**Introduction**

**to**

**Queues**

**and**

**Queueing Theory**

# Model of a Queue

Arrivals → | Departures

Server(s)

Waiting
Positions

## Input Specifications

- Arrival Process Description

- Service Process Description

- Number of Servers

- Number of Waiting Positions

- Special Queueing Rules, e.g. -

  - order of service (FCFS, LCFS, SIRO, etc.)

  - baulking, reneging, jockeying for queue position
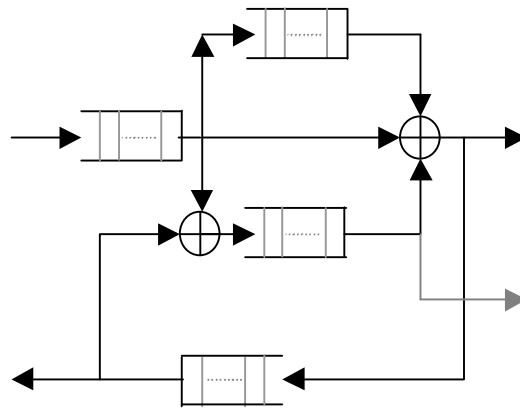
3

## Input Specifications

For networks of queues, one must provide additional information, such as -

- Interconnections between the queues

- Routing Strategy - deterministic, class based or probabilistic with given routing probabilities

- Strategy followed to handle blocking if the destination queue is one of finite capacity (i.e. with finite number of waiting positions)

4

# An Open Queueing Network

5
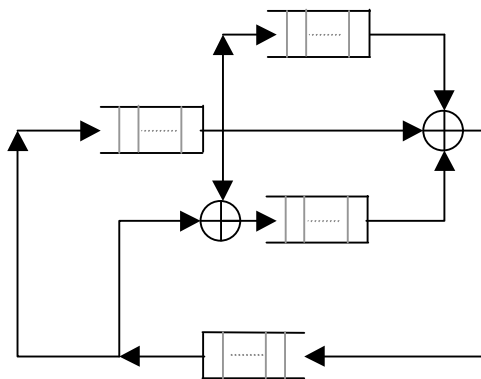
# A Closed Queueing Network

6

A Queue or a Queueing
Network may be studied in
different ways

The results may be provided
from different points of view

Analysis

and/or

Simulation

*Consider analytical
approaches here*

That of a customer entering
the system for service

That of a service provider
who provides the resources
(servers, buffers etc.)

7

---

Parameters of interest for a customer
arriving to the queue for service

(Service Parameters)

• Queueing delay

• Total delay

• Number waiting in queue

• Number in the system

• Blocking probability (*for finite capacity queues*)

• Probability that the customer has to wait for service

Transient Analysis

or

Equilibrium Analysis

Mean Results

or

Probability Distributions

8

## Parameters of interest for the Service Provider

## (Service Parameters)

- Server Utilization/Occupancy
- Buffer Utilization/Occupancy
- Total Revenue obtained
- Total Revenue lost
- Customer Satisfaction (Grade of Service)

Transient Analysis

or

Equilibrium Analysis

Mean Results

or

Probability Distributions

9

---

## Our approach to the study of queues and queueing networks

*"Subject to appropriate modelling assumptions, obtain exact analytical results for the mean performance parameters under equilibrium conditions"*

In some special cases, we can also obtain results on higher moments (variance etc.) or probability distributions and/or their transforms.

Transient analysis is not generally feasible, except for some very simple cases. For this, simulation methods are preferred.
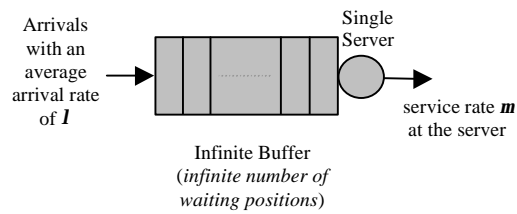
In some case, especially for queueing networks, exact analysis is not feasible but good approximate analytical methods are available.

10

# Analysis of a Simple Queue

*(with some simplifying assumptions)*

Arrivals with an average arrival rate of $l$

Single Server

service rate $m$ at the server

Infinite Buffer
(*infinite number of waiting positions*)

11

---

Assume that, as $Dt \to 0$

P{one arrival in time $Dt$} $= l\,D\,t$

P{no arrival in time $Dt$} $= 1 - l\,D\,t$

P{more than one arrival in time $Dt$} $= O((Dt)^2) = 0$

P{one departure in time $Dt$} $= m\,D\,t$

P{no departure in time $Dt$} $= 1 - m\,D\,t$

P{more than one departure in time $Dt$} $= O((Dt)^2) = 0$

P{one or more arrival and one or more departure in time $Dt$}
   $= O((Dt)^2) = 0$

Arrival Process

Mean Inter-arrival time $= \dfrac{1}{l}$

Service Process

Mean Service time $= \dfrac{1}{m}$

12

We have not really explicitly said it, but the implications of our earlier description for the arrivals and departures as $\Delta t \to 0$ is that -

• The arrival process is a Poisson process with exponentially distributed random inter-arrival times

• The service time is an exponentially distributed random variable

• The arrival process and the service process are independent of each other

13

---

The *state of the queue* is defined by defining an appropriate *system state* variable

System State at time $t = N(t) =$ Number in the system at $t$ (waiting and in service)

Let $p_N(t) =$ P{system in state $N$ at time $t$}

*Note that, given the initial system state at t=0 (which is typically assumed to be zero), if we can find $p_N(t)$ then we can actually describe probabilistically how the system will evolve with time.*

