# Research Report

## Performance Evaluation of a Tape Library System

I. Iliadis*, Y. Kim#, S. Sarafijanovic*, V. Venkatesan*

*IBM Research – Zurich
8803 Rüschlikon
Switzerland

#IBM Research - Netherlands

**IBM Research**
**Africa • Almaden • Austin • Australia • Brazil • China • Haifa • India • Ireland • Tokyo • Watson • Zurich**

# Performance Evaluation of a Tape Library System

Ilias Iliadis, Yusik Kim, Slavisa Sarafijanovic, Vinodh Venkatesan
IBM Research
ili@zurich.ibm.com, yusik.kim@nl.ibm.com, {sla,ven}@zurich.ibm.com

*Abstract*—Data with vastly different access characteristics is efficiently stored in multi-tiered storage systems. A cost-effective way to retain large volumes of infrequently accessed data is to store it on tape. Steady developments in tape technology deliver ever increasing storage capacities at low cost. This has established tape as a viable solution to cope with the extreme data growth in the context of Big Data. Assessing the performance of the various tiers is central to achieving appropriate tier dimensioning and storage provisioning. To that end, we develop an analytical model to evaluate the performance of a tape library system that considers various relevant aspects, such as the number of cartridges and tape drives as well as different mount/unmount policies. Closed-form expressions for the corresponding mean waiting times are derived. The validity of the model developed is confirmed by demonstrating that the predicted performance matches well with that obtained by simulation across a wide range of system parameter values.

## I. INTRODUCTION

Large data volumes are ubiquitous in modern enterprises with the demand for additional capacity continuing to grow. The Big Data explosion has led to the deployment of multi-tiered storage systems that make use of mass storage archives. In the past four decades, cost considerations have established tape as the natural choice for economically storing the bulk of data, that is, the data that is infrequently accessed [1]–[7]. The tape tier remains an economically attractive option, especially considering the potential growth of the capacity and bandwidth offered by tape devices [8]. Modern enterprise tape has reached a capacity of 10 TB per cartridge with data rates reaching 360 MB/s. A tape library system typically comprises a small number of tape drives for performing read/write operations on a large number of tape cartridges, which are located in a storage rack and are mounted in and unmounted from the tape drives through automation mechanisms (robot arms). To perform a read/write access to a tape cartridge, the cartridge should first be mounted in a free drive and then, after a *seek time*, be appropriately positioned to read/write the corresponding data. Also, the *unmount time* of a tape includes the time spent to rewind it before it is removed from the corresponding tape drive. Although the tape tier is economically attractive, it presents certain challenges as the access latencies are relatively high and can run into minutes even on lightly loaded systems. Furthermore, tape drives are very expensive compared to tape cartridges and consume more space. Therefore, the tape drive component is a scarce resource that needs to be properly provisioned for an efficient and cost-effective operation. Consequently, assigning a proper number of tape drives to a tape system is becoming increasingly important. The same applies to optical disk archives, which operate using optical disk readers and, generally, to tertiary

storage, which uses drives that accept removable media, such as tapes and optical disks [5], [6].

The complexity of storage systems has grown in terms of the heterogeneity required to satisfy user requirements. In this context, as well as in the context of cloud storage, multi-tiered storage systems store data with vastly different access characteristics in various types of storage devices. They typically store the least accessed data on tape, also referred to as tertiary storage [5], [6], and the most frequently accessed data on hard disk drives (HDD) and solid state drives (SSD). A scalable method for efficient tier provisioning and data placement that minimizes the system's mean response time for a given budget and workload was presented in [9]. It used an $M/G/1$ queueing model to capture the behavior of HDD and SSD devices and analytically assess the mean waiting times of the corresponding tiers. Note though, that the operation of tape libraries with the mounting, serving, and unmounting of cartridges is similar to that of polling systems, and therefore cannot be captured by a simple $M/G/1$ queueing model. Consequently, for a proper provisioning of a multi-tiered storage system that includes tape, an advanced tape tier model needs to be developed, which is the aim of this work.

In this article, we present an analytical queueing-based method for efficiently evaluating the performance of a tape library system and, more generally, of a storage system that operates using drives that accept removable media. The results obtained can then be used to provision the system such that desired performance guarantees are satisfied. Such performance guarantees are useful in appropriately provisioning tape systems for an expected workload in order to comply with service level agreements on the system performance. The model developed yields the mean waiting time analytically in closed-form as a function of the system parameters, including the number of tape cartridges, the number of tape drives, and the tape mount, unmount and seek times.

Performance evaluations of a tape library system reported in the literature were mainly conducted by means of simulation [1]–[3] or by actual measurements [6]. To properly dimension such a system, it is imperative to be able to assess the effect of the various parameters on the performance metric considered. Accomplishing this by simulation is time-consuming compared with analytical models that provide fast execution times and more insight. Initial analytical efforts were presented in [10], [11] where $M/M/c$ and $M^X/G/c$ models were used. These models, however, are simplistic because they do not capture the polling nature of operation of tape library systems. In this article, we present a comprehensive theoretical model that considers the resource contention between tapes and

drives, and provides accurate results, which cannot be directly derived from relevant previous works. Building upon $M/G/K$ and polling-system queueing results, we develop an enhanced model and subsequently assess the effect of two distinct tape mount/unmount policies. To the best of our knowledge, this is the first theoretical work to accurately evaluate the performance of tape libraries and, more generally, of systems with removable media, using closed-form expressions for the mean waiting time.

The remainder of the paper is organized as follows. Section II provides a survey of the relevant literature on performance evaluation of tape library systems. Sections III and IV describe the operation of a tape library system and its similarities to the operation of a polling system. Section V describes the tape system model along with the corresponding parameters. Section VI presents the analytical evaluation of the mean waiting time. Section VII shows numerical results for a typical tape system. Finally, we conclude in Section VIII.

## II. Related Work

The performance of tape libraries was studied primarily by means of simulation. The maximum throughput and the mean response time of an automated tape library were evaluated in [1]. An open queueing network model was developed, but it turned out to be analytically intractable. A performance evaluation was subsequently conducted by means of simulation. The effect of striping in large tape libraries was assessed in [2], [3] by simulating a closed system, where tape drives are always busy serving requests, and an open system, respectively. The effect of the various parameters on the maximum throughput and on the mean response time was subsequently evaluated.

To the best of our knowledge, there are only two initial efforts to assess the performance of a tape library analytically, both of which date back twenty years [10], [11]. The analysis in [10] considered an $M/M/c$ model applied together with an empirical expression, whereas that in [11] considered an $M^X/G/c$ model. As we will see in Section VII, the actual performance curve cannot be efficiently approximated using these models and this is due to the fact that none of them captures the polling nature of operation of a tape library system. At first glance, and given that there is a large body of queueing theoretical work on the performance evaluation of polling systems, it seems surprising that in the last two decades there has been no other analytical work on tape libraries. We proceed to review relevant prior work on polling systems and demonstrate that the corresponding results cannot be directly applied to analyze a tape library system. This could therefore be a reason for the lack of progress in this area.

## III. Operation of a Tape Library System

A tape library system consists of tape drives, automation mechanisms (robot arms), a storage rack for the tape cartridges, and a cartridge control unit. To serve a request, a robot arm (picker) fetches the appropriate tape cartridge from the storage rack and delivers it to a free tape drive. The tape drive unit mounts the tape, positions the head to the desired file and transfers the data. To free the tape drive, a robot arm unmounts the tape cartridge and returns it to the storage rack.

Tape read/write requests contain the following information: the cartridge id, the position in the cartridge where the corresponding data blocks reside, and the data size to be transferred. Requests submitted for cartridges are queued in the corresponding queues and are subsequently served according to a scheduling policy. In this article, we consider a hierarchical scheduling algorithm where at the upper level it employs a cyclic (round-robin) scheduling among the queues (mounting cartridges), and at the lower level it employs a first-come-first-served (FCFS) policy for serving requests within a queue (reading from a mounted cartridge). This algorithm ensures fairness and avoids starvation. The system's performance can be further improved by appropriately scheduling batches of random requests [12], [13]. This issue, however, is beyond the scope of this article. When all requests for a cartridge are served (exhaustive service), and there are still pending requests to some other, non-mounted cartridge, the cartridge is unmounted and another cartridge with pending requests is subsequently mounted. If, however, there are *no other pending requests to any other non-mounted cartridge*, the cartridge can either remain mounted in anticipation of future requests arriving for it or be unmounted so as to save its corresponding unmount time when future requests arrive for other non-mounted cartridges that subsequently need to be mounted. We proceed by considering the following two mount/unmount policies deployed in this context:

1) *Always-Unmount (AU) policy*: a tape cartridge is immediately unmounted upon completion of all pending requests for it, in anticipation of the next request arriving for another non-mounted cartridge.
2) *Not-Unmount (NU) policy*: a tape cartridge remains mounted upon completion of all pending requests for it, in anticipation of the next request arriving for this same currently mounted, but idle cartridge.

