# Solros: A Data-Centric Operating System Architecture for Heterogeneous Computing

**Changwoo Min**, Woonhak Kang, Mohan Kumar, Sanidhya Kashyap, Steffen Maass, Heeseung Jo, Taesoo Kim

Virginia Tech, eBay, Georgia Tech, Chonbuk National University

April 26, 2018

Specialization of general-purpose processors

Specialization

Generalization of co-processors

Specialization of processors

Generalization of accelerators

Specialization of co-processors

Blazingly fast storage/memory

# Blazingly fast IO Devices



**Intel's new Optane memory drives will turn your PC into a beast**

Jamie McKane   10 April 2017   33 Comments

67 shares

Intel's cutting-edge Optane memory boasts ma
increases over solid state drives.

HPC wire

Since 1987 - Covering the Fastest Computers in the World and the People Who Run Them

- ⊘ Home
- ⊘ Technologies
- ⊘ Sectors
- ⊘ Exascale
- ⊘ Specials
- ⊘ Resource Library
- ⊘ Events
- ⊘ Job Bank
- ⊘ About

FULL D COMPLETE COVERAGE
Keynote Reviews, Analysts Write Ups, Booth Vides, Student Competition, Awards and so much more

SC17

**Ahead of SC17, Mellanox Launches Scalable 200G Switch Platforms**
By Doug Black

November 9, 2017

In the run-up to the annual supercomputing conference SC17 next week in Denver, Mellanox made a series of announcements today, including a scalable switch platform based on its HDR 200G InfiniBand technology and the first deployment of a 100Gb/s Linux kernel-based Ethernet switch.

The company touts its HDR (High Data Rate) 200G InfiniBand Quantum, which offers up to 800 ports of 200Gb/s or 1,600 ports 100Gb/s in one chassis, as the most scalable switch platform available.

The platform family includes:

- Quantum QM8700: 40-port 200Gb/s or 80-port 100Gb/s
- Quantum CS8510: modular 200-port 200Gb/s or 400-port 100Gb/s
- Quantum CS8500: modular 800-port 200Gb/s or 1,600-port 100Gb/s

Mellanox said the Quantum product line's switch density will enable space and power consumption optimization, reducing network equipment cost by 4X, electricity costs by 2X and improving data transfer time by 2X.

Blazingly fast storage/m

Blazingly fast network

# Blazingly fast IO Devices

**Intel's new Optane memory drives will turn your PC into a beast**

Jamie McKane  10 April 2017  33 Comments

HPC wire

Since 1987 : Covering the Fastest Computers
in the World and the People Who Run Them

> How to exploit the full potential of such hardware devices without pain?
> - System-wide performance
> - Ease of programming

67 shares

Intel's cutting-edge Optane memory boasts ma
increases over solid state drives.

**FULL D COMPLETE COVERAGE**
Keynote Reviews, Analysts Write Ups,
Booth Video, Student Competition,
Awards and so much more

SC17

The platform family includes:

- Quantum QM8700: 40-port 200Gb/s or 80-port 100Gb/s
- Quantum CS8510: modular 200-port 200Gb/s or 400-port 100Gb/s
- Quantum CS8500: modular 800-port 200Gb/s or 1,600-port 100Gb/s

Mellanox said the Quantum product line's switch density will enable space and
power consumption optimization, reducing network equipment cost by 4X,
electricity costs by 2X and improving data transfer time by 2X.

Blazingly fast storage/r

Blazingly fast network

# Outline

# Host-Centric Approach

- Host OS controls co-processors and IO devices
- Examples: OpenCL, CUDA

# Host-Centric Approach

- Host OS controls co-processors and IO devices
- Examples: OpenCL, CUDA

# Host-Centric Approach

- Host OS controls co-processors and IO devices
- Examples: OpenCL, CUDA

# Host-Centric Approach

- Host OS controls co-processors and IO devices
- Examples: OpenCL, CUDA

# Host-Centric Approach

- Host OS controls co-processors and IO devices
- Examples: OpenCL, CUDA



## Problem

Redundant data communication
Complex to program and hard to optimize

# Coprocessor-Centric Architecture

- Co-processors control IO devices
- Examples: Xeon Phi (Linux), GPUfs [ASPLOS13], GPUNet [OSDI14]

# Coprocessor-Centric Architecture

- Co-processors control IO devices
- Examples: Xeon Phi (Linux), GPUfs [ASPLOS13], GPUNet [OSDI14]

