



HAL
open science

A note on the accuracy of several existing approximations for M/Ph/m queues

Thomas Begin, Alexandre Brandwajn

► **To cite this version:**

Thomas Begin, Alexandre Brandwajn. A note on the accuracy of several existing approximations for M/Ph/m queues. IEEE HSNCE 2013, Jul 2013, Kyoto, Japan. pp.5. hal-00864323v2

HAL Id: hal-00864323

<https://hal.science/hal-00864323v2>

Submitted on 3 Apr 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A note on the accuracy of several existing approximations for $M/Ph/m$ queues

Thomas Begin* and Alexandre Brandwajn†

*Inria, ENS Lyon, UCB Lyon 1, UMR 5668, 46 Allée d’Italie, 69007 Lyon, France

†Baskin School of Engineering, University of California Santa Cruz Santa Cruz, US

Email: thomas.begin@ens-lyon.fr, alexb@soe.ucsc.edu

Abstract—High variability of system parameters is a complicating factor in the modeling of the performance of big data systems. In this paper, we assess the potential inaccuracy of several existing approximations for evaluating the mean number of jobs queued in a parallelized device that can be represented as an $M/Ph/m$ queue. Unlike existing studies, we consider the effect of the third moment of the service time, or equivalently, its skewness.

We show that the approximations accuracy can be poor even for “easy” examples with a low coefficient of variation of the service time. Our examples demonstrate the important influence of the skewness of the service time distribution on the accuracy of the approximations. None of the approximations accounts for this property. We provide recommendations for the choice of the approximation that allow the user to choose the best suited approximation based on the actual queue parameters.

Index Terms— $M/Ph/m$ queue; Phase type distribution; Skewness; Approximate solution; Relative error.

I. INTRODUCTION

With the rapid emergence of cloud computing, the proliferation of social networking sites and the continuous growth of storage systems, practitioners are increasingly facing the challenge of dealing with extremely large datasets, commonly referred to as *big data*. The processing of these datasets raises many new issues related to their capture, storage, search, sharing, analysis, visualization and performance evaluation. For this latter purpose, the process of modeling systems, which include big data, is generally seen as complex due to, among other things, the high variability of some system parameters. For instance, in the case of a high-performance disk controller, the time needed to process an I/O request is likely to exhibit a skewed distribution because of the underlying multi-tier storage architecture.

Thus, when modeling a system, practitioners are often inclined to replace complex variables with more compact but approximate descriptions. By doing so, they tend to simplify the model and hopefully its solution. However, the resulting cost in terms of (in)accuracy is often unclear. In this paper, we assess the potential inaccuracy of several well-established approximations for evaluating the performance of a parallelized device. We focus our study on the case of highly variable workloads as it is likely to be the case in the context of big data.

The remainder of the paper is organized as follows. In Section II we present the simplest model to represent the behavior of a parallelized system. We briefly review existing

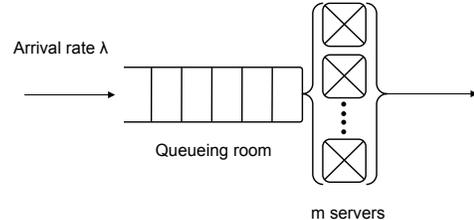


Fig. 1. Illustration of an $M/Ph/m$ queue.

approximations to evaluate its steady-state performance, and we detail our methodology to assess the applicability domain of these approximations. Section III presents the numerical results and discussions. Section IV concludes the paper.

II. MODEL OF PARALLELIZED DEVICE AND ITS SOLUTION

The simplest yet meaningful model for a parallelized device with a skewed service time distribution is the so-called $M/Ph/m$ queue. This queue represents a system with m homogeneous servers, where the time between arrivals are exponentially distributed and job service times have a phase type distribution¹. The buffer (i.e., queueing room) size is assumed to be infinite, so there is no limit on the number of jobs it can contain. The arrivals occur at a rate λ . The service time distribution has a mean of $1/\mu$ and a coefficient of variation of c_B . Recall that the coefficient of variation is defined as the ratio of the standard deviation to the mean, and can be seen as a normalized second order moment. The servers utilization is given by $\rho = \lambda/(m.\mu)$. Of course, if $\rho < 1$, i.e., $\lambda < m.\mu$, the model is stable and the number of jobs in the system has a stationary distribution. Figure 1 depicts the $M/Ph/m$ queue considered while Table I summarizes the principal notation used in our paper.

