

# Networks of Queues: Myth and Reality

Guang-Liang Li \*

Victor O. K. Li

Department of Electrical and Electronic Engineering

The University of Hong Kong

Room 601, Chow Yei Ching Bldg.

Pokfulam Road, Hong Kong

{glli,vli}@eee.hku.hk

Phone: (852)2857 8495 Fax: (852)2559 8738

## Abstract

*Jackson's theorem, which is the basis of the mathematical theory of networks of queues, does not hold.*

## 1 Introduction

Networks of queues are basic models for the analysis and design of computer networks [5], and represent, in their own right, an important research field, originated by the seminal work of Jackson [3, 4]. Various networks of queues, proposed after [3], are different generalizations or variations of a class of fundamental models, referred to as Jackson networks of queues, studied in [3]. For this reason, the classical result of Jackson, known as Jackson's theorem [3], is considered the cornerstone of the mathematical theory of networks of queues. However, Jackson's theorem does not hold.

In Section 2, we revisit Jackson's theorem. In Section 3, we disprove Jackson's theorem with simple counterexamples. In Section 4, we show that the limitation of the existing theory of stochastic modeling may explain why Jackson's proof and all other proofs of Jackson's theorem are flawed. After a brief discussion, we conclude in Section 5 by pointing out the implication of our result to networking studies.

## 2 Jackson's Theorem Revisited

We first introduce the notations used, and the assumptions made, by Jackson [3]. There are  $M$  queues (referred to by Jackson as "departments") in a Jackson network, labeled by  $m = 1, 2, \dots, M$ . At queue  $m$ , there are  $n_m$  servers; the service time required

by a customer is exponentially distributed, with mean  $1/\mu_m$ ; customers arrive from outside the network according to a Poisson process at rate  $\lambda_m \geq 0$ .

Once served at queue  $m$ , a customer may either go to queue  $k$  in the network with probability  $\theta_{km}$ , or leave the network with probability  $\theta_m^* = 1 - \sum_k \theta_{km}$ , where  $k = m$  is allowed. The total arrival rate at queue  $m$  is denoted by  $\Gamma_m$ . It is assumed

$$\Gamma_m < \mu_m n_m, \quad m = 1, 2, \dots, M. \quad (1)$$

The number of customers waiting and in service at queue  $m$  is  $k_m$ . At time  $t$ , the state of the network is defined by a vector  $(k_1, k_2, \dots, k_M)$ , which is further assumed to be a Markov process, with probability  $P_{k_1, \dots, k_M}(t)$ . Finally, write  $\alpha_i(k) = \min\{k, n_i\}$  and  $\delta_i = \min\{k_i, 1\}$ .

Jackson's theorem says, in essence, that "so far as steady states are concerned," queue  $m$  for  $m = 1, 2, \dots, M$  in a Jackson network behaves as if it was an M/M/ $n_m$  queue in isolation, independent of all other queues in the network. The proof given by Jackson begins with the following equation [3].

$$\begin{aligned} P_{k_1, \dots, k_M}(t+h) = & \left\{ 1 - \left( \sum \lambda_i \right) h - \left[ \sum \alpha_i(k_i) \mu_i \right] h \right\} \\ & P_{k_1, \dots, k_M}(t) \\ & + \sum \alpha_i(k_i+1) \mu_i \theta_i^* h P_{k_1, \dots, k_i+1, \dots, k_M}(t) \\ & + \sum \lambda_i \delta_i h P_{k_1, \dots, k_i-1, \dots, k_M}(t) \\ & + \sum \sum \alpha_j(k_j+1) \mu_j \theta_{ij} h \\ & P_{k_1, \dots, k_j+1, \dots, k_i-1, \dots, k_M}(t) + o(h) \end{aligned} \quad (2)$$

where  $h$  is a small time increment, and  $o(h)$  tends to zero faster than  $h$ , as  $h \rightarrow 0$ . Although not mentioned explicitly in [3], it should be clear that the condition  $j \neq i$  must be imposed on the double summation in (2).

\*The corresponding author.

### 3 Counterexamples

We give two simple counterexamples to Jackson's theorem. More counterexamples can be constructed similarly.

#### Counterexample 1: Single-Server Queue with Feedback

The simplest Jackson network, consisting of only one single-server queue with feedback, is a counterexample of Jackson's theorem.

