

by Ehssan Ghashim and Patrick Boily

Queueing theory is a branch of mathematics that studies and models the act of waiting in lines, or queues. As a topic in operational research, it combines elements of a variety of quantitative disciplines, but it is not often part of the data analyst’s toolbox. In this chapter, we introduce the terminology and basic framework of queueing models (including Kendall-Lee notation, birth-death processes, and Little’s formula), as well as the most commonly-used queueing system: $M/M/c$.

24.1 Background

Queueing theory is a branch of mathematics that studies and models the act of waiting in lines. The seminal paper on queueing theory [3] was published in 1909 by Danish mathematician A.K. Erlang; in it, he studied

the problem of determining how many telephone circuits were necessary to provide phone service that would prevent customers from waiting too long for an available circuit. In developing a solution to this problem, he began to realize that the problem of minimizing waiting time was applicable to many fields, and began developing the theory further. Erlang’s **switchboard problem** laid the path for modern queueing theory [1].

Queueing theory boils down to answering simple questions:

- How likely is it that objects/units/persons will queue up and wait in line?
- How long will the line be?
- How long will the wait be?
- How busy will the system be?
- How much capacity is needed to meet an expected level of demand?

Knowing how to think about these kinds of questions will help analysts and stakeholder anticipate **bottlenecks**. As a result, they will build systems and teams to be more efficient and more scalable, to have higher performance and lower costs, and to ultimately provide better service to their customers and end users.

Queueing theory also allows for the quantitative treatment of bottlenecks and effect on performance. For instance, a question such as “how long will the wait be, on average?” will have an answer, but so will other

24.1 Background	1573
24.2 Terminology	1575
Input/Arrival Processes	1579
Output/Service Processes	1579
Queue Discipline	1581
Joining a Queue	1581
24.3 Theoretical Framework	1581
Kendall-Lee Notation	1582
Birth-Death Processes	1583
Little’s Queueing Formula	1584
24.4 $M/M/1$ Queueing Systems	1585
Basics	1585
Limited Capacity	1587
24.5 $M/M/c$ Queueing Systems	1589
24.6 Exercises	1592
Chapter References	1592

questions concerning the variability of wait times, the distribution of wait times, and the likelihood that a customer will receive extremely poor service, and so on [5].

Let us consider a simple example. Suppose a grocery store has a single checkout line and a single cashier. If, on average, one shopper arrives at the line to pay for their groceries every 5 minutes and if scanning, bagging, and paying takes 4.5 minutes, on average, would we expect customers to have to wait in line?

When the problem is presented this way, our intuition says that there should be no waiting in line, and that the cashier should be idle, on average, 30 seconds every 5 minutes, only being busy 90% of the time. No one ever has to wait before being served!

If you have ever been in a grocery store, however, you know that this is not what happens in reality; many shoppers will wait in line, and they will have to wait a long time before being processed.

Fundamentally, **queueing** happens for three reasons:

- **irregular arrivals** – shoppers do not arrive at the checkout line on a regular schedule; they are sometimes spaced far apart and sometimes close together, so they **overlap** (an overlap automatically causes queueing and waiting);
- **irregular job sizes** – shoppers do not all get processed in 4.5 minutes; someone shopping for a large family will require much more time than someone shopping only for themselves, for instance (when this happens, overlap is again a problem because new shoppers will arrive and be ready to check out while the existing ones are still in progress), and
- **waste** – lost time can never be regained; shoppers overlap because the second shopper arrived too soon, before the first had the time to finish being served, but looking at it the other way, perhaps it's not the second shopper's fault; perhaps the first shopper should have arrived earlier, but they wasted time reading a magazine while the cashier was idle! They missed their chance for quick service and, as a result, made the second shopper have to wait.

Irregular arrival times and job sizes are guaranteed to cause queueing. The only time there is no queueing is when the job sizes are uniform, the arrivals are timed evenly, and there is little enough work for the cashier to keep up with the arrival. Even when the cashier is barely busy, irregular arrivals or arrivals **in bursts** will cause some queueing.

In general, queueing gets worse when the following hold:

- **high utilisation** – the busier the cashier is, the longer it takes to recover from wasted time;
- **high variability** – the more variability in arrivals or job sizes, the more waste and the more overlap (queueing) occurs, and
- **insufficient number of servers** – fewer cashiers means less capacity to absorb arrival spikes, leading to more wasted time and higher utilisation.

In order to describe queues, we must first know and understand some useful probability distributions, as well as input and output processes.

24.2 Terminology

Queueing theory studies processes in terms of three key concepts:

- **customers** are the units of work that the system serves – a customer can be a real person, or it can be whatever the system is supposed to process and complete: a web request, a database query, a part to be milled by a machine, etc.;
- **servers** are the objects that do the processing work – a server might be the cashier at the grocery store, a web server, a database server, a milling machine, etc., and
- **queues** are where the units of work wait if the server is busy and can not start the work as they arrive – a queue may be a physical line, reside in memory, etc.

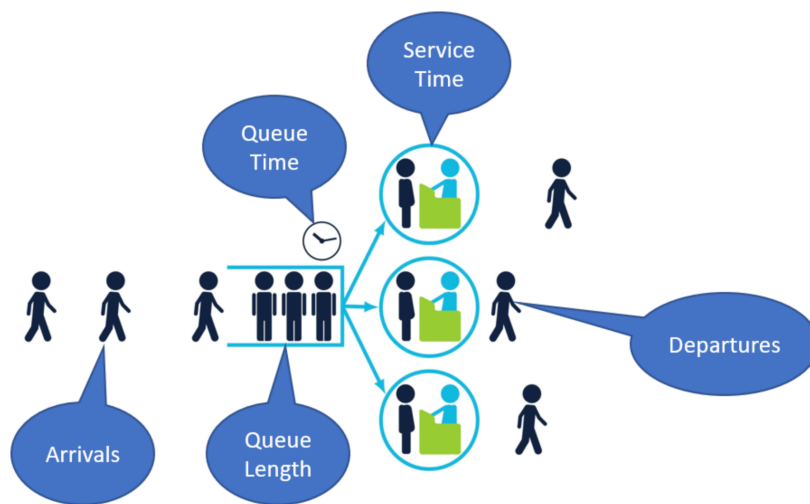


Figure 24.1: Components of a generic queueing system, by D. Hare [↗](#).

Useful Distributions Three distributions play a central role in queueing theory: Poisson, exponential, and Erlang distributions.

