

## Kendall Notation:

A/D/S/C/SD

A - arrival process

D - departure process

S - # of servers

C - buffer capacity

SD - service discipline

## Little's Law:

$$E(L) = \lambda E(S)$$

$$E(L_q) = \lambda E(W)$$

$$E(B) = \lambda E(\text{service time})$$

## M/M/1:

$$\rho = \frac{\lambda}{\mu}$$

$$\bar{\pi}_0 = 1 - \rho$$

$$\bar{s} = \frac{1}{\mu(1-\rho)}$$

$$\bar{L}_q = \frac{\rho^2}{1-\rho}$$

$$\bar{L} = \frac{\rho}{1-\rho}$$

$$\bar{W} = \frac{\rho}{\mu \cdot \lambda}$$

$$E(T) = \frac{1}{\lambda}$$

$$\bar{\pi}_n = (1-\rho) \rho^n$$

$$DBE: \bar{\pi}_n \lambda = \bar{\pi}_{n+1} \mu$$

$$P(X=i) = \frac{\lambda^i e^{-\lambda}}{i!}$$

$$P(T \leq t) = 1 - e^{-\lambda t}$$

## M/M/1/K:

$$\tilde{\pi}_0 = \frac{1-p}{1-p^{K+1}}$$

$$\tilde{\pi}_n = \frac{p^n(1-p)}{1-p^{K+1}}$$

$$P_{\text{Loss}} = \tilde{\pi}_K = \frac{p^K(1-p)}{1-p^{K+1}}$$

$\tilde{\pi}_K$  is the state of packet loss.

## M/M/C:

$$p = \frac{\lambda}{c\nu}$$

$$\tilde{\pi}_n = \frac{(cp)^n}{n!} \tilde{\pi}_0$$

$$\tilde{\pi}_{c+n} = p^n \tilde{\pi}_c = p^n \frac{(cp)^c}{c!} \tilde{\pi}_0$$

$$\tilde{\pi}_0 = \left[ \sum_{n=0}^{c-1} \frac{(cp)^n}{n!} + \frac{(cp)^c}{c!} \cdot \frac{1}{1-p} \right]^{-1}$$

queuing prob: prob. an incoming packet has to wait in the buffer

$$\tilde{\pi}_w = \tilde{\pi}_c + \tilde{\pi}_{c+1} + \dots = \frac{\tilde{\pi}_c}{1-p} = \frac{(cp)^c}{c!} \frac{1}{1-p} \tilde{\pi}_0$$

mean queue length

$$E(L_q) = \sum_{n=1}^{\infty} n \cdot \bar{M}_{C+n} = \bar{N}_w \frac{\rho}{1-\rho}$$

mean waiting time

$$E(W) = \frac{E(L_q)}{\lambda} = \bar{N}_w \frac{1}{1-\rho} \frac{1}{\bar{C}w}$$

## Continuous Time Markov Chain (CTMC)

**Def (CTMC):** A stochastic process  $\{X(t) : t \geq 0\}$  with discrete state space  $S$  is called CTMC if

$$\Pr \{ X(t) = j \mid X(h) = i, X(t_{n-1}) = i_{n-1}, \dots, X(t_1) = i_1 \} \\ = \Pr \{ X(t) = j \mid X(h) = i \}$$

where  $0 \leq t_1 \leq \dots \leq t_{n-1} \leq h \leq t$  and  
 $i_1, i_2, \dots, i_{n-1}, i_j, j \in S, \forall n \geq 1.$

Again, given the current state, the future state is independent from the past state.

Time Homogeneity: says CTMC is time homogeneous  
 if For any  $h \leq t$  and any state  $i, j \in S$ .

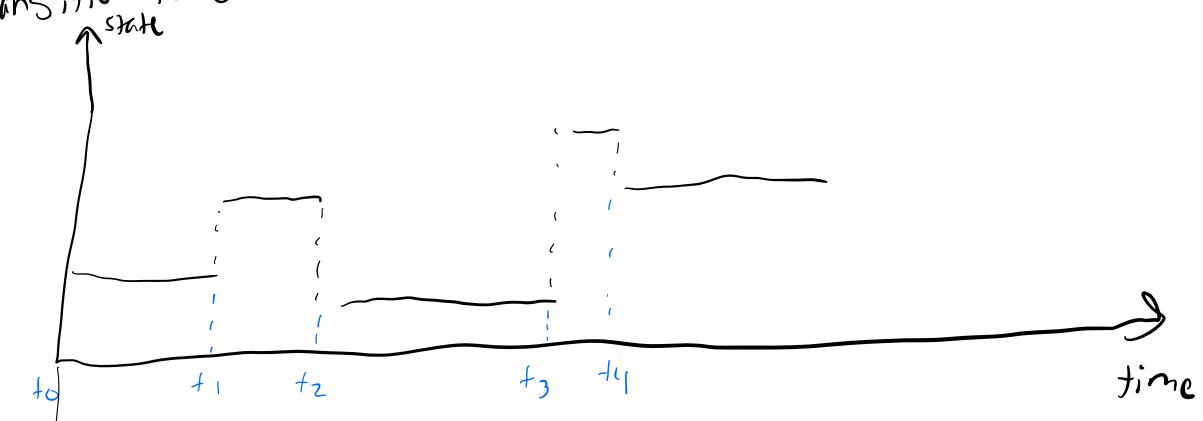
$$\Pr[X(t) = j \mid X(h) = i] = \Pr[X(t+h) = j \mid X(h) = i] = \Pr_{ij}(t+h)$$

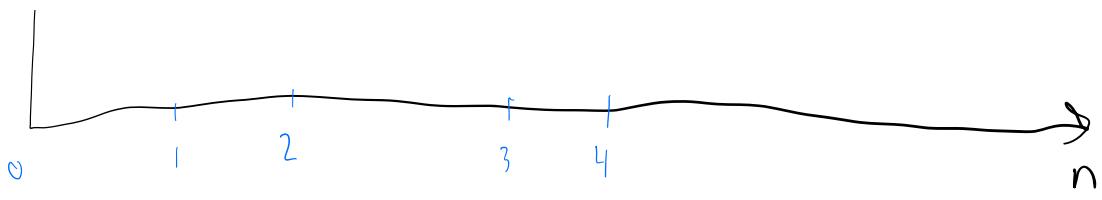
↑  
transition probability

10/6/25

By time homogeneity, whenever the process enters state  $i$ , the way it evolves probabilistically from that point is the same as if the process started in the state  $i$  at time 0. When the process enters state  $i$ , the time it spends there before it leaves state  $i$  is called the **holding time** in State  $i$ . The Markov property dictates that the holding time must have memoryless property and thus are exponentially distributed. Since an exponential distribution is completely determined by its rate  $\lambda_i > 0$ , we can conclude that for each  $i \in S$ , there must exist a constant rate  $\lambda_i > 0$ , such that the chain, when entering state  $i$  remains there for an amount of time  $T_i \sim \text{exp}(\lambda_i)$ .

**Embedded DTMC**: Discretize the CTMC by the moments when the process leaves a state and the state right after the transition is defined as the state of the DTMC.





every time the CTMC leaves state  $i$ , its transition to a state  $j$  is governed by a DTMC according to a transition probability matrix  $[P_{ij}]$ . ( $P_{ii} = 0$ )

Now, let's consider the first-order derivative of the continuous-time transition probability,  $P_{ij}(t)$ , evaluated at  $t=0$ ; i.e.,  $P_{ij}'(t)|_{t=0} = P_{ij}'(0)$   
i.e., the infinitesimal transition rate =

$$P_{ij}'(0) = \lim_{h \rightarrow 0} \frac{P_{ij}(h) - P_{ij}(0)}{h}$$

① Case 1: when  $i \neq j$        $P_{ij}(0) = 0$

$$P_{ij}'(0) = \lim_{h \rightarrow 0} \frac{P_{ij}(h)}{h} = \lim_{h \rightarrow 0} \frac{\Pr\{N_i(h)=1\}}{h} \cdot P_{ij}$$

$\checkmark$  CDF or exp. dist.

$$\Pr\{N_i(h)=1\} = \Pr\{\tau_i \leq h\} = 1 - e^{-\alpha_i h}$$

Do the Taylor expansion

$$\Pr\{N_i(h)=1\} = a \cdot h + O(h)$$

$$P'_{ij}(0) = \lim_{h \rightarrow 0} \frac{[a_i h + O(h)] \cdot P_{ij}}{h} = a_i \cdot P_{ij} \text{ for } i \neq j.$$

only need to know result

② Case 2: when  $i = j$

$$P_{ii}(0) = 1$$

$$P'_{ii}(0) = \lim_{h \rightarrow 0} \frac{P_{ii}(h) - 1}{h}$$

$$\begin{aligned} P_{ii}(h) &= 1 - \Pr\left\{\Pr(x/h) \geq i \mid X(0) = i\right\} \\ &\leq 1 - [1 - a_i h + O(h)] \\ &= a_i h + O(h) \end{aligned}$$

$$P'_{ii}(0) = \lim_{h \rightarrow 0} \frac{P_{ii}(h) - 1}{h}$$

10/13/25

The matrix  $Q = \{P'_{ij}(0)\}$  is called the transition rate matrix, or infinitesimal generator of the CTMC.

