# Legal Disclaimer

# WHO IS IN THE AUDIENCE?

- Fabric administrators
- MPI developers
- Knowledge of
  - InfiniBand
  - True Scale (aka QLogic InfiniBand)

# RELEVANCE TO THIS AUDIENCE?

- Omni-Path is available on Cray CS Series Clusters
- Future?

# Intel Fabrics over time



2014    2015    2016    FUTURE

**HPC Fabrics**

Intel® True Scale QDR40/80  40/2x40Gb

Intel® Omni-Path gen1

Future Intel® Omni-Path

Ethernet    10Gb    40Gb    100Gb

**Enterprise & Cloud Fabrics**

Forecast and Estimations, in Planning & Targets

Potential future options, subject to change without notice. Codenames.
All timeframes, features, products and dates are preliminary forecasts and subject to change without further notification.

# Omni-Path (OPA): a quick introduction

- It's like InfiniBand?
    - Yes, very much like InfiniBand
    - Switches, adapters and cables
    - Verbs/RDMA and PSM APIs
    - LIDs and GUIDs, Fat-trees and Subnet Managers
    - Similar admin commands to True Scale

- But it's not InfiniBand
    - Link Layer is different – More functionality
    - So cannot be directly connected to InfiniBand

# Omni-Path gen1 Architecture
## *OPA technology at a glance*

- Enhanced Intel® True Scale host stack on new 100Gb hardware
  - Link layer from Cray* Aries:
    - Packet pre-emption and interleaving minimises the impact of large storage packets on latency sensitive MPI traffic
    - Error correction optimised for latency
    - Enhanced to 100Gb
    *Significant scalability benefits over InfiniBand* roadmap.*
  - Host stack from True Scale
    - PSM: Connectionless tag-matching protocol
    - Proven scalable HPC platform

- Integration
  - Developing over time with each CPU generation

*Other names and brands may be claimed as the property of others.

# INTEL® OMNI-PATH ARCHITECTURE
## EVOLUTIONARY APPROACH, REVOLUTIONARY FEATURES, END-TO-END SOLUTION

| HFI Adapters | Edge Switches | Director Switches | Silicon | Software | Cables |
|---|---|---|---|---|---|
| *Single port* **x8 and x16** | *1U Form Factor* **24 and 48 port** | *QSFP-based* **192 and 768 port** | *OEM custom designs* **HFI and Switch ASICs** | *Open Source* **Host Software and Fabric Manager** | *Third Party Vendors* **Passive Copper Active Optical** |
| **x16 Adapter (100 Gb/s)** | **48-port Edge Switch** | **768-port Director Switch (20U chassis)** | **HFI silicon Up to 2 ports (50 GB/s total b/w)** | | |
| **x8 Adapter (58 Gb/s)** | **24-port Edge Switch** | **192-port Director Switch (7U chassis)** | **Switch silicon up to 48 ports (1200 GB/s total b/w)** | | |

## Building on the industry's best technologies

- Highly leverage existing Aries and Intel® True Scale fabric
- Adds innovative new features and capabilities to improve performance, reliability, and QoS
- Re-use of existing OpenFabrics Alliance* software

## Robust product offerings and ecosystem

- End-to-end Intel product line
- >100 OEM designs[1]
- Strong ecosystem with 70+ Fabric Builders members

# CPU-FABRIC INTEGRATION
## WITH THE INTEL® OMNI-PATH ARCHITECTURE

**INTEGRATRION**

**TIME**

## KEY VALUE VECTORS

✓ **Performance**

✓ **Density**

✓ **Cost**

✓ **Power**

✓ **Reliability**

**Next generation**
Additional integration, improvements, and features

**Tighter Integration**

Intel® OPA

Intel XEON PHI inside

**Multi-chip Package Integration**

Twinax Cable
Twinax Cable
Connector
Intel® OPA

intel XEON PHI inside

intel XEON inside

**Intel® OPA HFI Card**

Next Intel® Xeon® Phi™ processor (Knight Hill)

