

Design of Scalable Continuous Media Servers

Cheng-Fu Chou*

Leana Golubchik[†]

John C.S. Lui[‡]

I-Hsin Chung[§]

Abstract

Multimedia applications place high demands for quality-of-service (QoS), performance, and reliability on systems. These stringent requirements make design of cost-effective and scalable systems difficult. Therefore efficient *adaptive* and *dynamic* resource management techniques in conjunction with data placement techniques can be of great help in improving performance, scalability, and reliability of such systems. This is the focus of our paper.

Keywords: continuous media systems, dynamic resource management, data placement, storage servers.

1 Introduction

Multimedia applications place high demands for quality-of-service (QoS), performance, and reliability on storage servers and communication networks. These, often stringent, requirements make end-to-end design of cost-effective and *scalable* continuous media (CM) servers difficult. The scalability of a CM server's architecture depends on its ability to: (a) expand as user demand and data sizes grow; (b) maintain performance characteristics under degradation of system resources; and (c) maintain performance characteristics under growth or re-configuration. In particular, the choice of data placement techniques can have a significant effect on the scalability of a CM server and its ability to utilize resources efficiently. Existing data placement techniques in conjunction with scheduling algorithms address two major inefficiencies in such systems: (1) various overheads in reading data from storage devices and (2) load imbalance, e.g., due to *skews* in data access patterns. In this work, we focus on the latter issue and specifically on its bearing on the *scalability* characteristics of a *distributed* CM server.

An important issue in distributed designs is the placement of objects on the nodes of the CM server which directly affects its load balancing characteristics. In the recent past, a great deal of CM server designs, e.g., as in [9, 11], have focused on “wide” data striping techniques, where each object is striped across all the disks of the system as an approach to dealing with load imbalance problems due to skews in data access patterns (which change over time). An advantage of wide data striping is that it “implicitly” achieves load balance by decoupling an object's storage from its bandwidth requirements. However, wide data striping also suffers from a number of shortcomings which we detailed in [5, 4].

Another approach to dealing with the load imbalance problem arising from skews in data access patterns is replication, i.e., creating a sufficient number of copies of a (popular) object so as to meet the demand for that object, where the disadvantages of replication include: (1) a need for additional storage space, and (2) a need for techniques that adjust the number of replicas as the access patterns change. Some of these issues are addressed in [17], and in our earlier work [13, 6].

*Department of Computer Science, University of Maryland, chengfu@cs.umd.edu.

[†]Dept. of Computer Science & UMIACS, University of Maryland, leana@cs.umd.edu. This research was supported in part by the NSF CCR-98-96232 and ANI-00-70016 grants.

[‡]Dept. of Computer Science & Engineering, The Chinese University of Hong Kong, cslui@cs.cuhk.edu.hk. This research was supported in part by the RGC Earmarked Grant and the CUHK Mainline Research Grant.

[§]Department of Computer Science, University of Maryland, ihchung@cs.umd.edu.

In this paper, we consider a *hybrid* approach, where the main focus is on the *tradeoffs* between striping and replication. Essentially, striping is a good approach to load balancing while replication is a good approach to “isolating” nodes from being dependent on other (“non-local”) system resources. That is, the wider we stripe in a distributed CM system, the more we are dependent on the availability of network capacity as well as resources not local to a node. Furthermore, replication has the benefit of increased reliability in terms of: (a) longer mean time to loss of data from the disk sub-system, and (b) dealing with lack of network resources, including network partitioning. The downside of replication is that it increases storage space requirements and hence cost of storage. However, as storage costs decrease (fairly rapidly) and the need for scalability grows, replication becomes a more attractive technique. In summary, the appropriate compromise between the degree of striping and the degree of replication is key to the design of a *scalable* distributed CM server.

Much research has been done on design of CM storage servers, e.g., [9, 11]. To the best of our knowledge, most of these designs employ wide data striping techniques and the corresponding existing successful implementations employ only tens of disks. In contrast, the use of replication for the purpose of addressing changes in data access patterns has been less explored. In [16] the authors consider skews in data access patterns but in the context of a static environment. In [17], the authors address various questions arising in the context of load imbalance problems due to skews in data access patterns, but in a less dynamic environment (than we investigate here)¹. In [8, 7], the authors also consider dynamic replication as an approach to load imbalance, and in our previous work [13, 6], we study a taxonomy of dynamic replication schemes. However, all of these works, except our work in [5, 6], either (a) assume some *knowledge of frequencies of data access* to various objects in the system, and/or (b) do *not* provide users with *full use of VCR functionality*, and/or (c) consider less dynamic environments than the one considered here. Our motivation in doing away with such assumptions in our work is largely due to considerations of applicability of dynamic replication techniques in more general settings and to a wider range of applications of CM servers.

Lastly, to the best of our knowledge, previous works do not consider *alternative design characteristics* that affect the *scalability* of CM servers in an *end-to-end* setting under changes in access patterns (i.e., taking into consideration both the network and the storage resource constraints). The *quantitative* study of such issues and the cost/performance and reliability characteristics that distributed designs exhibit under growth, reconfiguration, degradation of resources, and changes in workloads are essential to assessing the *scalability* of proposed architectures and to the development of *large-scale* CM servers, in general.

Based on the end-to-end cost/performance and reliability study presented in this paper we believe that *hybrid* designs in conjunction with dynamic replication schemes result in large *scalable*, reliable, high performance, and cost-effective CM systems.

2 Hybrid CM System Architecture

A hybrid system architecture is illustrated in Figure 1(a). It consists of a set of nodes connected by, what we termed, a high speed *global* switch, which is a high bandwidth resource that can, for instance, correspond to a high capacity WAN or an ATM-type infrastructure. Each node i , as depicted in Figure 1(b), contains one or more processing units (PUs) and one, what we termed, *local* switch (e.g., switched Ethernet) which is used to connect all local PUs as well as local clients. Each client connects to its local switch. Requests from this client which are serviced by a PU from node i are termed “local”. When a request from a client cannot be serviced by its local node i , it is forwarded to a remote node j , which

¹We believe that the policies used in this paper can be complementary to the techniques developed in [17].

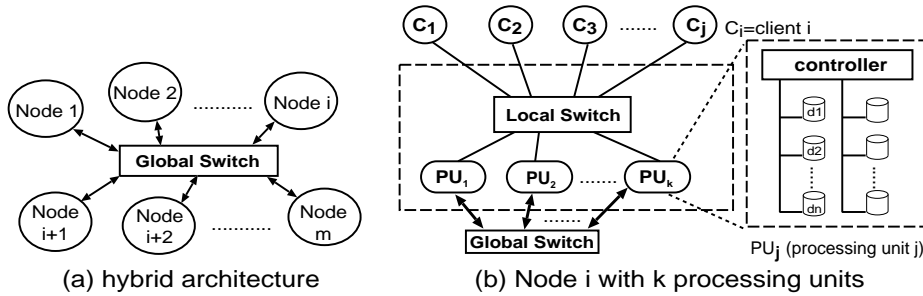


Figure 1: Hybrid System Architecture.

contains a replica of the requested object. We term this request “global”, as its service requires some capacity of the global switch, i.e., to deliver data from the remote node, through the global switch to the local node and subsequently to the client.

Each PU has one or more CPUs, memory, and an I/O sub-system (e.g., a cluster/array of disks), and it is also connected to the global switch. Each node $x \in S$, where S is the set of nodes in the system, has a finite storage capacity, D_x (*in units of CM objects*), as well as a finite service capacity, B_x (*in units of CM access streams*). Likewise, we measure the global and local switch capacities *in units of access streams*. In general, different nodes in such a *hybrid* system may differ in their storage, I/O bandwidth, and local switch service capacity. This flexibility of the hybrid architecture should result in a scalable system which can grow on a node by node basis.

Each CM object resides on one or more nodes of the system depending on its current popularity. An object is striped *only* across local disks which belong to the same node. Objects that require more than a single node’s service capacity (to support the corresponding demand) are replicated on multiple nodes. The number of replicas needed to support requests for a CM object is a function of demand, and therefore this number should change as the demand changes.

Let $R_i(t)$ denote the set of nodes containing replicas of object i at time t . Upon a customer’s arrival at time t , there is a probability $p_i(t)$ that the corresponding request is for object i and a probability $q_j^i(t)$ that this request is generated by a client local to node j . The admission of this customer into the system proceeds as follows. If at time t object i resides on node j and there is service capacity available at node j , then the system admits and serves this new request at node j , i.e., locally. Let $L_x(t)$ be the load on node x at time t . If at time t object i does not reside on node j or there is no service capacity available at node j , then the system examines the load information on each node in $R_i(t)$, and if there is sufficient capacity (on at least one of these node and in the global switch), to service the newly arrived request, the system assigns this request to the least-loaded node in $R_i(t)$. Otherwise, the customer cannot be served immediately. In this case we consider two alternatives : (1) the customer is immediately rejected (referred to as the no queueing case), or (2) the customer joins a FIFO queue and awaits service where there is no restriction on the size of the queue (referred to as the queueing case). These are two extreme cases (i.e., allowing no queueing at all or allowing infinite queueing). And, although these are not necessarily representative of how a real system should operate², they are useful in our performance evaluation study, i.e., these cases provide the necessary insight for the design of CM server.

Note that a “stream migration” approach to dealing with more “short-term” fluctuations in access patterns

²In a real system a reasonably sized finite queue is used, whose size depends on the required performance characteristics; the actual size of the queue should be a function of the performance characteristics required by an application.

is given in [17]; this is an orthogonal approach and can be combined with techniques presented in this paper. However, in the interest of isolating the performance effects of dynamic replication we do not consider this here further. We also note that the problem of determining whether or not there is “sufficient capacity” (under either CBR or VBR stream models) is orthogonal to the problems studied in this work; much literature exists on this topic (refer to [10]), and such solutions can be used in conjunction with policies developed in this paper.

Full VCR functionality (i.e., fast-forward, rewind, and pause/resume) is available to all admitted customers, with fast-forward and rewind provided at $n_{speed} > 1$ times the rate of normal playback. We assume that the user *views* the data as he/she is searching (e.g., fast-forwarding or rewinding) through it, and thus n_{speed} is *finite*. Let T_{np}^i be the mean amount of time that a customer spends in the normal playback mode, before entering some VCR function mode. And, let T_{ff}^i , T_{rw}^i , and T_{pause}^i be the the mean amount of time a customer spends in fast-forward, rewind, and pause modes, respectively, before returning to the normal playback mode. We also assume that the use of VCR functionality (such as fast-forward and rewind) does not require additional service capacity on the part of the CM server. This can be accomplished, for instance, by using techniques proposed in [2].

Note that, in a hybrid system we need to maintain load information on remote nodes and other bookkeeping information, which will require (a relatively small amount) of communication capacity; the exact amount depends on a particular implementation, and we leave these considerations to future work. Note also that in the case of wide data striping, the bookkeeping information must be exchanged between nodes to schedule *each* newly arrived request, whereas in hybrid architectures, we can tradeoff relying more on local (rather than remote) information for some loss in performance. Quantitative assessment of this tradeoff is left for future work.

