

Implementation and Evaluation of iSCSI over RDMA

Ethan Burns and Robert Russell
University of New Hampshire InterOperability Laboratory
121 Technology Drive, Suite 2
Durham, NH 03824-4716
{eaburns,rdr}@iol.unh.edu

Abstract

The Internet Small Computer Systems Interface (iSCSI) is a storage network technology that allows block-level access to storage devices, such as disks, over a computer network. Because iSCSI runs over the ubiquitous TCP/IP protocol, it has many advantages over proprietary alternatives. Due to the recent introduction of 10 gigabit Ethernet, storage vendors are interested in the benefits this large increase in network bandwidth could bring to iSCSI.

To make full use of the bandwidth provided by a 10 gigabit Ethernet link, specialized Remote Direct Memory Access (RDMA) hardware is being developed to offload processing and reduce the data-copy-overhead found in a standard TCP/IP network stack. This paper focuses on the development of an iSCSI software implementation capable of supporting this new hardware, and a preliminary evaluation of its performance.

We describe an approach used to implement iSCSI Extensions for Remote Direct Memory Access (iSER) with the UNH iSCSI reference implementation. This involved a three step process: moving UNH-iSCSI software from the Linux kernel to user-space, adding support for the iSER extensions to the user-space iSCSI, and finally moving everything back into the Linux kernel. Results are given that show improved performance of the completed iSER-assisted iSCSI implementation on RDMA hardware.

1. Introduction

1.1. The iSCSI Protocol

The *Internet Small Computer Systems Interface* (iSCSI) [11] protocol is a *Storage Area Network* (SAN) technology that uses the *Internet Protocol* (IP, specifically TCP/IP) as its underlying transport. Since iSCSI uses the TCP/IP protocol instead of a specialized network fabric, it has a lower total cost of ownership than other SAN technologies [6]. In

addition to being less expensive than its competitors, iSCSI also has the ability to provide networked storage over Local Area Networks (LANs) or across the wide area Internet [15] without requiring any specialized hardware.

1.2. Remote Direct Memory Access

The new 10 Gbit/s Ethernet (10GigE) standard allows IP networks to reach competitive speeds, but without offloading onto specialized *Network Interface Cards* (NICs) standard TCP/IP is not sufficient to make use of the entire 10GigE bandwidth. This is due to data copying, packet processing and interrupt handling on the CPUs at each end of a TCP/IP connection. In a traditional TCP/IP network stack, an interrupt occurs for every packet received, data is copied at least once in each host computer's memory, and the CPU is responsible for processing multiple nested packet headers for all incoming and outgoing packets. To eliminate these inefficiencies, specialized *Remote Direct Memory Access* (RDMA) hardware can be used. One type of RDMA hardware uses the *Internet Wide Area Remote Protocol* (iWARP) protocol suite [3, 12, 10] to move data directly from the memory of one computer to the memory of a remote computer without extra copying or CPU intervention at either end. The entire iWARP protocol suite is offloaded from the CPU onto the *RDMA Network Interface Card* (RNIC). The RNIC reduces the number of data copies that are required by the TCP/IP network to one per host (between main memory and the Ethernet wire), and offloads the bulk of the network processing from the CPU. This means applications utilizing RNICs have lower CPU usage than those utilizing standard NICs, and can achieve throughput close to the full capacity of 10GigE links.

1.3. iSCSI Extensions for RDMA

A difficulty with RDMA is that its usage differs from traditional TCP/IP, and applications need some redesigning to support it. For iSCSI, a new *Internet Engineering*

Task Force (IETF) standard has been created that defines the *iSCSI Extensions for RDMA* (iSER) [7]. iSER describes a set of operational primitives that must be provided by an iSER implementation for use by an iSCSI implementation, and a description of how a conforming iSER-assisted iSCSI session must operate. A design goal of iSER was to require minimal changes to the SCSI architecture and minimal changes to the iSCSI infrastructure, while retaining minimal state information [7]. With iSER, an iSCSI implementation can make use of RDMA hardware for very high speed transfers.

