

# Queueing theory

From Wikipedia, the free encyclopedia  
(Redirected from Queueing theory)



This article is **missing information about criticism of the scope of applicability of queueing theory to real problems**. Please expand the article to include this information. Further details may exist on the talk page. *(July 2014)*

**Queueing theory** is the **mathematical** study of waiting lines, or **queues**.<sup>[1]</sup> In queueing theory a model is constructed so that queue lengths and waiting time can be predicted.<sup>[1]</sup> Queueing theory is generally considered a branch of **operations research** because the results are often used when making business decisions about the resources needed to provide a service.

Queueing theory has its origins in research by **Agner Krarup Erlang** when he created models to describe the Copenhagen telephone exchange.<sup>[1]</sup> The ideas have since seen applications including **telecommunication**, **traffic engineering**, **computing**<sup>[2]</sup> and the design of factories, shops, offices and hospitals.<sup>[3][4]</sup>

Contents [hide]

1

Single queueing nodes

2

Service disciplines

3

Queueing networks

3.1

Example of M/M/1

3.2

Routing algorithms

4

Mean field limits

5

Fluid limits

6

Heavy traffic/diffusion approximations

7

Software for simulation/analysis

8

See also

9

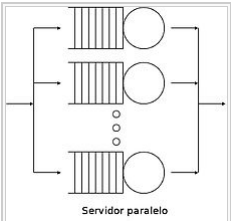
References

10

Further reading

11

External links



**Queue networks** are systems in which single queues are connected by a routing network. In this image servers are represented by circles, queues by a series of rectangles and the routing network by arrows. In the study of queue networks one typically tries to obtain the **equilibrium distribution** of the network, although in many applications the study of the **transient state** is fundamental.

## Single queueing nodes [edit]

Single queueing nodes are usually described using **Kendall's notation** in the form *A*/*S*/*C* where *A* describes the time between arrivals to the queue, *S* the size of jobs and *C* the number of servers at the node.<sup>[5][6]</sup> Many theorems in queueing theory can be proved by reducing queues to mathematical systems known as **Markov chains**, first described by **Andrey Markov** in his 1906 paper.<sup>[7]</sup>

**Agner Krarup Erlang**, a **Danish** engineer who worked for the Copenhagen Telephone Exchange, published the first paper on what would now be called queueing theory in 1909.<sup>[8][9][10]</sup> He modeled the number of telephone calls arriving at an exchange by a **Poisson process** and solved the **M/D/1 queue** in 1917 and **M/D/k queueing** model in 1920.<sup>[11]</sup> In Kendall's notation:

- M stands for Markov or memoryless and means arrivals occur according to a **Poisson process**
- D stands for deterministic and means jobs arriving at the queue require a fixed amount of service
- k* describes the number of servers at the queueing node (*k* = 1, 2,...). If there are more jobs at the node than there are servers then jobs will queue and wait for service

The **M/M/1 queue** is a simple model where a single server serves jobs that arrive according to a Poisson process and have **exponentially distributed** service requirements. In an **M/G/1 queue** the *G* stands for general and indicates an arbitrary **probability distribution**. The M/G/1 model was solved by **Felix Pollaczek** in 1930.<sup>[12]</sup> a solution later recast in probabilistic terms by **Aleksandr Khinchin** and now known as the **Pollaczek–Khinchine formula**.<sup>[11]</sup>

After **World War II** queueing theory became an area of research interest to mathematicians.<sup>[11]</sup> In 1953 **David Kendall** solved the G/M/k queue<sup>[13]</sup> and introduced the modern notation for queues, now known as **Kendall's notation**. In 1957 Pollaczek studied the G/G/1 using an **integral equation**.<sup>[14]</sup> **John Kingman** gave a formula for the **mean waiting time** in a **G/G/1 queue**: **Kingman's formula**.<sup>[15]</sup>

The **matrix geometric method** and **matrix analytic methods** have allowed queues with **phase-type distributed** inter-arrival and service time distributions to be considered.<sup>[16]</sup>

Problems such as performance metrics for the **M/G/k queue** remain an open problem.<sup>[11]</sup>

## Service disciplines [edit]

Various scheduling policies can be used at queueing nodes:

### First in first out

This principle states that customers are served one at a time and that the customer that has been waiting the longest is served first.<sup>[17]</sup>

### Last in first out

This principle also serves customers one at a time, but the customer with the shortest **waiting time** will be served first.<sup>[17]</sup> Also known as a **stack**.