This single-server queue with feedback has been used in the literature to explain Jackson's theorem, and is directly obtained from Jackson's model [3], by setting  $M = 1$  and  $n_1 = 1$ . In addition, write  $k_1 = k$ ,  $\lambda_1 = \lambda > 0$ ,  $\mu_1 = \mu > 0$ , and  $\theta_{11} = \theta > 0$ . Now for  $k > 1$ , (2) becomes

$$P_k(t+h) = (1 - \lambda h - \mu h)P_k(t) + \mu(1 - \theta)hP_{k+1}(t) + \lambda hP_{k-1}(t) + o(h). \quad (3)$$

Denote by  $Q(t)$  the number of customers waiting and being served in the single-server queue at time  $t$ . Equation (3) is actually

$$\begin{aligned} P\{Q(t+h) = k\} &= P\{Q(t+h) = k | Q(t) = k\}P\{Q(t) = k\} \\ &+ P\{Q(t+h) = k | Q(t) = k+1\}P\{Q(t) = k+1\} \\ &+ P\{Q(t+h) = k | Q(t) = k-1\}P\{Q(t) = k-1\} \\ &+ o(h). \end{aligned} \quad (4)$$

For simplicity, we follow the convention of ignoring higher order of infinitesimal. This will not affect our conclusion. Compare (3) with (4). We see immediately

$$P\{Q(t+h) = k | Q(t) = k\} = 1 - \lambda h - \mu h.$$

Replace  $k$  respectively by  $k-1$  and  $k+1$  in (3) and (4), and then compare the obtained equations. We see

$$P\{Q(t+h) = k-1 | Q(t) = k\} = \mu(1 - \theta)h$$

and

$$P\{Q(t+h) = k+1 | Q(t) = k\} = \lambda h.$$

We arrive at a contradiction: The sum of the above probabilities is not equal to one. Since (3) is not valid, a steady state derived based on (3) does not exist, and (1) is not meaningful. The contradiction shows that  $Q(t)$  is not a Markov process. This can also be seen without diagnosing (3).

Two time epochs,  $r$  and  $s$ , where  $r < s$ , define a time interval  $[r, s)$ , such that  $Q(r-) \neq k$ , where

$k > 0$  is given,  $Q(t) = k$  for  $t \in [r, s)$ , and  $Q(s) = k-1$ . Clearly,  $s-r$  is a sum of a random number of exponential random variables of rate  $\mu$ . Denote by  $Y_j$  the  $j$ th random variable in  $[r, s)$ .

$$s - r = \sum_{j=1}^N Y_j \quad (5)$$

where  $N$  is a geometric random variable.

For  $Q(t)$  to be a Markov process,  $s-r$  should be an exponential random variable. In an effort to salvage Jackson's theorem, it is claimed that  $s-r$  is exponential with parameter  $\mu(1-\theta)$ . Since  $s-r$  is a directly observable quantity, if the claim is true, in a sample path of  $Q(t)$ , realized values of the exponential random variable of rate  $\mu(1-\theta)$  should be directly observed. So the claim can be tested by the following experiment.

For given  $\lambda$ ,  $\mu$ , and  $\theta$ , let us simulate Jackson's feedback queue, and observe the realized values of  $s-r$ . If a sampled  $r-s$  is a realized value of the exponential random variable with parameter  $\mu$ , we color this  $s-r$  red. If the sampled  $r-s$  consists of at least two realized values of the exponential random variable with parameter  $\mu$ , we color it blue.

Such realized values of  $r-s$  are what we can see in the experiment. The different colors indicate that  $r-s$  is not a random variable with an unvarying distribution. Realized values of  $r-s$  colored red and blue are clearly not sampled from the same distribution. It is not even necessary for realized values of  $r-s$  colored blue to be sampled from the same distribution.

As shown by the experiment, realized values of the exponential random variable of rate  $\mu(1-\theta)$  do not exist in any sample path of  $Q(t)$ . In other words, the claim that  $r-s$  is exponential with rate  $\mu(1-\theta)$  is false. On the contrary, since  $r-s$  is not a random variable with any fixed distribution,  $Q(t)$  is not a Markov process, i.e., the experiment confirms the flaw in Jackson's proof.

#### Counterexample 2: Two Queues in Tandem

To show that Jackson's theorem may not necessarily hold, even if customers are restricted to visiting any queue in a Jackson network at most once, we impose this restriction in the following analysis.

Consider a simple Jackson network of two single-server queues in tandem, such that all customers arrive at the first queue, go to the second queue after service, and leave the network from there.

