

Silo, Rainbow, and Caching Token: Schemes for Scalable, Fault Tolerant Stream Caching

Youngsu Chae, Katherine Guo, Milind M. Buddhikot, Subhash Suri, and Ellen W. Zegura

Abstract— In the current Internet, web content is increasingly being cached closer to the end-user to reduce network and web server load and improve performance. Existing web caching systems [1] typically cache entire web documents and attempt to keep them consistent with the origin server. This approach works well for text and images; for bandwidth intensive multimedia data such as audio and video, caching entire documents is not cost effective and does not scale. An alternative approach is to cache parts of the multimedia stream on different caches in the network and coordinate stream playback from these independent caches. From the perspective of the clients, the collection of cooperating distributed caches act as a single fault tolerant, scalable cache [4]. In this paper, we focus on data placement and replacement techniques for such cooperating distributed caches. Specifically, we propose the following new schemes that work together: (1) A family of distributed layouts, consisting of two layouts, namely RCache and Silo. The RCache layout is a simple, randomized, easy-to-implement layout that distributes constant length segments of a clip among caches and provides modest storage efficiency. The Silo scheme improves upon RCache; it accounts for long term clip popularity and intra-clip segment popularity metrics and provides parameters to tune storage efficiency, server load, and playback switch-overs; (2) Rainbow, a local data replacement scheme based on the concept of segment access potential that accurately captures the popularity metrics. (3) Caching Token, a dynamic global data replacement or redistribution scheme that exploits existing data in distributed caches to minimize data distribution overhead. Our schemes optimize storage space, start-up latency, server load, network bandwidth usage, and overhead from playback switch-overs. Our analytical and simulation results show that the Silo scheme provides 3 - 8 times higher cache hit ratio than a comparable traditional web caching system that has the same amount of storage space.

Keywords— multimedia streaming, caching, VoD, cache replacement.

I. INTRODUCTION

In recent years, an exponential increase in the number of Internet users has strained the Internet infrastructure. Scalable, cluster-based web servers and smart data caching have been two of the many tools used to alleviate this problem. However, with the proliferation of new Internet services that use multimedia data such as audio and video, this

performance problem will worsen. Extensive research on multimedia streaming services and servers has attempted to address this problem and devise solutions that scale in terms of number of users and stream bandwidth. This research includes: (1) Scalable, high performance server and file system or operating system architectures [2], [8], [12], [24]. (2) Network and application level multicast techniques [6], [7], [18], [19], [20], [22], [28] that trade bandwidth and buffer space for interactivity and personalization. (3) Limited caching schemes such as stream patching [17] and prefix caching [23] that improve latency and playback quality perceived by the end-user.

Recent investigation of audio and video data on the Internet and in large test-beds [3], [5] reveals some important properties: (1) Bandwidth, duration, and therefore size of video files are gradually increasing. (2) Audio and video files are WORMS, that is, Write Once Read Many files; once created they seldom change and have a long life. (3) Access to such data exhibits high temporal locality which means over a period of time, large fraction of received requests are for same or similar objects. These properties suggest that in addition to the techniques mentioned above, efficient in-network caching of these multimedia objects can reduce network and server load, and improve start-up latency and quality perceived by the end-user.

Existing web caching systems [1] are typically stand-alone systems that independently cache entire multimedia clips in response to user requests. This approach of demand-driven replication of same content in different caches is only adequate for small sized clips and does not work well for the large, high bandwidth files that are becoming common [5]. For example, a single, two-hour long MPEG-2 movie requires about 4.5 GB of disk space. Such large objects pose two problems: (1) Replicating the movie at different caches leads to space wastage. (2) For a fixed investment in storage space at each cache, only a few media objects can be stored, and therefore, the hit ratio and the efficiency of the cache is limited.

An incremental improvement to this approach is *selective replication* which distributes an entire clip to a selected number of proxies in the network. However, optimal replication is difficult and can lead to severe load-imbalance, high server load, and storage inefficiency.

Manuscript received May 8, 2001; revised Jan. 15, 2002.

Y. Chae and E. Zegura are with the College of Computing, Georgia Institute of Technology, Atlanta, GA 30332 USA ({yschae,ewz}@cc.gatech.edu).

K. Guo and M. Buddhikot are with the Center for Networking Research, Bell Laboratories, Lucent Technologies, Holmdel, NJ 07733 USA ({kguo,milind}@dnrc.bell-labs.com).

S. Suri is with the Department of Computer Science, University of California, Santa Barbara, CA USA (suri@cs.ucsb.edu).

An alternate solution is to break media objects into smaller segments and distribute one or more copies of each segment in a system of caches connected by a storage, local or metro area network. This idea of data striping has been used widely in disk arrays such as RAIDs [16], and distributed cluster-based video-on-demand (VoD) servers [2], [8], [9], [12]. Such a scheme can reduce the required storage space compared to replication schemes and provide the added advantages of load-balance, fault tolerance, large parallelism, high cache hit ratio, and reduced network and server load. The scheme used for segment layout and replacement directly affects the efficiency of the caching system. In this paper, we focus on design and analysis of such a distributed caching system.

A. Research Contributions

In this paper we propose a novel distributed caching architecture and associated schemes that work together. Specifically, we describe the following: (1) A family of distributed layouts, consisting of two layouts, namely *RCache* and *Silo*. The *RCache* layout is a simple, randomized, easy-to-implement layout that distributes constant length segments of a clip among caches and provides modest storage efficiency. The *Silo* scheme improves upon *RCache*; it accounts for long term clip popularity and intra-clip segment popularity metrics and provides parameters to tune storage efficiency, server load, and playback switch-overs; (2) *Rainbow*, a local data replacement scheme based on the concept of *segment access potential* that accurately captures the popularity metrics. (3) A dynamic *global* data replacement or redistribution scheme called *Caching Token* that exploits existing data in distributed caches to reduce distribution overhead. Our schemes optimize storage space, start-up latency, server load, network bandwidth usage, and overhead from playback switch-overs. Our analytical and simulation results show that the *Silo* scheme provides 3 - 8 times as high cache hit ratio as that of traditional web caching system while utilizing the same amount of storage space.

B. Outline

The rest of this paper is organized as follows: Section II compares and contrasts relevant work in the areas of cooperative distributed caching and distributed data layouts in video-on-demand (VoD) servers. Section III describes our streaming cache architecture and introduces popularity measures we use in our design. In Section IV, we present design and analysis of our *RCache* and *Silo* distributed data layouts and provide guidelines for choosing various associated parameters. In Section V, we describe cache replacement policies that exploit the special structure inherent

in data layouts to optimize storage performance. In Section VI, we discuss mechanisms for redistributing a clip in the caching system in the event clip popularity changes. We propose a new *Caching Token* mechanism that reduces network traffic and server load. We describe the performance results from analytical and simulation models in Section VII. Finally, Section VIII presents the conclusions.

II. RELATED WORK

In general, two bodies of work are related to streaming cache design, namely web caching for regular objects and distributed video server design.

In cooperative caching systems, instead of operating independently, multiple caches cooperate to exploit cached content at neighboring caches. Such systems require protocols to organize caches into a topology and also to discover nearest caches from whom to retrieve content. Harvest [10], [14] represents one of the first *hierarchical caching* systems wherein the caches are *statically* configured in a hierarchy such that on a cache miss, the request is forwarded to the parent cache. The obvious drawback of such a system is the additional latency incurred in traversing the hierarchy when the requested object could be present in a neighboring (sibling) cache. This problem is fixed in Squid [1] using *static cooperative caching* where a cache always checks a fixed set of neighboring caches for any missing object before requesting it from its parent. A more sophisticated approach called *dynamic cooperative caching* used by Akamai [1] and Lucent *imminet* [1] extends the Squid ideas. It uses a centralized server that constantly computes a global map of the caching system. Using a cost function and DNS server based redirection, the client's request is then sent to a cache with the least cost. Our data layout schemes are designed for such cooperative caching systems and can be easily implemented in their architectural frameworks.

The resource-based caching algorithm in [25] is designed for web servers and proxies, and deals with storage of static and dynamic content on a disk system based on their bandwidth, storage requirements, and potential caching gain in an environment that has a single writer and multiple readers. The algorithm divides dynamic content into time intervals and stores the content by the unit of these intervals. We look at caching for streaming media only, therefore our algorithms are quite different than these.

The Middleman proxy architecture described in [4] has some similarities to our work but also, significant drawbacks as follows: (1) Like our architecture, it operates a collection of proxies as a scalable cluster cache. It consists of two kinds of proxies; storage proxies that store segments of clips and local proxies that serve the clients. The collection of these proxies is coordinated by a central-

ized coordinator which requires backup cohorts to ensure fault tolerance. Our architecture does not use a centralized coordinator. (2) Similar to our layouts, Middleman splits clips into smaller-sized segments and caches only one copy of each segment in the proxy cluster. In our layouts, we account for global and local popularity of a clip and its segments and store variable number of copies of each segment. Compared to Middleman, our approach guarantees better load balancing, fault tolerance and parallelism. (3) Middleman uses demand-driven data fetch and a modified version of LRU-K local cache replacement policy called *HistLRUpick*. We use a novel Rainbow data replacement policy based on the concept of access potential that accurately captures popularity of clip segments.

Several research efforts, such as Project MARS [12], Microsoft’s TIGER file system [9], Server Array [8], and research reported in the literature [13], [16], [24], [26] have addressed the problem of distributed or striped layouts. The primary focus of these schemes however, has been on video servers constructed using disks or tightly-coupled clusters. Also, they focus on optimizing data layouts for high concurrency and balanced operation of clusters under normal playout and interactive playback operations such as fast-forward and rewind and reduce origin server and network load. Our layouts are meant for loosely-coupled caching systems and are quite different.

Several researchers [6], [18], [19], [20], [22], [28] have proposed broadcast or multicast schemes that transmit variable or constant length segments of media clips on multiple broadcast or multicast channels at fixed or variable rate and offer near-video-on-demand (nVoD) services with limited interactivity. In contrast to a service model that uses a dedicated session per user, such schemes aim to use a fixed amount of network bandwidth and large client buffer to offer limited interactivity. Our layouts use variable length segments common with these schemes.

III. DISTRIBUTED CACHING ARCHITECTURE

Figure 1 illustrates our proposed distributed cache architecture which consists of a collection of streaming caches (SCs) interconnected together, a origin server where the content is created, and a set of multimedia clients (MCs). The origin server in our architecture is a streaming server in the content provider’s infrastructure such as *cnn.com* or a VoD server in a regional data center of a Metro Service Provider (MSP). The requests issued by a multimedia client (MC) for multimedia clips on an origin server are always sent to a designated Streaming Cache (SC) or proxy¹ associated with it. This association may be established

¹Here after, we use the terms SC, cache, and proxy interchangeably.

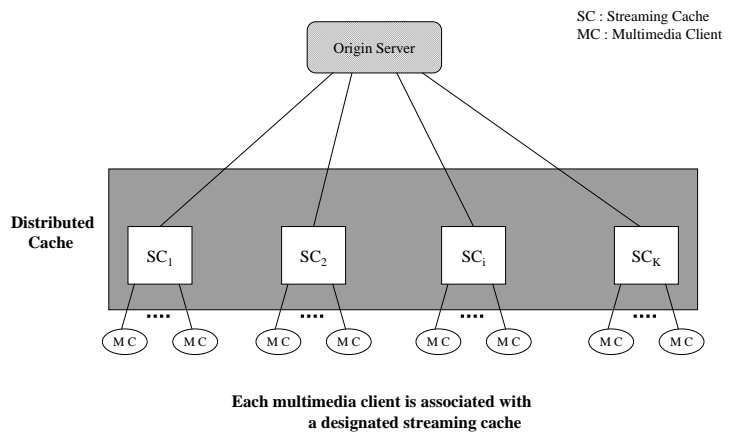


Fig. 1. Distributed streaming cache architecture.

manually by entering the address of SC in a “preferences” block in the web-browser or transparently using a layer-4 switch or DNS server based redirection as in [1].

