

Chapter 13 : M/M/1 to G/P/S

13.1 Generalized Processor Scheduling

Generalized Processor Scheduling (G/P/S): generalization of a resource sharing scheme with ideal properties.

When a single job A needs to be processed, it is allocated the full processing capability of the server. (just like M/M/1) + However, as soon as a second job B arrives and requests service, instead of queuing the job, we distribute the capacity in two portion.

Equal distribution을 원하면 assign 50% of server capacity to A and rest 50% to B. The quantities w_B and w_A can be seen as "weight" and are such that $w_A, w_B \in [0, 1]$, and that $w_A + w_B \leq 1$ 100%

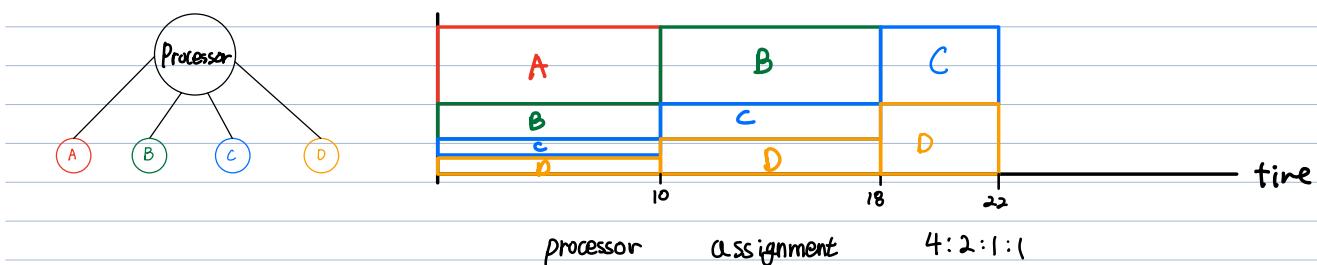
N jobs, $i \in \{1, 2, \dots, N\}$, and that $w_1 + w_2 + w_3 + \dots + w_N = \sum_{i=1}^N w_i \leq 1$. So, if the server under consideration has capacity C_S , then the capacity C_i assigned to each job be.

$$C_i = \frac{w_i}{\sum_{j=1}^N w_j} \cdot C_S$$

server capacity

그동안 얼마만큼 차지하나 이거야 / 생각해보면 꼭 맞는
걸내는 i들이 있으니, 작업에
있을 확률을 증가시킬 계산.

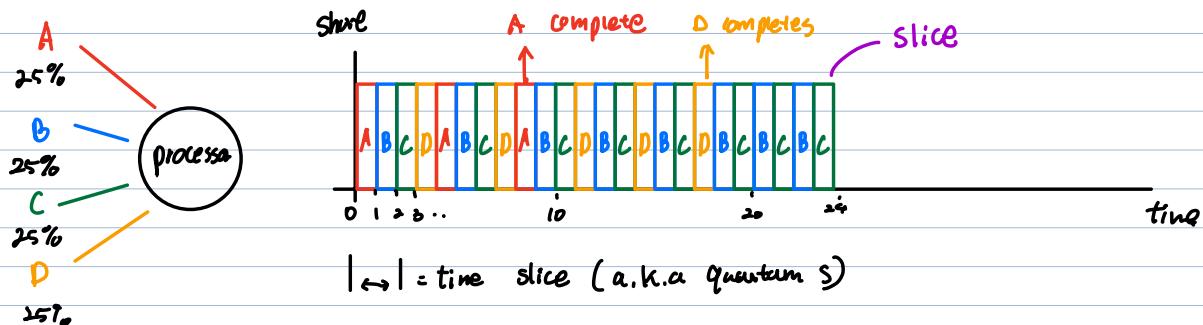
Class A	Class B	Class C	Class D
50%	25%	12.5%	12.5%



13.2 Round Robin (RR)

- No perfect way to split 50,50

So we implement time quantum | 시간대를 나누어서 돌아가면서 각각임을 수행



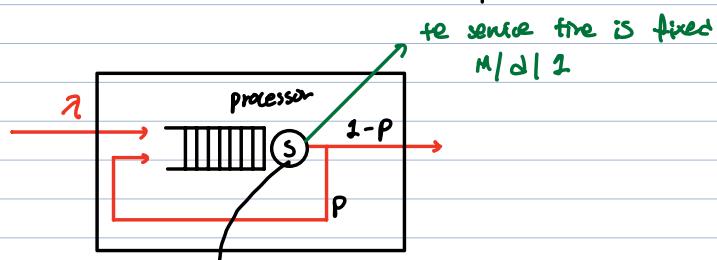
Round-Robin Scheduling

No job is picked twice until the "round" is completed

How long is each job? We don't know in principle.

→ But once it completes a slice, it may be done or not.

P	The probability of the process is not done after the slice	$1 - P$	"done"
---	--	---------	--------



끝난 확률이 $1 - P$, 다시 들어갈 확률이 P , 바로 둘의 법칙을 따른다.

$f(m) : m$ 번 이동이 다시 일어나는 re-enter : $P^m(1 - P)$ 하고 끝난 확률

그러니까, 총 일은 $m+1$ 번 해야 한다.

$t = (m+1) \cdot S \rightarrow m$ 번 차지되어 가고 꼭 번 일의 총동량
 ↓ Processor
 (execute process)

Geometric distribution에서 평균 값은

$$E[m] = \frac{p}{1-p}$$



이걸로 service time, T_s 구할 수 있다.

$$T_s = E[t] = E[m+0 \cdot s] = (E[m]+1) \cdot s = \left(\frac{p}{1-p} + 1\right) \cdot s = \left(\frac{1}{1-p}\right) \cdot s = \left(\frac{s}{1-p}\right)$$

13. 3 - System Utilization

1) 무한대에 세밀하게 분할 할 수 있는 시스템에서 $S \rightarrow 0$ 에 가깝게 사용하면 m수가 급격하게 늘어나, 작업을 완료하는데 엄청나게 많은 라운드가 필요합니다.

2) p 는 1 (각각 동작률)에 가까워 전자 $(1-p) \rightarrow 0$ 에 가까워지고 geometric distribution이 exponential과 매우 비슷해집니다. 이때 exponential은 λ 와 $\frac{1}{\mu}$ 를 가집니다.

$$T_s (\text{Service Time}) = E[(m+0) \cdot s] = E[m \cdot s] + s$$

$$\text{won } S \rightarrow 0, T_s = E[m \cdot s] + 0 = \frac{1}{\mu}$$

$$\text{그리고 결국 } T_s = \frac{1}{\mu} = \frac{S}{1-p}$$

여기 알파인 $\lambda/0/1 \rightarrow M/M/1$ 시스템으로 바운다. \Rightarrow Jackson theorem 중에 다른 법칙을 사용해보자

어떻게?

slice s 을 exponentially distributed로 정하여 적용.

시즈온 드록률 λ'

$$\lambda' \cdot (1-p) = \lambda \quad \text{In other words: } \lambda' = \frac{\lambda}{(1-p)}$$

$$p = \lambda' \cdot T_s$$

$$= \lambda' \cdot s = \frac{\lambda \cdot s}{1-p}, \text{ since } \mu \cdot s = 1-p \quad \frac{\lambda \cdot s}{\mu \cdot s} = \frac{\lambda}{\mu}$$

말에 내용은 모르겠어

Chapter 14 Basic Resource Scheduling

common strategies to perform scheduling for a single resource

14.1 Average vs Actual Case

여전히 Average를 찾으나, process-centric에서 보려면 "resource management"을 보자.

Resource management은 시스템에 나에게 적용

1) Resources are fairly Utilized

2) Requests are classified according to some metric of importance.

) 특히 heavily utilized
요인있으면 변화를
적각보록 있다.

14.2 Resource and Scheduling

항상 최고의 성능을 가질 수 있는 Scheduling은 "불가능"

성능측정 기준과 관리 차원에 따라 다르게 판斷

단일 작업

1) Stateful & Stateless

2) preemptive & non-preemptive

- Stateless

- Amount of time to service a given request is independent from the history request
- no resource is truly stateless. → 전부 요청이 새 요청에 무시할 수 있을 정도로 영향이 적을 때

- Stateful

- the time it takes to satisfy a request does depend on previously serviced.

CPU : stateless. Cache: Stateful.

