# Tecnologie e progetti innovativi in HPC e Cloud: alcune soluzioni e esperienze in Lenovo

Marco Briscolini - mbriscolini@lenovo.com

Workshop di CCR

LNGS,  22-26 Maggio, 2017

# Agenda

- HPC segment and trends
- Solution components
- Technology trends
- Over 20PF@CINECA

# Target Segments - Key Requirements
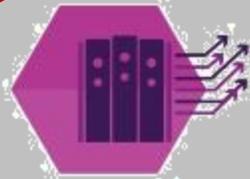
## Cloud Computing

**Key Requirements:**
- Mid-high bin EP processors
- Lots of memory (>256GB/node) for virtualization
- 1Gb/10/25/40 Ethernet

## Data Center Infrastructure

**Key Requirements:**
- Low-bin processors (low cost)
- Smaller memory (low cost)
- 1/10Gb Ethernet

## High Performance Computing

**Key Requirements:**
- High bin EP processors for maximum performance
- High performing memory
- Infiniband
- GPU support

## Data Analytics

**Key Requirements:**
- Mid-high bin EP processors
- Lots of memory (>256GB per node)
- 1Gb / 10Gb Ethernet
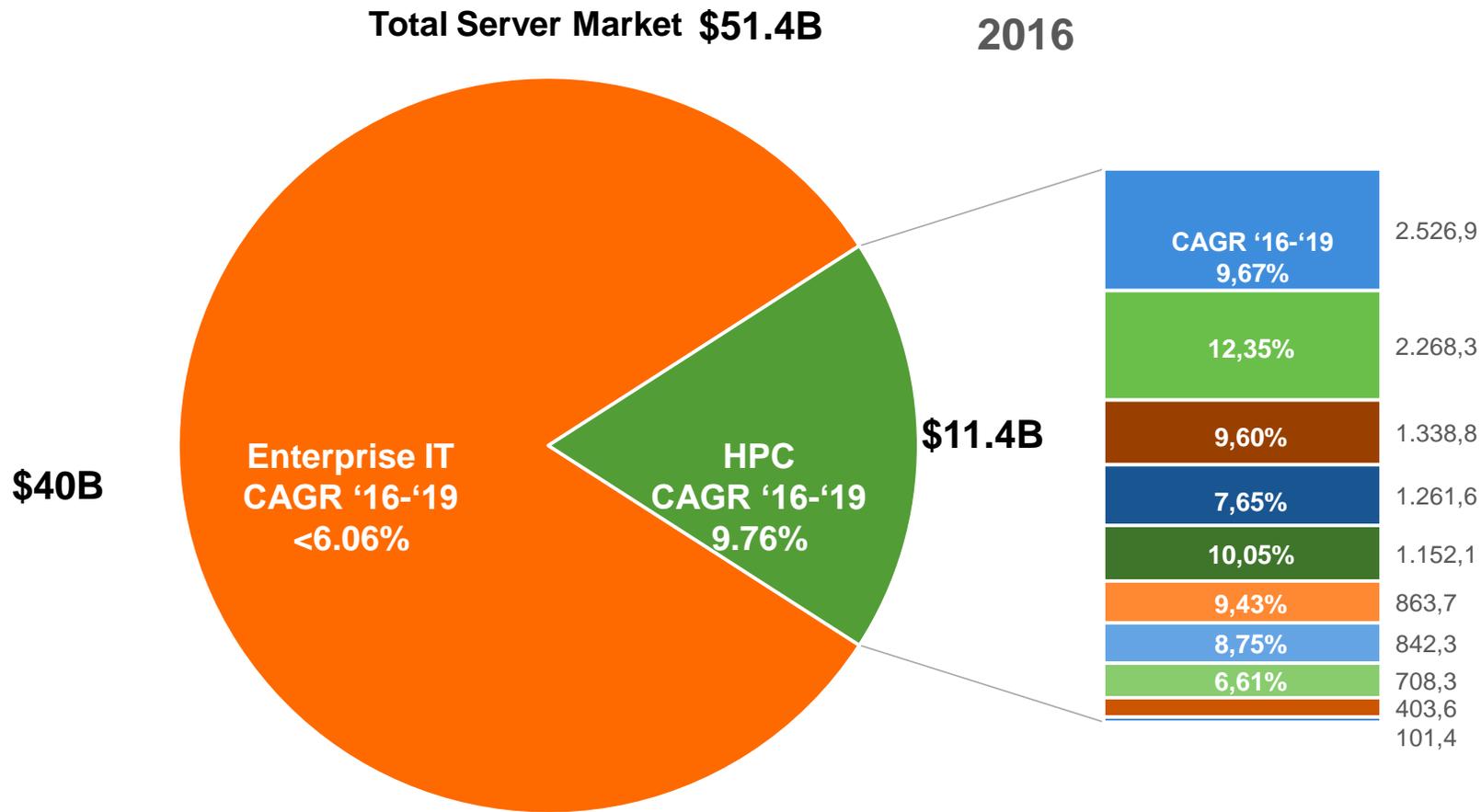- 1-2 SS drives for boot

## Virtual Desktop

**Key Requirements:**
- Lots of memory (> 256GB per node) for virtualization
- GPU support

Lenovo

3

# HPC – Fast Growing Opportunities

$11.4B Opportunity with 9.76% CAGR

**Total Server Market  $51.4B**

**2016**

**HPC is a Value Attach business**
**> 40% Storage**
**> 16% Service**

**$40B**

Enterprise IT
CAGR '16-'19
<6.06%

HPC
CAGR '16-'19
9.76%

**$11.4B**

| | CAGR '16-'19 | |
|---|---|---|
| | 9,67% | 2.526,9 |
| | 12,35% | 2.268,3 |
| | 9,60% | 1.338,8 |
| | 7,65% | 1.261,6 |
| | 10,05% | 1.152,1 |
| | 9,43% | 863,7 |
| | 8,75% | 842,3 |
| | 6,61% | 708,3 |
| | | 403,6 |
| | | 101,4 |

