



On the $M/M/c/N$ Call Center Queue Modeling and Analysis

By Nwonye Chukwunoso

University of Port Harcourt, Nigeria

Abstract – The $M/M/c/c$ model is the most widely applied queueing model in the mathematical analysis of call centers. The $M/M/c/c$ model is also referred to as the Erlang Loss System. The Erlang loss model does not take into consideration system attributes such as blocking and busy signals, balking and reneging, retrials and returns. Although, the Erlang loss model is analytically tractable, it is not easy to obtain insight from its results.

The need to develop a more accurate call center model has necessitated the modification of the Erlang loss model. In this research, we model and analyze a call center using $M/M/c/N$ the model. The goal of this paper is to extend existing results and prove new results with regards to the monotonicity and limiting behaviour of the $M/M/c/N$ model with respect to the system capacity N .

GJCST-C Classification : D.2.0



Strictly as per the compliance and regulations of:



On the $M/M/c/N$ Call Center Queue Modeling and Analysis

Nwonye Chukwunoso

Abstract - The $M/M/c/c$ model is the most widely applied queueing model in the mathematical analysis of call centers. The $M/M/c/c$ model is also referred to as the Erlang Loss System. The Erlang loss model does not take into consideration system attributes such as blocking and busy signals, balking and reneging, retrials and returns. Although, the Erlang loss model is analytically tractable, it is not easy to obtain insight from its results.

The need to develop a more accurate call center model has necessitated the modification of the Erlang loss model. In this research, we model and analyze a call center using the $M/M/c/N$ model. The goal of this paper is to extend existing results and prove new results with regards to the monotonicity and limiting behaviour of the $M/M/c/N$ model with respect to the system capacity N .

1. INTRODUCTION

The call center industry has grown explosively in the recent past and that has aroused the interest of researchers from different disciplines. Mandelbaum [11] have provided a comprehensive research bibliography with abstracts in diverse disciplines such as Operations research, Statistics, Engineering, and so on. Call Center research has been reviewed in the tutorial and survey paper by Gans et al. [6]. In this paper, our focus is on the computational rigor of the call center performance metrics using the $M/M/c/N$ model.

a) Description of a Call Center

A call center is a department of an establishment that attends to customers via telephone conversation often for the purpose of sales and product support, or that makes outgoing telephone calls to customers usually for the purpose of advertisement or telemarketing. Suppose the department also attends to e-mails, faxes, letters, and other similar written correspondence, then, it is called a contact center.

Inbound call center only handle incoming telephone calls initiated by customers while out bound call centers only make outgoing telephone calls to customers. There are call centers that deal with both types of calls. In majority of the call centers, inbound calls form the bulk of contacts with customers. In

addition, inbound calls are more time consuming compared to other types of contacting options (e.g. e-mails, faxes, or letters) in terms of waiting times in the telequeue or sojourn times. Hence, we will only focus on inbound call centers. In an inbound call center, there is a group of agents (Customer Sales Representatives, CSRs) who provide the needed service through talking to customers on phones. In this paper, we shall use the terms "agents" and "CSRs" interchangeably. Agents are equipped with equipment, such as a Private Automatic Branch Exchange (PABX or PBX), an Interactive Voice Response Unit (IVRU or VRU), an Automatic Call Distributor (ACD), and computers [16]. See Figure 1.1 for details on the operational process and components of an inbound call center.

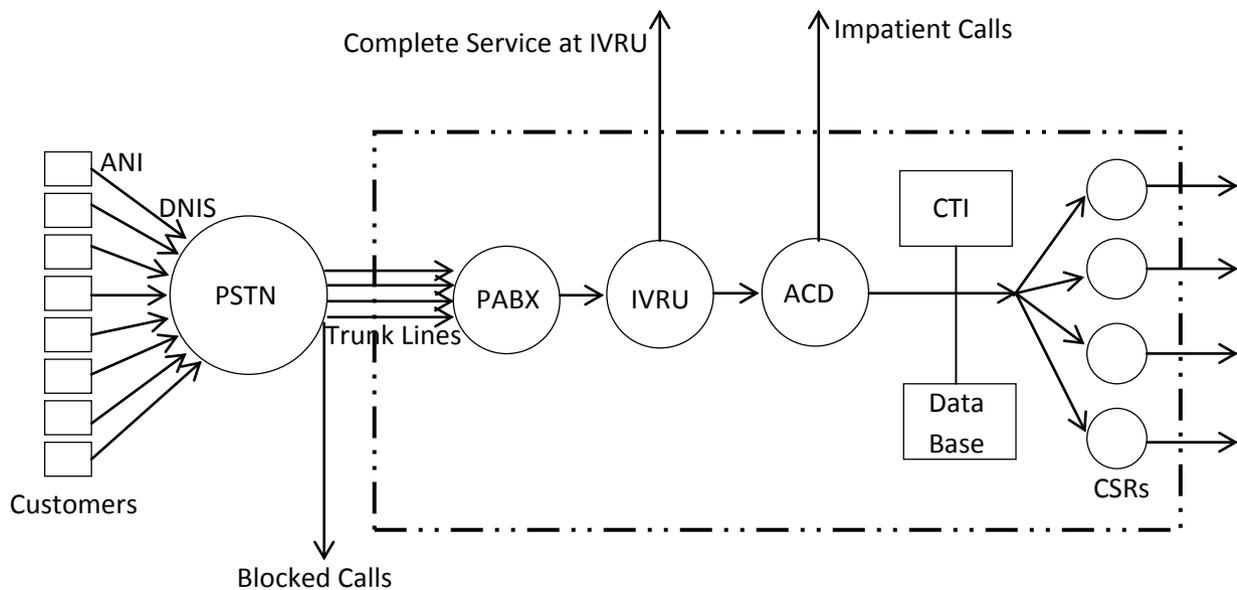


Figure 1.1 : Operational process of an inbound call center

b) *The Operational Process of an Inbound Call Center*

At some point in our lives, we have all called a call center. We will describe the operational process and components of an inbound call center in line with the description in [6, 17]. The process is depicted in Figure 1.1. Customers wanting to receive service from a call center, dial a special number provided by the call center. The Public Service Telephone Network (PSTN) company then uses the Automatic Number Identification (ANI) number (the phone number from which the customer dials) and the customer's Dialed Number Identification Service (DNIS) number (the special number being dialed) to connect the customer to the PABX privately-owned by the call center. The telephone lines (usually called trunk lines) connect the PABX to PSTN. If a trunk line is available, the customer seizes it; else the customer receives a busy signal and will be rejected. Hence, this customer is said to be blocked. Once the call is accepted, the customer will be connected through the PABX to the IVRU. The IVRU provides some automatic service for customers as well as several options for customers to choose from. Upon service completion at the IVRU, some customers leave the system and release the trunk lines. If the customer requires the service of an agent, the call will be passed from the IVRU to the ACD. The ACD is a sophisticated instrument designed to route calls to agents based on the specific needs of calls. If no appropriate CSRs are available, the customer is informed to wait and join a queue at the ACD. The customer is said to be delayed. The ACD decides the next customer to get service according to some preprogrammed queueing discipline (usually First Come First Served, FCFS). Delayed customers may decide to hang up and abandon (or renege) before they are served if they perceive that the

