



WCFS: a new framework for analyzing multiserver systems

Isaac Grosfot¹ · Mor Harchol-Balter¹ · Alan Scheller-Wolf²

Received: 1 March 2022 / Revised: 9 June 2022 / Accepted: 13 June 2022

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2022

Abstract

Multiserver queueing systems are found at the core of a wide variety of practical systems. Many important multiserver models have a previously-unexplained similarity: identical mean response time behavior is empirically observed in the heavy traffic limit. We explain this similarity for the first time. We do so by introducing the work-conserving finite-skip (WCFS) framework, which encompasses a broad class of important models. This class includes the heterogeneous $M/G/k$, the Limited Processor Sharing policy for the $M/G/1$, the Threshold Parallelism model and the Multiserver-Job model under a novel scheduling algorithm. We prove that for all WCFS models, scaled mean response time $E[T](1 - \rho)$ converges to the same value, $E[S^2]/(2E[S])$, in the heavy-traffic limit, which is also the heavy traffic limit for the $M/G/1/FCFS$. Moreover, we prove additively tight bounds on mean response time for the WCFS class, which hold for all load ρ . For each of the four models mentioned above, our bounds are the first known bounds on mean response time.

Keywords Queueing · Response time · Bounds · Heavy traffic · Multiserver · $M/G/k$ · Scheduling

Mathematics Subject Classification 68M20

✉ Isaac Grosfot
igrosfof@cs.cmu.edu
Mor Harchol-Balter
harchol@cs.cmu.edu
Alan Scheller-Wolf
awolf@andrew.cmu.edu

¹ Computer Science Department, Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh, PA 15213, USA

² Tepper School of Business, Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh, PA 15213, USA

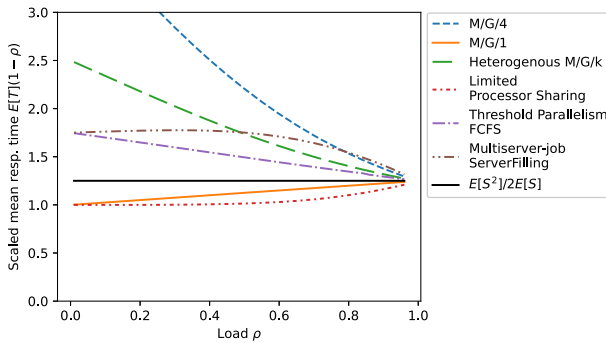


Fig. 1 Scaled mean response time of our four motivating models, as well as the related $M/G/k$ and $M/G/1$ models. Our four motivating models will be further defined in Sect. 3. In each case, the job size distribution S is distributed as Hyperexp($\mu_1 = 2, \mu_2 = \frac{2}{3}, p_1 = \frac{1}{2}$). The black line is $E[T](1 - \rho) = \frac{E[S^2]}{2E[S]}$, the heavy traffic behavior of $M/G/1/FCFS$ and each of our models of interest. 10^9 arrivals simulated. $\rho \in [0, 0.96]$ to ensure accurate results

1 Introduction

Consider the following four queueing models, which are each important, practical models, but which seem very different. We will refer to these models throughout the paper as our *four motivating models*:

- *Heterogeneous $M/G/k$* : A k -server system where servers run at different speeds. Jobs are held at a central queue and served in First-Come-First-Served (FCFS) order when servers become available. If multiple servers are vacant, a server assignment policy such as Fastest Server First is applied.
- *Limited processor sharing*: A single-server system where if at least k jobs are present, the k earliest arrivals each receive an equal fraction of service. If fewer than k jobs are present, the server is split equally among all jobs.
- *Threshold parallelism*: A multiserver system where jobs can run on any number of servers up to some threshold, with perfect speedup. We consider FCFS service, where each job is allocated a number of servers equal to its threshold, as long as servers are available. The final job served may be allocated fewer servers than its threshold.
- *Multiserver-jobs under the ServerFilling policy*: A multiserver system where the jobs are called “multiserver jobs,” because each job requires a fixed number of servers, which it holds concurrently throughout its service. We examine a service policy called *ServerFilling*, which always fills all of the servers if enough jobs are available.

We define these models in more detail in Sect. 3.

We will show that, while our four motivating models appear quite different, their mean response times, $E[T]$, are very similar, especially in the heavy-traffic limit. Specifically, we will show that their behavior in the heavy traffic limit is identical to that of the $M/G/1/FCFS$ model, and in fact the mean response time of each of these

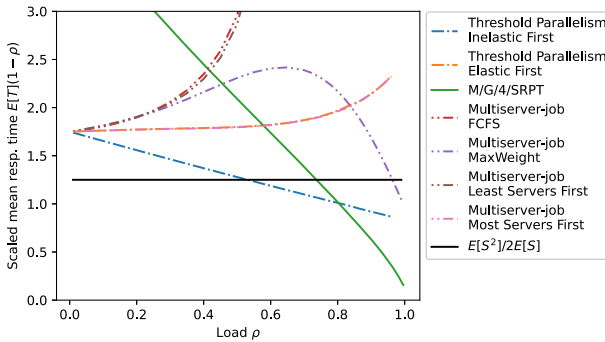


Fig. 2 Scaled mean response time of alternative models and policies. All of these models and policies will be explained in Sect. 6. $S \sim \text{Hyperexp}(\mu_1 = 2, \mu_2 = \frac{2}{3}, p_1 = \frac{1}{2})$. Black line is $E[T](1 - \rho) = \frac{E[S^2]}{2E[S]}$. 10^9 arrivals simulated, $\rho \in [0, 0.96]$ to ensure accurate results, except MaxWeight and $M/G/4/SRPT$: 10^{10} arrivals, $\rho \in [0, 0.99]$

disparate models only differs by an *additive* constant from that of $M/G/1/FCFS$ for all loads, a much stronger result than convergence in heavy traffic.

The similarity of these models is illustrated by Fig. 1, which shows mean response time, $E[T]$, scaled by a factor of $1 - \rho$, to help illustrate the asymptotic behavior in the $\rho \rightarrow 1$ limit. Observe that in each of our models of interest, as well as in the $M/G/1$ and the $M/G/4$, $E[T](1 - \rho)$ converges to $E[S^2]/2E[S]$, the mean of the equilibrium (excess) distribution, where S denotes the job size distribution and $\rho = \lambda E[S] < 1$ is the system load.

This similarity is striking—to see just how notable it is, consider a variety of alternative models and policies shown in Fig. 2. For these alternative models, scaled mean response time either does not converge at all, or converges to a different limit entirely.

This contrast poses an intriguing question:

Why do our four motivating models converge to $M/G/1/FCFS$ in heavy traffic?

To put it another way, we ask what crucial property our four motivating models share, that is not shared by the alternative models in Fig. 2.

To answer this question, we define the “work-conserving finite-skip” framework (WCFS), which applies to a broad class of models. The WCFS class contains our four motivating queueing models, as well many others. We demonstrate that for any model in the WCFS class (which we call a “WCFS model”), if the job size distribution S has bounded expected remaining size, then its scaled mean response time converges to the same heavy traffic limit as the $M/G/1/FCFS$. Specifically, we prove that

Theorem 1 *For any model $\pi \in WCFS$ with bounded expected remaining size¹,*

$$\lim_{\rho \rightarrow 1} E[T^\pi](1 - \rho) = \frac{E[S^2]}{2E[S]}.$$

¹ This assumption is defined in Sect. 2.3.

Theorem 1 follows from an even stronger result: We prove that the difference in mean response time between any WCFS model and $M/G/1/FCFS$ is bounded by an explicit additive constant, that may depend on the specific WCFS model.

Theorem 2 *For any model $\pi \in WCFS$ with bounded expected remaining size,*

$$E[T^\pi] \leq \frac{\rho}{1-\rho} \frac{E[S^2]}{2E[S]} + c_{upper}^\pi$$

$$E[T^\pi] \geq \frac{\rho}{1-\rho} \frac{E[S^2]}{2E[S]} + c_{lower}^\pi$$

for explicit constants c_{upper}^π and c_{lower}^π not dependent on load ρ .

Theorem 2 not only implies Theorem 1, it also guarantees rapid convergence of scaled mean response time to the heavy traffic limit specified in Theorem 1.

In summary, this paper makes the following contributions:

- We define the WCFS framework and our bounded expected remaining size assumption. (Sect. 2)
- We prove that each of the four motivating models is a WCFS model (Sect. 3).
- We discuss prior work on WCFS models (Sect. 4).
- We prove that all WCFS models with bounded expected remaining size have the same scaled mean response time as $M/G/1/FCFS$, and mean response time within an additive constant of $M/G/1/FCFS$ (Sect. 5).
- We empirically validate our results, contrasting heavy traffic behavior of WCFS models and non-WCFS models (Sect. 6).

2 The WCFS framework and WCFS models

In Sects. 2.1 and 2.2, we define the WCFS framework and resulting class of models. In Sect. 2.3, we define our “bounded expected remaining size” assumption. In Sect. 2.4, we define a few more concepts that will be used in the paper.

Job sizes are sampled *i.i.d.* from a job size distribution. Once sampled, job sizes are fixed: we assume preempt-resume service if a job is preempted while in service. Intuitively, the size of a job represents the amount of work associated with the job. Size will be defined in more detail in Sect. 2.1.2.

2.1 WCFS framework and WCFS models

The WCFS framework applies to the class of models with Poisson arrivals at rate λ , which satisfy the following properties:

1. Finite skip (Sect. 2.1.1),
2. Work conserving (Sect. 2.1.2)
3. Non-idling (Sect. 2.1.3).

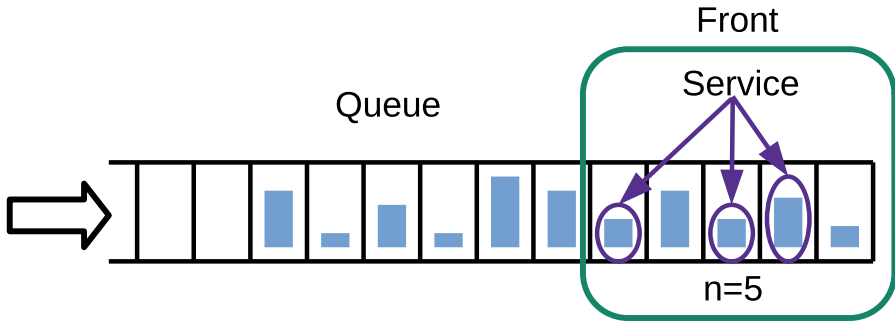


Fig. 3 Diagram of a finite-skip model

2.1.1 Finite skip

We first define the finite-skip property, which defines the class of finite-skip models. Consider the jobs in the system in arrival order. Associated with each finite-skip model, there is a finite parameter n . We partition the jobs in the system into two sets: the (up to) n jobs which arrived longest ago, which we call the *front*, and all other jobs, which we call the *queue*. The *finite-skip property* specifies that, among all of the jobs in the system, the server(s) only serve jobs in the front. In particular, no jobs beyond the first n jobs in arrival order receive any service. Figure 3 shows a generic finite-skip model.

