# Scheduling

Instructor: Youngjin Kwon

# Refine Process

## (Unix) Process

```
A(int tmp) {
  if (tmp<2)
    B();
  printf(tmp);
}
B() {
  C();
}
C() {
  A(2);
}
A(1);
…
```

Memory

I/O State
(e.g., file,
socket
contexts)

CPU state
(PC,
registers..)

Resources

Sequential
stream of
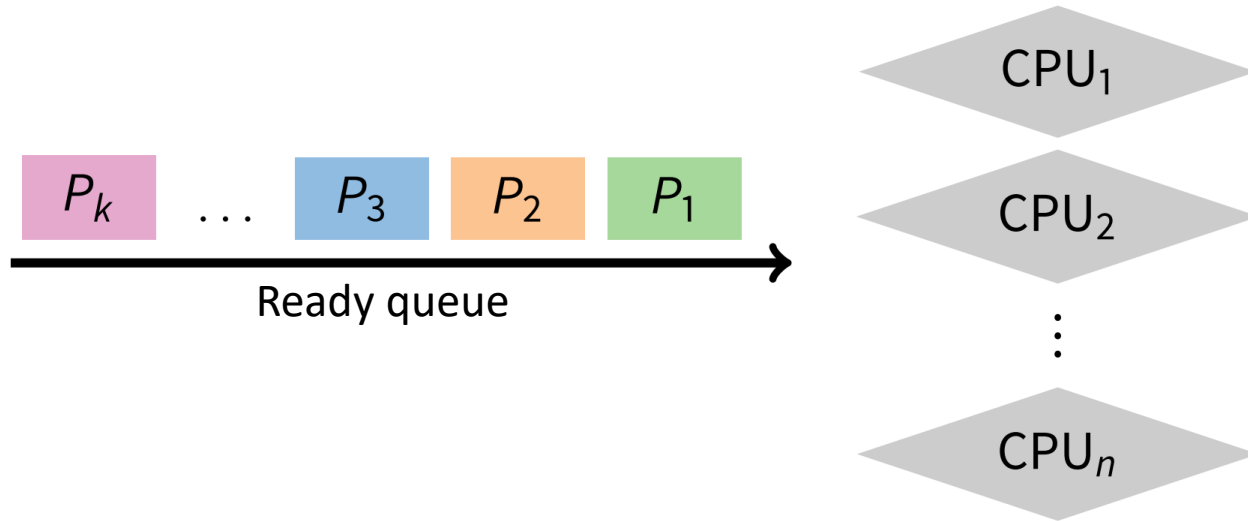instructions

# Processes

# Main Points

- **Scheduling policy**: ==what to do next==, when ==there are multiple tasks ready to run==
  - Or multiple packets to send, or web requests to serve, or …
- Uniprocessor policies
  - FIFO, round robin, optimal
  - multilevel feedback as approximation of optimal
- Multiprocessor policies
  - Affinity scheduling, gang scheduling

# Example

- You manage a web site, that suddenly becomes wildly popular.  Do you?
    - Buy more hardware?  *more processor*
    - Turn away some users?  Which ones?  *evict process*
    - Implement a different scheduling policy?

    *efficient policy*
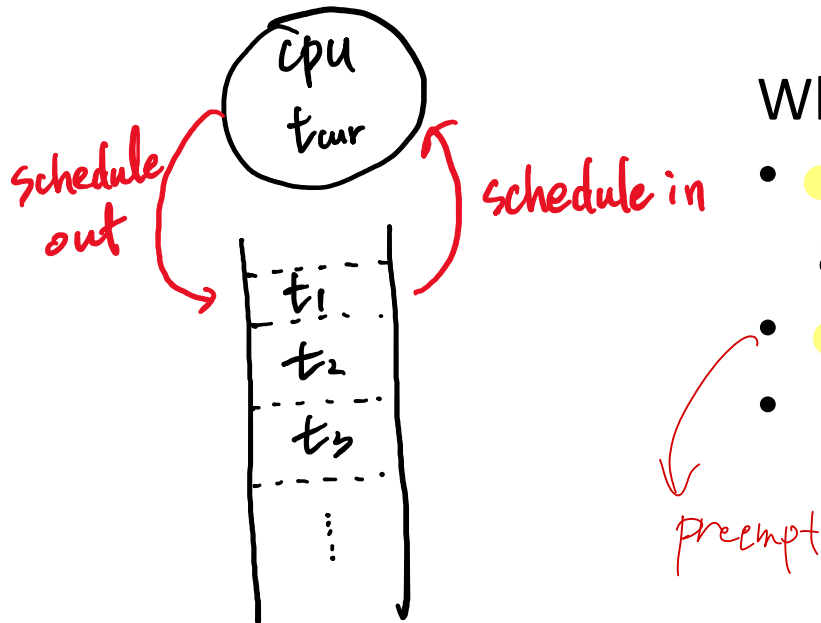
# Design: scheduling problem



- Scheduling algorithm
  - takes a workload as input
  - decides which tasks to do first
  - Performance metric (throughput, latency) as output

# Definitions

- Task/Job
  - User request: e.g., mouse click, web request, shell command, …
- Workload
  - Set of tasks for system to perform
- Overhead
  - How much extra work is done by the scheduler?
- Fairness
  - How equal is the performance received by different users?
- Predictability
  - How consistent is the performance over time?

# Scheduler concept



What scheduler does?
- Pick a task from run queue according to scheduler algorithm
- Kick out the running task from CPU
- Make the selected task run in CPU

# Scheduler design choice

- Preemptive scheduler
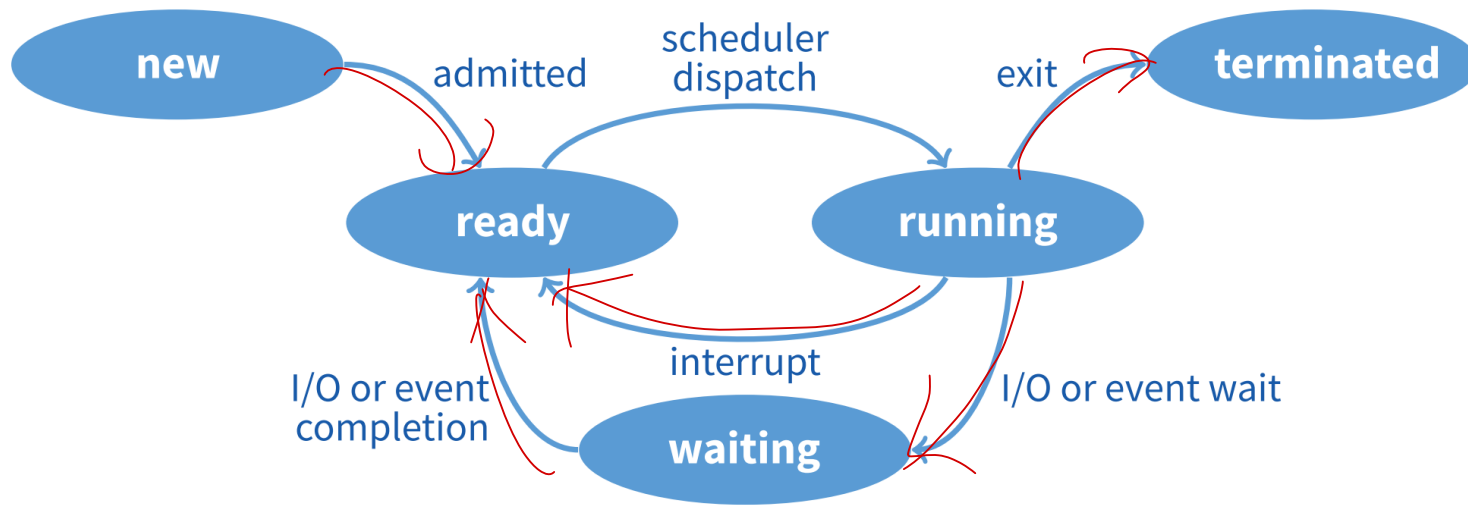  - If we can take resources away from a running task

→ thread가 리턴 상태에 있으면 CPU가 될 있음.