service is not worth the wait. Such customers are said to be impatient. Patient customers (who do not abandon service) will eventually be connected to an agent. In serving a customer, the CSR works with a PC furnished with Computer-Telephony Integration (CTI), which is the technology that allows interactions on a telephone and a computer to be integrated. CTI will help ACD to route the call, help the CSR to get the caller's information from the database and hence facilitate the service process. At the completion of service and exit of the customer, the CSR still needs some wrap-up time to finish the whole service process and then may be available for the next customer. The service time is the sum of talk time and wrap-up time. Customers who abandoned and were blocked may try to call again after some random amount of time and these calls are referred to as retrials. Customers who finished talking with an agent may also need further assistance and therefore call back. Hence they become return customers or feedback customers. Notice that these two types of customers are not shown in Figure 1.1.

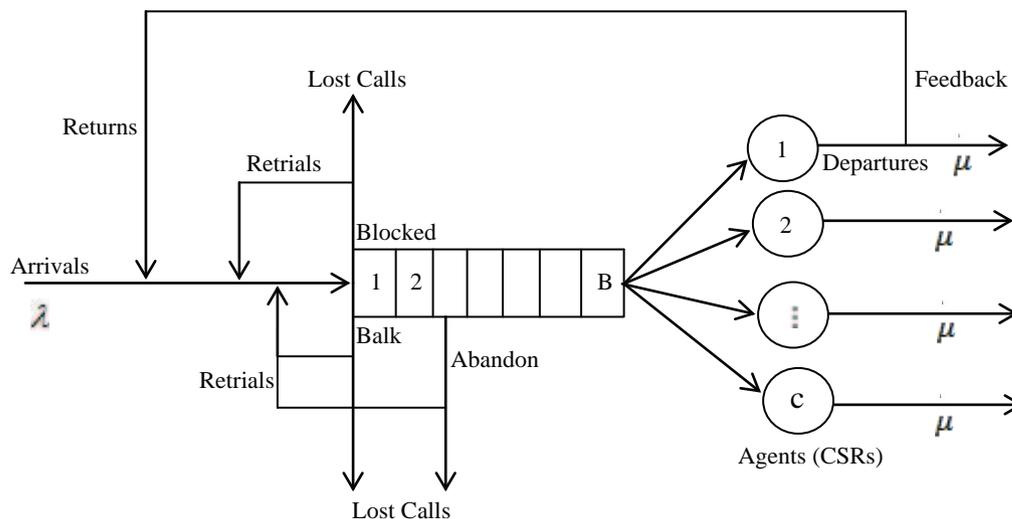


Figure 1.2: Call Center as a Queueing System

c) *The Call Center as a Queueing System*

Figure 1.2 depicts the call center as a queueing system [7]. The number of agents (CSRs) and waiting spaces are denoted respectively by c and B . Hence there are $N = B + c$ trunk lines at the call center with B waiting spaces. If an arriving call finds all N trunk lines occupied, it gets a busy signal and as such is blocked and cannot access the system. If there is an available trunk line, the call is either connected to the system and seizes one of the free trunk lines or it balks. Suppose there is an available trunk line and at least a free agent, then the call is immediately serviced. Otherwise the call experiences delay and has to wait in a queue at the ACD for a CSR to become available. Calls at the ACD may become impatient and abandon (renege) the system before being served and thus release the trunk line. The ACD usually implements the FCFS queueing discipline. Upon service completion by a CSR, the call leaves the system and then releases both the trunk line and the CSR and these resources become available to other arriving calls. Return (or feedback) calls are calls that return after been served by an agent. Some of those calls who do not get served (blocked, abandon or balk) may call again and they become retrials. The remaining calls become lost calls.

Suppose that the call arrivals follow a Poisson process with mean rate λ and that the service times of the calls are independent and identically distributed (*i. i. d*) exponential random variables with mean $1/\mu$. Then we can model the system as a $M/M/c/N$ queueing system with features such as balking, abandonment, retrial, and feedback.

d) *Performance Evaluation of the Call Center Queueing Model*

In this paper, we will ignore features such as balking, abandonment, retrial, and feedback. Following the above assumptions, we will apply the $M/M/c/N$ model in analyzing the call center performance. The $M/M/c/N$ queueing system has a closed-form solution for the system state (number of calls in the system), the queue length (number of calls in the queue) distribution and waiting time distribution. Then we can obtain system performance metrics such as average waiting time, average queue length, and probability of blocking. We will apply the performance analysis of the $M/M/c/N$ queueing system to call center modeling and in turn show new results. The call center performance measures (metrics or indicators) provide useful information in the design and management of call centers. Performance measures are used in determining the service levels (or quality of service) in call centers.

Not all queueing models can be analyzed exactly to obtain performance measures as $M/M/c/N$ model. For instance, if we include additional features such as Non-Poisson time varying arrival process, balking, abandonment, retrial, feedback, and non-exponential service times, the model may become insolvable using traditional queueing techniques and other techniques have to be used to analyze the model such as simulation modeling.

II. MODELING CALL CENTERS AS SINGLE-NODE EXPONENTIAL QUEUEING MODELS

In this section, we provide a detailed review of relevant single-node multiserver Markovian queueing models of call centers. Table 2.1 provides a list some main Markovian queueing models and their

performance indicators. Our emphasis is on the computational rigor of the exact performance measures

of these well-known models as well deriving new results.

Table 2.1 : Some Multiserver Markovian Queueing Models

Notation	Models	Performance Indicators
$M/M/c$	Delay Model (Erlang C)	$P(\text{delay}) = P(W_q > 0)$ (Erlang C formula) TSF; ASA
$M/M/c/c$	Blocking/Loss Model (Erlang B)	$P(\text{blocking}) = p_c$ (Erlang B formula)
$M/M/c/N$	Blocking and Delay Model	$P(\text{delay}) = P(W_q > 0); P(\text{blocking}) = p_N$ TSF; ASA

TSF : Telephone Service Factor ; ASA: Average Speed to Answer; AWT: Acceptable Waiting Time

a) Model Assumptions

The $M/M/c/N$ is a generalization of the $M/M/c$ and $M/M/c/c$ models. In order to analyze a call center using the $M/M/c/N$ Markovian queueing model, we assume that the inter arrival and service times are exponentially distributed random variables.

