

Unit-3 Queuing System

DATE: _____

3.1. Characteristics and structure of Basic Queuing System & Models of queuing system

* Concept of Basic Queuing System

- A queuing system is a facility consisting of one or several servers designed to perform certain tasks or process certain jobs and a queue of jobs waiting to be processed.
- Jobs arrive at the queuing system, wait for an available server, get processed by the server, and leave.
- Ex: An Isp whose customers connect to the internet, browse & disconnect.

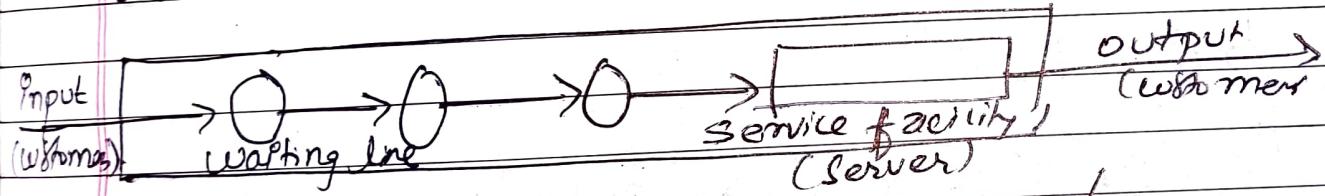


fig: structure of a basic queuing system (queuing model)

* Characteristics or elements of the queuing system

The elements or characteristics of queuing system are

- Calling population
- System capacity
- Arrival process
- Queue behaviour / Queue discipline.
- Service process

Note: Customer - people, machine, patients,

Server - receptionist, mechanics, doctors

1. Calling population.

It is the population of potential customers. The calling population may be finite or infinite population model. The assumption of an infinite population such that the rate of arrival of customers is not affected by the number of customers that have already joined the system. It means, the rate of arrival is constant throughout time.

2. System capacity

In many queuing systems, there is a limit to the no. of customers that may be in the waiting line or site. An arriving customer who finds the system full does not enter but returns immediately to the calling population. However, there are other systems that have infinite capacity. So, system capacity is either limited or unlimited.

3. Arrival process.

The arrival process for infinite popn models is usually characterized in terms of inter-arrival times of successive customers. Arrival may occur at scheduled times or random times. When at random times, the inter-arrival time is usually characterized by a probability distribution.

4. Queue Behaviour & Queue Discipline.

Queue behaviour refers to the actions of customers while in a queue waiting for service to begin. There are possibly walk, renege or jockey queue behaviour.

Queue discipline refers to the logical ordering of customers in a queue & determines which customer will be chosen for service when a server becomes free. Common queue disciplines are FIFO, LIFO, Service in Random Order (SRO), Priority and so on. SPT (Shortest process time)

5. Service Process.

Once the entities/customers have entered the system, they must be served. The physical meaning of "service" may vary from system to system. From the modeling point of view, we care about whether customers are processed in FCFS order or according to some kind of priority, etc.

Markov/exponential service process: It is a special service process in which entities are processed one at a time in FCFS order and serve time is independent & exponential. It is a memoryless service process.