- Work-conserving
  - Resource is used whenever there is a task to run
  - When is non work-conserving scheduler useful?

→ Network system waiting for full packets to arrive

# When does OS invoke scheduler?



Preemptive scheduler:
1. Waiting → Ready
2. Running → Waiting
3. Running → Ready
4. New/waiting → Ready
5. Exit

Non-preemptive scheduler:

New → Ready

exit

# Scheduler performance metric

- Throughput
  - How many tasks can be done per unit of time?
  - # of jobs / time
- Turnaround time
  - How long does a task take to complete?
  - $T_{finish}$ - $T_{arrival}$
- Response time
  - Time from request to "first" response
  - $T_{response}$ - $T_{arrival}$
- Waiting time
  - Waiting time of a task = $\sum$ Time spent in ready & wait states
  - Average waiting time = Avg. (waiting time of tasks in system)

# Contents

- Uniprocessor policies
  - FIFO, round robin, optimal
  - multilevel feedback as approximation of optimal

- Multiprocessor policies
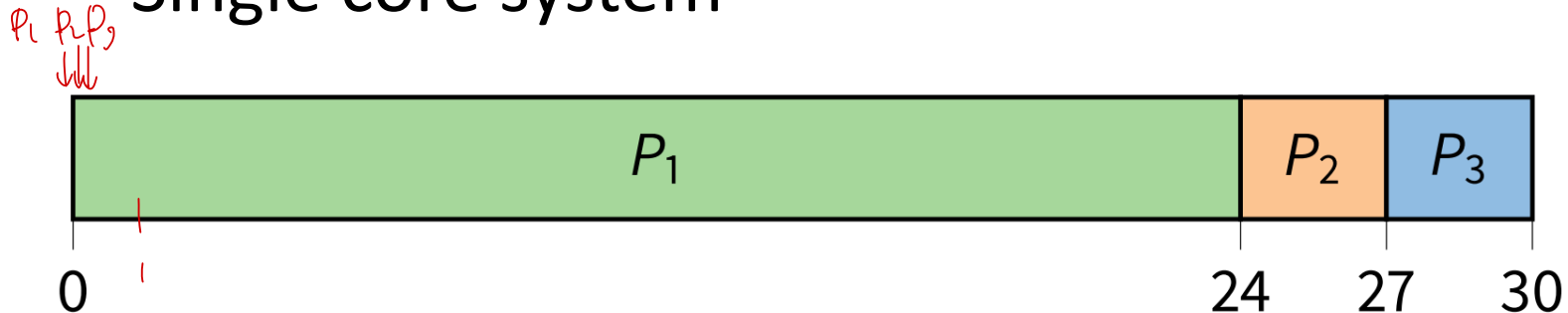  - Affinity scheduling, gang scheduling

# First In First Out (FIFO)

*Non preemptive scheduling*

- Schedule tasks in the order they arrive
  - Run tasks in order that they arrive

- Example: caching server
  - Facebook: cache of friend lists, image blobs etc


- On what workloads are FIFO particularly bad?

# FIFO scheduling example

- P1 needs 24 sec, P2 and P3 needs 3 sec
- Arrival order: P1, P2, P3
- Single core system



- Avg. Throughput: 3 jobs/30sec = 0.1 jobs/sec
- Avg. Turnaround time : (24 + 27 + 30) / 3 = 27
- Avg. wait time : (0 + 24 + 27) / 3 = 17

## Can we do better?

# Beyond FIFO scheduling

- T1 needs 24 sec, T2 and T3 needs 3 sec
- Changing scheduler order: P2, P3, P1



- Avg. Throughput: 3 jobs/30sec = 0.1 jobs/sec
- Avg. Turnaround time: (3 + 6 + 30) / 3 = 13    <27
- Avg. wait time : (0 + 3 + 6) / 3 = 3    <17

**Lesson: schedule algorithm**
**can reduce turnaround time and wait time**

# Convoy effect



Img source:
https://cs.jhu.edu/~huang/cs318/fall18/lectures/lec4_sched.pdf



image source:
http://web.cs.ucla.edu/classes/fall14/cs111/scribe/7a/convoy_effect.png

# Convoy effect
# in CPU-bound vs. IO-bound jobs

- CPU-bound jobs will hold CPU until exit or I/O

  – But I/O burst for CPU-bound job is small

  – Long periods where no I/O requests issued, and CPU held

  – Result: poor I/O device utilization

| job 1 | job 2 |
|---|---|
| CPU burst | I/O idle |
| I/O burst | CPU burst |
| | I/O burst |
| CPU burst | |
| | I/O idle |
| I/O burst | CPU burst |
| CPU-bound | I/O-bound |

# Shortest Job First (SJF)

*→ min ANT!*

*Non preemptive vs Preemptive*

- Always do the task that has the shortest remaining amount of work to do
  - Often called Shortest Remaining Time First (SRTF)

FIFO:

*24*                              *3*     *3*

| | | |
|---|---|---|
| $P_1$ | $P_2$ | $P_3$ |

0                              24   27   30

SJF:

