



# Dell HPC Omni-Path Fabric: Supported Architecture and Application Study

Deepthi Cherlopalle  
Joshua Weage  
Dell HPC Engineering  
June 2016

# Revisions

Date	Description
June 2016	Initial release – v1

THIS WHITE PAPER IS FOR INFORMATIONAL PURPOSES ONLY, AND MAY CONTAIN TYPOGRAPHICAL ERRORS AND TECHNICAL INACCURACIES. THE CONTENT IS PROVIDED AS IS, WITHOUT EXPRESS OR IMPLIED WARRANTIES OF ANY KIND.

Copyright © 2016 Dell Inc. All rights reserved. Dell and the Dell logo are trademarks of Dell Inc. in the United States and/or other jurisdictions. All other marks and names mentioned herein may be trademarks of their respective companies.



# Table of contents

Revisions.....	2
Executive Summary.....	5
Audience.....	5
1 Introduction.....	7
2 Dell Networking H-Series Fabric.....	9
2.1 Intel® Omni-Path Host Fabric Interface (HFI).....	9
2.1.1 Server Support Matrix.....	9
2.2 Dell H-Series Switches based on Intel® Omni-Path Architecture.....	9
2.2.1 Dell H-Series Edge Switches.....	10
2.2.2 Dell H-Series Director-Class Switches.....	11
3 Intel® Omni-Path Fabric Software.....	12
3.1 Available Installation Packages.....	12
3.1.1 Supported Operating Systems.....	12
3.1.2 Operating System Prerequisites.....	12
3.2 Intel® Omni-Path Fabric Manager GUI.....	13
3.2.1 Overview of Fabric Manager GUI.....	13
3.3 Chassis Viewer.....	14
3.4 OPA FastFabric.....	15
3.4.1 Chassis setup.....	15
3.4.2 Switch setup.....	15
3.4.3 Host Setup.....	15
3.5 Fabric Manager.....	15
3.5.1 Embedded Subnet Manager.....	16
3.5.2 Host-Based Fabric Manager.....	17
4 Test bed and configuration.....	18
5 Performance Benchmarking Results.....	20
5.1 Latency.....	20
5.2 Bandwidth.....	20
5.3 Weather Research Forecast.....	21
5.4 NAMD.....	22
5.5 ANSYS® Fluent®.....	22
5.6 CD-adapco® STAR-CCM+®.....	24
6 Conclusion and Future Work.....	26
7 References.....	27



## List of Figures

Figure 1	Intel® Omni-Path Host Fabric Interface Adapter 100 series 1 port PCIe x16 .....	9
Figure 2	Dell H-Series Switches.....	10
Figure 3	Dell H1024-OPF Edge Switch .....	10
Figure 4	H1048-OPF Edge Switch .....	10
Figure 5	H9106-OPF ( <i>left</i> ) and H9124-OPF ( <i>right</i> ) Director-Class Switches .....	11
Figure 6	Intel® Omni-Path Fabric Suite Fabric Manager Home page overview.....	13
Figure 7	Intel® Omni-Path Chassis Viewer Overview .....	14
Figure 8	Fastfabric Tools .....	15
Figure 9	Starting Embedded Subnet Manager .....	17
Figure 10	OSU Latency values based on Intel® Xeon® CPU E5-2697 v4 processor .....	20
Figure 11	OSU Bandwidth values based on Intel® Xeon® CPU E5-2697 v4 processor.....	21
Figure 12	Weather Research Forecasting model performance graph.....	21
Figure 13	NAMD performance graph.....	22
Figure 14	ANSYS® Fluent® Relative Performance Graph (1/2).....	23
Figure 15	ANSYS® Fluent® Relative Performance (2/2).....	23
Figure 16	STAR-CCM+® Relative Performance Graph (1/2) .....	24
Figure 17	STAR-CCM+® Relative Performance (2/2).....	25

## List of tables

Table 1	Server Configuration .....	18
Table 2	Application and benchmarks Details.....	19



## Executive Summary

In the world of High Performance Computing (HPC), servers with high-speed interconnects play a key role in the pursuit to achieve exascale performance. Intel's Omni-Path Architecture (OPA) is the latest addition to the interconnect universe and is a part of Intel Scalable System framework. It is based on innovations from Intel's True Scale technology, Cray's Aries interconnect, internal Intel® IP and from several other open source platforms. This paper provides an overview of the OPA technology, its features, Dell's support and walks through the OPA software ecosystem. It dwells into the performance aspects of OPA ranging from micro benchmarks to many commonly used HPC applications (both proprietary and open source) at a 32 node (~900 core) scale. To summarize the experience, the learning curve for OPA wasn't steep because of its roots in True Scale architecture and it satisfies the low latency and high bandwidth requirements needed for HPC applications.

*A big thank you to our team members Nishanth Dandapanthula, Alex Filby and Munira Hussain for their day-to-day effort assisting with the setup and configurations necessary for this paper, and we would like to thank James Erwin from Intel who has been helping us through our journey with Omni-Path.*

## Audience

This document is intended for people who are interested in learning about the key features and application performance of the new Intel® Omni-Path Fabric technology.





# 1 Introduction

The High Performance Computing (HPC) domain primarily deals with problems which surpass the capabilities of a standalone machine. With the advent of parallel programming, applications can scale past a single server. High performance interconnects provide low latency and high bandwidth which are needed for the application to divide the computational problem among multiple nodes, distribute data and then merge partial results from each node to a final result. As the computation power increases with greater number of nodes/cores added to the cluster configuration, the need for efficient and fast communication has become essential to continue to improve system performance. Applications may be sensitive to throughput and/or latency capabilities of the interconnect depending upon their communication characteristics.

Intel® Omni-Path Architecture (OPA) is an evolution of Intel® True Scale Fabric Cray Aries interconnect [1] and internal Intel® IP. In contrast to Intel® True Scale Fabric edge switches that support 36 ports of InfiniBand QDR-40Gbps performance, the new Intel® Omni-Path fabric edge switches support 48 ports of 100Gbps performance. The switching latency for True Scale edge switches is 165ns-175ns. The switching latency for the 48-port Omni-Path edge switch has been reduced to around 100ns-110ns. The Omni-Path Host Fabric Interface (HFI) MPI messaging rate is expected to be around 160 Million messages per second (Mmps) and a link bandwidth of 100Gbps.

The OPA technology includes a rich feature set. A few of those are described here <sup>[1]</sup>:

**Dynamic Lane Scaling:** When one or more physical lanes fail, the fabric continues to function with the remaining available lanes and the recovery process is transparent to the user and application. This allows jobs to continue and provides the flexibility of troubleshooting errors at a later time.

**Adaptive Routing:** This monitors the routing paths of all the fabrics connected to the switch and selects the least congested path to balance the workload. This implementation is based on the cooperation between Application-specific integrated circuits (ASIC) and Fabric Manager. The Fabric Manager performs the role of initializing the fabrics and setting up routing tables, once this is done the ASICs actively monitor and manage the routing by identifying fabric congestion. This feature helps the fabric to scale.

**Dispersive Routing:** Initializing and configuring the routes between the neighboring nodes of the fabric is always critical. Dispersive routing distributes the traffic across multiple paths as opposed to sending them to the destination via a single path. It helps to achieve maximum communication performance for the workload and promotes optimal fabric efficiency.

**Traffic Flow Optimization:** Helps in prioritizing packets in mixed traffic environments like storage and MPI. This helps to ensure that the high priority packets will not be delayed and there will be less/no latency variation on the MPI job. Traffic can also be shaped during run time by using congestion control and adaptive routing.

