# Emerging Ethernet standards
# &
# their impact on Storage

## Anupama B N
## NetApp

# Agenda

- Ethernet Technology Landscape
- Ethernet Standards and Technology
- Connector Standards and types
- RDMA (RoCE and iWARP)
- RoCE vs iWARP
- RoCE over long distance
- References

# Ethernet  Technology Landscape

- FCoE
- RDMA
    - iWARP
    - RoCE
- SDN
    - Extension **into** the VM environment vSphere/OpenVswitch/Nexus
    - Provisioning and orchestration tools, focus on **Overlays** – VXLAN, NVGRE, GENEVE, **WAN**
- NVMf
    - NVMe over Fabrics
- iSCSI
- iSER

2019 Storage Developer Conference India ©  All Rights Reserved.

# Data Center

- Bigger – 1km cable runs common
- Fill as you go, leave in place
- Manage via API (remote), ports set up on demand via API
- Leaf/Spine Clos (vs. Tree)



Old

New

# High Speed Interconnects

Low latency access **+** High speed transports

Storage Class Memory (SCM) as Cache

NVMe over Fabrics

Persistent Memory (PMEM) in Server

Hybrid
SAS

All-Flash
NVMe/SAS

Flash At All Tiers
QLC-SSD/SCM

**NetApp**

# Ethernet Technology and Standards

| CY16 | CY17 | CY18 | CY19 | CY20 | CY21 | CY22 |
|------|------|------|------|------|------|------|

**Technology**

- 50G PAM4 SERDES
- Silicon Photonics
- 100G PAM4 SERDES
- Silicon Photonics in NIC/Switch ASICs
- 50/100/200/400G Early NIC and Switch ASICs
- 100/200/400/800G Early NIC and Switch ASICs

**Standards**

- IEEE 802.3bq 25/40GBase-T
- IEEE 802.3bs 200G/400G (4/8 lane)
- IEEE 802.3cd 50/100/200G
- IEEE 802.3ck 100/200/400G
- IEEE 802.3by 25G
- IEEE 802.3cc 25G SMF
- IEEE 802.3cz 800G
- CWDM4 PSM4 MSA
- OSPF MSA
- IEEE 802.1AS-REV time sensitive Ethernet
- QSFP-DD

SDC 19 SNIA INDIA

6

NetApp

# 802.3bs/cd Signaling

**NRZ to PAM4**

☐ **PAM4**

☐ **1, 2, 4, and 8**

Source: Mellanox blog, neophotonics



**NRZ** =1,0; one-bit/clock pulse (Non-Return to Zero)

**PAM4** = 00,01,10,11; 2-bits per clock pulse
Pulse Amplitude Modulation u- levels

+ Enables twice the data transferred while using lower 25G
clock rate to keep component costs down.

# Connector Types

## MSA Mainstream: New Double Density Connectors



**QSFP-DD**
8-Channel
12W

**QSFP28**
4-Channel
3.5W

**SFP-DD**
2-Channel
3.5W

Source  SFP-DD consortium, QSFP-DD consortium

# What is RDMA

- RDMA – Remote Direct Memory Access
- Benefits
    - Very low latency, very high throughput, ≈ zero CPU
    - Bypasses traditional network stacks (TCP/IP)
    - Provides a Fibre Channel-equivalent solution at a lower cost
- Three hardware technologies
    - RoCE
    - iWARP
    - Infiniband
- Traditional protocols (SMB, NFS, iSCSI) can operate over RDMA

# Ethernet RDMA Stack



Blue content defined by the IBTA

Green content defined by IEEE / IETF

**Software**

**Typically Hardware**

| RDMA Application / ULP | | | |
|---|---|---|---|
| RDMA API (Verbs) | | | |
| RDMA Software Stack | | | |
| IB Transport Protocol | IB Transport Protocol | IB Transport Protocol | iWARP Protocol |
| IB Network Layer | IB Network Layer | UDP | TCP |
| | | IP | IP |
| IB Link Layer | Ethernet Link Layer | Ethernet Link Layer | Ethernet Link Layer |
| **InfiniBand** | **RoCE v1** | **RoCE v2** | **iWARP** |
| InfiniBand Management | Ethernet / IP Management | Ethernet / IP Management | Ethernet / IP Management |

# iWARP

- Delivers RDMA on top of Pervasive TCP/IP
- Runs over all Ethernet Infrastructure
- TCP provides Flow control and Congestion Management
- Highly routable and scalable Implementation
- Extensions eliminate TCP/IP stack process, mem copies and application contexts switches .
- iWARP addresses n/w bottlenecks of high speed Ethernet and provides high-throughput and low-latency with low-CPU utilization for data communication .

