

Clustering Support and Replication Management for Scalable Network Services

Kai Shen, Tao Yang, *Member, IEEE*, and Lingkun Chu

Abstract—The ubiquity of the Internet and various intranets has brought about widespread availability of online services and applications accessible through the network. Cluster-based network services have been rapidly emerging due to their cost-effectiveness in achieving high availability and incremental scalability. This paper presents the design and implementation of the *Neptune* middleware system that provides clustering support and replication management for scalable network services. Neptune employs a loosely connected and functionally symmetric clustering architecture to achieve high scalability and robustness. It shields the clustering complexities from application developers through simple programming interfaces. In addition, Neptune provides replication management with flexible replication consistency support at the clustering middleware level. Such support can be easily applied to a large number of applications with different underlying data management mechanisms or service semantics. The system has been implemented on Linux and Solaris clusters, where a number of applications have been successfully deployed. Our evaluations demonstrate the system performance and smooth failure recovery achieved by proposed techniques.

Index Terms—Network services, programming support, replication management, failure recovery, load balancing.

1 INTRODUCTION

RECENT years have seen widespread availability of online services and applications accessible through the network. Examples include document search engines [6], [19], online digital libraries [1], data mining on scientific data sets [8], online discussion forums [25], and electronic commerce [15]. Cluster-based network services have become particularly popular due to their cost-effectiveness in achieving high availability and incremental scalability, especially when the system experiences high growth in service evolution and user demands. Within a large-scale complex service cluster, service components are usually partitioned, replicated, and aggregated to fulfill external requests. Despite its importance, service clustering to achieve high scalability, availability, and manageability remains a challenging task for service designers. And, this is especially true for services with frequently updated persistent data. This paper recognizes the importance of service clustering support and its challenges. The main goal of this study is to answer the following question. *Given a service application running on a single machine with a modest amount of data, how can such a service be expanded quickly to run on a cluster environment for handling a large volume of concurrent request traffic with large-scale persistent service data?*

This paper examines the design and implementation of the *Neptune* middleware system, which provides clustering support and replication management for scalable network services. The main design goal of the Neptune clustering middleware is to support a simple, flexible, yet efficient model in aggregating and replicating network service

components. Neptune employs a loosely connected and functionally symmetric clustering architecture to achieve high scalability and robustness. This architecture allows Neptune service infrastructure to operate smoothly in the presence of transient failures and service evolutions. Additionally, Neptune provides simple programming interfaces to shield application developers from clustering complexities.

Replication of persistent data is crucial to load sharing and achieving high availability. However, it introduces the complexity of consistency management when service requests trigger frequent updates on the persistent data. Though large-scale network services often do not require strong ACID consistencies [9], [21], it is important to maintain *eventual consistency* [17] and prevent the loss of any accepted updates. Besides, users of time-critical services (e.g., auctions) demand certain guarantees to guard against accessing stale data. This paper investigates techniques in building a replication management framework as part of the Neptune clustering middleware. This work is closely related to a large body of previous research in network service clustering and data replication. In particular, recent projects have developed scalable replication support for specific cluster-based network services, including distributed hashtable [22] and the e-mail service [28]. The key contribution of our work is the support of replication management at the clustering middleware level. Such support can be easily applied to a large number of applications with different underlying data management mechanisms or service semantics.

The rest of this paper is organized as follows: Section 2 presents the Neptune clustering architecture. Section 3 describes Neptune's multilevel service replication support. Section 4 discusses the system implementation and service deployments. Section 5 presents the performance evaluation. Section 6 discusses the related work and Section 7 concludes this paper.

- K. Shen is with the Department of Computer Science, University of Rochester, Rochester, NY 14627. E-mail: kshen@cs.rochester.edu.
- T. Yang and L. Chu are with the Department of Computer Science, University of California, Santa Barbara, CA 93106. T. Yang is also with Ask Jeeves. E-mail: {tyang, lkchu}@cs.ucsb.edu.

Manuscript received 8 Dec. 2002; revised 31 July 2003; accepted 3 Aug. 2003. For information on obtaining reprints of this article, please send e-mail to: tps@computer.org, and reference IEEECS Log Number 118753.

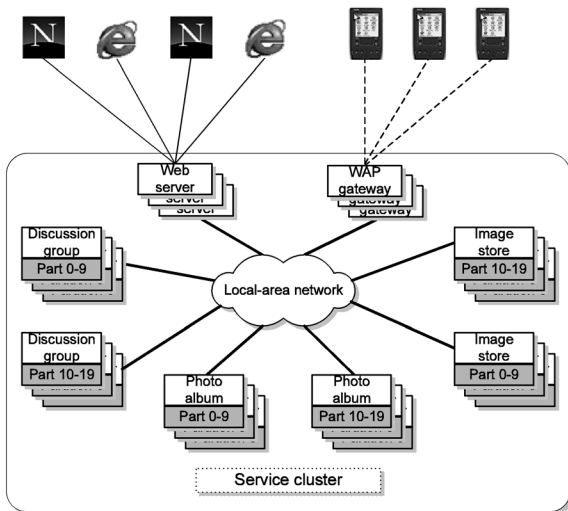


Fig. 1. Architecture for a discussion group and a photo album service.

2 NEPTUNE CLUSTERING ARCHITECTURE

Generally speaking, providing standard system components at the middleware level to support service clustering and replication tends to decrease the flexibility of service construction. Neptune demonstrates that it is possible to achieve these goals by taking advantage of the following characteristics existing in many network services:

1. *Information independence.* Network services tend to host a large amount of information addressing different and independent categories. For example, an auction site hosts different categories of items. Every bid only accesses data concerning a single item, thus providing an intuitive way to partition the data.
2. *User independence.* Information accessed by different users tends to be independent. Therefore, data may also be partitioned according to user accounts. E-mail service and Web page hosting are two such examples.

With these characteristics in mind, Neptune targets *partitionable* network services in the sense that persistent data manipulated by such a service can be divided into a large number of independent partitions and each service access can be delivered independently on a single partition or each access is an aggregate of a set of subaccesses, each of which can be completed independently on a single partition.

Fig. 1 illustrates the service cluster architecture for an online discussion group and a photo album service, similar to the MSN Groups service [25]. The discussion group service allows a group of users to post, read, and comment on messages. The photo album service provides a shared repository for a group of users to store, view, and edit photo images. The service cluster delivers these services to wide-area browsers and wireless clients through Web servers and WAP gateways. Note that the image store is an internal service in that it is only invoked by the photo album service components to fulfill external photo album service requests. In this example, all persistent service data is divided into 20 partitions according to user groups. All the components

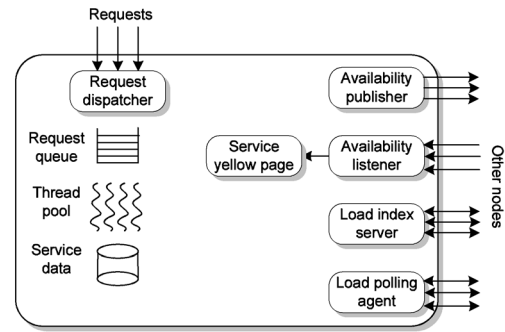


Fig. 2. Component architecture of a Neptune node.

(including protocol gateways) are replicated in the service cluster.

Previous work has recognized the importance of providing software infrastructures for cluster-based network services [17], [32]. Most of these systems rely on central components to maintain the server runtime workload and service availability information. In comparison, Neptune employs a loosely connected and functionally symmetric architecture in constructing the service cluster. This architecture allows the Neptune service infrastructure to operate smoothly in the presence of transient failures and service evolutions. The rest of this section examines the overall Neptune clustering architecture, its load balancing support, and programming interfaces.