Ex: if  $S = \{0, 1, 2, 3, 4\}$ , then,

$$Q = \begin{matrix} & 0 & 1 & 2 & 3 & 4 \\ 0 & -a_0 & a_{01} & a_{02} & a_{03} & a_{04} \\ 1 & a_{10} & -a_1 & a_{12} & a_{13} & a_{14} \\ 2 & a_{20} & a_{21} & -a_2 & a_{23} & a_{24} \\ 3 & a_{30} & a_{31} & a_{32} & -a_3 & a_{34} \\ 4 & a_{40} & a_{41} & a_{42} & a_{43} & -a_4 \end{matrix}$$

Note that each row adds to 0.

\* how to compute  $P_{ij}(t)$  → Kolmogorov backward equation

LBE: For CTMC with transition rate matrix  $Q = P'(0)$ , where  $P(t) = \{P_{ij}(t)\}$  the following set of linear differential equation is satisfied by  $P(t)$ .

$$P'(t) = Q P(t), \text{ where } t \geq 0, \quad P(0) = I.$$

$$\text{that is } P'_{ij}(t) = -a_i P_{ij}(t) + \sum_{k \neq i} q_{ik} P_{kj}(t)$$

$$= -a_i P_{ij}(t) + \sum_{k \neq i} a_k P_{ik} P_{kj}(t), \quad i, j \in S$$

The solution is thus of the exponential form.

$$P(t) = e^{Qt}, \quad t > 0$$

where for square matrix  $M$ ,

*don't have to calculate this*

$$e^M = \sum_{n=0}^{\infty} \frac{M^n}{n!}$$

**Stationary distribution:** Let  $\{x(t) : t \geq 0\}$  be CTMC with state space  $S$ , generator  $Q$  and transition probability matrix  $P(t)$ . A row vector  $\vec{\pi} = (\pi_i)_{i \in S}$  with  $\pi_i \geq 0$  for  $i \in S$  and  $\sum_{i \in S} \pi_i = 1$ . is a stationary distribution of the CTMC if  $\vec{\pi} = \vec{\pi} P(t)$  for all  $t \geq 0$ .

$$\vec{\pi} \text{ is a stationary distribution} \iff \vec{\pi} = \vec{\pi} P(t), \forall t \geq 0$$

$$\iff \vec{\pi} = \vec{\pi} \sum_{n=0}^{\infty} \frac{(-Q)^n}{n!} \text{ for all } t \geq 0$$

$$\iff \vec{\pi} = \vec{\pi} + \vec{\pi} \sum_{n=1}^{\infty} \frac{(-Q)^n}{n!}$$

$$\iff \vec{\sigma} = \vec{\pi} \sum_{n=1}^{\infty} \frac{(-Q)^n}{n!}$$

$$\dots \vec{\sigma} = \sum_{n=1}^{\infty} \frac{(-Q)^n}{n!} \vec{\pi} Q^n \text{ for all } t \geq 0$$

$$\dots \vec{\sigma} = \vec{\pi} Q^n \text{ for all } n \geq 1$$

$$\dots \vec{0} = \pi Q$$

Expanded form: The  $j^{\text{th}}$  equation in  $\vec{0} = \vec{\pi} Q$  is:

$$0 = -\pi_j + \sum_{i \leq j} q_{ij} \pi_i$$

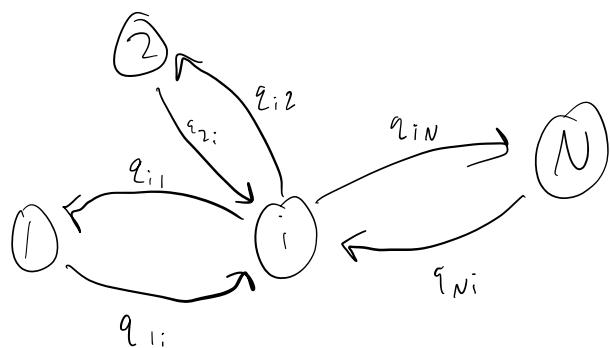
or

$$\pi_j q_j = \sum_{i \leq j} \pi_i q_{ij}$$

LHS (left): the rate leaving state  $j$  weighted by the probability of state  $j$ .

RHS:  $q_{ij}$  is the transition rate going into state  $j$  from  $i$  weighted by the stationary prob. of state  $i$  ( $\pi_i$ ).

This is called global balance equation. The long term rate out of state  $j$  has to be equivalent of the long term rate going into state  $j$ .



at stationary  $\pi_i (q_{i1} + q_{i2} + \dots + q_{in}) = \pi_i q_{1i} + \pi_2 q_{i2} + \dots + \pi_N q_{Ni}$

entering

leaving

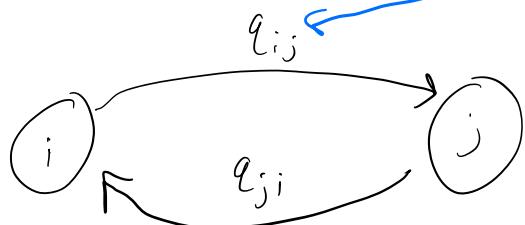
Normalization condition  $\sum_{i \in S} \pi_i = 1$ .

10/15/25

Detailed Balance Equation:

For a CTMC, if the solution  $\vec{\pi}$  to the following equation set exists

$$\pi_i q_{ij} = \pi_j q_{ji}, \text{ for all } i, j \in S$$



first in subscript is source & second is destination

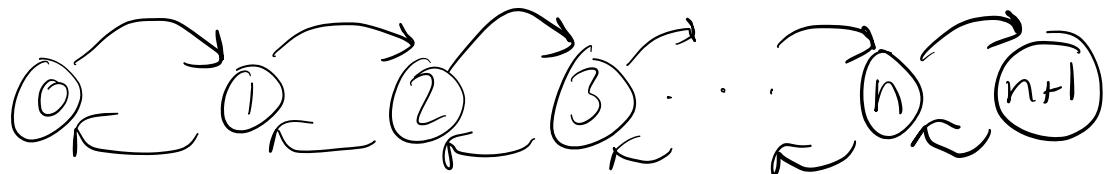
Then  $\vec{\pi}$  is the stationary distribution of the states.  
The above equation set is called detailed balance equation.

It can be shown that the solution to the detailed balance equation must be the solution to the global balance equation, but not vice versa, which means that a general CTMC should always satisfy the global balance equation but may not satisfy the detailed balance equation.

Detailed is sufficient but not necessary.  
Always there are Global.

It can be shown that in birth-death process satisfies the DBE.

Birth-death process : A markov chain with transition only between adjacent/neighboring states  
 $q_{ij}=0$  if  $|i-j| > 1$



In this case instead of looking at each state (node), we look at the edges and apply the DBE.

$$\pi_n q_{n,n+1} = \pi_{n+1} q_{n+1,n}$$

DBE is much easier to solve than GBE.