2. Implementation

In order to compare iSER-assisted and traditional software iSCSI implementations, we first needed to build a working iSCSI-iSER-iWARP stack. Due to the complexity of each layer in this stack, we chose to base our work on preexisting implementations instead of creating our own, whenever possible. This process involved selecting from various software projects to use as a base, and then reworking them to fit our needs.

2.1. Protocol Layers

The software layers developed for this project were based on the following prior work:

- The University of New Hampshire iSCSI Reference target and initiator (UNH-iSCSI) [8].
- The OpenFabrics Alliance (OFA) RDMA stack (for communication with the RNIC hardware).
- An old iSER implementation from the OpenIB target project [14].

In order to assist with development, debugging and testing before RDMA hardware became available, we also made use of a software-only implementation of the iWARP protocol, developed at the Ohio Super-Computer Center (OSC), which provides an OFA-compatible interface [5]. However, this software was not incorporated into our final solution.

Figure 1 shows how the protocol layers fit together with iSCSI and iSER. The shaded areas in this figure are implemented in the hardware RNIC; the unshaded areas are implemented in software. The following sections describe the choices made for each of the three software layers, and some of the problems encountered in fitting them together. Additionally, a description of some required modifications is given.

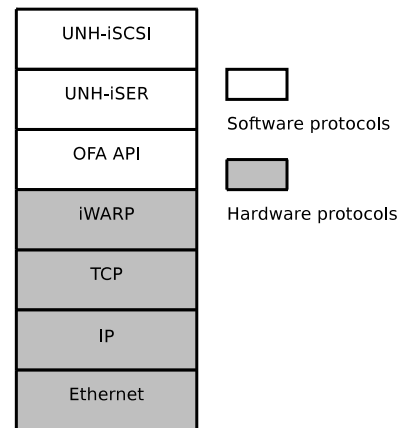


Figure 1. iSCSI-iSER-RDMA Stacks.

2.2. iSCSI

In order to create a working iSER solution, we first had to choose an iSCSI implementation to which we would add iSER support. There are a handful of freely available iSCSI implementations for the Linux operating system. For a number of reasons, we chose one created here at the University of New Hampshire InterOperability Lab (UNH-iSCSI) [8]:

- The UNH-iSCSI implementation was created and is maintained here at the InterOperability Lab (IOL) at the University of New Hampshire, which makes it very easy to get internal support. Most other choices would require a lot of support from external organizations.
- UNH-iSCSI contains both an iSCSI initiator and target together in one project. Both are well tested and supported by the IOL. Most other projects have built either an iSCSI initiator or a target, not both, which could lead to interoperability issues between separate implementations.
- Previous work had already been done to add iSER support to UNH-iSCSI [9]. While other projects are developing, or are thinking about developing, iSER, the UNH-iSCSI project already had some hooks for iSER that were added during this previous, incomplete attempt.

2.3. Three Steps

The UNH-iSCSI project was originally developed as a Linux kernel module, which means the UNH-iSCSI code runs in kernel-space, beneath the standard Linux SCSI layer (which also runs in kernel-space). The OFA stack, however, provides both a kernel-space API and a user-space API

that could be used by iSER-assisted iSCSI. We therefore decided to modify the existing UNH-iSCSI implementation so that it could run in both kernel-space and user-space. This would allow a direct comparison between the two modes of operation. Furthermore, by moving the existing implementation into user-space and then using that to develop the new iSER features, we expected to have a much easier time debugging the new code. This is because a mistake made in kernel-space can often bring down an entire system, whereas a mistake made in user-space will, at most, abort the process that caused the error. In addition, user-space code can be easily run in a debugger, which can provide line-by-line stepping and can give stack-traces when errors are encountered.

Since the Open Fabrics Alliance RDMA stack, which we chose to use for communicating with the RDMA hardware, provides a user-space API in addition to a kernel-space API, and since user-space programs are much easier to debug, and therefore to develop, than kernel modules, we decided to use three steps in our implementation.

