# Analysis of Scheduling Policies under Correlated Job Sizes

Varun Gupta[*], Michelle Burroughs, Mor Harchol-Balter

*Computer Science Department, Carnegie Mellon University, Pittsburgh, PA 15213, USA.*

## Abstract

Correlations in traffic patterns are an important facet of the workloads faced by real systems, and one that has far-reaching consequences on the performance and optimization of the systems involved. However, all the existing analytical work on understanding the effect of correlations between successive service requirements (job sizes) is limited to First-Come-First-Served scheduling. This leaves open fundamental questions: How do various scheduling policies interact with correlated job sizes? Can scheduling be used to mitigate the harmful effects of correlations?
In this paper we take the first step towards answering these questions. Under a simple model for job size correlations, we present the first asymptotic analysis of various common size-independent scheduling policies when the job size sequence exhibits high correlation. Our analysis reveals that the characteristics of various scheduling policies, as well as their performance relative to each other, are markedly different under the assumption of i.i.d. job sizes versus correlated job sizes. Further, among the class of size-independent scheduling policies, there is no single scheduling policy that is optimal for all degrees of correlations and thus any optimal policy must learn the correlations. We support the asymptotic analysis with numerical algorithms for exact performance analysis under an arbitrary degree of correlation, and with simulations. Finally, we verify the lessons from our correlation model on real world traces.

*Keywords:* Scheduling, Correlation, MMAP, $M/G/1$, Asymptotic analysis, Fluid analysis

## 1. Introduction

*Motivation*
The $M/G/1$ single-server queue has been used as a guiding model for performance analysis of widely varying systems, such as buffers for network switches, web server downlinks, and the CPU scheduler. There is a large body of work on the analysis of different scheduling policies and

---
[*]Corresponding author
 *Email addresses:* `varun@cs.cmu.edu` (Varun Gupta ), `mburroug@andrew.cmu.edu` (Michelle Burroughs), `harchol@cs.cmu.edu` (Mor Harchol-Balter)

their effects on response times of jobs (defined to be the time from the arrival to the completion of a job) [6]. However, almost all of the exact analysis has been performed under the assumptions of (*i*) Poisson arrival process and (*ii*) independent and identically distributed (*i.i.d.*) job sizes. Long ago, the need was recognized to relax these assumptions, as real systems workloads exhibit significant correlation patterns, and these patterns tend to greatly affect the accuracy of the traditional results [9, 20]. Primarily, there are three kinds of correlations that exist in real workloads: (*i*) correlations between consecutive interarrival times (e.g., network traffic [12], and web server traffic [8, 14, 25]), (*ii*) correlations between interarrival times and the subsequent service requirements (e.g., [3, 5, 12, 14]), and (*iii*) correlations between consecutive service requirements (e.g. packet sizes over network [12], supercomputing jobs [10, 16, 24], and disk request sizes [19]). In this paper we focus on studying the effects of correlations of type (*iii*).

While there has been a lot of *analytical* work studying the effect of all three types of correlation on mean response time in single server queues, all of this work has assumed First-Come-First-Served (FCFS) queues only. Fendick et al. [12] study all three types of correlation via a Brownian approximation and propose a stationary workload approximation based on heavy traffic limits. Adan and Kulkarni [2] also use analysis to study autocorrelation and cross-correlation of interarrival and service times in a MAP/G/1/FCFS queue. Riska et al. [22] use matrix-analytic methods to numerically calculate the mean response time in a MAP/PH/1/FCFS queue with correlated arrival stream. Ghosh and Squillante [14] propose a refinement to the Fendick et al. [12] approximation for FCFS queues, and propose approximations for a multi-class priority system with FCFS scheduling within each class. Cidon et al. [5] study correlations of type (*ii*) by deriving the Laplace transform of the workload using the theory of linear functional equations in a queue with an Interrupted Poisson arrival process.

The effect of correlation has also been studied via *simulation*, see for example [17–19, 23, 28]. In all except [18], FCFS scheduling was assumed. In [18] the authors examine an approximation of Shortest-Job-First (SJF) scheduling, which the authors call SWAP, and compare it against FCFS scheduling via simulation.

In summary, all the prior work dealing with correlations in successive job sizes has almost exclusively dealt with FCFS scheduling. Important questions have remained unanswered: How do different scheduling policies react to correlations in job sizes? Can scheduling be used to allay the detrimental effect of correlated job sizes on the performance?

In this paper, we take an important first step by analyzing the mean response time under various scheduling policies in the presence of correlated job sizes (see Table 1 for a list of policies analyzed in this paper). We restrict ourselves to the class of *size-independent* policies. That is, we consider policies which know the generative correlation model, but not the actual realizations of the sizes (or the size-class) of jobs. In most applications, including scheduling of CPU, IP flows, database queries etc., the job sizes are often not known a priori, and hence size-independent policies are more realistic. We consider the question of how the optimality of size-independent policies is affected by the presence or absence of correlation in the job sizes.

*The MMAP Correlation Model*

We assume the following simple *Markov Modulated Arrival Process (MMAP)* model for job-size correlations: jobs belong to one of two classes called little (L) and huge (H), where jobs of class L (respectively H) are Exponentially distributed with mean $\frac{1}{\mu_L}$ (respectively $\frac{1}{\mu_H} > \frac{1}{\mu_L}$) [1]. Therefore,

---

[1]Note that the mean sizes of the two classes can in fact be close. We have chosen the names of the classes to map to low (L) and high (H) load, respectively, in Section 2.

| Scheduling Policy | Description |
|---|---|
| FIRST-COME-FIRST-SERVED (FCFS) | Jobs are served in the order of arrival. |
| LAST-COME-FIRST-SERVED (LCFS) | Whenever a job completes service, the next job to be served is the one that arrived last. |
| PREEMPTIVE LCFS WITH RESUME (P-LCFS) | New arrivals immediately begin service by preempting the job at the server. On a service completion, the next job to resume service is the one that arrived last. |
| LEAST-ATTAINED-SERVICE (LAS) | The job with the least amount of received service (age) gets to serve. |
| PROCESSOR SHARING (PS) | If there are $n$ jobs in the system, each job gets $\frac{1}{n}$th of the server's capacity. |
| RANDOM-ORDER-OF-SERVICE (ROS) | Whenever a job completes service, the next job to be served is picked uniformly at random from amongst the jobs currently in the queue. |
| OPTIMAL OMNISCIENT (OPT) | A hypothetical optimal scheduling scheme that knows the class of all jobs, and gives preemptive priority to class L jobs. |

Table 1: A glossary of scheduling policies analyzed in this paper.

our jobs belong to a 2-phase hyperexponential ($H_2$) distribution. The system operates under a 2-state Markovian environment process with states L and H: while the environment process is in state L all arrivals are of class L, and while in state H all arrivals are of class H. The arrivals occur according to a Poisson process with rate $\lambda$ independent of the environment process. The times spent in state L during each visit are *i.i.d.* Exponentially distributed with mean $\frac{1}{\alpha_L}$, and those in state H are *i.i.d.* Exponentially distributed with mean $\frac{1}{\alpha_H}$. Denote $\alpha = \alpha_L + \alpha_H$, and $p = \frac{1/\alpha_L}{1/\alpha_L + 1/\alpha_H} = \frac{\alpha_H}{\alpha}$. Thus the time-average probability of an arrival belonging to class L is $p$, and of belonging to class H is $1 - p$. We will use $\rho = \lambda \cdot \left( \frac{p}{\mu_L} + \frac{1-p}{\mu_H} \right)$ to denote the long run fraction of time the system is busy. If we fix the job size distribution and arrival rate (i.e. $\mu_L, \mu_H, p, \lambda$) and set $\alpha = \infty$, then the job sizes form an *i.i.d.* stream. As we decrease $\alpha$ and thereby increase the mean residence time per sojourn of L and H states, we increase the correlation among successive job sizes, since the probability that a class L job is followed by another class L job ($p_{L,L} = p + \frac{\lambda(1-p)}{\lambda+\alpha}$) increases. By expressing $p_{L,L} = \frac{\alpha}{\lambda+\alpha} p + \frac{\lambda}{\lambda+\alpha}$, we can alternately visualize the correlation model as: with probability $\frac{\lambda}{\lambda+\alpha}$ the class of a job is the same as the class of the immediately preceding job, otherwise it is an independent sample from the $H_2$ distribution.

Let $\cdots, X_{-2}, X_{-1}, X_0, X_1, X_2 \cdots$ represent the sequence of job sizes. An appealing property of the above correlation model is the simple closed-form autocorrelation function (acf). In particular, the lag $n$ correlation for $n \geq 1$ is given by:

$$cor(X_m, X_{m+n}) = \frac{\mathbf{E}[X_m X_{m+n}] - \mathbf{E}[X_m]\mathbf{E}[X_{m+n}]}{\sqrt{var(X_m)}\sqrt{var(X_{m+n})}} = \frac{\left(\frac{\lambda}{\lambda+\alpha}\right)^n \left[ \frac{p}{\mu_L^2} + \frac{1-p}{\mu_H^2} \right] + \left(1 - \left(\frac{\lambda}{\lambda+\alpha}\right)^n\right)\mathbf{E}[X_0]^2}{var(X_0)} - \frac{\mathbf{E}[X_0]^2}{var(X_0)}$$

$$= \left(\frac{\lambda}{\lambda+\alpha}\right)^n \frac{\frac{\mathbf{E}[X_0^2]}{2} - \mathbf{E}[X_0]^2}{var(X_0)} = \frac{1}{2}\left(\frac{C^2-1}{C^2}\right)\left(\frac{\lambda}{\lambda+\alpha}\right)^n$$

where $C^2 = \frac{var(X_0)}{\mathbf{E}[X_0]^2} > 1$ denotes the squared coefficient of variation (SCV) of the $H_2$ job size distribution.

*Scope of the MMAP correlation model:.* The MMAP correlation model analyzed in this paper is similar to the model used in [2]. While MMAP models with more than 2 phases (e.g., [19]) or lo-
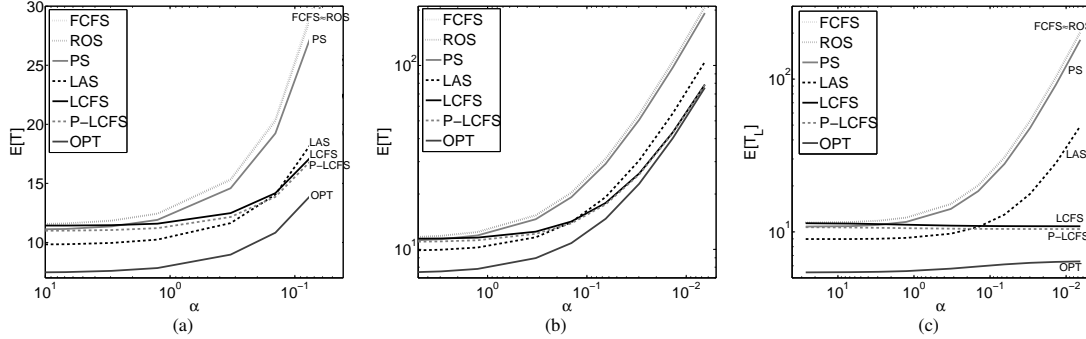
Figure 1: An example of the effect of job-size correlation on scheduling policies: (a) mean response time versus $\alpha$ for low to medium correlation; (b) mean response time versus $\alpha$ for medium to high correlation; (c) mean response time of the "little" (L) jobs versus $\alpha$. Here $\rho = 0.97$ and $C^2 \approx 1.08$. Note that the $\mathbf{E}[T]$ ordering changes from FCFS=ROS=LCFS>PS=P-LCFS>LAS>OPT at $\alpha = \infty$ (*i.i.d.* job sizes) to FCFS≈ROS>PS>LAS>LCFS=P-LCFS=OPT as $\alpha \to 0$ (high correlation).

cal sampling based models [11] are capable of modeling more general auto-correlation functions, the goal of this paper is to use an analytically tractable correlation model to explore *qualitative behavior* of different scheduling policies in the presence of correlated job sizes, and to gain insights for these behaviors and the effect of various system parameters on the performance. We believe that the qualitative behavior of scheduling policies discovered in this paper would extend to more general correlation structures, and we partially test this via real-world traces in Section 3.

