# DWS: Demand-aware Work-Stealing in Multi-programmed Multi-core Architectures

Quan Chen, Long Zheng, Minyi Guo

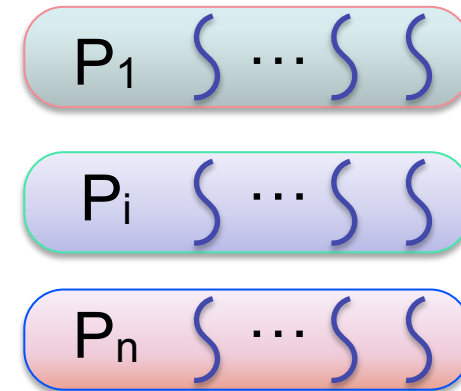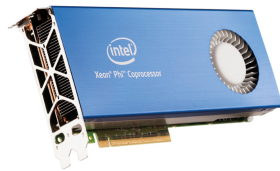Shanghai Jiao Tong University, China

1

# Outline

- Background

- Problem & Motivation

- Demand-aware Work-Stealing (DWS)

- Evaluation

- Conclusions

# Background

- Hardware: Multi-core/Many-core Architectures

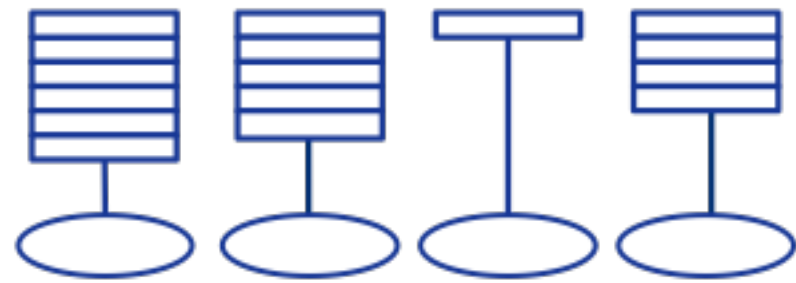- Scenario: Multiple parallel programs

$P_1$ … 

$P_i$ …

$P_n$ …

# Background-parallel programs

- ## Traditional parallel programs
  - *Hard* to adjust the number of threads at runtime
- ## Task-based parallel programs
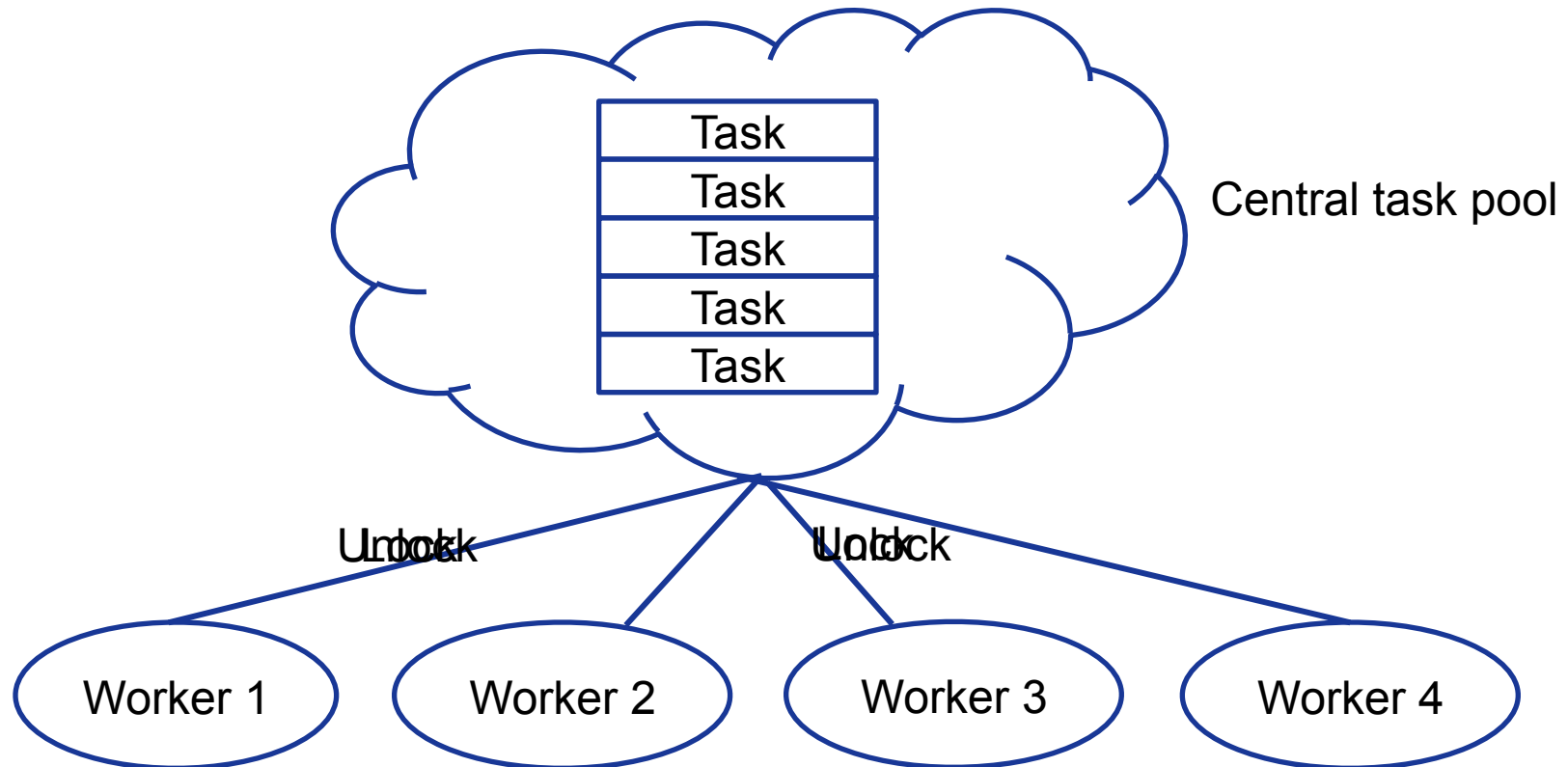  - Dynamic task scheduling