In the general case, no simple derivation exists to compute the exact stationary queue length distribution of an $M/Ph/m$ queue, even for its first moment. However, several approximate solutions have been proposed in the literature. Several of these approximations presents an easily implementable closed-form formula to evaluate the mean queueing time for a job².

¹A phase type distribution can approximate arbitrarily closely any positive valued distribution [9].

²By Little’s law [1], the mean queueing time $\bar{W}_{M/Ph/m}$ and the mean queue length $\bar{L}_{M/Ph/m}$ are related as follows: $\bar{L}_{M/Ph/m} = \lambda.\bar{W}_{M/Ph/m}$.

| | |
|--------------------------------|---|
| m | Number of servers |
| λ | Rate of jobs arrivals |
| $1/\mu$ | Mean service time for a job |
| σ | Standard deviation of the service time distribution |
| c_B | Coefficient of variation of the service time distribution |
| s_B | Skewness of the service time distribution |
| $\rho = \lambda/(m \cdot \mu)$ | Server utilization |
| $\bar{W}_{M/Ph/m}$ | Mean queueing time of a job in the $M/Ph/m$ queue |
| $\bar{L}_{M/Ph/m}$ | Mean queue length in the $M/Ph/m$ queue |

TABLE I
PRINCIPAL NOTATION USED IN THIS PAPER.

Possibly inspired by the distributional dependence factor in the Pollaczek-Khinchine formula for the $M/GI/1$ queue³, most existing studies seem to concentrate on the influence of the coefficient of variation of the service time distribution. However, a few papers [17], [16], and more recently [8], [6], report that even the common performance parameters like mean queueing delay depend, sometimes to a large extent, on more than the first two moments of the service time distribution. Hence, it would seem unrealistic to expect an approximation based on only two moments to perform well over a large spectrum of service time distribution. Most of the existing approximations base their solutions on the interpolation and extrapolation of two special cases of the $M/Ph/m$ queue, namely the $M/D/m$ and the $M/M/m$ queues, for which the mean queue length can be easily computed (exactly or approximately). We now present the detailed formula for six of these approximations, often found in textbooks on performance evaluation, to compute the mean queueing time $\bar{W}_{M/Ph/m}$.

- 1) Martin's approximation [13].

$$\bar{W}_{M/Ph/m} \simeq \frac{P_m/\mu}{1-\rho} \cdot \frac{1+c_B^2}{2m}$$

where P_m is the waiting probability for an arriving job estimated by the corresponding probability in an $M/M/m$ queue.

- 2) Cosmetatos's approximation [7] (also proposed by Björklund and Elldin [2]).

$$\bar{W}_{M/Ph/m} \simeq c_B^2 \bar{W}_{M/M/m} + (1-c_B^2) \bar{W}_{M/D/m}$$

- 3) Boxma, Cohen and Huffels referred to as BCH's approximation [5].

$$\bar{W}_{M/Ph/m} \simeq \frac{1+c_B^2}{2} \frac{2\bar{W}_{M/D/m}\bar{W}_{M/M/m}}{2a\bar{W}_{M/D/m} + (1-a)\bar{W}_{M/M/m}}$$

where

$$a = \begin{cases} 1 & \text{if } m = 1 \\ \frac{1}{m-1} (c_B^2 + 1 - m + 1) & \text{if } m > 1 \end{cases}$$

³The well-known Pollaczek-Khinchine formula states that the average queueing time in an $M/GI/1$ queue depends only on the first two moments of the service time.

and

$$\gamma_1 = \frac{1-c_B^2}{m+1} + \frac{c_B^2}{m}$$

- 4) Tijms's approximation [15].

$$\bar{W}_{M/Ph/m} \simeq \left((1-\rho)\gamma_1 m + \frac{\rho}{2}(c_B^2 + 1) \right) \bar{W}_{M/M/m}$$

where γ_1 is defined as in the BCH approximation.