- Academia & Research
- Government Lab
- Manufacturing & Construction
- Life Science & Health Care
- Security & Defense
- Natural Resources
- Silicon & Software
- Agriculture, Retail & Transportation
- Finance & Insurance
- Entertainment & Communication

**>$1 of every $5 x86 spend is HPC**

Lenovo

IDC 2015/16

4

# HPC Market Trends and our Strategy

## Trends

## Strategy

**Resurgence of Specialization**
Max performance for an expanding set of workloads

**Deliver a modular platform with easy to use management stack**
Allowing clients to optimize what they have today and easily adapt new technologies

**Open Everything**
Renewed Interest in Open HW and SW Globally

**Exceed client expectations for Openess with open SW and via deep collaboration**
That results in innovation and open IP

**Co-Design is Mandatory**
Truly optimized and holistic results based designs

**Design the best solution for any given workload, budget or constraint**
Using deep skills, partnership and flexibility

**Limited Budgets; Higher Demands**
Continued demand for best performance/$ + TCO/ECO/OPEX

**Use the power of our Global Scale of Economic and Data Center experience**
To maximize impact per spend

# LENOVO IS A FULL MEMBER ON THE EUROPEAN TECHNOLOGY PLATFORM 4 HPC

ETP4HPC will define research priorities for the development of a globally competitive HPC technology ecosystem in Europe. It will propose and help to implement a Strategic Research Agenda, while acting as the "one voice" of the European HPC industry in relations with the European Commission and national authorities.
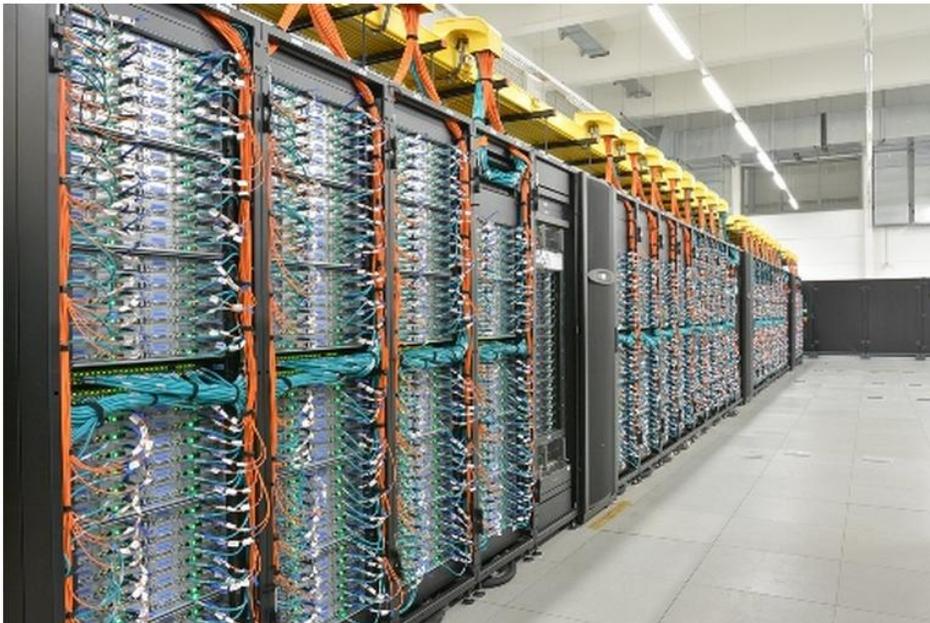
**ETP 4 HPC**
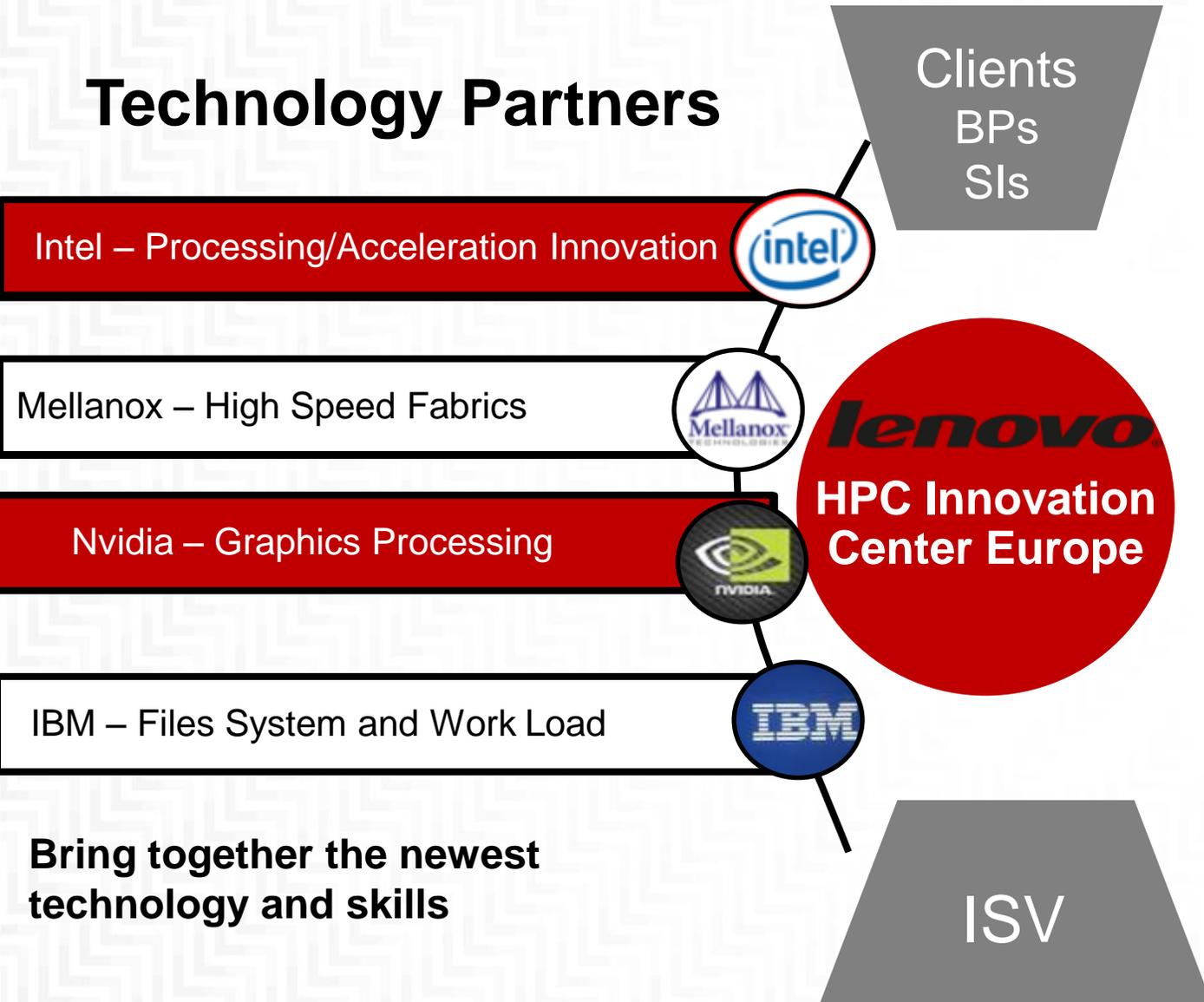THE EUROPEAN TECHNOLOGY PLATFORM FOR HIGH PERFORMANCE COMPUTING

