# MULTI-CHANNEL BI-LEVEL HETEROGENEOUS SERVERS BULK ARRIVAL QUEUEING SYSTEM WITH ERLANGIAN SERVICE TIME

Ahmed M. M. Sultan
Egypt Air Force, Cairo, Egypt
amsultan_52@yahoo.com

**Abstract**- An easily applicable algorithm to solve problems involving bulk-arrival queues with a breakdown of one of the heterogeneous servers in case of steady state is introduced. A Monte Carlo study for numerically finding the limiting distribution of the number in the system for the bulk arrival, multi-server queueing model (M[x]/EK/C; C-1/FCFS) with heterogeneous servers is presented. The system consists of servers of varying efficiency. This paper presents multi-channel queue with Poisson arrivals, Erlangian service time distributions in which all servers have equal breakdown chance. Measures of system performance including mean queue length, mean waiting time, and blocking probability are reported. Numerical results are obtained by simulation of the entire system. Examples of extensive numerical results for certain measures of efficiency are presented in tabular and chart form. In all cases, the proposed method is computationally efficient, accurate and reliable for both high and low values of the model parameters.

**Key Words**- Bulk Arrival; Multi-server; Heterogeneous Servers; Server Breakdown; Queueing System Performance; Erlangian Service Time

## 1.INTRODUCTION

We discuss a multi-channel, first-come, first-served situation in which arrivals occur in groups or in bulk and the system is in one of two levels of operation. In one level of operation all C servers are available and in the other level only C-1 servers are available. The interruptions in the service process are due to breakdown of one of the servers, to scheduling policy, or to one of the servers leaving the system temporarily.
Bulk arrival, multi-server queueing systems have been studied from different aspects (see Chaudhry and Briere [3], Chaudhry, et. al. [4], Chaudhry and Kim [5], Chanke [2], and Sultan, et. al. [1].)

This paper considers the multi-server Erlangian queueing system (M[x]/E$^k$/C; C-1/FCFS) with heterogeneous servers. Such a system has the following features:
Such a queue is considered with C heterogeneous servers. Only one of the servers can break down. When one of the servers breaks down, the system operates with the remaining (C-1) servers (system is in level-1). After the broken server is repaired and put back into service the system re-operates with C servers again (the system is in level-2). Thus a system alternates between two modes of system operations. This is due to breakdown of one of the servers, to scheduling policy, or to one of the servers leaving

the system temporarily. Service is on a first come, first served basis. There is no limit on queue length.

The paper presents different numerical measures of the system performance such as the expected number of customers in the system and in the queue, the expected waiting time per customer in the system and in the queue, the blocking probability and the probability that no customers in the system.

The paper is organized in 5 sections. Section 2 presents a description of the study design and introduces notation and basic assumptions for the model under study. Section 3 describes the procedure and steps used for the analysis of the system together with code and routines used. Section 4 discusses extensive Monte Carlo results obtained from analyzing the system. The results include different tables and graphs. Finally, section 5 concludes the paper.

## 2. MODEL DESIGN AND NOTATION

The model under study is ($M^{[x]}/E_k/C$; C-1/FCFS) where groups of customers (bulk) arrive at random times with mean bulk arrival rate $\lambda$. The group size is namely positive Poisson distributed. The queueing system under study has heterogeneous servers where service time has an Erlang type k distribution. In analyzing such a model, it is convenient to consider the Erlang as being made up of k exponential phases, each with mean 1/ ($k\mu_s$). Our system alternates between two levels of system operations, this is due to breakdown of one of the servers. This means that only one server is allowed to breakdown randomly according to a discrete uniform distribution that assigns one of the servers to be out of service where the server breakdown has equal chance over all servers in the system. The mean time that the system operates with C servers and (C-1) servers is $1/\alpha$ and $1/\beta$ respectively.

The queue discipline for groups is FCFS while the service discipline within the bulks is based on randomly choosing one of the customers mentioned earlier. The conditional probability of the customer waiting for d departures before his service commences given the state of the system n just before the arrival of the bulk is given by:

$$\sum_{r=c-n+d}^{\infty} \frac{ra_r}{m} \frac{1}{r} = \frac{1}{m} \sum_{r=c-n+d}^{\infty} a_r \qquad , d = 1, 2,\ldots$$

$$0 \le n \le d+C-1 \qquad (1)$$

where the bulk size X is an r.v. with distribution given by $a_r$ =P (X =r), r $\geq$1, X has mean m , $0 \le m = \sum_{r=1}^{\infty} ra_r < \infty$ and variance $\sigma_a^2, 0 < \sigma_a^2 < \infty$.

