

### Paving the Road to Exascale

**Gilad Shainer** 

August 2015, MVAPICH User Group (MUG) Meeting



# Mellanox Connect. Accelerate. Outperform."

### The Ever Growing Demand for Performance

### Performance



#### **Technology Development**



SMP to Clusters



**Single-Core to Multi-Core** 





### The Road to Exascale Computing





# Co-Design Architecture – From Discrete to System Focused













### Exascale will be Enabled via Co-Design Architecture



# Standard, Open Source, Eco-System Programmable, Configurable, Innovative





## Software-Hardware Co-Design? Example: Breaking the Latency Wall



- Today: Network devices are in 100ns latency today
- Challenge: How to enable the next order of magnitude performance improvement?
- Solution: Co-Design mapping the communication frameworks on all active devices
- Result: reduce HPC communication frameworks latency by an order of magnitude

### **Co-Design Architecture Paves the Road to Exascale Performance**

#### © 2015 Mellanox Technologies

![](_page_5_Picture_8.jpeg)

![](_page_5_Picture_9.jpeg)

6

### The Road to Exascale – Co-Design System Architecture

- The road to Exascale requires order of magnitude performance improvements
- Co-Design architecture enables all active devices to become co-processors

![](_page_6_Figure_3.jpeg)

![](_page_6_Picture_6.jpeg)

7

## **Co-Design Architecture**

![](_page_7_Figure_1.jpeg)

![](_page_7_Picture_3.jpeg)

### The Elements of the Co-Design Architecture

![](_page_8_Figure_4.jpeg)

# **Co-Design Implementation Via Offloading Technologies**

© 2015 Mellanox Technologies

![](_page_8_Picture_7.jpeg)

9

### Mellanox Co-Design Architecture (Collaborative Effort)

#### **Communication Frameworks (MPI, SHMEM/PGAS)**

![](_page_9_Figure_2.jpeg)

#### The Only Approach to Deliver 10X Performance Improvements

![](_page_9_Picture_5.jpeg)

![](_page_9_Figure_6.jpeg)

### Mellanox InfiniBand Proven and Most Scalable HPC Interconnect

![](_page_10_Picture_1.jpeg)

# Paving the Road to Exascale

![](_page_10_Picture_4.jpeg)

![](_page_10_Picture_5.jpeg)

![](_page_10_Picture_7.jpeg)

# High-Performance Designed 100Gb/s Interconnect Solutions

![](_page_11_Figure_1.jpeg)

![](_page_11_Picture_3.jpeg)

# InfiniBand Adapters Performance Comparison

Mellanox Adapters Single Port Performance	ConnectX-4 EDR 100G	Connect-IB FDR 56G
Uni-Directional Throughput	100 Gb/s	54.24 Gb/s
<b>Bi-Directional Throughput</b>	195 Gb/s	107.64 Gb/s
Latency	0.61 us	0.63 us
Message Rate (Uni-Directional)	149.5 Million/sec	105 Million/sec

![](_page_12_Picture_3.jpeg)

## ConnectX-3 Pro FDR 56G

#### 51.1 Gb/s

#### 98.4 Gb/s

#### 0.64 us

#### 35.9 Million/sec

### **EDR InfiniBand Performance – Commercial Applications**

**OptiStruct Performance** (Engine\_Assy.fem)

![](_page_13_Figure_2.jpeg)

FDR InfiniBand EDR InfiniBand

![](_page_13_Figure_4.jpeg)

**LS-DYNA Performance** 

(neon\_refined\_revised)

FDR InfiniBand

![](_page_13_Picture_7.jpeg)

![](_page_13_Figure_8.jpeg)

FDR InfiniBand EDR InfiniBand

![](_page_13_Picture_12.jpeg)

#### **RADIOSS 13.0 Performance** (NEON1M11, MPP)

Number of Nodes

EDR InfiniBand

![](_page_13_Picture_16.jpeg)

## **EDR InfiniBand Performance – Weather Simulation**

- Weather Research and Forecasting Model
- Optimization effort with the HPCAC
- EDR InfiniBand delivers 28% higher performance
  - 32-node cluster
  - Performance advantage increase with system size

![](_page_14_Figure_6.jpeg)

FDR InfiniBand

![](_page_14_Picture_8.jpeg)

![](_page_14_Picture_9.jpeg)

THE WEATHER RESEARCH & FORECASTING MODEL

<sup>2000</sup> Report to the second s second sec

![](_page_14_Picture_11.jpeg)

![](_page_14_Picture_13.jpeg)

![](_page_14_Picture_15.jpeg)

#### WRF Performance (conus12km)

EDR InfiniBand

### Lenovo EDR InfiniBand System (TOP500)

- "LENOX" EDR InfiniBand connected system at the Lenovo HPC innovation center
- EDR InfiniBand provides >20% higher performance versus over FDR on Graph500
  - At 128nodes

![](_page_15_Picture_4.jpeg)

![](_page_15_Picture_5.jpeg)

![](_page_15_Picture_7.jpeg)

![](_page_15_Picture_8.jpeg)

![](_page_16_Picture_0.jpeg)

# Unified Communication – X Framework (UCX)

The Next Generation HPC Software Framework To Meet the Needs of Future Systems / Applications