Work-sharing

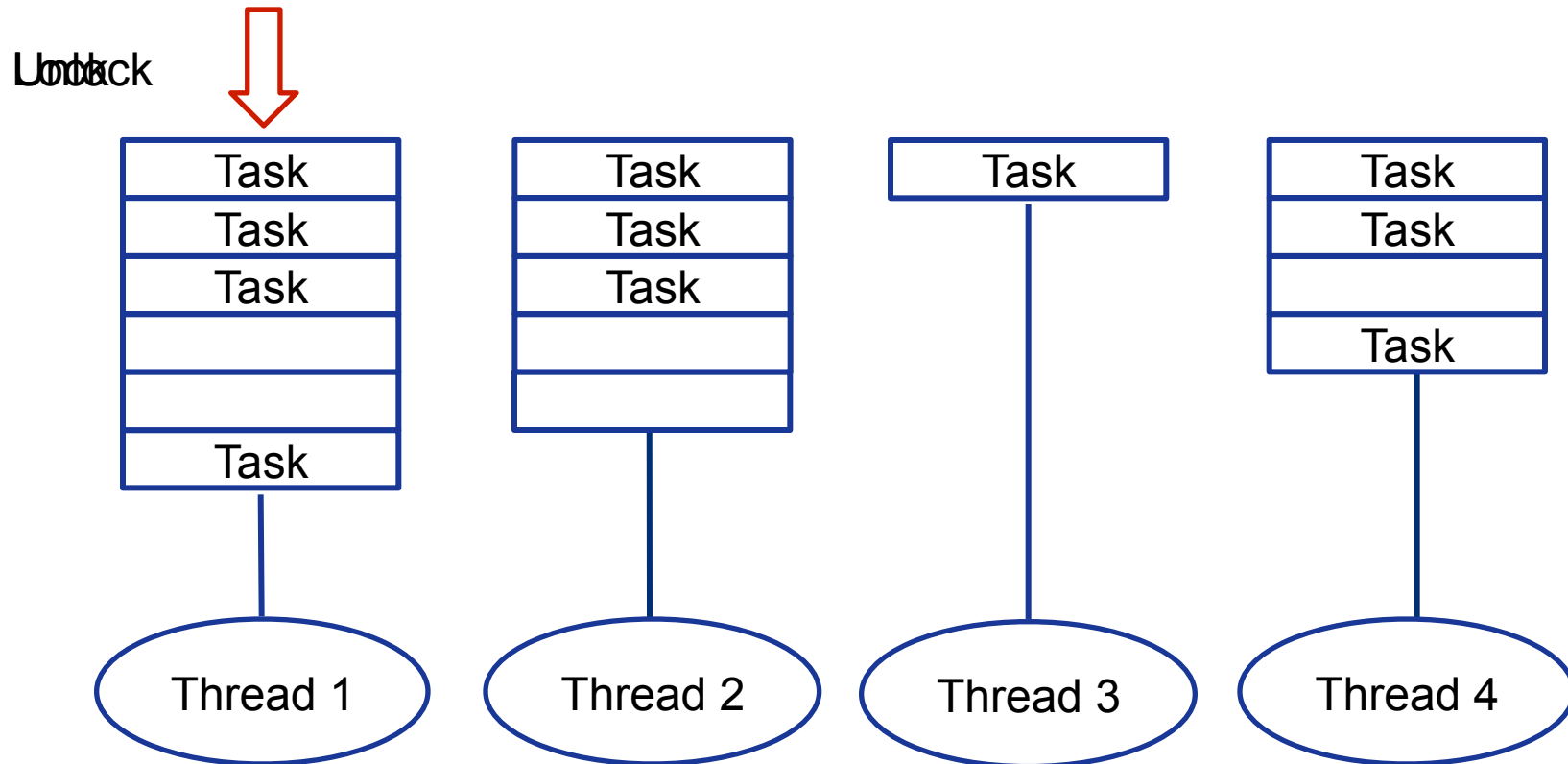Work-stealing

**Easy to adjust the number of workers**

# Work-sharing

Task
Task
Task
Task
Task

Central task pool

Unlock Lock

Unlock Lock

Worker 1  Worker 2  Worker 3  Worker 4

Lock the central task pool when getting a task

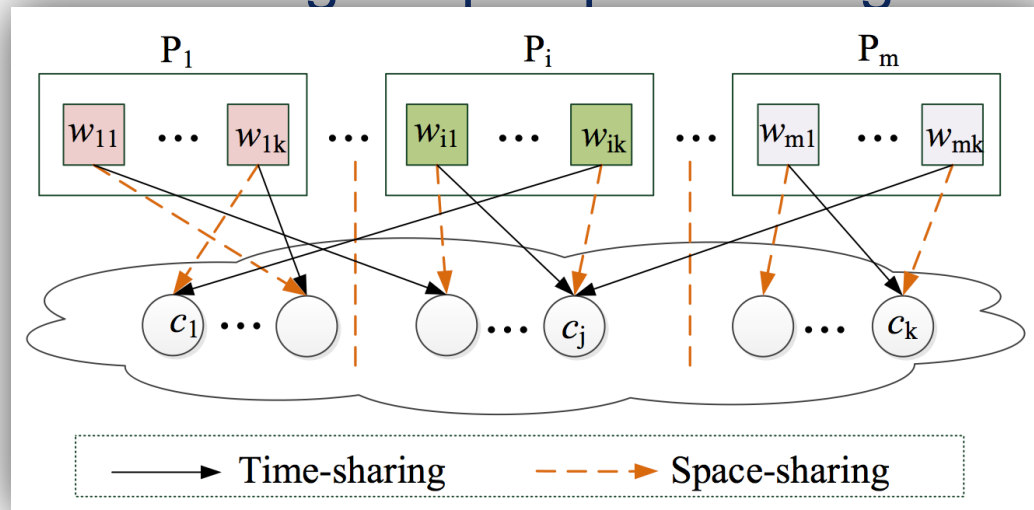# Work-stealing

