# OFA Workshop 2021 Session Abstracts

## Keynote

### Evolution of Interconnects and Fabrics to support future Compute Infrastructure

*Debendra Das Sharma, Intel Corporation*

High-performance workloads demand heterogeneous processing, tiered memory architecture including persistent memory, and infrastructure accelerators such as smart NICs to meet the demands of emerging applications in Artificial Intelligence, Machine Learning, Analytics, Cloud Infrastructure, Cloudification of the Network and Edge, communication systems, and high-performance Computing. Interconnect is a key pillar in this evolving computational landscape. Recent advances in I/O interconnects such as PCI Express and Compute Express Link with its memory and coherency semantics has made it possible to pool computational and memory resources at the rack level using low latency, higher throughput and memory coherent access mechanisms. We are already at a point where networking fabrics and tightly coupled load-store interconnects have overlapping coverage with some common attributes like standardized fabric manager for managing resources, low-latency message passing across nodes, and shared memory across multiple independent nodes. The traditional I/O interconnects are making changes to their load-store semantics to provide efficient access mechanisms for fabrics with advanced atomics, acceleration, smart NICs, persistent memory support etc. In this talk we will explore how synergistic evolution across interconnects and fabrics can benefit the compute infrastructure of the future.

## Technical Sessions

### Accelerating HPC Runtimes such as OpenSHMEM with COPA

*David Ozog, Intel Corp.; Andriy Kot, Intel Corp.; Michael Blocksome, Intel Corp.; Venkata Krishnan, Intel Corp.*

The Configurable Network Protocol Accelerator (COPA) framework enables FPGAs to incorporate custom inline/lookaside accelerators and attach directly to a 100 Gigabit/s Ethernet network. In addition to enabling FPGAs to be used as autonomous nodes for distributed computing, COPA also serves as a proof-of-concept for features intended for hardening in silicon.

The hardware component of COPA provides the necessary networking/accelerator infrastructure allowing custom accelerator modules to be integrated. On the software front, COPA abstracts the FPGA hardware by providing support for the OpenFabrics Interfaces (OFI). The COPA OFI provider also supports an enhanced OFI interface that exposes the acceleration and networking capabilities to upper-layer middleware/applications.

The performance of HPC middleware for distributed programming models (such as MPI and OpenSHMEM) would benefit from the acceleration capabilities provided by COPA's enhanced OFI, thereby improving overall application performance.  HPC middleware may require runtime and/or interface modifications to fully exploit COPA's underlying acceleration features, but there is also opportunity to design COPA's accelerator modules in a way that would require no changes to the middleware/application software to reap the performance benefits.

In this regard, we have taken up support for OpenSHMEM , a partitioned global address space (PGAS) programming model, by  enabling the Sandia OpenSHMEM library to run  with the COPA OFI provider. Additionally, the COPA framework also enables exploration of new features for accelerating OpenSHMEM runtimes. As part of that

approach, we have introduced new support for controlling hardware counters on the network device to track pertinent fabric events via libfabric provider software.

This presentation will provide a short overview of COPA and will describe the OpenSHMEM enablement on COPA OFI provider. OpenSHMEM makes use of OFI/libfabric fabric counters to track pending communication operations and to ensure remote completion and synchronization. This presentation will also provide details on the implementation of OpenSHMEM counters and how they can be utilized to support custom fine-grained communication tracking, which can support operations such as a PE-specific fence and quiet. The performance of OpenSHMEM (baseline and accelerated) will also be described in the context of OpenSHMEM kernels/benchmarks.

## Designing a High-performance MPI library using In-Network Computing

*Hari Subramoni, The Ohio State University; Dhabaleswar Panda, The Ohio State University; Mohammadreza Bayatpour, The Ohio State University; Bharath Ramesh, The Ohio State University; Kaushik Kandadi Suresh, The Ohio State University*
Emerging applications that use high-performance computing systems have large volumes of data being exchanged between workers present on different nodes. Based on these trends, to reach Exascale performance for end applications, it is beneficial to move the compute capabilities to the data instead of bringing the data to the compute elements. Thus, having support for in-network computing is critical for efficient scale-out of HPC and AI applications. High-performance networks such as InfiniBand have the support to offload more computational tasks to the network compared to portions of what was traditionally done on the host to the network. Excellent examples of modern in-network computing features in InfiniBand networks are the support for Scalable Hierarchical Aggregation and Reduction Protocol (SHARP) technology, Hardware Tag Matching, and Scatter-Gather List (SGL). These features can improve the overlap of communication and computation of point-to-point and collective operations while reducing the communication latency by lowering the traffic in the network. In this presentation, we discuss our various designs for these two technologies in MVAPICH2-X, and we provide benchmarking results on Frontera at large scale. Our experimental evaluation of the SHARP-based designs shows up to 5.1X for MPI_Allreduce, at a full system scale of 7,861 nodes over a host-based solution, up to 41% improvement for MPI_Iscatterv using the Hardware Tag Matching on 1024 nodes, and up to 30% improvement on 3D-Stencil kernel using SGL based Non-Contiguous Data transfers on 512 nodes.