To assess the *scalability* characteristics of the potential designs in an environment where data access patterns change over time, we consider the following cost/performance and reliability metrics: (1) the system’s *acceptance rate*, in the no queueing case, which is defined as the percentage of all arriving customer requests that are accepted by the system; (2) the mean and variance of the waiting time in the queueing case, i.e., the mean amount of time a customer spends in the queue before he/she is served (the variance is defined accordingly and is used as a measure of QoS); (3) the capacity of the global switch, the capacities as well as the number of local switches, the amount of disk storage, all required to support a particular architecture and corresponding acceptance rate; (4) the mean time to failure (MTTF) of a particular architecture. Table 1 summarizes the main notation used in this paper.

3 Dynamic Replication

Since the number of copies of object i *partly* determines the amount of resources available for servicing requests for that object, we adjust the number of replicas maintained by the system *dynamically*. The system’s performance depends on its ability to make such adjustments *rapidly* and *accurately*.

3.1 General Approach and Main Tradeoff

Given the distributed server described in Section 2, such a dynamic replication approach gives rise to several interesting design issues, including: (a) *when* is the right time for the system to reconfigure the number of replicas, i.e., when to create an additional copy of an object and when to remove a copy; (b) *to which* node should a (new) replica be added or from *which* node should a no longer (deemed) useful replica be removed; and (c) what are *proper policies* for actually creating a new replica (or removing a no longer useful one). In the remainder of this paper, we discuss and evaluate techniques that address these

| | |
|----------------|---|
| S | set of all nodes in the system |
| N | number of nodes in the system; $N = S $ |
| K | number of distinct objects in the system |
| B_x | maximum service capacity of node x (in streams) |
| \bar{B} | average service capacity of the nodes in the system (in streams) |
| C_x | maximum storage capacity of node x (in streams) |
| $L_x(t)$ | load on node x at time t (in streams) |
| $A_i(t)$ | available service capacity for object i at time t ; $A_i(x) = \sum_{x \in R_i(t)} (B_x - L_x(t))$ |
| $ReTh_i$ | replication threshold, i.e., the threshold for adding another copy of object i |
| $DeTh_i$ | de-replication threshold, i.e., the threshold for removing a copy of object i |
| D_i | difference between the replication and the de-replication thresholds, i.e., $D_i = ReTh_i - DeTh_i$ |
| T_{length}^i | length of object i |
| λ | average arrival rate to the system |
| $R_i(t)$ | set of replica nodes for object i at time t |
| $p_i(t)$ | probability of an arriving request being for object i at time t |
| $q_j^i(t)$ | probability of a request for object i being made at time t by a client local to node j |
| n_{speed} | ratio between the speed of fast forward (or rewind) and the speed of normal playback |
| T_{np}^i | mean amount of time a customer spends in normal playback when it enters that mode, for obj i |
| T_{ff}^i | mean amount of time a customer spends fast forwarding when it enters that mode, for obj i |
| T_{rw}^i | mean amount of time a customer spends rewinding when it enters that mode, for obj i |
| T_{pause}^i | mean amount of time a customer spends pausing when it enters that mode, for obj i |

Table 1: Summary of notation.

issues in an *efficient* manner. We note that unless otherwise stated, the following discussion is given in the context of the no queueing case³.

In general, a replication process of a CM object has a *source* node (which currently contains a copy of object i) and a *target* node (on which we are placing a new copy of object i). One simple approach to performing the replication is to “inject” a single replication stream into each of the source and target nodes, for reading and writing of the replica, respectively. We refer to this strategy as *sequential* replication. The sequential replication policy results in a relatively small increase in load on the source and target nodes, i.e., equal to the bandwidth requirements of a single user stream. However, such a policy results in a relatively long replication time (i.e., the replication time is equal to the playout time, at the normal display rate, of the object being replicated), and consequently many customers may be rejected/queued during the replication period due to lack of resources for *that* object, i.e., a lack of *other nodes* in the system, that can service requests for *that* object.

Clearly, one approach to reducing the replication time would be to increase the rate at which the replication is performed, i.e., to read (write) the CM object from (to) the source (target) node at M times the rate of a single stream. This, of course, requires M times the bandwidth of a single user stream on both the source and the target nodes. We refer to this strategy as *parallel* replication, i.e., conceptually this is equivalent to using M single streams in parallel to do the replication. Although this approach reduces the replication time, it also creates an additional load on both, the source and the target nodes, which could result in rejection of customer requests, possibly for CM objects *other* than the one being replicated, due to lack of resources on the *source* and/or *target nodes*, which are being used by the replication process. Thus, we essentially have conflicting goals of (a) using as few resources as possible to perform the replication (in order not to interfere with “normal” system operation) while (b) trying to complete the replication process as soon as possible.

³In the interest of isolating performance characteristics of data layout policies and dynamic replication schemes, we first consider the no queueing case. Later, we also consider a simple queueing policy (i.e., FIFO).

3.2 Early Acceptance

In an attempt to reach a compromise between the conflicting goals stated in the previous section, we consider “early acceptance” of customers, where admitted customers are allowed to use *incomplete* replicas (while the replication is in process). That is, once the system completes replication of the first T_{ea}^i time units of a new replica of an object i , it will treat it as a “virtually” complete copy and begin using it in servicing customer requests for object i . We are motivated by the continuous nature of objects, and essentially it is this property that we exploit here. Note that, for ease of presentation, in the remainder of the paper, we measure the amount of replication completed in time units of *normal playback time of that object*, from the beginning of the object, rather than in storage size, e.g., bytes. Furthermore, for simplicity and clarity of exposition of ideas, in the remainder of this section, much of the discussion is in terms of a specific object being replicated, and thus we drop the superscript i from our notation (the meaning should still be clear from the context of the discussion).

The issue that we need to consider is that a user might attempt to access a portion of an incomplete copy which has not been replicated yet, e.g., by fast-forwarding past the replication point. To allow customers to have full use of VCR functionality, we need to determine a “safe” value for T_{ea} . Clearly, one safe value is $T_{ea} = T_{length}$ (full length of the CM object). However, the intuition is that a smaller value of T_{ea} should result in a higher (at least in the “short term”) acceptance rate of customer requests.

In order to lower T_{ea} (and improve performance) we construct a model of user behavior which allows us to compute a “safe” (but lower than T_{length}) value of T_{ea} while still providing the desired QoS (with a high probability).

3.2.1 Deterministic Model

Given that the replication process constructs a new copy of an object, from the beginning to the end of the object (i.e., in a “linear” fashion), and using only knowledge of the ratio between normal playback and fast-forward, i.e., n_{speed} , we can construct a very simple model which will allow us to compute T_{ea} , as illustrated in Figure 2(a). Specifically, if a newly arrived customer is allowed to use an incomplete replica *after* $T_{ea} = T_{length} \frac{(n_{speed}-1)}{n_{speed}}$ time units of the object have been replicated, then he/she will *not* access data beyond the replication point with probability 1. Thus, if we use this *deterministic* model for admission of customers to the new replica, then (under a *sequential* policy, refer to Section 3.1) the “virtual replication completion time” of an object becomes T_{ea} as compared to T_{length} .

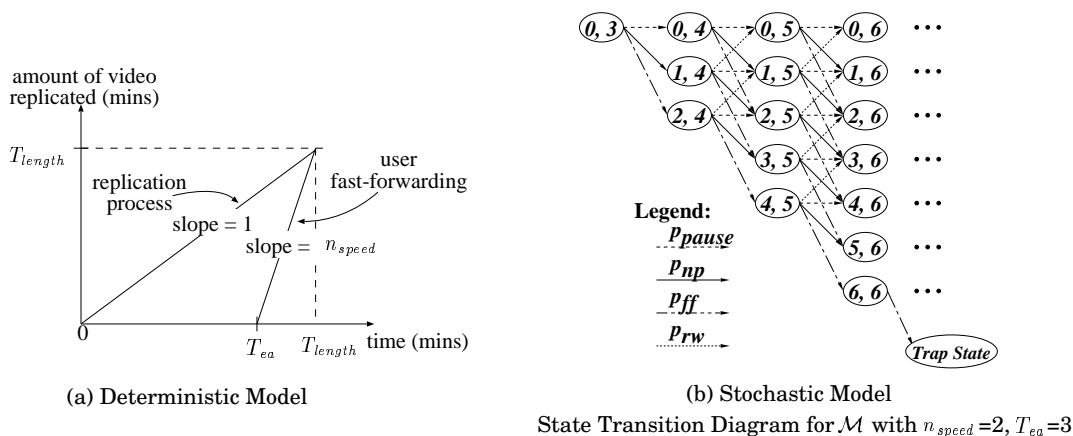


Figure 2: Mathematical Models of User Behavior.

3.2.2 Stochastic Model

To further lower the value of T_{ea} , we employ a stochastic model of user behavior, at the cost of lowering the probability that the user will not access data beyond the replication point (of course, this probability still has to be high, but less than 1). Specifically, we model the combination of the behavior of a user watching a display of a partially replicated object and the corresponding replication process using a Discrete Time Markov Chain (DTMC)⁴, \mathcal{M} , with the following state space \mathcal{S} :

$$\mathcal{S} = \{(V, R) \mid (0 \leq V \leq T_{length}) \wedge (T_{ea} \leq R \leq T_{length}) \wedge (V \leq R) \wedge ((R - V) \leq T_{length} \frac{n_{speed} - 1}{n_{speed}})\} \cup \{(\text{Trap State})\}$$

where V is the current viewing position of the customer and R is the current replication position of the *partial* copy being viewed by that customer.