2.1 Overall Clustering Architecture

Neptune encapsulates an application-level network service through a service access interface which contains several RPC-like access methods. Each service access through one of these methods can be fulfilled exclusively on one data partition. Neptune employs a *functionally symmetric* approach in constructing the server cluster. Every node in the cluster contains the same Neptune components and is functionally equal to each other. Each cluster node can elect to host service components with certain data partitions and it can also access service components hosted at other nodes in the cluster. Within a Neptune cluster, all the nodes are *loosely connected* through a publish/subscribe channel. This channel can be implemented using IP multicast or through a highly available well-known central directory.

Fig. 2 illustrates the architecture of a Neptune node. Since Neptune service clusters are functionally symmetric, all the nodes are based on the same architecture. The key components are described as follows: The *request dispatcher* directs an incoming request to the hosted service instance which contains a request queue and a pool of worker threads or processes. When all worker threads or processes are busy, subsequent requests will be queued. This scheme allows a Neptune server to gracefully handle spikes in the request volume while maintaining a desired level of concurrency. The *availability publisher* periodically announces the locally hosted service component, the data partitions, and the access interface to other nodes in the cluster. The *availability listener* monitors those announcements from other nodes and maintains a local *service yellow page* containing available service components and their locations in the cluster. The availability information in the yellow page is kept as soft state such that it has to be refreshed repeatedly to stay alive. This allows node failures

and system upgrades to be handled smoothly without complex protocol handshakes. For the load balancing purpose, each node can poll the load indexes of other nodes through a local *load polling agent*. The polls are responded to by *load index servers* at the polled servers. Neptune employs a random polling-based load balancing policy, described in the next section. The request dispatcher, request queue, thread/process pool, availability publisher, and load index server can be considered as service-side components in a service invocation, and they are in an aggregate called the *Neptune service module*. Similarly, the availability listener, service yellow page, and the load polling agent together are called the *Neptune client module*.

Basically, each service access is initiated at the requesting node with a service name, a data partition ID, a service method name, and a read/write access mode.¹ Then, the Neptune client module transparently selects a service node based on the service/partition availability, access mode, consistency requirement, runtime workload, and the load balancing policy. Upon receiving the access request, the Neptune server module in the chosen node dispatches a service instance to fulfill the service access. In the case of service replication, further request propagations may occur due to consistency management. Details on this will be discussed in Section 3.

In a Neptune-enabled service cluster, the application-level service programmer only needs to implement the stand-alone application components. The aggregation of multiple data partitions or application components as well as data replication are supported transparently by Neptune clustering modules. For the service cluster shown in Fig. 1, each photo album service instance locates and accesses an image store service instance through the local Neptune client module. In addition, each gateway node relies on its own Neptune client module to export the discussion group and photo album services to external users. Overall, Neptune’s loosely connected and functionally symmetric architecture provides the following advantages:

- *Scalability and robustness*. The functionally symmetric architecture removes scaling bottlenecks from the system. It also exhibits strong robustness in the face of random or targeted failures because a failure of one node is no more disastrous than that of any other.
- *Manageability*. Neptune allows service infrastructure to operate smoothly through service evolutions. For instance, new nodes can be added into the cluster by simply announcing locally hosted service components, data partitions, and the access interfaces. Therefore, the process of increasing system capacity or upgrading service features can be completed transparently to the on-going operations in the system.

2.2 Cluster Load Balancing

Neptune employs a random polling-based load balancing policy for the cluster-wide request distribution. For every service access, the service client polls several randomly selected servers for load information and then directs the

service access to the server responding with lightest load. The load index we use is the total number of active and queued service accesses at each server. Prior studies have suggested the resource queue lengths can be an excellent predictor of service response time [16], [34]. The random polling policy fits seamlessly into Neptune’s functionally symmetrical architecture because it does not require any centralized components. The load balancing is conducted by polling agents and load index servers, both of which exist at each node. An important parameter for a random polling policy is the poll size. Mitzenmacher demonstrated, through analytical models, that a poll size of two leads to an exponential improvement over pure random policy, but a poll size larger than two leads to much less substantial additional improvement [24]. By limiting the poll size, we can ensure the polling overhead for each service access does not increase with the growth of the system size, which is important for the system scalability.

In practice, Neptune’s cluster load balancing support is based on a random polling policy with a poll size of three. On top of the basic random polling policy, we also made an enhancement by discarding slow-responding polls. Through a ping-pong test on two idle machines in a Linux cluster connected by a switched 100Mbps Ethernet, we measured that a UDP-based polling roundtrip cost is around $290\mu\text{s}$. However, it may take much longer than that for a busy server to respond to a polling request. Such polling delays can affect the service response time and this impact is more severe for fine-grain services with small service processing time. With this in mind, we enhanced the basic polling policy by discarding polls not responded to within 10ms, which is the smallest timeout granularity supported by the `select` system call in Linux. Intuitively, this results in a trade off between spending less time waiting for polls and acquiring more polled information. In addition, slow-responding polls deliver inaccurate load information due to long delays. Discarding those slow-responding polls can avoid using such stale load information, serving as an additional advantage. This again tends to be more significant for fine-grain services.

2.3 Neptune Programming Interfaces

Neptune supports two communication schemes between clients and service instances inside the cluster: a request/response scheme and a stream-based scheme. In the request/response scheme, the client and the service instance communicate with each other through a request message and a response message. For the stream-based scheme, Neptune sets up a bidirectional stream between the client and the service instance as a result of the service invocation. Stream-based communication can be used for asynchronous service invocation and it also allows multiple rounds of interaction between the client and the service instance. Neptune only supports stream-based communication for read-only service accesses due to the difficulty of replicating and logging streams. Below, we briefly describe the interfaces between Neptune and service modules at both the client and the service sides for the request/response communication scheme. The complete programming interface can be found in the Neptune programming guide [12].

- At the client side, Neptune provides a unified interface to service clients for seeking location-transparent request/response service access. It is

1. We classify a service access as a *read access* (or *read*, in short) if it does not change the persistent service data, or as a *write access* (or *write*, in short) otherwise.

shown below in a language-neutral format:
 NeptuneRequest (NeptuneHandle, ServiceName, PartitionID, ServiceMethod, AccessMode, RequestMsg, ResponseMsg);

A NeptuneHandle should be used in every service request for a client session. It maintains the information related to each client session and we will discuss it further in Section 3.2. The meanings of other parameters are straightforward.

- At the service side, all the service method implementations need to be registered at the service deployment phase. This allows the Neptune service module to invoke the corresponding service instance when a request is received.

3 REPLICATION MANAGEMENT

In this section, we investigate techniques in providing scalable replication support for cluster-based network services. The goal of this work is to provide simple, flexible, and yet efficient service replication support for network services with frequently updated persistent data. This model should make it simple to deploy existing applications and shield application programmers from the complexities of replication consistency and fail-over support. It also needs to have the flexibility to accommodate a variety of data management mechanisms that network services typically rely on. In order to support applications with different service semantics and consistency requirements, Neptune employs a multilevel replication consistency scheme. The rest of this section describes the underlying assumptions of our work, Neptune's replication consistency support, and its failure recovery model.

3.1 Assumptions for Replication Management

We assume all hardware and software system components follow the fail-stop failure model and network partitions do not occur inside the service cluster. Following this principle, each component will simply terminate itself when an unexpected situation occurs. Nevertheless, we do not preclude catastrophic failures in our model. In other words, persistent data can survive through a failure that involves a large number of components or even all nodes. The replication consistency will be maintained after the recovery. This is important because software failures are often not independent. For instance, a replica failure triggered by high workload results in even higher workload in the remaining replicas and may cause cascading failures of all replicas.

