

STOCHASTIC PROCESSES

LECTURE 19: PERFORMANCE MEASURES,  
LITTLE'S LAW

Hailun Zhang@SDS of CUHK-Shenzhen

April 7, 2021

# Time-Average Performance Measures

$f(i)$  = “cost” or “reward” for being in state  $i$

What's the long-run average cost/reward?

## THEOREM (STRONG LAW OF LARGE NUMBERS)

If the CTMC  $\{X(t), t \geq 0\}$  with state space  $S$  is irreducible and positive recurrent, then for any  $f : S \rightarrow [0, \infty)$ ,

$$\mathbb{P} \left\{ \lim_{T \rightarrow \infty} \underbrace{\frac{1}{T} \int_0^T f(X(t)) dt}_{\text{long-run average}} = \sum_{i \in S} \pi_i f(i) \right\} = 1,$$

where  $\pi$  denotes the stationary distribution of the CTMC.

## Example: M/M/1 Queue

$$\lambda < \mu$$

Some Performance Measures:

- $f(i) = i \xrightarrow{\text{SLLN}}$  with probability 1,

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

long-run *average number of customers in sys.*  $= \sum_{i=0}^{\infty} i \pi_i = \frac{\lambda}{\mu - \lambda}$

- $f(i) = \mathbf{1}\{i > 0\} \xrightarrow{\text{SLLN}}$  with probability 1,

utilization

long-run *fraction of time the server is busy*  $= \sum_{i=1}^{\infty} \pi_i = \frac{\lambda}{\mu} \triangleq \rho$

- $f(i) = \mathbf{1}\{i = j\} \xrightarrow{\text{SLLN}}$  with probability 1,

long-run *fraction of time there're  $j$  customers in the system*  $= \pi_j$

# Headcount average performance measures

- $S_i$  be the time in system (waiting + service) of the  $i$ th customer. S<sub>1</sub>, S<sub>2</sub>, S<sub>3</sub>, ...
- average time in system n=100

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n S_i.$$

- SLLN for the arrival process: for a Poisson arrival process with rate  $\lambda > 0$ ,

$$\mathbb{P}\left\{\lim_{t \rightarrow \infty} \frac{N(t)}{t} = \lambda\right\} = 1. \quad (1)$$

- SLLN for  $X = \{X(t), t \geq 0\}$

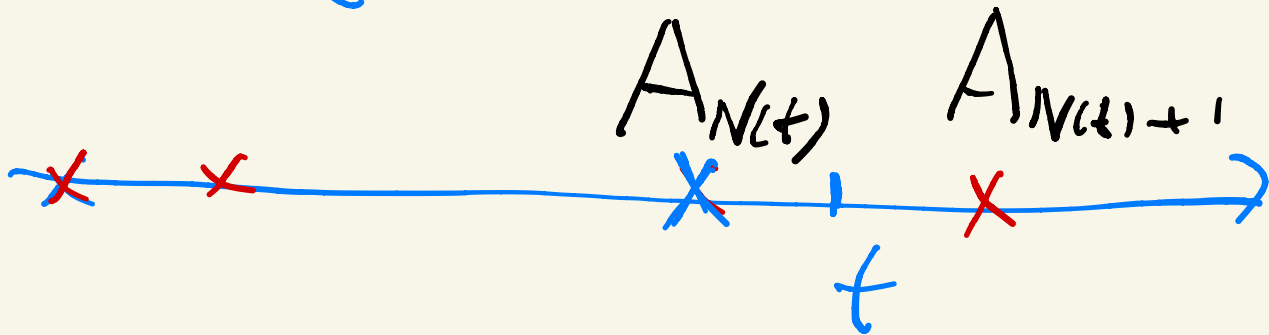
$$\mathbb{P}\left\{\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t X(s) ds = \frac{\rho}{1 - \rho}\right\} = 1. \quad (2)$$

- We claim with probability 1,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n S_i = \frac{1}{\lambda} \frac{\rho}{1 - \rho}.$$

=  $\frac{\rho}{\lambda(1-\rho)}$   
why?