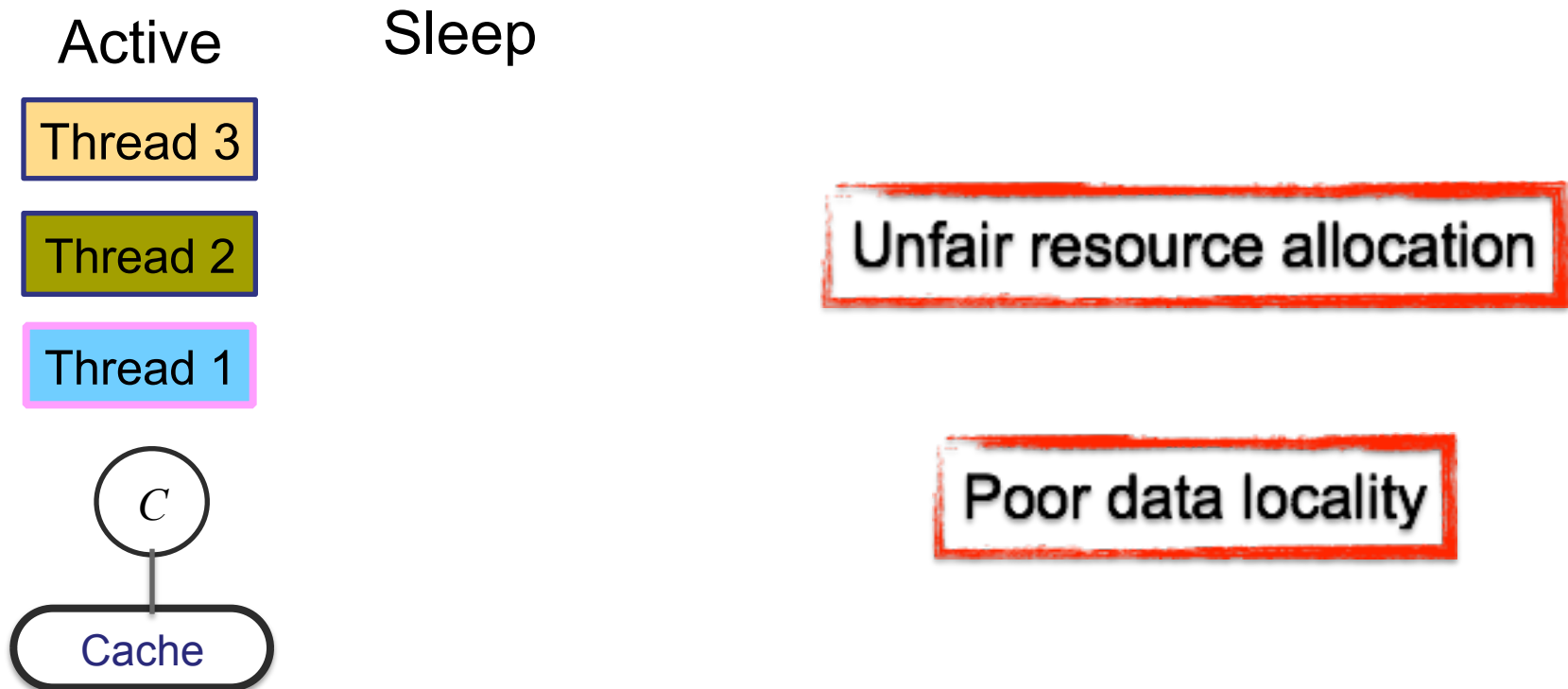- # Aggressive feature of work-stealing
  - On a $k$-core computer, $k$ threads/workers are launched

- # Existing solutions
  - Time-sharing - ABP yielding mechanism
  - Space-sharing - Equal-partitioning

⊛ **ABP yielding mechanism**

