



NVMe/TCP implementation on F-stack

Shinik Park ECE, Seoul National University

Network Becomes The Bottleneck of NVMeoF



SATA with AHCI

(Advanced Host Controller Interface)



PCIe with NVMe

(Non-Volatile Memory Express)

	SATA I	SATA II	SATA III	PCle
Max Bandwidth	1.5 Gbps	3 Gbps	6 Gbps	32 Gł

 PCle 3.0 x4
 PCle 4.0 x4
 PCle 5.0 x4

 32 Gbps
 64 Gbps
 128 Gbps

Fortunately, DAS (Directed Attached Storage) could be a feasible option at the datacenter.

NVMeoF (NVMe over Fabric) is becoming a trend which enables end nodes to access remote NVMe storage pool.



with combination of 100G Ethernet card (e.g., Intel E810, uses PCIe 4.0 x16) NVMeoF: Techniques to access remote NVMe Storage

□ NVMeoF could be supported by multiple options.

- FC (Fibre Channel)
- Infiniband
- iWARP
- RoCE (RDMA over Converged Ethernet)

Pros: Low latency, high bandwidth Cons: Custom host bus adapters and drivers are required, difficult to deploy, maintain & cost inefficient

TCP

NVMe/TCP (NVMe over TCP) is compatible with standard ethernet adapters. → Cost effective & Easy to deploy and maintain!

Possible Design Strategy for NVMe/TCP

□ Kernel level TCP implementation (NVMe/TCP)

- Directly modifying the raw kernel source code (i.e. kernel module implementation)
- \rightarrow Due to the hard deployment, it would not be handled.
- □ User-level implementation (NVMe/TCP)
 - DPDK (Data Plane Development Kit) based implementation (i.e. mTCP , F-stack)
 - Frameworks that implements TCP-stack into user-level.
 - ② eBPF (extended Berkeley Packet Filter) based implementation
 - Pushing packet processing code black into kernel hooking point.

User-level TCP is cost effective option to support high performance NVMe/TCP.

User-level TCP Can Mitigate Kernel Overhead

□ mTCP [NSDI `14] conducts experiment to see Linux kernel overhead.

Experiment setup

- 8-core Intel Xeon CPU (2.9 GHz, E5-2690) / 32 GB memory
- 10 Gbps NIC (Intel 82599 chipsets)
- Lighttpd v1.4.32 web server
- Multiple clients repeatedly download a 64 KB file per connection.



processed per CPU cycle in the kernel]

mTCP reduce CPU utilization of kernel from 80 % to 30 %. → mTCP can 4.3 times more effectively than linux.

User-level TCP Varients

- mTCP [NSDI `14]: One of the famous user-level TCP implementation. mTCP is basically framework for research so that has no compatibility to mature application.
- F-stack: Open source network framework developed by Tencent. F-stack include an user space TCP/IP stack(port FreeBSD 11.0 stable). It support Nginx, Redis and other mature applications.





Testbed Setup: Kernel TCP Stack Performance in 100G Environment

□ Simple one-hop testbed is organized to test the performance of F-stack and legacy TCP.



□ First, we test **default kernel TCP** to check the effect of kernel overhead.

- Simple data transmission from server to client by *iperf* tool
- \rightarrow Records around 30-35 Gbps when NIC and Switch can support up to 100 Gbps.



ntel@intel-S2600WFT:~\$ iperf -c 10.10.10.1						
Client connecting to 10.10.10.1, TCP port 5001 CCP window size: 85.0 KByte (default)						
3] local 10.10.10.2 port 35676 connected with 10.10.10.1 port 5001						
[3] 0.0-10.0 sec 37.4 GBytes 32.1 Gbits/sec						
		•				

[Snapshot of client-side log]

DPDK Speed Test with *testpmd* in 100G Environment

□ *testpmd*: TestPMD is the reference application distributed with the DPDK.

- Can measure performance by forwarding packets between Ethernet ports on a network interface.
- □ Results
 - \sim 98 Gbps is recorded \rightarrow DPDK itself can fully support 100G hardware with appropriate parameters!



Does F-stack can fully utilize the great performance of DPDK?

Performance Test Tool for F-stack

Performance comparison with kernel level TCP-stack

- SpeedTest source code should be implemented on F-stack.
- Since F-stack does not provide iperf-like application.

Implementation of SpeedTest application

- Done by using kqueue. (event driven I/O handling system call, similar with epoll in Linux)
- Send Large file from F-stack server to F-stack Client to evaluate maximum throughput under the given environment.

Goal:

Max throughput comparison & improving F-stack performance

[Core code block of SpeedTest]



Poor Performance of F-stack in 100G Environment

□ On the same one hop 100G testbed, **default F-stack**

Socket 2048 connected	
	SpeedTest Result
Goodput: 1351.981 Mbps,	Total Bytes: 100052700, Total elapsed time: 0.592036



- Records only ~1 Gbps... Poor performance compared to DPDK performance test.
- \rightarrow This implies that performance bottleneck placed between DPDK and user-level TCP stack!