Future Intel® Xeon® processor (14nm)

Intel® Xeon Phi™ processor (Knights Landing)

Next-Generation Intel® Xeon® processor

Intel® Xeon® processor E5-2600 v3

9

# What integration looks like

## Xeon Phi: KNL-F



(2) Internal-to-Faceplate Processor (IFP) cable supporting two-ports

Bottom view of card

PCIe carrier board, 2-port version (sideband cable and IFT connectors and cages on underside of the card)

EACH port requires:

(1) Internal Faceplate Transition (IFT) Connector

(1) IFT Cage

Top view of card

QSFP sideband header (cable not shown) Connects to header on motherboard

# The host sw stack, and PSM

# INTEL® OPA LINK LEVEL INNOVATION STARTS HERE
# LAYER 1.5: LINK TRANSFER LAYER

## InfiniBand*

**Application generates messages**

**Message segmented in packets of up to Maximum Transfer Unit (MTU) size**

**MPI Message**

256 B → 4 KB[1]

⋮

256 B → 4 KB[1]

**Packets sent until entire message is transmitted**

## Intel® Omni-Path Fabric

**New MTU Sizes Added**

**8K/10K**

The HFI segments data into 65-bit containers called **Fl**ow Control Dig**its** or **"Flits"**

Link Transfer Packets (LTPs) are created by assembling 16 Flits together (plus CRC)

LTPs send Flits sent over the **FABRIC** until the entire message is transmitted

CRC

Data Flit(s)

Control Flit(s) (optional)

