# Are Ethernet Attached SSDs Happening?

## NVMF-302B-1

Organizer/Chair: Rob Davis, Mellanox

Presenters:

Ilker Cebeli, Samsung

John Kloeppner, NetApp

Balaji Venkateshwaran, Toshiba

Khurram Milak, Netronome

Woo Suk Chung, SK Hynix

# Session Agenda

- **Ilker, Samsung – 15 minutes**
- **John - NetApp, Balaji -Toshiba, Khurram - Marvell – 40 minutes**
- **Woo, SK-Hynix – 15 minutes**
- **Q&A – 10 minutes**

# Are Ethernet Attached SSDs Happening?

## Disaggregated NVMe-oF Storage

Ilker Cebeli

Sr. Director of Planning

Samsung

# Disclaimer

# Data Center Evolution

## Traditional Data Center

**Manageability**

Console Console Console

**Compute** **Storage** **Networking**

**Stand Alone Component**
**Suited for Enterprise Applications**
**1GbE Networking**

## Hyper-converged Virtualized

**Manageability**

Console
Console
Console

VM VM VM
VM VM VM

**Resource Pools**

**Converged Management**
**Virtualized Computing/ Networking**
**10GbE Networking**

## Software-Defined Composable

**Manageability**

Console

VM VM
VM VM

**Rack Scale Software Defined**
**Disaggregated Compute and Storage**
**Composable**
**25-100GbE Networking**

**Evolution**

# Why Disaggregation?

## Converged

**Compute + Storage**



**NIC**  **CPU**  **DRAM**

**NVMe SSD**

**NETWORK**

## Disaggregated Compute & Storage

**Compute**

**Storage**



**NIC**  **CPU**  **DRAM**

**NETWORK**

**NVMe SSD**

❑ Pros:
  ✓ Scale Compute and Storage linearly
  ✓ Managed resources and storage services

❑ Cons
  ➤ Resources under-utilization
  ➤ Storage and Compute on the same network

❑ Pros:
  ✓ Compute and Storage scale independently
  ✓ Shared resources
  ✓ Improved utilization
  ✓ Grow as you go model based on workload demand
  ✓ Centralized storage services

❑ Cons
  ➤ Requires efficient storage protocols and latency
  ➤ Low latency and high bandwidth networking

# Some of the Use Cases for NVMe Over Fabrics



**NVMf Capable Network e.g. RDMA**

NVMf Bridge RNIC

CPU   CPU

Pcie Switch

Samsung NVMe SSD activated   Samsung NVMe SSD activated

NVMf Bridge RNIC

CPU   CPU

Pcie Switch

Samsung NVMe SSD activated   Samsung NVMe SSD activated

**Hyper-Converged
Most Common**

---

Server

CPU   CPU

PCie

Server

CPU   CPU

PCie

**NVMf Capable Network e.g. RDMA**

NVMf Target

Pcie Switch

Samsung NVMe SSD activated   Samsung NVMe SSD activated   Samsung NVMe SSD activated

Samsung NVMe SSD activated   Samsung NVMe SSD activated   Samsung NVMe SSD activated

**Disaggregated
JBOF Storage**

---

Server

CPU   CPU

PCie

NVMf Bridge RNIC

NVMf Target

Pcie Switch

Samsung NVMe SSD activated   Samsung NVMe SSD activated   Samsung NVMe SSD activated

Samsung NVMe SSD activated   Samsung NVMe SSD activated   Samsung NVMe SSD activated

**Direct Attached JBOF
SAS DAS Replacement**

# NVMe-oF JBOF



**2015 NVMe-oF JBOF**

Samsung All-Flash Array Reference Design

**NVMf Capable Network**

**2x100GbE ➜ ~24GB/s**

**Hyper-Converged Most Common**

**Typical NVMe-oF JBOF**

**24 NVMe SSDs ➜ ~24GB/s**

**2015 Platform Balanced Bandwidth between IO and NVMe**

# Future NVMe Bandwidth

**Throughput (GB/s)**
**Seq. Read 128KB**



| | |
|---|---|
| Future 24x Gen4x4 NVMe SSDs | FUTURE — 150+ |
| *2019 24x x4 Gen3 NVMe SSDs | 84 |
| 4x 100GbE | 44 |
| 4x 40GbE | 18 |
| 4x 10GbE | 4.5 |