## IV. Tape Library and Polling System Operation

In typical polling systems, which include computer communication, production, traffic and transportation systems [14], a server serves multiple queues in a given order. In such systems, it is important to optimize the strategy for serving the queues so as to minimize the average waiting and response time. Several policies were considered for serving jobs in a queue, such as exhaustive, gated, and limited. Under the exhaustive service discipline, a server continues to work serving a queue until it becomes empty. The order in which a server visits the queues is determined by a routing scheme, such as cyclic (round-robin), random, or first-come-first-served. Most of the work on polling systems assumes that the server continues polling successive queues even when the system is empty [15]. The transitions between successive queues may be instantaneous or may require a *switchover time*. If a polled queue is not empty, the server spends an additional *setup time* before beginning its service. This is a *state-dependent* setup time operation as opposed to a system operating under a *state-independent* setup time that is incurred in all queues visited, regardless of whether there is work to be done or the queues are empty.

We proceed by noting that the operation of a tape library bears similarities to that of a queueing polling system. Requests arriving for tape cartridges correspond to jobs arriving at queues. The mounting of cartridges to tape drives to serve requests corresponds to the servers visiting queues to serve jobs. In the remainder, we consider a cyclic mounting (visiting) policy. We now examine whether the cartridge mount/unmount times can be mapped to switchover and setup times in the context of polling models. In particular, when all requests of a cartridge are served while there are still pending requests in the system, the cartridge is unmounted and the next cartridge in the sequence that has pending requests is mounted. The time elapsed before starting to serve a new request is equal to the sum of the unmount and mount times, regardless of the position of the cartridge within the polling sequence. Thus, a moments' reflection reveals that the corresponding "switchover time" has to be equal to zero, for otherwise the time elapsed would depend on the cartridge's position. Furthermore, for the NU policy, the corresponding "setup time" is equal to the sum of the unmount and mount times, and is state-dependent as it is incurred only on non-empty queues. For the AU policy, however, the state-dependent setup time involves the unmount and mount times when the system is busy, but only the mount time when the system is idle. More details on further similarities and differences are provided next.

Polling systems have been extensively studied [16]. The vast majority of work has dealt with the case of a single server, and even in this case, the resulting models tend to be very complex with few explicit results. For instance, relevant to the NU policy is the work presented in [17] where a *patient server* halts at a queue as opposed to continuing to visit empty queues. This work considers state-dependent setup times and assesses system performance using an iterative procedure that involves a discrete-Fourier-transform numerical technique. Obtaining explicit closed-form expressions for the mean waiting time seems to be a daunting task [15]. The additional complexity of multiple (two or more) servers renders the polling models intractable [18]. In particular, there are no analytical results even for mean waiting times, and hence one must resort to either simulations or approximations.

We have shown that the operation of a tape library corresponds to that of a multi-server state-dependent polling system. An approximate analysis of a state-dependent polling system with multiple servers was presented in [19], but this work does not consider the exhaustive service discipline; it only provides results for the cases where a server visiting a queue serves either all present jobs (gated service) or a limited number of jobs (limited service). The exhaustive and gated service disciplines were considered in [20], but results were obtained under the assumption that multiple, and possibly all servers may simultaneously serve a queue. In the case of a tape library, this assumption does not hold as there can be at most one tape served by a drive, which implies that only one server may serve a queue at any given time. Also, for this reason, prior work on polling systems with multiple coupled servers [21], where all servers simultaneously visit a queue, is not applicable.

In this article, we present a model which is suitable for the

| Parameter | Definition |
|---|---|
| $c$ | number of cartridges |
| $d$ | number of tape drives |
| $a$ | number of arms |
| $M$ | mount time |
| $U$ | unmount time |
| $s$ | seek time |
| $\lambda_{ct}$ | arrival rate of requests for a cartridge |
| $Q$ | request size |
| $b_w$ | bandwidth |
| $\lambda$ | arrival rate of requests to the tape system ($\lambda = c\,\lambda_{ct}$) |
| $B$ | service time of a request ($B = s + Q/b_w$) |
| $\rho$ | system load ($\rho = \lambda\,\overline{B}$) |
| $n$ | number of cartridges per tape drive ($n = c/d$) |

analysis of the AU and NU policies in a common framework. Our approach follows the same direction as the one presented in [22], and the references therein, that efficiently obtain an expression for the expected delay by combining known expressions for the light and heavy traffic cases through interpolation. However, for the case of a multi-drive system, there are no known expressions for the light and heavy traffic cases. In our article, the mean waiting time of a multi-drive system in the heavy-load region is obtained through a non-trivial extension of previously published work for single-drive systems [23], whereas the one for the light-load region is obtained in a novel way through a virtual $M/G/K$ queueing model. Subsequently, we present a method to perform an interpolation between the light- and heavy-load regions.

## V. TAPE SYSTEM MODEL

The notation used is summarized in Table I. The parameters are divided according to whether they are independent or derived and are listed in the upper and the lower part of the table, respectively. The tape library system considered comprises $c$ cartridges, $d$ tape drives, and $a$ (typically one or two) robot arms. Note that the contention for the robot arm(s) can be neglected because at low and high loads it is expected to be negligible, and also because nowadays the time required by robot arms to fetch tapes is significantly smaller than the tape mount and unmount times, denoted by $M$ and $U$, respectively [24]. This assertion is indeed confirmed in Section VII. We therefore proceed by considering $a = d$, which implies that there is no contention for the robot arm(s). The workload is assumed to be symmetric with the requests assumed to arrive for the $c$ cartridges according to independent and identical Poisson processes at a rate of $\lambda_{ct}$ requests per unit of time. This implies that the arrival process of requests to the tape system is Poisson with rate $\lambda$, where $\lambda = c\,\lambda_{ct}$. The request size is denoted by $Q$, with the sizes of requests assumed to be independent and identically distributed. Each request incurs a seek time of $s$ and a transfer time of $Q/b_w$, where $b_w$ denotes the transfer bandwidth. Thus, the total time to serve a request, $B$, is equal to $s + Q/b_w$. The first and second moments of a random variable $X$ are denoted by $\overline{X}$ and $\overline{X^2}$, respectively, such that $\overline{B}$ and $\overline{B^2}$ denote the first and second moments of the service time $B$, respectively, given by

$$\overline{B} = \overline{s} + \overline{Q}/b_w \;, \tag{1}$$

and
$$\overline{B^2} = \overline{s^2} + \overline{Q^2}/b_w^2 + 2\,\overline{s}\,\overline{Q}/b_w \ . \qquad (2)$$

Owing to the assumptions, it follows that the service times are independent and identically distributed random variables.

## VI. System Analysis

The performance of the system depends on its load, which takes values in the [0,1) interval. In the following sections we will show that the performance of the system can be accurately modeled in the light-load and heavy-load regions. Subsequently, we will present a method to accurately assess the performance in the medium-load region based on an interpolation between the light-load and heavy-load results.

