# Operating System: Multiprocessor Scheduling

Sang Ho Choi (shchoi@kw.ac.kr)

School of Computer & Information Engineering

KwangWoon University
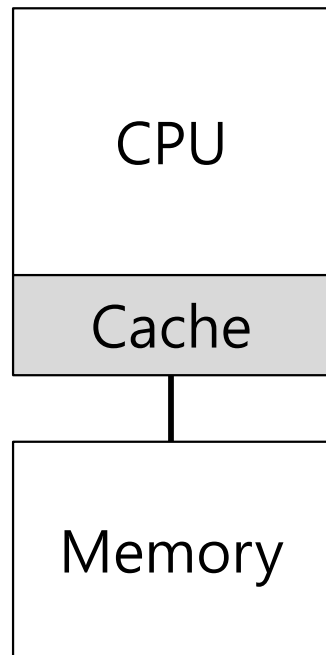
# Multiprocessor Scheduling

- The rise of the <span style="color:red">multicore processor</span> is the source of multiprocessor-scheduling proliferation 멀티프로세서의 scheduling을 어떻게 구성할것인가?

  – **Multicore**: Multiple CPU cores are packed onto a single chip

> **How to schedule jobs on <span style="color:red">Multiple CPUs</span>?**

- Adding more CPUs <u>does not</u> make that single application run faster 단일 어플리케이션의 속도가 빨라지는게 X

  ➢ You'll have to rewrite application to run in parallel, using **threads** 병렬로 스레드를 사용

광운대학교
KwangWoon University

# Single CPU with cache



SRAM

**Cache**
- Small, fast memories
- Hold copies of <u>popular</u> data that is found in the main memory 많이사용하는것은 캐쉬해둠.
- Utilize *temporal* and *spatial* <mark>locality</mark> 시간지역성 공간지역성.

DRAM

**Main Memory**
- Holds all of the data
- Access to main memory is slower than cache

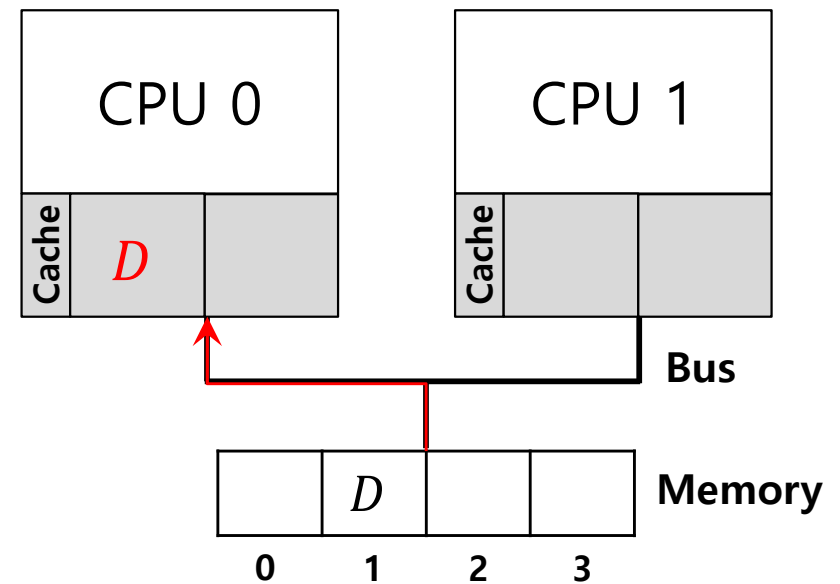**By keeping data in cache, the system can make slow memory appear to be a fast one**

