

# A unifying conservation law for single-server queues

Urtzi Ayesta  
LAAS-CNRS  
7, Avenue Colonel Roche  
31077, Toulouse  
e-mail: urtzi@laas.fr

## Abstract

We develop a conservation law for a multi-class  $GI/GI/1$  queue operating under a general work-conserving scheduling discipline. For single-class single-server queues, conservation laws have been obtained for both non-anticipating and anticipating disciplines with general service time distributions. For multi-class single-server queues, conservation laws have been obtained for (i) non-anticipating disciplines with exponential service time distributions and (ii) non-preemptive non-anticipating disciplines with general service time distributions. The unifying conservation law we develop generalizes already existing conservation laws. In addition it covers popular non-anticipating multi-class time-sharing disciplines such as Discriminatory Processor Sharing (DPS) and Generalized Processor Sharing (GPS) with general service time distributions. As an application we show that the unifying conservation law can be used to compare the expected unconditional response time under two scheduling disciplines.

## 1 Introduction

The so-called *work-conserving* property is fundamental to single-server (multi-class) systems. Let us consider a single-server queue with  $M$  job classes. Let  $U_j(t)$  be the unfinished work at time  $t$  of class- $j$  jobs,  $j = 1, \dots, M$ , and let  $U(t) = \sum_{j=1}^M U_j(t)$  denote the total unfinished work in the system. The unfinished work in the system,  $U(t)$ , is a function that has vertical jumps at arrival epochs equal in size to the corresponding service requirements of the job and remains constant when it hits the horizontal axis. We say that the scheduling discipline is work-conserving if  $U(t)$  decreases at rate  $1(sec/sec)$  whenever  $U(t) > 0$ . A sample path argument shows that the unfinished work in the system,  $U(t)$ , is the same regardless of the work-conserving scheduling discipline being deployed.

In this paper we focus on conservation laws for the time average unfinished work. We refer to Green and Stidham [11] and Sigman [23] for the derivation and application of sample-path conservation laws. Let  $\bar{U}_j$  denote the time average unfinished work of class  $j$ ,  $j = 1, \dots, M$ . The work-conservation property implies that the total time average unfinished work,

$$\bar{U} = \sum_{j=1}^M \bar{U}_j, \quad (1)$$

is a constant that depends only on the inter-arrival and service time distributions, and its value is independent of how the server's capacity is shared among the jobs of the various classes.

The work-conserving property has led to the development of so-called work-conservation laws. In the case of a single-class queue, Kleinrock [16, Section 4] proved that the expected conditional response time must satisfy an integral equation. Kleinrock's original result was obtained for the subset of non-anticipating scheduling disciplines. A scheduling discipline is said to be non-anticipating if the scheduling decision is independent of the actual service requirements of the jobs. O'Donovan [17] generalized this result by deriving a conservation law for the set of scheduling disciplines that are anticipating, that is, disciplines that may use information on the (remaining) service time when deciding which job will be served. For the multi-class case, the work-conserving property allows to obtain a linear relation that the expected unconditional response times of the various classes must satisfy. Such a linear relation has been obtained for (i) non-anticipating disciplines with exponential

service time distributions [4] and (ii) non-preemptive non-anticipating disciplines with general service time distributions [16, Section 3.4]. For more information on work-conservation laws we refer to the textbooks: Gelenbe and Mitrani [10, Chapter 6], Heyman and Sobel [13, Sections 11.4-5], Wolff [26, Chapter 10] and Baccelli and Brémaud [2, Section 3.2].

The application of work-conservation laws has proven extremely successful in the design of optimal control policies of queueing systems. For instance it has led to the development of the Achievable Region approach, see seminal work by Coffman and Mitrani [4] and Federgruen and Groenevelt [8], Shantikumar and Yao [22], Dacre, Glazebrook and Niño-Mora [6] and Green and Stidham [11].