**{ 1 Flit = 65 bits }**

**{ 16 Flits = LTP }**

[1] Intel® OPA supports up to 8KB for MPI Traffic and 10KB MTU for Storage

CRC: Cyclic Redundancy Check

**Goals:** Improved resiliency, performance, and consistent traffic movement

# Layer Innovation:
## Traffic Flow Optimization (TFO) – <u>Enabled</u>

Host Ports

ISL Ports

Intel® Omni-Path Architecture

48 Radix Switch

**Traffic Enters the Fabric**

**Same Priority Packet C Transmits after Packet A Completes**

**Packet A Suspended to Send High Priority Packet B then Packet A Resumes to Completion**

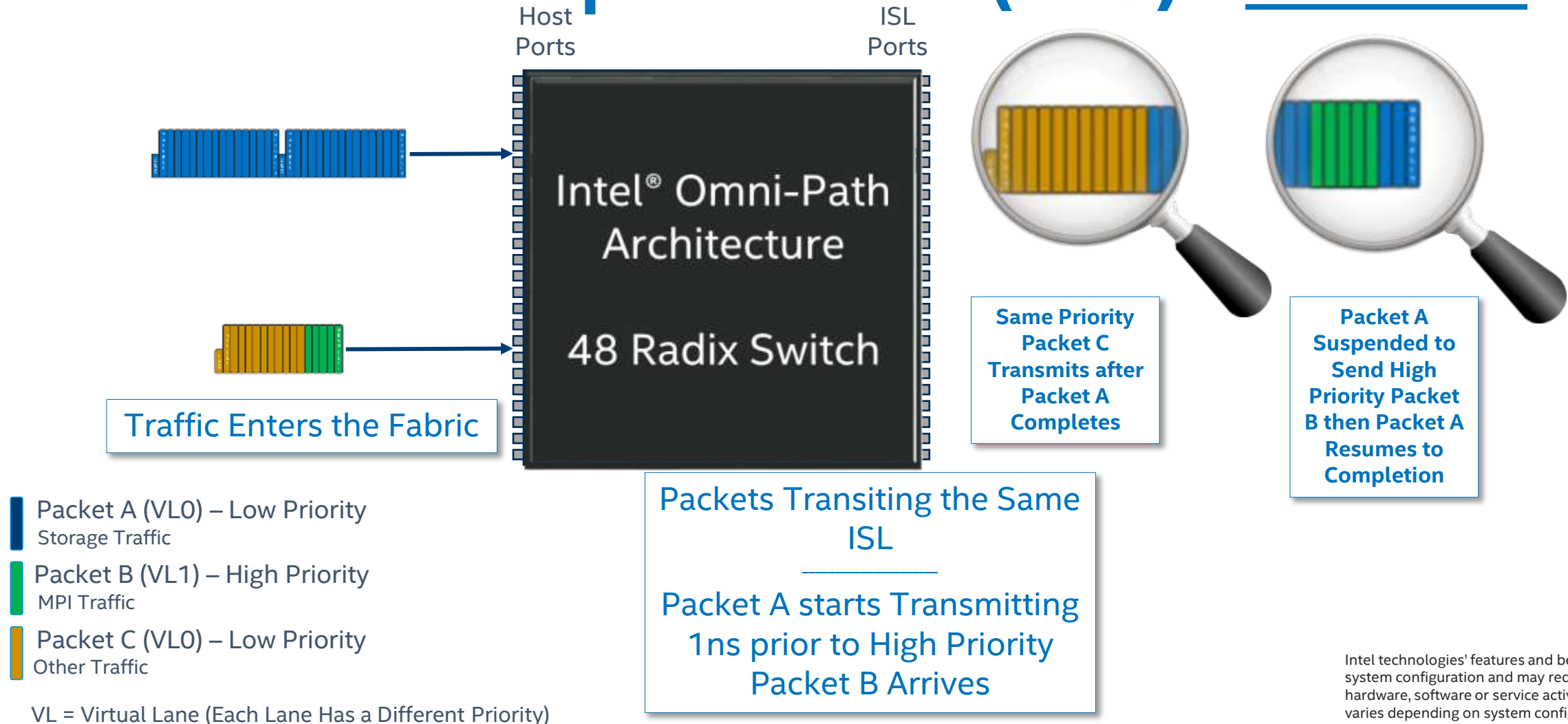**Packets Transiting the Same ISL**

_____

**Packet A starts Transmitting 1ns prior to High Priority Packet B Arrives**

**Packet A (VL0) – Low Priority**
Storage Traffic

**Packet B (VL1) – High Priority**
MPI Traffic

**Packet C (VL0) – Low Priority**
Other Traffic

VL = Virtual Lane (Each Lane Has a Different Priority)

