



# STICS: SCSI-to-IP cache for storage area networks

Xubin He<sup>a,\*</sup>, Ming Zhang<sup>b</sup>, Qing (Ken) Yang<sup>b</sup>

<sup>a</sup>Electrical and Computer Engineering, Tennessee Technological University, Cookeville, TN 38505, USA

<sup>b</sup>Electrical and Computer Engineering, University of Rhode Island, Kingston, RI 02881, USA

Received 31 July 2003; received in revised form 9 April 2004

## Abstract

Data storage plays an essential role in today's fast-growing data-intensive network services. New standards and products emerge very rapidly for networked data storage. Given the mature Internet infrastructure, the overwhelming preference among the IT community recently is using IP for storage networking because of economy and convenience. iSCSI is one of the most recent standards that allow SCSI protocols to be carried out over IP networks. However, there are many disparities between SCSI and IP in terms of protocols, speeds, bandwidths, data unit sizes, and design considerations that prevent fast and efficient deployment of storage area network (SAN) over IP. This paper introduces SCSI-to-IP cache storage (STICS), a novel storage architecture that couples reliable and high-speed data caching with low-overhead conversion between SCSI and IP protocols. A STICS block consists of one or several storage devices and an intelligent processing unit with CPU and RAM. The storage devices are used to cache and store data while the intelligent processing unit carries out the caching algorithm, protocol conversion, and self-management functions. Through the efficient caching algorithm and localization of certain unnecessary protocol overheads, STICS can significantly improve performance, reliability, and scalability over current iSCSI systems. Furthermore, STICS can be used as a basic plug-and-play building block for data storage over IP. Analogous to "cache memory" invented several decades ago for bridging the speed gap between CPU and memory, STICS is the first-ever "cache storage" for bridging the gap between SCSI and IP making it possible to build an efficient SAN over IP. We have implemented software STICS prototype on Linux operating system. Numerical results using popular benchmarks such as vxbench, IOzone, PostMark, and EMC's trace have shown a dramatic performance gain over the current iSCSI implementation.

© 2004 Elsevier Inc. All rights reserved.

**Keywords:** Cache; Networked storage; NAS; SAN; iSCSI; Performance evaluation.

## 1. Introduction

As we enter a new era of computing, data storage has changed its role from secondary with respect to CPU and RAM to primary importance in today's information world [13]. Online data storage doubles every 9 months [7] due to an ever-growing demand for networked information services [8,25,51]. In general, networked storage architectures have evolved from network-attached storage (NAS) [11,17,35,37], storage area network (SAN) [23,39,42], to most recent storage over IP (IP SAN) [17,44]. NAS

architecture allows a storage system/device to be directly connected to a standard network, typically via Ethernet. Clients in the network can access the NAS directly. A NAS-based storage subsystem has built-in file system to provide clients with file system functionality. SAN technology, on the other hand, provides a simple block level interface for manipulating nonvolatile magnetic media. Typically, a SAN consists of networked storage devices interconnected through a dedicated fibre channel (FC-4 protocol) network. The basic premise of a SAN is to replace the "point-to-point" infrastructure of server-to-storage communications with one that allows "any-to-any" communications. A SAN provides high connectivity, scalability, and availability using a specialized network protocol: FC-4 protocol. Deploying such a specialized network usually introduces additional cost

\* Corresponding author

E-mail addresses: [hexb@tntech.edu](mailto:hexb@tntech.edu) (X. He), [mingz@ele.uri.edu](mailto:mingz@ele.uri.edu) (M. Zhang), [qyang@ele.uri.edu](mailto:qyang@ele.uri.edu) (Q. (Ken) Yang).

for implementation, maintenance, and management. Internet SCSI (iSCSI) [4,20,30,45,49] is the most recently emerging technology with the goal of implementing the IP SAN.

Compared to FC-4, implementing SAN over IP (IP SAN) has several advantages [34,52]:

- IP SAN can run over standard off-the-shelf network components, such as switched Ethernet, which reduces the cost. One can extend and expand the switched network easily and quickly while riding the cost/performance improvement trends of Ethernet.
- IP SAN can exploit existing IP-based protocols, and IP SANs using iSCSI can be managed using existing and familiar IP-based tools such as SNMP, while Fibre Channel SANs require specialized management infrastructure.
- A network that incorporates IP SANs need use only a single kind of network infrastructure (Ethernet) for both data and storage traffic, whereas use of fibre channel protocol (FCP) requires a separate kind of infrastructure (fibre channel) for storage.

IP SAN brings economy and convenience whereas it also raises performance issues, which is the main downside of current IP SAN as compared to FC-SAN. Currently, there are basically two existing approaches to implement IP SAN using iSCSI: one carries out SCSI and IP protocol conversion at a specialized switch [39] and the other encapsulates SCSI protocol in TCP/IP at the host bus adapter (HBA) level [45]. Both approaches have severe performance limitations. Converting protocols at a switch places an additional special burden on an already-overloaded switch and requires specialized networking equipment in a SAN. Such a specialized switch is not only costly, as compared to off-the-shelf Ethernet switches, but also complicates installation, management, and maintenance. To encapsulate SCSI protocol over IP requires significant amount of overhead traffic for SCSI commands transfers and handshaking over the Internet. On a typical iSCSI implementation, we have measured around 58% of TCP/IP packets being less than 127 bytes long, implying an overwhelming quantity of small packets transferring SCSI commands and status (most of them are only one byte). A majority of such small packet traffic over the net is not necessary because of the reliable and connection-oriented services provided by underlying TCP/IP. Our experiments using the PostMark benchmark [22] have shown that efficient caching can reduce the total number of packets transferred over the network from 3,353,821 to 839,100 for same amount of remote storage data, a 75 percent reduction!

In addition to the above-mentioned protocol disparities between SCSI and IP, packet transfer latency exists over the network, particularly over long distances. Such latency does not decrease linearly with an increase in network bandwidth. For example, we measured average network latencies over 100Mbit and 1Gbit Ethernet switches to be 128.99 and 106.78  $\mu$ s, respectively. These results indicate that even though the bandwidth of Ethernet switches has increased to

gigabit or tens of gigabits, network latencies resulting from packet propagation delays are still there.

Protocol disparities and network latencies motivate us to introduce a new storage architecture: SCSI-to-IP cache storage (STICS). The purpose of STICS is to bridge the disparities between SCSI and IP so that an efficient SAN can be built over the Internet. A typical STICS block consists of a disk and an intelligent processing unit with an embedded processor and sufficient RAM. It has two standard interfaces: a SCSI interface and a standard Ethernet interface. The disk is used as a nonvolatile cache that caches data coming from possibly two directions: block data from the SCSI interface and network data from the Ethernet interface. In addition to standard SCSI and IP protocols running on the intelligent processing unit, it also implements a special caching algorithm controlling a two-level cache hierarchy that writes data very quickly. Besides caching storage data, STICS also localizes SCSI commands and handshaking operations to reduce unnecessary traffic over the Internet. In this way, it acts as a storage filter to discard a fraction of the data that would otherwise move across the Internet, reducing the bottleneck problem imposed by limited Internet bandwidth and increasing storage data transfer rate. Apparent advantages of the STICS are:

- It provides an iSCSI network cache to smooth out the traffic and improve overall performance. Such a cache or bridge is not only helpful but also necessary to a certain degree because of the different nature of SCSI and IP, such as speed, data unit size, protocols, and requirements. Wherever there is a speed disparity, cache helps. Analogous to “cache memory” used to cache memory data for a CPU [36], STICS is a “cache storage” used to cache networked storage data for a server host.
- It utilizes the techniques in a Log-structured file system [46,53] to quickly write data into magnetic media for caching data coming from both directions. A disk is used in caching, which is extremely important for caching data reliably since once data is written to a nonvolatile storage, it is considered to be safe.
- By localizing part of SCSI protocol and filtering out some unnecessary traffic, STICS can reduce the bandwidth required to implement a SAN.
- Active disks [1,29,43] are becoming feasible and popular. STICS represents another specific and practical implementation of active disks.
- It is a standard plug-and-play building block for SAN over the Internet. If ISTORE [7] is standard “brick” for building storage systems, then STICS can be considered as a standard “beam” or “post” that provides interconnect and support for the construction of SANs.

Overall, STICS adds a new dimension to the networked storage architectures. To quantitatively evaluate the performance potential of STICS in a real network environment, we have implemented the STICS under the Linux OS over

an Ethernet switch. We have used popular benchmark programs, such as PostMark [22], IOzone [40], vxbench,<sup>1</sup> and EMC's trace to measure system performance. Benchmark results show that STICS provides up to 4.3 times performance improvement over iSCSI implementation in terms of average system throughput. For EMC's trace measurement, STICS can be 6 times as fast as iSCSI in terms of average response time.

The paper is organized as follows. Section 2 summarizes related work. Section 3 presents the STICS architecture, followed by detailed descriptions of the design and implementation in Section 4. Section 5 presents our performance evaluation methodology and numerical results. Finally, Section 6 concludes the paper.

## 2. Related work

Existing research that is most closely related to STICS is NAS [11-12,43]. The NAS technology provides direct network connection for hosts to access through network interfaces. It also provides file system functionality. NAS-based storage appliances range from terabyte servers to a simple disk with an Ethernet plug. The main difference between NAS and SAN is that NAS provides storage at the file system level while SAN provides storage at the block device level. Another difference is that NAS is attached to the same LAN as the one connecting servers accessing storage, while SAN has a dedicated network connecting storage devices without competing for network bandwidth with the servers. STICS provides a direct SCSI connection to a server host to allow the server to access a SAN implemented over the Internet at the block level. In addition to being a storage component of the SAN, STICS performs network cache functions for a smooth and efficient SAN implementation over IP network.

