

Simulation and Modelling



NATIONAL UNIVERSITY
of Computer & Emerging Sciences

Spring 2023
CS4056

Muhammad Shahid Ashraf

Shahid Ashraf

CS4056

1 / 35

Queueing Models

Queueing Models

Shahid Ashraf

CS4056

2 / 35

Queueing Models

Characteristics of Queueing Systems

Characteristics of Queueing Systems



NATIONAL UNIVERSITY
of Computer & Emerging Sciences

- The Calling Population
- System Capacity
- The Arrival Process
- Queue Behavior and Queue Discipline
- Service Times and the Service Mechanism

Shahid Ashraf

CS4056

3 / 35

Queueing Models

Characteristics of Queueing Systems

Key elements of queueing systems



NATIONAL UNIVERSITY
of Computer & Emerging Sciences

- Customer: refers to anything that arrives at a facility and requires service, e.g., people, machines, trucks, emails, packets, frames.
- Server: refers to any resource that provides the requested service, e.g., repairpersons, machines, runways at airport, host, switch, router, disk drive, algorithm.

System	Customers	Server
Reception desk	People	Receptionist
Hospital	Patients	Nurses
Airport	Airplanes	Runway
Production line	Cases	Case-packer
Road network	Cars	Traffic light
Grocery	Shoppers	Checkout station
Computer	Jobs	CPU, disk, CD
Network	Packets	Router

Shahid Ashraf

CS4056

4 / 35

Notes

Notes

Notes

Notes

Calling Population



- Calling population: the population of potential customers, may be assumed to be finite or infinite.
 - Finite population model: if arrival rate depends on the number of customers being served and waiting, e.g., model of one corporate jet, if it is being repaired, the repair arrival rate becomes zero.



- Infinite population model: if arrival rate is not affected by the number of customers being served and waiting, e.g., systems with large population of potential customers.

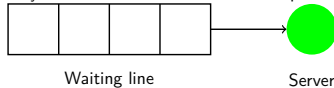


System Capacity

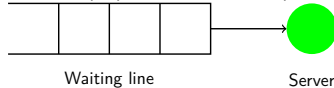


- System Capacity: a limit on the number of customers that may be in the waiting line or system.

- Limited capacity, e.g., an automatic car wash only has room for 10 cars to wait in line to enter the mechanism.
- If system is full no customers are accepted anymore



- Unlimited capacity, e.g., concert ticket sales with no limit on the number of people allowed to wait to purchase tickets.



Arrival Process



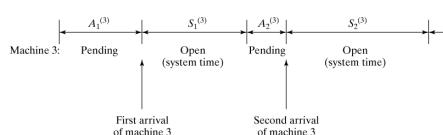
- For infinite-population models:
 - In terms of interarrival times of successive customers.
- Arrival types:
 - Random arrivals: interarrival times usually characterized by a probability distribution.
 - Most important model: Poisson arrival process (with rate λ), where a time represents the interarrival time between customer $n-1$ and customer n , and is exponentially distributed (with mean $\frac{1}{\lambda}$).
 - Scheduled arrivals: interarrival times can be constant or constant plus or minus a small random amount to represent early or late arrivals.
 - Example: patients to a physician or scheduled airline flight arrivals to an airport
- At least one customer is assumed to always be present, so the server is never idle, e.g., sufficient raw material for a machine.

Arrival Process



For finite-population models:

- Customer is pending when the customer is outside the queueing system, e.g., machine-repair problem: a machine is "pending" when it is operating, it becomes "not pending" the instant it demands service from the repairman.
- Runtime of a customer is the length of time from departure from the queueing system until that customer's next arrival to the queue, e.g., machine-repair problem, machines are customers and a runtime is time to failure (TTF).
- Let $A_1^{(i)}, A_2^{(i)}, \dots$ be the successive runtimes of customer i , and $S_1^{(i)}, S_2^{(i)}$ be the corresponding successive system times.



