

Queuing Theory

1.  $\lambda$  : arrival rate of job per unit time
2.  $\mu$  : service rate of job per unit time
3.  $m$  : # servers

Stability:  $\lambda < m\mu$

$$E[n] = E[n_q] + E[n_s]$$

$$E[r] = E[w] + E[s]$$

$$E[n] = \lambda E[r]$$

Little Law  $E[n_q] = \lambda E[w] \leftrightarrow$  Relation btwn jobs and time

$$E[n_s] = \lambda E[s]$$

Traffic intns  $\rho = \lambda E[s] = E[n_s]$  (one server)

$$P_0 = \text{Prob that server is idle}$$

$$= 1 - \rho$$

$$= 1 - \lambda E[s]$$

\*\* ??

$$E[r] = \frac{E[s]}{1 - \rho}$$

$$E[w] = \frac{\rho E[s]}{1 - \rho}$$

$$E[n] = \frac{\rho}{1 - \rho}$$

↗ service time distribution

$M/M/1 \leftarrow$  # server



inter-arrival time distribution

$$E[n_q] = \frac{\rho^2}{1 - \rho} \quad (M/M/1)$$

Operational Laws

## 1. Utilization Law

$$U_i = X_i S_i = \lambda_i S_i = X D_i \quad (D_i = V_i S_i)$$

## 2. Forced Flow

$$X_i = X V_i$$

## 3. Little's Law

$$Q = X R$$

$$Q_i = \lambda_i R_i = X_i R_i$$

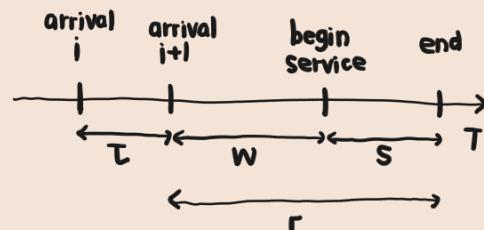
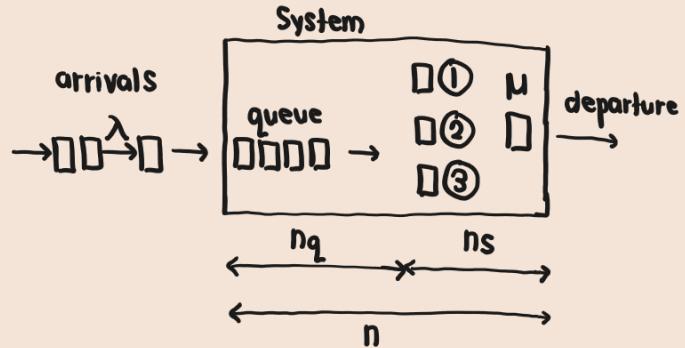
## 4. General Response Time

$$R = \sum_{i=1}^M V_i R_i$$

$$R_i = \frac{S_i}{1 - U_i}$$

## 5. Interactive Response Time

$$R = \frac{N}{X} - Z$$



## Asymptotic Bound

Thruput :  $\frac{N}{ND+Z} \leq x(N) \leq \min\left\{\frac{1}{D_{\max}}, \frac{N}{D+Z}\right\}$

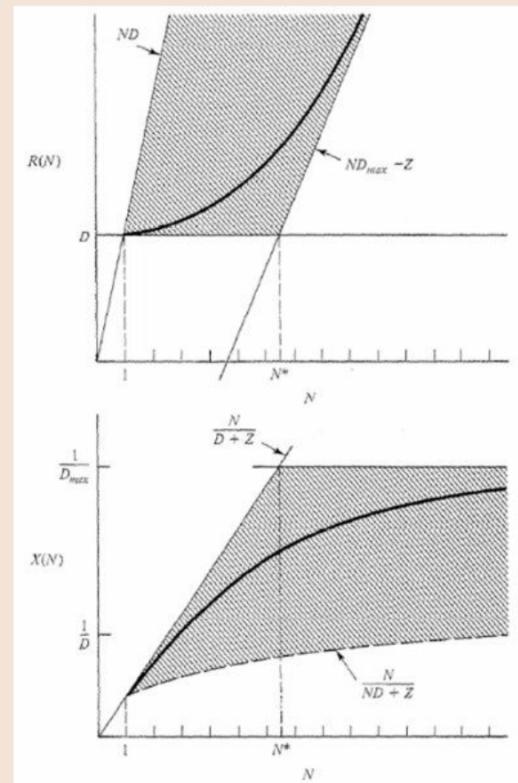
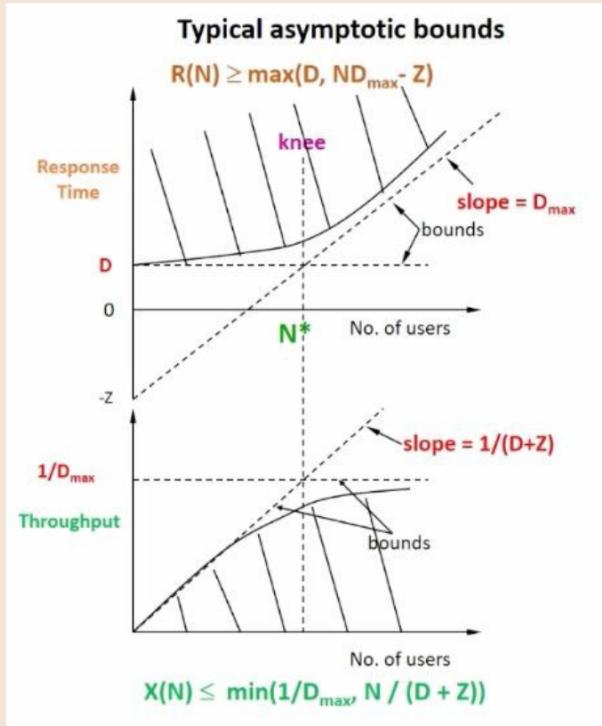
Response time:

$$\max\{D, ND_{\max} - Z\} \leq R(N) \leq ND$$

↑  
bottlenecked device

# users at knee

$$N^* = \frac{D+Z}{D_{\max}}$$



## Poisson $\lambda$

$\lambda$  = mean # events in duration unit time

pmf:  $f(x) = P(X=x) = \lambda^x \frac{e^{-\lambda}}{x!}$

Mean = Variance =  $\lambda$

## M: Exponential / Memoryless

$$Pr[M(t)=n] = \frac{e^{-\lambda t} (\lambda t)^n}{n!} \quad (t \geq 0, n=0,1,2\dots)$$

$$\xrightarrow{A_1}$$

$$Pr[A_1 \leq t] = 1 - e^{-\lambda t} \quad (\text{CDF})$$

$$Pr[A_1 > t] = Pr[M(t)=0] = e^{-\lambda t}$$

pdf:  $f(x) = \lambda e^{-\lambda x}$

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

$$0 \leq x \leq \infty$$

$$\text{Mean: } \frac{1}{\lambda}$$

$$\text{Var: } \frac{1}{\lambda^2}$$

$$E(x) = \begin{cases} \sum_{i \in I} x_i p(x_i) \\ \int_{-\infty}^{\infty} x f(x) dx \end{cases}$$

$$\begin{aligned} V(x) &= \sigma^2 = E[(x - E[x])^2] \\ &= E[x^2] - E[x]^2 \end{aligned}$$

## Operational Analysis

- Metrics that can be measured
- mean values measured during a finite duration (T or observation period)
- no assumption on distribution



- # job arrivals ( $A_i$ )
- # busy time of device i ( $B_i$ ) over T
- # job departures ( $C_i$ )



- Arrival rate  $\lambda_i = \frac{A_i}{T}$
- Thruput  $X_i = \frac{C_i}{T}$
- Utilization  $U_i = \frac{B_i}{T}$

$i=0 \Rightarrow$  system

$i>0 \Rightarrow$  devices

- Avg service time per job visit

$$S_i = \frac{B_i}{C_i}$$

- Avg # visits to device i during a job

$$V_i$$

- Utilization law

$$U_i = \frac{B_i}{T} \times 1 = \frac{B_i}{T} \times \frac{C_i}{C_i} = \frac{C_i}{T} \times \frac{B_i}{C_i} = X_i S_i$$

Device with the highest utilization is the bottleneck device

- Job Flow Balance Assumption

$$A_i = C_i \Leftrightarrow \frac{A_i}{T} = \frac{C_i}{T} \Leftrightarrow U_i = \lambda_i S_i$$

- Visit Ratio

Each job makes  $V_i$  requests for the  $i^{\text{th}}$  device

Balance  $\Rightarrow A_0 = C_0$

$$\Rightarrow C_i = C_0 V_i$$

## 11. Forced Flow Law

Relate sys throughput to individual device throughput

Assume each device has job flow balance ( $A_i = C_i$ )

$X$  is sys throughput

$$X = \frac{C_0}{T}$$

$$X_i = \frac{C_i}{T} \cdot \frac{C_0}{C_i} = \frac{C_0}{T} \cdot \frac{C_i}{C_0} = X V_i$$

let  $D_i$  be total service demand by a job at device  $i$

$$D_i = V_i S_i$$

combine forced flow and utilization law

$$U_i = X_i S_i = X V_i S_i \quad (\text{Forced Flow}) \\ = X D_i$$

$P_{ij}$  : probability that job moves from  $i$  to  $j$

$$C_j = \sum_{i=0}^M C_i P_{ij} \quad \Big) \div C_0 \\ V_j = \sum_{i=0}^M V_i P_{ij} \quad \text{visit ratio eqns} \\ V_0 = 1$$

## 12. Little's Law

Assume job flow balance

$Q_i$  := avg q size

$R_i$  := avg response time

$$Q_i = \lambda_i R_i = X_i R_i$$

## 13. General Response Time Law

Apply Little's Law to each component

$$Q = X R = Q_1 + Q_2 + \dots + Q_M$$

$$X R = \sum_{i=1}^M X_i R_i \quad \Big) \div X \text{ and forced flow (which imply Little)} \\ R = \sum_{i=1}^M V_i R_i$$

## 14. Interactive Response Time Law



$R :=$  sys response

$Z :=$  avg think time

Each user generates  $\frac{T}{R+Z}$  jobs in  $T$  time

$$\# \text{ jobs completed} = N * \frac{T}{R+Z}$$

$$x = \left( \frac{NT}{R+Z} \right) \cdot \frac{1}{T}$$

$$= \frac{N}{R+Z}$$

$$R = \frac{N}{x} - Z$$

## 15. Assume $N_i$ jobs at device $i$

$$\begin{aligned} \text{Then } R_i &= N_i S_i + S_i \\ &= R_i x_i S_i + S_i \\ &= R_i * u_i + S_i \end{aligned}$$

★  $R_i = \frac{S_i}{1-u_i}$  (M/M/1)