## Designing a ROCm-aware MPI Library for AMD GPUs over High-speed Networks

*Hari Subramoni, The Ohio State University; Dhabaleswar Panda, The Ohio State University; Ching-Hsiang Chu, The Ohio State University; Kawthar Shafie Khorassani, The Ohio State University*
GPU-aware MPI libraries have been the driving force behind scaling large-scale applications on GPU-enabled HPC and cloud systems. GPU-enabled communication ecosystems thus far have been dominated by NVIDIA GPUs due to features like GPU-to-GPU RDMA communication (GPUDirect) working in tandem with InfiniBand high-speed networking. While AMD GPUs had offered somewhat competitive compute performance through OpenCL, the lack of support for high-performance communication stacks made its adoption in large-scale HPC deployments very limited. In recent years, AMD has developed an open-source suite of libraries called Radeon Open Compute (ROCm) that is tailored towards achieving efficient computation and communication performance for applications running on AMD GPUs. It provides features like ROCmRDMA (similar to GPUDirect RDMA) to support efficient GPU-to-GPU communication over RDMA networks. The availability of such high-performance communication features (e.g., PeerDirect, ROCm IPC, etc.) has made AMD hardware a viable choice for designing large-scale HPC systems. In fact, two upcoming DOE exascale systems (Frontier and El Capitan) will be driven by AMD GPUs. Since MPI is the back-bone of any large-scale supercomputing system, it is essential to have MPI libraries optimized for ROCm to exploit the capabilities of AMD GPUs and provide scalable application-level performance. In this work, we present early results on designing a ROCm-aware MPI library (MVAPICH2-GDR) that utilizes high-performance communication features like ROCmRDMA to achieve scalable multi-GPU performance. We discuss various implementation challenges and their solutions in detail. The performance evaluation is carried out using micro-benchmarks (point-to-point and collective) and application kernels on up to 64 GPUs across 16 nodes. The results show that our ROCm-aware MPI library can achieve good performance and scalability. Compared to other ROCm-aware communication stacks, we demonstrate a speedup of up to 3-6X for point-to-point bandwidth and up to 10X for various collectives benchmarks.

### Direct PSM2 support for NCCL

*Brendan Cunningham, Cornelis Networks; Dennis Dalessandro, Cornelis Networks*

NCCL is the NVIDIA Collective Communications Library. It includes optimizations for collective-heavy GPU workloads. NCCL supports operations spanning multiple nodes using both conventional networks and high-performance interconnects. However NCCL does not include a high performance transport for OPA. There have been attempts to create network plugins for other architectures making use of OFI. This is not optimal for Omni-Path Architecture (OPA) based fabrics as it precludes the use of GPU Direct. In order to rectify this situation we present a plug-in for NCCL that will allow direct use of PSM2 which is the highest performing message passing middleware available for OPA. This presentation will highlight the technical challenges faced and performance results obtained.

### Efficient MPI Offloading Designs on Modern BlueField Smart NICs

*Dhabaleswar Panda, The Ohio State University; Hari Subramoni, The Ohio State University; Mohammadreza Bayatpour, The Ohio State University; Nick Sarkauskas, The Ohio State University*

In the state-of-the-art production quality MPI (Message Passing Interface) libraries, communication progress is either performed by the main thread or a separate communication progress thread. Taking advantage of separate communication threads can lead to a higher overlap of communication and computation as well as reduced total application execution time. However, such an approach can also lead to contention for CPU resources leading to sub-par application performance as the application itself has fewer available cores for computation. Recently, Mellanox has introduced the BlueField series of adapters which combine the advanced capabilities of traditional ASIC-based network adapters with an array of ARM processors. In this presentation, we discuss our proposed designs that can be used to offload different collective and point-to-point operations from the host CPU to the Smart NIC. Our designs provide full overlap of communication and computation while providing on-par pure communication latency to CPU based on-loading designs. Our experiments show that our designs can improve the performance at the benchmark level and as well as application-level up to 50%.