For this network, Jackson's theorem is a corollary of Burke's theorem [1]: By applying Burke's theorem

to the first queue, the output from the first queue is considered a Poisson process [1], and the second queue is considered an M/M/1 queue. This is, however, incorrect.

Let us simulate the first queue, or do the simulation as a thought experiment, and observe the inter-departure time sequence of this stable M/M/1 queue in steady state.

Let  $s$  and  $t$ , where  $s < t$ , be two arbitrary, consecutive departure epochs of the queue. There are two cases: The server is exclusively either (a) busy, or (b) idle, at time  $s$ . From a straightforward observation, the inter-departure time  $t - s$  is distributed as a service time in case (a), and as the sum of an idle time of the server and a service time in case (b).

An inter-departure time  $t - s$  is sampled in either case. In case (a), we color a line segment of length  $t - s$  red, and use “R” to represent the segment. In case (b), we color a segment of length  $t - s$  blue, and use “B” to represent the segment.

Corresponding to every sample path of the inter-departure time sequence, we now have a sequence of colored segments, which may look like

*RRRBRRRRBBBBRRRRBBRR...*

We can see a tendency for the segments with the same color to aggregate. This is because the segments are colored according to the state of the queue, which is a Markov process. So, the colored segments are Markov dependent. Due to the same reason, the corresponding inter-departure times in the physical departure process are also Markov dependent, not mutually independent.

As visualized by the colored segments, the inter-departure times in the physical departure process depend on the state of the queue, and follow two different distributions. Clearly, such inter-departure times cannot form an i.i.d. exponential sequence. The physical departure process corresponding to such inter-departure times is evidently not a Poisson process.

Burke discovered the Poisson output, by averaging out the impact of the state of the queue [1]. Taking the sequence of colored segments as an example, we illustrate how the impact of the state of the queue is averaged out, and hence the Poisson output is obtained based on the physical departure process. We begin with constructing a histogram, using the lengths of the colored segments as raw materials.

Consider interval  $(0, H)$  on the real line. Divide the interval into  $N(H)$  consecutive, disjoint small subintervals of equal length. As  $H$  increases, the number of subintervals  $N(H)$  increases, and the length of the subinterval decreases.

For all segments with lengths less than  $H$ , calculate the relative frequencies that the lengths of the segments are in the small subintervals, *regardless of their colors*. Increase  $H$  gradually. As  $H \rightarrow \infty$ , we obtain a probability density function (pdf), which is exponential with parameter identical to the arrival rate of the queue.

Sample random variables *independently, regardless of the state of the queue*, from so constructed exponential pdf. The sampled random variables are then mutually independent, follow an identical exponential distribution, and hence form an i.i.d. exponential sequence. Clearly, the Markov dependence observed in the sequence of colored segments does not exist in the exponential sequence, and the random variables of the exponential sequence are distributed neither as the lengths of the red segments, nor as the lengths of the blue segments. The Poisson process, corresponding to the exponential sequence, is the Poisson output constructed analytically by Burke in [1].

Actually, Burke’s theorem was motivated by isolating queues (referred to by Burke also as “stages”) in a tandem system, rather than analyzing them jointly, as said by Burke himself [1]:

“It is intuitively clear that, in tandem queuing processes of the type mentioned above, if the output distribution of each stage was of such character that the queuing system formed by the second stage was amenable to analysis, then the tandem queue could be analyzed stage-by-stage insofar as the separate delay and queue-length distributions are concerned. Such a stage-by-stage analysis can be expected to be considerably simpler than the simultaneous analysis heretofore necessary. Fortunately, under the conditions stated below, it is true that the output has the required simplicity for treating each stage individually.”

If the two queues in tandem are considered jointly as a network, the output from the first queue is the physical departure process. Since the physical departure process is not a Poisson process, the second queue *in the network* is not an M/M/1 queue.

Furthermore, the only output that can make the second queue amenable to a *steady-state* analysis is the constructed Poisson process, since the physical departure process is not stationary. As can be seen from the experiment, the inter-departure times in the physical departure process always follow two different distributions, even if the first queue is in steady state, and hence they cannot form a stationary random sequence. If the physical departure process is used as the input of the second queue, (1) is again not meaningful, since the arrival rate of the non-stationary arrival process at the second queue is not a constant. Thus, driven

by the non-stationary input process, the state of the second queue is not a stationary process, and hence the state of the tandem system as a whole cannot be modeled by a stationary process.

