# Approximate Description of Dynamics of a Closed Queueing Network Including Multi-servers

Svetlana Anulova$^{(\boxtimes)}$

V.A. Trapeznikov Institute of Control Sciences of Russian Academy of Sciences,
65 Profsoyuznaya Street, Moscow, Russia
anulovas@ipu.rssi.ru
http://www.ipu.ru/

**Abstract.** We investigate a closed network consisting of two multi-servers with $n$ customers. Service requirements of customers at a server have a common cumulative distribution function. The state of the network is described by the following state parameters: for each multi-server and for the queue empirical measures of the age of customers being serviced/waiting in the queue multiplied by $n^{-1}$. The approximation of a single multi-server dynamics is currently studied by famous scientists H. Kaspi, K. Ramanan, W. Whitt et al. We find approximation for a network, but only in discrete time.

A motivation for studying such systems is that they arise as models of computer data systems and call centers.

**Keywords:** Multi-server queues · GI/G/n queue · Fluid limits · Mean-field limits · Strong law of large numbers · Measure-valued processes · Call centers

## 1 Introduction

### 1.1 Review of Investigated Contact Centers Models

In the last ten years an extensive research in mathematical models for telephone call centers has been carried out, cf. [2–6,8–17]. The object has been expanded to more general customer contact centers (with contact also made by other means, such as fax and e-mail). One of important relating questions is the dynamics of multi-server queues with a large number of servers. In order to describe the object efficiently the state of the model must include: (1) for every customer in the queue the time that he has spent in it and (2) for every customer in the multi-server the time that he has spent after entering the service area, that is being received by one of the available servers.

The focus of research was on multi-server queues with a large number of servers, because it is typical of contact centers. For such queues were found fluid limits with the number of servers tending to infinity. Notice that such a limit is a deterministic fumction of time with values in a certain measure space, or in a space containing such a component. These developed deterministic fluid models provided simple first-order performance descriptions for multi-server queues under heavy loads.

## 1.2    A New Model for Contact Centers

We suggest here a more suitable model for contact centers. The number of customers is fixed. Customers may be situated in two states: normal and failure. There is a multi-server which repairs customers in the failure state. The repair time/the time duration of a normal state is a random variable, independent and identically distributed for all customers. Now "the arrival process" in the multi-server does not correspond to that of the previous $G/GI/s + GI$ model. For a large number of customers and a suitable number of servers calculate the number of current failures, so much as an approximation. This is a continuation of our work [1], where a single multi-server in a network was functioning.

We confine ourselves to a discrete time model. W. Whitt has written a very interesting seminal article [15], in a simple discrete case. About 150 authors have cited it and made generalizations to the continuous time. But their results do not enclose Whitt's discrete time ones. Walsh Zuñiga [12] with his results most close to discrete time admits only discrete time service but his arrival process is continuous.

In Whitt's article [15] the idea of the convergence proof is true and very lucid, and the proof is clearly presented. But Whitt has not covered all cases in his proofs, and at that Whitt does not accurately point to it. Only the main case is examined: the number of serviced customers is not zero and does not exceed the number of customers in the queue. Also when he proves convergence of the number of customers being served for a given time $b(t, k)$ in (6.33) he omits the case $b_s(n - 1, k - 1)$ tends to zero. We have transferred his proof technique to our new network model and filled all these misses.

## 1.3    Problem Origin

Consider a closed network consisting of $n$ customers. They may be situated in two states: normal and failure. A multi-server repairs customers in the failure state. The repair time (resp., the time duration of a normal state) is a random variable, independent and identically distributed for all customers. For a large number of customers and a suitable number of servers we shall calculate the number of current failures, so much as an approximation.

Now we give a rigorous description of this model.

Consider a closed network consisting of $n$ customers and two multi-servers. Multi-server 1 (further denoted MS1) consists of $n$ servers (for the customers in

the normal state), the time they service a customer has distribution $G^1$. Multi-server 2 (further denoted MS2) consists of $s_n n$ servers with a number $s_n \in (0, 1)$ (for the customers in the failure state), the time they service a customer has distribution $G^2$. The distributions $G^1, G^2$ are discrete: they are concentrated on $\{1, 2, \ldots\}$. Service times are independent for both servers and all customers. We will investigate the behavior of the net as $n \to \infty$, namely, we shall establish a stochastic-process fluid limit. It will be done only in a special case: discrete time $t = 0, 1, 2, \ldots$

We begin with a simple example of functioning of this network.