- 5) Lee's approximation [12].

$$\bar{W}_{M/Ph/m} \simeq \frac{1+c_B^2}{2} \bar{W}_{M/M/m}$$

- 6) Kimura's approximation [10]

$$\bar{W}_{M/Ph/m} \simeq \frac{1+c_B^2}{\frac{2c_B^2}{\bar{W}_{M/M/m}} + \frac{1-c_B^2}{\bar{W}_{M/D/m}}}$$

Following guidelines of [4], we used in the above approximation

$$\bar{W}_{M/D/m} \simeq \frac{1}{2} \bar{W}_{M/M/m} \left(1 + (1-\rho)(m-1) \frac{\sqrt{4+5m}-2}{16\rho m} \right).$$

We implemented other options for $\bar{W}_{M/D/m}$ (e.g. see eq. (4.17) and (4.29) in [11]) but they produce virtually identical results in our study. For more details on the rationale behind Martin, Cosmetatos, BCH and Tijms approximations, please refer to [4], and to [11] for the others, though their respective domain of applicability is often left unclear.

Overall, little seems to be known about the actual accuracy of these approximations. We briefly review the related state of the art. Not surprisingly, these approximations tend to be fair when the server utilization is low (say less than 0.5). The reason here is that under such condition there is little queue built up in the system, and its performance are mainly driven by the first moment of the service time distribution. Conversely, if the server utilization is high (close to 1), then an $M/Ph/m$ queue is not an adequate model anymore. In such a congestion regime, one needs a queueing model with a finite buffer to assess the system behavior. Existing literature suggests that these approximations will give excellent results when the variability of the service time distribution is low, which means its coefficient of variation ranges from 0 to not much higher than 1. This is quite expected since, by design, most of these approximations rely on the specific solutions of $M/D/m$ and $M/M/m$ queues. Furthermore, statistical distributions with low variability tend to be alike, and thus their expected queueing times tend to be close to those of $M/D/m$ and $M/M/m$ queues. On the other hand, with increasing values of the coefficient of variation of the service time, it is generally known that the approximations accuracy will gradually decrease, although it is generally not stated at what rate. In [4], the authors compare several approximations introduced above, namely Martin, Cosmetatos, BCH and Tijms, with the results obtained from discrete-events simulation. They consider an $M/Ph/m$ queue with 5 servers, a server utilization rate of 0.7 and values of c_B ranging from 0 to 3.2. They conclude that "all

approximations are good for $c_B < 1.4$, and that for higher values of c_B the approximation due to Cosmetatos is very good and the others are fair.”

In this paper, we aim to better characterize the domain of applicability for each of these six approximations. First, we show that their accuracy can be very poor, including for supposedly “easy” examples with a low coefficient of variation of the service time. Second, our results tend to suggest that one can get significantly improved approximate results for the $M/Ph/m$ queue by picking the right approximation according to the actual value of the third moment of the service time, or equivalently, its skewness.