Poisson Distribution The **Poisson distribution** counts the number of discrete events occurring in a fixed time period; it is closely connected to the exponential distribution, which measures the time between arrivals of the events. The Poisson distribution is a discrete distribution; the random variable can only take non-negative integer values. The exponential distribution can take any (nonnegative) real value.

Consider the problem of determining the probability of n arrivals being observed during a time interval of length t , where the following assumptions are made:

- the probability that an arrival is observed during a small time interval (say of length ν) is proportional to the length of interval; let the proportionality constant be λ , so that the probability is $\lambda\nu$;
- the probability of two or more arrivals in a small interval is zero;
- the number of arrivals in any time interval is independent of the number in non-overlapping time interval – for example, the number of arrivals occurring between times 5 and 25 does not provide

information about the number of arrivals occurring between times 30 and 50.

Let $P(n; t)$ be the probability of observing n arrivals in a time interval of length t . Then, for some $\lambda > 0$,

$$P_\lambda(n; t) = \frac{(\lambda t)^n}{n!} e^{-\lambda t}, \quad n = 0, 1, 2, \dots$$

is the probability mass function of the Poisson distribution for the discrete random variable n – the number of arrivals – for a given length of time interval t (see Figure 24.2).

Example: on average, 50 customers arrive in a coffee shop every hour. What is the probability that exactly 20 customers will arrive in a 30-minute period, if the arrivals follow a Poisson distribution?

Solution: given $\lambda = 50$ customers per hour, $t = 30 \text{ min} = 0.5 \text{ hr}$ and $n = 20$, we have

$$P_{50}(20; 0.5) = \frac{(50 \cdot 0.5)^{20}}{20!} e^{-50 \cdot 0.5} \approx 5.2\%.$$

We can evaluate the probability directly in R *via*

```
n=20
lambda=50
t=0.5
dpois(n, lambda*t)
```

```
[1] 0.05191747
```

In a queueing system, such arrivals are referred to as **Poisson arrivals**. The time between successive arrivals is called the **inter-arrival time**.

Exponential Distribution If the number of arrivals in a given time interval follows a Poisson distribution with parameter λt , the inter-arrival times follow an **exponential distribution** with probability density function

$$f_\lambda(t) = \lambda e^{-\lambda t}, \quad \text{for } t > 0,$$

and the probability $P(W \leq t)$ that a customer's waiting time W is smaller than the length of the time interval t is

$$P(W \leq t) = 1 - e^{-\lambda t}$$

(see Figure 24.2). We would write $W \sim \text{Exp}(\lambda)$.

Example: a manager of a fast food restaurant observes that an average of 9 customers are served by a waiter in a one-hour time period. Assuming that the service time follows an exponential distribution, what is the probability that a customer will be served within 15 minutes?

Solution: let w be the average waiting time. Given $\mu = 9$ customers per hours, $t = 15 \text{ min} = 0.25 \text{ hr}$, we have

$$P(w \leq 15 \text{ min}) = 1 - e^{-9 \times 0.25} \approx 89.5\%.$$

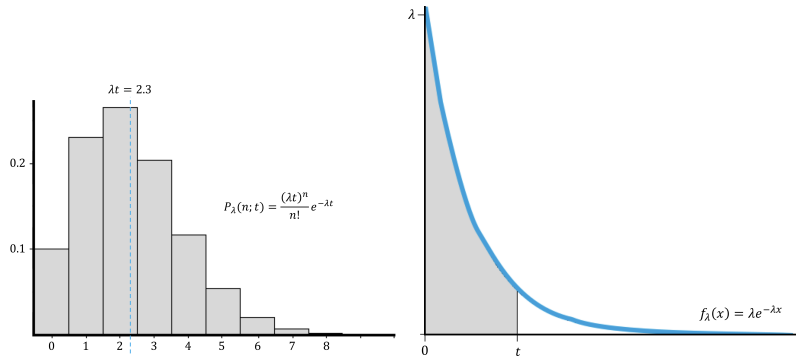


Figure 24.2: Poisson (with $\lambda t = 2.3$) and exponential distributions (with parameter λ). The shaded area (right) represents the probability that a customer will wait up to the length of the time interval t .

We can evaluate the probability directly in R.

```
t=0.25
mu=9
pexp(t, rate=mu)
```

```
[1] 0.8946008
```

In general, if the arrival rate is **stationary**, if **bulk** arrivals (two or more simultaneous arrivals) cannot occur, and if past arrivals do not affect future arrivals, then inter-arrival times follow an exponential distribution with parameter λ , and the number of arrivals in any interval of length t is Poisson with parameter λt .

One of the most attractive features of the exponential distribution relating to inter-arrival times is that it is **memoryless** – if X follows an exponential distribution, then for all non-negative values of t, h ,

$$P(X \geq t + h | X \geq t) = P(X \geq h).$$

No other density function satisfies this property [8].

The memoryless property of the exponential distribution is important because it implies that the probability distribution of the time until the next arrival is independent of the time since the last arrival. This is clearly not always the case – imagine if that was so when waiting for public transportation!

For instance, if we know that at least t time units have elapsed since the last arrival, then the distribution of the time h until the next arrival is independent of t . If $h = 4$, say, then we must have

$$P(X > 9 | X > 5) = P(X > 7 | X > 3) = P(X > 4).$$

Example: The time W a customer spends waiting in a bank queue is exponentially distributed with mean $\lambda = 10$ min, say. If they've already waited 10 minutes, what is the probability that they will have had to wait more than 15 minutes in total, when all is said and done?

Solution: thanks to the memory-less property of the exponential distribution, we have

$$P(W > 15 | W > 10) = P(W > 15 - 10 = 5) = \exp(-5/\lambda) = \exp(-1/2) \approx 60.6\%.$$

We can evaluate the probability directly in R.

```
w=5
lambda=10
1-pexp(w, rate=1/lambda)
```

```
[1] 0.6065307
```

Erlang Distribution The exponential distribution is not always an appropriate model of inter-arrival times, however (perhaps the process should not be memoryless, say).

A common alternative is to use the **Erlang** distribution $\mathcal{E}(R, k)$, a continuous random variable with **rate** and **shape** parameters $R > 0$ and $k \in \mathbb{Z}^+$, respectively, whose probability density function is

$$f_{R,k}(t) = \frac{R(Rt)^{k-1}e^{-Rt}}{(k-1)!}, \quad t \geq 0.$$

If $k = 1$, the Erlang distribution reduces to an exponential distribution with parameter R . It can further be shown that if $X \sim \mathcal{E}(R, k)$, where $R = k\lambda$, then $X \sim X_1 + X_2 + \dots + X_k$, where each $X_i \sim \text{Exp}(R)$ is an independent random variable.