Calls arriving at the call center are of a single type following a homogeneous Poisson process with rate λ . Callers are assumed to be patient and there is no form of impatience (balking or reneging). All agents (CSRs) are assumed to be statistically identical (i.e., equally skilled and provide service at the same rate). The service times are assumed to follow the exponential distribution with mean $1/\mu$. Services are rendered according to the First-Come-First-Serve queueing discipline. There are $N - c$ waiting spaces in the $M/M/c/N$ queueing system.

Let $L(t)$ be the number of calls in the system and $L_q(t)$ be the number of calls waiting in queue at time t . Let $W(t)$ and $W_q(t)$ be the steady-state sojourn time and waiting time in queue respectively. Since the models are Markovian, $L(t), L_q(t), W(t), \text{ and } W_q(t)$ can be obtained using the birth-death processes. Our focus is on steady-state

distribution of $L(t), L_q(t), W(t), \text{ and } W_q(t)$ with corresponding variables $L, L_q, W, \text{ and } W_q$, respectively.

Let $p_n = P(L = n), n = 0, 1, 2, \dots$ denote the steady-state probability (if it exists) of the system being in state n (i.e. having n calls in the system). Applying the modeling techniques of the birth-death processes, we can obtain some interesting system performance measures such as $E(L), E(L_q), E(W), \text{ and } E(W_q)$.

Due to the PASTA property, we have for the $M/M/c$ model, $P(\text{waiting}) = \sum_{n=c}^{\infty} p_n$ and in the cases of the $M/M/c/N$ and $M/M/c/c$, we have $P(\text{blocking}) = p_N$ and $P(\text{blocking}) = p_c$ respectively.

b) Review of the $M/M/c/c$ Model and the Erlang B Formula

In this section of the paper we will review the $M/M/c/c$ Erlang B model paying attention to the aspects that are relevant to call center modeling. The $M/M/c/c$ queue models a single-node system with c truck lines and no waiting spaces. Figure 2.1 depicts the $M/M/c/c$ queue and figure 3.2, its state transition.

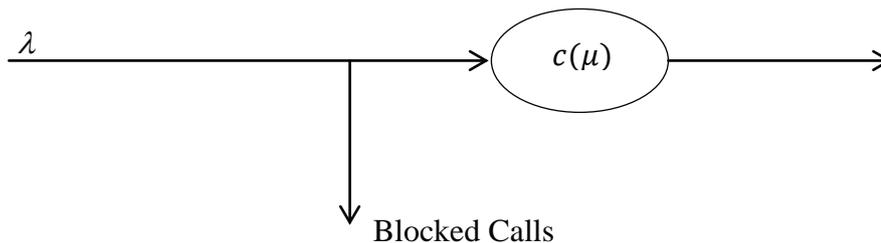


Figure 2.1 : Description of the $M/M/c/c$ Model and its parameters

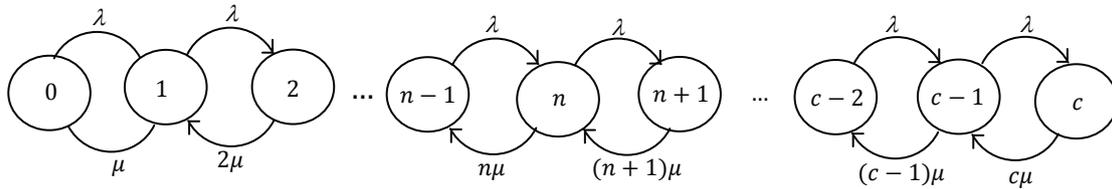


Figure 2.2 : $M/M/c/c$ Flow Rate [Multiple-Server Case($c \geq 2$)]

Considering figures 2.1 and 2.2, it is obvious that $L(t)$ is a finite birth-death process with birth rate

$$\lambda_n = \begin{cases} \lambda, & \text{if } n < c \\ 0, & \text{if } n \geq c \end{cases}$$

and state-dependent death rate $\mu_n = n\mu, n = 0, 1, 2, \dots, c$.

By the application of the fundamental equation in queueing theory, the steady-state solution of the $M/M/c/c$ model using birth-death process is given by

$$p_n = p_0 \prod_{i=1}^n \frac{\lambda_{i-1}}{\mu_i}$$

where p_0 is computed from $\sum_{n=0}^N p_n = 1$. The solution is given by

$$p_n = \frac{(\lambda/\mu)^n/n!}{\sum_{i=0}^c (\lambda/\mu)^i/i!} = \frac{a^n/n!}{\sum_{i=0}^c a^i/i!}, \quad 0 \leq n \leq c$$

where $a := \lambda/\mu$ is called the offered load which describes the demand made on the system.

Let

- $L(t)$ = Number of customers in the system at time t
- $p_n(t) = P(\text{System is in state } n \text{ at time } t) = P(L(t) = n)$
- $a_n(t) = P(\text{Arrival at time } t \text{ finds system in state } n)$
- $A(t, t + \delta t]$ = The event of an arrival in $(t, t + \delta t]$

Then

$$\begin{aligned} a_n(t) &= \lim_{\delta t \rightarrow 0} P(L(t) = n | A(t, t + \delta t]) \\ &= \lim_{\delta t \rightarrow 0} \frac{P(L(t) = n \text{ and } A(t, t + \delta t])}{P(A(t, t + \delta t])} \\ &= \lim_{\delta t \rightarrow 0} \frac{P(A(t, t + \delta t] | L(t) = n) P(L(t) = n)}{P(A(t, t + \delta t])} \\ &= \lim_{\delta t \rightarrow 0} \frac{P(A(t, t + \delta t]) P(L(t) = n)}{P(A(t, t + \delta t])} \\ &= P(L(t) = n) = p_n(t) \end{aligned}$$

i. *PASTA: Poisson Arrivals See Time Averages*

An important feature of the Markovian queueing models is that the arrival process follows a Poisson process. Considering the Poisson arrival process, the distribution of customers seen by an arrival to a queueing facility is, stochastically the same as the limiting distribution of customers at that facility. In other words, once the queueing system has reached steady state, each arrival from a Poisson process finds the system at equilibrium. If p_n is the probability that the system contains n customers at equilibrium and a_n denotes the probability that an arriving customer finds n customers already present, then PASTA states that $a_n = p_n$. This implies that the Poisson process sees the same distribution as a random observer, i.e., at equilibrium, Poisson arrivals take a random look at the system. This result is a direct consequence of the memoryless property of the interarrival time distribution of customers to a queueing system fed by a Poisson process. In particular, it does not depend on the service time distribution. To prove the PASTA property, we proceed as follows.