## 1 / MEMBERS

Members | Become a member | Membership benefits

< >

**Lenovo**™

http://www.lenovo.com

# Technology Partners

**Core Client Partners**

Intel – Processing/Acceleration Innovation

Mellanox – High Speed Fabrics

Nvidia – Graphics Processing

IBM – Files System and Work Load

Clients
BPs
SIs

## HPC Innovation Center Europe

**FZJ**

**HPC Storage**

ISV

**LRZ**

**Energy Efficient**

**Systems and Software**

**RZG**

Advancing Material Science

**CINECA**

Big Data and Many Cores

**BSC**

Extreme Application Scaling

**STFC** Hartree

HPC Software and Optimization

**Bring together the newest technology and skills**

**Focused knowledge and deep skills advance the science of HPC**

**http://news.lenovo.com/news+releases/first-global-hpc-innovation-centre.htm**

# ✚ 2 X 3 PFlops SuperMUC systems at LRZ Phase 1 and Phase 2

## Phase 1

Ranked 20 and 21 in Top500 June 2015

- Fastest Computer in Europe on Top 500, June 2012
  - 9324 Nodes with 2 Intel Sandy Bridge EP CPUs
  - HPL = 2.9 PetaFLOP/s
  - Infiniband FDR10 Interconnect
  - Large File Space for multiple purpose
    - 10 PetaByte File Space based on IBM GPFS with 200GigaByte/s I/O bw



- Innovative Technology for Energy Effective Computing
  - Hot Water Cooling
  - Energy Aware Scheduling

## Phase 2

- Most Energy Efficient high End HPC System
  - PUE 1.1
  - Total Power consumption over 5 years to be reduced by ~ 37% from 27.6 M€ to 17.4 M€
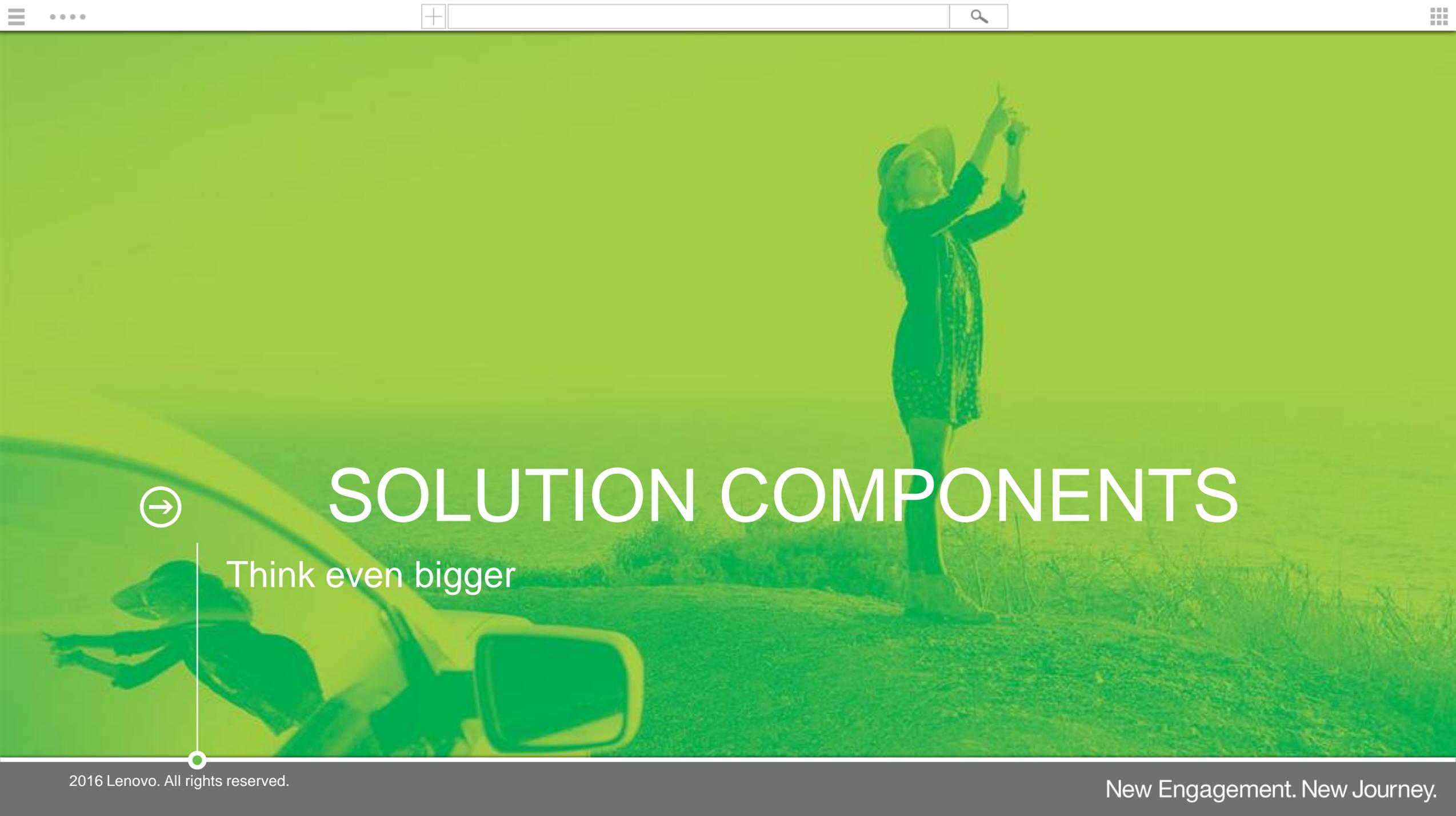
- Acceptance completed
  - 3096 nx360m5 compute nodes Haswell EP CPUs
  - HPL = 2.8 PetaFLOP/s
  - Direct Hot Water Cooled, Energy Aware Scheduling
  - Infiniband FDR14
  - GPFS, 10 x GSS26, 7.5 PB capacity , 100 GB/s IO bw