### A. Light-Load Analysis

Here we assess the mean waiting time in the light-load region, that is, when $\lambda$ is relatively small. A precise determination of where this region ends will be given in Section VI-C. We proceed by noting that when the load is light, most likely there is at most one request pending in each queue. Thus, every time a request of a cartridge is served, there is no other request pending for the same cartridge, and therefore a robot arm unmounts the cartridge and mounts another one to serve its request. In this context, the set of all outstanding requests form a virtual queue that is served by the $d$ tape drives. Furthermore, the behavior of this virtual queue can be analyzed by considering a fictitious service time of each request, $S_f$, consisting of three components: the unmount time $U$, the mount time $M$, and the service time $B$, that is,

$$S_f = U + M + B \ . \qquad (3)$$

In fact, the only difference between the AU and NU policies is that the sequence of appearance of these times varies; the AU policy has a mount-serve-ummount sequence whereas the NU policy has an unmount-mount-serve sequence. Consequently, from the virtual queue's perspective, the two policies have a service time, $S_f$, distributed identically with its first two moments given by

$$\overline{S_f} = \overline{U} + \overline{M} + \overline{B} \ , \qquad (4)$$

and
$$\overline{S_f^2} = \overline{U^2} + \overline{M^2} + \overline{B^2} + 2\,(\overline{U}\,\overline{M} + \overline{U}\,\overline{B} + \overline{M}\,\overline{B}) \ . \quad (5)$$

Consequently, the operation of the virtual queue can be captured by an $M/G/K$ queueing model with $K(=d)$ servers, where the arrival rate is $\lambda$ and the service time is $S_f$. The mean waiting of an $M/G/K$ queue is obtained approximately by $\mathrm{E}[W^{M/G/K}] = \frac{1}{2}\,(1 + C^2)\,\mathrm{E}[W^{M/M/K}]$ where $\mathrm{E}[W^{M/M/K}]$ is the mean waiting time of a corresponding $M/M/K$ queue with the same mean service time $\overline{S_f}$, and $C$ is the coefficient of variation of $S_f$ [25]. The fictitious mean waiting time, $\overline{W_f}$, can then be approximated by

$$\overline{W_f}(\rho_f) \approx \frac{\overline{S_f^2}\,d^{d-1}\,\rho_f^d}{2\,\overline{S_f}\,d!\,(1-\rho_f)^2}\,\pi_0 \ , \qquad (6)$$

where $\rho_f$ denotes the load of the virtual queue, given by

$$\rho_f = \frac{\lambda\,\overline{S_f}}{d} \ , \qquad (7)$$

and $\pi_0$ is the probability that the system is empty, given by

$$\pi_0 = \frac{1}{\frac{(d\,\rho_f)^d}{d!\,(1-\rho_f)} + \sum_{n=0}^{d}\frac{(d\,\rho_f)^n}{n!}} \ . \qquad (8)$$

It holds that the utilization (load) of the system is given by

$$\rho = \frac{\lambda\,\overline{B}}{d} \ . \qquad (9)$$

Combining (7) and (9) yields

$$\rho_f = \frac{\rho}{\rho^*} \ , \qquad (10)$$

where

$$\rho^* \triangleq \frac{\overline{B}}{\overline{S_f}} \stackrel{(4)}{=} \frac{\overline{B}}{\overline{U} + \overline{M} + \overline{B}} \ . \qquad (11)$$

The virtual queue saturates when $\rho_f = 1$, which, according to (10), occurs at load $\rho = \rho^*$. Note also that for a single-server queuing system ($K = d = 1$), approximation (6) is exact, yielding the Pollaczek–Khinchine formula.

As mentioned above, the AU and NU policies have the same fictitious service time distribution and, consequently, the same queueing time distribution. However, as the actual waiting times of requests are defined between the arrival and service initiation times, the actual waiting times are longer than the fictitious ones. In particular, parts of the fictitious service times need to be counted as actual waiting times. For the AU policy, the first component of the mount-serve-ummount fictitious service sequence, that is, the mount time, is in fact part of the actual waiting time. Similarly, for the NU policy, the first two components of the unmount-mount-serve fictitious service sequence, that is, the unmount and mount times, are in fact part of the actual waiting time. The actual waiting times, $W_{\mathrm{AU,l}}$ and $W_{\mathrm{NU,l}}$, corresponding to the *Always-Unmount* and *Not-Unmount* policies are therefore given by,

$$W_{\mathrm{AU,l}} = W_f + M \ , \qquad (12)$$

$$W_{\mathrm{NU,l}} = W_f + U + M \ , \qquad (13)$$

which in turn yields the mean waiting times as follows:

$$\overline{W_{\mathrm{AU,l}}} = \overline{W_f} + \overline{M} \ , \qquad (14)$$

$$\overline{W_{\mathrm{NU,l}}} = \overline{W_f} + \overline{U} + \overline{M} \ . \qquad (15)$$

Note that in the case of the NU policy, when the number of cartridges $c$ is not very large, there is a non-negligible possibility that a request may arrive for a cartridge that is already mounted. As for the NU policy there are always $d$ cartridges mounted, the probability of this event is equal to $d/c$. In this case, the service of this request can immediately start, which implies that the fictitious service of an arbitrary request is given by

$$S_f = \begin{cases} B, & \text{with prob. } \frac{d}{c} \\ B + U + M, & \text{with prob. } 1 - \frac{d}{c} \end{cases} , \quad \text{for NU. (16)}$$

Also, the actual waiting time, $W_{\mathrm{NU,l}}$, is now given by

$$W_{\mathrm{NU,l}} = \begin{cases} W_f, & \text{with prob. } \frac{d}{c} \\ W_f + U + M, & \text{with prob. } 1 - \frac{d}{c} \end{cases}, \qquad (17)$$

4

which in turn yields the mean waiting time $\overline{W_{\text{NU,l}}}$ as follows:

$$\overline{W_{\text{NU,l}}} = \overline{W_f} + \left(1 - \frac{d}{c}\right)(\overline{U} + \overline{M}).  \quad (18)$$

Note that (18) is also valid when $c$ is large because it reduces to (15). It can therefore be used for all values of $c$. Summarizing (14) and (18), and using (10), yields

$$\overline{W_l}(\rho) = \overline{W_f}(\rho/\rho^*) + H,  \quad (19)$$

where $\overline{W_f}(.)$ is given by (6), $\rho^*$ by (11), and $H$ by

$$H \triangleq \begin{cases} \overline{M}, & \text{for AU} \\ \left(1 - \frac{d}{c}\right)(\overline{U} + \overline{M}), & \text{for NU}. \end{cases}  \quad (20)$$

### B. Heavy-Load Analysis

We proceed by noting that when the load is high, the tape drives are busy serving requests most of the time, while newly arriving requests involve a large number of cartridges. Consequently, when all requests for a given cartridge are served, it is very unlikely that there are no other pending requests in the system, which in turn implies that the difference between the AU and NU policies will be negligible. It is also expected that there are requests to be served in all the queues visited, which implies that the setup times turn out to be state-independent, and therefore the performance could then be assessed by applying a state-independent polling model. Such models, though, were developed for the case of a single server, not for multiple servers. To cope with this issue, we consider the system as being partitioned in $d$ domains, with each domain containing a set of $n = c/d$ cartridges that are served by a single tape drive. Therefore, the arrival process of requests to each of the domains is Poisson with rate $\lambda_{\text{dm}}$, where $\lambda_{\text{dm}} = n\lambda_{\text{ct}} = \lambda/d$. Subsequently, we can apply the result obtained in [23]. In this case it holds that $N = n$. Also, for $j = 1, \ldots, n$, it holds that $\lambda_j = \lambda_{\text{ct}}$, $\rho_j = \lambda_{\text{ct}}\overline{B}$, (which implies that $\rho = \sum_{j=1}^n \rho_j = n\lambda_{\text{ct}}\overline{B} = \lambda_{\text{dm}}\overline{B} = \lambda\overline{B}/d$), $\overline{h_j^2} = \overline{B^2}$, $s_j = U_j + M_j$, with the variables $U_j$ and $M_j$ distributed according to $U$ and $M$, respectively, and $S = \sum_{j=1}^n s_j$. By substituting the preceding expressions into Eq. (4) of [23], we obtain the mean waiting time, $\overline{W_h}$, as follows:

$$\overline{W_h} = \frac{\lambda\overline{B^2}}{2d(1-\rho)} + \frac{\overline{s_{\text{um}}}}{2}\left[\frac{n-\rho}{1-\rho} + \frac{\overline{s_{\text{um}}^2} - \overline{s_{\text{um}}}^2}{\overline{s_{\text{um}}}^2}\right], \quad (21)$$

where $\overline{s_{\text{um}}} = \overline{U} + \overline{M}$ and $\overline{s_{\text{um}}^2} = \overline{U^2} + \overline{M^2} + 2\overline{U}\,\overline{M}.$ (22)