- Preemptive (설정 가능한 작업): Resources that can be temporarily interrupted to serve a different request

- Non-preemptive (비설정 가능한 작업): Resources that cannot be interrupted while servicing a request.
ex) printer, modem (GPU)

- Scheduling Overhead: 스케줄링 결정에 필요한 "Extra computation"

- Scheduling Event: scheduler가 도록 하는 시점, 다음 Schedule이 시작되는 것

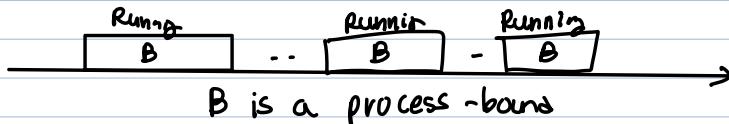
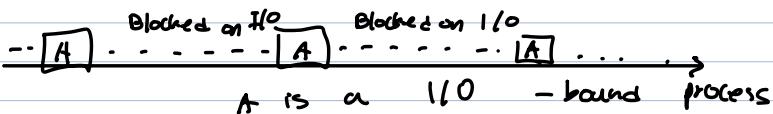
Scheduling policy: 어떤 작업을 다음에 실행할지 결정

Process states

- 1) Ready
- 2) blocked waiting

Bound processes — job (process)

- 1) Processor bound: 같은 job과 같은 I/O 상호작용, 대수로 프로세서에서 작업을 수행 ex) 수학계산, 메모리 처리
- 2) I/O-bound: 같은 job과 같은 I/O 대기 시간을 갖는다. 읽기, 쓰기, 네트워크 통신



Job and Task

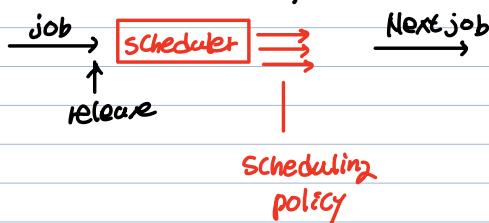
- Job: An interval of execution that can be entirely conducted on the processor.



- Task: a **sequence** of jobs

* process-bound, I/O-bound는 상대적인 것. GPU와 GPU와 CPU로 상대적인 것

Scheduler Anatomy



Scheduler
Innovation

When it wakes up and takes a decision

Scheduling policy

How the decision is taken.

process-centric metrics

— Response time (turnaround time)

$$R = t_f - t_r$$

↓
job has release time of job
Completed 2nd task t_f
작업이 완료되었을 때

— Fairness (공정성)

- job 같은 종류의 일마다 같은 대우
- fair scheduler would treat the two classes equally

System-centric metric

— Utilization : Scheduler utilization

- 스케줄러 평가의 일반적인 방향
- Scheduler can maximize processor

— Throughput (처리율)

- Rate at which requests are completed

— Complexity

- How difficult is to correctly implement and analyse a scheduling algorithm
- Implementation complexity impacts how much memory the scheduler will consume to compute its policy. 예를 들어 큐를 발생시킬 때, 큐 버퍼를 사용함.

우리가 work-conserving off는 일을 막론

14.4 General purpose Scheduling policies



$j_{i,k}$ i Task kth job

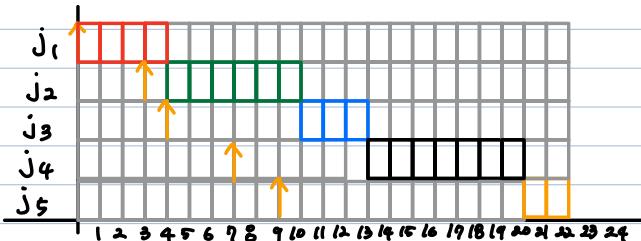
job	Arrival Time	length
j ₁	0	4
j ₂	3	6
j ₃	4	3
j ₄	7	7
j ₅	9	2

14.4.1 First come First serve (FCFS) NP First come (Earliest Arrival time)

1. Invocation: No job in the system. a new job has been released **NON-preemptive**
 2. policy: jobs are sorted in ascending order by arrival time. **The job that has been waiting longer has selected.**

↑ job Arrival

job	Arrival Time	length
j ₁	0	4
j ₂	3	6
j ₃	4	5
j ₄	7	7
j ₅	9	2



Average Response time

$$R_1 = 4 - 0 = 4$$

$$R_2 = 10 - 3 = 7$$

$$R_3 = 13 - 4 = 9$$

$$R_4 = 20 - 7 = 13$$

$$R_5 = 22 - 9 = 13$$

$$\bar{R} = \frac{1}{5} \sum_{i=1}^5 R_i = \frac{1}{5}(46) = 9.2$$

Fairness?

How fair FCFS is for the shortest job?

USE Slowdown to Reason about fairness.

$$\text{Slow down for } j_5 = \frac{R_5}{C_5} = \frac{\frac{13}{2}}{2} = 6.5$$

↑ response time
↓ service time

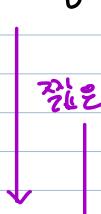
Could we Achieve a better average response time

$$\text{Ex)} \text{ Swap } j_4 \text{ and } j_5, R_4 = 22 - 7 = 15, R_5 = 15 - 9 = 6 \quad \bar{R}' = 8.2$$

Short coming of FCFS:

Once a long job is given control of the processor, all the other jobs need to wait a long time before having the chance to execute, even if their demand is small.

질문은 작업이 너무 오래 기다릴 수 있으니 Quantum [시간 차] 를 주는 건 어떤가?



14.4.2 Round Robin

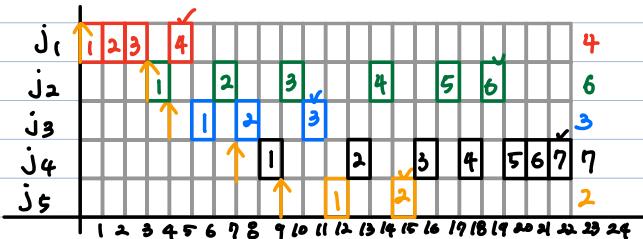
P

- Invocation: invoked after Expiration of a fixed time interval, quantum, q.
- Policy: the job that has been waiting the longest since the last time a processor quantum was assigned to it is picked next for execution.

↑ job Arrival

$$q = 1 \quad (\text{한 번 돌아갈 때})$$

job	Arrival Time	length
j ₁	0	4
j ₂	3	6
j ₃	4	5
j ₄	7	7
j ₅	9	2



Average Response Time

$$R_1 = 5 - 0 = 5$$

$$R_2 = 19 - 3 = 16$$

$$R_3 = 11 - 4 = 7$$

$$R_4 = 22 - 7 = 15$$

$$R_5 = 15 - 9 = 6$$

$$\bar{R} = \frac{1}{5} \sum_{i=1}^5 R_i = \frac{1}{5} (4a) = 9.8 \rightarrow \text{FCFS 보다 별로 편하지 않나?}$$

Fairness

How fair RR is for the shortest job?

USE slowdown to reason about fairness.

$$\text{Slow down for } j_5 = \frac{R_5}{C_5} = \frac{6}{2} = 3$$

RR Fairness (3) < FCFS Fairness (6.8)

Main Property of RR: USE preemptive **NP** to attempt fair assign of processor time to ready tasks. - prevent processor-hungry jobs from taking over the processor for long period of time.

$q \rightarrow 0$ 으로 갈수록 perfect fairness. q 가 커질수록 FCFS 랑 가까워진다. q 가 longest job in the System 보다 작으면 보통 FCFS 랑 할리된다.

Round Robin: 저분할도, 저고비헤드, 높은 Implementation. job을 작은 고정길이로 나누어
효과적

How about RR

First comes first serve

$\frac{q}{t} \rightarrow \infty$ (non-preemptive) FCFS scheduler
amount time slack for each Round Robin

$q \rightarrow 0$ perfectly fair GFS... but overhead?

Overhead Control

Say have measured the Overhead: O_m

And you want no more than 1% overhead

$$q | q \geq 100 \cdot O_m$$

↳ q 에서 O_m 의 배수이 1% 보다 적게 유지

Short Coming of RR

- 작업이 큼비되어 할당권 (q) 동안 프로세서가 실행, 끝나면 전환
- I/O bound 작업일 경우 q 가 끝나기 전에 완료, 대기시간이 아깝다.
- short job은 inefficient.



14.4.3 Virtual Round Robin → I/O bound를 조금 더 효율적으로 처리하기 위한 알고리즘

Invocation: Invoked at the expiration of a fixed time interval (q), Also invoked if the currently running job completes.

Policy: the scheduler keeps two queues.