Where could be the bottleneck in F-stack...?

F-stack Break Down for Finding the Bottleneck



□ **3X Performance enhancement** compared to the default F-stack.



- ① Main bottleneck is come from F-stack library to use DPDK functions.
- ② Development of congestion control algorithm for supporting NVMe/TCP

Congestion Control: Available Options in F-stack

- □ F-stack supports multiple congestion control algorithms.
 - I picked up some of algorithms that which could support NVMe/TCP.
 - TCP CUBIC [SIGOPS `08]
 - Loss-based algorithms used by default TCP of Linux(`06), MacOS(`14), Windows(`17).
 - DCTCP [SIGCOMM `10]
 - ECN (Explicit Congestion Notification)-based algorithms to support high bandwidth datacenter networks.
 - BBR [Queue `16]
 - Model-based congestion control developed by Google.
 - The performance of F-stack is highly dependent on the type of congestion control.
 - ① Limitations of these methods to make hard to support NVMe/TCP with the best performance.
 - ② Finding a suitable congestion control candidates to support NVMe/TCP.

Congestion Control: Limitation of introduced CC

□ The problem of TCP CUBIC and DCTCP which is implemented in F-stack

- These algorithm can only check maximum available bandwidth by "Max probing"
- \rightarrow Max probing implies that these algorithms **MUST accompany queueing** at the bottleneck link.
- $\hfill\square$ The problem of BBR
 - Model-based congestion control algorithm has its limitation on slow adaptation.
 - → Because it can respond to the network change **after satisfying the state transition condition**.



Mean and 95-th percentile RTT against the avg. throughput over a real network.

In actual performance evaluation, The performance of these CCs are far from BDP.

We decide to bridge the gap by using the core idea from HPCC [SIGCOMM `19].

HPCC: Direct Access to Queue Information via INT

Goal: To directly access the queue-related information of the bottleneck switch at the end node.

Adjusting flow rates per ACK Sender INT Overhead (42 bytes for 5 hops)

INT (In-Network Telemetry) Overview

ппор	pathiD	1 st Hop(64 bits)			2 nd Hop	
(4 bits)	(12 bits)	В	TS	txBytes	qLen	(64 bits)

- **B**: The type of speed of the egress port
- **TS**: The timestamp when the packet is emitted from its egress port
- txBytes: the accumulative total bytes sent from the egress port
- **qLen**: the current queue length of the egress port

Key Idea: Using In-Network Telemetry (INT), delivering queue information to control cwnd and sending rate



(Prototype) NVMe CC: HPCC with less overhead INT

- □ NVMe CC inherits the core idea of HPCC.
 - By using INT at the programmable switch.
- However, prototype of NVMe CC has multiple advantages with HPCC.
 - ① U calculation is done at the switch side.
 - ② INT overhead becomes 7 or 8 bytes. (HPCC: $2 + 8 \times n$ bytes, n: # of hops)
- Prototype implementation is now in progress.



NVMe CC Prototype Implementation on F-stack



TCP Header Option Implementation: addoptions

□ HPCC Option : Option length [server: 8 bytes, client: 7 bytes]

1 Byte	1 Byte	2 Bytes	2 Bytes	2 Bytes
Opt-	Opt-	Opt-	Opt-	Opt-data
kind	length	ExID	Magic#	(U and U-echo)

[Server-side option field reservation]

1 Byte	1 Byte 1 Byte		2 Bytes	1 Byte
Opt-	Opt-	Opt-	Opt-	Opt-data
kind	length	ExID	Magic#	(U-echo)

[Client-side U-echo relaying]

□ In <u>freebsd/netinet/tcp_output.c</u>, **tcp_addoptions()** modification.