Notes

Notes

Notes

Notes

Queue Behavior and Queue Discipline



- Queue behavior: the actions of customers while in a queue waiting for service to begin, for example:
 - Balk: leave when they see that the line is too long
 - Reneg: leave after being in the line when its moving too slowly
 - Jockey: move from one line to a shorter line
- Queue discipline: the logical ordering of customers in a queue that determines which customer is chosen for service when a server becomes free, for example:
 - First-in-first-out (FIFO)
 - Last-in-first-out (LIFO)
 - Service in random order (SIRO)
 - Shortest processing time first (SPT)
 - Service according to priority (PR)

Shahid Ashraf

CS4056

9 / 35

Notes

Service Times and Service Mechanism



- Service times of successive arrivals are denoted by S_1, S_2, S_3 .
- May be constant or random.
- $\{S_1, S_2, S_3, \dots\}$ is usually characterized as a sequence of independent and identically distributed (IID) random variables, e.g., Exponential, Weibull, Gamma, Lognormal, and Truncated normal distribution.
- A queueing system consists of a number of service centers and interconnected queues.
- Each service center consists of some number of servers (c) working in parallel, upon getting to the head of the line, a customer takes the 1st available server.

Shahid Ashraf

CS4056

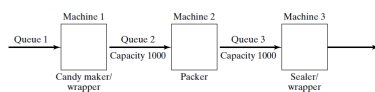
10 / 35

Notes

Examples of Queueing Systems



A candy manufacturer has a production line that consists of three machines separated by inventory in-process buffers. The first machine makes and wraps the individual pieces of candy, the second packs 50 pieces in a box, and the third machine seals and wraps the box. The two inventory buffers have capacities of 1000 boxes each. As illustrated by figure, the system is modeled as having three service centers, each center having $c = 1$ server (a machine), with queue capacity constraints between machines. It is assumed that a sufficient supply of raw material is always available at the first queue. Because of the queue capacity constraints, machine 1 shuts down whenever its inventory buffer (queue 2) fills to capacity, and machine 2 shuts down whenever its buffer empties.



Shahid Ashraf

CS4056

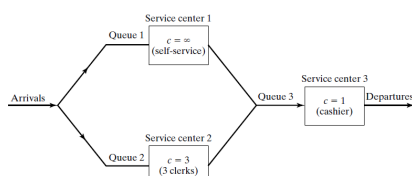
11 / 35

Notes

Examples of Queueing Systems



Consider a discount warehouse where customers may either serve themselves or wait for one of three clerks, then finally leave after paying a single cashier. The system is represented by the flow diagram in figure. The subsystem, consisting of queue 2 and service center 2. Other variations of service mechanisms include batch service (a server serving several customers simultaneously) and a customer requiring several servers simultaneously. In the discount warehouse, a clerk might pick several small orders at the same time, but it may take two of the clerks to handle one heavy item.



Shahid Ashraf

CS4056

12 / 35

Notes

A/B/c/N/K

These letters represent the following system characteristics:

- A represents the interarrival-time distribution.
- B represents the service-time distribution. Common symbols for A and B include
 - M (exponential or Markov),
 - D (constant or deterministic),
 - Ek (Erlang of order k),
 - PH (phase-type),
 - H (hyperexponential),
 - G (arbitrary or general), and
 - GI (general independent).
- c represents the number of parallel servers.
- N represents the system capacity.
- K represents the size of the calling population

Notes

A/B/c/N/K

For example,

- M/M/1/∞/∞ indicates a single-server system that has unlimited queue capacity and an infinite population of potential arrivals. The interarrival times and service times are exponentially distributed.
- When N and K are infinite, they may be dropped from the notation. For example, M/M/1/∞/∞ is often shortened to M/M/1.
- G/G/1/5/5: Single-server with capacity 5 and call-population 5.
- M/M/5/20/1500/FIFO: Five parallel server with capacity 20, call-population 1500, and service discipline FIFO