2019 Storage Developer Conference India ©  All Rights Reserved.

# RDMA over Converged Enhanced Ethernet (RoCE)

**Same RDMA, different L2 transport**

- Remote Direct Memory Access
    - Accelerates data exchange between servers
    - Bypass CPU & typical network stack
    - Reduced latency
- Converged Enhanced Ethernet
    - Priority Flow Control
    - Enhanced Transmission Selection
    - Lossless Ethernet fabric

2019 Storage Developer Conference India ©  All Rights Reserved.

# RoCE(V2)

- Well known on InfiniBand
- Works well on a lossless network
- Lower latency than alternative Transport protocols (TCP)
- Significantly lower overhead when offloaded to adapter

..BUT

- Ethernet is not lossless by design
- PFC is required to achieve lossless Ethernet fabric
- PFC (Part of DCB)has a high configuration and management overhead – VLANs, Priorities
- PFC is Layer 2 only

**SDC** 19
**SNIA INDIA**

**NetApp**

# RDMA  Pros and Cons

| Transport | Pros | | Cons |
|---|---|---|---|
| Non-RDMA Ethernet | | • TCP/IP-based protocol<br>• Works with any Ethernet switch<br>• Wide variety of vendors and models<br>• Support for in-box NIC teaming | • High CPU Utilization under load<br>• High latency |
| iWARP | Low CPU Utilization under load Low latency | • TCP/IP-based protocol<br>• Works with any Ethernet switch<br>• RDMA traffic routable<br>• Offers up to 100 Gbps per NIC port today* | • Requires enabling firewall rules |
| RoCE | | • Ethernet-based protocol<br>• Works with Ethernet switches<br>• Offers up to 100 Gbps per NIC port today*<br>• Routable with RoCEv2 | • Requires DCB switch with Priority Flow Control (PFC) |
| InfiniBand | | • Switches typically less expensive per port*<br>• Switches offer high speed Ethernet uplinks<br>• Commonly used in HPC environments<br>• Offers up to 54Gbps per NIC port today* | • Not an Ethernet-based protocol<br>• RDMA traffic not routable via IP infrastructure<br>• Requires InfiniBand switches<br>• Requires a subnet manager (typically on the switch) |

# RoCE vs IWARP differences

| | RoCE | iWARP |
|---|---|---|
| Underlying Network | UDP | TCP |
| Congestion Management | Rely on DCB | TCP handles with flow control |
| Adapter Offload | Full DMA | Full DMA w/TCP/IP |
| Routability | Yes | Yes |
| Cost | Need DCB enabled Switch Infra | Depends on the deployment , no requirement of Switch conf |

# Ethernet RDMA NIC Implementation

NVMe Native     NFS, CIFS, SAN

**Application**

NVMeoFabric

RDMA     FC     Ethernet

IB    RoCE    iWARP     FC

Ethernet/IP     Switches

**NVMe Target**

**Storage Compute**

HA PAIR

**NVMe Target**

**NVMe Initiator**

**NVMe Initiator**

Switches

NVMeoF

RDMA-ROCE

**Storage Drives**

Converged Shelves

NVMe SSDs

Integrated

NVMe SSDs

**RDMA NICs**

2019 Storage Developer Conference India © All Rights Reserved.

SDC 19
SNIA INDIA
16

NetApp

# PFC –Priority Flow Control

☐ By nature Ethernet is a lossy network

☐ Ethernet provides flow control mechanism which makes it lossless – 2 options:

• Applied FC over the whole port (Priority Flow Control - 802.3x)

• Applied FC over specific priority (Priority Flow Control -802.1Qbb)

☐ PFC negotiation between switch-host can be done by DCB (Data Center Bridging)

• Using Data Center Bridging Exchange (DCBX) negotiation

• End points (switch & host) exchange information about their capabilities

• If PFC is supported, it will be used

• If PFC is not supported, Global FC will be used

• If DCBX is not supported or the PFC capability is not supported, manual configuration is required