### Processor sharing

Service capacity is shared equally between customers.<sup>[17]</sup>

### Priority

Customers with high priority are served first.<sup>[17]</sup> Priority queues can be of two types, non-preemptive (where a job in service cannot be interrupted) and preemptive (where a job in service can be interrupted by a higher priority job). No work is lost in either model.<sup>[18]</sup>

### Shortest job first

The next job to be served is the one with the smallest size

### Preemptive shortest job first

The next job to be served is the one with the original smallest size<sup>[19]</sup>

### Shortest remaining processing time

The next job to serve is the one with the smallest remaining processing requirement.<sup>[20]</sup>

### Service facility

- Single server:customers line up and there is only one server
- Parallel servers:customers line up and there are several servers
- Tandem queue:there are many counters and customers can decide going where to queue

### Customer's behavior of waiting

- Balking:customers deciding not to join the queue if it is too long
- Jockeying:customers switch between queues if they think they will get served faster by so doing
- Reneging:customers leave the queue if they have waited too long for service

## Queueing networks [edit]

Networks of queues are 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 and queue for service, or leave the network. For a network of *m* the state of the system can be described by an *m*–dimensional vector (*x*<sub>1</sub>,*x*<sub>2</sub>,...,*x*<sub>*m*</sub>) where *x*<sub>*i*</sub> represents the number of customers at each node.

The first significant results in this area were **Jackson networks**,<sup>[21][22]</sup> for which an efficient **product-form stationary distribution** exists and the **mean value analysis**<sup>[23]</sup> which allows average metrics such as throughput and sojourn times to be computed.<sup>[24]</sup> If the total number of customers in the network remains constant the network is called a closed network and has also been shown to have a product–form stationary distribution in the **Gordon–Newell theorem**.<sup>[25]</sup> This result was extended to the **BCMP network**<sup>[26]</sup> where a network with very general service time, regimes and customer routing is shown to also exhibit a product-form stationary distribution. The **normalizing constant** can be calculated with the **Buzen's algorithm**, proposed in 1973.<sup>[27]</sup>

Networks of customers have also been investigated, **Kelly networks** where customers of different classes experience different priority levels at different service nodes.<sup>[28]</sup> Another type of network are **G-networks** first proposed by **Erol Gelenbe** in 1993:<sup>[29]</sup> these networks do not assume exponential time distributions like the classic Jackson Network.

### Example of M/M/1 [edit]

#### Birth and Death process

- A/B/C



21. <sup>^</sup> Jackson, J. R. (1957). "Networks of Waiting Lines". *Operations Research* **5** (4): 518–521. doi:10.1287/opre.5.4.518 ⓘ. JSTOR 167249 ⓘ.

22. <sup>^</sup> Jackson, James R. (Oct 1963). "Jobshop-like Queueing Systems". *Management Science* **10** (1): 131–142. doi:10.1287/mnsc.1040.0268 ⓘ. JSTOR 2627213 ⓘ.

23. <sup>^</sup> Reiser, M.; Lavenberg, S. S. (1980). "Mean-Value Analysis of Closed Multichain Queueing Networks". *Journal of the ACM* **27** (2): 313. doi:10.1145/322186.322195 ⓘ.

24. <sup>^</sup> Van Dijk, N. M. (1993). "On the arrival theorem for communication networks". *Computer Networks and ISDN Systems* **25** (10): 1135–2013. doi:10.1016/0169-7552(93)90073-D ⓘ.

25. <sup>^</sup> Gordon, W. J.; Newell, G. F. (1967). "Closed Queueing Systems with Exponential Servers". *Operations Research* **15** (2): 254. doi:10.1287/opre.15.2.254 ⓘ. JSTOR 168557 ⓘ.

26. <sup>^</sup> Baskett, F.; Chandu, K. Mani; Muntz, R.R.; Palacios, F.G. (1975). "Open, closed and mixed networks of queues with different classes of customers". *Journal of the ACM* **22** (2): 248–260. doi:10.1145/321879.321887 ⓘ.

27. <sup>^</sup> Buzen, J. P. (1973). "Computational algorithms for closed queueing networks with exponential servers" (PDF). *Communications of the ACM* **16** (9): 527. doi:10.1145/362342.362345 ⓘ.

28. <sup>^</sup> Kelly, F. P. (1975). "Networks of Queues with Customers of Different Types". *Journal of Applied Probability* **12** (3): 542–554. doi:10.2307/3212869 ⓘ. JSTOR 3212869 ⓘ.