## Basic Queueing Theory

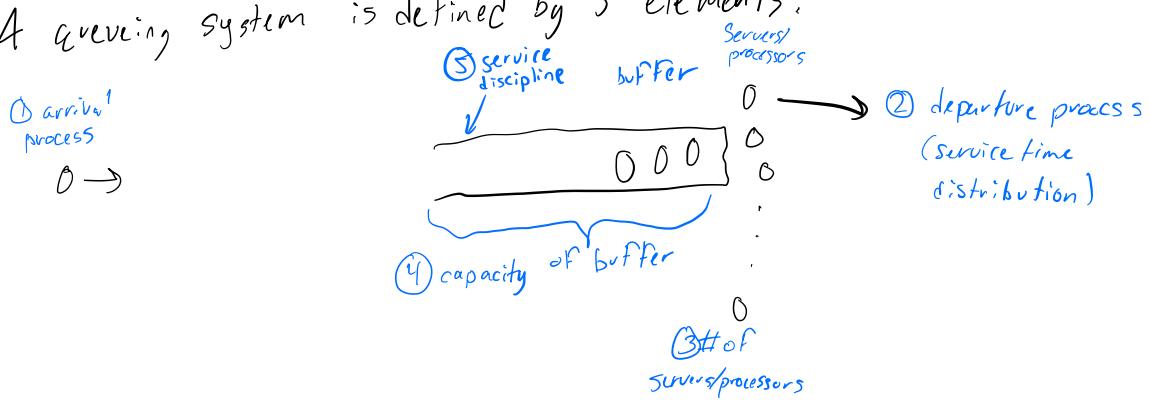
\* Kendall notation

\* Basic types of queues

\* How to compute important performance matrix

\* Queuing networks

A queuing system is defined by 5 elements.



Kendall notation: A/D/S/C/SD

A: arrival process (Poisson)

D: departure process (service time distribution)

S: # of servers

C: buffer capacity

SD: service discipline (FIFO, FILO)

M/M/1/∞

If no SD then it is FIFO

Also don't have to put ∞

More concise: M/M/1

→ M = Markovian → Poisson with rate,  $\lambda$

→ M<sub>s</sub> = Markovian → Exponential service time with rate  $\mu$

# of server = 1

queue capacity = ∞

SD = FIFO (first in first out)

10/17/25

M/M/1

Define the # of customers (packets/jobs) in the system (including the customers that are waiting in the queue and the one being serviced) as the state of the stochastic process. Because the customer interarrival times and the service times are all independent exponentially distributed r.v.s, the process

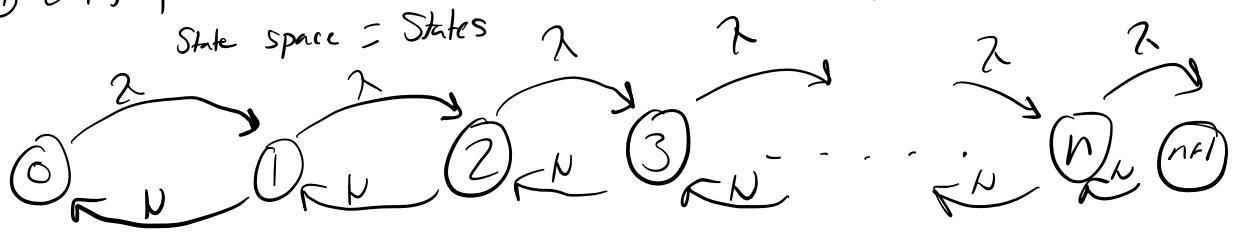
$\{X(t), t \geq 0\}$  is a CTMC with state space  $S = \{0, 1, 2, \dots, \infty\}$



only 1 event (arrival or departure) can happen, multiple event cannot happen at the same time.

To solve for the stationary distribution for  $X(t)$

① Let's plot the state transition diagram. this represents # of customers



Cannot have more than 1 event happening at 1 time

$\lambda$  = arrival rate

$\nu$  = departure rate

$$\Pr \{ \text{a new arrival within } h \}$$

$$= \Pr \{ \tau < h \}$$

$$= \lambda h + O(h)$$

$$\lim_{h \rightarrow 0} \frac{\lambda h + O(h)}{h} = \lambda$$

② Apply the GBE

$$\pi_0 \cdot \lambda = \pi_1 \cdot \nu \Rightarrow \pi_1 = \frac{\lambda}{\nu} \pi_0$$

cancel each other  
//

$$\begin{aligned} \pi_1(\lambda + \nu) &= \pi_0 \lambda + \pi_1 \nu \Rightarrow \pi_1 \nu = \underline{\pi_0 \lambda + \pi_1 \nu} \geq \pi_1 \lambda = \pi_2 \nu = \\ &\quad \pi_2 = \frac{\lambda}{\nu} \pi_1, \\ \vdots \quad \pi_n(\lambda + \nu) &= \pi_{n-1} \lambda + \pi_n \nu \\ &\quad \uparrow \quad \pi_3 = \left(\frac{\lambda}{\nu}\right)^3 \pi_0 \\ &\quad \pi_n = \left(\frac{\lambda}{\nu}\right)^n \pi_0 \end{aligned}$$

Let  $p \triangleq \frac{\lambda}{\nu}$  we have  $\pi_n = p^n \pi_0 \dots n=1, 2 \dots$   
 $0 \leq p \leq 1$

$$\begin{aligned} \pi_0 + \pi_1 + \dots &= 1 \\ \pi_0 (1 + p + p^2 + \dots) &= 1 \\ \pi_0 \frac{1}{1-p} &= 1 \Rightarrow \pi_0 = 1-p. \end{aligned}$$

Alternatively, apply DBE:

$$\begin{aligned} \pi_0 \lambda &= \pi_1 \nu \Rightarrow \pi_1 = \frac{\lambda}{\nu} \pi_0 \\ \pi_1 \lambda &= \pi_2 \nu \Rightarrow \pi_2 = \frac{\lambda}{\nu} \pi_1 = \left(\frac{\lambda}{\nu}\right)^2 \pi_0 \\ \pi_2 \lambda &= \pi_3 \nu \Rightarrow \pi_3 = \left(\frac{\lambda}{\nu}\right)^3 \pi_0 \\ \pi_{n-1} \lambda &= \pi_n \nu \Rightarrow \pi_n = \left(\frac{\lambda}{\nu}\right)^n \pi_0 \end{aligned}$$

It's clear that the stationary distribution exists if and only if  $p < 1$ .

$p$  is also called the load of the system (or occupation rate or service utilization).

Let's calculate a set of important performance matrix:

① Average # of customers in the system  $n \cdot p^{n+1} = (p^n)'$

$$E(L) = \sum_{n=0}^{\infty} n \pi_n = \sum_{n=0}^{\infty} n (1-p) \cdot p^n = \frac{p}{1-p}$$

② What's the prob. that the server is idle?

$$\pi_0 = 1-p \quad \text{so server utilization} = 1 - \pi_0 = p.$$

10/20/25

arrival rate - 2

service rate - N

$$\pi_0 = 1-p \quad \pi_n = p^n \pi_0 = p^n (1-p)$$

$$\text{average } L \quad E(L) = \frac{p}{1-p}$$

average number of customers being serviced: B

$$E(B) = 0 \cdot \pi_0 + 1 \cdot (1, \pi_0) = p \quad E(B) = p$$

Queue length:  $L^q$

$$L = L^q + B \Rightarrow E(L) = E(L^q) + E(B)$$

$$\Rightarrow E(L^q) = E(L) - E(B) = \frac{\rho^2}{1-\rho}$$

**Waiting time:** time spent in the queue  $\rightarrow$  from the moment of arrival (entering the buffer) to the moment of leaving the queue to get service. ( $w$ ).

**Sojourn time:** the waiting time plus service time  $\rightarrow$  from the moment of arrival to the moment of leaving the system (departure). (S)  
 (Also called delay, packet gets to router to leaving the router)

**Little's law:** specifies the relationship between  $E(L)$  and  $E(S)$ .

$$E(L) = \lambda E(S),$$



The above applies to any arrival and / departure process and any service order and even subsequences.

$$E(L^q) = \lambda E(w)$$

$$E(B) = \lambda E(\text{service time}) = E(\beta) = \lambda \frac{1}{\mu} = \rho$$

So For M/M/1 queue:

$$E(S) = \frac{E(L)}{\lambda} = \frac{1}{\mu(1-\rho)}$$

**PASTA property:** For queueing system with Poisson arrivals (M/M/1 system), the special property holds that arriving customers finds on average the same situation in the queueing system as an outside observer looking at the system at an arbitrary point in time.