High-Priority Queue

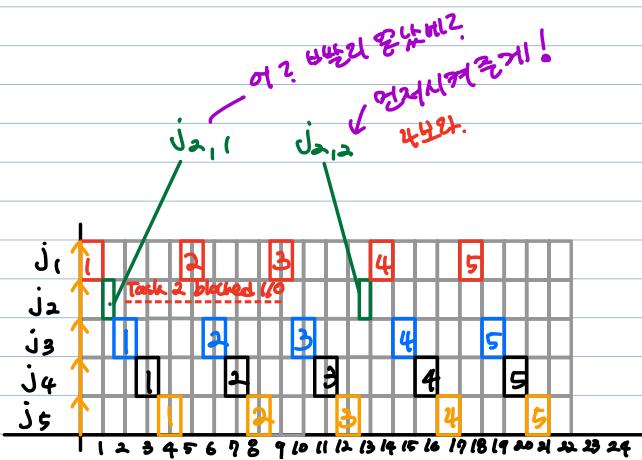
- A task enters HPQueue
- If a job belonging to that task has completed before the expiration of its assigned quantum.
- If a job in HP queue is ready then do this job first since these jobs have a high chance to become I/O bound

Low-Priority Queue

나머지는 그동안 jobs. (기장을)

↑ job Arrival

job	Arrival Time	length
j ₁	0	5
j ₂	0	1
j ₃	0	5
j ₄	0	5
j ₅	0	5



Chapter 15: Length-Aware job Scheduling

Still single resource!, but take account the length of scheduled jobs.

- 1) estimate,
- 2) predict the length of a job.

예제 1) 현재, 몇번나 죽었어 알고, 길이 몇갈지 3), 기한도 몇갈지.

이제는 Arrival time, length of a job 등과 설정을 바꿔 FCFS로 네비 시뮬레이션.

C_i | length of a generic job

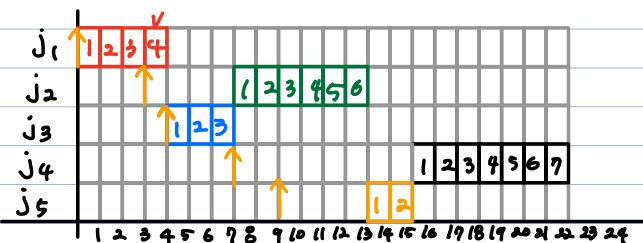
15.1.1 Shortest Job Next (SJN) FCFS 정복 4전 NP

Invocation: Completion of a job (Non-preemptive)

Policy: Read job with the **Shortest length** is picked next to execute on the processor

↑ job Arrival

job	Arrival Time	length
j ₁	0	4
j ₂	3	6
j ₃	4	5
j ₄	7	7
j ₅	9	2



Average Response Time

$$R_1 : 4 - 0 = 4$$

$$R_2 : 13 - 3 = 10$$

$$R_3 : 7 - 4 = 3$$

$$R_4 : 22 - 7 = 15$$

$$R_5 : 15 - 9 = 6$$

$$\bar{R} = \frac{1}{5} \sum_{i=1}^5 R_i = \frac{1}{5} (38) = 7.6$$

SJN : OPTimal non-preemptive strategy to minimize average response time

Fairness

Slowdown for shortest job

$$j_5 : \frac{R_5}{C_5} = \frac{6}{2} = 3$$

Short Coming of SJN: Starvation:

길이가 짧은 job이 계속 들어오면 길이가 긴 job은 영영 실행 못 시킬 수도 있다.

Blocking

j_2 = (길이)가 엄청 길리고 하면, j_5 는 짧은, 빠른 다른 job이라고 해도 j_2 아래에 blocking 되는 상황이다.

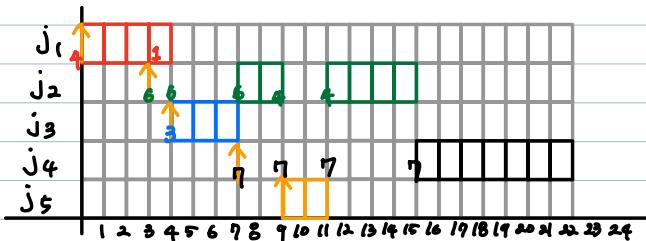
15.1.2 Shortest Remaining Time (SRT)

Invocation: at the completion of the job, or At the arrival of a new job (pre-emptive)

policy: Ready job with the shortest remaining execution time is picked next to execute on the processor

↑ job Arrival

Job	Arrival Time	length
j_1	0	4
j_2	3	6
j_3	4	5
j_4	7	7
j_5	9	2



Average Response Time

$$R_1 : 4 - 0 = 4$$

$$R_2 : 15 - 3 = 12$$

$$R_3 : 7 - 4 = 3$$

$$R_4 : 22 - 7 = 15$$

$$R_5 : 11 - 9 = 2$$

$$\bar{R} = \frac{1}{5} \sum_{i=1}^5 R_i = \frac{1}{5} (36) = 7.2$$

Fairness

Slowdown for shortest job

$$\underline{R_5} = \underline{2} - 1$$

Fairness improved.

$$J_S = C_S = \sigma - f$$

Blocking: 차단되었다.

But Starvation is Not fixed.

긴 job의 높은 실행요因地 short job만 계속 도착하면 안고져진다.

15.1.3 Highest Slowness Next (HSN) (NP)

Invocation: invoked at the completion of a job (non-preemptive) (NP)

Policy: Ready job with the **largest slowness** is picked next to execute on the processor.

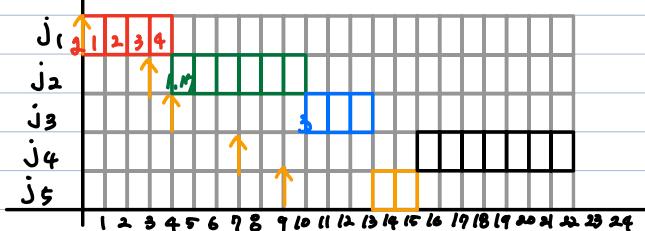
At a scheduling event, the slowness for ready tasks is calculated as the slowness that the task would have if it was executed next.

↳ 실행될 때의 slowness를 기준으로 쿠기

↑ job Arrival



job	Arrival Time	length
j ₁	0	4
j ₂	3	6
j ₃	4	5
j ₄	7	7
j ₅	9	2



4월 1주

Average Response Time

$$j_2 = 10 - 3 = 7 \quad \frac{7}{6} = 1.17$$

$$j_3 = 17 - 4 = 3 \quad \frac{3}{3} = 1$$

$$R_1 = 4 - 0 = 4$$

$$R_2 = 10 - 3 = 7$$

$$R_3 = 17 - 4 = 9$$

$$R_4 = 20 - 7 = 13$$

$$R_5 = 15 - 9 = 6$$

10월 1주

$$j_3 = 13 - 4 = 9 \quad \frac{9}{3} = 3 \quad \checkmark$$

$$j_4 = 17 - 7 = 10 \quad \frac{10}{7} = 1.42$$

$$j_5 = 12 - 9 = 3 \quad \frac{3}{3} = 1.5$$

13월 1주

$$j_4 = 20 - 7 = 13 \quad \frac{13}{7} = 1.85$$

$$j_5 = 15 - 9 = 6 \quad \frac{6}{6} = 1$$

Fairness

Slowness for shortest job

$$j_5 = \frac{5}{2} = 3$$

Slowdown 계산법 앞으그 햇을 깨끗이도 시점

$$\text{Slowdown} = \frac{t + \tilde{C}_i - a_i}{C_i}$$

t : Scheduling event

C_i : Task length.

a_i : Arrival time

HSN의 특성?

15.2 Estimating Job lengths. 작업 길이를 사용하여 스케줄러의 성능을 쟁답 시킬 수 있다.

It can be derived in few ways

1. Static analysis:

Traditional compiler application, job length analyze at compile-time or after compilation.

Static analysis rely on heuristics.

2. Offline measurements:

If the set of applications that will request execution on a given system/processor is known, it is possible to measure the individual execution time.

3. Online estimation:

Usually, application produce jobs that exhibit regular patterns in their parameter.
We can observe and predict the length of subsequent jobs

most practical way

15.2.1 Simple Average

Observe $C_{i,k}$ | $\frac{C_{1,1}, C_{2,2}, \dots \text{for Task } T_i,}{\text{Collected}}$ | $\bar{C}(n) - \text{predicted length of } n$

length of job 1
Actual

$$\bar{C}(n) = \frac{1}{n} \sum_{k=1}^n C_{i,k}$$

↓
try to keep to each $C_{i,k}$ value. to measure,

If a large number of samples, the equation above would require that we keep in memory the observed value of all the completed n jobs.