Definition 1 We call the front *full* if at least n jobs are present in the system, and therefore exactly n jobs are at the front.

The intuition behind the term “finite skip” comes from imagining moving through the jobs in the system in arrival order, skipping over some jobs and serving others. In a finite-skip model, only the first n jobs can be served, so only finitely many jobs can be skipped.

2.1.2 Work conserving

Now, we will specify what we mean by “work conserving,” which is a different concept here than in previous work.

First, we normalize the total system capacity to 1, regardless of the number of servers in the system. For instance, in a homogeneous k -server system, we think of each server as serving jobs at rate $1/k$.

Whenever a job is in service, it receives some fraction of the system’s total service capacity, which we call the job’s *service rate*. Let $B(t) \leq 1$ denote the total service rate of all jobs in service at time t , and let B be the stationary total service rate, assuming for now such a quantity exists.

We define a job’s *age* at time t to be the total amount of service the job has received up to time t : a job’s age increases at a rate equal to the job’s service rate whenever the job is in service. Each job has a property called its *size*. When the job’s age reaches its size, the job completes.

In particular, we assume that every job j has a size s_j and a class c_j drawn i.i.d. from some general joint distribution. Let (S, C) be the random variables denoting a job's size and class pair. A job's class is static information known to the scheduler, while a job's size is unknown to the scheduler. For instance, in the Threshold Parallelism model defined in Sect. 3.3, a job's parallelism threshold is its class.

Definition 2 We call the system *maximally busy* if the entire capacity of the system is in use, namely if the total service rate of jobs in service is 1.

We define a finite-skip model to be *work conserving* if whenever the front is full, the system is also maximally busy.

In other words, a finite-skip model is work conserving if, whenever there are at least n jobs in the system, the total service rate is 1.

Now that we have defined a job's size, we can also define the load of the system: $\rho = \lambda E[S]$. Load ρ is the time-average service rate, or equivalently the time-average fraction of capacity in use. Specifically, $\rho = E[B]$. We assume $\rho < 1$ to ensure stability.

2.1.3 Non-idling

We also assume that the total service rate $B(t)$ is bounded away from zero whenever a job is present. Specifically, whenever a job is present, we assume that $B(t) \geq b_{\text{inf}}$, for some constant $b_{\text{inf}} > 0$.

This assumption is key to bounding mean response time under low load. For an example, see the batch-processing system in Sect. 2.2.

2.2 Examples and non-examples

To clarify which models fit within the WCFS framework, we give several examples, both positive and negative.

- *M/G/k/FCFS*: This is a WCFS model with $n = k$.
- *M/G/∞*: This model is not finite skip. All jobs are in service, regardless of the number of jobs in the system: there is no finite bound on the number of jobs in service.
- *M/G/k/SRPT*: In this model, the k jobs with smallest remaining size are served at rate $1/k$. This model is not finite skip because the jobs with smallest remaining size can be arbitrarily far back in the arrival ordering.
- *Multiserver-job model*: Consider a multiserver system with $k = 2$ servers, and where each job requires either 1 or 2 servers. Let the front size $n = 2$.

If jobs are served in FCFS order, with head-of-the-line (HOLB) blocking, this policy is finite skip, but not work conserving. If the front consists of a job requiring 1 server followed by a job requiring 2 servers, under HOLB the system will only utilize one server. In this case, the front is full, because $n = 2$ jobs are present in the system, but the system is not maximally busy.

In contrast, consider a service policy which serves a 2 server job if either of the jobs in the front are 2 server jobs, or else serves each of the 1 server jobs at the

front. This policy is a special case of the ServerFilling policy, depicted in Fig. 1 and defined in general in Sect. 3.4.2. This policy is finite skip and work conserving.

- *Batch-processing M/G/k*: If there are at least k jobs present, the oldest k jobs in the system are each served at rate $\frac{1}{k}$. Otherwise, no service occurs. This model is finite skip and work conserving, but is not non-idling. To see why the non-idling property is necessary for our main results, specifically Theorem 2, one can show that in the $\lambda \rightarrow 0$ limit, response times will grow arbitrarily large in the batch-processing $M/G/k$. To rule out systems where $E[T]$ diverges in the $\lambda \rightarrow 0$ limit, we assume the non-idling property.
- *Red and blue M/G/k*: Imagine an $M/G/k$ with red and blue jobs. Only one color of jobs is allowed to be in service at a time. To determine which jobs to serve, the scheduler counts off jobs in arrival order until it finds k red jobs or k blue jobs and serves all k of the appropriate color (if fewer than k jobs are found for both colors, the system serves the more populous color). This scheduling policy is WCFS with $n = 2k - 1$.

2.3 Bounded expected remaining size: finite rem_{sup}

At a given point in time, let the *state* of a job j consist of its class c_j and its age a_j . Within our WCFS framework, we allow service to be based on the states of the jobs in the front, but not on the number or states of jobs in the queue.

A key assumption we make is that jobs have bounded expected remaining size from an arbitrary state. Let S_c be the job size distribution for jobs of class $c \in C$. We define $\text{rem}_{\text{sup}}(S, C)$ to be the supremum over the expected remaining sizes of jobs, taken over all states:

$$\text{rem}_{\text{sup}}(S, C) := \sup_{c \in C, a \in \mathbb{R}^+} E[S_c - a \mid S_c > a].$$

When size S is independent of class C , or when a model has no class information, we simply write $\text{rem}_{\text{sup}}(S)$.

In this paper, we focus on job size distributions for which $\text{rem}_{\text{sup}}(S, C)$ is finite. To better understand the finite $\text{rem}_{\text{sup}}(S, C)$ assumption, let us walk through a couple of examples. In all of these examples, let us suppose that the class information is independent of the job size distribution S , so we can simply write $\text{rem}_{\text{sup}}(S)$.

Consider a job size distribution S that is hyperexponential:

$$S = \begin{cases} \text{Exp}(\mu_1) & \text{w.p. } p_1 \\ \text{Exp}(\mu_2) & \text{w.p. } p_2 \\ \text{Exp}(\mu_3) & \text{w.p. } p_3 \end{cases}$$

For all ages a , the expected remaining size is bounded:

$$E[S - a \mid S > a] \leq \frac{1}{\min(\mu_1, \mu_2, \mu_3)} = \text{rem}_{\text{sup}}(S).$$

More generally, an arbitrary phase type job size distribution S' must have finite rem_{sup} .

On the other hand, Pareto job size distributions do not have finite rem_{sup} . Let $S'' \sim \text{Pareto}(\alpha = 3, x_{\text{min}} = 1)$, which has finite first and second moments.

$$E[S'' - a \mid S'' > a] = \frac{a}{2}, \quad \forall a \geq 1$$

$$\lim_{a \rightarrow \infty} E[S'' - a \mid S'' > a] = \infty$$

$$\text{rem}_{\text{sup}} = \sup_a E[S'' - a \mid S'' > a] = \infty$$

In general, finite rem_{sup} roughly corresponds to service time having an exponential or sub-exponential tail, though there are some subtleties. For instance, a Weibull distribution with $P(S \geq a) = a^{-k}$ for some $k < 1$ has infinite rem_{sup} , while for $k \geq 1$, rem_{sup} is finite.

As a final example, suppose the WCFS scheduling policy is a known-size policy, such as a policy which serves the job with least remaining size among the n jobs in the front, at rate 1. Because we require that service is based only on the age and class of a job, we model this situation by saying that a job's class is its original size. In this case, $S = C$, and the distribution S_x is simply the constant x . As a result, $\text{rem}_{\text{sup}}(S, C) = \sup(S)$. Therefore, in a known-size setting, rem_{sup} is finite only if S is bounded.

2.4 Work, number, response time

Let the *work* in the system be defined as the sum of the remaining sizes of all jobs in the system. Let $W(t)$ be the total work in the system at time t . Let $W_Q(t)$ and $W_F(t)$ be the work in the queue and the work at the front, respectively, at time t . (We will generally use the subscripts Q and F to denote the queue and the front.) Let W , W_Q , and W_F denote the corresponding time-stationary random variables.

Recall from Sect. 2.1.2 that $B(t)$ is the total service rate at time t . Note that $\frac{d}{dt} W(t) = -B(t)$, except at arrival instants.

Let $N(t)$ be the number of jobs in the system at time t . Note that $N_F(t) = n$ whenever $N(t) \geq n$, because the front is full, and $N_F(t) = N(t)$ otherwise.

Let T be a random variable denoting a job's time-stationary response time: the time from when a job arrives to when it completes.

3 Important WCFS models

Here, we define in more detail the four motivating models mentioned in the introduction and depicted in Fig. 1.

3.1 Heterogeneous $M/G/k$

The heterogeneous $M/G/k/FCFS$ models multiserver systems where servers have different speeds. This scenario commonly arises in datacenters, which are often composed of servers with a wide variety of different types of hardware [1, 2]. In the mobile device setting, the big.LITTLE architecture employs heterogeneous processors to improve battery life [3].

Let each server i have speed $v_i > 0$, scaled so that $\sum_i v_i = 1$. While a job is being served by server i , the job's age increases at a rate of v_i .

If there are multiple servers idle when a job arrives, a server is chosen according to an arbitrary server assignment policy. Jobs may also be migrated between servers when a job completes. We only assume that jobs are served in FCFS order, and that no job is left waiting while a server is idle. Under these assumptions, all assignment policies fit within the WCFS framework.

As an example, in Fig. 1 we show the scaled mean response time of a heterogeneous $M/G/4$ with server speeds 0.4, 0.3, 0.2, 0.1, and the Preemptive Fastest Server First assignment policy.

3.2 Limited processor sharing

The Processor Sharing policy for the $M/G/1$ is of great theoretical interest, and has been extensively studied [4]. However, in real systems, running too many jobs at once causes a significant overhead. A natural remedy is to utilize a policy known as Limited Processor Sharing (LPS) [5–8].

The LPS policy is parameterized by some Multi-Programming Level k . If at most k jobs are present in the system, then the policy is equivalent to Processor Sharing, serving all jobs at an equal rate, with total service rate 1. When more than k jobs are present, the k oldest jobs in FCFS order are each served at rate $1/k$. LPS is a WCFS model with $n = k$.

As an example, in Fig. 1 we show the scaled mean response time of a LPS system with MPL 4.

3.3 Threshold parallelism

In modern datacenters, it is increasingly common for jobs to be parallelizable across a variety of different numbers of servers, where the level of parallelism is chosen by the scheduler [9, 10]. Under Threshold Parallelism, a job j has two characteristics: its size s_j and its parallelism threshold ℓ_j , where ℓ_j is some number of servers. Job j may be parallelized across up to ℓ_j servers, with linear speedup. The pair (s_j, ℓ_j) is sampled i.i.d. from some joint distribution (S, L) . Note that ℓ_j is the class of the job j .