*Example 1.* Let at time $t = 0$ all $n$ customers be in a normal state. Each customer switches over to the failure state according to the distribution function $G^1$ and tries to enter multi-server 2. The early failure customers can do it, but with time growing multi-server 2 may become fully occupied. Then the failure customers create a queue, waiting for the first available server in multi-server 2. Recall that a server becomes afresh available with time distribution $G^2$.

In this example and everywhere further we demand:

**Assumption 1.** *Customers are served in order of their arrival to MS2 or to its queue (FCFS) by the first available server.*

Denote the number of customers at a moment $t = 0, 1, \ldots$ in MS1 (resp., MS2) by $B_n^1(t)$ (resp., $B_n^2(t)$) and the number of customers in the queue $Q_n(t) = n - B_n^1(t) - B_n^2(t)$. These quantities must be defined more exactly. Namely,

$$B_n^i(t) = \sum_{k=0}^{\infty} b_n^i(t, k), \; i = 1, 2, \text{ and } Q_n(t) = \sum_{0}^{\infty} q_n(t, k)$$

with $b_n^i(t, k)$ being the number of customers in the multi-server $i$ at the moment $t$ who have spent there time $k$, $i = 1, 2$, and $q_n(t, k)$ being the number of customers in the queue at the moment $t$ who have been there precisely for time $k$. $b_n^i(t, k)$ may also be interpreted as the number of busy servers at time $t$ in the multi-server i that are serving customers that have been in service precisely for time $k$, $i = 1, 2$.

At the same time moment $t \in \{1, 2, \ldots\}$ multiple events can take place, so we have to specify their order.

We must create a fictitious queue for the MS1—in fact this multi-server is so large ($n$ servers), that any customer of the whole quantity $n$ trying to enter the MS1 at once finds a free server in it.

At the time moment $t$ the parameters $b^1, b^2, q$ are taken from the previous time $t - 1$ and processed to the current situation.

For both multi-servers:

– first, customers in service are served;
– second, the served customers move to another multi-server queue, to the end of it;

– third, waiting customers in queue move into service of the multi-server according to Assumption 1.

Customers enter service in MS2 whenever a server is available, so that the system is work-conserving; i.e. we assume that $Q_n(t) = 0$ whenever $B_n^2(t) < s_n n$, and that $B_n^2(t) = s_n n$ whenever $Q_n(t) > 0$, $t = 0, 1, 2, \ldots$. This condition can be summarized by the equation

$$(s_n - B_n^2(t)/n)Q_n(t) = 0 \text{ for all } t \text{ and } n.$$

**Notations.** Let $\sigma_n^i(t)$ be the number of service completions in MSi at time moment $t = 1, 2, \ldots$, $i = 1, 2$. Denote for $k = 1, 2, \ldots$ $G^{i;c}(k) = 1 - G^i(k)$ and $g^i(k) = G^i(k) - G^i(k-1)$, $i = 1, 2$. Symbol $\Rightarrow$ means convergence of the network state characteristics to a constant in probability as the index $n$ denoting the number of cusomers tends to infinity.

**Theorem 1 (The Discrete-Time Fluid Limit).** *Suppose that for each $n$, the system is initialized with workload characterized by nonnegative-integer-valued stochastic processes*

$$b_n^i(0, k), \ i = 1, 2, \ and \ q_n(0, k), \ k = 0, 1, 2, \ldots,$$

*satisfying*
$$B_n^1(0) + B_n^2(0) + Q_n(0) = n, \tag{1}$$

$$B_n^2(0) \le s_n n, \ and \ (s_n n - B_n^2(0))Q_n(0) = 0 \tag{2}$$

*for each $n$ w.p.1. Suppose that $s_n \to s \in (0,1)$ and*

$$\frac{b_n^i(0, k)}{n} \Rightarrow b^i(0, k), \ i = 1, 2, \ and \ \frac{q_n(0, k)}{n} \Rightarrow q(0, k) \ for \ k = 0, 1, 2, \ldots \tag{3}$$