Depending on the type of network that connects the proxies, different scenarios are possible. For example, (1) **Cluster Caching**, where each proxy can be a storage node in a large cluster cache interconnected by a Storage Area Network (SAN). (2) **Metro Caching**, where each proxy can be an independent cache deployed in buildings interconnected by a Metro Area Network (MAN). (3) **Enterprise Caching**, where each proxy can be an independent cache deployed in different departments in an enterprise building and interconnected by a fast enterprise network. (4) **Wide Area Network Caching**, where each proxy can be an independent cache in the Internet. Recent results in the literature show that cooperative caching is useful in small-to-medium scale networks, principally in LAN and MAN environments [30]. Therefore, we are primarily interested in scenarios 1, 2, and 3 above. In all these cases, the proxies together form a *proxy or cache cluster* with K members and each proxy is aware of the other members or siblings. Also, the proxies exchange state among themselves using an inter-cache protocol, for example a variant of ICP [29].

In our architecture, each media clip is divided into relatively small segments that are distributed among different proxies in the network. The design of our layouts combined with demand-driven caching, guarantees that each segment is stored in *at least* one proxy. The state exchange among proxies guarantees that each proxy knows the clips present in the cluster and has a *global segment map* that describes which proxies store segments of every clip. In the event that a proxy fails, other proxies can detect the failure using one or more mechanisms such as ICP, remote session management, global cache replacement mechanism, etc. and selectively migrate the active sessions to other proxies that

have copies of the required segments. Clearly, our replicated segmentation approach enables high availability and increased level of fault tolerance similar to the RAID file system [16]. The mechanics of creation of data layout and cache behavior are described next.

Cold start for clip Q (Data layout creation): Suppose one of the proxies SC_i receives a request for a clip Q which is not currently present in the cache cluster. The proxy undertakes three steps: (1) It queries the origin server to find out popularity information for Q and uses the design of our distributed *RCache* (Section IV-A) or *Silo* layout (Section IV-B) to compute a *local segment map* that describes which segments in the clip it should cache. (2) It fetches the first segment of the clip from the origin server, stores it locally, and concurrently streams it to the client. (3) It informs the other proxies of the appearance of this new clip. They each compute their own *local segment map* based on the data layout design. The proxies exchange this map immediately to construct a consistent *global segment map* that describes the *supposed presence* of segments in the cluster.

For later segments, if the local segment map at proxy SC_i says it should not store the particular segment, it consults the global segment map and redirects the client to another proxy SC_j that is supposed to have the segment. If SC_j does not have the segment in its local storage, it first checks to see if any of its siblings have it. If they do, it fetches the segment, stores it locally, and also streams it to the client. If no sibling proxies have the segment, SC_j fetches the segment from the origin server, stores a local copy, and also streams it to the client. Whenever the segment can be requested from multiple proxies, a proxy server selection algorithm can be used to balance load in the proxy cluster. Our layout design virtually guarantees that each segment is fetched from the origin server only once. Also, the demand-driven mechanism of layout creation guarantees that storage is utilized just-in-time when there is a request for the clip.

Cache hit for clip Q : When the proxy SC_i receives a request for a clip Q that is already present in the cluster, one of the following scenarios happens for each segment j of the clip Q : (1) Requested segment j is present in the local storage: In this case, SC_i streams that segment to the client. (2) Segment j is supposed to be present in the local storage ($LocalSegmentMap[Q, j] = 1$), but it has not been cached yet or has been flushed out due to local replacement: In this case, SC_i consults our *Rainbow* local replacement policy to allocate space for this segment and fetches the segment from a sibling proxy or the origin server and streams it to the client. (3) Segment j is not supposed

to be cached locally ($LocalSegmentMap[Q, j] = 0$), but the global segment map says it is present in one or more sibling proxies. In this case, SC_i redirects the request to the next proxy S_{next} that is selected by consulting a proxy server selection algorithm that attempts to balance the load in the proxy cluster.

Cache miss for clip Q : A clip cache miss happens either when the clip is requested for the first time, or when all of its segments are flushed out of the proxy cluster. The first case is identical to cold start, and the second case can be handled segment by segment as in case (2) or (3) of cache hit. We later use analysis and simulations to show that our architecture provides excellent cache miss performance.

Note that when the popularity of a clip changes, data layout and therefore the total amount of storage used for the clip in the caching system should reflect this change. Our *Caching Token* scheme guarantees that this is accomplished just-in-time with minimal redistribution overhead and origin server load. Periodically global cache map and local cache map are updated deleting information about clips that are no longer cached anywhere in the proxy cluster.

When responding to client requests, the cluster of proxies together act as a single, loosely-coupled streaming cache. One undesirable artifact of data segment distribution is that two consecutive segments may have to be retrieved from different proxies, requiring signaling to switch between two streams coming on two different connections. Such switching is a potential source for disruption of client playback and deterioration of quality. However, using techniques such as double buffering or fetching data at a higher rate, playback from different proxies can be pipelined and disruption can be minimized. Also, our *Silo* data layout reduces the number of switch-overs by appropriate design of segment lengths.

In this architecture, we use the following performance metrics to assess our schemes: (1) total storage requirement, (2) playback start-up latency, (3) number of switch-overs, (4) number of active streams at each proxy, (5) origin server hit ratio and (6) cluster cache hit ratio, also called system-wide cache hit ratio. The goal for designing a data layout algorithm is to balance the tradeoffs between these metrics.

A. Clip and Segment Popularity

The popularity of clips and segments plays a critical role in our algorithms. We use two popularity metrics: (1) *Global Clip Hotness Rating (GCHR)*, a macro-metric that captures the global popularity of each clip, and (2) *Local Segment Hotness Rating (LSHR)*, a micro-metric that captures the popularity of segments of the clips.

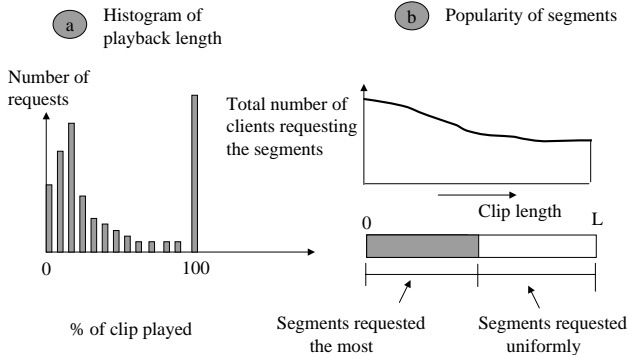


Fig. 2. Bimodal distribution of clip segment requests.

Global Clip Hotness Rating (GCHR) of a clip is defined as the total number of requests for the clip per unit time in the proxy cluster. At each proxy, the total number of requests to the clip per unit time is called *Local Clip Hotness Rating (LCHR)*. When a new media clip is stored on the origin server, the proxy cluster may set its GCHR value based on a priori knowledge of the expected popularity of the clip. Each proxy also maintains a LCHR for each clip, and the proxy cluster can compute the GCHR for a clip by adding its LCHR value at different proxies.

The GCHR is a measure of long term popularity of multimedia clips. Current understanding of long term access patterns of streaming media on the web is rather meager and often the Zipf distribution is used to model such long term access patterns [5], [11]. Zipf’s law relates clip access probability p and the popularity of the clip described by a rank metric i as $p = k(\frac{1}{i^\gamma})$ where $k = \sum_{j=1}^M (1/j^\gamma)$, M is the total number of clips in the proxy cluster and γ is in the range of $0.8 \leq \gamma \leq 1.2$. In our setup, we can determine exact value of γ by solving a curve fitting problem of mapping the series of GCHR for all the clips in the system to the series of their ranks. In other words, we calculate γ from GCHR of all the clips.

The popularity of segments within a clip strongly depends on the user access behavior. A recent study shows that playback length of media clips on the web follows a bimodal distribution [5]. Specifically, as shown in Figure 2(a), a fraction of the requests for a clip complete the playback, whereas the rest terminate after partial playback. Among the partial playback requests, the probability of playback exceeding a threshold drops rapidly for longer duration playbacks. This suggests that different parts of the clip are requested with different probability; segments in the beginning of clips are requested more often than later segments as in Figure 2(b). The popularity of a segment includes all playback requests that last long enough to access the segment. Therefore, the segment popularity shown in Figure 2(b) can be derived as a complementary cumula-

tive distribution function of the playback time distribution in Figure 2(a). Two important implications of such access patterns are: (1) A clip does not need to be cached in its entirety at each proxy; even partial segment caching can provide caching gain [4]. (2) Segment layout must account for *segment popularity* in addition to clip popularity.

Each proxy in our architecture maintains a Local Segment Hotness Rating (LSHR) for each segment, which is the total number of requests for the segment during unit time. We use the GCHR and LSHR in the design of our data placement and replacement schemes discussed in the following.

IV. DATA LAYOUTS IN DISTRIBUTED CACHES

In this section we first describe a simple random data layout scheme called *RCache*, and motivate the need for better layouts. We then present the details of our *Silo* layout.

A. *RCache: Random Caching of Media File Segments*

The basic *RCache* data layout scheme randomly distributes segments of media clips among the cluster of K SCs. It divides a media clip of length L into N equal-length segments and lets each SC independently decide which segment to cache. The goal of this scheme is to store all the segments within the proxy cluster, so that requests for the clip can be served entirely out of the cluster, without involving the origin server. Specifically, for a constant a ($0 \leq a < K$), each SC independently decides with a fixed probability a/K whether to store a particular segment. Notice the segment caching probability is proportional to $1/K$ and is independent of both the segment number and the proxy location. The algorithm is as follows:

```

/* Given constant  $1 \leq a < K$ , compute  $p = a/K$  */
FOR (all segments in the clip )
  r = UniformRandomNumGen(IPAddr)
  IF ( $0 < r < p$ ) THEN
    Cache the segment
  ELSE
    Ignore the segment
  END
END

```

The estimated total number of segments of the clip stored at each SC is then $(a \times N)/K$. We need to determine the constant a such that the probability that the entire clip is stored in the proxy cluster is sufficiently large. In other words, we are interested in the probability of the event – “there exists no segment which is not cached in the cluster of K SCs” or alternatively stated event – “each

segment is cached at least once in the cluster of K SCs”.

$$\begin{aligned}
& Pr(\text{Media clip is cached in the cluster}) \\
&= Pr(\text{Each segment is cached at least once in } K \text{ SCs}) \\
&= 1 - Pr(\text{At least one segment is not cached in any SC}) \\
&= 1 - Pr\left(\bigcup_{1 \leq i \leq N} E_i\right) \tag{1}
\end{aligned}$$

where E_i is the event that segment i ($1 \leq i \leq N$) is not cached at any of K SCs. Using the Inclusion-Exclusion Principle of probability: $Pr(\bigcup_{1 \leq i \leq N} E_i) \leq NPr(E_i)$, we can rewrite Equation 1 as follows:

$$Pr(\text{Media clip is cached in the cluster}) \geq 1 - NPr(E_i) \tag{2}$$

The probability that segment i is not stored by SC_j is $P_{uncache}(i, j) = 1 - \frac{a}{K}$. The probability that segment i is not stored by any SC is then

$$Pr(E_i) = P_{uncache}(i) = \left(1 - \frac{a}{K}\right)^K \tag{3}$$

Therefore, Equation 2 can be rewritten as

$$\begin{aligned}
& Pr(\text{Media clip is cached in the cluster}) \\
&\geq 1 - N\left(1 - \frac{a}{K}\right)^K \tag{4}
\end{aligned}$$

Since $\left(1 - \frac{1}{K}\right)^K \rightarrow e^{-1}$, as $K \rightarrow \infty$, we have $N\left(1 - \frac{a}{K}\right)^K \leq Ne^{-a}$ or $1 - N\left(1 - \frac{a}{K}\right)^K \geq 1 - Ne^{-a}$. Thus, if we choose $a = c \ln N$ where c is a constant, we get

$$\begin{aligned}
& Pr(\text{Media clip is cached in the cluster}) \\
&\geq 1 - \frac{N}{e^{c \ln N}} = 1 - \frac{1}{N^{c-1}} \tag{5}
\end{aligned}$$

As per Equation 5, the probability that the entire clip is cached in the cluster is *at least* $\left(1 - \frac{1}{N^{c-1}}\right)$, which is a desirable *lower* bound. Clearly, we can control this lower bound by changing c and N , and *virtually* guarantee that the entire clip is cached in the cluster. Consider, $N = 50$, $K = 10$, and $c = 2$. Using Equation 4, and $a = c \ln N$ we can see that $Pr(\text{Media clip is cached in the cluster}) \geq 1 - 50 * \left(1 - \frac{2 * \ln 50}{10}\right)^{10} \approx 0.99998$. Using the conservative lower bound, same can be calculated as $1 - (1/50)$ or 0.98. Clearly, the conservative lower bound is attractive as it is independent of K and gives an idea how well RCache does for a given number of segments.