III. NUMERICAL RESULTS, DISCUSSION AND CONCLUSION

A. Experimental protocol

To assess the accuracy of the approximate solutions, and more specifically the influence of the $M/Ph/m$ queue parameters, we proceed as follows. We consider values for the number of servers m ranging from 2 to 12, and for the server utilization ρ ranging from 0.1 to 0.95. For sake of simplicity, we keep the mean service time constant at $1/\mu = 1$. Taking into account the state-of-the-art, we focus our study on values of the coefficient of variation of the service time c_B significantly larger than 1, i.e., ranging from 2 to 10. Finally, unlike most of the existing literature, we pick values for the skewness parameter of the service time s_B ranging from c_B to 100. Recall that the skewness of a random variable X is defined as follows: $s \triangleq \mathbb{E}[(\frac{X-1/\mu}{\sigma})^3]$ where $1/\mu$ and σ are the mean and the standard deviation, respectively. s is thus the normalized third-order central moment, and it is a measure of the asymmetry of the distribution. The greater the skewness, the longer the tail. Note that a skewness of 0 occurs if the distribution is nearly symmetric. Although in theory negative values are possible for s_B , in the case of non-negative distributions and $c_B > 1$ only positive values of s_B can occur, meaning that the tail of the distribution is always on the right (see Appendix A). We rely on an easy computable algorithm [3] to find a phase type distribution that matches specific values of c_B and s_B . Then, we compute the exact value of the mean number of jobs in the corresponding $M/Ph/m$ queues using an iterative solution, and evaluate the corresponding relative error for each of the six approximate solutions.

We explore the obtained results along 4 parameters: m , ρ , c_B and s_B splitting our analysis in two. First we let m and ρ vary while c_B and s_B are kept constant, and then we do it reverse.

B. Influence of the number of servers and the server utilization

To assess the influence of m and ρ on the approximate solutions accuracy, we consider two distinct examples which we refer to as examples (A) and (B). Their description is given in Table II.

In example (A), the coefficient of variation of the service time is set to 8 while the skewness is equal to 30. Figure 2 shows the degree of accuracy attained by each approximate

| Example | m | ρ | c_B | s_B |
|---------|---------|-------------|---------|--------------|
| (A) | [0; 12] | [0.1; 0.95] | 8 | 30 |
| (B) | [0; 12] | [0.1; 0.95] | 5 | 50 |
| (C) | 4 | 0.7 | [2; 10] | $[c_B; 100]$ |
| (D) | 8 | 0.6 | [2; 10] | $[c_B; 100]$ |

TABLE II
THE $M/Ph/m$ QUEUE PARAMETERS.

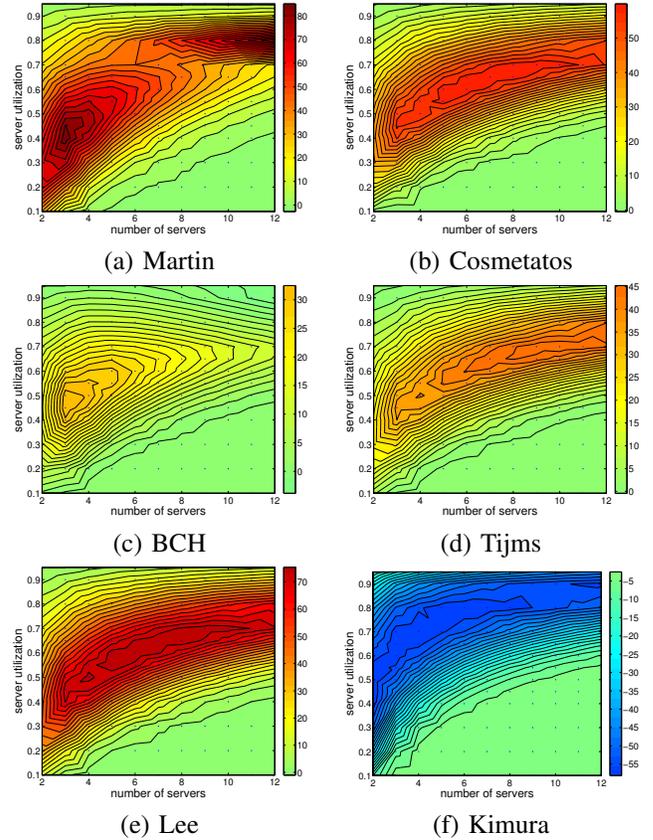


Fig. 2. Relative errors of the approximate solutions on the mean number of jobs in the $M/Ph/m$ queue of example (A).

solution on this example. It turns out that, with the exception of Kimura’s solution, virtually all approximations tend to overestimate the queue length, whatever the values of ρ and m . Deviations from the exact solution tend to be higher when ρ is close to 0.5 for small values of m , whereas they peak around ρ equal to 0.7 for larger values of m . Based on this example, the BCH solution with its relative error often less than 10% and at most of 35% significantly outperforms the others whose errors attain and exceed 50%.