In this paper, we derive a conservation law for a multi-class  $GI/GI/1$  queue with a general work-conserving scheduling discipline and general service time distributions. Provided the scheduling discipline is work-conserving, the scheduling discipline may be anticipating or non-anticipating, preemptive or non-preemptive. We will further show that already existing conservation laws for multi-class and single-class queues can be obtained as particular cases of our work-conservation law. Thus, the conservation law developed provides a unifying view on the already existing conservation laws for single-server queues. We note that our unifying conservation law covers popular multi-class disciplines such as Discriminatory Processor Sharing (DPS) and Generalized Processor Sharing (GPS) with general service time distributions, a case which was not covered by existing conservation laws. It is worthwhile to note that for the case of a single-server single-class queue with an anticipating service discipline, we obtain an alternative (equivalent) expression for the conservation law that O'Donovan [17] developed.

The remainder of the paper is organized as follows. In Section 2 we introduce the notation and assumptions. In Section 3 we develop the unifying conservation law. In Section 4 we use the new conservation law to compare the expected unconditional response time of a non-anticipating discipline with that of an anticipating discipline that favors short jobs. In Section 5 we show that existing conservation laws can be obtained as a particular case of our conservation law.

## 2 Notation and Assumptions

Throughout the paper we consider a  $GI/GI/1$  queue operating under a work-conserving discipline. Thus we have a single server with identically distributed inter-arrivals times and identically distributed service times. We assume that these random variables are mutually independent. Since the scheduling discipline is assumed to be work-conserving, the total unfinished work in the system at time  $t$ ,  $U(t)$ , is independent of the discipline being deployed.

Let  $A$  denote the inter-arrival time distribution and let  $\lambda = 1/E[A]$  be the mean arrival rate of jobs. With probability  $p_j$  an arrival is a class- $j$  job (independent of the inter-arrival times and classification of previous jobs). Let  $\lambda_j = \lambda p_j$  be the mean arrival rate of class- $j$  jobs. The service time distribution of class  $j$  is denoted by  $F_j(\cdot)$ , and let  $\bar{F}_j(\cdot) = 1 - F_j(\cdot)$  be its complementary distribution. Let  $E[X_j]$  and  $E[X_j^2]$ ,  $j = 1, \dots, M$ , denote the first and second moments of the service time distributions. The load of class  $j$  is given by  $\rho_j = \lambda_j E[X_j]$ , and the total load is  $\rho = \lambda \sum_{j=1}^M p_j E[X_j] = \sum_{j=1}^M \rho_j$ . We assume to be in the stable regime, i.e.,  $\rho < 1$ . As a direct consequence of the renewal assumption, the system regenerates itself at the beginning of each busy period. Throughout the paper we assume that the service time distributions of the various classes have a finite second moment, i.e.,  $E[X_j^2] < \infty$ ,  $j = 1, \dots, M$ . This ensures that, under any work-conserving discipline, the expected unfinished work at both arrival epochs (denoted by  $\bar{V}$ ) and random epochs (denoted by  $\bar{U}$ ) is finite, see [14] and [5, Section II.5.6], respectively. In the context of single-class systems, the subscript denoting the class will be dropped from all variables.

Recall that the total expected unfinished work at a random epoch,  $\bar{U}$ , is independent of the scheduling discipline, hence we have that  $\bar{U} = \bar{U}^{FCFS}$ , where  $\bar{U}^{FCFS}$  denotes the expected unfinished work under the First Come First Serve policy. In the case of Poisson arrivals, by the Pollaczek-Khinchin formula we get

$$\bar{U} = \bar{U}^{FCFS} = \frac{\sum_{j=1}^M \lambda_j E[X_j^2]}{2(1 - \rho)}. \quad (2)$$

Let  $T_j(x)$  be the expected conditional response time of a class- $j$  job with service time  $x$ . In addition, and bearing in mind the analysis of anticipating disciplines, let  $T_j(u; x)$  denote the expected conditional time that a class- $j$  job with total service time  $x$  spends in the system in order to obtain  $u$  units of service,  $u \leq x$ . In particular  $T_j(x; x) = T_j(x)$  denotes the regular expected conditional

response time. We note that for the set of disciplines that are non-anticipating,  $T_j(u; x) = T_j(u)$ , for all  $u \leq x$ . We denote by  $\bar{T}_j$  the expected unconditional response time of class- $j$  jobs, that is,  $\bar{T}_j = \int_{x=0}^{\infty} T_j(x) dF_j(x)$ . In the analysis we make the following assumption.