1. Move UNH-iSCSI into user-space.
2. Add support for iSER-assisted iSCSI in user-space.
3. Move everything back into the Linux kernel.

2.4. Moving to User-Space

In order to properly mimic kernel-space libraries in user-space, we attempted to use as much code directly from the Linux kernel as possible. Specifically, we were able to directly use Linux kernel source code for doubly-linked lists and the `atomic_t` (a structure that defines a multiprocessor atomic counter) data type. We also used Linux bit-wise operation functions and some other utility routines. We were not able to pull other higher level structures, such as threads and semaphores, directly out of the kernel source code, so we needed to use some user-space functionality. To handle threads, for example, we used the POSIX threads library to implemented functions that provide a behavior similar to the Linux kernel threads API.

The UNH-iSCSI target was the first piece moved to user-space. The kernel-space target is able to run in three different modes: *DISKIO* mode, *FILEIO* mode and *MEMORYIO* mode. In *DISKIO* mode, SCSI commands are sent directly to a hardware disk driver via the standard Linux SCSI mid-level software. In *FILEIO* mode, SCSI commands are mapped into operations on a file stored within the local file system – essentially the blocks of the file are used as if they constituted a disk, but the *FILEIO* software uses file system reads and writes rather than dealing directly with a disk driver. In *MEMORYIO* mode, no permanent storage medium is used – data is read into and written from a primary memory buffer, with no attempt made to retain the

data between operations. *MEMORYIO* mode is intended solely to test the iSCSI subsystem, not to store data.

Of these three modes, only *DISKIO* mode is required to be in kernel-space, because it causes the target to send SCSI commands directly to the Linux SCSI mid-level, an operation that can only happen within the kernel. In the other two modes (*FILEIO* and *MEMORYIO*) the target does not need to be in kernel-space. This design allows the target code to be moved to user-space without significant changes as long as it is not compiled in *DISKIO* mode.

The UNH-iSCSI initiator required more work to get running in user-space. Since the kernel-space initiator receives commands directly from the Linux SCSI mid-level¹ it required some redesigning to move it out of the kernel. In order for the initiator to act like a real SCSI initiator device, we needed to devise a way to pass it SCSI commands that it could then transfer to the target. To accomplish this, we wrote a simple interpreter to drive the initiator. This interpreter reads commands from standard input and performs basic actions, such as logging into a target with a new iSCSI connection, sending a SCSI command to the target, and performing a SCSI READ or WRITE operation. This simple interpreter allows us to write scripts that emulate a real SCSI session on top of the user-space SCSI initiator. In addition, these scripts collect timing information.

2.5. Adding iSER Support

There are a handful of RDMA interfaces that provide POSIX socket-like abstractions to communicate with RDMA hardware. This type of abstraction is advantageous because most programmers are familiar with the socket communication paradigm. A group called the OpenFabrics Alliance (OFA) [1] provides a free software stack with a socket-like abstraction layer called the Communication Manager Abstraction (CMA), along with a verbs layer that is used to perform data transfers over RDMA hardware (together these pieces are referred to as the *OFA stack* or the *OFA API*). The OFA stack has been growing in popularity due to its active development community and inclusion in the Linux kernel. In addition, the OFA stack also provides a user-space library that allows user-space applications to use its API to interface directly with RDMA hardware. We chose to use the OFA API for this project because of its availability, its applicability, and its growing popularity.

The iSER layer we modified was taken from the OpenIB iSER target project [14]. This code was given to the OpenIB group under a dual BSD/GPLv2 license, and therefore was freely available. In addition, an older version of this same iSER implementation was used in a previous, unsuccessful, attempt at adding iSER to UNH-iSCSI [9], so UNH-iSCSI

¹The kernel-space UNH-iSCSI initiator registers itself with the operating system as a Host Bus Adapter (HBA).

already had some support for this iSER interface². In order to use this OpenIB iSER layer, we needed to make many modifications, including moving it to user-space and modifying its lower layer interface to use the OFA stack instead of the kDAPL interface [4].