광운대학교
KwangWoon University

# Cache coherence

- Consistency of shared resource data stored in multiple caches

0. Two CPUs with caches sharing memory

| CPU 0 | CPU 1 |

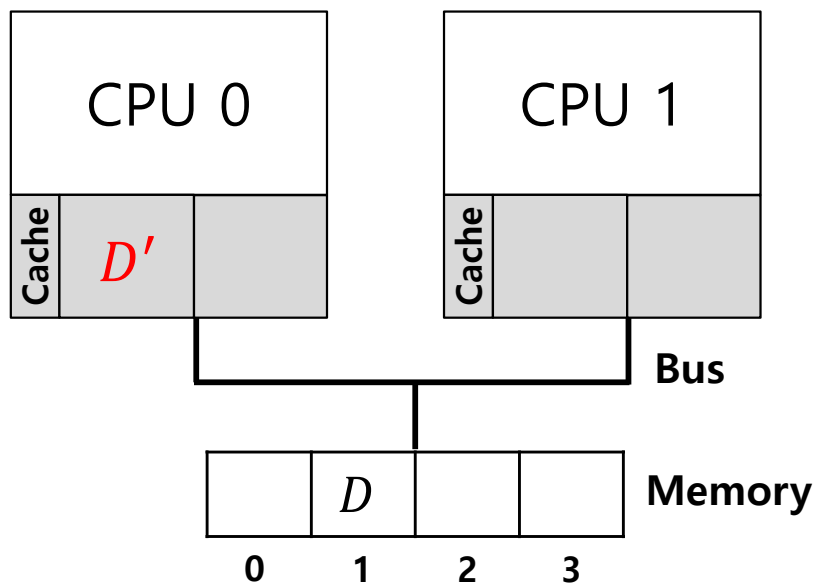Cache

Cache

Bus

| | D | | | Memory

0  1  2  3

1. CPU0 reads a data at address 1.

| CPU 0 | CPU 1 |

Cache  *D*

Cache

Bus

| | D | | | Memory

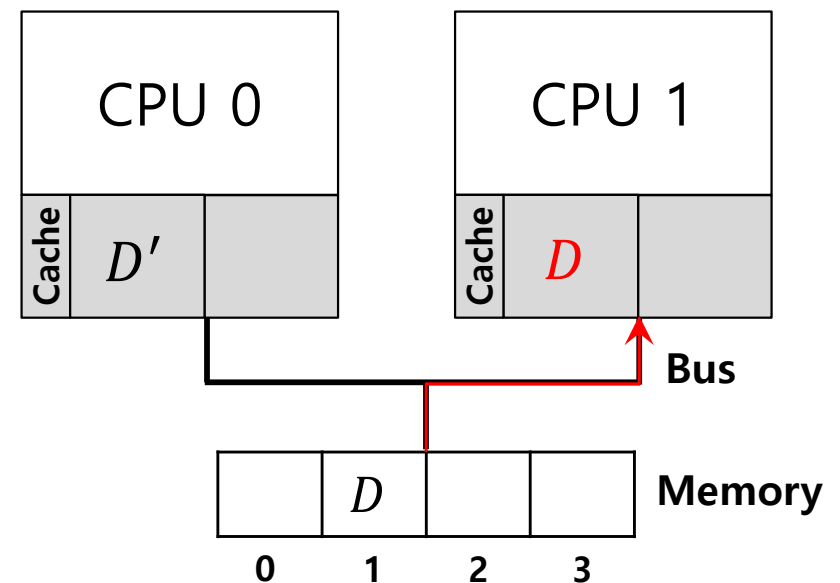0  1  2  3

광운대학교
KwangWoon University

# Cache coherence (Cont.)

2. $D$ is updated and CPU1 is scheduled.



3. CPU1 re-reads the value at address A



**CPU1 gets the old value $D$ instead of the correct value $D'$**

# Cache coherence solution

- Bus snooping
  - Each cache pays attention to memory updates by **observing the bus**

  - When a CPU sees an update for a data item it holds in its cache, it will notice the change and either <u>invalidate</u> its copy or <u>update</u> it

# Don't forget synchronization

- When accessing shared data across CPUs, mutual exclusion primitives should likely be used to guarantee correctness

```
1          typedef struct __Node_t {
2                  int value;
3                  struct __Node_t *next;
4          } Node_t;
5
6          int List_Pop() {
7                  Node_t *tmp = head;            // remember old head ...
8                  int value = head->value;       // ... and its value
9                  head = head->next;             // advance head to next pointer
10                 free(tmp);                     // free old head
11                 return value;                  // return value at head
12         }
```

