

# **OPEN** Compute Engineering Workshop March 9, 2015 San Jose



# Mellanox OCP Contribution Overview

# OCP Engineering Workshop, Mar 2015 Kevin Deierling

# Vice President, Mellanox Technologies

Mellanox Overview

# Public company (Nasdaq: MLNX)

- Founded in 1999
- ~1700 employees

# End to end provider of interconnect solutions Leading provider of InfiniBand & Ethernet solutions

- Silicon, Adapters/NICs, Switches
- State of the art Ethernet solutions
  - Dominant provider of 40GbE adapters
- Cables and QSFP Modules - Copper, VCSEL Optics, Silicon Photonic

# Mellanox and Open Compute Project

## 2011

Contribution of 10GbE OCP ConnectX-3 specification

## 2012

- UEFI supports for OCP platform
- Mechanical adjustment of OCP NIC to triplet servers

### 2013

Enablement of Shared NIC functionality with ConnectX-3 device

### 2014

- 1'st 40GbE OCP Mezz card with ConnectX-3 PRO
- Switch design specification contribution
  - 36x 40GbE or 48x 10G and 12x 40G switch
- ONIE and SAI contributions

### 2015



### 1st 25GbE, 50GbE and 100GbE Network Adapter Multi-Host design contributions and announcement OpenOptics WDM technical contribution



# Exponential Data Growth – Requires Platform Innovation

# We Live in a World of Data



# Data Needs to be Accessible Always and in Real-Time





### **More Data**

# Big, Fast, Real-Time Data Needs Innovation Too!











# Data Center Evolution Over Time

tot 1



# From Compute Centric to Data Centric Data Center (DCDC)

### **Compute-centric architecture**

- CPU at the center with attached peripherals
- Developed for transactional processing
  - Small, slow, fixed-format data
- Data is an afterthought!
- Not equipped for Big-Fast-Unstructured Data

## Focus is on server-level optimization

- Compute-centric optimization focus is the server
- Secondary focus is the storage chassis

## A higher level view is huge advantage!

- From compute to data centric architecture
- Explicitly considers Big-Fast-Unstructured Data
- Higher efficiency and better CapEx and OpEx



## **Compute Centric Center Architecture**

### Networking

**I/O** 

### Storage

# CPU

### Compute

# Data Centric Data Center Enables Rack Level Optimization



### Data centric view allows rack and data center level optimization

- Higher level means better optimization possible than at the server level
- Disaggregation of server resources
  - Allows sharing of CPU, storage, memory, and I/O resources
  - Flexibility to scale each element independently

Intelligent interconnect is at the heart of this transformation

### Storage



# **Rack & Data Center Level Optimization**

# **Rack Level Optimization** For The Data Centric Data Center



Three key requirements for rack & data center level optimization

- Server disaggregation enabled by network sharing 1.
- 2. Efficient data movement with RDMA and virtualization acceleration
- 3. High speed data connectivity100Gb/s copper, optical, & silicon photonics

# Server Disaggregation Enabled by Resource Sharing



- Enables efficient sharing of network & compute resources & efficient data transport
- Single socket CPU significantly reduces costs

1.1

Symmetric data access means deterministic performance under all circumstances

# Efficient Data Movement: eSwitch, RDMA, Network Virtualization





**Embedded Switch** Hardware OVS Switch Virtual Overlay Network Acceleration



## Efficient Data Handling

- eSwitch: Embedded hardware OVS switch flow steering engine
- 2. Virtual network acceleration
- 3. RDMA Efficient Data Exchange Low Latency, Low CPU Overhead

### Efficient Data Movement With RDMA

# **OpenOptics WDM Contribution to OCP**



End to end connectivity allows innovation and optimization at rack and data center level

- Standard QSFP form factors for adapters, switches, and cables
- Copper & fiber cable, single & multi-mode, VCSEL & silicon photonics
  - Use the best technology
- Advanced silicon photonics platform offers a future-proof roadmap
  - WDM, higher speed modulation