**Assumption 1** *The function  $T_j(u; x)$ ,  $j = 1, \dots, M$ ,  $u \leq x$ , has a continuous partial derivative with respect to  $x$ .*

This assumption does not seem very restrictive. For instance, for non-anticipating disciplines we have  $T_j(u; x) = T_j(u)$ , and hence  $\frac{\partial T_j(u; x)}{\partial x} = 0$  for all  $u \leq x$ . For SRPT, which is the most popular anticipating discipline, an expression for  $T(u; x)$  was provided in [17]. Taking the derivative with respect to  $x$ , it is easy to see that a sufficient condition for  $T^{SRPT}(u; x)$  to have a continuous derivative with respect to  $x$  is that the service time distribution is a continuous function from the right.

### 3 A Unifying Conservation Law

In this section we state the main result of the paper. In Theorem 1 we develop a conservation law for a  $GI/GI/1$  multi-class system operating under a general scheduling discipline and with general service time distributions. In spite of the apparent simplicity of the derivation, our result generalizes already existing conservation laws and in Section 5 we will show that previous laws can be obtained as particular cases of our unifying law.

Before stating our main result, we briefly mention the most important set of disciplines that Theorem 1 covers. For single-class systems, the set of non-anticipating scheduling disciplines includes among others FCFS, Processor-Sharing (PS), Foreground-Background (FB)<sup>1</sup> and Last Come First Served (LCFS). Important examples of anticipating disciplines are the Shortest Remaining Processing Time (SRPT) [21, 19], Shortest Job First and Fair Sojourn Protocol (FSP) [9]. In multi-class systems, the most popular disciplines are non-anticipating, for example Generalized Processor Sharing (GPS) [18, 24], Discriminatory Processor Sharing (DPS) [15, 7] and Priority Disciplines.

**Theorem 1** *Consider a  $GI/GI/1$  multi-class queue under a work-conserving scheduling discipline. Provided that Assumption 1 is satisfied, the expected conditional response times of the various classes satisfy*

$$\sum_{j=1}^M \lambda_j \int_{x=0}^{\infty} \bar{F}_j(x) \left( T_j(x) + \int_{u=0}^x \frac{\partial T_j(u; x)}{\partial x} du \right) dx = \bar{U}, \quad (3)$$

where  $\bar{U}$  is the total expected unfinished work in the system which depends only on the inter-arrival and service time distributions. If in addition we assume that class- $j$  jobs,  $j = 1, \dots, M$ , arrive according to a Poisson process, then

$$\bar{U} = \frac{\sum_{j=1}^M \lambda_j E[X_j^2]}{2(1 - \rho)}.$$

**Proof.** We consider class  $j$ ,  $j = 1, \dots, M$ . Let  $W_j^i$ ,  $i = 1, 2, \dots$ , be the cumulative burden of the  $i$ -th class- $j$  job to the unfinished work of the system (shaded region in Figure 1). Formally,  $W_j^i := \int_{t=0}^{d_j^i} R_j^i(a_j^i + t) dt$ , where  $a_j^i$  is the arrival time of the  $i$ -th class- $j$  job,  $d_j^i$  its response time and  $R_j^i(t)$  its remaining service requirement at time  $t$ . In particular,  $R_j^i(a_j^i) = x$  is the total service requirement of this job and  $R_j^i(a_j^i + d_j^i) = 0$ . Since  $\rho < 1$ , the busy period has a finite length with probability 1. Furthermore, since the superposed arrival process is a renewal process, the begin points of the busy periods constitute regeneration points and, as a consequence, the sequence  $\{W_j^n\}_{n=1}^{\infty}$  is a regenerative process with finite cycle lengths. Hence, the process  $\{W_j^n\}_{n=1}^{\infty}$  is stationary and ergodic. Applying the Palm inversion formula [3, 2] to the unfinished work of class  $j$ , we obtain  $\bar{U}_j = \lambda_j E[W_j]$ . This equation can also be obtained by the generalized Little's law (see Brumelle [3]). Informally, the equation  $\bar{U}_j = \lambda_j E[W_j]$  states that the time average of a stochastic process ( $\bar{U}_j$ ) is equal to the arrival rate ( $\lambda_j$ ) times the average contribution of each job to the process ( $E[W_j]$ ).