| | | |
|---|---|---|
| $P_2$ | $P_3$ | $P_1$ |

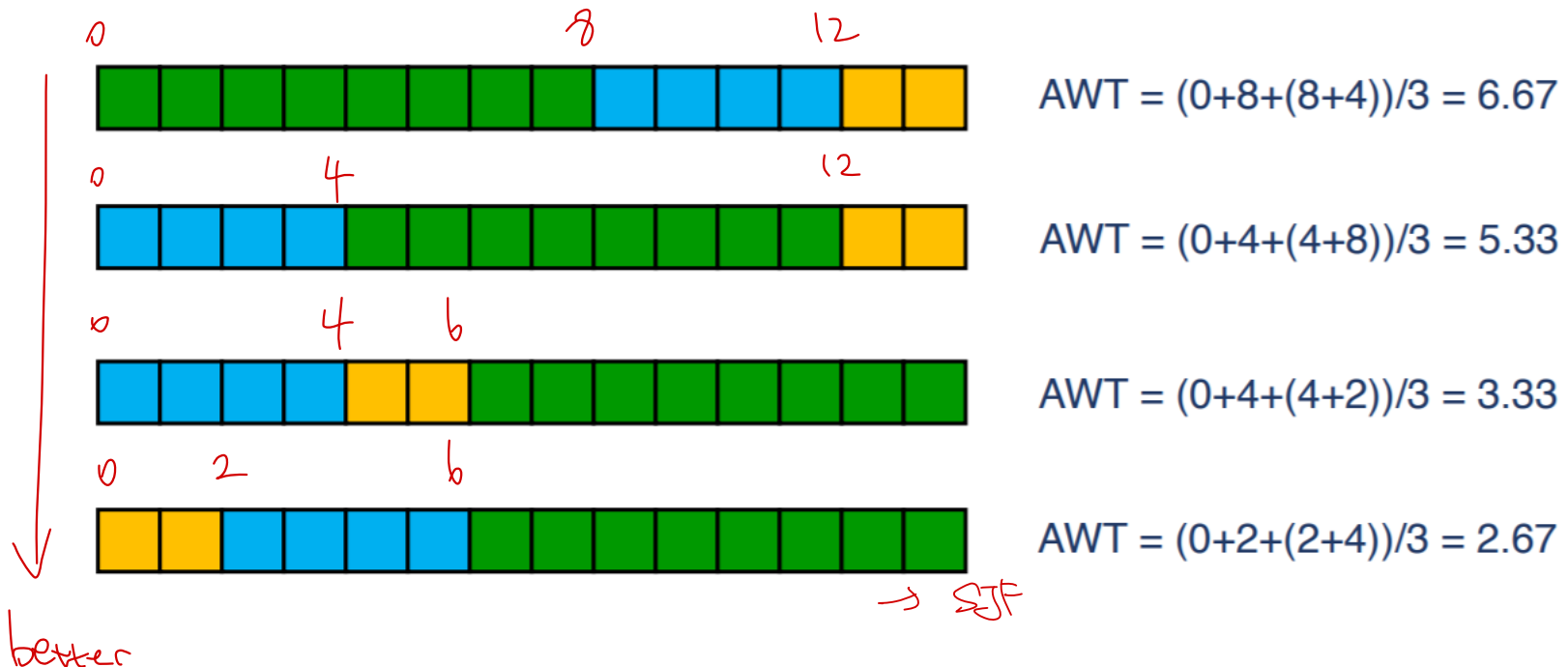0    3    6                                      30

# Shortest Job First (SJF)

- Provably optimal <mark>minimum average waiting time (AWT)</mark>

Choose the job with the smallest expected CPU burst



$AWT = (0+8+(8+4))/3 = 6.67$

$AWT = (0+4+(4+8))/3 = 5.33$

$AWT = (0+4+(4+2))/3 = 3.33$

$AWT = (0+2+(2+4))/3 = 2.67$

→ SJF

better

# Two schemes of SJF

- ## Non-preemptive SJF
  - Once CPU is given to a process, it cannot be preempted (kicked-out from currently using CPU) until it completes its work

- ## Preemptive SJF
  - If a new process arrives with shorter remaining time of current running process, preempt it

**What scheduling metric does SJF improve overs FIFO?**  *AWT, Response time*

**Non-preemptive vs preemptive.**  *preemptive.*
**Which has better response time? AWT?**  *?*

# Draw scheduling diagram

$TT_N = (7+4+10+11)/4 = 8$     $TT_P = (16+5+1+6)/4$
$= 7$

| Process | Arrival Time | Burst Time |
|---------|--------------|------------|
| $P_1$ | 0 | 7 |
| $P_2$ | 2 | 4 |
| $P_3$ | 4 | 1 |
| $P_4$ | 5 | 4 |

Response time
$P_1=0, P_2=6, P_7=3, P_4=7$

$AWT = (0+6+3+7)/4$
$= 16/4 = 4$

P1   P2   P3  P4                           P4 tie, choose by arrival time.

**Non-preemptive:**

| $P_1$ | $P_3$ | $P_2$ | $P_4$ |

0 ........... 7  8 ......... 12 .......... 16

Remaining time:   $P_1=5, P_2=4$     $P_1=5, P_2=2$     $P_1:5, P_4:4$     $P_1:5$
                                     $P_3=1$

**Preemptive:**

| $P_1$ | $P_2$ | $P_3$ | $P_2$ | $P_4$ | $P_1$ |

0 ..... 2 ..... 4  5 ..... 7 ......... 11 ......... 16

Response time | AWT of $P_2$
$P_1=0$  $P_4=2$ | $P_2 = 0+1$
$P_2:0$ | $P_1 = 0+9$   $P_3 = 0$   $P_4 = 2$
$P_3=0$

$P_1=5, P_3=2, P_4=4$

**Compute and compare avg. response time and AWT**

$AWT = 12/4 = 3$

$\rightarrow$ 둘다 preemptive가 better

# SJF limitations

*Turnaround Time*

- Doesn't always minimize average TT
  - Only minimizes waiting time
  - How to improve TT of SJF?

- Sometime, impossible to know size of CPU burst ahead of time → *predicting..*
  - Like choosing person in line without looking inside basket/cart

# How to predict CPU burst time?

request $r_1 \rightarrow$

$r_2 \rightarrow$

$r_3 \rightarrow$

predict burst time of $r_1$

- Estimate CPU burst length based on the past
  - E.g., Exponentially Weighted Average
    - $t_n$ : length of a process's n-th CPU burst
    - $t_{n+1} = \alpha t_n + (1 - \alpha)t_{n-1}$     linear combination of
      - $\alpha$ is parameter (e.g., 0.5)     n-th and (n-1)-th requests
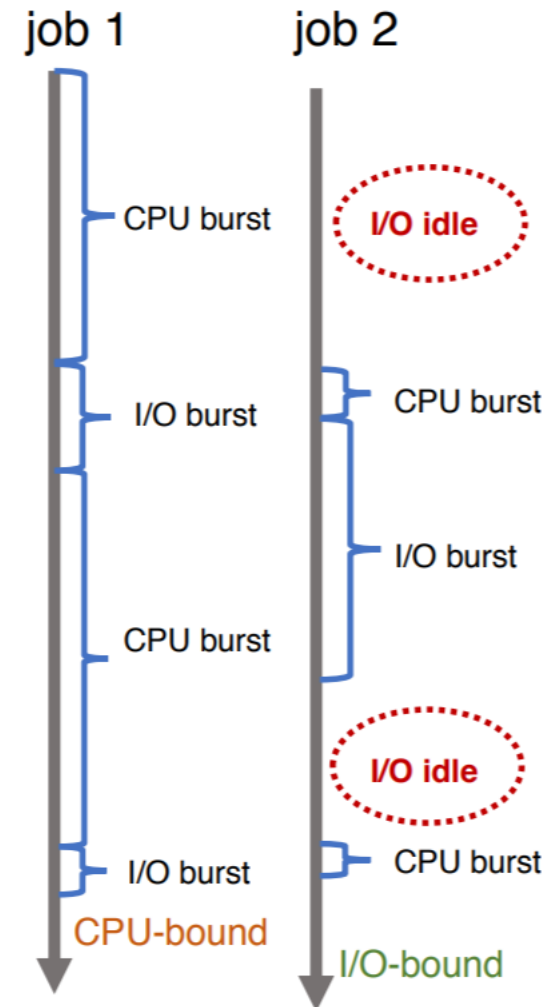
Measured value

# SJF limitations

- Doesn't always minimize average TT
  - Only minimizes waiting time
  - How to improve TT of SJF?