# Coprocessor-Centric Architecture

- Co-processors control IO devices
- Examples: Xeon Phi (Linux), GPUfs [ASPLOS13], GPUNet [OSDI14]

# Coprocessor-Centric Architecture

- Co-processors control IO devices
- Examples: Xeon Phi (Linux), GPUfs [ASPLOS13], GPUNet [OSDI14]



## Problem

Significant effort required for porting IO stack to co-processor
Not completely exploiting powerful host processors

# Outline

# Solros Goal

- Ease of programming
- Best use of processor architecture
- System-wide optimization

# Solros Goal

- Ease of programming
- Best use of processor architecture
- System-wide optimization

## Challenge

- Co-processor needs IO abstraction
- IO stacks is branch-divergent and difficult to parallelize
- It needs system-wide information

# Solros Architecture

## Split-Kernel Architecture

- Data-plane OS
  - Runs on a co-processor
  - Provides IO abstraction
  - Delegates actual IO operations to a control-plane OS
- Control-plane OS
  - Runs on a host processor
  - Runs actual IO stack
  - Performs system-wide coordination

# Solros Architecture

- Control-plane OS: actual OS service + system-wide coordination
- Data-plane OS: thin communication layer to host processor

# Solros Architecture

- Control-plane OS: actual OS service + system-wide coordination
- Data-plane OS: thin communication layer to host processor

# Solros Architecture

- Control-plane OS: actual OS service + system-wide coordination
- Data-plane OS: thin communication layer to host processor

# Solros Architecture

- Control-plane OS: actual OS service + system-wide coordination
- Data-plane OS: thin communication layer to host processor

# Solros Architecture

- Control-plane OS: actual OS service + system-wide coordination
- Data-plane OS: thin communication layer to host processor



- Co-processor has OS abstraction with minimal effort
- Best use of each of the fat and lean processors
- Efficient global coordination among devices (policy)

# Operating System Services

1. Transport service
2. Filesystem service
3. Network service

# Operating System Services

1. Transport service
2. Filesystem service

3. Network service

# Transport Service

High performance data transfer among devices are challenging:

- Uniform data transfer among devices
- High contention in massively-parallel co-processor
- Asymmetric performance between host processor and co-processor

# Transport Service

High performance data transfer among devices are challenging:

- Uniform data transfer among devices
- High contention in massively-parallel co-processor
- Asymmetric performance between host processor and co-processor

## Our approach

- Uniform data transfer $\Rightarrow$ system-mapped PCIe window
- High contention $\Rightarrow$ combining, replication, interleaving, etc.
- Asymmetric performance $\Rightarrow$ flexibly configurable (host DMA engine vs. co-processor DMA engine)

# Transport Service

High performance data transfer among devices are challenging:

- Uniform data transfer among devices
- High contention in massively-parallel co-processor
- Asymmetric performance between host processor and co-processor

## Our approach

- Uniform data transfer $\Rightarrow$ system-mapped PCIe window
- High contention $\Rightarrow$ combining, replication, interleaving, etc.
- Asymmetric performance $\Rightarrow$ flexibly configurable (host DMA engine vs. co-processor DMA engine)

See details in the paper

# Filesystem Service

- **Peer-to-peer operation**
- Buffered operation



Co-processor                        Host processor

| Application |                      | File system proxy |

          ① (File system stub)                    | File system |

PCIe

| DMA engine |      ·····▶ control
|    SSD     |      ——▶ data

Zero-copy of data between co-processor memory and SSD
Minimal data transfer

# Filesystem Service

- **Peer-to-peer operation**
- Buffered operation



Zero-copy of data between co-processor memory and SSD
Minimal data transfer

# Filesystem Service

- **Peer-to-peer operation**
- Buffered operation



Zero-copy of data between co-processor memory and SSD
Minimal data transfer

# Filesystem Service

- **Peer-to-peer operation**
- Buffered operation



Zero-copy of data between co-processor memory and SSD
Minimal data transfer

# Filesystem Service

- **Peer-to-peer operation**
- Buffered operation



Zero-copy of data between co-processor memory and SSD
Minimal data transfer

# Filesystem Service

- Peer-to-peer operation
- **Buffered operation**



Reduce storage IO by leveraging shared buffer cache among co-processors
Avoid performance anomaly of peer-to-peer communication over PCIe bus