Customers of a certain group are served randomly. If X is the size of a bulk then customer i, i =1,2,…,X has the same chance of joining service. In order to generate an equal chance to all customers in a given batch of size X, either a uniform assignment of the order in which they may be served is used or using a predefined permutation sequence. This permutation routine helps to recognize every customer in the batch.

Notation used for the model can be summarized by:

| | | | | |
|---|---|---|---|---|
| $\mu_s$ : | Constant service rate of a server number $s$; $1 \le s \le C$. | | $C$ : | Number of parallel servers. |
| $\overline{\mu}$ : | Expected service rate of servers. | | $\lambda$ : | Mean bulk arrival rate. |

| | | | | |
|---|---|---|---|---|
| $\alpha$ : | Transition rate from level -2 to level-1. | K : | The first phase of service. |
| | | 1 : | The last phase of service. |
| $P_B$ : | The blocking probability that defines the probability that all servers are busy. | m : | The average bulk size. |
| | | $L_Q$ : | The average queue length. |
| $\rho$ : | The traffic intensity. | $W_S$ : | The average waiting time per customer in the system. |
| $N$ : | The number of groups. | | |
| $\theta$ : | The parameter of the group size distribution. | $L_S$ : | The average number of customers in the system. |
| TOSC | TISC + ST | | |
| WTQ | Waiting time of a customer in the queue. | $W_Q$ : | The average waiting time per customer in the queue. |
| WTS | Waiting time of a customer in the system. | Level-1: | All $C$ servers are available for serving the customers in this level. |
| TBS | The busy time for a server. | Level-2: | All ($C$-1) servers are available for serving the customers in this level. |
| CUMWTQ | Cumulative waiting time in queue per customer. | $P_{0i}$ : | The probability of having no customers in the system when the system is in level-i, i= 1,2. |
| CUMWTS | Cumulative waiting time in system per customer. | | |
| TOB | Departure time for a batch. | $\beta$ : | Transition rate from level -1 to level-2. |

## 3. PROCEDURE

The procedure for studying the previously discussed system can be described in the following steps:

1- Generate 10000 interarrival times for 10000 different batches from exponential distribution with mean $1/\lambda$ with each batch size randomly generated from a positive Poisson distribution with parameter $\theta$.

2- The service times for each customer in the successive batches are generated from Erlang distribution with mean $1/\mu_S > 0$; $1 \leq s \leq C$ for each server. The system alternates between two modes of operations with equal breakdown chance for each of the C servers, using an exponential distribution with mean $1/\alpha$ and $1/\beta$ for the intervals of time in which the system operates in level-2 or level-1 successively.

3- Determine the event time of breakdown and repair for each server.

4- Determine the cumulative number of customers that enter the system as the sum of batch sizes that arrive to the system and determine the arrival time for each batch.

5- An arriving batch finds the system in either level-1 or level-2. If the arrival time of a batch is greater than the event time of breakdown and less than the event time of repair, the system will be in level-1. While if the arrival time of a batch is greater than the event time of repair and less than the event time of breakdown, the system will be in level-2.

6- Based on the system mode and the relation between the time of next batch arrival and the departure time of the previous batch, increment accordingly the number of batches for level-1 or level -2 by one.

7- An arriving customer in batch will immediately start service if one of the servers is free or wait until any server becomes free and this continues until all customers in a given batch are served.

8- Calculate TOSC, WTQ, WTS, TBS, CUMWTQ, CUMWTS, and TOB.
9- If a server is broken down during serving a certain customer, this customer will quit service and will start service at the first server available.
10- This continues till 10000 batches are generated.
11- The probability that the system is in level-i, i= 1,2 while there is no any customer in the system $P_{0,i}$, the probability that no customers are in system $P_0$, and the blocking probability $P_B$ are calculated.
12- Calculate $W_Q, W_S, L_Q$, and $L_S$.

