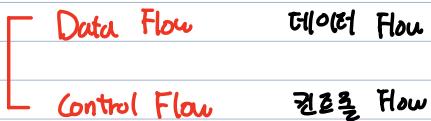


## Chapter 2

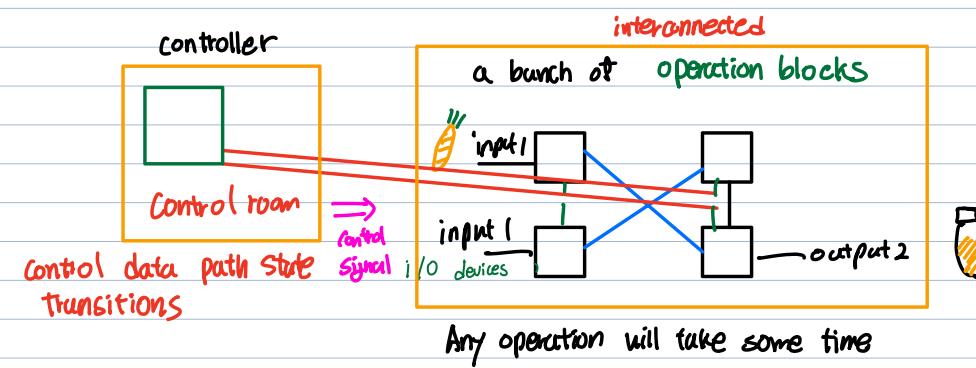
Two Design System

디자인 시스템은 2개



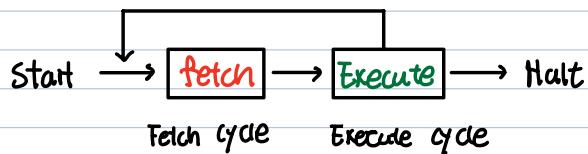
Building blocks

- Any system : beginning and end



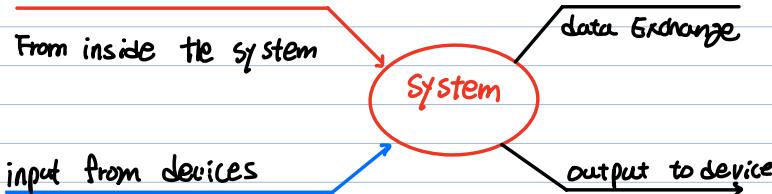
Data path (obstruction how you see set of resources data plane)

A simple controller: cycle through a sequence of states, in each state, send controls to data path.



Stuff outside of the controls?

Unexpected Events



**Synchronous** : 자원을 거치지 intervention

Exception : division by zero, misaligned instruction

Memory fault : wrong permission, wrong memory translation

SYS call : open file, create process

**Asynchronous** : 자원을 거치지 intervention

Timer fired : quantum expiration, packet send timeout : setup alarm

I/O Event : packet received at the NIC (Network interface controller), rendering completed at GPU

Hardware fault : CPU overheating, memory corruption detected

Control C

↓  
interrupt from other system

**Process** : Program in execution on a **system**, one of many

"State" is necessary

memory stored in, text, Heap, Stack

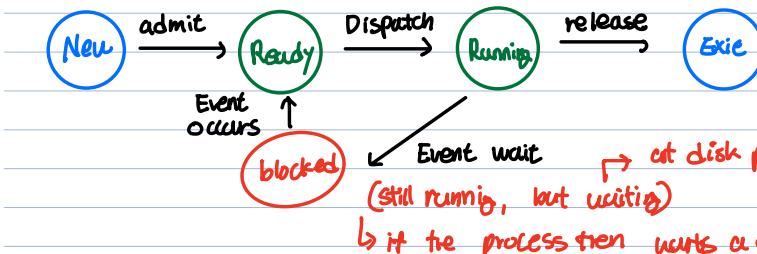
### Process

1. [Consume some resources] (e.g. CPU)
2. Request / wait for an **external** service / event
3. Do nothing until external service done (**wait**) : **A synchronous**
4. Back up!

Ready  
Running  
Block

### Transition State Diagram

System is the state of process.



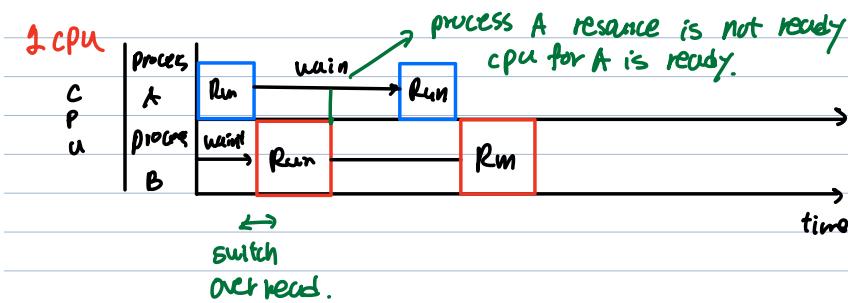
→ at disk part of diagram for text resource  
↳ if the process then wants a different resources

## Chapter 3: Processes as Resource Consumers

### 3.1 concurrency in computing system

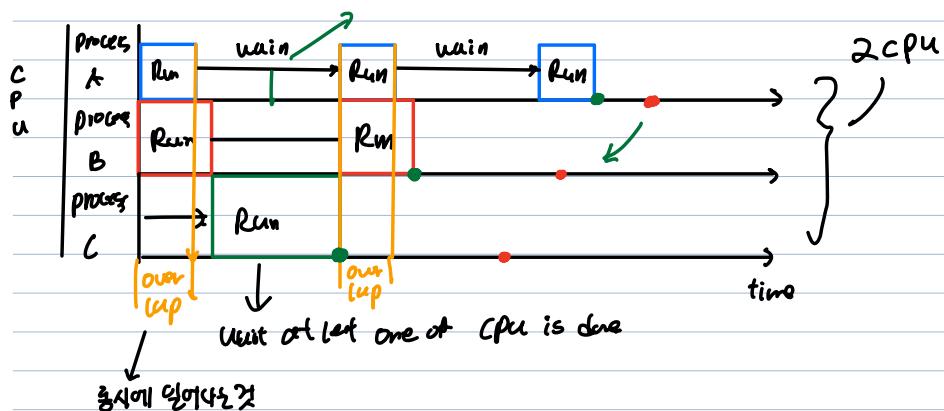
#### Multi programming (multi-tasking)