- • If a thread fails to steal a task, it goes to sleep

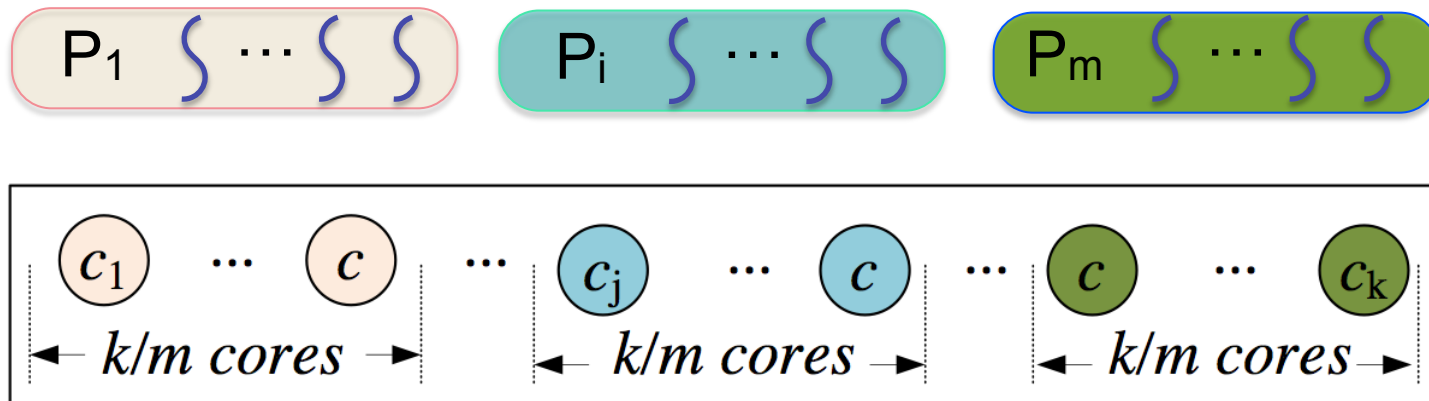Active  Sleep

Thread 3

Thread 2

Thread 1

$C$

Cache

Unfair resource allocation

Poor data locality

- Equal-partitioning mechanism
  - If $m$ programs co-run on a $k$-core computer, each program is allocated $k/m$ cores.