Let k be the total number of servers. Note that $\ell_j \in [1, k]$. If a job j is served on $q \leq \ell_j$ servers, then it receives service rate $\frac{q}{k}$ and will complete after $\frac{ks_j}{q}$ time in service. The number of servers a job is allocated can change over time, correspondingly changing its service rate.

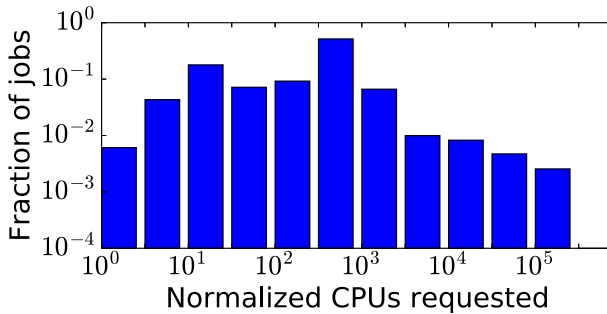


Fig. 4 The distribution of number of CPUs requested in Google’s recently published Borg trace [15]. Number of CPUs is normalized to the size of the smallest request observed, not an absolute value

We focus on the FCFS service policy. Under this policy, jobs are placed into service in arrival order until their total parallelism thresholds sum to at least k , or all jobs are in service. Each job j other than the final job in service is served by ℓ_j servers. The final job in service is served by the remaining servers. Under FCFS service, Threshold Parallelism fits the WCFS framework with $n = k$.

As an example, in Fig. 1 we show the scaled mean response time of a Threshold Parallelism model where the joint distribution (S, L) is $(\text{Exp}(2), 1)$ with probability $\frac{1}{2}$, and $(\text{Exp}(\frac{2}{3}), 4)$ with probability $\frac{1}{2}$, and with FCFS service.

As a comparison, in Fig. 2, we show Threshold Parallelism models with the same joint distribution (S, L) , but with different service policies: “Elastic First,” prioritizing jobs with $L = 1$, and “Inelastic First,” prioritizing jobs with $L = 4$. These policies do not fit within the WCFS framework, because a job may skip over an arbitrary number of jobs.

3.4 Multiserver-jobs under the ServerFilling policy

First, we will describe the Multiserver-Job setting. Then, we will specify the Server-Filling policy.

3.4.1 Multiserver-job setting

When we look at jobs in cloud computing systems [11] and in supercomputing systems [12–14], jobs commonly require an exact number of servers for the entire time the job is in service. To illustrate, in Fig. 4 we show the distribution of the number of CPUs requested by the jobs in Google’s recently published trace of its “Borg” computation cluster [15, 16]. The distribution is highly variable, with jobs requesting anywhere from 1 to 100,000 normalized CPUs².

The Multiserver-Job (MSJ) model is a natural model for these computing systems. In an MSJ model, a job j has two requirements: a number of servers v_j and an amount of time x_j , which are sampled i.i.d. from some joint distribution (V, X) . If job j

² The data was published in a scaled form [15]. We rescale the data so the smallest job in the trace uses one normalized CPU.

requires v_j servers, then it can only be served when exactly v_j servers are allocated to it. The job will complete after x_j time in service.

Let a job j 's size be defined as

$$s_j = \frac{v_j x_j}{k} \quad S = \frac{VX}{k}.$$

There are a wide variety of possible service policies for placing jobs at open servers, including FCFS, MaxWeight, Most Servers First and many others. (We formally define these policies in Sect. 6.) As examples, in Fig. 2, we show the scaled mean response time of Multiserver-Job models under a variety of service policies, where the joint distribution (V, X) is $(1, Exp(\frac{1}{2}))$ with probability $\frac{1}{2}$, and $(4, Exp(\frac{2}{3}))$ with probability $\frac{1}{2}$.

Unfortunately, no existing policies fit within the WCFS framework—all existing policies, including those shown in Fig. 2, are either non-finite-skip, such as Most Servers First, or non-work-conserving, such as FCFS. Correspondingly, in Fig. 2, we see that no existing policy has its scaled mean response time converge to the same limit as $M/G/1/FCFS$.

We therefore define a novel service policy called ServerFilling which yields a WCFS model. The scaled mean response time of this service policy is depicted in Fig. 1, with the same joint distribution (V, X) as the policies shown in Fig. 2.

3.4.2 ServerFilling

For simplicity, we initially define the Server Filling policy for the common situation in computer systems where all jobs require a number of servers which is a power of 2 (V is always a power of 2), and where k is also a power of 2. We discuss generalizations in Sect. 3.4.3.

First, ServerFilling designates a candidate set M , consisting of the minimal prefix (i.e., initial subset) of the jobs in the system in arrival order which collectively require at least k servers. If all jobs in the system collectively require fewer than k servers, then all are served. Note that $|M| \leq k$ because all jobs require at least 1 server.

For instance, if $k = 8$ and the jobs in the system require $[1, 2, 1, 1, 4, 2, 2, 1]$ servers, in arrival order (reading from left to right), then M would consist of the first 5 jobs: $[1, 2, 1, 1, 4]$, which collectively require 9 servers.

Next, the jobs in M are ordered by their server requirements v_j , from largest to smallest, tiebroken by arrival order. Jobs are placed into service in that order until no more servers are available. In our example, jobs requiring 4, 2, 1, and 1 server(s) would be placed into service.

To show that ServerFilling fits within WCFS with $n = k$, we must show that ServerFilling always utilizes all k servers if at least k jobs are in the system.

Lemma 1 *Let M be a set of jobs such that $\sum_{j \in M} v_j \geq k$, where each $v_j = 2^i$ for some i and $k = 2^{i'}$ for some i' . Label the jobs m_1, m_2, \dots in decreasing order of server requirement: $v_{m_1} \geq v_{m_2} \geq \dots$. Then, there exists some index $\ell \leq |M|$ such*

that

$$\sum_{j=1}^{\ell} v_{m_j} = k.$$

Proof Let $\text{REQ}(z)$ count the number of servers required by the first z jobs in this ordering:

$$\text{REQ}(z) = \sum_{j=1}^z v_{m_j}.$$

We want to show that $\text{REQ}(\ell) = k$ for some ℓ . To do so, it suffices to prove that:

$$\text{There exists no index } \ell' \text{ such that both } \text{REQ}(\ell') < k \text{ and } \text{REQ}(\ell' + 1) > k. \quad (1)$$

Equation (1) states that $\text{REQ}(z)$ cannot cross from below k to above k without exactly equaling k . Because $\text{REQ}(0) = 0$ and $\text{REQ}(|M|) \geq k$, $\text{REQ}(\ell)$ must exactly equal k for some ℓ .

To prove (1), let us examine the quantity $k - \text{REQ}(z)$, the number of remaining servers after z jobs have been placed in service. Because all v_j s are powers of 2, $k - \text{REQ}(z)$ carries an important property:

$$k - \text{REQ}(z) \text{ is divisible by } v_{m_{z+1}} \text{ for all } z. \quad (2)$$

We write $a|b$ to indicate that a divides b .

We will prove (2) inductively. For $z = 0$, $k - \text{REQ}(0) = k$. Because k is a power of 2, and v_{m_1} is a power of 2 no greater than k , the base case holds. Next, assume that (2) holds for some index z , meaning that $v_{m_{z+1}} | (k - \text{REQ}(z))$. Note that $\text{REQ}(z + 1) = \text{REQ}(z) + v_{m_{z+1}}$. As a result, $v_{m_{z+1}} | (k - \text{REQ}(z + 1))$. Now, note that $v_{m_{z+2}} | v_{m_{z+1}}$, because both are powers of 2, and $v_{m_{z+2}} \leq v_{m_{z+1}}$. As a result, $v_{m_{z+2}} | (k - \text{REQ}(z + 1))$, completing the proof of (2).

Now, we are ready to prove (1). Assume for contradiction that there does exist such an ℓ' . Then $k - \text{REQ}(\ell') > 0$, and $k - \text{REQ}(\ell' + 1) < 0$. Because $\text{REQ}(\ell' + 1) = \text{REQ}(\ell') + v_{m_{\ell'+1}}$, we therefore know that $v_{m_{\ell'+1}} > k - \text{REQ}(\ell')$. But from (2), we know that $v_{m_{\ell'+1}}$ divides $k - \text{REQ}(\ell')$, which is a contradiction. \square

As a result, the MSJ model under ServerFilling fits within the WCFS framework with $n = k$.

3.4.3 Generalizations of ServerFilling

The ServerFilling policy can be generalized, as long as all server requirements divide k . We describe the corresponding scheduling policy, which we call DivisorFilling, in the full version of this paper: [17, Appendix A].

DivisorFilling is the most general possible WCFS policy for the MSJ setting. If some server requirement does not divide k , then no policy fits within the WCFS framework, because the system is not work conserving if all jobs present require that non-divisible number of servers and more than n jobs are present.

4 Prior work

4.1 M/G/k

4.1.1 Fixed k

In this regime, the best-known bounds on response time either require much stronger assumptions on the job size distribution S than we assume [18], or prove much weaker bounds on mean response time [19, 20].

A paper by [18] bounds mean work in system in the $M/G/k$ to within an additive gap, under the strong assumption that the job size distribution S is bounded. While the paper mostly focuses on the overload regime ($\rho > 1$), their equations (9) and (10) apply in our setting ($\rho < 1$) as well. They couple the multiserver system with a single-server system on the same arrival sequence. They show that

$$0 \leq W^{M/G/k}(t) - W^{M/G/1}(t) \leq k \max_{1 \leq i \leq A(t)} S_i,$$

where $A(t)$ is the number of jobs that have arrived by time t . In the case of a bounded job size distribution S , one can therefore show that

$$0 \leq W^{M/G/k}(t) - W^{M/G/1}(t) \leq k \sup(S). \tag{3}$$

One could then use this workload bound to prove a bound on mean response time in the $M/G/k$.

These bounds are comparable to those in our Lemma 3 when S is bounded, but our bounds require a much weaker assumption on the job size distribution S .

Köllerström [19] proves convergence of queueing time to an exponential distribution in the $GI/GI/k$. Specialized to the $M/G/k$, the result states that in the $\rho \rightarrow 1$ limit, $T_Q^{M/G/k}$ converges to an exponential distribution with mean

$$\frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} - \frac{1}{\lambda} = E[T_Q^{M/G/1}] - \frac{1}{\lambda}.$$

Köllerström [20] improves upon [19] by characterizing the rate of convergence and thereby derives explicit moment bounds. However, unlike prior single-server results [21], these bounds are quite weak. Specialized to the $M/G/k$, [20]’s bounds state that

$$E[T_Q^{M/G/k}] - E[T_Q^{M/G/1}] \geq \frac{C_{\text{lower}}}{(1 - \rho)^{1/2}} \tag{4}$$

$$E[T_Q^{M/G/k}] - E[T_Q^{M/G/1}] \leq \frac{c_{\text{higher}}}{1 - \rho} \quad (5)$$

for constants $c_{\text{lower}}, c_{\text{higher}}$ not dependent on ρ .