To evaluate RCache storage requirement, the naive replication approach where each clip is stored in every proxy in the system is used as a baseline which requires a LK storage. The storage efficiency factor S_{eff} is defined as the ratio of the total storage of any scheme over that of the baseline scheme. RCache scheme will be a desirable scheme if its storage requirement is smaller than LK . That

is, S_{eff} should be less than 1. From the discussion above we can see that the storage requirement of RCache scheme is $K(aN/K)(L/N) = aL$. Therefore, the $S_{eff} = aL/(LK)$ or $(c \ln N)/K$. Clearly, as long as $c \ln N < K$, $S_{eff} < 1$ and the RCache strategy is superior. Observe that for a given c , the number of segments in the stream N that makes RCache superior is computed as $N < e^{K/c}$. Consider a simple example: if we have $K = 10$ proxies, and $c = 2$, then the number of clip segments should be less than $e^{10/2} = e^5 = 149$, in order for RCache to be better than the naive strategy. Notice with a fixed c , a decreasing N decreases S_{eff} , and therefore increases storage efficiency. Also notice, decreasing N decreases lower bound of the probability that the entire clip is cached in the cluster, $\left(1 - \frac{1}{N^{c-1}}\right)$. In practice, the value of N should be chosen to balance the tradeoff between these two factors.

In the basic RCache scheme above, segments of constant length are stored at proxies with constant probability. Each segment is virtually guaranteed to be in one of the proxies while keeping the storage requirement less than LK – the simple case where each segment is stored in every proxy. However, under RCache, the user perceived start-up latency can be large if the client’s designated SC does not store the first segment or even worse, the first segment is not stored in any of the SCs in the cluster. In the second scenario, the request has to be redirected to the origin server. We can eliminate these scenarios if we can guarantee that the first segment of the clip is always available on the designated proxy. This requires that the first segment be stored in every SC in the cluster.

Given that the segments in the early part of a clip are requested more frequently than those in the later part, it makes sense to distribute later segments to fewer SCs. Following this hypothesis, a layout that employs segments of increasing duration will decrease the average number of switch-overs for each user. In other words, instead of making each segment the same length and the probability to cache each segment the same, we can have a variety of RCache schemes, each with different segment lengths and different probability to cache each segment. We have two degrees of freedom in the design of RCache algorithm: one is how to divide a clip into segments, another is how to store the segments. In the following section, we discuss the *Silo* scheme for data layout, named after the shapes of the stacked-bar graphs of segment size versus segment number.

B. Silo Data Layouts

We first describe the *Silo* data layout scheme for a single clip and then describe the most general layout for multiple clips. We also provide guidelines for how relevant parameters for these data layouts should be selected.

Segment size variation over rounds

Probability variation over rounds

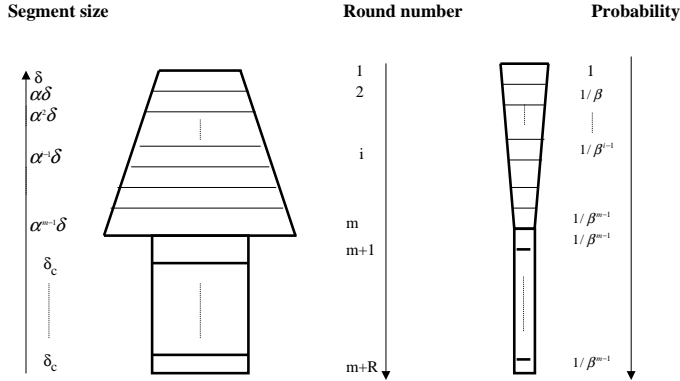


Fig. 3. Generalized Silo layout: segment size and probability variation.

B.1 Basic Silo Data Layout for a Single Clip

Consider a clip with length L , divided into N segments. The length of each segment and the probability with which an SC decides to cache a segment are controlled by two parameters: an increasing ratio of segment size (SSR) α and a decreasing ratio of intra-clip segment caching probability (SCPR) β . For the Silo layout, the size of the first segment is set to a small fraction of the entire clip, δ , and each proxy stores it with probability 1. The size of the second segment is $\alpha\delta$, and the probability to store the segment drops to $1/\beta$. In general, up to the m^{th} segment, the probability with which segment i of length $\alpha^{i-1}\delta$ is stored at each proxy is $1/\beta^{i-1}$.

At the m^{th} segment, the segment size reaches its maximum value of $\alpha^{m-1}\delta$. From that point on, the segment size is held constant at δ_c and the caching probability is fixed at $1/\beta^{m-1}$. We can either set (1) the constant segment size δ_c of the $(m+1)$ -th and all subsequent segments to the size of the m^{th} segment or (2) reduce it to smaller values as in Figure 3. The choice depends on the tradeoff between switch-overs and flexibility of storage management in cache replacement mechanism. As the segment size becomes large, the number of switch-overs is reduced at the expense of flexibility of cache replacement.

Table I illustrates the segment size and caching probability variation. The main principles used in this design are: (1) *Minimize playback start-up latency*: To this end, the layout always caches the first segment of a clip at each proxy. (2) *Minimize playback switch-over for most requests*: To this end, the layout uses increasingly longer segments for the portions of the clip that are most requested. The Silo shape matches the bimodal request distribution for segments. (3) *Achieve better load balance and fault tolerance*: Unlike Middleman [4] proxy architecture that stores only

one copy of constant length segments, Silo layout stores increasingly larger number of copies of popular segments and attempts to guarantee at least one copy of each segment. This requires modestly more storage but statistically guarantees better load distribution in the proxy cluster, better fault tolerance, and higher parallelism.

TABLE I
DETAILS OF BASIC SILO DATA LAYOUT

Segment number	Segment size	Pr(Seg Caching)	Expected number of proxies caching the segment	Expected total storage for the segment
1	δ	1	K	δK
2	$\alpha\delta$	$1/\beta$	K/β	$\frac{\alpha}{\beta}\delta K$
3	$\alpha^2\delta$	$1/\beta^2$	K/β^2	$(\frac{\alpha}{\beta})^2\delta K$
...				
m	$\alpha^{m-1}\delta$	$1/\beta^{m-1}$	K/β^{m-1}	$(\frac{\alpha}{\beta})^{m-1}\delta K$
$m+1$	δ_c	$1/\beta^{m-1}$	K/β^{m-1}	$\delta_c \frac{K}{\beta^{m-1}}$
$m+2$	δ_c	$1/\beta^{m-1}$	K/β^{m-1}	$\delta_c \frac{K}{\beta^{m-1}}$
...				
...				
$m+R$	δ_c	$1/\beta^{m-1}$	K/β^{m-1}	$\delta_c \frac{K}{\beta^{m-1}}$

Note that the sum of the size of first m segments is

$$\delta + \alpha\delta + \alpha^2\delta + \dots + \alpha^{m-1}\delta = \frac{\delta(\alpha^m - 1)}{\alpha - 1} \quad (6)$$

After round m , all the segments are of the same size δ_c . Therefore, for a clip of length L , the number of segments needed for the bottom half of the generalized Silo shape is

$$R = \frac{(L - \frac{\delta(\alpha^m - 1)}{\alpha - 1})}{\delta_c} \quad (7)$$

From the last column in the Table I, we can see that the total storage requirement for the basic Silo scheme is

$$T = K \left(R \frac{\delta_c}{\beta^{m-1}} + \delta \sum_{i=0}^{m-1} \left(\frac{\alpha}{\beta}\right)^i \right) \quad (8)$$

The storage efficiency factor is

$$\begin{aligned} S_{eff} &= \frac{T}{LK} \\ &= \frac{\left(R \frac{\delta_c}{\beta^{m-1}} + \delta \sum_{i=0}^{m-1} \left(\frac{\alpha}{\beta}\right)^i \right)}{L} \end{aligned} \quad (9)$$

Note that as long as the segment caching probability after the first segment is less than 1, S_{eff} is less than 1.

B.2 Silo Data Layout for Multiple Clips

If there are multiple identical-length media clips and all the clips are equally popular, each clip can be segmented in the same way using identical values of the parameter set $(\delta, \alpha, \beta, m)$. However, in reality, media clips differ in popularity and therefore, the caching system must allocate storage for clips in accordance to their popularity; more popular clips should be allocated more storage than less popular ones.

The Silo layout for a single clip accounts for segment popularity within a clip using a probability sequence

$$\begin{aligned} U &= (u_1, u_2, \dots, u_j, \dots) \\ &= (1, 1/\beta, 1/\beta^2, \dots, 1/\beta^{m-1}, \dots, 1/\beta^{m-1}) \end{aligned}$$

The layout caches more popular segments with higher probability. A natural way to extend this for multiple clips ordered by global popularity rating (GCHR) is as follows: given two clips with popularity ratings $GCHR_1$, $GCHR_2$ and associated segment caching probability sequences U_1 , U_2 , if $GCHR_1 > GCHR_2$, then set $U_1 > U_2$. Specifically, the probability sequence for a clip with GCHR i is

$$e(i) \otimes U = (e(i, 1)u_1, e(i, 2)u_2, \dots, e(i, j)u_j, \dots)$$

where $e(i)$ is called the *caching ratio* for a clip with rank i . It is defined as $e(i) = (e(i, 1), e(i, 2), \dots, e(i, j), \dots)$, and

$$e(i, j) = \begin{cases} 1 & \text{if } j = 1 \\ x, & 0 \leq x \leq 1 \text{ otherwise} \end{cases}$$

Using the Zipf's distribution, for a clip of rank i , we set $e(i, j) = 1/i^\gamma$, for $j \neq 1$, where γ ($0.8 \leq \gamma \leq 1.2$) controls decreasing ratio of inter-clip segment caching probability (ISCP) and the relative caching preference among clips with different ranks. Higher γ gives more preference to popular clips.

Table II shows the details of the configuration parameters and the cache layout for a clip of rank i . In this general layout, during the exponential increase phase, the segment size is increased by a factor α for successive segments until the maximum $\alpha^{m-1}\delta$ is reached. After this, the remaining segments have the same constant size of δ_c .