Another important related effort is *Petal* [27,48], a research project of Compaq's Systems Research Center. Petal uses a collection of NAS-like storage servers interconnected using a specially customized LAN to form a unified virtual disk space to clients at the block level. iSCSI [17,20,45], which emerged very recently, provides an ideal alternative to Petal's customized LAN-based SAN protocol. Taking advantage of existing Internet protocols and media, it is a natural way for storage to make use of TCP/IP, as demonstrated by the earlier research work of Meter et al., VISA [33] to transfer SCSI commands and data using IP protocol. iSCSI protocol is a mapping of the SCSI remote procedure invocation model over the TCP/IP protocol [45]. The STICS architecture attempts to localize some of SCSI protocol traffic by accepting SCSI commands and data from the host and filtering data block to be sent to the remote storage target. This SCSI-in block-out mechanism provides an immediate and transparent solution, both to the host and the storage,

eliminating some unnecessary remote synchronization. Furthermore, STICS provides a nonvolatile cache exclusively for SCSI commands and data that are supposed to be transferred through the network. This cache reduces latency from the host's point of view and avoids many unnecessary data transfers over the network, because many data are frequently overwritten.

The idea of using a disk-based log to improve system performance or to improve the reliability of RAM has been used in both file system and database systems for a long time. For example, the Log-structured File System (LFS [46,53]), Disk Caching Disk (DCD [18]), and other similar systems all use disk-based data/metadata logging to improve file system performance and speed-up crash recovery. Several RAID systems have implemented the LFS algorithm at the RAID controller level [19,32,54]. LFS collects writes in a RAM buffer to form large logs and writes large logs to data disks. While many implementation techniques are borrowed from existing work, the novelty of STICS is the new concept of caching between SCSI and IP.

## 3. Architecture

Essentially, STICS is a cache that bridges the protocol and speed disparities between SCSI and IP. Fig. 1 shows a typical SAN implementation over IP using STICS. Any number of storage devices or server computers can be connected to the standard Internet through STICS to form a SAN. Instead of using a specialized network or specialized switch, STICS connects a regular host server or a storage device to the standard IP network using off-the-shelf components. Consider STICS 1 in the diagram. It is directly connected to the SCSI HBA of Host 1 as a local storage device. It also acts as a cache and bridge to allow Host 1 to access, at the block level, any storage device connected to the SAN such as NAS, STICS 2, STICS 3, etc. In order to allow a smooth data access from Host 1 to the SAN, STICS 1

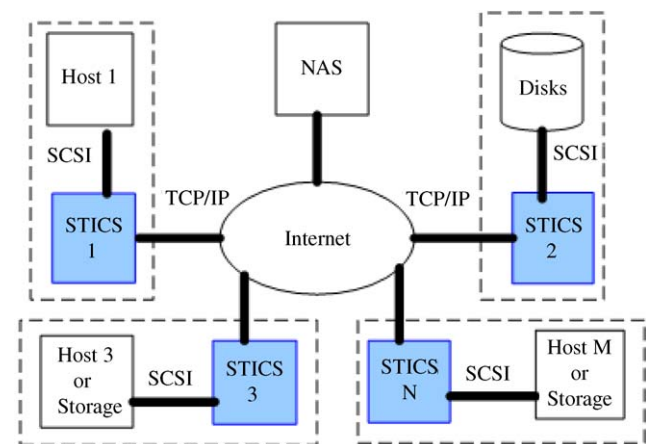


Fig. 1. System overview. A STICS connects to the host via SCSI interface and connects to other STICS' or NAS via Internet.

<sup>1</sup> An I/O benchmark developed by VERITAS Corp.

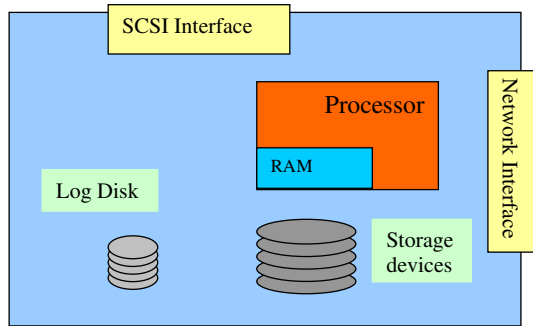


Fig. 2. STICS architecture.

provides SCSI service, caching service, naming service, and IP protocol service.

Fig. 2 shows the basic structure of STICS, which consists of five main components:

1. A *SCSI interface*: STICS supports SCSI communications with hosts and other extended storage devices. Via the SCSI interface, STICS may run under two different modes: initiator mode or target mode [55]. When a STICS is used to connect to a host, it runs in target mode receiving requests from the host, carrying out the I/O processing possibly through network, and sending back results to the host. In this case, the STICS acts as a directly attached storage device from the host's point of view. When a STICS is used to connect to a storage device such as a disk or RAID to extend storage, it runs in initiator mode, and it sends or forwards SCSI requests to the extended storage device. For example, in Fig. 1, STICS 1 runs in target mode while STICS 2 runs in initiator mode.
2. An *Ethernet interface*: Via the network interface, a STICS can be connected to the Internet and share storage with other STICS's or NAS.
3. An *intelligent processing unit*: This processing unit has an embedded processor and an amount of RAM. A specialized Log-structured file system, standard SCSI protocols, and IP protocols run on the processing unit. The RAM consists of a regular DRAM for read caching and a small (1–4 MB) nonvolatile RAM (NVRAM) for write caching. The NVRAM is also used to maintain the meta data such as hash table, LRU list, and the mapping information (STICS\_MAP). Alternatively, we can also use Soft Updates [10] technique to keep meta data consistency without using NVRAM.
4. A *log disk*: The log disk is a sequentially accessed device. It is used to cache write data along with the NVRAM above in the processing unit. The log disk and the NVRAM form a two-level hierarchical cache.
5. *Storage device*: The regular storage device is an optional component. It can be a disk, a RAID, or Just-Bunch-Of-Disks (JBOD). It is used as an extra and backup storage capacity. For example, it can be used as temporary offline

storage when the network fails. From the point of view of a server host to which the STICS is connected through the SCSI interface, this storage device can be considered as a local disk. From the point of view of the IP network through the network interface, it can be considered as a component of a SAN with an IP address as its ID.

## 4. Design and implementation

### 4.1. STICS naming service

To allow a true “any-to-any” communication between servers and storage devices, a global naming is necessary. In our design, each STICS is named by a global location number (GLN) which is unique for each STICS. Currently we assign an IP address to each STICS and use this IP as the GLN.

### 4.2. Cache structure of STICS

Each STICS has a read cache consisting of a large DRAM and a write cache consisting of a two-level hierarchy with a small NVRAM on top of a log disk. Frequently accessed data reside in the DRAM that is organized as LRU cache [21] for read operations. Write data are first stored in the small NVRAM. Whenever the newly written data in the NVRAM are sufficiently large or whenever the log disk is free, a log of data is written into the log disk sequentially. After the log write, the NVRAM becomes available to absorb additional write data. At the same time, a copy of the log is placed in the DRAM to speed up possible read operations of the data that have just been written to the log disk. Data in the log disk are organized in the format of *segments* similar to that in a LFS [46]. A segment contains a number of *slots* each of which can hold one data block. Data blocks in a segment are addressed by their *Segment IDs* and *Slot IDs*.

Fig. 3 shows the data structure in both DRAM and NVRAM. A Hash table is used to locate data in the RAM buffer including DRAM and NVRAM. DRAM and NVRAM can be differentiated through their addresses. A LRU list and a Free list are used to keep tracks of the most recently used data and the free slots, respectively.

Data blocks stored in the RAM buffer are addressed by their *Logical Block Addresses (LBAs)*. The Hash table contains location information for each of the valid data blocks in the buffer and uses LBAs of incoming requests as search keys. The slot size is set to be the size of a block. A slot entry consists of the following fields:

- An LBA of a cache line. It serves as the search key of hash table;
- Global Location Number (GLN) if the slot contains data from or to other STICS.



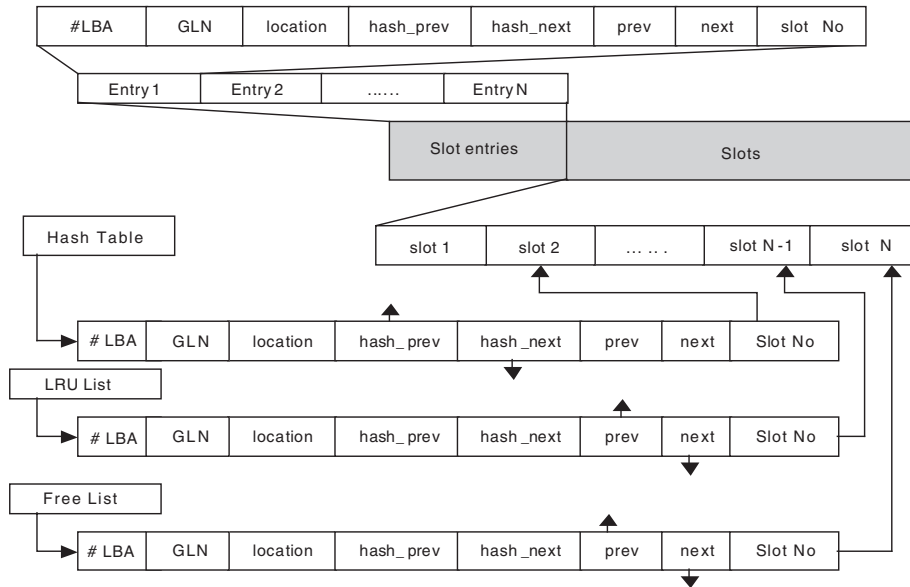


Fig. 3. RAM buffer layout. RAM buffer consists of slot entries and slots. The hash table, LRU list and Free list are used to organize the slot entries.

- A *location* field containing 2 parts:

- (1) a state tag (2 bits), used to specify where the slot data is: IN\_RAM\_BUFFER, IN\_LOG\_DISK, IN\_DATA\_DISK or IN\_OTHER\_STICS;
- (2) a log disk block index (30 bits), used to specify the log disk block number if the state tag indicates IN\_LOG\_DISK. The size of each log disk can be up to  $2^{30}$  blocks.

- Two pointers (*hash\_prev* and *hash\_next*) used to link the hash table.
- Two pointers (*prev* and *next*) used to link the LRU list and FREE list.
- A *Slot-No* used to describe the in-memory location of the cached data.