The $\Theta(\frac{1}{1-\rho})$ scaling in (5) is especially poor: this bound is too weak to give any explicit bound on the convergence rate of $E[T_Q^{M/G/k}](1 - \rho)$ to the previously established limit of $\frac{E[S^2]}{2E[S]}$.

Our bounds are tighter in that they are constants not depending on ρ , but we assume S has finite rem_{sup} , while [20] merely assumes that S has finite second moment.

4.1.2 Scaling k

Recent work has focused on regimes where both ρ and k scale asymptotically, such as the Halfin–Whitt regime. These results are not directly comparable to ours; they indicate that the limiting behavior in the Halfin–Whitt regime depends in a complex way on the job size distribution S [22–24].

Turning to the more general case of scaling k , in work currently under submission, [25] prove the first bounds on $E[T_Q]$ that scale as $\frac{c}{1-\rho}$ for an explicit constant c and arbitrary joint scaling of k and ρ . Unfortunately, the constant c is enormous, scaling as $10^{450} E[S^3]$. In contrast, we focus on the regime of fixed k , and prove tight and explicit bounds on mean response time. [25] also provide a highly detailed literature review on bounds on $E[T_Q]$ and related measures in the $M/G/k$ and related models.

4.2 Heterogeneous $M/G/k$

4.2.1 Heterogeneous $M/M/k$

Much of the previous work on multiserver models with heterogeneous service rates has focused on the much simpler $M/M/k$ setting, where jobs are memoryless [26–29]. In this model, one can analyze the preemptive Fastest-Server-First policy to derive a natural lower bound on the mean response time of any server assignment policy. One can similarly analyze the preemptive Slowest-Server-First policy to derive an upper bound. These two policies each lead to a single-dimensional birth-death Markov chain, allowing for straightforward analysis [27]. One can think of our bounds as essentially extending these bounds for the $M/M/k$ to the much more complex setting of the $M/G/k$.

4.2.2 Heterogeneous $M/H_m/k$

Van Harten and Sleptchenko [30] primarily study a *homogeneous* multiserver setting with hyperexponential job sizes. However, in their conclusion, they mention that their methods could be extended to a setting with heterogeneous servers, but at the cost of making their Markov chain grow exponentially. This exponential blowup seems inevitable when applying exact Markovian methods to a heterogeneous setting with differentiated jobs.

4.2.3 $M/(M+G)/2$ model

Another intermediate model is the $M/(M+G)/2$ model of [31]. In this model, jobs are not differentiated. Instead, the service time distribution is entirely dependent on the server. Server 1, the first server to be used, has an exponential service time distribution, while server 2 has a general service time distribution. [31] derive an implicit expression for the Laplace–Stieltjes transform of response time in this setting, which they are only able to make explicit when the general service time distribution has rational transform. Subsequent work has fully solved the $M/(M+G)/2$ model, under both FCFS service and related service disciplines [32–34].

Our results are not directly applicable to the $M/(M+G)/2$ setting, because the servers have different distributions of service time, not just different speeds. However, the slow progress on this two-server model illustrates the immense difficulty in solving even the simplest heterogeneous multiserver models. In contrast, our WCFS framework handles both differentiated jobs and an arbitrary number of servers with no additional effort.

4.3 Limited processor sharing

The Limited Processor Sharing policy has been studied by a wide variety of authors [5–8, 35–37], but none bound mean response time for all loads ρ .

4.3.1 Asymptotic regimes

A series of papers by Zhang et al. [7, 35, 36] derive the strongest known results on Limited Processor Sharing in a variety of asymptotic regimes. These authors derive the measure-valued fluid limit [35], the diffusion limit [36] and a steady-state approximation [7]. The most comparable of their results to our work is their steady-state approximation. When specialized to mean response time in the $M/G/1/LPS$, their approximation states that

$$E[T] \approx \frac{E[S]}{1 - \rho} (1 - \rho^k) + \frac{E[S^2]}{2E[S]} \frac{\rho^k}{1 - \rho}$$

They prove that this approximation is accurate in the heavy-traffic limit; they do not provide specific error bounds, but empirically show the approximation performs well at all loads ρ [7]. Our results therefore complement their results by proving concrete error bounds.

4.3.2 State-dependent server speed

To model the behavior of databases, [8] introduce a variant of the Limited Processor Sharing model, where the total server speed is a function of the number of jobs in service. In their setting, server speed increases to a peak, and then slowly declines as more jobs enter service. They derive a two-moment approximation for mean response

time, and use it to derive a heuristic policy for choosing the Multi-Programming Level (MPL). While this two-moment approximation is not known to be tight, it indicates that the optimal MPL for minimizing mean response time may be significantly larger than the service-rate-maximizing MPL, if job size variability is large and load is not too high.

Using our WCFS framework, it is possible to derive bounds on mean response time for the state-dependent server speeds setting. For MPL parameters less than or equal to the service-rate-maximizing MPL, both our upper and lower bounds apply, while if the MPL parameter is greater than the service-rate-maximizing MPL, only our upper bounds apply, because the system only partially fulfills our definition of work conservation.

Subsequently, [6] derive the Laplace–Stieltjes transform of response time in the LPS model with state-dependent server speed, under phase-type job sizes. Unfortunately, the transform takes the form of a complicated matrix equation, making it difficult to derive general insights across general job size distributions. Instead, the authors numerically invert the Laplace transform for a handful of specific distributions to derive empirical insights.

4.4 Threshold parallelism

Jobs with “speedup functions” are common in Machine Learning and other highly parallel computing settings. A job’s speedup function specifies the degree to which it can be parallelized. In [38–40], the authors study optimal allocation policies of servers to jobs when the arriving jobs have different speedup function. In many cases, a job’s speedup function takes the form of a “threshold” function: here the job receives perfect (linear) speedup up to some threshold number of servers and receives no additional speedup beyond that number of servers. We refer to this as the Threshold Parallelism model.

While understanding the response time in systems where jobs have speedup functions is generally intractable, [40] were able to approximately analyze response time in the case where every job is either “inelastic,” with parallelism threshold 1, or “elastic,” with parallelism threshold k . They also assume that inelastic jobs have size distributed as $\text{Exp}(\mu_I)$, and elastic jobs have size distributed as $\text{Exp}(\mu_E)$, with sizes unknown to the scheduler. They focus on two preemptive-priority service policies for this setting: Inelastic First (IF) and Elastic First (EF). In this setting, they approximate the mean response time of EF and IF within 1% error by using a combination of the Busy-Period Transitions technique and Matrix-Analytic methods to evaluate their multidimensional Markov chain.

The Threshold Parallelism model in our paper is far broader than that in the prior literature, and our bounds are tighter in the heavy-traffic limit.

4.5 Multiserver jobs

The Multiserver-Job model has been extensively studied, in both practical [12–14] and theoretical settings [11, 41–47]. It captures the common scenario in datacenters

and supercomputing where each job requires a fixed number of servers in order to run. Characterizing the stability region of policies in this model is already a challenging problem, and there were no bounds on mean response time for any scheduling policy, prior to our bound on ServerFilling.

4.5.1 FCFS scheduling

The most natural policy is FCFS, where the oldest jobs are placed into service until a job requires more servers than remain, at which point the queue is blocked. Therefore, the FCFS policy can leave a large number of servers idle even when many jobs are present. As a result, FCFS does not in general achieve an optimal stability region. Even worse, deriving the stability region of FCFS is an open problem, and has only been solved in a few special cases [41, 42].

One technique that may be useful for characterizing this stability region is the *saturated system* approach [48, 49]. The saturated system is a system in which additional jobs are always available, so the front is always full, only the composition of jobs in the front varies. The completion rate of the saturated system exactly matches the stability region of the equivalent open system, under a wide variety of arrival processes. Unfortunately, analyzing the general Multiserver-Job FCFS saturated system seems intractable.

Given the difficulty of proving results under FCFS scheduling, finding policies with better theoretical guarantees, such as ServerFilling, is desirable.

4.5.2 MaxWeight scheduling

One natural throughput-optimal policy is the MaxWeight policy [11]. Here, jobs are divided into classes based on their server requirements, with $N_i(t)$ denoting the number of jobs requiring i servers in the system at time t . Let the set $Z(t)$ denote all possible packings of jobs at time t onto servers. Let $z \in Z(t)$ be a particular packing, where z_i denotes the number of jobs requiring i servers that are served by packing z .

The MaxWeight service policy picks the packing z which maximizes

$$\max_z \sum_i N_i(t) z_i.$$

For example, if there are many jobs requiring 3 servers, we want to pick a packing that serves many 3-server jobs. While MaxWeight is throughput optimal, it is very computationally intensive to implement, requiring the scheduler to solve an NP-hard optimization problem whenever a job arrives or departs. For comparison, ServerFilling is also throughput-optimal given our assumptions on the server requirements V , but it is far computationally simpler, requiring approximately linear time as a function of k . Moreover, no bounds on mean response time are known for MaxWeight, due in part to its high complexity.

4.5.3 Nonpreemptive scheduling

In certain practical settings such as supercomputing, a nonpreemptive service policy is preferred. In such settings, a backfilling policy such as EASY backfilling or conservative backfilling is often used [12–14]. These start by serving jobs in FCFS order, until a job is reached that requires more servers than remain. At this point, jobs further back in the queue that require fewer servers are scheduled, but only if they will not delay older jobs, based on user-provided service time upper bounds. While these policies are popular in practice little is known about them theoretically, including their response time characteristics.

Finding any nonpreemptive throughput-optimal policy is a challenging problem. Several such policies have been designed [44, 45, 47], typically by slowly shifting between different server configurations to alleviate overhead. Because such policies can have very large renewal times, many jobs can back up while the system is in a low-efficiency configuration, which can empirically lead to very high mean response times. However, no theoretical mean response time analysis exists for any policy in the Multiserver-Job setting. As a result, there is no good baseline policy to compare against novel policies. Our bounds on the mean response time of ServerFilling can serve as such a baseline, albeit in the more permissive setting of preemptive scheduling.

5 Theorems and proofs

We perform a heavy traffic analysis within our WCFS framework, assuming finite $\text{rem}_{\text{sup}}(S, C)$. Specifically, we prove that the scaled mean response time of any WCFS model converges to the same constant as an $M/G/1/FCFS$:

Theorem 1 (Heavy Traffic response time) *For any model $\pi \in WCFS$, if $\text{rem}_{\text{sup}}(S, C)$ is finite,*

$$\lim_{\rho \rightarrow 1} E[T^\pi](1 - \rho) = \frac{E[S^2]}{2E[S]}.$$

To prove Theorem 1, we prove a stronger theorem, tightly and explicitly bounding $E[T^\pi]$ up to an additive constant, for any $\pi \in WCFS$.

Theorem 2 (Explicit response time bounds) *For any model $\pi \in WCFS$, if $\text{rem}_{\text{sup}}(S, C)$ is finite,*

$$\begin{aligned} E[T^\pi] &\leq \frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} + c_{\text{upper}}^\pi \\ E[T^\pi] &\geq \frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} + c_{\text{lower}}^\pi \end{aligned}$$

for explicit constants c_{upper}^π and c_{lower}^π not dependent on load ρ .