$$\frac{N(t)}{t} \rightarrow \lambda \quad t \rightarrow \infty$$



Interarrival  $Z_1, Z_2, \dots$

$A_n = n^{\text{th}} \text{ arrival time}$

$$= \sum_{i=1}^n Z_i$$

$$\frac{A_n}{n} \rightarrow \frac{1}{\lambda}$$

$$\underline{A_{N(t)} \leq t < A_{N(t)+1}}$$

$$\frac{N(t)}{A_{N(t)+1}} < \frac{N(t)}{t} \leq \frac{N(t)}{\underline{A_{N(t)}}}$$

$$\frac{N(t)}{A_{N(t)}} \rightarrow \lambda \quad (t \rightarrow \infty, N(t) \rightarrow \infty)$$

# Little's Law



$L$  = long-run *average number of customers* in the queue/system

$\lambda$  = long-run *average arrival rate* (or throughput of the system)

$W$  = long-run *average amount of time a customer waits* in the queue/system

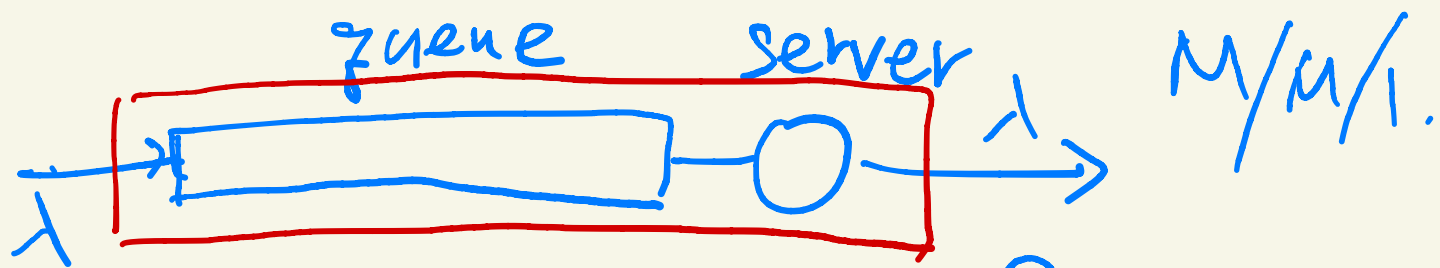
$$W = \frac{L}{\lambda}$$

## THEOREM (LITTLE'S LAW)

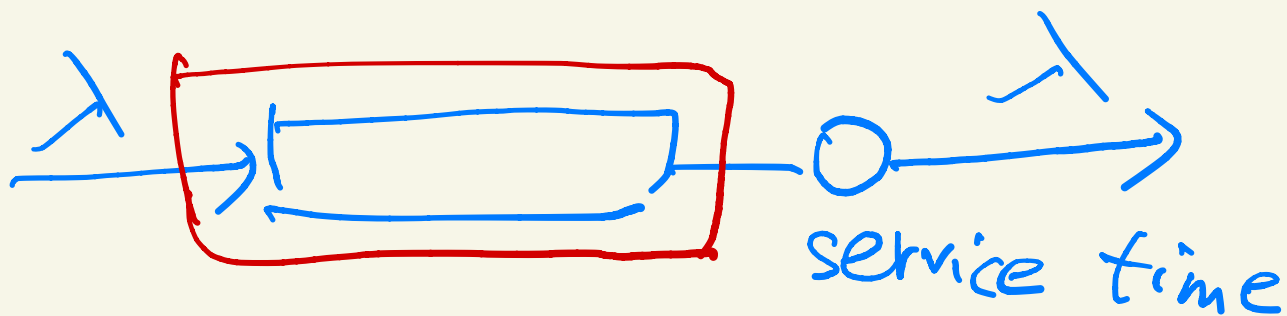
*If two quantities exist (well defined), the third quantity also exists.*