: Multiple programs must be executed on the same resource unit 노동력을 문제로

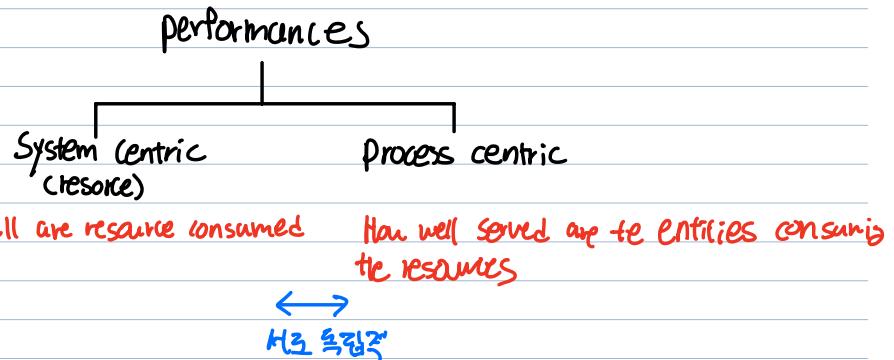


#### Multi processing

Overlapped concurrent execution



## Chapter 4: Performance Metrics and Perspectives



### System - Centric performance Metrics senior perspective

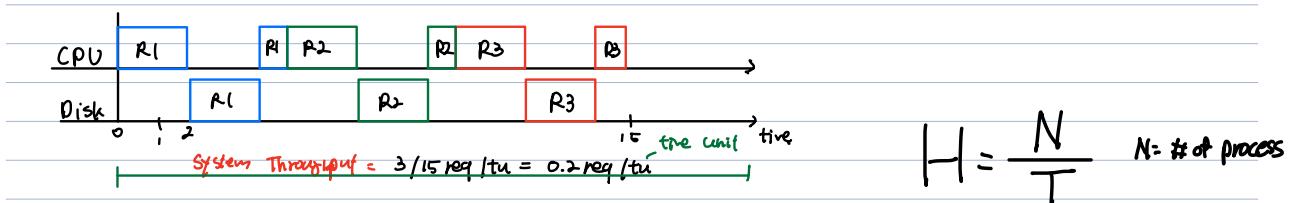
1 Utilization: a system resource S. over a period of Time T



Utilization is a **resource specific metric** → Utilization is  $\frac{\text{resource}}{\text{resource + waste}}$

2. Throughput: Number of completed requests over a given time window

To measure "holistic" performance of a system



Note: for ~~single resources~~, or for the whole system!

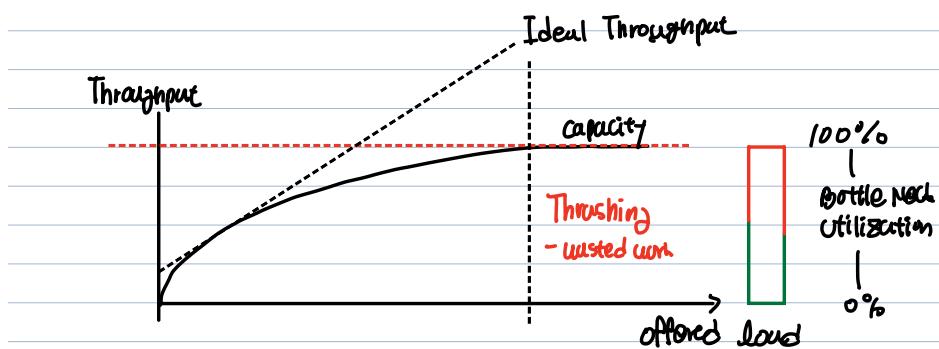
Increasing the MPL improves Utilization and throughput!

But then is this forever?

### 3. Capacity

Throughput은 request 증가에 따라 영향을 끌어내 (이유는 그림) 모델을 정확히 평가하기 어렵다.

→ Throughput이 높다고 다 안 좋은 모델이 아니다.



Demand이 늘어나면서 Throughput이 늘어나면서 ideal, 하지만 Bottleneck은 지나면 더 이상 늘어나지 않아. (Utilization이 100%가 되어있고, 더 이상 늘어날 수 없으면 bottleneck)

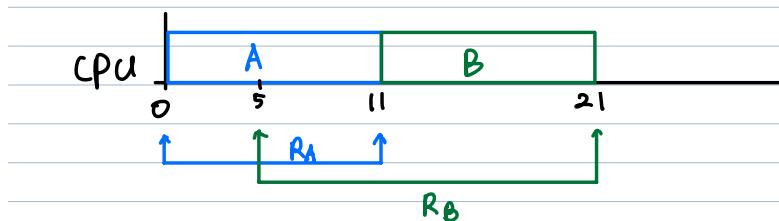
Culprit: when throughput reaches maximum

## Process-Centric Performance Metrics

n service requests. submitted  $t_1, t_2, t_s$ , finished  $t'_1, t'_2, t'_s$  respectively

### 1. Response Time

$$R_i = t'_i - t_i$$



$$R_A : 11 - 0 = 11$$

$$R_B : 21 - 11 = 10$$

$$KB: 21 - 5 = 16$$

## 2) Jitter (Response Time Variability), Variance, Standard deviation

Average Response Time  $\neq$  "real" response time

Exponential, Normal, Uniform, Gibrat distribution of  $\text{응답시간}$

## 3. Guarantee Ratio

Ex). if a deadline is imposed on the service requests, then the percentage of requests that end up meeting their deadlines is the guarantee ratio.

### 4.) Latency and Lateness

"How much later after the deadline service was completed" :  $\text{남은 시간 단위}$

### 5.) Fairness

requested execution times  $C_1, C_2, C_3$

response time:  $R_1, R_2, R_3$  fair system will keep the ratios  $R_i/C_i$

### 6.) Speed up and Slow down

System upgrade was beneficial or detrimental.

- Performance expressed in Time before the change of interest:  $T_{old}$ . ex) Average response Time

- Some metrics change under analysis :  $T_{new}$

Speed UP: How many times ( $\times$ ) faster is the new system compared to the old one

$$\text{Speed UP} = \frac{T_{old}}{T_{new}} \rightarrow \underline{\text{shorter}}$$

Slow down: Inverse of the speed up and returns a value greater than 1 indicates a performance loss

$$\text{Slowdown} = \frac{T_{new}}{T_{old}} \quad | \text{慢하므로 } \geq 1$$

Amdahl's Law: Overall performance improvements of the system.

f: fraction of time a resource is used in serving request

X: improving the performance of that resource X-fold.

Y: overall performance of the system by a factor  $\frac{T_{old}}{T_{new}}$  speed up

$$Y = \frac{1}{1-f \cdot (1-\frac{1}{X})}$$

/

parallelize proportion      number of processes

## Other Aspects of "System performance"

Real system, failures happens!

### 1 Reliability: Process-centric (User 角度)

measure of the probability that the system functions correctly continuously for a specified time period. → System 可用性

$$R(t) = e^{-t/MTBF}$$

### 2) Availability: Process-centric (Client 角度)

On-line solution is a measure of the probability that a system is available at an arbitrary point in time | → System 在线可用性

→ 99% availability, 99% reliability

1000 requests → 10% refuse 990 request → 3% failure 3% downtime.



$$\text{Availability} = \frac{\text{up Time}}{\text{Total time}} = \frac{\text{mean Time between failures (Available)}}{\text{mean Time between failures} + \text{Mean Time to Repair (dead)}}$$

$$= \frac{\text{MTBF}}{\text{MTBF} + \text{MTTR}}$$

$$MTBF = \frac{\text{Up Time (available total)}}{N \text{ fail}}$$

}

maintainability

$$MTTR = \frac{\text{down Time (System down total)}}{N \text{ fail}}$$

## Chapter 5 Performance Evaluation and Models

How do we obtain (utilization, response time, jitter, etc) from a given system?

시스템 설계자는 시스템을 모델링하고 분석할려면 자료 사용하는 데 “Queueing Analysis”

### Queueing Analysis

Customers: Individual requests for services

Queues: Waiting areas where requests for service wait.

Servers: Satisfying service requests

Distributions of Arrivals: Poisson process for the arrival of customers. The rate of arrivals  $\lambda$  # of customers coming into the system every period  $T$  is  $\lambda \cdot T$ . Arrival of requests are independent

Distribution of Service Time: Mean service time:  $T_s$  평균 서비스 시간에 걸리는 시간

$\mu$  (service rate): 단위 시간당 서비스가 처리할 수 있는 요청의 수  $T_s$ 와 반비례

$$\mu = \frac{1}{T_s}$$

Ex) 평균 서비스 시간이 2분이면, 서비스 평균적으로 1분당 0.5개의 요청을 처리할 수 있습니다. ( $\mu = \frac{1}{2} = 0.5$ )

#### 1) Single Queue, Single-Server System

Queue Size: Unlimited.