The average per-proxy cache size of the clip of rank i is

$$\begin{aligned} S_i &= \delta + \left(\delta \frac{\alpha}{\beta i^\gamma} + \delta \frac{\alpha^2}{\beta^2 i^\gamma} + \dots + \delta \frac{\alpha^{m-1}}{\beta^{m-1} i^\gamma} \right) + R \frac{\delta_c}{\beta^{m-1} i^\gamma} \\ &= \delta + \frac{(L - \frac{\delta(\alpha^m - 1)}{\alpha - 1})}{\beta^{m-1} i^\gamma} + \frac{\delta}{i^\gamma} \sum_{k=1}^{m-1} \left(\frac{\alpha}{\beta} \right)^k \end{aligned} \quad (10)$$

Therefore, the total size of the Silo cache with K proxies and M clips is

$$T = \sum_{i=1}^M K S_i \quad (11)$$

TABLE II

GENERAL SILO CACHE LAYOUT FOR A CLIP WITH RANK i

Parameters: (1) δ : Size of the first segment. (2) δ_c : Size of the constant segments. (3) α : Increasing Segment Size Ratio (SSR). (4) β : Decreasing Intra-clip Segment Caching Probability Ratio (SCPR). (5) γ : Exponent of the decreasing Inter-clip Segment Caching Probability Ratio (ISCP).				
Segment number	Segment size	Prob(Seg Caching)	Expected number of proxies caching the segment	Expected total storage for the segment
1	δ	1	K	δK
2	$\alpha\delta$	$\frac{1}{\beta i^\gamma}$	$\frac{K}{\beta i^\gamma}$	$\alpha\delta \frac{K}{\beta i^\gamma}$
3	$\alpha^2\delta$	$\frac{1}{\beta^2 i^\gamma}$	$\frac{K}{\beta^2 i^\gamma}$	$\alpha^2\delta \frac{K}{\beta^2 i^\gamma}$
\vdots				
m	$\alpha^{m-1}\delta$	$\frac{1}{\beta^{m-1} i^\gamma}$	$\frac{K}{\beta^{m-1} i^\gamma}$	$\alpha^{m-1}\delta \frac{K}{\beta^{m-1} i^\gamma}$
$m+1$	δ_c	$\frac{1}{\beta^{m-1} i^\gamma}$	$\frac{K}{\beta^{m-1} i^\gamma}$	$\delta_c \frac{K}{\beta^{m-1} i^\gamma}$
\vdots			\vdots	
$m+R$	δ_c	$\frac{1}{\beta^{m-1} i^\gamma}$	$\frac{K}{\beta^{m-1} i^\gamma}$	$\delta_c \frac{K}{\beta^{m-1} i^\gamma}$

The storage efficiency for a clip with rank i can then be computed as $S_{eff}^i = S_i/L$, and the total storage efficiency for M clips is computed as $S_{eff} = \sum_{i=1}^M \frac{S_{eff}^i}{M}$.

Note that the Middleman [4] proxy architecture stores only one copy of each segment and therefore always requires LM amount of storage for M clips. Clearly, Silo layout requires $\sum_{i=1}^M K S_i / LM = K S_{eff}$ times larger storage than Middleman.

C. Analytical Model

In this section, we introduce a simple analytical model of the static Silo configuration, providing a valuable means to understand the system dynamics and also providing the basis for our *local Rainbow* replacement scheme and proxy load balancing schemes.

The analytical model computes three important performance metrics: (1) *Local cache byte hit ratio (LHR)*, Φ_l , representing the percentage of stream data that is retrieved from local Silo caches. (2) *Remote cache byte hit ratio (RHR)*, Φ_r , representing the percentage of stream data that is retrieved from remote Silo caches. (3) *Server byte hit ratio (SHR)*, Φ_s , representing the percentage of stream

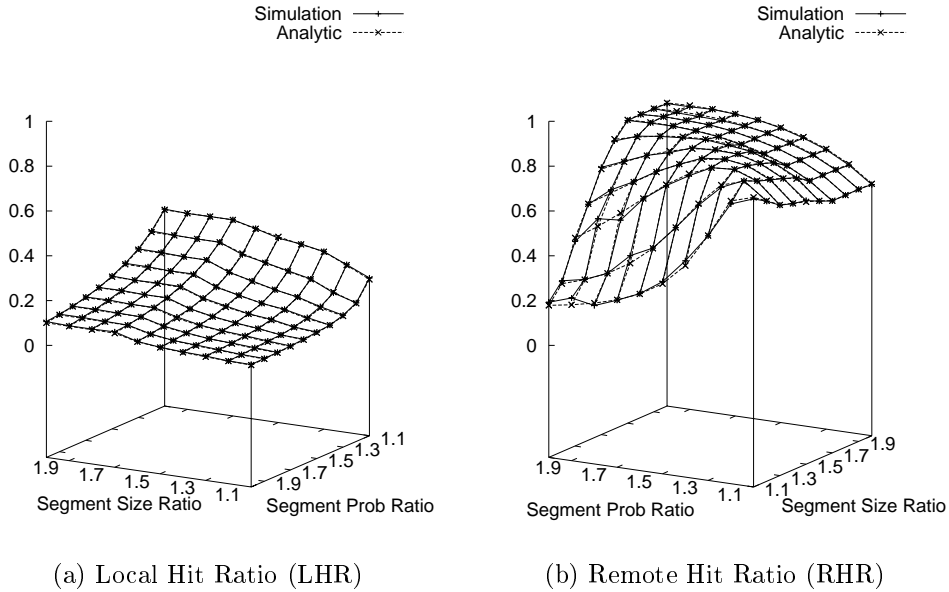


Fig. 4. Local and remote hit ratios with 100 Silo proxies ($\delta = 50$ MB).

data that is retrieved from the origin server. For the sake of brevity, we will refer to byte hit ratio as hit ratio from here on.

Assume that the number of Silo proxies is K and the number of clips is M . For the Silo cache layout shown in Table II, the segment caching probability of segment j of a clip with rank i , $P_{i,j}$, is

$$P_{i,j} = \begin{cases} 1 & \text{if } j = 1 \\ \frac{1}{\beta^{j-1}i^\gamma} & 2 \leq j \leq m \\ \frac{1}{\beta^{m-1}i^\gamma} & m < j \leq m + R \end{cases}$$

Given this, the probabilities that a segment can be accessed from local cache, origin server and remote caches, $P_{i,j}^l$, $P_{i,j}^s$, and $P_{i,j}^r$, respectively, are

$$P_{i,j}^l = P_{i,j}, \quad P_{i,j}^s = (1 - P_{i,j})^K, \\ P_{i,j}^r = (1 - P_{i,j})(1 - (1 - P_{i,j})^{K-1})$$

Also, let ρ_i be the probability that a clip with rank i is selected, then ρ_i is

$$\rho_i = \frac{\frac{1}{i^\gamma}}{\sum_{k=1}^M \frac{1}{k^\gamma}}$$

Given that a clip is selected, the probability that a segment j is selected, ψ_j , is

$$\begin{aligned} \psi_1 &= \text{Prob}\{\text{Playback time} \geq \text{StartTime}(1)\} \\ &= 1 \\ \psi_j &= \text{Prob}\{\text{Playback time} \geq \text{StartTime}(j)\} \end{aligned}$$

$$= \text{Prob}\{\text{Playback time} \geq \sum_{k=1}^{j-1} \tau_k\}, \quad 2 \leq j \leq m + R$$

where τ_k is the length of segment k and $\text{StartTime}(j)$ is the start time for segment j . The probability ψ_j can be calculated from either LSHR or the bimodal distribution [5].

The average number of bytes retrieved from local cache, remote caches, and origin server, θ_l , θ_r , and θ_s , respectively, are

$$\begin{aligned} \theta_l &= \sum_{k=1}^M \rho_k \left(\sum_{j=1}^{m+R} \eta_j \psi_j P_{k,j}^l \right), \quad \theta_r = \sum_{k=1}^M \rho_k \left(\sum_{j=1}^{m+R} \eta_j \psi_j P_{k,j}^r \right) \\ \theta_s &= \sum_{k=1}^M \rho_k \left(\sum_{j=1}^{m+R} \eta_j \psi_j P_{k,j}^s \right) \end{aligned}$$

where η_j is the size of the segment j . Finally, hit ratios are

$$\Phi_l = \frac{\theta_l}{\theta_l + \theta_r + \theta_s}, \quad \Phi_r = \frac{\theta_r}{\theta_l + \theta_r + \theta_s}, \quad \Phi_s = \frac{\theta_s}{\theta_l + \theta_r + \theta_s}$$

Figure 4 shows the results from (1) ns-2 [27] simulation and (2) our analytical model for a Silo system with 100 clips of 3 GB each stored on a cluster of 100 SCs interconnected by a mesh topology. Figure 4(a) shows the local hit ratio, fraction of data served from the designated local SCs, with varying segment size ratio (α) and segment caching probability ratio (β). Figure 4(b) shows the remote hit ratio, fraction of data served from the sibling remote SCs, with varying α and β parameters. It is clear that the results from our analytical model match very well with simulation, validating the accuracy of the model.

D. Guidelines for Choosing Parameters

Given the parameters L, K, S_{eff} , the parameters that need to be decided before the Silo layout for the clip can be realized are δ, α, β, m , and γ . The smallest segment size δ must be long enough to allow the client to start the pre-fetch pipeline for subsequent segments. Typically, an initial segment size of 1 to 2 minutes would be sufficient. For MPEG-2 streams with 5 Mbps average bandwidth, this results in the first segment size of roughly $\delta = 35 - 75$ MB. For α and m , the size of the Silo head is supposed to be large enough to cover the short beginning period that gets the maximum number of requests. Given the first ρ fraction of the clip is accessed most frequently, the size of Silo head structure from Equation (6) should satisfy the following condition:

$$\frac{\delta(\alpha^m - 1)}{\alpha - 1} \geq \rho L \quad (12)$$

where ρ can be obtained from either LSHR or the observed bimodal distribution of user playback time [5]. Equation (12) provides a set of candidate (α, m) tuples. Note also that a set of candidate (β, m) tuples can be calculated from either the bimodal distribution of playback time [5] or the LSHR. The final (α, β, m) can be determined from both candidate tuples of (α, m) and (β, m) . Figure 4 shows that there could be multiple (α, β, m) triples that provide similar cache hit ratio performance.

Finally, as γ becomes larger, S_{eff} improves while decreasing cache hit ratio for less popular clips. γ can be obtained from GCHR and adjusted in conjunction with S_{eff} with a tradeoff between S_{eff} and cache hit ratio for less popular clips.

V. LOCAL CACHE REPLACEMENT: THE RAINBOW ALGORITHM

Each component streaming cache in our distributed caching system has a fixed amount of storage space. In the event its storage space is exhausted, it requires a local segment replacement policy that removes data segments to create space for new segments. The purpose of a cache replacement algorithm is to retain segments that maximize the hit ratio. Clearly, such a replacement algorithm must account for both Global Clip Hotness Rating (GCHR) and Local Segment Hotness Rating (LSHR) metrics defined in Section III.

The Silo data layout has an inherent two-dimensional structure of relative ordering of importance among segments based on inter-clip and intra-clip popularity. In other words, clips themselves can be ordered in decreasing order of popularity (GCHR) and the segments within a clip can again be ordered based on the segment popularity (LSHR). Efficient selection of segments for replace-

ment requires that the two-dimensional ordering structure be transformed into a one-dimensional ordering based on a single quantity that captures both popularity metrics. To this end, we propose a novel cache replacement algorithm, *Rainbow*, that uses a single metric called *caching potential* to order segments.

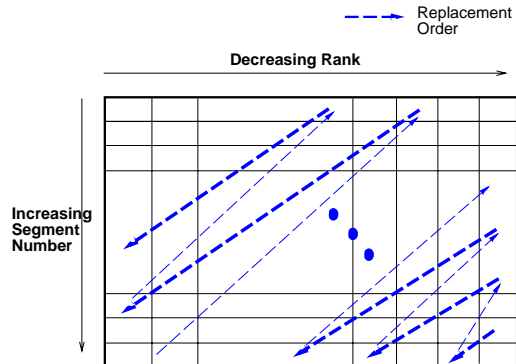


Fig. 5. Cache replacement order of Rainbow algorithm.