**30x**

* 24x Samsung PM1725b NVMe SSDs (3.5GB/s throughput each)

## Network links could throttle the storage throughput performance

# Evolution of Networking Speeds and 25Gb/s and Above



Source: Crehan Long-range Forecast - Ethernet Adapter forecast, January 2019 via Mellanox Q2'2019

# IO Bottleneck



**CPU and IO bottleneck for storage throughput performance**

# NVMe-oF SSD based EBOF



## Conventional NVMe JBOF

Storage Head Nodes
Or Application Servers

**Ethernet Switch**

PCIe Switch
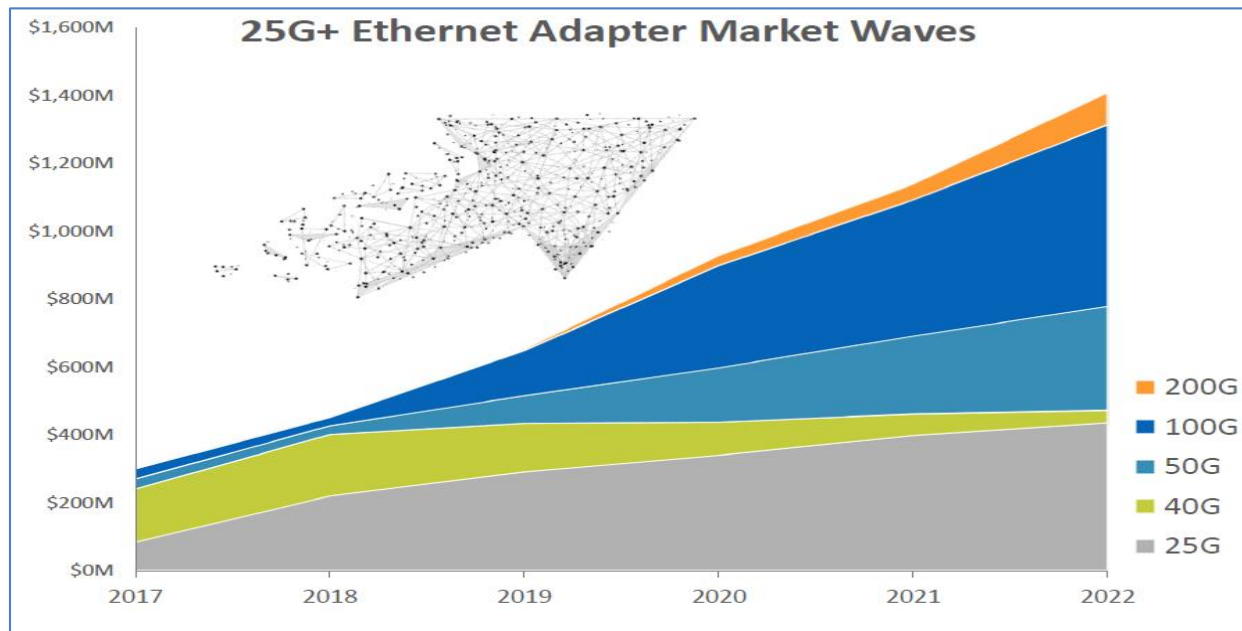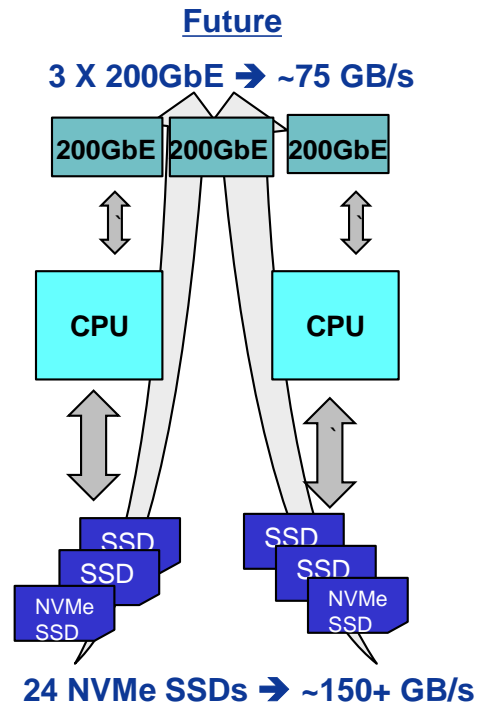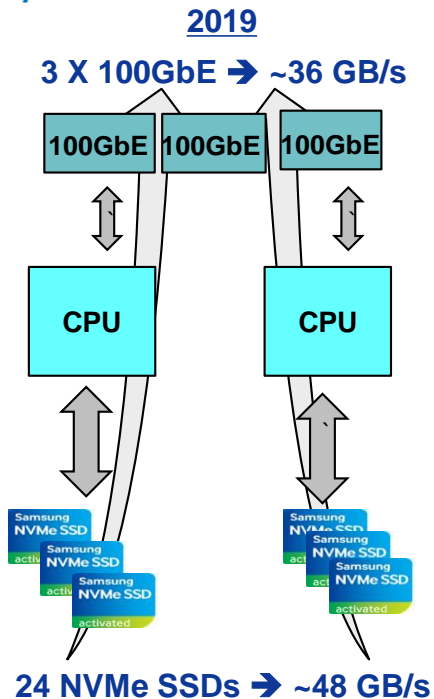
PCIe