Neptune supports atomic execution of service operations through failures only if each underlying service component can ensure atomicity in a stand-alone configuration. This assumption can be met when the persistent data is maintained in transactional databases or transactional file systems. To facilitate atomic execution, we assume that each service component provides a CHECK callback so that the Neptune service module can check if a previously spawned service instance has been successfully completed. The CHECK callback is very similar to the REDO and UNDO callbacks that resource managers provide in the transaction processing environment [9], [21]. It is only invoked during the node recovery phase and we will discuss its usage and a potential implementation further in Section 3.3.

We also assume that the service abort happens either at all replicas or not at all, which frees Neptune from coordinating service aborts. Note that a server failure does not necessarily cause the active service request to abort because it can be logged and successfully reissued when the failing server recovers. This assumption can be met when service aborts are only related to the state of service data, e.g., violations of data integrity constraints. In other words, we do not consider the situation in which a service aborts due to insufficient disk space. This assumption is not crucial to Neptune's fundamental correctness, but it greatly simplifies our implementation. Without such an assumption, a proper UNDO mechanism will be required to maintain replication consistency.

3.2 Multilevel Replication Consistency Model

In general, data replication is achieved through either *eager* or *lazy* write propagations [3], [13], [20]. Eager propagation keeps all replicas exactly synchronized by acquiring locks and updating data at all replicas within a globally coordinated transaction. In comparison, lazy propagation allows lock acquisitions and data updates to be completed independently at each replica. Previous work shows that synchronous eager propagation leads to high deadlock rates when the number of replicas increases [20]. The DDS project uses this synchronous approach and they rely on the timeout abort and client retry to resolve the deadlock issue [22]. Neptune's consistency model extends the previous work in lazy propagation with a focus on high scalability and runtime fail-over support. Lazy propagation introduces the problems of out-of-order writes and reading stale data versions. In order to support applications with different service semantics and consistency requirements, Neptune employs a flexible three-level replication consistency scheme to address these problems. In particular, Neptune's highest consistency level provides a staleness control which contains a quantitative staleness bound and a guarantee of progressive version delivery for each client's perspective.

Level 1: Write-anywhere replication for commutative writes. In this level, each write is initiated at any replica and is propagated to other replicas asynchronously. When writes are commutative, eventually the persistent data will converge to a consistent state after all outstanding writes are completed. The append-only discussion groups in which users can only append messages satisfy this commutativity requirement. Another example is a certain kind of e-mail service [28] in which all writes are total-updates, so out-of-order writes could be resolved by discarding all but the newest. The first level of replication consistency is intended to take advantage of application characteristics and to achieve high performance in terms of scalability and fail-over support.

Level 2: Primary-secondary replication for ordered writes. In this consistency level, writes for each data partition are totally ordered. A primary-copy node is assigned to each replicated data partition, and other replicas are considered as secondaries. All writes for a data partition are initiated at the primary, which asynchronously propagates them in an FIFO order to the secondaries. At each replica, writes for each partition are serialized to preserve the order. Serializing writes simplifies the write ordering for each partition, but it results in a loss of write concurrency within each partition. Since many network services have a large number

of data partitions due to the information and user independence, there should be sufficient write concurrency across partitions. Besides, concurrency among read operations is not affected by this scheme. Level two consistency provides the same client-viewed consistency support as level one without requiring writes to be commutative. As a result, it could be applied to more services.

Level 3: Primary-secondary replication with staleness control. Level two consistency is intended to solve the out-of-order write problem resulting from lazy propagation. This additional level is designed to address the issue of reading stale data versions. The primary-copy scheme is still used to order writes in this level. In addition, we assign a version number to each data partition and this number increments after each write. The staleness control provided by this consistency level contains two parts:

1. *Soft quantitative bound.* Each read is serviced at a replica that is at most a certain amount of time (x seconds) stale compared to the primary version. The quantitative staleness between two data versions is defined by the elapsed time between the two corresponding writes accepted at the primary. Thus, our scheme does not require a global synchronous clock. Currently, Neptune only provides a soft quantitative staleness bound and its implementation is described later in this section.
2. *Progressive version delivery.* From each client's point of view, the data versions used to service her read and write accesses should be monotonically non-decreasing. Both guarantees are important for services like online auction in which users would like to get as recent information as possible and they do not expect to see declining bidding prices.

We describe our implementation for the two staleness control guarantees in level three consistency. The quantitative bound ensures that all reads are serviced at a replica of at most x seconds stale compared to the primary version. In order to achieve this, each replica publishes its current version number as part of the service availability announcement and the primary publishes its version number at x seconds ago in addition. With this information, the Neptune client module can ensure that all reads are only directed to replicas within the specified quantitative staleness bound. Note that the published replica version number may be stale depending on the service publishing frequency, so it is possible that none of the replicas has a high enough version number to fulfill a request. In this case, the read is directed to the primary, which always has the latest version. Also, note that the " x seconds" is only a soft bound because the real guarantee depends on the latency, frequency, and intermittent losses of service availability announcements. However, these problems are insignificant in a low latency, reliable local area network.

The progressive version delivery guarantees that: 1) After a client writes to a data partition, she always sees the result of this write in her subsequent reads. 2) A user never reads a version that is older than another version she has seen before. In order to accomplish these guarantees, each service invocation returns a version number to the client side. For a read, this number stands for the data version used to fulfill this access. For a write, it stands for the latest data version as a result of this write. Each client keeps this

version number in a session handle (a.k.a., `NeptuneHandle`) and carries it in each service invocation. The Neptune client module can ensure that each client read access is directed to a replica with a published version number higher than any previously returned version number.

As mentioned in Section 2, Neptune targets partitionable network services in which service data can be divided into independent partitions. Therefore, Neptune's consistency model does not address data consistency across partition boundaries. Note that a consistency level is specified for each service and, thus, Neptune allows coexistence of services with different consistency levels.

3.3 Failure Recovery

In this section, we focus on the failure detection and recovery for the primary-copy scheme that is used in level two/three consistency schemes. The failure management for level one consistency is much simpler because the replicas are more independent of each other.

In order to recover lost propagations after failures, each Neptune service node maintains a REDO write log for each data partition it hosts. Each log entry contains the service method name, partition ID, the request message, along with an assigned *log sequence number (LSN)*. The write log consists of a committed portion and an uncommitted portion. The committed portion records those writes that are already completed, while the uncommitted portion records the writes that are received but not yet completed.

Neptune assigns a static priority for each replica of a data partition. The primary is the replica with the highest priority. When a node failure is detected, for each partition that the faulty node is the primary of, the remaining replica with the highest priority is elected to become the new primary. This election algorithm is the same as the classical Bully Algorithm [18] with the exception that each replica has a priority *for each data partition* it hosts. This fail-over scheme also requires that the elected primary does not miss any write that has committed in the failed primary. To ensure that, before the primary executes a write locally, it has to wait until all other replicas have acknowledged the reception of its propagation. If a replica does not acknowledge in a timeout period, this replica is considered to have failed due to our fail-stop assumption and, thus, this replica can only rejoin the service cluster after going through a recovery process.

When a node recovers after its failure, the underlying single-site service component first recovers its data into a consistent state. Then, this node will enter Neptune's three-phase recovery process as follows:

Phase 1: Internal synchronization. The recovering node first synchronizes its write log with the underlying service component. This is done by using registered CHECK callbacks to determine whether each write in the uncommitted log has been completed by the service component. The completed writes are merged into the committed portion of the write log and the uncompleted writes are reissued for execution.

