

**Original citation:**

Rizk, Amr, Poloczek, Felix and Ciucu, Florin. (2016) Stochastic bounds in fork-join queueing systems under full and partial mapping. Queueing Systems, 83 (3). pp. 261-291.

**Permanent WRAP URL:**

<http://wrap.warwick.ac.uk/79510>

**Copyright and reuse:**

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions. Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

**Publisher's statement:**

"The final publication is available at Springer via <http://dx.doi.org/10.1007/s11134-016-9486-x>

**A note on versions:**

The version presented here may differ from the published version or, version of record, if you wish to cite this item you are advised to consult the publisher's version. Please see the 'permanent WRAP URL' above for details on accessing the published version and note that access may require a subscription.

For more information, please contact the WRAP Team at: [wrap@warwick.ac.uk](mailto:wrap@warwick.ac.uk)

# Stochastic Bounds in Fork-Join Queueing Systems under Full and Partial Mapping

Amr Rizk · Felix Poloczek · Florin Ciucu

the date of receipt and acceptance should be inserted later

**Abstract** In a Fork-Join (FJ) queueing system an upstream fork station splits incoming jobs into  $N$  tasks to be further processed by  $N$  parallel servers, each with its own queue; the response time of one job is determined, at a downstream join station, by the maximum of the corresponding tasks' response times. This queueing system is useful to the modelling of multi-service systems subject to synchronization constraints, such as MapReduce clusters or multipath routing. Despite their apparent simplicity, FJ systems are hard to analyze.

This paper provides the first *computable* stochastic bounds on the waiting and response time distributions in FJ systems under full (bijective) and partial (injective) mapping of tasks to servers. We consider four practical scenarios by combining 1a) renewal and 1b) non-renewal arrivals, and 2a) non-blocking and 2b) blocking servers. In the case of non-blocking servers we prove that delays scale as  $\mathcal{O}(\log N)$ , a law which is known for first moments under renewal input only. In the case of blocking servers, we prove that the same factor of  $\log N$  dictates the stability region of the system. Simulation results indicate that our bounds are tight, especially at high utilizations, in all four scenarios. A remarkable insight gained from our results is that, at moderate to high utilizations, multipath routing “*makes sense*” from a queueing perspective for two paths only, i.e., response times drop the most when  $N = 2$ ; the technical

---

Amr Rizk  
University of Massachusetts Amherst, USA  
E-mail: arizk@umass.edu

Felix Poloczek  
University of Warwick, UK / TU Berlin, Germany  
E-mail: felix@inet.tu-berlin.de

Florin Ciucu  
University of Warwick, UK  
E-mail: F.Ciucu@warwick.ac.uk

explanation is that the resequencing (delay) price starts to quickly dominate the tempting gain due to multipath transmissions.

**Keywords** Fork-Join queue · Performance evaluation · Parallel systems · MapReduce · Multipath

## 1 Introduction

The performance analysis of Fork-Join (FJ) systems received new momentum with the recent wide-scale deployment of large-scale data processing that was enabled through emerging frameworks such as MapReduce [13]. The main idea behind these big data analysis frameworks is an elegant divide and conquer strategy with various degrees of freedom in the implementation. The open-source implementation of MapReduce, known as Hadoop [42], is deployed in numerous production clusters, e.g., Facebook and Yahoo [24].

The basic operation of MapReduce is depicted in Figure 1. In the *map phase*, a job is split into multiple tasks that are mapped to different workers (servers). Once a specific subset of these tasks finish their executions, the corresponding *reduce phase* starts by processing the combined output from all the corresponding tasks. In other words, the reduce phase is subject to a fundamental synchronization constraint on the finishing times of all involved tasks.

A natural way to model one reduce phase operation is by a *basic* FJ queueing system with  $N$  servers. Jobs, i.e., the input unit of work in MapReduce systems, arrive according to some point process. Each job is split into  $N$  (map) tasks (or *splits*, in the MapReduce terminology), which are simultaneously sent to the  $N$  servers. At each server, each task requires a random service time, capturing the variable task execution times on different servers in the map phase. A job leaves the FJ system when all of its tasks are served; this constraint corresponds to the specification that the reduce phase starts no sooner than when all of its map tasks complete their executions.

Concerning the execution of tasks belonging to different jobs on the same server, there are two operational modes. In the *non-blocking* mode, the servers are workconserving in the sense that tasks immediately start their executions once the previous tasks finish theirs. In the *blocking* mode, the mapped tasks of a job simultaneously start their executions, i.e., servers can be idle when their corresponding queues are not empty. The non-blocking execution mode prevails in MapReduce due to its conceivable efficiency, whereas the blocking execution mode is employed when the *jobtracker* (the node coordinating and scheduling jobs) waits for all machines to be ready to synchronize the configuration files before mapping a new job; in Hadoop, this can be enforced through the coordination service *zookeeper* [42].

In this paper we analyze the performance of the FJ queueing model in four practical scenarios by considering two broad arrival classes (driven by either renewal or non-renewal processes) and the two operational modes described above. The key contribution, to the best of our knowledge, are the

first non-asymptotic and computable stochastic bounds on the waiting and response time distributions in the most relevant scenario, i.e., non-renewal (Markov modulated) job arrivals and the non-blocking operational mode. Under all scenarios, the bounds are numerically tight especially at high utilizations. This inherent tightness is due to a suitable martingale representation of the underlying queueing system, an approach which was conceived in [27] for the analysis of GI/GI/1 queues, and which was recently extended to address multi-class queues with non-renewal arrivals [12, 34]. The simplicity of the obtained stochastic bounds enables the derivation of scaling laws, e.g., delays in FJ systems scale as  $\mathcal{O}(\log N)$  in the number of parallel servers  $N$ , for both renewal and non-renewal arrivals, in the non-blocking mode; more severe delay degradations hold in the blocking mode, and, moreover, the stability region depends on the same fundamental factor of  $\log N$ .

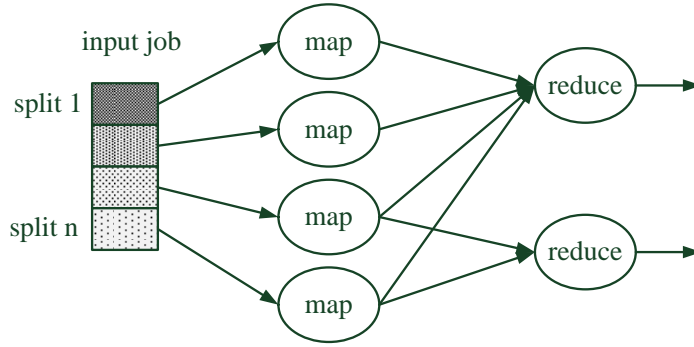
In addition to the direct applicability to the dimensioning of MapReduce clusters, there are other relevant types of parallel and distributed systems such as production or supply networks. In particular, by slightly modifying the basic FJ system corresponding to MapReduce, the resulting model suits the analysis of window-based transmission protocols over multipath routing. By making several simplifying assumptions such as ignoring the details of specific protocols (e.g., multipath TCP), we can provide a fundamental understanding of multipath routing from a queueing perspective. Concretely, we demonstrate that sending a flow of packets over two paths, instead of one, does generally reduce the steady-state response times. The surprising result is that by sending the flow over more than two paths, the steady-state response times start to increase. The technical explanation for such a rather counterintuitive result is that the  $\log N$  resequencing price at the destination quickly dominates the tempting gain in the queueing waiting time due to multipath transmissions.

The rest of the paper is structured as follows. We first discuss related work on FJ systems and related applications. Then we analyze full mapping, i.e., a mapping of jobs to  $N$  servers in Sections 3 and 4. We analyze both non-blocking and blocking FJ systems with renewal input in Section 3, and with non-renewal input in Section 4. The analysis of partial mapping, i.e., a mapping of jobs to  $H \leq N$  servers follows in Section 5. In Section 6 we apply the obtained results on the steady-state response time distributions to the analysis of multipath routing from a queueing perspective. Brief conclusions are presented in Section 7.

## 2 Related Work

We first review analytical results on FJ systems, and then results related to the two application case studies considered in this paper, i.e., MapReduce and multipath routing.

The significance of the Fork-Join queueing model stems from its natural ability to capture the behavior of many parallel service systems. The performance of FJ queueing systems has been subject of multiple studies such as



**Fig. 1** Schematic illustration of the basic operation of MapReduce.

[5, 31, 40, 25, 28, 6, 8]. In particular, [5] notes that an exact performance evaluation of general FJ systems is remarkably hard due to the synchronization constraints on the input and output streams. More precisely, a major difficulty lies in finding an exact closed form expression for the joint steady-state workload distribution for the FJ queueing system. However, a number of results exist given certain constraints on the FJ system. The authors of [15] provide the stationary joint workload distribution for a two-server FJ system under Poisson arrivals and independent exponential service times. For the general case of more than two parallel servers there exists a number of works that provide approximations [31, 40, 28, 29] and bounds [5, 6] for certain performance metrics of the FJ system. Given renewal arrivals, [6] significantly improves the lower bounds from [5] in the case of heterogeneous phase-type servers using a matrix-geometric algorithmic method. The authors of [28] provide an approximation of the sojourn time distribution in a renewal driven FJ system consisting of multiple G/M/1 nodes. They show that the approximation error diminishes at extremal utilizations. Refined approximations for the mean sojourn time in two-server FJ systems that take the first two moments of the service time distribution are given in [25]; numerical evidence is further provided on the quality of the approximation for different service time distributions. In a recent work, the authors of [30] establish Gaussian limits for the joint distributions of the service and waiting times for synchronization under general arrivals characterized by a limiting Brownian motion.

The closest related work to ours is [5], which provides computable lower and upper bounds on the expected response time in FJ systems under renewal assumptions with Poisson arrivals and exponential service times; the underlying idea is to artificially construct a more tractable system, yet subject to stochastic ordering relative to the original one. Our corresponding first order upper bound recovers the  $\mathcal{O}(\log N)$  asymptotic behavior of the one from [5], and also reported in [31] in the context of an approximation; numerically, our bound is slightly worse than the one from [5] due to our main focus on computing bounds on the whole distribution (first order bounds are secondarily obtained by integration). Moreover, we show that the  $\mathcal{O}(\log N)$  scaling law

also holds in the case of Markov modulated arrivals. In a parallel work [26] to ours, the authors adopt a network calculus approach to derive stochastic bounds in a non-blocking FJ system, under a strong assumption on the input; for related constructions of such arrival models see [20].

The work in [21,22] studies FJ systems where jobs leave the system when a subset  $H \leq N$  of its tasks are finished. This system is similar to the partial mapping FJ system that we study in Section 5, however, with subtle yet fundamental differences. The FJ system presented in [21,22] is based on the assumption that when  $H$  tasks finish execution, the finished job *purges* the unfinished  $N - H$  tasks out their corresponding queues. The authors of [21, 22] provide upper bounds for the mean response times in such systems under Poisson arrivals and general service distributions. In Section 5, we consider instead injective task mapping, i.e., jobs are *only* forked onto a subset of servers  $H \leq N$ . For this type of FJ systems we provide bounds on the steady state waiting and response time distributions under round-robin and random task placement.