*Furthermore, they satisfy*

$$L = \lambda W.$$



$$W = \frac{L}{\lambda} = \frac{\frac{\rho}{1-\rho}}{\lambda}$$

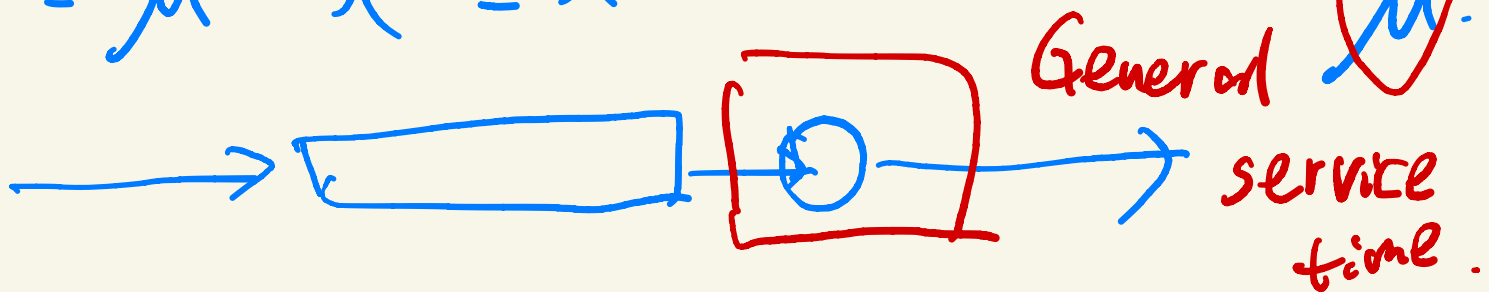


$$\lambda' = \lambda \quad \text{mean } \frac{1}{\mu}$$

$$W' = \text{average time in queue}$$

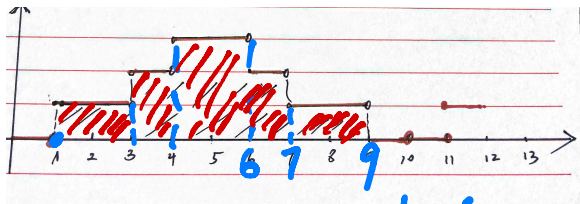
$$= \frac{L'}{\lambda'} = ?$$

$$W'' = \frac{1}{\mu} \quad \lambda'' = \lambda \quad L'' = \lambda'' W'' = \frac{\lambda}{\mu}$$



# An illustration

$$W = \frac{L}{\lambda}$$



- $t = 10$ ,  $N(t) = 3$ ,  $\lambda = N(t)/t = 3/10$ .
- $L$

$$L = \frac{1}{10} \int_0^{10} X(s) ds = \frac{1}{10} [1(8) + (4) + (2)] = \frac{14}{10}.$$

- $W$

$$W = \frac{W_1 + W_2 + W_3}{3}$$

$$\underline{W_1} = (6 - 1) = 5, \quad \underline{W_2} = 7 - 3 = 4, \quad \underline{W_3} = 9 - 4 = 5, \quad W = \frac{14}{3}.$$

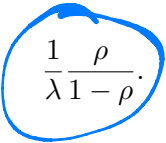


# Average time-in-system and waiting time in $M/M/1$ system

- Average number in system is

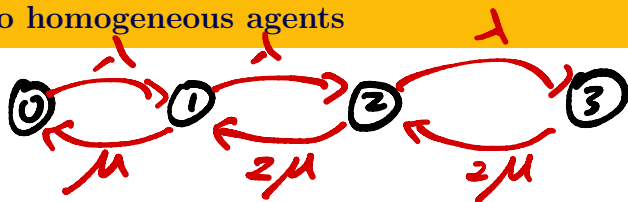
$$\frac{\rho}{1 - \rho}.$$

- Average time-in-system


$$\frac{1}{\lambda} \frac{\rho}{1 - \rho}.$$

Average time-in-queue ?