### Flatten the curve: Source Flow Control for sub-RTT management of network tail latency

*Jeongkeun Lee, Intel Corp.; Jeremias Blendin, Intel Corp.; Yanfang Le, Intel Corp.; Grzegorz Jereczek, Intel Corp.*

Tail latency is one of the most important performance metrics in HPC and modern datacenter workloads like distributed storage and RPC. Network congestion, especially caused by large fan-in application incast is a major contributing factor to the tail latency. Modern congestion control algorithms like Swift and HPCC have pulled the tail queueing latency down below 1ms. Pulling it down below 100us is much desired for HPC workloads but is challenging due to the inherent nature of end-to-end congestion control: the congestion information is carried by forward-direction data packets that are delayed by the ongoing congestion. Further, congestion algorithms take multiple RTTs to detect, react and converge. Meanwhile, as the network link speed goes up, the time needed to finish each flow/message transmission comes down, leaving not enough time for E2E congestion control to handle incast and sudden congestion. This talk will show that a simple L3 signaling between datacenter switches, combined with standard PFC flow control available in today's NICs, can enable instant (sub-RTT) reaction to incast at the traffic sources, protecting the scarce switch buffer and reducing the network tail latency. This Source Flow Control (SFC) can be deployed by upgrading dataplane program (such as P4) only at Top-of-Rack switches. No change to end host or NIC transport stack or congestion control protocol is required.

### How to efficiently provide software-defined storage using SmartNICs

*Jonas Pfefferle, IBM Research; Nikolas Loannou, IBM Research; Jose Castanos, IBM Research; Bernard Metzler, IBM Research*

In recent years SmartNICs have become the state-of-the-art solution for providing storage and network virtualization in cloud environments. Leading cloud providers are deploying SmartNIC-based services, like AWS Nitro or Azure SmartNIC, to provide isolation, security, increased performance and lower cost. By example of providing transparent block storage using SmartNICs as the storage clients and the distributed storage system Ceph as the backend, we discuss main aspects of the general design space of a SmartNIC deployment. We will focus on contradicting aspects like raw performance and flexibility/programmability, exemplified by two prototyped solutions: One where the Ceph client protocol is deployed in software directly on the SmartNIC, and another where a standard storage protocol client, NVMe-over-Fabrics RDMA, is deployed in hardware offload mode on the SmartNIC, while introducing an intermediate NVMeoF-to-Ceph gateway. Our initial performance evaluation of the two solutions shows that the NVMeoF offload-based solution is significantly faster, while giving in on flexibility. Both solutions were prototyped in a 100Gbps RoCE network.

## INEC: In-Network Erasure Coding

*Xiaoyi Lu, University of California Merced*

Erasure coding (EC) is a promising fault tolerance scheme that has been applied to many well-known distributed storage systems. The capability of Coherent EC Calculation and Networking on modern SmartNICs has demonstrated that EC will be an essential feature of in-network computing. This talk will introduce our proposed coherent in-network EC primitives, named INEC. Our analyses based on the proposed performance model demonstrate that INEC primitives can enable different kinds of EC schemes to fully leverage the EC offload capability on modern SmartNICs. We implement INEC on commodity RDMA NICs and integrate it into five state-of-the-art EC schemes. Our experiments show that INEC primitives significantly reduce 50th, 95th, and 99th percentile latencies, and accelerate the end-to-end throughput, write, and degraded read performance of the key-value store co-designed with INEC by up to 99.57%, 47.30%, and 49.55%, respectively.

## Infiniband Developments around the T7 Trading Engine at the Deutsche Boerse

*Christoph Lameter, Deutsche Boerse AG*

The Deutsche Boerse has been running Infiniband reliably in a production environment for almost a decade now. Deutsche Boerse has become one of the largest Stock Exchanges and runs a multiplicity of financial market places amoung them EUREX. In this talk we cover some of the recent improvements that were made to increase the scalability and performance of the T7 Trading Engine. - Basic design of the T7 infrastructure from an Infiniband and RDMA perspective - Benefits of the Sendonly Join Multicast feature to reduce data volume. - SubnetManager optimizations and the upgrade of the Subnet manager Infrastructure. - How to realize fast recovery for small failures - Challenges due to the scalability of reregistration events - A future perspective going forward to a mixed Ethernet and Infiniband environment.