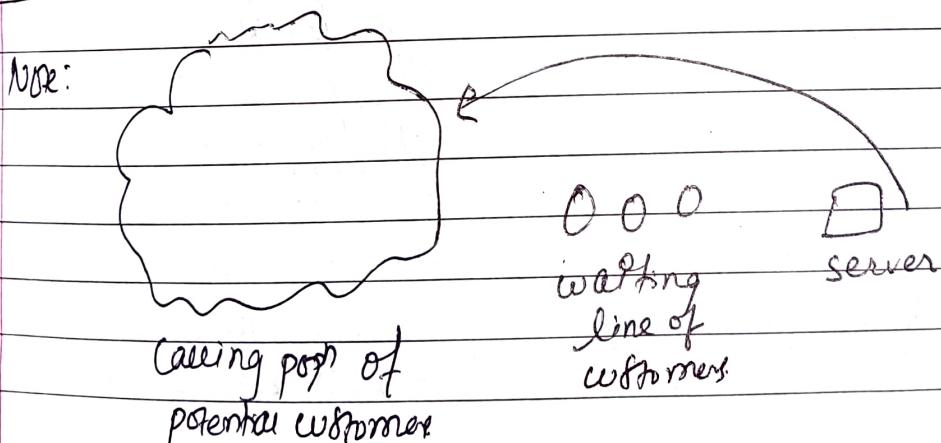


fig. Simple queuing model.

Example of queuing systems:

Systems	Customers	Servers
Reception desk	People	Receptionist
Garage	Trucks	Mechanic
Airport	Planes	Runway

8.2 Queuing Notation

X Kendall's notation for queuing system

Different notations are frequently used in queuing system and kendall's notation is one of them. kendall proposed a notational system for parallel servers systems which has been widely adopted. It can be represented in the form $A/B/c/N/K/D$ where.

$A \rightarrow$ interarrival time distribution

$B \rightarrow$ Service time distribution

$c \rightarrow$ No. of parallel servers

$N \rightarrow$ System capacity

$K \rightarrow$ Size of calling population

$D \rightarrow$ Queuing discipline.

Common symbols for $A \times B$ include M : Exponential/Markov

D : Constant or Deterministic

Ex: Erlang of order k

P_H : Phase-type

H : Hyper-exponential

G : Arbitrary or General

G_I : General Independent

When final three parameters are not specified

e.g. $M/M/1$, it is assumed to be $N \geq 0$, $K = \infty$, and $D = M$

Ex: $M/M/1$ or $M/M/1/\infty/\infty$ both same

$M \rightarrow$ Exponential/Markov distributed interarrival time

$M \rightarrow$ Exponential / Markov distributed service time

$1 \rightarrow$ Single server system

$\infty \rightarrow$ Unlimited system capacity

$\infty \rightarrow$ Infinite calling popⁿ with FIFO queuing discipline.

GI/M/1/100

- General time dist.
- Markov Service time dist.
- 5 parallel Servers system
- 100 system capacity
- ∞ calling popn
- FIFO queuing discipline

PH/GII/2/10/100/PR

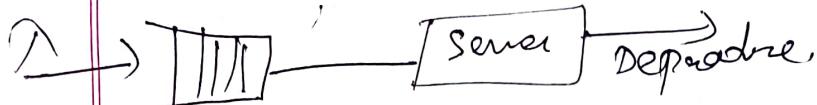
- Phase-type time dist. ana
- General Indep. service time dist.
- 2 parallel servers system
- 100 system capacity
- ∞ calling popn
- Priority queuing discipline

3.3 Single Server and Multiserver queuing system

Single Server queuing system

- It is a queuing system with only one server for any number of clients.
- It is a FIFO queuing system with Kendall notation MIM/I with poisson input, exponential service time and unlimited waiting positions.
- The model is based on following assumption.
 - The arrival follows poisson distribution with mean arrival rate λ
 - The service time has exponential distribution, average service rate μ
 - Arrivals are infinite population
 - Customers are served on FIFO basis.
 - There is only a single server.

In a single server queuing system, there is an infinite number of waiting positions in the queue. Hence, there can be any number of customers in the queue. So, it becomes a challenge to maintain the service rate in such a way to matchup with continuous arrival of customers.

DATE

$$Q = \begin{pmatrix} -\lambda & \lambda & & \\ u & -(\mu+u) & \lambda & \\ & u & -(\mu+u) & \lambda \\ & & u & -(\mu+u) \end{pmatrix}$$

The model is considered stable only if $\lambda < u$. We write $p = \frac{\lambda}{u}$ for the utilization of buffer & require $p < 1$ for the queue to be stable.