Notes

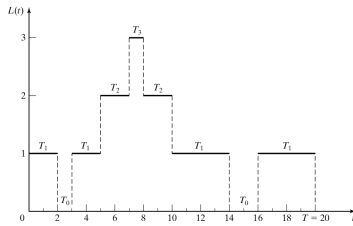
- Average waiting time: This is the average time a customer spends waiting in the queue before being served. It is an important measure of the quality of service provided by the system.
- Average queue length: This is the average number of customers waiting in the queue for service. It is an important measure of the efficiency of the system.
- Utilization: This is the proportion of time that the server is busy serving customers. A high utilization indicates that the system is being used efficiently, while a low utilization indicates that there is excess capacity.
- Throughput: This is the number of customers served by the system per unit of time. It is an important measure of the capacity of the system.
- Blocking probability: This is the probability that a customer is blocked (i.e., denied service) due to the system being at full capacity. It is an important measure of the effectiveness of the system in handling demand.

Notes

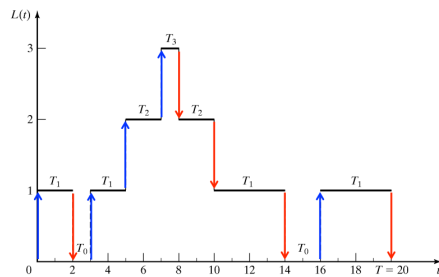
- Primary long-run measures of performance:
 - Long-run time-average number of customers in the system (L) and in the queue (L_Q)
 - Long-run average time spent in system (w) and in the queue (w_Q) per customer
 - Server utilization, or proportion of time that a server is busy (ρ)
 - "System" refers to the waiting line plus service mechanism; "queue" refers to waiting line alone
- Other measures of performance:
 - Long-run proportion of customers delayed in queue longer than t₀ time units
 - Long-run proportion of customers turned away due to capacity constraints
 - Long-run proportion of time waiting line contains more than k₀ customers
 - Defines measures of performance for general G/G/c/N/K queueing system
 - Discusses relationships and estimation using ordinary sample average or time-integrated sample average

Notes

Consider a queueing system over a period of time T , and let $L(t)$ denote the number of customers in the system at time t . A simulation of such a system is shown in Figure .

Figure: Simulation of a queueing system over a period of time T

Time-Average Number in System L



Time-Average Number in System L

- Consider a queueing system over a period of time T
 - Let T_i denote the total time during $[0, T]$ in which the system contained exactly i customers, the **time-weighted-average** number in the system is defined by:

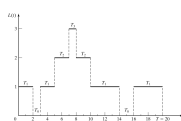
$$\hat{L} = \frac{1}{T} \sum_{i=0}^{\infty} iT_i = \sum_{i=0}^{\infty} i \left(\frac{T_i}{T} \right)$$

- Consider the total area under the function is $L(t)$, then,

$$\hat{L} = \frac{1}{T} \sum_{i=0}^{\infty} iT_i = \frac{1}{T} \int_0^T L(t) dt$$

- The long-run time-average number of customers in system, with probability 1:

$$\hat{L} = \frac{1}{T} \int_0^T L(t) dt \xrightarrow{T \rightarrow \infty} L$$

Figure: Simulation of a queueing system over a period of time T

Let T_i denote the total time during $[0, T]$ in which the system contained exactly i customers. In Figure, it is seen that $T_0 = 3$, $T_1 = 12$, $T_2 = 4$, and $T_3 = 1$. In general, the time-weighted-average number in a system is defined by:

$$L = \sum_{i=0}^{\infty} i \frac{T_i}{T}$$

Notice that T_i/T is the proportion of time the system contains exactly i customers. The estimator L is an example of a time-weighted average.