When we model the inter-arrival process as an Erlang distribution $\mathcal{E}(k\lambda, k)$, we are really saying that it is equivalent to customers going through k **phases** (each of which is memoryless) before being served.

For this reason, the shape parameter is often referred to as the number of phases of the Erlang distribution [7].

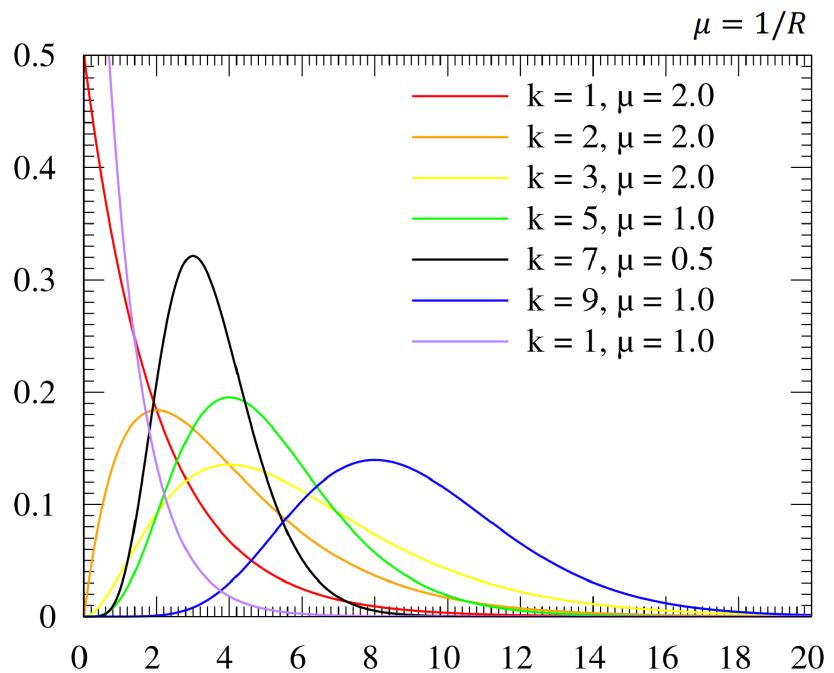


Figure 24.3: Probability distribution functions for various Erlang random variable [Wikipedia].

24.2.1 Input/Arrival Processes

The input process is usually called the **arrival process**. Arrivals are called **customers**. In the models under consideration, we assume that arrivals cannot be simultaneous (this might be unrealistic when modeling arrivals at a restaurant, say). If simultaneous arrivals are possible (in theory and/or in practice), we say that **bulk arrivals are allowed**.

Usually, we assume that the arrival process is **unaffected by the number of customers** in the system. In the context of a bank, this would imply that whether there are 500 or 5 people at the bank, the process governing arrivals remains unchanged. There are two common situations in which the arrival process may depend on the number of customers present. The first occurs when arrivals are drawn from a small population – the so-called **finite source models** – if all members of the population are already in the system, there cannot be another arrival!

Another such situation arises when the rate at which customers arrive at the facility decreases when the facility becomes too crowded. For example, when customers see that a restaurant's parking lot is full, they might very well decide to go to another restaurant or forego eating out altogether. If a customer arrives but fails to enter the system, we say that the customer has **balked**.

24.2.2 Output/Service Processes

To describe the output process (often called the **service process**) of a queuing system, we usually specify a probability distribution – the **service time distribution** – which governs the customers' service time.

In most cases, we assume that the service time distribution is independent of the number of customers present in the system. This implies, for example, that the server does not work faster when more customers are present. We can distinguish two types of servers: **in parallel** and **in series**.

Servers are **in parallel** if they all provide the same type of service and a customer only needs to pass through one of them to complete their service. For example, the tellers in a bank are usually arranged in parallel; typically, customers only need to be serviced by one teller, and any teller can perform the desired service.

Servers are **in series** if a customer must pass through several servers before their service is complete. An assembly line is an example of such a queuing system. Input and output processes occur in a variety of situations:

- **situation:** purchasing Blue Jays tickets at the Rogers Centre
input: baseball fans arrive at the ticket office
output: tellers serve the baseball fans
- **situation:** pizza parlour
input: requests for pizza delivery are received
output: pizza parlour prepares and bakes pizzas, and sends them to be delivered

- **situation:** government service centre
input: citizen/residents enter the service centre
output: receptionist assigns them to a specific queue based on their needs:
 - *input:* citizen/residents enter a specific queue based on their needs
 - output:* public servant addresses their needs
- **situation:** hospital blood bank
input: pints of blood arrive
output: patients use up pints of blood
- **situation:** garage
input: cars break down and are sent to the garage for repairs
output: cars are repaired by mechanics and sent back on the streets

The relevant computations are fairly easy to execute, as the following examples demonstrate.

Example: On average, 4.6 customers enter a coffee shop each hour. If the arrivals follow a Poisson process, what is the probability that at most two customers will enter in a 30 minute period?

Solution: since 30 min = 0.5 hr, we have

$$\begin{aligned}
 P_{\lambda=4.6}(n \leq 2; t = 0.5) &= P_{4.6}(0, 0.5) + P_{4.6}(1, 0.5) + P_{4.6}(2, 0.5) \\
 &= e^{-4.6 \cdot 0.5} \left[\frac{(4.6 \cdot 0.5)^0}{0!} + \frac{(4.6 \cdot 0.5)^1}{1!} + \frac{(4.6 \cdot 0.5)^2}{2!} \right] \\
 &\approx 0.5960;
 \end{aligned}$$

the corresponding Poisson distribution is shown in Figure 24.2.

We can evaluate the probability directly in R.

```

n=2
lambda=4.6
t=0.5
ppois(n, lambda*t)

```

[1] 0.5960388

Example: in a fast food restaurant, a cashier serves on average 9 customers in a one-hour time period. If the service time follows an exponential distribution, what percentage of customers will be served in 10 minutes or less? After 30 minutes?