Performance of M/M/1 system: (measurement of queue system per unit time)

Given, λ = Arrival rate of customers
 u = Service rate of the server

- P_0 = Probability that there are no customers in the system. (Server idle ratio)

$$P_0 = 1 - \frac{\lambda}{u}$$

- P_n = Probability that there are "n" customers in the system

$$P_n = \left[1 - \frac{\lambda}{u}\right] \left[\frac{\lambda}{u}\right]^n$$

- L = Average no. of customers in the system

$$L = \frac{\lambda}{u-\lambda}$$

- L_q = Average no. of customers in the queue

$$L_q = \left(\frac{\lambda}{u-u} \right) \left(\frac{\lambda}{u} \right)$$

- W = Average time a customer spends in the system

$$W = \frac{1}{u-\lambda}$$

- $W_q = \text{Average time a customer spends in queue}$

$$W_q = \frac{1}{\lambda - \mu} (\lambda / \mu)$$

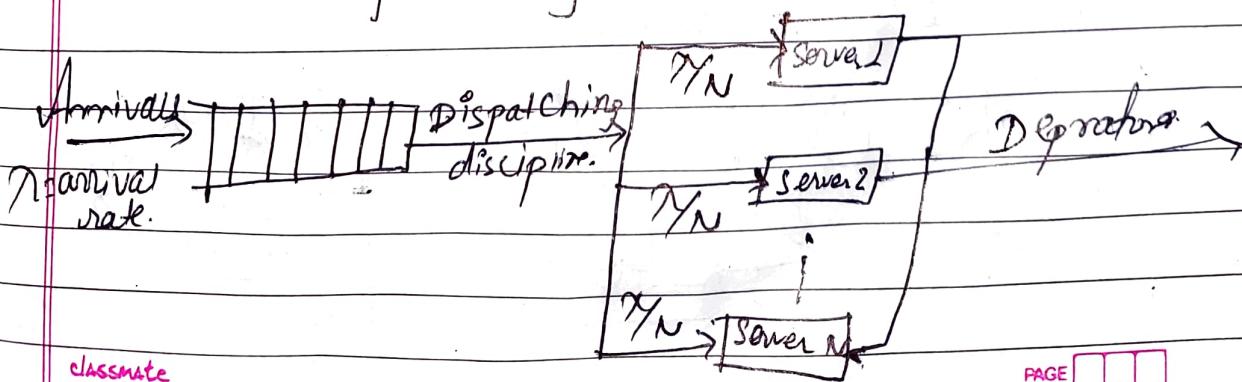
- $P_w = \text{probability that an arriving customer must wait for service}$

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

- $\rho = \text{Server utilization ratio} = \lambda / \mu$

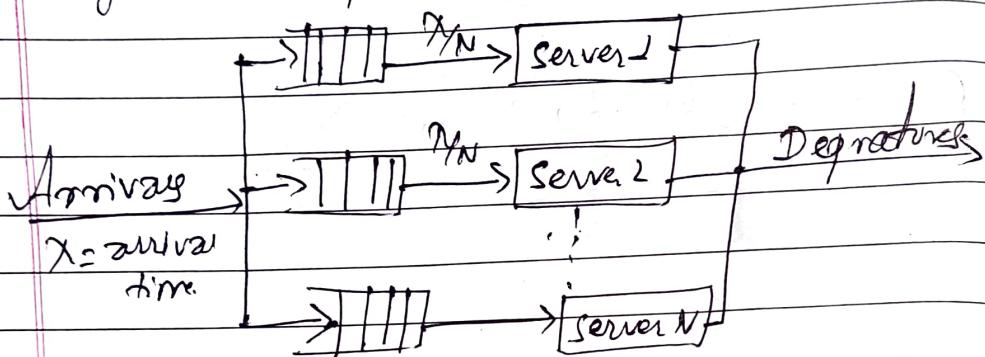
Multiserver Queuing System

- It is the queuing system with more than one server.
- In this system, all share a common queue.
- If an item arrives and at least one server is available, then the item is immediately dispatched to the server.
- It is assumed that all servers are identical, it makes no difference which server is chosen for an item.
- If all servers are busy, a queue begins to form.
- Multi-server queuing system is represented by M/M/C where arrivals form single queue & are governed by a poison process, there are C-servers & job source times are exponentially distributed.