## 4 Theoretical Explanation

Consider again  $t - s$ , the inter-departure time in the physical departure process of the first queue in counterexample 2. If we do not assume the state of the first queue at  $s$ , what is  $t - s$  then? According to the existing literature,  $t - s$  may only be modeled by (i) a random variable conditioned on the state of the queue at  $s$ , or (ii) an unconditional random variable, i.e., an exponential random variable following the inter-arrival time distribution [1].

Option (i) is not suitable, since we do not assume the state of the queue. Option (ii) contradicts the fact observed in the experiment. Thus,  $t - s$  exemplifies a phenomenon that cannot be described by a random variable. To describe such phenomena, we define a double random element, based on a random variable  $X$ , a collection  $E$  of random variables, and a map from  $A$ , the set of all possible values of  $X$ , to  $E$ . The random variable  $X$  and the set  $A$  are respectively the *index* and the *index set* of the double random element. Denote such an object by  $X|E$ . We may read  $X|E$  as “double random element  $X, E$ ,” or “double random element indexed by  $X$ ,” if omission of  $E$  will not cause confusion. The members of  $E$  are *components* of  $X|E$ .

**Definition 1** A double random element  $X|E$  is a bijection from  $A$  to  $E$ , where  $A$  is the set of all possible values of a random variable  $X$ , and  $E$  is a collection of random variables. Any realization of  $X|E$  begins with a realization of  $X$ .

Corresponding to  $x \in A$ , the random variable in  $E$  is denoted by  $x|E$ . We call  $x|E$  “the component of  $X|E$  indexed by  $x$ ,” or “the component indexed by  $x$ ,” if  $X$  and  $E$  can be omitted without causing confusion. The theorem below is an immediate result of Definition 1, and can be used as an alternative definition of double random elements.

**Theorem 1** (a) Any realization of  $X|E$  follows two steps:  $X$  must be realized first, and a realization of  $x|E$  follows, where  $x$  is the realized value of  $X$ . (b)  $z$  is a realized value of  $X|E$ , if and only if  $z$  is a realized value of some  $x|E$ , where  $x \in A$ .

In this paper, we shall focus on a special type of double random elements, such that  $X$  is a discrete random variable, and  $x|E$  is a continuous random variable for each  $x \in A$ .

**Definition 2** The marginal version of a double random element  $X|E$ , denoted by  $[X|E]$ , is a random variable with a distribution

$$P\{[X|E] \leq z\} = \sum_{x \in A} P\{x|E \leq z\}P\{X = x\}.$$

If  $A$  is a set of nonnegative integers, we denote  $X$  and  $X|E$  by  $N$  and  $N|E$ , respectively. The distribution of the marginal version of  $N|E$  is

$$P\{[N|E] \leq z\} = \sum_{n \in A} P\{n|E \leq z\}P\{N = n\}.$$

Consider the feedback queue again. As indicated clearly by (5), any realization of  $s - r$  follows two steps:  $N$  must be realized first, and a realization of  $\sum_{j=1}^n Y_j$  follows, where  $n$  is a realized value of  $N$ . Moreover, any realized value of  $s - r$  must be a realized value of  $\sum_{j=1}^n Y_j$  for some realized value  $n$  of  $N$ .

From Theorem 1, we recognize that  $s - r$  is a double random element. We use  $N|E$  to represent this double random element, where  $E$  is a collection of Erlang random variables with parameters  $\mu$  and  $n$ , denoted by  $n|E, n = 1, 2, \dots$ .

**Theorem 2** The state of Jackson’s feedback queue, i.e.,  $Q(t)$ , is not a Markov process.

**Proof:** Since  $s - r$  is a double random element  $N|E$ , whose components  $n|E, n = 1, 2, \dots$  follow different distributions,  $s - r$  is not a random variable with an unvarying distribution. The result to be proved then follows. □

It can be seen similarly that,  $t - s$ , the inter-departure time in the physical departure process of the first queue in counterexample 2, is also a double random element. Denote by  $Y$  a service time, and  $Z$  the idle time of the server of the first queue. Write

$$I = \begin{cases} 1, & \text{the server is busy at } s \\ 0, & \text{otherwise.} \end{cases}$$

We recognize now that  $t - s$  is a double random element  $I|E$ , which has two components:  $0|E = Z + Y$ , and  $1|E = Y$ .