Priorities: consider a M/M/1 queueing system servicing 2 types of user (packets) type 1 & type 2. Type 1 & Type 2 customers arrive according to independent Poisson processes with rates  $\lambda_1$  &  $\lambda_2$  respectively. The service times of all customers are iid exponential with mean  $\frac{1}{\mu}$ . We assume that  $\rho_1 + \rho_2 < 1$ , where  $\rho_i = \frac{\lambda_i}{\mu}$ . Type 1 has higher priority than Type 2.

### ① Preemptive Priority:

Def. In preemptive-resume rule, when type 2 job is in service, and a type 1 arrives, the type 2 service is interrupted (suspended) the server proceeds with the type 1 job. Once there is no more type 1 jobs in the system, the server resumes the service of the type 2 job at the point where its interrupted.

Let  $L_i$  denote the number of type  $i$ -jobs in the system and  $S_i$  denote the sojourn time of a type  $i$  job. Please calculate the expectation  $E(L_i)$  and  $E(S_i)$  for  $i=1, 2$ .

10/22/25

M/M/1 - queue 2 types of jobs

Type 1 > Type 2

Poisson rates  $\lambda_1, \lambda_2$

Service time  $\mu$

## ① Preemptive Priority

$$L_i, S_i, E(L_i), E(S_i)$$

For type 1, the type 2 jobs do not exist. So we have:

$$E(L_1) = \frac{P_1}{1-P_1} \quad E(S_1) = \frac{1}{\mu(1-P_1)}$$

Average delay for type 1 job

Since the service time is exponential, the total # of jobs in the system does not depend on the service order. So

$$E(L_1 + L_2) = \frac{P_1 + P_2}{1 - P_1 - P_2}$$

$$P_i = \frac{\lambda_i}{\mu}$$

$$E(L_1 + L_2) = E(L_1) + E(L_2) \Rightarrow E(L_2) = E(L_1 + L_2) - E(L_1)$$

$$= \frac{P_1 + P_2}{1 - P_1 - P_2} - \frac{P_1}{1 - P_1}$$

$$= \frac{P_2}{(1 - P_1)(1 - P_1 - P_2)}$$

↑  
average # of type 2

Applying Little's law on type 2 traffic:

$$E(L_2) = \lambda_2 \cdot E(S_2) \Rightarrow E(S_2) = \frac{E(L_2)}{\lambda_2} = \frac{\frac{P_2}{(1 - P_1)(1 - P_1 - P_2)}}{\lambda_2}$$

## (2) Non preemptive priority

Type 1 job is not allowed to interrupt the service of type 2 job.

there is 1

then is 0

$$P_2 \frac{1}{N} + (1-P_2) 0$$

$$E(S_1) = E(L_1) \cdot \frac{1}{N} + \frac{1}{N} P_2 \frac{1}{N}$$

↑ average service time  
 Service for type 1  
 jobs in front of me  
 ↑ Service for first firm if

By Little's Law:  $E(L_1) = \lambda_1 E(S_1)$

$$E(L_1) = \frac{(1+P_2)P_1}{1-P_1}$$

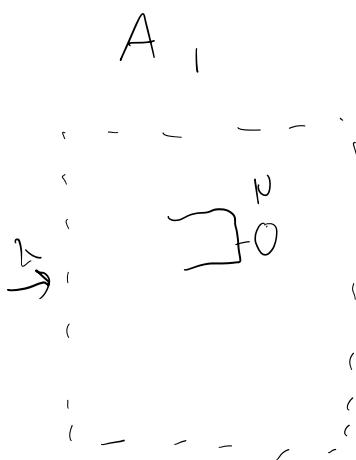
$$E(S_1) = \frac{1+P_2}{N(1-P_1)}$$

$$E(L_2) = E(L_1 + L_2) - E(L_1)$$

$$= \frac{P_1 + P_2}{1-P_1 - P_2} - \frac{(1+P_2)P_1}{1-P_1} = \frac{[1-P_1(1-P_1+P_2)]P_2}{(1-P_1)(1-P_1-P_2)}$$

$$E(S_2) = \frac{E(L_2)}{\lambda_2} = \frac{[1-P_1(1-P_1+P_2)]}{N(1-P_1)(1-P_1-P_2)}$$

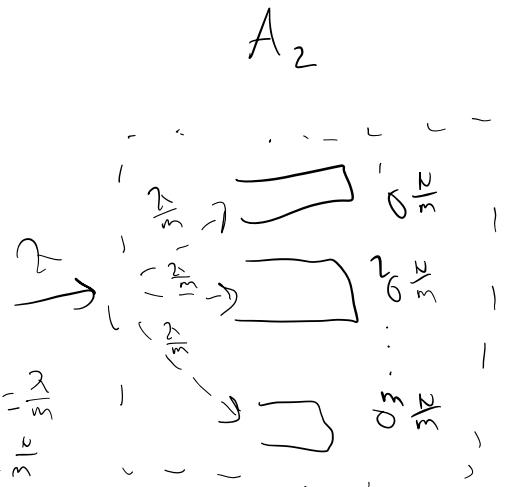
Let's compare the following 2 architectures:



$$\bar{L} = \frac{\rho}{1-\rho} = \frac{\lambda}{\mu - \lambda}$$

$$\bar{S} = \frac{\bar{L}}{\lambda} = \frac{1}{\mu - \lambda}$$

$$\bar{W} = \frac{\rho}{\mu - \lambda}$$



$$\lambda' = \frac{\lambda}{m}$$

$$\mu' = \frac{\mu}{m}$$

$$\rho' = \rho$$

$$\bar{L}' = \frac{\rho'}{1-\rho'} = \bar{L}$$

$$\bar{S}' = \frac{\bar{L}'}{\lambda'} = \frac{\bar{L}}{\frac{\lambda}{m}} = m\bar{S}$$

$$\bar{W}' = \frac{\rho'}{\mu' - \lambda'} = m\bar{W}$$

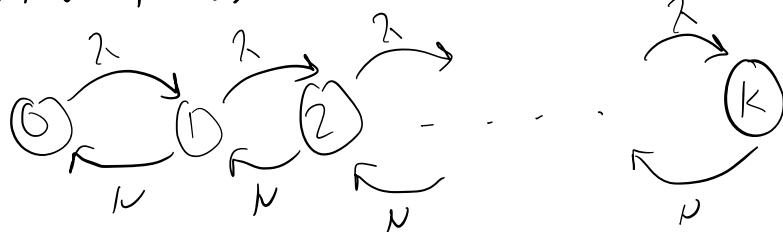
Key observation: A<sub>2</sub> slows down the arrival & service rate by a factor of  $m$   $\Rightarrow$  same queue length, but delay is  $m$  fold longer.

