

# Scalable Consistency Maintenance in Content Distribution Networks Using Cooperative Leases\*

Anoop Ninan, Purushottam Kulkarni<sup>§</sup>, Prashant Shenoy<sup>§</sup>,  
Krithi Ramamritham<sup>†</sup> and Renu Tewari<sup>‡</sup>

Software Engineering Division  
EMC Corporation  
Hopkinton, MA 01748  
ninan\_anoop@emc.com

<sup>§</sup>Department of Computer Science,  
University of Massachusetts  
Amherst, MA 01003  
{purukulk,shenoy,krithi}@cs.umass.edu

<sup>‡</sup>IBM Research Division  
Almaden Research Center  
San Jose, CA 95120  
tewarir@us.ibm.com

## Abstract

In this paper, we argue that cache consistency mechanisms designed for stand-alone proxies do not scale to the large number of proxies in a content distribution network and are not flexible enough to allow consistency guarantees to be tailored to object needs. To meet the twin challenges of scalability and flexibility, we introduce the notion of *cooperative consistency* along with a mechanism, called *cooperative leases*, to achieve it. By supporting  $\Delta$ -consistency semantics and by using a single lease for multiple proxies, cooperative leases allows the notion of leases to be applied in a flexible, scalable manner to CDNs. Further, the approach employs application-level multicast to propagate server notifications to proxies in a scalable manner. We implement our approach in the Apache web server and the Squid proxy cache and demonstrate its efficacy using a detailed experimental evaluation. Our results show a factor of 2.5 reduction in server message overhead and a 20% reduction in server state space overhead when compared to original leases albeit at an increased inter-proxy communication overhead.

**Index Terms** — Dynamic data, Data consistency, Data dissemination, World Wide Web, Scalability, Leases, Push, Pull, Content Distribution Networks

## 1 Introduction

### 1.1 Motivation

The past decade has seen a dramatic increase in the popularity and use of the World Wide Web. Numerous studies have shown that web accesses tend to be non-uniform in nature, resulting in

---

\*This research was supported in part by a NSF Career award CCR-9984030, NSF grants CCR-0098060, EIA-0080119, IBM, Intel, EMC, Sprint and the University of Massachusetts.

<sup>†</sup>Also affiliated with the Department of Computer Science and Engg., Indian Institute of Technology, Powai, Bombay.

(a) hot-spots of server and network load and (b) increases in the latency of web accesses. Content distribution networks have emerged as a possible solution to these problems. A *Content Distribution Network (CDN)* consists of a collection of proxies that act as intermediaries between the origin servers and the end users. Proxies in a CDN cache frequently accessed data from origin servers and serve requests for these objects from the proxy closest to the end-user. By doing so, a CDN has the potential to reduce the load on origin servers and the network and also improve client response times.

Architectures employed by a CDN can range from tree-like hierarchies [30] to clusters of cooperating proxies that employ content routing to exchange data [13]. From the perspective of endowing proxies with content, proxies within a CDN can either pull web content on-demand, prefetch popular content, or have such content pushed to them [11]. Mechanisms for locating the best proxy to service a user request range from Anycast [16] to DNS-based selection [22]. Regardless of the exact architecture and mechanisms, an important issue that must be addressed by a CDN is that of *consistency maintenance*. Since web pages tend to be modified at origin servers, cached versions of these pages can become inconsistent with their server versions. Using inconsistent (stale) data to service user requests is undesirable, and consequently, a CDN should ensure the consistency of cached data with the server by employing suitable techniques.

The problem of consistency maintenance is well studied in the context of a single proxy and several techniques such as time-to-live (TTL) values [5], client-polling, server-based invalidation [4], adaptive refresh methods [23, 25], and leases [27] have been proposed. In the simplest case, a CDN can employ these techniques at each individual proxy — each proxy assumes responsibility for maintaining consistency of data stored in its cache and interacts with the server to do so independently of other proxies in the CDN. Since a typical content distribution network consists of hundreds or thousands of proxies (e.g., the Akamai CDN has a footprint of more than 13,500 servers [1]), requiring each proxy to maintain consistency independently of other proxies is not scalable from the perspective of the origin servers (since the server will need to individually interact with a large number of proxies). Further, consistency mechanisms designed from the perspective of a single proxy (or a small group of proxies) do not scale well to large CDNs. The leases approach, for instance, requires the origin server to maintain per-proxy state for each cached object. This state space will grow with the number of objects a proxy caches and with the number of proxies that cache an object. These arguments motivate the need for designing novel consistency mechanisms that scale to large CDNs and is the focus of this paper.

## 1.2 Research Contributions

In this paper, motivated by the need to reduce the load at origin servers and to scale to a large number of proxies, we (a) argue that  $\Delta$ -consistency semantics are appropriate for CDNs because they allow

the tailoring of consistency guarantees to the nature of objects and their usage, and (b) introduce the notion of *cooperative consistency* along with a mechanism, called *cooperative leases*, to achieve it. Cooperative consistency enables proxies to cooperate with one another to reduce the overheads of consistency maintenance. By supporting  $\Delta$ -consistency semantics and by using a single lease for multiple proxies, our cooperative leases mechanism allows the notion of leases to be applied in a scalable manner to CDNs. Another advantage of our approach is that it employs application-level multicast to propagate server-notifications of modifications to objects, which reduces server overheads. We address the various design issues that arise in a practical realization of cooperative leases and then show how to implement the approach in the Apache web server and the Squid proxy cache using HTTP/1.1. Our work focuses more on cache consistency mechanisms and semantics, and less on the protocol details (i.e., message formats) used for sending invalidations. Finally, we experimentally demonstrate the efficacy of our approach using trace-driven simulations and the prototype implementation. Our results show that cooperative leases can reduce the number of server messages by a factor of 2.5 and the server state by 20% when compared to original leases, albeit at an increased proxy-proxy communication overhead.

The rest of this paper is structured as follows. Section 2 defines the problem of consistency maintenance in CDNs and presents our cooperative leases approach. We examine various design issues in instantiating cooperative leases in Section 3. Section 4 discusses the details of our prototype implementation. Section 5 presents our experimental results. Section 6 discusses related work, and finally, Section 7 presents some concluding remarks.

## **2 Cache Consistency: Semantics, Mechanisms, and Architecture**

### **2.1 $\Delta$ -Consistency: Consistency Semantics for Cached Objects**

Objects cached within a content distribution network need different levels of consistency guarantees depending on their characteristics and user preferences. For instance, users may be willing to receive slightly outdated versions of objects such as news stories, but are likely to demand the most up-to-date versions of “critical” objects such as financial information and sports scores. Typically, the stronger the desired consistency guarantee for an object, the higher the overheads of consistency maintenance. For reasons of flexibility and efficiency, rather than providing a single consistency semantics to all cached objects, a CDN should allow the consistency semantics to be tailored to each object or a group of related objects.

One possible approach for doing so is to employ  $\Delta$ -consistency semantics [25].  $\Delta$ -consistency requires that a cached version of an object is never out-of-date by more than  $\Delta$  time units with its server version. The value of  $\Delta$  determines the nature of the provided guarantee — the larger the value of  $\Delta$ , the weaker the consistency guarantee (since the object could be out of date by up

to  $\Delta$  time units at any instant). An advantage of  $\Delta$ -consistency is that it provides a quantitatively characterizable guarantee by virtue of providing an upper bound on the amount by which a cached object could be stale (unlike certain mechanisms that only provide qualitative guarantees). Another advantage is the flexibility of choosing a different value of  $\Delta$  for each object, allowing the guarantee to be tailored on a per-object basis. Finally, strong consistency — a guarantee that a cached object is never out-of-date with the server version — is a special case of  $\Delta$ -consistency with  $\Delta = 0$ .<sup>1</sup>

Due to the above advantages, in this paper, we assume a CDN that provides  $\Delta$ -consistency semantics. Next, we present a consistency mechanism to provide  $\Delta$ -consistency and then discuss its implementation in a CDN.

## 2.2 Cooperative Leases: A Cache Consistency Mechanism for CDNs

A consistency mechanism employed by a CDN should satisfy two key requirements: (i) *scalability*: the approach should scale to a large number of proxies employed by the CDN and should impose low overheads on the origin servers and proxies, and (ii) *flexibility*: the approach should support different levels of consistency guarantees. We now present a cache consistency mechanism that satisfies these requirements. Our approach is based on a generalization of *leases* [12].

In the original leases approach [12], the server grants a lease to each request from a proxy. The lease denotes the interval of time during which the server agrees to notify the proxy if the object is modified. After the expiration of the lease, the proxy must explicitly renew the lease in order to receive further notifications. Formally, a lease is a tuple  $\{O, p, d\}$  maintained by the server, where the server agrees to notify proxy  $p$  of all updates to an object  $O$  during time interval  $d$ .

The leases approach has two drawbacks from the perspective of a CDN. First, leases provide strong consistency semantics by virtue of notifying a proxy of *all* updates to an object. As argued earlier, not all objects cached within a CDN need such stringent guarantees. Second, leases require the server to maintain state for each proxy caching an object; the resulting state space overhead can be excessive for large CDNs. Thus, leases do not scale well to busy servers and large CDNs.

To alleviate these drawbacks, we generalize leases along two dimensions:

1. We add a *notification rate parameter*  $\Delta$  to leases that indicates the rate,  $1/\Delta$ , at which the server agrees to notify a proxy of updates to an object. This enhancement allows a server to relax the consistency semantics provided by leases from strong consistency to  $\Delta$ -consistency — a proxy is notified of updates at most once every  $\Delta$  time units (instead of after every update) and no later than  $\Delta$  time units after an update. Using  $\Delta = 0$  reverts to the original leases approach (i.e., strong consistency), while a positive value of  $\Delta$  allows the server to provide weaker consistency guarantees (and correspondingly reduces the number of notifications sent to a proxy).

---

<sup>1</sup>Implementing true strong consistency requires the server to first send notifications to all proxies caching the object to be updated; the write at the server commits only after receiving acknowledgments from these proxies.

2. We allow a server to grant a single lease collectively to a group of proxies, instead of issuing a separate lease to each individual proxy.<sup>2</sup> For each cached object, the proxy group designates an invalidation proxy, referred to as the *leader*, that is responsible for all lease-related interactions with the server. The leader of a group manages the lease on behalf of all the proxies in the group. Since a leader is selected per object, no single proxy becomes the bottleneck. Moreover, the server only notifies the leader upon an update to the object; the leader is then responsible for propagating this notification to other proxies in the group that are caching the object. Such an approach has two significant advantages: (i) it reduces the the amount of state maintained at a server (by using a single lease to represent a proxy group instead of an individual proxy); and (ii) it reduces the number of notifications that need to be sent by the server (by offloading some of notification burden to leader proxies).