A more efficient way of computing the quantity is only keep track of

i) the last average value $\bar{C}(n-1)$; and n

$$\bar{C}(n) = \frac{\bar{C}(n-1) \times (n-1) + C_{i,n}}{n},$$

Problem: "long age" → have the same weight of recent observations.

15.2.2. Averages over a Sliding Window.

Window of length w. The last w observations contribute to the current estimation

$$\bar{C}(n) = \frac{1}{w} \sum_{k=n-w+1}^n C_{i,k}$$

마지막 w 만 사용

a new observation is available,

$$\begin{aligned} \bar{C}(n) &= \bar{C}(n-1) - \frac{C_{i,n-w}}{w} + \frac{C_{i,n}}{w} \\ \text{New } &\downarrow \rightarrow \text{旧の } + \text{新しい } \end{aligned}$$

Sliding window Average





Actual $\rightarrow w=3$ $w=10$ $\rightarrow w=15$

w 값이 작을수록: 가능한 것 만큼 크게 반응. 예상치랑의 차 비슷

15.2.3 Exponentially weighted Moving Averages (EWMA)

지속 가중 이동 평균

$\alpha \in [0, 1]$ $0 < \alpha \leq 1$ 의 weight / keep track of the system's past but recent ones

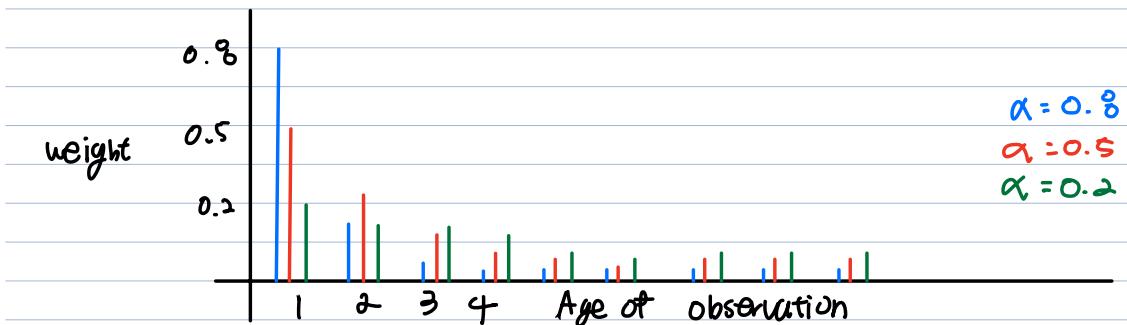
$$\bar{C}(n) = \alpha C_{i,n} + (1-\alpha) \bar{C}(n-1)$$

Why this called exponential?

$$\begin{aligned} \bar{C}(n) &= \alpha C_{i,n} + (1-\alpha) [\alpha (C_{i,n-1} + (1-\alpha) \bar{C}(n-2))] \\ &= \alpha [C_{i,n} + (1-\alpha) C_{i,n-1}] + (1-\alpha)^2 \bar{C}(n-2) \end{aligned}$$

After k steps,

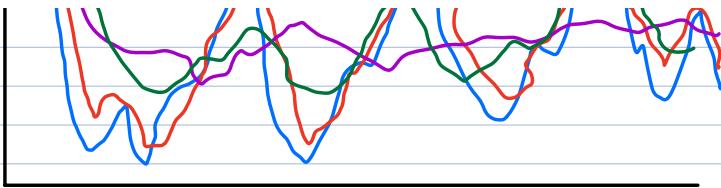
$$\begin{aligned} \bar{C}(n) &= \alpha [C_{i,n} + (1-\alpha) C_{i,n-1} + (1-\alpha)^2 C_{i,n-2} + \dots + (1-\alpha)^k C_{i,n-k}] + \\ &\quad + (1-\alpha)^{k+1} \bar{C}(n-(k+1)) \end{aligned}$$



Higher α follow better the observations, but are more influenced by temporary fluctuations
Smaller α yield to predictions that lag behind the real values

α 가 높으면 더 최근의 observation, but 투자 outlier에 더 반응
 α 가 작으면 더 오래전 과거까지 반영





Actual — $\alpha = 0.8$ — $\alpha = 0.5$ — $\alpha = 0.2$

- Priority Based Scheduling

Priority

FCFS: Inversely proportional to arrival time, higher the arrival timestamp, the lower priority

SJN: Inversely proportional to C_i

HSN: Priority of j_1 is directly proportional to the $\frac{t + C_i - a_i}{C_i}$

Chapter 16 State-Aware Resource Management

resource strategies for resources are stateful.

CPU is **stateless resource** because the time to perform context-switch / scheduling in between jobs is typically almost constant and much smaller than the length of jobs.

Stateful: 1) Traditional magnetic disk, 2) DRAM memory subsystem

16.1.1 System-centric Metrics

In context of stateless resources, utilization depends on **arrival rate & process time**. Utilization in stateless only depends on supply and demands

How about stateful Resources?

Stateless

- i) fetch new request
- ii) send request, repeat

- iii) fetch new request
- iv) recognize resource (overhead) → **Stateful Utilization**
→ overhead to recognize the resource
- v) send request

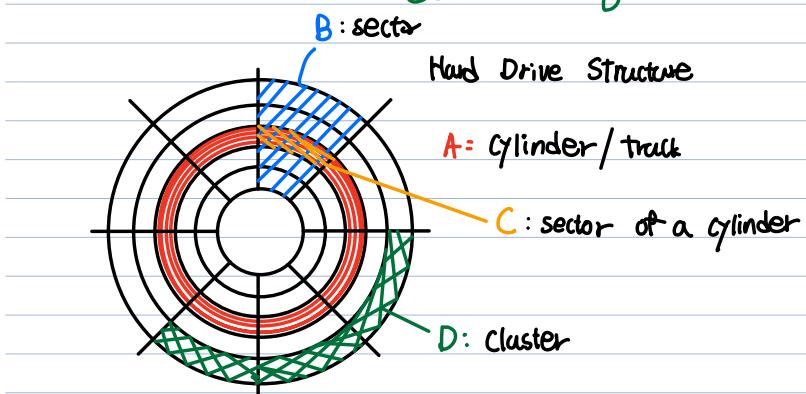
16.1.2 Process-centric Metrics

use Slabdown for fairness metrics

For a job j_i , Slabdown $\frac{R_i}{C_i}$

Key: Schedule Reordering

16.2 Hard disk Scheduling



- seek delay: Find Right **Cylinder or**
- Rotational delay: Find Right **sector.**
- Transfer time: Read/write time in the right position.

16.2.1 FCFS Scheduling

In order of arrival

Implementation ✓

Simple

Fairness ✓

No re-ordering, fairness is guaranteed

Head Movement X

when a and b requests are coming, but they are far. this case, efficiency is low

16.2.2 Shortest Scan First - SSF (가장 짧은 스캔우선)

Requests are kept in Queue, Then the position of the head after the last processed request is considered, and from the queue SSF extracts the pending request for the disk block that's closer to the disk head position

Implementation Halt

Geometry has to be known., keep track / predict of the position of the disk head to disk rotation

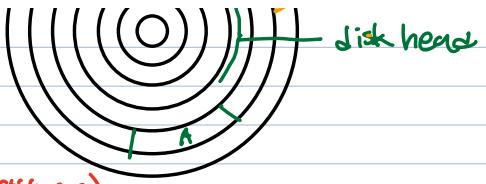


$C \rightarrow A \rightarrow B$

→ position of the disk head in the layer.