The total Server Utilization in case of Multi-server queue for N server queue is $S = \lambda t_{\text{avg}}$ where μ is the service rate & λ is the arrival rate.

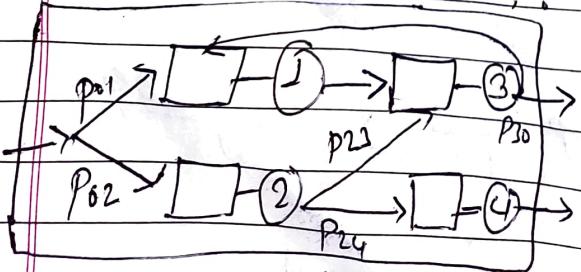
→ There is another concept which is called multiple Single server queue system



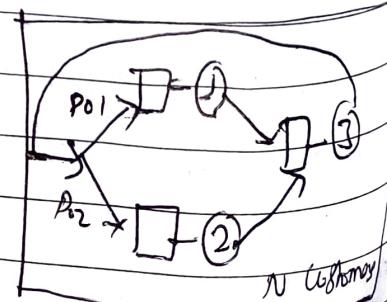
3.3.3 Networks of queuing (Networks of queues)

- Networks of queues are the systems in which a number of queues are connected by customer routing.
- When a customer is serviced at one node, it can join another node & queue for service, or leave the network.
- A queuing network is a system composed of several interconnected stations, each with a queue.

There are two networks, :



classmate Open network



Closed network

The following results assume a stable system with infinite calling population & no limit on system capacity.

- Provided that no customers are created or destroyed in the queue, the departure rate of queue P_j is same as arrival rate in the queue.
- If customers arrive to queue i^o at rate λ_i^o , & a fraction $0 \leq p_{ij} \leq 1$ of them are routed to queue j^o upon departure, then arrival rate from queue i^o to queue j^o is $\lambda_i^o p_{ij}$.
- The overall arrival rate into queue j^o , λ_j^o , is the sum of arrival rate from all sources. If customers arrive from outside the network at rate α_i , then, $\lambda_j^o = \alpha_j^o + \sum_{all i} \lambda_i^o p_{ij}$
- If queue j^o has $c_j < \infty$ parallel servers, each working at rate μ_j , then the long-run utilization of each server is $\rho_j = \frac{\lambda_j^o}{c_j \mu_j}$.

3.3.4 Applications of queuing system

- Used in designing & operating transportation systems such as airports, freeways, ports, and subways.
- It is used in determining min. no. of servers needed at a service center.
- Used in analysis of production & material handling systems.
- Used in the telephone & communication system.
- Used in analysis of telecommunications, computer networks, predicting computer performance, traffic, etc.
- Used in commercial organizations Ex: airlines, bank, ATM, gas stations etc. for serving external customers.

Numericals :

Q. Customers arrive at Mary's Shoes every 12 minutes on the average, according to a Poisson process. Service time is exponentially distributed with an average of 8 min per customer. Management is interested in determining the performance measures for this service system.

$$\Rightarrow \text{Given, } \lambda = \frac{1}{12} \text{ customers per minute} = \frac{60}{12} = 5 \text{ per hour}$$

$$\mu = \frac{1}{8} \text{ customer per minute} = \frac{60}{8} = 7.5 \text{ per hour}$$

$$\text{Now, } P_0 = 1 - \frac{\lambda}{\mu} = 1 - \frac{5}{7.5} = 0.333$$

$$L = \frac{\lambda}{\mu - \lambda} = 2$$

$$L_q = \frac{\lambda}{\mu - \lambda} \cdot P_{q0} = 1.333$$

$$W = \frac{1}{\mu - \lambda} = 0.4 \text{ hrs} = 24 \text{ min}$$

$$W_q = P_{q0} \cdot \frac{1}{\mu - \lambda} = 0.26667 \text{ hours} = 16 \text{ minutes}$$

Q. Customers arrive in a bank according to a Poisson's process with mean interarrival time of 10 minutes. Customers spend an average of 5 minutes on the single available counter, and leave. Discuss.