Concerning concrete applications of FJ systems, in particular MapReduce, there are several empirical and analytical studies analyzing its performance. For instance, [44,3] aim to improve the system performance via empirically adjusting its numerous and highly complex parameters. The targeted performance metric in these studies is the job response time, which is in fact an integral part of the business model of MapReduce based query systems such as [32] and time priced computing clouds such as Amazon's EC2 [1]. For an overview on works that optimize the performance of MapReduce systems see the survey article [33]. Using a similar idea as in [5], the authors of [37] derive asymptotic results on the response time distribution in the case of renewal arrivals; such results are further used to understand the impact of different scheduling models in the reduce phase of MapReduce. Using the model from [37] the work in [38] provides approximations for the number of jobs in a tandem system consisting of a map queue and a reduce queue in the heavy traffic regime. The work in [41] derives approximations of the mean response time in MapReduce systems using a mean value analysis technique and a closed FJ queueing system model from [39].

Concerning multipath routing, the works [4,19] provided ground for multiple studies on different formulations of the underlying resequencing delay problem, e.g., [18,43]. Factorization methods were used in [4] to analyze the disordering delay and the delay of resequencing algorithms, while the authors of [19] conduct a queueing theoretic analysis of an  $M/G/\infty$  queue receiving a stream of numbered customers. In [18,43] the multipath routing model comprises Bernoulli thinning of Poisson arrivals over  $N$  parallel queueing stations followed by a resequencing buffer. The work in [18] provides asymptotics on the conditional probability of the resequencing delay conditioned on the end-to-end delay for different service time distributions. For  $N = 2$  and exponential interarrival and service times, [43] derives a large deviations result on the resequencing queue size. Our work differs from these works in that we consider a model of the basic operation of window-based transmission protocols over

multipath routing, motivated by the emerging application of multipath TCP [35]. We point out, however, that we do not model the specific operation of any particular multipath transmission protocol. Instead, we analyze a generic multipath transmission protocol under simplifying assumptions, in order to provide a theoretical understanding of the overall response times comprised of both queueing and resequencing delays.

Relative to the existing literature, our key theoretical contribution is to provide *computable* and non-asymptotic bounds on the *distributions* of the steady-state waiting and response times under both *renewal* and *non-renewal* input in non-blocking FJ systems. These bounds can be found in Theorem 1, Theorem 3, and Theorem 5 – Theorem 7. The consideration of non-renewal input is particularly relevant, given recent observations that job arrivals are subject to temporal correlations in production clusters. For instance, [11, 23] report that job, respectively, flow arrival traces in clusters running MapReduce exhibit various degrees of burstiness. We augment the scope of the main contributions in this work by considering *blocking* FJ systems that essentially correspond to GI/G/1 queueing systems. Here, we recover and extend prominent results, e.g., from [2, 16] in Theorem 2 and Theorem 4, respectively. Note that non-blocking FJ systems behave fundamentally different from blocking FJ systems, thus requiring adapted mathematical tools for the analysis.

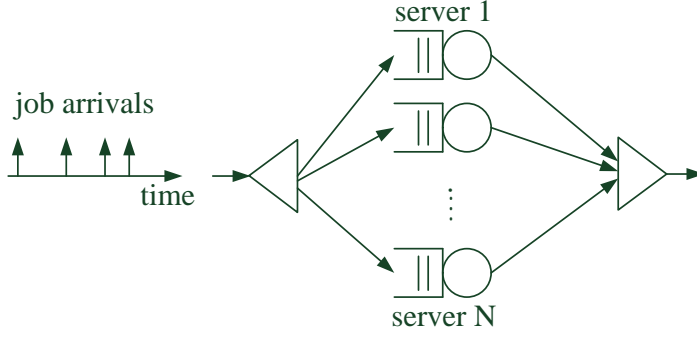
### 3 FJ Systems with Renewal Input

We consider a FJ queueing system as depicted in Figure 2. Jobs arrive at the input queue of the FJ system according to some point process with interarrival times  $t_i$  between the  $i$  and  $i + 1$  jobs. Each job  $i$  is split into  $N$  tasks that are mapped through a bijection to  $N$  servers. A task of job  $i$  that is serviced by some server  $n$  requires a random service time  $x_{n,i}$ . A job leaves the system when all of its tasks finish their executions, i.e., there is an underlying synchronization constraint on the output of the system. We assume that the families  $\{t_i\}$  and  $\{x_{n,i}\}$  are independent.

In the sequel we differentiate between two cases, i.e., *a*) non-blocking and *b*) blocking servers. The first case corresponds to workconserving servers, i.e., a server starts servicing a task of the next job (if available) immediately upon finishing the current task. In the latter case, a server that finishes servicing a task is blocked until the corresponding job leaves the system, i.e., until the last task of the current job completes its execution. This can be regarded as an additional synchronization constraint on the input of the system, i.e., all tasks of a job start receiving service simultaneously. We will next analyze *a*) and *b*) for renewal arrivals.

#### 3.1 Non-Blocking Systems

Consider an arrival flow of jobs with renewal interarrival times  $t_i$ , and assume that the waiting time of the first job is  $w_1 = 0$ . Given  $N$  parallel servers, the



**Fig. 2** A schematic Fork-Join queueing system with  $N$  parallel servers. An arriving job is split into  $N$  tasks, one for each server. A job leaves the FJ system when all of its tasks are served. An arriving job is considered waiting until the service of the last of its tasks starts, i.e., when the previous job departs the system.

waiting time  $w_j$  of the  $j$ th job is defined as

$$w_j = \max \left\{ 0, \max_{1 \leq k \leq j-1} \left\{ \max_{n \in [1, N]} \left\{ \sum_{i=1}^k x_{n, j-i} - \sum_{i=1}^k t_{j-i} \right\} \right\} \right\}, \quad (1)$$

for all  $j \geq 2$ , where  $x_{n,j}$  is the service time required by the task of job  $j$  that is mapped to server  $n$ . We count a job as waiting until its last task starts receiving service. Similarly, the *response times* of jobs, i.e., the times until the last corresponding tasks have finished their executions, are defined as  $r_1 = \max_n x_{n,1}$  for the first job, and for  $j \geq 2$  as

$$r_j = \max_{0 \leq k \leq j-1} \left\{ \max_{n \in [1, N]} \left\{ \sum_{i=0}^k x_{n, j-i} - \sum_{i=1}^k t_{j-i} \right\} \right\}, \quad (2)$$

where by convention  $\sum_{i=1}^0 t_i = 0$ ; for brevity, we will denote  $\max_n := \max_{n \in [1, N]}$ .

We assume that the task service times  $x_{n,j}$  are independent and identically distributed (iid). The stability condition for the FJ queueing system is given as  $E[x_{1,1}] < E[t_1]$ . By stationarity and reversibility of the iid processes  $x_{n,j}$  and  $t_j$ , there exists a distribution of the steady-state waiting time  $w$  and steady-state response time  $r$ , respectively, which have the representations

$$w =_{\mathcal{D}} \max_{k \geq 0} \left\{ \max_n \left\{ \sum_{i=1}^k x_{n,i} - \sum_{i=1}^k t_i \right\} \right\} \quad (3)$$

and

$$r =_{\mathcal{D}} \max_{k \geq 0} \left\{ \max_n \left\{ \sum_{i=0}^k x_{n,i} - \sum_{i=1}^k t_i \right\} \right\}, \quad (4)$$

respectively. Here,  $=_{\mathcal{D}}$  denotes equality in distribution. Note that the only difference in (3) and (4) is that for the latter the sum over the  $x_{n,i}$  starts at  $i = 0$  rather than at  $i = 1$ .



The following theorem provides stochastic upper bounds on  $w$  and  $r$ . The corresponding proof will rely on submartingale constructions and the Optional Sampling Theorem (see Lemma 1 in the Appendix).

**Theorem 1** (RENEWALS, NON-BLOCKING) *Given a FJ system with  $N$  parallel non-blocking servers that is fed by renewal job arrivals with interarrivals  $t_j$ . If the task service times  $x_{n,j}$  are iid, then the steady-state waiting and response times  $w$  and  $r$  are bounded by*

$$\mathbf{P}[w \geq \sigma] \leq N e^{-\theta_{nb}\sigma} \quad (5)$$

$$\mathbf{P}[r \geq \sigma] \leq N \mathbf{E}[e^{\theta_{nb}x_{1,1}}] e^{-\theta_{nb}\sigma}, \quad (6)$$

where  $\theta_{nb}$  (with the subscript ‘nb’ standing for non-blocking) is the (positive) solution of

$$\mathbf{E}[e^{\theta_{nb}x_{1,1}}] \mathbf{E}[e^{-\theta_{nb}t_1}] = 1. \quad (7)$$

We remark that the stability condition  $\mathbf{E}[x_{1,1}] < \mathbf{E}[t_1]$  guarantees the existence of a positive solution in (7) (see also [34]).

*Proof* Consider the waiting time  $w$ . We first prove that for each  $n \in [1, N]$  the process

$$z_n(k) = e^{\theta_{nb} \sum_{i=1}^k (x_{n,i} - t_i)}$$

is a martingale with respect to the filtration

$$\mathcal{F}_k := \sigma\{x_{n,m}, t_m \mid m \leq k, n \in [1, N]\}.$$

The independence assumption of  $x_{n,j}$  and  $t_j$  implies that

$$\begin{aligned} \mathbf{E}[z_n(k) \mid \mathcal{F}_{k-1}] &= \mathbf{E}\left[e^{\theta_{nb} \sum_{i=1}^k (x_{n,i} - t_i)} \mid \mathcal{F}_{k-1}\right] \\ &= \mathbf{E}\left[e^{\theta_{nb}(x_{n,k} - t_k)}\right] e^{\theta_{nb} \sum_{i=1}^{k-1} (x_{n,i} - t_i)} \\ &= e^{\theta_{nb} \sum_{i=1}^{k-1} (x_{n,i} - t_i)} \\ &= z_n(k-1), \end{aligned} \quad (8)$$

under the condition on  $\theta_{nb}$  from the theorem. Moreover,  $z_n(k)$  is obviously integrable by the condition on  $\theta_{nb}$  from the theorem, completing thus the proof for the martingale property.

Next we prove that the process

$$z(k) = \max_n z_n(k) \quad (9)$$

is a submartingale w.r.t.  $\mathcal{F}_k$ . Given the martingale property of each of the  $z_n$  and the monotonicity of the conditional expectation we can write for  $j \in [1, N]$ :

$$\mathbf{E}\left[\max_n z_n(k) \mid \mathcal{F}_{k-1}\right] \geq \mathbf{E}[z_j(k) \mid \mathcal{F}_{k-1}] = z_j(k-1),$$

where the inequality stems from  $\max_n z_n(k) \geq z_j(k)$  for  $j \in [1, N]$  a.s., whereas the subsequent equality stems from the martingale property (8) for  $z_n(k)$  for all  $n \in [1, N]$ . Hence, we can write

$$\mathbb{E}[z(k) \mid \mathcal{F}_{k-1}] \geq \max_n z_n(k-1) = z(k-1), \quad (10)$$

which proves the submartingale property.

To derive a bound on the steady-state waiting time distribution, let  $\sigma > 0$  and define the stopping time

$$K := \inf \left\{ k \geq 0 \mid \max_n \sum_{i=1}^k (x_{n,i} - t_i) \geq \sigma \right\}, \quad (11)$$

which is also the first point in time  $k$  where  $z(k) \geq e^{\theta_{nb}\sigma}$ . Note that with the representation of  $w$  from (3):

$$\{K < \infty\} = \{w \geq \sigma\}.$$

Now, using the Optional Sampling Theorem (see Lemma 1 from the Appendix) for submartingales with  $k \geq 1$ :

$$\begin{aligned} N &= \sum_{n \in [1, N]} \mathbb{E} \left[ e^{\theta_{nb} \sum_{i=1}^k (x_{n,i} - t_i)} \right] \\ &\geq \mathbb{E} \left[ \max_n e^{\theta_{nb} \sum_{i=1}^k (x_{n,i} - t_i)} \right] \\ &= \mathbb{E}[z(k)] \geq \mathbb{E}[z(K \wedge k)] \geq \mathbb{E}[z(K) 1_{K < k}] \\ &\geq e^{\theta_{nb}\sigma} \mathbb{P}[K < k], \end{aligned} \quad (12)$$

where we used the condition on  $\theta_{nb}$  from the theorem in the first line, the union bound in the second line, and the submartingale property in the third line. In the last line we used the definition of the stopping time  $K$ ; note that we use the notation  $K \wedge n := \min\{K, n\}$ . The proof completes by letting  $k \rightarrow \infty$ .