| NVMe SSD | NVMe SSD | NVMe SSD | NVMe SSD |

**NVMe JBOF**

❑ **Pros**
- ✓ Enables disaggregation of NVMe SSDs
- ✓ Management & Storage Services
- ✓ Utilizing existing storage & server architectures

❑ **Cons**
- ➢ Non-scalable Storage Controller - PCIe single root constraint
- ➢ Bandwidth Limitation
  - CPU, PCIe, Networking Constraints
- ➢ Power and Thermals

## NVMe-oF EBOF

| Storage Controllers | App. Server Web | App. Server Database | App. Server Analytics |

**Hypervisor, Container**

**Ethernet Switch**

| Ethernet Switch | Ethernet Switch |

**Ethernet**

| NVMe-oF SSD | NVMe-oF SSD | ••• | NVMe-oF SSD | NVMe-oF SSD |

**NVMe-oF EBOF**

❑ **Pros**
- ✓ High Bandwidth
- ✓ Scaled Linearly (Ethernet)
- ✓ Sharable via NVMe-oF
- ✓ Less power
- ✓ Lower latency

❑ **Cons**
- ➢ New platform architecture
- ➢ Management of Storage Services & Network Devices

## NVMe-oF EBOF can address bandwidth, scalability, and flexibility

# Example Datacenter Storage Disaggregation



**Where Storage Services and Network Devices managed**

ilker.cebeli@samsung.com

Thank You

# Session Agenda

- **Ilker, Samsung – 15 minutes**
- ***John - NetApp*, Balaji -Toshiba, Khurram - Marvell – 40/3 minutes**
- **Woo, SK-Hynix – 15 minutes**
- **Q&A – 10 minutes**

# NVMe -> NVMe over Fabrics



**Servers with embedded NVMe Storage**

Servers without Storage

NVMe-oF

Data Management

NVMe-oF (direct connected)

JBOF

- **Local high performance / low latency access**
- **Isolated Storage**
- **Under-utilized SSD Performance and Capacity**

- **Shared Storage, better utilization of storage**
- **Similar NVMe Performance**

# Disaggregated Compute/ Data Management / Storage

- **Scaling Compute, Data Management and Storage Independently**
- **Full Shared Storage**

Compute

NVMe-oF

Storage Controller | Storage Controller | Storage Controller | Data Management

NVMe-oF

NVMe-oF JBOF

RNIC

PCIe Switch | CPU | Memory

NVMe SSDs

JBOF | JBOF | JBOF | Data Storage

**NetApp**

# NVMe-oF JBOF Limitations

- **Performance**
  - Throughput - PCIe Gen3 -> PCIe Gen4 -> PCIe Gen5, SCM, limit by existing infrastructure
  - Latency - Store and Forward architecture
- **Cost – CPU, SOC/RNICs, Switches, Mem don't scale well to match increasing SSD performance**