We refer to the resulting approach as *cooperative leases*. Formally, a cooperative lease is a tuple  $\{O, G, L, d, \Delta\}$  where the server agrees to notify the leader  $L$  representing proxy group  $G$  of any updates to the object  $O$  once every  $\Delta$  time units for an interval  $d$ . While leases is a pure server-based approach to cache consistency, cooperative leases require both the server and the proxy (especially the leader) to participate in consistency maintenance. Hence this approach is more scalable when compared to original leases, and thus, more suited to CDN environments.

### 2.3 System Model of Cooperative Leases

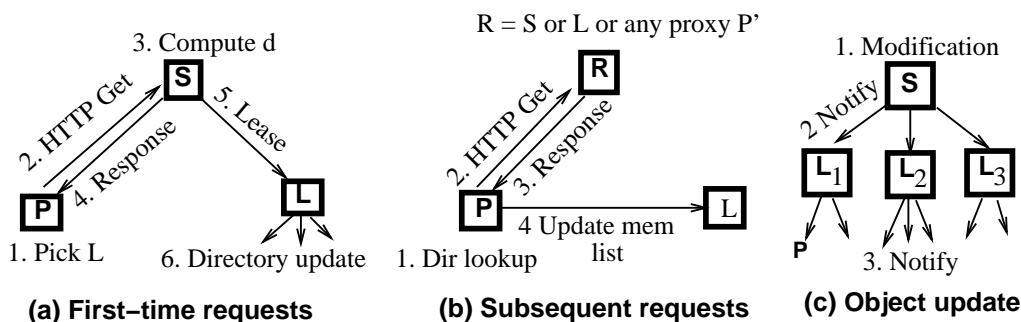
Before discussing the implementation of cooperative leases in CDNs, we present the system model assumed in this paper. A content distribution network is defined to be a collection of proxies that cache content stored on origin servers. For the purposes of maintaining consistency, proxies within the CDN are assumed to be partitioned into non-overlapping groups referred to as *regions* (issues in doing so are beyond the scope of this paper). Proxies within a region are assumed to cooperate with one another for maintaining consistency of cached objects. Cooperative consistency *is orthogonal to cooperative caching* — whereas the latter involves sharing of cached data to service user requests, the former involves cooperation solely for maintaining consistency of data cached by proxies within a region. Further, the organization of proxies into regions is limited to consistency maintenance; a different overlay topology can be used for exchanging data and meta-data within the CDN. Each proxy in a region is assumed to maintain a directory of mappings between the cached object and its corresponding leader (and possibly other information required by the CDN). Several directory schemes such as hint caches [24] and bloom filters [8] have been proposed to efficiently maintain such information. Another approach is to use a simple consistent hashing [15] based scheme to

---

<sup>2</sup>In addition, it is also possible for a lease to collectively represent multiple objects. Techniques for doing so are studied in [28].

**Table 1:** Design Considerations

<i>Event</i>	<i>Design Decision</i>	<i>Refer</i>
Issue a new lease	Choose a leader Choose $d$ and $\Delta$	Sec 3.1 Sec 3.2
Lease expiry	Lease renewal policy	Sec 3.3
Object changes	Send update/invalidate	Sec 3.4
Global cache miss	Issue a new lease	Sec 2.4
Local cache miss	Update membership list	Sec 2.4



**Figure 1:** Interactions between servers (S), leaders (L) and proxies (P) in cooperative leases

determine the mapping between an object and the proxy that acts as the leader. Here a hashing function is used to hash on both the unique object identifier and the list of proxy identifiers to determine the best match. Although this approach reduces the flexibility in assigning the leader for an object, it reduces the space and the message exchange overhead. Any of the above schemes suffices for our purpose.

## 2.4 Operations of Cooperative Leases

Cooperative leases can be instantiated as follows (see Table 1 and Figure 1).

*First-time requests:* When an object is requested for the first time within a region (i.e., upon a global cache miss), a leader needs to be chosen for the object. The proxy receiving the request executes a leader selection algorithm to pick a leader. Different cached objects can have different leaders — the cooperative leases approach attempts to distribute leader responsibilities across proxies in the region for load balancing purposes. Specific techniques for leader selection are discussed in Section 3.1. After choosing a leader, the proxy issues a HTTP request to the server and piggy-backs the leader information with the message; the message can also include optional information such as the desired notification rate parameter  $\Delta$ . The requested object is then sent to the proxy and the lease is sent to the leader and optionally, a copy of the object. As will be clear later, the presence

of a copy of the object at the leader enables us to perform certain optimizations. The leader proxy then broadcasts a directory update to all proxies in the region indicating it is the designated leader for the object. The leader also maintains a *membership list* consisting of all proxies caching the object; the list is initialized to the proxy that requested the object. Figure 1(a) depicts these interactions.

From this point on, the leader is responsible for renewing the lease on behalf of proxies in the region and for terminating the lease when proxies are no longer interested in the object. Policies for doing so are discussed in Section 3.3.

A minor modification of the protocol in Figure 1(a) is to have the proxy communicate with the server to request the object and communicate with the leader to obtain a lease on its behalf. If the leader does not already have a lease for that object it forwards the request to the server. The protocol described in Figure 1(a) has the advantage of lower message overhead for popular objects and integrates well with cooperative caching.

An alternate approach would have been to place the leader in the HTTP request path from the proxy to the server. This, however, suffers from the drawback of adding to the cache miss latency and increases the load on the leader as the entire object needs to be stored and forwarded by the leader. Secondly, it does not scale well in a multi-level proxy organization.

*Subsequent requests:* For each subsequent request to the object within the region, a proxy first examines its local cache. In the event of a cache hit, the proxy services the request using locally cached data. In the event of a local cache miss, the proxy can pursue one of several possible alternatives. It can either fetch the object from the server or consult its directory for a list of proxies caching the object and fetch the object from one such proxy (the exact proxy that is chosen may depend on the information in the directory and metrics such as proximity). Since the focus of our work is on consistency maintenance, the cooperative leases approach does not mandate the use of cooperative caching or require a particular policy for cooperative caching — the proxy is free to fetch the object from any entity that has the object, including the server. The only requirement imposed by cooperative leases is that the proxy notify the leader of its interest in the object. The leader then updates the membership list for the object and starts forwarding any subsequent notifications from the server to this proxy. Figure 1(b) depicts these interactions.

Observe that a proxy can optimize the overheads of the above operations by just fetching the object from the leader. If the leader cached the most recent version of the object (recall that a copy of the object could be optionally sent to the leader), this eliminates the need to send two different messages, one to fetch the object and the other to notify the leader of this fetch.

*Updates to the Object:* In the event the object is modified at the server, each proxy caching the object needs to be notified of the update. To do so, the origin server first notifies the leader of each region caching the object, subject to the notification rate parameter  $\Delta$ . The notification consists of either a cache invalidate or a new version of the object (see Section 3.4 for details).

Each leader in turn propagates this notification to every proxy in the region caching the object (i.e., to all proxies in the membership list). Depending on the type of notification, proxies then either invalidate the object in the cache or replace it with the newer version. Our approach is equivalent to using *application-level multicast* for propagating notifications; the membership list and the leader constitute the “multicast group”. Figure 1(c) depicts these interactions.

For simplicity of exposition, the above discussion assumed that the application-level multicast tree within a region is only two levels deep, spanning from servers to leaders and from leaders to proxies. Whereas a one-level hierarchy suffices for small regions (likely to be the common case), large CDNs are likely to constitute multiple regions and multi-level proxy hierarchies. Cooperative leases can be easily extended to multi-level proxy hierarchies and multi-region CDNs; techniques for doing so are discussed in Section 3.5 and Section 3.6.

### 3 Design Considerations for Cooperative Leases

In this section, we discuss design issues that arise when implementing cooperative leases in a CDN. These include leader selection, selecting lease duration and notification rate, policies for lease renewal, sending invalidations versus updates (see Table 1). We also present techniques to instantiate cooperative leases when CDN proxies are organized in multi-level hierarchies and/or multiple regions.

#### 3.1 Leader Selection

We consider two different policies for choosing a leader when an object is accessed for the first time within the region. In the simplest case, the proxy that receives this request can become the leader for the object. Since many web objects tend to be accessed by only one user [3], an advantage of this approach is that only one proxy is involved in consistency maintenance for such objects (since the proxy caching the object is also the leader). This results in lower communication overheads. A drawback, however, is that the approach has poor load balancing properties — leader responsibilities can become unevenly distributed if a small subset of proxies receive a disproportionate number of first-time requests. Additionally, if several proxies receive simultaneous first-time requests to an object, it is possible for multiple proxies to declare themselves the leader. Such duplication can be prevented using tie-breaking rules or by having the server perform additional error checks before issuing a new lease to a region.

An alternate approach is to employ a hashing function to determine the leader for an object. To illustrate, the leader could be determined based on the MD5 hash of the object URL (i.e.,  $L = MD5(URL) \bmod N$ , where  $N$  is the number of proxies in the region). More complex hashing functions can take other factors, such as the current load on proxies, into account in addition to the



URL [15]. An advantage of the hash-based approach is that it has good load balancing properties and results in a more uniform distribution of leader responsibilities across proxies. A limitation though is that it can impose a larger communication overhead than our first approach. Since the leader can be potentially different from proxies caching the object, additional directory updates, server notifications and lease management messages need to be exchanged between these proxies, which increases communication overheads. Section 5 quantitatively evaluates the tradeoffs of these two policies.

### 3.2 Choosing the Lease Duration and Notification Rate

Two key factors that influence the performance of cooperative leases are the lease duration  $d$  and the notification rate parameter  $\Delta$ . In a recent work, we investigated techniques for determining the lease duration for the original leases approach and proposed policies for computing  $d$  based on parameters such as object popularity, write frequency, and server/network load [7]. Since similar policies can be employed for computing the lease duration  $d$  in CDNs, we do not consider this issue any further.