- Sometime, impossible to know size of CPU burst ahead of time
  - Like choosing person in line without looking inside basket/cart

- What is a critical problem of SJF?

Can potentially lead to "Starvation"

Some jobs never get CPU
ex. 1 long task
+ many short tasks coming

# Revisit convoy effect in CPU-bound vs. IO-bound jobs

- CPU-bound jobs will hold CPU until exit or I/O
  - But I/O burst for CPU-bound job is small
  - Long periods where no I/O requests issued, and CPU held
  - Result: poor I/O device utilization
- **Simple solution: run process whose I/O completed**
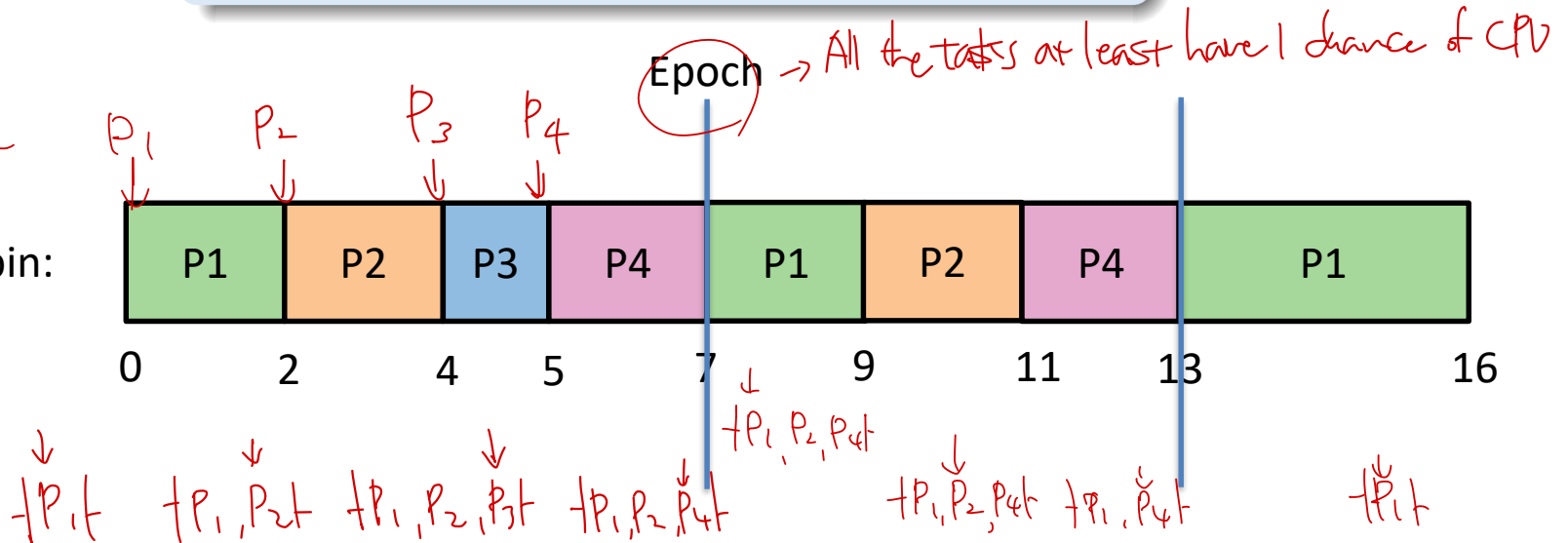  - What is a potential problem?

# Round Robin

- Solution to ==avoid the starvation== problem in SJF

- Each task gets resource for a ==fixed period of time== (time quantum or time slice)
  - If task doesn't complete, it goes back to FIFO queue

# Round Robin

Time slice: 2

| Process | Arrival Time | Burst Time |
|---------|--------------|------------|
| $P_1$ | 0 | 7 |
| $P_2$ | 2 | 4 |
| $P_3$ | 4 | 1 |
| $P_4$ | 5 | 4 |

Epoch → All the tasks at least have 1 chance of CPU

Arrive   $P_1$   $P_2$   $P_3$   $P_4$

Round Robin:

| P1 | P2 | P3 | P4 | P1 | P2 | P4 | P1 |
|----|----|----|----|----|----|----|----|

0   2   4   5   7   9   11   13   16

$\uparrow P_1 \uparrow$   $\uparrow P_1, P_2 \uparrow$   $\uparrow P_1, P_2, P_3 \uparrow$   $\uparrow P_1, P_2, P_4 \uparrow$   $\uparrow P_1, P_2, P_4 \uparrow$   $\uparrow P_1, P_4 \uparrow$   $\uparrow P_1 \uparrow$
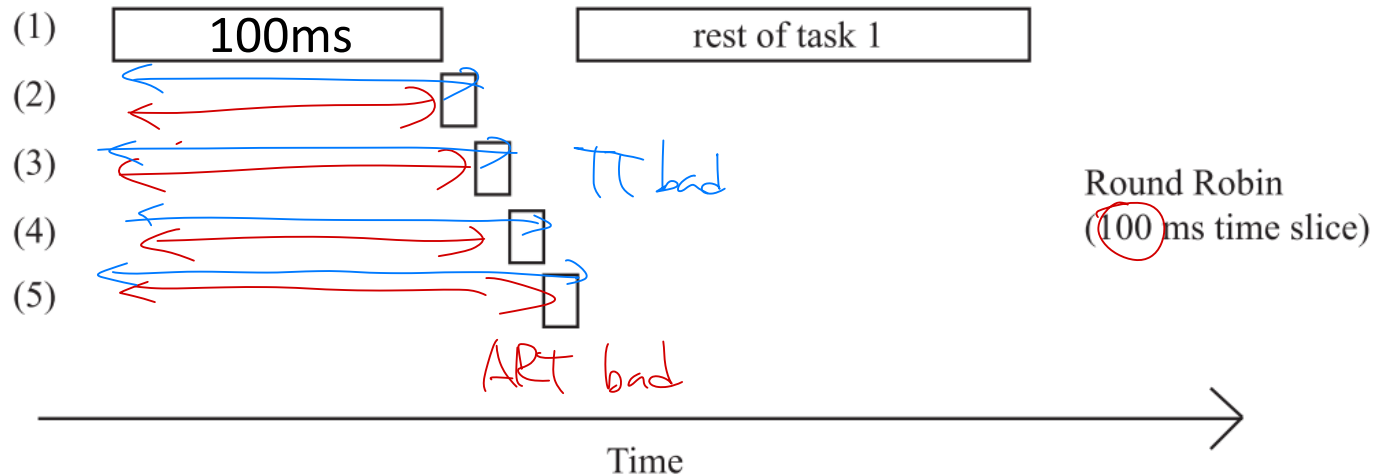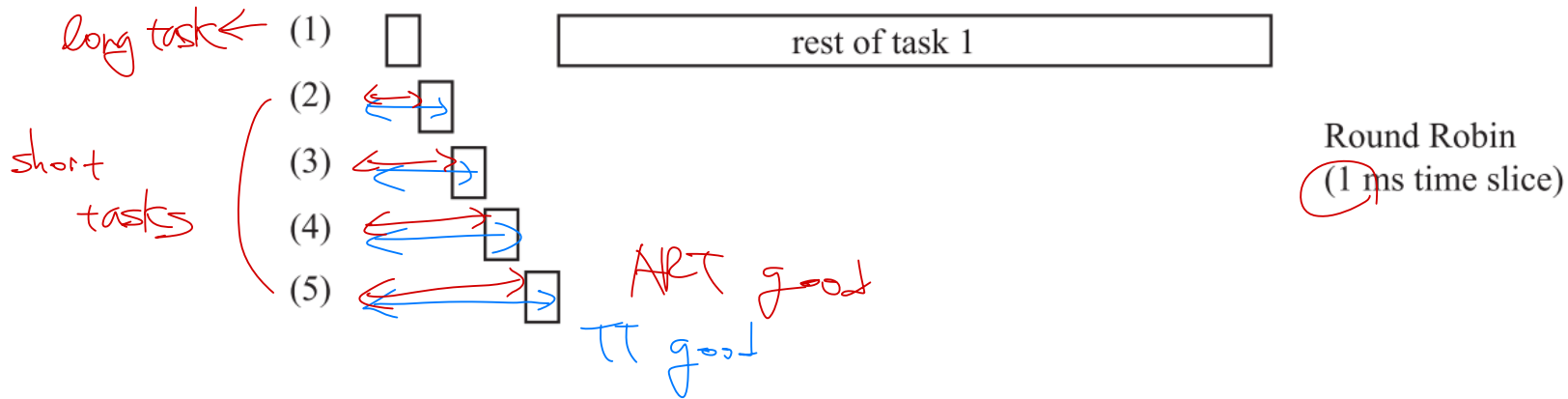
# Round Robin