The crucial step in the above argument is

$$P(A(t, t + \delta t) | L(t) = n) = P(A(t, t + \delta t))$$

This results from the fact that, since interarrival times possess the memoryless property, $P(A(t, t + \delta t))$ is independent of the past history of the arrival process and hence independent of the current state of the queueing system. With the Poisson arrival process having a constant rate λ , the probability of having an arrival in $(t, t + \delta t]$ is equal to

$$p_c = a_c = P(\text{blocking}) = B(c, \lambda/\mu) = B(c, a) = \frac{(\lambda/\mu)^n/n!}{\sum_{n=0}^c (\lambda/\mu)^n/n!} = \frac{a^c/c!}{\sum_{n=0}^c a^n/n!} \quad (2.1)$$

Notice that the probability that an arrival is lost is equal to the probability that all channels are busy. Erlang loss formula is also valid for the $M/G/c/c$ queue. In other words, the steady-state probabilities are a function only of the mean service time, and not of the complete underlying cumulative distribution function. An efficient recursive algorithm for computing $B(c, a)$ is given by

$$B(0, a) = 1, \quad B(c, a) = \frac{aB(c-1, a)}{c + aB(c-1, a)} \quad (2.2)$$

Recall that $a = \lambda/\mu$ is the offered load, we define $a' := E(L_b) = a[1 - B(c, a)]$ as the carried load, where we obtain the last equality by Little's law applying to number of busy servers and L_b is a random

$$B(c, a) = \frac{aB(c-1, a)}{c + aB(c-1, a)} < \frac{aB(c-1, a)}{a[1 - B(c-1, a)] + aB(c-1, a)} = B(c-1, a) \quad (2.4)$$

since

$$a[1 - B(c-1, a)] < c - 1 < c$$

We do not consider performance measures relating to waiting time and queue length since there is no waiting space in the $M/M/c/c$ model.

c) Review of the $M/M/c$ Model and the Erlang C Formula

The $M/M/c$ queue can be used to model multiprocessor systems or devices that have several identical servers (or agents) and all jobs (or calls) waiting for these servers are kept in one queue. It is assumed that there are c agents each with a service rate of μ jobs per unit time. The arrival rate is λ calls per unit time. If any of the c agents are idle, the arriving call is serviced immediately. If all c agents are busy, the arriving calls wait in a queue. The state of the system is

$\lambda\delta t + o(\delta t)$ which does not depend on $L(t)$. Note that the PASTA property only holds for Poisson arrival processes.

The formula for p_c is called "Erlang Loss Formula" and is the fraction of time that all c servers are busy. It denotes the probability that an arrival call finds all the truck line busy, (i.e. the blocking probability, p_c). It is written as $B(c, \lambda/\mu)$ s and is called "Erlang B formula":

variable representing the number of busy servers in steady-state.

The utilization

$$v := \frac{a'}{c} = \frac{a[1 - B(c, a)]}{c} = \rho[1 - B(c, a)] < 1$$

is the fraction of time that a server is busy, where $\rho := \frac{a}{c}$ is called the traffic intensity.

Hence we have that

$$1 - \frac{1}{\rho} < B(c, a) \quad (2.3)$$

which defines a lower bound for $B(c, a)$

Next, we show the monotonicity property of the $B(c, a)$ with respect to c .

represented by the number of calls n in the system. The state transition diagram is shown in figure 2.4. It is easy to see that the number of jobs in the system is a birth-death process with the following correspondence:

Using the concept of total probability, we have that

$$P(W_q > t) = \sum_{n=c}^{\infty} P(W_q > t \mid \text{system in state } n \text{ upon arrival})P(\text{system in state } n \text{ upon arrival})$$

$$= \sum_{n=c}^{\infty} P(\text{service completion time of } n - c + 1 \text{ calls} > t)a_n$$

Using the PASTA property, we can write $P(\text{system in state } n \text{ upon arrival}) = p_n = a_n$ which is the steady-state probability of an arriving call meeting n calls in the system. Since the service times are exponentially distributed and $n \geq c$, the completion

time of $n - c + 1$ calls (denoted by) Y_i has an Erlang distribution $Er(n - c + 1, c\mu)$ with survival function given by

$$R(t) = P(Y_i > t) = P(\text{service completion time of } n - c + 1 \text{ calls} > t) = \sum_{i=0}^{n-c} \frac{(c\mu t)^i e^{-c\mu t}}{i!}$$

Then we have that

$$P(W_q > t) = \sum_{n=c}^{\infty} \sum_{i=0}^{n-c} \frac{(c\mu t)^i e^{-c\mu t}}{i!} p_n = \frac{\alpha^c}{c!(1-\rho)} p_0 e^{-(c\mu-\lambda)t} \tag{2.10}$$

$$E(W_q) = \int_0^{\infty} P(W_q > t) dt = \frac{C(c, \alpha)}{c\mu - \lambda} \tag{2.11}$$

Note that $P(W_q = 0) = 1 - C(c, \alpha)$. By the application of Little's law, we have

$$E(L_q) = \lambda E(W_q) = \frac{\rho C(c, \alpha)}{1 - \rho}$$

Because of the closed-form solutions of most the performance indicators of the $M/M/c$ model, it is commonly used in performance modeling and analysis of call centers. In the application of $M/M/c$ model in call center analysis, it is usually assumed that the arrival and service rate are piece-wise constant and time-independent. Using the parameters of each interval, the $M/M/c$ is applied to each time interval. The $M/M/c$ model is not a realistic tool for modeling call centers due to the following reasons:

- It assumes there is no blocking since it has infinite buffer capacity.

- It does not consider the impatience (balking and reneging) attributes of customers.

d) *Review of the $M/M/c/N$ Model*

When the waiting room in a queueing system has a capacity limit we get a finite queue. In most situations, a finite queue occurs more naturally than a queue with a waiting room of infinite size. However, as the capacity limit gets larger, the behavior of the system approximates that of an infinite-capacity system, and in such cases we are justified in ignoring the size limit. A call center with a finite buffer and several agents is a good example of a finite queueing system. In this section we will review the $M/M/c/N$ model and prove new monotonicity properties of performance measures with respect to N .