Phase 2: Missing write recovery. In this phase, the recovering node announces its priority for each data partition it hosts. For each partition for which the recovering node has a higher priority than the current primary, this node will bully the current primary into a secondary as soon as its priority announcement is heard. Then, it contacts

the deposed primary to recover the writes that it missed during its down time. For each partition for which the recovering node does not have a higher priority than the current primary, this node simply contacts the primary to recover the missed writes.

Phase 3: Operation resumption. After the missed writes are recovered, this recovering node resumes normal operations by publishing the services it hosts and accepting requests from the clients.

Note that, if the recovering node has the highest priority for some data partitions it hosts, there will be no primary available for those partitions during phase two of the recovery. This temporary blocking of writes is essential to ensure that the recovering node can bring itself up-to-date before taking over as the new primary. We will present the experimental study for this behavior in Section 5.3. We also want to emphasize that a catastrophic failure that causes all replicas for a certain partition to fail requires special attention. No replica can successfully complete phase two recovery after such a failure because there is no preexisting primary in the system to recover missed writes. In this case, the replica with newest version needs to be manually brought up as the primary and then all other replicas can proceed according to the standard three-phase recovery.

Before concluding the discussion on our failure recovery model, we describe a possible implementation of the CHECK callback facility. We require the Neptune service module to pass the LSN with each request to the service instance. Then, the service instance fulfills the request and records this LSN on persistent storage. When the CHECK callback is invoked with an LSN during a recovery, the service component compares it with the LSN of the latest completed service access and returns appropriately. As we mentioned in Section 3.1, Neptune provides atomic execution through failures only if the underlying service component can ensure atomicity on single-site service accesses. Such support can ensure the service access and the recording of LSN take place as an atomic action. A transactional database or a transactional file system can be used to achieve atomicity for single-site service accesses.

4 SYSTEM IMPLEMENTATION AND SERVICE DEPLOYMENTS

Neptune has been implemented on Linux and Solaris clusters. The publish/subscribe channel in this implementation is realized using IP multicast. Each multicast message contains the service announcement and node runtime CPU and I/O workload, acquired through the `/proc` file system in Linux or the `kstat` facility in Solaris. We limit the size of each multicast packet to be within an Ethernet maximum transmission unit (MTU) in order to minimize the multicast overhead. We let each node send the multicast message once every second and the published information is kept as soft state in the service yellow page of each node, expiring in five seconds. That means a faulty node will be detected when five of its multicast messages are not heard in a row. This “soft state”-based node aliveness information can be inaccurate at times, especially for servers that keep going up and down. As a result, service connection setup may fail due to false node aliveness information. Neptune attempts three retries in these cases, after which failing nodes are excluded from local

service yellow page. Those failing nodes will be added back when future availability announcements are heard.

Inside each service node, hosted services can be compiled into a dynamically linked library. They are linked into Neptune process space at runtime and run as threads. Alternatively, each service instance could run as a separate OS process, which provides better fault isolation and resource control at the cost of degraded performance. In choosing between thread and process-based deployment, the rule of thumb is to pick threads for simple and short-running services while using processes for complex and large services.

Overall, this implementation incurs moderate overhead. For instance, we implemented an echo service on top of a Neptune-enabled Linux cluster connected by a switched 100 Mbps Ethernet. The echo service simply responds to the client with a message identical to the request message. The response time of such an echo service is measured at $1,128\mu\text{s}$, excluding the polling overhead incurred in service load balancing. Taking away the TCP roundtrip time with connection setup and teardown, which is measured at $1,031\mu\text{s}$, Neptune is responsible for an overhead of $97\mu\text{s}$ in each service invocation. So far, we have deployed a document search engine, an online discussion group, an auction, and a persistent cache service on Neptune clusters, described as follows. Document search serves read-only user requests, thus it is not concerned with the issue of replication consistency.

- *Document search* takes in a group of encoded query words, checks a memory mapped index database, and returns the identifications of the list of documents matching the input query words. Neptune has been deployed at Web search engines Teoma [29] and Ask Jeeves [6]. As of Summer 2003, a Neptune-enabled Linux/Solaris cluster at Ask Jeeves contains over 1,200 SMP servers and maintains a search index of more than 1 billion Web documents.
- *Online discussion group* handles three types of requests for each discussion topic: viewing the list of message headers (`ViewHeaders`), viewing the content of a message (`ViewMsg`), and adding a new message (`AddMsg`). Both `ViewHeaders` and `ViewMsg` are read-only requests. The messages are maintained and displayed in a hierarchical format according to the reply-to relationships among them. The discussion group uses MySQL database to store and retrieve messages and topics.
- A prototype *auction* service is also implemented on MySQL database. Level three consistency with proper staleness bound and progressive version delivery is desirable for this service because auction users are sensitive to data staleness.
- *Persistent cache* supports two service methods: storing a key/data pair into the persistent cache (`CacheUpdate`) and retrieving the data for a given key (`CacheLookup`). The persistent cache uses an MD5-based hashing function to map the key space into a set of buckets. We use the `mmap` utility to keep an in-memory reference to the disk data and we purge the updates and the corresponding LSN into the disk at every 10th `CacheUpdate` invocation. The

LSN is used to support the CHECK callback that we discussed in Section 3.3. The persistent cache is typically used as an internal utility, providing a scalable and reliable data store for other services. The level two consistency can be used for this service, which allows high throughput at the cost of intermittent false cache misses.

5 PERFORMANCE EVALUATION

Our evaluation examines the performance-scalability, impact of workload variation, and fault tolerance for Neptune's replication management. In addition, we also provide a performance evaluation on Neptune's cluster load balancing support. All the evaluations are conducted on a rack-mounted Linux cluster with 30 dual 400MHz Pentium II nodes, each of which contains 512MB or 1GB memory. Each node runs Linux 2.2.15 and has two 100Mbps Ethernet interfaces. The cluster is connected by a Lucent P550 Ethernet switch with 22Gbps backplane bandwidth. MySQL 3.23.22-Beta is used for the discussion group and auction services. We note that MySQL database does not have full-fledged transactional support, but its latest version supports "atomic operations," which is enough for Neptune to provide cluster-wide atomicity. Our persistent cache is built on a regular file system without atomic recovery support. We believe such a setting is sufficient for illustrative purposes.

This paper is focused on the clustering support and replication management within local-area server clusters. Networking issues between protocol gateways and wide-area clients are out of the scope of this study. Therefore, our experiments target the performance evaluation between client nodes and service nodes inside a Neptune service cluster. All the experiments presented in this section use up to 16 service nodes and six client nodes.

We use synthetic workloads in our evaluations on replication management. Two types of workloads are generated for this purpose: 1) *Balanced workloads*, where service requests are evenly distributed among data partitions are used to measure the best case scalability. 2) *Skewed workloads* are used to measure the system performance when some particular partitions draw a disproportional number of service requests. In our performance report, the system throughput is measured as the maximum request rate at which more than 98 percent of client requests can be successfully completed within two seconds.

5.1 Replication Management under Balanced Workload

We use the discussion group to study the system scalability under balanced workload. We vary the replication degree, the number of service nodes, the write percentage, and consistency levels in this evaluation. The replication degree is the number of replicas for each partition. The write percentage indicates the proportion of writes in all requests. We use two write percentages in our evaluation: 10 percent and 50 percent. We measure all three consistency levels in this study. Level one consistency requires writes to be commutative and, thus, we use a variation of the original service implementation to facilitate it. For the purpose of performance comparison with other consistency levels, we keep the changes to a minimum. For level three consistency, we choose one second as the staleness bound. We also notice that the performance of level three consistency is

affected by the request rate of individual clients. This is because a higher request rate from each client means a higher chance that a read has to be forwarded to the primary node to fulfill progressive version control which, in turn, restricts the system load balancing capabilities. We choose a very high per-client request rate (one request per second) in this evaluation to measure the worst-case scenario.