- **System A:**
- **1512 Lenovo nx360M5 ( 2 Petaflops)**
  - 21 racks
  - 126 NeXtScale WCT Chassis
  - 3,024 Intel Broadwell-EP E5-2697v4 (2.3GHz, 145W)
  - 54.432Processor Cores
  - 12.096 16GB DIMMs
- **3600 Adamspass KNL nodes ( 11 Petaflops)**
  - 50 Racks with 72 KNL nodes in Each Rack
  - 3.600 120GB SSD's
  - 244.800 cores
  - 345.600 GB RAM in 21.600 16GB DIMMs
  - 1.680 Optical cables

- **1512 Lenovo Stark nodes (>4 Petaflops)**
  - 21 racks
  - 3,024 Intel SkyLake 24c@2,1GHz
- **Over 60.000m Optical Cables**
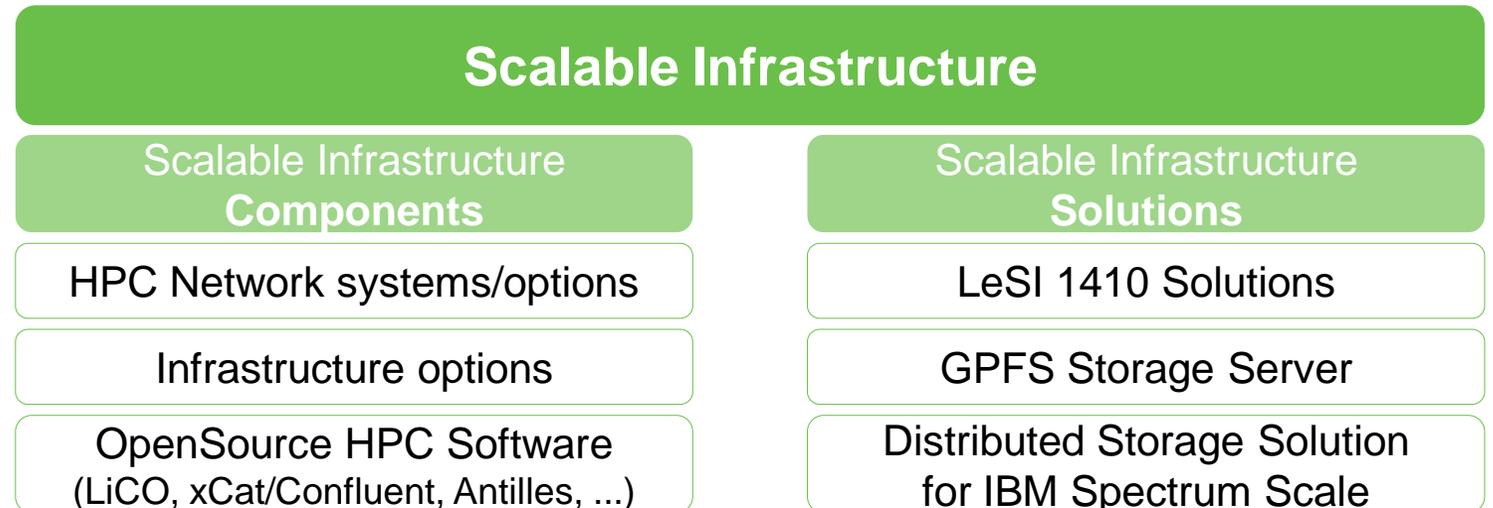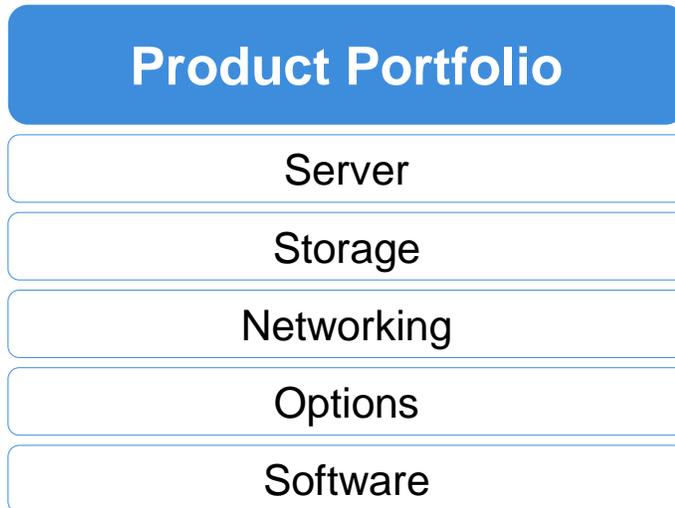- **6 GSS26 16PB raw in total >100GB/s**

# SOLUTION COMPONENTS

Think even bigger

New Engagement. New Journey.