*as $n \to \infty$, where $s$ is a constant and $b^i(0, k), i = 1, 2$, and $q(0, k)$ are deterministic functions. Moreover, suppose that for each $\epsilon > 0$ and $\eta > 0$, there exists an integer $k_0$ such that for $n = 1, 2, \ldots$*

$$\mathbf{P}(\sum_{k=k_0}^{\infty} \frac{b_n^i(0, k)}{n} > \epsilon) < \eta, \ i = 1, 2, \ and \ \mathbf{P}(\sum_{k=k_0}^{\infty} \frac{q_n(0, k)}{n} > \epsilon) < \eta. \tag{4}$$

*Then, as $n \to \infty$,*

$$\frac{b_n^i(t, k)}{n} \Rightarrow b^i(t, k), \ i = 1, 2, \tag{5}$$

$$\frac{q_n(t, k)}{n} \Rightarrow q(t, k), \tag{6}$$

$$\frac{\sigma_n^i(t)}{n} \Rightarrow \sigma^i(t), \ i = 1, 2, \tag{7}$$

*for each $t \ge 1$ and $k \ge 0$, where $(b^1, b^2, q, \sigma^1, \sigma^2)$ is a vector of deterministic functions (all with finite values).*

*Further, for each $t = 0, 1, \ldots$*

$$\frac{B_n^i(t)}{n} \equiv \frac{\sum_{k=0}^{\infty} b_n^i(t,k)}{n} \Rightarrow B^i(t) \equiv \sum_{k=0}^{\infty} b^i(t,k),\ i = 1, 2, \tag{8}$$

$$\frac{Q_n(t)}{n} \equiv \frac{\sum_{k=0}^{\infty} q_n(t,k)}{n} \Rightarrow Q(t) \equiv \sum_{k=0}^{\infty} q(t,k), \tag{9}$$

*with*

$$B^1(t), B^2(t), Q(t) \geq 0,\ B^1(t) + B^2(t) + Q(t) = 1, \tag{10}$$

$$B^2(t) \leq s,\ and\ (s - B^2(t))Q(t) = 0. \tag{11}$$

*The evolution of the vector $(b^1, b^2, q, \sigma^1, \sigma^2)(t)$, $t = 0, 1, 2 \ldots$, proceeds with steps of $t$ in the following way. As we go from time $t - 1$ to $t$, there are two cases, depending on whether $B^2(t-1) = s$ or $B^2(t-1) < s$.*

*Case 1. $B^2(t-1) = s$. In this first case, after moment $t - 1$ asymptotically all servers are busy and in general there may be a positive queue. In this case,*

$$\sigma^i(t) = \sum_{k=1}^{\infty} b^i(t-1, k-1) \frac{g^i(k)}{G^{i;c}(k-1)}, \tag{12}$$

$$b^i(t,k) = b^i(t-1, k-1) \frac{G^{i;c}(k)}{G^{i;c}(k-1)},\ k = 1, 2, \ldots,\ i = 1, 2, \tag{13}$$

$$b^1(t,0) = \sigma^2(t), \tag{14}$$

$$b^2(t,0) = \min\{\sigma^2(t), Q(t-1) + \sigma^1(t)\}, \tag{15}$$

*and finally $q$ is determined with the help of an intermediate queue $q'$,*

$$q'(t,0) = \sigma^1(t),\ q'(t,k) = q(t-1, k-1),\ k = 1, 2, \ldots: \tag{16}$$

$$if\ \sigma^2(t) = 0\ then\ q(t,k) = q'(t,k),\ k = 0, 1, \ldots, \tag{17}$$

$$if\ \sigma^2(t) \geq \sum_{k=0}^{\infty} q'(t,k)\ then\ q(t,k) = 0,\ k = 0, 1, \ldots, \tag{18}$$

$$if\ 0 < \sigma^2(t) < \sum_{k=0}^{\infty} q'(t,k)\ then\ with \tag{19}$$

$$c(t) = \min\{i \in \{0, 1, \ldots\} : \sum_{k=i}^{\infty} q'(t,k) \leq \sigma^2(t)\}, \tag{20}$$

$$q(t,k) = \begin{cases} 0\ for\ k \geq c(t), \\ \sum_{i=c(t)-1}^{\infty} q'(t,i) - \sigma^2(t)\ for\ k = c(t) - 1, \\ q'(t,k)\ for\ k < c(t) - 1. \end{cases}$$