By considering (9), (21) yields

$$\overline{W_h}(\rho) = \frac{\rho\overline{B^2}}{2\overline{B}(1-\rho)} + \frac{\overline{s_{\text{um}}}}{2}\left[\frac{n-1}{1-\rho} + \frac{\overline{s_{\text{um}}^2}}{\overline{s_{\text{um}}}^2}\right]. \quad (23)$$

### C. Medium-Load Analysis

Fig. 1 in Section VII shows the simulation-based mean waiting time results along with the theoretical curves obtained from (19) and (23). We observe that the theoretical light-load curve matches well with the simulation results in the light-load region $[0, \rho_l]$, and the theoretical heavy-load curve matches

well with the simulation results in the heavy-load region $[\rho_h, 1]$. In the medium-load region $[\rho_l, \rho_h]$, we observe that the mean-waiting-time curve increases almost linearly, which cannot be efficiently approximated by $M/M/c$ and $M^X/G/c$ queueing models such as those used in [10], [11]. Triggered by this observation, we propose, as an approximation in the medium-load region, the unique line that is tangent to the $\overline{W_l}(\rho)$ and $\overline{W_h}(\rho)$ curves at points $\rho_l$ and $\rho_h$, respectively. Note that the two curves $\overline{W_l}(\rho)$ and $\overline{W_h}(\rho)$ are convex and saturate at $\rho = \rho^* < 1$ and $\rho = 1$, respectively. It can also be shown that it holds that $\overline{W_l}(0) < \overline{W_h}(0)$, which implies that there is always a unique line that is tangent to the two curves, at points $\rho_l$ and $\rho_h$ with $\rho_l < \rho^* < \rho_h$, given by

$$\overline{W_m}(\rho) = \frac{\overline{W_h}(\rho_h) - \overline{W_l}(\rho_l)}{\rho_h - \rho_l}\rho - \frac{\rho_l\overline{W_h}(\rho_h) - \rho_h\overline{W_l}(\rho_l)}{\rho_h - \rho_l}. \quad (24)$$

Note that, owing to the complexity of expression (6), which yields the $\overline{W_f}$ component of the $\overline{W_l}$ function, the $\rho_l$ and $\rho_h$ values can be evaluated only numerically. This basic, partially numerical analysis yields the mean waiting time for the entire range of loads as follows:

$$\overline{W}(\rho) = \begin{cases} \overline{W_l}(\rho), & \text{for } \rho \in [0, \rho_l], \\ \overline{W_m}(\rho), & \text{for } \rho \in [\rho_l, \rho_h], \\ \overline{W_h}(\rho), & \text{for } \rho \in [\rho_h, 1), \end{cases} \quad (25)$$

with $\overline{W_l}(\rho)$, $\overline{W_m}(\rho)$ and $\overline{W_h}(\rho)$, given by (19), (24) and (23), respectively. It turns out that the tangent line, along with $\rho_l$ and $\rho_h$, can be analytically determined if instead of the virtual $M/G/K$ queue, we consider the corresponding single server $M/G/1$ queue, with the service time $B$ being $K = (d)$ times shorter. In this case, the approximate fictitious mean waiting time, $\overline{W_{\text{f,approx}}}$, is obtained by

$$\overline{W_{\text{f,approx}}} \approx \frac{\lambda\overline{S_f^2}}{2d^2(1-\rho_f)}. \quad (26)$$

Combining (7), (9), (11), and (26), (19) yields the approximate mean waiting time $\overline{W_l^{(a)}}$ as follows:

$$\overline{W_l^{(a)}}(\rho) = \overline{W_{\text{f,approx}}} + H \approx \frac{\overline{S_f^2}\rho}{2d\overline{S_f}(\rho^* - \rho)} + H. \quad (27)$$

From (23) and (27), after some manipulations, it follows that the tangent $\overline{W_m^{(a)}}$ is given by

$$\overline{W_m^{(a)}}(\rho) = \frac{\overline{W_h}(\rho_h^{(a)}) - \overline{W_l^{(a)}}(\rho_l^{(a)})}{\rho_h^{(a)} - \rho_l^{(a)}}\rho - \frac{\rho_l^{(a)}\overline{W_h}(\rho_h^{(a)}) - \rho_h^{(a)}\overline{W_l^{(a)}}(\rho_l^{(a)})}{\rho_h^{(a)} - \rho_l^{(a)}}, \quad (28)$$

where $\rho_l^{(a)} = (-Y + \sqrt{Y^2 - 4XZ})/(2X),$ (29)

$$\rho_h^{(a)} = 1 - R(1 - \rho_l/\rho^*), \quad (30)$$

with $R = \sqrt{(C+G)\rho^*/A},$ (31)

$$A = \overline{S_f^2}/(2d\overline{S_f}), \quad (32)$$

$$C = [n\overline{s_{\text{um}}} + (\overline{s_{\text{um}}^2} - \overline{s_{\text{um}}}^2)/\overline{s_{\text{um}}}^2]/2, \quad (33)$$

$$G = \overline{B^2}/(2\overline{B}) - \overline{s_{\text{um}}^2}/(2\overline{s_{\text{um}}}), \quad (34)$$

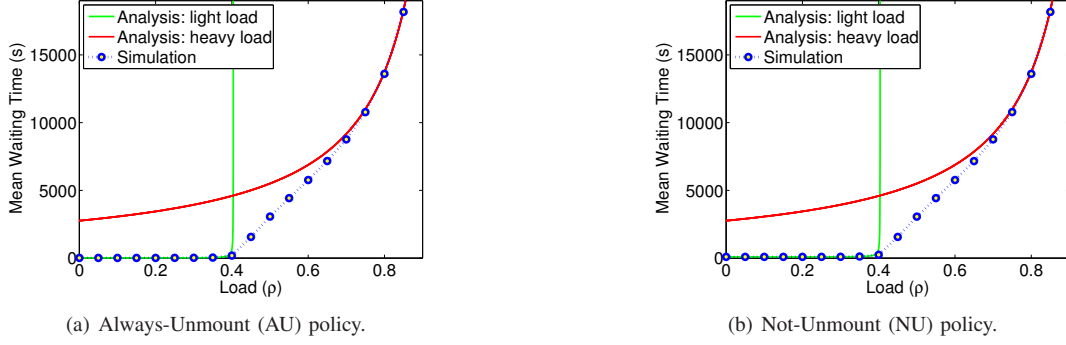(a) Always-Unmount (AU) policy.



(b) Not-Unmount (NU) policy.

Fig. 1. Theoretical mean waiting times for light- and heavy-load approximations vs. load along with simulation results; $c = 720$, $d = 12$.

$$X = (G + H - A) R \,, \qquad (35)$$

$$Y = A R^2 + [C + G - 2(G + H) R] \rho^* \,, \qquad (36)$$

$$Z = \{A R (1 - R) - [C + G - (G + H) R] \rho^*\} \rho^* \,, \qquad (37)$$

with $\rho^*$ given by (11), the moments of $S_f$ given by (4) and (5), the moments of $s_{\text{um}}$ given by (22), and $H$ given by (20).

Thus, this approximate analysis yields a closed-form expression for the mean waiting time $\overline{W^{(a)}}$ for the entire range of loads as follows:

$$\overline{W^{(a)}}(\rho) = \begin{cases} \overline{W_l^{(a)}}(\rho), & \text{for } \rho \in [0, \rho_l^{(a)}], \\ \overline{W_m^{(a)}}(\rho), & \text{for } \rho \in [\rho_l^{(a)}, \rho_h^{(a)}], \\ \overline{W_h}(\rho), & \text{for } \rho \in [\rho_h^{(a)}, 1) \,, \end{cases} \qquad (38)$$

with $\overline{W_l^{(a)}}(\rho)$, $\overline{W_m^{(a)}}(\rho)$, $\overline{W_h}(\rho)$, $\rho_l^{(a)}$ and $\rho_h^{(a)}$ given by (27), (28), (23), (29) and (30), respectively.