**Software Ecosystem:** It leverages the Open Fabric Alliance (OFA) <sup>[2]</sup> and uses a next generation Performance Scaled Messaging (PSM) layer called PSM2 which supports extreme scale but is still compatible with previous generation PSM applications. OPA also includes a software suite with extensive capabilities for monitoring and managing the fabric.

**On-load and Offload Model:** Intel Omni-Path can support both on-load and offload models depending on the data characteristics. There are two methods of sending data from one host to another.



### Sending the data

- Programmed I/O (PIO): This supports the on-load model. The host can be used to send small messages since these can be sent by the CPU faster than an RDMA setup time.
- Send DMA (SDMA): For larger messages the CPU sets up a RDMA send and then the 16 SDMA engines in the HFI transfer the data to the receiving host without CPU intervention.

### Receiving the data

- Eager receive: The data is delivered to the host memory and then copied to the application memory. This protocol is faster for smaller messages and does not require any responses
- Expected Receive: Data flows directly from HFI to application memory without CPU intervention.

Each data transfer method is independent. For example, SDMA can be used from Sender's side and Eager receive method can be used on the other side. This can be used for medium size messages since it does not require full RDMA setup. However, the packet size threshold for small, medium and large have default values but are configurable.

All these features make OPA an ideal component for HPC workloads. In the upcoming sections, this white paper details Dell's support of Intel OPA and dives deeper into the software ecosystem. Finally, it discusses the performance characterization of several HPC workloads at scale using OPA.





## 2 Dell Networking H-Series Fabric

Dell Networking H-Series Fabric is a comprehensive fabric solution that includes host adapters, edge and director class switches, cabling and complete software and management tools.

### 2.1 Intel® Omni-Path Host Fabric Interface (HFI)

Dell provides support for Intel® Omni-Path HFI 100 series cards <sup>[3]</sup> which are PCIe Gen3 x16 and capable of 100Gbps per port. Each Intel® Omni-Path HFI card has 4 lanes supporting 25Gbps each and can deliver up to 25GBps bidirectional bandwidth per port. Intel® OPA supports QSFP28 quad small form factor with pluggable passive and optical cables.



Figure 1 Intel® Omni-Path Host Fabric Interface Adapter 100 series 1 port PCIe x16

#### 2.1.1 Server Support Matrix

The following Dell servers support Intel® Omni-Path Host Fabric Interface cards

- PowerEdge R430
- PowerEdge R630
- PowerEdge R730
- PowerEdge R730XD
- PowerEdge R930
- PowerEdge C4130
- PowerEdge C6320
- PowerEdge FC830

## 2.2 Dell H-Series Switches based on Intel® Omni-Path Architecture

Dell networking provides H-series Edge and Director Class switches which are based on the Intel Omni-Path architecture and are targeted for HPC environments ranging from small to large scale clusters.





Figure 2 Dell H-Series Switches

### 2.2.1 Dell H-Series Edge Switches

Dell H-Series Edge Switches based on the Intel® Omni-Path Architecture consist of two models supporting 100Gbps for all ports: an entry-level 24-port switch targeting entry-level/small clusters and a 48-port switch which can be combined with other edge switches and directors to build larger clusters.



Figure 3 Dell H1024-OPF Edge Switch



Figure 4 H1048-OPF Edge Switch

A complete description and specifications for these switches can be found on the [Dell Networking H-Series Edge Switches](#) product page.

## 2.2.2 Dell H-Series Director-Class Switches

Dell H-Series Director-Class Switches based on the Intel® Omni-Path Architecture consist of two models supporting 100Gbps for all ports: a 192-port switch and a 768-port switch. These switches support HPC clusters of all sizes, from mid-level clusters to supercomputers.



Figure 5 H9106-OPF (*left*) and H9124-OPF (*right*) Director-Class Switches

A complete description and specifications for these switches can be found on the [Dell Networking H-Series Director-Class Switches](#) product page.

## 3 Intel® Omni-Path Fabric Software

### 3.1 Available Installation Packages

The following packages <sup>[4]</sup> are available for an Intel® Omni-Path Fabric:

- **Intel® Omni-Path Fabric Host Software** – This is the basic installation package that provides Intel® Omni-Path Fabric Host components needed to set up compute, I/O and service nodes with drivers, stacks and basic tools for local configuration and monitoring. This package is usually installed on compute-nodes.
- **Intel® Omni-Path Fabric Suite (IFS) Software** – This package installs all the components that are included in basic and adds additional fabric management tools like FastFabric and Fabric Manager.
- **Intel® Omni-Path Fabric Suite Fabric Manger GUI** – This package can be used to monitor and manage one or more fabrics and it can be installed on a computer outside of the fabric.

#### 3.1.1 Supported Operating Systems

The following list of operating systems are supported by Omni-Path Fabric Host/Suite Software with IFS 10.0.1.1.5(*pre-release*):

- Red Hat Enterprise Linux 7.1, 7.2
- SUSE Linux Enterprise Server 12, 12 SP1
- CentOS 7.1
- Scientific Linux 7.1

The Dell HPC release is based on Red Hat Enterprise Linux 7.2 kernel version 3.10.0-327.el7.x86\_64.

#### 3.1.2 Operating System Prerequisites

The following packages must be installed before installing the Intel® Omni-Path Software.

##### RHEL 7.x

The following packages can be installed using yum from the RHEL distribution:

- libibmad
- libibverbs
- infinipath-psm
- sysfsutils
- expect
- atlas

Please check the [Intel Omni-Path software installation Guide](#) for the list of packages to be installed on other Operating Systems.

Note: The IFS package conflicts with the MLNX OFED packages. Hence it is important that MLNX OFED not be installed in conjunction with IFS packages to avoid any conflicts or error messages.



## 3.2 Intel® Omni-Path Fabric Manager GUI

Fabric Manager GUI [5] provides a set of analysis tools for graphically monitoring fabric status and managing fabric health. This package is open source and can be run on a Linux or Windows system with TCP/IP connectivity to the Fabric Manager. To use the Fabric Manager GUI enable the following switch settings and ensure the opafm.xml files are identical on all switches.

- ismChassisSetFmEnabled
- smPmStart (SM is enabled, PM is enabled, FE is enabled)

When FMGUI is launched for the first time, a short setup wizard will be used to configure the GUI.

### 3.2.1 Overview of Fabric Manager GUI

Fabric Manager GUI can be used to review subnet performance, topology and manage the fabric on the switches.

- Subnet Summary- shows the master SM, master SM uptime, standby SM information, and Inter-switch links information.
  - Provides total number of active nodes, switches, and ports not being used in the fabric
  - Shows switch status, Host Fabric interface status and health trend in graphical form.
- Subnet Performance – shows the bandwidth and packet rate information on all the nodes, current traffic on the fabric, the top 10 worst performing nodes, etc.
- Fabric topology gives a snapshot of how the nodes are connected to the switches. Double clicking on a particular node or switch gives detailed device information and routing information for that device. Detailed information for FM GUI can be accessed from this [link](#).