Fairness X

Unfair, as a large burst of spatially co-located requests can delay requests for far-away blocks indefinitely.
(증거로 인접한 요청들) 대량으로 인접하여 Stuck, (head stickiness)
디스크 장애로



Head Movement ✓ but could be better

- Tries to minimize head moves
- But related to arrival time. 도착 시간이 짧아도 FCFS를 수행한다.
a1, b1, a2, b2

16.2.3 Scan Scheduling

Elevator Algorithm | Constraint on the direction

Pick the request that is closer in space to the position of the head, but only one direction

Ex) 29, 27, 25, 40, At 30 up ↑ → 40 → 29, 27, 25
down

Implementation → Save as SSF Half

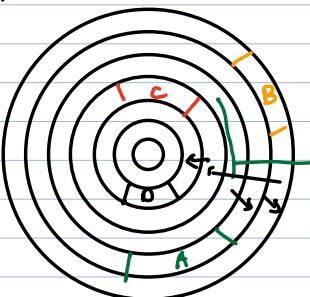
Need to keep track of current position,
predict / compute to desired head location

A → B → C → D

A B C D →

Fairness: X Still need to fix

Mono-polize & = F1 한 번만, 2번은 무시해라
증거로 Head Stickiness possible



16.2.4 C-SCAN scheduling

SCAN의 가장 큰 문제는 0에서 각각을 찾고, 다른쪽은 0이 될 때다. Worst case

0 → N, N → 0

limiting the direction of the head always on one direction

Request at 0 → 0 → 0 → 0

Other than mitigating the worst-case, identical to SCAN

16.2.5 Complete Fair Queuing Scheduler - CFQ (Linux)

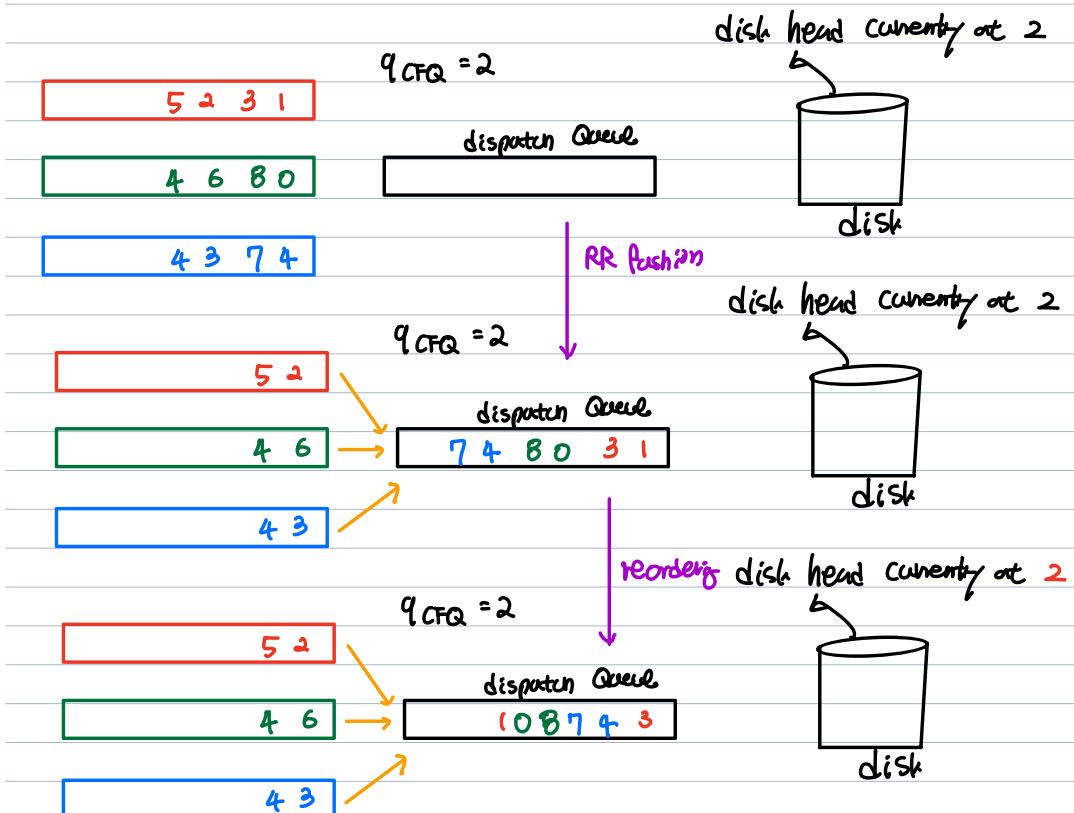
Distinguished between synchronous requests and asynchronous requests
(Requests need to be satisfied) (Can be done later without impacting performance)

Synchronous requests

2 stages

I) One queue, 9ms (Quantum) RR fashion, if empty ten wait and take ie.

2) Reorder, minimize disk head movement



Implementation X

Not Easy, maintain a set of perprocess Queue, a dispatch Queue , and perform re-ordering

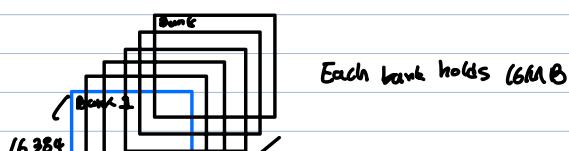
Fairness ✓

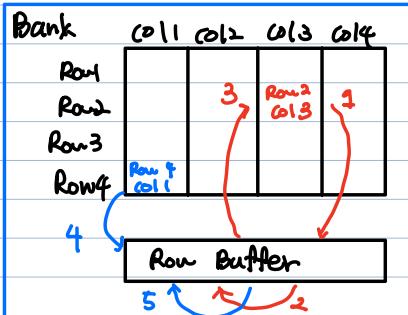
Round Robin 通过了。monopole는 아니고, reordering은 fairness와无关을 증명한 것이다.
QFCQ · N (# of tasks with pending requests)

Head Movement ✓

- Comprises minimization and fairness.
- In real words if per-application queue has less than QFCQ then CFQ non-work-conserving

16.3 DRAM Subsystem (Dynamic Random Access Memory)





At $t=0$, row buffer is empty



Operations:

1. Activate Row2
2. Read/write Col3 in Row2
3. Precharge Row2
4. Activate Row4
5. Read/write Col1 in Row4

Type A \rightarrow Chip X, bank 1, row 2, col 3

(FR, FC, FS)

16.3.1 First Ready, First come First served - FR-FCFS

A different queue is defined for each bank. FR-FC-FS keep track of which one is the row currently in the row buffer.

Implementation (in hardware) ✓

Fairly easy to implement.

keep track of currently open Row

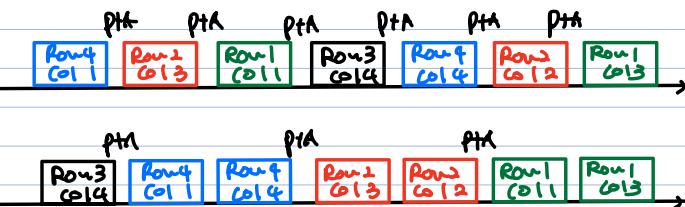
Fairness ✗

Starvation could happen

Memory throughput ✓

maximum throughput for outstanding requests

Least amount of PTA operations
(Precharge + Activ.)



16.3.2 Capped First Ready, First come First served - FR FCFS - Capped (Cap K, or reordering). No more than K younger row-hits can be re-ordered (effectively jump tie line) before older requests.

K=2 if Cut

Implementation (in hardware) Half

Tricker to implement / track of re-ordering and sufficing K

Fairness ✓



duration now fixed

Memory throughput ✓

maximum throughput for outstanding requests
Least amount of P+A operations
(Prefetch + Activ.)

Chapter 17: Real-Time Scheduling

17.1 Cyber-Physical Systems

Computation tightly with the physical world (CPS)
ex) Thermostat in home

Washing machine check temperature

Charlie card reader

Air bus system. → malfunc can lead to catastrophic

절대로 날로 작동하면 안되요.

17.2 Real Time Applications

Hard-Real Time

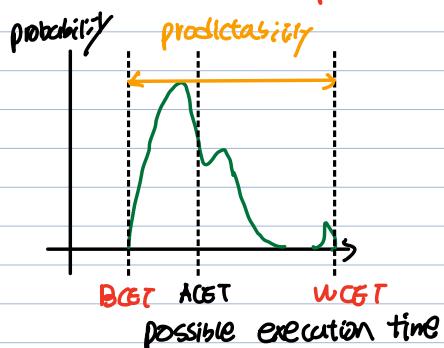
Deadlines

Absolutely **Not** okay to miss a deadline for a hard real-time task.

Predictability

Time difference between (A) best possible case and (B) worst-case behavior is

small



BCET: best-case Execution time

WCET: worst-case Execution Time

ACET: Average-case Execution time

It's okay to trade in some performance to bring BCET closer to WCET.

When analyzing real-time system, always reason in terms of WCET

17.2.2 Deadlines, Tardiness, and Utility

Soft Real-time System : Utility gradually decreases.

Hard Real-time System : Utility immediately 0 after deadline or -

17.3 Real-Time Scheduling

In Real-time Systems, it does not matter how long it will take for a job to complete as long as it is before deadlines.