10/24/25

\* **M/M/1/K** queue (at most  $K$  packets in the system)

$K$ -buffer capacity, # of jobs, capacity

birth/death b/c M/M/1



$\lambda$  - arrival rate

$\mu$  - service rate

Exactly the same as in M/M/1 queue, but state space is truncated!

$$\pi_n = \rho^n \pi_0, n=1, 2, \dots, K$$

Finite # of  $\pi_n$

normalization condition:  $\sum_{n=0}^K \pi_n = 1 \Rightarrow \pi_0 = \frac{1-\rho}{1-\rho^{K+1}}$

$$\pi_n = \frac{\rho^n(1-\rho)}{1-\rho^{K+1}}$$

- stationary distribution

Now there will be packet loss b/c overflow can happen

packet loss rate / buffer overflow rate  $P_{\text{loss}} = \pi_K = \frac{\rho^K(1-\rho)}{1-\rho^{K+1}}$

$E = \lambda S \rightarrow M/M/1$

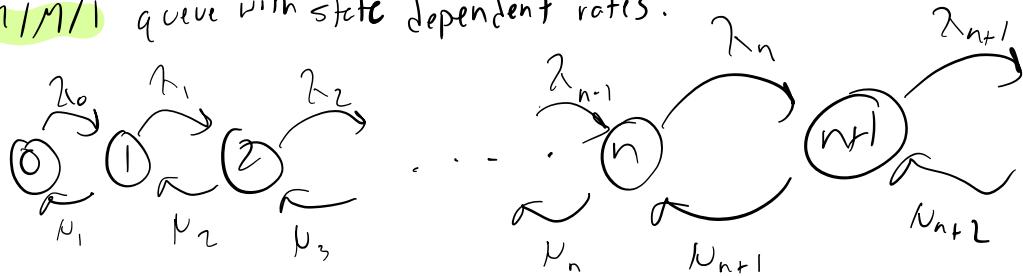
↑  
average # of  
customers in  
system

$E = \lambda_{\text{eff}} S \rightarrow M/M/1/K$

$$\lambda_{\text{eff}} = (1 - P_{\text{loss}}) \lambda$$

↑ total arrival

\* **M/M/1** queue with state dependent rates.



the rates  $\lambda_i$  depend on the underlying state of the system

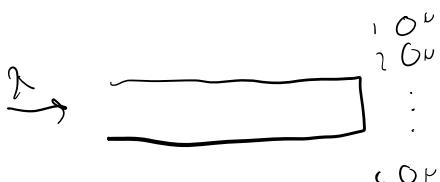
still birth/death process, so we use DBE;

$$\pi_n = \pi_0 \prod_{i=0}^{n-1} \frac{\lambda_i}{\mu_{i+1}}, \quad n \geq 1$$

$$\pi_0 = \left[ 1 + \sum_{n=1}^{\infty} \prod_{i=0}^{n-1} \frac{\lambda_i}{\mu_{i+1}} \right]^{-1}.$$

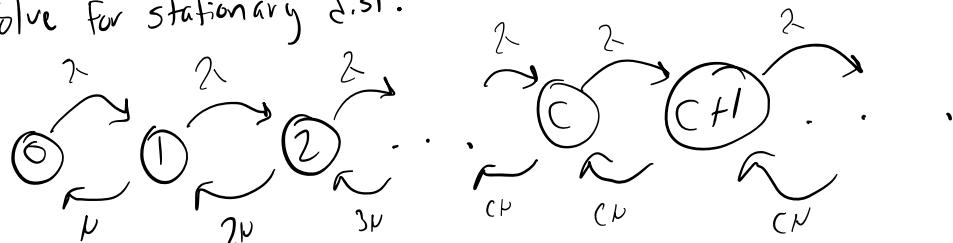
**M/M/c queue**

exponential interarrival times with rate  $\lambda$ , exponential service times with rate  $\mu$  for each server, and there are  $c$  parallel identical servers.



The total occupation rate/load  $P = \frac{\lambda}{c\mu} < 1$ .

Solve for stationary dist.



state transition diagram

We then get

$$\pi_n = \frac{(cp)^n}{n!} \pi_0, n=0, 1, \dots, c$$

$$\text{and } \pi_{cn} = p^n \pi_c = p^n \frac{(cp)^c}{c!} \pi_0, n=0, 1, \dots, \cancel{c}$$

$$\pi_{cn} = p^n \pi_c = p^n \frac{(cp)^c}{c!} \pi_0, \text{ normalization condition}$$

$$\pi_0 \text{ is calculated using the normalization condition}$$
$$\pi_0 = \left[ \sum_{n=0}^{c-1} \frac{(cp)^n}{n!} + \frac{(cp)^c}{c!} \cdot \frac{1}{1-p} \right]^{-1}$$

**Queueing probability** (the probability than an incoming packet/job/customer has to wait in the buffer).

$$\pi_w = \pi_c + \pi_{c+1} \dots = \frac{\pi_c}{1-p} = \frac{(cp)^c}{c!} \frac{1}{1-p} \pi_0$$

mean queue length:

$$E(L^q) = \sum_{n=1}^{\infty} n \cdot \pi_{cn} = \pi_w \frac{\rho}{1-p}$$

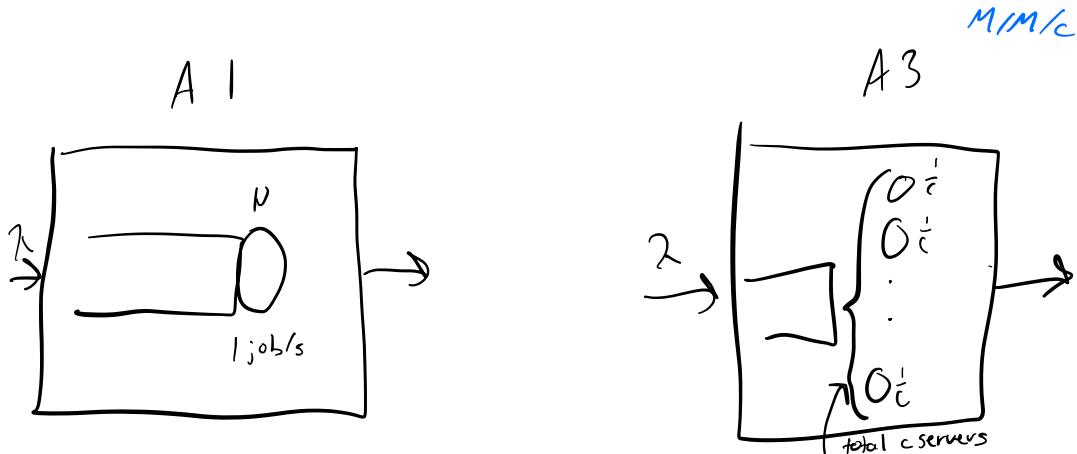
$n=1$   
↑  
ignore  
 $n=0$     ↑  
start at  
c+1 b/c now  
have backlog

mean waiting time:

$$E(W) = \frac{E(L^q)}{\lambda} = \pi_w \frac{1}{1-p} \frac{1}{\lambda p}$$

10/24/25

Mid 2 1/4



① Let's fix the incoming traffic rate  $\lambda = 0.9$  arrivals per second. and the load  $P = 0.9$  (ratio arrival/service, no rate). We will change the service rate of the servers in A3 with  $c$ . i.e., using less # of higher performance servers vs. using more # of lower performance servers.  $\leftarrow \rightarrow$  you have a choice on the type of servers.

# of servers	avg. waiting time	avg. queue length	avg. # of customers in the system	avg. sojourn time, total delay
--------------	-------------------	-------------------	-----------------------------------	--------------------------------

$c$	$E(W)$	$E(L^a)$	$E(L)$	$E(S)$
1	9	8.1	9	10
2	8.53	7.67	9.47	10.52
5	7.63	6.86	11.36	12.63
10	6.69	6.02	15.02	16.69

decreasing      decreasing      increasing      increasing  
 Clearly, putting more under performing servers is not desirable as the total delay will be higher. We prefer higher performance server.

(2) Let's fix  $n$  of each server. Say  $n=1$ . We fix  $P=0.9$ . (same kind of server and the same traffic load). Means we change  $\lambda$  with  $c$ .

prob of queueing

most important factor

$C$	$\Pi_w$	$E(W)$	$E(L^*)$	$E(L)$	$E(S)$
1	0.9	9	8.1	9	10
2	0.85	4.26	7.67	9.47	5.26
5	0.76	1.83	6.86	11.36	2.52
10	0.67	0.67	6.02	15.02	1.67
20	0.55	0.28	4.96	22.96	1.24

decreases      decreases      decreases      increases      decreases,  
↓                  ↓                  ↓                  ↑                  ↓                  ↓

Would rather have bigger team      smaller delay

Clearly, putting more servers leads to smaller queuing delay and overall delay. We prefer more servers.

Overall, the higher the performance of the server the better. The more servers you can have (when you cannot choose the type) the better.