An example state space for \mathcal{M} with $n_{speed} = 2$ is illustrated in Figure 2(b). There are 4 types of state transitions between adjoining states, which are attributed to: (1) normal playback, (2) fast-forward, (3) rewind, and (4) pause, which occur with probabilities p_{np} , p_{ff} , p_{rw} , and p_{pause} , respectively. A more formal specification of the state transitions in \mathcal{M} with a corresponding one step transition matrix, \mathbf{P} , is as follows:

$$\begin{aligned} (V, R) &\longrightarrow (min(V + n_{speed} * t_u, T_{length}), R + t_u) \\ &\text{Prob} = p_{ff} \mathbf{1}\{ ((V + n_{speed} * t_u) \leq (R + t_u)) \wedge ((R - V) < T_{length} \frac{n_{speed} - 1}{n_{speed}}) \wedge (R < T_{length}) \} \\ (V, R) &\longrightarrow (max(V - n_{speed} * t_u, 0), R + t_u) \\ &\text{Prob} = p_{rw} \mathbf{1}\{ (R < T_{length}) \wedge ((R - V) < T_{length} \frac{n_{speed} - 1}{n_{speed}}) \} \\ (V, R) &\longrightarrow (V, R + t_u) \\ &\text{Prob} = p_{pause} \mathbf{1}\{ (R < T_{length}) \wedge ((R - V) < T_{length} \frac{n_{speed} - 1}{n_{speed}}) \} \\ (V, R) &\longrightarrow (V + t_u, R + t_u) \quad \text{Prob} = p_{np} \mathbf{1}\{ (R < T_{length}) \} \\ (V, R) &\longrightarrow (\text{Trap State}) \quad \text{Prob} = p_{ff} \mathbf{1}\{ ((V + n_{speed} * t_u) > (R + t_u)) \wedge (R < T_{length}) \} \\ (V, R) &\longrightarrow (V, R) \quad \text{Prob} = \mathbf{1}\mathbf{1}\{ ((R - V) = T_{length} \frac{n_{speed} - 1}{n_{speed}}) \vee (R = T_{length}) \} \\ (\text{Trap State}) &\longrightarrow (\text{Trap State}) \quad \text{Prob} = 1 \end{aligned}$$

where $\mathbf{1}\{x\}$ is an indicator function (i.e., $\mathbf{1}\{x\} = 1$ if x is true and 0 if x is false), and t_u is the “granularity” of our model, i.e., the number of time units in an object’s display (under normal playback) corresponding to a unit of time⁵ in the DTMC \mathcal{M} . Finally, $p_{ff} = \frac{T_{ff}}{(T_{np} + T_{rw} + T_{pause} + T_{ff})}$, $p_{rw} = \frac{T_{rw}}{(T_{np} + T_{rw} + T_{pause} + T_{ff})}$, $p_{pause} = \frac{T_{pause}}{(T_{np} + T_{rw} + T_{pause} + T_{ff})}$, $p_{np} = 1 - p_{ff} - p_{rw} - p_{pause}$, where T_{ff} , T_{np} , T_{rw} , and T_{pause} are application-dependent model parameters which were defined in Section 2, and the “Trap State” in \mathcal{M} is a state corresponding to $V > R$, which represents a user’s attempt to access data which has not been replicated yet.

Our goal then is to determine a value of T_{ea} for which the probability of entering the “Trap State” *before the time the replication process completes* (i.e., before $R = T_{length}$) is sufficiently low. Or, conversely, given a value of T_{ea} , we need to compute the probability of entering the “Trap State” by time $t_n < T_{length} - T_{ea}$, which can be accomplished through a *transient* analysis of \mathcal{M} [15], i.e., by solving the following set of equations⁶:

$$\boldsymbol{\pi}(t_n) = \boldsymbol{\pi}(0) * \mathbf{P}^{t_n} \quad \text{and} \quad \sum_{j \in \mathcal{S}} \pi_j(t_n) = 1 \quad (1)$$

⁴Although by using a DTMC as our *model* we make the assumption that the amount of time a user spends in a particular playback mode has a memoryless distribution, we show in Section 4 that the performance of the *system* is fairly insensitive to either the *parameters* of the model or to the *distributional* assumptions.

⁵For instance, if the object is a video clip, then a “natural” time unit in its display would be the amount of time corresponding to the normal playback time of a single frame (on the order of $(\frac{1}{30})^{th}$ of a sec). However, in order to maintain a reasonable size of the DTMC state space, we allow t_u to take on larger time scales, e.g., on the order of minutes — essentially, performing (in general, approximate) aggregation of states.

⁶More sophisticated methods for computing transient results exist [15], but are not the subject of this paper.

where $\pi(t_n)$ is the vector of transient state probabilities at time t_n , $\pi(0) = \mathbf{e}_{(0, T_{ea})}$ is the initial state vector which is equal to a row vector of 0's in all components except for a 1 in the component corresponding to state $(0, T_{ea})$. Our interest then is in $\pi_{(\text{Trap State})}(t_n)$, which is the probability that the user will attempt to access data which has not been replicated yet.

Clearly, the efficiency of solving the above solution depends on the size of \mathcal{M} , which is finite but can be quite large. For instance, with $T_{length} = 90$, $t_u = 1$ min, and $n_{speed} = 2$, the size of \mathcal{M} 's state space is on the order of 3000 states. We can trade off computational cost for the system's performance by using higher values of t_u , e.g., for $t_u = 2$ min, the state space can be reduced to approximately 750 states. This modification can result in higher values of T_{ea} , due to a "coarser granularity" of the model, and hence the (potential) loss in the system's performance⁷. We will illustrate in Section 4 that it is not necessary to obtain extremely low values for $\pi_{(\text{Trap State})}(t_n)$ in order to provide a reasonable QoS — this is due to the fact that the model tends to be conservative, especially with higher values of t_u .

In general, this is an acceptable approach since it only needs to be performed rarely, on a per application basis. That is, a set of statistics or measurements corresponding to interactivity characteristics of an application intended to be run on the CM server can be used to compute the model parameters (i.e., T_{np} , T_{ff} , T_{rw} , and T_{pause}), needed to solve \mathcal{M} . We will show in Section 4 that the model is not very sensitive to the accuracy of the input parameters and thus is of reasonably practical use — this is partly due to its conservative nature. Therefore, we expect that the need for "re-solving" of the model with new parameters would be quite rare and occur only after significant changes in the *interactive nature* of an application have been detected.

However, if the state space of \mathcal{M} is still unacceptable or a more "run-time" approach to computing T_{ea} is desirable, instead of increasing the value of t_u , we can reduce the size of the state space by decreasing the amount of information included in the model about the user's behavior. Again, this reduction in computational cost results in more *conservative* estimates of T_{ea} , and thus we would be trading off system's performance for cost of the solution (for reasons similar to the ones stated above). We elaborate on this approach next.

We note that simple Markov chain models of user behavior have been employed in previous works on video servers, e.g., the two state Markov chain in [12]; however, these have been used for a somewhat different purpose and to the best of our knowledge, with interest in *steady state* characteristics only.

3.2.3 Reduction of the Stochastic Model

We can reduce the size of the state space and the number of transitions by not including all the information about user behavior in the DTMC. For instance, the state space and the number of transitions can be reduced by not including explicitly pause and rewind functionalities in the DTMC but rather "grouping" rewind and pause with the normal playback mode. This is still "safe" since pause and rewind can not cause the viewer to access an unreplicated portion of the data. That is, the reduced DTMC, \mathcal{M}^r , would have 2 types of state transitions between adjoining states, which are attributed to: (1) normal playback (with rewind and pause "grouped" here) and (2) fast-forward, which occur with probabilities p_{np}^r and p_{ff}^r , respectively. Similarly, we could have created another DTMC with only one of, pause or rewind, not explicitly included. We omit these variations since they are very similar in form to the one presented in this section. We omit a formal specification due to lack of space and point the interested reader to [4].

⁷In any case, a simple approach to determining the value of T_{ea} would be to solve the model (possibly) multiple times (e.g., using binary search), with different values of T_{ea} , until a desired value for $\pi_{(\text{Trap State})}(t_n)$ is obtained, which, corresponds to the desired QoS to be provided by the system.

Clearly, the cost of the solution will be reduced, as compared to \mathcal{M} , given the reduction in the state space and the number of transitions. For instance, with $T_{length} = 90$, $t_u = 1$ min, and $n_{speed} = 2$, the size of \mathcal{M}^r 's state space is on the order of 500 states (as compared to ≈ 3000 states in \mathcal{M}), where a brute force solution of Equation (1) takes several minutes to compute using MATLAB numerical solutions package on a Sun Ultra-1 machine. As before, we can trade off computational cost for system's performance (and obtain a more conservative solution) by using higher values of t_u , e.g., for $t_u = 2$ min, the size of the state state is reduced to approximately 130 states. In this case, a brute force solution of Equation (1), requires less than one minute.

3.3 Threshold-based Activation Scheme

We use a *threshold-based* approach to trigger object replication and de-replication, both of which are only triggered at customer *arrival* and/or *departure* instances. Threshold-based techniques for reacting to changes in workload are employed often for improving the cost/performance ratio of a system. Here, as in other systems, the main motivation for using a threshold-based scheme is that there is a non-negligible cost for creating or removing a replica⁸, and thus it should be done "sparingly".

Furthermore, in such an environment, having the amount of service capacity proportional to the access probabilities (even if we knew them) would not necessarily insure acceptance of newly arrived customers. An important factor in the performance of the system is the mixture of requests that arrives and is ultimately serviced by the nodes of the CM server. That is, we may reject requests for object i on node j due to an influx of requests for other objects residing on node j , i.e., other than object i .

Thus, in this paper we study dynamic data replication and de-replication techniques which do *not* assume knowledge of access probabilities. Without such information, one simple approach is to increase (decrease) the amount of service capacity allocated to an object when the amount of available resources left in the system to service that object falls below (above) some threshold value.

More formally, when a customer request for object i arrives to the system at time t , replication of object i is initiated if and only if *all* of the following criteria are satisfied: (1) $A_i(t) < ReTh_i$, where $ReTh_i$ is the replication threshold and $A_i(t) = \sum_{x \in R_i(t)} (B_x - L_x(t))$; (2) object i is not currently under replication; (3) there is sufficient available service capacity on the source node; (4) there is sufficient available storage space capacity **and** service capacity on the target node; (5) there is sufficient available service capacity in local switches as well as the global switch (i.e., interconnection network).

In the case of de-replication, it should be performed before the system runs out of storage space. Basically, we do not want to leave this decision until the time the system actually *needs* the space for creating a new replica. This is due to the fact that there might be customers using the copy that we would like to delete, and either we will have to wait for them to complete their display, or we will have to relocate them. "Planning ahead" for removing copies of "cold" objects before the space is actually needed should improve the system's performance.

De-replication is invoked at both the customer request arrival and departure instances. More formally, a replica of object i at node x will be removed at time t if and only if the following conditions are satisfied:

1. $A_i(t) = \max_{j \in S} \{A_j(t) > ReTh_i\}$. The motivation for this condition is that the number of replicas for object i at time t is more than its current workload demand and at this time it has the greatest excess of replicas among all relatively "cold" objects.

⁸Even removal time can be significant, since the copy may be utilized by users that, e.g., would have to be migrated to other nodes, at the time of removal, again due to real-time constraints on data delivery.

2. i has “crossed” the de-replication threshold, i.e.,

$$A_i(t) - (B_x - L_x(t)) - C_{ix}(t) > DeTh_i \quad (2)$$

where $C_{ix}(t)$ denotes the number of customers viewing object i at node x at time t . With the deletion of object i at node x , $A_i(t)$ would be decreased by $(B_x - L_x(t))$. Since a customer viewing object i at node x will have to be migrated to other replica nodes in $R_i(t)$, $A_i(t)$ would be further decreased by $C_{ix}(t)$.

3. $\sum_{x \in S} D_x(t) < D_S$, where D_S is the storage space threshold for activating de-replication.
4. In the case of the Delayed Migration (DM) de-replication policy *only* (see Section 3.4.2), there is an additional condition, namely that $C_{ix}(t)$ must be equal to 0.

To prevent the system from oscillating between replication and de-replication, a difference of D_i is introduced between $ReTh_i$ and $DeTh_i$, i.e., $DeTh_i = ReTh_i + D_i$. That is, we introduce *hysteresis* into the system.