Solution: since 1 hr = 60 mins, we have $\mu = 9$ customers/60 minutes, and so

$$P(W \leq 10/60) = 1 - e^{-9 \cdot 10/60} \approx 0.7769$$

$$P(W > 30/60) = e^{-9 \cdot 30/60} \approx 0.0111.$$

24.2.3 Queue Discipline

To describe a queuing system completely, we must also describe the **queue discipline** and the manner in which customers **join lines**. The queue discipline describes the method used to determine the order in which customers are served:

- the most common queue discipline is the **first come, first served** (FCFS) discipline, in which customers are served in the order of their arrival, as one would expect to see in an Ottawa coffee shop;
- under the **last come, first served** (LCFS) discipline, the most recent arrivals are the first to enter service; for example, if we consider exiting from an elevator to be the service, then a crowded elevator illustrates such a discipline;
- sometimes the order in which customers arrive has no effect on the order in which they are served; this would be the case if the next customer to enter service is randomly chosen from those customers waiting for service, a situation referred to as **service in random order** (SIRO) discipline; when callers to an inter-city bus company are put on hold, the luck of the draw often determines which caller will next be serviced by an operator;
- finally, **priority** discipline classifies each arrival into one of several categories, each of which is assigned a priority level (a **triage** process); within each priority level, customers enter the queue on a FCFS basis; such a discipline is often used in emergency rooms to determine the order in which customers receive treatment, and in copying and computer time-sharing facilities, where priority is usually given to jobs with shorter processing times.

24.2.4 Method Used by Arrivals to Join Queue

Another important factor for the behaviour of the queuing system is the **method** used by customers to determine which line to join. For example, in some banks, customers must join a single line, but in other banks, customers may choose the line they want to join.

When there are several lines, customers often join the shortest line. Unfortunately, in many situations (such as at the supermarket), it is difficult to define the shortest line. If there are several lines at a queuing facility, it is important to know whether or not customers are allowed to **switch**, or jockey, between lines. In most queuing systems with multiple lines, jockeying is permitted, but jockeying at a custom inspection booth would not be recommended (if it is even allowed), for instance.

24.3 Queueing Theory Framework

There is a standard notation that is used to describe large families of queueing systems: the **Kendall-Lee notation** [4].

24.3.1 Kendall-Lee Notation

Queueing systems can be described *via* six characteristics:

$$x_1/x_2/x_3/x_4/x_5/x_6.$$

The 1st characteristic x_1 specifies the nature of the **arrival process**. The following standard abbreviations are used:

M	inter-arrival times are independent identically distributed (iid) exponentials
D	inter-arrival times are iid and deterministic
E_k	inter-arrival times are iid Erlangs with shape parameter k
G	inter-arrival times are iid and governed by some general distribution

The 2nd characteristic x_2 specifies the nature of the **service times**:

M	service times are iid and exponential
D	service times are iid and deterministic
E_k	service times are iid Erlang with shape parameter k
G	service times are iid and follow some general distribution

The 3rd characteristic x_3 represents the **number of parallel servers**.

The 4th characteristic x_4 describes the **queue discipline**:

FCFS	first come, first served
LCFS	last come, first served
SIRO	service in random order
GD	general queue discipline

The 5th characteristic x_5 specifies the **maximum allowable number of customers in the system**.¹

1: Including customers who are waiting and customers who are in service.

The 6th characteristic x_6 gives the **size of the population** from which customers are drawn. Unless the number of potential customers is of the same order of magnitude as the number of servers, the population size is considered to be infinite.

2: When that is the case, the string is often omitted.

In many important models $x_4/x_5/x_6$ is $GD/\infty/\infty$.² As an example, $M/M/3/FCFS/20/\infty$ could represent a bank with 3 tellers, exponential arrival times, exponential service times, a “first come, first served” queue discipline, a total capacity of 20 customers, and an infinite population pool from which to draw. The situation is illustrated in Figure 24.4.

Examples: here are some commonly-used/studied queueing systems:

Name	Notation	Example
simple system	$M/M/1$	customer service desk in a small store
multi-server system	$M/M/c$	airline ticket counter
constant service	$M/D/1$	automated car wash
general service	$M/G/1$	auto repair shop
limited capacity	$M/M/1/N$	barber shop with N waiting seats

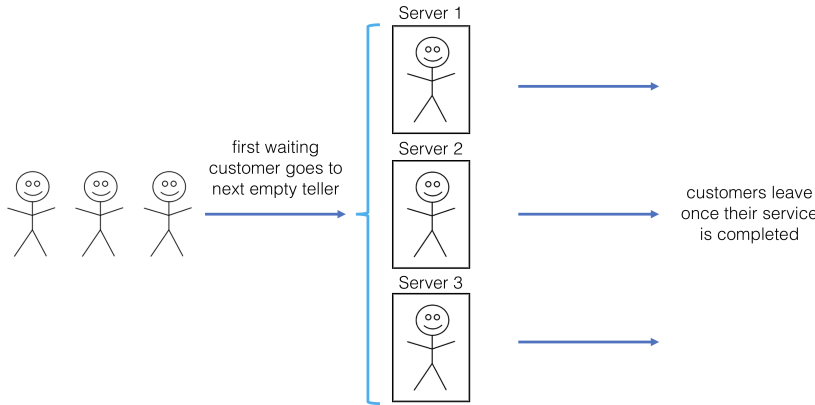


Figure 24.4: Single line at bank with three tellers – $M/M/3/FCFS/20/\infty$.

24.3.2 Birth-Death Processes

The **state** of a queueing system at time t is defined to be the number of customers in the queueing system, either waiting in line or in service, at time t . At $t = 0$, the state of the system is the initial number of customers in the system. This state is worth recording because it clearly affects the state at future times t .

Knowing this, we define $P_{i,j}(t)$ as the probability that the state at time t is j , given that the state at $t = 0$ was i . For large t , $P_{i,j}(t)$ becomes independent of i and approaches a limit π_j . This limit is known as the **steady-state** of state j .

It is generally quite difficult to determine the steps of arrivals and services that lead to a steady-state π_j . Likewise, starting from an early t , it is difficult to determine exactly when a system will reach its steady state π_j , if such a state even exists.

For simplicity’s sake, when a queueing system is studied, we begin by assuming that the steady-state has already been reached. A **birth-death process** is a Markov process in which states are indexed by non-negative integers, and transitions are only permitted between “neighbouring” states. After a “birth”, the state increases from n to $n + 1$; after a “death”, the state decreases from m to $m - 1$.

Typically, we denote the set of birth rates and death rates by λ_n and μ_m , respectively (see Figure 24.5).