Performance: Maximum work to process not passing deadline.

Capped work: Utilization bound, Utilization of ≤ 1 , response time \leq $O(1)$. / cap.

T_i | New job $j_{i,k}$ release period

C_i | WCET of a job Task T_i

D_i | deadline of a job T_i , if Job is released at t , then deadline is $t + D_i$

17.4 Static Priority Scheduling

Rate Monotonic Scheduling (RMS)

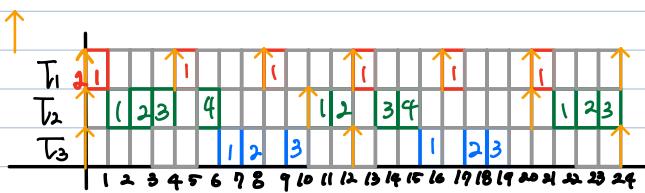
Priority: Inversely proportional to tie length

Invocation: completion of job or when a new job is released (P)

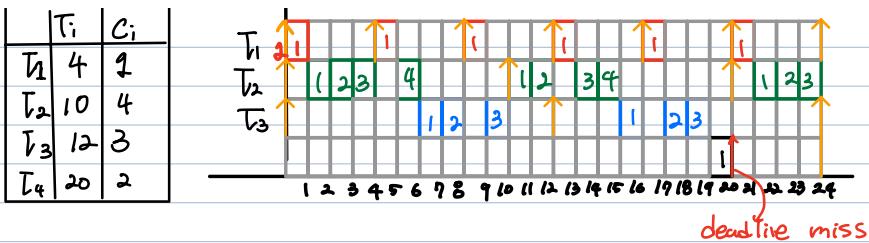
Policy: Ready job with the shortest period (higher priority) executes next on the processor.

$T_1 \geq T_2 \geq T_3$

	T_i	C_i
T_1	4	2
T_2	10	4
T_3	12	3



If we add one more task T_4 $T_4 = 20$, $C_4 = 2$.



Not schedulable under RM

RM Schedulability Test

Test whether or not a task-set is schedulable by just looking at it.

$$\sum_{i=1}^m \frac{C_i}{T_i} \leq m(2^{\frac{1}{m}} - 1)$$

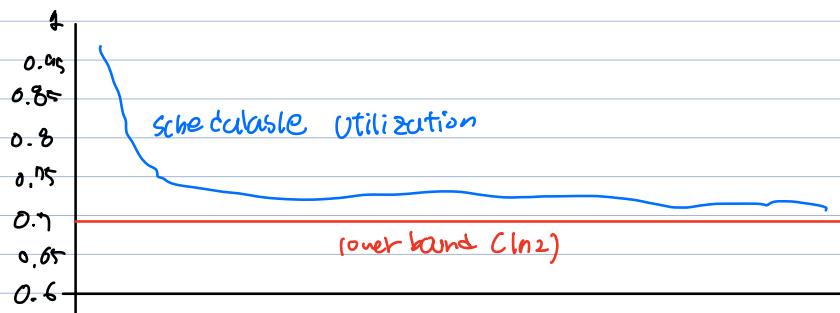
C_i : WCET

T_i : Arrival period

m : # of Tasks

※ 3개의 Task가 있던 경우, $\frac{1}{4} + \frac{4}{10} + \frac{3}{12} = \frac{1}{4} + \frac{2}{5} + \frac{1}{4} = 0.9 > 3(2^{\frac{1}{3}} - 1) = 0.18$ 방정식을 성립하지 않는다.
이면 충분조건이면 필요조건이 아님. 정확한건 그걸로 고려해야 알수있다.

$$m = \infty \quad m(2^{\frac{1}{m}} - 1) \rightarrow \ln(2) = 0.69 \quad \text{이용률이 } 69\% \text{ 이하인 경우는 RM 합격!}$$



17.5 Dynamic Priority Scheduling

프리온 방식은 Priority 를 도착할때마다 다시 계산.

Earliest Deadline First (EDF)

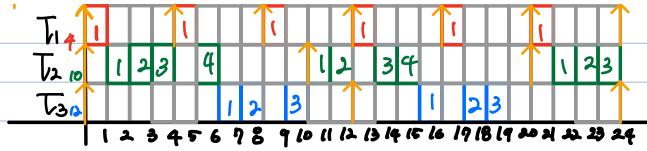
How close to deadline First

Invocation: preemptive P

Policy: Ready job with the closest deadline (highest priority) executes next on processor.

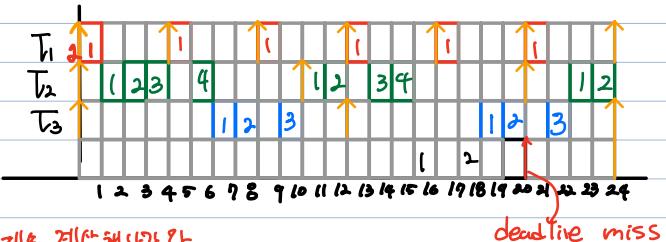


	T_i	C_i
T_1	4	1
T_2	10	4
T_3	12	3



If we add one more task T_4 $T_i=20$, $C_i = 2$.

	T_i	C_i
T_1	4	1
T_2	10	4
T_3	12	3
T_4	20	2



남은 deadline을 계속 계산해나간다.

deadline miss

Schedulability Test

$$\sum_{i=1}^m \frac{C_i}{T_i} \leq 1$$

Able to schedule a task-set with 100% utilization

17.6 Practical Scheduling

과연 EDF가 RM 보다 항상 좋을까?

1. 복잡성

2. 실제 시스템의 프로세서 사용률

3. Utilization이 100% 일 때 WCET가 조금으로 업기면 EDF에서 더 악화된다. (RM)도 영향을 받긴 한다.

$$RM \leq 69\% \quad | \quad EDF \leq 100\%$$

What would you choose for your hard real-time system?

Obviously... $RM =$

- Complexity
- Workload composition 30%를 사용하게 되면을 수 있지만, 69%를 넘어서는 경우로 스케줄링이 불가능하다.
- Robustness to tardiness, 가능한 최적화, non-critical 사용 가능하다.

실제로는 RM을 사용하는데.

Chapter 18 : Multi - Resource Scheduling

Multiple Resources scheduled

18.1 Multiple Resources

Global multiprocessor scheduling

- job sent indistinguishably to any of the available cores

partitioned multi-processor scheduling

- some jobs need to be done at certain processor.

18.2 Multiple parallel Resources

1) Compute the priority from time to time of the ready job

2) Decide which processor each job executes

18.2.1 Partitioned Scheduling

N / M / M / 2

Global

Multiprocessor scheduling

N processors / servers

1 ready queue

work-conserving

Partitioned

Multi-processor Scheduling

N processors / servers

N Ready queues

work-conserving

Partitioned

1) Offline, consider parameter and assign tasks to processors | 이전에 보낸거 계속

2) Apply single-core scheduling on each core independently | 각각의 processor 각각

in RT, schedulable is the most important!!

RM-based partitioning strategy

1) Try to assign a task, 2) check Schedulability

Tack	T_i	C_i	V_i			
T_1	10	3	0.3			
T_2	15	6	0.4			
T_3	20	6	0.3			
T_4	32	18	0.56			
T_5	8	2	0.25			

The diagram shows the assignment of tasks to processors. Task T_1 is assigned to processor P_1 . Task T_2 is also assigned to P_1 , indicated by a red arrow pointing to it. Task T_3 is assigned to processor P_2 . Task T_5 is assigned to P_2 , indicated by a blue arrow pointing to it. Task T_4 is assigned to processor P_3 .

$$T_1 + T_2 = 0.3 + 0.4 \geq 0.69 \text{ NO!} \quad \text{그치면 이런 종류조건이지 필수조건이 아님으로}$$

$$2(2^2 - 1) = 0.83 \text{ Yes!}$$

RM - FF (First Fit)

For the whole set: $U = \sum_{i=1}^m \frac{C_i}{T_i}$

A Task-set is

Schedulable

under RM-FF if

$$U \leq N \cdot (\sqrt{2} - 1)$$

processor

Tack	T_i	C_i	V_i
T_1	10	3	0.3
T_2	15	6	0.4
T_3	20	6	0.3
T_4	32	18	0.56
T_5	8	2	0.25

$\left. \right\} U = 1.81, \quad \leq 3 \cdot (\sqrt{2} - 1) = 1.24$
안可行하지.
 $1.81 \leq 1.24$