**Simple List Delete Code**

# Don't forget synchronization (Cont.)

- Solution

```
1          pthread_mtuex_t m;
2          typedef struct __Node_t {
3                    int value;
4                    struct __Node_t *next;
5          } Node_t;
6
7          int List_Pop() {
8                    lock(&m)
9                    Node_t *tmp = head;          // remember old head ...
10                   int value = head->value;     // ... and its value
11                   head = head->next;           // advance head to next pointer
12                   free(tmp);                   // free old head
13                   unlock(&m)
14                   return value;                // return value at head
15         }
```

**Simple List Delete Code with lock**

# Cache Affinity

*(한자성.)*

*프로세스 하나는 CPU 하나가 처리하기 기존에 관련있는 캐시데이터를 재활용 학습 있으나.*

- Keep a process on <span style="color:red">the same CPU</span> if at all possible
  - A process builds up a fair bit of state <u>in the cache</u> of a CPU
  - The next time the process run, it will run faster if some of its state is *already present* in the cache on that CPU

> **A multiprocessor scheduler should consider <span style="color:red">cache affinity</span> when making its scheduling decision**

# Single queue Multiprocessor Scheduling (SQMS)
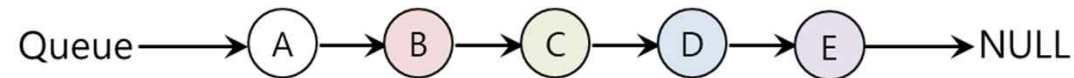
- Put all jobs that need to be scheduled into a single queue
    - Each CPU simply picks the next job from the globally shared queue
    - Cons:
        1) Some form of **locking** have to be inserted → Lack of scalability 확장성.
        2) Cache affinity
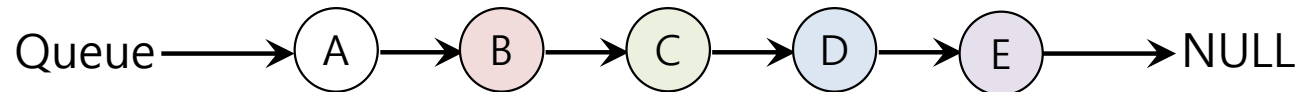            ✓ Ex) Possible job scheduler across CPUs:

Queue → A → B → C → D → E → NULL

| CPU0 | A | E | D | C | B | ... (repeat) ... |
|------|---|---|---|---|---|------------------|
| CPU1 | B | A | E | D | C | ... (repeat) ... |
| CPU2 | C | B | A | E | D | ... (repeat) ... |
| CPU3 | D | C | B | A | E | ... (repeat) ... |

광운대학교
KwangWoon University

# Scheduling Example with Cache affinity

캐시친화성을 적용

Queue → A → B → C → D → E → NULL

|        |   |   |   |   |   |            |
|--------|---|---|---|---|---|------------|
| CPU0   | A | E | A | A | A | … (repeat) … |
| CPU1   | B | B | E | B | B | … (repeat) … |
| CPU2   | C | C | C | E | C | … (repeat) … |
| CPU3   | D | D | D | D | E | … (repeat) … |

– <u>Preserving affinity</u> for most

  ➢ Jobs A through D are not moved across processors

  ➢ Only job e Migrating from CPU to CPU

– Implementing such a scheme can be **complex**

광운대학교
KwangWoon University
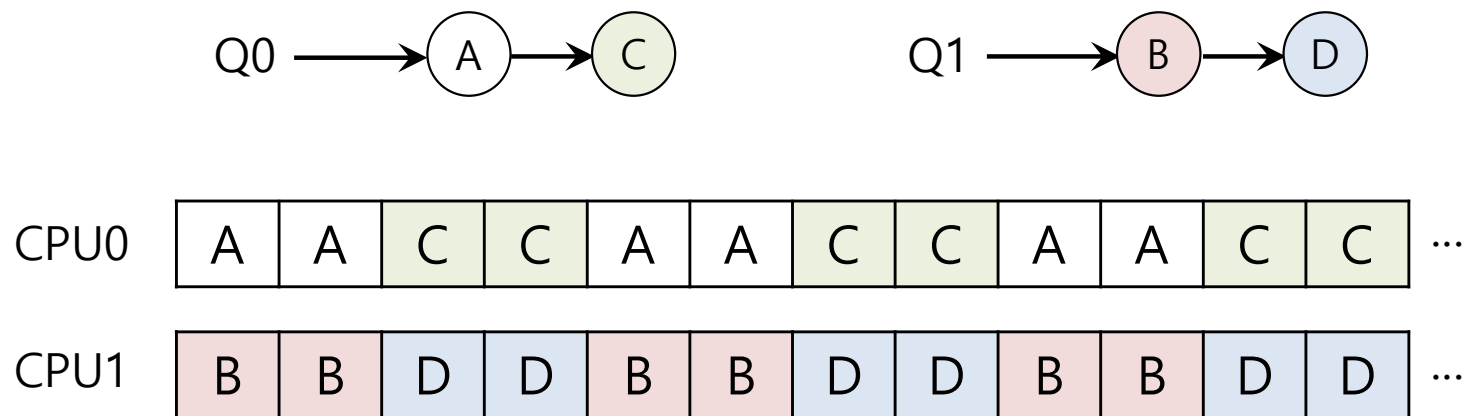
# Multi-queue Multiprocessor Scheduling (MQMS)