- Need to pick a time quantum
  - What if time quantum is too long?
    - Infinite? → Same as FIFO
  - What if time quantum is too short?
    - One instruction? → Context switch overhead gets heavy

# Round Robin

| Time Slice / Metric | Short | Long |
|---|---|---|
| TT | better | |
| ART | better | |

Tasks

| Task | | |
|---|---|---|
| (1) | □ | rest of task 1 |

long task ←

short tasks

(2) □
(3) □
(4) □
(5) □

ART good

TT good

Round Robin
(1 ms time slice)

| (1) | 100ms | | rest of task 1 |
|---|---|---|---|

(2)
(3) TT bad
(4)
(5)

ART bad

Round Robin
(100 ms time slice)

Time

**In terms of scheduling performance metric, what is different?**

# Time slice summary

Longer or shorter?

CPU bound task prefer ( *longer* ) time slices

*long CPU burst*

*→ Usually server*
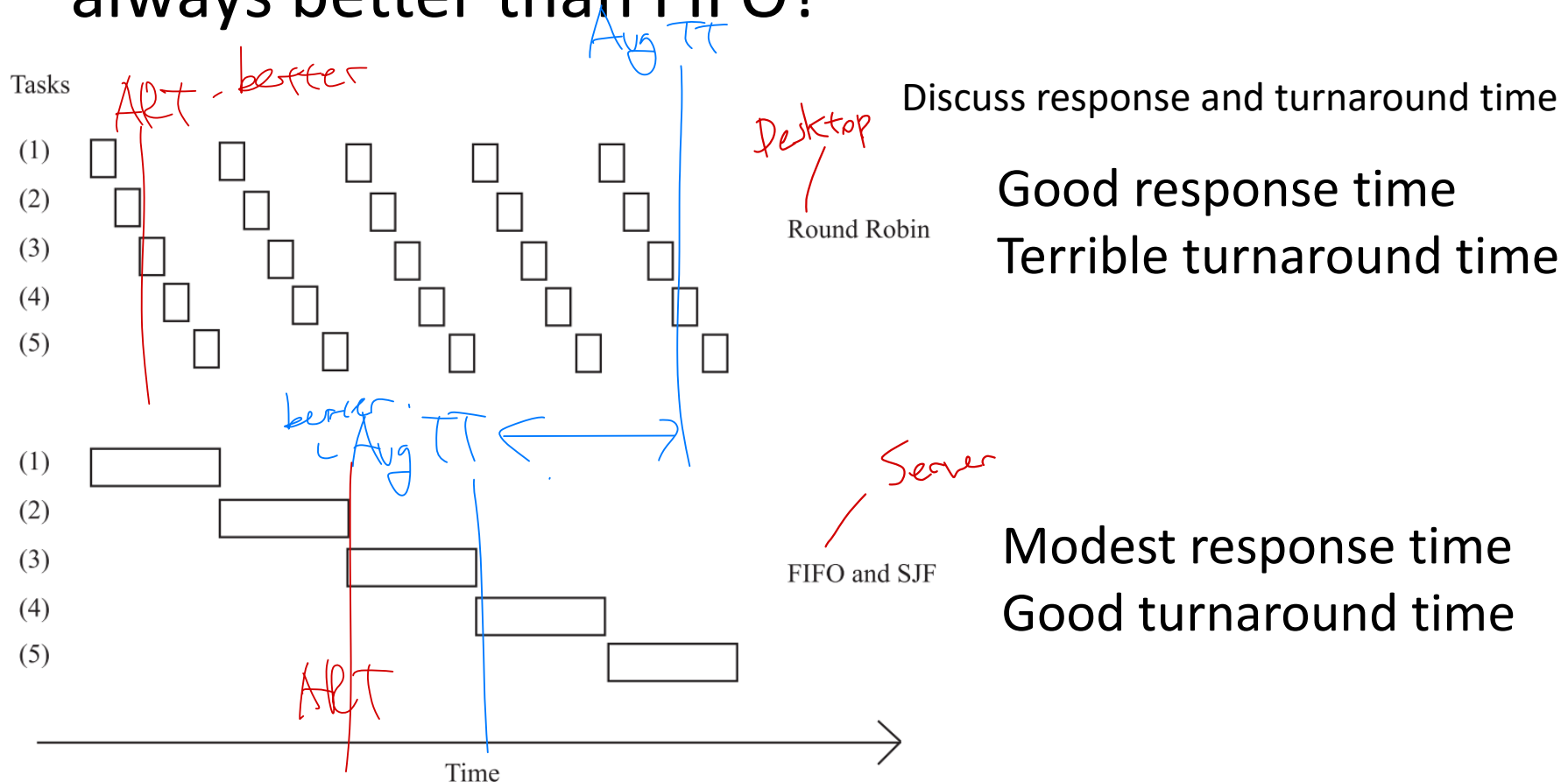
IO bound task prefer ( *shorter* ) time slices

*Short CPU burst*

*→ More responsive*

*→ Usually Desktop*

# Round Robin vs. FIFO

*Handwritten (top):* TT → FIFO better    Response Time → Round Robin

- Assuming zero-cost time slice, is Round Robin always better than FIFO?