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration.
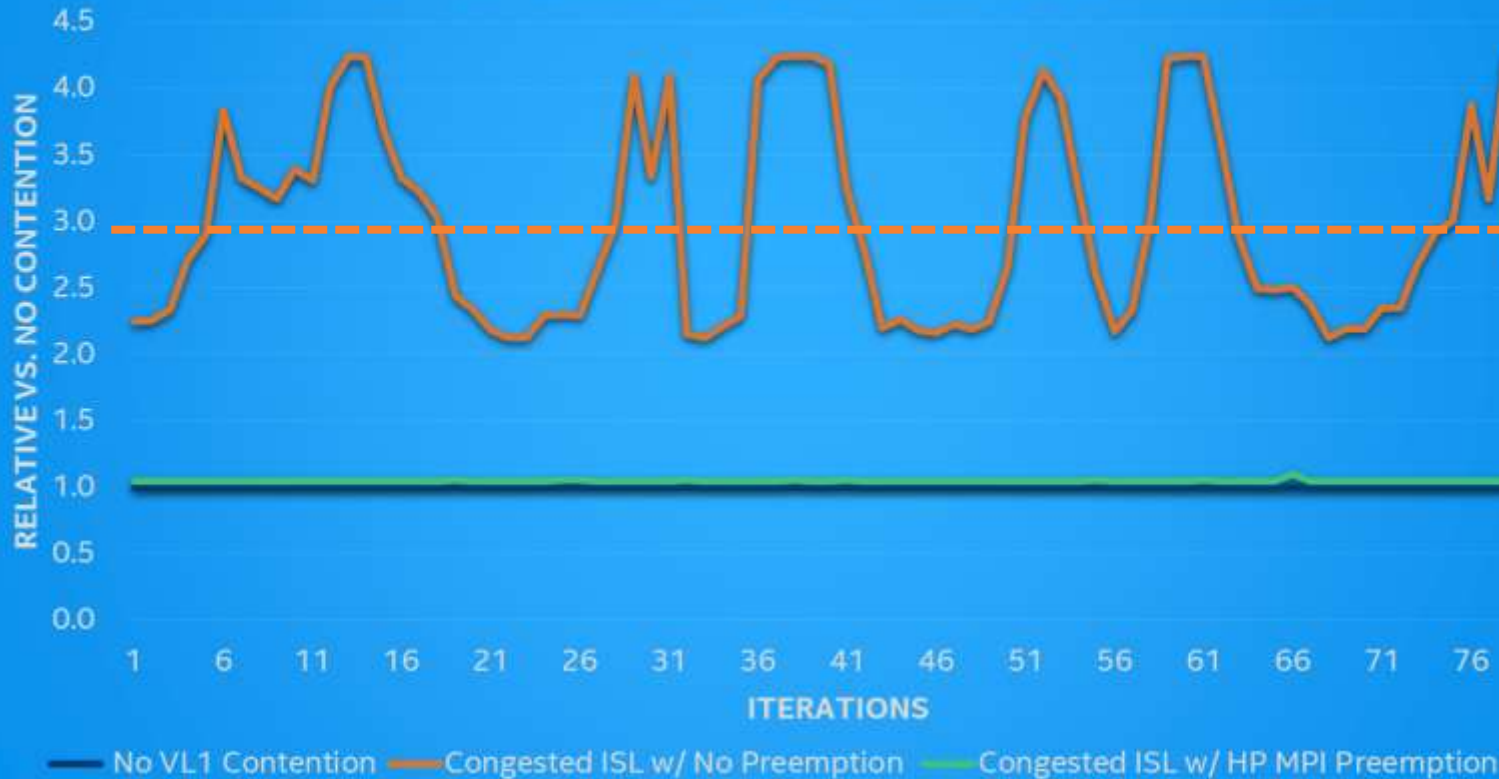
# Traffic Flow Optimization (TFO): MPI Performance Results

## – Qos Under Congested Link Conditions



VL1 Preemption: Latency Variation Results
Bi-directional Bulk Data Bandwidth

**MPI Job Running over an ISL Containing Bi-directional Bulk Data ~24.7GB on a Separate VL**

**No Prioritization with Data Contention**

**TFO Off (No Preemption) Average Latency**

**High Priority MPI Traffic with Contention**

**TFO Enabled (Preemption)**

**Relative Base MPI Latency – No Congestion**

Based on preliminary Intel internal testing using two pre-production Intel® OPA edge switch (A0)es with one inter-switch link, comparing MPI latency over multiple iterations with varying bandwidth allocations for storage and MPI traffic over multiple virtual lanes, both with Traffic Flow Optimization enabled and disabled.

BW allocation 10%/80% - (Avg. 80 Iterations)
**See Test Setup:   - Server Configuration:** Intel(R) Xeon(R) CPU E5-2699 v3 @ 2.30GHz, Turbo Disabled, Intel OPA 10.0.0.990.48 Software, ,RHEL 7.0, Kernel 3.10.0-123.el7.x86_64

# STORAGE: CONNECTING TO NEW AND EXISTING SYSTEMS

NEW systems:

- Key HPC storage vendors will deliver Intel® OPA-based storage devices

Accessing storage in EXISTING systems:

- Multi-homed solution
  - Direct-attach Intel® OPA to existing file system server along with the existing fabric connection
- Router solution
  - Lustre: Supported via LNET Router
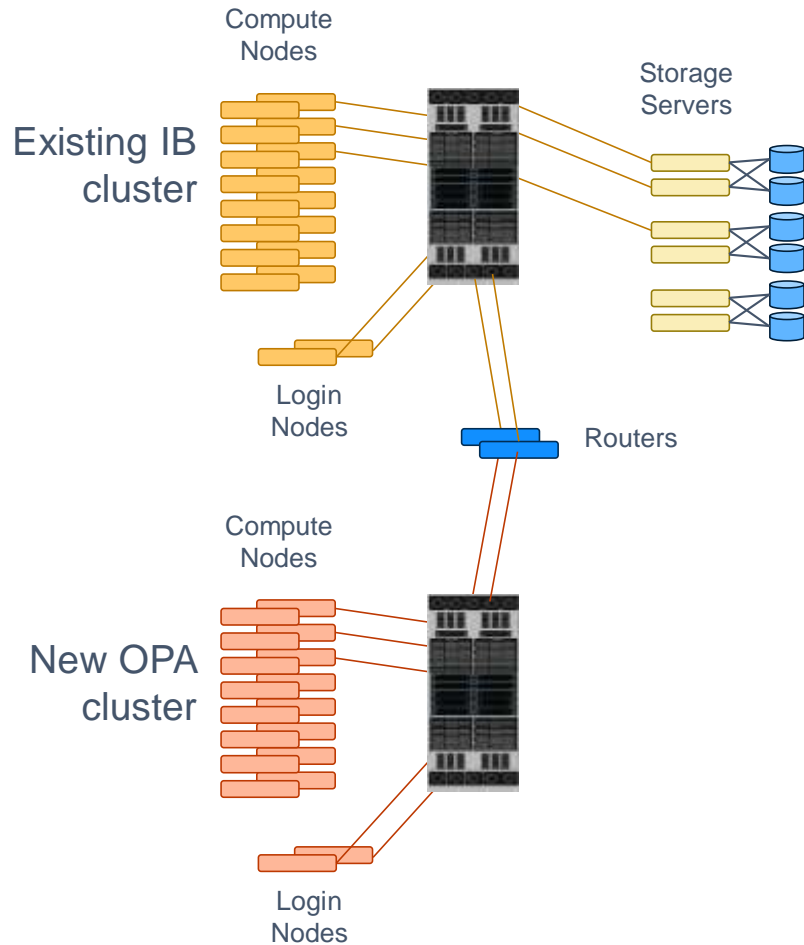  - GPFS/NAS/Other: Supported via IP Router

"Implementing Storage in Intel® Omni-Path Architecture Fabrics" white paper available now (*public link*)
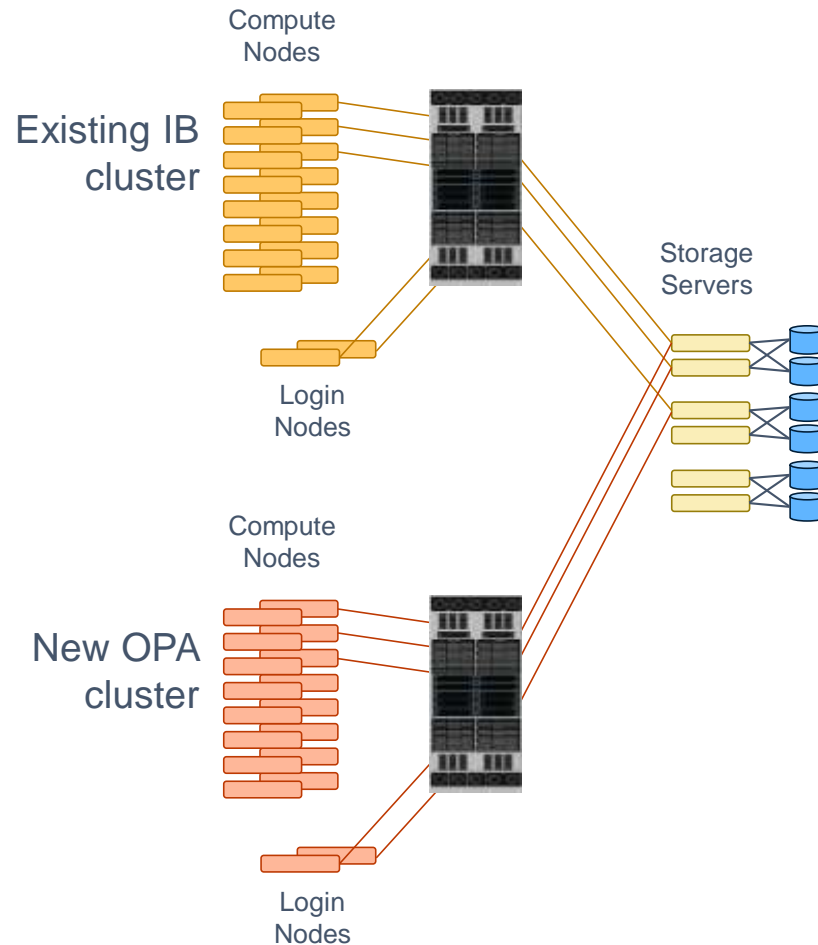"Intel® Omni-Path Storage Router Design Guide" available now (ask for access)