## 4. SIMULATION RESULTS AND PERFORMANCE ANALYSIS

Now, the results from the extensive Monte Carlo study described earlier will be presented and analyzed. The model performance is tested extensively for values of ($\rho$, C, X) with $0.1 \leq \rho \leq 0.9$, $1 \leq C \leq 100$ and batch size $X \leq 100$. Based on the methodology explained in the previous section, the input data includes the number of servers C, the number of batches N, the parameter of size of batches $\theta$, mean service time $1/\mu_s$, $1 \leq S \leq C$, mean interarrival time $1/\lambda$, the mean time that the system operates with C servers $1/\alpha$ and the mean time that the system operates with (C-1) servers $1/\beta$. Different performance measures are calculated. These measures include $L_Q$, $L_S$, $W_Q$, $W_S$, $P_0$, and $P_B$.. Two models are considered:

- In the first model the queueing model ($M^{[x]}/ E_2 /5; 4$/FCFS) is considered.
- In the second model the queueing model ($M^{[x]}/ M /3; 2$/FCFS) is considered.

The first model assumes that the system operates initially with five heterogeneous servers (system is in level-1). When one of the servers is broken down, the system operates with four heterogeneous servers (system is in level-2). After the repair of the broken server and putting it back into service, the system re-operates with five heterogeneous servers again. Hence, the system alternates between two modes of system operations. The service time has Erlang type 2 distribution with parameter $\mu_s > 0$, where $\mu_s$ is the service rate of server number $S$; $1 \leq S \leq C$. The size of groups followed positive Poisson distribution with parameter $\theta > 0$. There is no limit on system capacity. FCFS is the queue discipline (first-come, first-served).

In order to carry out the extensive Monte Carlo experimentation different input values are needed. These input values are considered as the input parameters for the designed computer routine. The input parameters include $1/\mu_1$, $1/\mu_2$, $1/\mu_3$, $1/\mu_4$, $1/\mu_5$, $1/\lambda$, $1/\alpha$, $1/\beta$, C, N and $\theta$.

Tables (1- 6) give different system performance measure variations with traffic intensities $\rho = m\lambda/C\bar{\mu}$ and the relative transition rate $\alpha/(\alpha+\beta)$ (see also [1] for more details).

Results from tables 1, 2, 3, and 4 are shown graphically in figures 1,2, 3, and 4 respectively. For example, figure 1 indicates the relation between traffic intensity $\rho$ and average number of customer in the queue $L_Q$ with the change of relative transition rate $\alpha/(\alpha+\beta)$ denoted REL in the figure.

**Table 1**
*Average number in the queue*
*(Number of servers: 5)*
*(Mean size of bulk=8)*

| ρ | α/(α+β) | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.200 | 0.290 | 0.332 | 0.388 | 0.445 | 0.506 |
| 0.300 | 0.748 | 0.952 | 1.134 | 1.351 | 1.483 |
| 0.400 | 1.338 | 1.646 | 1.943 | 2.196 | 2.478 |
| 0.500 | 2.160 | 3.074 | 3.468 | 4.029 | 4.611 |
| 0.600 | 3.113 | 5.283 | 6.586 | 8.018 | 9.943 |
| 0.700 | 4.446 | 9.367 | 12.395 | 16.663 | 27.132 |
| 0.800 | 6.009 | 20.134 | 28.036 | 49.735 | 151.714 |

**Table 2**
*The probability that no customers are in*
*the system at the arrival batch*
*(Mean size of bulk=8)*
*(Number of servers: 5)*

| ρ | α/(α+β) | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.200 | 0.6979 | 0.6976 | 0.6960 | 0.6949 | 0.6943 |
| 0.300 | 0.5887 | 0.5803 | 0.5680 | 0.5615 | 0.5494 |
| 0.400 | 0.4788 | 0.4678 | 0.4620 | 0.4507 | 0.4482 |
| 0.500 | 0.3872 | 0.3569 | 0.3489 | 0.3394 | 0.3244 |
| 0.600 | 0.3121 | 0.2634 | 0.2484 | 0.2214 | 0.1970 |
| 0.700 | 0.2451 | 0.1749 | 0.1537 | 0.1263 | 0.0926 |
| 0.800 | 0.1852 | 0.1017 | 0.0792 | 0.0453 | 0.0100 |