We derive now the value of  $E[W_j]$ . Let  $\tau_j^i(u; x)$  denote the amount of time that the  $i$ -th class- $j$  job, which has a size equal to  $x$ , needs to obtain  $u \leq x$  units of service. In particular we note that

<sup>1</sup>Also known as Least Attained Service (LAS)

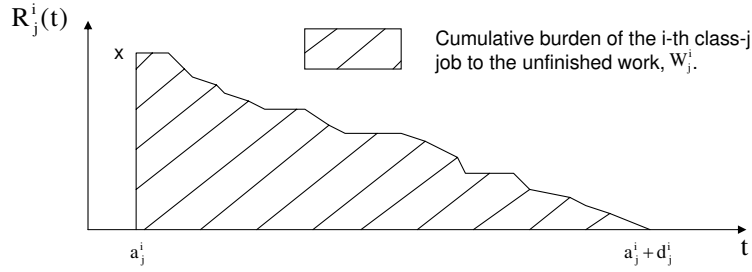


Figure 1: Cumulative contribution over time of the  $i$ -th class- $j$  job with size  $x$  on the unfinished work.

$\tau_j^i(x) := \tau_j^i(x; x) = d_j^i$  is equal to the response time of the  $i$ -th class- $j$  job. Note that  $E[\tau_j^i(u; x)] = T_j^i(u; x)$ . Then

$$\begin{aligned} E[W_j] &= E\left[\int_{x=0}^{\infty} \int_{t=0}^{\tau_j^i(x)} R_j^i(a_j^i + t) dt dF_j(x)\right] \\ &= E\left[\int_{x=0}^{\infty} \int_{u=0}^x \tau_j^i(x - u; x) du dF_j(x)\right]. \end{aligned}$$

This corresponds to integrating the shaded area in Figure 1 either horizontally (first equation) or vertically (second equation). By a simple change of variables and interchanging the order of the integrals we obtain

$$\begin{aligned} E[W_j] &= E\left[\int_{x=0}^{\infty} \int_{u=0}^x \tau_j^i(u; x) du dF_j(x)\right] \\ &= \int_{x=0}^{\infty} \int_{u=0}^x T_j^i(u; x) du dF_j(x). \end{aligned}$$

By Assumption 1 the function  $T_j^i(u; x)$  has a continuous partial derivative with respect to  $x$ , hence by partial integration we obtain

$$E[W_j] = -\bar{F}_j(x) \int_{u=0}^x T_j^i(u; x) du \Big|_{x=0}^{\infty} + \int_{x=0}^{\infty} \bar{F}_j(x) \left( T_j^i(x; x) + \int_{u=0}^x \frac{\partial T_j^i(u, x)}{\partial x} du \right) dx. \quad (4)$$

The second moment of the service time distribution satisfies  $E[X_j^2] = \int_{x=0}^{\infty} x^2 dF_j(x) = \int_{x=0}^{\infty} 2x \bar{F}_j(x) dx$ . By partial integration we obtain

$$\int_{x=0}^{\infty} x^2 dF_j(x) = \int_{x=0}^{\infty} 2x \bar{F}_j(x) dx + \lim_{x \rightarrow \infty} x^2 \bar{F}_j(x).$$