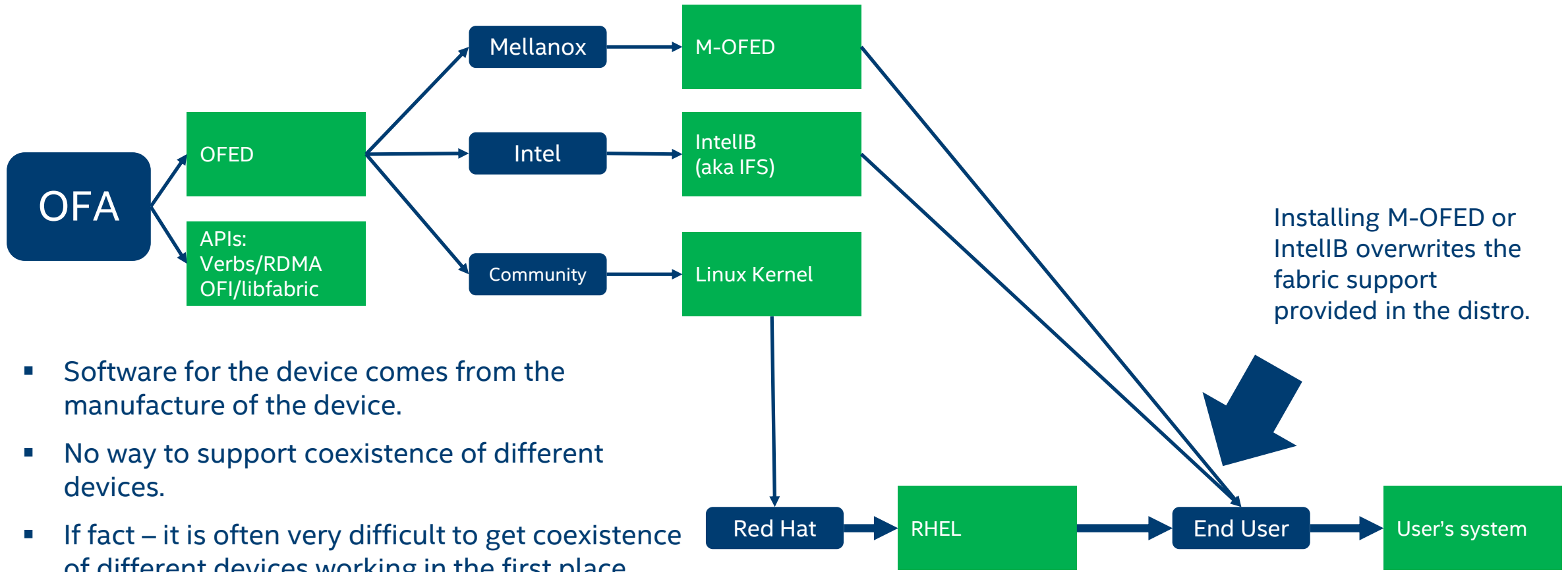# Accessing Existing Storage



Router solution

Multi-Homed Solution

Key enabler:
Interfaces must co-exist

Compute Nodes

Storage Servers

Existing IB cluster

Login Nodes

Routers

Compute Nodes

New OPA cluster

Login Nodes

Compute Nodes

Existing IB cluster

Login Nodes

Storage Servers

Compute Nodes

New OPA cluster

Login Nodes

# How It Used To Be...

*Install manufacture's software, developed from OFED*

Organizations

Deliverables

OFA → OFED → Mellanox → M-OFED

OFED → Intel → IntelIB (aka IFS)