# Filesystem Service

- Peer-to-peer operation
- **Buffered operation**



Reduce storage IO by leveraging shared buffer cache among co-processors
Avoid performance anomaly of peer-to-peer communication over PCIe bus

# Filesystem Service

- Peer-to-peer operation
- **Buffered operation**



Co-processor

Application

① 

File system stub

② 

Host processor

File system proxy

Buffer cache

③ 

File system

④ 

PCIe

······▶ control

——▶ data

DMA engine

SSD

Reduce storage IO by leveraging shared buffer cache among co-processors
Avoid performance anomaly of peer-to-peer communication over PCIe bus

# Filesystem Service

- Peer-to-peer operation
- **Buffered operation**



Reduce storage IO by leveraging shared buffer cache among co-processors
Avoid performance anomaly of peer-to-peer communication over PCIe bus

# Filesystem Service

- Peer-to-peer operation
- **Buffered operation**



Reduce storage IO by leveraging shared buffer cache among co-processors
Avoid performance anomaly of peer-to-peer communication over PCIe bus

# Implementation

- Host: 2-socket Xeon processor (12 cores each)
- Co-processor: 4 Xeon Phi (KNC, 61 cores, Linux, PCIe Gen 3x16)
- Storage device: 4 NVMe SSD
- NIC: 100 Gbps Ethernet

| Module | | Lines of code | |
|---|---|---|---|
| | | **Added lines** | **Deleted lines** |
| Transport service | | 1,035 | 365 |
| File system Service | Stub | 5,957 | 2,073 |
| | Proxy | 2,338 | 124 |
| Network Service | Stub | 2,921 | 79 |
| | Proxy | 5,609 | 34 |
| NVMe device driver | | 924 | 25 |
| SCIF kernel module | | 60 | 14 |
| **Total** | | **18,844** | **2,714** |

Questions:

- Performance of Solros services
- Impact on real-world applications

# Performance of Solros Services



(a) file random read on SSD

GB/sec — Phi-Solros, Phi-Linux

x-axis: 32KB, 64KB, 128KB, 256KB, 512KB, 1MB, 2MB, 4MB

(b) TCP latency: 64-byte message

Percentage of requests (%) vs Latency(usec) — Phi-Solros, Phi-Linux

File IO performance: 19x faster than the stock Linux on Xeon Phi
TCP latency (99 percentile): 7x shorter than the stock Linux on Xeon Phi

(a) file random read on SSD

(b) TCP latency: 64-byte message

Significant performance gain in data transport
Running IO stack on co-processor is slower

# Real-world Application - Image Search

- Image search engine is running on Xeon Phi
- Image database is on NVMe SSD (shared read-only)
- Image search queries are from network



Solros performs 2x faster than stock Linux on Xeon Phi

# Conclusion

- Solros, a new operating system architecture for co-processors and fast IO devices
- Control-plane, data-plane architecture allow:
    - Supporting high-level OS abstraction on co-processor
    - Efficient global coordination among devices
    - Near ideal IO performance from co-processor
- We will release source code soon

# Image Search - Scalability

- Increase the number of Xeon Phi to 4



Solros load balancing mechanism achieves near linear scaling

# Related work

- Control-plane/data-plane OS: Arrakis [OSDI'14], IX [OSDI'14]
- OS for heterogeneous systems: Helios [SOSP]09], M3 [ASPLOS'16], Hydra [ASPLOS'08]
- IO support for GPU: PTask [SOSP'11], GPUfs [ASPLOS'13], GPUnet [OSDI'14]

# Transport service

Host's virtual
address space

App 1

App 2

...

*mapped to*

*(in-host)*

Host's physical
address space

Host RAM

PCIe Window

mmio

*mapped to*

*(across devices)*

Device's physical
address sapce

DRAM    mmio

DRAM    mmio

DRAM    mmio

NIC

Xeon Phi

NVMe

# Network Service (TCP)

- Outbound operation
- **Inbound operation**



A load balancer on a host distributes incoming TCP connections to one of least-loaded co-processors.

See details in the paper.

# Discussion

- Hardware support other than Xeon Phi
    - Two atomic instructions: transport service
    - MMU: isolation among co-processor applications
- Scalability of control-plane OS
    - Limited by scalability of OS service, PCIe interconnect, and performance of IO devices

# Real-world Application - Text Search

- CLucene text indexing engine running on Xeon Phi
- Text data is on NVMe SSD



Solros performs 19x faster than stock Linux (ext4/virtio) on Xeon Phi