Since the service time distribution has a finite second moment we get that  $\lim_{x \rightarrow \infty} x^2 \bar{F}_j(x) = 0$ . Let  $B(y)$  be the expected length of the busy period initiated by a job of size  $y$ . Note that  $B(\cdot)$  refers to the regular busy period in the  $GI/GI/1$  (and not the sub-busy period of a particular class). Then it follows that

$$T_j^i(u; x) \leq T_j^i(x; x) \leq B(x + \bar{V}),$$

where  $\bar{V}$  is the expected total unfinished work in the system at an arrival epoch. Let  $L$  be a constant such that  $E[\min(A, L)] > \sum_{j=1}^M p_j E[X_j]$ , that is, we truncate the inter-arrival times such that the system is still stable. Using [12, Theorem III.3.1] and Wald's Lemma it is easy to show that

$$E[B(y)] \leq \frac{y + L}{\frac{E[\min(A, L)]}{E[A]} - \rho}.$$

Together with the fact that  $\lim_{x \rightarrow \infty} x^2 \bar{F}_j(x) = 0$  we obtain

$$\lim_{x \rightarrow \infty} \bar{F}_j(x) \int_{u=0}^x T_j^i(u; x) du \leq \lim_{x \rightarrow \infty} \bar{F}_j(x) \frac{x^2 + (\bar{V} + L)x}{\frac{E[\min(A, L)]}{E[A]} - \rho} = 0. \quad (5)$$

Thus, equation (4) becomes

$$E[W_j] = \int_{x=0}^{\infty} \bar{F}_j(x) \left( T_j(x) + \int_{u=0}^x \frac{\partial T_j(u; x)}{\partial x} du \right) dx.$$

Recall that  $\bar{U}_j = \lambda_j E[W_j]$ ,  $j = 1, \dots, M$ . Now the result follows after summing over all the classes and invoking the work-conservation property (1). When the arrival processes of the various classes are Poisson, the time average total unfinished work is given by equation (2). ■

In the context of a multi-class queue, the most important disciplines are non-anticipating. In the following corollary we specialize Theorem 1 to this set of disciplines.

**Corollary 1** *In addition to the conditions of Theorem 1, assume that the scheduling discipline is non-anticipating. Then*

$$\sum_{j=1}^M \lambda_j \int_{x=0}^{\infty} \bar{F}_j(x) T_j(x) dx = \bar{U}. \quad (6)$$

where  $\bar{U}$  is the expected unfinished work in the system which depends only on the inter-arrival and service time distributions. If the arrival processes of the various classes are Poisson, then  $\bar{U} = \frac{\sum_{j=1}^M \lambda_j E[X_j^2]}{2(1-\rho)}$ .

**Proof.** The proof follows readily from Theorem 1 after noting that for a non-anticipating discipline we have  $T_j(u; x) = T_j(u)$  and hence  $\frac{\partial T_j(u; x)}{\partial x} = 0$ , for all  $0 \leq u \leq x$ . ■

We note that Corollary 1 covers important multi-class disciplines such as DPS and GPS. For instance, Corollary 1 was used in [1] to study the asymptotics of the expected conditional response time in a DPS queue when the service time grows to infinite.

## 4 Performance of Anticipating Disciplines

In this section we show that the conservation law derived in Theorem 1 may be useful in evaluating the effect that deploying an anticipating service time discipline has on the expected unconditional response time.

**Proposition 1** *Consider a single-class queue with exponentially distributed service times. Let  $\pi_1$  be a work-conserving non-anticipating discipline and let  $\pi_2$  be a work-conserving anticipating discipline such that for all  $0 \leq u \leq x$*

$$\frac{\partial T^{\pi_2}(u, x)}{\partial x} \geq 0. \quad (7)$$

Then

$$\bar{T}^{\pi_1} \geq \bar{T}^{\pi_2},$$

where  $\bar{T}^{\pi_1}$  and  $\bar{T}^{\pi_2}$  denote the expected unconditional response time obtained under  $\pi_1$  and  $\pi_2$  respectively.