*Handwritten annotations over diagram:* ART - better, Avg TT, Desktop, better Avg TT, Server, ART

Tasks
(1)
(2)
(3)
(4)
(5)

Round Robin

FIFO and SJF

Time

Discuss response and turnaround time

Good response time
Terrible turnaround time

Modest response time
Good turnaround time

# Round Robin = Fairness?

- Is Round Robin always fair?

$P_1$   4 CPU burst

$P_2$   2 CPU burst

- What is fair?
  - FIFO?  →선착?
  - Equal share of the CPU?  → 같은 시간 사용?
  - What if some tasks don't need their full share?
  - ==Minimize worst case divergence?==
    - Time task would take if no one else was running
    - Time task takes under scheduling algorithm

# Mixed Workload



**How to define fairness in mixed workload?**

# Max-Min Fairness

- How do we balance a mixture of repeating tasks:
  - Some I/O bound, need only a little CPU
  - Some compute bound, can use as much CPU as they are assigned
- One approach: *maximize the minimum allocation given to a task*
  - If any task needs less than an equal share, schedule the smallest of these first
  - Split the remaining time using max-min
  - If all remaining tasks need at least equal share, split evenly

# Max-Min Fairness example

*Widely used in network scheduling*

- Demands of 4 tasks = {1.9, 2.5, 4, 5}, capacity = 10
  - Equal share = 10/4 = 2.5
  - Share: {**2.5**, 2.5, 2.5, 2.5}
  - A task only needs 1.9 → 0.6 extra *Put first*
  - Equally distribute 0.6 to the rest three tasks
  - Each of them could have 0.6/3 = 0.2 extra
  - Share: {1.9, **2.7**, 2.7, 2.7}
  - 0.2 extra: distribute 0.2/2 = 0.1 to the rest two tasks
  - Share: {1.9, 2.5, 2.8, 2.8} Done!

# Priority scheduling

- Give CPU to the process with highest priority
  - Preemptively or non-preemptively
  - What is the priority in SJF? → Burst time

- A critical problem of priority scheduling

→ Starvation

# Priority scheduling

- Give CPU to the process with highest priority
  - Preemptively or non-preemptively
  - What is the priority in SJF?

- Priority scheduling can cause starvation
  - What if higher priority tasks keep coming to system?

- Solution
  - Aging; increase a process's priority as it waits

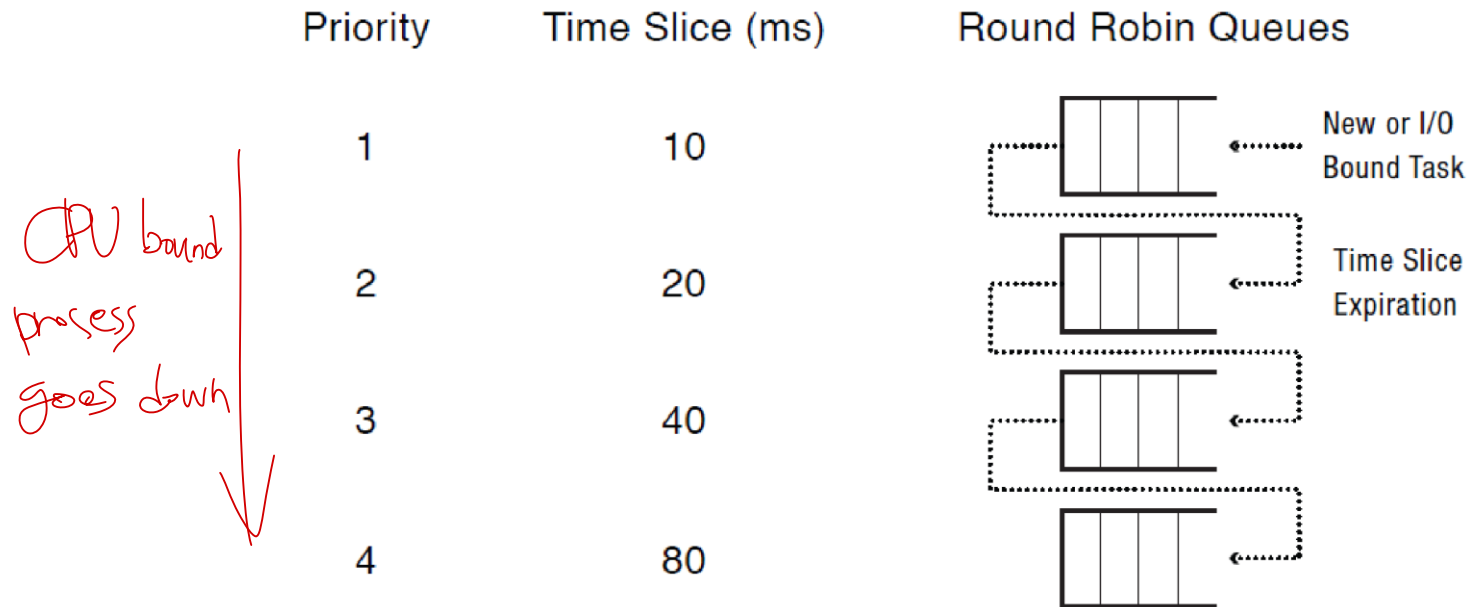# Multi-level Feedback Queue (MLFQ)

*└ Linux default scheduler*

- Goal #1: Optimize job turnaround time for "batch" jobs
  - Shorter jobs run first → *Long time slice for batch jobs*
  - Why not SJF? → *Starvation?*
- Goal #2: Minimize response time for "interactive" jobs

  → *Short time slice for interactive jobs*

- Challenge: - No a priori knowledge of what type a job is, what the next burst is, etc.

- Idea: - Change a process's priority based on how it behaves in the past ("feedback")

# MLFQ

*Good TT & Response Time*

- Set of Round Robin queues
  - Each queue has a separate priority
- High priority queues have *short time slices*  — *Interactive*
  - Low priority queues have *long time slices* — *batch*
- Scheduler picks the first thread in highest priority queue
- Tasks start in the highest priority queue
  - If time slice expires, task drops one level

# MLFQ example

| Priority | Time Slice (ms) | Round Robin Queues |
|----------|-----------------|--------------------|
| 1 | 10 | New or I/O Bound Task |
| 2 | 20 | Time Slice Expiration |
| 3 | 40 | |
| 4 | 80 | |

*CPU bound process goes down*

*Response time ↓ for interactive jobs*

Shorter time slice in higher priority queues. Why?
MLFQ ensures IO bound tasks to be schedule quickly. How?
The MLFQ algorithm is starvation-free? *No.*

*IO bound tasks가 계속 생겨 priority 낮은 상태 → Need aging*

# Contents

- Uniprocessor policies
  - FIFO, round robin, optimal
  - multilevel feedback as approximation of optimal

- Multiprocessor policies
  - Affinity scheduling, gang scheduling

# Multiprocessor Scheduling

- Must decide on more than which processes to run
    - Decide on which CPU to run which process