*1Q — 1CPU*

- MQMS consists of <span style="color:red">multiple scheduling queues</span>
  - Each queue will follow a particular scheduling discipline
  - When a job enters the system, it is placed on **exactly one** scheduling queue
  - Avoid the problems of <u>information sharing</u> and <u>synchronization</u>
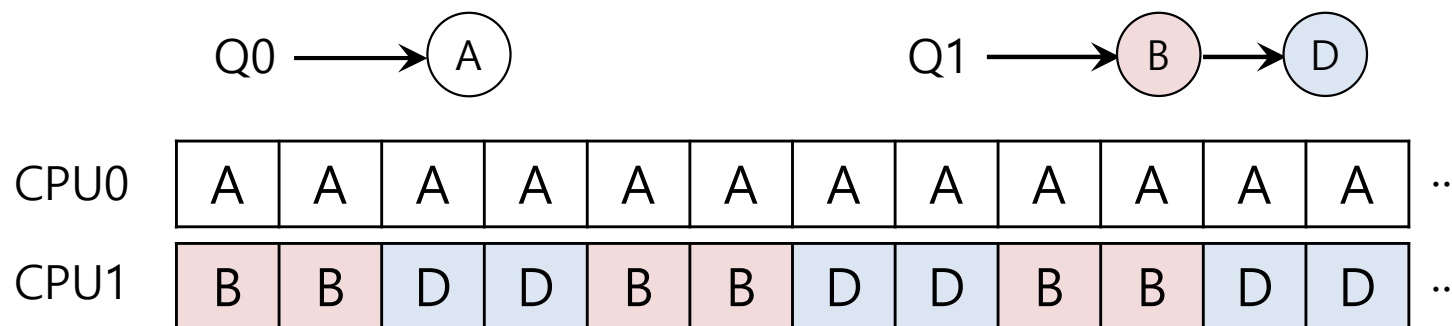
# MQMS Example

- With **round robin**, the system might produce a schedule that looks like this:

| Q0 → A → C | Q1 → B → D |
|---|---|

| CPU0 | A | A | C | C | A | A | C | C | A | A | C | C | … |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| CPU1 | B | B | D | D | B | B | D | D | B | B | D | D | … |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**MQMS provides more scalability and cache affinity**

# Load Imbalance issue of MQMS

- ## After job C in Q0 finishes:

Q0 ⟶ A          Q1 ⟶ B ⟶ D

| CPU0 | A | A | A | A | A | A | A | A | A | A | A | A | ... |
|------|---|---|---|---|---|---|---|---|---|---|---|---|-----|
| CPU1 | B | B | D | D | B | B | D | D | B | B | D | D | ... |

**A gets twice as much CPU as B and D**

- ## After job A in Q0 finishes:

Q0 ⟶          Q1 ⟶ B ⟶ D

| CPU0 |   |   |   |   |   |   |   |   |   |   |   |   | ... |
|------|---|---|---|---|---|---|---|---|---|---|---|---|-----|
| CPU1 | B | B | D | D | B | B | D | D | B | B | D | D | ... |