## Infiniband reliability engineering - stories from a public cloud

*Vladimir Chukov, 1&1 IONOS*

This presentation is focused on the fabric reliability in an ever-changing and growing cloud environment and is based on tests, user stories and RCAs. We will talk about components that fail (HCAs, switches, software), possible impact and available mitigation.

## Introduction to the OFA Fabric Software Development Platform

*Tatyana Nikolova, Intel Corp.; Doug Ledford, Red Hat, Inc.*

This year the OFA created a new cluster we call the Fabric Software Development Platform. This cluster is available for OFA members and upstream community members to use for developing and testing fabric related software. This presentation will cover the intended audience and use cases for the FSDP, the rationale behind it, what we hope to accomplish with it, and most importantly, how people can gain access to the FSDP cluster.

## Omni-Path Accelerated IP (AIP)

*Mike Marciniszyn, Cornelis Networks; Dennis Dalessandro, Cornelis Networks*

The Accelerated IP (AIP) feature of the hfi1 driver that was added to the Linux kernel in recent releases is what allows IPoIB to achieve near line rate throughput. Using advanced features of the hfi1 chip in OPA adapters we have enabled a method to reach throughput not otherwise possible. This presentation will look at the technical issues we faced as well as covering the design and implementation. We will highlight customer use cases and touch on the journey to have this feature included in the current Linux kernel.

## Libfabric Intranode Device Support

*Alexia Ingerson, Intel, Corp.*

With device memory usage becoming more common, ongoing work is being done to add device support to libfabric, the shared memory provider being one target area of where to add this support. The provider currently uses shared memory for local communication but needs extensions to support copying to and from device memory. This talk will provide an overview of the current design of the provider and how it will be adapted to support a variety of device types and transfers, including additions and modifications of its current protocols.

## Overview of Gen-Z Linux Subsystem, Fabric Manager, and Bridge Driver Development

*Jim Hull, Gen-Z Consortium*

A technical overview of the Gen-Z Linux Sub-System, including subsystem interface descriptions for the Gen-Z Bridge Driver, Fabric Manager, and Linux Local Management Services. Additionally, a discussion on the motivation to implement a Gen-Z subsystem in the first place.

## Performance Scaled Messaging v3 (PSM3) Architecture Overview

*Todd Rimmer, Intel Corp.*

PSM3 is a new OFI provider which supports a variety of devices and protocol mechanisms, including both on-load and off-load data movement strategies using devices such as Intel Columbiaville RoCE NICs. This session will discuss the features, architecture and internal components which comprise PSM3, including PSM3's rendezvous kernel module. The Rendezvous Module is a new kernel module designed to provide highly scalable communications support for PSM3.

## Progress of upstream GPU RDMA support

*Jianxin Xiong, Intel Corp.*

The ability to use device memory directly in RDMA transfers is important for scaling out computation systems that utilize GPU and other accelerators. Although such feature has been available in commercial systems for a while, it was missing from the upstream Linux kernel. One year ago, we started the effort of adding GPU RDMA support to upstream Linux kernel using DMA-Buf as the buffer sharing mechanism. Good progress have been made with the support from the community. This talk will present an update on the latest status, the issues we have encountered during the period and the major changes we have made since the original prototype.

## PSM3 Performance Studies

*James Erwin, Intel Corp.*

This session will review in-depth performance studies and results for Performance Scaled Messaging (PSM3) on Intel(R) Ethernet Network Adapter E810 Series NICs with RDMA, as well as some examples of PSM3 in other environments. A range of HPC applications at a variety of scales will be shown, as well as the latest best known methods for profiling and tuning PSM3.

## RDMA Spark Meets OSU INAM: Performance Engineering Big Data Applications on HPC Clusters

*Aamir Shafi, The Ohio State University; Dhabaleswar Panda, The Ohio State University; Hari Subramoni, The Ohio State University; Mansa Kedia, The Ohio State University; Pouya Kousha, The Ohio State University*