*Summary of Contributions*

Most of our results look at the effect of the parameter $\alpha$ on mean response time, $\mathbf{E}[T]$. We prove that, although all scheduling policies we consider are hurt by increasing the correlation, the degree to which correlation affects different policies varies widely. We consider two regimes: *(i)* $\mu_L > \mu_H > \lambda$, where the server is never in overload, and *(ii)* $\mu_L > \lambda > \mu_H$, where the system is in overload during bursts of H jobs, although it is still stable on average. For the no-overload regime, we prove that, as $\alpha$ decreases (correlation increases), all size-independent scheduling policies become the same with respect to mean response time. For the transient-overload regime, we prove that as correlation decreases, there can be a large (up to a factor of $\frac{\mu_L}{\mu_H}$) difference in $\mathbf{E}[T]$ between the policies. Also, the ordering of policies from "best" to "worst" mean response time changes a lot under correlation. An example of performance of the various scheduling policies under the transient-overload regime is shown in Figure 1(a). Some particularly interesting **findings include**:

- LAS is provably sub-optimal among size-independent policies when $\alpha \to 0$, while it has provably the best mean response time when $\alpha \to \infty$ for an $H_2$ job size distribution (due to its decreasing failure rate [21]).
- LCFS is provably best when $\alpha \to 0$, while it is worst (along with FCFS, ROS) when $\alpha \to \infty$.
- P-LCFS is also provably best when $\alpha \to 0$, which is interesting because under $\alpha \to \infty$ (*i.i.d.* case) LCFS and P-LCFS can be far apart for high variability job size distributions.
- PS can be arbitrarily worse than P-LCFS as $\alpha \to 0$, while they are provably equal as $\alpha \to \infty$.

The effect of correlation on the mean response time of the L jobs, $\mathbf{E}[T_L]$, is even more pronounced. In particular, we prove that:

- While $\mathbf{E}[T_L]$ increases for most policies, as $\alpha$ decreases (correlation increases), $\mathbf{E}[T_L]$ always *decreases* for P-LCFS and for LCFS. An example is shown in Figure 1(c).
- LAS performs poorly for $\mathbf{E}[T_L]$ compared to OPT, and even worse for $\mathbf{E}\left[T_L^2\right]$. Thus, while LAS is designed to help the little jobs by biasing towards jobs with least attained service, it fails to do this under correlation, and policies like LCFS which are entirely oblivious to job size distribution can actually help the little jobs.

The above results are primarily obtained by using fluid analysis and looking at asymptotic behavior of response time as $\alpha \to 0$, see Section 2. However, the effect of correlation under moderate $\alpha$ is also interesting. To study the moderate $\alpha$ regime, we derive numerical algorithms to analyze LCFS, OPT, P-LCFS, and FCFS. [2] For the other policies, we resort to simulations, see Section 3. These numerical and simulation results are useful for understanding the behavior of scheduling policies for intermediate $\alpha$ values and to explore how quickly scheduling policies converge to their asymptotically-limiting ($\alpha \to 0$) behavior. To see how our messages carry through to real-world scenarios, we end Section 3 with trace-driven simulation studies.

## 2. Asymptotic Analysis of Scheduling Policies as $\alpha \to 0$

Our goal in this section is to obtain an understanding of the "first-order effect" of correlations in the job sizes by considering the limiting case where the correlation approaches its maximum value under our model, that is, $\alpha \to 0$.[3] While this extremal case implies arbitrarily long consecutive streaks of only L and only H arrivals, an understanding of the behavior of the various scheduling policies under this asymptote gives us insights into why different scheduling policies react differently to correlation in job sizes, and should help guide the design of policies which are robust to correlation.

In Section 2.1, we present the asymptotic results for the simpler case $\mu_H > \lambda$. The non-trivial case of $\mu_H < \lambda$ is analyzed in Sections 2.2-2.5. A large number of scheduling policies that we will analyze will involve asymptotic analysis of busy periods. We have chosen to present the main results on busy period analysis in Appendix B and focus on the messages in the main body. For ready reference, we have summarized the notation used in this section in Table 2.

*Note on scaling and asymptotic notation:.* The asymptotic analysis of the scheduling policies is performed by considering a sequence of systems, indexed by the parameter $\alpha$. The system with index $\alpha$ is obtained by setting the switching rates of the environment process as $\alpha_H = p \cdot \alpha$ and $\alpha_L = (1 - p)\alpha$, where $p, \mu_L, \mu_H$ and $\lambda$ are held constant. We are interested in seeing the behavior of the scheduling policies in the asymptote $\alpha \to 0$, and hence the expressions for mean response times presented in this section will be written in the *asymptotic notation*: We say that a function $g(\alpha)$ is of a 'smaller order' than $h(\alpha)$ (and make the limit $\alpha \to 0$ implicit), denoted $g(\alpha) = o(h(\alpha))$, when $\frac{g(\alpha)}{h(\alpha)} \to 0$ when $\alpha \to 0$ (see Table 2). When we write the expressions for the mean response time under the $\alpha$th system, we only identify the dominant term in the expression, expressing the remaining terms which become negligible in comparison as $\alpha \to 0$ as

---

[2]Due to lack of space, the asymptotic analysis of PS and ROS, and the results on exact numerical analysis of LCFS, OPT, P-LCFS and FCFS are presented in the extended version [15].

[3]The analysis of the asymptote $\alpha \to 0$ should be seen analogously to heavy traffic analysis where the traffic intensity $\rho$ is allowed to approach 1 to observe the "first order" effect of system parameters (variance, cross-correlations) on the system performance.

| Notation | Meaning | Notation | Meaning |
|---|---|---|---|
| $\mathbf{E}[T_L^\pi], \mathbf{E}[T_H^\pi], \mathbf{E}[T^\pi]$ | mean response time of a class {L, H, avg} job under policy $\pi$ | $\mathbf{E}[D_L^\pi], \mathbf{E}[D_H^\pi], \mathbf{E}[D^\pi]$ | mean delay of a class {L, H, avg} job under scheduling policy $\pi$ |
| $\mathbf{E}[T_L^\pi(x)], \mathbf{E}[T_H^\pi(x)]$ | mean response time of a class L, H job of size $x$ under policy $\pi$ | $W_L, W_H$ | stationary workload conditioned on being in state L, H |
| $r_L$ | $= 1 - \frac{\lambda}{\mu_L}$ | $r_L(x)$ | $= 1 - \lambda s_L(x)$ |
| $r_H$ | $= 1 - \frac{\lambda}{\mu_H}$ | $r_H(x)$ | $= 1 - \lambda s_H(x)$ |
| $\rho$ | $= \lambda(p/\mu_L + (1-p)/\mu_H)$ | $\rho(x)$ | $= \lambda(p s_L(x) + (1-p)s_H(x))$ |
| $s_L(x)$ $s_H(x)$ | $= \mathbf{E}[\min\{\mathrm{Exp}(\mu_L), x\}] = \frac{1-e^{-\mu_L x}}{\mu_L}$ $= \mathbf{E}[\min\{\mathrm{Exp}(\mu_H), x\}] = \frac{1-e^{-\mu_H x}}{\mu_H}$ | $W_L^*, W_H^*$ | stationary fluid workload in a system with flow rates $r_L$ and $r_H$, conditioned on being in state L, H |
| $g(x) = \Theta(h(x))$ as $x \to x_0$ | $0 < \liminf_{x \to x_0} \frac{g(x)}{h(x)} \le \limsup_{x \to x_0} \frac{g(x)}{h(x)} < \infty$ | $W_L^*(x), W_H^*(x)$ | stationary fluid workloads in a system with flow rates $r_L(x), r_H(x)$ |
| $g(x) = o(h(x))$ as $x \to x_0$ | $\lim_{x \to x_0} \frac{g(x)}{h(x)} = 0$ | $\widetilde{X}(s) = \mathbf{E}[e^{-sX}]$ | Laplace transform of r.v. $X$ |

Table 2: Notation used in Section 2.

being of a smaller order than the dominant term. Similarly, we say $g(\alpha)$ is of 'the same order' as $h(\alpha)$ (again with the limit $\alpha \to 0$ implicit), denoted $g(\alpha) = \Theta(h(\alpha))$ when intuitively $\frac{g(\alpha)}{h(\alpha)}$ is eventually bounded between two strictly positive constants. Thus, for example, a $\Theta(1)$ function is eventually bounded between two strictly positive constants as $\alpha \to 0$. In proving theorems about response time, it will often suffice to just argue about the asymptotic order of busy period durations, probabilities and related quantities.

### 2.1. Analysis for case $\mu_H > \lambda$

Let $T_L^\pi$ and $T_H^\pi$ denote the random variables for response time of class L and class H jobs, respectively, under scheduling policy $\pi$ (see Table 2). When $\mu_H > \lambda$, the system is stable during both L and H states, and we have the following intuitive result which we state without proof.

**Theorem 1.** *Let $\pi$ be any work-conserving, size-independent policy. When $\mu_H > \lambda$,*

$$\lim_{\alpha \to 0} \mathbf{E}[T_L^\pi] = \frac{1}{\mu_L - \lambda} \quad ; \quad \lim_{\alpha \to 0} \mathbf{E}[T_H^\pi] = \frac{1}{\mu_H - \lambda}.$$

**Remark 1**: Theorem 1 says that as job sizes become more and more correlated, the behavior of all work-conserving, size-independent scheduling policies will tend to become the same, provided $\mu_H > \lambda$. This is because the system behaves as a mixture of two stable $M/M/1$ systems, and all size-independent scheduling policies have the same mean response time for an $M/M/1$ system. The same argument does not apply when $\mu_H < \lambda$ because the $M/M/1$ during the H states is unstable and the workload built up during the H states results in significant transient effects.

**Remark 2**: Since LAS is optimal (among size-independent policies) at each extreme, we intuitively expect LAS to be near-optimal through the entire range of $\alpha$, and thus for all levels of correlation. We verify that this is indeed true in Section 3, Figure 2.

## 2.2. Preliminaries: Workload analysis via Fluid model for the case $\mu_H < \lambda$

We begin our study of the case $\mu_H < \lambda$ by finding the distribution of stationary workload during the L and H states, respectively. To do this, we first introduce the *fluid model* of our MMAP correlation model.

**Definition 1.** *Under the fluid model, we assume that the workload increases at a constant rate of $-r_H$ during the H states (see Table 2), and decreases at a constant rate of $r_L$ during the L states as long as the workload is positive.*

**Lemma 1.** *Let $W_L^*$ and $W_H^*$ denote the random variables for the stationary workload during L and H states under the fluid model, respectively (we will superscript fluid model random variables by $*$). Let $\widetilde{W_L^*}(s) = \mathbf{E}\left[e^{-sW_L^*}\right]$ and $\widetilde{W_H^*}(s) = \mathbf{E}\left[e^{-sW_H^*}\right]$ denote their Laplace transforms. Then,*

$$\widetilde{W_H^*}(s) = \frac{\gamma_H - \gamma_L}{s + (\gamma_H - \gamma_L)} \; ; \quad \widetilde{W_L^*}(s) = \left(1 - \frac{\gamma_L}{\gamma_H}\right) + \frac{\gamma_L}{\gamma_H} \cdot \frac{\gamma_H - \gamma_L}{s + (\gamma_H - \gamma_L)}$$

*where $\gamma_L = \frac{\alpha_L}{r_L}$ and $\gamma_H = -\frac{\alpha_H}{r_H}$.*
*Thus the workload during the H states, $W_H^*$, is distributed as an $\mathrm{Exp}(\gamma_H - \gamma_L)$ random variable, and the workload during the L states, $W_L^*$, is a mixture of an $\mathrm{Exp}(\gamma_H - \gamma_L)$ random variable and an atom at $0$. Further, the mean of $W_L^*$ and $W_H^*$ are of the order $\Theta\left(\frac{1}{\alpha}\right)$. Thus, as $\alpha \to 0$, the fluid workload diverges at a rate of $\frac{1}{\alpha}$.*

**Lemma 2.** $W_L \stackrel{d}{=} W_L^* + o(\alpha^{-1})$ , $W_H \stackrel{d}{=} W_H^* + o(\alpha^{-1})$.

**Remark 3**: Lemma 2 says that, asymptotically as $\alpha \to 0$, the stationary workload, $W_L$ and $W_H$, of the stochastic system converge in distribution to the stationary workload, $W_L^*$ and $W_H^*$, under the fluid model. While a convergence of workloads on a sample path basis was proved in [4], we are unaware of a proof of the convergence of stationary workloads.

**Proof of Lemma 1:** We first note that by conditional PASTA [27], $W_L^*$ and $W_H^*$ are equal in distribution to the stationary workload at the end of L and H states respectively. Let $\tau_L$ and $\tau_H$ be Exponentially distributed random variables with mean $\frac{1}{\alpha_L}$ and $\frac{1}{\alpha_H}$, respectively. We have the following stochastic fixed point equations:

$$W_H^* \stackrel{d}{=} W_L^* - r_H\tau_H \; ; \quad W_L^* \stackrel{d}{=} \max\{W_H^* - r_L\tau_L, 0\}$$

Taking Laplace transforms of the above equations, we get the following fixed point equations:

$$\widetilde{W_H^*}(s) = \widetilde{W_L^*}(s) \cdot \frac{\alpha_H/r_H}{\alpha_H/r_H - s}; \quad \widetilde{W_L^*}(s) = \frac{s\widetilde{W_H^*}(\alpha_L/r_L) - (\alpha_L/r_L)\widetilde{W_H^*}(s)}{s - \alpha_L/r_L},$$

which yield the expressions in Lemma 1. ∎

**Proof of Lemma 2:** The lemma is proven by starting with Theorem 5 (Appendix A) which gives the exact expressions for the Laplace transforms of $W_L$ and $W_H$. According to Theorem 5:

$$\widetilde{W_L}(s) = \frac{(1-\rho)\alpha m_L m_H - sm_L g_H \pi_L(0)}{\alpha_L g_H m_L + \alpha_H g_L m_H - sg_L g_H} \tag{1}$$

where, $m_L = \mu_L + s$, $m_H = \mu_H + s$, $g_L = \mu_L - \lambda + s$, $g_H = \mu_H - \lambda + s$, $\pi_L(0) = \frac{(1-\rho)\alpha(\mu_H+\xi)}{\xi(\mu_H-\lambda+\xi)}$, and $\xi$ denotes the unique root of the denominator of (1) (viewed as a cubic in $s$) in the interval $(0, +\infty)$.