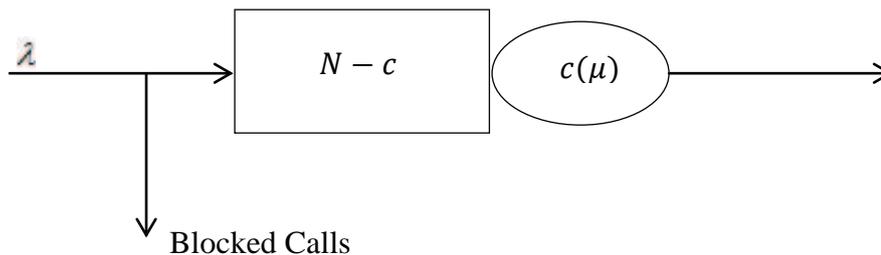


Figure 2.5: Description of the $M/M/c/N$ Model and its parameters

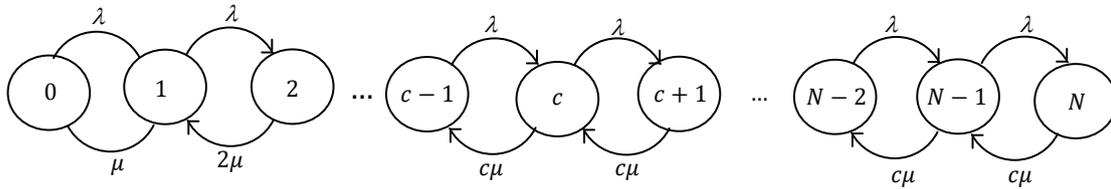


Figure 2.6 : $M/M/c/N$ Flow Rate [Multiple-Server Case ($c > 1$)]

The $M/M/c/N$ queue is similar to the $M/M/c$ queue except that the number of buffers is finite. After $B = N - c$ buffers are full, all arrivals are lost. We assume that B is greater than or equal to c ; otherwise, some servers will never be able to operate due to a lack of buffers and the system will effectively operate as a $M/M/B/B$ queue.

The state transition diagram for a $M/M/c/N$ queue is shown in Figure 2.6. The system can be modeled as a birth-death process using the following respective arrival and service rates:

$$\lambda_n = \lambda, \quad \text{if } 0 \leq n \leq N - 1$$

$$\mu_n = \begin{cases} n\mu, & \text{if } 0 \leq n \leq c \\ c\mu, & \text{if } c \leq n \leq N \end{cases}$$

Solving the balance equations derived from the state diagram, we obtain the following state probabilities.

$$P(\text{waiting}) = \sum_{n=c}^{N-1} p_n; \quad P(\text{blocking}) = p_N; \quad \text{and} \quad P(\text{no - waiting}) = 1 - P(\text{waiting}) - P(\text{blocking})$$

i. *The $M/M/c/N$ Waiting Time Distribution*

In this section, we shall provide a mathematical derivation of the waiting time distribution of the $M/M/c/N$ model. Due to the finiteness of the capacity of the $M/M/c/N$ system, deriving the waiting time distribution of the $M/M/c/N$ model is complicated because it results to finite series and also the arrival

$$p_n = \begin{cases} \frac{a^n}{n!} p_0, & \text{if } 0 \leq n \leq c \\ \frac{a^n}{c! c^{n-c}} p_0, & \text{if } c \leq n \leq N \end{cases} \quad (2.12)$$

with

$$p_0 = \begin{cases} \left(\sum_{n=0}^{c-1} \frac{a^n}{n!} + \frac{a^c}{c!} \frac{1 - \rho^{N-c+1}}{1 - \rho} \right)^{-1}, & \text{if } \rho \neq 1 \\ \left(\sum_{n=0}^{c-1} \frac{a^n}{n!} + \frac{a^c}{c!} (N - c + 1) \right)^{-1}, & \text{if } \rho = 1 \end{cases}$$

and

process is truncated by the system size N . The arrival process no longer follows the Poisson process and has necessitated the need to derive the arrival point probabilities, q_n since $p_n \neq q_n$. In this derivation of q_n , we shall apply the well-known Bayes' theorem.

$$q_n = P(\text{system is in state } n \mid \text{an arrival in } (t, t + \delta t])$$

$$= P(L(t) = n \mid A(t, t + \delta t]) = \frac{P(L(t) = n; A(t, t + \delta t])}{P(A(t, t + \delta t])}$$

$$= \frac{P(A(t, t + \delta t] \mid L(t) = n)P(L(t) = n)}{\sum_{n=0}^N P(A(t, t + \delta t] \mid L(t) = n)P(L(t) = n)} = \frac{P(A(t, t + \delta t] \mid L(t) = n)p_n}{\sum_{n=0}^N P(A(t, t + \delta t] \mid L(t) = n)p_n}$$

Taking limits of both sides and using the fact that the probability of an arrival in $(t, t + \delta t]$ is $\lambda \delta t + o(\delta t)$ we have that

$$\begin{aligned} \lim_{\delta t \rightarrow 0} q_n &= q_n = \frac{P(A(t, t + \delta t))p_n}{\sum_{n=0}^N P(A(t, t + \delta t))p_n} = \lim_{\delta t \rightarrow 0} \left(\frac{(\lambda \delta t + o(\delta t))p_n}{\sum_{n=0}^{N-1} (\lambda \delta t + o(\delta t))p_n} \right) \\ &= \lim_{\delta t \rightarrow 0} \left(\frac{\left(\lambda + \frac{o(\delta t)}{\delta t}\right)p_n}{\sum_{n=0}^{N-1} \left(\lambda + \frac{o(\delta t)}{\delta t}\right)p_n} \right) = \frac{\lambda p_n}{1 - p_N}, \quad \text{for } 0 \leq n \leq N - 1 \end{aligned}$$

which defines the probability of a call meeting n calls in the system upon arrival given that it is not blocked. Here we have used the fact that

$$p_n = \begin{cases} \frac{c!}{n! a^{c-n}} p_c, & \text{if } 0 \leq n \leq c \\ p_c \rho^{n-c}, & \text{if } c \leq n \leq N \end{cases}$$

$$\lim_{\delta t \rightarrow 0} \frac{o(\delta t)}{\delta t} = 0$$

32

Using equation (2.12), we can write $p_c = \frac{a^c}{c!} p_0$ which implies that $p_0 = \frac{c!}{a^c} p_c$. Then we can express p_n in terms of p_c as follows:

Then for $\rho \neq 1$,

$$\begin{aligned} 1 &= \sum_{n=0}^c \frac{c!}{n! a^{c-n}} p_c + \sum_{n=c+1}^N p_c \rho^{n-c} \\ p_c &= \left(\sum_{n=0}^c \frac{c!}{n! a^{c-n}} + \sum_{n=c+1}^N \rho^{n-c} \right)^{-1} = \left(\frac{1}{B(c, a)} + \frac{\rho(1 - \rho^{N-c})}{1 - \rho} \right)^{-1} \\ p_c &= \frac{(1 - \rho)B(c, a)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})} \end{aligned}$$

For $\rho \neq 1$,

$$P(\text{blocking}) = p_N = p_c \rho^{N-c} = \frac{(1 - \rho)B(c, a)\rho^{N-c}}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})}$$

In same way, for $\rho \neq 1$,

$$P(\text{waiting}) = \sum_{n=c}^{N-1} p_n = \sum_{n=c}^{N-1} p_c \rho^{n-c} = \frac{1 - \rho^{N-c}}{1 - \rho} p_c = \frac{1 - \rho^{N-c}}{1 - \rho} \frac{(1 - \rho)B(c, a)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})}$$

$$P(\text{waiting}) = \frac{(1 - \rho^{N-c})B(c, a)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})}$$

and

$$P(\text{no - waiting}) = 1 - P(\text{waiting}) - P(\text{blocking}) = \frac{(1 - B(c, a))(1 - \rho)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})}$$

Now, let us consider computing $P(\text{blocking})$, $P(\text{waiting})$ and $P(\text{no - waiting})$ in the case where $\rho = 1$.