**CPU0 will be left idle!**

*CPU utilization이 떨어짐. 특정큐에 몰리게 되면.*

# How to deal with load imbalance?

- ## The answer is to move jobs (<u>Migration</u>)
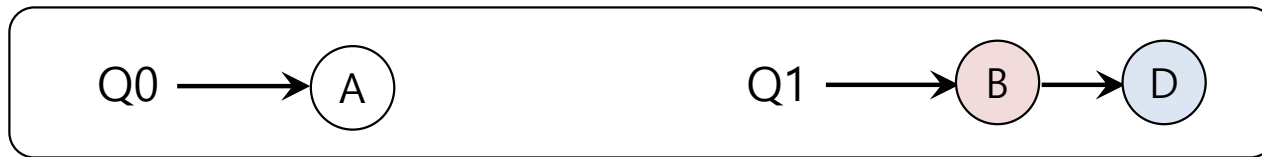
  - Example:
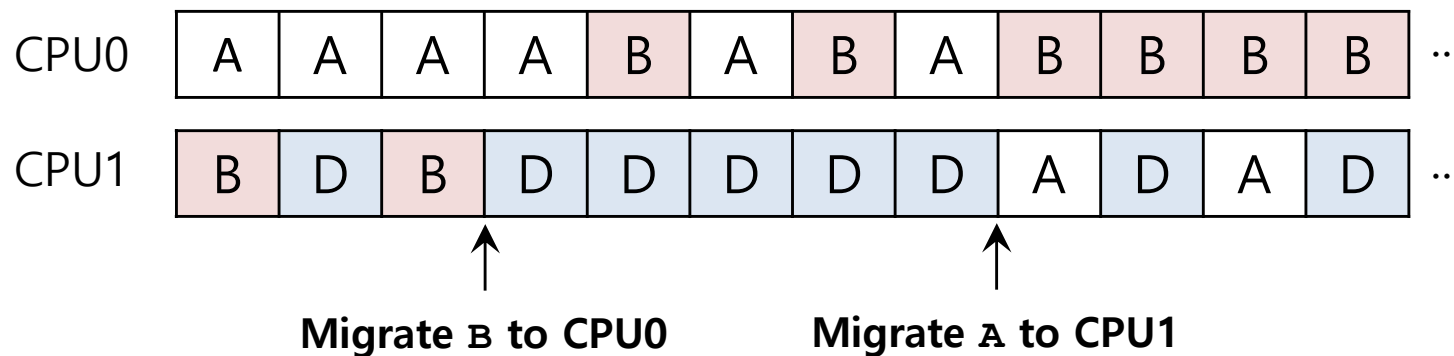


The OS moves one of B or D to CPU 0

Or

하나가 비면
그냥이동하면 되는데

# How to deal with load imbalance? (Cont.)

- A more tricky case:



- A possible migration pattern:

  - Keep switching jobs

# Work Stealing

- Move jobs between queues

  - Implementation:

    - A source queue that is <u>low on jobs</u> is picked   *job의 개수가 가장적은 큐를 source queue.*

    - The source queue occasionally peeks at another target queue   *노스큐가 다른큐를 바라보다가*

    - If the target queue is <u>more full than</u> the source queue, the source will "**steal**" one or more jobs from the target queue   *소스큐가 뺏어감*
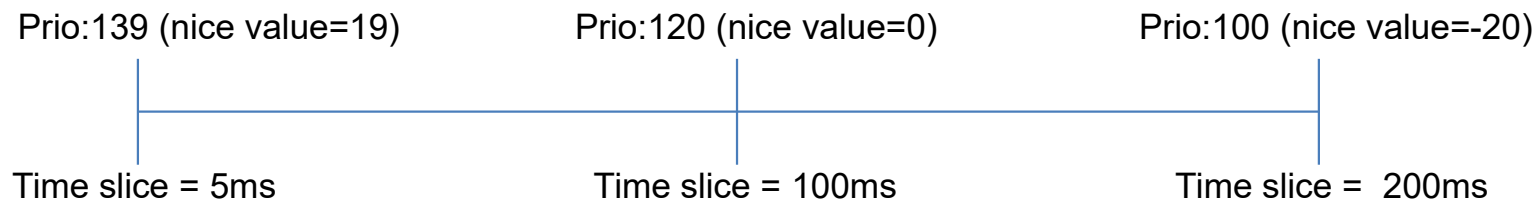
  - Cons:

    - *High overhead* and trouble *scaling*

      *steal 과정에서 생기는 overhead 나 확장성 문제가 있을수 있어 설계시 주의해야함.*

# Linux Multiprocessor Schedulers

- Time-sharing with time slice
  - Only processes with time slice can execute
  - Time slice is subtracted when timer interrupt occurs
  - When time slice = 0, another process is chosen
  - When all processes have time slice = 0, recalculation is performed

- Real-time with priority
  - Highest priority process always runs first
  - 140 priority levels
    - The lower the value, the higher priority
    - E.g. priority level 110 will have a higher priority than 130

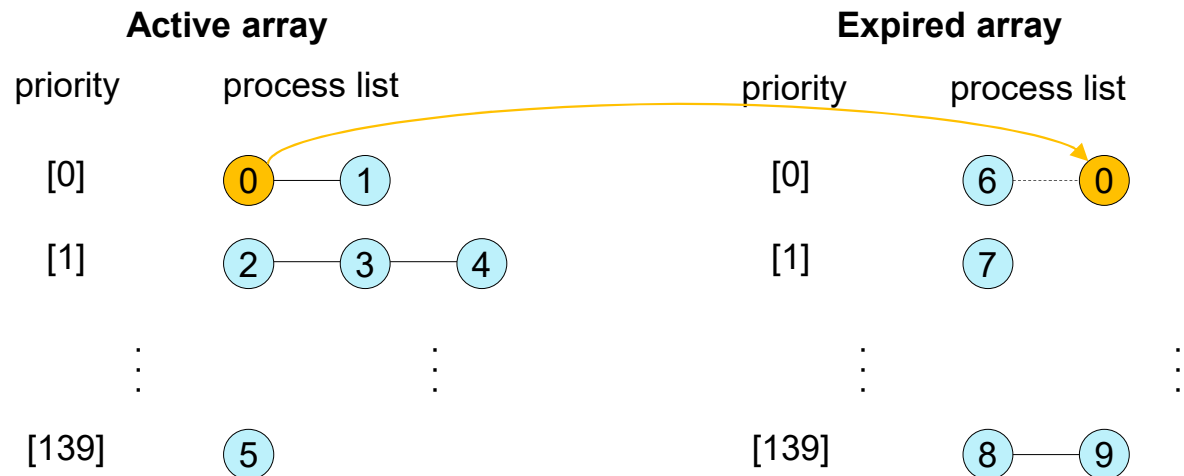- Linux schedules the highest priority process with time slice

# Linux Multiprocessor Schedulers (Cont.)

- O(1) scheduler
  - It selects a highest priority process in a ready queue within a constant amount of time, regardless of how many processes are running

  - It divides the ready queue into multiple queues according to priority
    - 0~99 for real time processes
    - 100-139 for non-real time processes

  - Time slice is given based on the priority

Prio:139 (nice value=19)      Prio:120 (nice value=0)      Prio:100 (nice value=-20)

Time slice = 5ms       Time slice = 100ms      Time slice = 200ms

# Linux Multiprocessor Schedulers (Cont.)

- O(1) scheduler
  - Active array - list of processes with time remaining in their time slices
  - Expired array - list of expired processes
  - The scheduler chooses the process with the highest priority from the active array

| Active array | | Expired array | |
|:---:|:---:|:---:|:---:|
| priority | process list | priority | process list |
| [0] | 0 — 1 | [0] | 6 ---- 0 |
| [1] | 2 — 3 — 4 | [1] | 7 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| [139] | 5 | [139] | 8 — 9 |

  - When the active array is empty, the two arrays are exchanged

*[handwritten Korean: 모두 종료하면 Active array가 비내가 이게 expired array로 사용하는거지, expired는 Active로 쓰고.]*

# Linux Multiprocessor Schedulers (Cont.)

- CFS (Complete Fair Scheduling) scheduler
  - Since kernel 2.6.23.
  - Proposed to overcome some oddity of O(1) scheduler
  - Time slice based on **priority** in O(1) scheduler