The quantity $\pi_L(0)$ denotes the long run fraction of time that the system is empty conditioned on being in state L. Taking the limit $\alpha \to 0$, we get

$$\xi = (\lambda - \mu_H) + \frac{p\alpha\lambda}{\lambda - \mu_H} + \Theta(\alpha^2)$$

and thus,

$$\pi_L(0) = \frac{(1-\rho)\alpha(\mu_H + \xi)}{\xi(\mu_H - \lambda + \xi)} = \frac{(1-\rho)\alpha(\lambda + \Theta(\alpha))}{(\lambda - \mu_H + \Theta(\alpha))\left(\frac{p\alpha\lambda}{\lambda - \mu_H} + \Theta(\alpha^2)\right)} = \frac{1-\rho}{p} + \Theta(\alpha)$$

Note that the above is not in disagreement with the result $\mathbf{Pr}\left[W_L^* = 0\right] = \left(1 - \frac{\gamma_L}{\gamma_H}\right)$ as the latter is only equivalent to $\mathbf{Pr}\left[W_L = o\left(\frac{1}{\alpha}\right)\right]$. The other roots of the denominator of (1) in the limit $\alpha \to 0$ are given by:

$$\chi = (\lambda - \mu_L) - \frac{p\alpha\lambda}{\mu_L - \lambda} + \Theta(\alpha^2) \quad \text{and} \quad \eta = -\frac{\alpha\mu_L\mu_H(1-\rho)}{(\mu_L - \lambda)(\lambda - \mu_H)} + \Theta(\alpha^2).$$

Canceling the common factor $(s - \xi)$, and noting that $\frac{\alpha\mu_L\mu_H(1-\rho)}{(\mu_L-\lambda)(\lambda-\mu_H)} = (\gamma_H - \gamma_L)$, we can rewrite:

$$\widetilde{W_L}(s) = \pi_L(0) + K_1 \frac{-\chi}{s - \chi} + K_2 \frac{-\eta}{s - \eta} = \frac{1-\rho}{p} + K_1 \frac{\mu_L - \lambda + \Theta(\alpha)}{s + (\mu_L - \lambda + \Theta(\alpha))} + K_2 \frac{\gamma_H - \gamma_L}{s + (\gamma_H - \gamma_L)}.$$

Matching the coefficients of $s$, we get $K_1 = \frac{1 - r_L}{r_L}\left(\frac{1-\rho}{p}\right) + \Theta(\alpha)$, and $K_2 = 1 - \frac{1-\rho}{pr_L} + \Theta(\alpha) = \frac{\gamma_L}{\gamma_H} + \Theta(\alpha)$. Thus we have proved that, as $\alpha \to 0$, the distribution of $W_L$ is a mixture of an Exponential distribution with mean $\frac{1}{\gamma_H - \gamma_L}$ with probability $\sim \frac{\gamma_L}{\gamma_H}$, and with the remaining probability the stationary distribution of an $M/M/1$ with arrival rate $\lambda$ and service rate $\mu_L$. ∎

*Goals of asymptotic analysis.* Since we are interested in analyzing work-conserving policies, the stationary workload, $W$, is the same across policies. What differs from one policy to another is what types of jobs make up that work. Since we restrict ourselves to size-independent policies, we can bound the mean remaining size of any job under our $H_2$ job size distribution between $\frac{1}{\mu_L}$ and $\frac{1}{\mu_H}$. This gives bounds on $\mathbf{E}[N^\pi]$ – the mean number of jobs in the system for any work-conserving policy $\pi$ – as $\mu_H\mathbf{E}[W] \leq \mathbf{E}[N^\pi] \leq \mu_L\mathbf{E}[W]$. Finally, by applying Little's law, we get $\frac{\mu_H}{\lambda}\mathbf{E}[W] \leq \mathbf{E}[T^\pi] \leq \frac{\mu_L}{\lambda}\mathbf{E}[W]$. Since $\mathbf{E}[W]$ diverges as $\frac{1}{\alpha}$ as $\alpha \to 0$, we have the following.

**Lemma 3.** *When $\mu_H < \lambda$ in the MMAP model, the mean response time of any work-conserving size-independent scheduling policy $\pi$ grows as $\mathbf{E}[T^\pi] = \frac{K^\pi}{\alpha} + o(\frac{1}{\alpha})$, for some constant $K^\pi$ which depends only on the scheduling policy and the parameters $\mu_H, \mu_L, p$ and $\lambda$.*

Our goal is to identify the $K^\pi$ for different policies. This is analogous to heavy traffic analysis, where space (response time, number of jobs in system, etc.) is scaled by $(1 - \rho)$ and analyzed in the limit $\rho \to 1$.

### 2.3. FCFS

**Theorem 2.** *In the regime $\mu_H < \lambda$,*

$$\mathbf{E}\left[D_L^{FCFS}\right] = \frac{(1-p)}{p(1-\rho)}\left(\frac{\lambda}{\mu_H} - 1\right)^2 \frac{1}{\alpha} + o\left(\frac{1}{\alpha}\right)$$

$$\mathbf{E}\left[D_H^{FCFS}\right] = \frac{1}{(1-\rho)}\left(1 - \frac{\lambda}{\mu_L}\right)\left(\frac{\lambda}{\mu_H} - 1\right)\frac{1}{\alpha} + o\left(\frac{1}{\alpha}\right)$$

**Proof:** By conditional PASTA, the delay of class L jobs is distributed as $W_L$, and that of class H as $W_H$. Applying Lemmas 2 and 1, the result is immediate. ∎

**Remark 4**: We already see a divergence in the behavior of scheduling policies when job sizes become correlated. When $\alpha \to \infty$ (*i.i.d.* case), and under a Poisson arrival process, the mean delay under FCFS depends only on the first two moments of the job size distribution. However, as $\alpha \to 0$, it depends on all the parameters of the $H_2$ job size distribution.

### 2.4. OPT, P-LCFS and LCFS

While it is hard to characterize the optimal size-independent policy when job sizes are correlated since the optimal policy might (and will) exploit the correlation structure to predict classes of future jobs based on observed history of job sizes, a trivial lower bound is obtained by considering an omniscient scheduler – that is, a scheduler that knows the *class* (L,H) of each job in the system, but not the exact size, and gives preemptive priority to class L jobs. We call this policy OPT.

**Theorem 3.** *When $\mu_H < \lambda$, we have for each policy $\pi \in \{OPT, P\text{-}LCFS, LCFS\}$:*

$$\mathbf{E}[D_L^\pi] = \Theta(1)$$

$$\mathbf{E}[D_H^\pi] = \left\lfloor \frac{\mu_H}{\lambda(1-p)} \right\rfloor \frac{(1-p)\lambda}{(1-\rho)} \left( \frac{1}{\mu_H} - \frac{1}{\mu_L} \right) \left( \frac{\lambda}{\mu_H} - 1 \right) \frac{1}{\alpha} + o\left( \frac{1}{\alpha} \right)$$

**Corollary 1.** *For $\pi \in \{LCFS, P\text{-}LCFS, OPT\}$, when $\mu_H < \lambda$:* $\lim_{\alpha\to0} \frac{\mathbf{E}[T^{FCFS}]}{\mathbf{E}[T^\pi]} = \frac{\lambda}{\mu_H}$.

**Proof of Theorem 3:** We first consider class L jobs. Under OPT, class L jobs get priority, and hence their response time is stochastically upper bounded by that of an $M/M/1$ with arrival rate $\lambda$ and service rate $\mu_L$, and is $\Theta(1)$. Under P-LCFS, the response time of class L jobs is the busy period started by $\text{Exp}(\mu_L)$ work in state L. By Theorem 6, Case 2 (see Appendix B), this is $\Theta(1)$. Under LCFS, the delay of class L jobs is a busy period started either by $\text{Exp}(\mu_L)$, $\text{Exp}(\mu_H)$ or 0 work. Again, by Theorem 6, Case 2, this is $\Theta(1)$.

To understand the delay of class H jobs, note that the above implies that the mean number of class L jobs in the system, and hence their contribution to the total workload is $\Theta(1)$. However, the stationary average workload is $\Theta(\alpha^{-1})$, and hence this must be composed (aside from a $\Theta(1)$ term) of class H jobs alone. Since, all scheduling policies are size-independent, the mean residual size of these class H jobs is $\frac{1}{\mu_H}$, yielding the mean number of class H jobs of $\frac{p\mathbf{E}[W_L]+(1-p)\mathbf{E}[W_H]}{1/\mu_H}$.

By Little's law, we obtain the mean delay of class H jobs as $\frac{p\mathbf{E}[W_L]+(1-p)\mathbf{E}[W_H]}{\lambda(1-p)/\mu_H}$. ∎

**Remark 5**: The proof does not extend to other policies in Table 1 as their $\mathbf{E}[T_L]$ is not $\Theta(1)$.

**Remark 6**: For the metric of $\mathbf{E}[T]$, all three policies – OPT, P-LCFS and LCFS – are asymptotically optimal. However, $\mathbf{E}[T_L]$ under the three policies is different, although always $\Theta(1)$, and given by the following lemma, whose proof we omit.

**Lemma 4.** *When $\mu_H < \lambda$, $\mathbf{E}[T_L]$ under OPT, LCFS and P-LCFS are given by:*

$$\mathbf{E}\left[T_L^{OPT}\right] = \frac{1}{\mu_L - \lambda} + o(1)$$

$$\mathbf{E}\left[T_L^{P\text{-}LCFS}\right] = \mathbf{E}\left[B_L^L\right] + o(1) = \frac{1-\rho_H}{\mu_L(1-\rho)} + o(1)$$

$$\mathbf{E}\left[T_L^{LCFS}\right] = \theta_H(1 - \frac{\lambda}{\mu_L})\mathbf{E}\left[B_L^H\right] + \frac{\lambda}{\mu_L}\mathbf{E}\left[B_L^L\right] + \frac{1}{\mu_L} + o(1)$$

*where $\theta_H = \frac{(1-p)(\lambda-\mu_H)}{(1-p)\lambda+(p-\rho)\mu_H}$, and $\mathbf{E}\left[B_L^L\right]$ and $\mathbf{E}\left[B_L^H\right]$ are given in Corollary 2 (see Appendix B).*

**Remark 7**: Comparing with $\mathbf{E}\left[T_L^{P-LCFS}\right]_{\alpha\to\infty} = \frac{1}{\mu_L} \cdot \frac{1}{1-\rho}$, we see that the extreme correlated $\mathbf{E}[T_L]$ for P-LCFS is always lower than the uncorrelated $\mathbf{E}[T_L]$. We can prove a similar result for LCFS.

**Remark 8**: A further difference between the three policies emerges if one looks at higher order metrics, such as $\mathbf{E}\left[(T_L^\pi)^2\right]$. As a byproduct of the proof of Theorem 6 (Case 2), we can see that $\mathbf{E}\left[(T_L^{P-LCFS})^2\right] = \Omega\left(\frac{1}{\alpha}\right)$, while it is $\Theta(1)$ for OPT. Thus, while simple policies such as P-LCFS and LCFS are asymptotically optimal for $\mathbf{E}[T]$, learning-based scheduling policies might be preferred when one cares about more fine-grained metrics.

## 2.5. LAS

The asymptotic analysis of LAS presented below builds on the analysis under *i.i.d.* arrivals given in [7]. In short, to analyze the response time of a tagged arrival of size $x$, we consider a modified system where jobs of original size $s$ are truncated to size $\min\{s, x\}$ when they enter the system. Under LAS, the response time of the tagged arrival is given by the busy period generated by the work it sees on arrival in this modified system.