## Three lines, two homogeneous agents



Consider a call center with two homogeneous agents and 3 phone lines. Arrival process is Poisson with rate  $\lambda = 2$  calls per minute. Processing times are iid exponentially distributed with mean 4 minutes.  $\mu = \frac{1}{4}$

- What is the long-run fraction of time that there are no customers in the system?  $\pi_0$
- What is the long-run fraction of time that both agents are busy?  $\pi_2 + \pi_3$
- What is the long-run fraction of time that all three lines are used?

$$\pi_3$$

$$\pi = (\pi_0, \dots, \pi_3)$$

## Solution

$$\begin{aligned} \tau &= 100 \\ n &= 100 \end{aligned}$$

$$\begin{aligned} 40 \\ 80 \end{aligned}$$

$$\begin{aligned} 90\% \\ 80\% \end{aligned}$$

- $X(t)$  is the number of calls in the system at time  $t$ .  $S = \{0, 1, 2, 3\}$ .
- flow in = flow out in each state.

$$2\pi_0 = \frac{1}{4}\pi_1, \quad 2\pi_1 = \frac{1}{2}\pi_2, \quad 2\pi_2 = \frac{1}{2}\pi_3, \quad \pi_0 + \pi_1 + \pi_2 + \pi_3 = 1$$

- Solving this by setting  $\pi_0 = 1$  and normalizing the result, we obtain

$$\pi = (1, 8, 32, 128) \Rightarrow \pi = \left( \frac{1}{169}, \frac{8}{169}, \frac{32}{169}, \frac{128}{169} \right).$$

- What is the long-run fraction of time that there are no customers in the system?  $\pi_0 = 1/169$
- What is the long-run fraction of time that both agents are busy?  
 $\pi_2 + \pi_3 = 160/169$
- What is the long-run fraction of time that all three lines are used?  
 $\pi_3 = 128/169$

## Other performance measures



- The number of calls lost per minute is  $\lambda\pi_3 = 2(128/169)$  which seems to be quite high.
- The throughput of the system is  $\lambda(1 - \pi_3)$ .
- The long-run fraction of calls that are lost is  $\pi_3$ ?
- PASTA property

Head count

$$n=3$$

*fraction of customers lost =  $\pi_3$*

### THEOREM (POISSON ARRIVALS SEE TIME AVERAGES)

*Suppose customers arrive at a queueing system according to a Poisson process. Then for any  $n \in \{0, 1, \dots\}$ , the*




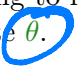



*long-run fraction of arrivals that see  $n$  customers in the system equals the*

*Headcount*

*long-run fraction of time that there are  $n$  customers in the system.*

*Time Average*

## Three lines, two non-homogeneous agents

- 3 *phone lines*, 2 *agents* (Alice & Bob)
- Incoming calls are routed to Alice if possible.  

- Calls *arrive* according to a **Poisson process** with rate  $\lambda$ .
- Alice's *processing times* are iid **exponential** with rate  $\mu_A$ .  

- Bob's *processing times* are iid **exponential** with rate  $\mu_B$ .  

- Times that callers are willing to hold (i.e., their *patience times*) are iid **exponential** with rate  $\theta$ .  
  
  
  