Example (B) deals with a smaller value of c_B but a larger s_B . Figure 3 presents the corresponding results. We observe that, again, all but Kimura’s solution tend to overestimate the mean number of jobs in the $M/Ph/m$ queue. Similarly to example (A), it shows that approximate solutions are the most likely to fail when ρ is close to 0.5 for m ranging from 2 to 4. For larger values of m the region of lowest accuracy gradually shifts to higher values of ρ around 0.8. However, unlike example (A), the solution of Kimura yields the best

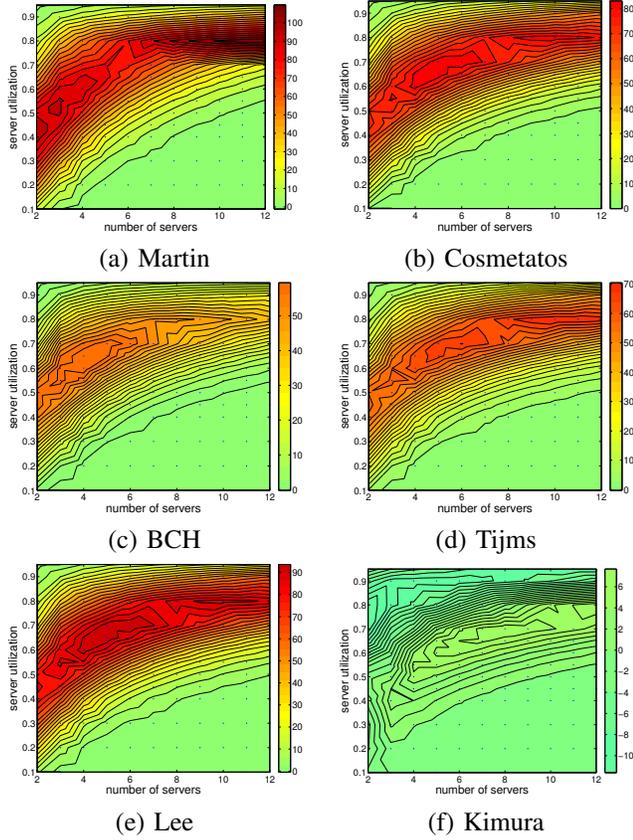


Fig. 3. Relative errors of the approximate solutions on the mean number of jobs in the $M/Ph/m$ queue of example (B).

results with a relative error not exceeding 10%.

Based on examples (A) and (B), it seems that when the number of servers is either much larger than the server utilization (say m verifies $m > 33\rho - 4.66$), or conversely much smaller than the server utilization (say $m < 50\rho - 33$), all approximate solutions lead to excellent results. On the other hand, the region where approximate solutions tend to fail occurs for moderate server utilization. When m is low (say less than 4) this region tends to cover server utilizations close to 0.5. As m grows, this region gradually shifts towards server utilizations around 0.8. It is also worth noting that, in this example, BCH's solution always outperforms Martin, Cosmetatos, Tijms and Lee's solution. Finally, it appears that depending on the actual service time distribution (here its c_B and s_B), BCH or Kimura's solution emerges as the best choice. In the next section, we study in more detail the influence of the service time distribution.

C. Influence of the coefficient of variation and the skewness

In order to investigate the influence of c_B and s_B , we now hold constant the number of servers m and the server utilization ρ .