**Theorem 4.** *When $\mu_H < \lambda$, the mean response time of a job of size $x$ under the LAS scheduling policy is given by:*
***Case** $\lambda s_H(x) > 1$:*

$$\mathbf{E}\left[T_L^{LAS}(x)\right] = \frac{\mathbf{E}\left[W_L^*(x)\right]}{1-\rho(x)} + o\left(\frac{1}{\alpha}\right); \quad \mathbf{E}\left[T_H^{LAS}(x)\right] = \frac{1}{\alpha_H} + \frac{\mathbf{E}\left[W_H^*(x)\right] + \frac{\lambda s_H(x)-1}{\alpha_H}}{1-\rho(x)} + o\left(\frac{1}{\alpha}\right)$$

***Case** $\lambda s_H(x) < 1$:*

$$\mathbf{E}\left[T_L^{LAS}(x)\right] = \mathbf{E}\left[T_L^{M/M/1/LAS}(x)\right] + o(1); \quad \mathbf{E}\left[T_H^{LAS}(x)\right] = \mathbf{E}\left[T_H^{M/M/1/LAS}(x)\right] + o(1)$$

*where $\mathbf{E}\left[T_L^{M/M/1/LAS}(x)\right]$ and $\mathbf{E}\left[T_H^{M/M/1/LAS}(x)\right]$ denote the mean response time of a job of size $x$ under LAS scheduling in $M/M/1$ queues with arrival rate $\lambda$, and job size distribution $\text{Exp}(\mu_L)$ and $\text{Exp}(\mu_H)$, respectively.*

**Proof:** **Case $\lambda s_H(x) > 1$:** In this case, the modified system with truncated job sizes is in transient overload during the H states. Theorem 6, Case 1 (see Appendix B), gives us the expression for the required mean busy period.

**Case $\lambda s_H(x) < 1$:** In this case, the modified system with truncated job sizes is stable during the H states. As $\alpha \to 0$, the system looks like a mixture of two independent stable $M/G/1$ queues with the modified job size distributions (similar to Theorem 1). The mean response time of a type L job of size $x$ in this modified system thus converges to the mean response time of a job of size $x$ under an $M/M/1/LAS$ system with arrival rate $\lambda$ and job sizes *i.i.d.* $\text{Exp}(\mu_L)$. A similar argument applies to type H jobs of size $x$. ∎

**Remark 9**: Under *i.i.d.* $H_2$ job sizes, LAS is the optimal size-independent policy for minimizing $\mathbf{E}[T]$ because it isolates the class L jobs from class H jobs. Intuitively we expect this behavior to carry over when correlations are introduced, *but this is not the case*. Not only does LAS perform suboptimally, but $\mathbf{E}[T_L]$ under LAS grows as $\Theta\left(\frac{1}{\alpha}\right)$, while it is $\Theta(1)$ under LCFS and P-LCFS. The reason for this counter-intuitive behavior lies in the fraction of L jobs that do not get isolation and hence experience $\Theta\left(\frac{1}{\alpha}\right)$ mean response time. Under LCFS and P-LCFS, this fraction is $\Theta(\alpha)$ with a net effect of $\Theta(1)$. Under LAS, however, all L jobs with a size bigger than $\frac{1}{\mu_H} \log\left(\frac{\mu_H}{\lambda-\mu_H}\right)$, which is a $\Theta(1)$ fraction, experience $\Theta\left(\frac{1}{\alpha}\right)$ mean response time.

## 3. Evaluation via Simulations

While Section 2 provided fluid asymptotics as $\alpha \to 0$ for a wide range of size-independent scheduling policies, we are only able to perform exact numerical analysis of the case $0 < \alpha < \infty$ for a smaller subset (FCFS, LCFS, P-LCFS, OPT) via algorithms proposed in the supplement [15]. This section studies the full range of policies for all $\alpha$ via numerical techniques for the policies mentioned above, and via simulation for the remaining policies in Table 1. We start with results for our MMAP model and then present results for trace-based experiments.
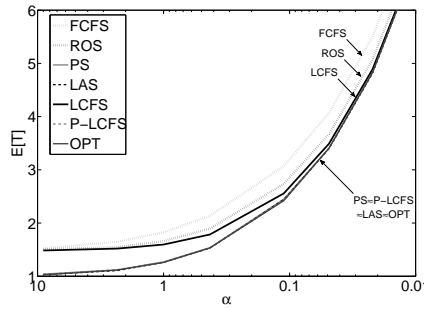


Figure 2: Effect of job size correlation when $\mu_H > \lambda$. The parameters chosen were $\mu_L = 50.73, \mu_H = 1.0055, p = 0.5073, \lambda = 1$ ($\rho = 0.5, C^2 \approx 2.9$).

**MMAP under No transient overload:** In Figure 2, we see the effect of correlation on scheduling policies when $\mu_H > \lambda$, so that there is no transient overload in H states. We see that for moderate $\alpha$, $\mathbf{E}[T]$ of the different scheduling policies range from $\mathbf{E}[T] = 1$ to about $\mathbf{E}[T] = 1.5$, with FCFS being the worst and LAS being the best. As $\alpha$ decreases, we see that the relative performance difference between scheduling policies begin to vanish ($\mathbf{E}[T]$ ranges from 6.9 to 7.5 for $\alpha \approx 0.01$). This behavior as $\alpha \to 0$ is consistent with Theorem 1. Observe also that while FCFS, ROS and LCFS are equal at the two extremes ($\alpha \to \infty$ and $\alpha \to 0$), for $0 < \alpha < \infty$ they are ordered as FCFS>ROS>LCFS with respect to $\mathbf{E}[T]$.

**MMAP under Transient overload:** Figure 3 shows the effect of correlation in the more interesting case of $\mu_H < \lambda$, implying that there is transient overload during the $H$ states. Figure 3(a) shows the $\mathbf{E}[T]$ vs. $\alpha$ curves for the different scheduling policies. We see that FCFS is the worst policy and LAS is optimal or close to optimal throughout the range of $\alpha$ shown. On the other hand, P-LCFS starts out equal to PS when $\alpha \to \infty$ and is clearly suboptimal; yet for low $\alpha$ (high correlation), P-LCFS approaches and even overtakes LAS, and becomes optimal. This is consistent with Theorem 3. Similarly, LCFS starts out equal to FCFS when $\alpha \to \infty$ and is worst in performance, but becomes optimal as $\alpha \to 0$, again confirming Theorem 3.

A major difference between Figure 3(a) (transient overload) and Figure 2 (no overload) is that the policies clearly do not converge to each other in Figure 3 as $\alpha \to 0$, whereas they do in Figure 2. Furthermore, for each policy $\pi$ in Figure 3(a), the $\mathbf{E}[T]$ curve asymptotes to a line on the plotted scale, which corresponds to $\mathbf{E}[T^\pi] \sim \frac{K^\pi}{\alpha}$ as in Lemma 3. Thus the mean response times grow unboundedly as $\alpha \to 0$, unlike in Figure 2.

Figure 3(b) verifies the expressions for $K^\pi$ obtained from our asymptotic analysis by showing $\left(\frac{\alpha}{1+\alpha}\right)\mathbf{E}[T]$ as a function of $\frac{1}{1+\alpha}$. We choose to scale $\mathbf{E}[T]$ by $\frac{\alpha}{1+\alpha}$ (instead of $\alpha$) to show the results for $\alpha \to \infty$ asymptote and the $\alpha \to 0$ asymptote in the same plot. In the former case, $\lim_{\alpha\to\infty} \frac{\alpha}{1+\alpha}\mathbf{E}[T^\pi] = \lim_{\alpha\to\infty} \mathbf{E}[T^\pi]$ and in the latter case $\lim_{\alpha\to 0} \frac{\alpha}{1+\alpha}\mathbf{E}[T^\pi] = \lim_{\alpha\to 0} \alpha\mathbf{E}[T^\pi] =$

$K^\pi$. The *x*-axis shows $\frac{1}{1+\alpha}$ which is bounded between 0 and 1 (unlike $\alpha$). The $\alpha \mathbf{E}[T^\pi]$ curves clearly converge to the analytically obtained values of $K^\pi$ marked with a small **x**. In the limit $\alpha \to 0$, $\mathbf{E}[T]$ for the different policies follows the order LCFS = P-LCFS < LAS < PS < ROS < FCFS. Due to the parameter settings, the difference between LAS and LCFS = P-LCFS as $\alpha \to 0$ is very slight; this contrasts with Figure 1 where the difference was significant.
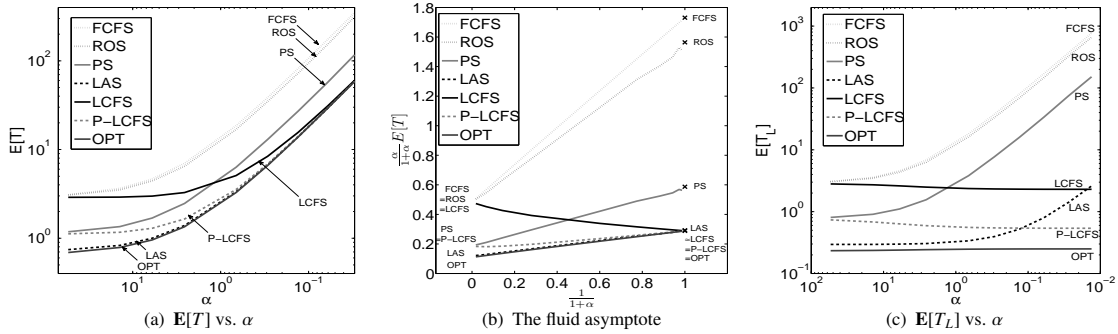


Figure 3: Effect of job size correlation when $\mu_H < \lambda$. The parameters chosen were $\mu_L = 10, \mu_H = 1, p = 0.95, \lambda = 6$ ($\rho = 0.87, C^2 \approx 4.66$).

Figure 3(c) shows mean response time for "little" (class L) jobs, denoted $\mathbf{E}[T_L]$, versus $\alpha$. For the L jobs, there is a wide difference (several orders of magnitude) in performance across policies. Several policies (FCFS, ROS, PS, LAS) show $\mathbf{E}[T_L]$ increasing in proportion to $\frac{1}{\alpha}$ (though this is less obvious in the case of LAS); however, other policies (LCFS, P-LCFS) show a decrease in $\mathbf{E}[T_L]$ as $\alpha$ decreases, as pointed out in Remark 7. Under the first group of polices, $\mathbf{E}[T_L]$ suffers from increased correlation, because L jobs are affected by H jobs. For LCFS and P-LCFS, this is not the case, since an L job is only affected by H jobs if they arrive during the L job's busy period. This happens with probability proportional to $\alpha$, which becomes zero as $\alpha \to 0$.

**Trace-based experiments:** While we garnered useful intuition by analyzing the MMAP correlation model, it is not obvious to what extent our results would extend to real-world applications. To investigate this, we consider two very different traces, one involving packets sizes (Bellcore) and a second involving supercomputing job sizes (SHARCNET). We have simulated FCFS, ROS, PS, LCFS and P-LCFS policies. In addition, we simulate PRIO-P, which gives preemptive priority to class L jobs, where class L jobs are defined as jobs with size below some threshold. Hence the PRIO-P policy is similar to the OPT policy, but is not necessarily the optimal size-independent policy because class L and H jobs are no longer Exponentially distributed. We also simulate SRPT (Shortest Remaining Processing Time) policy, and our plots show $\mathbf{E}[T]$ under the simulated policies normalized by the mean response time under SRPT scheduling.

Figures 4(a)-4(d) show the results of our experiments with a trace of packet sizes seen on the Bellcore Ethernet [13]. The autocorrelation function of packet sizes (Figure 4(a)) shows significant sequential job size correlation – the lag-1 correlation is approximately 0.45 with correlation persisting even at lags of up to 100 (unlike MMAP model where the correlation decreases exponentially in lag). Figure 4(b) shows the job size distribution which is almost a trimodal distribution. To perform the simulations, we modify the base trace as follows: In the first set of experiments (Figure 4(c)), we scale the interarrival times from the trace to vary the 'load'. In the second set of experiments (Figure 4(d)), we keep the same sequence of job sizes as the original trace, but create a new Poisson arrival process to eliminate correlations in the arrival process
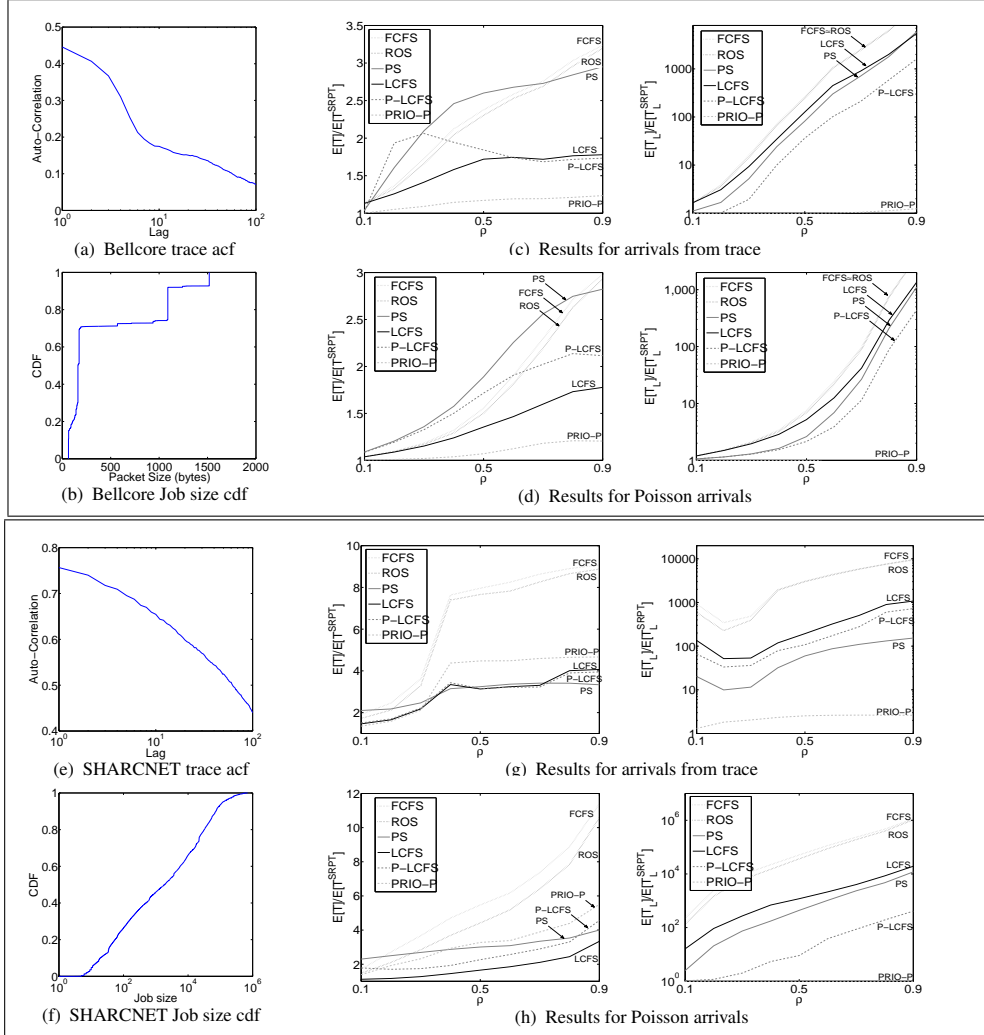
Figure 4: Trace-based experiments. Simulation results for the Bellcore trace are shown in the top box, and for the SHARCNET trace in the bottom box. For each set of traces, the top-left plot shows the autocorrelation function for job size sequence; the bottom-left plot shows the cdf of the job size distribution; the two top-right plots show the performance (as the ratio of $\mathbf{E}[T]$ to $\mathbf{E}\left[T^{SRPT}\right]$, and of $\mathbf{E}[T_L]$ to $\mathbf{E}\left[T_L^{SRPT}\right]$, respectively) when the interarrival times are taken from the trace; the two bottom right plots show the performance obtained by creating a synthetic Poisson arrival process.