The two main ideas in Rainbow are: (1) Assignment of caching potential to each segment. (2) Quantization of caching potential into “bands” of colors (potentials) to minimize the overhead of sorting the segments.

The *caching potential function* $F(i, j)$ is proportional to the probability that a segment j of a clip with rank i is accessed when a user request for clip playback arrives. The higher the probability is, the higher the potential that the segment is requested for playback. Clearly, $F(i, j)$ is based both on clip popularity distribution and playback time distribution and successfully captures both inter-clip and intra-clip orderings of importance. The potential function $F(i, j)$ can be formally written as

$$\begin{aligned} F(i, j) &= Prob\{\text{Clip rank} = i\} \times Prob\{\text{Segment played} = j\} \\ &= \frac{1}{i^\gamma} \times Prob\{\text{Playback time} \geq \text{StartTime}(j)\} \\ &= \rho_i \times \psi_j \end{aligned}$$

Using the analytical model in Section IV-C, we can compute ρ_i, ψ_j , and therefore, $F(i, j)$. Using this potential function, the cache replacement algorithm can sort all the segments in its store in decreasing order of the potential and replace segments from the end of the sorted list whenever it needs more space. Figure 5 shows the new diagonal ordering of the Rainbow cache replacement algorithm. Segments in the lower right corner have the lowest potential and segments in the upper left corner have the highest potential.

Clearly, with a large number of clips, the number of segments in the proxy can be very large and the overhead in-

curred to sort caching potentials every time ranks change or new clips are introduced can be prohibitively large. However, note that the exact ordering of the segments is not strictly required. Therefore, we adopt a quantization mechanism that is based on coarse grain sorting. Instead of using the exact values of the potential function and the exact ordering, the entire range of the values of the potential function is divided into subranges, each of which is assigned a color code, hence the name ‘‘Rainbow’’. The quantization can be either linear or non-linear. Non-linear quantization can give more flexibility by dividing higher potential regions with coarser granularity and lower potential regions with finer granularity. Because the most active region for cache replacement is the region with lower access potential, it is desirable to have more control over this region. On the other hand, a linear quantization gives simplicity. The Rainbow scheme stores segments into bins with matching color and does not maintain segment ordering within a bin. When adequate space cannot be found to cache a new segment with an assigned potential and color band, the Rainbow algorithm attempts to replace segments from the lowest color bin until it finds enough space for the new segment.

The Rainbow algorithm creates a novel one-dimensional replacement ordering of segments by combining both dimensions of the Silo layout. It works seamlessly under various configurations of Silo cache layout and is therefore quite flexible.

VI. GLOBAL REPLACEMENT: DYNAMIC CLIP REDISTRIBUTION

In this section, we present a new scheme, called *caching token* to efficiently implement global data redistribution required to reflect changes in clip popularity.

Recall that in our caching architecture, the rank of a clip plays an important role since it affects segment caching probabilities and decides the number of cached copies of segments. As the rank of a clip changes over time with changing client preference, data layout must be changed accordingly to account for the new rank information. If the rank of a clip goes up (down), the number of cached segments of the clip should be increased (decreased).

A naive way to deal with rank change is to redistribute the entire clip and rebuild a new cache layout. However, this approach has obvious drawbacks. Complete redistribution generates considerable load on the network and the origin server. During such a data redistribution phase, each proxy builds a new cache layout from scratch. Until this process completes, the Silo cache maps are inconsistent and a lot of client requests for the clip will be served from the origin server. Also, this process can trigger increased syn-

chronized ICP message exchanges among proxies. In the following we present an efficient solution to these problems.

A. Lazy Caching and Caching Token

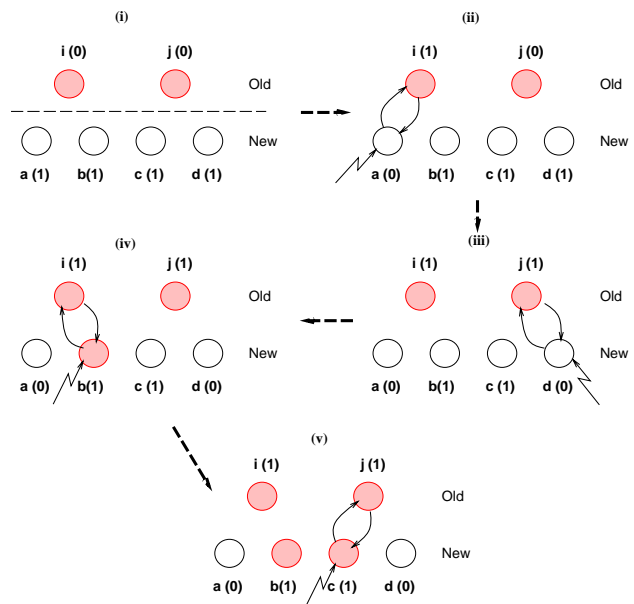


Fig. 6. Silo cache layout change: rank up.

We introduce two mechanisms – *lazy caching* and *caching token* – to deal with the aforementioned problems associated with the simple redistribution method. The caching token mechanism is based on the concept of moving the ‘‘pointer to data location’’ rather than the data itself. To achieve this basic objective, each proxy maintains two bits of state information for each segment: (1) *Segment presence bit (C)*: indicates if the segment is currently present in the local storage. (2) *Token bit (T)*: indicates if the proxy is ‘‘legally’’ holding the segment. If this bit is set, the segment is supposed to be cached in the new cache layout. When a proxy has rights to the segment but does not have the segment, that is, when $T = 1$, and $C = 0$, it attempts to pass the token to a proxy that may have the segment but not the token, that is, $T = 0$, and $C = 1$. The goal of such token exchange is to maintain the correct number of copies of the segment in the Silo cache system, while moving data as little as possible. We will elaborate how this is achieved in the following.

When a Silo proxy receives new rank information from the server, it recomputes the segment caching probabilities of the clip and determines a new cache layout. Note that recomputing the data layout results in two sets of proxies: (1) an old set of proxies which currently cache the segment, and (2) the new set of proxies which are supposed to cache the segment. Any proxy belonging to both sets does not need to change its cache layout. Only the remaining

proxies are involved in token exchanges to build the new cache layout. However, the new layout takes effect lazily, on an as-needed basis when there is a client request for the segment. The changes are realized using a caching token exchange on a segment-by-segment basis instead of the entire clip. Thus, exchange of a few bits of information per segment, instead of segments themselves, reduces network traffic and load on Silo proxies.

When a client request is received for a segment for which rank information and therefore layout has changed recently, two cases can arise: (1) no proxy has cached the requested segment, and (2) at least one proxy has cached the segment. For the first case, the proxy reroutes the client request to the origin server and also obtains a copy of the segment. For the second case, the caching token exchange mechanism attempts to utilize currently cached copies of the segment.

We illustrate these concepts with an example when clip rank is altered. Figures 6(i)-(v) illustrate the token exchange dynamics for one segment of a clip, when the clip rank goes up. In these figures, each circle corresponds to a proxy and the number $x(y)$ represents value of the token bit y at node x for the segment under consideration. In the old layout (Figure 6(i)), before the rank change, the two proxies i and j shown by shaded circles store the segment. After the rank change, a new set of Silo proxies, a , b , c and d , determine they should have a copy of the segment in their local storage and therefore, set the token bit T for this segment to 1 indicating their right to hold the segment. The proxies i and j on the other hand conclude that they should not keep this segment and therefore set the token bit to zero indicating they are illegally holding the segment. The complete change of layout does not begin until there is a client request for the segment. As in Figure 6(ii), when proxy a receives a request for the segment, its local state says that it is supposed to have this segment, but it does not have it in the storage. Also, it knows from the previous layout information that proxy i has this segment. So instead of fetching the segment, proxy a decides to hand over the “legal rights” for this segment to proxy i . To achieve this, it sends proxy i a segment request with the token bit set, indicating a token exchange (See Figure 6(ii)).

Upon arrival of a segment request with the token bit set, proxy i performs two checks: First, if the segment presence bit C is not set, indicating the segment is not stored in the proxy and therefore an error is returned to the proxy requesting token exchange. Otherwise, the proxy checks the token bit T of the segment; if it is clear, the proxy grabs the token, else an error is returned. In our example, proxy i has the segment with a clear token bit, so it grabs the token and sets the token bit of the segment to $T = 1$.

By doing this, proxy i gains the right to retain the segment in its cache. It also sends a segment reply message with token bit cleared to proxy a , which then concludes that the token exchange is successful and clears its token bit for the segment. Proxy a redirects the client to proxy i which streams the requested segment to the client. Notice the segment reply message from a does not contain the actual data segment. This terminates a successful cycle of token exchange. In Figure 6(iii), proxies d and j have exchanged the token successfully. In Figure 6(iv), proxy b tries to exchange its token with proxy i . However, the token exchange fails since proxy i already has the token for the segment. Thus, proxy i sends the segment with token bit set to proxy b . Proxy b learns the token exchange has failed and caches the segment with token bit set to 1. The same thing happens to proxy c in Figure 6(v). In the end, proxies i , j , b and c cache the segment instead of the new set of proxies, a , b , c , and d . By exchanging tokens, proxies i and j can keep the segment in their caches, eliminating the movement of the segment to proxies a and d .

The caching token mechanism works similarly for the rank-down case, removing unnecessary segments from local storage. The local replacement scheme lazily removes unnecessary segments when additional space is needed.

Clearly, compared to the simple redistribution mechanism, our new mechanism has several advantages: (1) The origin server is relieved from the redistribution of an entire clip and participates in reconstruction only if there is no proxy that caches a segment. (2) The token exchange method minimizes the required segment copying to build up a new cache layout. (3) The on-demand nature of token exchange eliminates synchronized ICP message exchanges and reduces instances of inconsistent cache map.

Note if there are N segments in a clip, in the worst case N bits of information need to be exchanged for token exchanges per clip. Clearly, a single maximum size IP packet can exchange token bits for a large number segments of multiple clips and thus, the caching token mechanism requires negligible overhead.

VII. PERFORMANCE ANALYSIS

In this section, we describe the simulation study we undertook to analyze the performance of Silo distributed cache system. We developed a simulation model using the ns-2 [27] simulation package. In the case of static configuration of Silo, we assume that the entire cache layout is built at start-up and remains the same through the end of simulation. The static configuration would be useful enough to capture the important performance characteristics of Silo caching system.

The important parameters that control the Silo caching

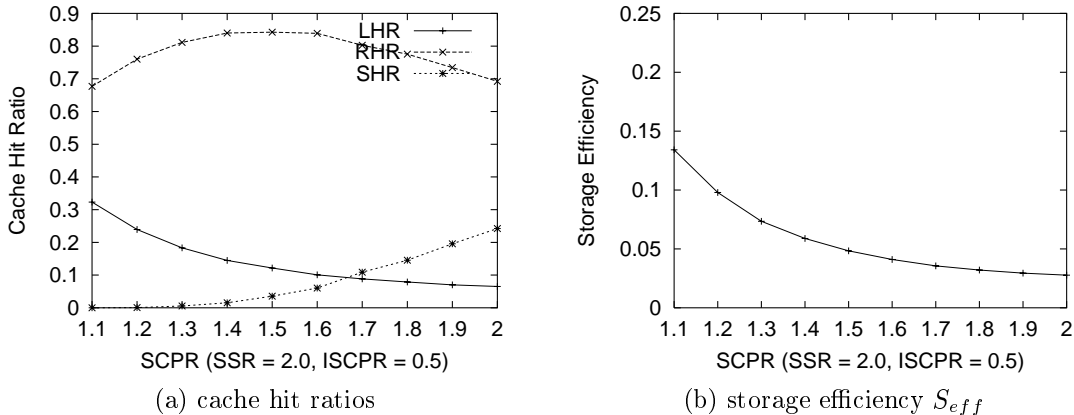


Fig. 7. Cache hit ratios and storage efficiency with varying segment caching probability ratio.

system are the following: (1) Segment Caching Probability Ratio (SCPR) β , (2) Segment Size Ratio (SSR) α , (3) Inter-Clip Segment Caching Probability Ratio (ISCPR) γ . In this section, we discuss how the control parameters affect the performance of Silo caching system. We measure the performance metrics defined in Section IV-C, namely: (1) Local hit ratio (LHR), (2) Remote hit ratio (RHR), (3) Server hit ratio (SHR), and (4) Number of switch-overs. Start-up latency is not an issue for Silo since SCs always cache the first segment in their local cache. We divide a clip in the Silo caching system into two regions: (1) *Silo roof* is the part where segment size is increasing, and (2) *Silo body* is the part where segment size is constant. We limit the maximum segment size in the Silo roof region to control the Silo layout. Table III shows the set of parameter values used throughout the simulation study. We assume that the SCs are interconnected in a simple mesh topology using network links of unlimited capacity. We set the aggregate arrival rate of client requests at each SC to 1 request per minute.