As we will see in Section VII, the light-load approximation given by (27) is not very accurate, leading to a large deviation of the mean waiting times in the light-load region. However, it turns out that the derived region limits $\rho_l^{(a)}$ and $\rho_h^{(a)}$ are very good approximations of the actual limits $\rho_l$ and $\rho_h$. To improve the accuracy of the model, we propose to use a variant of (25) with $\rho_l$ and $\rho_h$ replaced by $\rho_l^{(a)}$ and $\rho_h^{(a)}$, respectively. Consequently, the mean waiting time is now obtained by the following enhanced closed-form expression:

$$\overline{W^{(e)}}(\rho) = \begin{cases} \overline{W_l}(\rho), & \text{for } \rho \in [0, \rho_l^{(a)}], \\ \overline{W_m}(\rho), & \text{for } \rho \in [\rho_l^{(a)}, \rho_h^{(a)}], \\ \overline{W_h}(\rho), & \text{for } \rho \in [\rho_h^{(a)}, 1) \,, \end{cases} \qquad (39)$$

with $\overline{W_l}(\rho)$, $\overline{W_m}(\rho)$, $\overline{W_h}(\rho)$, $\rho_l^{(a)}$ and $\rho_h^{(a)}$ given by (19), (24), (23), (29) and (30), respectively. Because the $\rho_l^{(a)}$ and $\rho_h^{(a)}$ values are very close to those of $\rho_l$ and $\rho_h$, the enhanced $\overline{W^{(e)}}$ curve is almost identical to the partially numerically obtained $\overline{W}$ curve given by (25), with albeit the advantage of being expressed analytically in closed form.

## VII. Numerical Results

Here we assess the performance of a tape library by using both theoretical predictions and event-driven simulations. We have confirmed the validity of the model by considering scenarios for a range of distributions and parameter values. We begin by presenting the specific results for the IBM® TS4500
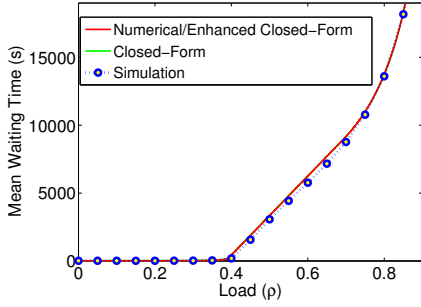
TABLE II
Medium-Load Region Boundaries

| Policy | $\rho_l$ | $\rho_h$ | $\rho_l^{(a)}$ | $\rho_h^{(a)}$ |
|--------|----------|----------|----------------|----------------|
| AU | 0.3945 | 0.6953 | 0.3944 | 0.6925 |
| NU | 0.3944 | 0.6920 | 0.3944 | 0.6912 |

tape library system [24]. The system comprises $c = 720$ cartridges and $d = 12$ servers, such that $n = 60$. The unmount and mount times are considered to be fixed, equal to $U = 77$ s and $M = 15$ s, respectively. Requests incur a fixed seek time of $s = 60$ s and, from a study that observes uniform access across files in archival data movement workloads [26], along with the assumption that files are accessed in their entirety in this context, we deduce that the distribution of I/O request sizes is the same as that of the file sizes in the archive. In particular, we consider the request size distribution to be the same as the file size distribution of CERN [27], whose mean $\overline{Q}$ is equal to 843 MB, the standard deviation to 2.8 GB and the second moment $\overline{Q^2}$ to 8.5 GB$^2$. The transfer bandwidth is assumed to be $b_w = 360$ MB/s.
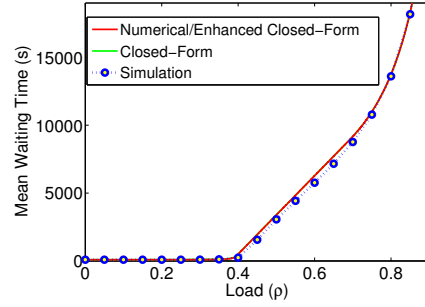
Fig. 1 illustrates the theoretical mean waiting times obtained from (19), which is derived by the light-load analysis, and (23), which is derived by the heavy-load analysis, indicated by the green and red curves, respectively, and shown for the entire range of loads. Figs. 1 (a) and (b) show the results for the AU and NU policies, respectively. Note that for both policies, the light-load green curves $\overline{W_l}(\rho)$ saturate at $\rho^* = 0.4039$. The heavy-load red curves do not depend on the policy and therefore are the same. The simulation results are indicated by the blue circles, with the 95% confidence intervals being extremely narrow and therefore not shown.

Fig. 2 shows the theoretical predictions for the AU and NU policies obtained by using (25) and indicated by the red curves, along with the simulation results. The values of $\rho_l$ and $\rho_h$ for the AU and NU policies are obtained numerically and listed in Table II. Fig. 2 also shows the approximate closed-form theoretical predictions obtained by making use of (38) as green curves. These curves are barely visible because they lie just below the red ones. The values of $\rho_l^{(a)}$ and $\rho_h^{(a)}$ are obtained by (29) and (30), respectively, and listed in Table II. We observe that the values of $\rho_l^{(a)}$ and $\rho_h^{(a)}$ are very close to those of $\rho_l$ and $\rho_h$, respectively.

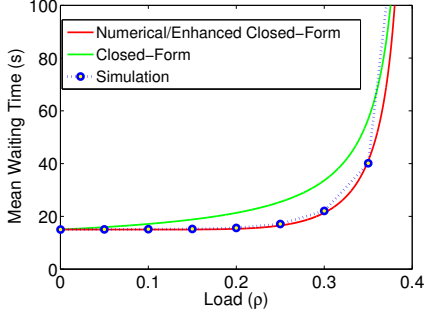Fig. 3 shows the results in the light-load region for the AU and NU policies. We observe that the red curves, obtained
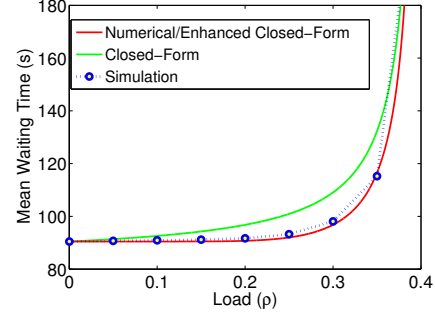
6

(a) Always-Unmount (AU) policy.



(b) Not-Unmount (NU) policy.

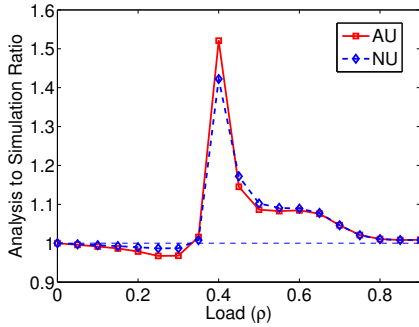Fig. 2.   Mean waiting times for AU and NU vs. load; $c = 720$, $d = 12$.



(a) Always-Unmount (AU) policy.


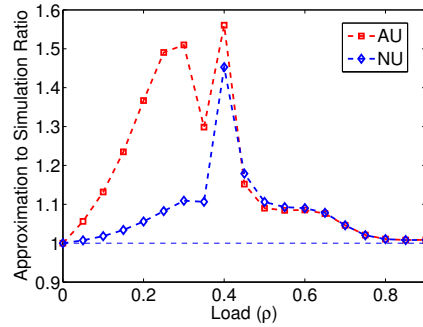
(b) Not-Unmount (NU) policy.

Fig. 3.   Mean waiting times for AU and NU vs. load in the light-load region; $c = 720$, $d = 12$.



(a) Partially numerical / Enhanced closed-form approximation.



(b) Closed-form approximation.

Fig. 4.   Ratio of analytical estimations to simulation results as a function of the load; $c = 720$, $d = 12$.

partially numerically by (25), match well with the simulation results as opposed to the green curves, obtained by the closed-form approximation (38), which exhibit a significant deviation. This is reflected in Fig. 4, which shows the ratio of the analytical predictions to the corresponding simulation values. Clearly, Fig. 4 (a) reveals that the theoretical results match well with the simulation ones in the light- and heavy-load regions, with their worst-case deviations being equal to 2% and 3%, respectively. In the medium-load region, the deviation is on the average about 10% and exhibits a spike with a peak of about 50% only in a narrow region around the saturation load $\rho^* = 0.4039$ of the virtual queue. Nevertheless, despite the drastic change in the system's mean waiting time, which is between one and two orders of magnitude, for loads close to $\rho^* \simeq 0.4$, the analytical results turn out to properly capture this behavior (cf. Figs. 2 and 3). In contrast, Fig. 4 (b) shows a significant deviation of the closed-form approximation results
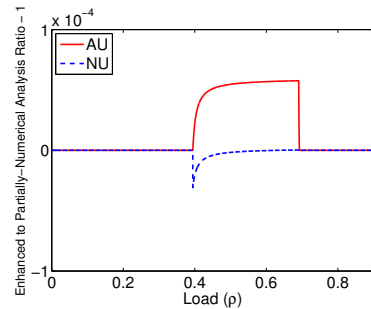


Fig. 5.   Deviation of the enhanced mean waiting times from the partially numerically ones for AU and NU vs. load; $c = 720$, $d = 12$.

in the light-load region, obtained by (38), from the simulation results.