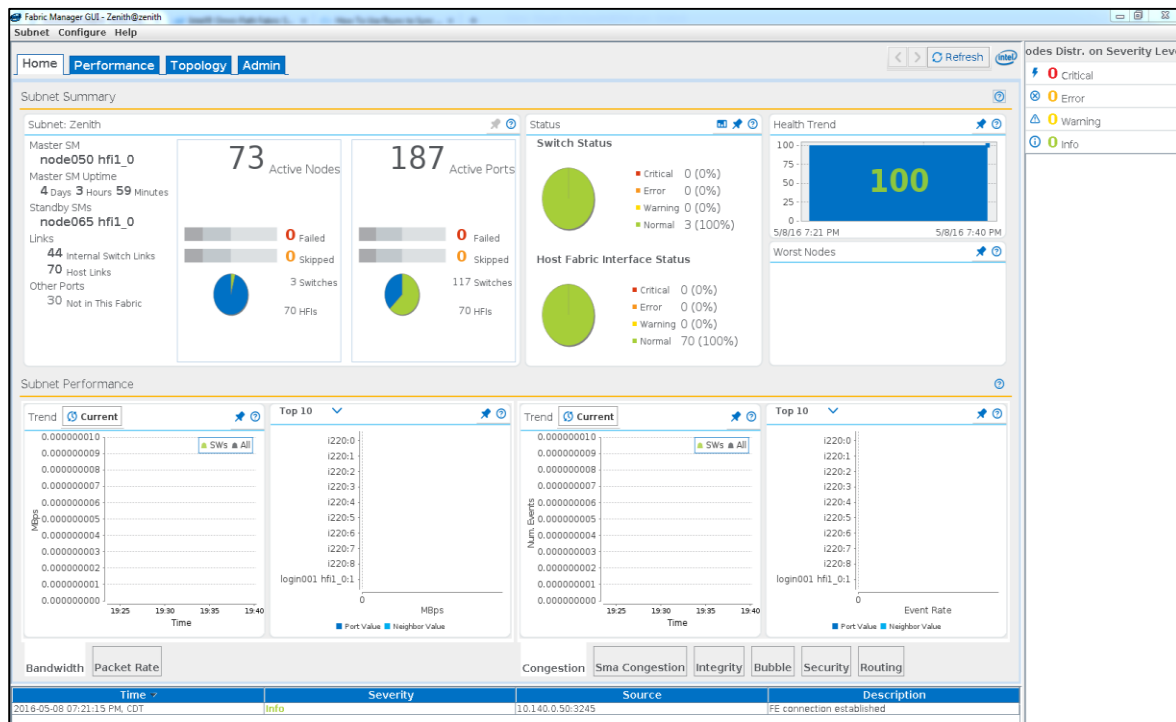


Figure 6 Intel® Omni-Path Fabric Suite Fabric Manager Home page overview

### 3.3 Chassis Viewer

Chassis viewer is a web interface which can be used to manage basic functionalities on edge switches with and without management cards. The following picture exemplifies the basic layout of Chassis Viewer:

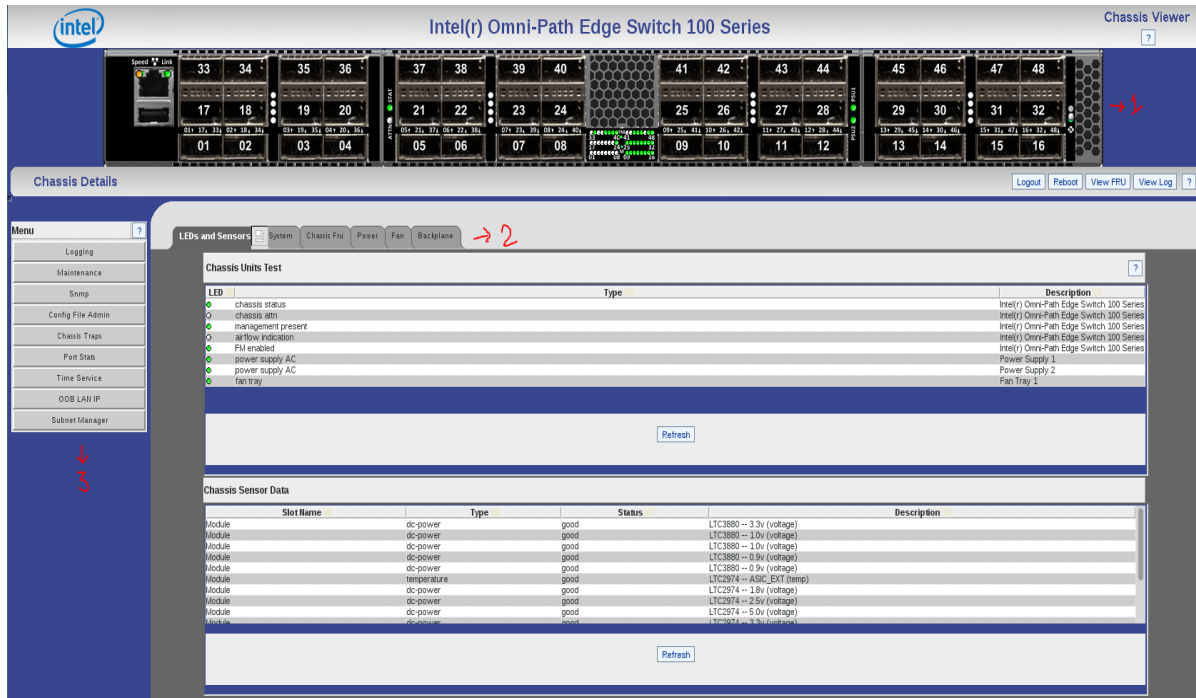


Figure 7 Intel® Omni-Path Chassis Viewer Overview

1. The LED lights at the center are “green” if the ports are active and “white” when in the polling state or the cables are not connected. It also shows the port status for PSU 1 & 2 and these LED indicators go red if any of the PSU’s need attention.
2. The main tab shows the LED and sensor information for chassis status, chassis attention, airflow indication, FM enabled, power supply units and fan trays. Toggling between the tabs shows detailed data on the chassis such as: Out of band IP information, FRU information, Power supply information, Fan Information and Backplane Information.
3. The main menu to the left shows options to manage the logging levels, maintenance of the switch, SNMP, Config File Admin, Chassis Traps, Port Stats, Time Service, OOB LAN IP, Subnet Manager



## 3.4 OPA FastFabric

FastFabric [6] is a set of fabric management tools used for fabric deployments, switch management and host management. This can be accessed if Intel Fabric Suite is installed on the nodes. The following screen appears when the “opafastfabric” command is given.

```
Intel FastFabric OPA Tools
Version: 10.0.1.1.2

 1) Chassis Setup/Admin
 2) Externally Managed Switch Setup/Admin
 3) Host Setup
 4) Host Verification/Admin
 5) Fabric Monitoring

 X) Exit
```

Figure 8 Fastfabric Tools

### 3.4.1 Chassis setup

This option helps the user to manage or monitor the chassis like setting up the Out-of-Band IP address on the switches, configuring syslog server, configuring NTP server, setting up time zone, OPA node description, rebooting the chassis, configuring the fabric manager, viewing firmware version, etc.

### 3.4.2 Switch setup

This can be used to retrieve firmware information, update firmware, and generate switch reports to aide in debugging. Additionally this allows running commands in parallel across all switches, as well as rebooting them.

### 3.4.3 Host Setup

This can be used to setup password-less SSH between hosts, upgrade the OPA software, configure IP over Fabric, running a command in parallel on all the hosts, etc.

A detailed explanation on opafastfabric can be accessed from the following [link](#).