Thus, we use dynamic replication and de-replication techniques which do *not* assume knowledge of access probabilities. To improve threshold estimation in the absence of access statistics, we use the last interarrival time for object i to (coarsely) “approximate” $p_i(t)$ and compute threshold values as follows:

1. For each object i , we record its last request access time, at_i . At arrival time, t , of a request for object i , we compute its latest interarrival time, $(t - at_i)$, and use it as a coarse “approximation” of $p_i(t)$. Whenever a new request for object i arrives, we record at_i and update the thresholds for the object i accordingly.
2. Then, $ReTh_i = \lceil \frac{T_{ea}^i}{t - at_i} \rceil$. That is, the replication threshold, $ReTh_i$, represents the amount of workload, corresponding to requests for object i , that are expected to arrive in the next T_{ea}^i time units, which is the amount of time needed to create a new (virtual) replica of object i (should we deem it necessary). Thus, the motivation for this setting of the replication threshold values is that if we have fewer resources than are estimated to be needed in a period representing the amount of time needed to create another replica (i.e., add resources), then we increase the number of replicas. That is, we attempt to match the available resources with anticipated workload.
3. Lastly, $DeTh_i = ReTh_i + H_i$, where $H_i = \lceil \frac{T_{length}^i}{t - at_i} \rceil$, i.e., we introduce a hysteresis. The motivation for this setting of the hysteresis value is similar to the motivation given above for $ReTh_i$, except that $\frac{T_{length}^i}{t - at_i}$ corresponds to the expected number of requests for object i that might arrive in the next T_{length}^i time units, i.e., during the entire duration (in normal playback) of the display of object i . That is, T_{length}^i corresponds to an estimate of the amount of time that will elapse before some of the currently allocated resources, which can be used to service requests for object i , are released. That is, we can release resources that are in excess of what we anticipate is needed in that time period.

Note that we dynamically adjust threshold values based on minimal amount of information, i.e., the last inter-arrival time of a request for object i . We used static threshold values in [6]; however, extensive simulations showed that the system’s performance (using the metrics described earlier) is significantly improved through the use of dynamic threshold values. The tradeoff is the need to adjust thresholds and the need to collect some information (i.e., the latest inter-arrival time of requests for each object in the system), which we note is fairly small. Hence, in the remainder of this paper, we only consider our replication policies under dynamic threshold values.

3.4 Policies

3.4.1 Selection Policies

Firstly, the choice of a source node for replication of object i is simple: we select the least-loaded node in the set $R_i(t)$. For the target node, we choose the node which has the highest estimated residual service capacity (in streams) and has available storage capacity. More formally, we choose the node x such that $x \notin R_i(t)$ and $L_x(t) = \max_{y \in (S - R_i(t))} \left\{ \frac{B_y - L_y(t)}{1 + \gamma_y(t)} \right\}$, and the remaining storage capacity on x is sufficient for the new replica, where $\gamma_y(t)$ corresponds to the number of replication processes already in progress on node y at time t . Intuitively, such a choice should avoid replication of multiple relatively popular objects on the same target node (which may later compete for that node's capacity).

3.4.2 Replication Policies

We now describe the replication policies. Recall that T_{ea}^i , as determined in Section 3.2, represents the initial amount of data that must be replicated before customers are allowed to use a *partial* replica, and it is measured in units of *normal playback time* of that object. Thus, the *value* of T_{ea}^i is *independent* of the replication policy used, but how *long* it takes to copy T_{ea}^i time units worth of data is a function of the replication policy. For example, if replication proceeds at the same rate as playback (as in the sequential policies below), then the replication time will be equal to T_{ea}^i , but if replication proceeds at a faster rate (as in the parallel policies below), then the replication time will be smaller than T_{ea}^i .

Sequential Replication (SR): The replication is performed “sequentially” (as described in Section 2), i.e., the system replicates at the rate of normal playback of a single stream by injecting a single read stream at the source node and a single write stream at the target node — each of these requires the same capacity as a single user stream. Thus replication of object i takes T_{length}^i time units, and users are not admitted to the new replica until the *entire* copy is complete. This policy is considered for comparison purposes *only*.

Sequential Replication + Early Acceptance (SREA): The replication is performed as in the SR policy, except that newly arrived users can be admitted to the new (incomplete) replica as soon as T_{ea}^i time units of that object have been replicated on the target node. Furthermore, this “virtual” replication completion time is used in checking the satisfaction of condition (2) in the decision of when to create a new replica (see Section 3.3).

Parallel Replication (PR): The system replicates at M times the rate of a normal display of a single user stream, where $M = \min((B_{source} - L_{source}(t)), (B_{target} - L_{target}(t)))$ at time t , when replication begins. Thus the “real” replication time of object i is reduced to $\frac{T_{length}^i}{M}$, and users are not admitted to the new replica until the *entire* copy is complete. This policy is considered in order to show a contrast in performance between policies that do and do not utilize the *early acceptance* technique.

Parallel Replication + Early Acceptance (PREA): The replication is performed in the same manner as in the PR policy, except that users are admitted to the new (incomplete) replica after the first T_{ea}^i time units of the replica are completed. Furthermore, as in the case of the SREA policy, this “virtual” replication completion time is used in checking the satisfaction of condition (2) in the decision of when to create a new replica.

Mixed Parallel, Early Acceptance + Sequential Replication (MPEA): The first T_{ea}^i time units of the object are replicated as in policy PREA and the remainder of the object is replicated as in policy SREA. Users are admitted to the new (incomplete) replica after the first T_{ea}^i time units of the replica are complete. And, as in the other policies using early acceptance, the “virtual” replication completion time

is used in checking condition (2) in the decision of when to create a new replica.

3.4.3 De-replication Policies

The decision process of *which* replica to remove will be described in Section 3.3. What remains to determine is the choice of the node from which to remove it. Part of the difficulty is in considering the customers that would have to be migrated from the node where the removal occurs. We consider the following de-replication policies.

Delayed Migration (DM): This policy removes a replica of object i only after the last customer finishes viewing the movie. That is, we only remove the replica of object i at node x at time t when $C_{ix}(t) = 0$ in Equation (2). No new users are admitted to this replica after the de-replication decision is made. This is motivated by the (possible) implementation complexity of migrating customers from one node to another.

Immediate Migration Minimum Overhead (IMMO): This policy chooses the node on which fewest customers are currently viewing object i . The motivation here is to reduce the (possible) system overheads associated with user migration. That is, at time t the replica of object i is removed from node x where $C_{ix} = \min_{y \in R_i(t)} \{C_{iy}(t)\}$ in Equation (2); The C_{ix} customers are distributed evenly among the *remaining* nodes in $R_i(t)$.

Immediate Migration Maximum Capacity (IMMC): This policy selects the node which could provide the greatest estimated (residual) service capacity after the replica of object i is removed. That is, at time t the replica of object i is removed from node x where $C_{ix} = \max_{y \in R_i(t)} \{C_{iy}(t) + (B_y - L_y(t))\}$ in Equation (2); the C_{ix} customers are distributed evenly among the *remaining* nodes in $R_i(t)$.

4 Discussion of Results

In this section, we first present a scalability study of CM server designs in the context of data placement techniques where the main concern is the system’s load balancing characteristics and the subsequent system performance. Then, we focus on a performance study of dynamic replication schemes, used in conjunction with the data placement techniques, with emphasis on their sensitivity to user model parameters, workload characteristics, and skewness of data access patterns, as well as applicability to various CM applications. (Table 2 lists parameters along with their default values/distributions and alternatives as used in the remainder of this section. All values are given in units of minutes, unless otherwise specified.)

4.1 Scalability Study

In this section we present results of our simulation study using the cost/performance and reliability metrics given in Section 2. The arrival process (of requests for objects) is Poisson with a mean arrival rate of $\lambda = \frac{a\bar{B}N}{T_{length}^i}$, where $0 \leq a \leq 1$ is the “relative arrival rate”. For ease of presentation, we discuss the results in terms of a , i.e., relative to the *total* service capacity of the system (e.g., $a = 1.0$ corresponds to the maximum service capacity of the system). We consider the design of a CM server with the following capacity requirements: (1) a total service capacity of $N * \bar{B} = 1600$ streams; (2) a total storage capacity of $K = 400$ distinct objects; and (3) each object is of length $T_{length}^i = 90$ minutes.

Since the main motivation for using *dynamic* replication policies is the need to react to changes in data access patterns, we consider the performance of these policies as a function of such changes. That is, the workload will have the characteristic that every “rotation time period” of X minutes, $p_i(t)$ s change. One change in access probabilities is described by Equation (3), which is intended to emulate a relatively

| Parameter | Default | Alternatives |
|--|---|---|
| Arrival process | Poisson with $a = 1.0$ | (1) Poisson with $a = 0.8$ (2) “time of day” based Poisson with $a = 0.9$ for 7 hrs, $a = 0.5$ for 17 hrs (3) on-off source with $a = 0.8$, $\frac{1}{\beta} = \frac{1}{2}$ (min), $\frac{1}{\alpha} = 6$ (min) (4) on-off source with $a = 1.0$, $\frac{1}{\beta} = \frac{1}{2}$ (min), $\frac{1}{\alpha} = 6$ (min) |
| User Behavior Model (used in computing T_{ea}^i) | Stochastic with no reduction in state space (DTMC \mathcal{M} with $\pi_{\text{Trap State}}(t_n) = 0.1$) | |
| T_{length}^i | 90 $\forall i$ | 10 $\forall i$ |
| Playback Mode distribution (NP,FF,RW,PAUSE) | Uniform [0.95 \times mean, 1.05 \times mean] | |
| Interactivity Parameters | NP:FF:RW:PAUSE = 19 : 1 : 1 : 1 $T_{length}^i = 90$, $n_{speed} = 4$ $T_{np}^i = 9.5$, $T_{ff}^i = 0.5$, $T_{rw}^i = 0.5$, $T_{pause}^i = 0.5$ $T_{ea}^i = 12$ | (1) NP:FF:RW:PAUSE=19 : 1 : 1 : 1 $T_{length}^i = 10$, $n_{speed} = 4$ $T_{np}^i = 1.9$, $T_{ff}^i = 0.1$, $T_{rw}^i = 0.1$, $T_{pause}^i = 0.1$ $T_{ea}^i = 3$ (2) NP:FF:RW:PAUSE=4 : 1 : 1 : 1 $T_{length}^i = 90$, $n_{speed} = 4$ $T_{np}^i = 2$, $T_{ff}^i = 0.5$, $T_{rw}^i = 0.5$, $T_{pause}^i = 0.5$ $T_{ea}^i = 18$ (3) NP:FF:RW:PAUSE=4 : 1 : 1 : 1 $T_{length}^i = 10$, $n_{speed} = 4$ $T_{np}^i = 2$, $T_{ff}^i = 0.5$, $T_{rw}^i = 0.5$, $T_{pause}^i = 0.5$ $T_{ea}^i = 3$ |
| Replication policy | SREA | PR, MPEA, PREA |
| De-replication policy | DM | |
| Access Probability change | “gradual” | “abrupt” |
| Skewness distribution | Zipf, $\theta = 0.0$ | Geometric, $\chi = 0.618$ |
| $q_j^i(t)$ | uniformly distributed between 1 and N , for each object i , $\forall t \geq 0$ | |
| D_S | 5 | |
| Architecture | (1) arch2.0: $B_x = 80 \forall x$, $C_x = 28 \forall x$, $N = 20$ | (2) arch1.0: $B_x = 20 \forall x$, $C_x = 7 \forall x$, $N = 80$ (3) arch2 group, (4) arch3 group, (5) arch4 group, (6) arch5 group, (7) arch2w (see Table 3) (8) heterogeneous arch (see Section 4.1.3) |
| K | 400 | |
| Request Queue Size | zero | infinite |
| Network constraint | $nc = 1.0$ | $nc = 0.1, \dots, 0.9$ |
| Rotation time period | 200 (min) | 50,100,400,600,800,1000,1200 (min) |

Table 2: Parameters.