During this project, we encountered two major problems implementing iSER as defined in the IETF iSER standard [7] with the currently available RNIC hardware and the OFA stack:

1. According to the iSER standard, an iSCSI Login should take place in “normal” TCP streaming mode in order to negotiate the use of the RDMA extensions for the new connection. If both the target and initiator agree, this connection should transition into RDMA mode between the end of the iSCSI Login phase and the start of Full Feature Phase. However, current RNIC hardware does not support such a transition. Therefore, we needed to add new iSER operational primitives to bring up a connection in RDMA mode. This makes moot the Login phase negotiations about the use of RDMA extensions.³
2. iSER for iWARP assumes that a buffer can be advertised with two pieces of information (the length of the buffer and an opaque handle that grants access to it). The OFA stack and current RNIC hardware require that a buffer advertisement actually convey three pieces of information (the length, the opaque handle and a base address). Unfortunately, the standard iSER header for iWARP does not allow space for this extra address. Therefore, we needed to add header fields to accommodate this 64-bit value.

There were also some minor mismatches between the iSER standard and the functionality provided by the OFA stack and iWARP hardware. These are due to the fact that the OFA stack was originally designed for InfiniBand [2] and was only later adapted to iWARP. The differences between the underlying InfiniBand definitions and the iWARP definitions have not been completely reconciled.

2.6. Moving Back into the Kernel

The advantage of having our code run in kernel-space is that we can send and receive SCSI commands to and from the Linux SCSI mid-level. This means that our iSCSI target can use a real disk, and the iSCSI initiator can receive SCSI commands passed down from the Linux virtual file system.

²The iSER version that was used by the previous attempt was not used in this project because of licensing issues.

³When opening a new RDMA connection, the RDMA standard allows for one exchange of private data, but this facility is inadequate to handle a full iSCSI Login phase negotiation.

Unlike the user-space version, the kernel-space iSCSI initiator registers itself with the operating system as a Host Bus Adaptor and can provide the user with access to a real disk attached to a target. This setup allows us to explore our implementation in a real iSCSI environment (with no emulation). This also gives us a real working product that people can freely utilize for high-performance storage networking.

Since we chose to implement a Linux-kernel compatible API for our user-space application, we did not need to change most of our code to compile in the Linux kernel. There are, however, some differences between the OFA user-space API and the OFA kernel-space API. Therefore, we created an abstraction layer on top of these two OFA interfaces so that our iSER layer could use most of the same code for both user-space and kernel-space. Although this small layer is needed mainly to abstract function and type names, there is some more complexity added for event notification and memory registration (which are handled differently in the kernel).

We encountered a large number of difficulties with some of differences between the two OFA interfaces when building our user-space/kernel-space abstraction layer. The kernel-space interface uses call-back functions to notify the user of events, whereas the user-space interface uses completion queues and polling. The kernel-space interface also handles memory registration differently, because the kernel does not use virtual user-space addresses for memory buffers. Since kernel code is not offered as much protection as user-space applications, these differences were difficult to find and fix.

3. Results

After we completed our iSER implementation and got it running in the Linux kernel, we performed a few experiments to evaluate the throughput with the following configurations, all utilizing MEMORYIO mode:

- iSER-assisted kernel-space target and initiator over iWARP RNICs.
- Traditional (unassisted) iSCSI kernel-space target and initiator over unmodified TCP/IP.
- iSER-assisted user-space target and initiator over iWARP RNICs.
- Traditional (unassisted) iSCSI user-space target and initiator over unmodified TCP/IP.