Population: independent.

1)  $W$ : # of customers waiting in the line

2)  $q$ : Total # of customers in the system

3)  $T_c$ : Computed Time in Computer

↳ 15. waiting time in queue

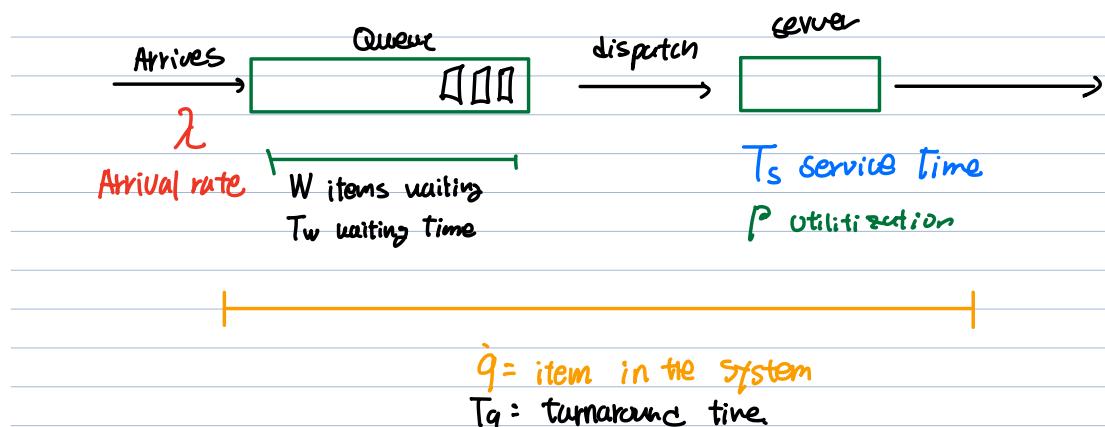
4)  $T_w$ : Waiting Time in Queue

5)  $T_q$ : turnaround time [response time]  $\frac{\text{req/sec} \times \text{sec/req}}{1} = \frac{\lambda}{\mu}$

6)  $P$ : Utilization of the system  $P = \lambda \cdot T_s = \frac{\lambda}{\mu}$  단위가 같아.

In steady state, the rate at which cannot exceed the rate at which the server is able to serve them.

$$\lambda < \mu \Rightarrow P < 1$$



## Little's Law

True for any "steady-state" queuing system.

$$q = \lambda \cdot T_q \quad \text{and} \quad w = \lambda \cdot T_w$$

## Relationships of Little's Law

$$T_q = T_s + T_w$$

$$q = w + \lambda \cdot T_s = w + P = w + \frac{\lambda}{\mu}$$

$$W = \lambda \cdot T_w$$

$$Q = \lambda \cdot (T_w + T_s) = \lambda \cdot T_q$$

- Request to a web server experiences average response time of 0.5 seconds and that the average rate of arrivals for requests is 100 requests per second., the average web server would be managing 50

$$T_q = 0.5 \text{ sec/req} \quad \lambda = 100 \text{ req/s} \quad Q = 100 \times 0.5 = 50 \text{ req}$$

- Multiprogramming Level (number of process in system) is 30 on average and that the ratio with which processes are created is 2 per minutes. then you may conclude that the average "lifet ime" of a process is  $30/2 = 15$

$$Q = 30 \quad \lambda = 2 \quad T_q = 30/2 = 15 \text{ req/min}$$

- Rate of arrival of packets to a router is 1000 per second and that it takes 0.5 msec to service each one of these packets and that the average number of packets queued at the router is 3, then "delay through this router" will be

$$\lambda = 1000 \text{ req/sec} \quad T_s = 0.5 \text{ msec/req} = 0.0005 \text{ sec/req}$$

$$W = 3 \text{ req}$$

$$T_q = T_w + T_s \quad T_w = \frac{W}{\lambda} = 3 \times \frac{1}{1000} = \frac{3}{1000} = 0.003 \text{ sec/req}$$

$$T_s = 0.0005 \text{ sec/req}$$

$$T_q = 0.003 + 0.0005 = 0.0035 \text{ sec/req} = 3.5 \text{ msec/sec}$$

$$\lambda = 100 \text{ req/sec}$$

# Chapter 7 : Basic Probability Analysis

Probability = quantification of (probable)

## Modeling Probabilistic Outcomes

Probability helps in modeling the "unknown" because outcomes are "unknown"

Random Variables

Systems Examples

- throughput, arrival rate, utilization.

- response time.

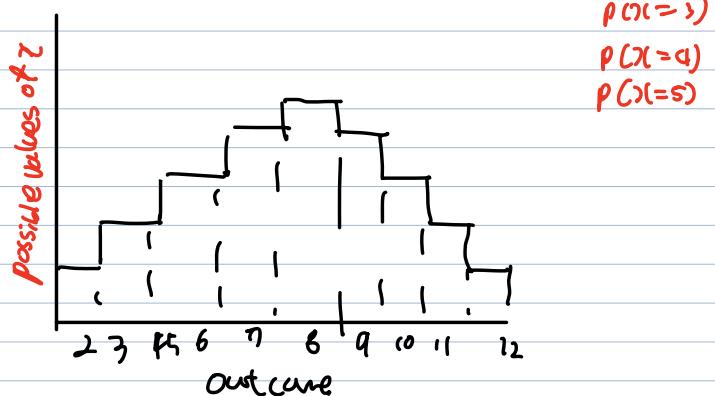
# of processes blocked at a given time

The time interval between two failures — Exponential

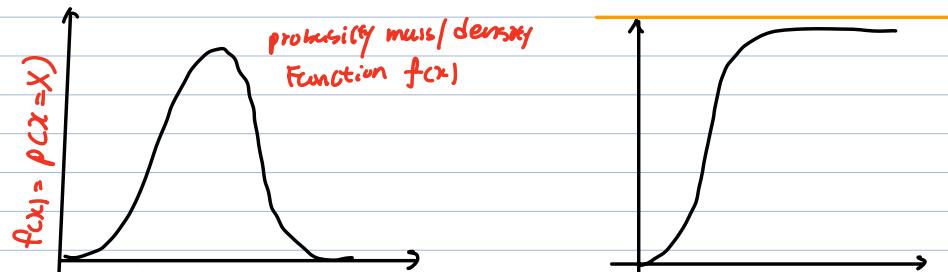
Time taken to access a hierarchical memory

## Probability distribution

### Probability Mass function



## How about continuous distribution



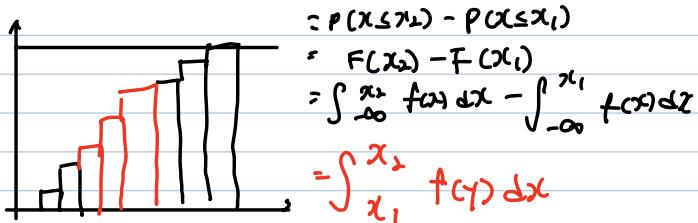
Outcome of experiment

Cumulative

Extracting More From  $F(x)$

$f(x)$  to find the probability

$$P(X_1 < x \leq X_2) =$$



## Chapter 8 : Probability Distributions as Modeling Tools

### 1. Uniform distribution

$$f(x) = \begin{cases} c & \text{if } x \geq a \text{ and } x \leq b \\ 0 & \text{if } x < a \text{ or } x > b \end{cases} \quad a, b \text{ and } c$$