Fair but inefficient
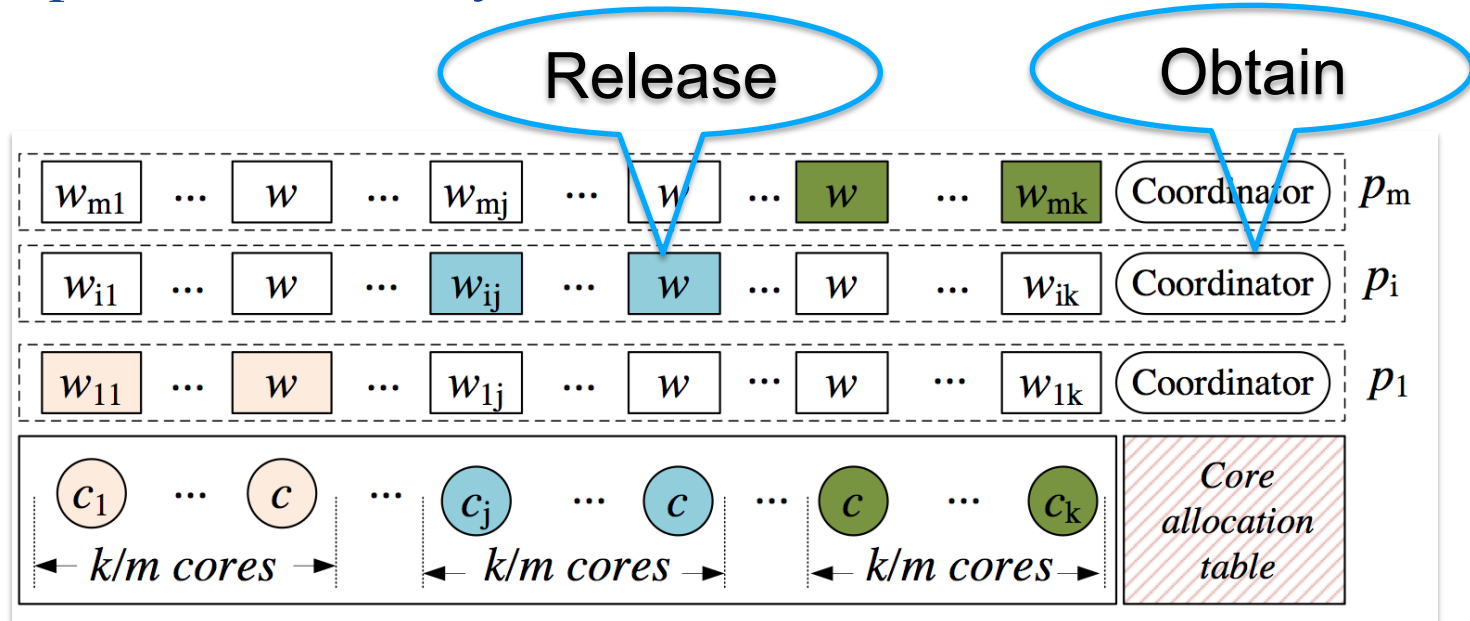
- **Start from Equal-partitioning**

- **Dynamically balance cores at runtime**
  - If $p_i$ cannot fully-utilized a core, it release the core
  - If $p_i$ has too many tasks, it tries to obtain more cores

Release

Obtain

| $w_{m1}$ | ... | $w$ | ... | $w_{mj}$ | ... | $w$ | ... | $w$ | ... | $w_{mk}$ | Coordinator | $p_m$ |
| $w_{i1}$ | ... | $w$ | ... | $w_{ij}$ | ... | $w$ | ... | $w$ | ... | $w_{ik}$ | Coordinator | $p_i$ |
| $w_{11}$ | ... | $w$ | ... | $w_{1j}$ | ... | $w$ | ... | $w$ | ... | $w_{1k}$ | Coordinator | $p_1$ |

$c_1$ ... $c$ ... $c_j$ ... $c$ ... $c$ ... $c_k$ | Core allocation table

← $k/m$ cores → ← $k/m$ cores → ← $k/m$ cores →

Runtime Arch. of DWS

- A worker decides whether to release its core by itself

**Algorithm 1:** Work-stealing algorithm in DWS

**Input:** $w$: current worker

```
 1  int failed_steals = 0;          // num of failed steals
 2  while work is not done do
 3      if w is free then
 4          if its task pool is not empty then
 5              w obtains a task t from its own task pool ;
 6              failed_steals = 0 ;
 7          else
 8              w randomly selects v as victim worker ;
 9              if v has a non-empty task pool then
10                  w steals t from v ;
11                  failed_steals = 0 ;
12              else
13                  failed_steals ++ ;
14                  if failed_steals > T_SLEEP then
15                      w goes to sleep ;
16                      w waits to be woken up ;
17                  end if
18              end if
19          end if
20          if t then
21              w executes t ;
22          end if
23      end if
24  end while
```

If a worker fails too many times (T_SLEEP) to steal a new task, it goes to sleep

● The coordinator decides whether to obtain more cores

- If a program has too many queued tasks, it should try to get some free cores

How Many?