The notification rate can either be specified by the user (or proxy), or computed by the server. In the former approach, the end-user or the proxy specifies a tolerance  $\Delta$  based on the desired consistency guarantee. The server then grants a lease with this  $\Delta$  if it has sufficient resources to meet the desired tolerance. In the latter approach, the server computes an appropriate notification rate based on various system parameters while issuing a new lease. For instance, the server could compute  $\Delta$  based on the server or network load. Rather than rejecting a request for a lease during periods of heavy load, the server could continue to grant leases but provide weaker guarantees (i.e., use a larger  $\Delta$ ). To illustrate,

$$\Delta = \begin{cases} 0 & \text{load} < LWM \\ c \cdot \text{load} & LWM \leq \text{load} < HWM \\ d & \text{load} \geq HWM \end{cases} \quad (1)$$

where  $c$  is a constant (based on the values of the load and the watermarks) and  $LWM$  and  $HWM$  denote low and high watermarks (thresholds), respectively. Here the server notifies leaders of all updates at low loads.  $\Delta$  is increased linearly with the load at moderate utilizations and is finally set to the lease duration at high loads ( $d$  is the least possible notification rate, since at least one update should be sent in each lease duration).

To analyze the message overhead due to a particular choice of these parameters, consider an object with a lease duration  $d$  and notification rate  $1/\Delta$ . Let  $\hat{W}$  denote the non-consecutive write frequency for the object (a non-consecutive write is defined to be a write followed by a read, while a consecutive write is one followed by another write; the non-consecutive write frequency is computed by counting all consecutive writes to an object as a single write—only one notification is necessary

for each set of consecutive writes). In such a scenario, the number of notifications sent by the server within each lease duration is  $\min(\hat{W}, \frac{1}{\Delta}) \cdot d$ . In the event the notification is an invalidate, each such message will trigger an HTTP GET message upon a subsequent read at one of the proxies in the region, resulting in another  $\min(\hat{W}, \frac{1}{\Delta}) \cdot d$  messages. Hence, the total number of messages processed by the server for the object is  $2 \cdot \min(\hat{W}, \frac{1}{\Delta}) \cdot d$ .

### 3.3 Eager versus Lazy Lease Renewals

Another important issue in cooperative leases is the policy for lease renewals. Since the leader manages the lease on behalf of all proxies in the region, i.e., the proxies do not maintain the lease directly with the server, it needs to decide whether and when to renew a lease. Two different renewal policies are possible:

**Eager renewals:** In this policy, the leader continuously renews the lease upon each expiration until it is explicitly notified by proxies not to do so. This approach requires each proxy to track its interests in locally cached objects and send a “terminate lease” message to the leader when it is no longer interested in an object. For instance, a proxy can send such a message if it has not received a request for an the object for a long time period. Upon receiving such a message, the leader removes that proxy from its membership list and stops forwarding server notifications to the proxy. Consequently, a “terminate lease” message is equivalent to a “leave” message from the application-level multicast group. When the membership list becomes empty (i.e., all proxies caching the object send terminate messages), the leader stops renewing the lease. It then broadcasts a directory update to all proxies indicating that it has relinquished leader responsibilities for the object. Eager renewals are beneficial in scenarios where the objects that are being modified are also the most popular objects.

To analyze the overhead of eager renewal, consider a proxy region that caches an object. Assume that  $P$  proxies within the region cache the object. Let  $d$  denote the lease duration for the object and let  $R$  and  $\hat{W}$  denote the read and the non-consecutive write frequency, respectively, for the object at the server. For simplicity, assume that  $\Delta = 0$ , implying that every update triggers a notification. Then the server sends notifications at the rate of  $\hat{W}$  to the leader, which in turn forwards these notifications to the remaining  $P - 1$  proxies. Each proxy issues a HTTP GET message upon a subsequent read. Thus, the total message frequency due to invalidates and GETs is  $2 \cdot \hat{W} \cdot P$ . In addition to these messages, the leader sends an eager renewal message once every  $d$  time units, resulting in a message frequency of  $1/d$  for renewal messages. Finally, so long as the object is popular (large  $R$ ), no terminate messages are sent. When the object becomes cold, each proxy sends a terminate message to the leader, which then terminates the lease when no proxies are interested in the object. Thus, there can be at most  $P$  terminate messages within a lease duration. Consequently,

the per-object control message frequency  $F$  in a proxy region that employs eager renewals is

$$F = \begin{cases} 2 \cdot \hat{W} \cdot P + \frac{1}{d} & \text{for popular objects with } R \gg 0 \\ 2 \cdot \hat{W} \cdot P + \frac{1}{d} + \frac{P}{d} & \text{for unpopular objects; } R \ll 1 \end{cases} \quad (2)$$

**Lazy renewals:** Here, the leader does not renew a lease upon expiration. Instead it sends a “lease expired” message to all proxies caching the object; proxies in turn flag the object as “potentially stale”. This message is required as member proxies do not maintain lease information. Upon receiving a subsequent request for this object, a proxy sends an if-modified-since (IMS) request to the server. The server then issues a new lease for the object, if one has not already been issued, and responds to the IMS request by sending a new version of the object, if the object was modified in the interim. The lease, if one is issued, is sent to the leader. In the lazy approach, proxies do not need to track their interest in each cached object. Moreover, since leases are renewed lazily and only when an object is accessed, the approach is efficient for less popular objects (e.g., “one-timers”). The drawback, though, is that each request received after a lease expiration involves an additional interaction with the server (in the form of an IMS request). In contrast, the eager approach only involves leader-server interactions after lease expiry; individual proxies do not need to interact with the server, which reduces server load.

To analyze the overhead of lazy renewals, consider a region where  $P$  proxies are caching the object. Like in eager renewals, the message frequency due to invalidate and GET messages is  $2 \cdot \hat{W} \cdot P$ . If the object is popular, each of the  $P$  proxies will receive a read request after lease expiry, which triggers an IMS request and an appropriate (server) response. Thus, at most  $\frac{2P}{d}$  IMS requests and responses are exchanged. On the other hand, if the object is unpopular, no messages are sent by proxies upon lease expiry. Consequently, the per-object message frequency  $F$  in a proxy region that employs lazy renewals is

$$F = \begin{cases} 2 \cdot \hat{W} \cdot P + \frac{2P}{d} & \text{for popular objects with } R \gg 0 \\ 2 \cdot \hat{W} \cdot P & \text{for unpopular objects; } R \ll 1 \end{cases} \quad (3)$$

From the above analysis, we observe that eager renewals yield a lower message overhead for popular objects, while lazy renewals yield a lower message overhead for unpopular objects.

### 3.4 Propagating Invalidates versus Updates

When an object is modified, the server notifies each leader proxy with an active lease (subject to the notification rate parameter  $\Delta$ ). As explained earlier, this notification consists of either a cache invalidate or an updated (new) version of the object. On receiving an invalidate message for an object, a proxy marks the object as *invalid*. Subsequent requests for the object requires the proxy to fetch the object from the server (or from another proxy in the region if that proxy has already fetched

the updated object). Thus, each request after a cache invalidate incurs an additional delay due to this remote fetch. No such delay is incurred if the server sends out the new version of the object upon a modification.<sup>3</sup> In such a scenario, subsequent requests can be serviced using locally cached data. A drawback, however, is that sending updates incurs a large network overhead (especially for large objects). This extra effort is wasted if the object is never subsequently requested at the proxy. Consequently, cache invalidates are better suited for less popular objects, while updates can yield better performance for frequently requested objects. Observe that sending invalidates is equivalent to a *lazy update* policy at proxies, while sending new versions of objects amounts to *eager updates*.

A server can dynamically decide between invalidates and updates based on the characteristics of an object. One policy is to send updates for objects whose popularity exceeds a threshold and to send invalidates for all other objects. Although a server does not have access to the actual request stream at proxies to compute object popularities, it can estimate the popularity based on lease renewals. A continuously renewed lease is an indication that the object is popular within a region. Hence, the server can send updates for objects whose leases have been renewed at least  $\tau$  consecutive times ( $\tau$  is a threshold). Using  $\tau = 0$  causes only updates to be sent, whereas  $\tau = \infty$  causes only invalidates to be sent; an intermediate value of  $\tau$  allows the server to dynamically choose between the two based on the object popularity. A more complex policy is to take both popularity and object size into account. Since large objects impose a larger network transfer overhead, the server can use progressively larger thresholds for such objects (the larger a object, the more popular it needs to be before the server starts sending updates).

We analyze the overheads of propagating invalidates and updates to understand when each of these options should be used. Let  $R$  and  $W$  denote the read and write frequency for an object cached in a region. For simplicity, assume that  $\Delta = 0$ , implying that the proxy needs to be notified of all updates.

To analyze the overhead of sending invalidates, observe that a invalidate message needs to be sent for each write that follows a read request. After each invalidate, the leader needs to download a new version of the object upon a subsequent read (via a HTTP GET). No request needs to be sent to the server for consecutive reads (since the cached version is up-to-date). Let all consecutive reads be grouped as a single read and let the resulting frequency of non-consecutive reads be denoted by  $\hat{R}$ . Hence, in the worst case, the bandwidth required for message exchanges is the overhead due to invalidates, the overhead of GET messages and the overhead of transferring the object, which is given by  $2Wc + \hat{R}S$ ,

where  $c$  denotes the size of a control message and  $S$  denotes the size of the object.

As an optimization, we note that each set of consecutive writes requires only a single notification to be sent by the server. Thus, in the best case, only the first write after a read request triggers an

---

<sup>3</sup>Security and authentication issues in doing so are beyond the scope of this paper.

invalidate, provided the server maintains enough state about the leader’s last read time. In this case, the above expression for bandwidth usage simplifies to  $2\hat{W}c + \hat{W}S$ , where  $\hat{W}$  is the frequency of non-consecutive writes. Note that the non-consecutive write frequency and non-consecutive read frequency become equal after consecutive reads and consecutive writes are grouped together (i.e.,  $\hat{R} = \hat{W}$ ).

In contrast, if the server propagates updates, every write at the server results in the object being pushed to the leader. No additional control messages need to be exchanged. This results in a bandwidth requirement of  $WS$ .