## 3.5 Fabric Manager

Fabric Manager (FM) [7] consists of several elements that are required to manage the fabric. The elements are listed as follows:

### Subnet Manager (SM)

- Subnet Manager usually assigns Local Identifiers (LIDs) to each port connected to the fabric
- Manages the routing table of the fabric
- Link and port initialization

- Sweeping the fabric to discover topology changes, managing those changes when nodes are added or deleted, etc.
- Adaptive routing
- Congestion Control

### **Subnet Administration**

Subnet Administrator (SA) actively engages with Subnet Manager (SM) to store and retrieve the fabric information. The SM/SA is single unified entity. With the help of SA messages, nodes connected to the fabric can have node-to-node path information, fabric topology and configuration, event notifications etc.

### **Performance Manager**

Performance Manager (PM) communicates with Performance Manager Agent (PMA) on each node of the fabric to gain statistical information like link utilization bandwidth, link packet rates, link congestion, error statistics (packet discards, packet routing errors etc.).

### **Performance Administration**

Performance Administration (PA) actively engages with Performance Manager (PM) to store and retrieve fabric performance information. The PM/PA is single unified entity. With the help of PA messages, management nodes on the fabric can gain access to fabric information like overall health, fabric utilization, packet information, counters and status for specific port etc.

### **Fabric Executive**

Fabric Executive (FE) provides out-of-band access to the FM. It communicates with Fabric Manager GUI over TCP/IP.

## **3.5.1 Embedded Subnet Manager**

Embedded Subnet Manager (ESM) is a switch-based Subnet Manager which deploys all the FM components as an embedded solution in an internally managed switch. ESM can manage small clusters (up to 128 nodes). Embedded FM can be managed through CLI of the switch or through the chassis viewer of the switch.

### **Controlling Subnet Manager using Chassis Viewer**

From the chassis viewer of the switch, click on "subnet manager" -> "control" to start/stop/restart the subnet manager.





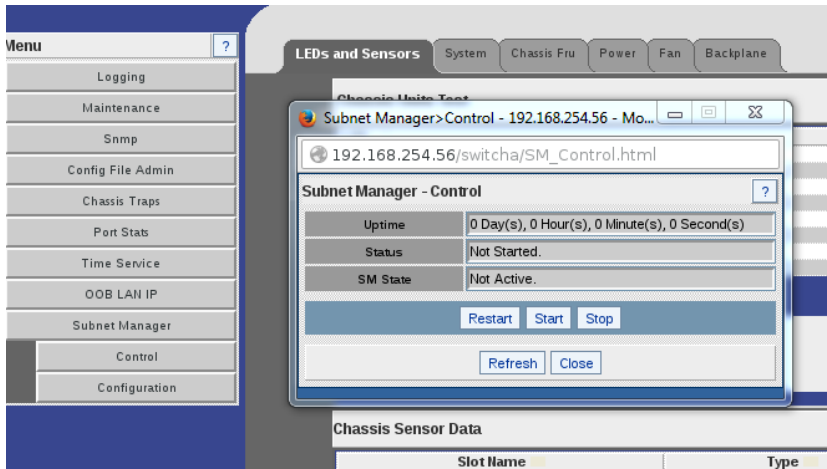


Figure 9 Starting Embedded Subnet Manager

### Controlling Subnet Manager using Switch CLI

The default login credentials for the switch is admin/adminpass. From the switch CLI, use the following commands to manage the SM.

- smcontrol start
- smcontrol stop
- smcontrol restart
- smcontrol status

## 3.5.2 Host-Based Fabric Manager

Host-based Fabric Manager can run on a compute node/host. This package is available in the Intel Fabric Suite. It can be used to manage both small and large scale clusters; however, host-based fabric manager is recommended for cluster sizes greater than 128 nodes. Upon installation, the configuration file is located at /etc/sysconfig/opafm.xml. The opafm service can be managed using the standard Linux systemctl command. Fabric manager needs at least 1GB of physical memory and one CPU core for each Fabric Manager instance. The physical memory needed varies with the size of the cluster. For example, when managing a cluster of 10,000 nodes or more, 5GB of memory per fabric instance is required. For very large clusters with more than 16,000 nodes, 15GB of memory is required per fabric manager instance. When running multiple Fabric Manager Instances on a single management node, the physical memory requirements should be multiplied by the number of fabric manager instances. When multiple Fabric managers are running on the fabric, master subnet manager will be assigned to the host with lowest Host Fabric Interface GUID.

### Controlling Fabric Manager

- systemctl start opafm.service
- systemctl stop opafm.service
- systemctl restart opafm.service
- systemctl status opafm.service

## 4 Test bed and configuration

This section explains the server configuration, BIOS options and application versions used for the application performance study. This study utilized the Zenith cluster located in the [Dell HPC Innovation Lab](#).

Component	Details
<b>Server</b>	32 PowerEdge R630
<b>Processor</b>	Intel® Xeon® CPU E5-2697 v3 @ 2.60GHz No. of cores: 14 Processor Base Freq: 2.6GHz AVX Base Freq: 2.2GHz
<b>Memory</b>	8*8 GB @ 2133MHz
<b>Operating System</b>	Red Hat Enterprise Linux Server release 7.2 (Maipo)
<b>Kernel-release</b>	3.10.0-327.el7.x86_64
<b>BIOS version</b>	2.0.2
<b>System Profile</b>	Performance Profile <ul style="list-style-type: none"><li>• Turbomode : Enabled</li><li>• Cstates : disabled</li><li>• Nodeinterleave : disabled</li><li>• Logical processor : disabled</li><li>• Snoop mode : Cluster-On-Die</li><li>• IO-NonpostedPrefetch: Disabled</li></ul>
<b>HFI card</b>	Intel® OPA HFI Adapter 100 series
<b>Switch Firmware</b>	Dell Networking H1048 OPF Edge Switch 10.0.1.0.21
<b>Fabric Manager</b>	Host Based FM
<b>Intel OPA IFS</b>	10.0.1.1.5 ( <i>pre-release</i> )

Table 1 Server Configuration

The Intel® Scalable Solutions Framework is a combination of technologies and recommendations that together provide a set of building blocks and reference architectures for designing flexible HPC systems. The Zenith cluster, housed in Dell's HPC innovation Center, leverages the Intel® SSF in its design. Currently at 256 nodes, Zenith utilizes Intel® Xeon® Processors, Intel Omni-Path fabric, Intel® Enterprise Edition Lustre®, and various Intel software such as Intel® Parallel Studio.



Application	Version	MPI	Benchmark
<b>OSU</b>	4.4.1	OpenMPI-hfi-1.10	osu_latency osu_bw osu_bibw
<b>NAMD</b>	2.11	Intel MPI 5.1.3	Apoa1 F1atpase Stmv
<b>WRF</b>	3.8	Intel MPI 5.1.3	Conus 2.5km
<b>ANSYS® Fluent®</b>	17.0	Platform MPI 9.1.3.1	eddy_417k pump_2m aircraft_wing_2m ice_2m fluidized_bed_2m rotor_3m sedan_4m oil_rig_7m combustor_12m truck_poly_14m aircraft_wing_14m landing_gear_15m lm6000_16m
<b>CD-adapco® STAR-CCM+®</b>	11.02.010	Platform MPI 9.1.4.0	EglinStoreSeparation KcsWithPhysics TurboCharger Reactor_9m HIMach10 EmpHydroCyclone13m lemans_poly_17m civil_trim_20m

Table 2 Application and benchmarks Details

Note: OSU latency, bandwidth results are based on the Intel® Xeon® CPU E5-2697 v4 @ 2.30GHz Broadwell processor and rest of the application results are as mentioned in Table 1.