Which?

C1: The more queued tasks in a program, the more cores should the program obtain

C2: A program can take its allocated cores back

C3: A program cannot obtain the busy cores

- C1: The more queued tasks in a program, the more cores should the program obtain

How many: $$N_w = \frac{N_b}{N_a}$$

| | |
|---|---|
| *Num of active workers* | $N_a$ |
| *Num of queued tasks* | $N_b$ |
| *Num of free cores* | $N_f$ |
| *Num of released cores* | $N_r$ |
| *Num of cores expected* | $N_w$ |

- $N_w <= N_f$
  - Randomly select $N_w$ free cores

- $N_f < N_w <= N_f + N_r$      **(C2)**
  - Select $N_f$ free cores + its $(N_w - N_f)$ released core

- $N_w > N_f + N_r$      **(C3)**
  - $N_f$ free cores+its $N_r$ released cores

| | |
|---|---|
| *Num of active workers* | $N_a$ |
| *Num of queued tasks* | $N_b$ |
| *Num of free cores* | $N_f$ |
| *Num of released cores* | $N_r$ |
| *Num of cores expected* | $N_w$ |

# Evaluation platform

- A Dual-socket Quad-core computer with Hyper-Threading Technology

- Each socket is a Quad-Core Intel Xeon E5620
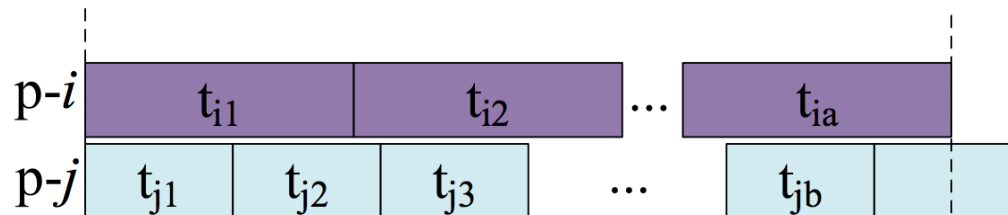
| Hardware & Configuration | Size/Version |
|---|---|
| *L1/L2 cache size (each core)* | *256 KB/1MB* |
| *L3 cache size  (each socket)* | *12 MB* |
| *Main memory size* | *32 GB* |
| *Operation system* | *Linux 2.6.32-38* |

