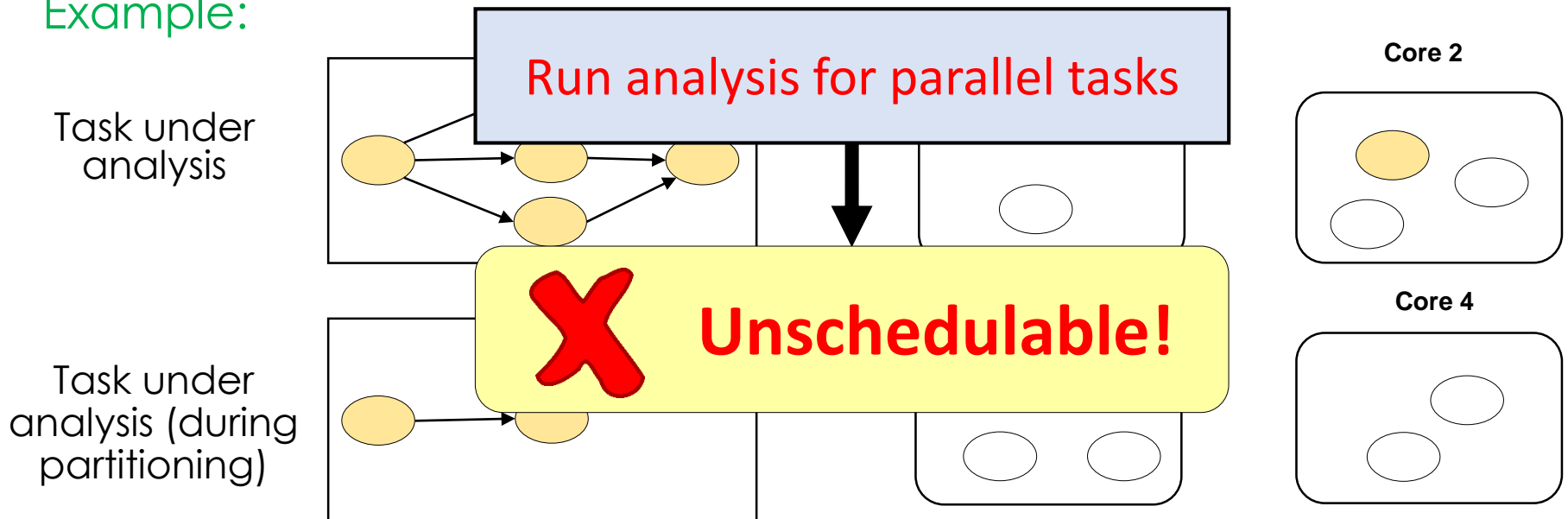
1. Node partitioning

Example:

Task under analysis

Task under analysis (during partitioning)

Core 1

Core 2

Core 3

Core 4

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:
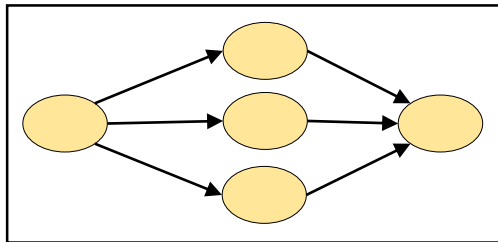
1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:

Task under analysis

Run analysis for parallel tasks

Core 2

Task under analysis (during partitioning)

Core 3

Core 4

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

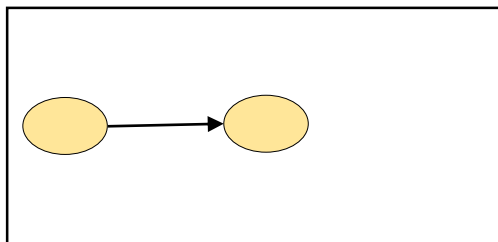1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:

Task under analysis

Task under analysis (during partitioning)