### 4.3. Basic operations

#### 4.3.1. Write

Write requests may come from one of two sources: the host via the SCSI interface and another STICS via the Ethernet interface.

*Write requests from the host via SCSI interface:* After receiving a write request, the STICS first searches the Hash table by the LBA address. If an entry is found, the entry is overwritten by the incoming write and is moved to the NVRAM if it is in DRAM. If no entry is found, a free slot entry in the NVRAM is allocated from the Free list, the data are copied into the corresponding slot, and its address is recorded in the Hash table. The LRU list and Free list are then updated. When enough data slots (128 in our preliminary implementation) are accumulated or when the log disk

is idle, the data slots are written into log disk sequentially in one large write. After the log write completes successfully, STICS signals the host that the request is complete and the log is moved from the NVRAM to DRAM.

*Write requests from another STICS via Ethernet interface:* A packet coming from the network interface may turn out to be a write operation from a remote STICS on the network. After receiving such a write request, STICS gets a data block with GLN and LBA. It then searches the Hash table by the LBA and GLN. The same writing process as above is then performed.

#### 4.3.2. Read

Similar to write operations, read operations may also come either from the host via the SCSI interface or from another STICS via the Ethernet interface.

*Read requests from the host via SCSI interface:* After receiving a read request, the STICS searches the Hash table by the LBA to determine the location of the data. Data requested may be in one of four different places: the RAM buffer, the log disk(s), the storage device in the local STICS, or a storage device in another STICS on the network. If the data are found in the RAM buffer, the data are copied from the RAM buffer to the requesting buffer. The STICS then signals the host that the request is complete. If the data are found in the log disk or the local storage device, the data are read from the log disk or storage device into the requesting buffer. Otherwise, the STICS encapsulates the request including LBA, current GLN, and destination GLN into an IP packet and forwards it to the corresponding STICS.

*Read requests from another STICS via Ethernet interface:* When a read request is found after unpacking an incoming IP packet, the STICS obtains the GLN and LBA from the

packet. It then searches the Hash table by the LBA and the source GLN to determine the location of the data. It locates and reads data from that location. Finally, it sends the data back to the source STICS through the network.

#### 4.3.3. Destages

The operation of moving data from a higher-level storage device to a lower level storage device is defined as *destage* operation [50]. Two levels of destage operations are defined in STICS: destaging data from the NVRAM buffer to the log disk (*Level 1 destage*) and destaging data from the log disk to a storage device (*Level 2 destage*). We implement a separate kernel thread, *LogDestage*, to perform the destaging tasks. The *LogDestage* thread is registered during system initialization and monitors the *STICS* states. *Level 1 destage* activates whenever the log disk is idle and there are data to be destaged in the NVRAM. *Level 2 destage* activates whenever one of the following events occurs: (1) the *STICS* detects a CPU idle period; (2) the size of data in the log disk exceeds a threshold value. *Level 1 destage* has higher priority than *Level 2 destage*. Once the *Level 1 destage* starts, it continues until a log of data in the NVRAM buffer is written to the log disk. *Level 2 destage* may be interrupted if a new request comes in or until the log disk becomes empty. If the destage process is interrupted, the destage thread would be suspended until the *STICS* detects another idle period. For extreme burst writes, where the log disk is full, *Level 1 destage* forces subsequent writes to the addressed network storage to bypass the log disk to avoid cache overflow [50].

As for *Level 1 destage*, the data in the NVRAM buffer are written to the log disk sequentially in large size (64 KB). At the same time, the data are moved from NVRAM to DRAM. The log disk header and the corresponding in-memory slot entries are updated. All data are written to the log disk in “append” mode, which ensures that every time the data are written to consecutive log disk blocks.

For *Level 2 destage*, we use a “first-write-first-destage” algorithm according to the LRU List. Currently, we are using the LRU replacement algorithm, and other algorithms [56] are in consideration for the future implementation. Each time 64 KB data are read from the consecutive blocks of the log disk and written to the addressed network storage. The LRU list and free list are updated subsequently.

#### 4.4. Cache coherence

There are three ways to configure a distributed storage system using *STICS*, placing *STICS* near the host, target storage, or both. If we place a *STICS* near the host, the corresponding *STICS* building block is a private cache. If we place a *STICS* near the storage, we have a shared cache system. There are tradeoffs between shared cache and private cache configurations. From the point of view of cache efficiency, we would like to place cache as close to a host as possible to minimize latency [38]. Such a private cache sys-

tem allows multiple copies of a shared storage data to reside in different caches giving rise to the well-known cache coherence problem [2,3,6,9,16,24,26,28,41]. Shared caches, on the other hand, do not have such cache coherence problem because each cache is associated with target storage. However, each request has to go through the network to obtain data at the target storage side. We have considered both private and shared cache configurations. Shared cache configuration is relatively simple. For private cache configuration, a coherence protocol is necessary. One possible way to implement a cache coherence protocol in private cache system is using the local consistency (LC) model [3], which helps to minimize meta-data network traffic pertaining to coherence protocol. A shared-read/exclusive-write lock (token) can be used to implement the necessary synchronization [5,47]. The details of the cache coherence protocol are out of scope of this paper. Interested readers are referred to [14].

#### 4.5. Implementation

There are several ways to implement *STICS*. A software *STICS* is a device driver or kernel module that controls and coordinates SCSI HBA and network interface card (NIC). It uses a part of host’s system RAM and part of disk to form the cache. *STICS* can also be implemented at HBA controller level as a *STICS* card. Such a card has sufficient intelligence with RAM, IDE or SCSI interface, and Ethernet interface. The IDE or SCSI interface is used to connect to a log disk for caching. Finally, *STICS* can be implemented as a complete cache box with built-in controller, log disks, and local storage.

Currently, we have implemented a software prototype of *STICS* on Linux kernel 2.4.2, and it is implemented as kernel module that can be loaded and unloaded dynamically. We simulate NVRAM using part of system RAM. This part of system RAM is reserved when the system is boot up, and it is not accessible to other applications. So this part of RAM is “immune” to application-level software crash and more reliable than regular RAM. The log disk is an additional IDE hard disk for caching function. There is no local storage and all I/O operations are remote operations going through the network.

### 5. Performance evaluations

#### 5.1. Methodology

For the purpose of performance evaluation, we have implemented a *STICS* prototype and deployed a software iSCSI. For a fair performance comparison, both iSCSI and *STICS* have exactly the same CPU and RAM size. This RAM includes read cache and write buffer used in *STICS*. All I/O operations in both iSCSI and *STICS* are forced to be remote operations to target disks through a switch.

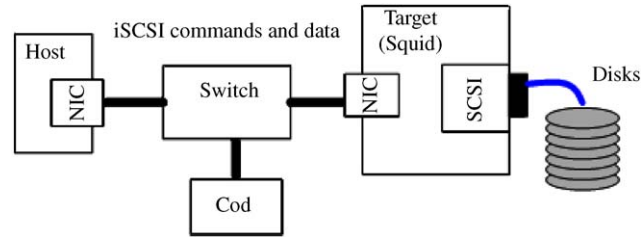


Fig. 4. iSCSI configuration. The host Trout establishes connection to target, and the target Squid responds and connects. Then the Squid exports hard drive and Trout sees the disks as local.

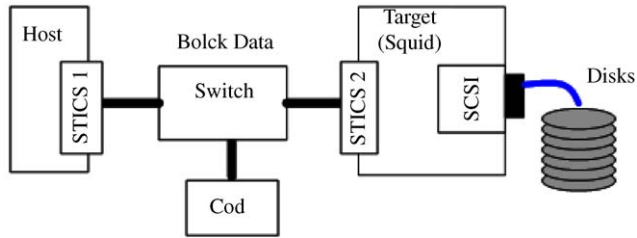


Fig. 5. STICS configuration. The STICS cache data from both SCSI and network.

Our experimental settings are shown in Figs. 4 and 5. Three PCs and a SUN workstation are involved in our experiments, namely *Trout*, *Cod*, *Squid*, and *Clam*. *Trout* and *Clam* serves as the host and *Squid* as the storage target. *Cod* serves as a switch console to monitor the network traffic. For STICS experiment, a software STICS is loaded as kernel module. To focus on the STICS performance measurement, we disabled STICS2's cache function on the target machine. All cache function is performed by STICS1 on the initiator side. STICS2 only receives and forwards I/O requests in our experiments. All these machines are interconnected through an 8-port Gigabit switch (Intel NetStructure 470 T) to form an isolated LAN. Each machine is running Linux kernel 2.4.2 with a Netgear GA622 T Gigabit NIC and an Adaptec 39160 high-performance SCSI adaptor. The network cards and switch can be tuned to gigabit and 100Mbit dynamically. The configurations of these machines are described in Table 1 and the characteristics of individual disks are summarized in Table 2.

For iSCSI implementation, we compiled and run the Linux iSCSI developed by Intel Corporation [20]. The iSCSI is compiled under Linux kernel 2.4.2 and configured as shown in Fig. 4. There are 4 steps for the two machines to establish communications via iSCSI. First, the host establishes connection to target; second, the target responds and connects; third, the target machine exports its disks and finally the host sees these disks as local. All these steps are finished through socket communications. After these steps, the iSCSI is in "full feature phase" mode where SCSI commands and data can be exchanged between the host and the target. For each SCSI operation, there will be at least 4 socket communications as follows: (1) the host encapsulates the SCSI command into packet data unit (PDU) and

sends this PDU to the target; (2) the target receives and decapsulates the PDU. It then encapsulates a response into a PDU and sends it back to the host; (3) the host receives and decapsulates the response PDU. It then encapsulates the data into a PDU and sends it to the target if the target is ready to transfer; (4) the target receives the data PDU and sends another response to the host to acknowledge the finish of the SCSI operation. iSCSI supports both solicited and unsolicited writes, but we found in current iSCSI implementation, the performance difference between solicited and unsolicited writes are less than 10%, which is very small compared to our STICS performance gain, so in our experiments, we configured iSCSI with solicited writes which is the default setting.

