



100G JBOF using Chelsio Offloads

Using FADU Delta SSDs & Chelsio T6 Adapter

Executive Summary

This paper presents the storage performance results of FADU Delta SSDs Just a Bunch of Flash (JBOF) target connected to multiple hosts/initiators with the following network configurations using Chelsio T6 100GbE adapters:

S.No	Target	Hosts/Initiators
1	Kernel NVMe/TCP	Kernel NVMe/TCP
2	Kernel NVMe/TCP Offload	Kernel NVMe/TCP
3	SPDK NVMe/TCP Offload	Kernel NVMe/TCP
4	LIO iSCSI Offload	Open iSCSI
5	Kernel NVMe/iWARP Offload	Kernel NVMe/iWARP Offload
6	SPDK NVMe/iWARP Offload	Kernel NVMe/iWARP Offload

The test result proof points in this paper show that T6 Offload:

- Delivers line-rate 94 Gbps READ throughput.
- Reaches 2.9 Million IOPs at 4K I/O size.
- Provides local-like access to remote storage.
- Provides significant CPU savings compared to NVMe/TCP (no-offload).

Overview

FADU is a company developing advanced flash storage technology to meet the explosively increasing data storage demands placed on hyperscale, enterprise and cloud data centers. The FADU innovative SSD solutions are based on industry-standard specifications, designed with FADU's proprietary Flash Memory Controller architecture, and compatibility with multiple industry NAND suppliers. FADU's storage platform design addresses all aspects of SSD storage requirements - very low power, ultra-high performance, rich feature sets, solid reliability, and superior QOS. FADU is currently delivering PCIe Gen3 (Bravo) and Gen4 (Delta) E1.S and U.2 SSD solutions.

Chelsio Terminator 6 (T6) 1/10/25/40/50/100GbE Unified Wire adapters offer Storage Protocol and TCP/IP offload (TOE) capabilities. Chelsio adapters enable enterprise storage systems to deliver optimized NVMe/TCP and iSCSI performance for various application workloads in mission-critical virtualized and private cloud environments. Chelsio adapters are designed for industry-leading performance, efficiency, and have a unique ability to offload multiple storage protocols (iSCSI, iSER, NVMe/iWARP, NVMe/TCP, FCoE, S2D, SMB) using a single ASIC and firmware. Chelsio adapters unburden communication responsibilities and processing overhead from host servers and storage systems resulting in a dramatic increase in application performance with a minimum of CPU cycles.

The combination of FADU SSDs with Chelsio's Unified Wire adapter solution delivers compelling performance, power savings, and total cost of ownership (TCO) advantages. This enables innovative topologies and networked computing models to address the most demanding processing needs.



Test Overview

The following tests were conducted showing the FADU Delta SSDs JBOF Target performance results:

1. [Random READ Throughput \(Bandwidth\) and %CPU/Gbps](#) at I/O sizes from 8 to 256 Kbytes
2. [4K Random READ IOPs](#)
3. [4K Random READ and WRITE latency](#) for local (direct-attached using the PCIe bus) versus remote (Ethernet network-based) access