**Table 3**
*Average waiting time in the queue*
*(Number of servers: 5)*
*(Mean size of bulk=8)*

| ρ | α/(α+β) | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.200 | 1.452 | 1.659 | 1.937 | 2.248 | 2.519 |
| 0.300 | 3.724 | 4.739 | 5.643 | 6.724 | 7.381 |
| 0.400 | 6.660 | 8.219 | 9.669 | 10.929 | 12.335 |
| 0.500 | 10.750 | 15.304 | 17.264 | 20.055 | 22.968 |
| 0.600 | 15.495 | 26.296 | 32.783 | 39.909 | 49.494 |
| 0.700 | 22.133 | 46.635 | 61.697 | 82.950 | 135.074 |
| 0.800 | 29.909 | 100.225 | 139.564 | 247.627 | 756.348 |

**Table 4**
*The probability that no servers*
*are idle in the system*
*(Mean size of bulk=8)*
*(Number of servers: 5)*

| ρ | α/(α+β) | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.100 | 0.393 | 0.422 | 0.449 | 0.478 | 0.505 |
| 0.200 | 0.450 | 0.477 | 0.507 | 0.535 | 0.562 |
| 0.300 | 0.516 | 0.549 | 0.582 | 0.613 | 0.643 |
| 0.400 | 0.578 | 0.614 | 0.641 | 0.669 | 0.695 |
| 0.500 | 0.639 | 0.683 | 0.709 | 0.740 | 0.766 |
| 0.600 | 0.691 | 0.752 | 0.783 | 0.817 | 0.851 |
| 0.700 | 0.741 | 0.824 | 0.856 | 0.890 | 0.925 |
| 0.800 | 0.789 | 0.895 | 0.923 | 0.958 | 0.992 |



**Fig1:** $\rho$ versus $L_Q$ for ($\theta$=7.997309, $m$=8, C=5)



**Fig. 2:** $\rho$ versus $P_0$ for ($\theta$=7.997309, m=8, C=5)



**Fig. 3:** $\rho$ versus $W_Q$ for ($\theta$=7.997309, m=8, C=5)



**Fig. 4:** $\rho$ versus $P_B$ for ($\theta$=7.997309, m=8, C=5)

A. M. M. Sultan

### Table 5
*Average number in the system*
*(Number of servers: 5)*
*(Mean size of bulk=8)*

| $\rho$ | $\alpha/(\alpha+\beta)$ | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.200 | 1.133 | 1.163 | 1.221 | 1.274 | 1.344 |
| 0.300 | 2.131 | 2.346 | 2.541 | 2.782 | 2.926 |
| 0.400 | 3.234 | 3.515 | 3.798 | 4.033 | 4.311 |
| 0.500 | 4.576 | 5.464 | 5.829 | 6.381 | 6.960 |
| 0.600 | 5.991 | 8.179 | 9.485 | 10.939 | 12.887 |
| 0.700 | 7.827 | 12.777 | 15.810 | 20.113 | 30.616 |
| 0.800 | 9.895 | 24.040 | 31.929 | 53.663 | 155.650 |

### Table 6
*Average waiting time in the system*
*(Number of servers: 5)*
*(Mean size of bulk=8)*

| $\rho$ | $\alpha/(\alpha+\beta)$ | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.200 | 5.665 | 5.815 | 6.106 | 6.43 | 6.692 |
| 0.300 | 10.608 | 11.677 | 12.648 | 13.848 | 14.562 |
| 0.400 | 16.099 | 17.553 | 18.903 | 20.076 | 21.461 |
| 0.500 | 22.778 | 27.200 | 29.017 | 31.763 | 34.642 |
| 0.600 | 29.819 | 40.714 | 47.215 | 54.449 | 64.147 |
| 0.700 | 38.961 | 63.612 | 78.699 | 100.121 | 152.415 |
| 0.800 | 49.256 | 119.663 | 158.943 | 267.185 | 775.969 |

The second model studies the case with heterogeneous servers ($M^{[x]}/ M /C$; C-1/FCFS) which is a special case of the model (($M^{[x]}/ E_k /C$; C-1/FCFS) with heterogeneous servers which occurs when K=1.

In this model, the system is assumed to operate initially with three heterogeneous servers (system is in level-2) and when one of the three servers is broken down, the system operates with two heterogeneous servers (system is in level-1). After the repair of the broken server and putting it back into service, the system re-operates with three heterogeneous servers again. Hence, the system alternates between two modes of system operations. The service time has Exponential distribution with parameter $\mu_s > 0$, where $\mu_s$ is the service rate of server number $s$, $1 \le s \le 3$. The size of bulks followed positive Poisson distribution with parameter $\theta > 0$. There is no limit on system capacity. The first-come, first-served is the queue discipline.