    *우선순위차이는 비소형 차이가멸로 안나.*

    | process 1 (prio=120) | process 2 (prio=121) |
    |---|---|

    ➔ Time slice = 100ms, 95ms   *차이가멸로 X*

    | process 3 (prio=138) | process 4 (prio=139) |
    |---|---|

    ➔ Time slice = 10ms, 5ms   *2배*   ➔ relatively 2 times

  - Time slice based on **weight** in CFS scheduler
    - ➢ Roughly 1.25 times for each priority level difference

      | process 5 (prio=n) | process 6 (prio=n+1) |
      |---|---|

      ➔ Time slice = 56ms, 44ms   ➔ relatively 1.25 times

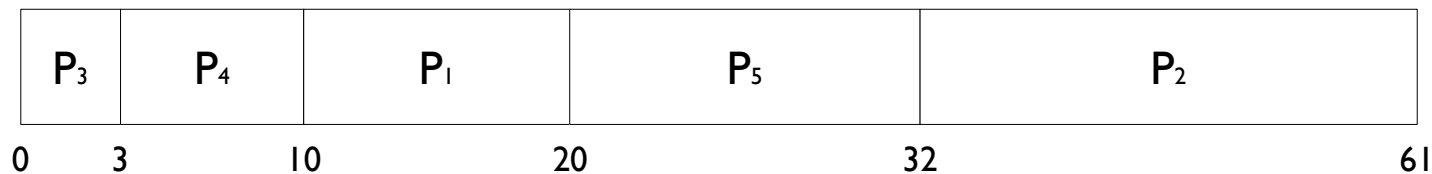      | process 7 (prio=m) | process 8 (prio=m+5) |
      |---|---|

      ➔ Time slice = 75ms, 25ms   ➔ relatively 3 times

# Algorithm Evaluation

- Deterministic modeling
  - takes a particular predetermined workload, and then
  - defines the performance of each algorithm for that workload

  - E.g. SJF with Gantt Chart, when given a set of processes

| $P_3$ | $P_4$ | $P_1$ | $P_5$ | $P_2$ |
|-------|-------|-------|-------|-------|

0    3        10           20           32           61

  - Simple and fast
  - Useful in describing the algorithm and providing examples

광운대학교
KwangWoon University
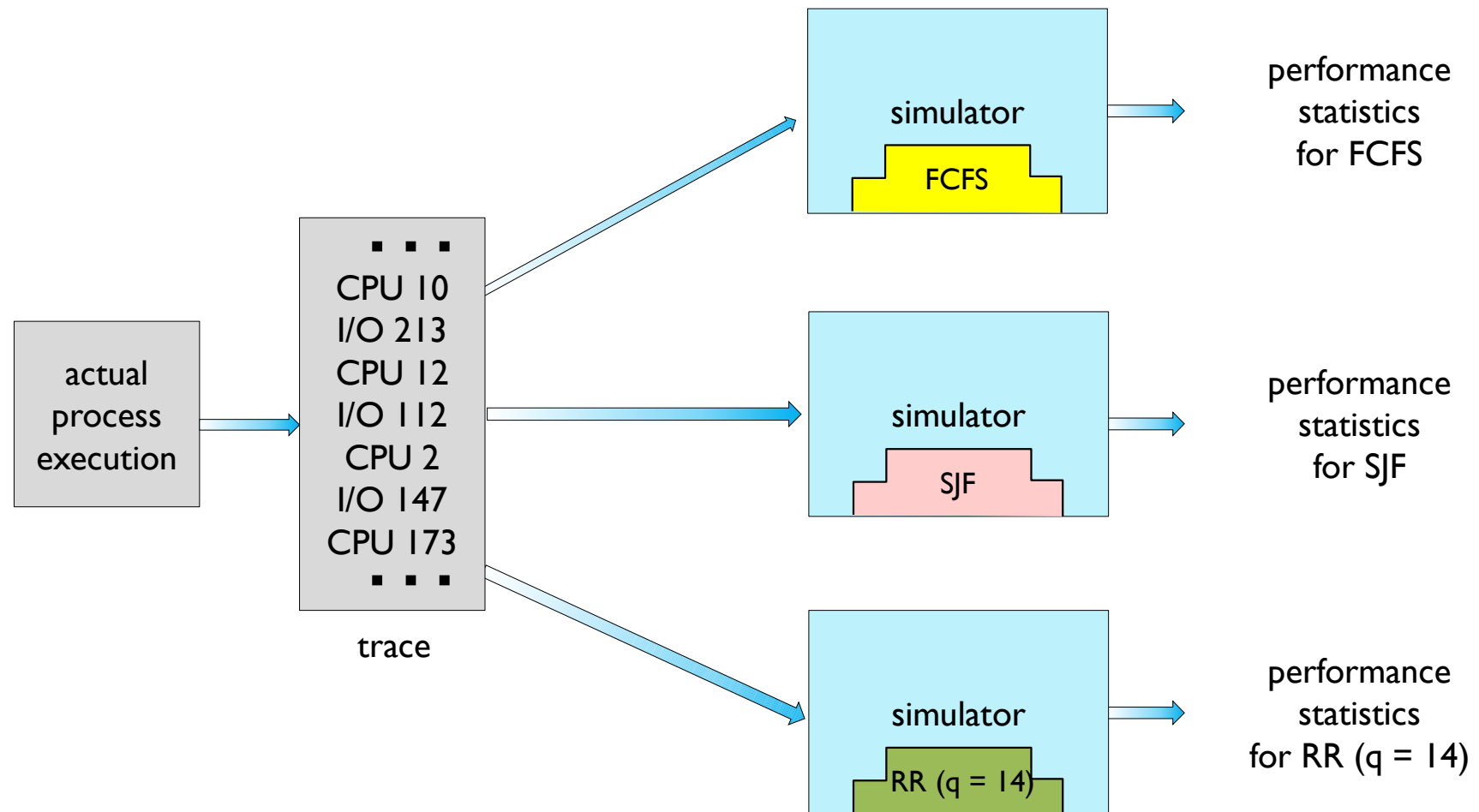
# Algorithm Evaluation