# Lenovo Scalable Infrastructure (LeSI)

**Lenovo Scalable Infrastructure (LeSI) is a framework for development, configuration, build, delivery and support of integrated data center solutions**

- Complete HPC data center portfolio with the best-of-breed partner technology
- Collaborate on OpenSource HPC software in true commitment to Openess
- End-to-end expert-designed, tested, integrated and supported HPC solutions

| Product Portfolio | Scalable Infrastructure | |
| --- | --- | --- |
| | Scalable Infrastructure **Components** | Scalable Infrastructure **Solutions** |
| Server | HPC Network systems/options | LeSI 1410 Solutions |
| Storage | Infrastructure options | GPFS Storage Server |
| Networking | OpenSource HPC Software (LiCO, xCat/Confluent, Antilles, ...) | Distributed Storage Solution for IBM Spectrum Scale |
| Options | | |
| Software | | |

# The Combined x86 Portfolio – Delivering more choice

## High-end systems

4 socket+ enterprise-class x86 performance, resiliency, security



## Converged/Blade systems

Integration across Lenovo assets in systems and SW for maximum client optimization and value

Flex System

## Storage

Simple, Efficient, Reliable storage solutions : DAS, SAN, Tapes

## Dense systems

Optimize space-constrained data centers with extreme performance and energy efficiency

NEXTSCALE

## 1P & 2P Rack & Tower systems

Broad rack and tower portfolio to meet a wide range of client needs from infrastructure to technical computing

ThinkServer          System x

## Switches

System Networking & SAN switches for Data Centers & Virtualization needs

## Services

Warranty upgrade, maintenance, installation services, SW support, …

## SOLUTIONS

| Cloud | Analytics | Technical Computing |
|-------|-----------|---------------------|

**Management Standalone or Integration with VMware and Microsoft**
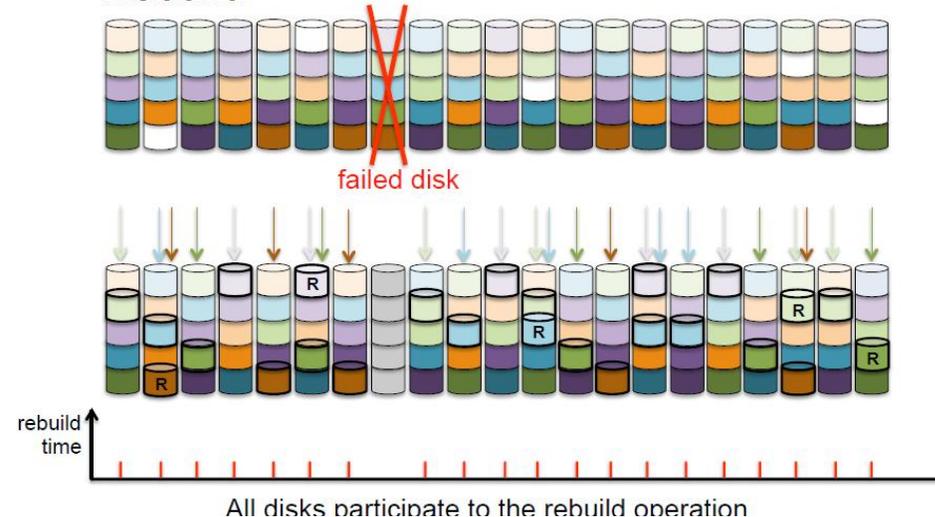
Lenovo

13

# HPC Storage

## Lenovo DSS-G

## Solution design

- **Embedded GPFS filesystem**
- **RAID support at filesystem level**
- **Fast data reconstruction by declustered RAID**
- **40GbE, FDR, EDR, OPA support**
- **Up-to 5PB raw in a system**
- **2 to 6 high density Jbod attached to two servers**
- **Reduced maintenance costs due to HW semplification**

### Declustered RAID – How it works

- Rebuild

failed disk

rebuild time

All disks participate to the rebuild operation

# DSS–G Storage Available in Either D3284 or D1224

- Lenovo D3284 JBODs (5U84)

- Lenovo D1224 JBODs (2U24)

- Two x3650M5 servers
  - SAS p2p connections to JBODs (12Gbps)
  - HPC interconnect: Ethernet, IB, OPA

- 2, 4 or 6 D3284 JBODs (5U84, 12Gbps)
  - 3.5" choice of 4,6,8,**10** TB NL-SAS disks
  - Up to **5 PB raw capacity**

- 1, 2, 4 or 6 D1224 JBODs (2U24, 12Gbps)
  - 2.5", choice of:
    - o 15K - 600GB or 300GB
    - o 10K - 1.8TB , 1.2TB, 900GB,600GB
    - o 7.2K - 2TB NL-SAS, 1TB NL-SAS
    - o SSD - 1.6TB 3 DWD, 800GB 3 DWD, 400GB 3 DWD