Spark is a popular Big Data framework that enables high-performance data science in multiple languages including Java, Spark, and Python. RDMA Spark is an enhanced version of the Apache Spark with support for high-performance RDMA networks like InfiniBand. It is important for data scientists and system administrators to be able to profile their Big Data applications executed using Spark in order to understand bottlenecks and ultimately optimize their codes. In this context, Spark provides a Web User Interface (Web UI) exposing basic application statistics, which are used today to monitor application performance. This information, however, is inadequate since it does not provide a holistic view especially on modern HPC systems equipped with high-performance interconnect like InfiniBand. The most communication intensive part of the Spark software is the so-called shuffle operation that is triggered to re-distribute distributed data between various stages of the application execution. It is critical for end-users to gain insights about the shuffle operation in order to mitigate sources of performance overheads by taking appropriate corrective measures. We take up this challenge by introducing the capability of performance engineering Big Data applications, executed by the RDMA Spark framework, in OSU INAM that is a network profiling, monitoring, and analysis tool. Unlike the Spark Web UI, OSU INAM is designed to provide real-time scalable performance insight to network traffic on HPC interconnects through tight integration with middleware runtimes and job schedulers. This allows end-users to conduct performance optimization and workload characterization. This is possible because OSU INAM provides a holistic view of application performance by collecting and correlating network communication at system level with data exchanges carried out during the shuffle phase.

## Remote Persistent Memory Access as Simple as Local Memory Access

*Tomasz Gromadzki, Intel Corp.*
Remote PMem (RPMem) is the fastest distributed storage available on the market, while it is the most challenging aspect of persistent memory deployment in distributed systems. All needed components to use RPMem are ready and available in every modern Linux distribution. But many potential users experience the situation where two technologies must be incorporated simultaneously: persistent memory and RDMA, as many applications still use TCP/IP based data exchange methods. Recently PMDK team developed the librpma library to address this issue. The library provides an easy to integrate, simple but powerful API for remote persistent memory operations. The library comes with several examples and, together with the dedicated Fio engine, provides a very powerful environment for remote persistent memory programming.

## sRDMA: secure transport for Remote Direct Memory Access

*Konstantin Taranov, ETH Zurich; Benjamin Rothenberger, ETH Zurich; Adrian Perrig, ETH Zurich; Torsten Hoefler, ETH Zurich*
RDMA architectures such as RoCE and InfiniBand were designed for high-performance computing and private networks, and have neglected security in their design in favor of focusing on high performance. Therefore, current RDMA protocols lack any form of cryptographic authentication or encryption, making them a powerful attack vector for an adversary. During the presentation, we demonstrate multiple vulnerabilities in the design of InfiniBand-based architectures and implementations of RDMA-capable network interface cards and exploit those vulnerabilities to enable powerful attacks such a packet injection, unauthorized memory access, and Denial-of-Service attacks. Given these threats, we propose sRDMA, a protocol that provides authentication and encryption for RDMA networking. sRDMA uses symmetric cryptography and requires minimal changes to the current InfiniBand architecture. Additionally, sRDMA utilizes efficient dynamic key derivation to remove the need for storing security keys for each connection and to extend memory protection mechanisms enabling the delegation of memory accesses.

## Standards-based scalable fabrics management now and into the future

*Richelle Ahlvers, Intel Corp.; Phil Cayton, Intel Corp.; Rajalaxmi Angadi, Intel Corp.*
The SNIA Swordfish™ specification helps provide a unified approach for the management of storage in hyperscale and cloud infrastructure environments, making it easier for IT administrators to integrate scalable solutions into data centers. Swordfish now ensures efficient management of NVMe and NVMe-oF technology environments. The Swordfish ecosystem includes a family of open-source tools and a vendor-neutral conformance test suite both designed to accelerate implementations and adoptions of the standard.
SNIA is partnering with other industry organizations to expand the scope covered in Swordfish management to meet client workloads and use-cases covering everything from direct attached storage (DMTF) to QOS-based resource orchestration (SODA), to fabric technology management and administration (OFA).

## Status of OpenFabrics Interfaces (OFI) Support in MPICH

*Yanfei Guo, Argonne National Laboratory*
This session will give the audience an update on the OFI integration in MPICH. MPICH underwent a large redesign effort (CH4) in order to better support high-level network APIs such as OFI. We will show the benefits realized with this design, as well as ongoing work to utilize more aspects of the API and underlying functionality. This talk has a special focus on how MPICH is using Libfabric for GPU support and the development updates on GPU fallback path in Libfabric.