**Proof.** Since  $\pi_1$  is non-anticipating and noting that the exponential assumptions implies that  $\bar{F}(x)dx = E[X]dF(x)$ , from Corollary 1 we get

$$\begin{aligned} \bar{U} &= \lambda \int_{x=0}^{\infty} \bar{F}(x) T^{\pi_1}(x) dx = \lambda E[X] \int_{x=0}^{\infty} T^{\pi_1}(x) dF(x) \\ &= \rho \bar{T}^{\pi_1}. \end{aligned}$$

In the case of  $\pi_2$ , the conservation law (3) can be written as

$$\bar{U} = \rho \bar{T}^{\pi_2} + \lambda \int_{x=0}^{\infty} \bar{F}(x) \int_{u=0}^x \frac{\partial T^{\pi_2}(u; x)}{\partial x} du dx.$$

Taking the difference we obtain

$$\rho \left( \overline{T}^{\pi_1} - \overline{T}^{\pi_2} \right) = \lambda \int_{x=0}^{\infty} \overline{F}(x) \int_{u=0}^x \frac{\partial T^{\pi_2}(u; x)}{\partial x} du dx.$$

The final result follows as a direct consequence of inequality (7). ■

Proposition 1 shows that  $\frac{\partial T^{\pi_2}(u, x)}{\partial x} \geq 0$  for all  $0 \leq u \leq x$  is a sufficient condition for an anticipating discipline  $\pi_2$  to have a smaller expected unconditional response time compared to any non-anticipating discipline, i.e., discipline  $\pi_2$  discriminates against large jobs. The set of scheduling disciplines that satisfy equation (7) is large. Using O'Donovan's expression [17] it is easy to verify that relation (7) is indeed satisfied for SRPT. We expect that policies such as FSP, Shortest-Job-First (preemptive and non-preemptive) and SMART [25] will also satisfy (7) since they all discriminate against large jobs.

In future work we plan to investigate whether a precise evaluation of the term  $\frac{\partial T(u, x)}{\partial x}$  allows one to obtain precise bounds on the expected unconditional response time of anticipating disciplines.

## 5 Relation with previously obtained Conservation Laws

The derivations of existing conservation laws for the single-class and multi-class systems were obtained by different approaches. In this section we show that these conservation laws can all be obtained as a particular case of the unifying conservation law as stated in Theorem 1.

### 5.1 Single-class queue

#### 5.1.1 Non-anticipating discipline

It is straightforward to derive a work-conservation law for a single-class, non-anticipating scheduling discipline. Setting  $M = 1$  in Corollary 1 we obtain

$$\overline{U} = \lambda \int_{x=0}^{\infty} \overline{F}(x) T(x) dx,$$

which is precisely the conservation law for the single-class system as stated in [17],[16, Section 4.9] and [2, Section 2.3].

#### 5.1.2 Anticipating discipline

In this section we show how the conservation law for a general anticipating scheduling discipline obtained by O'Donovan [17, equation (9)] can be retrieved from the unifying conservation law. Setting  $M = 1$  in Theorem 1 we obtain

$$\begin{aligned} \overline{U} &= \lambda \int_{x=0}^{\infty} \overline{F}(x) \left( T(x) + \int_{u=0}^x \frac{\partial T(u; x)}{\partial x} du \right) dx \\ &= \lambda \int_{x=0}^{\infty} \overline{F}(x) \frac{d}{dx} \left( \int_{u=0}^x T(u; x) du \right) dx \\ &= \lambda \int_{x=0}^{\infty} \int_{u=0}^x T(u; x) du dF(x), \end{aligned}$$

where the last inequality is obtained by partial integration and (5). Making the change of variable  $u = x - r$ , and integrating by parts the inner integral we get

$$\overline{U} = \lambda \int_{x=0}^{\infty} \int_{r=0}^x (-r \partial_r T(x - r; x)) dF(x),$$

where the notation  $\partial_r$  denotes a differential with respect to the variable  $r$ . Now by interchanging the order of the integrals we obtain

$$\overline{U} = \lambda \int_{r=0}^{\infty} r \int_{x=r}^{\infty} (-\partial_r T(x - r; x)) dF(x),$$

which is precisely the conservation law as stated in O'Donovan [17, equation (9)]. Note that due to the definition of  $T(x - r; x)$ , we have  $-\partial_r T(x - r; x) \geq 0$ .