Having derived the overheads, we now analyze scenarios under which each approach outperforms the other. Consider the case where the object is small and is approximately equal to the size of a control message ( $c \approx S$ ). Then  $2\hat{W}c + \hat{W}S \approx 3\hat{W}S = 3\hat{R}S$ . Hence, updates are better if  $WS \leq 3\hat{R}S$ , i.e., if the write frequency is less than thrice the non-consecutive read frequency. In this scenario, updates will incur a smaller overhead than invalidates. Intuitively, updates are a better option if the object is small and if lots of reads amortize the overhead of pushing an update (due to a write). In contrast, for large objects (where  $\frac{c}{S} \rightarrow 0$ ),  $2\hat{W}c + \hat{W}S \rightarrow \hat{W}S (= \hat{R}S)$ . So, invalidates use lower bandwidth as long as write frequency is greater than the non-consecutive read (write) frequency, which is always the case. Thus, invalidates are better for large objects and when there are more writes than reads.

### 3.5 Multi-Level Proxy Organization

If the number of proxies in a region is large, then a single level proxy hierarchy will not scale since the leader needs to propagate notifications to a large number of proxies and hence the leader can become a potential bottleneck in the system. To reduce the burden on the leader and improve the scalability of the system, we propose a generic technique for extending cooperative leases to a  $N$ -ary,  $L$ -level proxy hierarchy, where  $N$  is an input to the system and  $L$  is determined by the number of proxies interested in an object in a region. The leader for an object becomes the root of the hierarchy which is  $L$ -levels deep and has a fan-out of  $N$ .  $N$  corresponds to the maximum length of the membership list for any proxy in the region. Several protocols exist that can be used for dynamic tree formation among proxies; however the simple tree construction protocol described below suffices for our purpose. In what follows, we describe how the cooperative leases approach presented in Section 2.4 can be extended to multi-level hierarchies.

*First Request:* The operation of cooperative leases for this event remains the same as that described in Section 2.4. Although any leader selection scheme suffices for our purpose, for simplicity, we assume the *first proxy is leader* scheme for the rest of this section. To make cooperative leases generically applicable to a multi-level hierarchical, the leader needs to maintain additional state about the structure of the tree and each member proxy needs to maintain information about

its children. Doing so enables generation of a well-balanced tree. We note that, like in single level hierarchies, a multi-level hierarchy is generated on a per-object basis; different objects may have different leaders and different tree structures.

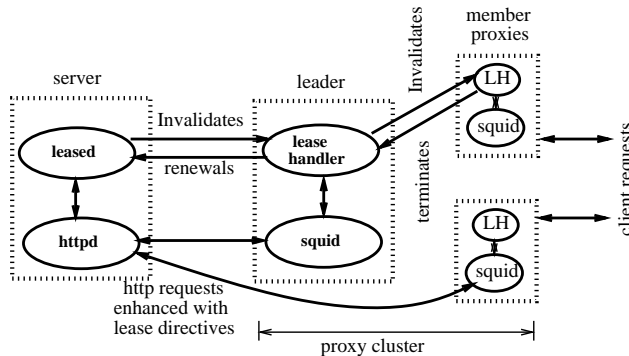
*Subsequent Requests:* This is similar to what is described in Section 2.4 with the following exception. If a proxy experiences a cache miss, it forwards this request to the leader; the leader adds such a proxy as a member (an immediate child) as long as it has less than  $N$  immediate children. For the formation of further levels, the proxy is added as a child of an existing member with the least number of children; since the leader knows the structure of the entire tree, determining this information is trivial. If all members have an equal number of children, then a member is picked at random. If the selected member already has  $N$  immediate children it selects one of its children, to be the parent in a similar fashion. This process continues recursively until a parent is found. The number of children at each level of the tree gets updated as the parent selection process continues. We call the messages sent by the leader for identifying a parent, *AddMember* messages. Once a parent proxy is found, the parent adds the new proxy as its child (using a *JoinMe* message) and hence the new proxy joins the hierarchy. The hierarchy creation procedure of selecting a member with the smallest cardinality ensures that the tree created is balanced and has the smallest depth for a given number of proxies, thereby, reducing the propagation delay from leader to leaf-proxies.

*Updates to the Object:* Since the basis for multi-level proxy organization is sharing the task for propagating updates, in addition to that stated in Section 2.4, every proxy in the hierarchy maintains a membership list. When a proxy in the hierarchy receives an update, it forwards this only to proxies in its membership list. Hence a leader proxy multicasts updates to its children and the process continues recursively until all leaf-proxies in the hierarchy receive the updates.

*Terminate Requests:* When a proxy loses interest in an object and if it has no children, it sends a terminate request to its parent and this recurses all the way to the root (leader). This is necessary so that state information at all proxies required for maintaining balance in the tree gets updated. To differentiate these messages from *Terminate Requests* we see in one-level proxy hierarchies, we call these *UpdateChildrenState* messages. As new proxies join the hierarchy the tree will dynamically re-balance itself.

### 3.6 Multiple regions in a CDN

If proxies in a CDN are organized into multiple regions, the server has to maintain additional state information, which includes leader proxy and lease information on a per-region basis. As a result, if proxies in multiple regions express interest in an object, the server will be responsible for propagating updates to multiple leaders (one per region), and they in turn will be responsible for catering to member proxies in their regions. While it is clear that such a scheme requires maintenance of extra state, the maintenance of a lease for an object per-region has its advantages. This, in general,



**Figure 2:** Implementation architecture. The figure depicts the architecture of a single region and a single-level region. Each CDN will have a number of such regions.

adds flexibility in independently managing critical resources for each region. For instance, CDN administrators may choose different lease duration computation techniques [7] for leases for different regions based on specific region characteristics. Note that the operations of cooperative leases within each individual region remains as stated in Section 2.4 and Section 3.5.

#### 4 Implementation Issues

We have implemented the cooperative leases algorithm in the Squid proxy cache and the Apache web server<sup>4</sup>. Our implementation is based on HTTP/1.1, which allows user defined extensions as part of the request/response header. We use these header extensions to enable proxies to request and renew leases from a server. To do so, lease requests and responses are piggybacked onto normal HTTP requests and responses. Lease renewals and invalidation requests are also sent as request header extensions. The exact HTTP grammar for lease requests, renewals and invalidations is described in [20].

For simplicity and modularity, our implementation separates lease management functionality from the serving of web requests. Lease management at the server is handled by a separate lease server (`leased`). Such an architecture results in a clean separation of functionality between the Apache server, which handles normal HTTP processing, and the lease server which handles lease processing and maintains all the state information (see Figure 2). Whenever the Apache server receives a lease grant/renewal request piggybacked on a HTTP request, it forwards the former to the lease server for further processing. The lease duration  $d$  and the notification rate parameter  $\Delta$  are computed using policies listed in [7] and Section 3.2. The HTTP response is then sent back to the client (proxy), while the lease is sent to the leader. Invalidation requests are handled similarly — the web server forwards the request to the lease server, which then sends invalidations to all leaders with active leases. Leaders forward the invalidations to all proxies caching the object as described below.

<sup>4</sup>Source code for our implementation can be downloaded from <http://lass.cs.umass.edu/software/cdn>

Analogous to the web server architecture, our implementation in Squid consists of two components — the proxy cache and the lease handler — that separate the caching functionality from lease management. The lease handler (*LH*) can either act as a leader or as a member. In the former case, the lease handler maintains a membership list of all proxies caching the object and forwards notifications from the server to this list. The lease handlers at member proxies are responsible for tracking object popularities and sending lease terminate messages to the leader for cold objects. Server failures and/or network partitions can be handled at the leader by exchanging heartbeat messages [18] or by maintaining a persistent TCP connection with the server — a broken connection indicates a failure and requires cached objects to be invalidated within  $\Delta$  time units. The heartbeat interval should, in this case, be smaller than the notification rate parameter.

## 5 Experimental Evaluation

In this section, we demonstrate the efficacy of cooperative leases by (i) comparing the approach with the original leases from the perspective of scalability, (ii) evaluating the tradeoffs of various policies described in section 3 and (iii) quantifying the implementation overheads of cooperative leases. We employ a combination of trace-driven simulation and prototype evaluation for our experiments. We use simulations to explore the parameter space along various dimensions and use our prototype to measure implementation overheads (an aspect that simulations do not reveal). In what follows, we first present our experimental methodology and then our experimental results.

### 5.1 Experimental Methodology

#### 5.1.1 Simulation Environment

We have designed an event-based simulator to evaluate the efficacy of cooperative leases. The simulator simulates one or more proxy regions within a CDN. Each proxy is assumed to receive requests from a large number of clients. Cache hits are serviced using locally cached data. Cache misses involve a remote fetch and are serviced by fetching the object from the leader (if one exists) or from the server. The directory maintained by the proxy is used to make this decision. Our simulator supports all policies discussed in Section 3 for leader selection, server notifications, lease renewals and rate computations.

Our experiments assume that each proxy maintains a disk-based cache to store objects. We assume each proxy cache is infinitely large — a practical assumption, since disk capacities today are in tens of gigabytes and a typical proxy can employ multiple disks. Data retrievals from disk (i.e., cache hits) are modeled using an empirically derived disk model with a fixed operating system overhead added to each request. For cache misses, data retrieval over the network are modeled using the round trip time, available network bandwidth and the object size. The network latency and bandwidth



Trace	Num Requests	Duration (secs)	Unique Objects	Num Writes
DEC	750000	42031	276914	17126
NLANR	750000	56518	393853	14385

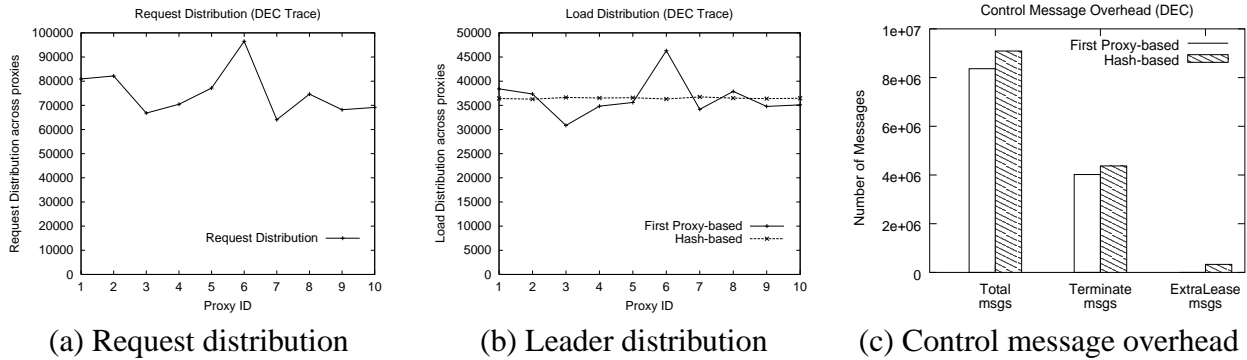
**Table 2:** Trace Characteristics. The DEC traces are from Dec. 1996 while the NLANR traces are from June 2001. The number of writes is the total number across all objects. These were synthetically introduced with the write frequency following a tri-modal distribution based on earlier studies.

between proxies and leaders is assumed to be  $75ms$  and  $500KB/s$ , while that between proxies and origin servers is  $250ms$  and  $250KB/s$ . Although actual network latencies and bandwidths vary with network conditions, the use of this simple network model suffices for our purpose (due to our focus on consistency maintenance rather than end-user performance).