Run analysis for parallel tasks

❌ **Unschedulable!**

Core 2

Core 4

**IDEA:** Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

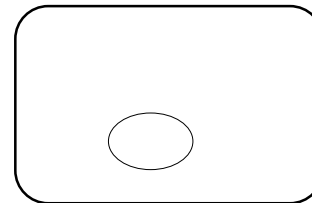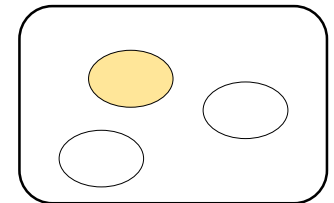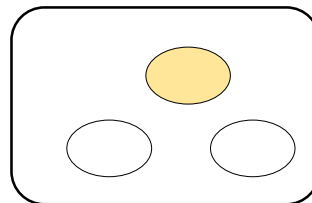1. Node partitioning

Example:



Task under analysis
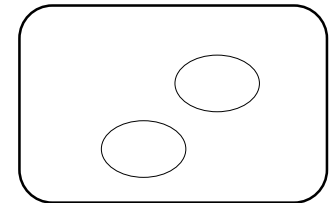
Task under analysis (during partitioning)

Core 1

Core 2

Core 3

Core 4

# Partitioning (meta-)algorithm

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:
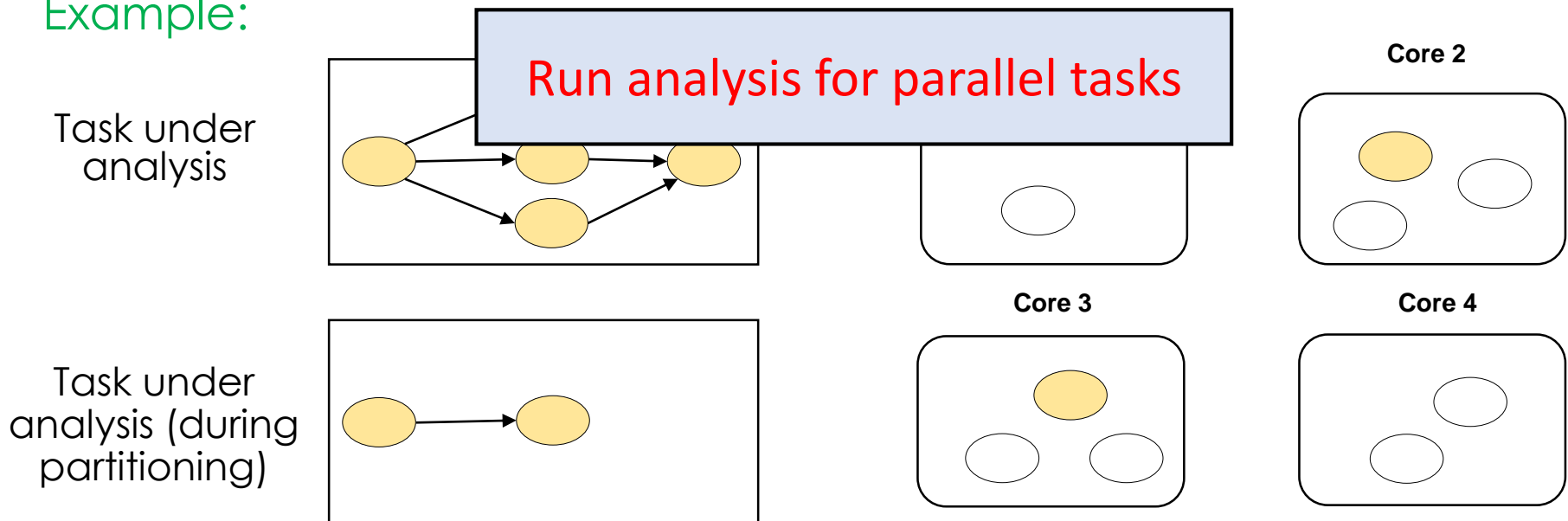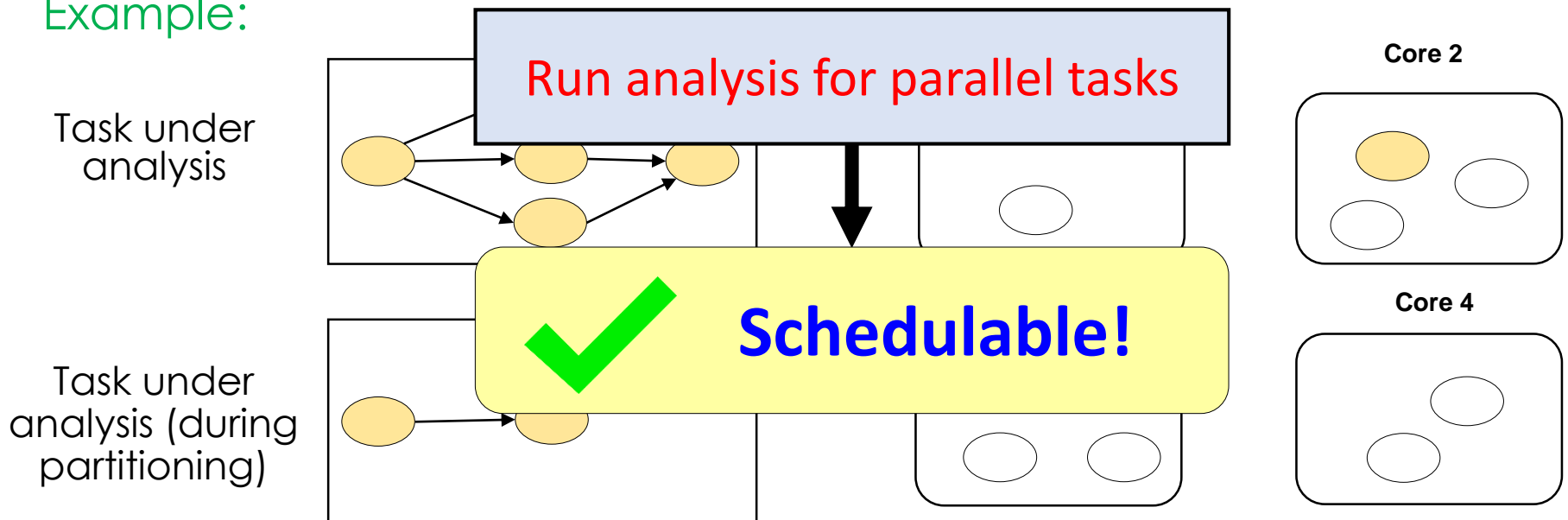
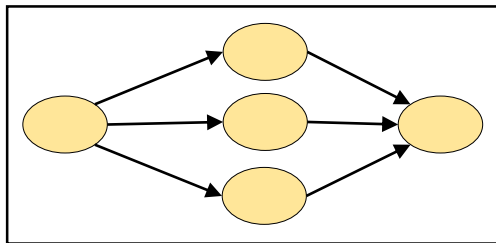1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:



Run analysis for parallel tasks

Task under analysis

Task under analysis (during partitioning)

Core 2

Core 3

Core 4

**IDEA:** Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

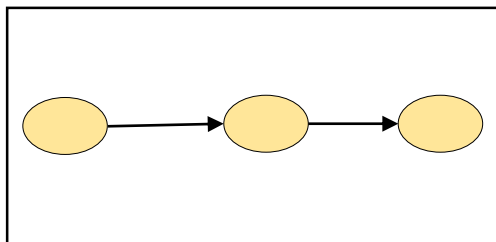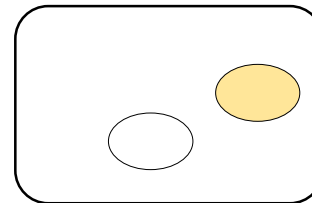1. Strategy for ordering tasks

2. Strategy for ordering cores

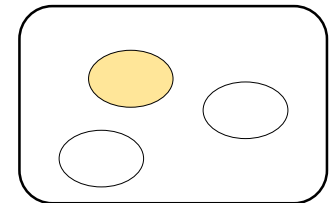Output:

1. Node partitioning

Example:

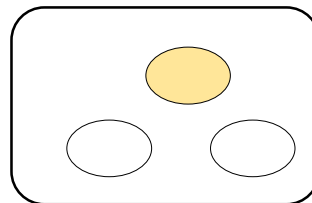Task under analysis

Task under analysis (during partitioning)

Run analysis for parallel tasks

✔ **Schedulable!**

**Core 2**

**Core 4**

# Partitioning (meta-)algorithm

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:

**IDEA:** Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:



Task under analysis
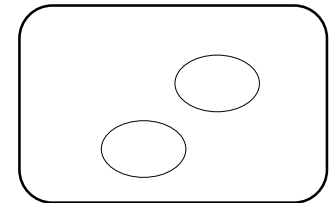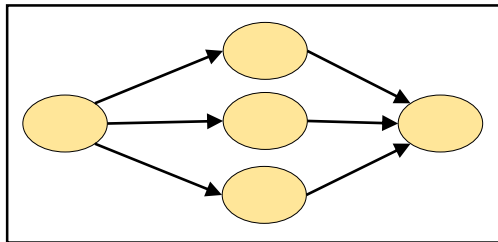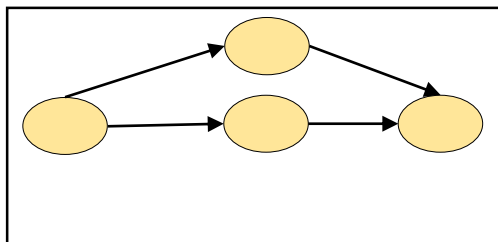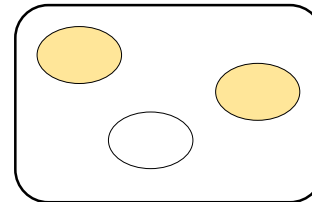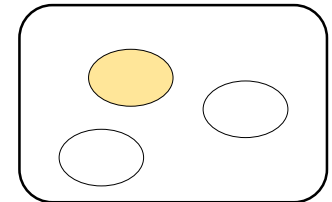
Task under analysis (during partitioning)
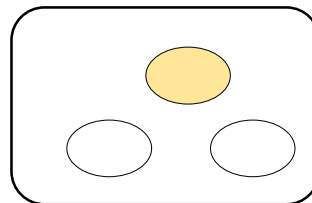
Core 1

Core 2

Core 3

Core 4

# Partitioning (meta-)algorithm

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:

IDEA: Analyzing schedulability incrementally, adding one node at a time, and perform schedulability analysis on a subgraph

Inputs:

1. Strategy for ordering tasks

2. Strategy for ordering cores

Output:

1. Node partitioning

Example:



Task under analysis
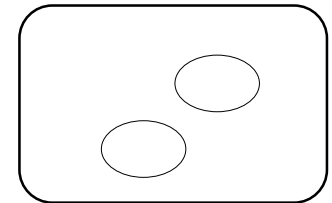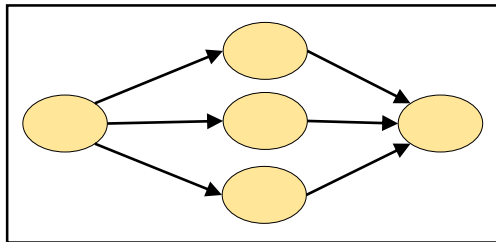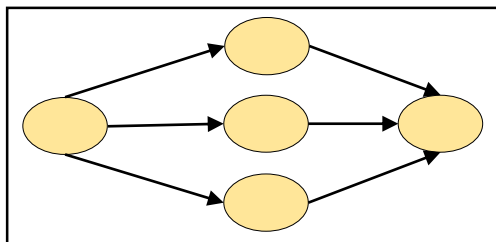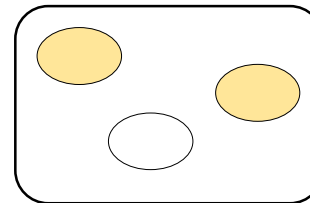
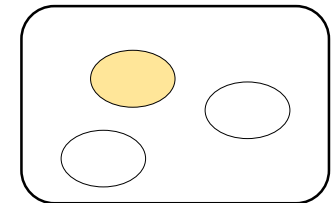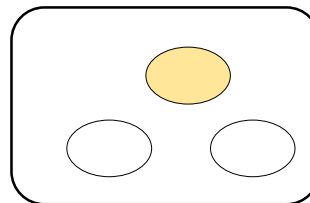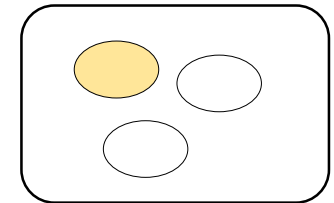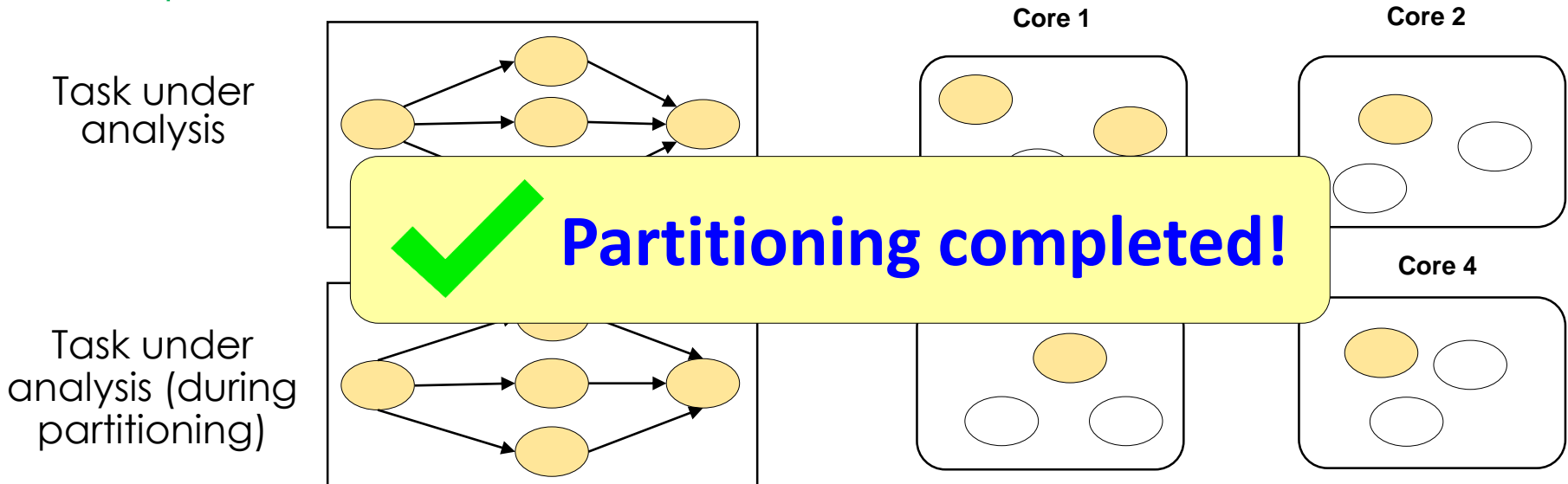Task under analysis (during partitioning)

Core 1

Core 2

Core 4

✅ **Partitioning completed!**

# Experimental Results

# Experimental Study

- Experimental study based on synthetic workload

    - We compared against the only previous work targeting non-preemptive scheduling of parallel tasks, which targets global scheduling (Serrano et al. 2017)

    - Same DAG generator used in [Serrano et al. 2017]

        - WCETs randomly generated in (0,100] with uniform distribution

        - Tasks utilizations obtained with U-Unifast

        - Tasks periods computed as $T_i = U_i \sum_{nodes} C_{i,j}$

12 tasks, 16 processors

The higher the better

## 12 tasks, 16 processors



Increasing task-set utilization

## 12 tasks, 16 processors



Improvement up to 100 percentage points over [Serrano et al. 2017]

Legend: GLOBAL, WF_UTIL, FF_ALGO, BF_ALGO, WF_ALGO, PARTITIONED

Our experimental study revelead a similar trend varying the number of tasks and processor, e.g.,