All the graphs shown in this section used the following negotiated iSCSI parameter values:

- MaxRecvDataSegmentLength=512KB

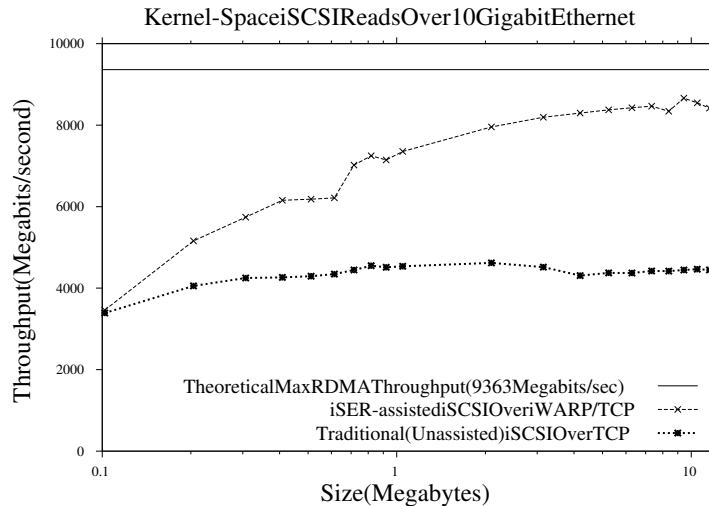


Figure 2. Kernel-space iSCSI Read Throughput

- MaxBurstLength=512KB
- InitiatorRecvDataSegmentLength=512KB
- TargetRecvDataSegmentLength=512KB
- InitialR2T=Yes
- ImmediateData=No

The latter two parameter values were chosen to prevent use of iSCSI immediate and unsolicited data, because the current iSER implementation does not properly handle immediate and unsolicited data. The first four parameter values were chosen because 512KB is the largest SCSI data transfer size generated by the Linux SCSI mid-level when it segments large data transfers into multiple SCSI Read or Write CDBs. Smaller values were not tried because previous work on iSCSI without iSER [13] demonstrated that large values for these parameters produced better performance. Note that InitiatorRecvDataSegmentLength and TargetRecvDataSegmentLength are applicable only in iSER-assisted mode.

It should be noted that when using iSER, the iSCSI R2T, DataIn and DataOut PDUs never appear on the wire between the two communication endpoints. Rather, the iSER-assisted iSCSI target maps each DataIn PDU onto an RDMA Write operation that moves the data (that would have been contained in the DataIn PDU) directly into the initiator's memory via the RNIC. Similarly, the iSER-assisted iSCSI target maps each R2T PDU onto an RDMA Read operation, which moves all the data (that would have been requested in the R2T) directly from the initiator's memory into the target's memory via the RNIC. In non-iSER-assisted mode, the data transferred by this RDMA read

operation would have been carried in DataOut PDUs sent by an initiator.

We used MEMORYIO mode to demonstrate the throughput of the iSER protocol without the overhead of a disk drive, since for this project, we were mainly interested in demonstrating the benefits of using RDMA hardware to assist an iSCSI connection. Overhead introduced by a disk drive is independent of the performance of the iSCSI protocol itself. In any case, all disk drives available to us at present are too slow to demonstrate the throughput possible with 10GigE RNIC hardware.

The computers used to perform these evaluations contained four Intel 2.6GHz, 64-bit processor cores with a total of four gigabytes of main memory. We used the Red Hat Enterprise Linux operating system version 5 with version 1.3 of the OpenFabrics Enterprise Distribution. The kernel version of the systems, as reported by the `uname -a` command, was "Linux 2.6.18-8.el5" with Symmetric Multiprocessing support enabled. The hardware RNIC on each machine was a 10GigE Chelsio R310E-CXA that plugs in to a PCI Express 8X slot and fully offloads the iWARP stack from the CPU. No intervening switch was used – the RNICs in the two machines were connected back-to-back with CX4 copper cable.