For the response time  $r$ , define the processes

$$\tilde{z}_n(k) = e^{\theta_{nb}(\sum_{i=0}^k x_{n,i} - \sum_{i=1}^k t_i)},$$

which differs from the  $z_n$  only in the range of the sum of the service times  $x_{n,i}$ . Then we proceed as for the derivation of the bound on the waiting time  $w$ . The only difference in the derivation is that inequality (12) translates to

$$N \mathbb{E} \left[ e^{\theta_{nb} x_{1,1}} \right] \geq \mathbb{E} \left[ \max_n e^{\theta_{nb}(\sum_{i=0}^k x_{n,i} - \sum_{i=1}^k t_i)} \right].$$

Fixing the right hand sides in (5) and (6) to  $\varepsilon$ , we find that the corresponding quantiles on the waiting and response times grow with the number of parallel servers  $N$  as  $\mathcal{O}(\log N)$ , a law which was already demonstrated in the special case of Poisson arrival and exponential service times, and for first moments, in [31], and more generally in [5]. This scaling result is essential for dimensioning FJ systems such as MapReduce computing clusters, as it explains

the impact of a MapReduce server pool size  $N$  on the job waiting/response times.

We note that the bound in Theorem 1 can be computed for different arrival and service time distributions as long as the MGF (moment generating function) and Laplace transform from (7) are computable. Given a scenario where the job interarrival process and the task size distributions in a MapReduce cluster are not known a priori, estimates of the corresponding MGF and Laplace transforms can be obtained using recorded traces, e.g., using the method from [17].

Next we illustrate two immediate applications of Theorem 1.

*Example 1: Exponentially distributed interarrival and service times*

Consider that the interarrival times  $t_i$  and service times  $x_{n,i}$  are exponentially distributed with parameters  $\lambda$  and  $\mu$ , respectively; note that when  $N = 1$  the system corresponds to the M/M/1 queue. The corresponding stability condition becomes  $\mu > \lambda$ . Using Theorem 1, the bounds on the steady-state waiting and response time distributions are

$$\mathbb{P}[w \geq \sigma] \leq N e^{-(\mu-\lambda)\sigma} \quad (13)$$

and

$$\mathbb{P}[r \geq \sigma] \leq \frac{N}{\rho} e^{-(\mu-\lambda)\sigma}, \quad (14)$$

where the exponential decay rate  $\mu - \lambda$  follows by solving  $\frac{\mu}{\mu-\theta} \frac{\lambda}{\lambda+\theta} = 1$ , i.e., the instantiation of (7). Here, we use  $\rho$  to denote the utilization  $\lambda/\mu$ .

Next we briefly compare our results to the existing bound on the mean response time from [5], given as

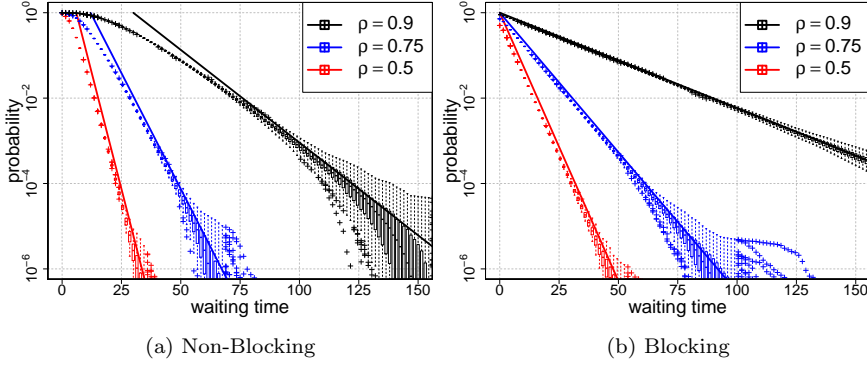
$$\mathbb{E}[r] \leq \frac{1}{\mu - \lambda} \sum_{n=1}^N \frac{1}{n}. \quad (15)$$

By integrating the tail of (14) we obtain the following upper bound on the mean response time

$$\mathbb{E}[r] \leq \frac{\log(N/\rho) + 1}{\mu - \lambda}.$$

Compared to (15), our bound exhibits the same  $\log N$  scaling law but is numerically slightly looser; asymptotically in  $N$ , the ratio between the two bounds converges to one. A key technical reason for obtaining a looser bound is that we mainly focus on deriving bounds on distributions; through integration, the numerical discrepancies accumulate.

For the numerical illustration of the tightness of the bounds on the waiting time distributions from (13) we refer to Figure 3.(a); the numerical parameters and simulation details are included in the caption.



**Fig. 3** Bounds on the waiting time distributions vs. simulations (renewal input): (a) the non-blocking case (13) and (b) the blocking case (22). The system parameters are  $N = 20$ ,  $\mu = 1$ , and three utilization levels  $\rho = \{0.9, 0.75, 0.5\}$  (from top to bottom). Simulations include 100 runs, each accounting for  $10^7$  slots.

*Example 2: Exponentially distributed interarrival times and constant service times*

We now consider the case of iid exponentially distributed interarrival times  $t_i$  with parameter  $\lambda$ , and deterministic service times  $x_{n,i} = 1/\mu$ , for all  $i \geq 0$  and  $n \in [1, N]$ ; note that when  $N = 1$  the system corresponds to the M/D/1 queue.

The condition on the asymptotic decay rate  $\theta_{nb}$  from Theorem 1 becomes

$$\frac{\lambda}{\lambda + \theta_{nb}} = e^{-\frac{\theta_{nb}}{\mu}},$$

which can be numerically solved; upper bounds on the waiting and response time distributions follow then immediately from Theorem 1.

### 3.2 Blocking Systems

Here, we consider a blocking FJ queueing system, i.e., the start of each job is synchronized amongst all servers. We maintain the iid assumptions on the interarrival times  $t_i$  and service times  $x_{n,i}$ . The waiting time and response time for the  $j$ th job can then be written as

$$w_j = \max \left\{ 0, \max_{1 \leq k \leq j-1} \left\{ \sum_{i=1}^k \max_n x_{n,j-i} - \sum_{i=1}^k t_{j-i} \right\} \right\}$$

$$r_j = \max_{0 \leq k \leq j-1} \left\{ \sum_{i=0}^k \max_n x_{n,j-i} - \sum_{i=1}^k t_{j-i} \right\}.$$

Note that the only difference to (1) and (2) is that the maximum over the number of servers now occurs inside the sum. Note that this blocking system corresponds to a GI/GI/1 queue which is analyzed, e.g., in [2].

It is evident that the blocking system is more conservative than the non-blocking system in the sense that the waiting time distribution of the non-blocking system is dominated by the waiting time distribution of the blocking system. Moreover, the stability region for the blocking system, given by  $\mathbb{E}[t_1] > \mathbb{E}[\max_n x_{n,1}]$ , is included in the stability region of the corresponding non-blocking system (i.e.,  $\mathbb{E}[t_1] > \mathbb{E}[x_{1,1}]$ ).

Analogously to (3), the steady-state waiting and response times  $w$  and  $r$  have now the representations

$$w =_{\mathcal{D}} \max_{k \geq 0} \left\{ \sum_{i=1}^k \max_n x_{n,i} - \sum_{i=1}^k t_i \right\} \quad (16)$$

$$r =_{\mathcal{D}} \max_{k \geq 0} \left\{ \sum_{i=0}^k \max_n x_{n,i} - \sum_{i=1}^k t_i \right\} . \quad (17)$$

The following theorem provides upper bounds on  $w$  and  $r$ .

**Theorem 2** (RENEWALS, BLOCKING) *Given a FJ queueing system with  $N$  parallel blocking servers that is fed by renewal job arrivals with interarrivals  $t_j$  and iid task service times  $x_{n,j}$ . The distributions of the steady-state waiting and response times are bounded by*

$$\begin{aligned} \mathbb{P}[w \geq \sigma] &\leq e^{-\theta_b \sigma} \\ \mathbb{P}[r \geq \sigma] &\leq \mathbb{E}[e^{\theta_b \max_n x_{n,1}}] e^{-\theta_b \sigma} , \end{aligned} \quad (18)$$

where  $\theta_b$  (with the subscript ‘b’ standing for blocking) is the (positive) solution of

$$\mathbb{E}[e^{\theta \max_n x_{n,1}}] \mathbb{E}[e^{-\theta t_1}] = 1 . \quad (19)$$

Before giving the proof we note that, in general, (19) can be numerically solved. Moreover, for small values of  $N$ ,  $\theta_b$  can be analytically solved.

*Proof* Consider the waiting time  $w$ . We proceed similarly as in the proof of Theorem 1. Letting  $\mathcal{F}_k$  as above, we first prove that the process

$$y(k) = e^{\theta_b \sum_{i=1}^k (\max_n x_{n,i} - t_i)}$$

is a martingale w.r.t.  $\mathcal{F}_k$  using a technique from [27]. We write

$$\begin{aligned} \mathbb{E}[y(k) | \mathcal{F}_{k-1}] &= \mathbb{E}\left[e^{\theta_b \sum_{i=1}^k (\max_n x_{n,i} - t_i)} \middle| \mathcal{F}_{k-1}\right] \\ &= e^{\theta_b \sum_{i=1}^{k-1} (\max_n x_{n,i} - t_i)} \mathbb{E}\left[e^{\theta_b (\max_n x_{n,k} - t_k)}\right] \\ &= e^{\theta_b \sum_{i=1}^{k-1} (\max_n x_{n,i} - t_i)} \\ &= y(k-1) , \end{aligned}$$

where we used the independence and renewal assumptions for  $x_{n,i}$  and  $t_i$  in the second line, and finally the condition on  $\theta_b$  from (19).

In the next step we apply the Optional Sampling Theorem (45) to derive the bound from the theorem. We first define the stopping time  $K$  by

$$K := \inf \left\{ k \geq 0 \left| \sum_{i=1}^k \left( \max_n x_{n,i} - t_i \right) \geq \sigma \right. \right\} . \quad (20)$$

Recall that  $\mathbf{P}[K < \infty] = \mathbf{P}[w \geq \sigma]$ . We can next write for every  $k \in \mathbb{N}$

$$\begin{aligned} 1 &= \mathbf{E}[y(0)] \\ &= \mathbf{E}[y(K \wedge k)] \\ &\geq \mathbf{E}[y(K \wedge k) 1_{K < k}] \\ &= \mathbf{E} \left[ e^{\theta_b \sum_{i=1}^K (\max_n x_{n,i} - t_i)} 1_{K < k} \right] \\ &\geq e^{\theta_b \sigma} \mathbf{P}[K < k] . \end{aligned}$$

Taking  $k \rightarrow \infty$  completes the proof. The proof for the response time  $r$  is analogous.