**10 tasks, 8 processors**

# Conclusions

**1**    **Methodology** for analyzing non-preemptive parallel tasks as a set of self-suspending tasks
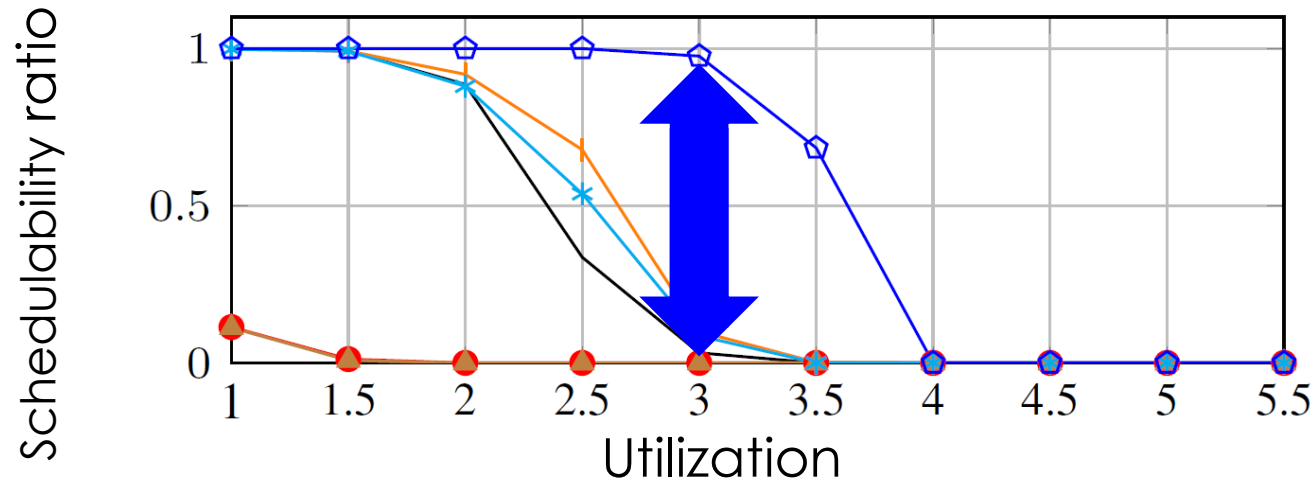
**2**    **Analysis** for non-preemptive self-suspending tasks which analytically dominates the only previous result

**3**    **Partitioning algorithm** to allocate nodes to the available processors
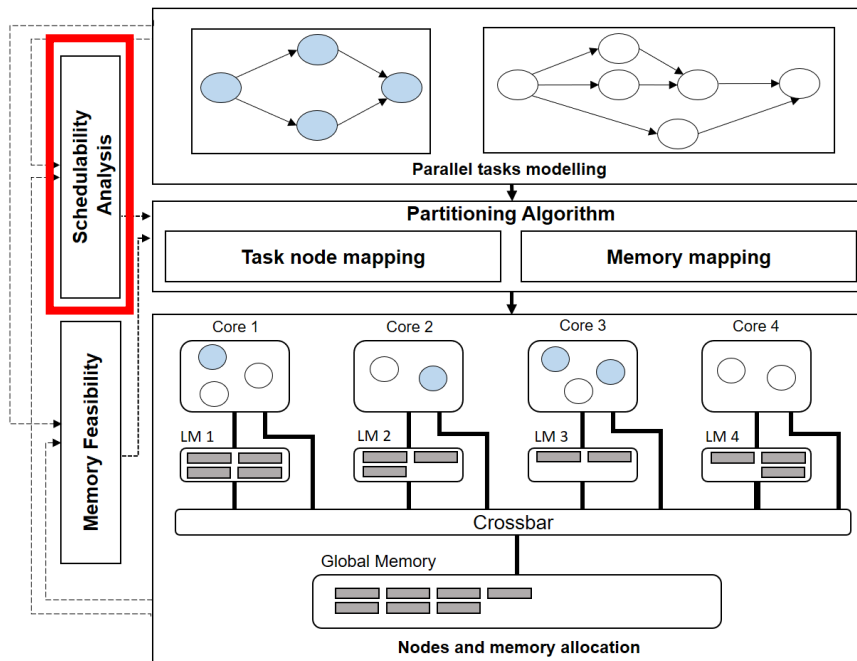
**4**    **Experimental study** to assess the improvement in terms of schedulability – **up to 100 p.p.** w.r.t. the only existing previous work for global scheduling

Deeper investigation of partitioning strategies

Improvement in the analysis precision

Integration of **communication delays** in the analysis



Parallel tasks modelling

Partitioning Algorithm

Task node mapping | Memory mapping

Core 1 | Core 2 | Core 3 | Core 4

LM 1 | LM 2 | LM 3 | LM 4

Crossbar

Global Memory

Nodes and memory allocation

Schedulability Analysis

Memory Feasibility

Memory Feasibility Analysis of Parallel Tasks Running on Scratchpad-Based Architectures

*Daniel Casini, Alessandro Biondi, Geoffrey Nelissen and Giorgio Buttazzo*

*This morning @ **RTSS***

# Thank you!

Daniel Casini
daniel.casini@sssup.it