☐ Routers rebuild the layer 2 header

• Among it the routers rebuild the PCP filed using a DSCP to PCP mapping

SDC 19
SNIA INDIA

NetApp

# PFC contd..



PFC Priorities | Receive Buffers

| Default | 0 | 35% |
| Cluster | 1 | 10% |
| | 2 | |
| | 3 | |
| | 4 | |
| Storage | 5 | 40% |
| | 6 | |
| | 7 | |

Darryl
van der Peijl

https://www.darrylvanderpeijl.com/wp-content/uploads/part1_PFC.gif

# PFC-ETS

# RoCE for Long Distance

- ❒ Minimize the recovery impact from lost packets
    - ❒ Congestion, faulty networking components, alpha particles, etc.
- ❒ Congestion

  Can not use normal congestion control

    - PFC and ECN latency is too great because of distance

    - ❒ Solution options - Packet Pacing(NIC and Application) - Prioritize flows(QPs) through local networks(NIC and Switches)
        - ❒ ECN, PFC, other QOS
- ❒ Enhance recovery for lost packet
    - ❒ Resilient RoCE
    - ❒ Create a lot of small flows (Application)
    - ❒ Minimize the latency of retry

# Routable RoCE

- Routable RoCE requires a higher level congestion mechanism
    - ECN – Explicit Congestion Notification
- ECN can slow down traffic to prevent congestion
- ECN configuration overhead is lower than PFC, simple and easy



1. Switch detects onset of congestion
2. Switch marks packets' CE (Congestion Experienced) bit
3. Receiver sees CE bit on packets and notifies sender
4. Sender slows transmission rate until congestion risk is gone

Source:  Mellanox web

NetApp

# Resilient RoCE

- Resilient RoCE can cope with packet loss and Out of Order packets
- ECN is suggested but not required
- Out of Order packets are held in buffer to fill the gaps. Re-ordered packets are then written to memory
- Missing packets are requested from the sender

So..

- No loss – everything is fast
- Some loss – slows down, but stays in working order
- Still significantly better than TCP/IP

NetApp

# References

- [http://conferences.sigcomm.org/sigcomm/2015/pdf/papers/p523.pdf](http://conferences.sigcomm.org/sigcomm/2015/pdf/papers/p523.pdf)
- [https://community.mellanox.com/s/article/understanding-qos-configuration-for-roce](https://community.mellanox.com/s/article/understanding-qos-configuration-for-roce)
- [https://www.snia.org/sites/default/files/ESF/RoCE-vs.-iWARP-Final.pdf](https://www.snia.org/sites/default/files/ESF/RoCE-vs.-iWARP-Final.pdf)
- [http://files.gpfsug.org/presentations/2017/Manchester/04_Mellanox.pdf](http://files.gpfsug.org/presentations/2017/Manchester/04_Mellanox.pdf)

SDC 19
SNIA INDIA

NetApp

# CONCLUSION

- ☐ High-Speed Ethernet is the new back-bone which could replace FC
- ☐  Different media/storage via network require reliable connectivity with High throughput and low latency .
- ☐ High Availability and Disaster Recovery solutions are On-Demand with high data re-locational capabilities across geographies.
- ☐ Transports for NVMe over Fabric  with Ethernet is gaining momentum .

# Ethernet NIC

| CY17 | CY18 | CY19 | CY20 | CY21 | CY22 |
|------|------|------|------|------|------|

**Broadcom**
- Cumulus 2x25G gen3
- Stratus 1x100G gen3
- Thor 1x200G gen4
- Stingray 1x100G gen3
- Stingray 2 1x200G gen4

**Mellanox**
- ConnectX-5 2x100G gen4
- ConnectX-6 2x200G gen4
- Bluefield 2x100G gen4 x32
- Bluefield 2Lx 2x100G gen4 x16

**Chelsio**
- T6 2x100G gen3
- T7 2x100G gen4 RoCE + iWARP
- T8 2x200G gen4 RoCE + iWARP

**Intel**
- Columbiaville 2 x 100G gen4 RoCE + iWARP
- Mt. Stellar 2x200G 2x gen4 x16

**Cavium**
- Arrowhead 2/4x25G gen3
- Elbrus (E5) 2x100G (50G serdes) 1x200G (50G serdes) gen4 RoCE + iWARP
- F1 2x200G (50G serdes) 1x100G (25G serdes) gen4 RoCE + iWARP
- FastLinQ
- Liquid I/O III 1x100G gen4

NetApp

**Questions ?**

**n NetApp**