In all the graphs shown below, the throughput rate is calculated for the user-level payload only and is shown on the y-axis; the user-level payload size is plotted in logarithmic scale on the x-axis. This payload is the data supplied by or delivered to a user-level application – all protocol header overhead at all levels is excluded from the payload. However, the time used in the rate calculations necessarily includes all overheads, so the maximum user-level payload throughput possible on a 10GigE link using standard 1500

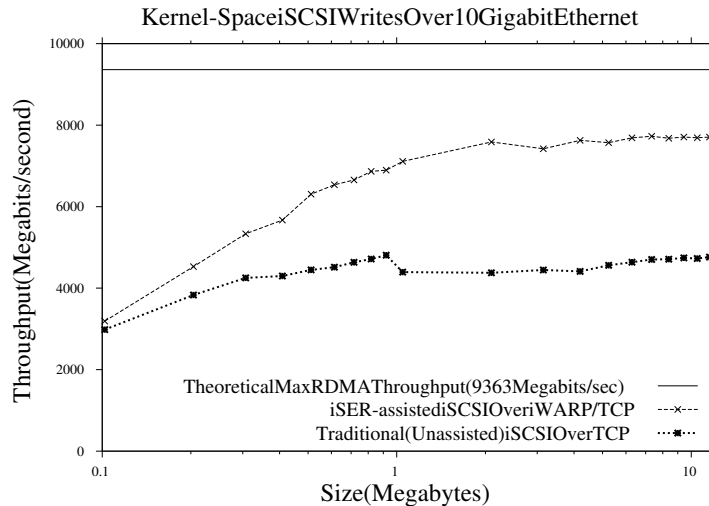


Figure 3. Kernel-space iSCSI Write Throughput

byte Ethernet frames is about 9.36 Gbits/s – the rest of the bandwidth is unavailable to user payload because it is used by headers added by each protocol layer in the network stack.

3.1. Kernel-Space Throughput

Figure 2 shows that the iSER-assisted iSCSI Read throughput surpasses that of a traditional iSCSI implementation for all payload sizes. The improvement due to iSER-assisted iSCSI increases rapidly with an increase in transfer size below 2MB, because, as the transfer sizes increase, the constant overhead involved in setting up an RDMA transfer becomes relatively less significant than the benefits of off-loading and zero-copy transfers. Transfer sizes greater than 2MB show a more constant throughput improvement. The throughput falls short of the maximum payload throughput possible with a 10GigE network due to the following sources of extra overhead:

- Ethernet, TCP/IP and RDMA header bytes sent on the wire use a portion of the 10GigE bandwidth, reducing the amount that is available for payload data.
- iSCSI Command and Response PDUs that are transmitted for each 512KB of data also use a portion of the bandwidth, further reducing that available for payload data.
- CPU processing of each iSCSI Command and Response PDU on both the target and initiator increases the latency of the transfer and correspondingly reduces the throughput.

Figure 3 shows that the iSER-assisted iSCSI Write throughput surpasses that of traditional iSCSI. In this figure, we also see a steep climb in throughput until about 2MB. This happens as the overhead of setting up an RDMA transfer becomes less significant with respect to the size of the data transfer. For transfer sizes greater than 2MB the throughput flattens out and does not achieve the maximum capacity possible with a 10GigE connection. Additionally, this figure shows that Write throughput is not quite as good as Read throughput. One potential reason for this is because iSER-assisted iSCSI Writes are mapped to RDMA Read operations. Each RDMA Read operation consists of a Read-Request/Read-Response message exchange. Since the Linux kernel will not build a SCSI command that contains more than 512KB of data, these Read-Request/Read-Response exchanges happen twice per megabyte transferred, in addition to all of the other overhead mentioned previously.

3.2. User-Space Throughput

Figure 4 shows that the throughput of user-space, iSER-assisted iSCSI Writes ramps up steeply for transfer sizes of under 1MB, but then flattens out at just over 8800 Megabits/sec, which is about 94% of the 9363 Megabits/sec maximum payload throughput possible with a 10GigE network. The user-space results shown in Figure 4 demonstrate noticeably better performance than the kernel-space results shown in Figure 3 (8800 Megabits/sec vs 7700 Megabits/sec).