Test Results

Throughput and %CPU/Gbps Results

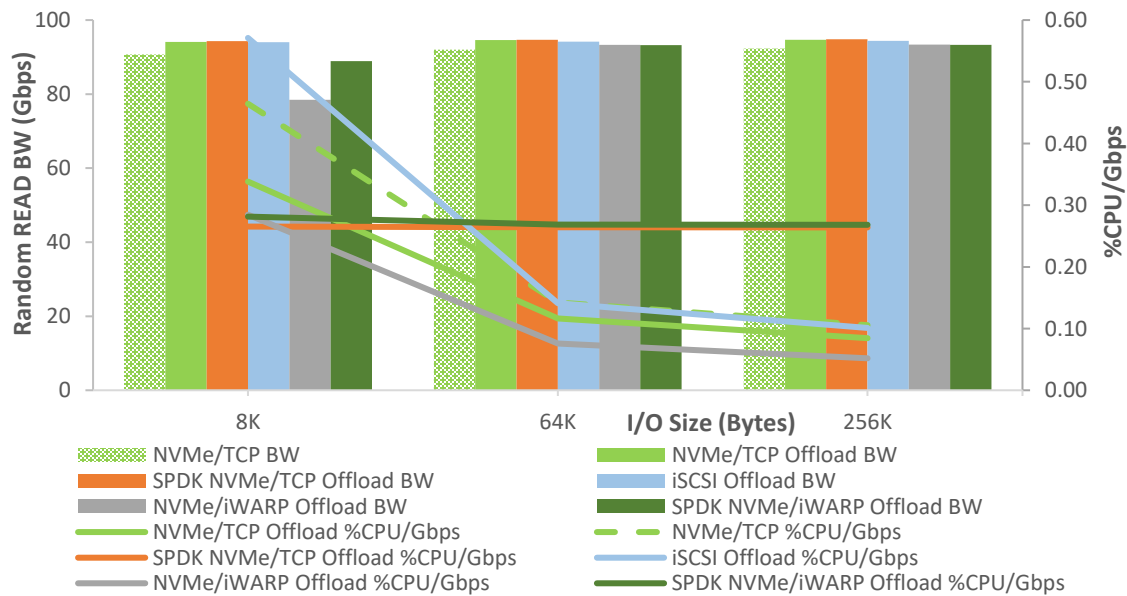


Figure 1 – Random READ Throughput (BW) and %CPU/Gbps vs. I/O size

Random READ		8K		64K		256K	
S.No	Target	BW	%CPU/Gbps	BW	%CPU/Gbps	BW	%CPU/Gbps
1	NVMe/TCP	90.70	0.46	92.00	0.14	92.30	0.11
2	NVMe/TCP Offload	94.10	0.34	94.60	0.12	94.70	0.08
3	SPDK NVMe/TCP Offload	94.30	0.27	94.70	0.26	94.80	0.26
4	iSCSI Offload	94.00	0.57	94.20	0.14	94.40	0.10
5	NVMe/iWARP	78.5	0.28	93.3	0.08	93.4	0.05
6	SPDK NVMe/iWARP	88.90	0.28	93.20	0.27	93.30	0.27

Table 1 – BW (Gbps) and %CPU/Gbps vs. I/O size

The above graph and table show how the T6 delivers line-rate READ throughput for all the target modes of the FADU solution. T6 Kernel NVMe Offloads consume significantly less server CPU compared to the Kernel NVMe/TCP (no-offload).



4K Random READ IOPs Results

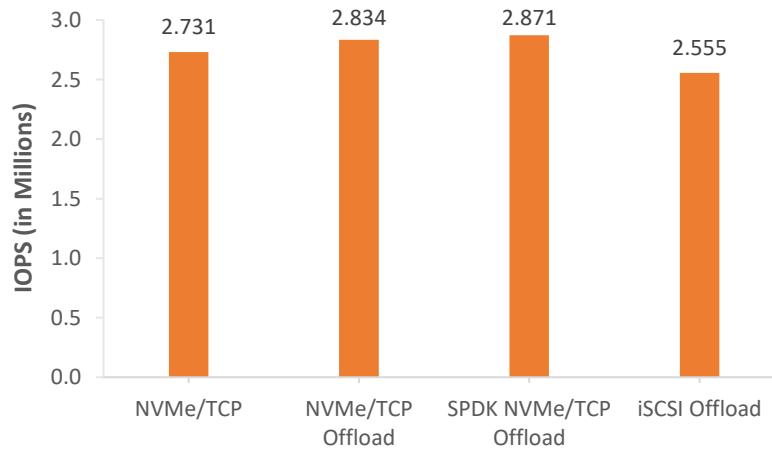


Figure 2 – 4K Random READ IOPs

The T6 NVMe Offload enabled solution delivers 2.9M random READ IOPs at 4K IO Size.

4K Latency Results

Local refers to the latency for accessing the direct-attached SSDs (i.e., using the PCIe bus). *Remote* refers to the latency for accessing the SSDs across the Ethernet network.

4K Random Latency		Read			Write		
S.No	Target	Local	Remote	Delta	Local	Remote	Delta
1	NVMe/TCP	62.08	89.14	27.06	10.44	36.19	25.75
2	NVMe/TCP Offload	62.08	87.59	25.51	10.44	34.44	24
3	SPDK NVMe/TCP Offload	62.08	83.64	21.56	10.44	30.49	20.05
4	iSCSI Offload	62.08	88.17	26.09	10.44	36.13	25.69
5	NVMe/iWARP Offload	62.08	75.91	13.83	10.44	22.9	12.46
6	SPDK NVMe/iWARP Offload	62.08	74.23	12.15	10.44	21.38	10.94

Table 2 – 4K Random Latency

T6 Offload (with SPDK) provides the least delta latency between remote and local storage access. This demonstrates the local like performance of remote distributed storage using T6 Offload enabled and SPDK solution.