(the trace arrival process is bursty) and correlations between interarrival times and job sizes (the correlations between a job size and immediately following interarrival time is −0.15). We see that with respect to $\mathbf{E}[T]$, the ordering of the policies largely obeys FCFS≈ROS≈PS>LCFS≈P-LCFS>PRIO-P>SRPT. This is consistent with the ordering we obtained via analysis using the MMAP correlation model. We also see that $\mathbf{E}\left[T^{FCFS}\right]$ is up to 1.8 times worse than $\mathbf{E}\left[T^{LCFS}\right]$ which contrasts with the uncorrelated case where they are equal. We also investigate the effect of scheduling on the little jobs by classifying packets of size less than 400 bytes as L. Under our criterion, the L jobs make up 70% of the packets, and 25% of the total bytes. We find that

$\mathbf{E}\left[T_L^{FCFS}\right]$ is up to 3 to 4 times worse than $\mathbf{E}\left[T_L^{LCFS}\right]$ and almost 10 times worse than $\mathbf{E}\left[T_L^{P-LCFS}\right]$. We also see that PS outperforms LCFS but not P-LCFS in terms of $\mathbf{E}[T_L]$. This can be explained by the fact that under the uncorrelated case PS and P-LCFS have identical performance and outperform LCFS which suffers due to job size variability. Under moderate correlation, we see a behavior that is the mixture of uncorrelated and high-correlation cases: job size variability is still hurting class L jobs under LCFS and thus gives them worse performance than PS, however due to correlation P-LCFS is able to perform better than PS (our MMAP simulation results also suggest that for moderate correlations, PS still outperforms LCFS). The same observations hold under a Poisson arrival process, but the gains are more moderate. This suggests that in the presence of cross-correlations and bursty arrivals, the effect of scheduling will be more pronounced.

Figures 4(e)-4(h) show the results for the SHARCNET trace [1], which is a supercomputing workload. Here job size is defined as the run time of jobs submitted to the server, and the correlation in the sequence of job sizes is very high (lag-1 autocorrelation is over .7, and even lag-100 correlation is over .4). The ordering of policies with respect to $\mathbf{E}[T]$ largely obeys FCFS>ROS>PRIO-P>PS≈P-LCFS≈LCFS>SRPT. The gains of utilizing LCFS instead of FCFS for the SHARCNET trace are even more significant, as the ratio of $\mathbf{E}\left[T^{FCFS}\right]$ to $\mathbf{E}\left[T^{LCFS}\right]$ can be over 2. For the SHARCNET trace, we defined L jobs as those smaller than 54000 seconds ( 86% jobs, 25% of total load). There is again a significant difference between $\mathbf{E}\left[T_L^{FCFS}\right]$ and $\mathbf{E}\left[T_L^{LCFS}\right]$, up to 4X when scaling the original interarrival times, and 15X to 20X when the arrival process has been converted to a Poisson process. Comparing $\mathbf{E}[T_L]$ for LCFS, PS and P-LCFS, we see that PS does better than LCFS which can be explained by the presence of job size variability. However the ordering of PS and P-LCFS under arrival times from the SHARCNET trace switches when a Poisson arrival process is considered. While under a Poisson arrival process, PS performs worse than P-LCFS as predicted by our analysis of the MMAP correlation model, under the arrival sequence from the SHARCNET trace, PS outperforms P-LCFS. This suggests that the correlation between the arrival times (the SHARCNET arrival sequence has extremely bursty and variable interarrival times compared to the Bellcore trace) is also an important aspect to consider to fully understand the effect of scheduling under correlated traffic pattern.

## 4. Conclusions

To the best of our knowledge, this is the first paper to study analytically how common scheduling policies, like PS, LAS, ROS, P-LCFS, LCFS, etc. are affected by correlation among consecutive job sizes. We find the ranking of scheduling policies, from highest to lowest mean response time ($\mathbf{E}[T]$), changes dramatically under correlation: LCFS which performs poorly under no correlation becomes optimal among size-independent policies under high correlation; the optimal size-independent policy for i.i.d. job sizes, LAS, becomes sub-optimal under high correlation; the mean response times of policies which are insensitive to job-size variability when job sizes are i.i.d., like PS and P-LCFS, now depend on the entire job-size distribution, to cite a few examples. When examining the mean response time of "little" jobs only ($\mathbf{E}[T_L]$), the change in ranking is even more dramatic, with correlation actually making some policies like LCFS and P-LCFS perform better, and making other policies like LAS perform far worse.

We have only scratched the surface of how correlation in job sizes affects performance. First, our correlation model is very simple, chosen for analytical tractability and to gain insights; extending the results presented here to richer models is left for future work. Second, while this paper shows that P-LCFS and LCFS perform optimally among size-independent policies under very

high correlation, the paper does not answer the question of which policy is best under moderate correlation. Furthermore, we have not even explored policies which might exploit the correlation structure to improve performance. Third, our model only captures correlations in consecutive job sizes, but we believe that the techniques introduced herein can be applied to understanding the effect of all three types of correlation on the performance of scheduling policies.

## Appendix A. Transforms for stationary workload

**Theorem 5.** *Let $\widetilde{W_L}(s)$ and $\widetilde{W_H}(s)$ denote the transform for the stationary workloads during the L and H states, respectively, under the MMAP model. Then:*

$$\widetilde{W_L}(s) = \frac{(1 - \rho)\alpha m_L m_H - s m_L g_H \pi_L(0)}{\alpha_L g_H m_L + \alpha_H g_L m_H - s g_L g_H} \tag{A.1}$$

*where,*

$$
\begin{array}{llll}
m_L = \mu_L + s & ; & m_H = \mu_H + s & \\
g_L = \mu_L - \lambda + s & ; & g_H = \mu_H - \lambda + s &
\end{array}
; \quad \pi_L(0) = \frac{(1 - \rho)\alpha(\mu_H + \xi)}{\xi(\mu_H - \lambda + \xi)}
$$

*and $\xi$ denotes the unique root of the denominator of* (A.1) *in the interval* $(0, +\infty)$. *The quantity $\pi_L(0)$ denotes the long run fraction of time that the system is empty conditioned on being in state L. The expression for $\widetilde{W_H}(s)$ is obtained by flipping $\mu_H$ and $\mu_L$, and flipping $\alpha_L$ and $\alpha_H$.*

**Proof:** The first step is analysis of the transient workload in an $M/G/1$. Consider an $M/G/1$ with arrival rate $\lambda$, *i.i.d.* job sizes $X_1, X_2, \ldots$ with Laplace transform of the job size distribution given by $\mathbf{E}\left[e^{-sX_1}\right] = \widetilde{X}(s)$. We can write the following equation for the evolution of the workload $W(t)$ in this $M/G/1$:

$$W(t + \delta t) = W(t) - \delta t \mathbf{1}_{W(t)>0} + \sum_n X_n \mathbf{1}_{n\text{th arrival in } (t,t+\delta t)}$$

Let $\widetilde{W_t}(s) = \mathbf{E}\left[e^{-sW(t)}\right]$. Taking Laplace transforms in the above equation, and then letting $\delta t \to 0$,

$$\frac{d}{dt}\widetilde{W_t}(s) = \widetilde{W_t}(s)\left(s - \lambda(1 - \widetilde{X}(s))\right) - s\mathbf{Pr}[W_t = 0]$$

Let $T$ be an $\mathrm{Exp}(\nu)$ random variable and $\widetilde{W_T}(s) = \mathbf{E}\left[e^{-sW(T)}\right]$. Using integration by parts, we get:

$$
\begin{aligned}
\widetilde{W_T}(s) &\equiv \int_{u=0}^{\infty} \widetilde{W_u}(s)\nu e^{-\nu u}du = \left[\frac{\widetilde{W_u}(s)\nu e^{-\nu u}}{-\nu}\right]_{u=0}^{\infty} + \int_{u=0}^{\infty} \frac{d\widetilde{W_u}(s)}{du}e^{-\nu u}du \\
&= \widetilde{W_0}(s) + \frac{1}{\nu}\int_{u=0}^{\infty}\left(\widetilde{W_u}(s)\left[s - \lambda(1 - \widetilde{X}(s))\right] - s\mathbf{Pr}[W_u = 0]\right)\nu e^{-\nu u}du \\
&= \widetilde{W_0}(s) + \frac{1}{\nu}\left(\widetilde{W_T}(s)[s - \lambda(1 - \widetilde{X}(s))] - s\mathbf{Pr}[W(T) = 0]\right)
\end{aligned}
$$

Specializing to our problem, we obtain the following two relations by applying the above equation during L and H states, and noting that by PASTA $\widetilde{W_L}(s)$ and $\widetilde{W_H}(s)$ also denote the stationary workloads at the *ends* of L and H states, respectively:

$$\widetilde{W_L}(s) = \widetilde{W_H}(s) + \frac{s}{\alpha_L}\left[\widetilde{W_L}(s)\left(1 - \frac{\lambda}{\mu_L + s}\right) - \pi_L(0)\right]$$

$$\widetilde{W_H}(s) = \widetilde{W_L}(s) + \frac{s}{\alpha_H}\left[\widetilde{W_H}(s)\left(1 - \frac{\lambda}{\mu_H + s}\right) - \pi_H(0)\right]$$

Eliminating $\widetilde{W_H}(s)$, and $\pi_H(0)$ by using the fact $\frac{\pi_L(0)}{\alpha_L} + \frac{\pi_H(0)}{\alpha_H} = (1 - \rho)\left(\frac{1}{\alpha_L} + \frac{1}{\alpha_H}\right)$, we obtain the expression for $\widetilde{W_L}(s)$ shown in the Theorem. It now remains to determine the unknown $\pi_L(0)$. To obtain this, we note that the polynomial in the denominator of $\widetilde{W_L}(s)$ is a cubic in $s$ which approaches $-\infty$ as $s \to \infty$. Further, the denominator is positive at $s = 0$ but negative at $s = \lambda - \mu_L < 0$. Therefore there is exactly one root of the denominator in the interval $(0, +\infty)$, which we denote by $\xi$, at which there is a degeneracy in the denominator. Since the transform must converge in $Re(s) > 0$, the numerator must share this root, yielding the unknown $\pi_L(0)$. ∎

## Appendix B. Asymptotic Expressions for Mean Busy Periods

Busy periods form the core of the analysis for scheduling policies, and therefore we deal with the problem of analyzing busy periods in as much generality as possible.

We consider a system with an environment controlled by a 2-state Markov chain with states L and H. The time spent in state L during each visit is $\text{Exp}(\alpha_L)$ and time spent in state H is $\text{Exp}(\alpha_H)$. Let $\alpha = \alpha_L + \alpha_H$, $p = \frac{\alpha_H}{\alpha}$. The arrivals occur at a rate $\lambda$ in each state. The arrivals during an L state have *i.i.d.* general job sizes and are denoted by random variable $S_L$. Similarly, the arrivals during an H state have *i.i.d.* general job sizes denoted by random variable $S_H$. We will assume $\mathbf{E}[S_L] < \mathbf{E}[S_H]$. We index this system by $\alpha$.