Figure 5 plots SCSI Write throughput performance using the four combinations of user-space and kernel-space with target and initiator. This figure shows that the most limiting agent is the kernel-space iSCSI initiator, as both of the runs

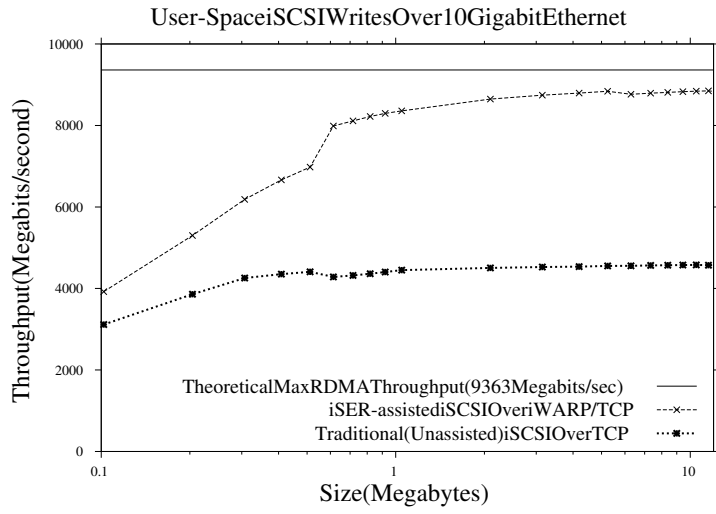


Figure 4. User-space iSCSI Write Throughput

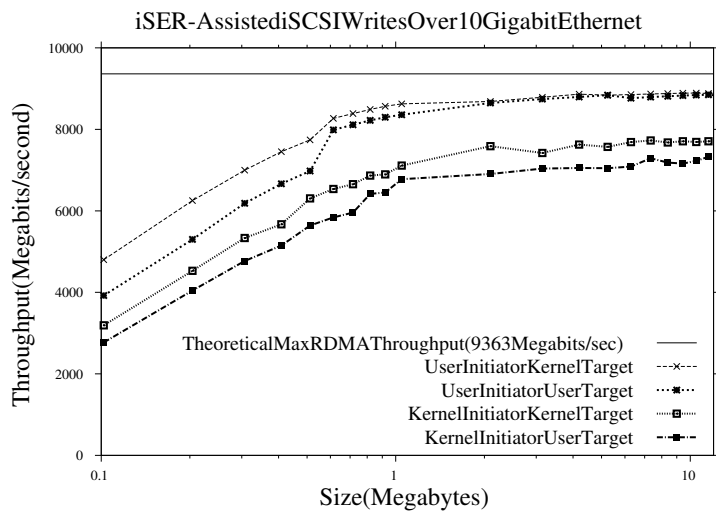


Figure 5. Hybrid User/Kernel-Space iSER Write Operation Throughput

in which it was used exhibit significantly lower throughput than the two runs with the user-space initiator.

Our initiator performs significantly better in user-space than in kernel-space because in user-space it does not use multiple SCSI commands for each data transfer. In kernel-space, the `read` and `write` system calls are passed down from the Linux file-system layer and eventually to the Linux SCSI mid-level. The SCSI mid-level builds a set of SCSI CDBs that it queues up for the UNH-iSCSI initiator to perform the actual data transfer operations. The maximum size SCSI CDB that the Linux SCSI mid-level generates is 512KB. Thus for each megabyte transferred, the mid-level generates two SCSI CDBs, the initiator transmits two iSCSI Command PDUs (one for each CDB), and the target performs two transfer operations and sends two iSCSI Response PDUs. It is also important to note that the iSCSI initiator only sends a single command at a time to the iSCSI target, so multiple commands are not being handled in parallel.

In user-space, the UNH-iSCSI initiator script performs data transfers by building SCSI CDBs that are all the exact size of the transfer. This means that for each transfer, regardless of its size, there is only a single SCSI CDB, a single iSCSI Command PDU, a single transfer operation, and a single iSCSI Response PDU. Thus, for a 12MB transfer, the kernel-space implementation will transmit 24 SCSI CDBs, perform 24 data transfers and send 24 responses, whereas the user-space implementation will only use a single SCSI CDB, data transfer and response.