*Example 3: Exponentially distributed interarrival and service times*

Consider interarrival and service times  $t_i$  and  $x_{n,i}$  that are exponentially distributed with parameters  $\lambda$  and  $\mu$ , respectively. In [36] it was shown that

$$\max_n L_n =_{\mathcal{D}} \sum_{n=1}^N \frac{L_n}{n}$$

for iid exponentially distributed random variables  $L_n$ , so that the stability condition  $\mathbf{E}[t_1] > \mathbf{E}[\max_n x_{n,1}]$  becomes

$$\frac{1}{\lambda} > \frac{1}{\mu} \sum_{n=1}^N \frac{1}{n} . \quad (21)$$

By applying Theorem 2, the bounds on the steady-state waiting and response time distributions are

$$\mathbf{P}[w \geq \sigma] \leq e^{-\theta_b \sigma} \quad (22)$$

and

$$\mathbf{P}[r \geq \sigma] \leq \frac{\mu}{\mu - \theta_b} e^{-\theta_b \sigma} ,$$

where  $\theta_b$  can be numerically solved from the condition

$$\prod_{n=1}^N \frac{n\mu}{n\mu - \theta_b} \frac{\lambda}{\lambda + \theta_b} = 1 .$$

For quick numerical illustrations we refer back to Figure 3.(b).

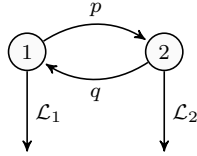
The interesting observation is that the stability condition from (21) depends on the number of servers  $N$ . In particular, as the right hand side grows in  $\log N$ , the system becomes unstable (i.e., waiting times are infinite) for sufficiently large  $N$ . This shows that the optional blocking mode from Hadoop should be judiciously enabled.

*Example 4: Exponentially distributed interarrival and constant service times*

If the service times are deterministic, i.e.,  $x_{n,i} = 1/\mu$  for all  $i \geq 0$  and  $n \in [1, N]$ , the representations of  $w$  and  $r$  from (16) and (17) match their non-blocking counterparts from (3) and (4) and hence the corresponding stability regions and stochastic bounds are equal to those from Example 2.

#### 4 FJ Systems with Non-renewal Input

In this section we consider the more realistic case of FJ queueing systems with non-renewal job arrivals. This model is particularly relevant given the empirical evidence that clusters running MapReduce exhibit various degrees of burstiness in the input [11, 23]. Moreover, numerous studies have demonstrated the burstiness of Internet traces, which can be regarded in particular as the input to multipath routing.



**Fig. 4** Markov modulating chain  $c_k$  for the job interarrival times.

We model the interarrival times  $t_i$  using a Markov modulated process. Concretely, consider a two-state modulating Markov chain  $c_k$ , as depicted in Figure 4, with a transition matrix  $T$  given by

$$T = \begin{pmatrix} 1-p & p \\ q & 1-q \end{pmatrix}, \quad (23)$$

for some values  $0 < p, q < 1$ . In state  $i \in \{1, 2\}$  the interarrival times are given by iid random variables  $L_i$  with distribution  $\mathcal{L}_i$ . Without loss of generality we assume that  $L_1$  is stochastically smaller than  $L_2$ , i.e.,

$$\mathbf{P}[L_1 \geq t] \leq \mathbf{P}[L_2 \geq t],$$

for any  $t \geq 0$ . Additionally, we assume that the Markov chain  $c_k$  satisfies the burstiness condition

$$p < 1 - q, \quad (24)$$

i.e., the probability of jumping to a different state is less than the probability of staying in the same state.

Subsequent derivations will exploit the following exponential transform of the transition matrix  $T$  defined as

$$T_\theta := \begin{pmatrix} (1-p) \mathbb{E} \left[ e^{-\theta L_1} \right] & p \mathbb{E} \left[ e^{-\theta L_2} \right] \\ q \mathbb{E} \left[ e^{-\theta L_1} \right] & (1-q) \mathbb{E} \left[ e^{-\theta L_2} \right] \end{pmatrix},$$

for some  $\theta > 0$ . Let  $\Lambda(\theta)$  denote the maximal positive eigenvalue of  $T_\theta$ , and the vector  $h = (h(1), h(2))$  denote a corresponding eigenvector. By the Perron-Frobenius Theorem,  $\Lambda(\theta)$  is equal to the spectral radius of  $T_\theta$  such that  $h$  can be chosen with strictly positive components.

As in the case of renewal arrivals, we will next analyze both non-blocking and blocking FJ systems.

#### 4.1 Non-Blocking Systems

We first analyze a non-blocking FJ system fed with arrivals that are modulated by a stationary Markov chain as in Figure 4. We assume that the task service times  $x_{n,j}$  are iid and that the families  $\{t_i\}$  and  $\{x_{n,i}\}$  are independent. Note that both the definition of  $w_j$  from (1) and the representation of the steady-state waiting time  $w$  in (3) remain valid, due to stationarity and reversibility; the same holds for the response times.

The next theorem provides upper bounds on the steady-state waiting and response time distributions in the non-blocking scenario with Markov modulated interarrivals.

**Theorem 3** (NON-RENEWALS, NON-BLOCKING) *Given a FJ queueing system with  $N$  parallel non-blocking servers, Markov modulated job interarrivals  $t_j$  according to the Markov chain depicted in Figure 4 with transition matrix (23), and iid task service times  $x_{n,j}$ . The steady-state waiting and response time distributions are bounded by*

$$\mathbb{P}[w \geq \sigma] \leq N e^{-\theta_{nb} \sigma} \quad (25)$$

$$\mathbb{P}[r \geq \sigma] \leq N \mathbb{E} \left[ e^{\theta_{nb} x_{1,1}} \right] e^{-\theta_{nb} \sigma}, \quad (26)$$

where  $\theta_{nb}$  is the (positive) solution of

$$\mathbb{E} \left[ e^{\theta x_{1,1}} \right] \Lambda(\theta) = 1.$$

(Recall that  $\Lambda(\theta)$  was defined as a spectral radius.)

We remark that the existence of a positive solution  $\theta_{nb}$  is guaranteed by the Perron-Frobenius Theorem, see, e.g., [34].



*Proof* Consider the filtration

$$\mathcal{F}_k := \sigma \{x_{n,m}, t_m, c_m \mid m \leq k, n \in [1, N]\} ,$$

that includes information about the state  $c_k$  of the Markov chain. Now, we construct the process  $z(k)$  as

$$\begin{aligned} z(k) &= h(c_k) e^{\theta_{nb}(\max_n \sum_{i=1}^k x_{n,i} - \sum_{i=1}^k t_i)} \\ &= \left( e^{\theta_{nb}(\max_n \sum_{i=1}^k x_{n,i} - kD)} \right) \left( h(c_k) e^{\theta_{nb}(kD - \sum_{i=1}^k t_i)} \right) \end{aligned} \quad (27)$$

with the deterministic parameter

$$D := \theta_{nb}^{-1} \log \left( \mathbb{E} \left[ e^{\theta_{nb} x_{1,1}} \right] \right) .$$

Note the similarity of  $z(k)$  to (9) except for the additional function  $h$ . Roughly, the function  $h$  captures the correlation structure of the non-renewal interarrival time process.

Next we show that both terms of (27) are submartingales. In the first step we note that by the definition of  $D$ :

$$\mathbb{E} \left[ e^{\theta_{nb}(\sum_{i=1}^k x_{n,i} - kD)} \mid \mathcal{F}_{k-1} \right] = e^{\theta_{nb}(\sum_{i=1}^{k-1} x_{n,i} - (k-1)D)} ,$$

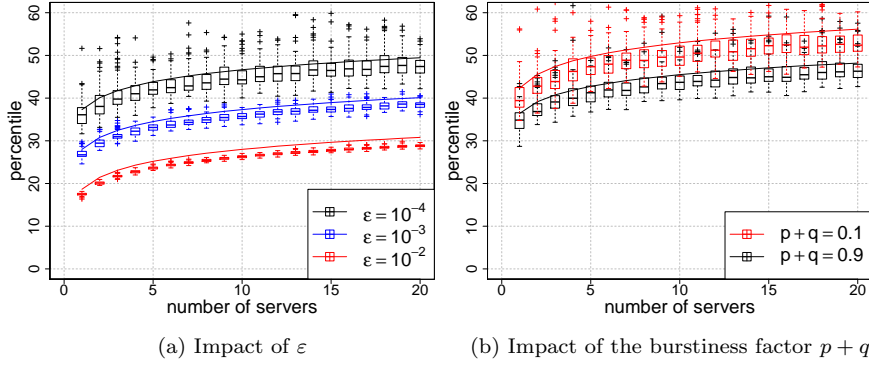
hence, following the line of argument in (10) the left factor of (27), which accounts for the additional  $\max_n$ , is a submartingale. The second step is similar to the derivations in [10, 14]. First, note that

$$\begin{aligned} \mathbb{E} \left[ h(c_k) e^{\theta_{nb}(D - t_k)} \mid \mathcal{F}_{k-1} \right] &= e^{\theta_{nb}D} T_{\theta_{nb}} h(c_{k-1}) \\ &= e^{\theta_{nb}D} \Lambda(\theta_{nb}) h(c_{k-1}) \\ &= h(c_{k-1}) , \end{aligned} \quad (28)$$

where the last line is due to the definitions of  $D$  and  $\theta_{nb}$ . Now, multiplying both sides of (28) by  $e^{\theta_{nb}((k-1)D - \sum_{i=1}^{k-1} t_i)}$  proves the martingale and hence the submartingale property of the right factor in (27). As the process  $z(k)$  is a product of two independent submartingales, it is a submartingale itself w.r.t.  $\mathcal{F}_k$ .

Next, we derive a bound on the steady-state waiting time distribution using the Optional Stopping Theorem. Here, we use the stopping time  $K$  defined in (11). Recall that  $\mathbb{P}[K < \infty] = \mathbb{P}[w \geq \sigma]$ . On the one hand we can write for every  $k \in \mathbb{N}$

$$\begin{aligned} \mathbb{E}[z(k)] &\geq \mathbb{E}[z(K \wedge k)] \\ &\geq \mathbb{E}[z(K \wedge k) 1_{K < k}] \\ &= \mathbb{E} \left[ \max_n h(c_K) e^{\theta_{nb}(\sum_{i=1}^K x_{n,i} - \sum_{i=1}^K t_i)} 1_{K < k} \right] \\ &\geq e^{\theta_{nb}\sigma} \mathbb{E}[h(c_K) 1_{K < k}] \\ &= e^{\theta_{nb}\sigma} \mathbb{E}[h(c_K) \mid K < k] \mathbb{P}[K < k] . \end{aligned} \quad (29)$$



**Fig. 5** The  $\mathcal{O}(\log N)$  scaling of waiting time percentiles  $w^\varepsilon$  for Markov modulated input (the non-blocking case (25)). The system parameters are  $\mu = 1$ ,  $\lambda_2 = 0.9$ ,  $\rho = 0.75$  (in both (a) and (b))  $p = 0.1$ ,  $q = 0.4$  (in (a)), three violation probabilities  $\varepsilon$  (in (a)),  $\varepsilon = 10^{-4}$  and only two burstiness parameters  $p + q$  (in (b)) (for visual convenience). Simulations include 100 runs, each accounting for  $10^7$  slots.