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

$$\int_{-\infty}^{\infty} f(x) \cdot dx = \int_a^b c \cdot dx = c \cdot b - ca = c(b-a) = 1$$

$$\therefore c = \frac{1}{b-a}$$

$$\text{Mean } \mu = \frac{b+a}{2} \quad \text{Variance } \sigma^2 = \frac{(b-a)^2}{12} \quad \sigma = \sqrt{\frac{(b-a)^2}{12}}$$

### 2) Bernoulli distribution

Coin, a random trial with binary output

$$f(x) = \begin{cases} p & \text{if } x=0 \text{ fail} \\ 1-p & \text{if } x=1 \end{cases} \quad 1-p \text{가 성공 확률} \quad \text{Geometric oil 더 쓸 줄 알기}$$

$$\mu = 0 \cdot p + 1 \cdot (1-p) = 1-p$$

$$\sigma^2 = E[(x-\mu)^2] = E[x^2] - \mu^2$$

$$\begin{aligned}
 & 0 \cdot p + 1^2(1-p) - (1-p)^2 \\
 &= 1 \cdot p - (1-2p+p^2) \\
 &= 1 - p - 1 + 2p - p^2 \\
 &= p - p^2 = p(1-p)
 \end{aligned}$$

### 3) Geometric distribution

$$PMF f(n) = p^n(1-p), \quad 0 \leq p \leq 1, n \geq 0$$

성공전  $n$  번의 실패

$$\text{mean: } M = \sum_{n=0}^{\infty} n \cdot p^n(1-p) = \frac{p}{1-p} \text{ converse}$$

$$\frac{1}{0.9} = \frac{10}{9} = 1.1111$$

하지만 평균적 그 것 성공을 하기까지의 시장 횟수는

$$f(n) = f(n-1) = p^{n-1}(1-p) \quad M = \frac{1}{1-p}$$

$$S^2 = \frac{p}{(1-p)^2}$$

### 4. Binomial Distribution

conduct Bernoulli trial  $n$  times

$$f(x) = \binom{n}{x} \cdot p^x \cdot (1-p)^{n-x}$$

$$\binom{n}{x} = \frac{n!}{(n-x)! x!}$$

$$\text{mean: } M = np \quad \text{Varience: } n \cdot p(1-p)$$

The probability that a sever is "up" is 0.7 what is the probability that 3 servers will be up?

$$p = 0.7$$

$$n = 3$$

$$f(0) = \binom{3}{0} (0.7)^0 (0.3)^0 = 1 \cdot 1 \cdot 1 = 0.027$$

$$f(1) = \binom{3}{1} (0.7)^1 (0.3)^2 = \frac{3!}{1! \cdot 2!} \cdot \frac{3^2 \cdot 1}{2^2 \cdot 1} = 3 \times 0.7 \times 0.09 = 0.189$$

$$f(2) = \binom{3}{2} (0.7)^2 (0.3)^1 = \frac{3!}{2! \cdot 1!} = \frac{3 \times 2 \times 1}{2 \times 1 \times 1} \times (0.7)^2 \cdot 0.3 = 0.441$$

$$f(3) = \binom{3}{3} (0.7)^3 = \frac{3!}{3! \cdot 0!} = (0.7)^3 = 0.343$$

## 5 Poisson Distribution

Consider binomial Process with large  $n$  and a small  $p$ .  $n \rightarrow \infty$ ,  $p \rightarrow 0$ . Product  $np = \lambda > 0$ .

$\lambda$ : 단위시간 또는 단위공간에서 특정 사건이 발생하는 평균 발생 횟수

a) 1시간 동안 평균적으로 5번의 전화:  $\lambda = 5$  / 1분에 10명의 방문객:  $\lambda = 10$

$$\text{PMF } f(x) = \frac{\lambda^x}{x!} e^{-\lambda}$$

$x$ 는 사건이 발생한 횟수

$$\mu = \sigma^2 = \lambda$$

$$f(x) = \lim_{n \rightarrow \infty, p \rightarrow 0} \binom{n}{x} p^x (1-p)^{n-x} = \frac{\lambda^x}{x!} e^{-\lambda}$$

CDF  
자리잡기.

$$P(2) = 2\text{번 일어나는 확률}$$

$$P(3) = 3\text{번 일어나는 확률}$$

## 6. Exponential distribution

사건이 발생할 두 시간 사이의 시간 간격의 모델링

$\lambda$ : 단위시간당 사건이 발생할 것으로 기대되는 평균 발생 횟수

$$f(x) = \lambda e^{-\lambda x} \quad x > 0$$

사건이 한시간에 5번 발생하는  
시간 사이의 평균 시간 간격은  $\frac{1}{\lambda}$ 시간

$$F(x) = 1 - e^{-\lambda x}$$

$$\mu = \frac{1}{\lambda} \quad \sigma^2 = \frac{1}{\lambda}$$

A Note on Reliability

7



It is active than fails every  $\text{MTBF}$   
 then its failure rate is  $1/\text{MTBF}$   
 ... hence probability of failing at  $t$  or earlier  
 $F(t) = 1 - e^{-\lambda t / \text{MTBF}}$   
 ... hence prob of being "up" for  $t$  or longer  
 $1 - F(t) = e^{-\lambda t / \text{MTBF}}$

$$\text{MTBF} \geq 5 \text{ minutes}$$

$$\lambda = \frac{1}{5}$$

Failure rate is  $\lambda = 0.2$

failure rate 두 시간 사이의 간격이 MTBF (mean Time between failure)

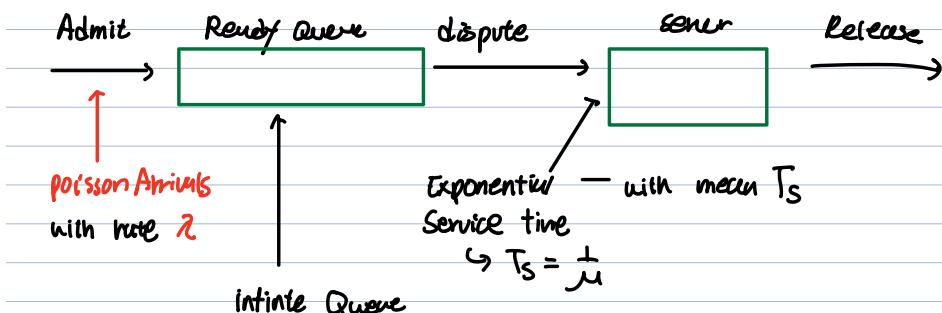
## Chapter 9 Elementary Queuing Analysis

Analysis of M/M/1 Queuing system

Single queueing single server

Arrival rate is poisson and service time is exponential

M/M/1 System



Consider a very small interval of time of length  $h$ . Assume that this interval time ( $h$ ) is so small that a max of one arrival can realistically occur in that period of time. Since the rate of arrival is  $\lambda$  per request/unit time the rate of arrival per interval is  $\lambda h$ . if  $\lambda = 100 \text{ req/sec}$  and  $h = 1 \text{ msec}$  rate of arrival =  $0.001$ ,  $0.1 \text{ req/msec}$ .