## 5 Performance Benchmarking Results

### 5.1 Latency

OSU Micro-benchmarks were used to determine latency. These latency tests were done in Ping-Pong fashion. HPC applications need low latency and high throughput. As seen in the graph below, the back to back latency is 0.77 $\mu$ s and switch latency is 0.9 $\mu$ s which is on par with industry standards.

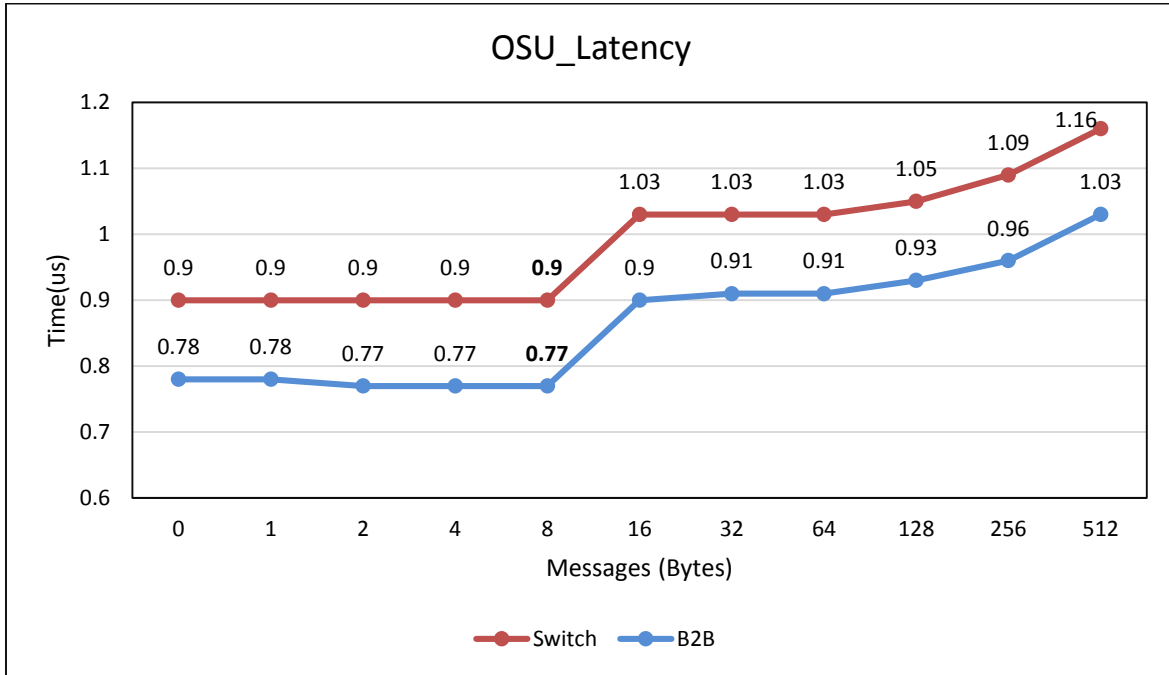


Figure 10 OSU Latency values based on Intel® Xeon® CPU E5-2697 v4 processor

### 5.2 Bandwidth

Figure 11 shows the OSU uni-directional and bi-directional bandwidth results with OpenMPI-1.10-hfi version. At 4MB uni-directional bandwidth is around 12.3 GB/s and bi-directional bandwidth is around 24.3 GB/s which is on par with the theoretical peak values.

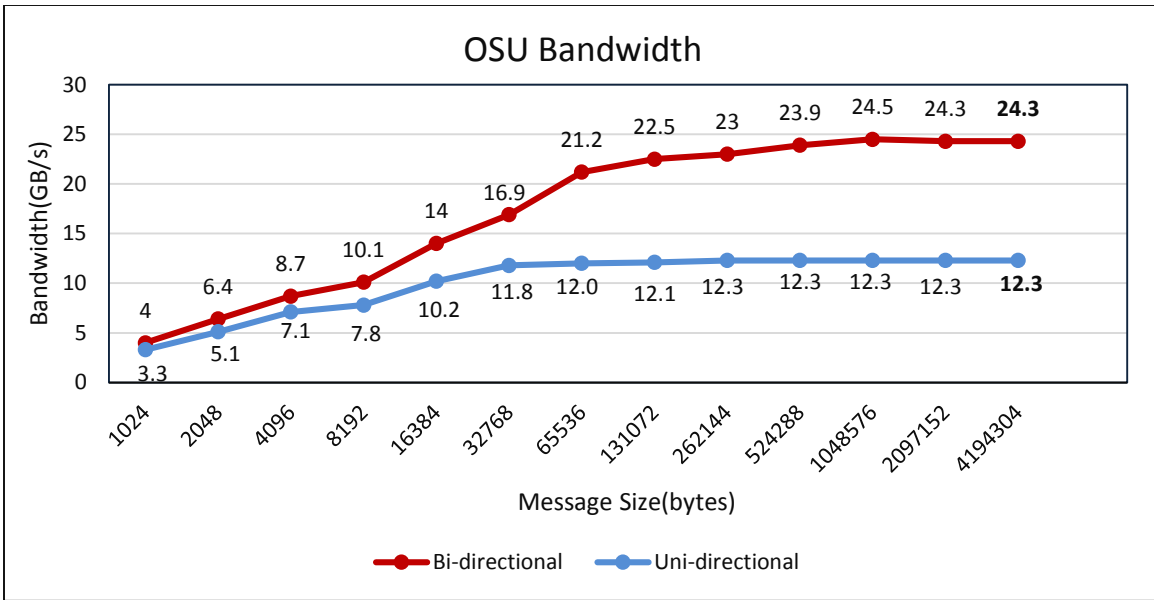


Figure 11 OSU Bandwidth values based on Intel® Xeon® CPU E5-2697 v4 processor.

### 5.3 Weather Research Forecast

Weather Research and Forecasting Model [8] is a weather prediction system designed for atmospheric research and operational forecasting needs. The WRF system contains two dynamic solvers, referred to as the ARW (Advanced Research WRF) core and NMM (Nonhydrostatic Mesoscale Model) core. The WRF-ARW code was used for testing.

Figure 12 illustrates the application performance for node counts scaling from 1 to 32. All the data points in the graph are relative to one node. As the node count increases most of the MPI time is spent on MPI\_Bcast for this dataset. This application scales linearly across all the nodes counts. Higher values represents better performance and at 32 nodes, a 44.2X improvement is seen over a single node.