**DSS G201** | **DSS G202** | **DSS G204** | **DSS G206** | **DSS G220** | **DSS G240** | **DSS G260**

SSD / SAS Option for High Performance / IOPS

*Low Cost of Entry*

*Performance optimized*

DSS G201:
- D1224
- x3650M5 HPIO
- x3650M5 HPIO

DSS G202:
- D1224
- x3650M5 HPIO
- x3650M5 HPIO
- D1224

DSS G204:
- D1224
- D1224
- x3650M5 HPIO
- x3650M5 HPIO
- D1224
- D1224

DSS G206:
- D1224
- D1224
- D1224
- x3650M5 HPIO
- x3650M5 HPIO
- D1224
- D1224
- D1224

DSS G220:
- D3284
- x3650M5 HPIO
- x3650M5 HPIO
- D3284
- 164 x NL-SAS

DSS G240:
- D3284
- D3284
- x3650M5 HPIO
- x3650M5 HPIO
- D3284
- D3284
- 334 x NL-SAS

DSS G260:
- D3284
- D3284
- D3284
- x3650M5 HPIO
- x3650M5 HPIO
- D3284
- D3284
- D3284
- 502 x NL-SAS

HPIO = High Performance I/O

15

Lenovo

# Lenovo Cloud Network Operating System (CNOS)

Enables Enterprise networks to scale in cloud environments

## Resilient

- Event driven Multi process architecture
- Fault isolation for control plane stability
- High availability features

## Cloud Scale

- State of the art routing protocol stack
- 32-way multipath scale out Clos fabric
- Multi-tenant aware

## Programmable

- Enable automation at large scale
- DevOps innovation
- Native Linux shell access for server/network tools integration

# SDN/CLOUD DATA CENTER ECO-SYSTEM

# Current Lenovo HPC Software Solutions

| Customer Applications | | |
|---|---|---|
| **Debuggers & Monitoring** | **Eclipse PTP + debugger, gdb,..** | **ICINGA** / **Ganglia** |
| **Compilers & Tools** | **Intel Parallel Studio, MKL** | **Open Source Tools: FFTW, PAPI, TAU, ..** |
| **Parallel Runtime** | **Intel MPI** | **Open MPI** / **MVAPICH, IBM PMPI** |
| **Workload & Resources** | **IBM LSF** HPC & Symphony | **Adaptive Moab** / **Maui/Torque Slurm** |
| **Parallel File Systems** | **IBM GPFS** | **Lustre** / **NFS** |
| **Systems Management** | **xCat Extreme Cloud Admin. Toolkit** | **IBM PCM** |

**OS VM**  **OFED**

**NEXTSCALE**

**LenovoSystem x Virtual, Physical, Desktop, Server**

**Compute  Storage  Network** (intel) **OmniPath**  Mellanox **UFM**

- **Building Partnerships to provide the "Best In-Class" HPC Cluster Solutions for our customers**
- Collaborating with software vendors to provide features that optimizes customer workloads
- Leveraging "Open Source" components that are production ready
- Contributing to "Open Source" (i.e. xCAT, Confluent, OpenStack) to enhance our platforms
- Providing "Services" to help customers deploy and optimize their clusters

# Future HPC Open Source Management Stack

**Web Console GUI**



AN OPEN SOURCE HPC SOFTWARE COMMUNITY

| Parallel File Systems | Lenovo GSS | Intel Lustre | NFS |
| --- | --- | --- | --- |

| Systems Management | xCAT | | Confluent |
| --- | --- | --- | --- |

**OS VM**   **OFED**

**NEXTSCALE**     Leovo System x
Virtual, Physical, Desktop, Server

**Compute   Storage   Network**  (intel) **OmniPath**   Mellanox  **UFM**

- **Adding new features to the stack**
  - **Web Console GUI**
  - **xCAT**
    - **Heat Map of servers/racks**
    - **Fluid Return Temperature /Flow rate of CDU**
  - **Energy Awareness**
    - **scheduler independent**

# TECHNOLOGY TRENDS

Think even bigger

New Engagement. New Journey.

# Intel processors Development Model



Tick-Tock Development Model:
Sustained Microprocessor Leadership

Previous Generation | Current generation | Next Generation

| Intel® Microarchitecture Codename Nehalem | | Intel® Microarchitecture Codename Sandy Bridge | | Intel® Microarchitecture Codename Haswell | | Intel® Microarchitecture Codename Skylake | |
|---|---|---|---|---|---|---|---|
| Nehalem 45nm New Micro-architecture | Westmere 32nm New Process Technology | Sandy Bridge 32nm New Micro-architecture | Ivy Bridge 22nm New Process Technology | Haswell 22nm New Micro-architecture | Broadwell 14nm New Process Technology | Skylake 14nm New Micro-architecture | Future Product |
| Tock | Tick | Tock | Tick | Tock | Tick | Tock | Tick |

Innovation delivers new microarchitecture with Skylake

# Intel processors Development Model

# details

# Knights Landing Architectural Diagram

Over 3 TF DP peak

Full Xeon ISA compatibility through AVX-512

~3x single-thread vs. compared to Knights Corner

Up to 72 cores

2D mesh architecture

Up to 16GB high-bandwidth on-package memory (MCDRAM)

Exposed as NUMA node

~500 GB/s sustained BW

2x 512b VPU per core (Vector Processing Units)

6 channels DDR4

Up to 384GB

MCDRAM  MCDRAM  MCDRAM  MCDRAM

DDR4

DDR4

DDR4

Up to 72 cores

DDR4

DDR4

DDR4

MCDRAM  MCDRAM  MCDRAM  MCDRAM

**Tile**

2 VPU   HUB   2 VPU

Core   1MB L2   Core

Wellsburg PCH

DMI

HFI

Common with Grantley PCH

1S (no QPI/KTI)

Connector

Based on Intel® Atom Silvermont processor with many HPC enhancements

Deep out-of-order buffers

Gather/scatter in hardware

Improved branch prediction

4 threads/core

High cache bandwidth & more

2 ports Intel® Omni-Path Integrated Fabric

On-package 50 GB/s bi-directional

Micro-Coax Cable (IFP)

Micro-Coax Cable (IFP)

PCIe Gen3 x36

Diagram is for conceptual purposes only and only illustrates a CPU and memory – it is not to scale and does not include all functional areas of the CPU, nor does it represent actual component layout.