On the other hand we can upper bound the term

$$\begin{aligned} \mathbb{E}[z(k)] &= \mathbb{E}\left[\max_n e^{\theta_{nb}(\sum_{i=1}^k x_{n,i} - kD)}\right] \mathbb{E}\left[h(c_k) e^{\theta_{nb}(kD - \sum_{i=1}^k t_i)}\right] \\ &\leq N \mathbb{E}[h(c_1)] . \end{aligned}$$

Letting  $k \rightarrow \infty$  in (29) leads to

$$\mathbb{P}[K < \infty] \leq \frac{\mathbb{E}[h(c_1)]}{\mathbb{E}[h(c_K) | K < \infty]} N e^{-\theta_{nb}\sigma} . \quad (30)$$

In Lemma 2 it is shown that the distribution of the random variable  $(c_K | K < k)$  is stochastically smaller than the stationary distribution of the Markov chain. Given the burstiness condition in (24) and that the function  $h$  is monotonically decreasing [9], we can further upper bound the prefactor in (30) as

$$\frac{\mathbb{E}[h(c_1)]}{\mathbb{E}[h(c_K) | K < \infty]} \leq 1 ,$$

which completes the proof. The proof for the response time  $r$  is analogous.

**Remark:** Note that, if the burstiness condition (24) is not fulfilled then we can still upper bound the prefactor in (30) using the trivial upper bound

$$\frac{\mathbb{E}[h(c_1)]}{\mathbb{E}[h(c_K) | K < \infty]} \leq \frac{\mathbb{E}[h(c_1)]}{\min_k h(c_k)} .$$

Figure 5 displays the bounds on the waiting time percentiles  $w^\varepsilon$ , for various violation probabilities  $\varepsilon$ , in the FJ system with non-renewal input. The bounds closely match the corresponding simulation results, shown as box-plots, while also exhibiting the  $\mathcal{O}(\log N)$  scaling behavior (which can be also derived from both (25) and (26), as in Section 3).

## 4.2 Blocking Systems

Now we turn to the blocking variant of the FJ system that is fed by the same non-renewal arrivals as in the previous section. In the following, we consider exponential distributions  $\mathcal{L}_m$  for  $m \in \{1, 2\}$ . The main result is:

**Theorem 4** (NON-RENEWALS, BLOCKING) *Given a FJ system with  $N$  blocking servers, Markov modulated job interarrivals  $t_j$ , and iid task service times  $x_{n,j}$ . The steady-state waiting and response time distributions are bounded by*

$$\begin{aligned} \mathbb{P}[w \geq \sigma] &\leq e^{-\theta_b \sigma} \\ \mathbb{P}[r \geq \sigma] &\leq \mathbb{E}[e^{\theta_b \max_n x_{n,1}}] e^{-\theta_b \sigma}, \end{aligned} \quad (31)$$

where  $\theta_b$  is the (positive) solution of

$$\mathbb{E}[e^{\theta \max_n x_{n,1}}] \Lambda(\theta) = 1.$$

We remark that the positive solution for  $\theta_b$  is guaranteed under the stronger stability condition  $\mathbb{E}[t_1] > \mathbb{E}[\max_n x_{n,1}]$  and the Perron-Frobenius Theorem.

*Proof* Let  $D := \theta_b^{-1} \log \mathbb{E}[e^{\theta_b \max_n x_{n,1}}]$  and define the process  $y$  by:

$$\begin{aligned} y(k) &= h(c_k) e^{\theta_b (\sum_{i=1}^k \max_n x_{n,i} - \sum_{i=1}^k t_i)} \\ &= (e^{\theta_b (\sum_{i=1}^k \max_n x_{n,i} - kD)}) (h(c_k) e^{\theta_b (kD - \sum_{i=1}^k t_i)}). \end{aligned}$$

Similarly to the proofs of Theorem 2 and Theorem 3 one can show that both the first and second factor of  $y$  are martingales, and hence  $y$  is a martingale. We use the stopping time  $K$  in (20) and write

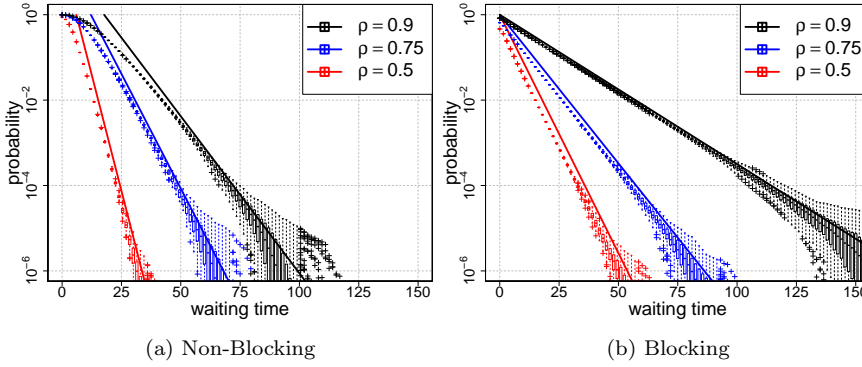
$$\begin{aligned} \mathbb{E}[h(c_1)] &= \mathbb{E}[y(0)] \\ &\geq \mathbb{E}[y(K \wedge k)] \\ &\geq \mathbb{E}[y(K \wedge k) 1_{K < k}] \\ &= \mathbb{E}\left[e^{\theta_b (\sum_{i=1}^K \max_n x_{n,i} - \sum_{i=1}^K t_i)} h(c_K) 1_{K < k}\right] \\ &\geq e^{\theta_b \sigma} \mathbb{E}[h(c_K) | K < \infty] \mathbb{P}[K < k]. \end{aligned}$$

Taking  $k \rightarrow \infty$  we obtain the bound

$$\mathbb{P}[K < \infty] \leq \frac{\mathbb{E}[h(c_1)]}{\mathbb{E}[h(c_K) | K < \infty]} e^{-\theta_b \sigma} \leq e^{-\theta_b \sigma},$$

where we used Lemma 2 for the last inequality. The proof for  $r$  is analogous.

A close comparison of the waiting time bound in the non-renewal case (31) to the corresponding bound in the renewal case (18) reveals that the decay factors  $\theta_b$  depend on similar conditions, whereby the MGF of the interarrival times in (18) is replaced by the spectral radius of the modulating Markov chain in (31). Moreover, given the ergodicity of the underlying Markov chain,



**Fig. 6** Bounds on the waiting time distributions vs. simulations (non-renewal input): (a) the non-blocking case (25) and (b) the blocking case (31). The parameters are  $N = 20, \mu = 1, p = 0.1, q = 0.4, \lambda_1 \in \{0.4, 0.72, 0.72\}$  and  $\lambda_2 \in \{0.9, 0.9, 1.62\}$  leading to utilizations  $\rho \in \{0.5, 0.75, 0.9\}$ . Simulations include 100 runs, each accounting for  $10^7$  slots.

the blocking system with non-renewal input is subject to the same degrading stability region (in  $\log N$ ) as in the renewal case (recall (21)).

For quick numerical illustrations of the tightness of the bounds on the waiting time distributions in both the non-blocking and blocking cases we refer to Figure 6.

So far we have contributed stochastic bounds on the steady-state waiting and response time distributions in FJ systems fed with either renewal and non-renewal job arrivals. The key technical insight was that the stochastic bounds in the non-blocking model grow as  $\mathcal{O}(\log N)$  in the number of parallel servers  $N$  under non-renewal arrivals, which extends a known result for renewal arrivals [31, 5]. The same fundamental factor of  $\log N$  was shown to drive the stability region in the blocking model. A concrete application follows next.

## 5 Partial Mapping

In this section we consider FJ queueing systems where jobs are mapped to a subset of  $H \leq N$  servers. This model captures a crucial aspect of the operation of parallel systems, i.e., the amount of resources provided to some job is not necessarily the entire amount of resources available. This corresponds, for example, to batch systems, where servers are grouped into resource pools and incoming jobs are assigned to one such pool. In general, partial mapping provides a basis for service differentiation and isolation within parallel systems. In the following we regard two contrasting types of partial mapping, i.e., a rigid round-robin mapping and a random partial mapping of jobs to  $H \leq N$  servers. The subsequent analysis of the fan-out ratio  $H/N$  on the system performance provides a reference for dimensioning such server pools. In the following, we

restrict the exposition to the more interesting case of non-blocking servers since most of the derivations rely on results from Sections 3 and 4.

### 5.1 Round-robin Partial Mapping, Dyadic System

We consider a dyadic FJ system where the number of servers is given as  $N = 2^W$  (with  $W \geq 1$ ) and a job is split into  $H = 2^V$  tasks (with  $1 \leq V \leq W$ ). The assignment of tasks to servers follows a round-robin scheme such that the first job is assigned to servers  $1, \dots, H$ , the second to the servers  $H+1, \dots, 2H$ , etc.

In the following, we consider job arrivals as renewal processes similar to Sect. 3. For the analysis it is sufficient to look only at an equivalent “FJ subsystem” that consists of only  $H$  servers and adjust the job interarrival times  $\bar{t}_k$  to that system accordingly:

$$\bar{t}_k := \sum_{i=1}^{2^{(W-V)}} t_{(k-1)2^{(W-V)}+i}.$$

Note that for the extremal case  $V = W$  we recover the scenario from Sect. 3, i.e.,  $\bar{t}_k = t_k$ .

The Laplace transform of the job interarrival times  $\bar{t}_k$  to one subsystem is obtained directly from the Laplace transform of the original job interarrival times  $t_k$  and the number of subsystems:

$$\mathbb{E} \left[ e^{-\theta \bar{t}_1} \right] = \mathbb{E} \left[ e^{-\theta t_1} \right]^{2^{W-V}} = \mathbb{E} \left[ e^{-\theta t_1} \right]^{\frac{N}{H}}.$$

The steady-state waiting time distribution now has the following representation:

$$w =_{\mathcal{D}} \max_{k \geq 0} \left\{ \max_{1 \leq n \leq H} \left\{ \sum_{i=1}^k x_{n,i} - \sum_{i=1}^k \bar{t}_i \right\} \right\} \quad (32)$$

and the response time:

$$r =_{\mathcal{D}} \max_{k \geq 0} \left\{ \max_{1 \leq n \leq H} \left\{ \sum_{i=0}^k x_{n,i} - \sum_{i=1}^k \bar{t}_i \right\} \right\}. \quad (33)$$

The next theorem provides upper bounds on the steady-state waiting and response time distributions in the non-blocking scenario with partial round-robin mapping and renewal interarrivals.

**Theorem 5** (ROUND-ROBIN MAPPING, RENEWALS, NON-BLOCKING) *Given a FJ queueing system with  $N = 2^W$  non-blocking servers and partial round-robin mapping of jobs to  $H = 2^V$  servers with  $1 \leq V \leq W$ . The system is fed by renewal job arrivals with interarrivals  $t_j$ . If the input job size is normalized such that the MGF of the task service time is given as  $\mathbb{E} [e^{\theta x_{n,i}/H}]$ , with the*

service times  $x_{n,i}$  being iid, then the steady-state waiting and response times  $w$  and  $r$  are bounded by

$$\begin{aligned} \mathbb{P}[w \geq \sigma] &\leq H e^{-\theta \sigma} , \\ \mathbb{P}[r \geq \sigma] &\leq H \mathbb{E}[e^{\theta x_{1,1}}] e^{-\theta \sigma} , \end{aligned}$$

where  $\theta$  is the solution of

$$\mathbb{E}[e^{\theta x_{1,1}/H}] \mathbb{E}[e^{-\theta t_1}]^{\frac{N}{H}} = 1 . \quad (34)$$

*Proof* The proof goes along the same arguments of the proof of Theorem 1, however, with modified MGF and Laplace transform for the task service times  $x_{n,i}$  and the job interarrival times  $t_i$ , respectively.