Proof deferred to Sect. 5.1 From Theorem 2, Theorem 1 follows via a simple rearrangement:

$$\frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} = \frac{\frac{E[S^2]}{2E[S]}}{1 - \rho} - \frac{E[S^2]}{2E[S]}.$$

□

Theorem 2 also implies rapid convergence of scaled mean response time to its limiting constant for any WCFS policy:

Corollary 1 For any model $\pi \in \text{WCFS}$, if $\text{rem}_{\text{sup}}(S, C)$ is finite,

$$E[T^\pi](1 - \rho) = \frac{E[S^2]}{2E[S]} + O(1 - \rho).$$

5.1 Outline of Proof of Theorem 2

We will prove Theorem 2 where

$$\begin{aligned} c_{\text{upper}}^\pi &= (n - 1)\text{rem}_{\text{sup}}(S, C) + \frac{nE[S]}{b_{\text{inf}}}, \\ c_{\text{lower}}^\pi &= -(n - 1)\text{rem}_{\text{sup}}(S, C) + E[S], \end{aligned}$$

where n denotes the size of the front, and where b_{inf} is defined in Sect. 2.1.3.

Our goal is simply to prove the bounds in Theorem 2 for some constants $c_{\text{upper}}^\pi, c_{\text{lower}}^\pi$ independent of ρ ; we have made no effort to optimize these constants, leaving that to future work. Specifically, for three of our four motivating models, the $\frac{n}{b_{\text{inf}}}$ term scales as $O(n^2)$. For these models, this term is unnecessarily loose, and could easily be lowered to an $O(n)$ bound by using a more detailed view.

Our approach is to split response time T into two pieces, queueing time T_Q and front time T_F , and bound the expectation of each separately. We first bound $E[T_Q]$, which forms the bulk of our proof. The two key ideas come from the intuition that a WCFS model behaves like a FCFS $M/G/1$ system. In Theorem 2, we prove that $E[T_Q] = E[W] + c$, for some constant c ; in a WCFS model, jobs progress through the system in essentially FCFS order, and as $\rho \rightarrow 1$ work is completed essentially at rate 1.

In Lemma 3, we prove that $E[W] = E[W^{M/G/1}] + c$, for some constant c . The key idea here is that in a WCFS model, if W is large, work arrives and completes in exactly the same way as in an $M/G/1$. Likewise, if the front is not full, then W cannot be large.

In Lemma 4, we combine Lemmas 2 and 3 to prove that $E[T_Q] = E[T^{M/G/1}] + c$ for some constant c .

In Lemma 5, we prove that work W is indeed stationary with finite mean. This is a technical lemma that rules out pathological scenarios, which is necessary because our WCFS class of models is very general. Lemma 5 is used by both Lemmas 2 and 3.

Finally, in Lemma 6, we bound $E[T_F]$, utilizing Little's law. Combining Lemmas 4 and 6 proves Theorem 2.

5.2 Two views

At several steps in our proof of Theorem 2, we will make use of two different views of the queueing system, corresponding to two different state descriptors:

Omniscient view In the omniscient view the state descriptor consists of the remaining size and class of all jobs in the system; we sample jobs' sizes and classes when the jobs enter the system. For a given system state, work is a deterministic quantity.

Limited view In the limited view, the state descriptor consists of the age and class of the jobs in the front, and the number of jobs in the queue. We sample jobs' classes when they enter the front, and determine whether jobs complete according to the hazard rate of the job size distribution, as the job ages. For a given system state, work is a random variable.

We will make it clear which view of the system we are using in each step of the proof. Generally, the omniscient view is useful when analyzing total work in the system, and the limited view is useful when analyzing work at the front.

5.3 Lemma 2: $E[T_Q]$ and $E[W]$

First, we prove that mean queueing time and mean work are similar:

Lemma 2 (Queueing time and work) *For any model $\pi \in WCFS$, if $\text{rem}_{\text{sup}}(S, C)$ is finite,*

$$E[W] - (n - 1)\text{rem}_{\text{sup}}(S, C) \leq E[T_Q] \leq E[W].$$

Proof Start by writing time in queue T_Q in terms of work in system. Let us consider the omniscient view of the system, so work W is a deterministic quantity given the system state. Consider an arbitrary tagged job j . When j arrives, let $W^A(j)$ be the amount of work j sees in the system. Let $W_F^F(j)$ be the amount of work j sees in the front other than j itself, when j leaves the queue and enters the front. In W_F^F , the subscript F indicates that we are looking at the amount of work at the front, and the superscript F indicates that we are looking at the moment when j enters the front.

Because the model is finite skip, jobs move from the queue to the front in arrival order, so all of the $W^A(j)$ work that was in the system when j arrived is either complete or in the front when j enters the front. As a result, the amount of work which is completed while j is in the queue is exactly $W^A(j) - W_F^F(j)$. Note that if j enters the front upon arrival to the system, $W^A(j) = W_F^F(j)$, and no work is completed while j is in the queue.

While j is in the queue, the front must be full; the system must be maximally busy during this time, completing work at rate 1. Job j is in the queue for $T_Q(j)$ time, so the system must complete $T_Q(j)$ work during that time. We can therefore conclude

that

$$W^A(j) - W_F^F(j) = T_Q(j).$$

Because j is an arbitrary job, we can write $W_F^F(j)$ as W_F^F , a random variable over all jobs that pass through the system. Likewise, $T_Q(j)$ is simply T_Q . Because Poisson arrivals see time averages, $W^A(j) \sim W$, the time-stationary amount of work in the system. Combining these equivalencies, we find that

$$W - W_F^F = T_Q. \tag{6}$$

Note that W is time-stationary, while W_F^F and T_Q are event-stationary.

To rigorously demonstrate (6), we need to prove that the system converges to a stationary distribution, which we prove in Lemma 5.

To give bounds on W_F^F , we switch to the limited view of the system, where the state of the front consists of the classes and ages of the jobs at the front. We have two simple bounds on W_F^F : First, $W_F^F \geq 0$. Next, because $W_F^F(j)$ is the work of at most $n - 1$ jobs, the jobs at the front when a given job enters the front, we know that

$$E[W_F^F] \leq (n - 1)\text{rem}_{\text{sup}}(S, C).$$

Combining these bounds with (6), we can bound $E[T_Q]$ in terms of $E[W]$:

$$E[W] - (n - 1)\text{rem}_{\text{sup}}(S, C) \leq E[T_Q] \leq E[W].$$

□

5.4 Lemma 3: bounding $E[W]$

Lemma 3 (Work bounds) *For any model $\pi \in \text{WCFS}$, if $\text{rem}_{\text{sup}}(S, C)$ is finite,*

$$\frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} \leq E[W] \leq \frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} + (n - 1)\text{rem}_{\text{sup}}(S, C).$$

Proof Consider the stationary random variable W^2 in the omniscient view, so work is a deterministic quantity at a given time on a given sample path. W^2 evolves in two ways: continuous decrease as work is completed, and stochastic jumps as jobs arrive. Because W^2 is a stationary random variable, the expected rate of decrease and increase must be equal, due to the rate conservation law [50] with respect to W^2 .

To calculate the expected rate of decrease, note that, ignoring moments where jobs arrive, $\frac{d}{dt} W(t) = -B(t)$, by definition, where $B(t)$ is the total service rate of the system at time t . As a result, $\frac{d}{dt} W(t)^2 = -2W(t)B(t)$, ignoring arrival epochs. This expected rate of decrease is a well-defined random variable, because the system converges to stationarity. Thus, the expected rate of decrease of W^2 is $2E[WB]$.

To calculate the expected rate of increase, let t^- be the time just before a job arrives to the system. When the job arrives, W^2 increases from $W(t^-)^2$ to $(W(t^-) + S)^2$, a change of $2W(t^-)S + S^2$. Note that $W(t^-)$ is distributed as W , by PASTA. Note also that W and S are independent, because S is sampled i.i.d.. As a result, the expected increase per arrival is $2E[W]E[S] + E[S^2]$. Arrivals occur at rate λ . As a result, the expected rate of increase is $2\lambda E[W]E[S] + \lambda E[S^2]$.

To show that these rates are equal, we must show that the rates are finite. This follows from the fact that $E[W]$ is finite, which we prove in Lemma 5.

As a result, the rates of increase and decrease of W^2 are equal:

$$\begin{aligned}
 2E[WB] &= 2\lambda E[W]E[S] + \lambda E[S^2] \\
 E[WB] &= \lambda E[W]E[S] + \frac{\lambda}{2} E[S^2] \\
 E[WB] &= \rho E[W] + \frac{\lambda}{2} E[S^2] \\
 E[W] - E[W(1 - B)] &= \rho E[W] + \frac{\lambda}{2} E[S^2] \\
 E[W](1 - \rho) &= E[W(1 - B)] + \frac{\lambda}{2} E[S^2] \\
 E[W] &= \frac{E[W(1 - B)]}{1 - \rho} + \frac{\lambda E[S^2]}{2(1 - \rho)} \tag{7}
 \end{aligned}$$

Now, we merely need to bound $E[W(1 - B)]$. We do so by switching to the limited view. Note that

$$\begin{aligned}
 E[W(1 - B)] &= E[W(1 - B)\mathbb{1}\{B = 1\}] + E[W(1 - B)\mathbb{1}\{B < 1\}] \\
 &= E[W(1 - B)\mathbb{1}\{B < 1\}]
 \end{aligned}$$

Because the model is work conserving, if $B < 1$, the front is not full, and there are at most $n - 1$ jobs in the system. Taking expectations over the future randomness of these jobs, at any time t for which $B(t) < 1$,

$$E[W(t)] \leq (n - 1)\text{rem}_{\text{sup}}(S, C)$$

Therefore,

$$\begin{aligned}
 E[W(1 - B)\mathbb{1}\{B < 1\}] &\leq (n - 1)\text{rem}_{\text{sup}}(S, C)E[(1 - B)\mathbb{1}\{B < 1\}] \\
 &= (n - 1)\text{rem}_{\text{sup}}(S, C)E[1 - B] \\
 &= (n - 1)\text{rem}_{\text{sup}}(S, C)(1 - \rho) \\
 E[W(1 - B)] &\leq (n - 1)\text{rem}_{\text{sup}}(S, C)(1 - \rho).
 \end{aligned}$$

Substituting this into (7), our equation for $E[W]$, we find that

$$E[W] \leq \frac{\lambda E[S^2]}{2(1 - \rho)} + (n - 1)\text{rem}_{\text{sup}}(S, C).$$

Dropping the first term of (7), we also get a lower bound:

$$E[W] \geq \frac{\lambda E[S^2]}{2(1 - \rho)}.$$

□

One might alternatively try to prove Lemma 3 via a coupling argument, by coupling the WCFS system to an $M/G/1$ with the same arrival process. Unfortunately, this proof strategy does not succeed, for a subtle reason.