Our STICS runs on Linux kernel 2.4.2 with target mode support and is loaded as a kernel module as shown in Fig. 5. Four MB of the system RAM is used to simulate STICS NVRAM buffer, another 16 MB of the system RAM is used as the DRAM read cache in our STICS, and the log disk is a standalone hard drive. When requests come from the host, the STICS first processes the requests locally. For write requests, the STICS writes the data to its write buffer. Whenever the log disk is idle, the data will be destaged to the log disk through level 1 destage. After data are written to the log disk, STICS signals host write complete and moves the data to DRAM cache. When data in the log disk exceeds a threshold or the system is idle, the data in log disk will be destaged to the remote target storage through the network. The hash table and LRU list are updated. When a read request comes in, the STICS searches the hash table, locates where the data are, and accesses the data from RAM buffer, log disk, or remote disks via network.

In our previous discussions, all STICS are configured in "report after complete" mode. This scheme has a good reliability because a write is guaranteed to be stored in a disk before the CPU is acknowledged. If the 4-MB RAM buffer is nonvolatile, "immediate report" mode can be used, where as soon as the data are transferred to the RAM buffer, STICS sends an acknowledgement of "write complete" to the host.

## 5.2. Benchmark programs and workload characteristics

It is important to use realistic workloads to drive our STICS for a fair performance evaluation and comparison.

Table 1  
Machines configurations

|       | Processor          | RAM (MB) | IDE disk         | SCSI disk   |
|-------|--------------------|----------|------------------|-------------|
| Trout | PII-450            | 128      | 2 Maxtor ASO10a1 | N/A         |
| Cod   | PII-400            | 128      | Maxtor ASO10a1   | N/A         |
| Squid | PII-400            | 128      | 2 Maxtor ASO10a1 | IBM 07N3200 |
| Clam  | Ultra SPARC II 450 | 256      | N/A              | SUN18G      |

Table 2  
Disk parameters

| Disk model | Interface     | Capacity (G) | Data buffer | RPM    | Latency (ms) | Transfer rate (MB/s) | Seek time (ms) | Manufacturer |
|------------|---------------|--------------|-------------|--------|--------------|----------------------|----------------|--------------|
| O7N3       | Ultra SCSI    | 36.7         | N/A         | 10,000 | 3.0          | 29.8                 | 4.9            | IBM          |
| AS010a1    | Ultra ATA/100 | 10.2         | 2MB         | 7200   | 4.17         | 16.6                 | 8.5            | Maxtor       |
| SUN18G     | Ultra SCSI    | 18.2         | 512KB       | 10,000 | 3.0          | 18                   | 8.5            | SUN          |

1 For this reason, we chose to use 3 popular benchmark programs and a real-world trace.

3 The benchmarks we used to measure system throughput  
 5 are PostMark [22] which is a popular file system bench-  
 7 mark developed by Network Appliance, IOzone [40], and  
 9 vxbench developed by VERITAS. PostMark measures per-  
 11 formance in terms of transaction rates in an ephemeral small-  
 13 file environment by creating a large pool of continually  
 15 changing files. “PostMark was created to simulate heavy  
 17 small-file system loads with a minimal amount of soft-  
 19 ware and configuration effort and to provide complete re-  
 21 producibility [22].” PostMark generates an initial pool of  
 23 random text files ranging in size from a configurable low  
 25 bound to a configurable high bound. This file pool is of  
 27 configurable size and can be located on any accessible file  
 system. Once the pool has been created, a specified number  
 of transactions occur. Each transaction consists of a pair of  
 smaller transactions, i.e. *Create file or Delete file*, and *Read  
 file or Append file*. Each transaction type and its affected  
 files are chosen randomly. The read and write block size can  
 be tuned. On completion of each run, a report is generated  
 showing some metrics such as elapsed time, transaction rate,  
 total number of files created and so on. IOzone and vxbench  
 generate and measure a variety of file operations based on a  
 big file. Vxbench can only run on Solaris operating system.  
 We used these two benchmarks to measure the I/O through-  
 put in terms of KB/Sec.

29 In addition to benchmark programs, we also used a real-  
 world trace obtained from EMC Corporation. The trace, re-  
 31 ferred to as *EMC-tel* trace hereafter, was collected by an  
 EMC Symmetrix system installed at a telecommunication  
 33 consumer site. The trace file contains 230,370 requests, with  
 a fixed request size of 4 blocks. The whole dataset size is  
 35 900 M bytes. The trace is write-dominated with a write ratio  
 of 89%. The average request rate is about 333 requests/s. In  
 37 order for the trace to be read by our STICS and the iSCSI  
 implementation, we developed a program called *ReqGener-*  
*ator* to convert the traces to high-level I/O requests. These

requests are then fed to our STICS and iSCSI system to 39  
 measure performance.

### 5.3. Measured results and discussions 41

#### 5.3.1. Postmark results

43 Our first experiment is to use PostMark to measure the  
 I/O throughput in terms of transactions per second. In our  
 45 tests, PostMark was configured in two different ways. First,  
 a small pool of 1000 initial files and 50,000 transactions;  
 47 and second a large pool of 20,000 initial files and 100,000  
 transactions. The total sizes of accessed data are 436 MB  
 (151.05 MB read and 285.08 MB write) and 740 MB  
 (303.46 MB read and 436.18 MB write), respectively. They  
 49 are much larger than host system RAM (128 MB). We left  
 all other PostMark parameters at their default settings. The  
 network is configured as a 100 MB network. 51

53 In Fig. 6, we plotted two separate bar graphs corre-  
 sponding to the small file pool case and the large one, 55  
 respectively. Each group of bars represents the system  
 throughputs of STICS with report after complete (STICS: 57  
 light blue bars), iSCSI (iSCSI: dark red bars) and STICS  
 with immediate report (STICS-Imm: light yellow bars) for 59  
 a specific data block size. It is clear from this figure that  
 STICS shows obvious better system throughput than the 61  
 iSCSI. The performance improvement of STICS over iSCSI  
 is consistent across different block sizes and for both small 63  
 pool and large pool cases. The performance gains of STICS  
 with report after complete over iSCSI range from 60% to 65  
 110%. STICS with immediate report outperforms iSCSI by  
 a factor of 2.69–4.18. 67

69 To understand why STICS provides such impressive per-  
 formance gains over the iSCSI, we monitored the network  
 activities at the Ethernet Switch through the console machine  
 71 *cod* for both STICS and iSCSI implementations. While both  
 implementations write all data from the host to the remote  
 storage, STICS transfers dramatically less packets over the 73



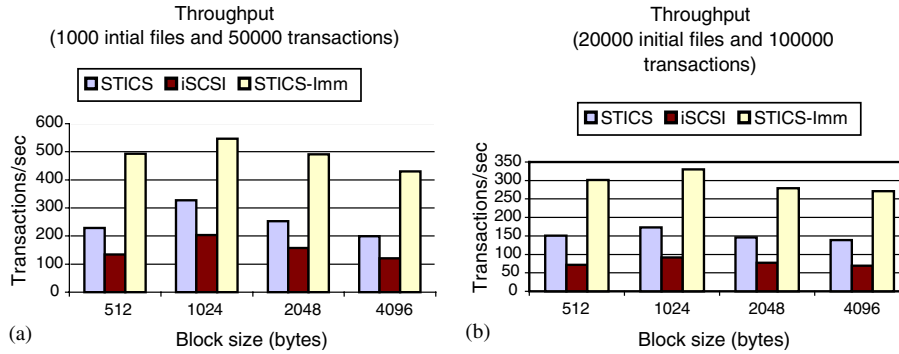


Fig. 6. PostMark measurements (100 Mb network).

Table 3  
Packet distribution

|       | No. of packets with different sizes |           |         |         |          |           |
|-------|-------------------------------------|-----------|---------|---------|----------|-----------|
|       | < 64 Bytes                          | 65–127    | 128–255 | 256–511 | 512–1023 | > 1024    |
| iSCSI | 7                                   | 1,937,724 | 91      | 60      | 27       | 1,415,912 |
| STICS | 4                                   | 431,216   | 16      | 30      | 7        | 607,827   |

Table 4  
Network traffic

|       | Total packets | Full/partial packet ratio | Bytes transferred | Average bytes/packet |
|-------|---------------|---------------------------|-------------------|----------------------|
| iSCSI | 3,353,821     | 0.73                      | 1,914,566,504     | 571                  |
| STICS | 839,100       | 1.41                      | 980,963,821       | 944                  |

1 network than iSCSI does. Tables 3 and 4 show the measured network activities for both STICS and iSCSI. Based  
 3 on our analysis of the numerical results, we believe that the performance gain of STICS over iSCSI can be attributed to  
 5 the following facts. First, the log disk along with the RAM buffer forms a large cache for the host and absorbs small  
 7 writes very quickly, which reduces the network traffic because many data are overwritten in the local log disk. As  
 9 shown in Table 4, the number of total bytes transferred over the network is reduced from 1,914,566,504 to 980,963,821  
 11 although the total data stored in the target storage is the same. Secondly, STICS eliminates many remote handshaking  
 13 caused by iSCSI, which in turn reduce the network traffic. We noticed in Table 3 that the small size packets, which  
 15 are mainly used to transfer iSCSI handshaking messages, are dramatically reduced from 1,937,724 to 431,216. Thirdly,  
 17 by combining small writes into large ones, STICS increases the network bandwidth utilization. If we define full packet  
 19 as the packet with size larger than 1024 bytes of payload data, and other packets are defined as partial packets. As  
 21 shown in Table 4, STICS improves the ratio of full packets to partial packets from 0.73 to 1.41, and average bytes per  
 23 packet is increased from 571 in iSCSI to 944 in STICS.

25 At this point, readers may wonder how many I/O requests are satisfied by the host Linux file system cache and how

many are satisfied by STICS cache. To answer this question, we profile the I/O requests by adding several counters as  
 follows to record the number of requests received at each layer.