CPU based
NVMe-oF JBOF

RNIC

PCIe Switch — CPU — Memory

24 - x4 PCIe Gen4 NVMe SSDs

SOC RNIC based
NVMe-oF JBOF

SOC RNIC — Memory

PCIe Switch

24 - x4 PCIe Gen4 NVMe SSDs

**NetApp**

# Native Ethernet / NVMe-oF SSDs

- **Optimize NVMe-oF performance at SSD**
- **Options for NVMe-oF SSDs**

# Solution with Native NVMe-oF SSDs



- **Lower Latency, Higher Throughput**
- **Lower Cost and overall TCO**

**NVMe-oF JBOF**

Ethernet Switch

Management & Discovery

**24 Native NVMe-oF SSDs**

## JBOF price comparison (Excluding SSD cost)

| PCIe Gen 4 CPU Based* | PCIe Gen 4 SOC Based** | NVMe-oF 25G SSD Based*** |

## JBOF Price per Gbit of performance (Excluding SSD cost)

| PCIe Gen 4 CPU Based* | PCIe Gen 4 SOC Based** | NVMe-oF 25G SSD Based*** |

**\* Supports one 2x200G RNIC connected with x16 PCIe Gen4**

**\*\* Supports one 2x200G SOC RNIC connected with x16 PCIe Gen4**

**\*\*\* Supports three 200G Host connected Ethernet ports**

**NetApp**

# Additional Benefits

- Additional Benefits
  - Performance/cost scales with SSDs
  - Lower Power, reduced TCO
  - Including Ethernet switching within JBOF … potential to reduce networking cost, footprint, cabling

**■ NetApp**

# Other Activities

- Industry Standardization / Enablement
  - Standardization – Work underway in SNIA to define Form Factor, Pinout, Management – Toshiba will cover
  - Enablement – Fabrico Interposer – Marvell will cover

**NetApp**

Thanks!

# Session Agenda

- **Ilker, Samsung – 15 minutes**
- **John - NetApp, *Balaji -Toshiba*, Khurram - Marvell – 40/3 minutes**
- **Woo, SK-Hynix – 15 minutes**
- **Q&A – 10 minutes**

# Ethernet SSD-based Storage Platforms

App Server

CPU  DRAM

SSD SSD SSD SSD

**Block storage over the network**

## AFA software
**Scale for higher performance**

Ethernet Fabric

**Ethernet JBOF (EBOF)**

**Scale for larger capacity**
Ethernet Switch

SSD SSD SSD SSD SSD SSD SSD SSD

**Ethernet SSD**

## Advantages

– Independent scaling between performance (controller node) and capacity (JBOF) for optimal HW deployment in large scale systems

– Manage NVMe™ -based pools for separate storage/caching tiers

# Enabling NVMe-oF™ Functionality in SSDs

- Connector
  - SFF 8639 connector predominant for NVMe™-based systems
  - SFF-TA-1002 (EDSFF) specification a future-proof option
  - Standardizing Ethernet pinout in the connector a must for industry adoption

- Management Framework
  - NVMe™ devices attached to a system get enumerated using OS resources
  - Ethernet-attached device enumeration needs equivalent network functionality
  - Potential candidates for easier manageability:
    - NVMe-MI – from a BMC (not network)
    - RedFish – works for scalability in a Datacenter Network
    - RSD – uses RedFish

# Considerations in Connector Standardization

- Ethernet-based pinout should ensure:
  - SSDs of different types can be interchanged without electrical damage
    - First look in the VPD via SMBus, then apply power and signals
  - Forward-compatible
    - Connector of choice should support 25G → 50G → 100G transitions
    - Multi-lane for dual-port connectivity
  - Backwards-compatible
    - Ethernet pinout-based SSD should share midplane with SAS/SATA/PCIe pinouts
- Discovery of SSD:
  - Use standardized discovery mechanisms to obtain IP address, slot location
  - Discover and manage through RedFish
- Partnering to solve these challenges
  - Comprehensive standard specification in development in SNIA

# Management Frameworks for Ethernet SSDs

- Some administration will be done in-band via NVMe™ Admin commands once attached to a host

- But allocation and attachment needs to happen first at scale
    - Drive parameters and health monitoring
    - Encryption / Decryption key management
    - Host usage Authentication and Authorization
    - Logical assignment of drive resources on demand to multiple hosts

- NVMe™ functionality being mapped to RedFish management schema for these purposes

# Other Advanced NVMe™ Features

- Data Path Functionality
  - Zoned Name Space Support
  - Key Value namespaces
  - Endurance Group /  NVM Set / Namespace Management
  - Future Computational Storage platform for FPGA, Accelerators, etc.