Unless noted otherwise, our experiments assume a default of 1 region with a one-level proxy hierarchy, a region size of 10 proxies, and a lease duration of 30 minutes. We simulate  $750K$  trace requests. The details of these traces are outlined in the section that follows. We also assume that a leader always caches a copy of the object and this copy is updated upon a modification.

### 5.1.2 Workload Characteristics

The workload for our experiments is generated using traces from actual proxies, each containing several hundred thousand requests. We use two different traces for our study; the characteristics of these traces are shown in Table 2. The same set of traces are used for our simulations as well as our prototype evaluation (which employs trace replay). Each request in the trace provides information such as the time of the request, the requested URL, the size of the object, the client ID, etc. We use the client ID to map each request in the trace to a proxy in the region — all requests from a client are mapped to the same proxy. To determine when objects were modified, we considered using the last modified times as reported in the trace. However, these values were not always available. Since the modification times are crucial for evaluating cache consistency mechanisms, we employ an empirically derived model to generate modification times. Based on observations in [2, 14], we assume that 90% of all web objects change very infrequently (i.e., have an average lifetime of 60 days). We assume that 7% of all objects are mutable (i.e., have an average lifetime of 20 days) and the remaining 3% objects are very mutable (i.e., have a lifetime of 5 days). We partition all objects in the trace into these three categories and generate write requests and last modified times using exponentially distributed lifetimes. Although the average lifetimes are in days, given the high variance in the modification times there were numerous writes within the sampling duration of the trace. The number of synthetic writes generated for each trace is shown in Table 2. In practice the



**Figure 3:** Comparison of leader selection schemes

server will rely on a publishing system or a database trigger to detect a modification, the details of which are beyond the scope of the paper.

Next, we describe our experimental results.

## 5.2 Impact of Leader Selection Policies

To evaluate leader selection policies, we simulated a region of 10 proxies that employed two different policies — the hash based policy and the “first proxy is leader” policy. Our experiment assumed eager lease renewals and notifications in the form of invalidations (leaders were sent updates, leaders forwarded invalidations). For each policy, we measured how evenly leader responsibilities were distributed across proxies in the region as well as the total control message overhead imposed. Figures 3(b) and (c) depict our results, while Figure 3(a) shows the number of requests processed by each proxy in the region (we only plot results for one of the traces due to space constraints. See [20] for complete results). As expected, the “first proxy is leader” scheme suffers from load imbalances since some proxies service a larger number of requests (and assume leader responsibilities for a correspondingly larger number of first-time requests). The figure also shows that there is a factor of 1.5 difference in load between the most heavily-loaded and the least-loaded proxy. In contrast, the hash-based policy shows better load balancing properties but imposes a larger communication overhead (since leaders can be different from proxies caching the object, requiring additional message exchanges). As shown in Figure 3(c), the total increase in control message overhead is about 10% and the increase is primarily due to the lease terminate messages sent from proxies to leaders. Since a small (10%) increase in message overhead is tolerable to correct a potentially large imbalance (factor of 1.5), our results indicate that the hash-based leader selection is a better policy than the “first proxy is leader” approach.

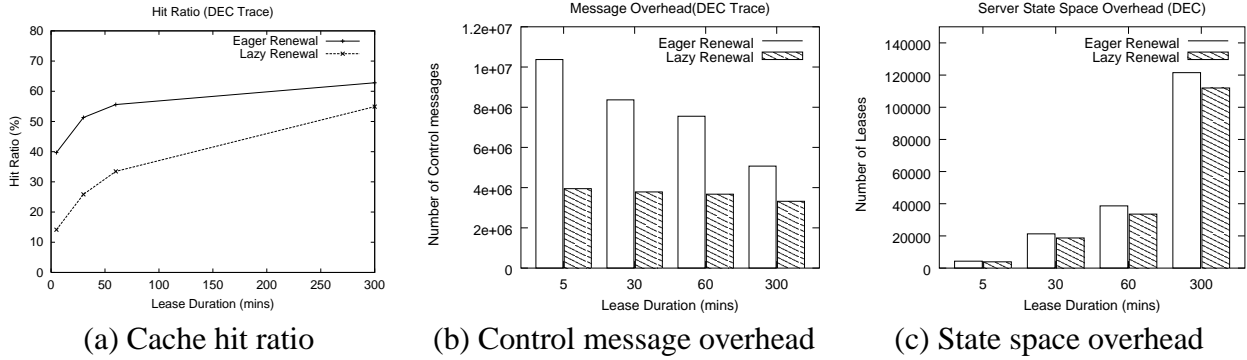


Figure 4: Comparison of lease renewal policies

### 5.3 Eager versus Lazy Renewals

Next, we evaluate the impact of eager and lazy lease renewals on performance. Like in the previous experiment, we assume a region of 10 proxies, each with an infinite cache. We vary the lease duration from 5 minutes to 5 hours and measure its impact on lazy and eager renewals. Figure 4 depicts our results. As shown in Figure 4(a), depending on the lease duration, eager renewals result in a 15–63% improvement in cache hit ratios; the hit ratio is lower for lazy renewals since requests arriving after a lease expiry trigger an IMS request to the server. The higher hit ratios for eager renewals are at the expense of an increased control message overhead (see Figure 4(b)). The message overhead is 33–175% higher and is primarily due to extra lease renew and terminate messages. The overhead for both policies decreases with increasing lease durations (since longer leases require fewer renewals). Finally, Figure 4(c) plots the state space overhead of the two policies; as expected, eager renewals result in a larger number of active leases at any instant, causing a 3–9% increase in state space overhead.

An important factor governing the performance of the eager renewals is the lease termination policy — the policy employed by member proxies to notify the leader that they are no longer interested in the object. As shown in Figure 5, the larger the period of inactivity before which a “terminate lease” message is sent, the larger the state space overhead at the server and the larger the control message overhead (since the leader continuously renews leases until such a message is received).

Thus, the two policies show a clear tradeoff — eager renewals yield better hit ratios and response times at the expense of a larger control message overhead and a slightly larger state space overhead. Depending on whether user performance or network/server overheads are the primary factors of interest, one policy can be chosen over the other.

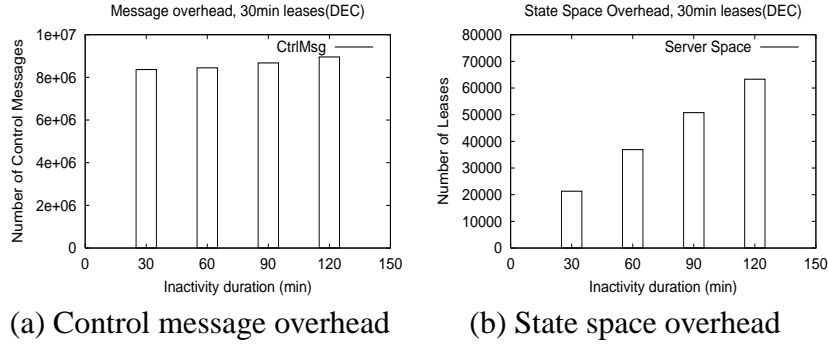


Figure 5: Evaluation of lease termination policies

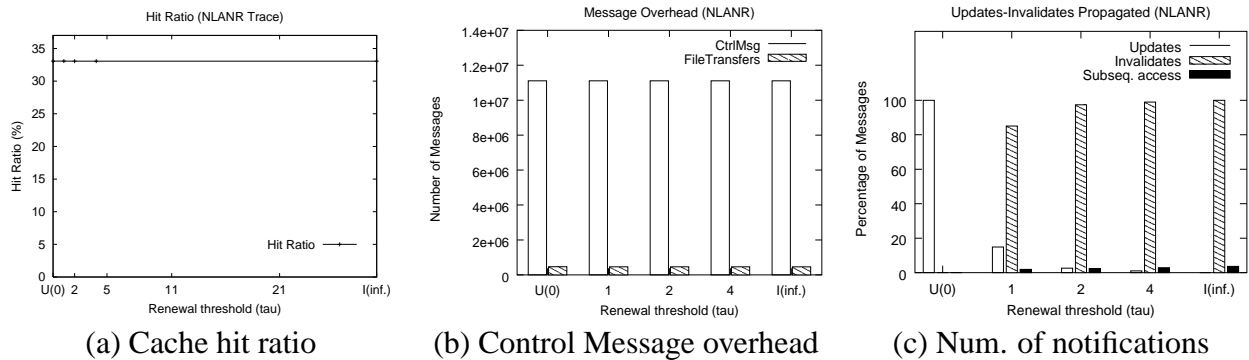
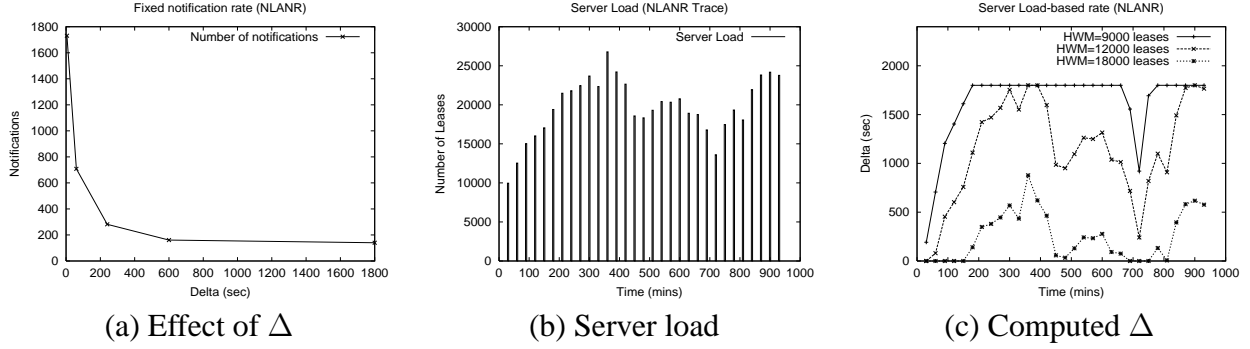