- i. What is the probability that a customer will not have to wait at the counter?
- ii. What is the expected no. of customers in the bank?
- iii. How much time can a customer expect to spend in the bank?

∴ Given,

$$\lambda = \frac{1}{10} \text{ customers per minute} = \frac{60}{10} \text{ customers per hr} = 6$$

$$\mu = \frac{1}{5} \text{ customers per minute} = \frac{60}{5} \text{ customers per hr} = 12$$

∴ The customers will not have to wait if there are no customers in the bank.

$$\therefore P_0 = 1 - \lambda/\mu = 1 - \frac{6}{12} = 0.5$$

$$\therefore \text{Expected number of customers} L = \lambda/\mu - \lambda = \frac{6}{12} = 1$$

$$\therefore \text{Expected time to be spent in the bank} = \frac{1}{\mu - \lambda} = 10 \text{ min}$$

Q. Customers arrive at a sales counter manned by a single person according to a Poisson process with a mean rate of 20 per hour. The time required to serve a customer has an exponential distribution with a mean of 100 seconds. Find the average waiting time of a customer for the system and queue.

$$\Rightarrow \text{Arrival rate } (\lambda) = 20 \text{ per hour}$$

$$\text{Service rate } (\mu) = \frac{1}{100} \text{ customers per second} = \frac{60 \times 60}{100} = 36 \text{ customers per hour.}$$

$$\text{The average waiting time of a customer in the queue} \\ = \frac{\lambda}{\mu(\mu-\lambda)} = \frac{20}{36(36-20)} = \frac{5}{36 \times 4} \text{ hrs} = 125 \text{ seconds}$$

The average waiting time of a customer in the system

$$N = \frac{1}{\mu - \lambda} = \frac{1}{36-20} = \frac{1}{16} \text{ hrs} = \frac{60 \times 60}{16} \text{ sec.}$$

Q.

Self service at a university cafeteria, at an average rate of 7 minutes per customer, is slower than attendant service, which has a rate of 6 minutes per student. The manager of the cafeteria wishes to calculate the average time each student spends waiting for service. Assume that customers arrive randomly at each time, at the rate of 5 per hour. Calculate the appropriate statistics for this cafeteria.

	Self Service	Attended
⇒ Arrival rate (λ)	5	5
Service rate (μ)	$6 \text{ per hr} = 6/60 = 0.1 \text{ per min}$	$10 \text{ per hr} = 10/60 = 0.1667 \text{ per min}$
Expected no. of customers in the system	$\frac{\lambda}{\mu - \lambda} = \frac{5}{0.1 - 0.0833} = 1.40$	$\frac{5}{0.1667 - 0.0833} = 1$
Expected no. of students waiting for service.	$\frac{\lambda^2}{\mu(\mu - \lambda)} = \frac{5^2}{0.1(0.1 - 0.0833)} = 0.82$	0.5
Average time in system	$\frac{1}{\mu - \lambda} = \frac{1}{0.1 - 0.0833} = 0.28 \text{ hr}$	$= 0.20 \text{ hr}$
Average time in queue.	$\frac{1}{\mu - \lambda} \times \frac{\lambda}{\mu} = \frac{1}{0.1 - 0.0833} \times \frac{5}{0.1} = 0.1667 \text{ hr}$	$= 0.10 \text{ hr}$