*Case 2. $B^2(t-1) < s$. In this second case, after the time moment $t - 1$ asymptotically all servers are not busy so that there is no queue. As in the first case, Eqs. (12), (13), and (14) hold. Instead of (15),*

$$b^2(t,0) = \min\{s - B^2(t-1) + \sigma^2(t), \sigma^1(t)\}. \tag{21}$$

*Then,*

$$q(t,k) = 0\ for all\ k > 0\ and\ q(t,0) = \sigma^1(t) - b^2(t,0). \tag{22}$$

**Proof.** For $t = 0$ conditions (3) and (4) imply the convergence presented in (8) and (9), and conditions (1), (2) provide properties (10), (11). Thus the proof of the Theorem is reduced to the following Lemma:

**Lemma 1.** *Suppose for a given $t \in \{1, 2, \ldots\}$ Eqs. (5), (6) hold for $t - 1$. Then all the statements of the Theorem hold for $t$.*

We divide the proof of this Lemma to several Lemmas below.

**Lemma 2.** *Fix $n \in \{1, 2, \ldots\}$. If for given $\epsilon > 0$ and $\eta > 0$ and integer $k_0$ holds*

$$\mathbf{P}(\sum_{k=k_0}^{\infty} \frac{b_n^i(0, k)}{n} > \epsilon) < \eta, \ i = 1, 2, \ and \ \mathbf{P}(\sum_{k=k_0}^{\infty} \frac{q_n(0, k)}{n} > \epsilon) < \eta, \quad (23)$$

*then for every $t = 1, 2, \ldots$*

$$\mathbf{P}(\sum_{k=k_0+t}^{\infty} \frac{b_n^i(t, k)}{n} > \epsilon) < \eta, \ i = 1, 2, \ and \ \mathbf{P}(\sum_{k=k_0+t}^{\infty} \frac{q_n(t, k)}{n} > \epsilon) < \eta. \quad (24)$$

*Proof.* Evidently, for any $t = 1, 2, \ldots$, $k = 0, 1, \ldots$ w.p.1

$$b_n^i(t, k + t) \leq b_n^i(0, k) \ and \ q_n(t, k + t) \leq q_n(0, k).$$

**Corollary 1.** *Condition (4) holds not only for time $t = 0$, but for any time $t = 1, 2, \ldots$, i.e., for each $t, \epsilon > 0$ and $\eta > 0$, there exists an integer $k_0$ such that for $n = 1, 2, \ldots$*

$$\mathbf{P}(\sum_{k=k_0}^{\infty} \frac{b_n^i(t, k)}{n} > \epsilon) < \eta, \ i = 1, 2, \ and \ \mathbf{P}(\sum_{k=k_0}^{\infty} \frac{q_n(t, k)}{n} > \epsilon) < \eta. \quad (25)$$

In a network consisting of $n \in \{1, 2, \ldots\}$ customers denote by $\sigma_n^i(t, k)$ the number of customers served in MSi, $i = 1, 2$, at time moment $t$ who had been in service for time $k$ at this time moment $t$, for $t, k \in \{1, 2, \ldots\}$.

**Lemma 3.** *If for given $t, k \in \{1, 2, \ldots\}$*

$$\frac{b_n^i(t - 1, k - 1)}{n} \Rightarrow b^i(t - 1, k - 1),$$

*then*

$$\frac{\sigma_n^i(t, k)}{n} \Rightarrow b^i(t - 1, k - 1) \frac{g^i(k)}{G^{i;c}(k - 1)}, \ k = 1, 2, \ldots, \ i = 1, 2,$$

$$\frac{b_n^i(t, k)}{n} \Rightarrow b^i(t, k) = b^i(t - 1, k - 1) \frac{G^{i;c}(k)}{G^{i;c}(k - 1)}, \ k = 1, 2, \ldots, \ i = 1, 2.$$

*Proof.* Set $i = 1, 2$. For each $n, t, k \geq 1$, we can represent $\sigma_n^i(t, k)$ and $b_n^i(t, k)$ as random sums of IID Bernoulli random variables; in particular,

$$\sigma_n^i(t, k - 1) = \sum_{i=1}^{b_n^i(t-1,k-1)} X_i \tag{26}$$

and

$$b_n^i(t, k) = \sum_{i=1}^{b_n^i(t-1,k-1)} (1 - X_i), \tag{27}$$

where $X_i$ assumes the value 1 if the $i$th customer among those in service at time moment $t - 1$ that have been in the system for time $k - 1$ is served at time moment $t$, and assumes the value 0 otherwise. Thus, $X_i, i \geq 1$, is a sequence of IID random variables with

$$\mathbf{P}(X_1 = 0) = \frac{g^i(k)}{G^{i;c}(k - 1)}, \ k = 1, 2, \ldots, i = 1, 2.$$

Apply now Appendix Lemma 5.