For the second model, similar results in tables and graphs are given in tables 7–10 and figures 5 – 8.

### Table 7
*Average number in the queue*
*(Number of servers: 3)*
*(Mean size of bulk=8)*

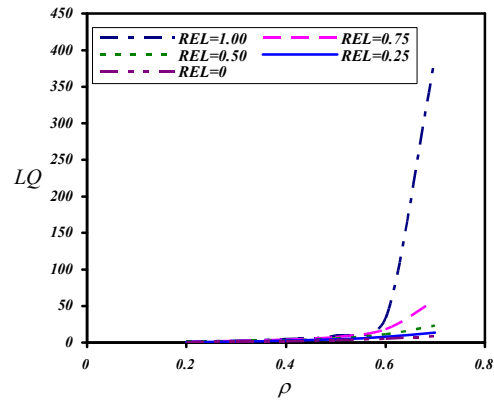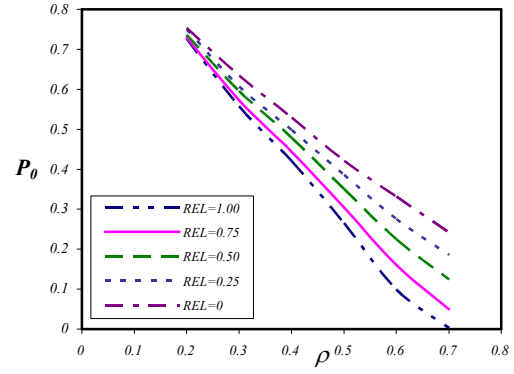| $\rho$ | $\alpha/(\alpha+\beta)$ | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.200 | 0.542 | 0.672 | 0.789 | 0.931 | 1.039 |
| 0.300 | 1.163 | 1.462 | 1.755 | 2.169 | 2.492 |
| 0.400 | 1.964 | 2.552 | 3.158 | 3.836 | 4.658 |
| 0.500 | 3.410 | 4.468 | 5.781 | 7.631 | 9.820 |
| 0.600 | 5.361 | 7.771 | 11.369 | 18.150 | 33.619 |
| 0.700 | 8.865 | 13.327 | 23.523 | 57.488 | 388.467 |
| 0.800 | 15.686 | 29.520 | 64.754 | 1447.44 | 4525.67 |



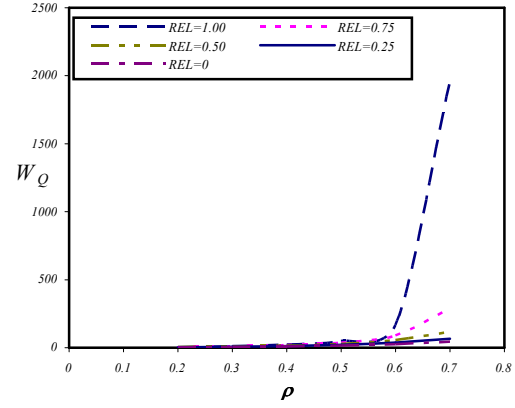**Fig. 5:** $\rho$ versus $L_Q$ for ($\theta$=7.997309, $m$=8, C=3)

**Table 8**
*The probability that no customers are in
the system at the arrival batch
(Mean size of bulk = 8)
(Number of servers: 3)*

| ρ | α/(α+β) | | | | |
|---|---|---|---|---|---|
| | **0.00** | **0.250** | **0.500** | **0.750** | **1.00** |
| **0.200** | *0.7533* | *0.7493* | *0.7359* | *0.7294* | *0.7262* |
| **0.300** | *0.6339* | *0.6072* | *0.5957* | *0.5723* | *0.5576* |
| **0.400** | *0.5290* | *0.4980* | *0.4790* | *0.4440* | *0.4201* |
| **0.500** | *0.4211* | *0.3860* | *0.3504* | *0.3047* | *0.2647* |
| **0.600** | *0.3317* | *0.2754* | *0.2252* | *0.1606* | *9.92E-02* |
| **0.700** | *0.2420* | *0.1859* | *0.1244* | *4.94E-02* | *2.50E-03* |
| **0.800** | *0.1567* | *9.66E-02* | *4.08E-02* | *7.00E-04* | *4.00E-04* |