4. Conclusion and Future Work

The goal of this project was to implement an iSER-assisted iSCSI solution for RDMA over a 10GigE network and to evaluate its performance. In order to accomplish this goal, we added support to the UNH-iSCSI reference implementation for iSER extensions to the iSCSI protocol. This paper has described the approach we took and some of the difficulties encountered. We also presented results of a set of throughput evaluations which show that the iSER extensions enable iSCSI software together with a hardware RNIC to make nearly full use of a 10GigE network.

The performance analysis in this paper is just the beginning of the possible analysis that can be done. Initially we concentrated only on throughput, but clearly latency and CPU utilization are also important. In addition, the many protocol layers offer various parameters that can be changed and the resulting performance compared. We only tested with 1500-byte Ethernet frames, but the software and hardware can equally well utilize 9000-byte Ethernet frames, which should show better performance since the overhead is cut by a factor of 6. A comparison between iWARP and InfiniBand would also be interesting, and one of the ben-

efits of designing software based on the OFA stack is that our implementation should run over both technologies. Finally, we need to run some comparisons against available ASIC-based TOE/iSCSI acceleration adaptors, which provide a fully offloaded, 10GigE-based hardware alternative to iWARP RNICs.

The source for this project, both user-space and kernel-space versions, is freely available under the GPL open-source license at <http://sourceforge.net/projects/unh-iscsi>.

References

- [1] O. F. Alliance. <http://www.openfabrics.org>.
- [2] I. T. Association. Infiniband Architecture Specification version 1.2.1.
- [3] P. Culley, U. Elzur, R. Recio, S. Bailey, and J. Carrier. Marker pdu aligned framing for tcp specification. RFC 5044 (Standards Track), Oct. 2007.
- [4] DAT-Collaborative. *kDAPL: Kernel Direct Access Programming Library*, June 2002.
- [5] D. Delessandro, A. Devulapalli, and P. Wyckoff. Design and Implementation of the iWarp Protocol in Software. PDCS, Nov. 2005.
- [6] J. L. Hufferd. *iSCSI The Universal Storage Connection*. Addison Wesley, 2003.
- [7] M. Ko, M. Chadalapaka, J. Hufferd, U. Elzur, H. Shah, and P. Thaler. Internet Small Computer System Interface (iSCSI) Extensions for Remote Direct Memory Access (RDMA). RFC 5046 (Standards Track), Oct. 2007.
- [8] A. Palekar, A. Chadda, N. Ganapathy, and R. Russell. Design and implementation of a software prototype for storage area network protocol evaluation. In *Proceedings of the 13th International Conference on Parallel and Distributed Computing and Systems*, pages 21–24, Aug. 2001.
- [9] M. Patel and R. Russell. Design and Implementation of iSCSI Extensions for RDMA, Sept. 2005.
- [10] R. Recio, B. Metzler, P. Culley, J. Hilland, and D. Garcia. A remote direct memory access protocol specification. RFC 5040 (Standards Track), Oct. 2007.
- [11] J. Satran, K. Meth, C. Sapuntzakis, M. Chadalapaka, and E. Zeidner. Internet Small Computer Systems Interface (iSCSI). RFC 3720 (Proposed Standard), Apr. 2004. Updated by RFCs 3980, 4850.
- [12] H. Shah, J. Pinkerton, R. Recio, and P. Culley. Direct data placement over reliable transports. RFC 5041 (Standards Track), Oct. 2007.
- [13] Y. Shastry, S. Klotz, and R. Russell. Evaluating the Effect of iSCSI Protocol Parameters on Performance. In *Proceedings of the IASTED International Conference on Parallel and Distributed Computing and Networks (PDCN 2005)*, pages 159–166, Feb. 2005.
- [14] Voltaire. Open source iSER code, at <https://svn.openfabrics.org/svn/openib/gen2/>.
- [15] F. Xu and R. Russell. An architecture for public internet disks. In *Proceedings of the 3rd International Workshop on Storage Network Architecture and Parallel I/O (SNAPI05)*, pages 73–80, Sept. 2005.