29. <sup>^</sup> Gelenbe, Erol (Sep 1993). "G-Networks with Triggered Customer Movement". *Journal of Applied Probability* **30** (3): 742–748. doi:10.2307/3214781 ⓘ. JSTOR 3214781 ⓘ.

30. <sup>^</sup> Bobbio, A.; Griboaldo, M.; Telek, M. S. (2008). "Analysis of Large Scale Interacting Systems by Mean Field Method". *2008 Fifth International Conference on Quantitative Evaluation of Systems*. p. 215. doi:10.1109/QEST.2008.47 ⓘ. ISBN 978-0-7695-3360-5.

31. <sup>^</sup> Bramson, M. (1999). "A stable queueing network with unstable fluid model". *The Annals of Applied Probability* **9** (3): 818. doi:10.1214/aop/1029962815 ⓘ. JSTOR 2667284 ⓘ.

32. <sup>^</sup> Chen, H.; Whitt, W. (1993). "Diffusion approximations for open queueing networks with service interruptions". *Queueing Systems* **13** (4): 335. doi:10.1007/BF01149260 ⓘ.

33. <sup>^</sup> Yamada, K. (1995). "Diffusion Approximation for Open State-Dependent Queueing Networks in the Heavy Traffic Situation". *The Annals of Applied Probability* **5** (4): 958. doi:10.1214/aop/1177004602 ⓘ. JSTOR 2245101 ⓘ.

Further reading [edit]

- Gross, Donald; Carl M. Harris (1998). *Fundamentals of Queueing Theory*. Wiley. ISBN 0-471-32812-X. Online ⓘ.
- Deitel, Harvey M. (1984) [1982]. *An introduction to operating systems* ⓘ (revisited first ed.). Addison-Wesley. p. 673. ISBN 0-201-14502-2. chap.15, pp. 380–412
- Lazowska, Edward D.; John Zahorjan; G. Scott Graham; Kenneth C. Sevcik (1984). *Quantitative System Performance: Computer System Analysis Using Queueing Network Models* ⓘ. Prentice-Hall, Inc. ISBN 0-13-746975-6.
- Zukerman, Moshe. *Introduction to Queueing Theory and Stochastic Teletraffic Models* (PDF).

External links [edit]

- Queueing theory calculator ⓘ
- Teknomo's Queueing theory tutorial and calculators ⓘ
- Virtamo's Queueing Theory Course ⓘ
- Myron Hlynka's Queueing Theory Page ⓘ
- Queueing Theory Basics ⓘ
- A free online tool to solve some classical queueing systems ⓘ

<span><span>v</span> · <span>t</span> · <span>e</span></span>	<b>Queueing theory</b> <span><span>[</span>hide<span>]</span></span>
<b>Single queueing nodes</b>	DM/1 queue · MD/1 queue · MD/c queue · MM/1 queue (Burke's theorem) · MM/c queue · MM/∞ queue · MG/1 queue (Pollaczek–Khinchine formula · Matrix analytic method) · MG/k queue · GM/1 queue · G/G/1 queue (Kingman's formula · Lindley equation) · Fork-join queue · Bulk queue
<b>Arrival processes</b>	Poisson process · Markovian arrival process · Rational arrival process
<b>Queueing networks</b>	Jackson network (Traffic equations) · Gordon–Newell theorem (Mean value analysis · Buzen's algorithm) · Kelly network · G-network · BCMP network
<b>Service policies</b>	FIFO · LIFO · Processor sharing · Shortest job first · Shortest remaining time
<b>Key concepts</b>	Continuous-time Markov chain · Kendall's notation · Little's law · Product-form solution (Balance equation · Quasireversibility · Flow-equivalent server method) · Arrival theorem · Decomposition method · Beneš method
<b>Limit theorems</b>	Fluid limit · Mean field theory · Heavy traffic approximation (Reflected Brownian motion)
<b>Extensions</b>	Fluid queue · Layered queueing network · Polling system · Adversarial queueing network · Loss network · Retrial queue
<b>Information systems</b>	Data buffer · Erlang (unit) · Erlang distribution · Flow control (data) · Message queue · Network congestion · Network scheduler · Pipeline (software) · Quality of service · Scheduling (computing) · Teletraffic engineering
<span><span><span></span></span></span> <span>Category</span>	
<b>Authority control</b>	NDL: 00567524 <span> </span> <span><span>ⓘ</span></span>

Categories: Stochastic processes | Production and manufacturing | Services management and marketing | Operations research | Formal sciences | Queueing theory | Rationing and licensing | Network performance | Markov models | Markov processes