“gradual” increase/decrease in popularities:

$$p_i(t') = \begin{cases} p_{i+2}(t) & \text{if } i \text{ is odd \& } 1 \leq i < K - 1 \\ p_K(t) & \text{if } i \text{ is odd \& } i = K - 1 \\ p_{i-2}(t) & \text{if } i \text{ is even \& } 2 < i \leq K \\ p_1(t) & \text{if } i \text{ is even \& } i = 2 \end{cases} \quad (3)$$

where t and t' refer to two consecutive rotation periods and for ease of presentation we assume that K is even. This is to illustrate that even under a relatively gradual change, dynamic policies are still useful. Furthermore, we believe this is a reasonable “emulation” of change in access patterns for many CM applications.

To test our policies further, we also emulate an “abrupt” increase in popularity of currently unpopular objects as well as a “gradual” decrease in popularity of the currently more popular objects as follows:

$$p_i(t') = \begin{cases} p_1(t) & \text{if } i = K \\ p_{i+1}(t) & \text{if } 1 \leq i \leq K - 1 \end{cases} \quad (4)$$

Here, we consider the Zipf distribution, given in Equation (5), for skewness of access probabilities.

$$\text{Prob}[\text{request for object } i] = \frac{c}{i^{(1-\theta)}} \quad \forall i = 1, 2, \dots, K \quad \text{and} \quad 0 \leq \theta \leq 1 \quad (5)$$

where $c = \frac{1}{H_K^{(1-\theta)}}$ and $H_K^{(1-\theta)} = \sum_{j=1}^K \frac{1}{j^{(1-\theta)}}$. We set $\theta = 0.0$, which corresponds to the *measurements* performed in [3] (for a movies-on-demand application). In Section 4.2 we also consider a finite geometric distribution, for skewness of access probabilities.

The architectural settings considered in this section are the default parameters of Table 2 together with Table 3⁹. Here arch2w corresponds to wide data striping, where a single copy of each object is striped across all nodes of the system, and arch2 through arch5 groups correspond to various configurations of a *hybrid* CM server. For the hybrid architectures we experiment with different amounts of per node storage space capacity, in order to illustrate the tradeoff between storage space capacity local to a node and the corresponding required capacity of the global switch.

| Arch type | No. of nodes | Srv cap/node Lcl switch cap (in streams) | Stor space/node (in objects) |
|-------------|--------------|--|---------------------------------|
| arch2w | 20 | 80 | 20 |
| arch2 group | 20 | 80 | 22; 24; 26; 28; 30 |
| arch3 group | 10 | 160 | 44; 48; 52; 56; 60 |
| arch4 group | 5 | 320 | 88; 92; 96; 100; 104 |
| arch5 group | 2 | 800 | 205;210;215;220;225 |

Table 3: Parameters for architecture groups.

Moreover, we consider the affect on the overall system performance of limitations of communication network resources. Let nc represent the ratio of the global switch and the storage system service capacities, i.e., $nc = 1.0$ represents equal service capacities in the storage and communication sub-systems. Then we vary the service capacity of the global switch, $0.1 \leq nc \leq 1$, and compute the subsequent degradation in performance experienced by the various architectures. The motivation for these experiments is to: (1) observe the performance degradation characteristics of possible CM server designs (as this is an indication of their scalability), and (2) assess whether reduction in overall required global switch capacity (which should lead to lower costs) is possible without significant loss in the overall system performance.

Lastly, below “upper bound” on the acceptance rate refers to the acceptance rate that a wide data striping system can achieve *without considering network capacity constraints*; thus this is the only curve in the following figures that is *not* a function of nc .

4.1.1 Wide Data Striping System vs. Hybrid System

Figures 3 and 4(b) illustrate that under lower network capacities, a hybrid system has better overall performance as well as performance degradation characteristics than the wide data striping system. More importantly, the hybrid architecture allows us to tradeoff capacities of the various system resources in order to achieve a more cost-effective system overall. Specifically, we can tradeoff local storage space and local switch capacities with global switch capacity and achieve nearly the same performance characteristics. For instance, for a particular architectural setting, the larger the local storage space capacity is, the smaller the global switch capacity need be, in order to achieve the same overall system performance, e.g., consider the “arch2 group” in Figure 3(a) — in the case of the 24 objects/node architecture, the corresponding

⁹For a hybrid system that requires more storage space than the corresponding wide data striping system we only increase the storage *space* per disk, *not* the number of disks, as that would also increase the service capacity and hence would not make for a fair comparison.

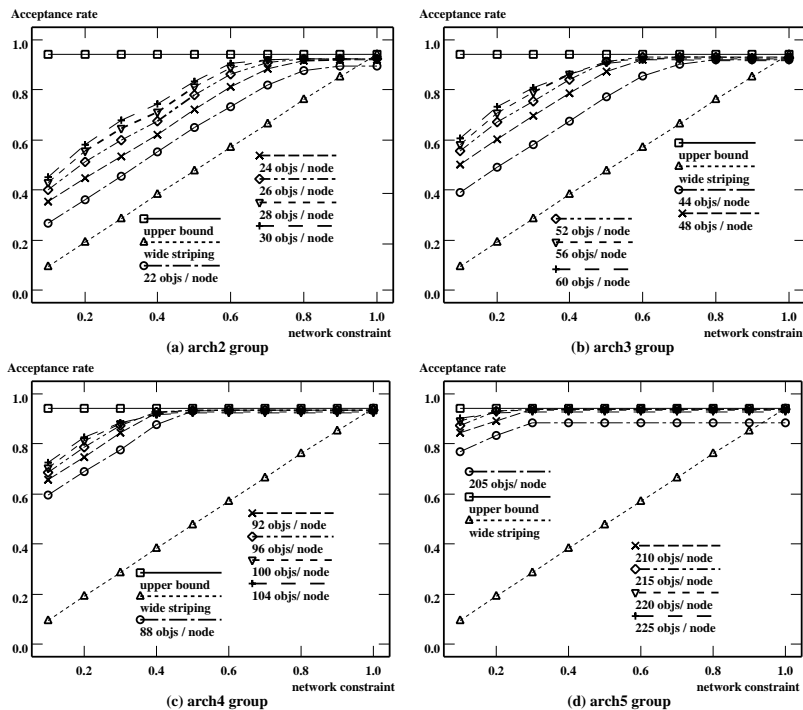


Figure 3: Different network constraints.

required service capacity of the global switch¹⁰ is 1280 streams, whereas in the case of the 30 objects/node architecture, it is only 960 streams.

Conversely, the larger the local switch is, the more we can reduce the storage space and global switch capacities, e.g., consider the “arch2 group” in Figure 3(a) — in this case with a local switch capacity of 80 streams¹¹, the corresponding required total storage space capacity is 600 objects (i.e., 30 objects/node) and the corresponding required service capacity of the global switch is 960 streams. Consider now the “arch4

¹⁰The needed global switch capacity is determined from Figure 3(a) by first fixing the acceptance rate that we would like to achieve. Here, we fix the required acceptance rate to be *at least* 0.95 * acceptance rate of the “upper bound result”, and then determine, using Figure 3(a), the smallest network constraint, nc , that satisfies that acceptance rate for the appropriate architecture curve; then, the required global switch capacity is $nc * 1600$ streams where 1600 streams corresponds to the maximum required system capacity.

¹¹These values can be determined from Tables 3 and 4.

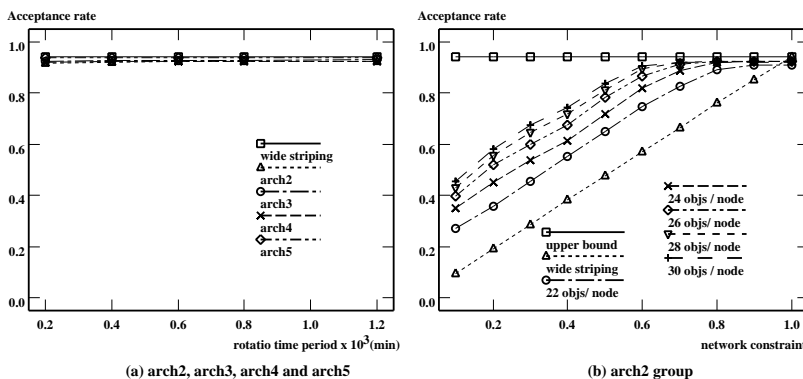


Figure 4: Abrupt increase/gradual decrease.

group” in Figure 3(c) — although in this case the local switch capacity increases to 320 streams, the corresponding required service capacity of the global switch drops down to 640 streams and the corresponding required total storage space capacity drops down to 460 objects (i.e., 92 objects/node). These results are due to the fact that with larger local switch capacities, we can service more customer requests locally (hence the corresponding smaller required global switch capacity). Furthermore, larger local switches and corresponding larger node service capacities also provide more opportunities to take advantage of the load balancing characteristics of striping *within* a node (hence the smaller required storage space capacity per node).

4.1.2 System Sizing Issues

Quantitatively evaluating the tradeoffs between local switch capacity, node storage space capacity, and global switch capacity is no easy task, as it is not immediately clear how to tradeoff one resource for another. Ideally, one would like to evaluate these tradeoffs based on cost. However, cost considerations are a complex issue, given that costs depend on many factors. Thus, next we instead evaluate the different hybrid designs based on the amount of each resource they require *relative* to the wide data striping system. Such an evaluation quantitatively illustrates to the designer the *relative* merits of the different architectures, without the need for choosing a specific technology¹². The purpose of these experiments is to illustrate how a CM server designer can deal with these (fairly complex) *system sizing* issues.