**Corollary 2.** *If for given $t \in \{1, 2, \ldots\}, i \in \{1, 2\}$*

$$\frac{b_n^i(t - 1, k - 1)}{n} \Rightarrow b^i(t - 1, k - 1), \ k = 1, 2, \ldots,$$

*then limit* (7) *holds:*

$$\frac{\sigma_n^i(t)}{n} \Rightarrow \sigma^i(t). \tag{28}$$

*Proof.* This is guaranteed by Lemma 2.

Now we shall calculate the queue after it has filled the free servers. For a given $n \in \{1, 2, \ldots\}$ we move $\sigma_n^2(t)$ customers from the end of queue, that is, the customers who have spent the longest time in the queue (if there are fewer customers in the queue, we move they all).

Define an intermediate queue $q_n'$:

$$q_n'(t, 0) = \sigma_n^1(t), \ q_n'(t, k) = q_n(t - 1, k - 1), \ k = 1, 2, \ldots \tag{29}$$

Obviously,

$$\text{if } \sigma_n^2(t) = 0 \text{ then } q_n(t, k) = q_n'(t, k), \ k = 0, 1, \ldots, \tag{30}$$

$$\text{if } \sigma_n^2(t) \geq \sum_{k=0}^{\infty} q_n'(t, k) \text{ then } q_n(t, k) = 0, \ k = 0, 1, \ldots, \tag{31}$$

$$\text{if } 0 < \sigma_n^2(t) < \sum_{k=0}^{\infty} q_n'(t, k) \text{ then with} \tag{32}$$

$$c_n(t) = \min\{i \in \{0, 1, \ldots\} : \sum_{k=i}^{\infty} q_n'(t, k) \leq \sigma_n^2(t)\}, \tag{33}$$

$$q_n(t,k) = \begin{cases} 0 \text{ for } k \geq c_n(t), \\ \sum\limits_{i=c_n(t)-1}^{\infty} q_n'(t,i) - \sigma_n^2(t) \text{ for } k = c_n(t)-1, \\ q_n'(t,k) \text{ for } k < c_n(t)-1. \end{cases}$$

*Remark 1.* As $q_n'(t,k) = q_n(t-1,k-1)$,

$$\frac{q_n'(t,k)}{n} \Rightarrow q'(t,k), \; k = 1, 2, \ldots$$

It is easy to understand other presentations of $q_n$ and $q$:

**Lemma 4.** *For $i = 0, 1, 2, \ldots$*

$$\sum_{k=i}^{\infty} q_n(t,k) = (\sum_{k=i}^{\infty} q_n'(t,k) - \sigma_n(t)) \vee 0, \tag{34}$$

$$\sum_{k=i}^{\infty} q(t,k) = (\sum_{k=i}^{\infty} q'(t,k) - \sigma(t)) \vee 0. \tag{35}$$

It follows straightforward:

$$\sum_{k=i}^{\infty} \frac{q_n(t,k)}{n} = (\sum_{k=i}^{\infty} \frac{q_n'(t,k)}{n} - \frac{\sigma_n(t)}{n}) \vee 0. \tag{36}$$

Setting in Appendix Lemma 7

$$a_n(k) = \frac{q_n'(t,k)}{n} \text{ and } \theta_n = \frac{\sigma_n(t)}{n},$$

$$a_n^-(k) = \frac{q_n(t,k)}{n}, \; a^-(k) = q(t,k) \text{ and } \theta = \sigma(t),$$

and using Remark 1 and Corollary 2, we obtain convergence in Eqs. (6), (7).