- Part of a Composable Infrastructure
  - Storage "stack" assembled on demand tailored to application needs
  - Drawn from pools of Ethernet Drives, then returned to the pool when finished

# World's First True Ethernet NVMe-oF™ SSD

**In-Form Factor Native NVMe-oF™ SSD (Ethernet SSD)**

- Standard 2.5" In-Form Factor
- No external components needed
- SFF 8639 / 9639 standardized connector with Ethernet pinout
- Dual-port 25Gbit Ethernet
- RDMA over Converged Ethernet ver. 2 (RoCEv2)
- 675K IOPS @ 4KB Random Read
  - Equivalent performance to PCIe® Gen3x4

# Visit the Toshiba Memory FMS Booth #307

**2.5" Ethernet SSD Prototype**

**Example of Ethernet SSD-based AFA architecture**

Demonstration of 2.5" Ethernet SSD prototype with native NVMe-oF™ support

Prototype of a possible AFA platform using EBOF (Ethernet SSD-based)

# Thanks!

# Session Agenda

- **Ilker, Samsung – 15 minutes**
- **John - NetApp, Balaji -Toshiba, *Khurram - Marvell* – 40/3 minutes**
- **Woo, SK-Hynix – 15 minutes**
- **Q&A – 10 minutes**

# Native NVMe-oF SSD

## Khurram Malik

## Sr. Product Marketing Manager, Marvell

# Current Challenges with NVMe-oF

- ## SSD Industry is Diverging:
  - Different interfaces (SATA, SAS, PCIe)
  - Different protocols/transports (NVMe-oF variants; NVMe; SCSi …)
  - Different form factor (U.2, U.3, EDSFF S, ESDFF L, EDSFF 3")

- ## Challenges:
  - **Standards are diverging instead of converging.**
  - **No clear direction which standard will eventually win.**
  - Selecting a right standard and enable NVMe-oF SSD.
  - Managing two different SSDs skews; NVMe and NVME-oF
  - **Managing two different midplanes; PCIe (NVMe) & Ethernet (NVMe-oF)**
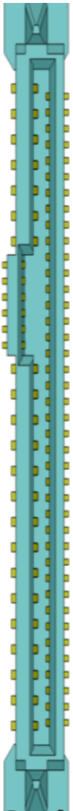  - **Designing a new chassis to use NVMe-oF SSDs.**

# OCP Kinetic and SNIA Ethernet Drive Pins

| | | SATA | SATA Express | SAS | MultiLink SAS | Quad PCIe | USB | OCP Kinetic | SNIA Ethernet Drive |
|---|---|---|---|---|---|---|---|---|---|
| S1 | Ground | GND | GND | GROUND | GROUND | Ground | GND | Ground | Ground |
| S2 | Rcvr+ | A+ | PETp0 | PR+ | RX0+ | | SSRX+ | RX0+ | RX0+ |
| S3 | Rcvr- | A- | PETn0 | PR- | RX0- | | SSRX- | RX0- | RX0- |
| S4 | Ground | GND | GND | GROUND | GROUND | Ground | GND | Ground | Ground |
| S5 | Xmtr- | B- | PERn0 | TP- | TX0- | | SSTX- | TX0- | TX0- |
| S6 | Xmtr+ | B+ | PETR0 | TP+ | TX0+ | | SSTX+ | TX0+ | TX0+ |
| S7 | Ground | GND | GND | GROUND | GROUND | Ground | GND | Ground | Ground |
| S8 | Ground | | GND | GROUND | GROUND | Ground | | Ground | Ground |
| S9 | Rcvr+ | | PETp1 | SR+ | RX1+ | | | RX1+ | RX1+ optional |
| S10 | Rcvr- | | PETn1 | SR- | RX1- | | | RX1- | RX1- optional |
| S11 | Ground | | GND | GROUND | GROUND | Ground | | Ground | Ground |
| S12 | Xmtr- | | PERn1 | ST+ | TX1- | | | TX1- | TX1- optional |
| S13 | Xmtr+ | | PERp1 | ST- | TX1+ | | | TX1+ | TX1+ optional |
| S14 | Ground | | GND | GROUND | GROUND | Ground | | Ground | Ground |

- U2 OCP Kinetic and SNIA Ethernet Drive pin assignments induce crosstalk between adjacent TX and RX pairs, which reduce the max supported channel length. Therefore we recommend different differential pin assignments for 25Gbps PAM2 or 50Gbps PAM4 two Lanes Ethernet application.