![](_page_16_Picture_3.jpeg)

![](_page_16_Picture_4.jpeg)

**Exascale Co-Design Collaboration** 

# **The Next Generation**

# **HPC Software Framework**

# Collaborative Effort Industry, National Laboratories and Academia

![](_page_17_Picture_4.jpeg)

![](_page_17_Picture_6.jpeg)

Lawrence Livermore National Laboratory

# **UCX Framework Mission**

- Collaboration between industry, laboratories, and academia
- Create open-source production grade communication framework for HPC applications
- To enable the highest performance through co-design of software-hardware interfaces
- To unify industry national laboratories academia efforts

![](_page_18_Figure_5.jpeg)

#### **Co-design of Exascale Network APIs**

![](_page_18_Picture_8.jpeg)

#### cations erfaces

#### **Production quality**

Developed, maintained, tested, and used by industry and researcher community

#### **Cross platform**

Support for Infiniband, Cray, various shared memory (x86-64 and Power), GPUs

### The UCX Framework

#### UC-S for Services

This framework provides basic infrastructure for component based programming, memory management, and useful system utilities

Functionality: Platform abstractions and data structures

#### UC-T for Transport

Low-level API that expose basic network operations supported by underlying hardware

Functionality: work request setup and instantiation of operations

High-level API uses UCT framework to construct protocols commonly found in applications

Functionality: Multi-rail, device selection, pending queue, rendezvous, tag-matching, software-atomics, etc.

![](_page_19_Picture_11.jpeg)

#### UC-P for Protocols

# UCX High-level Overview

![](_page_20_Figure_1.jpeg)

![](_page_20_Picture_3.jpeg)

# Collaboration

![](_page_21_Picture_1.jpeg)

- Mellanox co-designs network interface and contributes MXM technology
  - Infrastructure, transport, shared memory, protocols, integration with OpenMPI/SHMEM, MPICH

![](_page_21_Picture_4.jpeg)

ORNL co-designs network interface and contributes UCCS project • InfiniBand optimizations, Cray devices, shared memory

![](_page_21_Picture_6.jpeg)

NVIDIA co-designs high-quality support for GPU devices • GPU-Direct, GDR copy, etc.

![](_page_21_Picture_8.jpeg)

IBM co-designs network interface and contributes ideas and concepts from PAMI

![](_page_21_Picture_10.jpeg)

![](_page_21_Picture_11.jpeg)

![](_page_21_Picture_13.jpeg)

Lawrence Livermore National Laboratory

# UCX Performance

![](_page_22_Figure_1.jpeg)

![](_page_22_Figure_2.jpeg)

![](_page_22_Picture_4.jpeg)

![](_page_23_Picture_0.jpeg)

# Mellanox Multi-Host<sup>™</sup> Technology Next Generation Data Center Architecture

![](_page_23_Picture_2.jpeg)

### New Compute Rack / Data Center Architecture

![](_page_24_Figure_1.jpeg)

The Network is The Computer

![](_page_24_Picture_4.jpeg)

## Multi-Host Dramatically Reduces Server Cost

![](_page_25_Figure_1.jpeg)

Proprietary **Cache Coherent** Bus

> Modern CPUs with 8-20 cores don't require expensive SMP architectures.

> Additional parallelism achieved with higher level network based distributed programming techniques such as Hadoop Map-Reduce

#### **Multi-Host 4-Socket Architecture**

![](_page_25_Picture_6.jpeg)

#### Expensive 4-Way CPU

- Massive but unused cache-coherent domain
- High overhead but un-necessary CPU bus
  - High pin count ,high power, complex layout
- Asymmetric (NUMA) of data access

- Low cost single-socket CPU
  - Clean, simple, cost-effective, software transparent
- Cache coherent domain: Multi-Core CPU
  - Eliminates pins, Lower power, Simpler layout
- Symmetric Data Access

![](_page_25_Picture_18.jpeg)

## ConnectX-4 on Facebook OCP Multi-Host Platform (Yosemite)

![](_page_26_Picture_1.jpeg)

![](_page_26_Picture_2.jpeg)

![](_page_26_Picture_3.jpeg)

![](_page_26_Picture_4.jpeg)

The Next Generation Compute and Storage Rack Design

© 2015 Mellanox Technologies

![](_page_26_Picture_7.jpeg)

![](_page_26_Picture_8.jpeg)

#### **Compute Slots**

![](_page_27_Picture_0.jpeg)

# Summary

# Paving the Road to Exascale Computing

![](_page_27_Picture_3.jpeg)

### End-to-End Interconnect Solutions for All Platforms

#### **Highest Performance and Scalability for**

#### X86, Power, GPU, ARM and FPGA-based Compute and Storage Platforms

#### 10, 20, 25, 40, 50, 56 and 100Gb/s Speeds

![](_page_28_Figure_4.jpeg)

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

![](_page_28_Picture_7.jpeg)

### Technology Roadmap – One-Generation Lead over the Competition

![](_page_29_Figure_1.jpeg)

![](_page_29_Picture_3.jpeg)

![](_page_29_Picture_4.jpeg)

![](_page_30_Picture_0.jpeg)

# Thank You

![](_page_30_Picture_2.jpeg)

# Mellanox Connect. Accelerate. Outperform.™