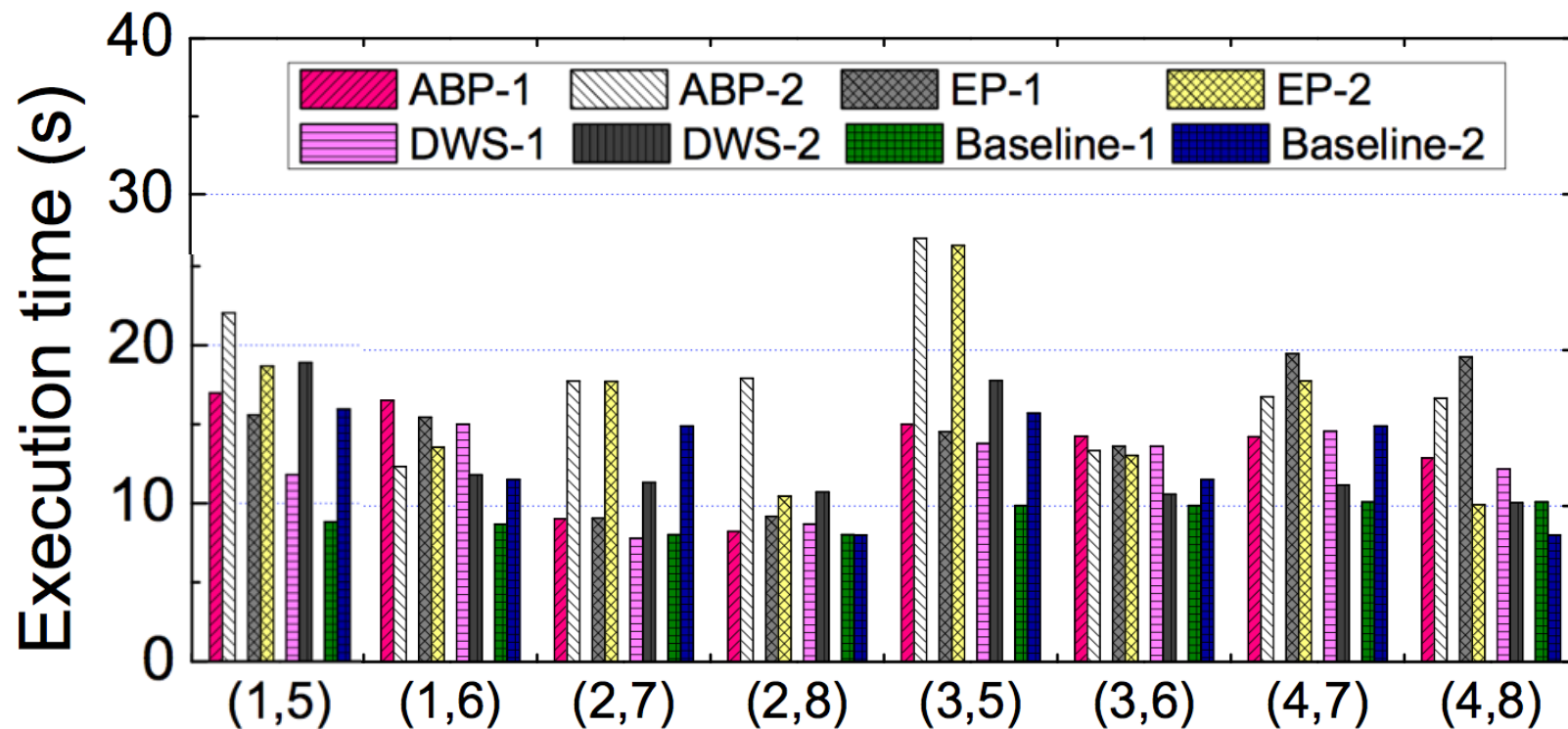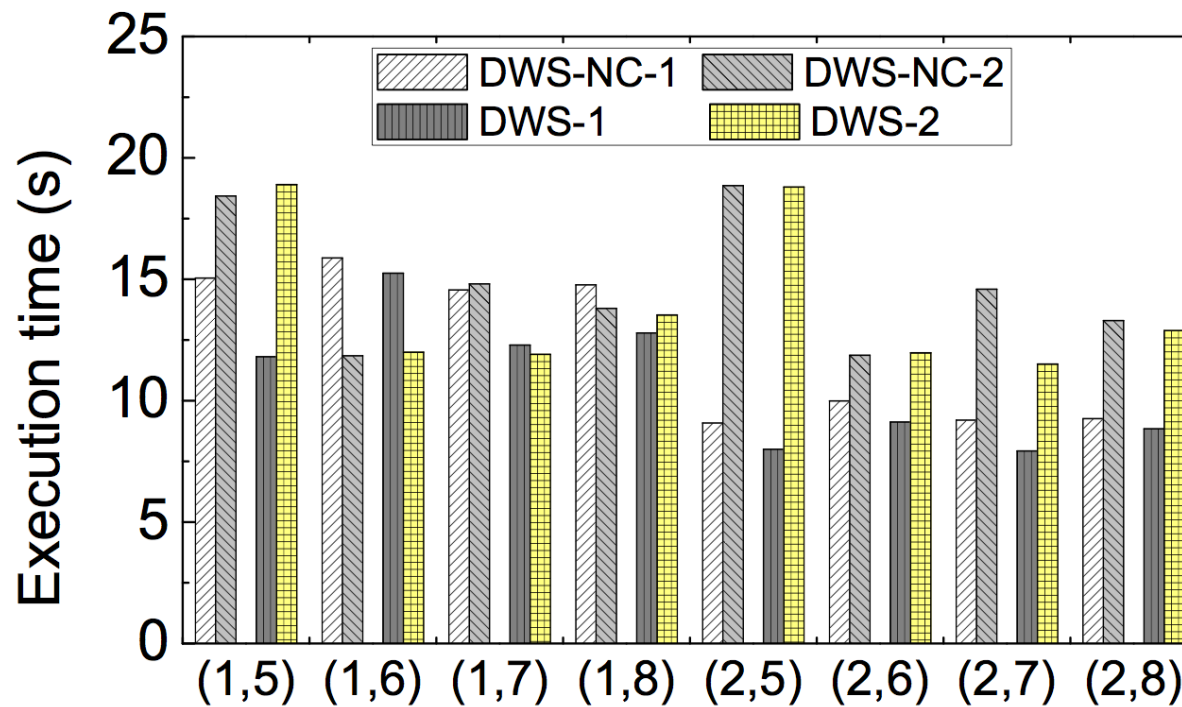| ID | Name | Description |
|---|---|---|
| p-1 | FFT | Fast Fourier Transform |
| p-2 | PNN | Polynomial Neural Network |
| p-3 | Cholesky | Cholesky decomposition |
| p-4 | LU | LU decomposition |
| p-5 | GE | Gaussian Elimination algorithm |
| p-6 | Heat | Five-point heat distribution |
| p-7 | SOR | 2D Successive Over-Relaxation |
| p-8 | Mergesort | Merge sort on 4E6 numbers |