- Implementation
  - Only completely accurate way to evaluate algorithm
  - is very difficult
    - algorithm coding
    - operating system modification

# Algorithm Evaluation

- Simulation
  - programs a model of computer system
    - ➢ Data structures represent the major components of system
    - ➢ E.g. a variable for clock

  - Random number generator is required
    - ➢ The number of processes, CPU burst times, arrivals, departures, etc
    - ➢ Uniform, exponential, poisson, etc

  - Can be expensive
    - ➢ Design, coding, and debugging of simulator is time consuming
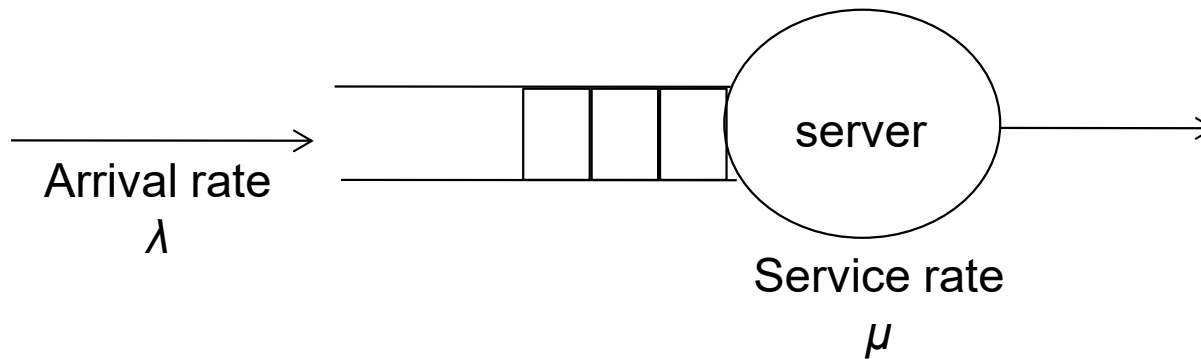    - ➢ Large storage space and time may be required

# Algorithm Evaluation

- ## Simulation example



actual process execution → trace (CPU 10, I/O 213, CPU 12, I/O 112, CPU 2, I/O 147, CPU 173) → simulator (FCFS) → performance statistics for FCFS; simulator (SJF) → performance statistics for SJF; simulator (RR (q = 14)) → performance statistics for RR (q = 14)

# Algorithm Evaluation

- Queueing models
  - with a mathematical formula
  - Simple queueing network model



  - Given arrival rates and service rates,
    - utilization, average queue length, and average waiting time can be derived
  - Limitedly used in some evaluation