# Introduction to Queueing Theory

Queueing theory is a set of mathematical tools for the analysis of probabilistic systems of customers and servers. Among the oldest of Management Science tools, queueing theory can be traced to the work of A. K. Erlang, a Danish mathematician who studied telephone traffic congestion in the first decade of the 20th century.

#### **Definitions:**

<u>Arrival Rate</u> refers to the average number of customers who require service within a specific period of time.

A <u>**Capacitated Queue</u>** is limited as to the number of customers who are allowed to wait in line.</u>

<u>Customers</u> can be people, work-in-process inventory, raw materials, incoming digital messages, or any other entities that can be modeled as lining up to wait for some process to take place.

A **<u>Queue</u>** is a set of customers waiting for service.

**Queue Discipline** refers to the priority system by which the next customer to receive service is selected from a set of waiting customers. One common queue discipline is first-in-first-out, or FIFO.

A <u>Server</u> can be a human worker, a machine, or any other entity that can be modeled as executing some process for waiting customers.

<u>Service Rate</u> (or Service Capacity) refers to the overall average number of customers a system can handle in a given time period.

<u>Stochastic Processes</u> are systems of events in which the times between events are random variables. In queueing models, the patterns of customer arrivals and service are modeled as stochastic processes based on probability distributions.

<u>Utilization</u> refers to the proportion of time that a server (or system of servers) is busy handling customers.

## Queueing Notation:

In the literature, queueing models are described by a series of symbols and slashes, such as A/B/X/Y/Z, where A indicates the arrival pattern, B indicates the service pattern, X indicates the number of parallel servers, Y indicates the

queue's capacity, and Z indicates the queue discipline. We will be concerned primarily with the M/M/1 queue, in which the letter M indicates that times between arrivals and times between services both can be modeled as being exponentially distributed. The number 1 indicates that there is one server; we will also study some M/M/n queues, where n is some number greater than 1.

### Symbols:

Performance Measure	Random Variable	Expected Value
Number of Customers in the System	Ν	L
Number of Customers in the Queue	Nq	$L_q$
Number of Customers in Service	Ns	Ls
Time Spent in the System	Т	W
Time Spent in the Queue	Tq	$W_q$
Time Spent in Service	Ts	Ws

System Parameters				
Number of Servers	S			
Arrival Rate (number per unit of time)	λ	(Greek letter lambda)		
Service Rate (number per unit of time)	μ	(Greek letter mu)		
Utilization Factor	ρ	(Greek letter rho)		

#### Formulas:

The utilization factor, or  $\rho$  (rho), also the probability that a server will be busy at any point in time:

$$\rho = \frac{\lambda}{S\mu} \tag{i}$$

Idle time, or the proportion of time servers are not busy, or the probability that a server will be idle at any given time:

$$1-\rho$$
 (ii)

The average time a customer spends in the system:

$$W = W_q + W_s = \frac{1}{\mu - \lambda}$$
(iii)

The average number of customers in the queue:

$$L_q = \frac{\lambda^2}{\mu(\mu - \lambda)} = \frac{\rho^2}{1 - \rho}$$
(iv)

The average number of customers in the system:

$$L = L_q + L_s \tag{v}$$

The average time spent waiting in the queue:

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

Those who would commune more directly with the Arab mathematician al Jebr (whose name is the source for our word algebra) will note that there is another formula for N:

$$N_q + N_s = \frac{\rho^2}{1 - \rho} + \rho = \frac{\frac{\lambda}{\mu} \left(\frac{\lambda}{\mu}\right)}{1 - \frac{\lambda}{\mu}} + \frac{\lambda}{\mu} = \frac{\left\{\mu \left[\frac{\lambda}{\mu} \left(\frac{\lambda}{\mu}\right)\right]\right\} + \left\{\left(1 - \frac{\lambda}{\mu}\right)(\lambda)\right\}}{\mu \left(1 - \frac{\lambda}{\mu}\right)} = \frac{\frac{\lambda}{\mu}}{1 - \frac{\lambda}{\mu}} = \frac{\rho}{1 - \rho} = N$$

Therefore:

$$L = \frac{\rho}{1 - \rho} \tag{vii}$$

The single most important formula in queueing theory is called Little's Law:

$$L = \lambda W$$
 (viii)

Little's Law can be restated in several ways; for example, here is a formula for W:

$$W = \frac{L}{\lambda}$$
(ix)

There are also some formulas that are based on the special characteristics of the exponential distribution. These generally are beyond the scope of this course, but one useful example is offered here. In an M/M/n system, the probability of a customer spending longer than some specific time period t in the system is given by the following formula. e is the base of the natural logarithms (roughly 2.718):

$$P\{T \ge t\} = e^{(\lambda - \mu)t} \tag{x}$$

#### What Every Manager Should Know about Queueing Theory

1. Little's Law describes the fundamental trade-off between customer service and server utilization. To reduce the amount of time a customer spends waiting in your system, you must increase the amount of time your server is idle. Consider the following graph, which uses Little's Law to illustrate the relationship between rho and the expected number of people in the system:



2. Great economies of scale with respect to waiting times can be achieved by joining multiple queues into a single queue with multiple servers. For example, supermarket customers spend less time waiting in a single queue served by 4 cashiers than in one of four separate queues, each served by one cashier.

This seeming paradox stems from the fact that, with four independent queues, it is possible for one cashier to be idle while the other three have long lines of people waiting. Although supermarkets are set up with separate queues for each cashier, the customers intuitively understand and make adjustments; they won't stand in line for one cashier when another cashier is idle.

3. Lambda (the arrival rate of customers) is rarely stationary. For example, customers arrive at a hospital emergency room more frequently from 12:00 to 1:00 AM than from 9:00 to 10:00 AM. A common error is to use a long-run average to calculate lambda, without accounting for changes in the arrival rate over time. Depending on the nature of the operation being planned, it may be best to use a "worst-case" estimate of lambda, based on the arrival rate during peak hours.

4. A simple rule of thumb for deciding how many servers to provide is given by Peter Kolesar. Let us suppose that you want only a 5% chance that any given customer will have to wait in line; call this probability of waiting  $\alpha$  (alpha). Using the normal table, we can estimate the number of servers needed to provide this level of service using the following formula:

$$s \approx \frac{\lambda}{\mu} + z_{1-\alpha} \sqrt{\frac{\lambda}{\mu}}$$

If, for example, we have a system in which we expect 75 customers per hour, and in which each server can handle about 10 customers per hour, we can achieve the 5% probability of waiting by providing 12 servers:

$$s \approx \frac{75}{10} + z_{1-.05} \sqrt{\frac{75}{10}}$$
$$= 7.5 + z_{1-.05} \sqrt{7.5}$$
$$= 7.5 + 1.64(2.74)$$
$$= 11.99$$