Figure 6: Comparison of server notification policies

#### 5.4 Server Notifications: Invalidates versus Updates

To understand the implications of sending invalidates versus updates, we considered a policy where the server sent updates for objects whose leases were renewed at least  $\tau$  times in succession; invalidates were sent for the remaining objects. We varied  $\tau$  from 0 to  $\infty$  and measured its impact on the cache hit ratio and the control message overhead. Figure 6 shows that the notification policy has a negligible impact on the cache hit ratio ( $< 1\%$  reduction as  $\tau$  increases from 0 to  $\infty$ ). The control message overhead increases slightly (by about 1%) with increasing  $\tau$ . This small increase is due to an increase in the number of invalidates, each of which triggers an HTTP request upon a subsequent user request. To better understand this behavior, Figure 6(c) plots the percentage of updates and invalidates sent for different values of  $\tau$ ; the percentage of objects accessed subsequent to a server notification is also shown. As shown, when  $\tau = \infty$  (i.e., the invalidate-only scenario), only 5% of the invalidated objects are accessed subsequently. Thus, our results show that updates should be sent only for those modified objects that are also popular, which can be achieved using a large  $\tau$ . More generally, our analysis of read and write frequencies has shown that updates are advantageous when the write frequency is (i) less than 3 times the read frequency for small objects and (ii) less than the



**Figure 7:** Impact of the notification rate

read frequency for large objects [20].

### 5.5 Impact of the Notification Rate

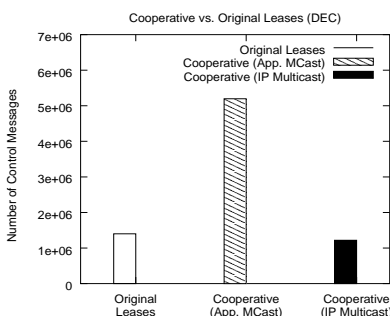
To understand the impact of the notification rate, we varied  $\Delta$  from 5 seconds to 30 minutes and measured the impact on the number of notifications (invalidates) sent by the server (the leases duration was fixed at 30 minutes). As shown in Figure 7(a), the number of notifications drops by an order of magnitude with increasing  $\Delta$ s. This indicates that an appropriate choice of  $\Delta$  can result in substantial savings at the server, albeit at the expense of weaker consistency guarantees. Next we considered a policy where the server computes  $\Delta$  based on the load as explained in Equation 1; the server state space overhead, measured by the number of concurrent leases, is used as an indicator of the load. Note that  $\Delta$  is computed based on the server load only at the beginning of a lease; once picked,  $\Delta$  does *not* change for that lease until lease expiry. We varied the high and low watermarks in Eq. 1 and measured its impact on  $\Delta$ . Figure 7(b) shows the variation in the server load over a 15-hour period, while Figure 7(c) plots the corresponding value of  $\Delta$  used for new leases and renewals. The figure shows that the value of  $\Delta$  closely matches the variation in the server load. Further, depending on the low and high watermarks used, the server uses  $\Delta = 0$  during periods of low load and increases  $\Delta$  to its maximum value (i.e., the lease duration) during periods of heavy load. Thus, an intelligent choice of  $\Delta$  helps provide the desired level of consistency guarantee while lowering server overheads.

### 5.6 Scalability Issues: Comparison with Original Leases

We use a more write-intensive workload to study scalability issues — in addition to the set up we described in Section 5.1.2. We split the 3% of very mutable objects into two categories — 2.5% of these change once every 480 minutes (*Type A*), and 0.5% of these change once every 32 minutes (*Type B*).

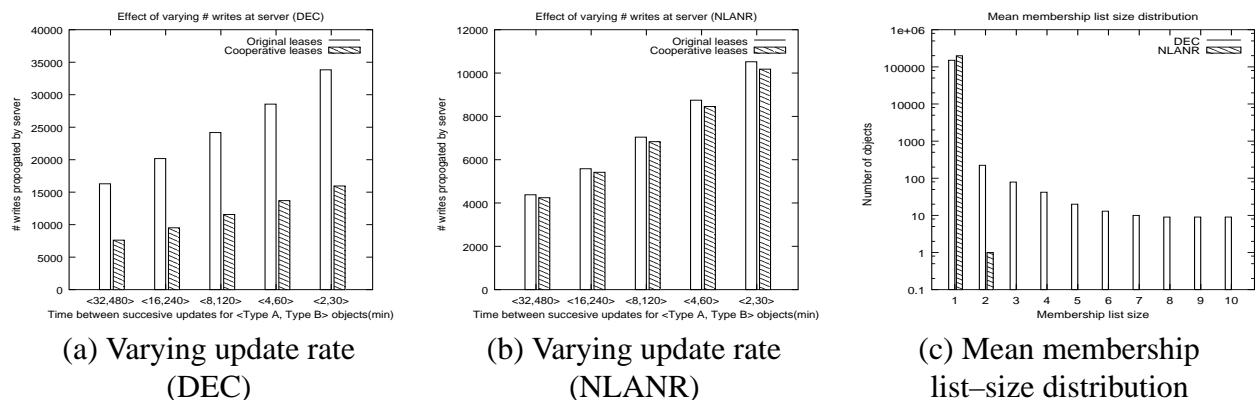
**Table 3:** Comparison with original leases (DEC)

	20 proxies		10 proxies	
	Coop. Leases	Leases	Coop. Leases	Leases
Active leases	23038	28653	23013	27477
Svr. notifications	5085	12618	5082	9813

**Figure 8:** Control message overhead for cooperative and original leases

### 5.6.1 State space and Control message overheads

To compare cooperative leases with original leases, we consider a region of 20 proxies. We also mention results for a 10-proxy region for comparison with results in prior sections. To permit a fair comparison, other than the cache consistency mechanism, all simulation parameters are kept identical across these experiments, the first involving cooperative leases and the second employing the original leases approach. The lease duration is set to 30 minutes and  $\Delta = 0$ . Due to resource constraints, we simulate only 500K read requests from the DEC trace (this represents a duration of 23565 seconds). Figure 8 and Table 3 depict our results. As expected, the number of leases managed by the server decreases when cooperative leases is used (since each lease represents multiple proxies, fewer leases are needed). The reduction in state space overhead is 20% (see Table 3); the reduction is smaller than expected since a large number of objects in the workload are requested by only one proxy and cooperative leases do not provide any benefits in such scenarios. Note however that the number of active leases in the region at any instant is only in the order of a few thousands. The number of server notifications is smaller by a factor of 2.5 indicating that cooperative leases successfully offloads the burden of sending notifications to leader proxies, thereby improving server scalability. These reductions come at the expense of having to maintain a directory of cached objects and an increased control message overhead due to directory updates. This results in an increase in the message overhead by a factor of 3.7, for a 20-proxy region — the directory update overhead is proportional to the number of proxies in the region when application-level multicast (i.e., unicast) is



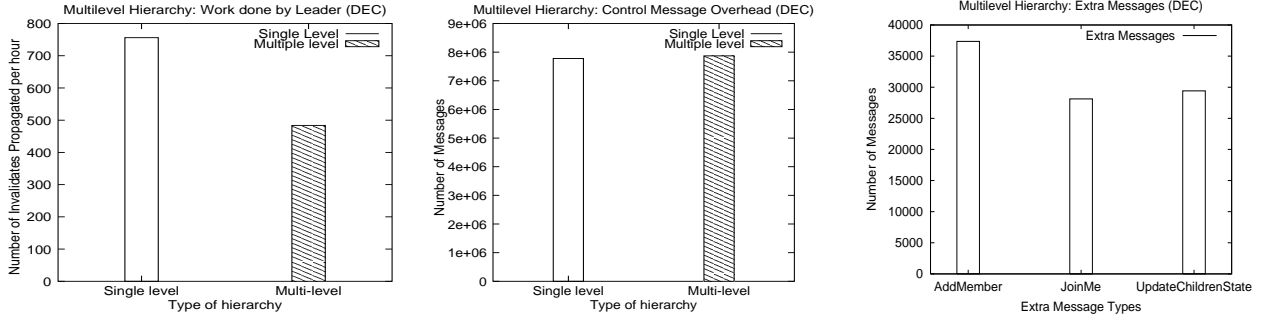
**Figure 9:** Effect of the update rate at the server

used (see Figure 8). The use of IP-multicast, instead of application-level multicast, to send directory updates can help lower this overhead (since IP-multicast is more efficient than unicast). Also note that each unique  $\Delta$  value associated with an object needs its own application level multicast group; a server can reduce the number of multicast groups by restricting itself to a small set of  $\Delta$ s. For a 10-proxy region, the reduction in state space overhead is 16%, the number of server notifications is smaller by a factor of 1.9 and the increase in the control message overhead is by a factor of only 2.2. Thus, we conclude that cooperative leases do indeed enhance scalability from the perspective of the server (in terms of the state space and server message overhead), albeit at the expense of increased inter-proxy communication overhead.

### 5.6.2 Effect of varying update rate at the server

A second dimension of comparing the cooperative leases and original leases mechanisms is by studying the effect of write frequency of objects at the server. In this section, we change *Type A* objects once every 30 to 480 minutes, and *Type B* objects once every 2 to 32 minutes. Experiments were run on both DEC and NLANR traces using 10 proxies for original leases and a single region of 10 proxies for cooperative leases. All other parameters are as described in Section 5.1.1.

Figures 9(a), 9(b) and 9(c) summarize our results. In Figure 9(a), cooperative leases consistently reduced the number of updates the server propagated to proxies by 50–53% for the DEC trace. However, Figure 9(b) shows that the corresponding gain is only 2.9–3.3% for the NLANR trace. We attribute this to majority of the objects in the NLANR trace not being accessed by multiple proxies. Figure 9(c) plots the distribution of membership list sizes at leader proxies for both DEC and NLANR workloads. As seen from the figure most leases have only one proxy in the membership list for the NLANR workload, whereas a sizable number of objects (popular objects) have greater



(a) Invalidate Propagation rate at Leaders      (b) Control Message Overhead      (c) The Extra Messages

**Figure 10:** Comparing single and multiple level hierarchical proxy organizations in a region

than one proxy in their membership lists for the DEC workload. In scenarios where objects are accessed by only one proxy, cooperative leases do not provide any benefits over normal leases.