with variable definition at tcpopt structure in <u>freebsd/netinet/tcp_var.h</u>

```
case TOF_INT:
                                                                                                                       ield reservation in data packet at the sender
                                                                                                           * Sender returns received U as the U-echo to the switch */
        /* Receiver-side U-echo */
        /* U-echo option generation (7 bytes) in ACK packet at the receiver */
                                                                                                          else {
       /* To do: Have to relay the data from the data packet to the U-echo field! */
if( to->uflags != 0 ) {
    while (optlen % 4) {
                                                                                                                   while (optlen % 4) {
                                                                                                                           optlen += TCPOLEN_NOP;
                                                                                                                            *optp++ = TCPOPT_NOP;
                         optlen += TCPOLEN_NOP;
                         *optp++ = TCPOPT_NOP;
                                                                                                                  optlen += TCPOLEN_SINT;
                                                                                                                   *optp++ = TCPOPT_SINT;
                optlen += TCPOLEN_RINT;
                *optp++ = TCPOPT_RINT;
                                                                                                                   /* Field Reservation */
                                                                                                                    *optp++ = TCPOLEN_SINT;
                *optp++ = TCPOLEN_RINT;
                                                                                                                   to \rightarrow to ExID = 2;
                to->to_ExID = htons(to->to_ExID);
                                                                                                                   to->to_ExID = htons(to->to_ExID);
                bcopy((u_char *)&to->to_ExID, optp, sizeof(to->to_ExID));
optp += sizeof(to->to_ExID);
                                                                                                                   bcopy((u_char *)&to->to_ExID, optp, sizeof(to->to_ExID));
                                                                                                                   optp += sizeof(to->to_ExID);
                to->to_magicnum = htons(to->to_magicnum);
                                                                                                                   to -> to_magicnum = 1;
                bcopy((u_char *)&to->to_magicnum, optp, sizeof(to->to_magicnum));
                                                                                                                   to->to_magicnum = htons(to->to_magicnum);
                                                                                                                   bcopy((u_char *)&to->to_magicnum, optp, sizeof(to->to_magicnum));
                optp += sizeof(to->to_magicnum);
                                                                                                                   optp += sizeof(to->to_magicnum);
                /* U-echo with received U value */
                 *optp++ = to->to_U;
                                                                                                                   /* U value reservation with 0 value */
                                                                                                                   *optp++ = 0;
                                                                                                                   /* Actual Uecho value return to the switch */
                                                                                                                   to->to_Uecho = to->to_U;
                                                                                                                   *optp++ = to->to_Uecho:
```

break:

TCP Header Option Implementation: doptions

□ Packet Reception : Option length [server: 7 bytes, client: 8 bytes]

Server receives 7 bytes option from ACK / Client receives 8 bytes option from the switch.

1 Byte	1 Byte	2 Bytes	2 Bytes	1 Byte
Opt-	Opt-	Opt-	Opt-	Opt-data
kind	length	ExID	Magic#	(U-echo)

[Server-side option reception from ack]

1 Byte	1 Byte	2 Bytes	2 Bytes	2 Bytes
Opt-	Opt-	Opt-	Opt-	Opt-data
kind	length	ExID	Magic#	(U and U-echo)

[Client-side option reception from the switch]

□ In <u>freebsd/netinet/tcp_input.c</u>, **tcp_dooptions()** modification.

with variable definition at tcpopt structure in <u>freebsd/netinet/tcp_var.h</u>



- Server have to save U-echo from the ack and relay to the HPCC algorithm.
- Client have to save U value from the switch and feedback to the server through the ACK.

NVMe CC Core Algorithm Implementation

□ HPCC module implementation is finished (Verified by unittest)

- By relaying U-echo from ack, implements whole HPCC algorithm on the f-stack.
- cc_ack_received() roles the main control loop for the HPCC algorithm.



[Core code block of HPCC algorithm]

hpcc_compute_wind(CCV(ccv, downlink_U), true);

hpcc_data->lastUpdate_seq = CCV(ccv, snd_nxt);

hpcc_compute_wind(CCV(ccv, downlink_U, false)

/*To Do: have to check actual value is precisely reflect intened value */ if(U >= (uint8_t)eta || hpcc_data->incStage >= hpcc_data->maxStage) {

hpcc_data->wnd = hpcc_data->wnd_c + hpcc_data->wnd_ai;

hpcc_data->wnd = hpcc_data->wnd_c * (u_long)(eta / U) + hpcc_data->wnd_ai

if(ccv->curack > hpcc_data->lastUpdate_seq) {

pcc ack received(struct cc var *ccv, uint16 t type)

CCV(ccv, snd_cwnd) = hpcc_data->wnd;

npcc_compute_wind(struct cc_var *ccv, uint8 t U, bool updateWc)

hpcc_data->incStage = 0; hpcc data->wnd c = hpcc data->wnd;

hpcc_data->incStage++; hpcc data->wnd c = hpcc data->wnd;

struct hpcc *hpcc_data; hpcc data = ccv->cc data;

if(type == CC ACK) {

else {

struct hpcc *hpcc_data;

hpcc data = ccv->cc data;

if(updataWc) {

if(updataWc) {

static void

else {

Conclusion and Future Works

□ Conclusion

- ① NVMe/TCP is new trend to support fast access to the remote storage pool.
 - With low cost & easy deployment
- ② User-level TCP is one of the best options for NVMe/TCP.
- ③ Performance bottleneck of user-level TCP is coming from F-stack library.
- ④ Congestion control algorithm is the key factor for the performance of NVMe/TCP.

□ Future works

- ① Main performance bottleneck is come from F-stack library.
 - We have to focus on F-stack library optimization to dramatically increase F-stack performance.
 - F-stack developers are also aware of this issue.
- ② Upgrade NVMe CC algorithm to utilize "the number of flows" which could be observed at the programmable switch.
 - By doing so, fairness performance is expected to be definitely improved.