**Fig 6:** ρ *versus* $P_0$ *for* (θ=7.997309, m=8, C=3)

**Table 9**
*Average waiting time in the queue
(Number of servers: 3)
(Mean size of bulk=8)*

| ρ | α/(α+β) | | | | |
|---|---|---|---|---|---|
| | **0.00** | **0.250** | **0.500** | **0.750** | **1.00** |
| **0.200** | *2.694* | *3.339* | *3.927* | *4.634* | *5.174* |
| **0.300** | *5.788* | *7.279* | *8.738* | *10.796* | *12.406* |
| **0.400** | *9.775* | *12.701* | *15.719* | *19.093* | *23.182* |
| **0.500** | *16.974* | *22.239* | *28.776* | *37.985* | *48.877* |
| **0.600** | *26.684* | *38.681* | *56.592* | *90.349* | *167.345* |
| **0.700** | *44.128* | *66.34* | *117.129* | *286.194* | *1963.37* |
| **0.800** | *78.078* | *146.965* | *322.494* | *7517.54* | *25851.3* |



**Fig. 7:** ρ versus $W_Q$ for (θ=7.997309, m=8, C=3)

**Table 10**
*The probability that no servers
are idle in the system
(Mean size of bulk = 8)
(Number of servers: 3)*

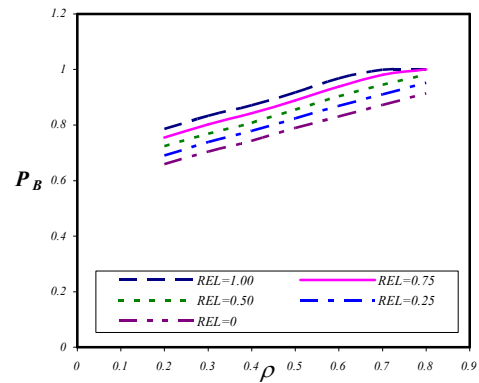| ρ | α/(α+β) | | | | |
|---|---|---|---|---|---|
| | **0.00** | **0.250** | **0.500** | **0.750** | **1.00** |
| **0.200** | *0.65925* | *0.69048* | *0.72386* | *0.75583* | *0.78627* |
| **0.300** | *0.70425* | *0.7388* | *0.76886* | *0.80216* | *0.83292* |
| **0.400** | *0.74433* | *0.77958* | *0.80848* | *0.84226* | *0.87091* |
| **0.500** | *0.79038* | *0.82322* | *0.85465* | *0.88773* | *0.91712* |
| **0.600** | *0.83092* | *0.86901* | *0.90312* | *0.93840* | *0.96837* |
| **0.700** | *0.87232* | *0.91029* | *0.94507* | *0.98028* | *0.99904* |
| **0.800** | *0.91375* | *0.95173* | *0.98142* | *0.99976* | *0.99987* |



**Fig. 8:** ρ *versus* $P_B$ *for* (θ=7.997309, m=8, C=3)

*Table 11*
*Average waiting time in the system*
*(Number of servers: 3)*
*(Mean size of bulk=8)*

| $\rho$ | $\alpha/(\alpha+\beta)$ | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.200 | 5.506 | 6.066 | 6.573 | 7.198 | 7.660 |
| 0.300 | 10.208 | 11.619 | 13.013 | 14.996 | 16.543 |
| 0.400 | 15.637 | 18.464 | 21.389 | 24.676 | 28.661 |
| 0.500 | 24.415 | 29.593 | 36.056 | 45.198 | 56.015 |
| 0.600 | 35.617 | 47.599 | 65.519 | 99.317 | 176.312 |
| 0.700 | 54.565 | 76.746 | 127.544 | 296.645 | 1973.60 |
| 0.800 | 89.928 | 158.698 | 334.146 | 7529.17 | 25862.8 |

*Table 12*
*Average number in the system*
*(Mean size of bulk=8)*
*(Number of servers: 3)*

| $\rho$ | $\alpha/(\alpha+\beta)$ | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.250 | 0.500 | 0.750 | 1.00 |
| 0.200 | 1.106 | 1.219 | 1.320 | 1.446 | 1.539 |
| 0.300 | 2.051 | 2.334 | 2.614 | 3.013 | 3.324 |
| 0.400 | 3.142 | 3.709 | 4.297 | 4.957 | 5.758 |
| 0.500 | 4.905 | 5.945 | 7.243 | 9.080 | 11.253 |
| 0.600 | 7.155 | 9.563 | 13.163 | 19.951 | 35.420 |
| 0.700 | 10.961 | 15.418 | 25.615 | 59.587 | 390.492 |
| 0.800 | 18.067 | 31.877 | 67.094 | 1449.68 | 4527.70 |

## 6. CONCLUSION

The final conclusion from the previous study shows that: -

1- The increase of relative transition $\alpha/(\alpha+\beta)$ has remarkable effect on the average number of customers in the queue and in the system, while the increase of the traffic intensity $\rho$ has a big noticeable effect on a performance measures. (See tables 1, 7 and 5, 12).