## 2  Appendix

**Lemma 5.** *$\xi_n, n = 1, 2, \ldots$, is a random variable with values in $\{0, 1, 2, \ldots\}$, and $\frac{\xi_n}{n} \Rightarrow \xi$ (a constant). Let $X_i, i = 1, 2, \ldots$, be IID random variables with values in $\{0, 1\}$. Then*

$$\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{\xi_n} X_i = \xi \mathbf{P}(X_1 = 1). \tag{37}$$

*Proof.* As $\sum_{i=1}^{\xi_n} X_i \leq \xi_n$, $n = 1, 2, \ldots$, the proposition is evident for $\xi = 0$. If $\xi > 0$, then $\lim_{n \to \infty} \xi_n = \infty$ and (weak LLN)

$$\frac{1}{n} \sum_{i=1}^{\xi_n} X_i = \frac{\xi_n}{n} \frac{\sum_{i=1}^{\xi_n} X_i}{\xi_n} \xrightarrow[n \to \infty]{} \xi \mathbf{E} X_1 = \xi \mathbf{P}(X_1 = 1).$$

**Assumption 2.** *For $n = 1, 2, \ldots$ there exists a sequence $a_n$ of nonnegative random variables $a_n(k)$, $k = 0, 1, 2, \ldots$ with finite deterministic limits $a_n(k) \Rightarrow a(k)$, $k = 0, 1, 2, \ldots$*

*For each $\epsilon > 0$ and $\eta > 0$, there exists an integer $k_0$ such that for $n = 1, 2, \ldots$*

$$\mathbf{P}(\sum_{k=k_0}^{\infty} a_n(k) > \epsilon) < \eta. \tag{38}$$

**Lemma 6.** *If Assumption 2 holds then*

$$\mathbf{P}(\sum_{k=0}^{\infty} a_n(k) < \infty) = 1, \; n = 1, 2, \ldots, \; and \; \sum_{k=0}^{\infty} a(k) < \infty. \tag{39}$$

*Furthermore, for $i=0,1,2,\ldots$*

$$\sum_{k=i}^{\infty} a_n(k) \Rightarrow \sum_{k=i}^{\infty} a(k). \tag{40}$$

*Proof.* Whitt has introduced this Assumption and derives from it partially the statement of the Lemma for $i = 0$ without proof (see the first paragraph in the proof of Theorem 6.1 [15]).

**Lemma 7.** *Suppose that Assumption 2 holds and nonnegative random variables $\theta_n$, $n = 1, 2, \ldots$, with deterministic limit $\theta$, $\theta_n \Rightarrow \theta$, are given.*

*Define new sequences[1] $a_n^-$, $n \in \{1, \ldots\}$, $a^-$ by the following equations sequences: for $i = 0, 1, 2, \ldots$*

$$\sum_{k=i}^{\infty} a_n^-(k) = (\sum_{k=i}^{\infty} a_n(k) - \theta_n) \vee 0, \tag{41}$$

$$\sum_{k=i}^{\infty} a^-(k) = (\sum_{k=i}^{\infty} a(k) - \theta) \vee 0. \tag{42}$$

*These $a_n^-$, $n \in \{1, \ldots\}$, $a^-$ sequences satisfy:[2]*

$$a_n^-(k) \Rightarrow a^-(k), \; k = 0, 1, 2, \ldots, \; and \; \sum_{k=0}^{\infty} a_n^-(k) \Rightarrow \sum_{k=0}^{\infty} a^-(k). \tag{43}$$

---

[1] They describe queues after customers leave them and go to the emptied servers.

[2] Even starting from an arbitrary $k$.

*Proof.* According to Lemma 6,

$$\text{for } i = 0, 1, 2, \ldots \quad \sum_{k=i}^{\infty} a_n(k) \Rightarrow \sum_{k=i}^{\infty} a(k). \tag{44}$$

As function $f(x, y) = (x - y) \vee 0$ is continuous, this and convergence of $\theta_n$, $n = 1, 2, \ldots$, guarantees convergence

$$\text{for } i = 0, 1, 2, \ldots \quad (\sum_{k=i}^{\infty} a_n(k) - \theta_n) \vee 0 \Rightarrow (\sum_{k=i}^{\infty} a(k) - \theta) \vee 0. \tag{45}$$

what is identical to

$$\sum_{k=i}^{\infty} a_n^-(k) \Rightarrow \sum_{k=i}^{\infty} a^-(k). \tag{46}$$

Finally, for $i = 0, 1, 2, \ldots$

$$a_n^-(i) = \sum_{k=i}^{\infty} a_n^-(k) - \sum_{k=i+1}^{\infty} a_n^-(i) \Rightarrow \sum_{k=i}^{\infty} a^-(k) - \sum_{k=i+1}^{\infty} a^-(k) = a^-(k). \tag{47}$$

## 3   Conclusion

We have investigated a model without abandonment in the queue, although the necessary details of the customers age in the queue are provided. Since such a behavior in the queue is universally recognized, we intend to consider it next. As customers in a closed network cannot abandon it, probably we shall choose instead a similar version of "nonpersistent customers", see [7].