The number of discussion groups in our synthetic workload is 400 times the number of service nodes. Those groups are in turn divided into 64 partitions. These partitions and their replicas are evenly distributed across service nodes in this evaluation. Each request is sent to a discussion group chosen according to an even distribution. The distribution of different requests (`AddMsg`, `ViewHeaders`, and `ViewMsg`) is determined based on the write percentage.

Fig. 3 shows the scalability of discussion group service at various configurations. Each subfigure illustrates the system performance under no replication (`NoRep`) and replication degrees of two, three, and four. The `NoRep` performance is acquired through running a stripped down version of Neptune that does not incur any replication overhead except logging. The single node performance under no replication is 152 requests/second for 10 percent writes and 175 requests/second for 50 percent writes. Note that a read is more costly than a write because `ViewHeaders` displays the message headers in a hierarchical format according to the reply-to relationships, which may invoke expensive SQL queries.

We can draw the following conclusions based on the results in Fig. 3:

1. When the number of service nodes increases, the throughput steadily scales across all replication degrees.
2. Service replication comes with an overhead because every write has to be executed more than once. Not surprisingly, this overhead is more significant under higher write percentage. In general, a nonreplicated service performs twice as fast as its counterpart with a replication degree of four at 50 percent writes. However, Section 5.2 shows that replicated services can outperform nonreplicated services under skewed workloads due to better load balancing.
3. All three consistency levels perform very similarly under balanced workloads. This means level one consistency does not provide a significant performance advantage and our staleness control does not incur significant overhead either.

5.2 Impact of Partition Imbalance for Replication Management

This section studies the performance impact of partition imbalance. Each skewed workload in this study consists of requests that are chosen from a set of partitions according to a Zipf distribution. Each workload is also labeled with a *partition imbalance factor*, which indicates the proportion of the requests that are directed to the most popular partition. For a service with 64 partitions, a workload with an imbalance factor of 1/64 is completely balanced. A workload with an imbalance factor of 1 is the other extreme in which all requests are directed to one single partition.

Fig. 4a shows the impact of partition imbalance on services with different replication degree for the discussion

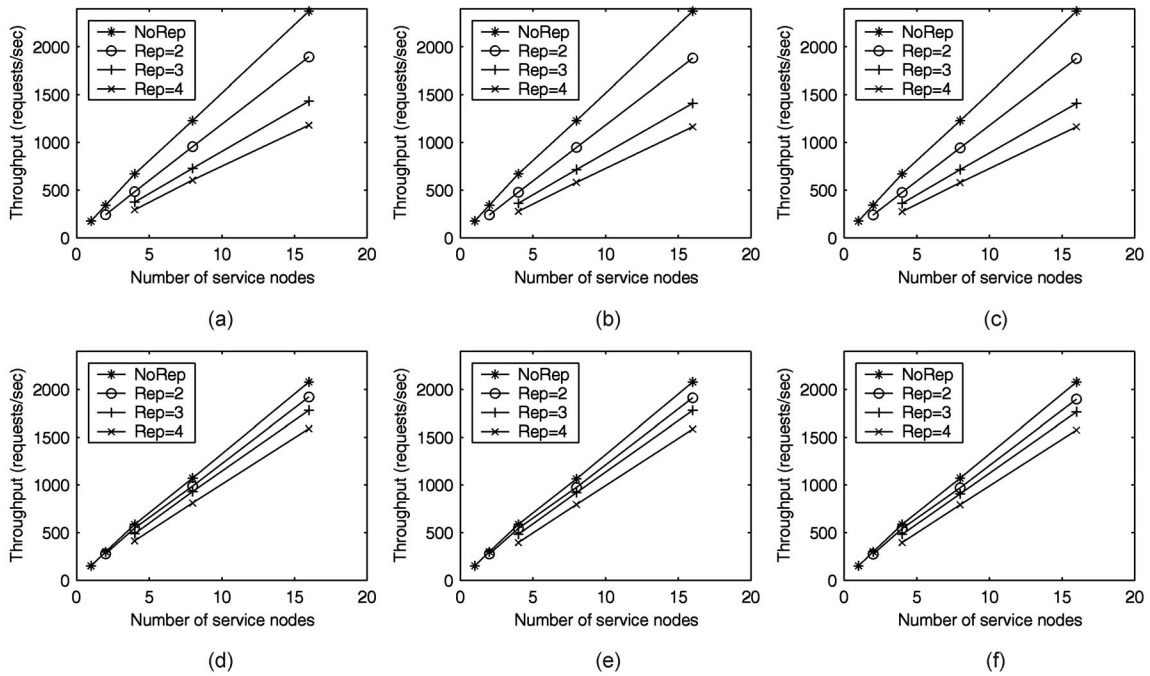


Fig. 3. Scalability of the discussion group service under balanced workload. (a) Level 1, 50 percent writes, (b) Level 2, 50 percent writes, (c) Level 3, 50 percent writes, (d) Level 1, 10 percent writes, (e) Level 2, 10 percent writes, and (f) Level 3, 10 percent writes.

group service. The 10 percent write percentage, level two consistency, and eight service nodes are used in this experiment. We see that, even though service replication carries an overhead under balanced workload (imbalance factor = $1/64$), replicated services can outperform non-replicated ones under skewed workload. Specifically, under the workload where all requests are directed to one single partition, the service with a replication degree of four performs almost three times as fast as its nonreplicated counterpart. This is because service replication provides better load sharing by spreading hot-spots over several service nodes, which completely amortizes the overhead of extra writes in achieving the replication consistency.

We learned from Section 5.1 that all three consistency levels perform very closely under balanced workload. Fig. 4b illustrates the impact of workload imbalance on different consistency levels. The 10 percent write percentage, a replication degree of four, and eight service nodes are used in this experiment. The performance difference among three consistency levels becomes slightly more prominent when the partition imbalance factor increases. Specifically, under the workload where all requests are directed to one single partition, level one consistency yields 12 percent better performance than level two consistency which, in turn, performs 9 percent faster than level three consistency with staleness control at one second. Based on these results, we learned that: 1) The freedom of directing writes to any replica in level one consistency only yields moderate performance advantage and 2) our staleness control scheme carries an insignificant overhead, even though it appears slightly larger for skewed workload.

5.3 Replication Management during Failure Recoveries

Fig. 5 depicts the behavior of a Neptune-enabled discussion group service during three node failures in a 200-second period. Eight service nodes, level two consistency, and a

replication degree of four are used in the experiments. Three service nodes fail simultaneously at time 50. Node 1 recovers 30 seconds later. Node 2 recovers at time 110 and node 3 recovers at time 140. It is worth mentioning that a recovery may take much longer than 30 seconds in practice, especially when large data files need to be loaded over the network as part of such recovery. However, we believe

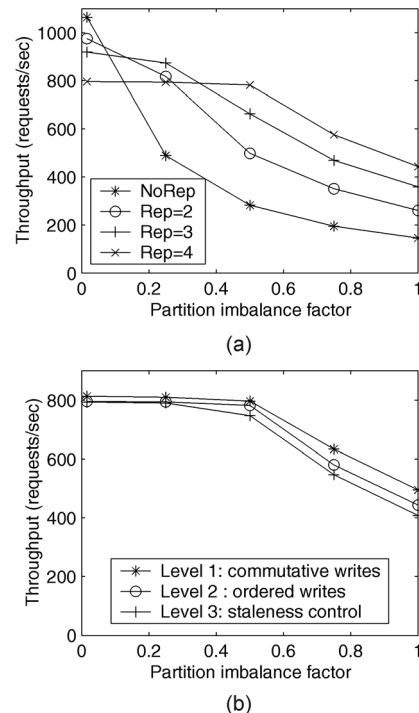


Fig. 4. Impact of partition imbalance on an 8-node system. (a) Impact on replication degrees and (b) impact on consistency levels.