Pure birth processes are those for which $\mu_m = 0$ for all m ; **pure death** processes those for which $\lambda_n = 0$ for all n . The **steady-state solution** of a birth-death process, i.e., the probability π_n of being in state n , can actually be computed:

$$\pi_n = \pi_0 \frac{\lambda_0 \lambda_1 \cdots \lambda_{n-1}}{\mu_1 \mu_2 \cdots \mu_n}, \quad \text{for } n = 1, 2, \dots$$

where π_0 is the probability of being in state 0 (i.e., without users). It can further be shown [5] that:

$$\pi_0 = \frac{1}{1 + \sum_{n=1}^{\infty} \prod_{j=0}^{n-1} \frac{\lambda_j}{\mu_{j+1}}}$$

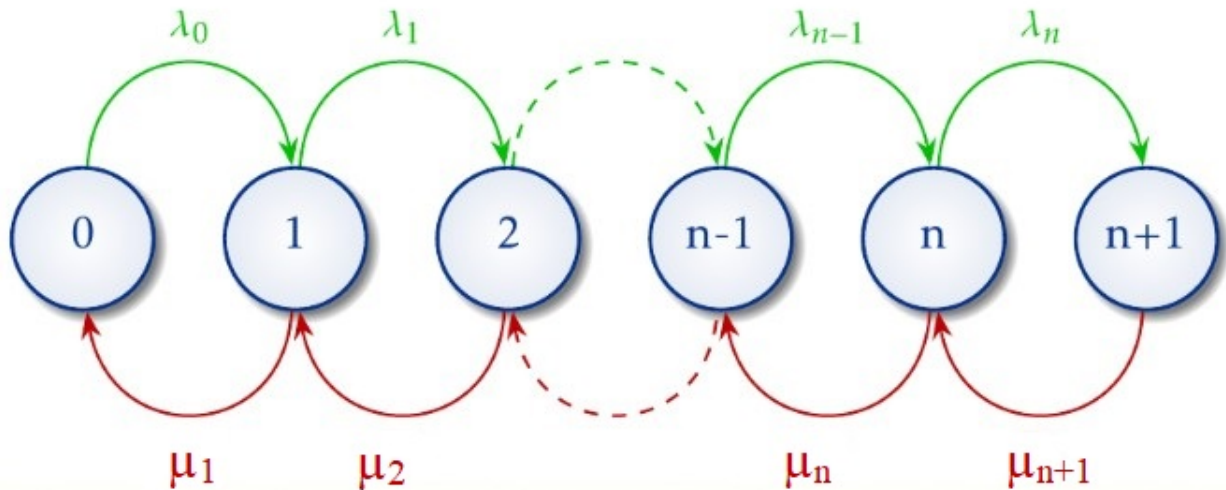


Figure 24.5: Birth-death process; queueing states indexed by integers; birth rates and death rates indicated by λ_n and μ_n , respectively (source unknown).

24.3.3 Little’s Queuing Formula

It is often the case that clients and end users are interested in determining the **amount of time** that a typical customer spends in the queuing system. Let W be the **expected waiting time** spent in the queuing system, including time in line plus time in service, and W_q be the **expected time a customer spends waiting in line**.

Both W and W_q are computed under the assumption that the steady state has been reached. By using a powerful result known as **Little’s queuing formula**, W and W_q are easily related to the number of customers in the queue and those waiting in line. For any queuing system (or any subset of a queuing system), consider the following quantities:

- λ = average number of arrivals entering the system per unit time;
- L = average number of customers present in the queuing system;
- L_q = average number of customers waiting in line;
- L_s = average number of customers in service;
- W = average time a customer spends in the system;
- W_q = average time a customer spends in line, and
- W_s = average time a customer spends in service.

Customers in the system can only be found in the queue or being serviced, so that $L = L_q + L_s$ and $W = W_q + W_s$. In these definitions, all averages are **steady-state averages**.

For most queuing systems in which a steady-state exists, **Little’s queuing formula** are summarized by:

$$L = \lambda W, \quad L_q = \lambda W_q, \quad \text{and} \quad L_s = \lambda W_s.$$

Example: if, on average, 46 customers enter a restaurant each hour it is opened, and if they spend, on average, 10 minutes ($1/6$ hours) waiting to be served, then we should expect $46 \cdot 1/6 \approx 7.7$ customers in the queue at all time (on average).

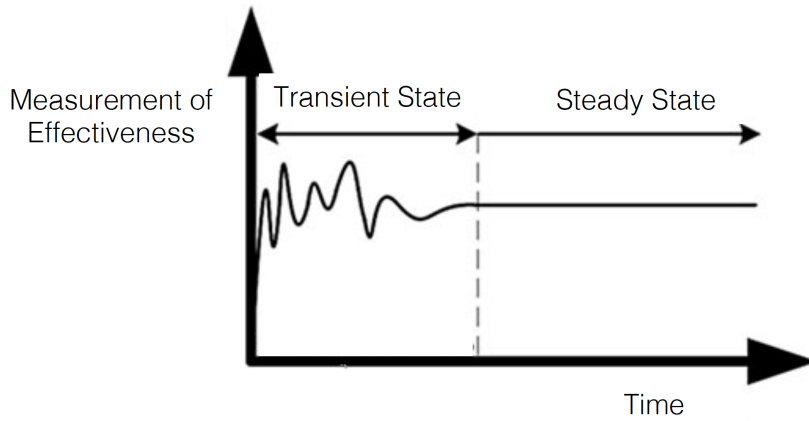


Figure 24.6: Schematics of steady state vs. transient behaviour (source unknown).

24.4 M/M/1 Queueing Systems

We now discuss the **simplest non-trivial** queueing system.

24.4.1 Basics

An $M/M/1/GD/\infty/\infty$ queueing system has exponential inter-arrival times, exponential service times, and a single server. It can be modeled as a birth-death process with

$$\begin{aligned} \lambda_j &= \lambda, \quad j = 0, 1, 2, \dots \\ \mu_0 &= 0 \\ \mu_j &= \mu, \quad j = 1, 2, 3, \dots \end{aligned}$$

Substituting these rates in the steady-state solution of a birth-death process yields

$$\pi_j = \frac{\lambda^j \pi_0}{\mu^j} = \rho^j \pi_0,$$

where $\rho = \lambda/\mu$ is the **traffic intensity** of the system.

Since the system has to be in exactly one of the states at any given moment, the sum of all probabilities is 1:

$$\pi_0 + \pi_1 + \pi_2 + \dots = \pi_0(1 + \rho + \rho^2 + \dots) = 1.$$

If $0 \leq \rho < 1$, the infinite series converges to $\frac{1}{1-\rho}$ from which we derive

$$\pi_0 \cdot \frac{1}{1-\rho} = 1 \implies \pi_0 = 1 - \rho \implies \pi_j = \rho^j \pi_0 = \rho^j (1 - \rho)$$

as the **steady-state probability of state j** .