Solution



- Let T_i denote the total time during $[0, T]$ in which the system contained exactly i customers.
- In Figure 6, $T_0 = 3$, $T_1 = 12$, $T_2 = 4$, and $T_3 = 1$.
- Using Eq. (1), the time-weighted-average number of customers in the system is:

$$\begin{aligned}
 L &= \frac{\sum_{i=0}^{\infty} iT_i}{T} \\
 &= \frac{0(3) + 1(12) + 2(4) + 3(1)}{20} \\
 &= \frac{23}{20} \\
 &= 1.15 \text{ customers.}
 \end{aligned}$$

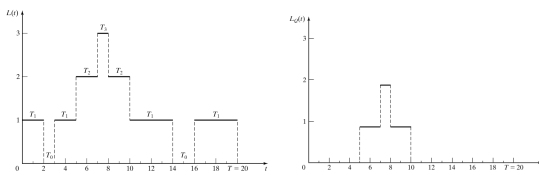
Navigation icons: back, forward, search, etc.

Shahid Ashraf

CS4056

21 / 35

Queueing System and Waiting Customers



Suppose that Figure represents a single-server queue—that is, a G/G/1/N/K queueing system ($N \geq 3$, $K \geq 3$). Then the number of customers waiting in queue is given by $L_Q(t)$, defined by

$$L_Q(t) = \begin{cases} 0, & \text{if } L(t) = 0, \\ L(t) - 1, & \text{if } L(t) \geq 1, \end{cases}$$

and shown in Figure .

Navigation icons: back, forward, search, etc.

Shahid Ashraf

CS4056

22 / 35

Queueing System and Waiting Customers



Suppose that Figure represents a single-server queue—that is, a G/G/1/N/K queueing system ($N \geq 3$, $K \geq 3$). Then the number of customers waiting in queue is given by $L_Q(t)$, defined by

$$L_Q(t) = \begin{cases} 0, & \text{if } L(t) = 0, \\ L(t) - 1, & \text{if } L(t) \geq 1, \end{cases}$$

and shown in Figure .

Thus, $T_{0,Q} = 5 + 10 = 15$, $T_{1,Q} = 2 + 2 = 4$, and $T_{2,Q} = 1$. Therefore,

$$L_Q = \frac{0(15) + 1(4) + 2(1)}{20} = 0.3 \text{ customers.}$$

Navigation icons: back, forward, search, etc.

Shahid Ashraf

CS4056

23 / 35

Average Time Spent in System Per Customer w



- The average time spent in system per customer, called the average system time, is:

$$\hat{w} = \frac{1}{N} \sum_{i=1}^N W_i$$

where W_1, W_2, \dots, W_N are the individual times that each of the N customers spend in the system during $[0, T]$.

- For stable systems: $\hat{w} \rightarrow w$ as $N \rightarrow \infty$
- If the system under consideration is the queue alone:

$$\hat{w}_Q = \frac{1}{N} \sum_{i=1}^N W_i^Q \xrightarrow{N \rightarrow \infty} w_Q$$

Navigation icons: back, forward, search, etc.

Shahid Ashraf

CS4056

24 / 35

Notes

Notes

Notes

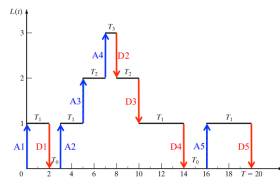
Notes

- $G/G/1/N/K$ example (cont.):

- The average system time is ($W_i = D_i - A_i$)

$$\hat{w} = \frac{W_1 + W_2 + \dots + W_5}{5} = \frac{2 + (8-3) + (10-5) + (14-7) + (20-16)}{5} = 4.6 \text{ time units}$$

- The average queuing time is $\hat{w}_Q = \frac{0 + 0 + 3 + 3 + 0}{5} = 1.2 \text{ time units}$



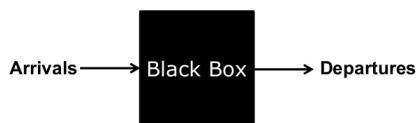
Queueing Notation

- General performance measures of queueing systems:

- P_n steady-state probability of having n customers in system
- $P_n(t)$ probability of n customers in system at time t
- λ arrival rate
- λ_e effective arrival rate
- μ service rate of one server
- ρ server utilization
- A_n interarrival time between customers $n-1$ and n
- S_n service time of the n -th arriving customer
- W_n total time spent in system by the n -th customer
- W_n^Q total time spent in the waiting line by customer n
- $L(t)$ the number of customers in system at time t
- $L_Q(t)$ the number of customers in queue at time t
- L long-run time-average number of customers in system
- L_Q long-run time-average number of customers in queue
- \hat{w} long-run average time spent in system per customer
- w_Q long-run average time spent in queue per customer

The Conservation Equation: Little's Law

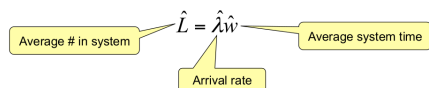
- One of the most common theorems in queueing theory
- Mean number of customers in system
- Conservation equation (a.k.a. Little's law)



average number in system = arrival rate \times average system time

The Conservation Equation: Little's Law

- Conservation equation (a.k.a. Little's law)



$$L = \lambda w \text{ as } T \rightarrow \infty \text{ and } N \rightarrow \infty$$

- Holds for almost all queueing systems or subsystems (regardless of the number of servers, the queue discipline, or other special circumstances).
- $G/G/1/N/K$ example (cont.): On average, one arrival every 4 time units and each arrival spends 4.6 time units in the system. Hence, at an arbitrary point in time, there are $(1/4)(4.6) = 1.15$ customers present on average.

Server Utilization



- Definition: the proportion of time that a server is busy.
 - Observed server utilization, $\hat{\rho}$, is defined over a specified time interval $[0, T]$.
 - Long-run server utilization is ρ .
 - For systems with long-run stability: $\hat{\rho} \rightarrow \rho$ as $T \rightarrow \infty$
- For $G/G/1/\infty/\infty$ queues:
 - Any single-server queueing system with
 - average arrival rate λ customers per time unit,
 - average service time $E(S) = 1/\mu$ time units, and
 - infinite queue capacity and calling population.
 - Conservation equation, $L = \lambda W$, can be applied.
 - For a stable system, the average arrival rate to the server, λ_s , must be identical to λ .
 - The average number of customers in the server is:

$$\hat{L}_s = \frac{1}{T} \int_0^T (L(t) - L_Q(t)) dt = \frac{T - T_0}{T}$$

Navigation icons: back, forward, search, etc.

Shahid Ashraf

CS4056

29 / 35

Notes

Server Utilization



- In general, for a single-server queue:

$$\hat{L}_s = \hat{\rho} \xrightarrow{T \rightarrow \infty} L_s = \rho$$

$$\text{and } \rho = \lambda \cdot E(S) = \frac{\lambda}{\mu}$$

- For a single-server stable queue: $\rho = \frac{\lambda}{\mu} < 1$
- For an unstable queue ($\lambda > \mu$), long-run server utilization is 1.

Navigation icons: back, forward, search, etc.

Shahid Ashraf

CS4056

30 / 35

Notes

Server Utilization



- For $G/G/c/\infty/\infty$ queues:
 - A system with c identical servers in parallel.
 - If an arriving customer finds more than one server idle, the customer chooses a server without favoring any particular server.
 - For systems in **statistical equilibrium**, the average number of busy servers, L_s , is:

$$L_s = \lambda E(S) = \frac{\lambda}{\mu}$$
 - Clearly $0 \leq L_s \leq c$
 - The long-run average server utilization is:

$$\rho = \frac{L_s}{c} = \frac{\lambda}{c\mu}, \text{ where } \lambda < c\mu \text{ for stable systems}$$

Navigation icons: back, forward, search, etc.