EDF - FF

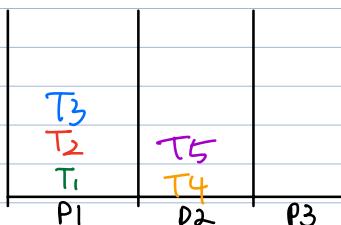
For the whole set: $U = \sum_{i=1}^m \frac{C_i}{T_i}$

A Task-set is

Schedulable

under EDF-FF if

$$U \leq \frac{\beta \cdot N + q}{\beta + 1}$$



$$\beta = \lfloor \frac{1}{\max_i \frac{C_i}{T_i}} \rfloor : \text{floor value, 정수}$$

Tack	T_i	C_i	V_i
T_1	10	3	0.3
T_2	15	6	0.4
T_3	20	6	0.3
T_4	32	18	0.56
T_5	8	2	0.25

$\left. \right\} U = 1.81, \quad \leq 2 \quad \text{Schedulable}$
 $\text{EDF-FF: } \frac{2 \cdot 3 + 1}{1+2} = \frac{4}{2} = 2$
 $\beta = \lfloor \frac{1}{\max_i \frac{C_i}{T_i}} \rfloor = \lfloor 1.78.. \rfloor = 1$

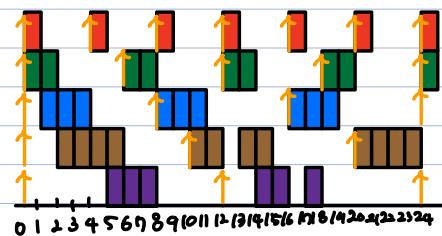
3) Partitioning scheduling 3 치고.

18.2.2 Global Scheduling

Try 가능한 수 빼기 정리.

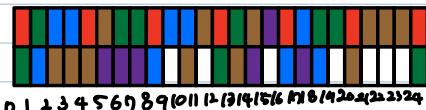
- 1) Offline, rank all the tasks/jobs in the shared ready queue
- 2) Select $\text{too } N$ to run on the N available processors.

Task	T_i	C_i	U_i
T_1	4	1	0.28
T_2	6	2	0.33
T_3	8	3	0.39
T_4	10	4	0.4
T_5	12	3	0.26

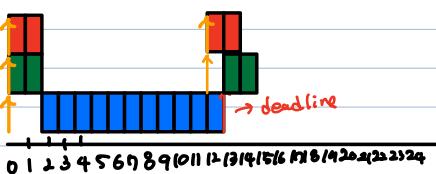


lower job ID wins

$$U \approx 1.6$$

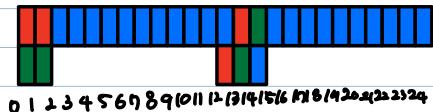


Task	T_i	C_i	U_i
T_1	12	2	0.17
T_2	12	2	0.17
T_3	13	12	0.92



lower job ID wins

$$U \approx 1.2$$



Task-set is Not Schedulable on 2 processor under Global EDF.

Some low-utilization task-sets are **not schedulable**, even with large N

The Dhall's Effect.

Chapter 19: Mutual Exclusion

19.1 Data Races

Multiprogramming

No real concurrency: interleaved execution at different resources.

Multiprocessing

real concurrency AND interleaved execution



Distributed processing

All of the above, but with possibly large communication

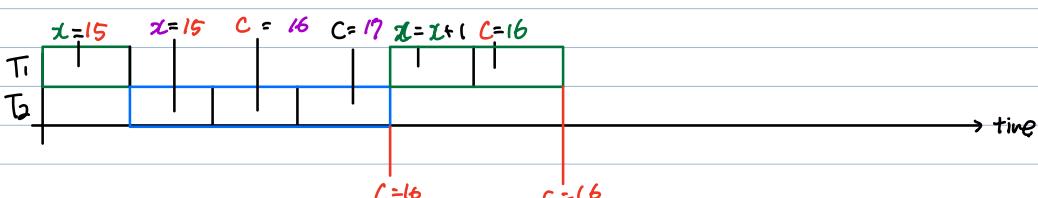
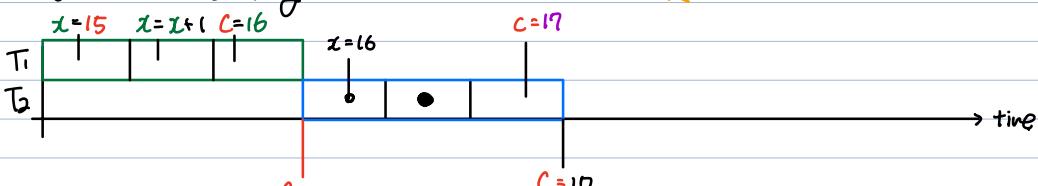
Concurrency

- 1) might need to share resources (ex) I/O devices, processor, interface
- 2) might need to share data a memory data structure, etc

```
void main() {  
    int x  
    x = C;  
    x = x + 1;  
    C = x;  
}
```

C is shared variable

Uncontrolled Sharing



Focus on Process Interleaving:

Step	Task 1	Task 2
1

Step	Task 1	Task 2
1

"Correctness" is wrong

	$x = 0$				$x = 0$	$x = C$	"Data Race" - Uncontrolled access by two or more concurrent processes to a shared variable.
2	$x = x + 1$			2	$x = x + 1$	$x = C$	
3	$C = x$			3	$x = x + 1$	$C = x$	
4		$x = C$		4			
5	$x = x + 1$			5	$x = x + 1$		
6	$C = x$			6	$C = x$		

$C = 17$ $C = 16$

19.2 Critical Sections & Code

- when a process manipulates shared resources / data, it is said to be in a **Critical section**

19.3 Mutual Exclusion (상호 배제 문제)

Sharing Resource **R** is safe if **at any time only one process is in a Critical section that manipulates R**

Assumptions about the processing Environment

- progress: None of process stops. It moves forward.
- Unknown Rate of progress: arbitrary speeds
- Concurrency exists: It could be any concurrency
- Unknown Interleaving: Any arbitrary way of interleaving
- Serialized Memory: Allow for shared-memory.

Generic Code

```
repeat:
    reminder section
    [Entry Section]
    Critical section!
    [Exit Section]
    reminder section
forever
```

Desirable properties

Mutual Exclusion

at most one process in a given critical section CS

Isolation

only processes trying to enter a CS get to decide which one enters

Progress (no deadlocks)

if no progress in CS and processes want to enter, then someone will

No starvation

no process waits an infinite amount of time before entering a CS

1)

2 Party Mutual Exclusion

Process i

repeat:

 reminder section

 [Entry section]

 Critical section!

 [Exit section]

 reminder section

 forever

index of the current process

Taking Turns

Process i

repeat:

 reminder section

 while (turn != i);

 Critical section!

 turn = j;

 reminder section

 forever

otherwise skip to

busy-loop until it's
process i's turn

once done, give turn
to process j

Solution is Not Satisfied

4)

Mutual Exclusion is satisfied, when P_j is not in used, P_i is always held by P_i .

Attempt 2: Signal Use of Critical section
Showing Signal (flag) that I'm in the Critical section

Locked Variable

Process i

repeat:

 while (locked);

 locked = true;

 Critical section!

 locked = false;

 forever

initialized to False

Step	Task 1	Task 2
1	while (locked);	
2		while (locked);
3		locked = true;
4		Critical section!
5	locked = true;	
6	Critical section!	

Mutual Exclusion is Not Satisfied

Task 1 and 2 can enter into locked at the same time since locked value is shared

giving house.

To fix this, **use one such flag per process**

Each Flag for processes

Global Variables

```
var flag[i] = false;  
var flag[j] = false;
```

Process i

```
repeat:  
    reminder section  
    while (flag[j] == true);  
    flag[i] = true;  
    Critical section!  
    flag[i] = false;  
    reminder section  
forever
```

NO, Tasks can entry to
critical section after read each others
flag and before change their flag

Step	Task 1	Task 2
1	while (flag[j] == true);	
2		while (flag[i] == true);
3	flag[i] = true;	flag[i] = true;
4		Critical section!
5		
6	Critical section!	

Attempt 3 : Signal intention to use Critical section
The mutual exclusion is invalid due to check others flag before setting my flag.