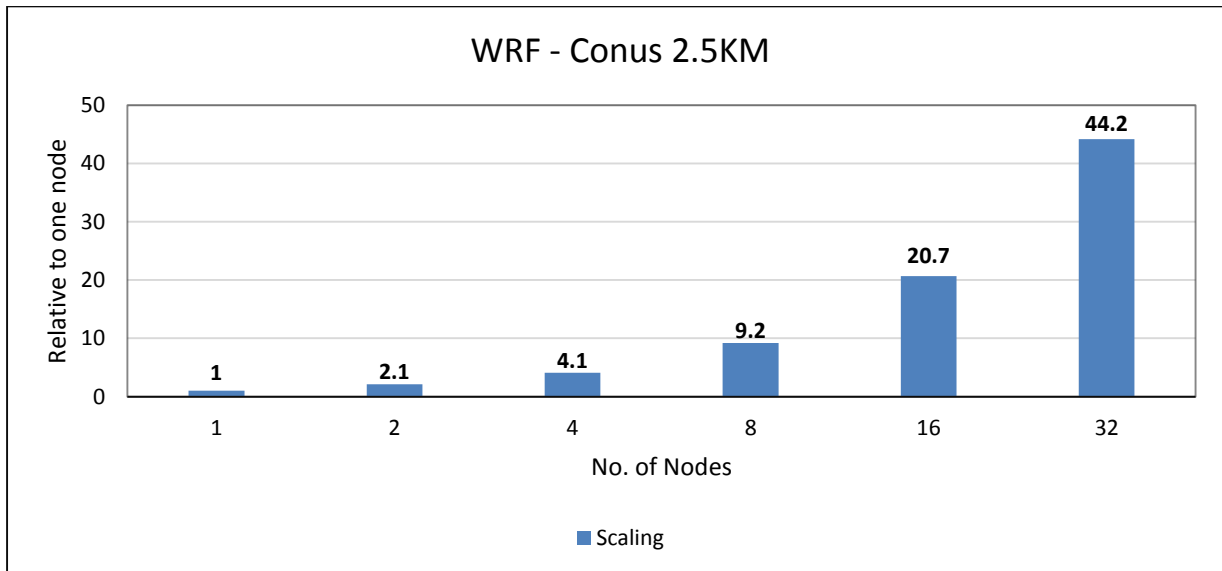


Figure 12 Weather Research Forecasting model performance graph



## 5.4 NAMD

NAMD [9] is a parallel molecular dynamics code designed for high-performance simulation of large biomolecular systems. Three proteins ApoA1 (92,224 atoms), F1ATpase (327,506 atoms) and STMV (1,066,628) are used in our study due to their relatively large problem size.

Figure 13 illustrates the performance of NAMD for three different datasets apoA1, f1atpase and stmv. All the datasets show results for node counts from 1 to 32. All the data points in the graph are relative to one node and a higher value represents better performance. Most of the MPI time is spent on broadcasting for these datasets and the performance scales linearly at higher node counts.

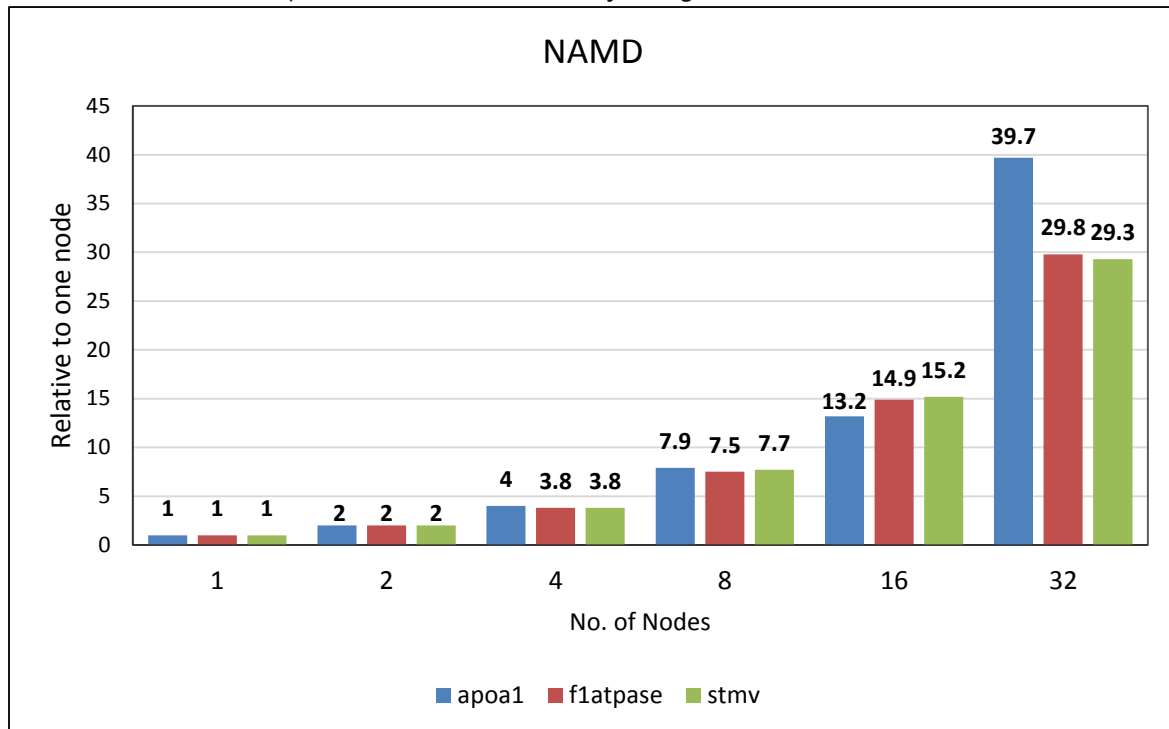


Figure 13 NAMD performance graph

## 5.5 ANSYS® Fluent®

At the time of publication of this whitepaper, Intel Omni-Path was not officially supported by ANSYS® Fluent®. In order to obtain preliminary performance data for this application, the MPI software was modified to use the appropriate Intel Omni-Path library.

Multiple cases from Fluent benchmark suites v15 and v16 were tested on the lab test system. The relative performance of twelve benchmark cases are presented in this section.

The graphs in Figure 14 and Figure 15 show the relative performance of the benchmarks on 1 to 32 nodes using 28 to 896 cores. Each data point on the graphs represents the relative performance of the specific benchmark data set using the number of cores marked on the x-axis in a parallel simulation. The results are presented as performance relative to the performance of a single node or 28 cores. A higher value represents better performance. The results are divided into two charts for easy readability. Figure 14 presents benchmarks that scale relatively well up to 32 nodes and Figure 15 presents benchmarks that do not scale as well, due to smaller model sizes.