    → 어떤 process는 어떤 CPU에서 돌릴지
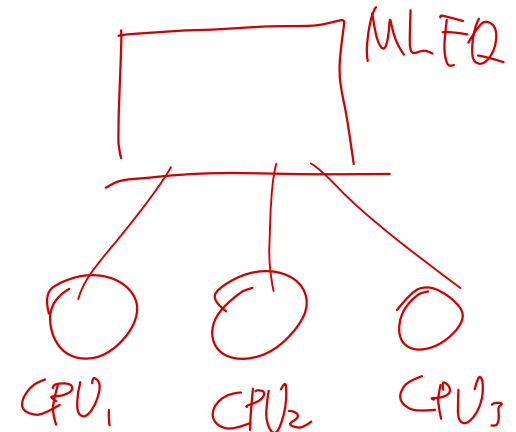
- Moving between CPUs has costs
    - Cache/TLB misses, Task migration costs
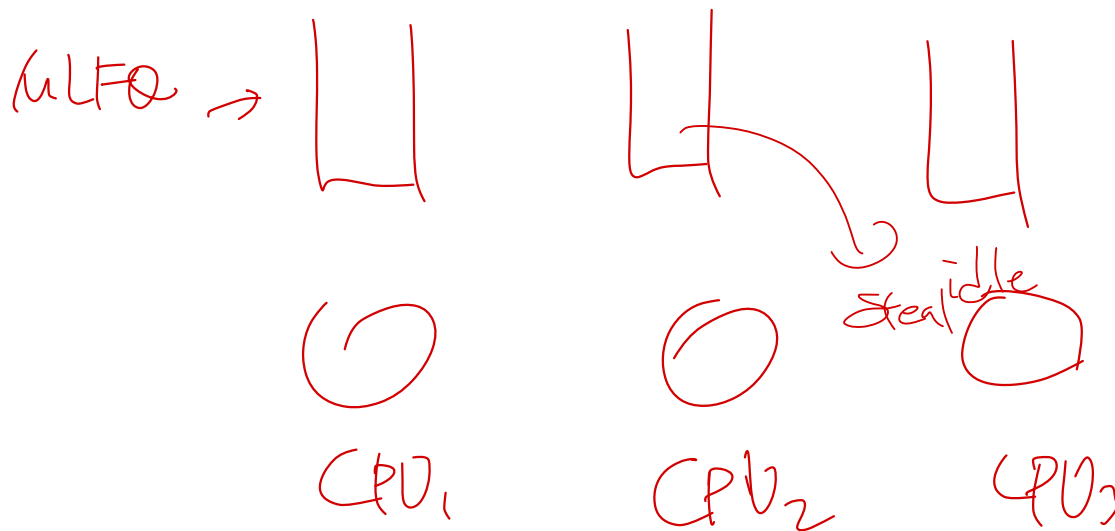
    ( Cold miss

# Multiprocessor Scheduling

- What would happen if we used a global MFQ on a multiprocessor?

  - Contention for scheduler spinlock

  - Cache slowdown due to ready list data structure pinging from one CPU to another

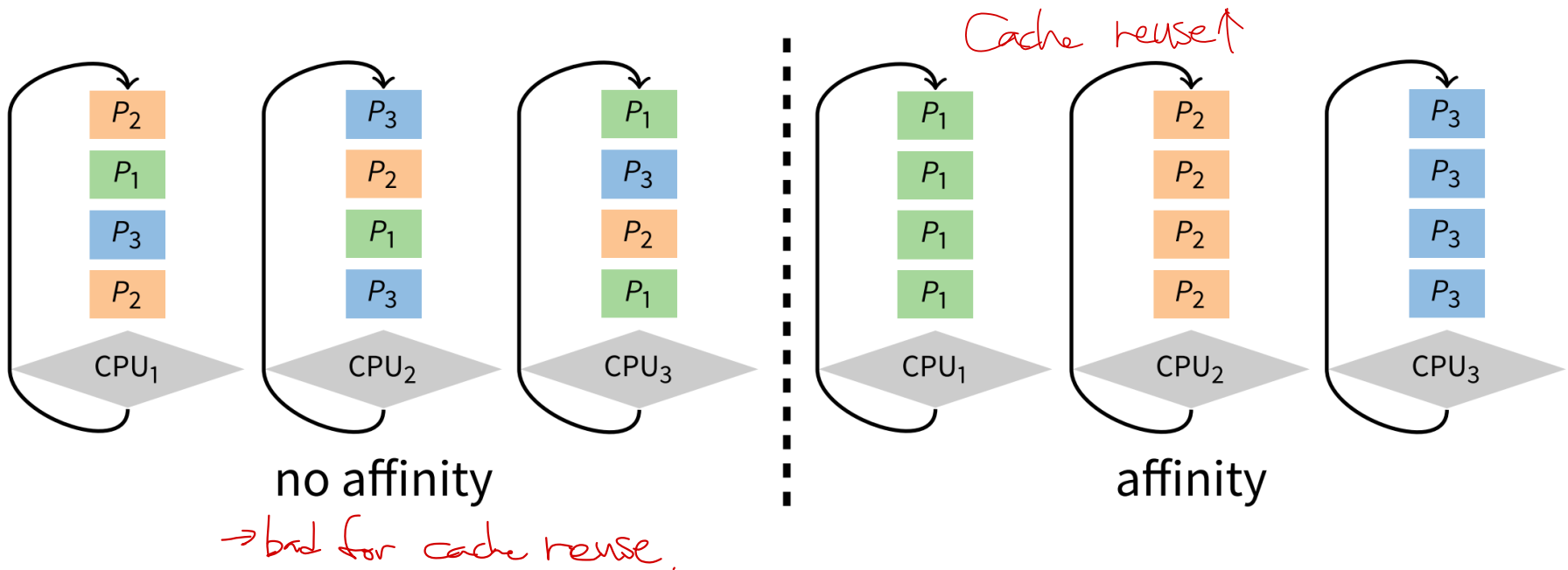  - Limited cache reuse: thread's data from last time it ran is often still in its old cache

*(handwritten annotations: "L", "Solution: → different MLFQ for each core", diagram of MLFQ connected to CPU₁, CPU₂, CPU₃)*

# Per-core ready list

- Each CPU has a per-core ready list
- Work-conserving: Idle processors can **steal** process from other processor