If $\rho \geq 1$, the infinite series diverges and no steady-state exists. Intuitively, this happens when $\lambda \geq \mu$, that is, if the arrival rate is greater than the service rate, then the state of the system grows without bounds and the queue is never cleared. From this point on, we assume $\rho < 1$ to guarantee that the steady-state probabilities π_j exist, from which we can determine several quantities of interest.

Assuming that the steady state has been reached, it can be shown that L , L_s , and L_q are given respectively by:

$$L = \frac{\lambda}{\mu - \lambda} = \frac{\rho}{1 - \rho}, \quad L_s = \rho, \quad L_q = \frac{\rho^2}{1 - \rho}.$$

Using Little's queuing formula, we can also solve for W , W_s , and W_q by dividing each of the corresponding L values by λ :

$$W = \frac{1}{\mu - \lambda}, \quad W_s = \frac{1}{\mu}, \quad W_q = \frac{\lambda}{\mu(\mu - \lambda)}.$$

Note that, as expected, both $W, W_q \rightarrow +\infty$ when $\rho \rightarrow 1$. On the other hand, $W_q \rightarrow 0$ and $W \rightarrow \frac{1}{\mu}$ (the **mean service time**) as $\rho \rightarrow 0$.

Example: (based on [9]) an average of 10 cars arrive at a single-server drive-in teller every hour. Assume that the average customer is served in 4 minutes, and that both inter-arrival times and service times are exponentially distributed.

1. What is the probability that the teller is idle?
2. Excluding the car that is being served, what is the average number of cars waiting in line at the teller?
3. What is the average amount of time a drive-in customer spends in the bank parking lot (including time in service)?
4. On average, how many customers per hour are served by the teller?

Solution: by assumption, we are dealing with an $M/M/1/GD/\infty/\infty$ queuing system for which $\lambda = 10$ cars/hr and $\mu = 15$ cars/hr, and as such $\rho = 10/15 = 2/3$.

1. The teller is idle one third of the time on average because $\pi_0 = 1 - \rho = 1/3$.
2. There are $L_q = \rho^2/(1 - \rho) = 4/3$ cars waiting in line for the teller.
3. We know that $L = \lambda/(\mu - \lambda) = 10/(15 - 10) = 2$, and so $W = L/\lambda = 0.2$ hr = 12 min.
4. If the teller were always busy, it would serve an average of $\mu = 15$ customers per hour. From part 1., we know that the teller is only busy two-thirds of the time, thus during each hour, the teller serves an average of $15 \cdot 2/3 = 10$ customers. This is reasonable since, in a steady-state, 10 customers are arriving each hour and 10 customers must leave the system every hour.

Example: (based on [6]) Suppose that all car owners fill up when their tanks are exactly half full. On average, 7.5 customers arrive every hour at a single-pump gas station. It takes an average of 4 minutes to fuel a car.

Assume that inter-arrival times and service times are both exponential.

1. What are the values of L and W in this scenario?
2. Suppose that a gas shortage occurs and panic buying takes place. To model this phenomenon, assume that all car owners now purchase gas when their tanks are exactly three-quarters full. Since each car owner is now putting less gas into the tank during each visit to the station, we assume that the average service time has been reduced to 10/3 minutes. How has panic buying affected the values of L and W ?

Solution: by assumption, we again have an $M/M/1/GD/\infty/\infty$ queueing system, with $\lambda = 7.5$ cars/hr and $\mu = 60/4 = 15$ cars/hr. Thus, $\rho = 7.5/15 = 1/2$.

1. By definition, $L = \lambda/(\mu - \lambda) = 7.5/(15 - 7.5) = 1$ and $W = 1/7.5 \approx 0.13$ hr = 7.8 min. Hence, in this situation, everything is under control, and long lines appear to be unlikely.
2. Under the panic buying scenario, $\lambda = 2(7.5) = 15$ cars/hr as each car owner now fills up twice as often, and $\mu = 60 \cdot 3/10 = 18$ cars/hr, so $\rho = \lambda/\mu = 5/6$. In that scenario,

$$L = \frac{\rho}{1 - \rho} = 5 \text{ cars, and } W = \frac{L}{\lambda} = \frac{5}{15} = 20 \text{ min.}$$

Thus, panic buying has more than doubled the wait time in line. In a $M/M/1$ queueing system, we have

$$L = \frac{\rho}{1 - \rho} = -1 + \frac{1}{1 - \rho},$$

and it is easy to see that $L \rightarrow \infty$ as $\rho \rightarrow 1$. The 5-fold increase in L when ρ jumps from $1/2$ to $5/6$ (with accompanying jumps in W) illustrates that fact.

ρ	L in a $M/M/1$ queue
0.30	0.43
0.60	1.50
0.80	4.00
0.90	9.00
0.95	19.00
0.99	99.00

24.4.2 Limited Capacity

In the real world, queues never become infinite – they are limited due to requirements of space and/or time, or service operating policy. Such a queueing model falls under the purview of **finite queues**.

Finite queue models restrict the number of customers allowed in the service system. Let N represent the maximum allowable number of customers in the system. If the system is at **capacity**, the arrival of a $(N + 1)^{\text{th}}$ customer results in a failure to enter the queue – the customer is assumed to balk and depart without seeking service.

Finite queues can also be modeled as a birth-death process, but with a slight modification in its parameters:

$$\begin{aligned} \lambda_j &= \lambda, \quad j = 0, 1, 2, \dots, N - 1 \\ \lambda_N &= 0, \quad \mu_0 = 0 \\ \mu_j &= \mu, \quad j = 1, 2, 3, \dots, N. \end{aligned}$$

The restriction $\lambda_N = 0$ is what sets this model apart from the $M/M/1/\infty$. It makes it impossible to reach a state greater than N . Because of this

restriction, a steady-state always exist because even if $\lambda \geq \mu$, there can never be more than N customers in the system.

Mathematically, this has the effect of replacing the infinite series linking the π_j 's by a finite geometric series, which always converges:

$$\pi_0 + \pi_1 + \dots + \pi_N = \pi_0(1 + \rho + \dots + \rho^N) = 1,$$

from which we can derive

$$\begin{aligned} \pi_0 \cdot \frac{1 - \rho^{N+1}}{1 - \rho} = 1 &\implies \pi_0 = \frac{1 - \rho}{1 - \rho^{N+1}} \\ &\implies \pi_j = \begin{cases} \rho^j \frac{1 - \rho}{1 - \rho^{N+1}} & \text{for } j = 0, \dots, N \\ 0 & \text{for } j > N \end{cases} \end{aligned}$$

Since $L = \sum_{j=0}^N j \cdot \pi_j$ (why?),

$$L = \frac{\rho[1 + N\rho^{N+1} - (N+1)\rho^N]}{(1 - \rho)(1 - \rho^{N+1})}$$

when $\lambda \neq \mu$. As in the $M/M/1/\infty$ queue, $L_s = 1 - \pi_0$, and $L_q = L - L_s$.