Test Setup

The Bandwidth, %CPU and IOPs test setup consists of a storage target machine (with 2 FADU Delta U.2 3.84 TB SSDs direct-attached) and 4 host/initiator machines connected through a 100GbE switch using a single port on each system. 2 hosts/initiators connect to 1 SSD using 20 connections each.

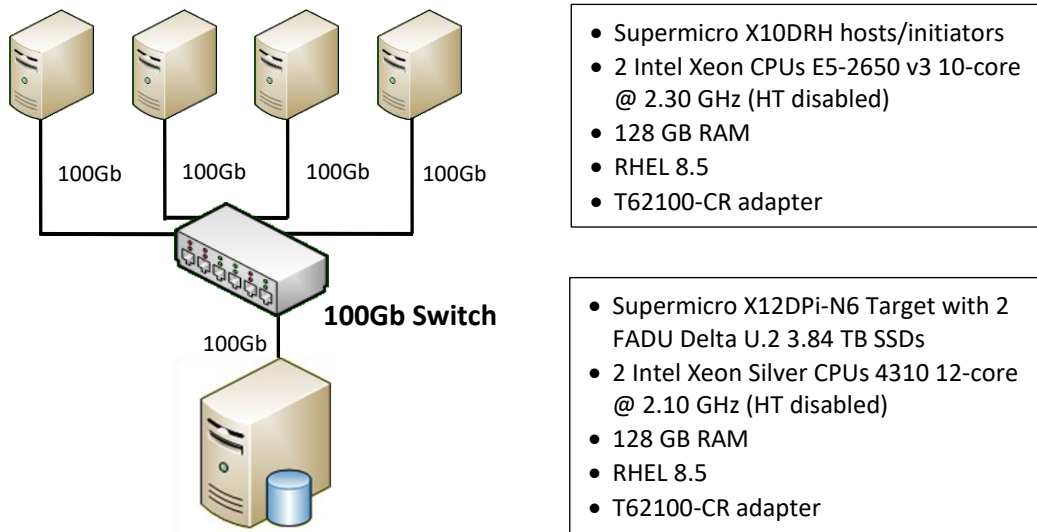


Figure 3 – BW, %CPU and IOPs Test Setup

The Latency test setup consists of a storage target machine (with 1 FADU Delta U.2 3.84 TB SSD direct attached) connecting directly (no switch) to a single host/initiator machine using a single port on each system. The host/initiator connects to the target using one connection.

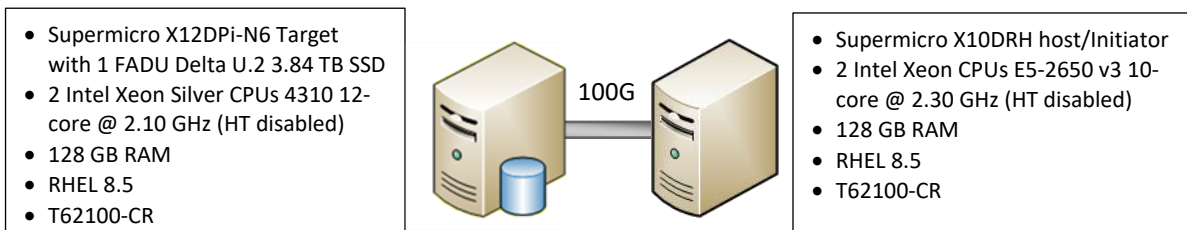


Figure 4 – Latency Test Setup

For all the tests, standard MTU of 1500 is used on the ports under test. The latest Chelsio Unified Wire driver for Linux is installed on all machines.

Set-up Configuration

General Configuration

Execute the below steps on target and all host/initiator machines.

- Disable virtualization, c-state technology, VT-d, Intel I/O AT, SR-IOV in system BIOS.
- Compile and install the latest Chelsio Unified Wire package and reboot the machine.

```
[root@host~]# cd ChelsioUwire-x.x.x.x
[root@host~]# make install
[root@host~]# reboot
```

- Add the below parameters to grub kernel command line.