- *ReqVFSRcv*: This counter is used to record how many I/O requests received at Linux file system layer. This is done by modifying Linux kernel file system *read\_write* function.
- *ReqToRaw*: This counter is used to record how many I/O requests are forwarded to low-level block layer. This is done by modifying Linux kernel block I/O *ll\_rw\_blk* function. The difference between *ReqToRaw* and *ReqBufferRcv* roughly reflects the number of requests satisfied by Linux file system cache (*ReqFSCache*).
- *ReqSTICSCache*: These counter records the number of requests satisfied by local STICS1 cache.

Table 5 shows the detail breakdown of I/O requests. It is obvious that during the STICS test, 39.58% I/O requests are satisfied by the host file system cache, and additional 29.2% requests are satisfied by STICS cache. We also found that compared to original iSCSI, the reservation of 20 MB system RAM for STICS did not dramatically reduce the hit ratio of file system cache (39.58% vs. 39.84%). The reason is that

Table 5  
Requests breakdown

|       | ReqVFSRcv | ReqToRAW  | ReqFSCache (percentage) | ReqSTICSCache (percentage) |
|-------|-----------|-----------|-------------------------|----------------------------|
| iSCSI | 3,104,257 | 1,867,529 | 1,236,728 (39.84%)      | N/A                        |
| STICS | 3,069,438 | 1,854,554 | 1,214,884 (39.58%)      | 896,582 (29.21%)           |

more requests are satisfied by STICS cache, the replacement possibility is reduced in the file system cache.

Besides the above requests, we have also measured the number of major SCSI commands [55] received on the target side as shown below in Table 6. We have observed dramatic reduction of data transfer commands (*READ* and *WRITE*). STICS filters out most *READ* commands (from 237,659 to 94,502) because many reads are satisfied by local STICS cache. We have also observed a reduction in the number of *INQUIRY* command (from 2009 to 1278) because of less number of data transfer commands (the host does not have to inquiry the storage device so often as seen in original iSCSI implementation).

Above results are measured under 100 MB network, when we configured the switch and network cards as gigabit network, we observed similar results as shown in Fig. 7. The performance gains of STICS with report after complete over iSCSI range from 51% to 80%. STICS with immediate report outperforms iSCSI by a factor of 2.49–3.07. The reason is as follows. When the network is improved from 100 MB to 1 GB, the network latency is not decreased linearly. In our test, we found the average latencies for 100 MB and 1 GB network are 128.99 and 106.78  $\mu$ s. The network performance is improved less than 20% in terms of latency from 100 MB to GB network.

### 5.3.2. IOzone results

Above experiment shows that STICS outperforms iSCSI for workloads consisting of a lot of small files. Our next experiment is to use IOzone to measure the behavior of STICS and iSCSI under a huge file. The network is configured to a 100 MB network. All I/O operations are set to “Synchronous” mode, which means the host is not acknowledged until data is written to a disk. The data set is 512 M bytes in our test.

In Fig. 8, we plotted bar graphs for random read and random write operations against request size for STICS and iSCSI, separately. The request sizes range from 1 KB to 32 KB. In all scenarios, STICS outperforms iSCSI. The performance gain of STICS over iSCSI ranges from 51% to a factor of 4.3.

### 5.3.3. Vxbench results

In order to verify that STICS works well under different host platforms, our next experiment is to test STICS under Solaris operating systems. In this test, we use vxbench to measure the performance, which is a popular file system benchmark program developed by VERITAS Corp. The net-

work is configured to a 100 MB network. The data set is set to 512 M bytes.

We measured the system throughput and host CPU utilization for two workload patterns: random write and mixed I/O as shown in Figs. 9 and 10, respectively. Under STICS, the host CPU utilization increased dramatically, which improved the system throughput as in Fig. 9. The performance gain of STICS over iSCSI ranges from 52% to a factor of 1.9.

### 5.3.4. Response times

Our next experiment is to measure and compare the response times of STICS and iSCSI under EMC trace. The network is configured as a Gigabit network. Response times of all individual I/O requests are plotted in Fig. 11 for STICS with immediate report (Fig. 11a), STICS with report after complete (Fig. 11b) and iSCSI (Fig. 11c). Each dot in a figure represents the response time of an individual I/O request. It can be seen from the figures that overall response times of STICS are much smaller than that of iSCSI. In Fig. 8b, we noticed 4 requests take up to 300 ms. These few peak points drag down the average performance of STICS. These excessive large response times can be attributed to the *destaging* process. In our current implementation, we allow the *level 2 destaging* process to continue until the entire log segment is empty before serving a new storage request. It takes a long time to move data in a full log segment to the remote data disk. We are still working on the optimization of the *destage* algorithm. We believe there is sufficient room to improve the *destaging* process to avoid the few peak response times of STICS.

We also plotted histogram of request numbers against response times in Fig. 12. In this figure, X-axis represents response time and Y-axis represents the number of storage requests finished at a particular response time. For example, a point  $(X, Y) = (10, 2500)$  means that there are 2500 requests finished within 10 ms. As shown in Fig. 12a, for the STICS with immediate report, most requests are finished within 2 ms, because STICS signals the complete of requests when the data are transferred to NVRAM buffer for write requests. The average response time is 2.7 ms. For the STICS with report after complete as shown in Fig. 12b, the response times of the majority of requests fall within the range of 2–5 ms. The rest of requests take longer time to finish but very few of them take longer than 40 ms. The average response time is 5.71 ms.

The iSCSI, on the other hand, has obvious larger response time. The response times of the majority of requests fall within the range of 6–28 ms as shown in Fig. 9c. No requests

Table 6  
Major SCSI commands breakdown

|       | READ    | WRITE   | INQUIRY | READ_CAPACITY | TEST_UNIT_READY |
|-------|---------|---------|---------|---------------|-----------------|
| iSCSI | 237,659 | 309,182 | 2009    | 3             | 53              |
| STICS | 94,502  | 200,531 | 1278    | 3             | 47              |

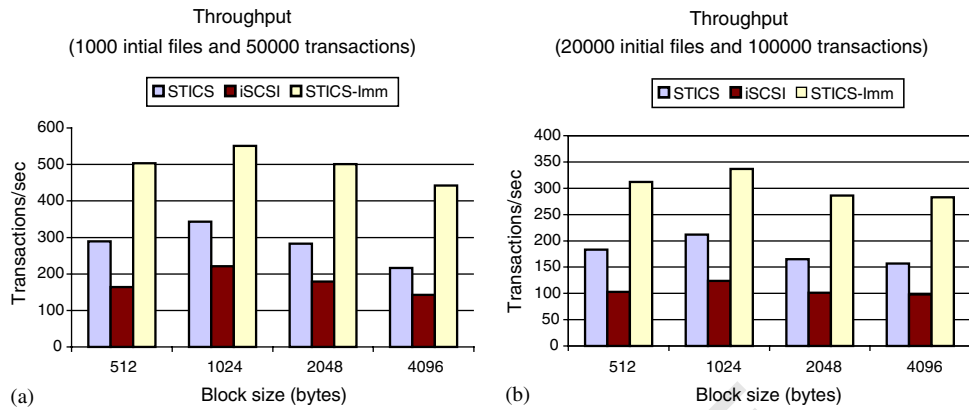


Fig. 7. PostMark measurements (gigabit network).

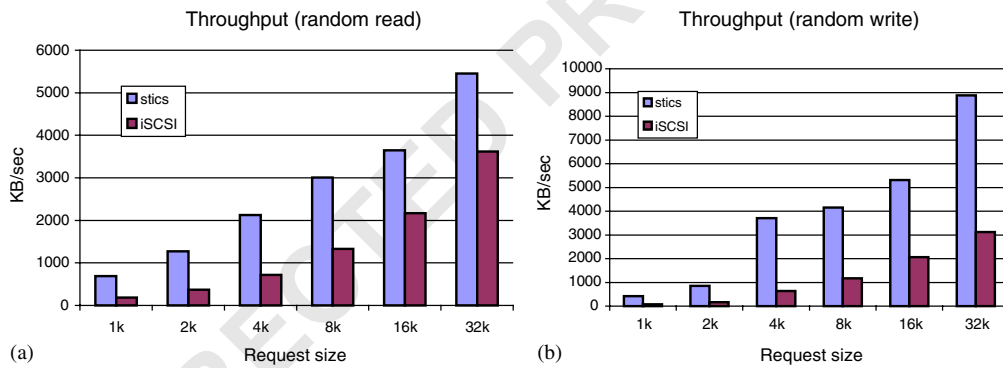


Fig. 8. IOzone measurements.

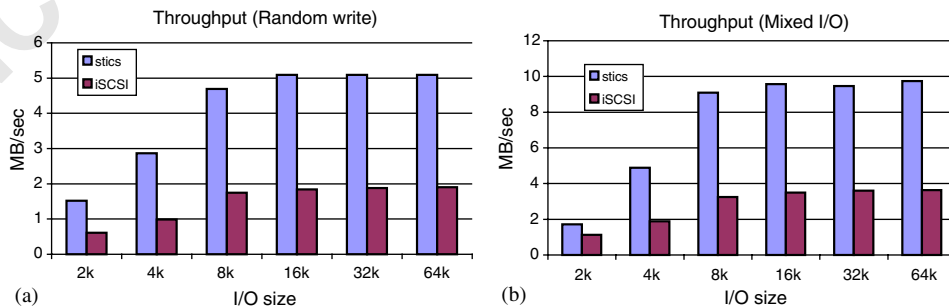


Fig. 9. vxbench Results (throughput).