A. Varying Segment Caching Probability Ratio (SCPR)

First, we consider the performance metrics with varying SCPR β . Figure 7(a) shows the local, remote, and server hit ratios with varying SCPR. Local hit ratio (LHR) decreases as the SCPR increases. This is because, as SCPR increases, segment caching probability decreases more quickly and becomes smaller for later segments, resulting in less chance for the segments to be cached.

Remote hit ratio (RHR), however, shows an interesting behavior. It increases as SCPR increases till a certain threshold value and decreases as SCPR increases further. This can be explained by considering the system-wide hit ratio – the sum of local and remote hit ratios.²

²It is clear that the server hit ratio becomes (1 - system-wide hit ratio).

TABLE III
DEFAULT VALUES OF SIMULATION PARAMETERS

Parameter	Default Value
Initial seg size δ	50 MB
Max seg size of Silo roof	400 MB
Seg size of Silo body δ_c	50 MB
Clip size L	3 GB
Clip type	5 Mbps MPEG-2
Number of clips M	100
Number of SCs K	100
Client request	Poisson process
Clients' preference to clips	Zipf's distribution with $\gamma = 1$
Playback time distribution	Bimodal distribution [5]
Local cache size	0.8 - 4.5 GB depending on parameters

In general, the system-wide hit ratio decreases as SCPR increases. However, for the Silo cache system, the system-wide hit ratio is almost 1 in some range of SCPR, where at least one SC inside the system caches each segment. In that range, the Silo cache system can serve client requests either through local designated SC or remote SCs, keeping the server hit ratio (SHR) almost zero. For example, in Figure 7, with less than 8% of S_{eff} at SCPR of 1.3, the Silo caching system services almost 100% of stream requests through either local SC or remote SCs. The only difference, as SCPR varies in that range, is whether the requests are served through the local designated SC or remote SCs. The ratio between local and remote hit ratios is strongly related to S_{eff} . To serve more requests from the local SC, more segments should be cached locally, increasing S_{eff} as in Figure 7(b). As SCPR increases beyond the threshold value, however, the chance that the segments are

being cached within the Silo cache system decreases, resulting in increase of the server hit ratio.

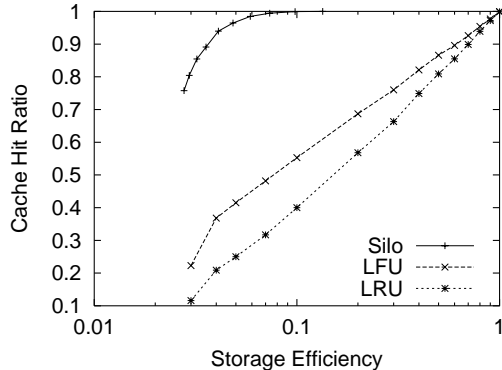


Fig. 8. Silo vs. traditional web cache with varying cache size.

In Figure 8, we compare the performance of Silo to that of traditional web caching system, where SCs independently cache entire clips upon client requests. For the traditional web caching system, we evaluate two different cache replacement algorithms: Least Recently Used (LRU) and Least Frequently Used (LFU). Compared to traditional web caching systems, Silo provides excellent cache hit ratio performance. At S_{eff} of 3%, Silo provides 3 – 8 times higher cache hit ratio than that of traditional web caching systems with LRU or LFU. As S_{eff} increases, cache hit ratio of Silo also grows much faster than that of traditional web cache systems, providing almost 100% cache hit ratio roughly at S_{eff} of 7%. To obtain near 100% of cache hit ratio, traditional web caching system needs ≈ 15 times more storage space than Silo. One more interesting point is LFU outperforms LRU over all range of S_{eff} . We also studied the impact of SCPR on the number of switch-overs averaged over all requests in the simulation. Figure 9 shows the average switch-over rate with varying SCPR. The average switch-over rate is the ratio of the total number of switch-overs to the total number of segment boundaries. Clearly, user requests that entirely fall into the first segment do not have any switch-overs. The average switch-over rate only accounts for user requests that require more than one segment. Figure 9 shows the interesting behavior of this metric. There are fewer switch-overs at both end regions and more in the middle region of SCPR. This can be explained by the behavior of the remote cache hit ratio (RHR). The region of the high rate of switch-overs matches with the region of the high remote hit ratio. In both extreme regions, where SCPR is either very small or very large, there are large fraction of cache hits from the local proxy or the origin server, increasing the chance of consecutive segments being fetched from either the local proxy or

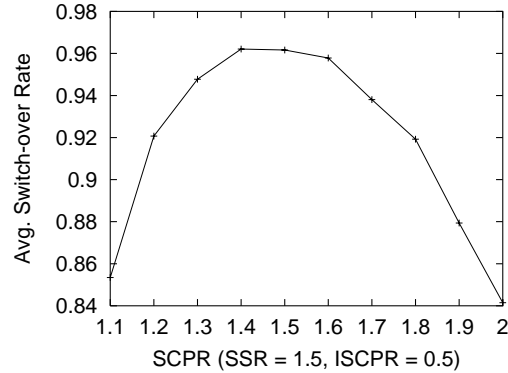


Fig. 9. Switch-overs with varying SCPR.

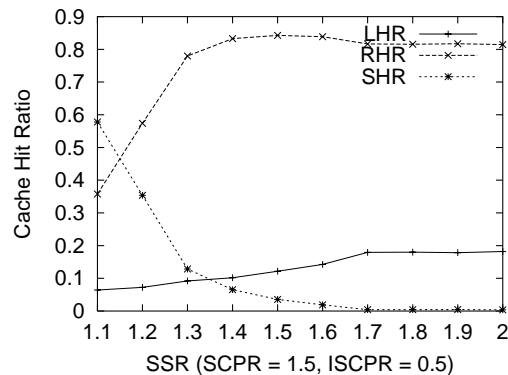


Fig. 10. Cache hit ratios with varying SSR.

the origin server. However, in the region where the remote cache hit ratio is high, there is less chance that consecutive segments are serviced from the same remote proxy, resulting in higher rate of switch-overs. One approach to solve this problem is segment clustering where more consecutive segments are clustered together in the same proxy.

In short, SCPR controls intra-clip segment replication by decreasing the ratio of the segment caching probability within a clip. In the following section, we discuss Segment Size Ratio (SSR) α that also controls intra-clip segment replication.

B. Varying Segment Size Ratio (SSR)

Figure 10 indicates that varying SSR α has the opposite effect of varying SCPR on the performance metrics. This is because, as SSR increases, the number of segments in the Silo roof region decreases. This results in smaller number of steps in decreasing segment caching probability and higher segment caching probability for the segments in the Silo body region. By controlling the number of decreasing steps in segment caching probability, SSR achieves intra-clip segment replication control. All the performance metrics show a plateau in the region of large values of SSR

(Figure 10). In this region, the number of segments in the Silo roof does not change since the increasing segment size reaches the maximum segment size with the same number of steps, although SSR changes.

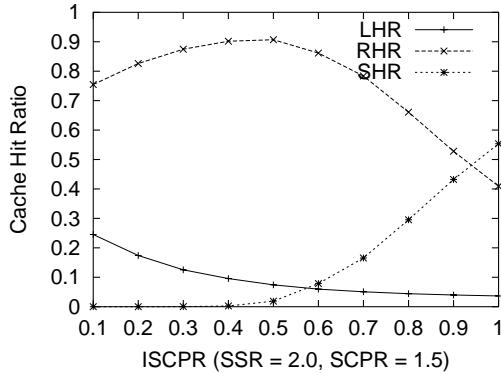


Fig. 11. Cache hit ratios with varying ISCPR.

C. Varying Inter-Clip Segment Caching Probability Ratio (ISCPR)

Although varying ISCPR γ shows similar performance results as the SSR and SCPR cases (Figure 11), the reason for this performance change is quite different from those cases. ISCPR controls inter-clip segment replication while SSR and SCPR control intra-clip segment replication. By varying ISCPR, the Silo system can control the preference to clips with different ranks. If the ISCPR is high, more storage space is allocated to popular clips. Since the inter-clip segment replication control through varying ISCPR is enforced for the entire clip, the performance metrics are more sensitive to ISCPR than SSR and SCPR as shown in Figures 7, 10 and 11.

The three control parameters (α, β, γ) can be used together to configure proper operational region of the Silo system, where the server hit ratio can be kept very small and the difference between the local and the remote hit ratios can be controlled with varying storage efficiency.

D. Load Balancing

We use the number of active streams at each proxy as the load balancing metric for the caching system. Figure 12 shows the average, maximum and standard deviation of the number of active streams per proxy. As discussed in Section IV-B, Silo replicates segments according to their popularity to achieve better load balancing performance. When there are multiple copies of a segment in Silo, one of the proxies needs to be selected.

We evaluate three remote proxy selection methods: *Best*, *Random*, and *Round-Robin*. The *Best* method selects the proxy that currently has the least load among the candidate

proxies. We use the number of current active streams as the indicator for load. We assume instant update of proxy state for the *Best* method for comparison purpose. The *Random* method randomly selects a target remote proxy from the candidate proxies. The *Round-Robin* method selects a target proxy from the candidate proxies in a round-robin manner. The *Best* method shows almost perfect load balancing performance as in Figure 12(a). However, the *Best* method requires prompt update of proxy states and its performance strongly depends on the freshness of state information. We use the *Best* method as an “ideal” case for comparison. The *Random* and the *Round-Robin* methods show comparable load balancing performance considering they do not use any state information. Figure 12(b) shows the standard deviation of the number of active streams on each proxy. The *Best* method shows the least variance of the number of active streams – a direct outcome of selecting the proxy with the minimum number of active streams. The *Random* and the *Round-Robin* methods show higher variance performance since they do not have any explicit load balancing mechanism.

E. Performance of Rainbow Cache Replacement Algorithm

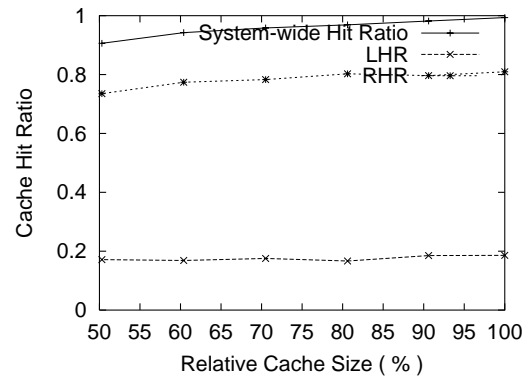


Fig. 13. Cache hit ratios with varying cache size limit.