For $\rho = 1$, implies that $a = c$ so that we have

$$p_c = \left(\sum_{n=0}^c \frac{c!}{n! a^{c-n}} + \sum_{n=c+1}^N \rho^{n-c} \right)^{-1} = \left(\frac{1}{B(c, a)} + N - c \right)^{-1} = \frac{B(c, c)}{1 + (N - c)B(c, c)}$$

$$P(\text{blocking}) = p_N = p_c \rho^{N-c} = p_c = \frac{B(c, c)}{1 + (N - c)B(c, c)}$$

$$P(\text{waiting}) = \sum_{n=c}^{N-1} p_n = \sum_{n=c}^{N-1} p_c \rho^{n-c} = \sum_{n=c}^{N-1} p_c = (N - c)p_c = \frac{(N - c)B(c, c)}{1 + (N - c)B(c, c)}$$

$$P(\text{no - waiting}) = 1 - P(\text{waiting}) - P(\text{blocking}) = \frac{1 - B(c, c)}{1 + (N - c)B(c, c)}$$

Theorem 2.1

Suppose $N = c$ then the $M/M/c/N$ model reduces to the $M/M/c/c$ model with

$$P(\text{blocking}) = B(c, a), P(\text{waiting}) = 0 \text{ and } P(\text{no - waiting}) = 1 - B(c, a).$$

Proof

If $N = c$,

$$P(\text{blocking}) = p_c = p_N = p_c \rho^{N-c} = \frac{B(c, a) \rho^{N-c}}{1 + (N - c)B(c, a)} = \frac{B(c, a)(1)}{1 + (0)B(c, a)} = B(c, a)$$

$$P(\text{waiting}) = \frac{(1 - \rho^{N-c})B(c, a)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})} = \frac{(0)B(c, a)}{1 - \rho + \rho B(c, a)(0)} = 0$$

$$P(\text{no - waiting}) = 1 - P(\text{waiting}) - P(\text{blocking}) = 1 - 0 - B(c, a) = 1 - B(c, a)$$

Theorem 2.2

In the limit, as $N \rightarrow \infty$, we have the following results:

$$1. \lim_{N \rightarrow \infty} P(\text{waiting}) = \begin{cases} \frac{1}{\rho}, & \text{if } \rho > 1 \\ 1, & \text{if } \rho = 1 \\ C(c, a), & \text{if } 0 < \rho < 1 \end{cases}$$

$$2. \lim_{N \rightarrow \infty} P(\text{blocking}) = \begin{cases} 1 - \frac{1}{\rho}, & \text{if } \rho \geq 1 \\ 0, & \text{if } 0 < \rho \leq 1 \end{cases}$$

$$3. \lim_{N \rightarrow \infty} P(\text{no - waiting}) = \begin{cases} 0, & \text{if } \rho \geq 1 \\ 1 - C(c, a), & \text{if } 0 < \rho \leq 1 \end{cases}$$

Proof

1. For $0 < \rho < 1$,

$$P(\text{waiting}) = \frac{(1 - \rho^{N-c})B(c, a)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})} \xrightarrow{\text{as } N \rightarrow \infty} \frac{B(c, a)}{1 - \rho + \rho B(c, a)} = C(c, a)$$

For $\rho = 1$,

$$P(\text{waiting}) = \frac{(N - c)B(c, c)}{1 + (N - c)B(c, c)} = \frac{B(c, c)}{\frac{1}{(N-c)} + B(c, c)} \xrightarrow{\text{as } N \rightarrow \infty} 1$$

For $\rho > 1$,

$$P(\text{waiting}) = \frac{(\rho^{N-c} - 1)B(c, a)}{\rho + \rho B(c, a)(\rho^{N-c} - 1) - 1} = \frac{\frac{(\rho^{N-c} - 1)B(c, a)}{\rho^{N-c} - 1}}{\frac{\rho - 1}{\rho^{N-c} - 1} + \frac{\rho B(c, a)(\rho^{N-c} - 1)}{\rho^{N-c} - 1}} = \frac{1}{\rho}$$

So that

$$\lim_{N \rightarrow \infty} P(\text{waiting}) = \frac{1}{\rho}, \quad \text{if } \rho > 1$$

2. For $0 < \rho < 1$,

$$P(\text{blocking}) = p_N = p_c \rho^{N-c} = \frac{(1 - \rho)B(c, a)\rho^{N-c}}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})} \xrightarrow{\text{as } N \rightarrow \infty} 0$$

Since $\rho^{N-c} \rightarrow 0$ as $N \rightarrow \infty$, for $0 < \rho < 1$.

For $\rho = 1$,

$$P(\text{blocking}) = p_N = \frac{B(c, c)}{1 + (N - c)B(c, c)} \xrightarrow{\text{as } N \rightarrow \infty} 0$$

For $\rho > 1$,

$$\begin{aligned} P(\text{blocking}) &= p_N = \frac{(\rho - 1)B(c, a)\rho^{N-c}}{\rho + \rho B(c, a)(\rho^{N-c} - 1) - 1} = \frac{\frac{(\rho - 1)B(c, a)\rho^{N-c}}{\rho^{N-c}}}{\frac{\rho - 1}{\rho^{N-c}} + \frac{\rho B(c, a)(\rho^{N-c} - 1)}{\rho^{N-c}}} \\ &= \frac{(\rho - 1)B(c, a)}{\frac{\rho - 1}{\rho^{N-c}} + \frac{\rho B(c, a)(\rho^{N-c} - 1)}{\rho^{N-c}}} \xrightarrow{\text{as } N \rightarrow \infty} \frac{(\rho - 1)B(c, a)}{\rho B(c, a)} = 1 - \frac{1}{\rho} \end{aligned}$$

3. For $\rho > 1$,

$$\begin{aligned} P(\text{no - waiting}) &= \frac{(1 - B(c, a))(\rho - 1)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})} = \frac{\frac{(1 - B(c, a))(\rho - 1)}{\rho^{N-c}}}{\frac{\rho - 1}{\rho^{N-c}} + \frac{\rho B(c, a)(\rho^{N-c} - 1)}{\rho^{N-c}}} \\ &= \frac{\frac{(1 - B(c, a))(\rho - 1)}{\rho^{N-c}}}{\frac{\rho - 1}{\rho^{N-c}} - \frac{\rho B(c, a)}{\rho^{N-c}} + \frac{\rho B(c, a)\rho^{N-c}}{\rho^{N-c}}} \xrightarrow{\text{as } N \rightarrow \infty} 0 \end{aligned}$$

For $0 < \rho < 1$,

$$\begin{aligned} (no - waiting) &= \frac{(1 - B(c, a))(1 - \rho)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})} \xrightarrow{as\ N \rightarrow \infty} \frac{(1 - B(c, a))(1 - \rho)}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})} \\ &= 1 - C(c, a) \end{aligned}$$

For $\rho = 1$,

$$P(no - waiting) = \frac{1 - B(c, c)}{1 + (N - c)B(c, c)} \xrightarrow{as\ N \rightarrow \infty} 0$$

Before we proceed to derive the formula for computing an important performance measure $E(W_q)$, we shall prove some new results that will be useful in the course of our derivations and computations.