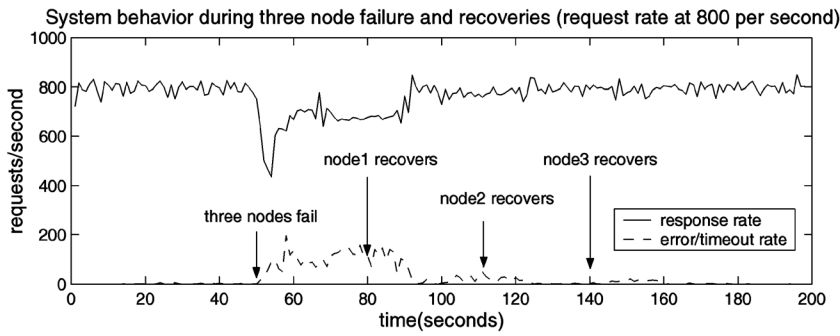


Fig. 5. Behavior of the discussion group service during three node failure and recoveries. Eight service nodes, level two consistency, and a replication degree of four are used in this experiment.

those variations do not affect the fundamental system behavior illustrated in this experiment. We observe that the system throughput goes down during the failure period. A sharp performance drop at the beginning of the failure period is caused by the temporary loss of primaries, which is remedied after the elections of new primaries. And, we also observe a tail of errors and timeouts trailing each node recovery. This is caused by the overhead of synchronizing lost updates as part of the recovery. Overall, the service quickly stabilizes and resumes normal operations after the recoveries.

5.4 Replication Performance of Auction and Persistent Cache

In this section, we present the performance of replication management for the Neptune-enabled auction and persistent cache service. Our test workload for auction is based on the published access trace from eBay between 29 May and 9 June 1999. Excluding the requests for embedded images, we estimate that about 10 percent of the requests are for bidding, and 3 percent are for adding new items. More information about this analysis can be found in our earlier study on dynamic Web caching [36]. We choose the number of auction categories to be 400 times the number of service nodes. Those categories are in turn divided into 64 partitions. Each request is made for an auction category selected from a population according to an even distribution. We choose level three consistency with staleness control at one second in this experiment. This consistency level fits the auction users' needs to acquire the latest information. Fig. 6a shows the performance of a Neptune-enabled auction service. In general, the results match the performance of the discussion group with 10 percent writes in Section 5.1. However, we do observe that the replication overhead is smaller for the auction service. The reason is that the trade off between the read load sharing and extra write overhead for service replication depends on the cost ratio between a read and a write. Most writes for the auction service are bidding requests which incur small overhead.

Fig. 6b illustrates the performance of the persistent cache service. Level two consistency and 10 percent write percentage are used in the experiment. The results show large replication overhead caused by extra writes. This is because CacheUpdate may cause costly disk accesses, while CacheLookup can usually be fulfilled with in-memory data.

5.5 Evaluation on Cluster Load Balancing

This section shows a performance evaluation on Neptune's cluster-wide load balancing support. We use read-only

services in this study to minimize the impact of replication management on the load balancing performance. This evaluation is based on two service traces from the Web search engine Teoma [29]. Both traces were collected across a one-week time span in late July 2001. One of the services provides the translation between query words and their internal representations. It has a mean service time of 22.2 ms and we call it the *Fine-Grain* trace. The other service supports a similar translation for Web page descriptions. It has a mean service time of 208.9 ms and we call it the *Medium-Grain* trace. In addition to the two traces, we also include a synthetic workload with Poisson process arrivals and exponentially distributed service times. We call this workload *Poisson/Exp*. The service processing for Poisson/Exp is emulated using a CPU-spinning microbenchmark that consumes the same amount of CPU time as the intended service time.

Fig. 7 shows the mean response time for using different load balancing policies. The random polling policies with

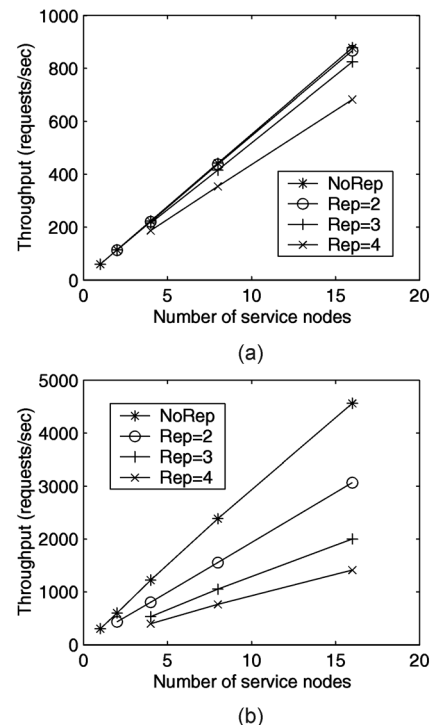


Fig. 6. Replication performance for auction and persistent cache. (a) Auction and (b) number of service nodes.

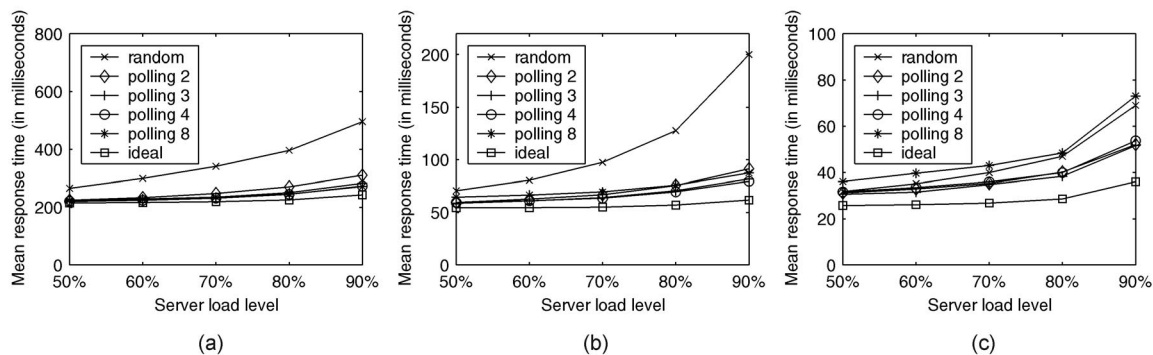


Fig. 7. Performance of cluster load balancing on a 16-server system. (a) Medium-grain trace, (b) Poisson/Exp with a mean service time of 50 ms, and (c) fine-grain trace.

the poll size of 2, 3, 4, and 8 are compared with the random and ideal approach in a 16-server system. The ideal performance in this study is measured through a centralized load index manager which keeps track of all server load indexes. In terms of server load level, we consider a server has reached full load (100 percent) when around 98 percent of client requests are successfully completed within two seconds. Then, we use this as the basis to calculate the client request rate for various server load levels. We observe that the results for Medium-Grain trace and Poisson/Exp workload largely confirm Mitzenmacher's analytical results [24]. For the Fine-Grain trace, however, we notice that a poll size of 8 exhibits far worse performance than policies with smaller poll sizes and it is even slightly worse than the pure random policy. This is caused by excessive polling overhead coming from two sources: 1) longer polling delays resulting from larger poll size and 2) less accurate server load information due to longer polling delay. And, those overheads are more severe for fine-grain services. Our conclusion is that a small poll size (e.g., 2 or 3) provides sufficient information for load balancing. And, an excessively large poll size may even degrade the performance due to polling overhead, especially for fine-grain services.