# The Corresponding CTMC

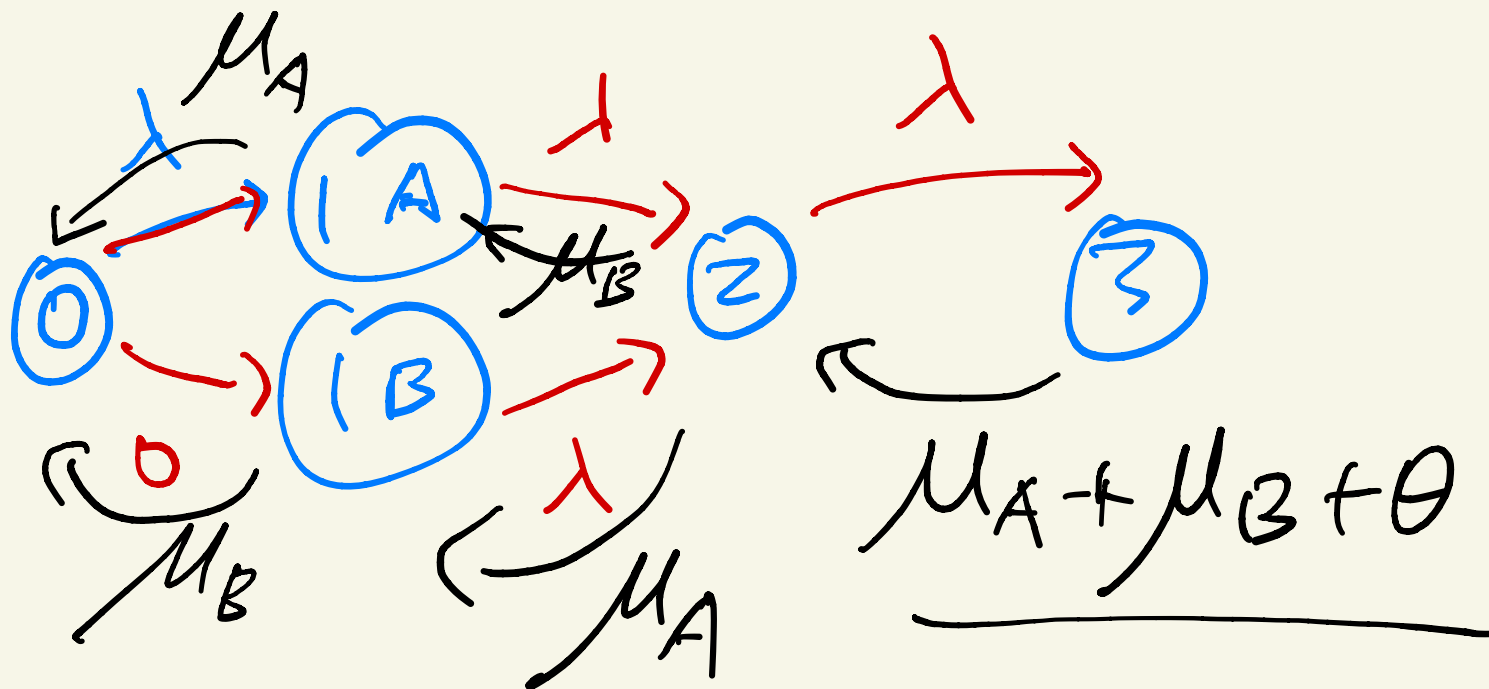
$\{0, 1, 2, 3\}$  ✗

*State Space* =  $\{0, \underline{1A}, \underline{1B}, 2, 3\}$

- 0 = no calls in the system
- 1A (resp. 1B) = 1 call in the system, with Alice (resp. Bob)
- 2 (resp. 3) = 2 (resp. 3) calls in the system

*Generator Matrix* (rows correspond to states in the order listed above)

$$G = \begin{bmatrix} -\lambda & \lambda & 0 & 0 & 0 \\ \mu_A & -(\lambda + \mu_A) & 0 & \lambda & 0 \\ \mu_B & 0 & -(\lambda + \mu_B) & \lambda & 0 \\ 0 & \mu_B & \mu_A & -(\lambda + \mu_A + \mu_B) & \lambda \\ 0 & 0 & 0 & \mu_A + \mu_B + \theta & -(\mu_A + \mu_B + \theta) \end{bmatrix}$$





# Stationary Distribution

The stationary distribution  $\pi = [\pi_0, \pi_{1A}, \pi_{1B}, \pi_2, \pi_3]$  satisfies

$$\pi G = 0 \quad \text{and} \quad \pi_0 + \pi_{1A} + \pi_{1B} + \pi_2 + \pi_3 = 1.$$

Use this to solve for  $\pi$ .

- e.g., write all the  $\pi_i$ 's in terms of  $\pi_{1B}$ , and use the fact that they should sum to 1.

# Some Performance Measures

What is the

- long-run fraction of time that both Alice and Bob are free?

$$\pi_0$$

- long-run fraction of time that Bob is free?

$$\pi_0 + \pi_{1,A}?$$

- long-run fraction of arrivals that get a busy signal?

Headcount

$$\pi_3$$

PASTA