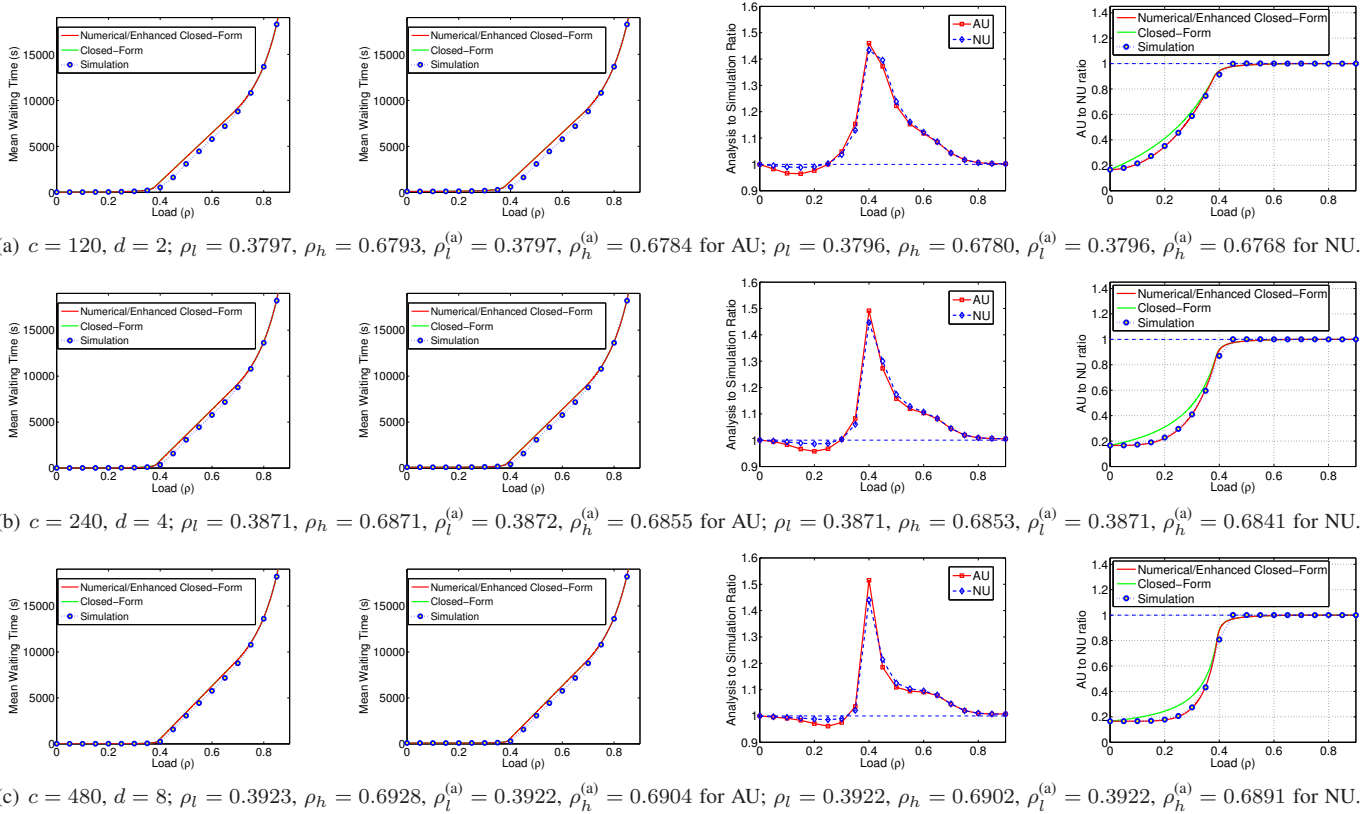Fig. 5 plots the $\overline{W^{(e)}}(\rho)/\overline{W}(\rho) - 1$ quantity, which reveals

(a) $c = 120$, $d = 2$; $\rho_l = 0.3797$, $\rho_h = 0.6793$, $\rho_l^{(a)} = 0.3797$, $\rho_h^{(a)} = 0.6784$ for AU; $\rho_l = 0.3796$, $\rho_h = 0.6780$, $\rho_l^{(a)} = 0.3796$, $\rho_h^{(a)} = 0.6768$ for NU.



(b) $c = 240$, $d = 4$; $\rho_l = 0.3871$, $\rho_h = 0.6871$, $\rho_l^{(a)} = 0.3872$, $\rho_h^{(a)} = 0.6855$ for AU; $\rho_l = 0.3871$, $\rho_h = 0.6853$, $\rho_l^{(a)} = 0.3871$, $\rho_h^{(a)} = 0.6841$ for NU.



(c) $c = 480$, $d = 8$; $\rho_l = 0.3923$, $\rho_h = 0.6928$, $\rho_l^{(a)} = 0.3922$, $\rho_h^{(a)} = 0.6904$ for AU; $\rho_l = 0.3922$, $\rho_h = 0.6902$, $\rho_l^{(a)} = 0.3922$, $\rho_h^{(a)} = 0.6891$ for NU.

Fig. 8. Various configurations with a ratio of the number of cartridges to the number of tape drives of $n = 60$.
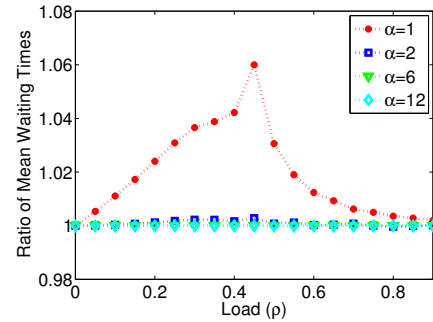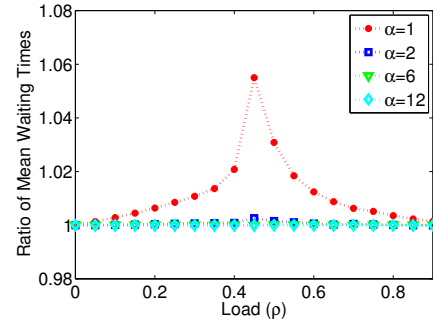


Fig. 6. Ratio of mean waiting times of AU to NU vs. load; $c = 720$, $d = 12$.

that the enhanced closed-form $\overline{W^{(e)}}$ curve given by (39) is almost identical to the partially numerically obtained $\overline{W}$ curve given by (25). We have verified that this holds in all cases presented next, and this is because the values of $\rho_l^{(a)}$ and $\rho_h^{(a)}$ are always very close to those of $\rho_l$ and $\rho_h$, respectively. Consequently, the red curves shown in all figures correspond to the values obtained by both the partially numerical model and the enhanced closed-form model.