For $\rho \neq 1$,

$$\begin{aligned} P(waiting|no - blocking) &= P(W_q > 0) = \sum_{n=c}^{N-1} q_n = \sum_{n=c}^{N-1} \frac{p_n}{1 - p_N} = \frac{1}{1 - p_N} \sum_{n=c}^{N-1} p_n \\ &= \frac{1}{1 - p_N} P(waiting) = \frac{P(waiting)}{1 - P(blocking)} \\ &= \frac{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c-1})} \frac{B(c, a)(1 - \rho^{N-c})}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c})} \end{aligned}$$

$$P(waiting|no - blocking) = P(W_q > 0) = \frac{B(c, a)(1 - \rho^{N-c})}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c-1})}$$

But

$$\begin{aligned} P(no - waiting|no - blocking) &= P(W_q = 0) = \sum_{n=0}^{c-1} q_n = 1 - P(waiting|no - blocking) \\ &= \frac{(1 - \rho)[1 - B(c, a)]}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c-1})} \end{aligned}$$

For $\rho = 1$,

$$\begin{aligned} P(waiting|no - blocking) &= \frac{P(waiting)}{1 - P(blocking)} = \frac{\frac{(N-c)B(c, c)}{1+(N-c)B(c, c)}}{1 - \frac{B(c, c)}{1+(N-c)B(c, c)}} \\ &= \frac{(N - c)B(c, c)}{1 + B(c, c)[N - c - 1]} \end{aligned}$$

Now, using the principles of conditional probability, we can write

$$P(W_q > t | W_q > 0) = \frac{P(W_q > t; W_q > 0)}{P(W_q > 0)} = \frac{P(W_q > t)}{P(W_q > 0)}$$

$$P(W_q > t) = P(W_q > t | W_q > 0)P(W_q > 0)$$

$$P(W_q > t) = P(W_q > t | W_q > 0)P(\text{waiting|no - blocking})$$

For $\rho \neq 1$,

$$P(W_q > t | W_q > 0) = P(W_q > t | c \leq Q \leq N - 1)$$

$$= \sum_{n=0}^{N-c+1} P(W_q > t | L = c + n; c \leq L \leq N - 1)P(L = c + n | c \leq L \leq N - 1)$$

$$= \sum_{n=0}^{N-c-1} \left(\sum_{k=0}^n \frac{(c\mu t)^k e^{-c\mu t}}{k!} \right) \frac{\rho^n}{1 + \rho + \dots + \rho^{N-c-1}}$$

$$= \sum_{k=0}^{N-c-1} \frac{(c\mu t)^k e^{-c\mu t}}{k!} \sum_{n=k}^{N-c-1} \frac{\rho^n}{1 + \rho + \dots + \rho^{N-c-1}}$$

$$P(W_q > t | W_q > 0) = \sum_{k=0}^{N-c-1} \frac{(c\mu t)^k e^{-c\mu t}}{k!} \frac{\rho^k - \rho^{N-c}}{1 - \rho^{N-c}} = \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N-c-k}}{1 - \rho^{N-c}}$$

Then for $\rho \neq 1$, we have that

$$P(W_q > t) = P(\text{waiting|no - blocking}) \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N-c-k}}{1 - \rho^{N-c}}, t \geq 0 \tag{2.13}$$

Where we have used the fact that[15]

$$P(L = c + n | c \leq L \leq N - 1) = \frac{\rho^n}{1 + \rho + \dots + \rho^{N-c-1}}, \quad \text{for } 0 \leq n \leq N - c - 1 \tag{2.14}$$

For $\rho = 1$, we also have that

$$P(W_q > t) = P(\text{waiting|no - blocking}) \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \left(1 - \frac{k}{N-c} \right) \tag{2.15}$$

In same line of reasoning, we derive the mathematical formula for computing the Average Speed to Answer (ASA) as follows:

$$\begin{aligned}
 ASA &= E(W_q) = \int_0^\infty P(W_q > t) dt \\
 &= P(\text{waiting|no - blocking}) \sum_{k=0}^{N-c-1} \frac{\rho^k - \rho^{N-c}}{1 - \rho^{N-c}} \int_0^\infty \frac{(c\mu t)^k e^{-c\mu t}}{k!} dt \\
 &= P(\text{waiting|no - blocking}) \sum_{k=0}^{N-c-1} \frac{\rho^k - \rho^{N-c}}{(1 - \rho^{N-c})c\mu} \\
 &= P(\text{waiting|no - blocking}) \frac{1 - \rho^{N-c}(1 + (1 - \rho)(N - c))}{(1 - \rho)(1 - \rho^{N-c})c\mu} \\
 &= \frac{B(c, a)(1 - \rho^{N-c})}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c-1})} \frac{1 - \rho^{N-c}(1 + (1 - \rho)(N - c))}{(1 - \rho)(1 - \rho^{N-c})c\mu} \\
 E(W_q) &= \frac{(1 - \rho^{N-c}(1 + (1 - \rho)(N - c)))B(c, a)}{(1 - \rho + \rho B(c, a)(1 - \rho^{N-c-1}))(1 - \rho)c\mu}
 \end{aligned}$$

By the application of Little’s law, we have that

$$E(L_q) = \lambda E(W_q)(1 - P(\text{blocking}))$$

III. LIMITING BEHAVIOUR OF THE $M/M/c/N$ MODEL PERFORMANCE INDICATORS

In this section of the paper we shall prove some limiting properties of the $M/M/c/$ model with respect to N .

Theorem 3.1

Given that c and other model parameters remain constant, $P(W_q > t)$ is an increasing function of N .