In a finite capacity model, only $\lambda - \lambda\pi_N = \lambda(1 - \pi_N)$ arrivals per unit time actually enter the system on average (λ arrive, but $\lambda\pi_N$ find the system full). With this fact,

$$W = \frac{L}{\lambda(1 - \pi_N)} \quad \text{and} \quad W_q = \frac{L_q}{\lambda(1 - \pi_N)}.$$

What does that look like in practice?

Example: consider a one-man barber shop with a total of 10 seats. Assume, as has always been the case so far (but need not be), that inter-arrival times are exponentially distributed with an average of 20 prospective customers arriving each hour at the shop. Those customers who find the shop full do not enter (perhaps they do not like standing). The barber takes an average of 12 minutes to cut each customer's hair; assume that haircut times are also exponentially distributed.

1. On average, how many haircuts per hour will the barber complete?
2. On average, how much time will be spent in the shop by a customer who enters?

Solution:

1. A fraction π_{10} of all arrivals will find the shop full, so that only an average of $\lambda(1 - \pi_{10})$ will actually enter the shop each hour. All entering customers receive a haircut, so the barber will give an average of $\lambda(1 - \pi_{10})$ haircuts per hour. In this scenario, $N = 10$, $\lambda = 20$ customers/hr, and $\mu = 60/12 = 5$ customers/hr. Thus $\rho = 20/5 = 4$ and we have

$$\begin{aligned} \pi_0 &= \frac{1 - \rho}{1 - \rho^{N+1}} = \frac{1 - 4}{1 - 4^{11}} \approx 7.15 \times 10^{-7} \text{ and} \\ \pi_{10} &= 4^{10} \pi_0 = \frac{3}{4} \text{ (from formula in opposite column).} \end{aligned}$$

In that case, an average of $20(1 - 3/4) = 5$ customers per hour will receive haircuts. This means that an average of $20 - 5 = 15$ prospective customers per hour will not enter the shop.

2. To determine W , we must first compute

$$L = \frac{4[1 + (10)4^{11} - (11)4^{10}]}{(1 - 4)(1 - 4^{11})} = 9.67.$$

Using the formulas described above, we obtain

$$W = \frac{L}{\lambda(1 - \pi_{10})} = \frac{9.67}{5} = 1.93 \text{ hr.}$$

This barber shop is quite crowded – the barber would be well-advised to hire at least one more barber!

But what *would* be the effect of hiring a second barber?

In order to answer this question, we need to look into *M/M/c* queueing systems.

24.5 *M/M/c Queueing Systems*

An *M/M/c/GD/∞* queueing system also has exponential inter-arrival and service times, with rates λ and μ , respectively. What sets this system apart is that there are now $c > 1$ servers willing to serve from a single line of customers, perhaps like one would find in a bank (see Figure 24.7).

If $j \leq c$ customers are present in the system, then every customer is being served and there is no wait time; if $j > c$ customers are in the system, then c customers are being served and the remaining $j - c$ customers are waiting in the queue. To model this as a birth-death process, we have to observe that the death rate is dependent on how many servers are actually being used.

If each server completes service at a rate of μ (which may not be the case in practice as there might be variations in servers, at least for human servers), then the **actual death rate** is $\mu \times$ the number of customers actually being served. The parameters for this process are

$$\lambda_n = \lambda, \quad n = 0, 1, 2, \dots$$

$$\mu_n = \begin{cases} n\mu, & n = 0, 1, 2, \dots, c \\ c\mu, & n = c + 1, c + 2, \dots \end{cases}$$

The traffic intensity for the *M/M/c* system is $\rho = \lambda/(c\mu)$ and the steady-state solution is

$$\pi_n = \begin{cases} \frac{(c\rho)^n}{n!} \pi_0, & 1 \leq n \leq c \\ \frac{c^c \rho^n}{c!} \pi_0, & n \geq c \end{cases}$$

where

$$\pi_0 = \left[1 + \frac{(c\rho)^c}{c!(1 - \rho)} + \sum_{n=1}^{c-1} \frac{c\rho^n}{n!} \right]^{-1}.$$

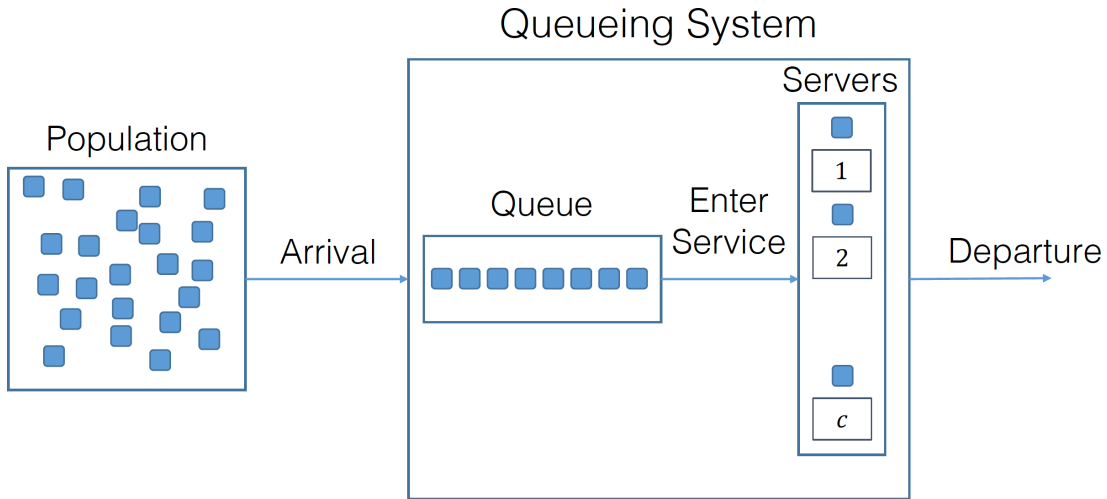


Figure 24.7: Generic $M/M/c$ queue.

Note that, as was the case in a $M/M/1$ system, if $\rho \geq 1$, there can be no steady state – in other words, if the arrival rate is at least as large as the maximum possible service rate ($\lambda \geq c\mu$), then the system “blows up”.