BW/IOPs test: intel_idle.max_cstate=0 processor.max_cstate=0 intel_pstate=disable

Latency test: idle=poll



- iv. Set cpupower governor to performance.

```
[root@host~]# cpupower frequency-set --governor performance
```

- v. Set the below tuned-adm profile for BW/IOPs test.

```
[root@host~]# tuned-adm profile network-throughput
```

Set the below tuned-adm profile for latency test.

```
[root@host~]# tuned-adm profile network-latency
```

- vi. Load the Chelsio NIC driver (*cxgb4*) and bring up interface with an IPv4 address.

```
[root@host~]# modprobe cxgb4  
[root@host~]# ifconfig ethX <IPv4 address> up
```

- vii. Precondition the SSDs.

```
[root@host~]# nvme format /dev/nvme0n1 -f  
[root@host~]# fio --rw=write --name=precon --norandommap --fsync_on_close=1  
--invalidate=1 --direct=1 --filename=/dev/nvme0n1 --iodepth=64 --numjobs=1 -  
-ioengine=libaio --bs=512K --ramp_time=5 --size=512G
```

Kernel NVMe-oF Target Configuration

- i. Load the Chelsio NVMe/TCP Offload drivers (not required for non-offload NVMe/TCP).

```
[root@host~]# modprobe t4_tom
```

Note: For NVMe/iWARP, load *iw_cxgb4* instead of *t4_tom*.

- ii. Apply the below cop policy for NVMe/TCP Offload.

```
[root@host~]# cat /root/cop_policy  
all => offload !nagle !ddp !coalesce  
[root@host~]# cop -d -o file /root/cop_policy  
[root@host~]# cxgbtool ethX policy file
```

- iii. Set the CPU affinity for BW/IOPs test.

```
[root@host~]# t4_perftune.sh -s -n -Q ofld -c 0-23
```

Set the CPU affinity for Latency test to use a single CPU (0 in this case).

```
[root@host~]# t4_perftune.sh -s -n -Q ofld -c 0
```

Note: For NVMe/TCP, please use *nic* instead of *ofld*. For NVMe/iWARP, please use *rdma* instead of *ofld*.

- iv. Load the NVMe drivers.

```
[root@host~]# modprobe nvmet  
[root@host~]# modprobe nvmet-tcp  
[root@host~]# modprobe nvmet-rdma
```



v. Configure the target using the below script.

```
#!/bin/bash
nvmectl clear > /dev/null 2>$1
mount -t configfs none /sys/kernel/config > /dev/null 2>$1

IPPORT="4420"      # 4420 is the reserved NVME/Fabrics port
IPADDR="0.0.0.0"
NAME="nvme-ssd"
DEV="/dev/nvme"

for i in `seq 0 1`; do
mkdir /sys/kernel/config/nvmet/subsystems/${NAME}${i}
mkdir -p /sys/kernel/config/nvmet/subsystems/${NAME}${i}/namespaces/1
echo -n "${DEV}${i}n1"
>/sys/kernel/config/nvmet/subsystems/${NAME}${i}/namespaces/1/device_path
echo 1 > /sys/kernel/config/nvmet/subsystems/${NAME}${i}/attr_allow_any_host
echo 1 > /sys/kernel/config/nvmet/subsystems/${NAME}${i}/namespaces/1/enable
done

mkdir /sys/kernel/config/nvmet/ports/1
# echo 8192 > /sys/kernel/config/nvmet/ports/1/param_inline_data_size
echo "ipv4" > /sys/kernel/config/nvmet/ports/1/addr_adrfam
echo "tcp" > /sys/kernel/config/nvmet/ports/1/addr_trtype
echo $IPPORT > /sys/kernel/config/nvmet/ports/1/addr_trsvcid
echo $IPADDR > /sys/kernel/config/nvmet/ports/1/addr_traddr

for i in `seq 0 1`; do
ln -s /sys/kernel/config/nvmet/subsystems/${NAME}${i}
/sys/kernel/config/nvmet/ports/1/subsystems/${NAME}${i}
done
```

Note: For NVMe/iWARP, please use *param_inline_data_size* and use *rdma* instead of *tcp*.