**The scaling**: We consider a sequence of systems, indexed by $\alpha$, obtained by setting the switching rates as $\alpha_L + \alpha_H = \alpha$, while fixing $p = \frac{\alpha_H}{\alpha}$. We start the $\alpha$th system in a prescribed state with initial workload (a random variable) denoted by $W_\alpha$. We will say that the workload sequence $W_\alpha$ is $\Theta(g(\alpha))$ if the sequence $\left\{\frac{W_\alpha}{g(\alpha)}\right\}$ is uniformly integrable and $\lim_{\alpha \to 0} \frac{W_\alpha}{g(\alpha)} \xrightarrow{d} \overline{W}$, where $\overline{W}$ is some non-degenerate random variable. Similarly, we say $W_\alpha = o(h(\alpha))$ if $W_\alpha = \Theta(g(\alpha))$ and $\lim_{\alpha \to 0} \frac{g(\alpha)}{h(\alpha)} = 0$, or $W_\alpha = \omega(h(\alpha))$ if $W_\alpha = \Theta(g(\alpha))$ and $\lim_{\alpha \to 0} \frac{h(\alpha)}{g(\alpha)} = 0$.

**Goal**: Let $B_L(W_\alpha)$ and $B_H(W_\alpha)$ denote the random variables for the busy periods started by work $W_\alpha$ in states L and H, respectively, in the $\alpha$th system. We will be interested in obtaining the mean busy period in the asymptotic regime $\alpha \to 0$. That is, we are interested in obtaining the dominant term in $\mathbf{E}[B_L(W_\alpha)]$ or $\mathbf{E}[B_H(W_\alpha)]$, as the switching rate $\alpha \to 0$.

**Notation**: $\widetilde{S_L}(s) = \mathbf{E}\left[e^{-sS_L}\right]$; $\widetilde{S_H}(s) = \mathbf{E}\left[e^{-sS_H}\right]$

$$r_L = 1 - \lambda\mathbf{E}[S_L]; \quad r_H = 1 - \lambda\mathbf{E}[S_H]; \quad \rho = \lambda(p\mathbf{E}[S_L] + (1 - p)\mathbf{E}[S_H])$$

We first present the theorems on asymptotic expressions for the mean busy periods. After presenting the theorems, we first present a brief proof sketch to elucidate how the theorems were derived, and then the detailed proofs. Theorem 6 considers the case $\lambda\mathbf{E}[S_H] > 1$, and Theorem 7 considers the case $\lambda\mathbf{E}[S_H] < 1$.

**Theorem 6.** *Let $r_H < 0$. That is, the system is under temporary overload during H states.*
**Case 1:** $W_\alpha = \omega(1)$, $\mathbf{Pr}\left[\overline{W} = 0\right] = 0$:

$$\mathbf{E}[B_L(W_\alpha)] = \frac{\mathbf{E}[W_\alpha]}{1 - \rho} + o(W_\alpha)$$

$$\mathbf{E}[B_H(W_\alpha)] = \frac{\mathbf{E}[W_\alpha] + \frac{1 - \rho - r_H}{\alpha_H}}{1 - \rho} + o(\max\left\{W_\alpha, \alpha^{-1}\right\})$$

**Case 2:** $W_\alpha = \Theta(1)$:

$$\mathbf{E}[B_L(W_\alpha)] = \frac{\mathbf{E}\left[\overline{W}\right]}{r_L} + p_{switch} \cdot (1 - Q_f)\frac{1 - \rho - r_H}{\alpha_H(1 - \rho)} + o(1)$$

$$\mathbf{E}[B_H(W_\alpha)] = (1 - P_f) \cdot \frac{\mathbf{E}\left[\overline{W}\right] + \frac{1-\rho-r_H}{\alpha_H}}{1 - \rho} + O(1)$$

*where, $p_{switch}$ denotes the probability that the environment state switches to H before the busy period started by $\overline{W}$ in state L ends. We call this event a 'switch'. The expression for $p_{switch}$ is given by $p_{switch} = \frac{\mathbf{E}\left[\overline{W}\right]\alpha_L}{r_L} + o(\alpha)$. The quantity $Q_f$ denotes the probability that, given a 'switch' occurs, the residual busy period is finite if the H state were to last indefinitely from then on:*

$$Q_f = \widetilde{V}(\lambda(1 - p_f)) + o(1)$$

*where $\widetilde{V}(\cdot)$ is given by [4]: $\widetilde{V}(s) = \frac{r_L \cdot \frac{1-\widetilde{\overline{W}}(s)}{\mathbf{E}[W]}}{s - \lambda(1 - \widetilde{S}_L(s))}$, and $p_f \in (0, 1)$ solves the fixed point equation [5]: $p_f = \widetilde{S_H}(\lambda(1 - p_f))$.*

*The quantity $P_f$ denotes the probability that the busy period started by $\overline{W}$ during an H state is finite if the H state were to last indefinitely and is given by $P_f = \widetilde{\overline{W}}(\lambda(1 - p_f))$.*

**Corollary 2.** *Consider the case $S_L \sim \text{Exp}(\mu_L)$ and $S_H \sim \text{Exp}(\mu_H)$, $\mu_L > \lambda > \mu_H$. Let $B_s^c$ $(c, s \in \{L, H\})$ denote the busy period duration started by a class c job in environment state s. Then,*

$$\mathbf{E}\left[B_L^L\right] = \frac{1}{\mu_L - \lambda}\left(1 + \frac{1 - p}{p} \cdot \frac{\lambda - \mu_H}{\mu_L - \mu_H} \cdot \frac{\frac{\lambda}{\mu_H} - \rho}{1 - \rho}\right) + o(1);$$

$$\mathbf{E}\left[B_L^H\right] = \frac{\mu_L}{\mu_H(\mu_L - \lambda)}\left(1 + \frac{1 - p}{p}(1 - Q_{f_H})\frac{\frac{\lambda}{\mu_H} - \rho}{1 - \rho}\right) + o(1)$$

*and:*

$$\mathbf{E}\left[B_H^H\right] = \left(1 - \frac{\mu_H}{\lambda}\right) \cdot \frac{\frac{\lambda}{\mu_H} - \rho}{\alpha_H(1 - \rho)} + o(\alpha^{-1})$$

$$\mathbf{E}\left[B_H^L\right] = \left(1 - \frac{\mu_L}{\mu_L + \lambda - \mu_H}\right) \cdot \frac{\frac{\lambda}{\mu_H} - \rho}{\alpha_H(1 - \rho)} + o(\alpha^{-1}).$$

*In the above, $1 - Q_{f_H} = 1 - \widetilde{V_H}(\lambda(1 - \phi_f))$, where $\phi_f = \frac{\mu_H}{\lambda}$, and $\widetilde{V_H}(s) = \frac{\left(1 - \frac{\lambda}{\mu_L}\right)\left(\frac{\mu_H}{\mu_H + s}\right)}{1 - \frac{\lambda}{\mu_L}\left(\frac{\mu_L}{\mu_L + s}\right)}$.*

---

[4] The function $\widetilde{V}(s)$ denotes the Laplace transform of the workload in the system just before the 'switch' event occurs. $\widetilde{V}(s)$ is obtained as the Laplace transform of the stationary workload conditioned on server being busy in an $M/G/1$ with repeated vacations, with service distribution $S_L$ and *i.i.d.* vacations distributed as $\overline{W}$.

[5] The quantity $p_f$ denotes the probability that a busy period started by an H job in an H state is finite if the H state were to last indefinitely.

**Theorem 7.** *Let $r_H > 0$. That is, the system is stable during H states.*
**Case 1:** $W_\alpha = \omega(\alpha^{-1})$

$$\mathbf{E}[B_L(W_\alpha)] = \frac{\mathbf{E}[W_\alpha]}{1-\rho} + o(W_\alpha) \quad ; \qquad \mathbf{E}[B_H(W_\alpha)] = \frac{\mathbf{E}[W_\alpha]}{1-\rho} + o(W_\alpha).$$

**Case 2:** $W_\alpha = \Theta(\alpha^{-1})$

$$\mathbf{E}[B_L(W_\alpha)] = \frac{\mathbf{E}[W_\alpha]}{1-\rho}(1-u_\alpha) + \frac{\mathbf{E}[W_\alpha]}{r_L}u_\alpha + o(\alpha^{-1})$$

$$\mathbf{E}[B_H(W_\alpha)] = \frac{\mathbf{E}[W_\alpha]}{1-\rho}(1-u_\alpha) + \frac{\mathbf{E}[W_\alpha]}{r_H}u_\alpha + o(\alpha^{-1})$$

*where* $u_\alpha \equiv \left[\frac{1-\widetilde{W_\alpha}\left(\frac{\alpha_L}{r_L}+\frac{\alpha_H}{r_H}\right)}{\mathbf{E}[W_\alpha]\left(\frac{\alpha_L}{r_L}+\frac{\alpha_H}{r_H}\right)}\right]$, $0 < u_\alpha < 1$, *and* $\lim_{\alpha\to 0} u_\alpha = u = \left[\frac{1-\widetilde{\overline{W}}\left(\frac{1-p}{r_L}+\frac{p}{r_H}\right)}{\mathbf{E}[\overline{W}]\left(\frac{1-p}{r_L}+\frac{p}{r_H}\right)}\right]$; *and recall*
$\overline{W} = \lim_{\alpha\to 0} \alpha W_\alpha$.
**Case 3:** $W_\alpha = o(\alpha^{-1})$

$$\mathbf{E}[B_L(W_\alpha)] = \frac{\mathbf{E}[W_\alpha]}{r_L} + o(W_\alpha); \quad \mathbf{E}[B_H(W_\alpha)] = \frac{\mathbf{E}[W_\alpha]}{r_H} + o(W_\alpha).$$

**Proof Sketch of Theorems 6 and 7:** Recall our fluid model, in which the workload decreases at deterministic rate $r_L$ during the L states, and increases at rate $-r_H$ during the H states. We would like to believe that given an initial workload $W_\alpha$, asymptotically the busy period started by $W_\alpha$ is the same as the duration of the busy period started by $W_\alpha$ under the fluid model. However, this is only partially true. When $W_\alpha = \Theta(\alpha^{-1})$, this asymptotic equivalence is justified by [4, Theorem 1(b)] which proves the convergence of workload sample paths of the stochastic and fluid systems (although one needs to do a little more work to convert it to convergence of busy periods). For the remaining cases, we must consider the tree of events that may occur until each leaf corresponds to an empty system, or one with workload that is $\Theta(\alpha^{-1})$ so that we can apply [4, Theorem 1(b)]. We describe this below.

**Case:** $W_\alpha = \omega(\alpha^{-1})$: In this case, the initial workload is of a higher order than the scale at which the system switches. Thus, asymptotically, the number of times the system switches states before $W_\alpha$ drains goes to $\infty$ as $\alpha \to 0$, and the workload sees the "average system" during its sojourn. Thus the mean busy period is $\frac{\mathbf{E}[W_\alpha]}{1-\rho} + o(W_\alpha)$.

**Case:** $W_\alpha = \Theta(\alpha^{-1})$: As noted above, in this case from [4, Theorem 1(b)] asymptotically the mean busy period is given by the busy period under the fluid model. The final expressions are obtained by setting up and solving recurrences for the mean busy period under the fluid model.

**Remark 10**: When $r_H > 0$, the mean busy period started in state $s$ is a convex combination of the busy period if the state $s$ were to last indefinitely, and the busy period of the "average system", with the coefficient being a function of the Laplace transform of the workload.

**Case:** $W_\alpha = o(\alpha^{-1})$, $r_H > 0$: In this case, the system is stable in both states. Consider a busy period starting in state L. If the L state were to last forever, the busy period would exactly be $\frac{\mathbf{E}[W_\alpha]}{r_L}$. However, since we may switch at rate $\Theta(\alpha)$, there is a $o(1)$ probability that the system switches to state H before the busy period finishes. If this switch were to happen, the remaining busy period would be stochastically bounded by a $\Theta(W_\alpha)$ random variable, as the system is always stable, thus giving a $o(W_\alpha)$ contribution to the overall busy period after multiplying by

the probability of switching. Thus asymptotically, the mean busy period started by $W_\alpha$ workload in state L would be $\frac{\mathbf{E}[W_\alpha]}{r_L} + o(W_\alpha)$.