### 5.1.3 Non-preemptive anticipating discipline

Baccelli and Brémaud [2, p. 163] develop a conservation law for the subset of anticipating scheduling disciplines that are non-preemptive (for example non-preemptive Shortest Job First). In this case the expected conditional sojourn time can be expressed as  $T(u; x) = u + \bar{V}(x)$ , for all  $0 \leq u \leq x$ , where  $\bar{V}(x)$  denotes the expected waiting time in the queue for a job of size  $x \geq 0$ , i.e., the elapsed time between the arrival time and the time at which the job starts to be served. Setting  $M = 1$  in (3) and substituting the expression for  $T(u; x)$ , we get

$$\begin{aligned} \bar{U} &= \lambda \int_{x=0}^{\infty} \bar{F}(x) \left( x + \bar{V}(x) + \int_{u=0}^x \frac{d\bar{V}(x)}{dx} du \right) dx \\ &= \frac{1}{2} \lambda E[X^2] + \int_{x=0}^{\infty} \bar{F}(x) \frac{d(x\bar{V}(x))}{dx} dx \\ &= \frac{1}{2} \lambda E[X^2] + \int_{x=0}^{\infty} x \bar{V}(x) dF(x), \end{aligned}$$

where the last equality is obtained by partial integration. This expression is the same as the conservation law stated in [2, p. 163].

## 5.2 Multi-class queue

### 5.2.1 Non-anticipating discipline and exponential service time distributions

A conservation law for non-anticipating disciplines with exponential service times was obtained in [4] (see also [10, Section 6.2]). From the assumption of exponential service time distributions, it follows that for  $j = 1, \dots, M$ ,

$$\bar{F}_j(x) dx = E[X_j] dF_j(x). \quad (8)$$

Plugging (8) into (6) we get

$$\bar{U} = \sum_{j=1}^M \lambda_j E[X_j] \int_{x=0}^{\infty} T_j(x) dF_j(x) = \sum_{j=1}^M \rho_j \bar{T}_j,$$

which is precisely the conservation law as stated in [4] and [10, Section 6.2].

### 5.2.2 Non-preemptive non-anticipating discipline and general service time distributions

Finally, let us consider a non-preemptive non-anticipating scheduling discipline with general service time distributions. Such a policy specifies which class to serve whenever a job leaves the system. For example, the server may visit the classes in some order: fixed, random or following a priority rule. The policy is non-anticipating if this decision is made based only on the past history and current state of the system. Within one class, the job that will be served can be determined with a non-preemptive non-anticipating policy like FCFS or Random Order of Service. Under a non-preemptive non-anticipating scheduling discipline the expected conditional response time for a class- $j$  job of size  $x$  satisfies  $T_j(x) = x + \bar{V}_j$ , for all  $x \geq 0$ , where  $\bar{V}_j$  denotes the expected waiting time in the queue for a class- $j$  job (which is independent of  $x$ ), i.e., the elapsed time between the arrival time and the time at which the job starts to obtain service. Substituting this into (6) we obtain

$$\begin{aligned} \bar{U} &= \sum_{j=1}^M \lambda_j \int_{x=0}^{\infty} T_j(x) \bar{F}_j(x) dx = \frac{1}{2} \sum_{j=1}^M \lambda_j E[X_j^2] + \sum_{j=1}^M \lambda_j \bar{V}_j \int_{x=0}^{\infty} \bar{F}_j(x) dx \\ &= \frac{1}{2} \sum_{j=1}^M \lambda_j E[X_j^2] + \sum_{j=1}^M \rho_j \bar{V}_j, \end{aligned}$$

which is equivalent to the conservation law as stated in [20] and [16, Section 3.4].

## Acknowledgments

The author is thankful to S. Aalto, K.E. Avrachenkov, S.C. Borst, R. Núñez-Queija, M. Nuyens and especially to I.M. Verloop for reading and providing valuable comments on the preliminary versions of the paper. The author is also thankful to A. Sapozhnikov and A.P. Zwart who pointed out reference [12] and to the anonymous referee.