# **ConnectX-4 Multi-Host Contribution**

# **ConnectX-4**

# Multi-Host Technology





# Mellanox ConnectX-4 100GbE Network Adapter

# ConnectX-4: Highest Performance Adapter in the Market

InfiniBand: QDR, FDR, EDR

Ethernet: 10 / 25 / 40 / 50 / 56 / 100GbE

100Gb/s, <0.7us latency

150 million messages per second (8B)

100 million packets per second (64B)

RDMA, RoCE

Multi-Host technology

CORE-Direct technology

**GPUDirect RDMA** 

**Overlay Networks offload** 







### ConnextX-4 OCP Adapter with **Multi-Host Technology**

# Mellanox ConnectX-4 100GbE Network Adapter

# ConnectX-4: Highest Performance Adapter in the Market

OCP 2.0 compliant

Automatic self-configuring for legacy and Multi-Host platforms

Host management support

- NC-SI compliant multi-instance BMC
- MCTP (MCTPoSMBus and MCTPoPCIe) compliant multiinstance **BMC**
- Dedicated management resources per each managed host Supports:
- Single 16x Gen 3.0 PCle
- Dual 8x Gen 3.0 PCIe
- Quad 4x Gen 3.0 PCIe









### ConnextX-4 OCP Adapter with **Multi-Host Technology**

# Mellanox Multi-Host<sup>™</sup> Technology

the the



### Multi-Host Technology

# New Compute Rack / Data Center Architecture



### **Scalable Data Center with Multi-Host**

- The future of data center design •
- •
- •
- •
- •

### **Traditional Data Center**

- Expensive design for scalable data centers
- Requires many ports on ToR switch
- Dedicated NIC / cable per server



Modular, share components CPU & NIC Flexible, configurable, application optimized **Optimized networking configurations** Enabled by high-throughput network

# New Scalable Rack Design



### Smart Interconnect to Unleash the Power of All Compute Architectures



**Complete Architectural Flexibility** 

**Highest Performance and Scalability** X86, Power, GPU, ARM and FPGA-based **Compute and Storage Platforms** 





10/25/40/50/100 Gb/s

Smart Interconnect Enables Single Platform to Support Broad Range of Workloads



# Higher Performance Data Center

# **Overcoming CPU-to-CPU Connectivity Bottlenecks Lower Application Latency**



### Smart Interconnect for a High Variety of Compute Architectures



# Multi Host and BMC Management

### Single BMC for Multi-host

- Most efficient and lowest cost solution. A single BMC operates as multiple instances. Each instance controls a separate host. Each instance has dedicated management resources on ConnectX-4.

### Dedicated BMC for each host

- Easiest migration to Multi-host platforms. Each BMC has dedicated management resources on ConectX-4

### Chassis MGMT BMC

 Additional BMC for managing common chassis resources. Chassis manager also has dedicated management resources on ConnectX-4.

### Independent medium migration

• Each BMC (or BMC instance) is fully independent. Each BMC can migrate between mediums (RBT, SMBus or PCIe) independently – allows optimizing management resources and system power for each host

### Multiple BMCs sharing a common interface

- Common medium can be shared by multiple BMCs (or BMC instances) - allows minimizing the number of connections



# Summary

Next generation compute and storage rack design

- Multi-host allows server disaggregation and resource sharing
- **Enables Data Centric Data Center (DCDC)**
- Enables scalable high performance, Cloud and Web 2.0 data centers
- Enables rack level optimization full disaggregation of system elements
- Complete flexibility: x86, ARM, Power, GPU data centers
- Independent Host management, individual QoS per server
- Reduces cabling reducing cost, easing deployment, simplifying maintenance
- **Reduces switch port count**

### ConnextX-4 OCP Adapter with Multi-Host Technology





# Thank You



# **OPEN** Compute Engineering Workshop March 9, 2015 San Jose