Probability density of the poisson distribution is

$$f(x) = \frac{\lambda^x}{x!} e^{-\lambda}$$

given that the rate of arrival per interval  $h$  is  $\lambda h$ , the probability of  $x$  arrivals per interval  $h$  is

$$f(x) = \frac{(\lambda h)^x}{x!} e^{-(\lambda h)}$$

$$f(0) = \text{no arrival is } f(0) = \frac{(\lambda h)^0}{0!} e^{-(\lambda h)} = e^{-\lambda h}$$

$$\text{so } f(\text{no arrival}) = e^{-\lambda h}$$

let's say arrival to the system is birth

$$P(\text{birth}) = 1 - P(\text{no arrival}) = 1 - e^{-\lambda h}$$

Rewrite using Taylor series

$$1 - e^{-\lambda h} = 1 - \left(1 - \lambda h + \frac{(\lambda h)^2}{2!} - \frac{(\lambda h)^3}{3!} + \dots\right)$$

$$\text{so } P(\text{birth}) = \lambda h \quad \text{for more arrival}$$

Similar to this, finish request is "Death"

rate of finish service is  $\mu$  req/sec. so for the interval of time  $h$ , the rate of finish service is  $\mu h$ .

so I can use poisson distribution

$$f(y) = \frac{(\mu h)^y}{y!} e^{-\mu h}$$

when  $y=0$   $f(0)$  : For a given time  $h$ , probability of finishing request = 0. so the probability of finishing request, "death" is

$$1 - f(0) = 1 - \frac{\lambda^0 \mu^0}{0!} e^{-\lambda \mu} = 1 - e^{-\lambda \mu} = \mu h$$

$$P(\text{death}) = \mu h$$

|

1 or more request finished

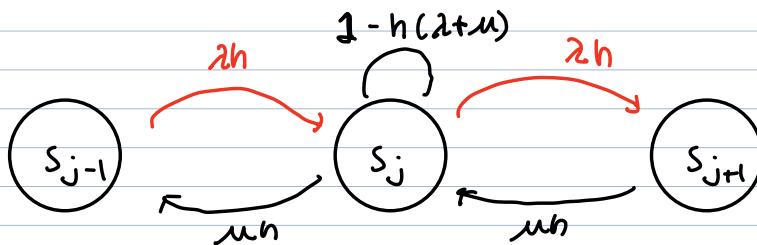
### State Transition diagram for M/M/1

State of the system by  $S_j$  j customers are in the system

1. The system was in  $S_{j-1}$  and a birth occurred, the probability of that happening is  $\lambda h$

2. The system was in  $S_{j+1}$  and a death occurred, the probability of that happening is  $\mu h$

3. The system was in state  $S_j$  and, neither death or birth occurred, The probability of that happening is  $1 - \lambda h - \mu h = 1 - h(\lambda + \mu)$



$$\text{Relationship: } P(S_j) = \lambda h \cdot P(S_{j-1}) + (1 - h(\lambda + \mu)) P(S_j) + \mu h \cdot P(S_{j+1})$$

Re arrange

$$0 = \lambda h \cdot P(S_{j-1}) + (1 - h(\lambda + \mu)) P(S_j) - P(S_j) + \mu h \cdot P(S_{j+1})$$

$$0 = \lambda h P(S_{j-1}) - h \lambda P(S_j) - \mu h P(S_j) + \mu h P(S_{j+1})$$

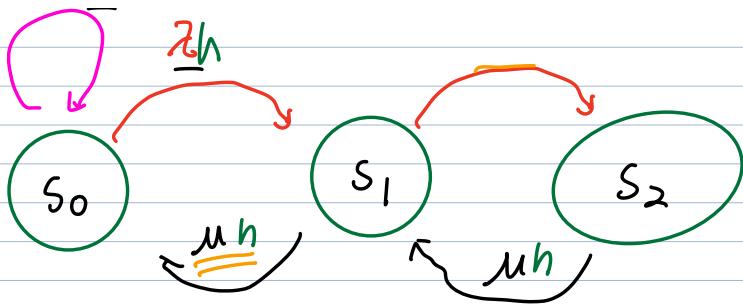
$$0 = \lambda P(S_{j-1}) - \lambda P(S_j) - \mu P(S_j) + \mu P(S_{j+1})$$

$$(\lambda + \mu) P(S_j) - \lambda P(S_{j-1}) = \mu P(S_{j+1})$$

$$P(S_{j+1}) = \frac{(\mu + \lambda) P(S_j)}{1 + \text{util}} - \frac{P(S_j)}{1 - \text{util}}$$

So : System state in which there are no requests

$$1 - \underline{\lambda h}$$



$$P(S_0) = (1 - \mu_h) P(S_0) + \mu_h P(S_1)$$

$$\alpha_h P(S_0) = \mu_h P(S_1)$$

$$\left(\frac{\lambda}{\mu}\right) = p$$

$$P(S_0) = P(S_1)$$

$P(S_1) = P(S_0)$

$$P(S_2) = P(S_2) = (1 + p) P(S_1) - P(S_0)$$

$$= (1 + p) \cdot P(S_0) - P(S_0)$$

$$= p^2 P(S_0)$$

$$\begin{aligned} P(S_3) &= P(S_3) = (1 + p) P(S_2) - P(S_0) \\ &= (1 + p) p^2 P(S_0) - P(S_0) \\ &= (p^2 + p^3 - p^1) P(S_0) \\ &= p^3 P(S_0) \end{aligned}$$

$$\therefore \text{so } P(S_j) = p^j P(S_0)$$

overall probability  $\sum_{i=0}^{\infty} P(S_i) = 1$

$$\sum_{n=0}^{\infty} P(S_n) = \sum_{n=0}^{\infty} p^n \cdot P(S_0) = 1$$

$$P(S_0) = \frac{1}{\sum_{i=0}^{\infty} p_i}$$

$\sum_{i=0}^{\infty} p_i$  is Geometric series.  $P \in [0, 1]$  (Probability)

$$M = \frac{1}{1-p} \text{ will converge.}$$

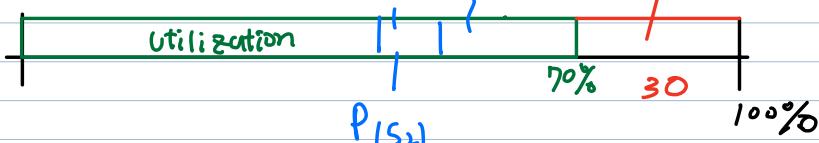
$$P(S_0) = \frac{1}{\frac{1}{1-p}} = 1-p$$

$$0.7 \times 0.3$$

$$P(S_1) = p \cdot (1-p)$$

$$\frac{1}{1-p} P(S_1)$$

$$P(S_0)$$



$$P(S_2)$$

$$0.7^2 \times 0.3$$

### Average number of customer in an M/M/1 system

How can we compute average number of customers  $q$  in an M/M/1 system.