**Case:** $W_\alpha = \Theta(1)$, $r_H < 0$: This case is the most non-trivial of all, and clearly explains the failure of fluid modeling of busy periods. First, consider a busy period started in state H by $W_\alpha = \Theta(1)$ work. The fluid model would imply that the workload keeps increasing at rate $-r_H$ until the system switches to L. At this point we have $\Theta(\alpha^{-1})$ workload built up, and we could apply [4, Theorem 1(b)]. *However, given that we start with $\Theta(1)$ workload in state H (which is in transient overload), there is still a constant ($\Theta(1)$) probability that the stochastic busy period started by the $\Theta(1)$ workload is finite*! This probability is denoted by $P_f$ in Theorem 6, and given that this event does not happen, we can use the fluid busy period expressions. Next, consider a busy period started in state L by $W_\alpha = \Theta(1)$ work. In this case, with $\Theta(\alpha)$ probability (given by $p_{switch}$), there is a class H arrival before the busy period ends. We are now in state H with $\Theta(1)$ workload (whose transform is given by $\widetilde{L}(s) \cdot \widetilde{S}_H(s)$). Given that a class H arrival happens, the residual busy period (from our argument above) is $\Theta(\alpha^{-1})$. After multiplying it with $p_{switch}$, we see that the contribution of this term to the overall busy period is $\Theta(1)$, and hence is of the same asymptotic order as the duration of the busy period started in state L conditioned on it ending in state L ($= \frac{\mathbf{E}[W_\alpha]}{r_L} + o(1)$). Therefore, we need to be precise with each of the terms involved, and applying the fluid method does not yield the correct expressions. ∎

**Proof of Theorem 6:**

**Case 1:** $W_\alpha = \omega(1)$, $\mathbf{Pr}\left[\overline{W} = 0\right] = 0$: We first show that under the fluid regime, the expressions for the busy periods are as given. Then we will argue that when $W_\alpha = \omega(1)$, the fluid approximation for the mean busy period is asymptotically the same as the stochastic busy period. Let $W_\alpha$ be deterministic $x$, and $\tau_L \sim \text{Exp } \alpha_L$. Then we can write the following recurrence relation for the fluid busy period started in L or H state by workload $x$.

$$\mathbf{E}[B_H(x)] = \frac{1}{\alpha_H} + \mathbf{E}\left[B_L\left(x - \frac{r_H}{\alpha_H}\right)\right]$$

$$\mathbf{E}[B_L(x)] = \mathbf{E}\left[\min\left\{\frac{x}{r_L}, \tau_L\right\}\right] + \mathbf{E}\left[B_H\left(x - r_L \min\left\{\frac{x}{r_L}, \tau_L\right\}\right) \cdot \mathbf{1}_{\{x > r_L \cdot \tau_L\}}\right]$$

Now we assume $\mathbf{E}[B_L(x)] = b_L x$ and $\mathbf{E}[B_H(x)] = a_H + b_H x$ for some constants $b_L, a_H, b_H$, and then verify that these forms are indeed correct by identifying the unknown constants. Under the assumed forms for fluid busy periods, the recurrences reduce to:

$$a_H + b_H x = \frac{1}{\alpha_H} + b_L x - b_L \frac{r_H}{\alpha_H}; \quad b_L x = \frac{1 - e^{-\frac{\alpha_L}{r_L} x}}{\alpha_L} + a_H(1 - e^{-\frac{\alpha_L}{r_L} x}) + b_H x - b_H r_L \frac{1 - e^{-\frac{\alpha_L}{r_L} x}}{\alpha_L}$$

Since the above equations should be satisfied for all $x$, we get $b_L = b_H = \frac{1}{1-\rho}$ and $a_H = \frac{1-\rho-r_H}{\alpha_H(1-\rho)}$ yielding the expressions in the theorem statement.

Now we verify that when $W_\alpha = \omega(1)$, the fluid busy period expressions are asymptotically correct. In the simple case $W_\alpha = \omega(\alpha^{-1})$, the system switches on a faster time-scale ($\Theta(\alpha^{-1})$) than the initial amount of work ($\omega(\alpha^{-1})$). Thus this workload sees the "average" system (rather than the transient system) and its busy period is simply $\frac{\mathbf{E}[W_\alpha]}{1-\rho} + o(W_\alpha)$.

When the workload is $\Theta(\alpha^{-1})$, then using [4], the sample paths of the stochastic system (scaled by $\alpha$) converge as $\alpha \to 0$ to the fluid sample path in the space $D[0, \infty)$. Thus, the mean busy period of the stochastic system is within $o(\alpha^{-1})$ of the mean busy period of the fluid system.

Now consider the case $W_\alpha = \Theta(g(\alpha))$ where $g(\alpha) = \omega(1)$, but $g(\alpha) = o(\alpha^{-1})$ (e.g., $g(\alpha) = \frac{1}{\sqrt{\alpha}}$).
**Subcase 1:** Busy period beginning in state H: We will show that even though the initial workload is $o(\alpha^{-1})$, since it is $\omega(1)$, with overwhelming probability, the sample paths will follow the fluid trajectory. Let $\widetilde{W_\alpha}(s) = \mathbf{E}\left[e^{-sW_\alpha}\right]$. Since reordering the jobs served in a busy period does not change the busy period duration, consider the case where the initial workload $W_\alpha$ is served first. If the H state were to last forever, the $z$−transform for the number of arrivals of class H jobs while workload $W_\alpha$ is served is given by $\widetilde{W_\alpha}(\lambda(1 - z))$. The main idea is to show that since the H state is in overload, with probability tending to 1, at least one of the class H jobs will start a busy period that lasts until the end of the H state, whereby by the Strong Law of Large Numbers the accumulated workload will be $\Theta(\alpha^{-1})$. Consider the busy period that one class H job starts, provided the H state continues forever. The Laplace transform of the busy period in an $M/G/1$ with only class H jobs, $\widetilde{B_H}(s)$, satisfies:

$$\widetilde{B_H}(s) = \widetilde{S_H}(s + \lambda(1 - \widetilde{B_H}(s))).$$

Since the $M/G/1$ is in overload, there is a constant probability that the busy period is infinite. The probability that the busy period is finite is obtained as

$$p_f = \lim_{s \to 0} \widetilde{B_H}(s).$$

Taking limit in the expression for $\widetilde{B_H}(s)$, we obtain:

$$p_f = \widetilde{S_H}(\lambda(1 - p_f))$$

The busy period started by $W_\alpha$, given the H phase lasts forever, is finite if and only if the busy period started by each H arrival while $W_\alpha$ was served is finite. This probability, then is given by

$$\mathbf{Pr}[\text{busy period started during H is finite}] = \sum_{i=0}^{\infty} \mathbf{Pr}[i \text{ arrivals during } W_\alpha]\cdot p_f^i = \widetilde{W_\alpha}(\lambda(1-p_f)) \to 0.$$

The last fact is true since $\frac{W_\alpha}{g(\alpha)} \to \overline{W}$, $\widetilde{W_\alpha}(s) \to \widetilde{\overline{W}}(s \cdot g(\alpha)) \to 0$ as $\alpha \to 0$ ($\widetilde{\overline{W}}(s)$ is a decreasing function from 1 to 0, and $g(\alpha) = \omega(1)$). The fact that $\lim_{s\to\infty} \widetilde{\overline{W}}(s) = 0$ follows from the assumption $\mathbf{Pr}\left[\overline{W} = 0\right] = 0$.
Therefore, with probability approaching 1, the busy period started by $W_\alpha$ in phase H (under the assumption that the H phase lasts forever) is not finite. In other words, during the H phase, the workload increases asymptotically along the fluid trajectory, and then the system switches to the L phase. Since the work built up during the H state is $\Theta(\alpha^{-1})$, the workload follows the fluid trajectory after switching to the L state. Therefore, the expression for the mean busy period started in H phase by $\omega(1)$ work is indeed given by the mean busy period under the fluid regime within a $o(\max\{W_\alpha, \alpha^{-1}\})$ term.
**Subcase 2:** Busy period beginning in state L: Now we consider the case where the busy period starts in the L phase. If the L phase were to last forever, the workload in the system, scaled by $g(\alpha)$, would follow the fluid trajectory, and hence the mean busy period would be the mean busy period under the fluid regime within a $o(W_\alpha)$ term. However, with probability $\Theta(\alpha \cdot g(\alpha))$ the system switches to H state before the fluid workload reaches 0. Conditioned on switching to the H state before the period ends, the workload at the beginning of the H state is again $\Theta(g(\alpha))$. We have already argued above that subsequently the workload follows the fluid trajectory – and

the residual busy period will be $\Theta(\alpha^{-1})$ within an $o(\alpha^{-1})$ term. Therefore, the mean busy period started in L phase will be the mean busy period under the fluid regime, within a $o(W_\alpha)$ term.

**Case 2:** $W_\alpha = \Theta(1)$: We first consider the case where the busy period begins in the H state by workload $\overline{W}$ with Laplace transform $\widetilde{\overline{W}}(s)$. As we have argued above, since the H state is in overload, there is a constant probability that the busy period does not end before the system switches to the L state. This probability is given by $1 - P_f$ where,

$$P_f = \widetilde{\overline{W}}(\lambda(1 - p_f))$$

and $p_f$ is the solution to the fixed point equation $p_f = \widetilde{S_H}(\lambda(1 - p_f))$. $P_f$ denotes the probability that a busy period started by work $\overline{W}$ in the $M/G/1$ under overload is finite, and $p_f$ is the probability that a busy period started by a single class $H$ job is finite.

Given that the busy period does not end before the system switches, the work that builds up in the system is given by $\tau_H(\frac{\lambda}{\mu_H} - 1) + o(\alpha^{-1})$ where $\tau_H$ denotes the duration of the H state and is $\Theta(\alpha^{-1})$. We can thus apply the previous case and conclude that the mean busy period in this case, that is with probability $1 - P_f$, is given by $\frac{1}{\alpha_H} - \frac{r_H}{\alpha_H(1-\rho)}$. In simpler terms, we are starting the busy period with $\Theta(1)$ work in the H state. With $\Theta(1)$ probability, the busy period does not end in the H phase, in which case we start the subsequent L state with $\Theta(\alpha^{-1})$ work, with an overall contribution to the mean busy period of $\Theta(\alpha)$. If however, the original busy period ends in the H state itself, then this event contributes a $\Theta(1)$ term and hence is asymptotically negligible compared to the contribution of the event where the busy period does not end in the H state.

Now we consider the case where the busy period begins in an L state. Again, we have two cases – either the busy period ends in the L state itself, or the system switches to an H state before the busy period ends. If the busy period ends in the L state, an event which happens with probability $1 - \Theta(\alpha)$, then the mean busy period conditioned on this event is given by $\frac{\mathbf{E}[\overline{W}]}{r_L}$. However, the system can switch with probability $\Theta(\alpha)$, and the contribution of the residual busy period conditioned on this event can be $\Theta(\alpha^{-1})$ (from the previous subcase). Therefore, this event also contributes a $\Theta(1)$ term to the mean busy period, and we handle this event next.

Consider an $M/G/1$ busy period started by work $\overline{W}$. We let this $M/G/1$ evolve in the L state, and consider an independent Poisson$(\alpha_L)$ marking process. Our aim is to find the workload in the $M/G/1$ when the first mark arrives during the busy period. The probability that no mark arrives is given by $1 - \frac{\mathbf{E}[\overline{W}]\alpha_L}{r_L}$, which we denote by $1 - p_{switch}$ in the theorem statement. Thus, with probability $p_{switch}$, at least one mark arrives, or equivalently, the environment processes switches before the busy period ends and hence the busy period now evolves in the H state.

The subsequent busy period (that which evolves after the system switches to H) is given by the the busy period that starts in H state with work $\widetilde{V}(s)$, where $\widetilde{V}(s)$ denotes the transform of the work that is seen by the Poisson$(\alpha_L)$ marking process *conditioned on being the first mark of a busy period*. We will now argue that this is asymptotically given by the stationary work in an $M/G/1$ conditioned on the server being busy, with exceptional service distribution for the job that starts the busy period given by $\overline{W}$, and service distribution $S_L$. We first note that if we have such an $M/G/1$ where we consider the distribution of work seen by all marks, then this is indeed the stationary work conditioned on the server being busy, and hence is given by the stationary delay seen by arrivals finding the server busy in an $M/G/1$ system with special first service (this expression, $\widetilde{V}(s) = \frac{r_L \cdot \frac{1-\widetilde{\overline{W}}(s)}{\mathbf{E}[\overline{W}]}}{s-\lambda(1-\widetilde{S_L}(s))}$, is given in the theorem statement; see [26] or [15, Appendix B] for proof). However, we are interested in the work that the first mark sees in a

busy period, call this $W_1$. We will argue that as the probability of marking goes to 0, the work seen by the first mark converges in distribution to the stationary work conditioned on the server being busy (and that this sequence of random variables remains uniformly integrable so that the Laplace transforms converge). We first note that the work seen by the first mark is stochastically bounded above by the supremum of the work in a busy period started by work $\overline{W}$, denote this by $W_1^*$. Further, conditioned on a second marked arrival, we can upper bound the work that this mark sees by the supremum of the work in the busy period started by $W_1^*$ (which is an upper bound on the work after the arrival of the first mark), denote this by $W_2^*$. Similarly, we can obtain an upper bound on the work seen by the $n$th marked arrival in a busy period. We also have the trivial lower bound of 0 on the work seen by the $n$th marked arrival in a busy period. Note that both these upper and lower bounds are independent of the marking probability. Let $p_i$ denote the probability that there are $i$ marked arrivals in a busy period. We can thus sandwich the stationary work of the $M/G/1$ conditioned on it being busy between $\frac{p_1 \cdot W_1}{\sum_{i=1}^{\infty} p_i}$ and $\frac{p_1 \cdot W_1 + \sum_{i=2}^{\infty} p_2 \cdot W_i^*}{\sum_{i=1}^{\infty} p_i}$. However, as the marking probability ($\Theta(\alpha)$) goes to 0, $p_i \sim \Theta(\alpha^i)$. Therefore, $W_1$ converges to the stationary work in the $M/G/1$ with special service, conditioned on server being busy. ∎

**Proof of Theorem 7:**

Recall that the work is decreasing during both the $L$ and $H$ states. There is a negative drift of $r_L = 1 - \frac{\lambda}{\mu_L}$ during the $L$ phase and a negative drift of $r_H = 1 - \frac{\lambda}{\mu_H}$ during the $H$ phase.