Erlang-B formula:

Consider a circuit switching telephone network e.g. satellite phone network. There are  $m$  circuit (lines) provided by the satellite. So at most  $m$  calls can be supported by the satellite at the same

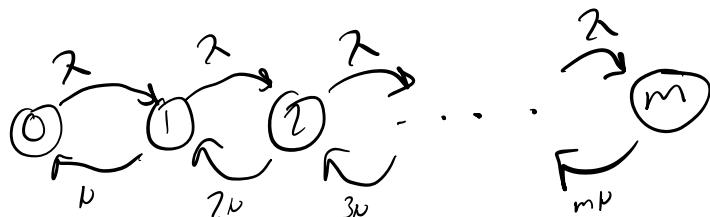
time. When these  $m$  lines are all being occupied, a dialog caller will be blocked, i.e., the caller will have a busy time when dials. What's the blocking probability of a user, suppose in total in the service area we have  $K$  users each user on average makes  $A$  calls a day and each call on average lasts  $B$  minutes?

$0$  buffer  
circuit is servers ( $n$ )

$M/M/m/0$ , the  $0$  is representing the buffer size

$$\lambda = \frac{KA}{24 \times 60} \text{ calls/min} \quad \mu = \frac{1}{B} \text{ calls/min}$$

$$\begin{matrix} 0 \\ \vdots \\ 0_m \end{matrix} \quad \begin{matrix} \text{DBE} \\ \text{birth/death process} \end{matrix}$$



$$\pi_n = \pi_0 \left( \frac{\lambda}{\mu} \right)^n \quad \text{For } n=0, 1, \dots, m$$

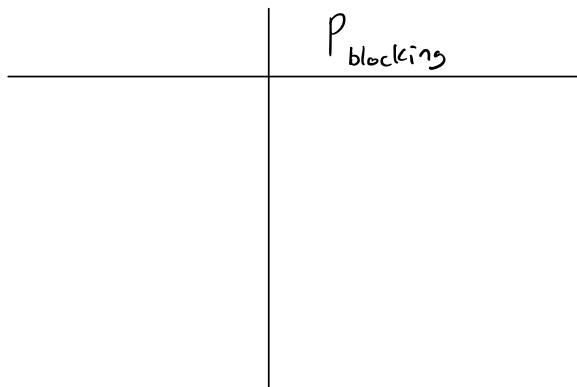
$$\pi_0 = \left[ \sum_{n=0}^m \frac{\left( \frac{\lambda}{\mu} \right)^n}{n!} \right]^{-1}$$

$$P_{\text{blocking}} = \pi_m = \frac{\left( \frac{\lambda}{\mu} \right)^m / m!}{\sum_{n=0}^m \left( \frac{\lambda}{\mu} \right)^n / n!} \quad \leftarrow \text{Erlang-B}$$

Ex.  $\lambda = 1$  calls/min

Finish notes

$$\mu = \frac{1}{3} \Rightarrow \frac{\lambda}{\mu} = 12$$



10/29/25

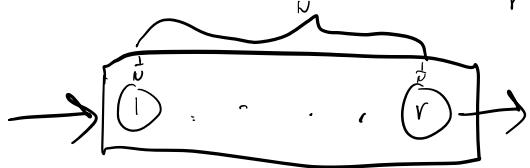
Midterm 2: 1/10

M/Erl/1

arrival - Poisson

service time - Erlong distributed

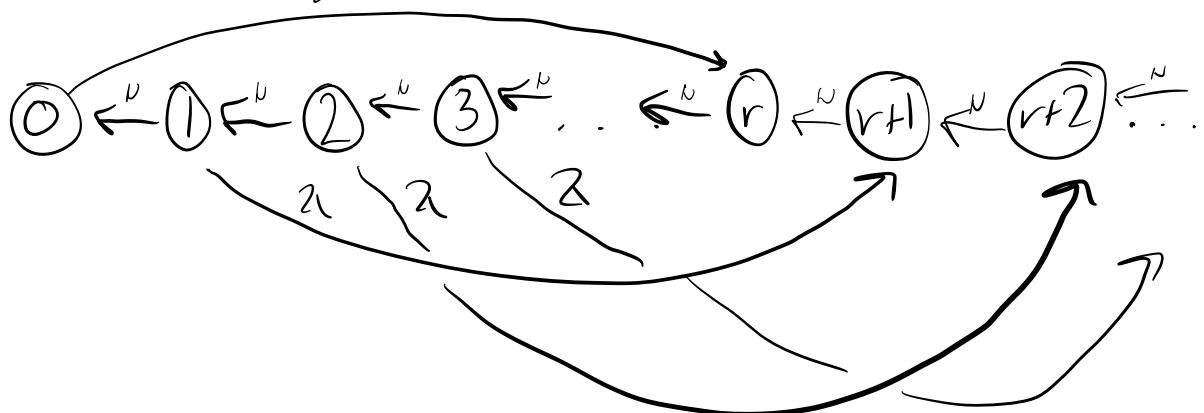
To service a job/packet, the server has to go through  $r$  stages sequentially, where each stage takes an exponentially distributed time ( $\mu$ ). The job arrival process is Poisson with rate  $\lambda$ . The load of the system  $P = \frac{r\lambda}{\mu}$ .



$$P = \frac{\text{arrival rate}}{\text{service rate}}$$

**Definition of the state:** the natural way to describe the state of nonempty system is by the pair  $(K, l)$ , where  $K$  denotes the # of jobs in the system and  $l$  denotes the remaining # of stages of the job that is in service. This is a 2-dimensional MC. An alternative 1-D way is to describe the state by counting the total # of uncompleted stages that the jobs in the system still need to finish. There is a one-to-one mapping between this number (say  $n$ ) and the pair  $(K, l)$ .

$n = (K-1) \cdot r + l$ . So we can use the following 1-D model.



Let  $\tilde{\pi}_n$  denote the stationary prob. of  $n$  stages in the system. Based on GBE, the flow going out of state  $i$  is equivalent to total flows going into that state  $i$ .

$$\tilde{\pi}_0 u = \tilde{\pi}_1 u$$

$$\tilde{\pi}_1 (\lambda + u) = \tilde{\pi}_2 u$$

$$\tilde{\pi}_n (\lambda + u) = \tilde{\pi}_{n+1} u \quad \text{for } n = 1, \dots, r-1.$$

$$\tilde{\pi}_n (\lambda + u) = \tilde{\pi}_{n-r} \lambda + \tilde{\pi}_{r+1} u \quad \text{for } n = r, r+1, \dots$$

Below is not tested on

If can be shown that

$$\bar{Y}_i = \sum_{k=1}^r c_k x_k^i, \text{ where } i=0, 1, 2, \dots$$

where  $x_1, \dots, x_r$  where the  $r$  are the distinct roots to the

following equation:

$$(\lambda + \mu) \cdot x^r = \lambda + \mu x^{r+1} + |x| \lambda \cdot$$

the coefficients

$$c_k = \frac{1 - \rho}{\prod_{j \neq k} \left(1 - \frac{x_j}{x_k}\right)}, \quad k=1, \dots, r$$

---

Distribution of # of jobs:  $K$   $q_k \stackrel{\Delta}{=} \Pr\{k \text{ jobs}\}$

$$q_0 = \bar{Y}_0$$

$$q_1 = \bar{Y}_1 + \dots + \bar{Y}_r$$

:

$$q_i = \sum_{n=(i-1)r+1}^{ir} \bar{Y}_n.$$

**M/G/1 queue**: the jobs arrive according to Poisson process with rate  $\lambda$ , the service time is iid with CDF  $F_B(t)$  and pdf  $f_B(t)$ . To be stable, we require  $\rho = \lambda E(B) < 1$ .

10/31/25

**Generating Function**:

Let the probability mass function of a discrete non-negative r.v.

$X$  be  $\Pr(X=n) = p(n)$ , where  $n=0, 1, 2, \dots$ . Then the moment generating function of  $P_X(z)$  or

$X$  is defined as:

$$P_X(z) = E(z^X) = \sum_{n=0}^{\infty} z^n p(n).$$

(Recap :  $z$ - transform of  $p(n)$  is  $P(z) = \sum_{n=0}^{\infty} p(n) z^{-n}$ )

**Properties:** (a)  $|P_X(z)| \leq 1$  for all  $|z| \leq 1$ .

(b)  $P_X(0) = p(0)$ ,  $P_X(1) = 1$ .