We further refine our test cases in Table 4, and choose the per node storage space and corresponding global switch capacity of each architecture based on the results of the previous section, i.e., we choose those architectures that can achieve an acceptance rate of *at least 0.95** acceptance rate of the “upper bound result” with reasonably small per node storage space and global switch capacities¹³. Figure 5

| Arch type | No of nodes | Srv cap/node Lcl switch cap (in streams) | Stor space per node (in obj's) | Glbl switch cap (in streams) |
|-----------|-------------|--|--------------------------------------|------------------------------------|
| arch2w | 20 | 80 | 20 | 1600 |
| arch2 | 20 | 80 | 26 | 1120 |
| arch2.1 | 20 | 80 | 30 | 960 |
| arch3 | 10 | 160 | 52 | 800 |
| arch4 | 5 | 320 | 92 | 640 |
| arch5 | 2 | 800 | 215 | 320 |
| arch5.1 | 2 | 800 | 225 | 160 |

Table 4: Parameters for architectures used in Section 4.1.

depicts the results of this comparison for each resource as $\frac{\text{resource requirement of arch } i}{\text{resource requirement of arch2w}} \quad \forall i \neq 2w$. Hence, the straight line at the value of 1.0 in each of the graphs of Figure 5 corresponds to the (“scaled”) resource requirement of “arch2w”.

As already stated, these results illustrate to the designer the relative merits of the different architectures by quantifying the tradeoffs between the various resources of the CM server. Based on these results and current costs and technology trends, the designer can make system sizing decisions.

¹²Characterizing a resource using only its capacity may result in a simplification for certain types of resources; however, this is still a good abstraction for evaluating cost-effectiveness of designs, without having to choose a specific technology for each system component.

¹³The upper bound is computed without considering network capacity constraints; since it is not always achievable by an architecture, we choose a performance goal that is reasonably close.

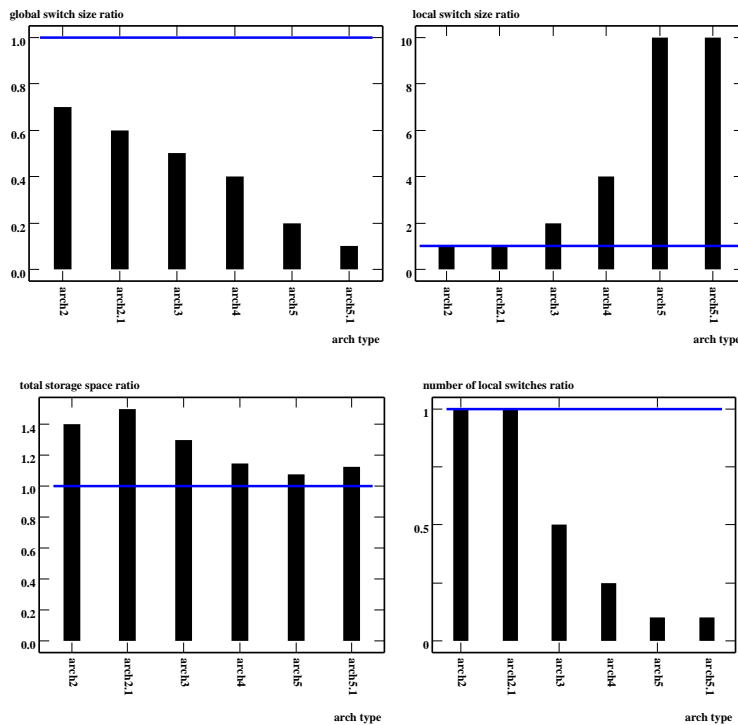


Figure 5: System sizing.

4.1.3 Heterogeneous Systems

Next, we illustrate the ease of dealing with heterogeneous systems when using hybrid CM server designs *without* loss of performance as compared to an equivalent homogeneous case¹⁴. For this purpose, we consider a hybrid CM architecture with 5 nodes and a *total* service capacity of 1600 streams. We use two test cases in the following experiments, both based on the homogeneous version of “arch4” with the storage space capacity of 104 objects per node (refer to Table 3). We introduce 5% and 10% differences in storage space and service capacities between the nodes of the system (as well as corresponding differences in local switch capacities), e.g., to emulate a system that gradually grows (as well as experiences replacements due to failures) and thus is forced to use heterogeneous resources¹⁵. In the results depicted in Figure 6(a) heterogeneous systems achieve performance that is comparable to homogeneous system performance.

4.1.4 Reliability

We use the mean time to failure (MTTF) as our reliability metric, defined as the mean time until some combination of disk failures results in loss of some data that can no longer be recovered through the use of redundant information. We compare the architectures in Table 4, using a *conservative* estimate for the *hybrid* system based on an assumption of only a *single* copy per object. The results are depicted in Figure 6(b) as the following *ratio* $\frac{\text{MTTF of arch } i}{\text{MTTF of arch } 2w} \quad \forall i \neq 2w$, where the straight line at 1.0 corresponds to the (“scaled”) MTTF of wide data striping. The derivations of the MTTF equations used to compute the

¹⁴We expect CM servers to have heterogeneous resources as a result of system growth (due to expansion) and disk failures (and subsequent replacement). We are not aware of wide data striping techniques for systems with heterogeneous resources which do not suffer from loss of either bandwidth or storage resources. Hybrid designs do not suffer from such a loss.

¹⁵We have one test case of a 5 node system, with 84, 94, 104, 114, and 124 objects/node, respectively and service capacities of 256, 288, 320, 352, and 384 streams, respectively; we have another test case of a 5 node system, with 94, 99, 104, 109, and 114 objects/node, respectively and service capacities of 288, 304, 320, 336, and 352 streams, respectively.

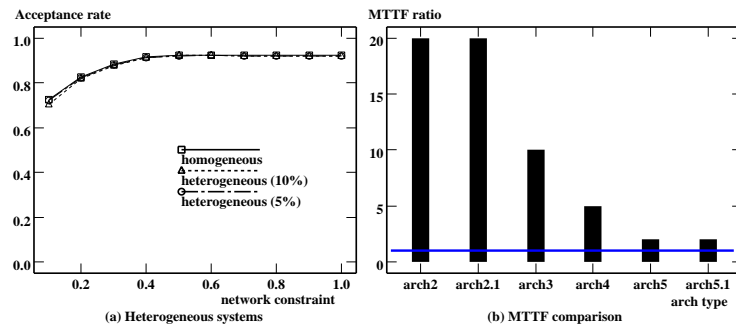


Figure 6: Heterogeneous systems and MTTF.

results in Figure 6(b) are given in [4].

These results clearly show that higher reliability can be achieved by hybrid systems, even for objects that only have a *single* copy, as compared to wide data striping. This increase in reliability is due to the “isolation” of fault affects, i.e., the wider we stripe an object, the more disk failures can affect it. Of course, the reliability is even higher for objects with multiple copies, as is natural in a system which employs data replication. Thus, in a hybrid system, we can provide significantly higher reliability for the *popular* objects, as there will always be multiple replicas of such objects in a hybrid system. Lastly, under network failures or high workload conditions at remote nodes, local nodes can at least deliver *some* objects¹⁶, which is not the case for wide data striping, as all nodes and network capacity must be available in order to serve a request for *any* object.

4.2 Policy Performance Evaluation

In this section we conduct a number of experiments¹⁷ in order to understand the sensitivity of the dynamic replication policies to differences in architectures, user model parameters, workload characteristics, skewness in data access patterns, as well as applicability to different CM applications.

There is a multitude of parameters that can be varied in studying performance of dynamic replication policies. And, we performed a great number of experiments [4]; we present a subset here to illustrates our points. Several entries in Table 2 require a few words of clarification. The default arrival process (of requests for objects) is Poisson¹⁸ with a mean arrival rate of $\lambda = \frac{aBN}{T_{length}^2}$, where $0 \leq a \leq 1$ is the “relative arrival rate”. As before, we discuss the results in terms of the *relative* arrival rate, a , i.e., relative to the *total* service capacity of the system. In order to further explore the benefits of early acceptance schemes, we consider another arrival process, specifically an on-off process which, with proper parameter settings, exhibits more bursty characteristics¹⁹. An on-off process is a two state Markov chain, i.e., one state is an “on” state and the other state is an “off” state. Let β be the transition rate from “on” state to “off” state and α be the transition rate from “off” state to “on” state. When the process is in the “on” state, one customer arrives every (deterministic) time interval $\frac{1}{\lambda}$. There are no arrivals when the process is in the

¹⁶For example, in a movies-on-demand application, even if a requested movie is not available, the user has the option to choose another movie that may be available.

¹⁷Most results presented here are obtained with $95 \pm 5\%$ confidence intervals; all results are within $95 \pm 10\%$ confidence intervals.

¹⁸We believe it is reasonable for us to consider a Poisson arrival process for purposes of this study, since user requests are essentially considered on a “per session” basis; refer to [14] for details that support the use of Poisson arrivals in this case.

¹⁹In our case, the on-off process with $\frac{1}{\beta} = 0.5$ minutes and $\frac{1}{\alpha} = 6$ minutes has the same mean arrival rate, λ , as the corresponding Poisson process. However, its variance is 12λ times the variance of the corresponding Poisson process.

“off” state. In order to make a meaningful comparison, the parameters of the on-off arrival process, i.e., the deterministic arrival rate A and transition rates β and α are set such that the on-off process has the same mean arrival rate as the Poisson process, i.e., $A\frac{\alpha}{\alpha+\beta} = \lambda$, but a higher variance, i.e., $A^2\frac{\alpha\beta}{(\alpha+\beta)^2}$.

At customer arrival and departure instances, as well as at completion of object replication instances, the system first adjusts the threshold values (as described in 3.3). Then, if excess service capacity is available, the system gives priority to replication of objects based on current threshold values where remaining capacity (in the queueing case) is given to currently “serve-able” requests in the queue, in a FIFO manner. A request for object i is “serve-able” at time t if $A_i(t) > 0$ and there is sufficient network capacity to serve this request.

In this performance study, we not only consider the Zipf distribution for the skewness of access probabilities (refer to Section 4.1), but also a finite geometric distribution, given in Equation (6), where we set $\chi = 0.618$.

$$\text{Prob}[\text{request for object } i] = \frac{(\chi)^{(i-1)}(1 - \chi)}{1 - \chi^K} \quad \forall i = 1, 2, \dots, K \quad (6)$$

The motivation being that some applications (other than movies-on-demand) may exhibit higher skewness in data access, e.g., news-on-demand. As we are not aware of measurements available for applications such as news-on-demand, we use a “generically” highly skewed distribution, i.e., the geometric²⁰.

Moreover, the interactivity entry in Table 2 refers to how interactive the users are, with NP: FF: RW: PAUSE referring to the ratio between normal playback (NP) and the various VCR functions (FF, RW, PAUSE/RESUME). These values are used as parameters of \mathcal{M} in the computation of T_{ea}^i . The default values are in agreement with the range of values used in [12]. Unless otherwise stated, in all figures below we use the default values given in Table 2.

Finally, we note that, although the evaluation of the replication policies presented in the remainder of this section is *quantitative*, the main focus of the following discussion is “trends” in the curves and relative performance of the policies, rather than absolute performance. This is due to the fact that our main motivation is to explore the above stated issues and tradeoffs, rather than to predict the (exact) performance of the system through simulation. To this end, we run the simulations at a very *high* load for the no queueing case²¹ in order to illustrate our points (since it almost does not matter what resource management techniques are used at low loads). This is *not* to say that we recommend that the system is *operated* at such high loads; e.g., clearly, under extremely frequent changes in access patterns²² the acceptance rate will be low under very high loads and thus, under such conditions the real system should be operated at lower loads. In the queueing case, we use somewhat lower loads to insure system stability.