$$q = E[\text{number of customers}] = \sum_{i=0}^{\infty} i \cdot P(S_i) = \sum_{i=0}^{\infty} i \cdot p^i \cdot (1-p)$$

$$= (1-p) \sum_{i=0}^{\infty} i \cdot p^i$$

$$q = (1-p) \sum_{i=0}^{\infty} \frac{p^i}{i!}$$

} poly logarithms

$$= (1-p) L_{i-1}(p) \quad L_{i-1}(p) = \frac{p}{(1-p)^2}$$

$$= (1-p) \times \frac{p}{(1-p)^2}$$

$$q = \frac{P}{1-P}$$

Average number of customers waiting for service in a M/M/1 System

$$w = q - P = \frac{P}{1-P} - P = \frac{P - P + P^2}{1-P} = \frac{P^2}{1-P}$$

Average time in an M/M/1 system

$$T_q = \frac{q}{\lambda} = \frac{P}{1-P} \div \lambda = \frac{P}{1-P} \times \frac{1}{\lambda} = \frac{\frac{\lambda}{\mu}}{1-P} \times \frac{1}{\lambda} = \frac{1}{\mu(1-P)} = \frac{1}{\mu(1-P)}$$

Average Time waiting in the M/M/1 system

$$Tw = \frac{w}{\lambda} = \frac{P^2}{1-P} = \frac{P}{\mu(1-P)}$$

### Slow down

Slow down is defined as average length of time it takes to get rendered (response time) over minimum amount of time in which can be rendered (service time)

$$\text{Slow down} = \frac{T_q}{T_s} = \frac{\frac{1}{\mu(1-P)}}{\frac{1}{\mu}} = \frac{1}{\mu(1-P)} \times \frac{\mu}{1} = \frac{1}{1-P}$$

### Variability measures

Standard deviation in the number of customers in an M/M/1 system as well as the total time in the system.

$$\sigma_q = \frac{\sqrt{P}}{1-P}$$

$$\sigma_{T_q} = \frac{1}{\mu(1-P)} = \frac{T_s}{1-P}$$

### Ex) M/M/1 Example

1. Poisson Arrival rate  $\lambda = 60 \text{ req/sec}$ , service exponential  $T_s = 9 \text{ msec/req}$

calculate response time  $T_q$

$$T_s = 0.009 \text{ sec/req} \quad \lambda = 10 \text{ req/sec}$$

$$p = 0.009 \times 10 = 0.09$$

$$q = \frac{p}{1-p} = \frac{0.09}{0.91} = 0.098901098$$

$$T_q = q \div \lambda = 0.098901098 \text{ sec/req}$$

2. Due to super bowl, hit to the webserver went up to 10 folds =  $\lambda \times 10$

$$\lambda = 100 \text{ req/sec}$$

response time  $T_q =$

$$p = 0.009 \times 10 = 0.9$$

$$q = \frac{0.9}{0.1} = 9$$

$$T_q = q \div \lambda = 9 \div 100 = 0.09$$

3. relative slow down

$$\text{Slow down} = \frac{T_q}{T_s}$$

$$\text{Before : slow down} : \frac{0.099}{0.009} = 1.099$$

After the super bowl : 10

4. How much faster the web server.

$$(0.009 \times x) \times 100 = p$$

$$1 - (0.009 \times 100) = \frac{p}{(0.009x)} = 1.09$$

$$\frac{p}{(0.009x)} = \frac{1.09}{1000}$$

$$1 - \alpha x \quad \frac{0.809x}{0.001}$$

$$\frac{0.809}{1 - 0.809} = 1.09 \frac{1}{0.001}$$

$$1.09 - 9.81x = 100.0$$

$$x = 101$$

## Chapter 11: Variations of the M/M/1 Queue System

Exponential distribution is memory less.

What if service time is not exponential?

### 11.1) Constant time M/D/1

Service time is constant. All customers require the same amount of service.

$$q = E[\text{number of customers}] = \frac{P^2}{2 \cdot \lambda(1-P)} + P \quad \text{ldk why}$$

$$W = q \cdot P = \frac{P^2}{2 \cdot \lambda(1-P)}$$

$$Tq = q \cdot \lambda = \frac{P}{2 \cdot \lambda(1-P)} + T_s \quad Wq = w \cdot \lambda = \frac{P}{2 \cdot \lambda(1-P)}$$

### 11.2) General Service Time - M/G/1

We do not know the distribution of service time.

But 1) Average service time,  $T_s$ , Service time standard deviation  $\sigma_{T_s}$

2) Normalized value of the service time standard deviation  $\sigma_{T_s}/T_s$

$$A \mid \frac{1}{2} \cdot \left[ 1 + \left( \frac{\sigma_{T_s}}{T_s} \right)^2 \right]$$

compute  $q$  and  $w$

$$\left. \begin{array}{l} q = \frac{P^2 \cdot A}{1-P} + P \\ w = \frac{P^2 \cdot A}{1-P} \end{array} \right\} \rightarrow \begin{array}{l} T_q = \frac{P \cdot A}{M(1-P)} + \frac{1}{\mu} = \frac{1}{\mu} \left( \frac{P \cdot A}{1-P} + 1 \right) \\ w_q = \frac{P \cdot A}{M(1-P)} \end{array}$$

### Sanity check

1) M/M/1.

For an exponential distribution, Standard is Equal to the mean  $\sigma_{T_S} = T_S$ ,  $\sigma_{T_S}/T_S = 1$

$$A \text{ for } M/M/1 = \frac{1}{2} [1 + (1)^2] = 1$$

$$q = \frac{P^2}{1-P} + P = \frac{P^2 + P - P^2}{1-P} = \frac{P}{1-P}$$

↓

$q$ 가 고정 constant 같다.

2) M/d/1

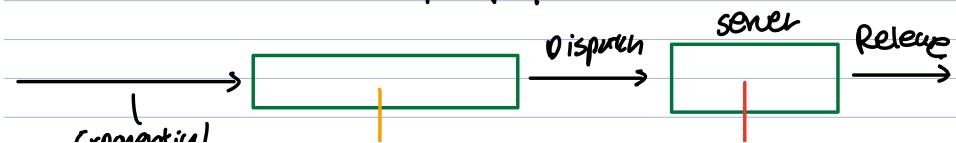
For an deterministic, Standard is 0

$$A \text{ for } M/d/1 = \frac{1}{2} [1 + 0^2] = \frac{1}{2}$$

$$q = \frac{P^2 \cdot \frac{1}{2}}{(1-P)} + P = \frac{\frac{P^2}{2}}{2(1-P)} + P$$

### II. 3 Finite Queue Systems - M/M/1/k

M/M/1/k:



Experiments

finite queue size

Exponential service time

Question

$M/M/1/(k)$ 에서  $P > 1$ 도 적용됩니다.

$P(S_j) \leq q \neq j$ 가 되도록 확률을 조정할 때,  $j=0$  일 때 마무리 요청이 없을 때,  $j=k$  일 때는 가족을 한 번에

$$P(S_j) = \begin{cases} \frac{(1-p) \cdot p^j}{1-p^{k+1}} & p \neq 1 \\ \frac{1}{k+1} & p = 1 \end{cases}$$

제속 흘러나가고 일도망 full  
p=1 최대로 요청을 내고 일도망 full

$$q = \begin{cases} \frac{p}{1-p} - \frac{(k+1) \cdot p^{k+1}}{(1-p^{k+1})} & \text{if } p \neq 1 \\ \frac{k}{2} & \text{if } p = 1 \end{cases}$$

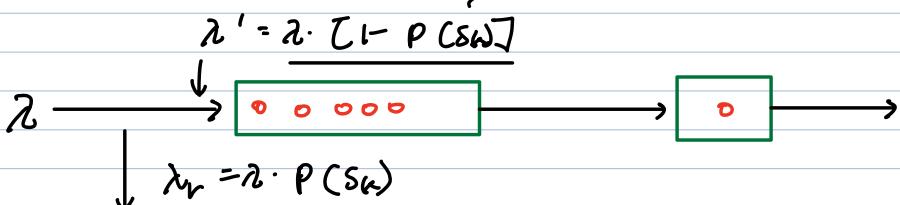
$M/M/1/k$

Rejection Rate

$P(S_k) = P(\text{rejection})$ . 즉  $k$ 가 full

$$P(S_k) = \begin{cases} \frac{(1-p)p^k}{1-p^{k+1}} & p \neq 1 \\ \frac{1}{k+1} & p = 1 \end{cases}$$