The rationale behind the normalization of the input job size such that the MGF of the task service time is given as  $\mathbb{E}[e^{\theta x_{n,i}/H}]$  is to compare different fan-out factors  $H$  such that the mean task service time is  $\mathbb{E}[x]/H$ .

*Example: Exponentially distributed interarrival and service times*

In the case of exponentially distributed interarrival times with parameter  $\lambda$  the job interarrival times at one subsystem have an Erlang  $E_{\frac{N}{H}}$  distribution. We assume the tasks are exponentially distributed with a mean  $1/H\mu$ . The condition (34) from Theorem 5 becomes

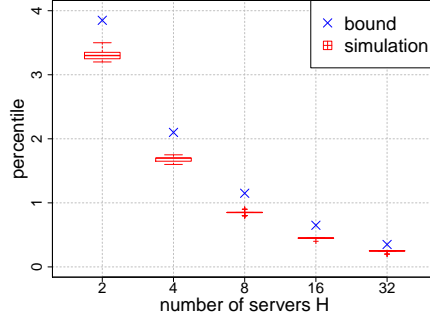
$$\left(\frac{H\mu}{H\mu - \theta}\right) \left(\frac{\lambda}{\lambda + \theta}\right)^{\frac{N}{H}} = 1 . \quad (35)$$

In Figure 7 we show simulation box-plots as well as corresponding bounds on the waiting time percentile  $w^\varepsilon$  from Theorem 5 for an increasing number of fan-out servers  $H$ . Observe the diminishing gain in terms of waiting time reduction with increasing the server fan-out.

## 5.2 Random Partial Mapping

Here, we consider a system that randomly maps a job to  $H$  out of  $N$  available servers based on a uniform distribution over the set  $\{A \subseteq \{1, \dots, N\} | |A| = H\}$  of server combinations with cardinality  $H$ . We bound the job waiting and response time in this system using the following abstraction which considers the probability of assigning a task to a specific server. Note that the probability for a task dedicated to a certain server is given by  $p_d = H/N$ . Now, if we focus on only one server of this FJ system, the task service times at that server can be represented by the compound distribution

$$\bar{x}_{n,i} = \begin{cases} x_{n,i} & \text{with probability } p_d \\ 0 & \text{with probability } 1 - p_d , \end{cases} \quad (36)$$



**Fig. 7** Round-robin partial mapping: Bound on the waiting time percentile  $w^\varepsilon$  for renewal arrivals and increasing number of servers (fan-out)  $H$ . The system parameters are  $\mu = 1$ ,  $\lambda = 0.75$ ,  $\varepsilon = 10^{-3}$  and the overall number of servers is  $N = 2^8$ .

since a job that is not assigned to this server can be considered to have a service time equal to 0. Hence, one server of this FJ system with random partial mapping can be modelled as if it is part of a FJ system with full mapping as in Sect. 3, but with the modified service times  $\bar{x}_{n,i}$ . Note that the MGF of  $\bar{x}_{n,i}$  can be computed as:

$$\mathbb{E} [e^{\theta \bar{x}_{n,i}}] = (1 - p_d) + p_d \mathbb{E} [e^{\theta x_{n,i}}] .$$

The representations for the waiting and response time, respectively, become

$$w =_{\mathcal{D}} \max_{k \geq 0} \left\{ \max_{1 \leq n \leq H} \left\{ \sum_{i=1}^k \bar{x}_{n,i} - \sum_{i=1}^k t_i \right\} \right\} , \quad (37)$$

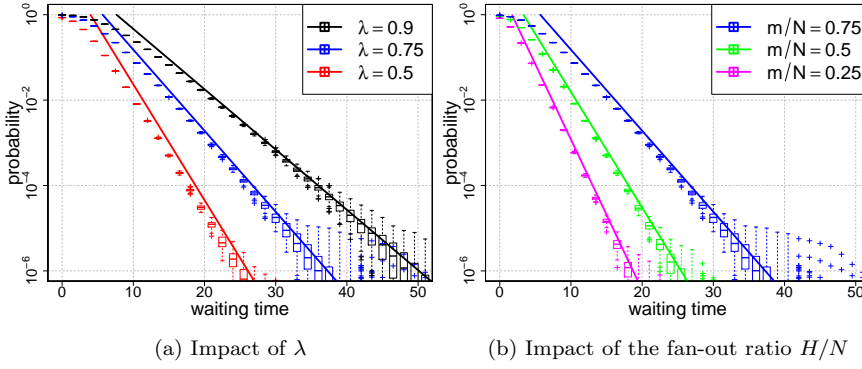
and

$$r =_{\mathcal{D}} \max_{k \geq 0} \left\{ \max_{1 \leq n \leq H} \left\{ x_{n,0} + \sum_{i=1}^k \bar{x}_{n,i} - \sum_{i=1}^k t_i \right\} \right\} . \quad (38)$$

Note the asymmetry for the response time in (38). For  $i \geq 1$  we consider the modified service times  $\bar{x}_{n,i}$  as the corresponding server is only selected with probability  $p_d$ . In turn, for  $i = 0$ , we need to consider the unmodified service time  $x_{0,i}$  as we only look at those servers which have been selected for mapping.

The following theorems provide upper bounds on the steady-state waiting and response time distributions in the non-blocking scenarios with partial random mapping for renewal and Markov-modulated interarrivals, respectively.

**Theorem 6** (RANDOM MAPPING, RENEWALS, NON-BLOCKING) *Given a FJ queueing system with  $N$  servers and random partial mapping of jobs to  $H \leq N$  servers based on a uniform distribution over the set  $\{A \subseteq \{1, \dots, N\} | |A| = H\}$  of server combinations with cardinality  $H$ . The system is fed with renewal job*



**Fig. 8** Bounds on the waiting time distributions vs. simulation box-plots for renewal input with random server mapping. The parameters are  $N = 16, \mu = 1$ . (a) Here, we fix the fan-out ratio to  $H = 12$  and change the job arrival rate  $\lambda \in \{0.5, 0.75, 0.9\}$  while in (b) we fix the arrival rate to  $\lambda = 0.75$  and vary the fan-out ratio  $H/N \in \{0.25, 0.5, 0.75\}$ . Simulations include 100 runs, each accounting for  $10^6$  slots.

arrivals. If the task service times  $x_{n,j}$  are iid, then the steady-state waiting and response times  $w$  and  $r$  are bounded by

$$\begin{aligned} P[w \geq \sigma] &\leq H e^{-\theta \sigma}, \\ P[r \geq \sigma] &\leq H E[e^{\theta x_{1,1}}] e^{-\theta \sigma}, \end{aligned}$$

where  $\theta$  is the solution of

$$((1 - p_d) + p_d E[e^{\theta x_{n,i}}]) E[e^{-\theta t_1}] = 1. \quad (39)$$

*Proof* The proof goes along similar steps as for Theorem 5, however, using the process

$$z_n(k) = e^{\theta \sum_{i=1}^k (\bar{x}_{n,i} - t_i)}$$

which is a martingale for each  $n \leq N$  under the criterion (39) on  $\theta$ .

Figure 8 shows a numerical illustration of the tightness of the bounds on the waiting time distribution from Theorem 6. The illustrated results are for the example of exponentially distributed interarrival and service times with parameters  $\lambda$  and  $\mu$ , respectively.

By combining the above consideration of the compound service time distribution with the results from Section 4, one can extend the analysis of random partial mapping to the case of non-renewal input.

**Theorem 7** (RANDOM MAPPING, NON-RENEWALS, NON-BLOCKING) *Given a FJ queueing system with  $N$  parallel non-blocking servers, Markov modulated*



job interarrivals  $t_j$  as in Section 4, and task service times  $\bar{x}_{n,i}$  that are described by Eq. (36). Jobs are randomly mapped to servers according to a uniform distribution over the set of server combinations with cardinality  $H$ . The steady-state waiting and response time distributions are bounded by

$$\begin{aligned} \mathbb{P}[w \geq \sigma] &\leq H e^{-\theta\sigma}, \\ \mathbb{P}[r \geq \sigma] &\leq H \mathbb{E}[e^{\theta x_{1,1}}] e^{-\theta\sigma}, \end{aligned}$$

where  $\theta$  is the solution of

$$((1 - p_d) + p_d \mathbb{E}[e^{\theta x_{1,1}}]) \Lambda(\theta) = 1.$$

(Recall that  $\Lambda(\theta)$  was defined as a spectral radius of  $T_\theta$  in Section 4).

*Proof* The proof follows analogously to the proof of Theorem 3 with the difference that  $x_{n,i}$  is replaced by  $\bar{x}_{n,i}$  and  $N$  by  $H$ , respectively.

**Remark: Random number of servers H:** One variation of the system that is considered in Sect. 5.2 is a random mapping of arriving jobs to a random number of servers  $1 \leq H \leq N$  based on a uniform distribution over the power set  $\{2^A \setminus \emptyset\}$  with  $A = \{1, \dots, N\}$ . In this case the steady state waiting and response times are bounded by

$$\begin{aligned} \mathbb{P}[w \geq \sigma] &\leq N e^{-\theta\sigma}, \\ \mathbb{P}[r \geq \sigma] &\leq N \mathbb{E}[e^{\theta x_{1,1}}] e^{-\theta\sigma}, \end{aligned}$$

where  $\theta$  is the solution of (39) with  $p_d = 2^{N-1}/(2^N - 1)$ .

## 6 Application to Window-based Protocols over Multipath Routing

In this section we slightly adapt and use the non-blocking FJ queueing system from Section 3.1 to analyze the performance of a *generic* window-based transmission protocol over multipath routing. While this problem has attracted much interest lately with the emergence of multipath TCP [35], it is subject to a major difficulty due to the likely overtaking of packets on different paths. Consequently, packets have to additionally wait for a *resequencing delay*, which directly corresponds to the synchronization constraint in FJ systems. We note that the employed FJ non-blocking model is subject to a convenient simplification, i.e., each path is modelled by a single server/queue only.

As depicted in Figure 9, we consider an arrival flow containing  $l$  batches of  $N$  packets, with  $l \in \mathbb{N}$ , at the fork node  $A$ . In practice, a *packet* as denoted here may represent an entire train of consecutive datagrams. The incoming packets are sent over multiple paths to the destination node  $B$ , where they need to be eventually reordered. We assume that the batch size corresponds to the transmission window size of the protocol, such that one packet traverses a single path only. For example, the first path transmits the packets  $\{1, N + 1, 2N + 1, \dots\}$ , i.e., packets are distributed in a round-robin fashion over the  $N$

paths. We also assume that packets on each path are delivered in a (locally-) FIFO order, i.e., there is no overtaking on the same path.

In analogy to Section 3.1, we consider a batch waiting until its last packet starts being transmitted. When the transmission of the last packet of batch  $j$  begins, the previous batch has already been received, i.e., all packets of the batch  $j - 1$  are *in order* at node  $B$ .

We are interested in the response times of the batches, which are upper bounded by the largest response time of the packets therein. The arrival time of a batch is defined as the latest arrival time of the packets therein, i.e., when the batch is entirely received. Formally, the response time of batch  $j \in \{lN + 1 \mid l \in \mathbb{N}\}$  can be given by slightly modifying (2), i.e.,

$$r_j = \max_{0 \leq k \leq j-1} \left\{ \max_n \left\{ \sum_{i=0}^k x_{n,j-i} - \sum_{i=1}^k t_{n,j-i} \right\} \right\} .$$

The corresponding steady-state response time has the modified representation

$$r =_{\mathcal{D}} \max_{k \geq 0} \left\{ \max_n \left\{ \sum_{i=0}^k x_{n,i} - \sum_{i=1}^k t_{n,i} \right\} \right\} .$$

The modifications account for the fact that the packets of each batch are asynchronously transmitted on the corresponding paths (instead, in the basic FJ systems, the tasks of each job are simultaneously mapped). In terms of notations, the  $t_{n,i}$ 's now denote the interarrival times of the packets transmitted over the same path  $n$ , whereas  $x_{n,i}$ 's are iid and denote the transmission time of packet  $i$  over path  $n$ ; as an example, when the arrival flow at node  $A$  is Poisson,  $t_{n,i}$  has an Erlang  $E_N$  distribution for all  $n$  and  $i$ .