4.2.1 Static vs. Dynamic

We first summarize the motivation for using *dynamic* replication policies, as opposed to static ones²³. Based on extensive simulations (refer to [4]), we believe that dynamic replication with early acceptance does result in significantly better performance as compared to a static scheme. With a reasonably large amount of per node resources dynamic policies perform better than the static policy. As the per node resources become fewer, while keeping the same amount of overall system resources, the benefits of the

²⁰Furthermore, applications with relatively little skew in access patterns should not, in a sense, present a performance problem in this case, and thus we do not consider such access patterns here.

²¹There is no stability issue here.

²²We include these for the sake of completeness.

²³The details of a static policy used and the results of the simulations are omitted due to lack of space and can be found in [4]; this comparison is a conservative one.

more conservative dynamic policies are significantly diminished. However, addition of early acceptance mitigates this problem. In general, these results are due to the fact that, dynamic replication schemes make good decisions about: (a) which CM objects are hot, and (b) when to replicate such objects. The problem with the more conservative policies is that the inability to make rapid adjustments in number of replicas is a more severe handicap when resources are few, which is alleviated through the use of early acceptance.

One advantage of the static policy, of course, is that it is easier to implement. Specifically, the need to migrate users from one node to another (in mid-stream) during de-replication may result in complications in the implementation. In the results summarized above as well as in the remainder of the paper, unless otherwise stated, we use the DM de-replication policy in an attempt to make a more fair comparison with the static policy.

4.2.2 Early Acceptance vs. No Early Acceptance

We now motivate the use of *early acceptance* techniques in conjunction with dynamic replication policies. To this end we compare performance of the dynamic policies with and without the use of early acceptance. This comparison is depicted in Figures 7, 8, 9, and 11, where the more important observations are as follows.

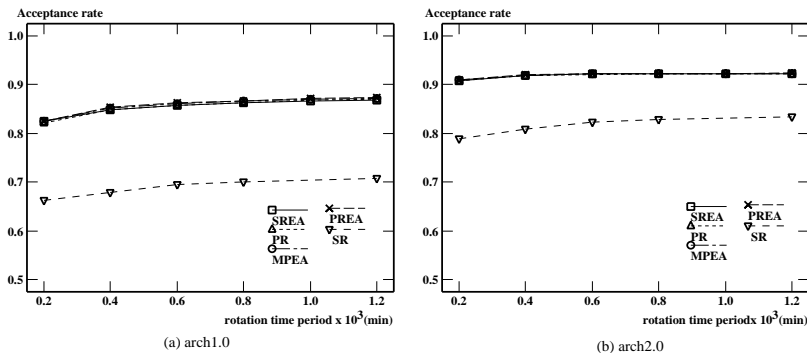


Figure 7: Default settings.

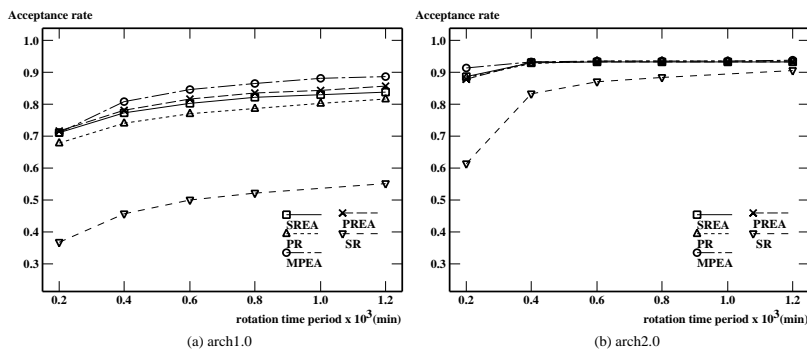


Figure 8: Default settings, with geometric distribution of access patterns and both architectures.

Firstly, based on extensive simulations, we conclude that early acceptance does result in a nice compromise between using resources for performing replication and using resources for servicing customer requests (as stated earlier). This point is best illustrated by considering the more conservative policy with early acceptance, i.e., SREA, and comparing it to the least conservative policy we have without early acceptance,

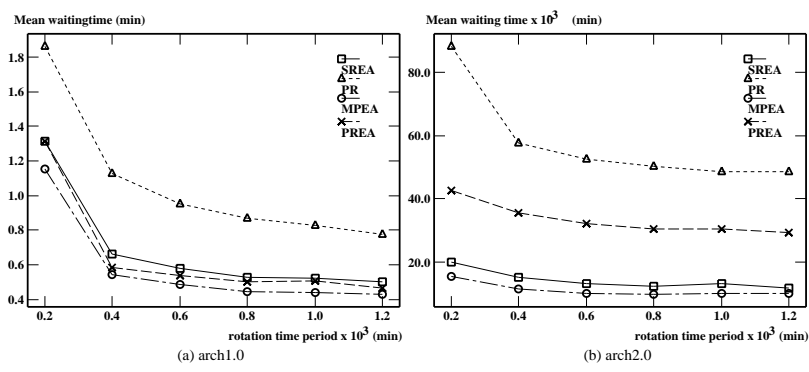


Figure 9: Default settings with $a=0.8$

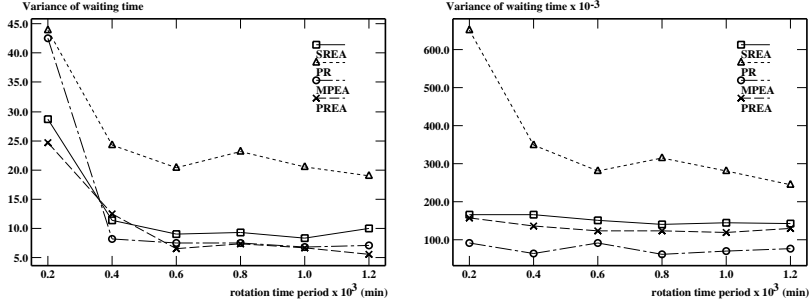


Figure 10: Default settings but with $a=0.8$

i.e., PR. SREA uses as few resources as possible for replication but still makes the new copy available to customers fairly quickly — this policy performs well consistently, i.e., it either results in the best or nearly the best performances in the test cases examined in Figures 7, 8, 9, and 11.

To further explore the effects of early acceptance, we consider (1) more bursty arrival processes, i.e., the on-off arrival process and (2) the queueing case. More specifically, we observe that all policies with early acceptance (such as SREA, MPEA, and PREA) as compared to the PR policy are: (a) less sensitive to the more bursty arrival process (compare Figure 7(b) with Figure 11(b)), and (b) result in better performance in the queueing case²⁴, i.e., smaller mean waiting time (see Figure 9) and better QoS, i.e., smaller waiting time variance (see Figure 10). Note that, in a system with a small mean waiting time but a large waiting

²⁴We already showed that these result in better performance in the no queueing case.

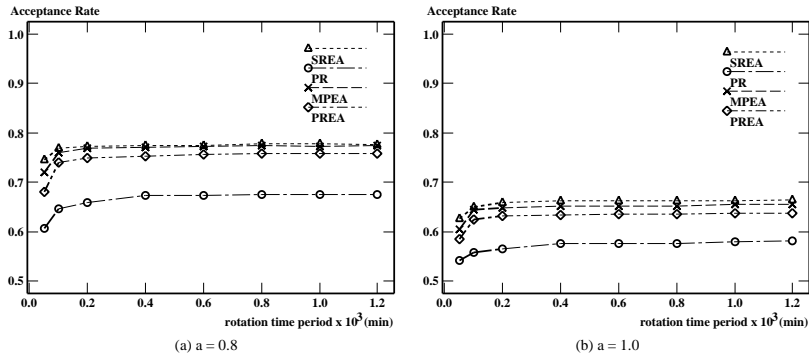


Figure 11: Default setting but with on-off process: $\frac{1}{\beta} = \frac{1}{2}$ (min), and $\frac{1}{\alpha} = 6$ (min)

time variance, it is still possible for a customer to suffer a fairly long waiting time period on occasion; this is clearly undesirable. Hence, we use the variance of the waiting time as our measure of quality of service.

4.2.3 Sensitivity to User Model

Next we show that the mathematical model of user behavior is not very sensitive to the precision of the model parameters (which need to be computed based on statistics or measurements collected about user behavior), and thus it is of reasonably practical use. To validate this conjecture, in our *simulation* we *deviate* on several points from the analytical model. Firstly, in our simulations the distribution of residence times in various user playback modes (NP, FF, RW, PAUSE) is **uniform** as compared to the exponential assumption made in the analytical model²⁵. For all cases where the interactivity model corresponds to NP:FF:RW:PAUSE = 19:1:1:1, the fraction of admitted customers that enter “Trap State” in our simulations is zero — recall that, in our computation of T_{ea}^i we chose $\pi_{\text{Trap State}}(t_n) = 0.1$ (refer to Table 2). This is partly due to the fact that our analytical model tends to be conservative (see Section 3.2)²⁶.

To further “stress test” our model we ran a set of simulations where T_{ff}^i was increased by 20% in the *simulation*, as compared to the parameter used in the *analytical* model. The result is that there is no change in the fraction of admitted customers entering the “Trap State” in the simulation, i.e., it is still zero. This supports our conjecture (made in Section 3.2), that the parameters used in the analytical model do not have to be exact, with respect to the “real” user behavior — that is, fairly large inaccuracies in the collected statistics about the user behavior can be tolerated and consequently the model is reasonable, and “re-solving” of the model with new parameters only needs to be performed “occasionally” (and *not* necessarily in real-time as explained in Section 3.2).

These results are due to: (1) the conservative nature of the analytical model; (2) good dynamic threshold adjustment methods, i.e., the system is able to distinguish between “hot” and “cold” objects sufficiently well; and (3) the fact that the level of interactivity (19:1:1:1) is relatively low (although reasonable for a movies-on-demand application [12]). Thus, in order to further “stress test” the analytical model, we consider a workload with a significantly higher level of interactivity (i.e., alternative (2) for interactivity settings in Table 2) — this may not necessarily correspond to a realistic workload but is useful for purposes of illustration. Here we consider the following observations in simulations: the fraction of admitted customers entering the “Trap State”, the mean amount of time a customer spent in the “Trap State”, given that he/she entered it, and the maximum amount of time a customer spent in the “Trap State”, given that he/she entered it.

Even with such high interactivity levels, the results show that the fraction of admitted customers which enter the “Trap State” in the simulation is quite small (on the order of 10^{-6}). That is, in our experiments, most admitted customers *do not* enter the “Trap State”. Even in the rare case when an admitted customer did enter the “Trap State”, the mean and the maximum time he/she spent in the “Trap State” in the simulation was less than 30 seconds. However, we would like to stress that this rarely occurs.