Fig. 6 shows the results for the ratio of the mean waiting times of the AU to those of the NU policy. As expected, for light loads, the AU policy is superior because it does not waste any time to unmount cartridges when requests arrive. As the load increases, the performance of the AU policy approaches that of the NU policy. Clearly, the theoretical results obtained by the partially numerical and enhanced models and shown
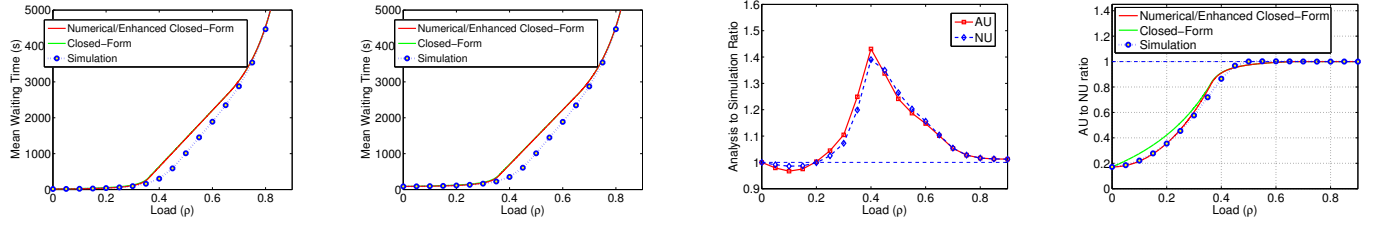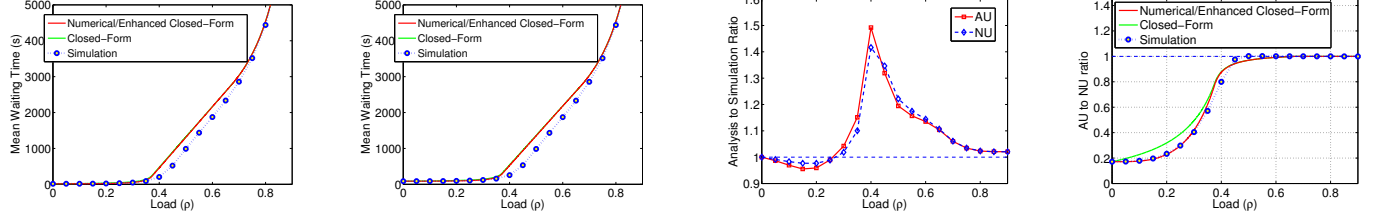


(a) Always-Unmount (AU) policy.
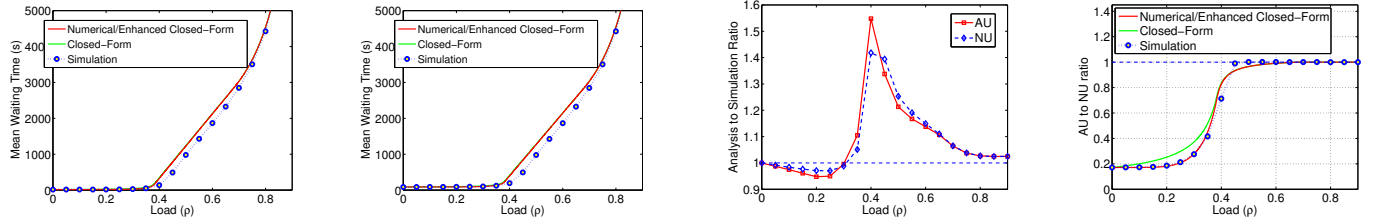


(b) Not-Unmount (NU) policy.

Fig. 7. Effect of number of arms on the mean waiting time as a function of the load; $c = 720$, $d = 12$.

(a) $c = 40$, $d = 2$; $\rho_l = 0.3595$, $\rho_h = 0.6615$, $\rho_l^{(a)} = 0.3592$, $\rho_h^{(a)} = 0.6597$ for AU; $\rho_l = 0.3588$, $\rho_h = 0.6562$, $\rho_l^{(a)} = 0.3585$, $\rho_h^{(a)} = 0.6543$ for NU.
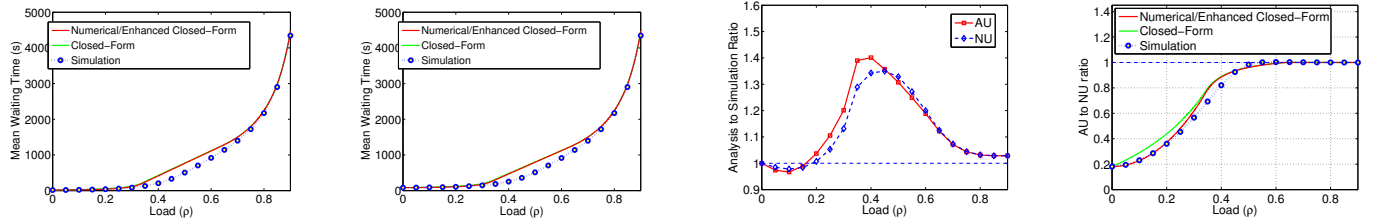
(b) $c = 80$, $d = 4$; $\rho_l = 0.3737$, $\rho_h = 0.6742$, $\rho_l^{(a)} = 0.3735$, $\rho_h^{(a)} = 0.6727$ for AU; $\rho_l = 0.3733$, $\rho_h = 0.6698$, $\rho_l^{(a)} = 0.3731$, $\rho_h^{(a)} = 0.6678$ for NU.
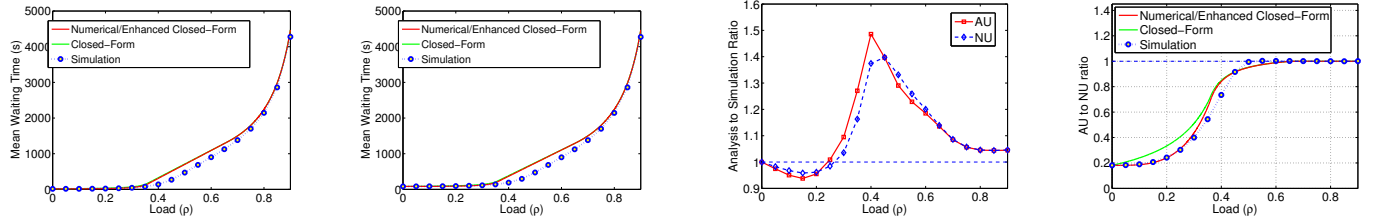
(c) $c = 160$, $d = 8$; $\rho_l = 0.3832$, $\rho_h = 0.6841$, $\rho_l^{(a)} = 0.3830$, $\rho_h^{(a)} = 0.6816$ for AU; $\rho_l = 0.3829$, $\rho_h = 0.6795$, $\rho_l^{(a)} = 0.3827$, $\rho_h^{(a)} = 0.6771$ for NU.
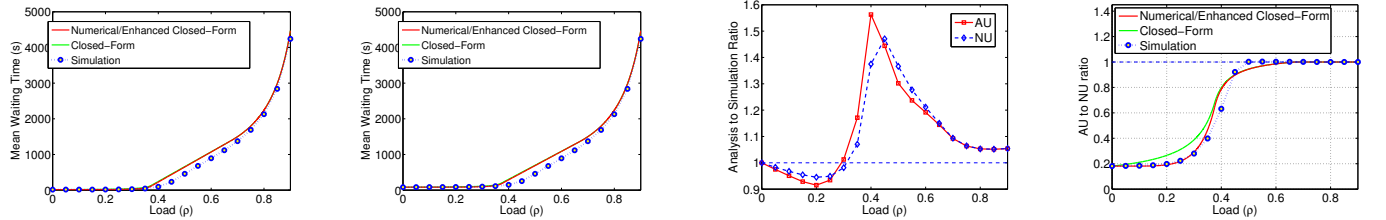
Fig. 9.   Various configurations with a ratio of the number of cartridges to the number of tape drives of $n = 20$.

(a) $c = 20$, $d = 2$; $\rho_l = 0.3374$, $\rho_h = 0.6435$, $\rho_l^{(a)} = 0.3365$, $\rho_h^{(a)} = 0.6401$ for AU; $\rho_l = 0.3352$, $\rho_h = 0.6316$, $\rho_l^{(a)} = 0.3342$, $\rho_h^{(a)} = 0.6276$ for NU.

(b) $c = 40$, $d = 4$; $\rho_l = 0.3595$, $\rho_h = 0.6627$, $\rho_l^{(a)} = 0.3588$, $\rho_h^{(a)} = 0.6593$ for AU; $\rho_l = 0.3582$, $\rho_h = 0.6527$, $\rho_l^{(a)} = 0.3574$, $\rho_h^{(a)} = 0.6486$ for NU.

(c) $c = 80$, $d = 8$; $\rho_l = 0.3737$, $\rho_h = 0.6752$, $\rho_l^{(a)} = 0.3732$, $\rho_h^{(a)} = 0.6724$ for AU; $\rho_l = 0.3729$, $\rho_h = 0.6665$, $\rho_l^{(a)} = 0.3724$, $\rho_h^{(a)} = 0.6629$ for NU.

Fig. 10.   Various configurations with a ratio of the number of cartridges to the number of tape drives of $n = 10$.

by the red curve are in agreement with the simulation ones, which establishes a confidence for the model presented.