OFED → Community → Linux Kernel

OFA → APIs: Verbs/RDMA OFI/libfabric

Linux Kernel → Red Hat → RHEL → End User → User's system

M-OFED → End User

IntelIB (aka IFS) → End User

Installing M-OFED or IntelIB overwrites the fabric support provided in the distro.

- Software for the device comes from the manufacture of the device.

- No way to support coexistence of different devices.

- If fact – it is often very difficult to get coexistence of different devices working in the first place.

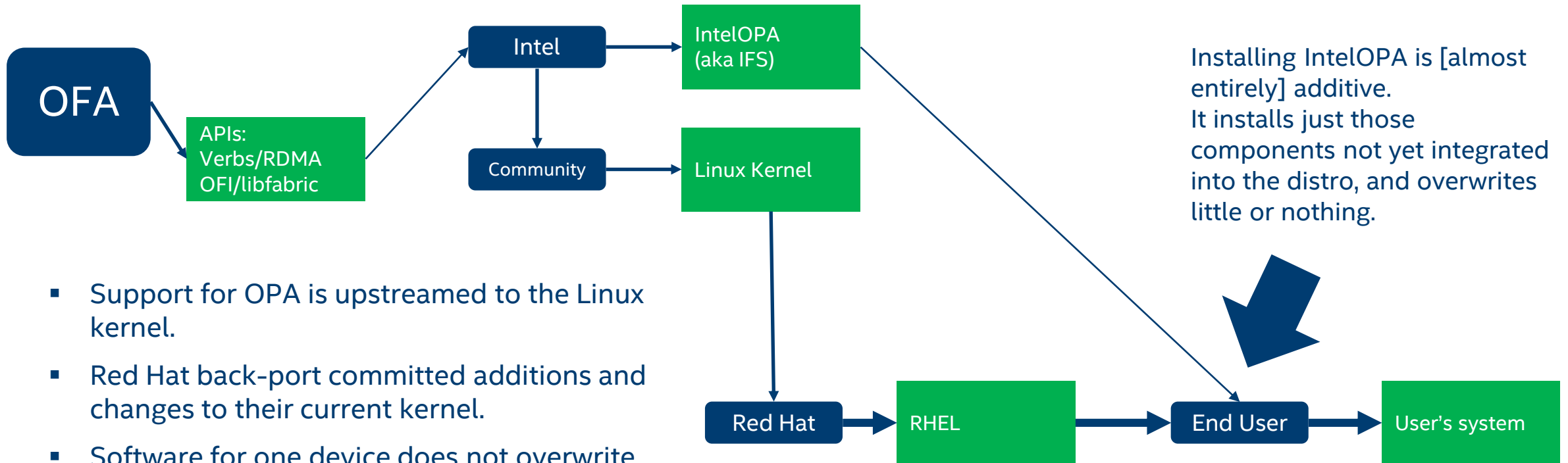- ***OS updates can break the device drivers!***

# ...How We Do It Now

## *Push support into the distro*

**Organizations**

**Deliverables**

**OFA** → APIs: Verbs/RDMA OFI/libfabric → **Intel** → IntelOPA (aka IFS)

**Intel** → **Community** → Linux Kernel

Linux Kernel → **Red Hat** → RHEL → **End User** → User's system

IntelOPA (aka IFS) → **End User**

Installing IntelOPA is [almost entirely] additive.
It installs just those components not yet integrated into the distro, and overwrites little or nothing.

- Support for OPA is upstreamed to the Linux kernel.

- Red Hat back-port committed additions and changes to their current kernel.

- Software for one device does not overwrite support for another.