We next analyze the performance of the considered multipath routing for both renewal and non-renewal input.

### *Renewal Arrivals*

Consider first the scenario with renewal interarrival times. Similarly to Section 3.1 we bound the distribution of the steady-state response time  $r$  using a submartingale in the time domain  $j \in \{lN + 1 \mid l \in \mathbb{N}\}$ . Following the same steps as in Theorem 1, the process

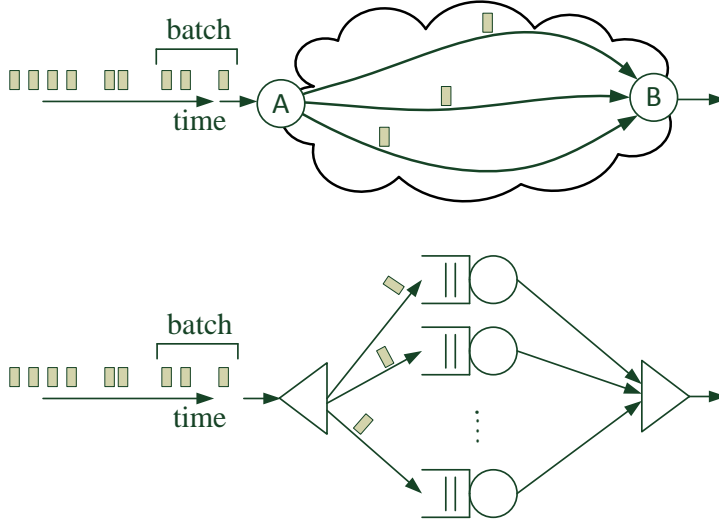
$$z_n(k) = e^{\theta(\sum_{i=0}^k x_{n,i} - \sum_{i=1}^k t_{n,i})}$$

is a martingale under the condition

$$\mathbb{E} [e^{\theta x_{1,1}}] \mathbb{E} [e^{-\theta t_{1,1}}] = 1 ,$$

where we used the filtration

$$\mathcal{F}_k := \sigma\{x_{n,m}, t_{n,m} \mid m \leq k, n \in [1, N]\} .$$



**Fig. 9** A schematic description of the window-based transmission over multipath routing; each path is modelled as a single server/queue.

Note that  $\mathbf{E} [e^{-\theta t_{1,1}}]$  denotes the Laplace transform of the interarrival times of packets transmitted over each path. The proof that  $\max_n z_n(k)$  is a submartingale follows a similar argument as in (10). Hence, we can bound the distribution of the steady-state response time as

$$\mathbf{P} [r \geq \sigma] \leq N \mathbf{E} [e^{\theta x_{1,1}}] e^{-\theta \sigma}, \quad (40)$$

with the condition on  $\theta$  from above.

#### Non-Renewal Arrivals

Next, consider a scenario with non-renewal interarrival times  $t_i$  of the packets arriving at the fork node  $A$  in Figure 9, as described in Section 4. On every path  $n \in [1, N]$  the interarrivals are given by a sub-chain  $(c_{n,k})_k$  that is driven by the  $N$ -step transition matrix  $T^N = (\alpha_{i,j})_{i,j}$  for  $T$  given in (23). Similarly as in the proof of Theorem 3, we will use an exponential transform  $(T^N)_\theta$  of the transition matrix that describes each path  $n$ , i.e.,

$$(T^N)_\theta := \begin{pmatrix} \alpha_{1,1}\beta_1 & \alpha_{1,2}\beta_2 \\ \alpha_{2,1}\beta_1 & \alpha_{2,2}\beta_2 \end{pmatrix},$$

with  $\alpha_{i,j}$  defined above and  $\beta_1, \beta_2$  being the elements of a vector  $\beta$  of conditional Laplace transforms of  $N$  consecutive interarrival times  $t_i$ . The vector  $\beta$  is given by

$$\beta := \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} \mathbf{E} [e^{-\theta \sum_{i=1}^N t_i} \mid c_1 = 1] \\ \mathbf{E} [e^{-\theta \sum_{i=1}^N t_i} \mid c_1 = 2] \end{pmatrix},$$

and can be computed given the transition matrix  $T$  from (23) via an exponential row transform [10] (Example 7.2.7) denoted by

$$\tilde{T}_\theta := \begin{pmatrix} (1-p)\mathbb{E}[e^{-\theta L_1}] & p\mathbb{E}[e^{-\theta L_1}] \\ q\mathbb{E}[e^{-\theta L_2}] & (1-q)\mathbb{E}[e^{-\theta L_2}] \end{pmatrix},$$

yielding  $\beta = (\tilde{T}_\theta)^N \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ .

Denote  $\Lambda(\theta)$  and  $h = (h(1), h(2))$  as the maximal positive eigenvalue of the matrix  $(\tilde{T}_\theta)^N$  and the corresponding right eigenvector, respectively. Mimicking the proof of Theorem 3, one can show for every path  $n$  that the process

$$z_n(k) = h(c_{n,k}) e^{\theta(\sum_{i=0}^k x_{n,i} - \sum_{i=1}^k t_{n,i})}$$

is a martingale under the condition on (positive)  $\theta$

$$\mathbb{E}[e^{\theta x_{1,1}}] \Lambda(\theta) = 1. \quad (41)$$

Given the martingale representation of the processes  $z_n(k)$  for every path  $n$ , the process

$$z(k) = \max_n z_n(k)$$

is a submartingale following the line of argument in (10). We can now use (30) and the remark at the end of Section 4.1 to bound the distribution of the steady-state response time  $r$  as

$$\mathbb{P}[r \geq \sigma] \leq \frac{\mathbb{E}[h(c_{1,1})]}{h(2)} N \mathbb{E}[e^{\theta x_{1,1}}] e^{-\theta \sigma}, \quad (42)$$

where we also used that  $h$  is monotonically decreasing and  $\theta$  as defined in (41).

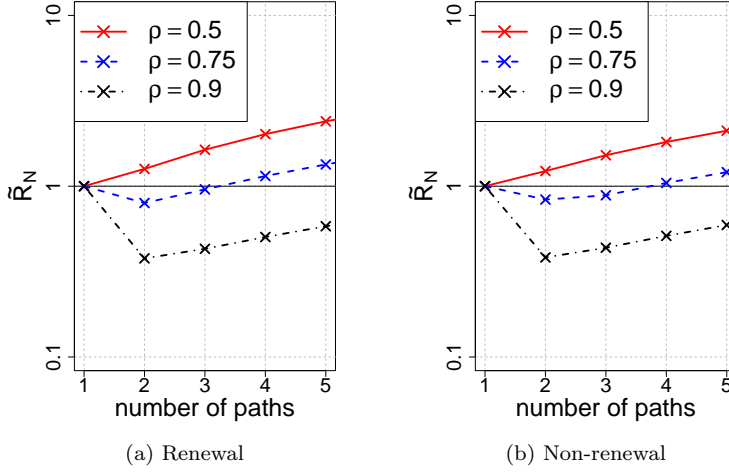
As a direct application of the obtained stochastic bounds (i.e., (40) and (42)), consider the problem of optimizing the number of parallel paths  $N$  subject to the batch delay (accounting for both queueing and resequencing delays). More concretely, we are interested in the number of paths  $N$  minimizing the overall average batch delay. Note that the path utilization changes with  $N$  as

$$\rho = \frac{\lambda}{N\mu},$$

since each path only receives  $\frac{1}{N}$  of the input. In other words, the packets on each path are delivered much faster with increasing  $N$ , but they are subject to the additional resequencing delay (which increases as  $\log N$  as shown in Section 3.1).

To visualize the impact of increasing  $N$  on the average batch response times we use the ratio

$$\tilde{R}_N := \frac{E[r_N]}{E[r_1]},$$



**Fig. 10** Multipath routing reduces the average batch response time when  $\tilde{R}_N < 1$ ; smaller  $\tilde{R}_N$  corresponds to larger reductions. Baseline parameter  $\mu = 1$  and non-renewal parameters:  $p = 0.1, q = 0.4, \lambda_1 = \{0.39, 0.7, 0.88\}, \lambda_2 = 0.95$ , yielding the utilizations  $\rho = \{0.5, 0.75, 0.9\}$  (from top to bottom).

where, with abuse of notation,  $E[r_N]$  denotes a bound on the average batch response time for some  $N$ , and  $E[r_1]$  denotes the corresponding baseline bound for  $N = 1$ ; both bounds are obtained by integrating either (40) or (42) for the renewal and the non-renewal case, respectively.

In the renewal case, with exponentially distributed interarrival times with parameter  $\lambda$ , and homogenous paths/servers where the service times are exponentially distributed with parameter  $\mu$ , we obtain

$$\tilde{R}_N = \left( \frac{\log(N\mu/(\mu - \theta)) + 1}{\log(1/\rho) + 1} \right) \left( \frac{\mu - \lambda}{\theta} \right), \quad (43)$$

where  $\theta$  is the solution of

$$\frac{\mu}{\mu - \theta} \left( \frac{\lambda}{\lambda + \theta} \right)^N = 1.$$

In the non-renewal case we obtain the same expression for  $\tilde{R}_N$  as in (43) except for the additional prefactor  $\frac{E[h(c_1(1))]}{h(2)}$  prior to  $N$ ; moreover,  $\theta$  is the implicit solution from (41).

Figure 10 illustrates  $\tilde{R}_N$  as a function of  $N$  for several utilization levels  $\rho$  for both renewal (a) and non-renewal (b) input; recall that the utilization on each path is  $\frac{\rho}{N}$ . In both cases, the fundamental observation is that at small utilizations (i.e., roughly when  $\rho \leq 0.5$ ), multipath routing increases the response times. In turn, at higher utilizations, response times benefit from multipath routing but only for 2 paths. While this result may appear as counterintuitive,

the technical explanation (in (a)) is that the waiting time in the underlying  $E_N/M/1$  queue quickly converges to  $\frac{1}{\mu}$ , whereas the resequencing delay grows as  $\log N$ ; in other words, the gain in the queueing delay due to multipath routing is quickly dominated by the resequencing delay price.

## 7 Conclusions

In this paper we have provided the first computable and non-asymptotic bounds on the waiting and response time distributions in Fork-Join queueing systems under full and partial server mapping. We have analyzed four practical scenarios comprising of either workconserving or non-workconserving servers, which are fed by either renewal or non-renewal arrivals. In the case of workconserving servers, we have shown that delays scale as  $\mathcal{O}(\log N)$  in the number of parallel servers  $N$ , extending a related scaling result from renewal to non-renewal input. In turn, in the case of non-workconserving servers, we have shown that the same fundamental factor of  $\log N$  determines the system's stability region. Given their inherent tightness, our results can be directly applied to the dimensioning of Fork-Join systems such as MapReduce clusters and multipath routing. A highlight of our study is that multipath routing is reasonable from a queueing perspective for two routing paths only.