One can show that the difference in work between the two systems during an interval when the WCFS system has a full front is bounded by the amount of work in the WCFS system at the beginning of the interval. This is analogous to the many-jobs interval argument used by Grosf et al. [51] to analyze relevant work in the $M/G/k/SRPT$. The key difference is that in the WCFS setting, we consider total work, not relevant work, meaning that job sizes are not bounded. As a result, while the expected work at the beginning of a full-front interval is bounded, the realization of that work may be arbitrarily large.

A coupling argument would therefore need to bound the relative length of full-front intervals started by different amounts of work, to prove a time-average bound on the gap between $E[W]$ and $E[W^{M/G/1}]$. This seems intractable, given the generality of WCFS policies.

Instead, by using a rate-conservation approach, formalized by Palm Calculus, we directly connect the small expected amount of work in a WCFS system with non-full front to a small expected difference in work between the two systems. We therefore prove Lemma 3, while avoiding all of the complications of a coupling-based argument.

5.5 Lemma 4: bounding $E[T_Q]$

Now, we can bound $E[T_Q]$ by combining Lemmas 2 and 3:

Lemma 4 (Queueing time bounds) *For any model $\pi \in WCFS$, if $\text{rem}_{\text{sup}}(S, C)$ is finite,*

$$E[T_Q^\pi] \leq \frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} + (n - 1)\text{rem}_{\text{sup}}(S, C)$$

$$E[T_Q^\pi] \geq \frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} - (n - 1)\text{rem}_{\text{sup}}(S, C)$$

5.6 Lemma 5: finite $E[W]$

Lemma 5 (Finite mean work) *For any model $\pi \in WCFS$, if $\text{rem}_{\text{sup}}(S, C)$ is finite, for any load $\rho < 1$, W is a well-defined stationary random variable and $E[W]$ is finite.*

Proof Recall that $W = W_F + W_Q$; we first focus on W_F . There are at most n jobs in the front at any time. In the limited view, each job has expected remaining size at most $\text{rem}_{\text{sup}}(S, C)$, so $E[W_F] \leq n\text{rem}_{\text{sup}}(S, C)$.

As for the stationarity of the state of the front, this follows from two assumptions we made in Sect. 2.3. First, we assumed that the service policy is dependent only on the state of the front. Second, the front must empty and thereby undergo renewals, because the service rate $B(t)$ is at least b_{inf} whenever the system is nonempty. As a result, W_F is stationary.

We now turn to W_Q . To prove that W_Q is stationary and well-defined with finite mean, we will apply the “inventory process” results of [52, 53]’s refinement of those results.

We upper bound W_Q by \mathcal{W} , which we will write as an inventory process.

$$\mathcal{W} := W \mathbb{1}\{W_Q > 0\}.$$

Here, we will use the omniscient view, so $\mathcal{W}(t)$ is a specific value. By proving \mathcal{W} is stationary and well-defined with finite mean, we also show the same is true of W_Q . Because $W_Q = (\mathcal{W} - W_F)^+$, the stationarity of \mathcal{W} also implies the stationarity of W_Q , given the stationarity of W_F .

To write \mathcal{W} as an inventory process as in [52], we must define a process $X(t)$ with stationary and ergodic increments, such that

$$\mathcal{W}(t) = X(t) + L(t),$$

where

$$L(t) := \sup_{0 \leq s \leq t} (-\min\{0, X(s)\}).$$

Here, $X(t)$ represents the potential workload process, and $L(t)$ corrects for the fact that the queue can empty.

We will apply [53, Theorem 2.2.1], for the special case of the first moment. Note by Remarks 1 and 3, for the first moment of an inventory process, it suffices to show:

- Negative drift: There exists an amount of work $w < \infty$ and a drift rate $\delta > 0$ such that conditioned on $\mathcal{W}(t) \geq w$,

$$\lim_{\epsilon \rightarrow 0} \frac{E_{\mathcal{F}_t}[X(t + \epsilon) - X(t)]}{\epsilon} \geq -\delta$$

where \mathcal{F}_t is the filtration defined by the behavior of the system up to time t .

- Finite second moment of positive jumps: There exists a constant $k_1 < \infty$ such that

$$\lim_{\epsilon \rightarrow 0} E_{\mathcal{F}_t}[(X(t + \epsilon) - X(t))^+]^2 \leq k_1$$

Now, we define the potential workload process $X(t)$ based on $W(t)$ and $W_Q(t)$.

During intervals when $W_Q(t) = 0$, $X(t)$ is constant. If t_0 is the beginning of an interval where $W_Q(t) > 0$, $X(t)$ jumps up by $W(t_0^+)$ at time t_0 . During an interval where $W_Q(t) > 0$, $X(t)$ mimics $W(t)$: $X(t)$ rises by S when a job arrives, and decreases

at rate 1. If t_1 is the end of an interval where $W_Q(t) > 0$, $X(t)$ jumps down by $W(t_1^-)$ at time t_1 .

By construction, $X(t)$ generates $\mathcal{W}(t)$ as an inventory process. For example, let t_1 be the end of an interval where $W_Q(t) > 0$. Assume that the desired relationship between $X(t)$ and $\mathcal{W}(t)$ holds up to time t_1^- . In particular, $\mathcal{W}(t_1^-) = W(t_1^-)$. Then, $\mathcal{W}(t_1^+) = 0$, as desired.

Next, we show that $X(t)$ has stationary and ergodic increments. $X(t)$ has two types of increments: First, Poisson arrivals cause increments sampled i.i.d. from S , which are clearly stationary and ergodic. Second, the beginning and end of intervals where $W_Q(t) = 0$ cause increments equal to $W_F(t)$. These increments are stationary and ergodic because the state of the front, and W_F in particular, are stationary. Thus, $X(t)$ has stationary and ergodic increments.

To demonstrate negative drift, let w be an arbitrary nonzero amount of work. Whenever $\mathcal{W}(t) \geq w$, $X(t)$ has two types of increments: jumps of size S occurring at rate λ , and continuous decrease at rate 1. As a result, the drift of $X(t)$ is $\rho - 1 < 0$.

To demonstrate finite second moment of positive jumps, note that $X(t)$ has two kinds of positive jumps: Jumps of size S , when $W_Q(t) > 0$, and jumps of size $W(t)$, at the beginning of a $W_Q > 0$ interval.

Switching back to the limited view, note that the latter kind of jump consists of the remaining size of at most n jobs. These remaining sizes are distributed as

$$R(a, c) \sim [S_c - a \mid S_c > a]$$

for some age a and class c .

It therefore suffices to show that there exists a constant r such that for all a, c ,

$$E[R(a, c)^2] \leq r < \infty.$$

To do so, we will write $R(a, c)_e$, the excess of the remaining size distribution, as a mixture of remaining size distributions for different ages. Note that for any distribution Y , the excess Y_e is equivalent to

$$Y_e \sim [Y - Y_e \mid Y > Y_e].$$

This holds because the forward and backwards renewal times are distributed identically [37, Chapter 23]. By applying this construction with $Y = R(a, c)$, we find that

$$\begin{aligned} R(a, c)_e &\sim [R(a, c) - R(a, c)_e \mid R(a, c) > R(a, c)_e] \\ &= [S_c - (a + R(a, c)_e) \mid S_c > a + R(a, c)_e]. \end{aligned}$$

As a result, $a + R(a, c)_e$ is the desired age distribution.

For any age a' , $E[R(a', c)] \leq \text{rem}_{\text{sup}}(S, C)$. Because $R(a, c)_e$ can be written as a mixture of remaining size distributions, $E[R(a, c)_e] \leq \text{rem}_{\text{sup}}(S, C)$, which is finite by assumption.

We can now bound $E[R(a, c)^2]$:

$$E[R(a, c)_e] = \frac{E[R(a, c)^2]}{2E[R(a, c)]}$$

$$E[R(a, c)^2] = 2E[R(a, c)]E[R(a, c)_e] \leq 2\text{rem}_{\text{sup}}(S, C)^2$$

Thus, the requirements of [53, Theorem 2.2.1] are satisfied, so both \mathcal{W} and W_Q are stationary and well-defined, and have finite mean. \square

5.7 Lemma 6: bounding $E[T_F]$

Lemma 6 (Front time bounds) *For any model $\pi \in \text{WCFS}$,*

$$E[S] \leq E[T^F] \leq \frac{nE[S]}{b_{\text{inf}}}$$

Proof First, to prove that $E[T^F] \geq E[S]$, note that if a job receives service at the maximum possible rate of 1 for the entire time it is in the front, then the job will complete in time S . As a result, $E[T^F] \geq E[S]$.

To prove the upper bound, recall that by the non-idling assumption from Sect. 2.1.3, in all states of the front s where $N_F(s) \geq 1$, the service rate $B(s) \geq b_{\text{inf}}$. Because $N_F(s) \leq n$, we can bound the ratio $B(s)/N_F(s)$ in all $N_F(s) \geq 1$ states:

$$\frac{B(s)}{N_F(s)} \geq \frac{b_{\text{inf}}}{n}.$$

Therefore, in all states,

$$B(s) \geq \frac{b_{\text{inf}}}{n} N_F(s).$$

In expectation, the same must hold:

$$E[B] \geq \frac{b_{\text{inf}}}{n} E[N_F].$$

Note that $E[B] = \rho$ and $E[N_F] = \lambda E[T_F]$ by Little’s Law. Thus,

$$\rho \geq \frac{b_{\text{inf}}}{n} \lambda E[T_F]$$

$$\frac{nE[S]}{b_{\text{inf}}} \geq E[T_F].$$

\square

Note that Lemma 6 proves a relatively weak bound on $E[T^F]$, because we have only made the weak assumption that b_{inf} is positive. In many models, one can prove

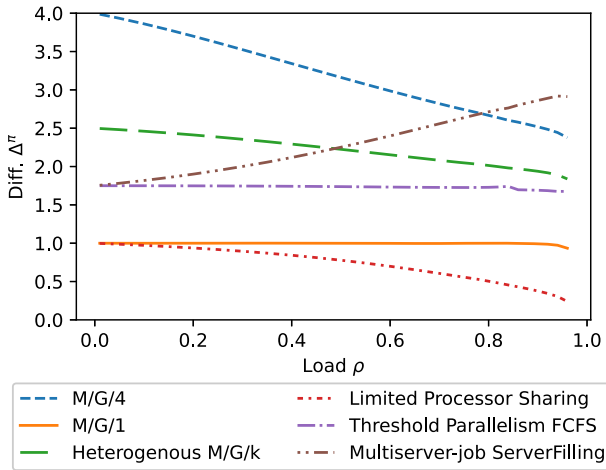


Fig. 5 Δ^π for WCFS models. Job size distribution S is hyperexponential: $\text{Exp}(2)$ w.p. $1/2$, $\text{Exp}(2/3)$ otherwise. 10^9 arrivals simulated. $\rho > 0.96$ omitted due to the large amount of random noise under high load. Specific settings: Heterogeneous $M/G/k$ with speeds $[0.4, 0.3, 0.2, 0.1]$. Limited Processor Sharing with Multi-Programming Level 4. Threshold Parallelism FCFS with joint random variable (S, L) of $(\text{Exp}(2), 1)$ w.p. $1/2$, $(\text{Exp}(2/3), 4)$ otherwise. Multiserver-Job ServerFilling with joint random variable (V, X) of $(1, \text{Exp}(1/2))$ w.p. $1/2$, $(4, \text{Exp}(2/3))$ otherwise

a stronger bound on $E[T^F]$ by using more information about the model’s dynamics when the front is not full.