Table 1 shows the performance improvement of discarding slow-responding polls. The experiments are conducted with a poll size of 3 and a server load level of 90 percent. The result on the Medium-Grain trace shows a slight performance degradation due to the loss of load information. However, the results on both Fine-Grain trace and Poisson/Exp workload exhibit sizable improvement. Overall, the enhancement of discarding slow-responding polls can improve the load balancing performance by up to 8.3 percent. Note that the performance results shown in Fig. 7 are not with discarding slow-responding polls.

TABLE 1
Performance Gains of Discarding Slow-Responding Polls

Workload	Mean response time		Performance gain
	Original	Optimized	
Medium-Grain trace	282.1ms	283.1ms	-0.4%
Poisson/Exp	81.8ms	79.2ms	3.2%
Fine-Grain trace	51.6ms	47.3ms	8.3%

6 RELATED WORK

Our work is built upon a large body of previous research in network service clustering and replication management in distributed systems. We describe the work related to this paper in the following categories.

6.1 Software Infrastructure for Cluster-Based Network Services

Previous work has addressed the scalability and availability issues in providing clustering support for network service aggregation and load balancing [17], [32]. For instance, the TACC project employs a two-tier architecture in which the service components called "workers" run on different backends, while incoming requests to workers are controlled by front-ends [17]. Most of these systems employ two-tier clustering architectures and rely on centralized components to maintain the server runtime workload and service availability information. In comparison, Neptune employs a loosely connected and functionally symmetric architecture to achieve high scalability and robustness. This architecture allows the Neptune service infrastructure to operate smoothly in the presence of transient failures and service evolutions.

6.2 Replication Consistency

Replication of persistent data is crucial to load sharing and achieving high availability. The earlier analysis by Gray et al. shows that the synchronous replication based on eager update propagations leads to high deadlock rates [20]. A recent study by Anderson et al. confirms this using simulations [4]. The asynchronous replication based on lazy propagations has been used in Bayou [27]. Adya and Liskov have studied a type of lazy consistency in which server data is replicated in a client cache [2]. The serializability for lazy propagations with the primary-copy method is further studied by a few other research groups [4], [13], and their work addressed causal dependence for accessing multiple objects. The recent work by Yu and Vahdat provides a tunable framework to exploit the trade off among availability, consistency, and performance [33]. Neptune's replication consistency maintenance differs from the previous work in providing flexible replication consistency at the clustering middleware level. In particular, Neptune's multilevel replication consistency model supports data staleness control at its highest level.

Neptune's data staleness control is similar to the session guarantees in the Bayou project [30]. Bayou's session guarantee supports the following properties for operations belonging to the same client session: "read your writes," "monotonic reads," "writes follow reads," and "monotonic writes." The first two guarantees are covered in Neptune's "progressive version delivery." The other two properties are inherently supported by level two and three of Neptune's replication consistency model. Additionally, our "soft quantitative bound" enforces a soft bound on the data staleness in time quantities. Both "soft quantitative bound" and "progressive version delivery" are important for staleness-sensitive services like auction.

6.3 Replication Support for Network Services and Database Systems

The Ivory system provides automatic state replication for dynamic Web services through Java bytecode rewriting [10] and the Porcupine project developed a large-scale cluster-based e-mail service with replicated mailboxes [28]. Replication models for these systems are intended for services with only commutative updates for which ordered updates are unnecessary. Recently, the DDS project addressed replication of persistent data with a carefully built data management layer that encapsulates scalable replication consistency and fail-over support [22]. While this approach has been demonstrated for services with simple processing logic like distributed hash tables, constructing such a data management layer could be difficult for applications with complex data management logic.

Traditional databases and transaction processing systems enforce ACID consistencies in replication management through global synchronous locking [9], [21], [31]. Recent database systems from Oracle, Sybase, and IBM support lazy updates for data replication and they rely on user-specified rules to resolve conflicts. Neptune differs from those systems by supporting replication management in the clustering middleware level. The key advantage is that such a support can be easily applied to a large number of applications with different underlying data management mechanisms. Our work demonstrates this wide applicability can be achieved along with flexible replication consistency, performance scalability, and failure recovery support for cluster-based network services.

6.4 Clustering Load Balancing

A large body of work has been done to optimize HTTP request distribution among a cluster of Web servers [5], [11], [23], [26]. Most load balancing policies proposed in such a context rely on the premise that all network packets go through a single front-end dispatcher or a TCP-aware (layer 4 or above) switch so that TCP level connection-based statistics can be accurately maintained. However, clients and servers inside the service cluster are often connected by high-throughput, low-latency Ethernet (layer 2) or IP (layer 3) switches, which do not provide any TCP level traffic statistics. Our study in this paper shows that an optimized random polling policy that does not require centralized statistics can deliver good load balancing performance.

Previous research has proposed and evaluated various load balancing policies for distributed systems [7], [14], [24], [35]. Those studies are mostly designed for coarse-grain

distributed computation and they often ignore fine-grain jobs by simply processing them locally. Built on top of these previous results, our study finds that the random polling policy with a small poll size performs well for cluster-based network services. In addition, we find that the technique of discarding slow-responding polls can further improve the performance for fine-grain services.

7 CONCLUDING REMARKS

Building large-scale network services with the ever-increasing demand on scalability and availability is a challenging task. This paper investigates techniques in building the Neptune middleware system. Neptune provides clustering support and replication management for scalable network services. The system has been implemented on Linux and Solaris clusters where a number of applications have been successfully deployed.

Neptune provides multilevel replication consistency for cluster-based network services with performance scalability and fail-over support. This work demonstrates that a flexible replication management framework can be built at the clustering middleware level. In terms of service guarantees, our level three consistency ensures that client accesses are serviced progressively within a specified soft staleness bound, which is sufficient for many network services. Nevertheless, the current three consistency levels are not meant to be complete. Additional consistency levels can be incorporated into Neptune's replication consistency model when needed.

ACKNOWLEDGMENTS

This work was supported in part by US National Science Foundation grants ACIR-0082666, ACIR-0086061, ACIR-0234346, and CCR-0306473.

REFERENCES

- [1] The Alexandria Digital Library Project, <http://www.alexandria.ucsb.edu>, 2003.
- [2] A. Adya and B. Liskov, "Lazy Consistency Using Loosely Synchronized Clocks," *Proc. ACM Symp. Principles of Distributed Computing*, pp. 73-82, Aug. 1997.
- [3] D. Agrawal, A. El Abbadi, and R.C. Steinke, "Epidemic Algorithms in Replicated Databases," *Proc. 16th Symp. Principles of Database Systems*, pp. 161-172, May 1997.
- [4] T. Anderson, Y. Breitbart, H.F. Korth, and A. Wool, "Replication, Consistency, and Practicality: Are These Mutually Exclusive?" *Proc. ACM SIGMOD Int'l Conf. Management of Data*, pp. 484-495, June 1998.
- [5] M. Aron, D. Sanders, P. Druschel, and W. Zwaenepoel, "Scalable Content-Aware Request Distribution in Cluster-Based Network Services," *Proc. USENIX Ann. Technical Conf.*, June 2000.
- [6] Ask Jeeves Search, <http://www.ask.com>, 2003.
- [7] A. Barak, S. Guday, and R.G. Wheeler, *The MOSIX Distributed Operating System: Load Balancing for UNIX*, Springer-Verlag, 1993.
- [8] D.A. Benson, I. Karsch-Mizrachi, D.J. Lipman, J. Ostell, B.A. Rapp, and D.L. Wheeler, "GenBank," *Nucleic Acids Research*, vol. 30, no. 1, pp. 17-20, 2002.
- [9] P. Bernstein and E. Newcomer, *Principles of Transaction Processing*. Morgan Kaufmann, 1997.
- [10] G. Berry, J. Chase, G. Cohen, L. Cox, and A. Vahdat, "Toward Automatic State Management for Dynamic Web Services," *Proc. Network Storage Symp.*, Oct. 1999.