## References

1. Amazon Elastic Compute Cloud EC2. <http://aws.amazon.com/ec2>
2. Abate, J., Choudhury, G.L., Whitt, W.: Exponential approximations for tail probabilities in queues, I: Waiting times. *Oper. Res.* **43**, 885–901 (1995)
3. Babu, S.: Towards automatic optimization of MapReduce programs. In: *Proc. of ACM SoCC*, pp. 137–142 (2010)
4. Baccelli, F., Gelenbe, E., Plateau, B.: An end-to-end approach to the resequencing problem. *J. ACM* **31**(3), 474–485 (1984)
5. Baccelli, F., Makowski, A.M., Schwartz, A.: The Fork-Join queue and related systems with synchronization constraints: Stochastic ordering and computable bounds. *Adv. in Appl. Probab.* **21**(3), 629–660 (1989)
6. Balsamo, S., Donatiello, L., Van Dijk, N.M.: Bound performance models of heterogeneous parallel processing systems. *IEEE Trans. Parallel Distrib. Syst.* **9**(10), 1041–1056 (1998)
7. Billingsley, P.: *Probability and Measure*, 3rd edn. Wiley (1995)
8. Boxma, O., Koole, G., Liu, Z.: Queueing-theoretic solution methods for models of parallel and distributed systems. In: *Proc. of Performance Evaluation of Parallel and Distributed Systems*. CWI Tract 105, pp. 1–24 (1994)
9. Buffet, E., Duffield, N.G.: Exponential upper bounds via martingales for multiplexers with Markovian arrivals. *J. Appl. Probab.* **31**(4), 1049–1060 (1994)
10. Chang, C.S.: *Performance Guarantees in Communication Networks*. Springer (2000)
11. Chen, Y., Alspaugh, S., Katz, R.: Interactive analytical processing in big data systems: A cross-industry study of mapreduce workloads. *Proc. VLDB Endow.* **5**(12), 1802–1813 (2012)
12. Ciucu, F., Poloczek, F., Schmitt, J.: Sharp per-flow delay bounds for bursty arrivals: The case of FIFO, SP, and EDF scheduling. In: *Proc. of IEEE INFOCOM*, pp. 1896–1904 (2014)
13. Dean, J., Ghemawat, S.: MapReduce: Simplified data processing on large clusters. *Commun. ACM* **51**(1), 107–113 (2008)

14. Duffield, N.: Exponential bounds for queues with Markovian arrivals. *Queueing Syst.* **17**(3–4), 413–430 (1994)
15. Flatto, L., Hahn, S.: Two parallel queues created by arrivals with two demands I. *SIAM J. Appl. Math.* **44**(5), 1041–1053 (1984)
16. Ganesh, A., O’Connell, N., Wischik, D.: Big queues. No. 1838 in *Lecture notes in mathematics*. Springer (2004)
17. Gibbens, R.J.: Traffic characterisation and effective bandwidths for broadband network traces. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* (1996)
18. Han, Y., Makowski, A.: Resequencing delays under multipath routing - Asymptotics in a simple queueing model. In: *Proc. of IEEE INFOCOM*, pp. 1–12 (2006)
19. Harrus, G., Plateau, B.: Queueing analysis of a reordering issue. *IEEE Trans. Softw. Eng.* **8**(2), 113–123 (1982)
20. Jiang, Y., Liu, Y.: *Stochastic Network Calculus*. Springer (2008)
21. Joshi, G., Liu, Y., Soljanin, E.: Coding for fast content download. In: *Proc. of the Allerton Conference on Communication, Control, and Computing*, pp. 326–333 (2012)
22. Joshi, G., Liu, Y., Soljanin, E.: On the delay-storage trade-off in content download from coded distributed storage systems. *IEEE J. Sel. Areas Commun.* **32**(5), 989–997 (2014)
23. Kandula, S., Sengupta, S., Greenberg, A., Patel, P., Chaiken, R.: The nature of data center traffic: Measurements & analysis. In: *Proc. of ACM IMC*, pp. 202–208 (2009)
24. Kavulya, S., Tan, J., Gandhi, R., Narasimhan, P.: An analysis of traces from a production MapReduce cluster. In: *Proc. of IEEE/ACM CCGRID*, pp. 94–103 (2010)
25. Kemper, B., Mandjes, M.: Mean sojourn times in two-queue Fork-Join systems: Bounds and approximations. *OR Spectr.* **34**(3), 723–742 (2012)
26. Kesidis, G., Urgaonkar, B., Shan, Y., Kamarava, S., Liebeherr, J.: Network calculus for parallel processing. In: *Proc. of the ACM MAMA workshop* (2015)
27. Kingman, J.F.C.: Inequalities in the theory of queues. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **32**(1), 102–110 (1970)
28. Ko, S.S., Serfozo, R.F.: Sojourn times in G/M/1 Fork-Join networks. *Naval Res. Logist.* **55**(5), 432–443 (2008)
29. Lebrecht, A.S., Knottenbelt, W.J.: Response time approximations in Fork-Join queues. In: *Proc. of UKPEW* (2007)
30. Lu, H., Pang, G.: Gaussian limits for a Fork-Join network with nonexchangeable synchronization in heavy traffic. *Math. Oper. Res.* **41**(2), 560–595 (2016)
31. Nelson, R., Tantawi, A.: Approximate analysis of Fork/Join synchronization in parallel queues. *IEEE Trans. Computers* **37**(6), 739–743 (1988)
32. Pike, R., Dorward, S., Griesemer, R., Quinlan, S.: Interpreting the data: Parallel analysis with Sawzall. *Sci. Program.* **13**(4), 277–298 (2005)
33. Polato, I., R, R., Goldman, A., Kon, F.: A comprehensive view of Hadoop research - a systematic literature review. *J. Netw. Comput. Appl.* **46**(0), 1 – 25 (2014)
34. Poloczek, F., Ciucu, F.: Scheduling analysis with martingales. *Perform. Evaluation* **79**, 56–72 (2014)
35. Raiciu, C., Barre, S., Pluntke, C., Greenhalgh, A., Wischik, D., Handley, M.: Improving datacenter performance and robustness with multipath TCP. *SIGCOMM Comput. Commun. Rev.* **41**(4), 266–277 (2011)
36. Rényi, A.: On the theory of order statistics. *Acta Math. Hungar.* **4**(3–4), 191–231 (1953)
37. Tan, J., Meng, X., Zhang, L.: Delay tails in MapReduce scheduling. *SIGMETRICS Perform. Eval. Rev.* **40**(1), 5–16 (2012)
38. Tan, J., Wang, Y., Yu, W., Zhang, L.: Non-work-conserving effects in MapReduce: Diffusion limit and criticality. *SIGMETRICS Perform. Eval. Rev.* **42**(1), 181–192 (2014)
39. Varki, E.: Mean value technique for closed Fork-Join networks. *SIGMETRICS Perform. Eval. Rev.* **27**(1), 103–112 (1999)
40. Varma, S., Makowski, A.M.: Interpolation approximations for symmetric Fork-Join queues. *Perform. Evaluation* **20**(1–3), 245–265 (1994)
41. Vianna, E., Comarela, G., Pontes, T., Almeida, J., Almeida, V., Wilkinson, K., Kuno, H., Dayal, U.: Analytical performance models for MapReduce workloads. *Int. J. Parallel Prog.* **41**(4), 495–525 (2013)
42. White, T.: *Hadoop: The Definitive Guide*, 1st edn. O’Reilly Media, Inc. (2009)

43. Xia, Y., Tse, D.: On the large deviation of resequencing queue size: 2-M/M/1 case. IEEE Trans. Inf. Theory **54**(9), 4107–4118 (2008)
44. Zaharia, M., Konwinski, A., Joseph, A.D., Katz, R., Stoica, I.: Improving MapReduce performance in heterogeneous environments. In: Proc. of USENIX OSDI, pp. 29–42 (2008)

## Appendix

We assume throughout the paper that all probabilistic objects are defined on a common filtered probability space  $(\Omega, \mathcal{A}, (\mathcal{F}_n)_n, \mathbb{P})$ . All processes  $(X_n)_n$  are assumed to be *adapted*, i.e., for each  $n \geq 0$ , the random variable  $X_n$  is  $\mathcal{F}_n$ -measurable.

**Definition 1** (MARTINGALE) An integrable process  $(X_n)_n$  is a *martingale* if and only if for each  $n \geq 1$

$$\mathbb{E}[X_n | \mathcal{F}_{n-1}] = X_{n-1} . \quad (44)$$

Further,  $X$  is said to be a sub-(super-)martingale if in (44) we have  $\geq$  ( $\leq$ ) instead of equality.

The key property of (sub, super)-martingales that we use in this paper is described by the following lemma:

**Lemma 1** (OPTIONAL SAMPLING THEOREM) Let  $(X_n)_n$  be a martingale, and  $K$  a bounded stopping time, i.e.,  $K \leq n$  a.s. for some  $n \geq 0$  and  $\{K = k\} \in \mathcal{F}_k$  for all  $k \leq n$ . Then

$$\mathbb{E}[X_0] = \mathbb{E}[X_K] = \mathbb{E}[X_n] . \quad (45)$$

If  $X$  is a sub-(super)-martingale, the equality sign in (45) is replaced by  $\leq$  ( $\geq$ ).

*Proof* See, e.g., [7].

Note that for *any* (possibly unbounded) stopping time  $K$ , the stopping time  $K \wedge n$  is always bounded. We use Lemma 1 with the stopping times  $K \wedge n$  in the proofs of Theorems 1 – 4.

**Lemma 2** Let  $c_k$  be the Markov chain from Figure 4 and  $K$  be the stopping time from (11). Then the distribution of  $(c_K | K < \infty)$  is stochastically smaller than the steady-state distribution of  $c_k$ , i.e.,

$$\mathbb{P}[c_K = 2 | K < \infty] \leq \mathbb{P}[c_1 = 2] ,$$

or, equivalently,

$$\mathbb{E}[h(c_K) | K < \infty] \geq \mathbb{E}[h(c_k)] ,$$

for all monotonically decreasing functions  $h$  on  $\{1, 2\}$ .

*Proof* Using Bayes' rule and the stationarity of the process  $c_k$ , it holds:

$$\begin{aligned} \mathbb{P}[c_K = 2 | K < \infty] &= \sum_{k=1}^{\infty} \mathbb{P}[c_k = 2 | K = k] \mathbb{P}[K = k] \\ &= \sum_{k=1}^{\infty} \mathbb{P}[K = k | c_k = 2] \mathbb{P}[c_k = 2] \\ &= \mathbb{P}[c_1 = 2] \sum_{k=1}^{\infty} \mathbb{P}[K = k | c_k = 2] . \end{aligned}$$



Since  $L_1$  is stochastically smaller than  $L_2$ , we have for any  $k \geq 1$

$$\begin{aligned}
 & \mathbb{P}[K = k \mid c_k = 2] \\
 &= \mathbb{P} \left[ t_k \leq \max_n \sum_{i=1}^k x_{n,i} - \sum_{i=1}^{k-1} t_i - \sigma, \max_n \sum_{i=1}^{k-1} (x_{n,i} - t_i) < \sigma \mid c_k = 2 \right] \\
 &\leq \mathbb{P} \left[ t_k \leq \max_n \sum_{i=1}^k x_{n,i} - \sum_{i=1}^{k-1} t_i - \sigma, \max_n \sum_{i=1}^{k-1} (x_{n,i} - t_i) < \sigma \right] \\
 &= \mathbb{P}[K = k] .
 \end{aligned}$$

Hence  $\sum_{k=1}^{\infty} \mathbb{P}[K = k \mid c_k = 2] \leq 1$ , which completes the proof.