**Case 1:** $W_\alpha = \omega(\alpha^{-1})$: As in the proof of Theorem 6, since the system switches at a faster time scale ($\Theta(\alpha^{-1})$) than the initial work ($\omega(\alpha^{-1})$), the work during its sojourn sees an average system, and hence the busy period is $\frac{\mathbf{E}[W_\alpha]}{1-\rho} + o(W_\alpha)$.

**Case 2:** $W = \Theta(\alpha^{-1})$: We begin by noting that since the initial work is $\Theta(\alpha^{-1})$, the workload trajectory of the stochastic system, scaled by $\alpha$, converges to the fluid trajectory. Hence the busy period of the stochastic system is given by the fluid busy period and an additional $o(\alpha^{-1})$ term.

We now set up the recurrences for busy periods started by deterministic work $x$ during the H and L phases under the fluid regime:

$$\mathbf{E}[B_H(x)] = \mathbf{E}\left[\min\left\{\frac{x}{r_H}, \tau_H\right\}\right] + \mathbf{E}\left[B_L\left(x - r_H \min\left\{\frac{x}{r_H}, \tau_H\right\}\right) \cdot \mathbf{1}_{\{x > r_H \cdot \tau_H\}}\right]$$

$$\mathbf{E}[B_L(x)] = \mathbf{E}\left[\min\left\{\frac{x}{r_L}, \tau_L\right\}\right] + \mathbf{E}\left[B_H\left(x - r_L \min\left\{\frac{x}{r_L}, \tau_L\right\}\right) \cdot \mathbf{1}_{\{x > r_L \cdot \tau_L\}}\right]$$

where $\tau_H$ is an Exp $(\alpha_H)$ random variable and $\tau_L$ is an Exp $(\alpha_L)$ random variable.

We now guess and verify that $\mathbf{E}[B_H(x)]$ and $\mathbf{E}[B_L(x)]$ have the following function form:

$$B_i(x) = a_i + b_i x + c_i e^{-\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x}$$

where $a_i, b_i$ and $c_i$, $i \in \{L, H\}$, are constants to be determined. The 'guess' is in fact an educated attempt arrived at by exact analysis of the mean busy period started by $n$ jobs in an alternate discrete system which is identical on fluid scale to the system we want to analyze, but with 0 arrival rate. Since $B_i(0) = 0$, we have $a_i = -c_i$. Since the Laplace transform for $x - r_i \min\left\{\frac{x}{r_i}, \tau_i\right\}$

is $\mathbf{E}\left[e^{-s\left(x - \min\left\{\frac{x}{r_i}, \tau_i\right\}\right)}\right] = \frac{se^{-\frac{\alpha_i}{r_i}x} - \frac{\alpha_i}{r_i}e^{-sx}}{s - \frac{\alpha_i}{r_i}}$ and $\mathbf{E}\left[\min\left\{\frac{x}{r_i}, \tau_i\right\}\right] = \frac{1 - e^{-\frac{\alpha_i}{r_i}x}}{\alpha_i}$, our recurrences become:

$$a_L + b_L x + c_L e^{-\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x} = \frac{1 - e^{-\frac{\alpha_L}{r_L}x}}{\alpha_L} + a_H + b_H\left(x - \frac{1 - e^{-\frac{\alpha_L}{r_L}x}}{\frac{\alpha_L}{r_L}}\right) + c_H\left(\frac{\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)e^{-\frac{\alpha_L}{r_L}x} - \frac{\alpha_L}{r_L}e^{-\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x}}{\frac{\alpha_H}{r_H}}\right)$$

$$a_H + b_H x + c_H e^{-\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x} = \frac{1 - e^{-\frac{\alpha_H}{r_H}x}}{\alpha_H} + a_L + b_L\left(x - \frac{1 - e^{-\frac{\alpha_H}{r_H}x}}{\frac{\alpha_H}{r_H}}\right) + c_L\left(\frac{\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)e^{-\frac{\alpha_H}{r_H}x} - \frac{\alpha_H}{r_H}e^{-\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x}}{\frac{\alpha_L}{r_L}}\right).$$

Since the above equalities hold for all $x$, together with $a_i = -c_i$, we get:

$$b_L = b_H = \left(\frac{\frac{r_L}{\alpha_L} + \frac{r_H}{\alpha_H}}{\frac{1}{\alpha_L} + \frac{1}{\alpha_H}}\right)^{-1} = \frac{1}{1 - \rho},$$

$$-a_L = c_L = \frac{r_L - r_H}{\alpha_L \alpha_H \left(\frac{r_L}{\alpha_L} + \frac{r_H}{\alpha_H}\right)^2} \cdot \frac{r_H}{\alpha_H},$$

$$-a_H = c_H = -\frac{r_L - r_H}{\alpha_L \alpha_H \left(\frac{r_L}{\alpha_L} + \frac{r_H}{\alpha_H}\right)^2} \cdot \frac{r_L}{\alpha_L}.$$

Therefore the expected busy period started by a work of size $x$ during L and H phases, respectively, can be expressed in the following convenient/intuitive form:

$$\mathbf{E}[B_L(x)] = \frac{x}{1 - \rho} - \left(\frac{x}{1 - \rho} - \frac{x}{r_L}\right) \cdot \left[\frac{1 - e^{-\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x}}{\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x}\right] \tag{B.1}$$

$$\mathbf{E}[B_H(x)] = \frac{x}{1 - \rho} - \left(\frac{x}{1 - \rho} - \frac{x}{r_H}\right) \cdot \left[\frac{1 - e^{-\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x}}{\left(\frac{\alpha_L}{r_L} + \frac{\alpha_H}{r_H}\right)x}\right] \tag{B.2}$$

which show that $\mathbf{E}[B_L(x)]$ and $\mathbf{E}[B_H(x)]$ are weighted averages of the busy periods of the $\alpha \to 0$ and $\alpha \to \infty$ cases. Taking expectation over $x$ (which is distributed as $W_\alpha$), we obtain the expressions given in the theorem.

**Case 3:** $W_\alpha = o(\alpha^{-1})$: Since the system is stable during both the L and H states, the busy period is $\Theta(W_\alpha)$ (being upper bounded by the busy period started by $W_\alpha$ in an $M/G/1$ with service distribution $S_H$). Suppose the busy period starts in the L state. If the L state were to last forever, the busy period would indeed be $\frac{\mathbf{E}[W_\alpha]}{r_L}$. Now either the system switches to the H state before this busy period ends, and this event happens with probability $1 - o(1)$. In this case, the length of the busy period conditioned on it being smaller than $\text{Exp}(\alpha_L)$ will be $\frac{\mathbf{E}[W_\alpha]}{r_L} + o(W_\alpha)$ since $W_\alpha = o(\alpha^{-1})$. However, if the system switches before the busy period ends, which happens with probability $o(1)$, the residual busy period is still $\Theta(W_\alpha)$. The overall contribution of the second event to the mean busy period started by $W_\alpha$ is $o(W_\alpha)$. By law of total probability, the mean busy started in L phase is $\frac{\mathbf{E}[W_\alpha]}{r_L} + o(W_\alpha)$.

The proof for busy periods started during H phases is identical. ∎

## References

[1] http://www.cs.huji.ac.il/labs/parallel/workload/.

[2] I. J. B. F. Adan and V. G. Kulkarni. Single-server queue with Markov-dependent inter-arrival and service times. *QUESTA*, 45:113–134, 2003.

[3] S. C. Borst, O. J. Boxma, and M. B. Combé. Collection of customers: a correlated *M/G/1* queue. In *SIGMETRICS/Performance*, pages 47–59, New York, NY, USA, 1992. ACM.

[4] G. L. Choudhury, A. Mandelbaum, M. I. Reiman, and W. Whitt. Fluid and diffusion limits for queues in slowly changing environments. *Stoch. Mod.*, 13:121–146, 1997.

[5] I. Cidon, R. Gurin, A. Khamisy, and M. Sidi. Analysis of a correlated queue in a communication system. In *INFOCOM'93*, pages 209–216, 1993.

[6] J. Cohen. *The single server queue*. North Holland, 1969.

[7] R. Conway, W. Maxwell, and M. Miller. *Theory of Scheduling*. Addison-Wesley, 1967.

[8] M. E. Crovella and A. Bestavros. Self-similarity in World Wide Web traffic: Evidence and possible causes. In *ACM SIGMETRICS'96*, pages 160–169, May 1996.

[9] A. Erramilli, O. Narayan, and W. Willinger. Experimental queueing analysis with long-range dependent packet traffic. *IEEE/ACM Trans. on Networking*, 4:209–223, 1996.

[10] D. G. Feitelson. Packing schemes for gang scheduling. In *IPPS*, pages 89–110, London, UK, 1996. Springer-Verlag.

[11] D. G. Feitelson. Locality of sampling and diversity in parallel system workloads. In *Proceedings of the 21st annual International Conference on Supercomputing*, pages 53–63, New York, NY, USA, 2007. ACM.

[12] K. Fendick, V. Saksena, and W. Whitt. Dependence in packet queues. *IEEE Trans. Commun.*, 37:1173–1183, 1989.

[13] H. J. Fowler, W. E. Leland, and B. Bellcore. Local area network traffic characteristics, with implications for broadband network congestion management. *IEEE Journal on Selected Areas in Communications*, 9:1139–1149, 1991.

[14] S. Ghosh and M. Squillante. Analysis and control of correlated web server queues. *Computer Communications*, 27(18):1771–1785, 2004.

[15] V. Gupta, M. Burroughs, and M. Harchol-Balter. Analysis of scheduling policies under correlated job sizes. Technical Report CMU-CS-10-107, School of Computer Science, Carnegie Mellon University, 2010.

[16] H. Li, D. Groep, and L. Wolters. Workload characteristics of a multi-cluster supercomputer. pages 176–193. Springer Verlag, 2004.

[17] M. Livny, B. Melamed, and A. K. Tsiolis. The impact of autocorrelation on queuing systems. *Manage. Sci.*, 39(3):322–339, 1993.

[18] N. Mi, G. Casale, and E. Smirni. Scheduling for performance and availability in systems with temporal dependent workloads. In *DSN'08*, pages 336–345, 2008.

[19] N. Mi, G. Casale, Q. Zhang, A. Riska, and E. Smirni. Autocorrelation-driven load control in distributed systems. In *MASCOTS'09*, 2009.

[20] V. Paxson and S. Floyd. Wide area traffic: the failure of Poisson modeling. *IEEE/ACM Transactions on Networking*, 3(3):226–244, 1995.

[21] R. Righter, J. G. Shanthikumar, and G. Yamazaki. On extremal service disciplines in single-stage queueing systems. *J. Appl. Probab.*, 27(2):409–416, 1990.

[22] A. Riska, M. Squillante, S.-Z. Yu, Z. Liu, and L. Zhang. Matrix-analytic analysis of a $MAP/PH/1$ queue fitted to web server data. *Matrix-Analytic Methods: Theory and Applications*, pages 335–356, 2002.

[23] E. Smirni, Q. Zhang, N. Mi, A. Riska, and G. Casale. New results on the performance effects of autocorrelated flows in systems. In *IEEE IPDPS'07*, pages 1–6, 2007.

[24] B. Song, C. Ernemann, and R. Yahyapour. Parallel computer workload modeling with markov chains. In *Proc. of the 10th Job Scheduling Strategies for Parallel Processing (JSSPP)*, pages 47–62. Springer, 2004.

[25] M. S. Squillante, D. D. Yao, and L. Zhang. Internet traffic: periodicity, tail behavior, and performance implications. *System performance evaluation: methodologies and applications*, pages 23–37, 2000.

[26] H. Takagi. *Queueing Analysis, Vol. 1: Vacation and Priority Systems*. North-Holland, 1991.

[27] E. van Doorn and J. Regterschot. Conditional PASTA. *Oper. Res. Lett.*, 7:229–232, 1988.

[28] Q. Zhang, N. Mi, A. Riska, and E. Smirni. Load unbalancing to improve performance under autocorrelated traffic. In *ICDCS'06*, Lisboa, Portugal, 2006.