## Supporting Live Migration of VMs communicating with bare-metal RDMA endpoints

*Jorgen Hansen, VMware; Bryan Tan, VMware; Rajesh Jalisatgi, VMware; Adit Ranadive, VMware; Vishnu Dasa, VMware; Georgina Chua, VMware*
For several years, VMware has been providing a paravirtual RDMA device (PVRDMA) allowing VMs in a cluster to communicate through Remote Direct Memory Access (RDMA). One major limitation of PVRDMA has been that all RDMA endpoints had to be PVRDMA in order to support advanced virtualization features such as snapshots and live migration. Through close collaboration with hardware partners, PVRDMA is now able to support communication with non-virtualized RDMA endpoints such as storage arrays, NVMe-OF targets and physical file servers with full virtualization support for the PVRDMA-enabled VMs. In this session, we provide an overview of the design to support VMs communicating with non-virtualized endpoints, discuss the major challenges in providing transparent live migration

without end to end control over the RDMA connections, and how additional hardware support has bridged this gap. In addition to this, we present performance measurements showing the benefits of the added hardware support when running storage and HPC workloads.

## True Zero memory copy Direct Memory Transfer for Clusters and High-performance Computing over RDMA using p2pmem API's

*Suresh Srinivasan, Intel Corp.; Phil Cayton, Intel Corp.*

High Performance Computing (HPC) and Cluster Computing largely depend on data sharing across nodes. HPC and Cluster Computing AI application (e.g., TensorFlow, PyTorch) performance largely relies on internode communication speed and data storage/retrieval rate. Communication links such as RDMA, highspeed ethernet etc., achieve high data rates such as 100Gbps but storage poses a significant bottleneck to these data rates. This is primarily due to additional kernel buffer copy and virtual memory translation before transferring data between memories across nodes. Though Linux kernel proposes zero memory copy, it fails to bypass virtual memory translations causing additional delays in storing and retrieving data. Peer to Peer memory (p2pmem) transfer is one of the key accelerators for enhancing data transfer rate. In true zero copy RDMA, remote data storage and retrieval are accessed directly bypassing the kernel. We will discuss exploiting the existing p2pmem API from the upstream kernel to achieve true zero copy RDMA data transfer across nodes.

## Understanding Compute Express Link 2.0: A Cache-coherent Interconnect

*Jim Pappas, CXL Consortium*

Compute Express Link™ (CXL™) is an industry-supported cache-coherent interconnect for processors, memory expansion, and accelerators. Datacenter architectures continue to evolve rapidly to support the ever-growing demands of emerging workloads such as Artificial Intelligence and Machine Learning. CXL is an open industry-standard interconnect offering coherency and memory semantics using high-bandwidth, low-latency connectivity between the host processor and devices such as accelerators, memory buffers, and smart I/O devices. CXL technology is designed to address the growing needs of high-performance computational workloads by supporting heterogeneous processing and memory systems for applications in Artificial Intelligence, Machine Learning, communication systems, and High-Performance Computing. These applications deploy a diverse mix of scalar, vector, matrix, and spatial architectures through CPU, GPU, FPGA, smart NICs, and other accelerators. At this session, attendees will learn about the next generation of CXL technology. The CXL 2.0 specification adds support for switching for fan-out to connect to more devices; memory pooling for increased memory utilization efficiency and providing memory capacity on demand; and support for persistent memory. This presentation will explore the new features of CXL 2.0 and share how CXL technology is keeping pace to meet the performance and latency demands of emerging workloads for Artificial Intelligence and Machine Learning.

## Birds of a Feather

## OpenFabrics Management Framework

*Michael Aguilar, Sandia National Laboratories; Jeff Hilland, HPE; Russ Herrell, HPE; Paul Grun*

This BoF will be both an illustration and discussion of the work being performed on behalf of the OpenFabrics Management Framework Working Group towards a new high-speed network Subnet Manager. The Management Framework is designed to provide on-demand fabric subnets and route management that match performance requirements for bandwidth, latency, and security. In our illustration, we will provide a sample of Use-Case descriptions that we are developing to outline the software functionality of the key components of the Fabric Manager.

## To logo program, or not to logo program?

*Tatyana Nikolova, Intel Corp.; Doug Ledford, Red Hat*

The OFA needs to decide whether or not we will run a logo program using the FSDP cluster as a means of doing logo testing. We are actively soliciting input from the larger ecosystem on this topic. Vendors - Do you want a logo program? If we have one, what requirements do you think it should meet? Consumers - Do you want a logo program? Do you see value to having logo certified hardware? If we have one, what requirements do you think it should meet? This is intended to be an active discussion from which the OFA can listen to your input and decide whether a logo program is something the OFA should undertake.