Intent my flag before read

Global Variables

```
var flag[i] = false;  
var flag[j] = false;
```

Process i

```
repeat:  
    reminder section  
    flag[i] = true;  
    while (flag[j] == true);  
    Critical section!  
    flag[i] = false;  
    reminder section  
forever
```

Might Enter a DEADLOCKS

Step	Task 1	Task 2
1	flag[i] = true;	
2		flag[j] = true;
3	while (flag[j] == true);	
4		while (flag[i] == true);
5	while (flag[j] == true);	
6		while (flag[i] == true);

Attempt 4 : Intermittently, Signal Intention to Use Critical Section
 when process is realizing that my process is not able to proceed, then set my flag free.

Set my flag free

Global Variables

```
var flag[i] = false;
var flag[j] = false;
```

Process i

```
repeat:
    remainder section
    flag[i] = true;
    while (flag[j] == true) {
        flag[i] = false;
        delay();
        flag[i] = true;
    }
    Critical section!
    flag[i] = false;
    remainder section
forever
```

Still NO DEADLOCK
 NO PROGRESS

Step	Task 1	Task 2
1	flag[i] = true;	
2		flag[j] = true;
3	while (flag[j] == true);	
4		while (flag[i] == true);
5	flag[i] = false;	
6	delay();	flag[j] = false;
7		delay();
8
k+1	flag[i] = true;	flag[i] = true;
k+2		
k+3	while (flag[j] == true);	while (flag[j] == true);
k+4	flag[i] = false;	flag[i] = false;
k+5		
k+6	delay();	delay();

Attempt 5 : Intermittently and Assertively, Signal Intention to Use Critical Section
 Previous was almost correct, but still possible for deadlock

One possibility is giving higher priorities. higher priority does not reset flag

Still Not Fair → could lead to starvation, bound waiting correctness criterion.

Dekker's Algorithm : Intermittently but Fairly, Signal Intention to Use Critical section.

대장을 깨는건 올바른 선택, 한쪽에만 대장을 깨고 있지도 않다.

Turn off flag

Dekker Algorithm

Global Variables

```
turn_flag = false;
```

```

var flag[i] = false,
var flag[j] = false;
var turn = i;

Process i:
repeat:
    reminder section
    flag[i] = true;
    while (flag[j] == true) {
        if (turn == j) {
            flag[i] = false;
            while (turn == j);
            flag[i] = true;
        }
    }
    turn or
    (if flag[i] == true)
    flag[i] = false;
    (if flag[i] == false)
    flag[i] = true;
    i += 1;
    if (i > n) i = 1;
}

Critical section!
turn = j;
flag[i] = false;
reminder section
forever

```

With Turn set to j, it's Task i's turn to back off!

After CS, release resource and switch to

Got All 4. mutual Exclusion
 (synchronization
 progress)
 No Starvation

How about Fairness.

peterson's Algorithm

which is also fair

```

process i:
repeat:
    flag[i] = true;
    turn = i; → 우선권 양보
    while (flag[j] && turn == j);
    Critical section! → j가 대기 중
    flag[i] = false;
forever

```

차이점: CS에 놓여 있는지 확인하는,
 우선권을 양보.

Mutual Exclusion: Lock / Unlock

lock

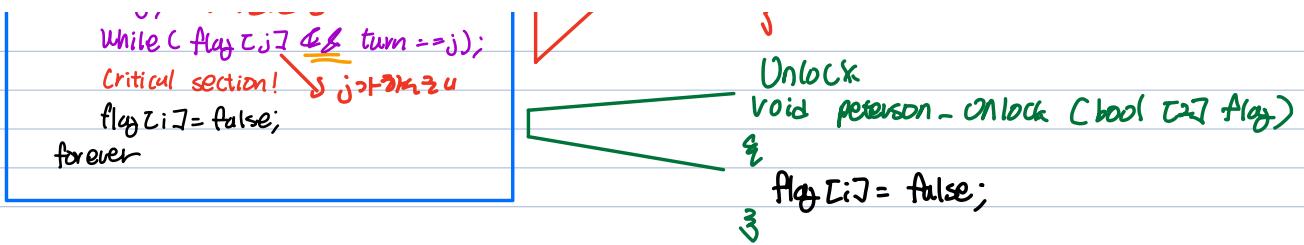
void peterson_lock (bool &flag, bool &turn)

```

process i:
repeat:
    flag[i] = true;
    turn = i; → 우선권 양보

```

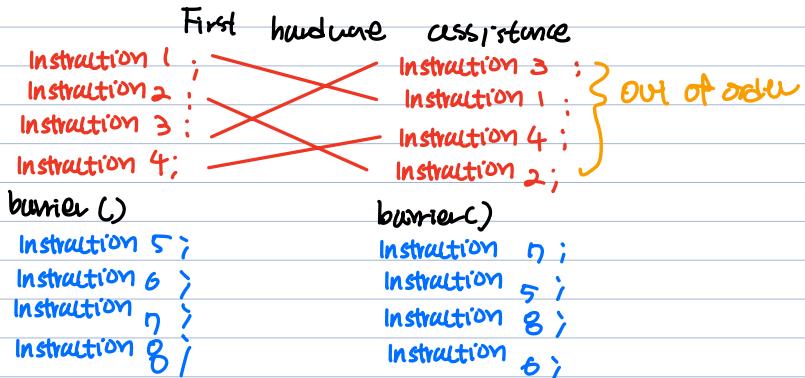
↑ flag[i] = true; → 내부 먼저 signal (intent)
 turn = j;
 while (flag[j] && turn == j);



In reality Out of Order (OOO) for optimization

Chapter 20: Hardware Assisted Synchronization

Out of order: 4개 이상의 명령을 순서로 진행된다.



Due to Out of Orders

20.1 Hardware Assisted Mutual Exclusion

There are two task:

- i) read (check of a flag)
 - ii) write / set the flag
- make it one
Atomic) do it together at once.

bool test_and_set (bool & flag)

{

 → 사용 안하는 원자연수

 if (flag == 0) {

 flag = 1;

 return true;

 } else {

 return false;

 }

Mutual Exclusion

void test_set_lock (bool & flag)

{

 while (test_and_set (flag) == false);

}

void test_set_unlock (bool & flag)

Hardware provides

" flag = 0;
3

20.2 Spinlocks

void acquire (bool & flag)

{

 while (test_and_set (flag) == false);
 barrier(); → critical section (flag)
3

void require (bool & flag)

{

 barrier(); → busy wait (flag)
 flag = 0;

3

/* some process */

Repeat;

 acquire (flag);

 critical section

 release (flag)

 remainder section

 forever

i) Using Busy wait so, it could lead to process waste.

2) For this reason, used for short period of CS

3) barrier() is critical for OOO.

20.3 Semaphores.

In order to fix processor busy wait waste from spinlock, lock the process to put in queue

- i) Attempt to acquire lock
- ii) succeeds, continue with execution
- iii) fails, blocked and enter a lock-specific queue

1 Count: int counter that captures the current status of semaphore;

2. queue: queue of processor blocked on the resources / critical section for which the semaphore is being used.

- Count > 0 : How many more tasks are allowed for concurrent usage.

- Negative value: # of tasks waiting in the queue.

- queue: reference to tasks blocks, Another job is released then Unblock waiting Task

definition

class semaphore {
 int count;
 queue queue;
 bool flag;

 void wait (semaphore & sem) {
 acquire (sem.flag);
 sem.count --;
 if (sem.count < 0) {
 sem.queue.enqueue (currProc);

3

```
# block process & unlock  
block( curr- proc , sem. flag );  
8 else,  
release( sem. flag );  
3 }
```

```

void Signal (semaphore & sem)
{
    acquire (sem.flag);
    sem.count++;
    if (sem.count < 0) {
        sem.queue.deQueue ();
        release (sem.flag);
        unblock (curr_proc);
    } else {
        release (sem.flag);
    }
}

```

{
 }
 {
 }
 }
 {
 }
 }
 {
 }
 }
}

unblock curr_proc

use queue until think goes

Synchronization Problem 1.

Simple waiting

proc 1's section 2 should only be run if proc 2's section 1 is completed.

1



2.



Proc

Initialization:

Semaphore Synch := 0

Proc. 1

repeal

Section 1

E-mail

PROC 2

repeat

Section 1
Signal search

www wwww
section 2
forever

optional copying
section 2
forever

Dating

①



②



Initialization:

Semaphore synch1 := 0

Semaphore synch2 = 0

Proc 1

repeat

 Section 1

 Signal (synch1)

 Wait (synch2)

 Section 2

forever

Proc 2

repeat

 Section 1

 Signal 1 (synch2)

 Wait (synch1)

 Section 2

forever