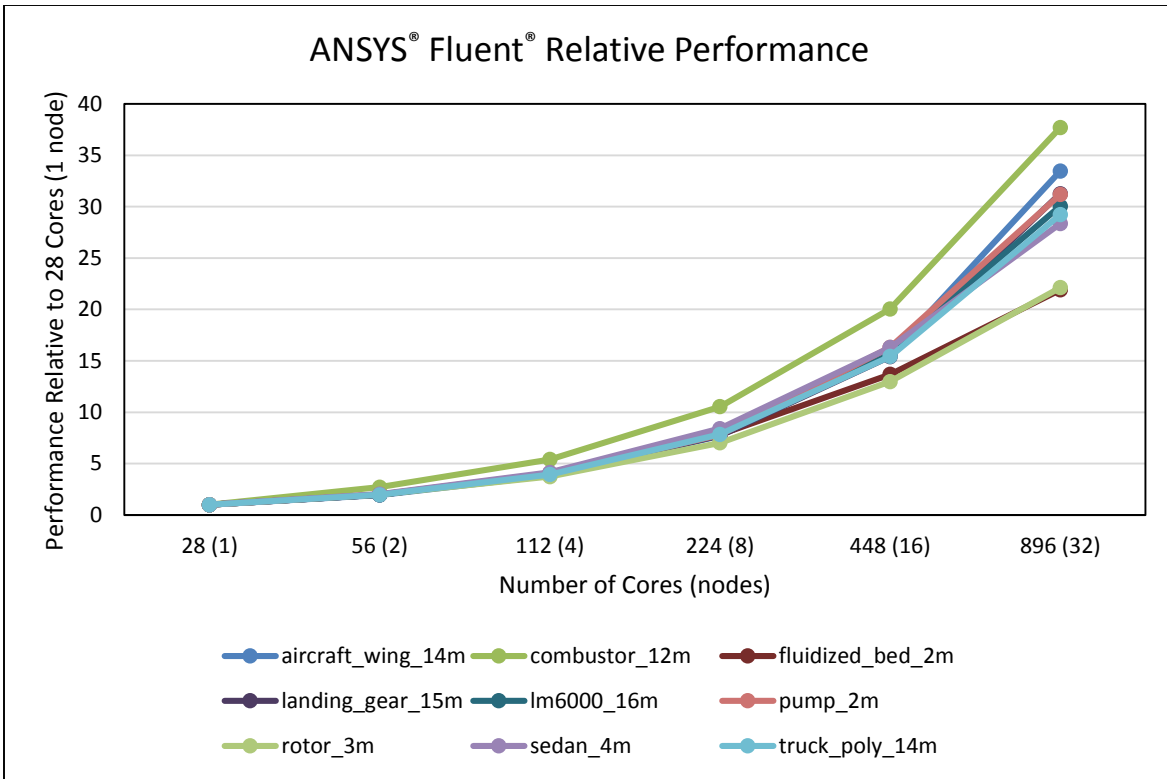


Figure 14 ANSYS® Fluent® Relative Performance Graph (1/2)

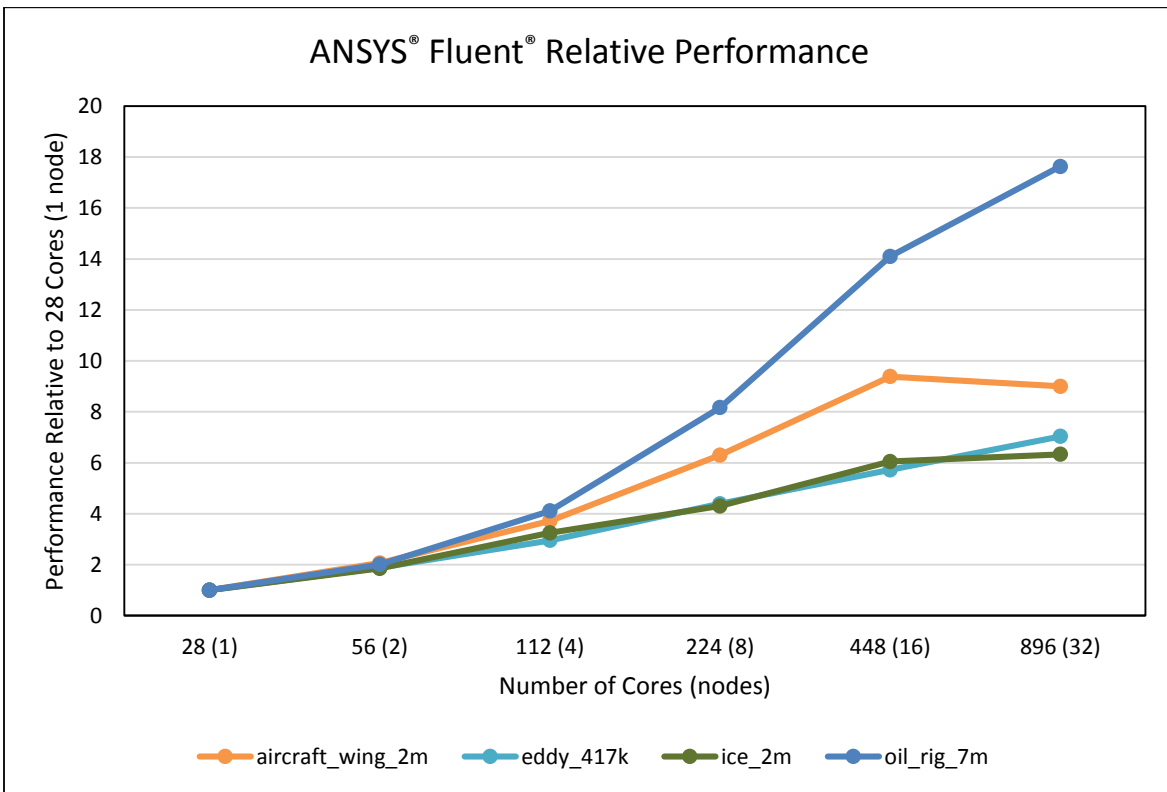


Figure 15 ANSYS® Fluent® Relative Performance (2/2)



## 5.6 CD-adapco® STAR-CCM+®

At the time of publication of this whitepaper, Intel Omni-Path was not officially supported by CD-adapco® STAR-CCM+®. In order to obtain preliminary performance data for this application, the MPI software was modified to use the appropriate Intel® Omni-Path library.

Multiple cases from the STAR-CCM+ benchmark suite were tested on the lab test system. The relative performance of eight benchmark cases are presented in this section.

The graphs in Figure 16 and Figure 17 show the relative performance of the benchmarks on 1 to 32 nodes using 28 to 896 cores. Each data point on the graphs represents the relative performance of the specific benchmark data set using the number of cores marked on the x-axis in a parallel simulation. The results are presented as performance relative to the performance of a single node or 28 cores. A higher value represents better performance. The results are divided into two charts for easy readability. Figure 16 presents benchmarks that scale relatively well up to 32 nodes and Figure 17 presents benchmarks that don't scale as well. For all datasets, the system scalability is as expected.

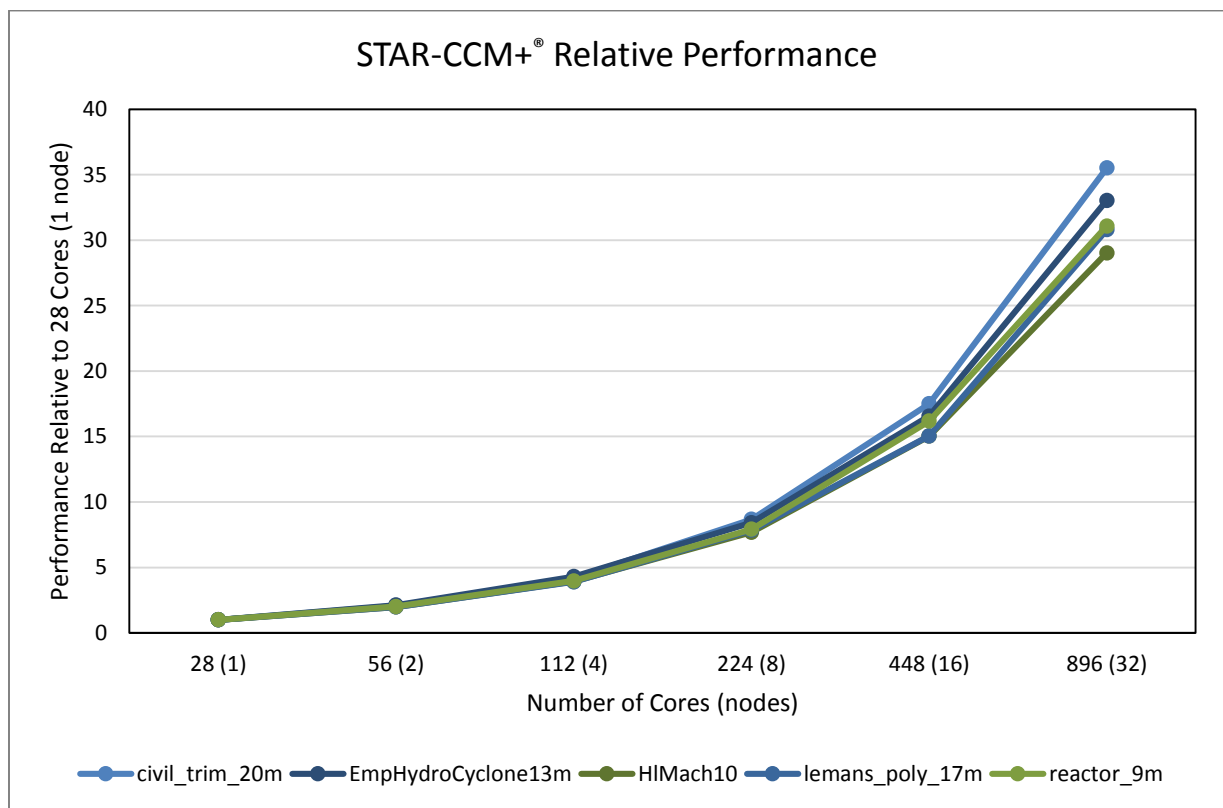


Figure 16 STAR-CCM+® Relative Performance Graph (1/2)