## Calculate execution time:



$$T_i = \frac{\sum_{r=1}^{a} t_{ir}}{a}, T_j = \frac{\sum_{r=1}^{b} t_{jr}}{b}$$

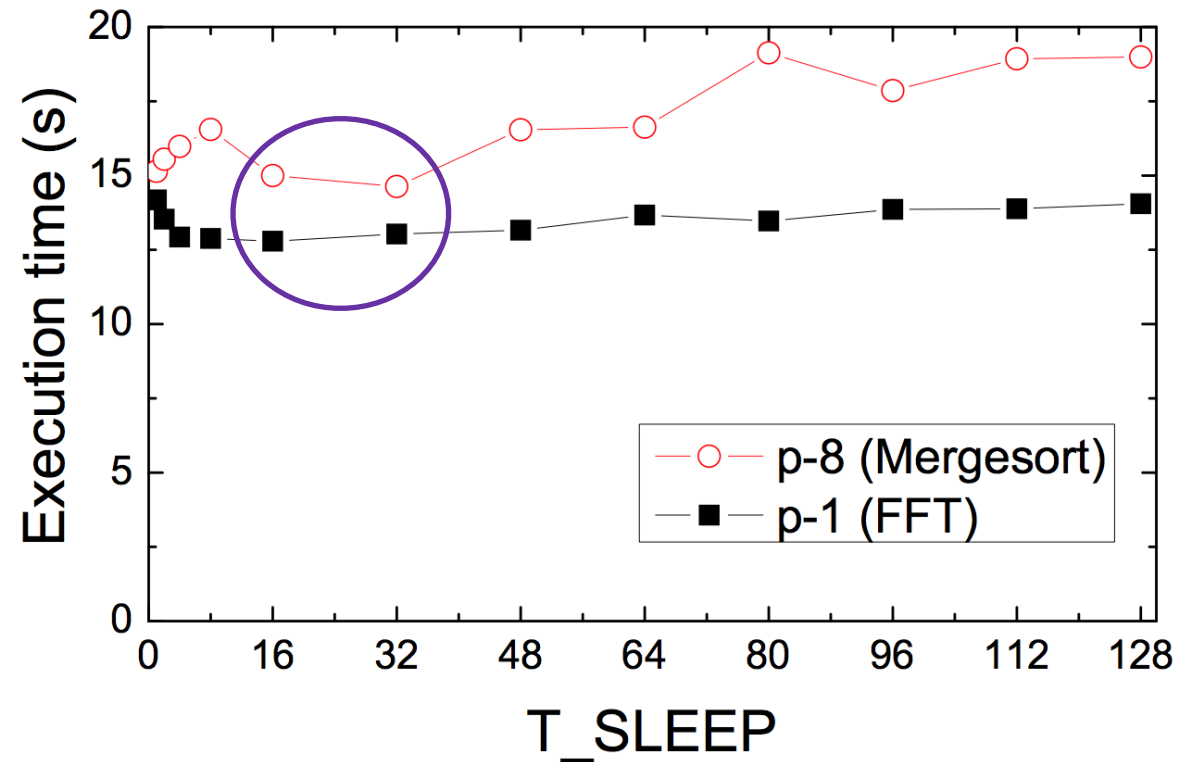DWS can significantly improve the performance of the benchmarks

Without the coordinator, the performance of the benchmarks is degraded

We should choose T_SLEEP = *k* or *2k* on a k-core computer

# Contributions & conclusions

- A modified work-stealing algorithm that enables a program to release the under-utilized cores.

- A coordinator to manage the workers. It enables a program to grab and use the under-utilized cores released by other programs.

- We have implemented DWS, which achieves a performance gain of up to 32.3% in the best cases compared to traditional work-stealing schedulers.

# Thanks!

# Questions?