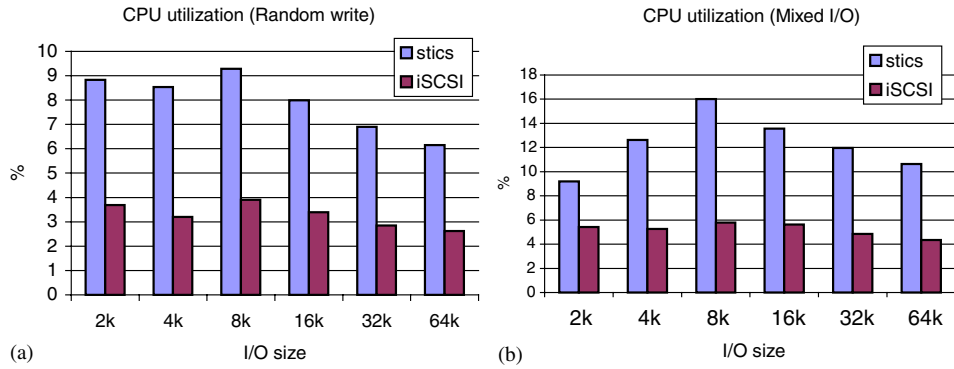


Fig. 10. vxbench Results (CPU utilization).

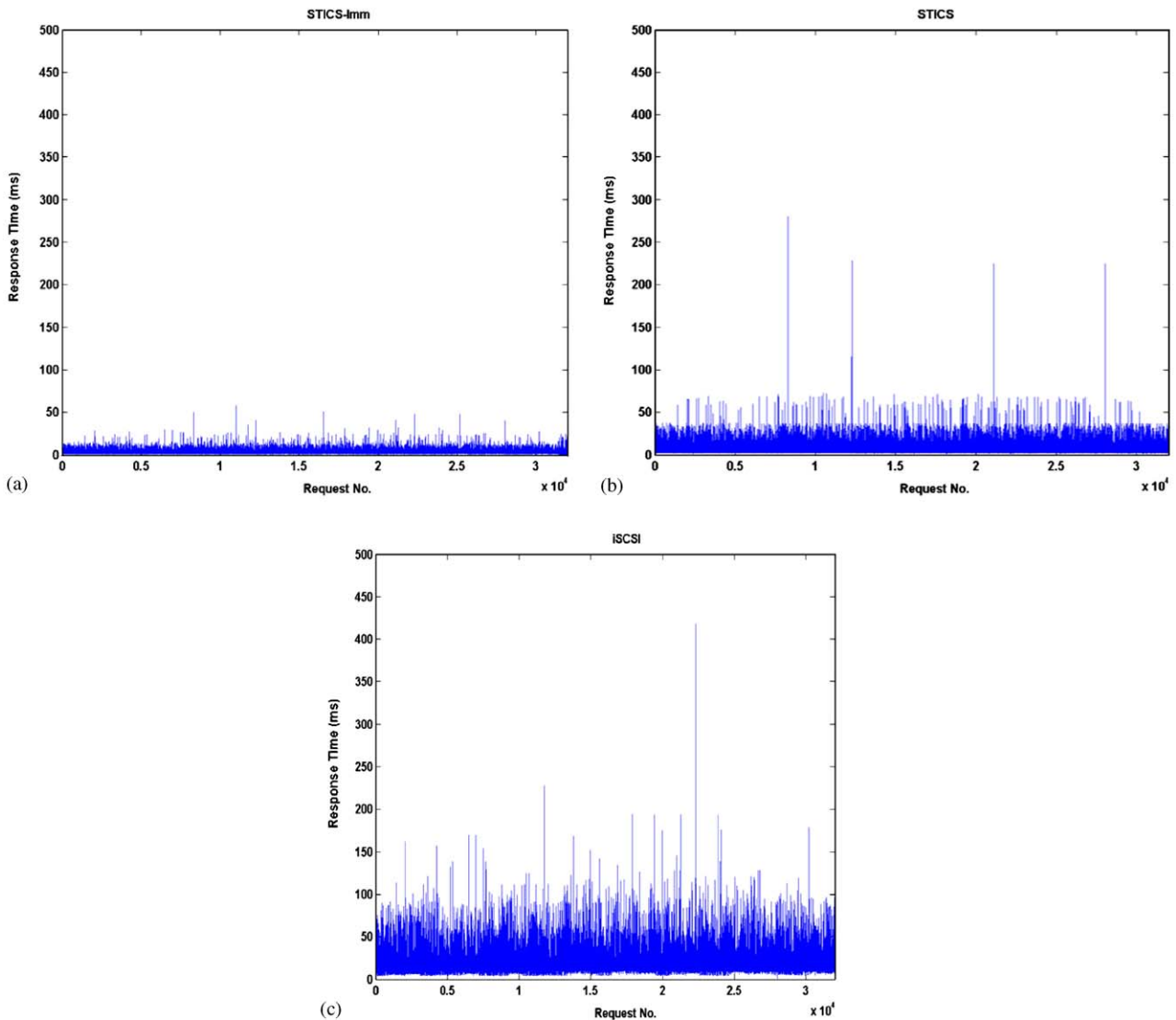


Fig. 11. Response times for EMC-tel trace. Each dot in this figure shows the response time of an individual I/O request. (a) STICS with immediate report (b) STICS with report after complete. (c) iSCSI.

1 are finished within 5 ms. Some of them even take up to  
 400 ms. The average response time is 16.73 ms, which is  
 3 2.9 times as much as STICS with report after complete and

6.2 times as much as STICS with immediate report. Such a  
 long response time can be mainly attributed to the excessive  
 network traffic of iSCSI.



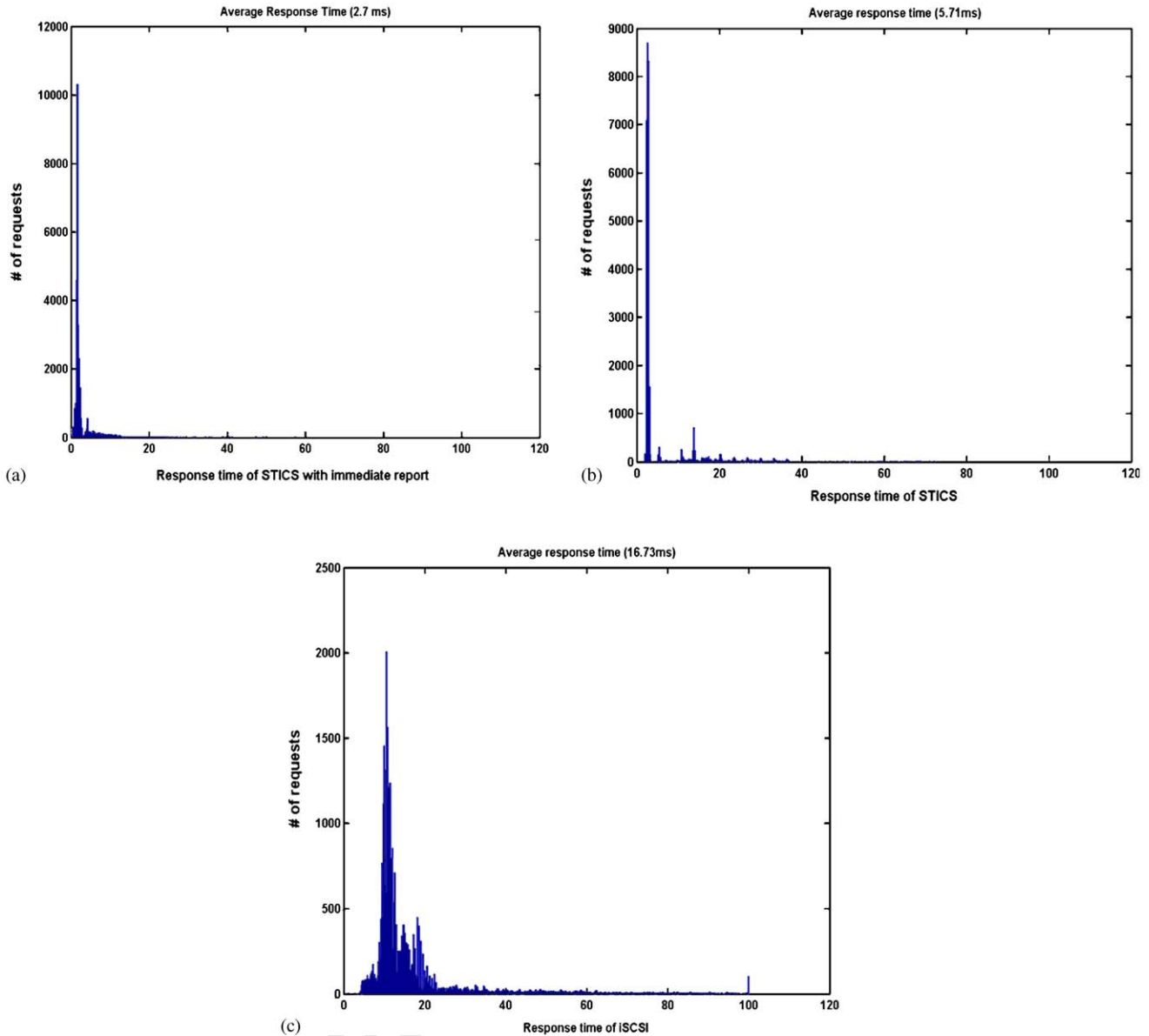


Fig. 12. Histograms of I/O response times for trace EMC-tel. (a) STICS with immediate report (b) STICS with report after complete. (c) iSCSI.

#### 1 5.4. Costs, reliability, and scalability analysis

3 As shown in the last subsection, STICS presents  
 4 significant performance gains over the standard iSCSI im-  
 5 plementation. One obvious question to ask is whether such  
 6 performance improvement comes at extra hardware cost. To  
 7 answer this question, we have carried out cost analysis as  
 8 compared to iSCSI. In our experimental implementations,  
 9 all hardware components such as CPU, RAM, cabling and  
 10 switches are exactly the same for both iSCSI and STICS  
 11 except for an additional disk in STICS for caching. With  
 12 rapid dropping of disk prices, such an additional disk is  
 13 easily justifiable. Typical cost of a 20 GB disk is just around  
 \$50 while a typical SAN costs over tens of thousands dol-

lars, implying a very small fraction of additional cost of  
 STICS.