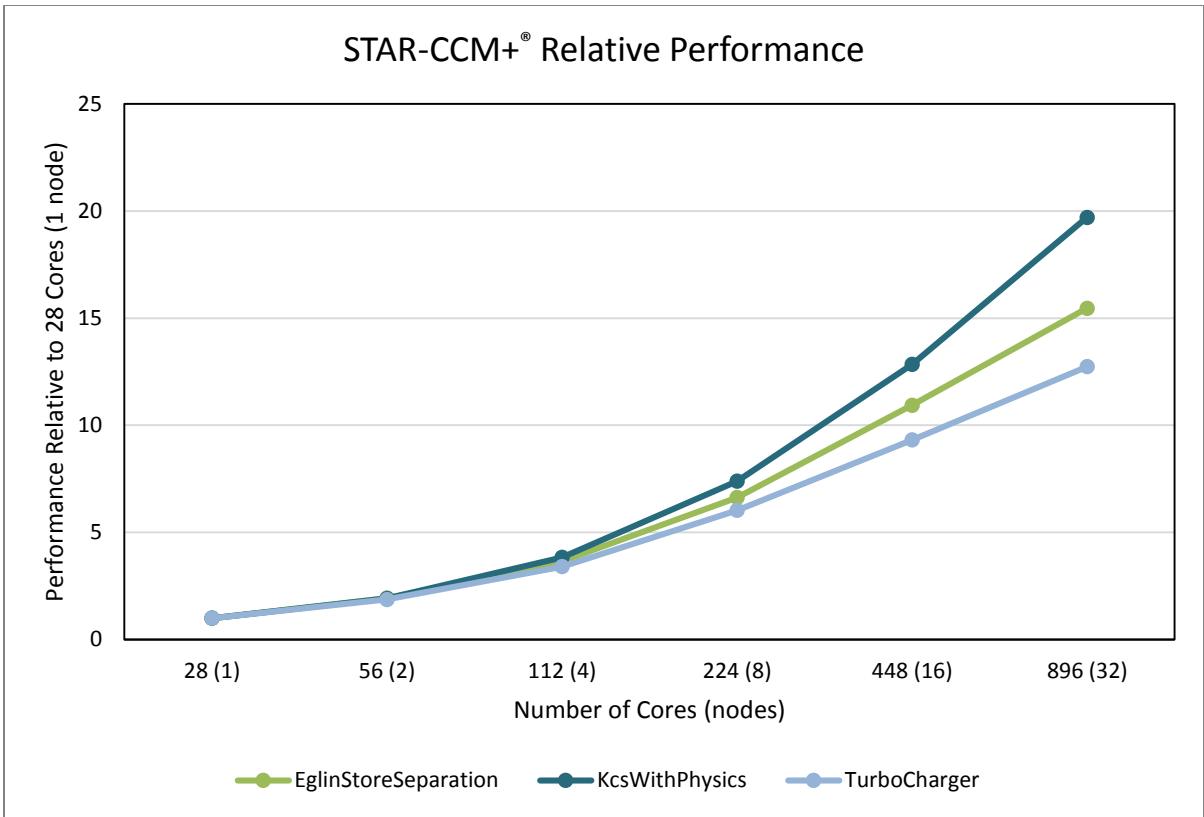


Figure 17 STAR-CCM+® Relative Performance (2/2)



## 6 Conclusion and Future Work

The Intel® Omni-Path Architecture is a new option available for low-latency, high-bandwidth cluster fabrics. The micro benchmark results prove that OPA is an ideal candidate for HPC workloads. Good application scalability is demonstrated with NAMD, WRF, STAR-CCM+, and ANSYS Fluent. Users of Intel® Omni-Path benefit from the freely available fabric monitoring and managing tools such as Fabric Manager GUI, chassis viewer, and opafastfabric.

As the drivers and software ecosystem of OPA mature further, we plan on evaluating additional fabric features such as dynamic lane scaling, packet integrity protection and traffic flow optimization, as well as continuing to evaluate the performance of various applications using Intel® Omni-Path and Intel® Broadwell Processors on Dell Servers.



## 7 References

- [1] [Online]. Available: <http://www.intel.com/content/dam/www/public/us/en/documents/product-briefs/transforming-economics-hpc-fabrics-opa-brief.pdf>.
- [2] [Online]. Available: <http://www.intel.com/content/www/us/en/high-performance-computing-fabrics/omni-path-architecture-fabric-overview.html>.
- [3] [Online]. Available: <http://www.intel.com/content/www/us/en/high-performance-computing-fabrics/omni-path-host-fabric-interface.html>.
- [4] [Online]. Available: [http://www.intel.com/content/dam/support/us/en/documents/network-and-i-o/fabric-products/Intel\\_OP\\_Fabric\\_Software\\_IG\\_H76467\\_v2.0.pdf](http://www.intel.com/content/dam/support/us/en/documents/network-and-i-o/fabric-products/Intel_OP_Fabric_Software_IG_H76467_v2.0.pdf).
- [5] [Online]. Available: [http://www.intel.com/content/dam/support/us/en/documents/network-and-i-o/fabric-products/Intel\\_OP\\_FabricSuite\\_Fabric\\_Manager\\_GUI\\_UG\\_H76471\\_v2\\_0.pdf](http://www.intel.com/content/dam/support/us/en/documents/network-and-i-o/fabric-products/Intel_OP_FabricSuite_Fabric_Manager_GUI_UG_H76471_v2_0.pdf).
- [6] [Online]. Available: [http://www.intel.com/content/dam/support/us/en/documents/network-and-i-o/fabric-products/Intel\\_OP\\_FabricSuite\\_FastFabric\\_CLI\\_RG\\_H76472\\_v2\\_0.pdf](http://www.intel.com/content/dam/support/us/en/documents/network-and-i-o/fabric-products/Intel_OP_FabricSuite_FastFabric_CLI_RG_H76472_v2_0.pdf).
- [7] [Online]. Available: [http://www.intel.com/content/dam/support/us/en/documents/network-and-i-o/fabric-products/Intel\\_OP\\_FabricSuite\\_Fabric\\_Manager\\_UG\\_H76468\\_v2\\_0.pdf](http://www.intel.com/content/dam/support/us/en/documents/network-and-i-o/fabric-products/Intel_OP_FabricSuite_Fabric_Manager_UG_H76468_v2_0.pdf).
- [8] [Online]. Available: <http://www.wrf-model.org/index.php>.
- [9] [Online]. Available: <http://www.ks.uiuc.edu/Research/namd/>.