If we want to further reduce the time a customer may spend in the “Trap State”, some possible solutions include: (1) migration of customers entering the “Trap State” to other nodes which contain a copy of

²⁵We experimented with other distributions as well, e.g., normal; the results were within a few percent of those given here (in the interests of brevity we do not present them).

²⁶We experimented with values of T_{ea}^i that were smaller than those predicted by the analytical model, to determine whether the model’s conservative nature is resulting in some loss of performance — this turned out *not* to be the case in our experiments, that is, beyond a certain value of T_{ea}^i one only obtains diminishing returns in performance gains, but the fraction of accepted customers which enter the “Trap State” continues to grow.

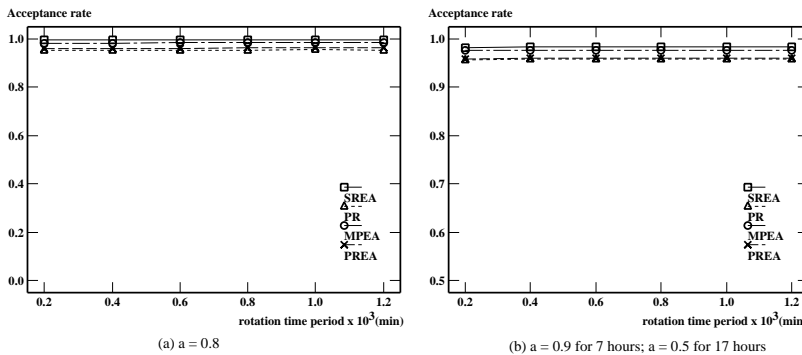


Figure 12: Default settings, but with alternatives (1) and (2) for arrivals.

the object they are viewing and have available service capacity, or (2) increasing T_{ea}^i (at the cost of some performance degradation).

4.2.4 Sensitivity to workload characteristics

Next, we show the lack of sensitivity to the workload characteristics, accomplished through the use of *early acceptance*. To this end we ran a set of simulations with three different modification to the workload characteristics (refer to Table 2), as compared to the default workload used thus far (i.e., as compared to a Poisson process with a constant rate and $a = 1.0$): (1) “time of day” based workload, which is still Poisson but with arrival rates based on time of day, i.e., with $a = 0.9$ over 7 hours of a day and $a = 0.5$ over the remaining 17 hours; (2) lower workloads, i.e., still Poisson with a constant arrival rate but with $a = 0.8$; and (3) more bursty workloads, i.e., an on-off arrival process with $a = 0.8$, $\frac{1}{\beta} = 0.5$ min and $\frac{1}{\alpha} = 6$ min. The results are depicted in Figures 11 and 12. More specifically, we observe that all policies with early acceptance (such as SREA, MPEA, and PREA) as compared to the PR policy are less sensitive to changes of workload characteristics (e.g., compare Figure 11(a) with Figure 12(a)).

4.2.5 Applicability to a variety of applications

Next, we consider applicability of dynamic replication with early acceptance to a range of applications of CM servers. To this end, we ran a set of simulations with: (a) smaller objects, i.e., shorter clips with $T_{length}^i = 10$ min $\forall i$, as well as *in addition* (to smaller clips) (b) higher levels of interactivity, NP:FF:RW:PAUSE=4:1:1:1; these cases correspond to interactivity alternatives (1) and (3) in Table 2. The results for cases (a) and (b) are illustrated in Figures 13 and 14, respectively. Qualitatively, all the conclusions made above, in the context of $T_{length}^i = 90$, still hold²⁷. We also considered the sensitivity of the mathematical model of user behavior to the new workload settings; the results are similar to the ones stated above (refer to [4] for details).

5 Conclusions

In this work, we focused on *scalability* of large CM end-to-end server designs (in the context of data placement issues) as a function of their cost/performance and reliability characteristics under various workloads and system constraints. We have presented a performance study where the main observations are as follows:

²⁷We observe that there is a significant reduction in the acceptance rate in Figure 14 as compared to Figure 13. This is due to the fact that with higher interactivity levels each customer stays in the system longer, on the average. That is, with higher interactivity levels each customer holds on to the resources allocated to him/her longer.

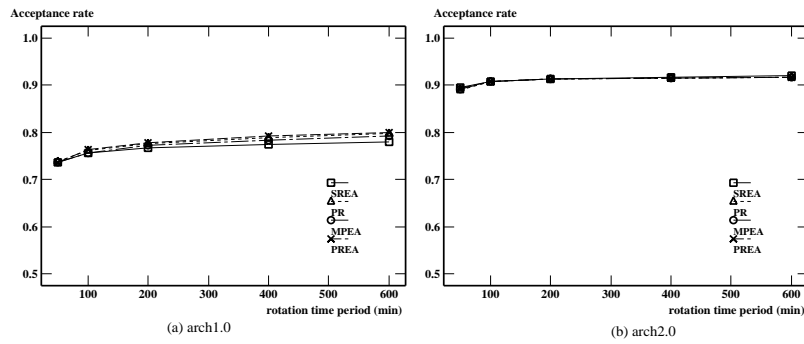


Figure 13: Default settings with alternative (1) for interactivity settings.

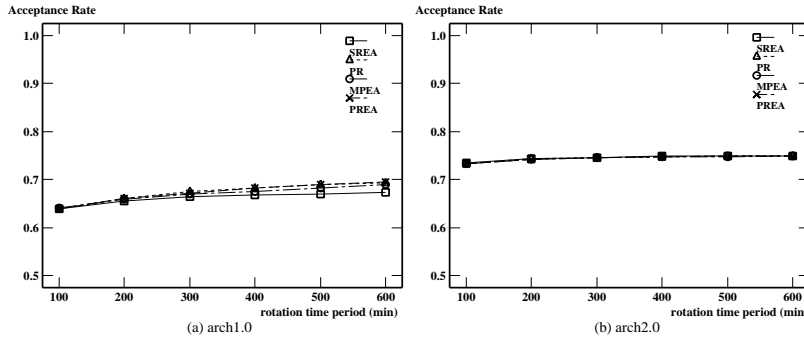


Figure 14: Default settings with alternative (3) for interactivity settings.

- The use of hybrid designs allows us to tradeoff resources and specifically we can tradeoff a node's local storage space, its local switch capacity, and the system's global switch capacity given a particular performance requirement. This should allow a system designer to make better system sizing decisions by making appropriate tradeoffs, i.e., it should allow for cost-effective designs while satisfying the required system performance.
- Hybrid designs naturally extend to systems with heterogeneous resources with little or no loss in performance as compared to the corresponding homogeneous systems. Since in practice we expect CM servers to have heterogeneous resources as a result of system growth (due to expansion) and disk failures (and subsequent replacement), this is an important characteristic.
- Hybrid designs result in higher system reliability even under the conservative assumption of have a single copy per object. Better reliability characteristics are of significant importance to CM servers since they service real-time workload and furthermore there is a need to recover from failure in real-time [1].
- Dynamic replication schemes in conjunction with early acceptance provide a good compromise between using resources for performing replication and using resources for serving customers which in turn results in good performance characteristics across a wide range of workloads, skewness data access patterns and frequency of changes in data access patterns.
- The mathematical model of user behavior, which is used to allow early acceptance in a "safe" manner, is not very sensitive to the precision of the model parameters, as illustrated by testing over a wide range of workloads as well as by using different assumptions in the simulations as compared to the model. Thus, this approach is of reasonably practical use.
- Dynamic adjustment of threshold values, in conjunction with dynamic replication and early accep-

tance, results in systems that are able to make a reasonably good distinction between “hot” and “cold” objects. This in turn results in improved system performance.

We believe that this suggests the usefulness of these techniques across a wide range of continuous media applications.

References

- [1] S. Berson, L. Golubchik, and R. R. Muntz. Fault Tolerant Design of Multimedia Servers. In *Proc. of the ACM SIGMOD Conf. on Management of Data*, pages 364–375, San Jose, CA, May 1995.
- [2] M.-S. Chen, D. D. Kandlur, and P. S. Yu. Support for Fully Interactive Playout in a Disk-Array-Based Video Server. *Proceedings of the 2nd ACM Intl. Conf. on Multimedia*, pages 391–398, October 1994.
- [3] A. L. Chervenak. Tertiary Storage: An Evaluation of New Applications. *Ph.D. Thesis, UC Berkeley*, 1994.
- [4] C. Chou, L. Golubchik, J. C.S. Lui, and I-Hsin Chung. Design of Scalable Continuous Media Servers with Dynamic Replication. Technical Report CS-TR-4232, University of Maryland, March 2001.
- [5] C.-F. Chou, L. Golubchik, and J. C.S. Lui. Striping doesn’t scale: How to achieve scalability for continuous media servers with replication. In *Proceedings of the International Conference on Distributed Computing Systems (ICDCS)*, Taipei, Taiwan, April 2000.
- [6] C.-F. Chou, L. Golubchik, and J. C.S. Lui. A performance study of dynamic replication techniques in continuous media servers. In *International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems*, San Francisco, CA, August 2000.
- [7] A. Dan, M. Kienzle, and D. Sitaram. A Dynamic Policy of Segment Replication for Load-Balancing in Video-on-Demand Servers. *ACM Multimedia Systems*, 3:93–103, 1995.
- [8] A. Dan and D. Sitaram. An Online Video Placement Policy Based on Bandwidth to Space Ratio (BSR). In *Proceedings of ACM SIGMOD’95*, 1995.
- [9] Bolosky et al. The Tiger Video Fileserver. Technical Report MSR-TR-96-09, Microsoft Research, 1996.
- [10] S. Ghandeharizadeh and R. R. Muntz. Design and Implementation of Scalable Continuous Media Servers. *Parallel Computing Journal*, pages 123–155, January 1998.
- [11] R.L. Haskin. Tiger Shark : A Scalable File System for Multimedia. Technical report, IBM Research, 1996.
- [12] K.D. Jayanta, J.D. Salehi, J.F. Kurose, and D. Towsley. Providing VCR Capacities in Large-Scale Video Servers. In *Proc. ACM Intl. Conf. on Multimedia*, pages 25–32, 1994.
- [13] P. W. K. Lie, J. C.-S. Lui, and L. Golubchik. Threshold-Based Dynamic Replication in Large-Scale Video-on-Demand Systems. *Multimedia Tools and Applications*, 11(1):35–62, 2000.
- [14] V. Paxson and S. Floyd. Wide-Area Traffic: The Failure of Poisson Modeling. *IEEE/ACM Transactions on Networking*, 3(3):266–244, June 1995.
- [15] W. J. Stewart. *Introduction to Numerical Solution of Markov Chains*. Princeton University Press, 1994.
- [16] N. Venkatasubramanian and S. Ramanathan. Load Management in Distributed Video Servers. In *Proceedings of ICDCS*, pages 528–535, Baltimore, MD, May 1997.
- [17] J. Wolf, H. Shachnai, and P. Yu. DASD Dancing: A Disk Load Balancing Optimization Scheme for Video-on-Demand Computer Systems. In *ACM SIGMETRICS/Performance Conf.*, 1995.