From Lemmas 4 and 6, Theorem 2 follows immediately, with explicit formulas for c_{upper}^π and c_{lower}^π .

6 Empirical comparison: WCFS and non-WCFS

We have proven tight bounds on mean response time for all WCFS policies. To quantify the tightness of our bounds, we define the *mean response time difference* Δ^π for a given policy π :

$$\Delta^\pi = E[T^\pi] - \frac{\rho}{1 - \rho} \frac{E[S^2]}{2E[S]} = E[T^\pi] - E[T_Q^{M/G/1}].$$

For instance, $\Delta^{M/G/1} = E[S]$.

This definition is useful because we have shown in Theorem 2 that for any load ρ , $\Delta^\pi \in [c_{\text{lower}}^\pi, c_{\text{upper}}^\pi]$, for constants $c_{\text{lower}}^\pi, c_{\text{upper}}^\pi$ not dependent on ρ , but potentially depending on the model π .

To investigate the behavior of Δ^π , we turn to simulation. We simulate both WCFS models, to confirm our results, as well as non-WCFS models, to show that non-WCFS models typically do not have constant Δ^π in the $\rho \rightarrow 1$ limit.

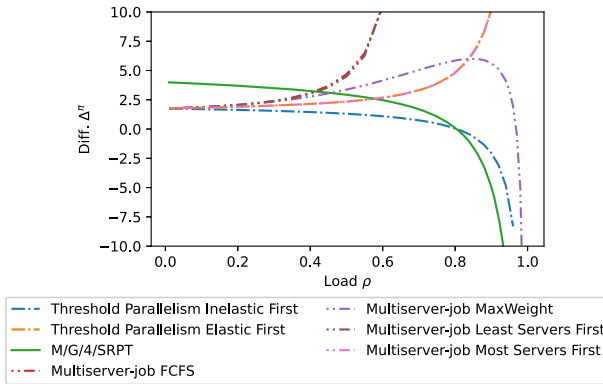


Fig. 6 Δ^π for non-WCFS models. Same job sizes and specific settings as in Fig. 5. Same number of arrivals and range of ρ except MaxWeight: 10^{10} arrivals, $\rho \in [0, 0.99]$

In Fig. 5, we simulate WCFS models: our four motivating models from Sect. 3, as well as the simpler $M/G/k$ and $M/G/1$ models. In each case, we find that Δ^π remains bounded quite close to 0, meaning that Theorem 2 holds with constants close to 0.

In Fig. 5, we see that for some models, Δ^π increases with ρ , while for others, Δ^π decreases with ρ . Intuitively, this depends on which jobs tend to be prioritized as $\rho \rightarrow 1$. Policies which serve many jobs at once, such as the $M/G/4$ and Limited Processor Sharing systems, typically have Δ^π decrease as $\rho \rightarrow 1$, because they allow small and large jobs to share service. As a result, small jobs can complete faster than in an $M/G/1$, lowering Δ^π if ρ is large enough that many jobs are typically in the system.

In contrast, policies which reorder large jobs ahead of small jobs typically have Δ^π increase as $\rho \rightarrow 1$, by the same principle. For example, Multiserver-Job ServerFilling prioritizes jobs in the front which require 4 servers. In the setting depicted in Fig. 5, such jobs have mean size $3/2$ in this system, compared to the overall mean size $E[S] = 1$.

In all of the settings simulated in Fig. 5, $\Delta^\pi > 0$. This is merely a coincidence, not a general rule, as can be seen in Fig. 7b.

Regardless of the different reordering behavior of these different WCFS policies, Δ^π does not diverge as $\rho \rightarrow 1$, as predicted by Theorem 2.

In contrast, in Fig. 6, we simulate several non-WCFS models, which we depicted earlier in Fig. 2. These models are:

- *Threshold parallelism inelastic first*: This is the Threshold Parallelism model from Sect. 3.3, but rather than serving jobs in FCFS order, we prioritize jobs j with smaller parallelism threshold p_j [38].
- *Threshold parallelism elastic first*: This is the Threshold Parallelism model from Sect. 3.3, but we prioritize jobs j with larger parallelism threshold p_j .
- *M/G/k/SRPT*: This is an $M/G/k$, where each of the k servers runs at speed $1/k$, and we prioritize jobs of least remaining size.
- *Multiserver-job FCFS*: This is the Multiserver-Job model from Sect. 3.4, but we serve jobs in FCFS order. If the next job to be served does not “fit” in the remaining

servers, those servers remain idle until other jobs complete, idling sufficient servers to allow the job to fit.

- *Multiserver-job least servers first*: This is the Multiserver-Job model from Sect. 3.4, but we prioritize jobs j with smaller server requirements v_j . Again, if the next job does not fit, the remaining servers remain idle until the job can fit.
- *Multiserver-job most servers first*: This is the Multiserver-Job model from Sect. 3.4, but we prioritize jobs j with larger server requirements v_j .
- *Multiserver-job MaxWeight*: This is the Multiserver-Job model from Sect. 3.4, but we serve jobs according to the “MaxWeight” policy which we describe in Sect. 4.5.2.

In all cases, prioritization is preemptive.

Our empirical results in Fig. 6 indicate that for these non-WCFS policies, Δ^π diverges as $\rho \rightarrow 1$. Specifically, for Threshold Parallelism Elastic First, Multiserver-Job FCFS, Multiserver-Job Least Servers First and Multiserver-Job Most Servers First, Δ^π appears to diverge in the positive direction. For Threshold Parallelism Inelastic First, $M/G/k/SRPT$, and Multiserver-Job ServerFilling, Δ^π appears to diverge in the negative direction. Note the expanded scale of Fig. 6 as compared to Fig. 5. For Multiserver-Job MaxWeight, we performed additional simulation, which indicated that Δ^π diverged in the negative direction as $\rho \rightarrow 1$.

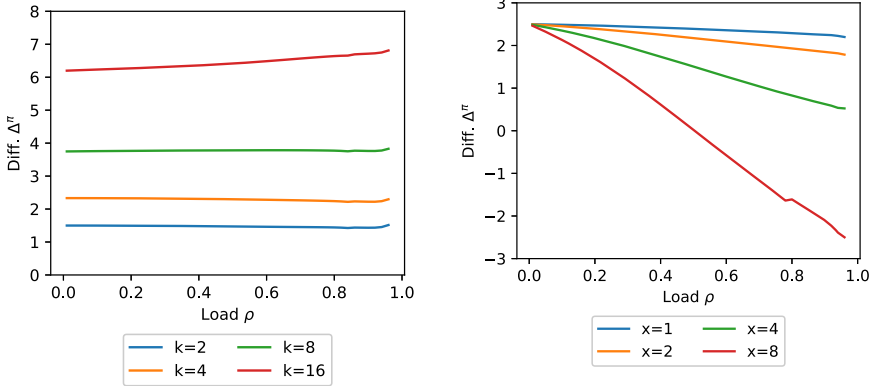
Next, we explore the behavior of Δ^π for WCFS models, as we vary the front size n and the job size distribution S .

First, in Fig. 7a, we investigate the effects of varying front size n on Δ^π for the Multiserver-Job model with our ServerFilling policy; under this model, the front size n is equal to the number of servers k . In this setting, the difference Δ^π empirically grows approximately linearly with the number of servers k , and is nearly constant as $\rho \rightarrow 1$. This matches the behavior of our bounds proven in Theorem 2, which expand linearly with n . Our simulations indicate that other WCFS policies similarly experience linear relationships between n and Δ^π .

In Fig. 7b, we investigate the effects of varying job size distribution S on Δ^π in the Heterogeneous $M/G/k$ where the job size distribution S is parameterized by a real value x . Each S is a hyperexponential distribution with $E[S] = 1$. At large ages a , the remaining size distributions $[S - a \mid S > a]$ of these job size distributions converge to $\text{Exp}(1/x)$, the larger exponential branch. From this, it is straightforward to show that $\text{rem}_{\text{sup}}(S) = x$.

In Fig. 7b, we see that as x increases, Δ^π at loads near 1 falls linearly, with more negative slope for larger x . However, for each specific x , it does not appear that Δ^π is diverging to positive or negative infinity. For instance, consider the red curve, $x = 8$: as $\rho \rightarrow 1$, Δ^π converges to a value near -3 , rather than diverging.

Broadly, Fig. 7b matches the behavior of our bounds proven in Theorem 2, which expand linearly with $\text{rem}_{\text{sup}}(S)$, which here is x . We have empirically found that other WCFS policies similarly experience linear relations between $\text{rem}_{\text{sup}}(S)$ and Δ^π , for hyperexponential job size distributions S , and we believe that similar behavior will occur for other common job size distributions.



(a) Varying front size n . Multiserver-job ServerFilling with $k = [2, 4, 8, 16]$. S distributed $Exp(1)$. Server requirement V distributed uniformly over all integer powers of $2 \leq k$.

(b) Varying job size distributions. Heterogeneous $M/G/4$ with speeds $[0.4, 0.3, 0.2, 0.1]$. S distributed hyperexponential: $Exp(1/x)$ with probability $1/2x$, else $Exp((2x - 1)/x)$, for $x \in [1, 2, 4, 8]$. $E[S] = 1, C^2 \approx [1, 1.67, 3.57, 7.53]$.

Fig. 7 a Varying front size n . Multiserver-Job ServerFilling with $k = [2, 4, 8, 16]$. S distributed $Exp(1)$. Server requirement V distributed uniformly over all integer powers of $2 \leq k$. b Varying job size distributions. Heterogeneous $M/G/4$ with speeds $[0.4, 0.3, 0.2, 0.1]$. S distributed hyperexponential: $Exp(1/x)$ with probability $1/2x$, else $Exp((2x - 1)/x)$, for $x \in [1, 2, 4, 8]$. $E[S] = 1, C^2 \approx [1, 1.67, 3.57, 7.53]$. Δ^π under WCFS models with varying conditions. Up to 10^9 arrivals simulated

7 Conclusion

We introduce the *work-conserving finite-skip* (WCFS) framework, and use it to analyze many important queueing models which have eluded analysis thus far. We prove that the scaled mean response time $E[T^\pi](1 - \rho)$ of any WCFS model π converges in heavy traffic to the same limit as $M/G/1/FCFS$. Moreover, we prove that the additive gap $\Delta^\pi = E[T^\pi] - E[T_Q^{M/G/1}]$ remains bounded by explicit constants at all loads ρ , proving rapid convergence to the heavy traffic limit.