For $0 < \rho < 1$,

$$\begin{aligned}
 P(\text{waiting|no - blocking}) &= \frac{B(c, a)(1 - \rho^{N-c})}{1 - \rho + \rho B(c, a)(1 - \rho^{N-c-1})} = \frac{B(c, a)}{\frac{1 - \rho}{(1 - \rho^{N-c-1})} + \rho B(c, a)} \\
 &< \frac{B(c, a)}{1 - \rho + \rho B(c, a)}, \quad \text{Since } 1 - \rho + \rho B(c, a) > \frac{1 - \rho}{(1 - \rho^{N-c-1})} + \rho B(c, a)
 \end{aligned}$$

Proof

First, we need to show that $P(\text{waiting|no - blocking})$ is an increasing function of N :

For $\rho = 1$,

$$P(\text{waiting}|\text{no - blocking}) = \frac{(N - c)B(c, c)}{1 + B(c, c)[N - c - 1]} = \frac{B(c, c)}{\frac{1}{N - c} + \frac{B(c, c)[N - c - 1]}{N - c}}$$

$$\cong \frac{B(c, c)}{\frac{1}{N - c} + B(c, c)} < 1, \quad \text{since } \frac{1}{N - c} + B(c, c) > B(c, c)$$

Now that we have established the fact that $P(\text{waiting}|\text{no - blocking})$ is an increasing function of N , we will proceed to show that $P(W_q > t)$

is an increasing function of N , given that c and other model parameters remain constant.

38

Recall from equation 2.13; for $0 < \rho < 1$

$$P(W_q > t) = P(\text{waiting}|\text{no - blocking}) \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N-c-k}}{1 - \rho^{N-c}}, t \geq 0$$

and from equation 2.15; for $\rho = 1$, we have

$$P(W_q > t) = P(\text{waiting}|\text{no - blocking}) \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \left(1 - \frac{k}{N - c}\right)$$

We are only left to show that

$$\gamma(N) := \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N-c-k}}{1 - \rho^{N-c}} \quad \text{and} \quad \varphi(N) := \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \left(1 - \frac{k}{N - c}\right)$$

are increasing functions of N .

$$\text{For } i \in \mathbb{N}, \quad \gamma(N + i) := \sum_{k=0}^{N+i-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N+i-c-k}}{1 - \rho^{N+i-c}}$$

$$\begin{aligned} \gamma(N + i) &= \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N+i-c-k}}{1 - \rho^{N+i-c}} + \sum_{k=N-c}^{N+i-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N+i-c-k}}{1 - \rho^{N+i-c}} \\ &> \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N-c-k}}{1 - \rho^{N-c}} = \gamma(N) \end{aligned}$$

since

$$\sum_{k=N-c}^{N+i-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \frac{1 - \rho^{N+i-c-k}}{1 - \rho^{N+i-c}} > 0 \quad \text{and} \quad \frac{1 - \rho^{N+i-c-k}}{1 - \rho^{N+i-c}} \geq \frac{1 - \rho^{N-c-k}}{1 - \rho^{N-c}}$$

In same way,

$$\text{For } i \in \mathbb{N}, \quad \varphi(N+i) := \sum_{k=0}^{N+i-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \left(1 - \frac{k}{N+i-c}\right)$$

$$\begin{aligned} \varphi(N+i) &= \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \left(1 - \frac{k}{N+i-c}\right) + \sum_{k=N-c}^{N+i-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \left(1 - \frac{k}{N+i-c}\right) \\ &> \sum_{k=0}^{N-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \left(1 - \frac{k}{N-c}\right) = \varphi(N) \end{aligned}$$

since

$$\sum_{k=N-c}^{N+i-c-1} \frac{(\lambda t)^k e^{-c\mu t}}{k!} \left(1 - \frac{k}{N+i-c}\right) > 0 \text{ and } 1 - \frac{k}{N+i-c} \geq 1 - \frac{k}{N-c}, 0 \leq k \leq N-c-1$$

IV. CONCLUSIONS

In this paper, we have discussed in detail the modeling of a call center as single-node using the Markovian queueing techniques. We considered the $M/M/c/c$ Erlang B Loss model and the $M/M/c/c$ Erlang C model as well as the more general $M/M/c/N$ model. Our emphasis is on the derivation of the exact performance measures of these well-known models. Considering the $M/M/c/N$ model, we expressed the system performance measures in terms of Erlang B formula, which facilitates the computation as well as the analysis. Using the results emanating from the analysis, we showed the monotonicity properties for performance measures with respect to N and c .

REFERENCES RÉFÉRENCES REFERENCIAS

- Adan I. and Jacques R. Queueing Theory, Eindhoven, Netherlands, 2002.
- Boucherie, R. Product-Form in Queueing Networks. Ph.D. Thesis, Free University, Amsterdam, The Netherlands, 1992.
- Bolch, G., Greiner, S., de Meer, H., and Trivedi, K. S. Queueing networks and Markov chains: modeling and performance evaluation with computer science applications. Wiley Interscience, New York, NY, USA, 1998.
- Borst, S., Mandelbaum, A., and Reiman, M. Dimensioning large call centers. Operations Research 52, 17-34, 2004.
- Chen, H., and Yao, D. Fundamentals of queueing networks: performance, asymptotics and optimization. Springer-Verlag, New York, 2001.
- Cooper, R. Introduction to Queueing Theory. North-Holland, New York, 1981.
- Gans, N., Koole, G., and Mandelbaum, A. Telephone call centers: Tutorial, review, and research prospects. Manufacturing and Service Operations Management 5, 79-141, 2003.
- Gnedenko, B., and Kovalenko, I. Introduction to Queueing Theory. Birkhauser Boston Inc., Cambridge, MA, 1968.
- Donald Gross, John F. Shortle, James M. Thompson, and Carl M. Harris. Fundamentals of Queueing Theory, third ed. Wiley, New York, NY, 1998.
- Kleinrock, L. Queueing Systems, Vol. 1. Wiley, New York, NY, 1975.
- Koole, G., and Mandelbaum, A. Queueing models of call centers: An introduction. Annals of Operations Research 113, 41-59, 2002.
- Mandelbaum, A. Call centers (centres) research bibliography with abstracts. Tech. rep., Technion, Haifa, Israel, Version 7: May 4, 2006.
- Mandelbaum, A., and Zeltyn, S. The Palm/Erlang-A queue, with applications to call centers. Tech. rep., Service Engineering Lecture Notes, 2004.
- Ross, S. M. Stochastic Processes. Second Edition. Wiley Inc, New York, 1996.
- Ross, S. M. Introduction to Probability Models, Tenth Edition. Academic Press, 2011.
- Srinivasan, R., Talim, J., and Wang, J. Performance analysis of a call center with interactive voice response units. TOP: An Official Journal of the Spanish Society of Statistics and Operations Research 12, 91-110, 2004.

17. Zhidong Zhang. Call centres with balking and abandonment: from queueing to queueing network models. Ph.D. Thesis, University of Saskatchewan, Saskatoon, Saskatchewan, 2010.