3- The average number of customers in the system and in the queue increases as the average group size increases. So the increase of relative transition $\alpha/(\alpha+\beta)$ has noticeable effect on $L_Q$ and $L_S$ but not as well when the average group size is relatively close to the value of the number of servers. (See tables 1, 7 and 5, 12).

1- The most effect on the average number of customers in the system and in the queue is happening when the traffic intensity $\rho$ is very close to unity. The absence of one server affects the average number of customers as traffic intensity gets near the unity (heavy traffic). (See tables 1, 7 and 5, 12).

5- The average waiting time in the queue and in the system is highly affected by the absence of one of the servers for both low and high traffic intensity. (See tables 3, 6 and 9, 11).

6- The probability of having no customers in the system $P_0$ at an arrival batch i.e. the percentage of time the system is idle is not highly affected by the absence of one of the servers, but it decreases as the traffic intensity increases. (See tables 2 and 8).

7- The probability that no servers are idle in the system $P_B$ is affected by the increase of traffic intensity, where it increases as the traffic intensity increases. So the blocking probability has remarkable effect resulting from absence of one of servers where it increases as the relative transition $\alpha/(\alpha+\beta)$ increases. (See tables 4 and 10).

8- The results show in tables 7 to 12 for the model ($M^{[x]}/M/3$; 2/FCFS) illustrate that the results with the absence of one server is much more noticeable than the results show in tables 1 to 6 for the model ($M^{[x]}/E_k/5$; 4/FCFS). Where, in the first model the service time has Exponential distribution and the number of servers is fewer than the number of servers in the second model.

9- The results show in tables 7 to 12 for the model ($M^{[x]}$/M /3; 2/FCFS) illustrate that the results with heterogeneous servers and the absence of one servers is much more noticeable than the results for the model ($M^{[x]}$/M/5; 4/FCFS) with homogeneous servers and increasing of mean of size of bulk which described in [1] (i.e. the performance measures for model ($M^{[x]}$/M/3; 2/FCFS) with heterogeneous servers are noticeable highly affected than the model ($M^{[x]}$/M/5; 4/FCFS) with homogeneous servers when the traffic intensity $\rho$ is very close to unity and an absence of one server.

Finally, this paper introduces an easily applicable algorithm to solve problems involving bulk-arrival queues with a breakdown of one of the heterogeneous server in case of steady state. This approach was preferred to producing large tables of exact results, varying the queueing parameters because of the endless list of possible combinations when applied to bulk queues and one of the servers break down. The performance measures are changed in response to the changes of the operating parameters. We documented the behavior of the system when one of the servers temporarily leaves the system with useful graphical representation to give the reader an opportunity to watch the system behavior over the traffic intensity and the relative transition rate.

## 7. REFERENCES

1. A.M. Sultan, N.A.Hassan, N.M. Elhamy, Computational analysis of a multi-server bulk arrival system with two modes server breakdown, *Journal of the Mathematical & Computational Applications,*2003.
2. J. Chanke, Two thresholds of an M/G/1 queueing system with server breakdowns and two vacation types. *Journal of National Taichung Institute of Technology*, **129**, 2001.
3. M.L. Chaudhry and G. Briere, Computational analysis of single-server bulk arrival queues. $M^{[x]}$/G/1, *Journal of Computer and Operations Research*, **15**, 283-292, 1988.
4. M.L. Chaudhry, J.G.C.  Templeton, and J. Medhi, Computational results of multi-server bulk arrival queues with constant service time $M^{[x]}$/D/C, *Journal of Opns. Res.*, **40**, 229-238, 1992.
5. M.L. Chaudhry, and N. K. Kim, A complete and simple solution for the discrete time multi-server queue with bulk arrivals and deterministic service times, Dept. Indust. Eng., KAIST, 2002.