↑ when Queue is full



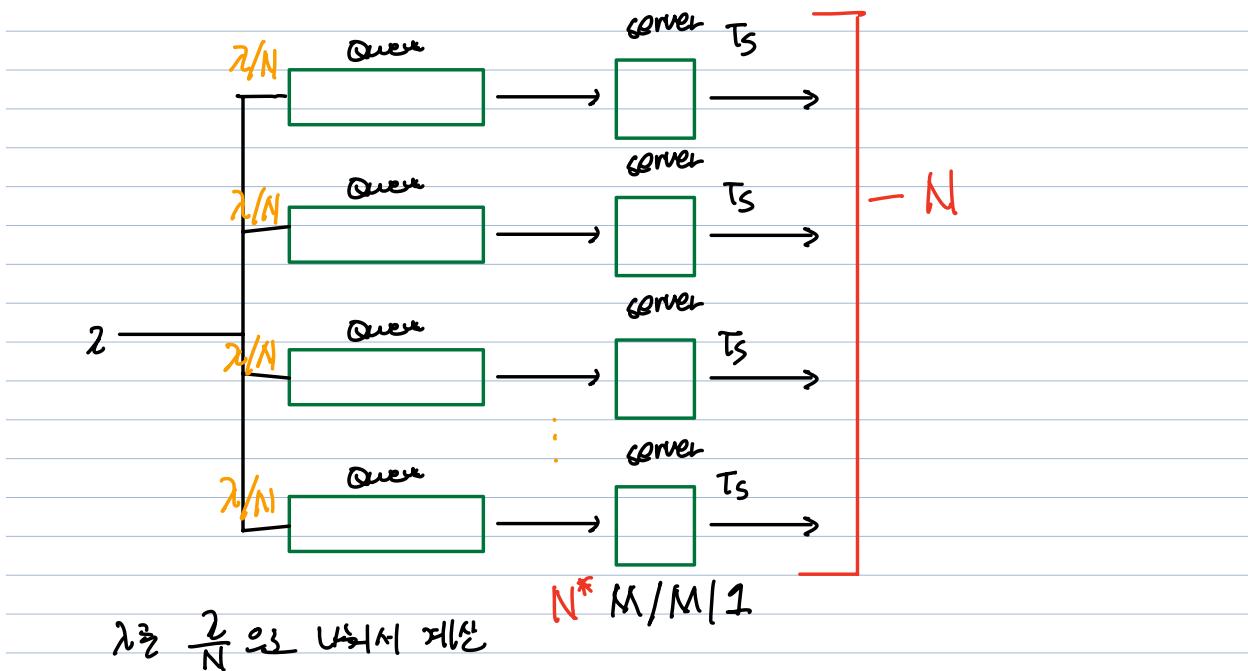
∴ For  $M/M/1/k$   $T_q = \frac{q}{\lambda}$

## Chapter 12: Networks of Queue

### 2) Multiple Single-server Queues

Customers choose one of  $N$  queues. Grocery stores star market

$N \times M/M/1$

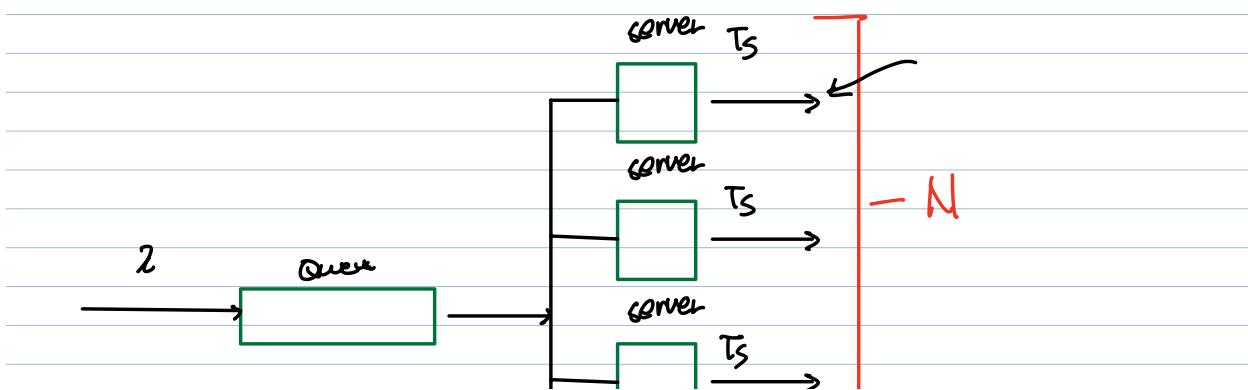


$$N \times q = q_{\text{tot}} \quad P' = \frac{\lambda}{N} \cdot T_s$$

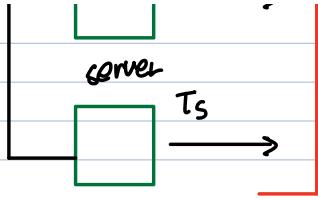
$$q' = \frac{P'}{1-P'} \quad N \times q' = q_{\text{total}}$$

### 2) Multi-server Queue

$M/M/N$



M/M/N



Probability N or more customers in the system (All servers are busy)

$$C = \frac{1 - k}{1 - p^k}$$

k is a poisson ratio function

$$k = \frac{\sum_{i=0}^{N-1} \frac{(NP)^i}{i!}}{\sum_{i=0}^N \frac{(NP)^i}{i!}} = 1 \cdot \frac{\sum_{i=0}^N \frac{(NP)^i}{i!}}{\sum_{i=0}^N \frac{(NP)^i}{i!}}$$

$$\text{erlang-c P} = \frac{\lambda T_s}{N} = \frac{\lambda}{\lambda N}$$

Utilization of  
single server eqn

$$P = \frac{\lambda}{N} T_s$$

$$Q | \frac{P}{1-P} \cdot C + NP \rightarrow \text{Requests being served}$$

$\downarrow$   
225 req x probability of all servers  
utilization are busy

Example

2 servers

$$\lambda = 20 \quad T_s = 80 \text{ msec /req} = 0.08 \text{ sec}$$

Gout: Compute Response Time T<sub>q</sub>

$$1) N M/M/1 \Rightarrow 2 * M/M/1$$

$$\lambda = 20 / 2 = 10 \text{ req/sec}$$

$$P = 10 \times 0.08 = 0.8$$

$$q' = \frac{0.8}{1-0.8} = 4 \quad q_{\text{total}} = 2 \times 4 = 8$$

$$Tq = 8 \div \lambda = 8 \div 20 = 0.4$$

2) M/M/2

$$p' = 10 \times 0.08 = 0.8$$

Get K first

$$K = 1 - \frac{(2 \times 0.8)^2}{2 \times 1}$$

$$i=0 \quad \frac{(2 \times 0.8)^0}{0!} = 1$$

$$i=1 \quad \frac{(2 \times 0.8)^1}{1!} = 1.6$$

$$i=2 \quad \frac{(2 \times 0.8)^2}{2 \times 1} = 1.28$$

3.88

$$k = 0.670$$

$$= 0.67$$

$$C = \frac{1 - 0.67}{1 - 0.8 \times 0.67} = \frac{0.33}{0.46} \approx 0.72$$

$$q = \frac{0.8}{0.2} \times 0.72 + 2 \times 0.8 = \\ 2.88 + 1.6 = 4.48$$

$$q = 4.48$$

$$Tq = 0.224$$

$$Q = \frac{P}{1-P} \cdot C + NP$$

Average Number of requests waiting for service in M/M/N

$$W = Q - \underline{NP} \rightarrow \text{현재 시기의 대기열의 평균값입니다.}$$

$$= \frac{P}{1-P} \cdot C + NP - NP = \frac{P}{1-P} \cdot C$$

Average Time in an M/M/N System

$$T_q = \frac{q}{\lambda} = \frac{NP'}{\lambda} + \frac{P'}{\lambda(C-P)} \cdot C$$