MLFQ →

steal idle

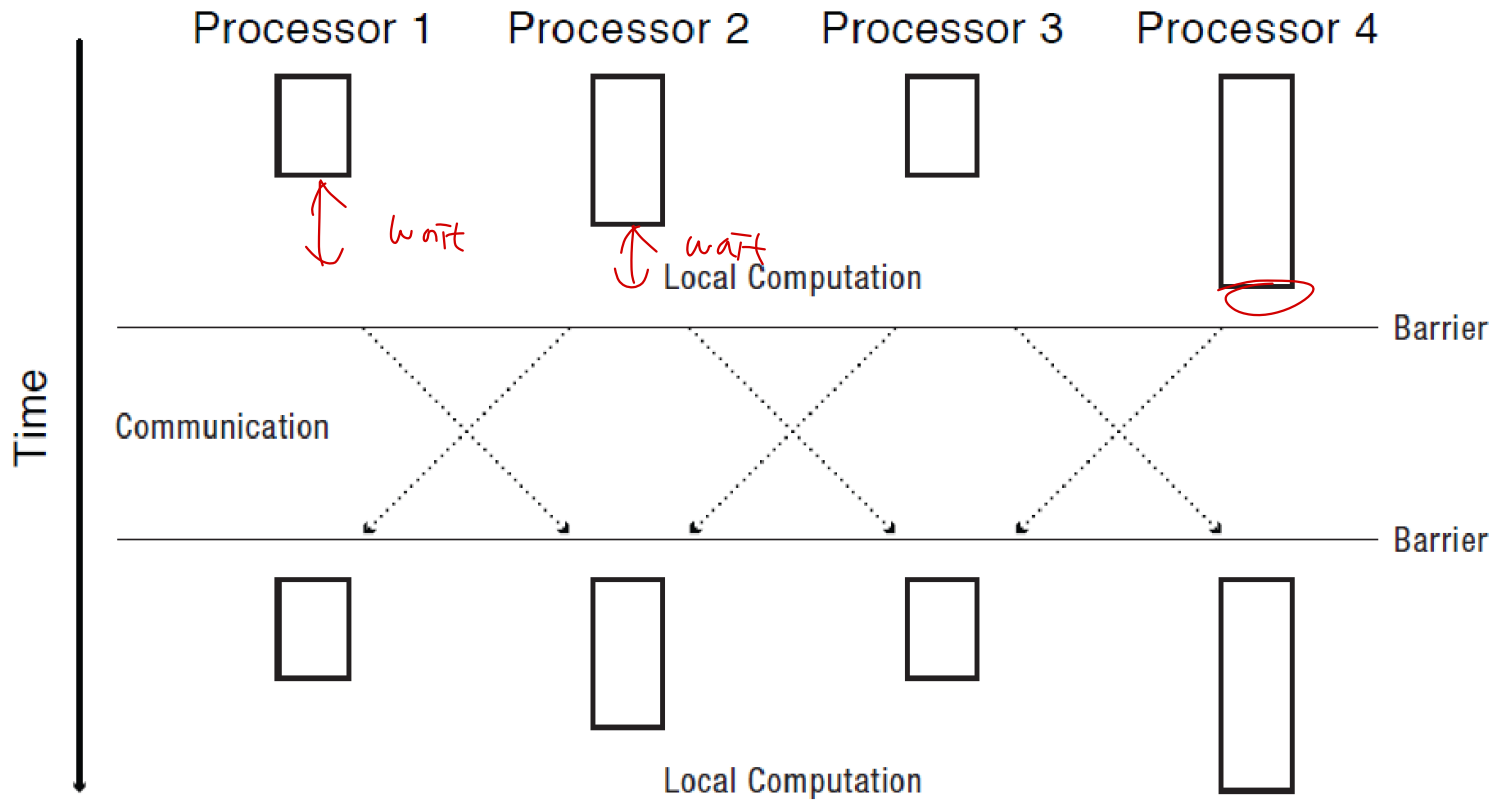CPU₁        CPU₂        CPU₃

# Per-Processor Affinity Scheduling

- Each processor has its own ready list
  - Protected by a per-processor spinlock

- Put threads back on the ready list where it had most recently run
  - Ex: when I/O completes, or on Condition->signal



no affinity

→ bad for cache reuse.

Cache reuse ↑

affinity

# Scheduling Parallel Programs

- What happens if one thread gets a short time slice while other threads have longer time slices?

  - Assuming program uses locks and condition variables, it will still be correct
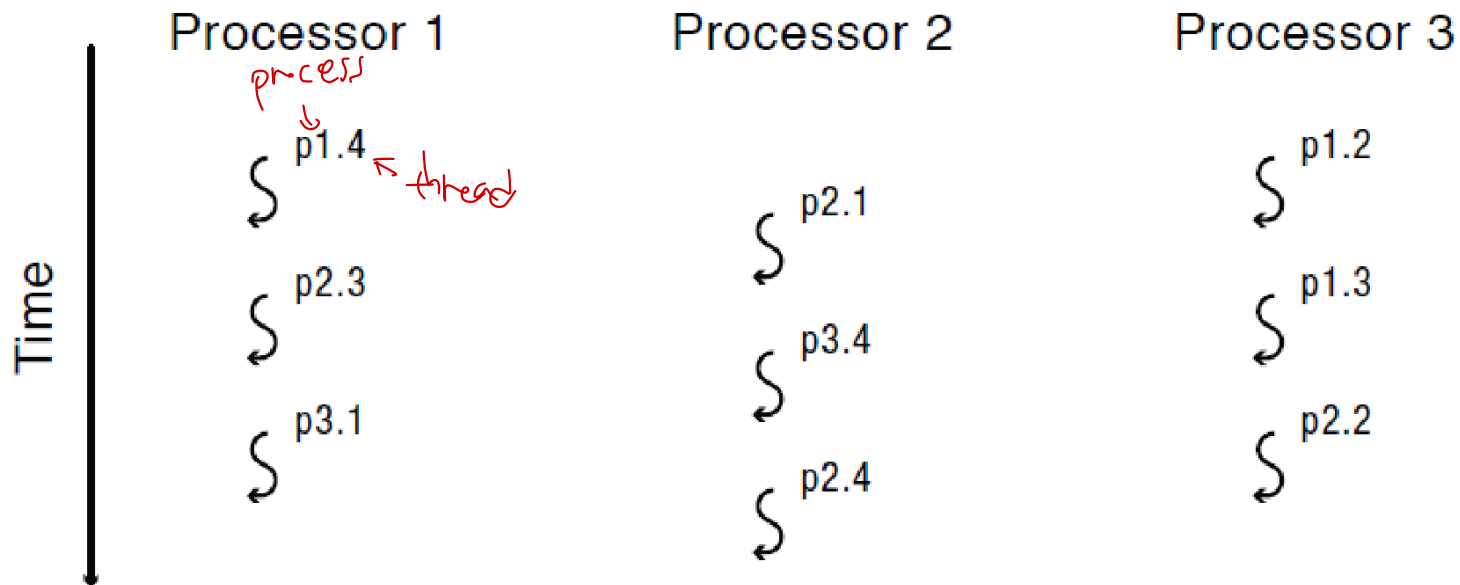
  - What about performance?

# Tail Latency
# by bulk synchronous pattern

# Scheduling Parallel Programs

의하지기 못하는

Oblivious: each processor time-slices its ready list independently of the other processors



| Processor 1 | Processor 2 | Processor 3 |
|---|---|---|
| p1.4 (process) (← thread) | p2.1 | p1.2 |
| p2.3 | p3.4 | p1.3 |
| p3.1 | p2.4 | p2.2 |

px.y = Thread y in process x

# Gang scheduling

- Schedule related processes/threads together
  - Based on communication, data sharing pattern

- Schedule all CPUs synchronously
  - With synchronized time slice, which is easier to scheduler related processes/threads together