We conclude that as long as objects are popular and accessed by clients associated with different proxies in a region, cooperative leases are effective in propagating server notifications.

## 5.7 Multi-level Hierarchical Proxy Organization

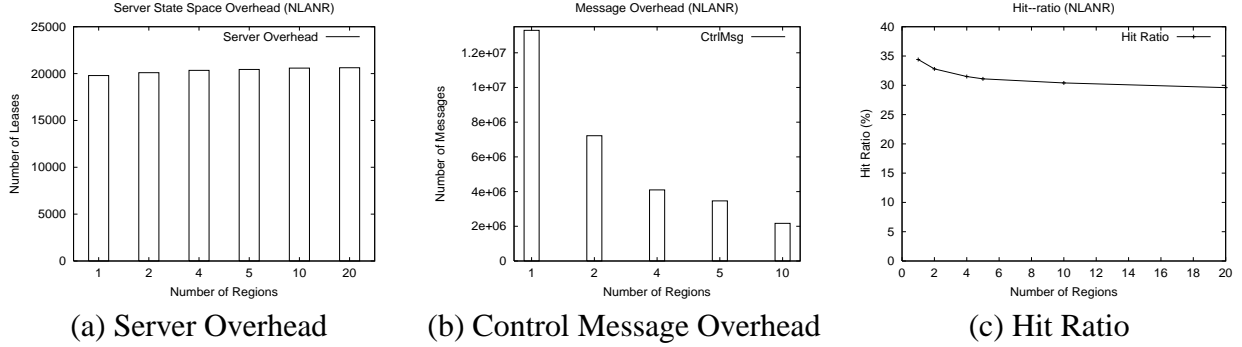
In this section, we study how organizing a region of CDN proxies into a multi-level hierarchy can help reduce some of the overheads of maintaining a one-level proxy hierarchy. As we mentioned earlier, the amount of work a leader proxy will have to do in order to propagate updates and/or invalidates to all interested members in the group can become prohibitive if the number of proxies in the region becomes large.

We ran this experiment with  $500K$  read requests of the DEC trace (a duration of 23565 seconds). While we keep all other parameters the same, as elaborated in Section 5.1.1, we use a 15-proxy region for a one-level proxy hierarchy (server to leader proxy to member proxies) and for a hierarchy that has a fixed span-out of 2 (four levels at most). Note that the one-level proxy hierarchy can have a maximum span-out of 14 (from leader to members) as opposed to the maximum span-out of 2 in the other case.

Figures 10(a), (b) and (c) summarize our results. Figure 10(a) shows a decrease by about 36%, in the number of invalidates propagated by leader proxies per hour, when proxies get organized into hierarchies with a span-out of 2. Figure 10(b) shows an increase in the number of control messages by about 1.2%. This is attributed to the additional messages (see Figure 10(c) and Section 3.5) involved in maintaining a proxy hierarchy of more than one level.

We conclude that cooperative leases scales well with the organization of proxies in a region into multiple levels, by sharing the responsibility of propagation of invalidate messages by the leader





**Figure 11:** Various metrics with increasing number of regions

proxies with the member proxies, at the cost of a very small increase in the control message overhead.

### 5.8 Organization of CDN proxies into varying number of regions

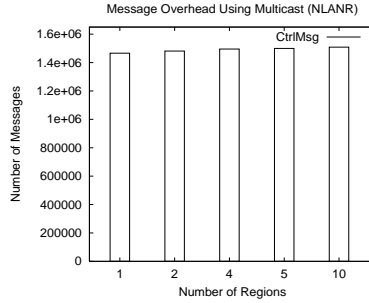
In this section we study how cooperative leases scale when a fixed number of proxies  $P$  is organized into varying number of regions,  $R$  (where  $R < P$ ;  $R = P$  represents the original leases [12] case).

This experiment was run using  $P = 20$  proxies and 490K requests of the NLANR trace (a duration of 36623 seconds). All other same parameters are as mentioned in Section 5.1.1. We ran experiments by uniformly distributing the 20 proxies across  $R = 1, 2, 4, 5$  and 10 regions. Let us now compare the results we got for single ( $R = 1$ ) versus multiple ( $R > 1$  and  $R < P$ ) regions.

As we increase the number of regions, since there can be only one leader per region, the server has to maintain more state than it does in the single region case (a list of leader proxies, one per region and a list of leases, one per leader, per object). As we go from  $R = 1$  to  $R = 10$ , we see a 20% increase in the mean state space overhead at the server (see Figure 11(a)). However, along with the increase in the number of regions, since the number of proxies per region ( $P/R$ ) decreases, the degree of cooperation each region decreases. As a result, we observe a decrease in message overhead by 80.7% (see Figure 11(b)) and a decrease in overall hit ratio by a factor of 1.34 (see Figure 11(c)).

Therefore, if server resources are critical, maintaining fewer number of regions will be more economical; and if network resources are critical, creating multiple regions and distributing CDN proxies across these is useful, albeit a decrease in overall hit ratio and increase in state space consumption at the server.

In addition to the above, we ran the same experiment but simulated IP multicast this time. From Figure 8 we observed that for a 20-proxy region, the large increase in message overhead when we use cooperative leases instead of the original leases, can be overcome by using IP Multicast. Figure



**Figure 12:** Control Message Overhead when using IP Multicast with varying number of regions

12 on the other hand shows how this message overhead varies with increasing number of regions, when IP multicast is used. We see only a 2.8% increase in the message overhead as we vary  $R$  from 1 to 10. This is because the increase in the number of regions decreases the benefits of IP multicast.

Hence, cooperative leases exhibits good scalability properties as we distribute CDN proxies across more regions when using either application-level or IP multicast.

### 5.9 Implementation Overheads

Whereas the preceding sections examined the efficacy of cooperative leases using simulations, in this section we study the overheads of various operations needed for consistency maintenance. The testbed for our experiments consists of the lease-enhanced Apache web server, a region consisting of four Squid proxy caches and a client workload generator, all of which run on a cluster of Linux PCs. Each PC in our experiment is a  $700MHz$  Pentium III with  $512MB$  RAM, interconnected by  $100Mb/s$  switched ethernet. The client workload generator employs trace replay and uses the traces described in Table 2. To do so, it maps each URL in the trace to a unique object stored on the server of approximately the same size. Further, like in our simulations, each end-host in the trace is bound to a fixed proxy cache using a hashing function. The proxy and the server maintain consistency using cooperative leases as described in Section 4. We measured the overhead of various lease management operations at the server and the proxies over the duration of the trace. Table 4 lists our results. As shown in the table, the overhead of granting and renewing leases is very small (order of milliseconds). Similarly directory updates and server notifications (invalidates) can be propagated efficiently to proxies in the region (clearly these overheads depend on the number of proxies in the region and number of proxies that cache an object, respectively).

These results indicate that cooperative leases can be implemented efficiently in web servers and CDN proxies.

**Table 4:** Implementation Overheads (NLANR)

Proxy Overheads		Server Overheads	
Event	Time (ms)	Event	Time (ms)
Dir. lookup	0.052	Grant lease	0.64
Dir. update	0.056	Renew lease	0.28
Dir. broadcast	2.7	Inv. to Leader	3.36
Renew lease	2.65		
Inv. to Proxy	0.565		

## 6 Related Work

Recently several cache consistency mechanisms have been developed for single proxies [4, 5, 7]; as argued earlier, these mechanisms do not scale well to proxies in a CDN. Three recent efforts have focused on the issue of scalability [18, 27, 30]. We discuss each in turn.

A cache consistency mechanism for hierarchical proxy caches was discussed in [30]. The approach does not propose a new consistency mechanism, rather it examines issues in instantiating existing approaches into a hierarchical proxy cache using mechanisms such as multicast. They argue for a fixed hierarchy (i.e., a fixed parent-child relationship between proxies), whereas we allow different proxies to be leaders for different objects. In addition to consistency, they also consider pushing of content from servers to proxies.

Mechanisms for scaling leases are studied in [27]. The approach assumes volume leases [28], where each lease represents *multiple objects* cached by a stand-alone proxy. In contrast, we employ cooperative leases where a lease can represent *multiple proxies*. They examine issues such as delaying invalidations until lease renewals; our work is based on a different consistency mechanism —  $\Delta$ -consistency — for propagating invalidations.  $\Delta$ -consistency allows a separation of the notification frequency from the lease duration, providing additional flexibility to the server. They also discuss prefetching and pushing of lease renewals. Our renewal policies are more complex, since leaders need to interact with member proxies to decide on renewals. We should note that if a large number of objects are serviced by only one proxy in a region and if several such objects originate from the same server, we could further optimize state and message overheads by employing a single volume lease to manage these objects.

Techniques for dynamically growing and shrinking consistency hierarchies are presented in [29]. Issues such as fault tolerance and performance of hierarchies are studied in this work. The study suggests that a promising configuration for providing strong consistency is a two-level hierarchy and dynamic hierarchies are almost always better than static hierarchies. In contrast to their focus on dynamic hierarchies and fault tolerance, we focus on issues such as leader selection and eager versus lazy lease renewals.

The web cache invalidation protocol (WCIP) is an attempt to standardize propagation of server invalidations using application-level multicast [18]. The focus is on a protocol for propagating invalidations; the approach is agnostic of the actual cache consistency mechanism employed by proxies. Like WCIP, our approach also employs leaders to propagate invalidations and manage lease renewals on behalf of proxies in a region. While we study specific cache consistency mechanisms and policies as well as their performance, their focus is on protocol issues (message formats, heartbeat messages etc.). Indeed, our prototype implementation could have employed WCIP instead of HTTP for sending invalidations.

The distributed object consistency protocol (DOCP) [6] proposes extensions to the current HTTP cache control mechanism for providing consistency guarantees. DOCP uses a publish and subscribe mechanism along with server invalidations to provide consistency guarantees. In DOCP, *master* proxies that publish content are discovered by *slave* proxies for subscription using an optimistic discovery mechanism. Proxies subscribe to master proxies for popular objects and directly interact with server for other objects. After subscription, the use of master and slave proxies is similar to our approach of member proxies joining a group after receiving first-time requests for objects. In both cases subsequent invalidates are always sent via the group’s leader/master proxy. In contrast to the DOCP approach of only using invalidates, we study both the effect of propagating updates and invalidates depending on object characteristics. While DOCP requires deployment of certain proxies as *masters*, our approach allows flexible leader selection.