We have also investigated by simulation the effect of the number of robot arms on the mean waiting time when arm operations take 3.3 s [24]. Fig. 7 shows the ratio of the mean waiting times for various values of $a$ to the optimal ones obtained when there is no contention for the robot arms, that is, when the number of robot arms is equal to the number of tape drives ($a = d = 12$). We observe that for $a \geq 2$, the performance is insensitive to the number of arms because there is practically no contention for them. For a single arm, the contention results in an increase of the mean waiting time by at most 6%, which occurs at medium loads; as expected at low and and high loads the corresponding increase is very small, at most 4% and 0.7%, respectively. These results establish that the effect of the number of arms is negligible, which in turn confirms the validity of the model presented.

Next, we examine the sensitivity of the accuracy of the theoretical results to the number of cartridges and the number of tape drives. First, we consider a ratio of $n = 60$ cartridges per tape drive. The theoretical and simulation results obtained for $d = 2, 4$ and 8 tape drives are shown in Fig. 8, with each of the corresponding three rows containing four figures. In each row, the first two figures show the mean waiting times for the AU and NU policies, respectively, similarly to Fig. 2. The values of $\rho_l$, $\rho_h$, $\rho_l^{(a)}$ and $\rho_h^{(a)}$ for the AU and NU policies are listed in the corresponding figure caption. The third figure shows the ratio of the analytical predictions of the enhanced theoretical model, obtained by (39), to the corresponding simulation values. The fourth figure shows the ratio of the mean waiting times of the AU to the NU policies, similarly to Fig. 6. From Fig. 8, we observe that as the number of drives increases, the accuracy of the theoretical approximation improves in that the spike, shown in the third figures for loads around the value of $\rho^* = 0.4039$, becomes narrower, with its height remaining roughly constant. Also, the red curves in the fourth figures match well with the simulation results and, as the number of drives increases, tend to become a step function for a load equal to $\rho^*$. This implies that for $\rho < \rho^*$, the ratio of AU to NU remains practically the same, and for $\rho > \rho^*$, it tends to one. Figs. 9 and 10 reveal that the same applies when $n$ takes smaller values. However, the width of the spike becomes wider, which implies that the accuracy of the model deteriorates at medium loads, but is quite satisfactory at light and heavy loads, with the deviation being less than 10%.

## VIII. CONCLUSIONS

The unrelenting growth of big data has fueled a high demand for tape storage. Predicting the performance of a tape library system is key to efficiently dimension it, and, in particular, to dimension multi-tiered storage systems that contain several type of devices, including tape. Our work is the first to provide an accurate analytical model for assessing the performance of tape library systems. Extending this model to include enhanced interpolation schemes, to study asymmetric workloads, and to incorporate advanced scheduling policies is the subject of further investigations.

## REFERENCES

[1] S. S. Lavenberg and D. R. Slutz, "Regenerative simulation of a queuing model of an automated tape library," *IBM J. Res. Dev.*, vol. 19, no. 5, pp. 463–475, Sep. 1975.

[2] A. L. Drapeau and R. H. Katz, "Striping in large tape libraries," in *Proc. ACM/IEEE Int'l Conf. on Supercomputing (SC)*, Nov. 1993, pp. 378–387.

[3] L. Golubchik, R. R. Muntz, and R. W. Watson, "Analysis of striping techniques in robotic storage libraries," in *Proc. 14th IEEE Symp. on Mass Storage Systems (MASS)*, Sep. 1995, pp. 225–238.

[4] O. I. Pentakalos, D. A. Menasce, M. Halem, and Y. Yesha, "Analytical performance modeling of hierarchical mass storage systems," *IEEE Trans. Comput.*, vol. 46, no. 10, pp. 1103–1138, Oct. 1997.

[5] P. Triantafillou and T. Papadakis, "On-demand data elevation in hierarchical multimedia storage servers," in *Proc. 23rd Int'l Conf. on Very Large Data Bases (VLDB)*, Aug. 1997, pp. 226–235.

[6] T. Johnson and E. L. Miller, "Performance measurements of tertiary storage devices," in *Proc. 24rd Int'l Conf. on Very Large Data Bases (VLDB)*, Aug. 1998, pp. 50–61.

[7] I. Koltsidas *et al.*, "Seamlessly integrating disk and tape in a multi-tiered distributed file system," in *Proc. IEEE 31st Int'l Conf. on Data Engineering (ICDE)*, Apr. 2015, pp. 1328–1339.

[8] M. Lantz *et al.*, "123 Gbit/in$^2$ recording areal density on barium ferrite tape," *IEEE Trans. Magn.*, vol. 51, no. 11, pp. 1–4, Nov. 2015.

[9] I. Iliadis, J. Jelitto, Y. Kim, S. Sarafijanovic, and V. Venkatesan, "Exa-Plan: Queueing-based data placement and provisioning for large tiered storage systems," in *Proc. 23nd Annual IEEE Int'l Symp. on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*, Oct. 2015, pp. 218–227.

[10] J. J. Gniewek, "Evolving requirements for magnetic tape data storage systems," in *Proc. 5th NASA Goddard Conf. on Mass Storage Systems and Technologies (MSST)*, Sep. 1996, pp. 477–491.

[11] T. Johnson, "An analytical performance model of robotic storage libraries," *Perform. Eval.*, vol. 27-28, pp. 231–251, Oct. 1996.

[12] B. K. Hillyer and A. Silberschatz, "Random I/O scheduling in online tertiary storage systems," in *Proc. ACM Int'l Conf. on Management of Data (SIGMOD)*, Jun. 1996, pp. 195–204.

[13] P. Triantafillou and I. Georgiadis, "Hierarchical scheduling algorithms for near-line tape libraries," in *Proc. 10th Int'l Workshop on Database and Expert Systems Applications (DEXA)*, Sep. 1999, pp. 50–54.

[14] M. A. A. Boon, R. D. van der Mei, and E. M. M. Winands, "Applications of polling systems," *Surveys in Operations Research and Management Science*, vol. 16, no. 2, pp. 67–82, Jul. 2011.

[15] M. P. Singh and M. M. Srinivasan, "Exact analysis of the state-dependent polling model," *Queueing Systems*, vol. 41, no. 4, pp. 371–399, Aug. 2002.

[16] H. Takagi, *Analysis of Polling Systems*. Cambridge, MA, USA: MIT Press, 1986.

[17] Y. Günalay and D. Gupta, "Polling systems with a patient server and state-dependent setup times," *IIE Trans.*, vol. 29, no. 6, pp. 469–480, 1997.

[18] W. S. Lai, D. J. Houck, and S. W. Fuhrmann, "Two-stage polling system with multiple servers," in *Proc. SPIE, Performance and Control of Network Systems*, vol. 3231, Oct. 1997, pp. 535–545.

[19] R. J. T. Morris and Y. T. Wang, "Some results for multi-queue systems with multiple cyclic servers," in *Proc. IFIP Int'l Symp. on the Performance of Computer Communication Systems*, Mar. 1984, pp. 245–258.

[20] S. C. Borst and R. D. van der Mei, "Waiting time approximations for multiple-server polling systems," *Perform. Eval.*, vol. 31, no. 3-4, pp. 163–182, Jan. 1998.

[21] S. C. Borst, "Polling systems with multiple coupled servers," *Queueing Systems*, vol. 20, no. 3-4, pp. 369–393, Sep. 1995.

[22] M. J. Fischer, C. M. Harris, and J. Xie, "An interpolation approximation for expected wait in a time-limited polling system," *Comput. Oper. Res.*, vol. 27, no. 4, pp. 353–366, Apr. 2000.

[23] D. Everitt, "Simple approximations for token rings," *IEEE Trans. Commun.*, vol. 34, no. 7, pp. 719–721, Jul. 1986.

[24] IBM® TS4500 Tape Library: Specifications. http://www-03.ibm.com/systems/storage/tape/ts4500/specifications.html, Jan. 2016. IBM is a registered trademark of International Business Machines Corporation, registered in many jurisdictions worldwide.

[25] A. M. Lee and P. A. Longton, "Queueing processes associated with airline passenger check-in," *Operations Research*, vol. 10, pp. 56–71, Sep. 1959.

[26] I. F. Adams *et al.*, "Usage behavior of a large-scale scientific archive," in *Proc. ACM/IEEE Int'l Conf. on Supercomputing (SC)*, Nov. 2012, pp. 1–11.

[27] G. Cancio *et al.*, "Tape archive challenges when approaching exabyte-scale," 2010, Presentation at CHEP 2010, available online.