There might be a desire to ensure that customers do not wait in line an inordinate amount of time, but there might also be a desire to minimize the amount of time for which at least one of the server is idle. In a $M/M/c$ queueing system, this steady-state probability is given by

$$P(n \geq c) = \frac{(c\rho)^c}{c!(1-\rho)}\pi_0.$$

This table shows the probabilities $P(n \geq c)$ that all servers are busy in an $M/M/c$ system for $c = 2, \dots, 7$ and $0.1 \leq \rho \leq 0.95$ [9, p.1088].

ρ	$c = 2$	$c = 3$	$c = 4$	$c = 5$	$c = 6$	$c = 7$
.10	.02	.00	.00	.00	.00	.00
.20	.07	.02	.00	.00	.00	.00
.30	.14	.07	.04	.02	.01	.00
.40	.23	.14	.09	.06	.04	.03
.50	.33	.24	.17	.13	.10	.08
.55	.39	.29	.23	.18	.14	.11
.60	.45	.35	.29	.24	.20	.17
.65	.51	.42	.35	.30	.26	.21
.70	.57	.51	.43	.38	.34	.30
.75	.64	.57	.51	.46	.42	.39
.80	.71	.65	.60	.55	.52	.49
.85	.78	.73	.69	.65	.62	.60
.90	.85	.83	.79	.76	.74	.72
.95	.92	.91	.89	.88	.87	.85

Cumbersome calculations, using $W_s = \frac{1}{\mu}$, yield

$$L_q = \frac{\rho}{1-\rho}P(n \geq c), \quad W_q = \frac{L_q}{\lambda}, \quad W = \frac{1}{\mu} + W_q, \quad L = \frac{\lambda}{\mu} + L_q.$$

Example: consider, for instance, a bank with two tellers. An average of 80 customers arrive at the bank each hour and wait in a single line for an idle teller. For this specific bank, the average service time is 1.2 minutes. Assume that inter-arrival times and service times are exponential. Determine:

1. The expected number of customers in the bank.
2. The expected length of time a customer spends in the bank.
3. The fraction of time that a particular teller is idle.

Solutions: we are dealing with an $M/M/2$ system with $\lambda = 80$ customers/hr and $\mu = 50$ customers/hr. Thus, $\rho = \frac{80}{2 \cdot 50} = 0.80 < 1$ and the steady-state exists.

1. From the above table, $P(n \geq 2) = 0.71$, from which we compute

$$L_q = P(n \geq 2) \cdot \frac{.8}{1 - .8} = 2.84 \text{ customers}$$

$$L = \frac{80}{50} + L_q = 4.44 \text{ customers.}$$

2. We know that $W = \frac{L}{\lambda} = \frac{4.44}{80} = 0.055 \text{ hr} = 3.3 \text{ min.}$
3. To determine the fraction of time that a particular server is idle, note that tellers are idle during all moments when $n = 0$, and half the time (by symmetry) when $n = 1$. The probability that a server is idle is thus given by $\pi_0 + 0.5\pi_1$. But

$$\pi_0 = \left[1 + \frac{(2 \cdot .8)^2}{2!(1 - .8)} + \sum_{n=1}^{2-1} \frac{2 \cdot .8^n}{n!} \right]^{-1} = \frac{1}{9}$$

and

$$\pi_1 = \frac{1.6}{1!} \pi_0 = 0.176$$

and so the probability that particular teller is idle is $0.111 + 0.5(0.176) = 0.199$.

Important Note: general queueing models are not understood to the same extent as $M/M/1$ (and $M/M/c$ to a lesser extent), and their given performance measurements may only be approximate and highly-dependent on the specifics of the problem at hand.

For this reason, $M/M/c$ models are sometimes used even when their use is not supported by the data (the situation is not unlike the widespread use of the normal distribution in a variety of probability and statistics problems).

In numerous applications, the empirical distributions of arrivals and service times are nearly Poisson and exponential, respectively, so that the assumption is not entirely off the mark, but numerical simulations should not be eschewed when departures from the $M/M/c$ model are too pronounced.

24.6 Exercises

The *Borealian Aeronautic Security Agency* (BASA) runs pre-board screening of passengers and crew for all flights departing the nation's airfields. There are 4 Major Airfields:

- Auckland
- Chebucto
- Saint-François
- Queenston

The screening process (PBS) is structurally similar at each airfield:

1. Passengers arrive at the beginning of the main queue
2. Boarding passes may or may not be scanned at S₁
3. Passengers enter the main queue
4. Boarding passes are scanned at S₂
5. Passengers are directed to a server entry position
6. Passengers and carry-on luggage are screened by a server

Some factors influence the PBS wait time, including:

- schedule intensity of departing flights
- passenger volume on these flights
- number of servers and processing rates at a given airfield, etc.

There might also be:

- yearly, seasonal, time-of-day, day-of-week interaction effects (among others) depending on the airfield, the flight destination, etc.
- trend level shifts in the number of passengers, flights, destinations, etc.

Datasets: [20262030.csv](#), [BASA_AUC_2028_912.csv](#), [dat_F_sub.csv](#), [dat_P_sub_c.csv](#).

1. Build a data dictionary for the datasets
2. Explore and visualize the datasets
3. Perform a queueing model analysis to predict the wait times at each airfield for which you have data.

Use the CATSA case study to inform your analysis [2].

Chapter References

- [1] R. Berry. *Queueing Theory and Applications*. 2nd. PWS/Kent Publishing, 2002.
- [2] P. Boily and J. Schellinck. *Introduction to Quantitative Consulting*. Quadrangle/Data Action Lab, 2025.
- [3] A.K. Erlang. 'The theory of probabilities and telephone conversations'. In: *Nyt Tidsskrift for Matematik B* (1909).
- [4] D.G. Kendall. 'Stochastic Processes Occurring in the Theory of Queues and their Analysis by the Method of the Imbedded Markov Chain'. In: *The Annals of Mathematical Statistics* 24.3 (1953), pp. 338–354. DOI: [10.1214/aoms/1177728975](https://doi.org/10.1214/aoms/1177728975).
- [5] L. Kleinrock. *Queueing Systems, Volume I*. Wiley-Interscience. Wiley, 1974.
- [6] 'Management Science and the Gas Shortage'. In: *Interfaces* 4.4 (Aug. 1974), pp. 47–51.
- [7] C. Newell. *Applications of Queueing Theory*. Ettore Majorana International Science Series. Springer Netherlands, 2013.
- [8] S.M. Ross. *Introduction to Probability Models*. 11th ed. San Diego, CA, USA: Academic Press, 2014.
- [9] W.L. Winston. *Operations Research: Applications and Algorithms*. Cengage Learning, 2022.