SPDK NVMe-oF Target Configuration

i. Load the Chelsio SPDK NVMe/TCP Offload driver.

```
[root@host~]# modprobe chtcp
```

Note: For NVMe/iWARP, please use *iw_cxgb4* instead of *chtcp*.

ii. Configure Huge Pages.

```
[root@host~]# echo 8192 > /proc/sys/vm/nr_hugepages
[root@host~]# echo 8192 > /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages
[root@host~]# cd ChelsioUwire-x.x.x.x/build/src/chspdk/user/spdk/
[root@host~]# CLEAR_HUGE=yes HUGENODE='nodes_hp[1]=8192,nodes_hp[0]=8192'
scripts/setup.sh config
```

iii. Start the target.

```
[root@host~]# cd ChelsioUwire-x.x.x.x/build/src/chspdk/user/spdk/
[root@host~]# ./build/bin/nvmf_tgt -m 3F000
```

iv. Configure the target with the SSDs.



- ii. Run *fio* tool on all 4 hosts at the same time.

```
[root@host~]# fio --rw=randwrite/randread --ioengine=libaio --name=random --norandommap --group_reporting --exitall --fsync_on_close=1 --invalidate=1 --direct=1 --runtime=60 --time_based --filename=<device list> --iodepth=16 --numjobs=20 --bs=<value> --unit_base=1 -kb_base=1000 --ramp_time=2
```

Latency test:

- i. Connect single host to the target.

```
[root@host~]# nvme connect -t tcp -a 10.2.2.136 -n nvme-ssd0 -i 1

nqn.2016-06.io.spdk:cnode0 should be used while connecting to the SPDK NVMe/TCP Offload Target.
```

Note: For NVMe/iWARP, please use *rdma* instead of *tcp*.

- ii. Run *fio* tool on the host using a single CPU (10 in this case).

```
[root@host~]# taskset -c 10 fio --rw=randwrite/randread --ioengine=libaio --name=random --invalidate=1 --direct=1 --runtime=60 --time_based --fsync_on_close=1 --group_reporting --filename=<device list> --iodepth=1 --numjobs=1 --bs=4K --unit_base=1 -kb_base=1000 --ramp_time=2 --size=512G --randrepeat=0
```

iSCSI Offload Target Configuration

- i. Load the Chelsio LIO PDU Offload Target driver.

```
[root@host~]# modprobe cxgbit
```

- ii. Set the CPU affinity for BW/IOPs test (to use *local_cpulist* of interface).

```
[root@host~]# t4_perftune.sh -s -n -Q iSCSIT -c 12-23
```

Set the CPU affinity for Latency test to use a single CPU (0 in this case).

```
[root@host~]# t4_perftune.sh -s -n -Q iSCSIT -c 0
```

- iii. Create 10 namespaces on each NVMe SSDs.

```
[root@host~]# nvme delete-ns /dev/nvme0 -n 1
[root@host~]# nvme delete-ns /dev/nvme1 -n 1
[root@host~]# for i in `seq 1 10` ; do nvme create-ns /dev/nvme0 --nsze=58605568 --ncap=58605568 --flbas=0 -dps=0 ; done
[root@host~]# for i in `seq 1 10` ; do nvme create-ns /dev/nvme1 --nsze=58605568 --ncap=58605568 --flbas=0 -dps=0 ; done
[root@host~]# for i in `seq 1 10` ; do nvme attach-ns /dev/nvme0 --namespace-id=${i} -controllers=0x1 ; done
[root@host~]# for i in `seq 1 10` ; do nvme attach-ns /dev/nvme1 --namespace-id=${i} -controllers=0x1 ; done
```