## References

1. Anulova, S.V.: Age-distribution description and "fluid" approximation for a network with an infinite server. In: International Conference "Probability Theory and its Applications", Moscow, pp. 219–220. 26–30 June 2012. (M.: LENAND)
2. Brown, L., Gans, N., Mandelbaum, A., Sakov, A., Shen, H., Zeltyn, S., Zhao, L.: Statistical analysis of a telephone call center: a queueing-science perspective. J. Am. Stat. Assoc. **100**(469), 36–50 (2005)
3. Dai, J., He, S.: Many-server queues with customer abandonment: a survey of diffusion and fluid approximations. J. Syst. Sci. Syst. Eng. **21**(1), 1–36 (2012). http://dx.doi.org/10.1007/s11518-012-5189-y, and http://link.springer.com/article/10.1007/s11518-012-5189-y
4. Gamarnik, D., Goldberg, D.A.: On the rate of convergence to stationarity of the M/M/n queue in the Halfin-Whitt regime. Ann. Appl. Probab. **23**(5), 1879–1912 (2013)
5. Gamarnik, D., Stolyar, A.L.: Multiclass multiserver queueing system in the Halfin-Whitt heavy traffic regime: asymptotics of the stationary distribution. Queueing Syst. **71**(1–2), 25–51 (2012)

6. Kang, W., Pang, G.: Equivalence of fluid models for Gt/GI/N+GI queues. ArXiv e-prints, February 2015. http://arxiv.org/abs/1502.00346
7. Kang, W.: Fluid limits of many-server retrial queues with nonpersistent customers. Queueing Systems (2014). http://gen.lib.rus.ec/scimag/index.php?s=10.1007/s11134-014-9415-9
8. Kaspi, H., Ramanan, K.: Law of large numbers limits for many-server queues. Ann. Appl. Probab. **21**(1), 33–114 (2011)
9. Koçağa, Y.L., Ward, A.R.: Admission control for a multi-server queue with abandonment. Queueing Syst. **65**(3), 275–323 (2010)
10. Pang, G., Talreja, R., Whitt, W.: Martingale proofs of many-server heavy-traffic limits for Markovian queues. Probab. Surv. **4**, 193–267 (2007). http://www.emis.ams.org/journals/PS/viewarticle9f7e.html?id=91&layout=abstract
11. Reed, J.: The $G/GI/N$ queue in the Halfin-Whitt regime. Ann. Appl. Probab. **19**(6), 2211–2269 (2009)
12. Zuñiga, A.W.: Fluid limits of many-server queues with abandonments, general service and continuous patience time distributions. Stochast. Process. Appl. **124**(3), 1436–1468 (2014)
13. Ward, A.R.: Asymptotic analysis of queueing systems with reneging: A survey of results for FIFO, single class models. Surv. Oper. Res. Manage. Sci. **17**(1), 1–14 (2012). http://www.sciencedirect.com/science/article/pii/S1876735411000237
14. Whitt, W.: Engineering solution of a basic call-center model. Manage. Sci. **51**(2), 221–235 (2005)
15. Whitt, W.: Fluid models for multiserver queues with abandonments. Oper. Res. **54**(1), 37–54 (2006). http://pubsonline.informs.org//abs/10.1287/opre.1050.0227
16. Xiong, W., Altiok, T.: An approximation for multi-server queues with deterministic reneging times. Ann. Oper. Res. **172**, 143–151 (2009). http://link.springer.com/article/10.1007/s10479-009-0534-3
17. Zhang, J.: Fluid models of many-server queues with abandonment. Queueing Syst. **73**(2), 147–193 (2013). http://link.springer.com/article/10.1007/s11134-012-9307-9