We define *relative cache size (RCS)* of a SC as the ratio of its *available cache size (ACS)* to the *minimum cache size (MCS)* for which no local cache replacement occurs. For example, 100 % RCS represents no cache replacement, whereas any value of RCS less than 100 % will result in cache replacement. In the following evaluation we use the same RCS for all SCs. Using the parameters listed in Table III, and the analysis of Silo layout in Section IV, we calculate the MCS for each SC. In the simulation, we vary RCS from 50% to 100%, compute corresponding ACS, and measure LHR, RHR and system-wide hit ratio. The results are shown in Figure 13.

The Rainbow algorithm performs very well over a wide range of RCS. When RCS changes from 100 % to 50 %, the

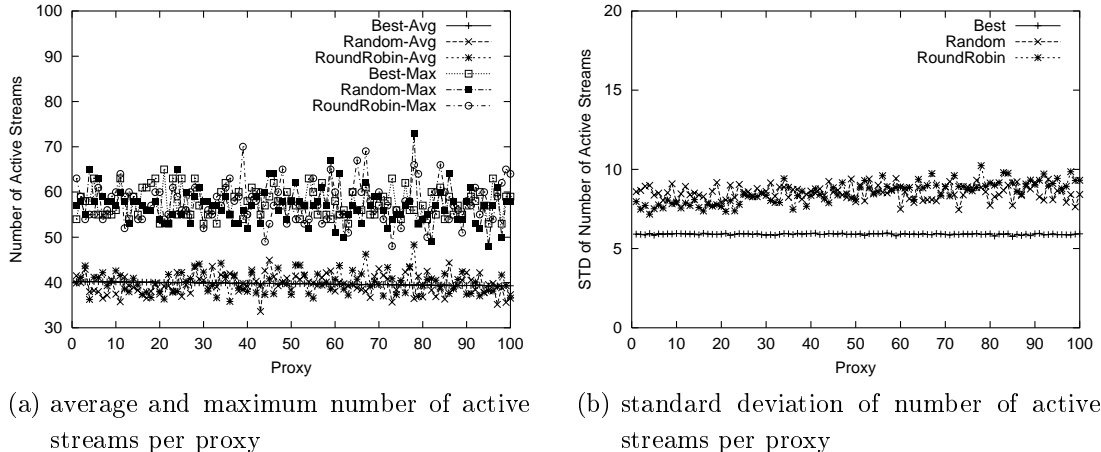


Fig. 12. Load balancing among proxies.

system-wide cache hit ratio only drops from 100 % to 90 %. In other words, 50 % reduction in cache size costs only 10 % loss of system-wide cache hit ratio.

The Rainbow algorithm also shows almost constant LHR value over the wide range of RCS from 50 % to 100 %. As RCS decreases, RHR behaves similar to system-wide cache hit ratio, decreasing very slowly until 78 % at RCS of 50 %. The constant LHR and the decreasing RHR results from the fact that the Rainbow algorithm replaces segments from the least popular ones in increasing order of segment caching potential. The replacement of unpopular segments does not affect LHR since these segments are cached only on a very small number of SCs and user requests to the segments are usually satisfied through remote cache hits. As soon as the unpopular segments are entirely replaced from the Silo system, user requests to the segments cannot be satisfied within the Silo system and have to be satisfied by the origin server. This results in decreasing RHR and increasing SHR. Clearly, our experiment shows that the Rainbow algorithm scales well with varying RCS.

F. Performance Comparison of RCache vs. Silo

To analyze the performance of RCache scheme, we consider a set of different RCache models: (1) RCache with first segment caching and inter-clip segment replication control (RCache-FI), where each SC always caches the first segment of all the clips and uses inter-clip segment replication control. (2) RCache with no first segment caching but with inter-clip segment replication control (RCache-NF), where each SC does not necessarily cache the first segment of all clips but uses inter-clip segment replication control. (3) RCache with no first segment caching and no inter-clip segment replication control (RCache-NFNI), where each SC does not necessarily cache the first segment

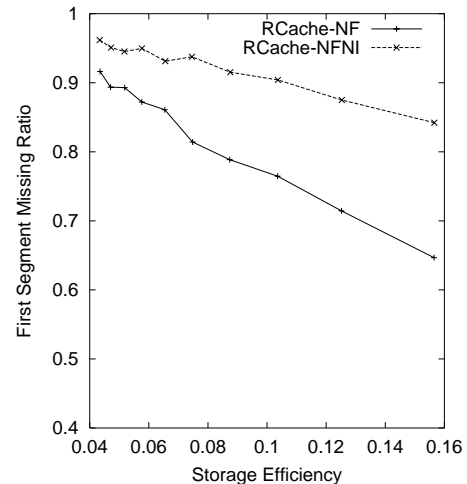


Fig. 15. First segment missing ratio.

of all clips and does not use inter-clip segment replication control.

Figure 14 shows the cache hit performance of various RCache configurations and Silo. For this comparison, we fix the segment size of RCache to be the same as the first segment size of Silo. For local cache hit ratio, Silo performs the best as shown in Figure 14(a). Both RCache-FI and RCache-NF schemes show the next best performance with slightly lower cache hit ratio, while RCache-NFNI shows the worst performance. Figure 14(a) shows an important fact that caching schemes with inter-clip segment replication control give higher local cache hit ratio. As in Figure 14(d), system-wide cache hit ratios for all the caching schemes are similar, while the schemes with first segment caching show slightly lower cache hit ratio. Since the first segment caching option takes more space for storing the first segments for all clips, less space is reserved for caching the remaining segments, resulting in slightly lower overall system-wide cache hit ratio. This shows a tradeoff

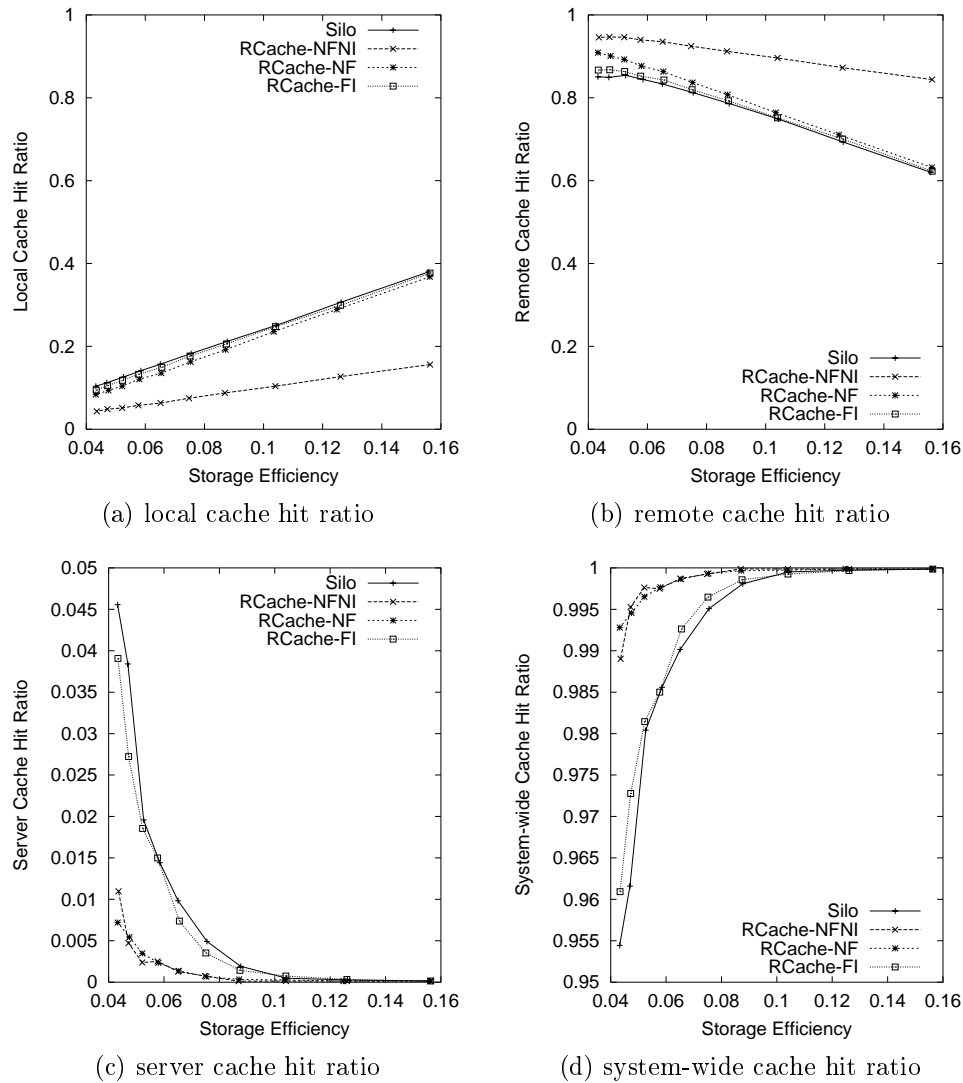


Fig. 14. Cache hit ratios with varying storage efficiency.

between start-up latency and system-wide cache hit ratio. Figure 15 shows that both RCache-NFNI and RCache-NF schemes experience high first segment missing ratio, resulting in high start-up latency to fetch the first segments from remote proxies. The caching schemes with first segment caching give better start-up latency performance at the expense of system-wide cache hit ratio. Figure 14(b) shows similar difference between the caching schemes with inter-clip segment replication control and the scheme without inter-clip segment replication control. Inter-clip segment replication control gives higher segment caching probability to the segments of popular clips, giving higher chance of local cache hit for the popular segments. Because of low local cache hit ratio, RCache-NFNI requires more segments from remote SCs.

Figure 16 compares the switch-over performance of these schemes. One big difference between Silo and RCache schemes is the number of switch-overs. Since Silo utilizes

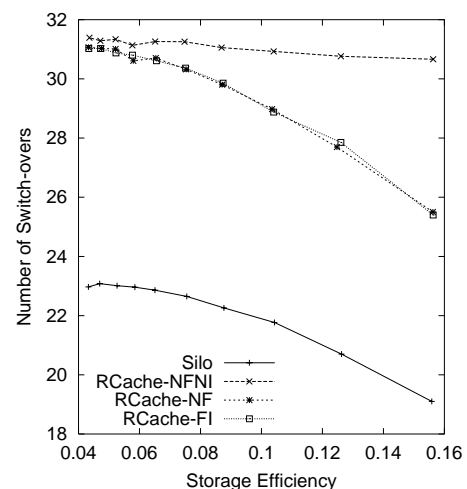


Fig. 16. Average number of switch-overs.

larger segment size for early segments, it allows less number of switch-overs for the large set of clip requests with small playback time [5], showing the smallest number of switch-overs. Figure 16 also shows that intra-clip segment replication control plays an important role in decreasing the number of switch-overs. RCache schemes with intra-clip segment replication control, RCache-FI and RCache-NF, result in smaller number of switch-overs than that of RCache-NFNI. The difference between RCache-NFNI and RCache-NF shows the effect of intra-clip segment replication control on the number of switch-overs.

In short, the first segment caching option gives better start-up latency performance while marginally sacrificing the system-wide cache hit ratio. The intra-clip segment replication control option selectively increases the chance of segments being cached in the system and thus, plays an important role in increasing the local cache hit ratio and in decreasing the number of switch-overs. Finally, increasing the segment size for the early segments in Silo also gives better switch-over performance.

VIII. CONCLUSIONS

In this paper, we focus on scalable, fault tolerant data placement and replacement techniques for caching multimedia streams in cooperating distributed caches. Specifically, we propose the following new schemes that work together: (1) *RCache*, a randomized, easy to implement distributed data layout scheme with good storage efficiency. (2) *Silo*, a novel, fault tolerant multimedia data layout scheme that accounts for clip popularity and intra-clip segment popularity. This scheme is superior to naive clip replication and RCache layout. Using various parameters associated with this layout, one can control its storage efficiency and tradeoff storage for performance. (3) *Rainbow*, a local replacement scheme based on the concept of segment access potential used for data replacement at each proxy. The Rainbow scheme captures the micro and macro popularity of clip segments and replaces only the least useful segments. (4) *Caching Token*, a dynamic *global* data replacement or redistribution scheme that exploits the existing data in the distributed cache to minimize distribution overhead. This scheme reduces network and origin server load dramatically and is straightforward to implement.