Back to Contents

intel | 27

23

# NVIDIA NVLink architecture

### 1st Generation

### 2nd Generation

# AMD Naples and multicores – 1P or 2P in HPC?



2 SOCKET
64 CORE
128 THREAD

| Component | AMD | INTEL |
|---|---|---|
| CPU model | "Naples" | E5-2699A V4 |
| Total CPUS | 2 | 2 |
| Total cores (SMT/HT on) | 128 | 88 |
| Total memory channels | 16 | 8 |
| Total memory capacity (16 GB DIMMS) | 512 | 384 |
| Memory frequency | 2400 | 1866 |
| Total PCIE gen3 lanes to CPUs | 8x16=128 | 2x40=80 |

o Intel server is a standard, commercially available server from a major OEM

https://www.nextplatform.com/2017/05/17/amd-disrupts-two-socket-server-status-quo/

# AMD Naples and multicores – 1P or 2P in HPC?

**ARM solution from mobile to server  to offer a solution at lower power consumption**

## Maximizing Throughput Density: per mm², per Watt



**ARM Solution Benefits:**

- Less than 1/3rd the power for equivalent performance*
- Allows power headroom for specialized
- computing or greater thread density

Comparison for equivalent number of threads
- Platforms used:
  - Xeon-E5 2660  10C20T platform (measured)
  - Xeon-E5 2650  10C20T platform  (measured)
  - Gcc compiler v4.9 with –o3 flag
  - TDP rating source: ark.intel.com

- Estimated result on example 20C ARM Cortex platforms with CCN-508, 28MB total L2+L3 cache
  - per-core measurements on RTL with relevant memory system
  - Gcc compiler v4.9 with –o3 flag
  - Scaled to 20T based on modelled and empirical results
  - Power estimated in 16nm based on ARM internal implementations for entire CPU+interconnect complex including 20xCPU, CCN-508, L2+L3 caches
  - Actual results on silicon platforms may vary
* A portion of Intel TDP power will be consumed by IO.
The Cortex-A72 and Cortex-A57 estimates exclude IO power

https://www.nextplatform.com/2017/03/21/new-arm-architecture-offers-dynamiq-response-compute/

# Performance trends in a server

- **Technology evolution determines a significant performance growth in the next 3yrs**
- **From 2015 to 2018 peak performances double at least on x86, X-Phi, GPUs**
- **Technology solutions to hundreds of PFs is not so evident and will depend by several conditions:**
  - **Peak performance vs cost**
  - **Peak performance vs power consumption (GFs/W)**
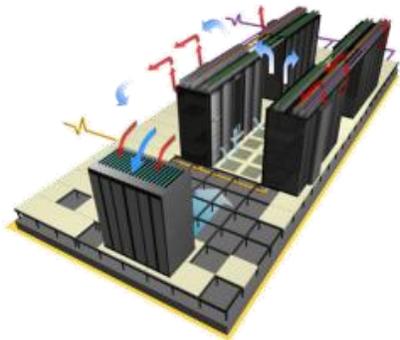  - **Sustained performances vs power consumption and TCO**

## Peak performance trends

| | 2014 | 2015 | 2016 | 2017 | 2018 |
|---|---|---|---|---|---|
| x86 2p | 0,7 | 1 | 1,5 | 3 | 4 |
| X-Phi | 0,8 | 1,2 | 3 | 3 | 7,2 |
| GPUs | 1,2 | 1,8 | 4 | 6 | 7 |

# COOLING TECHNOLOGY AND TCO

Think even bigger

New Engagement. New Journey.

# Choice of Cooling

**Air Cooled**



**Air Cooled with
Rear Door Heat Exchangers**



**Direct Water Cooled**



- Standard air flow with internal fans
- Fits in any datacenter
- Maximum flexibility
- Broadest choice of configurable options supported
- Supports Native Expansion nodes (Storage NeX, PCI NeX)

**PUE ~1.5**

**ERE ~ 1.5**

- Air cool, supplemented with RDHX door on rack
- Uses chilled water with economizer (18C water)
- Enables extremely tight rack placement

**PUE ~1.2**

**ERE ~ 1.2**

- Direct water cooling with no internal fans
- Higher performance per watt
- Free cooling (45C water)
- **Energy re-use**
- Densest footprint
- Ideal for geos with high electricity costs and new data centers
- Supports highest wattage processors