$$\begin{aligned} P' &= \frac{\lambda}{N} \cdot T_S \\ &= \frac{NT_S}{N} + \frac{T_S}{N(C-P)} \cdot C \\ &= T_S + \frac{T_S}{N(C-P)} \cdot C \end{aligned}$$

$$\sigma_{T_q} = \frac{T_S}{N(C-P)} \sqrt{C(C-1) + N^2(C-P)^2}$$

Average Waiting Time M/M/N

$$T_w = \frac{W}{2} = \frac{P}{1-P} \cdot C \times \frac{1}{2} = \frac{T_S}{N(C-P)} \cdot C$$

### 3) General Network of Queues



Considering the system at steady state:

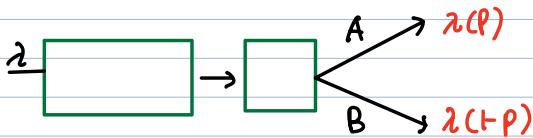
1) The rate in must be equal to the rate out of every queue

2) If the arrival process to a queue is Poisson, the departure process is also Poisson

Steady state arrival rate should be lower than the lowest service rate

Utilization  $\leq 100\%$ , then there will be steady state

## Splitting



1) Path A takes with probability:  $p$   
Path B taken with probability:  $1-p$

2) If the arrival process is Poisson, then a fixed "portion"  $p$  of arrivals is also Poisson.

## Merging



1) If two flows are joined, resulting flow has sum of the rates

2) If Arrival processes are Poisson, with  $\lambda_A$  and  $\lambda_B$ ,  
resulting process also Poisson with rate  $\lambda_A + \lambda_B$

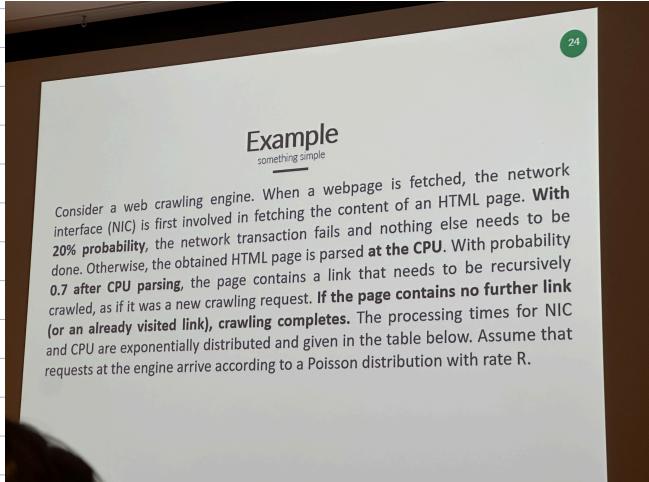
## Jackson Theorem

Assumption:

- Steady state
- Arrivals from outside are Poisson

- Service times for all servers Exponential
- No finite queues ✓

### Example slides



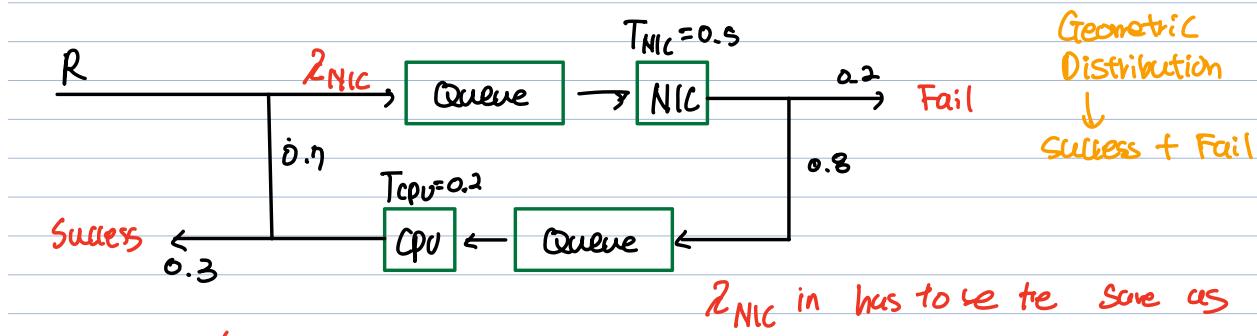
Given

Request processing times

CPU	0.2 sec
NIC	0.5 sec

Assume

- Steady State
- infinite Queue



$\lambda_{\text{NIC}}$  in has to be the same as  
 $\lambda_{\text{NIC}}$  coming out

$$\lambda_{\text{NIC}} = R + 0.8 \times 0.7 \times \lambda_{\text{NIC}}$$

$$\lambda_{\text{NIC}} = R + 0.56 \lambda_{\text{NIC}}$$

$$x = R + 0.56x \quad 0.44x = R$$

$$x = \frac{R}{0.44} = R \times \frac{100}{44} = \frac{25}{11} R$$

$$\lambda_{\text{NIC}} = \frac{25}{11} R$$

전체 시스템

$$\lambda_{\text{NIC}} = R + 0.8 \times 0.7 \times \lambda_{\text{NIC}}$$

$$\lambda_{\text{NIC}} = R + 0.56 \lambda_{\text{NIC}}$$

$$R = \lambda_{\text{NIC}} \cdot 0.2 + \lambda_{\text{NIC}} \cdot 0.8 \cdot 0.3$$

$$\lambda_{\text{NIC}} = \frac{R}{0.2 + 0.8 \cdot 0.3} = \frac{R}{0.44}$$

$$P_{NIC} = \lambda_{NIC} \cdot T_{NIC} = \frac{R}{0.44} \cdot 0.5 =$$

$$\lambda_{CPU} = \lambda_{NIC} \cdot 0.8$$

$$P_{CPU} = \lambda_{CPU} \cdot T_{CPU} = \frac{0.8R}{0.44} \cdot 0.2$$

$$q_{NIC} = \frac{P_{NIC}}{1 - P_{NIC}}$$

$$q_{CPU} = \frac{P_{CPU}}{1 - P_{CPU}}$$

$$q_{TOT} = q_{NIC} + q_{CPU}$$

$$T_{q,TOT} = \frac{q_{TOT}}{R}$$

$$\left[ \begin{array}{l} \frac{W_{total}}{R} = T_{W,tot} \\ \frac{T_{q,tot}}{T_{tot}} = T_{S,tot} \\ T_{q,tot} = T_{S,tot} + T_{W,tot} \\ \frac{T_{q,tot}}{T_{S,tot}} \approx \text{Slowdown} \end{array} \right]$$