A possible direction for future work would be to tighten the explicit constants on Δ^π . Doing so will likely require use of more detailed properties of the WCFS models being analyzed, but seems quite doable.

This paper considers models which are finite skip and work conserving relative to the FCFS service ordering. Another interesting direction would be to investigate policies which are “finite-skip” relative to other base service orderings. Hopefully, one could prove bounds on mean response time of models in this new class relative to an $M/G/1$ operating under the base service ordering.

Finally, one could try to characterize other metrics of response time for WCFS policies, such as tail metrics. One possible approach to doing so would be to generalize the rate-conservation technique used in Lemma 3.

Funding Funding was provided by National Science Foundation (Grant Nos. CMMI-1938909, CSR-1763701) and a Google 2020 Faculty Research Award.

References

1. Nathuji, R., Isci, C., Gorbato, E.: Exploiting platform heterogeneity for power efficient data centers. In: Fourth International Conference on Autonomic Computing (ICAC'07), p. 5 (2007)
2. Mars, J., Tang, L., Hundt, R.: Heterogeneity in “homogeneous” warehouse-scale computers: a performance opportunity. *IEEE Comput. Arch. Lett.* **10**(2), 29–32 (2011)
3. Cho, H.-D., Engineer, P.D.P., Chung, K., Kim, T.: Benefits of the big. LITTLE architecture. *EETimes*, (2012)
4. Yashkov, S., Yashkova, A.: Processor sharing: a survey of the mathematical theory. *Autom. Remote Control* **68**(9), 1662–1731 (2007)
5. Nuyens, M., Van Der Weij, W.: Monotonicity in the limited processor sharing queue. *Resource* **4**, 7 (2008)
6. Telek, M., Van Houdt, B.: Response time distribution of a class of limited processor sharing queues. *SIGMETRICS Perform. Eval. Rev.* **45**(3), 143–155 (2018)
7. Zhang, J., Zwart, B.: Steady state approximations of limited processor sharing queues in heavy traffic. *Queueing Syst.* **60**(3), 227–246 (2008)
8. Gupta, V., Harchol-Balter, M.: Self-adaptive admission control policies for resource-sharing systems. *SIGMETRICS Perform. Eval. Rev.* **37**(1), 311–322 (2009)
9. Delimitrou, C., Kozyrakis, C.: Quasar: Resource-efficient and QoS-aware cluster management. In: Proceedings of the 19th International Conference on Architectural Support for Programming Languages and Operating Systems. ASPLOS '14, pp. 127–144 (2014)
10. Peng, Y., Bao, Y., Chen, Y., Wu, C., Guo, C.: Optimus: an efficient dynamic resource scheduler for deep learning clusters. In: Proceedings of the Thirteenth EuroSys Conference. EuroSys '18 (2018)
11. Maguluri, S.T., Srikant, R., Ying, L.: Stochastic models of load balancing and scheduling in cloud computing clusters. In: 2012 Proceedings IEEE Infocom, pp. 702–710. IEEE, Orlando (2012)
12. Feitelson, D.G., Rudolph, L., Schwiegelshohn, U.: Parallel job scheduling—a status report. In: Workshop on Job Scheduling Strategies for Parallel Processing, pp. 1–16. Springer, New York (2004)
13. Srinivasan, S., Kettimuthu, R., Subramani, V., Sadayappan, P.: Characterization of backfilling strategies for parallel job scheduling. In: Proceedings of International Conference on Parallel Processing Workshop, pp. 514–519 (2002)
14. Carastan-Santos, D., De Camargo, R.Y., Trystram, D., Zrigui, S.: One can only gain by replacing easy backfilling: a simple scheduling policies case study. In: 2019 19th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID), pp. 1–10 (2019)
15. Tirmazi, M., Barker, A., Deng, N., Haque, M.E., Qin, Z.G., Hand, S., Harchol-Balter, M., Wilkes, J.: Borg: the next generation. In: Proceedings of the Fifteenth European Conference on Computer Systems. EuroSys '20 (2020)
16. Groszof, I., Harchol-Balter, M., Scheller-Wolf, A.: Stability for two-class multiserver-job systems. *arXiv preprint arXiv:2010.00631* (2020)
17. Groszof, I., Harchol-Balter, M., Scheller-Wolf, A.: WCFS: a new framework for analyzing multiserver systems (2022). <https://doi.org/10.48550/ARXIV.2109.12663>
18. Loulou, R.: Multi-channel queues in heavy traffic. *J. Appl. Probab.* **10**(4), 769–777 (1973)
19. Köllerström, J.: Heavy traffic theory for queues with several servers. I. *J. Appl. Probab.* **11**(3), 544–552 (1974)
20. Köllerström, J.: Heavy traffic theory for queues with several servers. II. *J. Appl. Probab.* **16**(2), 393–401 (1979)
21. Kingman, J.: Some inequalities for the queue $GI/G/1$. *Biometrika* **49**(3/4), 315–324 (1962)
22. Gamarnik, D., Momčilović, P.: Steady-state analysis of a multiserver queue in the Halfin-Whitt regime. *Adv. Appl. Probab.* **40**(2), 548–577 (2008)
23. Aghajani, R., Ramanan, K.: The limit of stationary distributions of many-server queues in the Halfin-Whitt regime. *Math. Oper. Res.* **45**(3), 1016–1055 (2020)
24. Dai, J., Dieker, A., Gao, X.: Validity of heavy-traffic steady-state approximations in many-server queues with abandonment. *Queueing Syst.* **78**(1), 1–29 (2014)
25. Goldberg, D.A., Li, Y.: Simple and explicit bounds for multi-server queues with universal $1/(1-\rho)$ scaling. *arXiv preprint arXiv:1706.04628* (2017)
26. Efrosinin, D.V., Rykov, V.V.: On performance characteristics for queueing systems with heterogeneous servers. *Autom. Remote Control* **69**(1), 61–75 (2008)

27. Alves, F., Yehia, H., Pedrosa, L., Cruz, F., Kerbache, L.: Upper bounds on performance measures of heterogeneous M/M/c queues. *Math. Probl. Eng.* **2011** (2011)
28. Efrosinin, D., Stepanova, N., Sztrik, J., Plank, A.: Approximations in performance analysis of a controllable queueing system with heterogeneous servers. *Mathematics* **8**(10), 1803 (2020)
29. Lin, W., Kumar, P.: Optimal control of a queueing system with two heterogeneous servers. *IEEE Trans. Autom. Control* **29**(8), 696–703 (1984)
30. Van Harten, A., Slepchenko, A.: On Markovian multi-class, multi-server queueing. *Queueing Syst.* **43**(4), 307–328 (2003)
31. Boxma, O.J., Deng, Q., Zwart, A.P.: Waiting-time asymptotics for the M/G/2 queue with heterogeneous servers. *Queueing Syst.* **40**(1), 5–31 (2002)
32. Keaogile, T., Fatai Adewole, A., Ramasamy, S.: Geo (λ)/Geo (μ)+ G/2 queues with heterogeneous servers operating under FCFS queue discipline. *Am. J. Appl. Math. Stat* **3**(2), 54–58 (2015)
33. Sani, S., Daman, O.A.: The M/G/2 queue with heterogeneous servers under a controlled service discipline: stationary performance analysis. *IAENG Int. J. Appl. Math.* **45**(1) (2015)
34. Ramasamy, S., Daman, O.A., Sani, S.: An M/G/2 queue where customers are served subject to a minimum violation of FCFS queue discipline. *Eur. J. Oper. Res.* **240**(1), 140–146 (2015)
35. Zhang, J., Dai, J.G., Zwart, B.: Law of large number limits of limited processor-sharing queues. *Math. Oper. Res.* **34**(4), 937–970 (2009)
36. Zhang, J., Dai, J.G., Zwart, B.: Diffusion limits of limited processor sharing queues. *Ann. Appl. Probab.* **21**(2), 745–799 (2011)
37. Harchol-Balter, M.: *Performance Modeling and Design of Computer Systems: Queueing Theory in Action*. Cambridge University Press, Cambridge (2013)
38. Berg, B., Dorsman, J.-P., Harchol-Balter, M.: Towards optimality in parallel scheduling. *Proc. ACM Meas. Anal. Comput. Syst.* **1**(2) (2017)
39. Berg, B., Harchol-Balter, M.: Optimal scheduling of parallel jobs with unknown service requirements. In: *Handbook of Research on Methodologies and Applications of Supercomputing*, pp. 18–40. IGI Global, Hershey (2021)
40. Berg, B., Harchol-Balter, M., Moseley, B., Wang, W., Whitehouse, J.: Optimal resource allocation for elastic and inelastic jobs. In: *Proceedings of the 32nd ACM Symposium on Parallelism in Algorithms and Architectures*. SPAA '20, pp. 75–87 (2020)
41. Brill, P.H., Green, L.: Queues in which customers receive simultaneous service from a random number of servers: a system point approach. *Manag. Sci.* **30**(1), 51–68 (1984)
42. Rumyantsev, A., Morozov, E.: Stability criterion of a multiserver model with simultaneous service. *Ann. Oper. Res.* **252**(1), 29–39 (2017)
43. Hong, Y., Wang, W.: Sharp zero-queueing bounds for multi-server jobs (2021)
44. Ghaderi, J.: Randomized algorithms for scheduling VMs in the cloud. In: *IEEE INFOCOM 2016—The 35th Annual IEEE International Conference on Computer Communications*, pp. 1–9 (2016)
45. Psychas, K., Ghaderi, J.: Randomized algorithms for scheduling multi-resource jobs in the cloud. *IEEE/ACM Trans. Netw.* **26**(5), 2202–2215 (2018)
46. Psychas, K., Ghaderi, J.: On non-preemptive VM scheduling in the cloud. *Proc. ACM Meas. Anal. Comput. Syst.* **1**(2), 35–13529 (2017)
47. Maguluri, S.T., Srikant, R.: Scheduling jobs with unknown duration in clouds. *IEEE/ACM Trans. Netw.* **22**(6), 1938–1951 (2014)
48. Baccelli, F., Foss, S.: On the saturation rule for the stability of queues. *J. Appl. Probab.* **32**(2), 494–507 (1995)
49. Foss, S., Konstantopoulos, T.: An overview of some stochastic stability methods. *J. Oper. Res. Soc. Jpn.* **47**(4), 275–303 (2004)
50. Miyazawa, M.: Rate conservation laws: a survey. *Queueing Syst.* **15**(1), 1–58 (1994)
51. Grosf, I., Scully, Z., Harchol-Balter, M.: SRPT for multiserver systems. *Perform. Eval.* **127–128**, 154–175 (2018)
52. Sigman, K., Yao, D.D.: Finite moments for inventory processes. *Ann. Appl. Probab.* **4**, 765–778 (1994)
53. Scheller-Wolf, A.: Finite moment conditions for stationary content processes with applications to fluid models and queues. PhD thesis, Columbia University (1996)