- Red Hat support any combination of devices whose software is in-distro.

- *OS updates can be applied safely*

# Acceptance: Intel® Fabric Builders



https://fabricbuilders.intel.com/

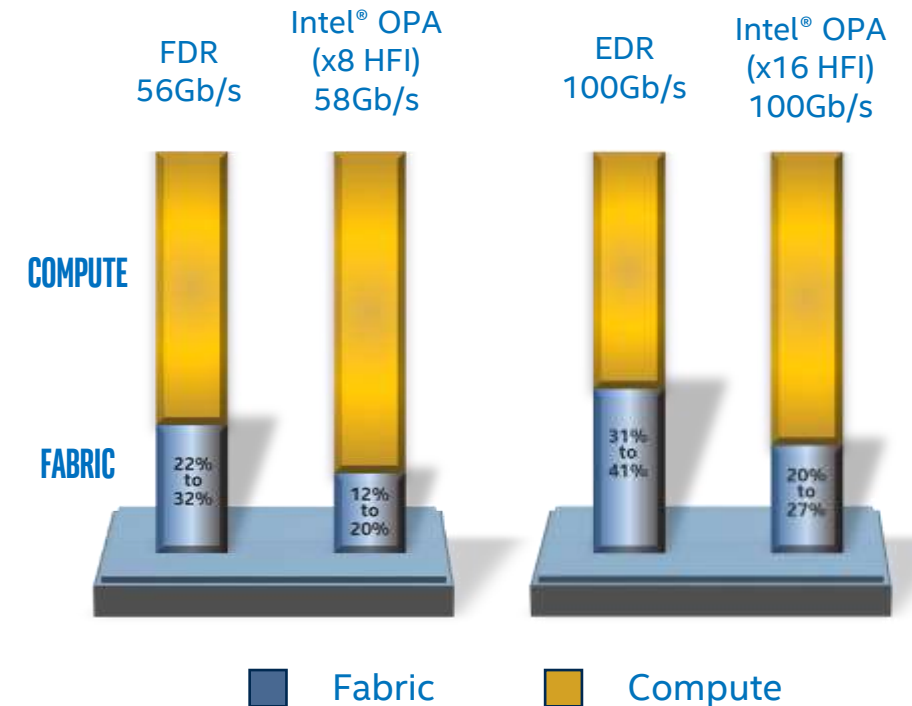*Last update November 13, 2015*

# Cost advantages

Compute / interconnect cost ratio has changed

- Compute price/performance improvements continue unabated
- Current corresponding fabric metrics unable to keep pace as a percentage of total cluster costs which includes compute and storage

Challenge: Keeping fabric costs in check to free up cluster $$$ for increased compute and storage capability

**Up to 10% lower cluster cost mix than either FDR or EDR (at similar bandwidths)**

**Hardware Cost Estimate[1]**

FDR 56Gb/s

Intel® OPA (x8 HFI) 58Gb/s

EDR 100Gb/s

Intel® OPA (x16 HFI) 100Gb/s

COMPUTE

FABRIC

22% to 32%

12% to 20%

31% to 41%

20% to 27%

■ Fabric    ■ Compute

**More Compute = More FLOPS**

# Intel® OPA Momentum is Building Quickly

## Worldwide design wins keep rolling in
## >100 OEM platform, switch, and adapters expected in 1H'16[1]

**Worldwide Wins**

- Penguin: US DoE CTS-1
- Dell: TACC Stampede 1.5
- HPE: Pittsburgh Supercomputing

- Inspur: Qingdao, Tsinghua University
- Dell: NCAR, NASA, Uni Colorado
- Sugon: Beijing Academy of Science and Technology

> >800 nodes deployed in 3 days!

> 280 nodes pre-stage "just worked"

**EMEA Wins**

- Clustervision: AEI Potsdam, University of Hull
- Dell: Wartsilla, University of Sheffield
- Lenovo: Cineca
- Cray: AWI, Juelich

- Plus more that cannot be named at this time

[1] Source: Intel internal information