14

By ignoring terms with $(\Delta t)^2$ and higher order terms, the probability of the system state at time $t+\Delta t$ may then be found as -

$$p_0(t + \Delta t) = p_0(t)[1 - \lambda \Delta t] + p_1(t)\mu \Delta t \qquad N=0 \qquad (1.1)$$

$$p_N(t + \Delta t) = p_N(t)[1 - \lambda \Delta t - \mu \Delta t] + p_{N-1}(t)\lambda \Delta t + p_{N+1}(t)\mu \Delta t$$

$$N>0 \qquad (1.2)$$

subject to the normalisation condition that $\displaystyle\sum_{\forall i} p_i(t) = 1$ for all $t \geq 0$

Taking the limits as $\Delta t \to 0$, and subject to the same normalisation, we get

$$\frac{dp_0(t)}{dt} = -\lambda p_0(t) + \mu p_1(t) \qquad\qquad N=0 \quad (1.3)$$

$$\frac{dp_N(t)}{dt} = -(\lambda + \mu) p_N(t) + \lambda p_{N-1}(t) + \mu p_{N+1}(t) \qquad N>0 \quad (1.4)$$

These equations may be solved with the proper initial conditions to get the *Transient Solution.*

If the queue starts with $N$ in the system, then the corresponding initial condition will be

$p_i(0)=0$          for $i \neq N$

$p_N(0)=1$          for $i=N$

For the *equilibrium solution,* the conditions invoked are -

$$\frac{dp_i(t)}{dt} = 0$$

and

$$p_i(t) = p_i \qquad \text{for } i=0, 1, 2.....\boldsymbol{\mu}$$

---

For this, defining $\boldsymbol{r} = \boldsymbol{l/m}$ *erlangs*, with $\boldsymbol{r}<1$ for stability, we get

$$p_1 = \boldsymbol{r}p_0$$
$$p_{N+1} = (1+\boldsymbol{r})p_N - \boldsymbol{r}p_{N-1} = \boldsymbol{r}p_N = \boldsymbol{r}^{N+1}p_0 \qquad\qquad N \geq 1 \quad\Bigg\} \ (1.5)$$

Applying the Normalization Condition $\displaystyle\sum_{i=0}^{\infty} p_i = 1$ we get

$$p_i = \boldsymbol{r}^i (1 - \boldsymbol{r}) \qquad\qquad i = 0, 1, ......,\boldsymbol{\mu} \qquad\qquad (1.6)$$

as the equilibrium solution for the state distribution when the arrival and service rates are such that $\boldsymbol{r} = \boldsymbol{l/m} < 1$

*Note that the equilibrium solution does not depend on the initial condition but requires that the average arrival rate must be less than the average service rate*

Mean Performance Parameters of the Queue

(a) Mean Number in System, $N$

$$N = \sum_{i=0}^{\infty} i p_i = \sum_{i=0}^{\infty} i r^i (1-r) = \frac{r}{1-r} \qquad (1.7)$$

(b) Mean Number Waiting in Queue, $N_q$

$$N_q = \sum_{i=1}^{\infty} (i-1) p_i = \frac{r}{1-r} - (1-p_0) = \frac{r}{1-r} - r = \frac{r^2}{1-r} \qquad (1.8)$$

19

---

Mean Performance Parameters of the Queue

(c) Mean Time Spent in System $W$

    This would require the following additional
    assumptions

• FCFS system though the mean results will hold for any queue where the server does not idle while there are customers in the system

• The equilibrium state probability $p_k$ will also be the same as the probability distribution for the number in the system as seen by an arriving customer

• The mean residual service time for the customer currently in service when an arrival occurs will still be $1/m$   Memory-less Property satisfied only by the exponential distribution

20

Mean Performance Parameters of the Queue *(continued)*

Using these assumptions, we can write

$$W = \sum_{k=0}^{\infty} \frac{(k+1)}{m} p_k = \frac{1}{m(1-r)} \qquad (1.9)$$

(d) Mean Time Spent Waiting in Queue $W_q$

This will obviously be one mean service time less than $W$

$$W_q = W - \frac{1}{m} = \frac{r}{m(1-r)} \qquad (1.10)$$

21

Mean Performance Parameters of the Queue *(continued)*

Alternatively, $W_q$ may be obtained using the same kind of arguments as those used to obtain *W* earlier. This will give

$$W_q = \sum_{k=0}^{\infty} \frac{k}{m} p_k = \frac{r}{m(1-r)}$$

which is the same result as obtained earlier.

(e) P{Arriving customer has to wait for service} = $1 - p_0 = r$

22

11

Mean Performance Parameters of the Queue *(continued)*

(f) Server Utilization    *"Fraction of time the server is busy"*

= P{server is not idle}

$= 1 - p_0 = \rho$

The queue we have analyzed is the single server M/M/1/∞ queue with Poisson arrivals, exponentially distributed service times and infinite number of buffer positions

23

---

The analytical approach given here may actually be applied for simple queueing situations where -

• The arrival process is Poisson, i.e. the inter-arrival times are exponentially distributed

• The service times are exponentially distributed

• The arrival process and the service process are independent of each other

24

Some other simple queues which may be similarly analyzed, under the same assumptions -

• Queue with Finite Capacity

• Queue with Multiple Servers

• Queue with Variable Arrival Rates

• Queue with "Balking"

25