Table 7 lists the practical cost of building a small (480 GB)  
 SAN configuration with 6 servers using iSCSI and STICS,  
 respectively.<sup>2</sup> As shown in this table, the cost difference  
 between the two is just around 5%. Considering software  
 cost (\$22,059) and maintenance cost (\$8,676) for the even  
 smaller SAN system (200 GB) [31], the cost difference be-  
 tween the two is much less than 2%. We believe trading 2%  
 of additional cost for six folds performance gain is certainly  
 worthwhile.

<sup>2</sup> Prices are as of 12/12/2003 at [www.dell.com](http://www.dell.com).

Table 7  
Hardware costs comparison

|                  | iSCSI |             |           | STICS |            |            |
|------------------|-------|-------------|-----------|-------|------------|------------|
|                  | Qty   | Cost        | Total     | Qty   | Cost       | Total      |
| HBA card         | 12    | \$336.95    | \$4043.40 | 12    | \$336.95   | \$4,043.40 |
| Switch           | 1     | \$1803.95   | \$1803.95 | 1     | \$1803.95  | \$1803.95  |
| GB NIC           | 12    | \$185.95    | \$2231.40 | 12    | \$185.95   | \$2231.40  |
| OS HDD           | 12    | \$52.16     | \$625.92  | 12    | \$52.16    | \$625.92   |
| SCSI storage HDD | 6     | \$573.26    | \$3439.56 | 6     | \$573.26   | \$3439.56  |
| Log disks        |       |             |           | 12    | \$52.16    | \$625.92   |
| Total            |       | \$12,144.23 |           |       | \$12770.15 |            |

We have also considered the cost of implementing iSCSI and STICS in hardware. For the same SAN configuration with 6 servers, iSCSI would need an iSCSI to SCSI converter costing \$5083 [31] or iSCSI cards. The additional hardware for each STICS would include an I/O processor with 4-MB NVRAM. We can conservatively estimate the total cost in addition to Table 6 for 12 STICS to be under \$5000. While at the same time, the cost of a Dell entry level SAN (CX200LC) with 360 GB is \$25,604.

High reliability of STICS is obvious as compared to traditional storage cache using large RAM because STICS uses disks for caching. The small NVRAM in our cache hierarchy is only up to 4 MB. Transient data stay in this NVRAM less than a few hundreds milliseconds. Majority of cached data are in disks that are made extremely reliable today with the mean time to failure of millions of hours. RAM, on the other hand, has much higher failure rate with mean time to failure of a few thousands hours. In addition, RAM cache is also vulnerable to hardware failures such as board failure, CPU failure, and so forth. Disks can be unplugged from a failed system and plugged to another good system with data intact. In addition, log disks can be mirrored to further increase the reliability.

STICS-based SAN systems are also highly scalable. Off-the-shelf Ethernet Switches can be used to connect as many STICS as possible without obvious bottleneck. Furthermore, the LAN connecting STICS can be a completely separate network from the LAN interconnecting servers. This is in contrast to NAS that is attached to the same LAN where servers are connected, competing for the same network resources with servers that access the NAS.

To show the scalability of STICS, we built a larger system as shown in Fig. 13, where STAR1 is the initiator, and STAR2 through 5 are targets. STAR6 is the network traffic monitor. STICS are installed in STAR1..5, where only STICS on the initiator is cache enabled. We also deployed an iSCSI system using the same 5 nodes, where STAR1 is the iSCSI initiator and STAR2 through 5 are iSCSI targets.

We measured the performance of both STICS with report after complete and iSCSI using PostMark benchmark under similar configuration as Section 5.3.1 under Gigabit network. The measurement results in terms of transactions

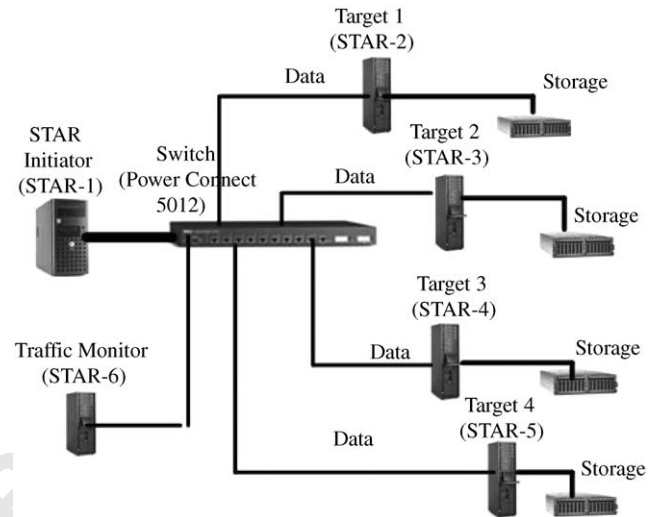


Fig. 13. A larger system area network utilizing STICS.

per second are shown in Table 8. We observed performance gains ranging from 52% to 74% for different block sizes, indicating a good scalability of STICS.

## 6. Conclusions

In this paper, we have introduced a new concept “*SCSI-to-IP cache storage*” (STICS) to bridge the disparities between SCSI and IP in order to facilitate implementation of SAN over the Internet. STICS adds a new dimension to networked storage architectures allowing any server host to efficiently access a SAN on Internet through a standard SCSI interface. Using a nonvolatile “*cache storage*”, STICS smoothes out the storage data traffic between SCSI and IP very much like the way “*cache memory*” smoothes out CPU-memory traffic. We have implemented a prototype STICS under the Linux operating system. We measured the performance of STICS as compared to a typical iSCSI implementation using popular benchmarks (PostMark, IOzone, and vxbench) and a real world I/O workload (EMC’s trace). PostMark, IOzone, and vxbench results have shown that STICS en-

Table 8

Performance comparison between STICS and iSCSI with 5 nodes (20,000 initial files and 10,000 transactions)

|             | Block Size (bytes) |      |      |      |
|-------------|--------------------|------|------|------|
|             | 512                | 1024 | 2048 | 4096 |
| iSCSI       | 318                | 367  | 354  | 283  |
| STICS       | 552                | 589  | 537  | 432  |
| STICS/iSCSI | 1.74               | 1.61 | 1.52 | 1.53 |

1 hances performance of iSCSI by a factor of 4.18, 4.3, and  
 3 1.9, respectively, in terms of average system throughput.  
 Numerical results under EMC's trace show a factor of  
 5 2.9–6.2 performance gain in terms of average response  
 time. Furthermore, STICS is a plug-and-play building block  
 for storage networks.

## 7 Acknowledgements

The first author's research was partially supported by the Center for Manufacturing Research and Research Office under Faculty Research Grant at Tennessee Technological University. The second and third authors' research was supported in part by National Science Foundation under grants CCR-0312613 (ITR) and CCR-0073377, and a Prototype Funding Award of University of Rhode Island Foundation. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the Foundations. The authors would like to thank EMC Corporation for providing trace files to us and thank Jian Li for his invaluable assistance in our experiments. The authors would like to thank Dr. Joe Anderson for proofreading the original manuscript. The authors also acknowledge many valuable comments and suggestions from the anonymous reviewers. A short and preliminary version of this paper was presented at the *International Conference on Parallel Processing (ICPP'2002)* [15], Vancouver, Canada, August 18–21, 2002.

## References

- 9 [1] A. Acharya, M. Uysal, and J. Saltz, Active disks: programming model,  
 algorithms and evaluation, Proceedings of the Eighth International  
 11 Conference on Architectural support for Programming Languages  
 and Operating Systems (ASPLOS'98), October 2-7, 1998, pp. 81–91.
- 13 [2] S. Adve, V. Adve, M. Hill, and M. Vernon, Comparison of  
 hardware and software cache coherence schemes, Proceedings of  
 15 the 18th Annual International Symposium on Computer Architecture  
 (ISCA'91), May 1991, pp. 298–308.
- 17 [3] M. Ahamad, R. Kordale, Scalable consistency protocols for  
 distributed services, IEEE Trans. Parallel Distrib. Systems 10 (1999)  
 19 888.
- 21 [4] S. Aiken, D. Grunwald, A. Pleszkun, J. Willeke, A performance  
 analysis of iSCSI protocol, IEEE Mass Storage Conference, 2003.
- 23 [5] T. Anderson, M. Dahlin, J. Neefe, D. Patterson, D. Roselli, R. Wang,  
 Serverless network file systems, ACM Trans. Comput. Systems 14  
 (1996) 41–79.
- [6] E. Bilir, R. Dickson, Y. Hu, M. Plakal, D. Sorin, M. Hill, D. Wood,  
 Multicast snooping: a new coherence method using a multicast  
 address network, Proceedings of the 26th Annual International  
 Symposium on Computer Architecture (ISCA'99), 1999, pp. 294–304.
- [7] A. Brown, D. Oppenheimer, K. Keeton, R. Thomas, J. Kubiawicz,  
 D. Patterson, ISTORE: introspective storage for data-Intensive  
 network services, Proceedings of the seventh Workshop on Hot Topics  
 in Operating Systems (HotOS-VII), March 1999.
- [8] Y. Chen, L. Ni, C.-Z. Xu, J. Kusler, P. Zheng, CoStore: a reliable  
 and highly available storage systems, Proc. of the 16th Annual  
 International Symposium on High Performance Computing Systems  
 and Applications, Canada, June 2002.
- [9] F. Dahlgren, P. Stenstrom, Using write caches to improve  
 performance of cache coherence protocols in shared-memory  
 multiprocessors, J. of Parallel Distrib. Comput. 26 (2) (April 1995)  
 193–210.
- [10] G. Ganger, M. McKusick, C. Soules, Y. Patt, Soft updates: a solution  
 to the metadata update problem in file systems, ACM Trans. on  
 Comput. Systems 18 (2) (2000) 127–153.
- [11] G. Gibson, R. Meter, Network attached storage architecture, Comm.  
 ACM 43 (11) (November 2000) 37–45.
- [12] G. Gibson, D. Nagle, W. Courtright, N. Lanza, P. Mazaitis, M.  
 Unangst, J. Zelenka, NASD scalable storage systems, USENIX99,  
 Extreme Linux Workshop, Monterey, CA, June 1999.
- [13] J. Gray, Storage bricks have arrived, Keynote in Conference on Fast  
 and Storage Technologies (FAST'2002), Monterey, CA, 2002 January  
 28–30.
- [14] X. He, Q. Yang, M. Zhang, Introducing SCSI-To-IP cache for storage  
 area networks, Technical Report, URL: [http://ele.uri.edu/hexb/  
 publications/STICS-Tech-200112.pdf](http://ele.uri.edu/hexb/publications/STICS-Tech-200112.pdf).
- [15] X. He, Q. Yang, M. Zhang, Introducing SCSI-To-IP cache for  
 storage area networks, The 2002 International Conference on Parallel  
 Processing (ICPP'2002), August 18–21, 2002.
- [16] J. Hennessy, M. Heinrich, A. Gupta, Cache-coherent distributed  
 shared memory: perspective on its development and future challenges,  
 Proc. IEEE 87 (1998) 418–429.
- [17] R. Hernandez, C. Kion, G. Cole, IP storage networking: IBM NAS  
 and iSCSI solutions, Redbooks Publications (IBM), SG24-6240-00,  
 June 2001.
- [18] Y. Hu, Q. Yang, DCD-disk caching disk: a new approach for boosting  
 I/O performance, 23rd Annual Intl. Symposium on Computer  
 Architecture (ISCA'96), May 1996, pp. 169–178.
- [19] Y. Hu, Q. Yang, T. Nightingale, RAPID-Cache—a reliable  
 and inexpensive write cache for disk I/O systems, In: The  
 fifth International Symposium on High Performance Computer  
 Architecture (HPCA-5), Orlando FL, January 1999.
- [20] Intel iSCSI project, URL: <http://sourceforge.net/projects/intel-iscsi>,  
 January 2003.
- [21] J. Jeong, M. Dubois, Cost-sensitive cache replacement algorithms,  
 Proceedings of the Ninth International Symposium on High  
 Performance Computer Architecture (HPCA-9), Anaheim, CA,  
 February 2003.
- [22] J. Katcher, PostMark: a new file system benchmark, Technical  
 Report TR3022, Network Appliance, URL: [http://www.netapp.  
 com/tech\\_library/3022.html](http://www.netapp.com/tech_library/3022.html).