(c)  $P_X'(z) \Big|_{z=1} = E(X)$ .

$$P_X'(z) = \sum_{n=0}^{\infty} n z^{n-1} p(n)$$

$$\text{So } P_X'(1) = \sum_{n=0}^{\infty} n p(n) = E(X)$$

$$P_X^{(2)}(z) = \sum_{n=0}^{\infty} n(n-1) z^{n-2} p(n)$$

$$P_X^{(2)}(1) = \sum_{n=0}^{\infty} n(n-1) p(n) = E(X(X-1)) = E X^2 - E X$$

$$\text{So in general: } P_X^{(k)}(1) = E[X(X-1) \cdots (X-k+1)]$$

$M/G/1$  queue The service time is iid with CDF  $F_B(\tau)$ , and pdf  $f_B(\tau)$ .

The moment generating function of the number of jobs/products in the system:

$$P_L(z) = \frac{(1-p)\tilde{B}(\lambda - \lambda z)(1-z)}{\tilde{B}(\lambda - \lambda z) - z} \text{ where } p = \lambda E(B)$$

$\tilde{B}(s)$  is the Laplace transform of the service time distribution.

$$\tilde{B}(s) = \int_0^\infty f_B(t) e^{-st} dt.$$

Distribution of the sojourn time:

$$\tilde{S}(s) = \frac{(1-p)\tilde{B}(s) \cdot s}{\lambda \tilde{B}(s) + s - \lambda}$$

Distribution of the waiting time:

$$\tilde{W}(s) = \frac{(1-p) \cdot s}{\lambda \tilde{B}(s) + s - \lambda}$$

All above are the Pollaczek-Khintchine formula

① mean waiting time in the buffer.

$$\bar{W} = \frac{\lambda E(x^2)}{2(1-p)} \quad Var(X) = E(x^2) - (E(x))^2$$

$$\Rightarrow E(x^2) = Var(X) + (E(x))^2$$

$$= \frac{p_f \lambda N Var(X)}{2(N-2)}$$

② mean sojourn time

$$\bar{S} = \bar{W} + E(x) = \frac{1}{N} + \frac{p_f \lambda N Var(X)}{2(N-2)}$$

service time

(3) mean queue length

$$\bar{L}_q = \lambda \bar{W} = \frac{\lambda^2 \bar{x}^2}{2(1-p)} - \frac{p^2 + \lambda^2 \text{Var}(x)}{2(1-p)}$$

(4) mean number of packets in the system

$$\bar{L} = p_r \bar{L}_q = p_r \frac{p^2 + \lambda^2 \text{Var}(x)}{2(1-p)}$$

Example: M/D/1 queuing system: Poisson arrival with rate  $\lambda$ , deterministic service time  $\frac{1}{\mu}$

$$\bar{L} = p + \frac{1}{2} \frac{p^2}{1-p}$$

$$\bar{L}_q = \frac{1}{2} \frac{p^2}{1-p}$$

M/M/1 queue same  $\lambda, \mu$

$$\bar{L}_q = \frac{p^2}{1-p}$$

*Not required*

$$\bar{S} = \frac{1}{\mu} + \frac{p}{2\mu(1-p)}$$

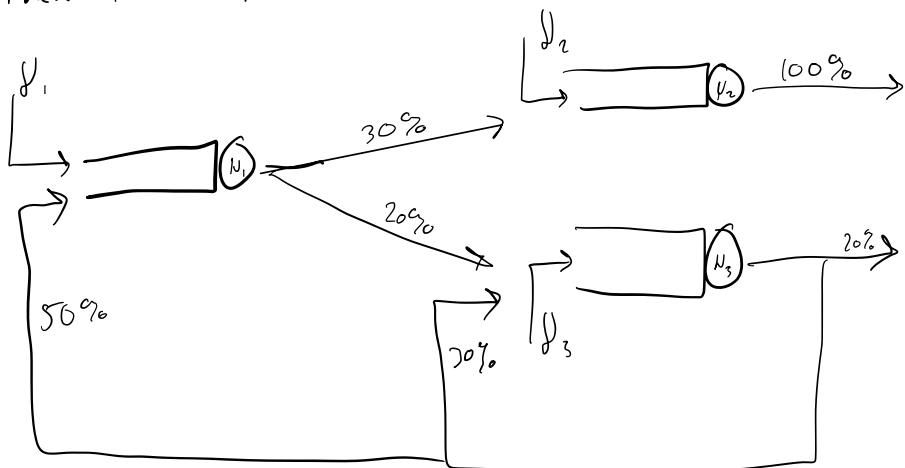
$$\bar{W} = \frac{p}{2\mu(1-p)} \leftarrow \text{half of } \bar{W} \text{ for M/M/1}$$

**Queuing networks**: a set of interconnected queues.

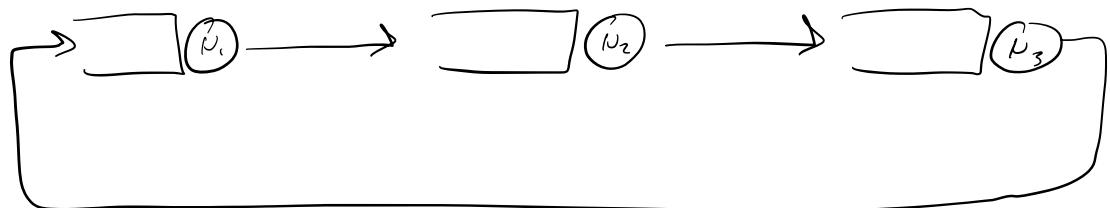
Types of queuing networks

- Open networks - external traffic
- Closed networks
- Hybrid networks

**open networks:** traffic comes from outside the system, and are serviced and then leave/depart from the network.



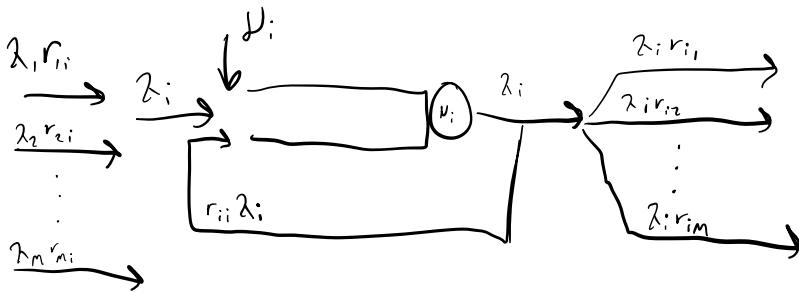
**Closed networks:** a fixed number of customers ( $k$ ) are trapped in the system and circulate among the queues.



11/3/25

Only worry about open networks

Consider an open network; assume an arbitrary network of  $M$  queues with infinite buffer capacity. Jobs arriving from outside the system are serviced and then depart. Note that a job may visit several queues before departing, including possibly visiting some queues more than once. Service time of queue is generally distributed with rate  $\mu_i$  for queue  $i$ . Arrivals from outside to queue  $i$  occur according to general iid process with mean rate  $\lambda_i$ . The total mean job arrival rate to queue  $i$  is denoted by  $\bar{\lambda}_i$ .



$r_{ij}$  - routing probability, fraction of packets, that a job completing service at queue  $i$  goes to queue  $j$ .

$r_{im}$  - routing probability that a job completing service at queue  $i$  leaves the network (goes to outside).

traffic conservation condition:  $\sum_{j=1}^{M+1} r_{ij} = 1$ .

Let  $\bar{\lambda}_i$  be the total mean arrival rate of queue  $i$ .

$$\bar{\lambda}_i = \lambda_i + \sum_{j=1}^M \lambda_j r_{ji}$$

Denote  $\vec{\lambda} \triangleq [\lambda_1 \ \lambda_2 \ \dots \ \lambda_M]$

$$\vec{\mu} \triangleq [\mu_1 \ \mu_2 \ \dots \ \mu_M] \quad \text{routing matrix}$$

$$\vec{R} \triangleq [r_{ij}] \quad (1 \leq i \leq M, 1 \leq j \leq M)$$

The flow conservation equation can be written in matrix form

$$\vec{\lambda} = \vec{d} + \vec{\lambda} \vec{R} \Rightarrow \vec{\lambda} [I - \vec{R}] = \vec{d}$$

$$\vec{\lambda} = \vec{d} [I - \vec{R}]^{-1}$$

identity matrix  
routing matrix

↑  
external traffic

It relates external arrival rates and routing to determine the total flow rate at each queue.