**Theorem 3** The physical inter-departure times of the first queue cannot form a stationary random sequence. The second queue and hence the network of the two queues in tandem cannot be stable.

**Proof:** The results follow directly from the fact that  $t - s$  is a double random element  $I|E$  with components  $0|E = Z + Y$  and  $1|E = Y$ .

□

The marginal version of the double random element  $N|E$  is an exponential random variable of rate  $\theta(1-\mu)$ . The marginal version of  $I|E$  is also an exponential random variable, with a parameter equal to the Poisson arrival rate of the M/M/1 queue. Since double random elements are not included in the existing literature, the double random elements  $N|E$  and  $I|E$  are confused with their marginal versions. This may explain why Jackson's proof, and all other proofs of Jackson's theorem, are flawed.

We have suggested two experiments for the feedback queue and the two queues in tandem. The reader is encouraged to perform the experiments, at least to perform them as thought experiments, to see firsthand the difference between the double random elements and their marginal versions.

## 5 Discussion and Conclusion

If a Jackson network prohibits customers to visit any queue in the network more than once, we may apply Burke's theorem in a "stage-by-stage" way [1], so as to isolate the queues in the network, and then treat each queue individually. Only after such treatment, the queues, which used to be connected but now stand alone, are M/M/m queues amenable to analysis. If any loop, such as direct or indirect feedback, exists in a Jackson network, the traffic within the network cannot be modeled by Poisson processes. In this case, the "stage-by-stage" method is no longer applicable.

The single-server queue with feedback and the two queues in tandem are basic components of more complicated Jackson networks. As demonstrated by the counterexamples, Jackson's theorem may not hold for a network constructed based on the two components, since the assumptions made by Jackson (and in other proofs of Jackson's theorem), such as the state of the network can be modeled as a stationary Markov process, may not be valid. Consequently, it may be impossible to describe the actual process of the network with Jackson's theorem.

Networks of queues, including those allowing customers to visit the same queue more than once, are widely used in networking studies, e.g., cf. [2, 6]. However, Jackson's theorem, which is the basis of the mathematical theory of networks of queues, does not hold in general. A reinvestigation of all related issues may be necessary.

For example, it is typical to assume that traffic processes in computer networks are (asymptotically) stationary. This assumption is, however, unlikely to

hold. Even if a computer network can be modeled by a Jackson network of queues without any loop, the traffic process within the network is not stationary, and will never approach any stationary process, since the output of a queue *in the network* is not a stationary process.

Consequently, a steady state may not exist for a network considered as a whole. So, it may not make sense to study "steady-state" behavior of the network, either by analysis or by simulation, and the results obtained in "steady-state" may not be meaningful approximation to the real process of the system.

In particular, if a network is studied by simulation, the simulated output from a link, modeled as a queue, is not stationary. Although a downstream queue driven by this non-stationary process does not have a steady state, a "steady-state" version of the queue may be constructed from the simulation. But such version differs essentially from the true behavior of the queue. To see this, suppose that the behavior of a system is described by  $X(n)$ , a function of discrete time  $n = 1, 2, \dots$ . For simplicity, let  $X(n) = 0$  if  $n = 1, 3, 5, \dots$ , and  $X(n) = 1$  otherwise. By doing a "steady-state" simulation of  $X(n)$ , a "steady-state" version  $X$  can be constructed, such that  $X$  may be either 0 or 1, with equal probability, for any  $n$ . Evidently, this "steady-state" version fails to capture the true behavior of  $X(n)$ .

## References

- [1] P. J. Burke, "The output of a queueing system," *Operations Research*, vol. 4, pp. 699-714, 1956.
- [2] S. B. Fredj, T. Bonald, A. Proutiere, G. Régnié, J. W. Roberts, "Statistical bandwidth sharing: a study of congestion at flow level," *Sigcomm'01*, August 27-31, 2001, San Diego, California, USA.
- [3] J. R. Jackson, "Networks of waiting lines," *Operations Research*, vol. 5, no. 4, 518-521, 1957.
- [4] J. R. Jackson, "How networks of queues came about," *Operations Research*, vol. 50, no. 1, 112-113, 2002.
- [5] L. Kleinrock, "Creating a mathematical theory of computer networks," *Operations Research*, vol. 50, no. 1, 125-131, 2002.
- [6] M. Schwartz, **Telecommunication Networks: Protocols, Modeling and Analysis**, Addison-Wesley, 1987.