- [23] R. Khattar, M. Murphy, G. Tarella, K. Nystrom, Introduction to storage area network, Redbooks Publications (IBM), SG24-5470-00, September 1999.
- [24] A. Klausner, R. Posch, Distributed caching in networked file systems, *J. Universal Comput Sci* 1 (6) (June 1995) 399–408.
- [25] J. Kubiatowicz, et al., OceanStore: an architecture for global-scale persistent storage, Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'2000), December 2000.
- [26] A. Lebeck, D. Wood, Dynamic self-invalidation: reducing coherence overhead in shared-memory multiprocessors, Proceedings of the 22nd Annual International Symposium on Computer Architecture (ISCA'95), 1995, pp. 48–59.
- [27] E. Lee, C. Thekkath, Petal: distributed virtual disks, Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS 1996), 1996, pp. 84–92.
- [28] D. Lilja, Cache coherence in large-scale shared memory multiprocessors: issues and comparisons, *ACM Comput Surveys* 25 (1993) 303–338.
- [29] H. Lim, V. Kapoor, C. Wighe, D. Du, Active disk file system: a distributed, scalable file system, Proceedings of the 18th IEEE Symposium on Mass Storage Systems and Technologies, April 2001, pp. 101–115.
- [30] Y. Lu, D. Du, Performance study of iSCSI-based storage systems, *IEEE Comm.* 41 (8) (2003).
- [31] B. Mackin, A Study of iSCSI total cost of ownership (TCO) vs. fibre channel and SCSI, Adaptec Co., URL: <http://www.adaptec.com>, Oct. 2001.
- [32] J. Menon, A performance comparison of RAID-5 and log-structured arrays, Proceedings of Fourth IEEE International Symposium on High Performance Distributed Computing, August, 1995, pp. 167–178.
- [33] R. Meter, G. Finn, S. Hotz, VISA: netstation's virtual internet SCSI adapter, Proceedings of the Eighth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), October 1998, pp.71–80.
- [34] K. Meth, J. Satran, Features of the iSCSI protocol, *IEEE Comm* 41 (8) (2003).
- [35] E. Miller, D. Long, W. Freeman, B. Reed, Strong security for network-attached storage, Proceedings of the Conference on Fast and Storage Technologies (FAST'2002), Monterey, CA, January 28–30, 2002.
- [36] V. Milutinovic, M. Valero, The evolution of cache memories, (Special Issue on Cache Memory) *IEEE Trans. Comput.* (February 1999), 97–99.
- [37] D. Nagle, G. Ganger, J. Butler, G. Goodson, C. Sabol, Network support for network-attached storage, Hot Interconnects'1999, August 1999.
- [38] W. T. Ng, et al., Obtaining high performance for storage outsourcing, Proceedings of the FAST'02, January 2002.
- [39] Nishan System white paper, Storage over IP (SoIP) Framework—The Next Generation SAN, URL: [http://www.nishansystems.com/techlib/techlib\\_papers.html](http://www.nishansystems.com/techlib/techlib_papers.html), June 2001.
- [40] W. Norcott, D. Capps, IOzone file system benchmark, URL: [www.iozone.org](http://www.iozone.org).
- [41] A. Nowatzky, G. Aybay, M. Browne, E. Kelly, M. Parkin, B. Radke, S. Vishin, Exploiting parallelism in cache coherency protocol engines, *Euro-Par'95*, 1995, pp. 269–286.
- [42] B. Phillips, Have storage area networks come of age, *IEEE Comput.* 31 (7) (1998).
- [43] E. Riedel, G. Faloutsos, G. Gibson, D. Nagle, Active disks for large-scale data processing, *IEEE Comput.* 34 (6) (June 2001).
- [44] P. Sarkar, S. Uttamchandani, K. Voruganti, Storage over IP: when does hardware support help? Proceedings of second USENIX Conference on File And Storage Technologies (FAST'2003), San Francisco, CA, March 2003.
- [45] J. Satran, et al., iSCSI draft standard, URL: <http://www.ietf.org/internet-drafts/draft-ietf-ips-iscsi-20.txt>, January 2003.
- [46] M. Seltzer, K. Bostic, M. McKusick, C. Staelin, An implementation of a log-structured file system for UNIX, Winter USENIX Proceedings, January 1993, pp. 201–220.
- [47] S. Soltis, T. Ruwart, M. O'Keefe, The global file system, in: Proceedings of the fifth NASA Goddard Space Flight Center Conference on Mass Storage Systems and Technologies, College Park, MD, 1996.
- [48] C. Thekkath, T. Mann, E. Lee, Frangipani: a scalable distributed file system, Proceedings of the 16th ACM Symposium on Operating Systems Principles (SOSP), October 1997, pp. 224–237.
- [49] University of New Hampshire Interoperability Lab iSCSI consortium, URL: <http://www.iol.unh.edu/consortiums/iscsi/>, October 2001.
- [50] A. Varma, Q. Jacobson, Destage algorithms for disk arrays with non-volatile caches, Proceedings of the 22nd Annual International Symposium on Computer Architecture (ISCA'95), 1995, pp. 83–95.
- [51] A. Veitch, E. Riedel, S. Towers, J. Wilkes, Towards global storage management and data placement, Technical Memo HPL-SSP-2001-1, HP Labs, March 2001.
- [52] P. Wang, et al., IP SAN—from iSCSI to IP-addressable ethernet disks, 20th IEEE Conference on Mass Storage Systems and Technologies, 2003.
- [53] R. Wang, T. Anderson, D. Patterson, Virtual log based file systems for a programmable disk, Proceedings of the Third USENIX Symposium on Operating Systems Design and Implementation (OSDI), February 22–25, 1999, pp. 29–43.
- [54] J. Wilkes, R. Golding, C. Staelin, T. Sullivan, The HP autoRAID hierarchical storage system, Proceedings of the Fifteenth ACM Symposium on Operating System Principles, December. 3-6, 1995, pp. 96–108.
- [55] Working Draft, Information Technology: The SCSI Architecture Model-2 (SAM-2), Revision 14, T10-1157-D, URL: <ftp://ftp.t10.org/t10/drafts/sam2/sam2r14.pdf>, Sept. 2000.
- [56] Y. Zhou, J.F. Philbin, K. Li, Multi-queue replacement algorithm for second level buffer caches, USENIX Annual Technical Conference 2001.



**Xubin He** received the Ph.D. in electrical engineering from the University of Rhode Island, USA, in 2002 and both the BS and MS in computer science from the Huazhong University of Science and Technology, China, in 1995 and 1997, respectively. He is an assistant professor of electrical and computer engineering at the Tennessee Technological University. His research interests include computer architecture, storage systems, computer security, and performance evaluation. He received the Ralph E. Powe Junior Faculty Enhancement Award in 2004.

He is a member of the IEEE Computer Society, Sigma Xi, and ASEE.



**Ming Zhang** is currently a Ph.D. candidate in Electrical Engineering at the University of Rhode Island. He received his Bachelor and Master degrees in computer science from the Huazhong University of Science and Technology, China, in 1997 and 2000, respectively. His research interests include computer architecture, networked storage systems, benchmarking, and performance evaluation. He is a student member of the IEEE Computer Society and ACM.



3  
5  
7  
9  
11  
13  
15  
17

**Qing (KEN) Yang** received his B.Sc. in computer science from Huazhong University of Science and Technology, Wuhan, China, in 1982, the M.A.Sc. in electrical engineering from University of Toronto, Canada, in 1985, and the Ph.D. in computer Engineering from the Center for Advanced Computer Studies, University of Louisiana at Lafayette, in 1988. Presently, he is a Distinguished Engineering Professor in the Department of Electrical and Computer Engineering at The University of Rhode Island where he has been a faculty member since 1988. His research interests include computer

architectures, memory systems, disk I/O systems, networked data storages, parallel and distributed computing, performance evaluation, and local area networks. He is a senior member of the IEEE Computer Society and a member of the SIGARCH of the ACM.

19  
21

UNCORRECTED PROOF