Finally, numerous studies have focused on specific aspects of cache consistency or content distribution. For instance, piggybacking of invalidations [17], the use of deltas for sending updates [19], an application-level multicast framework for internet distribution [10], the efficacy of sending updates versus invalidates [9], distribution strategies [21], and various schemes for prefetching content in CDNs [26] have also been studied. These efforts complement our work and can coexist with our approach.

## 7 Concluding Remarks

In this paper, we argued that existing consistency techniques are not suitable for CDN environments. To alleviate this drawback, we proposed the notion of cooperative consistency and a mechanism called cooperative leases to achieve it. Cooperative leases meets the twin goals of flexibility and scalability by (i) employing  $\Delta$ -consistency semantics, (ii) using a single lease to represent multiple proxies and (iii) using application-level multicast to propagate server notifications. We implemented our approach into a prototype web server and proxy cache and demonstrated its efficacy via an experimental evaluation. While we examined a single region with a one-level proxy hierarchy for most design considerations, we went on to underline that cooperative leases indeed scale well when

proxies in a region are organized into a multi-level hierarchy and when CDN proxies are organized into multiple regions. Although some results shown use the DEC trace and others use the NLANR trace, all corresponding results using both the traces have displayed identical trends and are reported in [20].

## References

- [1] Akamai Technologies, Inc. <http://www.akamai.com>.
- [2] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker. Web Caching and Zipf-like Distributions: Evidence and Implications. In *Proceedings of Infocom'99, New York, NY*, March 1999.
- [3] M. Busari and C. Williamson. On the Sensitivity of Web Proxy Cache Performance to Workload Characteristics. In *Proceedings of IEEE Infocom'01, Anchorage, Alaska*, April 2001.
- [4] P. Cao and C. Liu. Maintaining Strong Cache Consistency in the World-Wide Web. In *Proceedings of the Seventeenth International Conference on Distributed Computing Systems*, May 1997.
- [5] V. Cate. Alex: A Global File System. In *Proceedings of the 1992 USENIX File System Workshop*, pages 1–12, May 1992.
- [6] John Dille, Martin Arlitt, Stephane Perret, and Tai Jin. The Distributed Object Consistency Protocol. Technical report, Hewlett-Packard Labs Technical Reports, 1999.
- [7] V. Duvvuri, P. Shenoy, and R. Tewari. Adaptive Leases: A Strong Consistency Mechanism for the World Wide Web. In *Proceedings of the IEEE Infocom'00, Tel Aviv, Israel*, March 2000.
- [8] L. Fan, P. Cao, J. Almeida, and A Z. Broder. Summary Cache: A Scalable Wide-area Web Cache Sharing Protocol. In *Proceedings ACM SIGCOMM'98, Vancouver, BC*, pages 254 – 265, September 1998.
- [9] Z. Fei. A Novel Approach to Managing Consistency in Content Distribution Networks. In *Proceedings of the 6th Workshop on Web Caching and Content Distribution, Boston, MA*, June 2001.
- [10] P. Francis. Yoid: Extending the Internet Multicast Architecture. Technical report, AT&T Center for Internet Research at ICSI (ACIRI), April 2000.
- [11] S. Gadde, J. Chase, and M Rabinovich. Web Caching and Content Distribution: A View From the Interior. In *Proceedings of the 5th International Web Caching and Content Delivery Workshop*, 2000.
- [12] C. Gray and D. Cheriton. Leases: An Efficient Fault-Tolerant Mechanism for Distributed File Cache Consistency. In *Proceedings of the Twelfth ACM Symposium on Operating Systems Principles*, pages 202–210, 1989.
- [13] M. Gritter and D R. Cheriton. An Architecture for Content Routing Support in the Internet. In *Proceedings of the USENIX Symposium on Internet Technologies, San Francisco, CA*, March 2001.

- [14] J. Gwertzman and M. Seltzer. World-Wide Web Cache Consistency. In *Proceedings of the 1996 USENIX Technical Conference*, January 1996.
- [15] D. Karger, E. Lehman, T. Leighton, M. Levine, D. Lewin, and R. Panigrahy. Consistent Hashing and Random Trees: Distributed Caching Protocols for Relieving Hot Spots on the World Wide Web. In *Proceedings of the Twenty-ninth ACM Symposium on Theory of Computing*, 1997.
- [16] D. Katabi and J. Wroclawski. A Framework for Scalable Global IP-Anycast. In *Proceedings of ACM SIGCOMM, Stockholm, Sweden*, pages 3–15, August 2000.
- [17] B. Krishnamurthy and C. Wills. Study of Piggyback Cache Validation for Proxy Caches in the WWW. In *Proceedings of the 1997 USENIX Symposium on Internet Technology and Systems, Monterey, CA*, pages 1–12, December 1997.
- [18] D. Li, P. Cao, and M. Dahlin. WCIP: Web Cache Invalidation Protocol. IETF Internet Draft, November 2000.
- [19] J C. Mogul, F. Douglis, A. Feldmann, and B. Krishnamurthy. Potential Benefits of Delta Encoding and Data Compression for HTTP. In *Proceedings of ACM SIGCOMM Conference*, 1997.
- [20] A. Ninan. Maintaining Cache Consistency in Content Distribution Networks. Master’s thesis, Department of Computer Science, Univ. of Massachusetts, June 2001.
- [21] G. Pierre, M. van Steen, and A. Tanenbaum. Dynamically Selecting Optimal Distribution Strategies for Web Documents. *IEEE Transactions on Computers*, 51(6), pages 637–651, June 2002.
- [22] A. Shaikh, R. Tewari, and M. Agrawal. On the Effectiveness of DNS-based Server Selection. In *Proceedings of IEEE Infocom, Anchorage, Alaska*, April 2001.
- [23] R. Srinivasan, C. Liang, and K. Ramamritham. Maintaining Temporal Coherency of Virtual Warehouses. In *Proceedings of the 19th IEEE Real-Time Systems Symposium (RTSS98), Madrid, Spain*, December 1998.
- [24] R. Tewari, M. Dahlin, H M. Vin, and J. Kay. Beyond Hierarchies: Design Considerations for Distributed Caching on the Internet. In *Proceedings of the 19th International Conference on Distributed Computing Systems (ICDCS)*, June 1999.
- [25] B. Urgaonkar, A. Ninan, M. Raunak, P. Shenoy, and K. Ramamritham. Maintaining Mutual Consistency for Cached Web Objects. In *Proceedings of the 21st International Conference on Distributed Computing Systems (ICDCS-21), Phoenix, AZ*, pages 371–380, April 2001.
- [26] A. Venkataramani, P. Yalagandula, R. Kokku, S. Sharif, and M. Dahlin. Potential Costs and Benefits of Long-term Prefetching for Content Distribution. In *Proceedings of the Web Caching Worskhop, Boston, MA*, June 2001.

- [27] J. Yin, L. Alvisi, M. Dahlin, and A. Iyengar. Engineering Server-driven Consistency for Large-scale Dynamic Web Services. In *Proceedings of the 10th World Wide Web Conference, Hong Kong, May 2001*.
- [28] J. Yin, L. Alvisi, M. Dahlin, and C. Lin. Volume Leases for Consistency in Large-Scale Systems. *IEEE Transactions on Knowledge and Data Engineering*, January 1999.
- [29] Jian Yin, Lorenzo Alvisi, Mike Dahlin, and Calvin Lin. Hierarchical Cache Consistency in a WAN. In *Proceedings of the USENIX Symposium on Internet Technologies and Systems*, Boulder, Colorado, October 1999.
- [30] H. Yu, L. Breslau, and S. Shenker. A Scalable Web Cache Consistency Architecture. In *Proceedings of the ACM SIGCOMM'99, Boston, MA, September 1999*.

**Anoop George Ninan** received the B.E (Honors) degree in Computer Science and Engineering from the Motilal Nehru Regional Engineering College (MNREC), Allahabad, India in 1996 and the M.S degree in Computer Science from the University of Massachusetts, Amherst in 2001. Between MNREC and UMass, he worked on the MasterCraft project at the Tata Research Development and Design Center, Pune, India. He is currently positioned as a Senior Software Engineer at EMC Corporation, Hopkinton, Massachusetts and works on performance, scalability and capacity planning aspects of EMC's ControlCenter suite of storage management products.

His interests include distributed systems, storage systems/networking/management, web caching and model-driven architectures.

**Purushottam Kulkarni** received the B.E degree in Computer Science from the University of Pune, India in 1997 and the M.S degree in Computer Science from the University of Minnesota, Duluth in 2000. He is currently a Ph.D candidate in the Department of Computer Science at the University of Massachusetts, Amherst. Before joining the University of Minnesota, he worked as a Software Engineer at Tata Technologies India Limited, India.

His research interests include scalable data dissemination over the web, Internet technologies, distributed systems and mobile computing.

**Prashant Shenoy** received the B.Tech degree in Computer Science and Engineering from the Indian Institute of Technology, Bombay in 1993, and the M.S and Ph.D degrees in Computer Science from the University of Texas, Austin, in 1994 and 1998, respectively. He is currently an Assistant Professor in the Department of Computer Science, University of Massachusetts, Amherst.

His research interests are multimedia systems, operating systems, computer networks, and distributed systems. Dr. Shenoy is a member of the Association for Computing Machinery (ACM). He has been the recipient of the National Science Foundation Career Award, the IBM Faculty Development Award, the Lilly Foundation Teaching Fellowship, and the UT Computer Science Best Dissertation Award.

**Krithi Ramamritham** received the Ph.D degree in Computer Science from the University of Utah, Salt Lake City, and then joined the University of Massachusetts, Amherst. He is currently at the

Indian Institute of Technology, Bombay as the Vijay and Sita Vashee Chair Professor in the Department of Computer Science and Engineering.

His interests span the areas of real-time systems, transaction processing in database systems, and real-time databases systems. He is applying concepts from these areas to solve problems in mobile computing, e-commerce, intelligent internet and the Web. Dr. Ramamritham is a Fellow of the IEEE and a Fellow of the ACM.

**Renu Tewari** received the Ph.D degree in Computer Science from the University of Texas, Austin in 1998. She is currently a research staff member at the IBM Almaden Research Center.

Her main interests are web caching, content distribution networks, edge services, web server architecture and quality-of-service. She is actively involved in IETF standards related to caching and CDNs.