We consider example (C) where m is 4 and ρ is set to 0.7. The corresponding results are illustrated in Figure 4. All solutions except Kimura's approximation yield fair results even for large values of the coefficient of variation as long

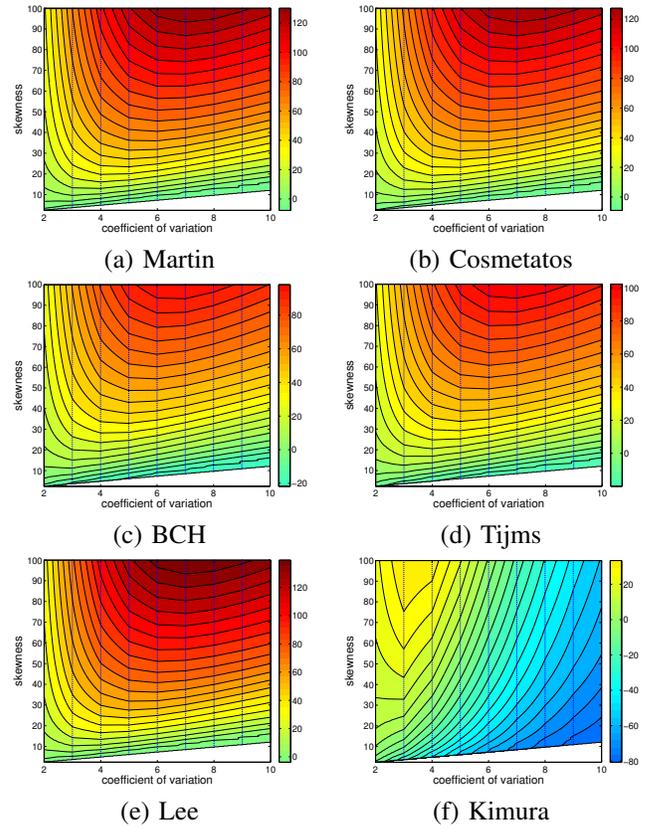


Fig. 4. Relative errors of the approximate solution on the mean number of jobs in the $M/Ph/m$ queue of example (C).

as the skewness is low (say less than 20). As the skewness increases and reaches 40 or 50, then these solutions lead to highly inaccurate results, unless c_B is less than 3. On the other hand, Kimura's solution tends to work best when the other solutions are failing. Except when c_B is high and s_B is kept low, its deviations from the exact solution are below 25%.

In example (D) we increase the number of servers to $m = 8$ and we set the server utilization ρ to 0.6. Figure 5 shows the corresponding results. Overall, all the approximate solutions perform better here than in the previous example. Kimura's solution yields excellent results unless c_B is high and s_B is very low. Conversely, in such a case, all other solutions produce fair results.

From examples of (C) and (D), we observe that if both c_B and s_B are low (say less than 4 and 30, respectively), any of the six approximate solutions provides accurate results. If the skewness is large (say 10 times greater than the coefficient of variation), but the coefficient of variation is low, our results favor the use of Kimura's solution. Conversely, when the skewness is low, and the coefficient of variation is high, then BCH solution becomes the best option. Finally, when both c_B and s_B are high, none of the approximations seem to be acceptable.

To conclude, based on our four examples, we recommend the following guidelines:

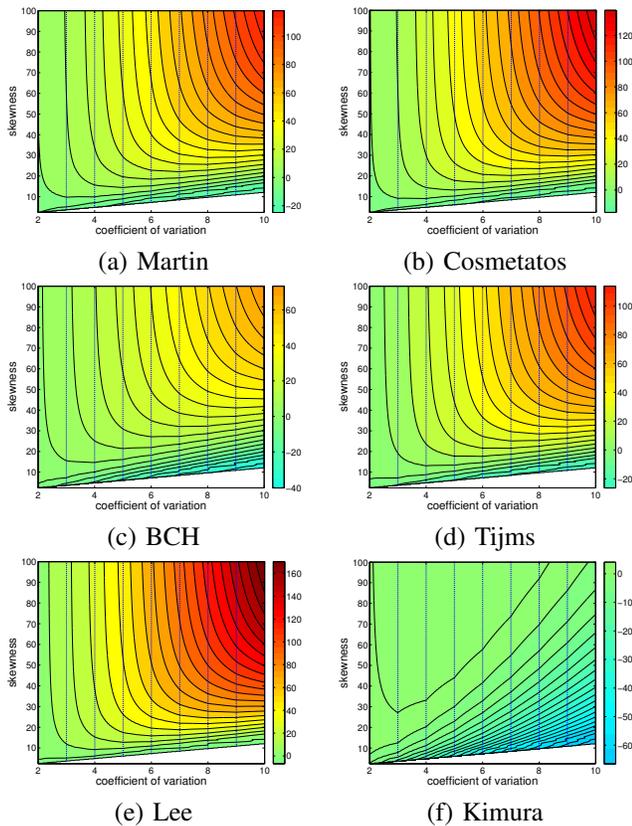


Fig. 5. Relative errors of the approximate solution on the mean number of jobs in the $M/Ph/m$ queue of example (D).

- if m is much larger than ρ (say $m > 33\rho - 4.66$), or conversely much smaller (say $m < 50\rho - 33$), then any approximate solutions yields excellent results;
- otherwise
 - if both c_B and s_B are low (say less than 4 and 30, respectively), any of the six approximate solutions provides accurate results;
 - if s_B is large (say 10 times greater than the coefficient of variation) and c_B is low, we recommend the use of the Kimura solution;
 - if s_B is low and c_B is large, we recommend the use of the BCH solution;
 - if both c_B and s_B are high, none of the approximations seem to be acceptable.

The use of these guidelines would ensure that, in almost all our numerous experiments, the discrepancy between the exact and the approximate mean number of jobs in the queue stays within a range of 25% relative error (except obviously when both c_B and s_B are high and ρ is moderate).

IV. CONCLUSION

The modeling and performance evaluation of systems that process big data is often hampered by difficulties arising due to the high variability of system parameters. In this paper, we assess the potential inaccuracy of several well-established approximations for evaluating the performance of a parallelized

device. We focus our study on the case of highly variable and skewed workloads as it is likely to be the case in the context of big data. We use an exact numerical solution to the $M/Ph/m$ queue and an extensive set of experiments (thousands of cases) to evaluate the difference between approximation results and the exact value of the mean queue length. Unlike existing studies, we explore the (in)accuracy of these approximations not only along the commonly explored coefficients of variation but also through the actual value of the third moment of the service time, or equivalently, its skewness.

Our work helps to better characterize the domain of applicability for each of these approximations. First, we show that their accuracy can be very poor, including for supposedly “easy” examples with a low coefficient of variation of the service time. Second, our examples emphasize the important influence of the skewness of the service time distribution on the approximations accuracy whereas none of them account for it. More precisely, our study shows that for 5 out of the 6 tested approximations the results are quite close to each other results (although some differences do exist) and they tend to overestimate the congestion level of the queue. Conversely, the sixth approximate solution is most likely to underestimate the mean length of the queue. Thus, we provide recommendations regarding the choice of the approximation based on the actual set of values for the queue parameters. Using these guidelines would ensure that, in most cases, the relative error of the approximation would remain below 25% relative error.

However, it is important to stress that our results only pertain to a subset of phase type distributions since they were all obtained using a matching algorithm for the first 3 moments of the service time distributions that minimizes the number of phases for the resulting distribution. Thus, a meaningful extension to this work would consist in studying the accuracy of the approximations for phase type distributions that are no more subject to this constraint. Preliminary results (not shown in this paper) seem to indicate that our recommendations still apply in the general case. Another option would be to extend the results of this study to the case of distributions with negative skewness (and hence coefficient of variation less than 1).