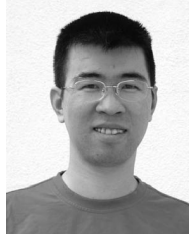
- [11] E.V. Carrera and R. Bianchini, "Efficiency vs. Portability in Cluster-Based Network Servers," *Proc. Eighth ACM Symp. Principles and Practice of Parallel Programming*, pp. 113-122, June 2001.
- [12] L. Chu, K. Shen, and T. Yang, "A Guide to Neptune: Clustering Middleware for Online Services," part of the Neptune software distribution, <http://www.cs.ucsb.edu/projects/neptune>, Apr. 2003.
- [13] P. Chundi, D.J. Rosenkrantz, and S.S. Ravi, "Deferred Updates and Data Placement in Distributed Databases," *Proc. 12th Int'l Conf. Data Eng.*, pp. 469-476, Feb. 1996.
- [14] D.L. Eager, E.D. Lazowska, and J. Zahorjan, "Adaptive Load Sharing in Homogeneous Distributed Systems," *IEEE Trans. Software Eng.*, vol. 12, no. 5, pp. 662-675, May 1986.
- [15] eBay Online Auctions, <http://www.ebay.com>, 2003.
- [16] D. Ferrari, "A Study of Load Indices for Load Balancing Schemes," Technical Report CSD-85-262, EECS Dept., Univ. of California Berkeley, Oct. 1985.
- [17] A. Fox, S.D. Gribble, Y. Chawathe, E.A. Brewer, and P. Gauthier, "Cluster-Based Scalable Network Services," *Proc. 16th ACM Symp. Operating System Principles*, pp. 78-91, Oct. 1997.
- [18] H. Garcia-Molina, "Elections in a Distributed Computing System," *IEEE Trans. Computers*, vol. 31, no. 1, pp. 48-59, Jan. 1982.
- [19] Google Search, <http://www.google.com>, 2003.
- [20] J. Gray, P. Helland, P. O'Neil, and D. Shasha, "The Dangers of Replication and a Solution," *Proc. ACM SIGMOD Int'l Conf. Management of Data*, pp. 173-182, June 1996.
- [21] J. Gray and A. Reuter, *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann, 1993.
- [22] S.D. Gribble, E.A. Brewer, J.M. Hellerstein, and D. Culler, "Scalable, Distributed Data Structures for Internet Service Construction" *Proc. Fourth USENIX Symp. Operating Systems Design and Implementation*, Oct. 2000.
- [23] G.D.H. Hunt, G.S. Goldszmidt, R.P. King, and R. Mukherjee, "Network Dispatcher: A Connection Router for Scalable Internet Services," *Proc. Seventh Int'l World Wide Web Conf.*, Apr. 1998.
- [24] M. Mitzenmacher, "On the Analysis of Randomized Load Balancing Schemes," *Proc. Ninth ACM Symp. Parallel Algorithms and Architectures*, pp. 292-301, June 1997.
- [25] MSN Groups Service, <http://groups.msn.com>, 2003.
- [26] V.S. Pai, M. Aron, G. Banga, M. Svendsen, P. Druschel, W. Zwaenepoel, and E. Nahum, "Locality-Aware Request Distribution in Cluster-Based Network Servers," *Proc. Eighth ACM Conf. Architectural Support for Programming Languages and Operating Systems*, pp. 205-216, Oct. 1998.
- [27] K. Petersen, M.J. Spreitzer, D.B. Terry, M.M. Theimer, and A.J. Demers, "Flexible Update Propagation for Weakly Consistent Replication," *Proc. 16th ACM Symp. Operating Systems Principles*, pp. 288-301, Oct. 1997.
- [28] Y. Saito, B.N. Bershad, and H.M. Levy, "Manageability, Availability, and Performance in Porcupine: A Highly Scalable, Cluster-Based Mail Service," *Proc. 17th ACM Symp. Operating Systems Principles*, pp. 1-15, Dec. 1999.
- [29] Teoma Search, <http://www.teoma.com>, 2003.
- [30] D.B. Terry, A.J. Demers, K. Petersen, M.J. Spreitzer, M.M. Theimer, and B.B. Welch, "Session Guarantees for Weakly Consistent Replicated Data," *Proc. Int'l Conf. Parallel and Distributed Information Systems*, pp. 140-149, Sept. 1994.
- [31] WebLogic and Tuxedo Transaction Application Server White Papers, <http://www.bea.com/products/tuxedo/papers.html>, 2003.
- [32] J. Robert von Behren, E.A. Brewer, N. Borisov, M. Chen, M. Welsh, J. MacDonald, J. Lau, S. Gribble, and D. Culler, "Ninja: A Framework for Network Services," *Proc. USENIX Ann. Technical Conf.*, June 2002.
- [33] H. Yu and A. Vahdat, "Design and Evaluation of a Continuous Consistency Model for Replicated Services," *Proc. Fourth USENIX Symp. Operating Systems Design and Implementation*, Oct. 2000.
- [34] S. Zhou, "An Experimental Assessment of Resource Queue Lengths as Load Indices," *Proc. Winter USENIX Technical Conf.*, pp. 73-82, Jan. 1987.
- [35] S. Zhou, "A Trace-Driven Simulation Study of Dynamic Load Balancing," *IEEE Trans. Software Eng.*, vol. 14, no. 9, pp. 1327-1341, Sept. 1988.
- [36] H. Zhu and T. Yang, "Class-Based Cache Management for Dynamic Web Contents," *Proc. IEEE INFOCOM*, Apr. 2001.



Kai Shen received the BS degree in computer science and engineering from Shanghai Jiaotong University, China, in 1996, and the PhD degree in computer science from the University of California, Santa Barbara, in 2002. He is currently an assistant professor in the Department of Computer Science at the University of Rochester, Rochester, New York. His research interests are in the areas of parallel and distributed systems, computer networks, and Web searching.



Tao Yang received the BS degree in computer science from Zhejiang University, China, in 1984, and the MS and PhD degrees in computer science from Rutgers University, in 1990 and 1993, respectively. He is an associate professor at the Department of Computer Science, University of California, Santa Barbara. His main research interests have been parallel and distributed systems, high performance scientific computing, cluster-based network services, and Internet search. He has published more than 70 refereed conference papers and journal articles. Dr. Yang has also been the chief scientist and vice president of Research and Development at Teoma Technologies from 2000 to 2001 and the chief scientist at Ask Jeeves since 2001, where he has directed research of large-scale Internet search and led the development of the Teoma search engine running on giant SMP clusters. Dr. Yang was an associate editor for *IEEE Transactions on Parallel and Distributed Systems* from 1999 to 2002, and served on the program committee of many high performance computing conferences. He received the Research Initiation Award from the US National Science Foundation in 1994, the Computer Science Faculty Teacher Award in 1995 from the UCSB College of Engineering, the CAREER Award from the US National Science Foundation in 1997, and the Noble Jeeviant Award from AskJeeves in 2002. He is a member of the IEEE.



Lingkun Chu received the BS and MS degrees in computer science from Zhejiang University, China, in 1996 and 1999, respectively. He is a PhD candidate in the Department of Computer Science, University of California, Santa Barbara. He is currently working with Professor Tao Yang in the area of parallel and distributed systems, with a particular focus on cluster-based network services.

► For more information on this or any other computing topic, please visit our Digital Library at <http://computer.org/publications/dlib>.