37

# U.2 connector pin assignment for Ethernet application



Fig1. U.2 pin assignment

**Notes:**

**Marvell has recommended two high speed signal pin assignment proposals for Ethernet application to minimize connector impacts on the overall Channel Operating Margin(COM).**

- **Proposal1: Maximize the distance from one differential pair to other signals; (Highlighted as red column)**
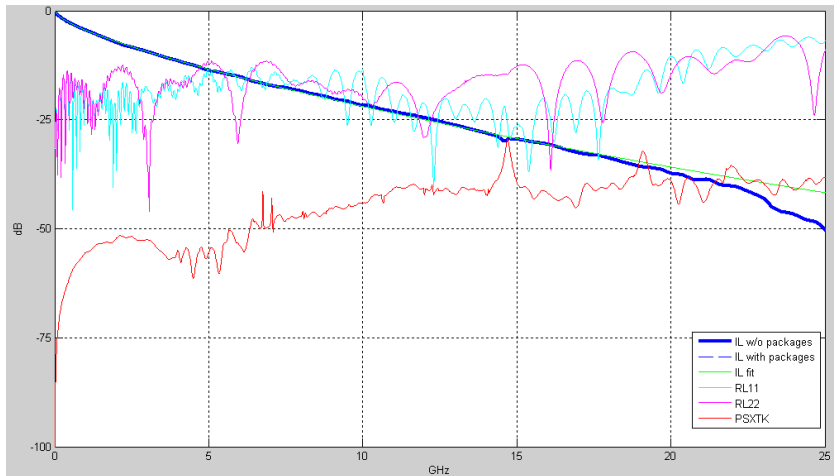- **Proposal2: Based on proposal1 concept, keep pin compatible with PCIe signals. (Highlighted as blue column)**

# MRVL COM simulation Setup and Results

Based on below long lossy channel, run end to end COM/ERL simulation with two proposed U.2 pin configurations.

- IEEE 802.3by 25GBASE-KR Channel Operating Margin(COM>3dB) without FEC.
- IEEE 802.3bs 50GBASE-KR Channel Operating Margin(COM>3dB,ERL>10dB)



| Operation mode | U.2 Pin proposal1 (SAS & Ethernet Signals) Proposal1 | | U.2 Pin proposal2 (PCIe & Ethernet Signals) Poposal2 | |
|---|---|---|---|---|
| | COM(dB) | ERL(dB) | COM(dB) | ERL(dB) |
| 25Gbps PAM2 | 3.52 | NA | 3.65 | NA |
| 50Gbps PAM4 | 3.25 | 14.08 | 3.20 | 14.18 |

# Convert NVMe SSD to NVMe-oF SSD



NVMe-oF Converter Controller interposer in a carrier



NVMe-oF Converter Controller  Interposer (SSD side)



NVMe-oF Converter Controller Interposer (network side)
*(*8639 is used to drive 2x25Gb Ethernet)*



NVMe-oF Converter Controller Interposer (profile)
Connected to U.2 (non-carrier)

# Enabling NVMe-oF

**Simple, low RBOM, low power backplane**

# Enabling NVMe-oF

- Marinating NVMe and NVMe-oF support
  - NVMe
  - NVMe-oF : ROCEv2 ; TCP

- NVMe-oF Converter Controller
  - Can fit interposer
  - Can fit inside U.2/EDSFF
  - Can be merged with SSD Controller

- Re use of backplane
  - Re use 8639/9639
  - No changes to mid plane
  - Swap IOM

- No extra enclosure expense (other than IOM)

- Single SSD can work both PCIe and Ethernet (Better inventory management)

# Thanks!

# Session Agenda

- **Ilker, Samsung – 15 minutes**
- **John - NetApp, Balaji -Toshiba, Khurram - Marvell – 40 minutes**
- **<u>Woo, SK-Hynix – 15 minutes</u>**
- **Q&A – 10 minutes**

- Woo slides

# Thanks!

# Session Agenda

- **Ilker, Samsung – 15 minutes**
- **John - NetApp, Balaji -Toshiba, Khurram - Marvell – 40 minutes**
- **Woo, SK-Hynix – 15 minutes**
- **<u>Q&A – 10 minutes</u>**

Q/A

# Thanks!