REFERENCES

- [1] Allen, A. O., *Probability, Statistics, and Queueing Theory with Computer Science Applications*. Academic Press, 2nd edition, 1990.
- [2] Björklund, M., and Elldin, A., A practical method of calculation for certain types of complex common control systems, *Ericsson Technics*, Vol. 20, 1964, pp. 3-75.
- [3] Bobbio, A., Horváth, A., and Telek, M., Matching three moments with minimal acyclic phase type distributions, *Stochastic Models*, Vol. 21, 2005, pp. 303-326.
- [4] Bolch, G., Greiner, S., Meer, H., and Trivedi, K., *Queueing Networks and Markov Chains*. Second Edition, Wiley-Interscience, 2005.
- [5] Boxma, O. J., Cohen, J. W., and Huffels, N., Approximations of the Mean Waiting Time in an $M/G/s$ -Queueing System, *Operations Research*, Vol. 27, No. 6., 1979, pp. 1115-1127.
- [6] Brandwajn, A. and Begin, T., Considerations in workload characterization for Parallel Access Volumes, *proceedings of CMG*, 2009.
- [7] Cosmetatos, G., Some Approximate Equilibrium Results for the Multi-server Queue ($M/G/r$), *Operations Research Quarterly*, USA, 1976, pp. 615-620.

- [8] Gupta, V., Harchol-Balter, M., Dai, J. and Zwart, B., The effect of higher moments of job size distribution on the performance of an $M/G/s$ queueing system, *Performance Evaluation Review*, Vol. 35 (2), 2007, pp. 12-14.
- [9] Johnson, M. A., and Taaffe, M. R. The denseness of phase distributions. School of Industrial Engineering, Purdue University. 1988.
- [10] Kimura, T., A two-moment approximation for the mean waiting time in the $GI/G/s$ queue, *Management Science*, Vol. 32, 1986, pp. 751-763.
- [11] Kimura, T., Approximations for multi-server queues: system interpolations, *Queueing Systems*, Vol. 17, 1994, pp. 347-382.
- [12] Lee, A.M. and Longton, P.A., Queueing process associated with airline passenger check-in, *Operations Research Quarterly*, Vol. 10, 1957, pp. 56-71.
- [13] Martin, J., *System Analysis for Data Transmission*. Prentice-Hall, Englewood Cliffs, NJ, 1972.
- [14] Osogami, T. and Harchol-Balter, M., Closed form solutions for mapping general distributions to quasi-minimal PH distributions, *Performance Evaluation*, Vol. 63 (6), 2006, pp. 524-552.
- [15] Tijms, H., *Stochastic Modelling and Analysis: A Computational Approach*. John Wiley, NY, 1986.
- [16] Whitt, W., The effect of variability in the $GI/G/s$ queue, *Journal of Applied Probability*, Vol. 17 (4), 1980, pp. 1062-1071.
- [17] Wolff, R.W., The Effect of Service Time Regularity on System Performance, *Computer Performance*, North Holland, 1977, pp. 297-304.

APPENDIX

A. Relation between coefficient of variation and skewness

Let X be a non-negative random variable. We denote by m_i its i -th non-central moment, i.e., $m_i \triangleq \mathbb{E}(X^i)$, and by n_k its k -th normalized moments. By definition, we have: $n_2 \triangleq m_2/m_1^2$ and $n_3 \triangleq m_3/(m_1 m_2)$. It is known (e.g., [14]) that for any non-negative distribution we have:

$$n_3 \geq n_2. \quad (1)$$

Using the definitions of n_2 and n_3 , (1) can be rewritten as:

$$s_B \geq (c_B - \frac{1}{c_B}), \quad (2)$$

where s_B and c_B represent the skewness and the coefficient of variation, respectively. Recall that $c_B \triangleq \frac{\sqrt{m_2 - m_1^2}}{m_1}$ and that s_B can be expressed as $s_B = \frac{m_3 - 3m_1 m_2 + 2m_1^3}{(m_2 - m_1^2)^{3/2}}$.

From (2), it follows:

$$s_B < 0 \iff c_B < 1 \quad (3)$$

since m_1 and c_B are both positive.

This simple derivation explains (for a non-negative distribution) why the skewness values can only be positive if we limit the coefficient of variation to values larger than 1.