## References

- [1] K.E. Avrachenkov, U. Ayesta, P. Brown, and R. Núñez-Queija. Discriminatory processor sharing revisited. In *Proceedings of IEEE INFOCOM*, 2005.
- [2] F. Baccelli and P. Brémaud. *Elements of Queuing Theory: Palm Martingale Calculus and Stochastic Recurrences*. Springer, 2003.
- [3] S.L. Brumelle. On the relation between customer and time average in queues. *Journal of Applied Probability*, 2:508–520, 1971.
- [4] E.G. Coffman and I. Mitrani. A characterization of waiting time performance realizable by single-server queues. *Operations Research*, 3(28):810–821, 1980.
- [5] J.W. Cohen. *The single server queue*. North-Holland, 1982.
- [6] M. Dacre, K. Glazebrook, and J. Niño-Mora. The achievable region approach to the optimal control of stochastic systems. *Journal of the Royal Statistical Society. Series B, Methodological*, 61(4):747–791, 1996.
- [7] G. Fayolle, I. Mitrani, and R. Iasnogorodski. Sharing a processor among many job classes. *Journal of the ACM*, 27(3):519–532, 1980.
- [8] A. Federgruen and H. Groenevelt. Characterization and optimization of achievable performance in general queueing systems. *Operations Research*, 36:733–741, 1988.
- [9] E. Friedman and S. Henderson. Fairness and efficiency in processor sharing protocols to minimize sojourn times. *Proceedings of ACM SIGMETRICS*, pages 229–337, 2003.
- [10] E. Gelenbe and I. Mitrani. *Analysis and Synthesis of Computer Systems*. London: Academic Press, 1980.
- [11] T.C. Green and S. Stidham. Sample-path conservation laws, with application to scheduling queues and fluid systems. *Queueing Systems*, 36:175–199, 2000.
- [12] A. Gut. *Stopped random walks: limit theorems and applications*. Springer-Verlag, 1988.
- [13] D.P. Heyman and M.J. Sobel. *Stochastic Models in operations research, Volume I: Stochastic Processes and operating characteristics*. McGraw-Hill, 1982.
- [14] J.Kiefer and J.Wolfowitz. On the theory of queues with many servers. *Transactions of the American Mathematical Society*, 78(1):1–18, 1955.
- [15] L. Kleinrock. Time-shared systems: A theoretical treatment. *Journal of the ACM*, 14(2):242–261, 1967.
- [16] L. Kleinrock. *Queueing Systems, vol. 2*. John Wiley and Sons, 1976.
- [17] T.M. O’Donovan. Distribution of attained service and residual service in general queueing systems. *Operations Research*, 22:570–575, 1974.
- [18] A.K. Parekh and R.G. Gallager. A generalized processor sharing approach to flow control in integrated services networks: the single-node case. *IEEE/ACM Transactions on Networking*, 1(3):344–357, 1993.



- [19] L.E. Schrage. The queue M/G/1 with feedback to lower priority queues. *Management Science*, 13:466–471, 1967.
- [20] L.E. Schrage. An alternative proof of a conservation law for the queue G/G/1. *Operations Research*, 18:185–187, 1970.
- [21] L.E. Schrage and L.W. Miller. The queue M/G/1 with the shortest remaining processing time discipline. *Operations Research*, 14:670–684, 1966.
- [22] J. Shanthikumar and D. Yao. Multiclass queueing systems: Polymatroidal structure and optimal scheduling control. *Operations Research*, 40(2):293–299, 1992.
- [23] K. Sigman. A note on a sample-path rate conservation law and its relationship with  $H=\lambda G$ . *Advances in Applied Probability*, 23:662–665, 1991.
- [24] M.J.G. van Uitert. *Generalized Processor Sharing Queues*. PhD thesis, Eindhoven University of Technology, 2003.
- [25] A. Wierman, M. Harchol-Balter, and T. Osogami. Nearly insensitive bounds on SMART scheduling. In *Proceedings of ACM SIGMETRICS*, 2005.
- [26] R.W. Wolff. *Stochastic Modeling and the theory of Queues*. Prentice-Hall, 1989.