- iv. Configure the target using the below script.

```
i=1
for j in `seq 1 10` ; do
```




```
        /usr/local/bin/targetcli /backstores/block/ create name=nvme0n${j}
dev=/dev/nvme0n${j} readonly=false
        /usr/local/bin/targetcli /iscsi create iqn.2022-11.org.linux-
iscsi.target${i}
        /usr/local/bin/targetcli /iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpg1/ set attribute authentication=0
demo_mode_write_protect=0 generate_node_acls=1 cache_dynamic_acls=1
        /usr/local/bin/targetcli /iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpg1/luns create lun=0
storage_object=/backstores/block/nvme0n${j}
        /usr/local/bin/targetcli /iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpg1/portals/ delete ip_address=0.0.0.0 ip_port=3260
        /usr/local/bin/targetcli /iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpg1/portals create ip_address=11.11.11.10 ip_port=3260
echo 1 > /sys/kernel/config/target/iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpgt_1/np/11.11.11.10\:3260/cxgbit
        let i=i+1
done

for j in `seq 1 10` ; do
        /usr/local/bin/targetcli /backstores/block/ create name=nvme1n${j}
dev=/dev/nvme1n${j} readonly=false
        /usr/local/bin/targetcli /iscsi create iqn.2022-11.org.linux-
iscsi.target${i}
        /usr/local/bin/targetcli /iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpg1/ set attribute authentication=0
demo_mode_write_protect=0 generate_node_acls=1 cache_dynamic_acls=1
        /usr/local/bin/targetcli /iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpg1/luns create lun=0
storage_object=/backstores/block/nvme1n${j}
        /usr/local/bin/targetcli /iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpg1/portals/ delete ip_address=0.0.0.0 ip_port=3260
        /usr/local/bin/targetcli /iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpg1/portals create ip_address=11.11.11.10 ip_port=3260
echo 1 > /sys/kernel/config/target/iscsi/iqn.2022-11.org.linux-
iscsi.target${i}/tpgt_1/np/11.11.11.10\:3260/cxgbit
        let i=i+1
done
```

iSCSI Initiator Configuration

- i. Enable Adaptive Rx for the Chelsio Interface.

```
[root@host~]# ethtool -C ethX adaptive-rx on
```

- ii. Set the CPU affinity for BW/IOPs test (to use *local_cpulist* of interface).

```
[root@host~]# t4_perftune.sh -s -n -Q nic -c 10-19
```

Set the CPU affinity for Latency test to use a single CPU (0 in this case).

```
[root@host~]# t4_perftune.sh -s -n -Q nic -c 0
```

BW/IOPs test:

- i. Connect all the 4 initiators to the targets using the below command.

```
[root@host1~]# for i in `seq 1 20` ; do iscsiadm -m node -T iqn.2022-
11.org.linux-iscsi.target${i} -p 10.10.10.10 -l ; done
```



- ii. Run *fiio* tool on all 4 initiators at the same time.

```
[root@host~]# fio --rw=randwrite/randread --ioengine=libaio --name=random --norandommap --group_reporting --exitall --fsync_on_close=1 --invalidate=1 --direct=1 --runtime=60 --time_based --filename=<device list> --iodepth=64 --numjobs=20 --bs=<value> --unit_base=1 -kb_base=1000 --ramp_time=2
```

Latency test:

- i. Connect single initiator to the target.

```
[root@host~]# iscsiadm -m node -T iqn.2022-11.org.linux-iscsi.target1 -p 10.10.10.10 -l
```

- ii. Run *fiio* tool on the host using a single CPU (10 in this case).

```
[root@host~]# taskset -c 10 fio --rw=randwrite/randread --ioengine=libaio --name=random --invalidate=1 --direct=1 --runtime=60 --time_based --fsync_on_close=1 --group_reporting --filename=<device list> --iodepth=1 --numjobs=1 --bs=4K --unit_base=1 -kb_base=1000 --ramp_time=2 --size=512G --randrepeat=0
```

Conclusion

This paper showcases the performance capabilities of Chelsio T6 100G Offload-enabled adapters with FADU Delta SSDs. With concurrent support for NVMe/TCP, NVMe/iWARP and iSCSI, users can create and maintain a true Converged Fabric-based server cluster for software-defined storage and other applications.

The Chelsio T6 enables the FADU SSDs to be shared, pooled, and managed more effectively across a low latency, high-performance, scalable, standard Ethernet network, with CPU server savings - delivering a highly cost-effective solution.

Related Links

[FADU Delta Gen4 SSD](#)

[100G JBOF using Chelsio Offloads](#)

[Best Practices for NVMe/TCP Deployment](#)

[“No-compromise” NVMe/TCP Deployment using Server Storage I/O Offload](#)

[100G Kernel and User Space NVMe/TCP Using Chelsio Offload](#)

[100G iSCSI Performance for AMD EPYC](#)

[Demartek Evaluation: Chelsio Terminator 6 \(T6\) Unified Wire Adapter iSCSI Offload](#)

[Windows iSCSI Performance at 100Gbps](#)