Shahid Ashraf

CS4056

31 / 35

Notes

Server Utilization and System Performance

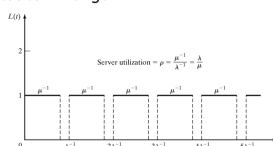


- System performance varies widely for a given utilization ρ .
 - For example, a $D/D/1$ queue where $E(A) = 1/\lambda$ and $E(S) = 1/\mu$, where:

$$L = \rho = \lambda/\mu, \quad w = E(S) = 1/\mu, \quad L_Q = W_Q = 0$$

- By varying λ and μ , server utilization can assume any value between 0 and 1.

- In general, variability of interarrival and service times causes lines to fluctuate in length.



Navigation icons: back, forward, search, etc.

Shahid Ashraf

CS4056

32 / 35

Notes

Server Utilization and System Performance

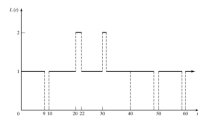
- Example: A physician who schedules patients every 10 minutes and spends S_i minutes with the i -th patient:

$$S_i = \begin{cases} 9 \text{ minutes with probability } 0.9 \\ 12 \text{ minutes with probability } 0.1 \end{cases}$$
- Arrivals are deterministic:

$$A_1 = A_2 = \dots = \lambda^{-1} = 10$$
- Services are stochastic
 - $E(S_i) = 9.3 \text{ min}$
 - $V(S_0) = 0.81 \text{ min}^2$
 - $\sigma = 0.9 \text{ min}$
- On average, the physician's utilization is

$$\rho = \lambda \mu = 0.93 < 1$$

- Consider the system is simulated with service times: $S_1 = 9, S_2 = 12, S_3 = 9, S_4 = 9, S_5 = 9, \dots$
- The system becomes:
- The occurrence of a relatively long service time ($S_2 = 12$) causes a waiting line to form temporarily.



Steady-State Behavior of Markovian Models

- Markovian models:
 - Exponential-distributed arrival process (mean arrival rate = $1/\lambda$).
 - Service times may be exponentially (M) or arbitrary (G) distributed.
 - Queue discipline is FIFO.
 - A queueing system is in **statistical equilibrium** if the probability that the system is in a given state is **not time dependent**:

$$P(L(t) = n) = P_n(t) = P_n$$

- Mathematical models in this chapter can be used to obtain approximate results even when the model assumptions do not strictly hold, as a rough guide.
- Simulation can be used for more refined analysis, more faithful representation for complex systems.

Steady-State Behavior of Markovian Models

- Properties of processes with statistical equilibrium
 - The state of statistical equilibrium is reached from any starting state.
 - The process remains in statistical equilibrium once it has reached it.

[illegible]

Q4 A service station can hold only three customers: one in service, and two waiting. Additional customers are turned away when the system is full. The offered load is a , namely $a = 2/3, N = \text{System Capacity}$, $\alpha = \lambda/\mu$, $\rho = \frac{a}{N}$

- Space only for 2 customers: one in service and two waiting
- First compute P_0

$$P_0 = \left[1 + \frac{c}{n+1} \frac{\rho^n}{c!} + \frac{\rho^n}{c!} \sum_{m=n+1}^{\infty} \rho^{m-n} \right]^{-1}$$
- P(system is full)
$$P_s = P_n = \frac{1}{n!} P_0 = P_0 \frac{\rho^n}{c! n!} P_0$$
- Average of the queue
$$L_q = \frac{\rho P_0}{c! - \rho P_0} = 1 - \rho^{c-1} - (N - c) \rho^{c-1} (1 - \rho)$$
- Effective arrival rate
$$\lambda_e = \lambda (1 - P_s)$$

- Queue time

$$w_q = \frac{E_q}{\lambda_s} =$$
- System time, time in shop

$$w = w_q + \frac{1}{\mu} =$$
- Expected number of customers in shop

$$L = \lambda_s w =$$
- Probability of busy shop

$$1 - p_0 = \frac{\lambda_s}{\mu} =$$

Notes

[illegible]

Notes

[illegible]

Notes

[illegible]

Notes

[illegible]