A special type of open network is the **Jackson network**.

- this is a network of  $M$  queues
- there is only one class of jobs in the network
- All service times are exponentially distributed with service rate  $\mu_i$  at queue  $i$ . The service discipline at all queues is FCFS (First come First serve). All external arrival processes are Poisson with rate  $\lambda_i$  at queue  $i$ .

Because the merge of independent Poisson processes is Poisson with rate equal to the sum of the individual rates. The departure process of an  $M/M/1$  queue is also Poisson with rate equal to the input rate of the queue (i.e.,  $\lambda_i$ ), and probabilistic splitting of Poisson process results in Poisson processes for each of the sub flows. Therefore, the input and output processes of each queue  $i$  in the network is a Poisson process.

**Solution for Stationary Distribution:** Let  $\tilde{n}_i(t)$  be the number of jobs at queue  $i$  at time  $t$ . The state of the network is defined by the vector  $(\tilde{n}_1(t), \tilde{n}_2(t), \dots, \tilde{n}_m(t))$ . Then  $\{(\tilde{n}_1(t), \tilde{n}_2(t), \dots, \tilde{n}_m(t)), t \geq 1\}$ .

is a  $M$ -dimensional Markov process.

The stationary distribution,  $\vec{n} = (n_1, n_2, \dots, n_m)$

$$\pi(\vec{n}) = \lim_{t \rightarrow \infty} \text{Prob}\{\tilde{n}_1(t) = n_1, \tilde{n}_2(t) = n_2, \dots, \tilde{n}_m(t) = n_m\}$$

If can be shown that the stationary state probability has the following form:

$$\pi(\vec{n}) = C P_1^{n_1} P_2^{n_2} \dots P_m^{n_m} = C \prod_{i=1}^m P_i^{n_i} \quad \text{where } P_i = \frac{x_i}{\mu_i}$$

and  $C$  is some constant

normalization condition:  $\sum_{n_1=0}^{\infty} \sum_{n_2=0}^{\infty} \dots \sum_{n_m=0}^{\infty} C P_1^{n_1} P_2^{n_2} \dots P_m^{n_m} = 1$

$$\Rightarrow C = \prod_{i=1}^M (1 - P_i)$$

Therefore, your  $\pi(\vec{n}) = \prod_{i=1}^M (1 - P_i) P_i^{n_i}$

recall that for an  $M/M/1$  queue  $\pi_n = (1 - p)p^n$

$\therefore \pi(\vec{n}) = \prod_{i=1}^M \pi_{n_i}$  ← the product of the stationary probability of  $M$  independent  $M/M/1$  queues.

1/5/25

**Jackson's Theory**: If in an open Jackson network  $\lambda_i < \mu_i$

holds for all the queue  $i = 1, \dots, M$ , then,

① The arrival rates  $\lambda_i$  can be computed by  $\vec{\lambda} = \vec{D}(\vec{I} - \vec{R})^{-1}$

② The stationary probability of the network can be expressed as the product of the stationary probability of the individual queues.

$$\pi(k_1, k_2, \dots, k_m) = \pi_1(k_1) \pi_2(k_2) \dots \pi_m(k_m)$$

③ The nodes of the network can be considered as independent M/M/1 queues with arrival rate (aggregate)  $\lambda_i$  and service rate  $\mu_i$ , respectively.

**Performance matrix:** Because each queue  $i$  in M/M/1 with

$$\overline{L}_i = \frac{P_i}{1-P_i}, \quad \overline{W}_i = (1-P_i)P_i^{-1}, \quad \overline{S} = \frac{1}{\mu_i - \lambda_i}$$

$\uparrow$   
 average # of packets  
 in the system

$\uparrow$   
 Average  
 sojourn time

$$\overline{W} = \frac{P_i}{\mu_i - \lambda_i}$$

$\uparrow$   
 Average  
 Waiting time  
 in buffer

For the network as a whole

$L_N$  - average # of packets in the networks

$W_N$  - average delay for a packet to go through the network

$$L_N = \sum_{i=1}^M \overline{L}_i$$

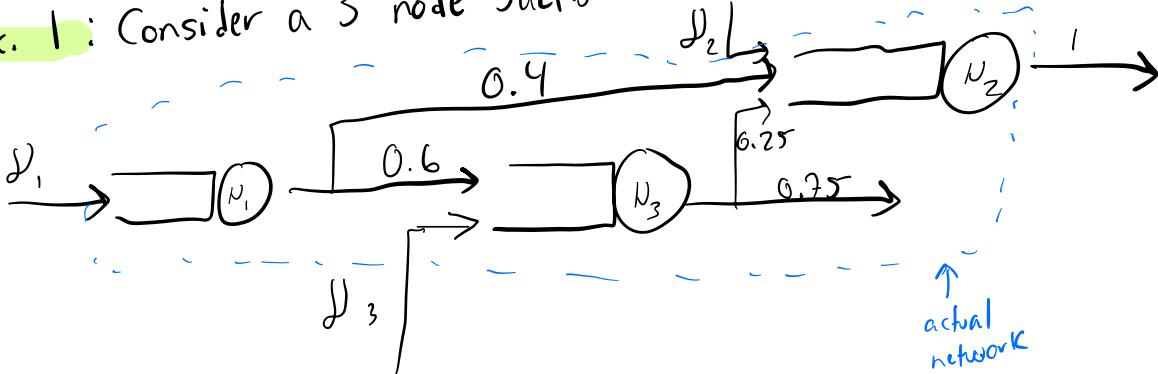
Use Little's law to find  $W_N$

$\lambda_N$  - total external traffic arrival rate to the network

$$\lambda_N = \sum_{i=1}^m \lambda_i$$

$$\overline{W}_N = \frac{\sum_{i=1}^m \frac{\rho_i}{1-\rho_i}}{\sum_{i=1}^m \lambda_i}$$

Ex. 1: Consider a 3 node Jackson network



Poisson external arrival with rates  $\lambda_1 = 0.5$ ,  $\lambda_2 = 0.25$ ,  $\lambda_3 = 0.25$ .

Exponential service times at each queue/server:  $\mu_1 = 1$ ,  $\mu_2 = 1$ ,  $\mu_3 = 1$ .

Average delay?

The routing matrix:  $r_{12} = 0.4$     $r_{11} = 0$     $r_{13} = 0.6$

$$\begin{matrix} r_{21} = 0 & r_{22} = 0 & r_{23} = 0 \\ r_{31} = 0 & r_{32} = 0.25 & r_{33} = 0 \end{matrix}$$

$$\begin{aligned} r_{11} &= 0 && \text{2 don't need} \\ r_{21} &= 1 \\ r_{31} &= 0.75 \end{aligned}$$

$$\vec{R} = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} - \begin{bmatrix} 0 & 0.4 & 0.6 \\ 0 & 0 & 0 \\ 0 & 0.25 & 0 \end{bmatrix}$$

$$\vec{\lambda} = \vec{J}(I-R)^{-1} \Rightarrow \vec{\lambda} = [0.5 \quad 0.5875 \quad 0.55]$$

$$P_1 = \frac{\lambda_1}{\nu_1} = 0.5 \quad P_2 = \frac{\lambda_2}{\nu_2} = 0.5875 \quad P_3 = \frac{\lambda_3}{\nu_3} = 0.55$$

$$J_N = J_1 + J_2 + J_3 = 0.5 + 0.25 + 0.25 = 1$$

$$\bar{W}_N = \frac{\bar{L}_N}{J_N} = \frac{1}{1} = 3.646 \text{ sec.}$$

$$\frac{0.5}{1-0.5} + \frac{0.5875}{1-0.5875} + \frac{0.55}{1-0.55} = 1 + \frac{0.5875}{1-0.5875} + \frac{0.55}{0.45}$$