Our schemes together optimize various performance metrics, namely: (1) storage space, (2) start-up latency, (3) overhead from playback switch-overs, and (4) network bandwidth. Our analytical and simulation results show excellent performance and provide a means to achieve optimal tradeoff between various design parameters such as the number of caches, origin server bandwidth, and layout parameters. Silo provides 3 - 8 times as high cache hit

ratio as that of traditional web caching system while utilizing the same amount of storage space. Silo also provides almost 100 % of cache hit ratio with less than 10 % of storage space required for traditional web cache to achieve the same cache hit ratio.

ACKNOWLEDGMENTS

A significant portion of the work reported in this paper was carried out during the first author's summer internships at the Center for Networking Research in Lucent Bell Labs during Summer 1999 and Summer 2000. The first author would like to thank Dr. Krishan Sabnani, Vice President, Center for Networking Research for these valuable opportunities and for continued help and encouragement. The authors would like to thank their colleagues Michael Vernick and Adishesu Hari for insightful discussions and comments on the early versions of this manuscript. Authors would also like to thank the editors and anonymous reviewers for their valuable feedback that has improved the final version.

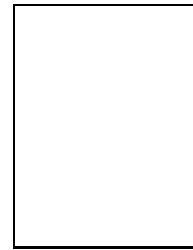
REFERENCES

- [1] www.lucent.com/businessunit/icdd.html, www.akamai.com, www.cacheflow.com, www.inktomi.com, www.infolibria.com, www.networkappliance.com.
- [2] <http://www.ncube.com>.
- [3] Acharya, S. and Smith, B. "An Experiment to Characterize Videos on the World Wide Web", *Multimedia Computing and Networking 1998*.
- [4] Acharya, S. and Smith, B., "Middleman: A Video Caching Proxy Server", *Proc. of NOSSDAV'00*, June, 2000.
- [5] Acharya, S., Smith, B., and Parnes, P., "Characterizing User Access to Videos on the World Wide Web", *Multimedia Computing and Networking 2000*.
- [6] Aggarwal, C., Wolf, J., L., and Yu., P. S., "A Permutation-based Pyramid Broadcasting Scheme for Video-on-demand Systems", *Proc. of IEEE International Conference on Multimedia Computing and Systems*, Jun, 1996.
- [7] Almeroth, K. and Ammar, M., "On the Use of Multicast Delivery to Provide a Scalable and Interactive Video-on-Demand Service", *Journal of Selected Areas of Communication*, Aug, 1996.
- [8] Bernhardt, C., and Biersack, E., "The Server Array: A Scalable Video Server Architecture", *High-Speed Networks for Multimedia Applications*, editors, Danthine, A., Ferrari, D., Spaniol, O., and Effelsberg, W., Kluwer Academic Press, 1996.
- [9] Bolosky, W., et al., "The Tiger Video File Server", *Proc. of NOSSDAV96*, Apr, 1996.
- [10] Bowman, C., M., Danzig, P., B., Hardy, D., R., et al. "Harvest: A Scalable, Customizable Discovery and Access System", Tech Report, CU-CS-732-94, University of Colorado, Boulder, USA, 1994.
- [11] L. Breslau, P. Cao, L. Fan, G. Phillips and S. Shenker, "Web Caching and Zipf-like Distributions: Evidence and Implications", *Proc. of IEEE Infocom'99*, Mar, 1999.
- [12] Buddhikot, M., Parulkar, G., and Cox, J., R., Jr., "Design of a Large Scale Multimedia Storage Server", *Journal of Computer Networks and ISDN Systems*, pp. 504-517, December 1994.
- [13] Chan, S. and Tobagi F., "Distributed Servers Architecture for Networked Video Services", *IEEE/ACM Transactions on Networking*, pp. 125-136, Apr, 2001.

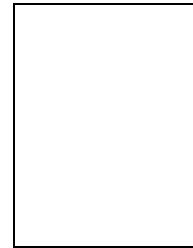
- [14] Chankhunthod, A., Danzig, P., B., Neerdaels, C., Schwartz, M., F., and Worrell, K., J., "A Hierarchical Internet Object Cache", Proc. of the 1996 Usenix Technical Conference, 1996.
- [15] Chen, M., Kandlur, D., and Yu, S., P., "Support for Fully Interactive Playout in a Disk-Array-Based Video Server", *Proc. of Second International Conference on Multimedia, ACM Multimedia'94*, 1994.
- [16] Chen, P., et al., "RAID: High-Performance, Reliable Secondary Storage," *ACM Computing Surveys*, Jun, 1994.
- [17] Hua, K., Chai, Y. and Sheu, S., "Patching: A Multicast Technique for True Video-on-Demand Services", *Proc. of ACM Multimedia'98*, Sept, 1998.
- [18] Hua, K., A., Sheu, S., "Skyscraper Broadcasting: A New Broadcasting Scheme for Metropolitan Video-on-Demand Systems", *Proc. of ACM SIGCOMM*, Sept, 1997.
- [19] Juhn, L-S., and Tseng, L-M., "Harmonic Broadcasting for Video-on-Demand Service", *IEEE Transaction on Broadcasting*, Vol. 43, No. 3, pp, 268-271, Sept, 1997.
- [20] Juhn, L-S, and Tseng, L-M., "Staircase Data Broadcasting and Receiving Scheme for Hot Video Service", *IEEE Transaction on Consumer Electronics*, Vol. 43, No. 4, pp. 1110-1117, Nov, 1997.
- [21] Partridge, C., Mendex, T., and Milliken, W., "Host Anycasting Service", *RFC 1546*, Nov, 1993.
- [22] Ramesh, S., Rhee, I., and Guo, K., "Multicast with Cache (Mcache): An Adaptive Zero Delay Video-on-Demand Service", *IEEE Transactions on Circuits and Systems for Video Technology*, Vol. 11, No. 3, Mar, 2001.
- [23] Sen, S., Rexford, J., and Towsley, D., "Proxy Prefix Caching for Multimedia Streams", *Proc. of IEEE Infocom'99*, Mar, 1999.
- [24] Shenoy, P., "Symphony: An Integrated Multimedia File System", *Doctoral Dissertation*, Dept. of Computer Science, University of Texas at Austin, Aug, 1998.
- [25] Tewari, R., Vin, H., Dan, A., and Dias, D., "Resource Based Caching for Web Servers", *ACM Multimedia Systems Journal*, 1999.
- [26] Tewari, R., Mukherjee, R., Dias, D., M., and Vin, H., M., "Design and Performance Tradeoffs in Clustered Video Servers", *Proc. of the IEEE ICMCS'96*, May, 1996.
- [27] UCB/LBNL/VINT, "Network Simulator, ns", <http://www.isi.edu/nsnam/ns/>.
- [28] Viswanathan, S., and Imielinski, T., "Metropolitan Area Video-on-Demand Service using Pyramid Broadcasting", *ACM Multimedia Systems*, Vol. 4, pp. 197-208, 1996.
- [29] Wessels, D. and Claffy, K., "Internet Cache Protocol (ICP)", version 2, *IETF RFC 2186*, Sept, 1997.
- [30] Wolman A., Voelker, G., Sharma N., Cardwell, N., Karlin A., and Levy, M., "On the Scale and Performance of Cooperative Web Proxy Caching", *Proceedings of ACM SOSP'99*, Dec, 1999.

Youngsu Chae (S 1994) received the B.S. (1994) and the M.S. (1996) degrees in computer science from Pohang University of Science and Technology, Pohang, Korea. He is currently a Ph.D. candidate in the College of Computing, Georgia Institute of Technology. His thesis work focuses on a distributed stream caching architecture and a scalable message delivery service that generalizes

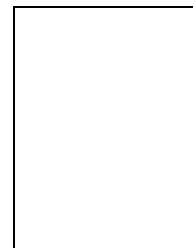
both anycast and multicast. His research interests include multimedia stream caching, active networking, anycast and multimedia routing, and multimedia streaming server architecture.



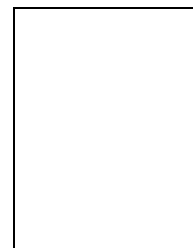
Katherine Guo received her B.A. in mathematics, B.S. in computer science from the University of Texas at Austin in 1992, M.S. and Ph.D. in computer science from Cornell University in 1998. She is currently a Member of Technical Staff in the Center for Networking Research at Bell Laboratories, Lucent Technologies. She has published in the areas of computer networks, distributed systems and multimedia networking. Her current interest includes systems and protocols for wireless networks, web caching and network support for streaming media.



Milind M. Buddhikot is a Member of Technical Staff in the Center for Networking Research at Lucent Bell Labs. His current research interests are in the areas of systems and protocols for public wireless networks, MPLS path routing, and multimedia messaging and stream caching. Milind holds a Doctor of Science (D. Sc.) in computer science (July 1998) from Washington University in St. Louis, and a Master of Technology (M.Tech.) in communication engineering (December 1988) from the Indian Institute of Technology (I.I.T), Bombay. He has authored over 20 research papers and 5 patent submissions in design of multimedia systems and protocols, layer-4 packet classification, and MPLS path routing. Milind served as a co-guest-editor of IEEE Network magazine's March 2001 Special issue on "Fast IP Packet Forwarding and Classification for Next Generation Internet Services".



Subhash Suri is a Professor of Computer Science at the University of California, Santa Barbara. Prior to joining UCSB, he was an Associate Professor at Washington University (1994-2000) and a Member of Technical Staff at Bellcore (1987-1994). Prof. Suri received a Ph.D. in computer science from the Johns Hopkins University in 1987. Professor Suri has published over 80 research papers in design and analysis of algorithms, Internet commerce, computational geometry, computer networking, combinatorial optimization, and graph theory. He serves on the editorial board of the Computational Geometry journal, and has served on numerous panels and symposium program committees. Prof. Suri has been awarded several research grants from the National Science Foundation.



Ellen W. Zegura (M 1993) received the B.S. degrees in computer science and electrical engineering (1987), the M.S. degree in computer science (1990) and the D.Sc. degree in computer science (1993) all from Washington University, St. Louis, Missouri. She is currently an Associate Professor in the College of

Computing at Georgia Institute of Technology. Her research interests include active networking, server selection, anycast and multicast routing, and modeling large-scale internetworks. Dr. Zegura is currently an editor for IEEE/ACM Transactions on Networking. Her email address is ewz@cc.gatech.edu.

LIST OF FIGURES

1	Distributed streaming cache architecture.	3
2	Bimodal distribution of clip segment requests.	5
3	Generalized Silo layout: segment size and probability variation.	7
4	Local and remote hit ratios with 100 Silo proxies ($\delta = 50$ MB).	9
5	Cache replacement order of Rainbow algorithm.	10
6	Silo cache layout change: rank up.	11
7	Cache hit ratios and storage efficiency with varying segment caching probability ratio.	13
8	Silo vs. traditional web cache with varying cache size.	14
9	Switch-overs with varying SCPR.	14
10	Cache hit ratios with varying SSR.	14
11	Cache hit ratios with varying IS CPR.	15
13	Cache hit ratios with varying cache size limit.	15
12	Load balancing among proxies.	16
15	First segment missing ratio.	16
14	Cache hit ratios with varying storage efficiency.	17
16	Average number of switch-overs.	17

LIST OF TABLES

I	Details of basic Silo data layout	7
II	General Silo cache layout for a clip with rank i	8
III	Default values of simulation parameters	13