**PUE <= 1.1**

**ERE ~ 0.3**  **with hot water**

Choose for broadest choice of customizable options

Choose for balance between configuration flexibility and energy efficiency

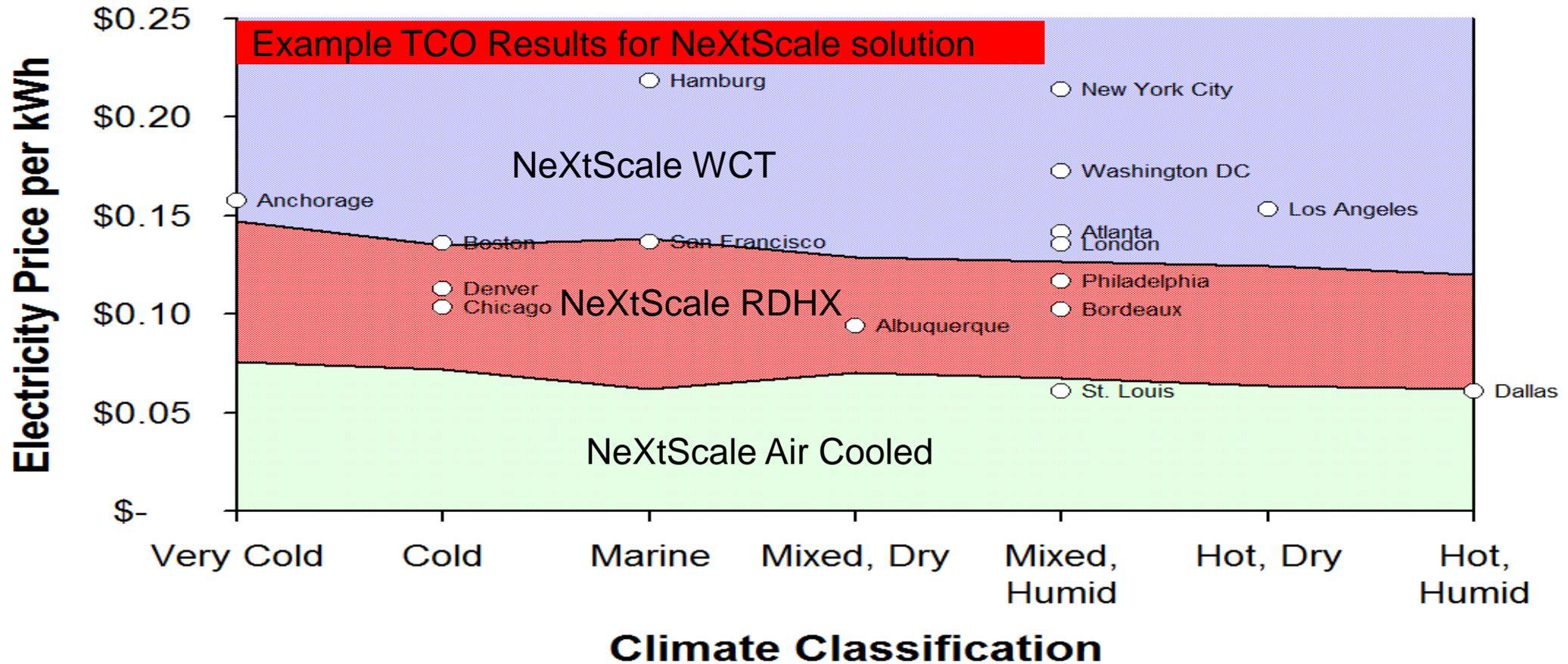Choose for highest performance and energy efficiency

**Power cooling with hybrid W+A solution: Tinlet air 25°C and water on RDHX at 20°C and 8gpm**



% heat removal as function of water temperature and flow rate for given rack power, rack inlet temperature, and rack air flow rate
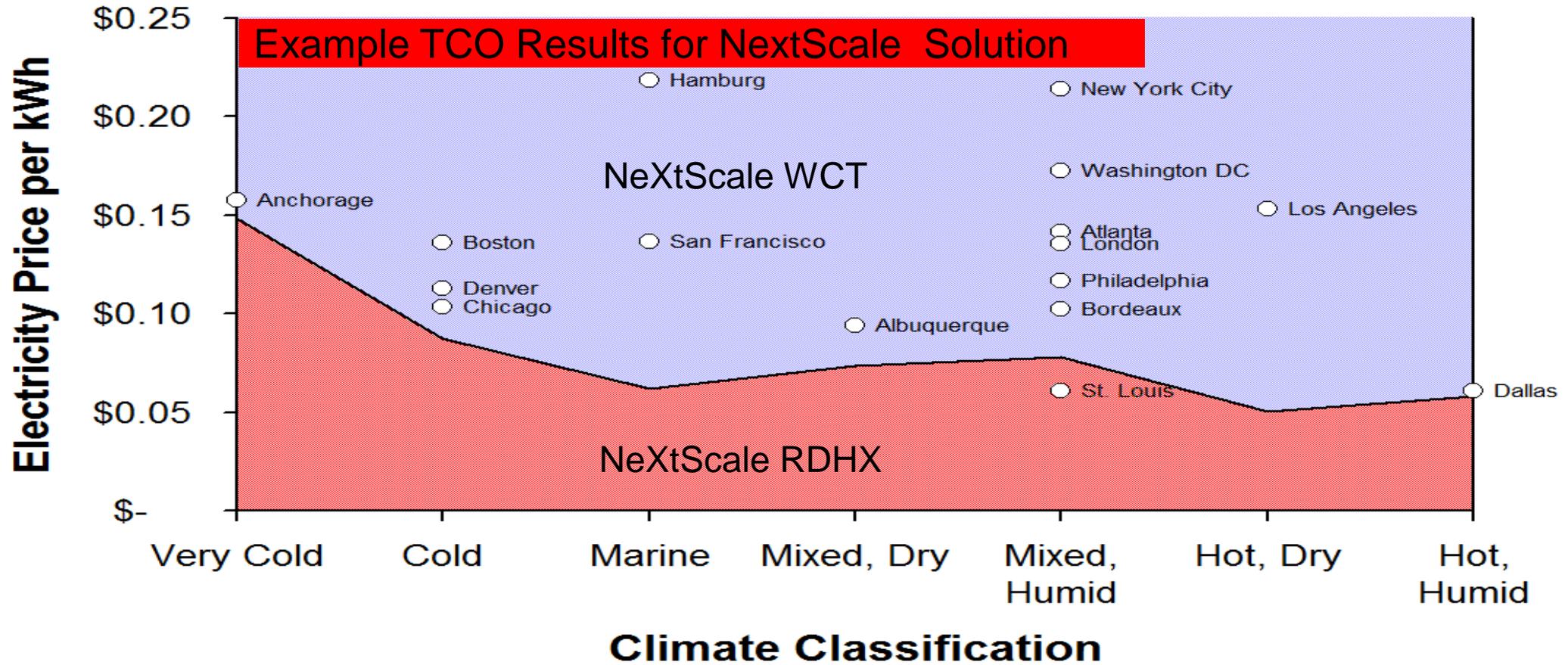
# Technology Selection for an <u>Existing</u> Data Center Installation



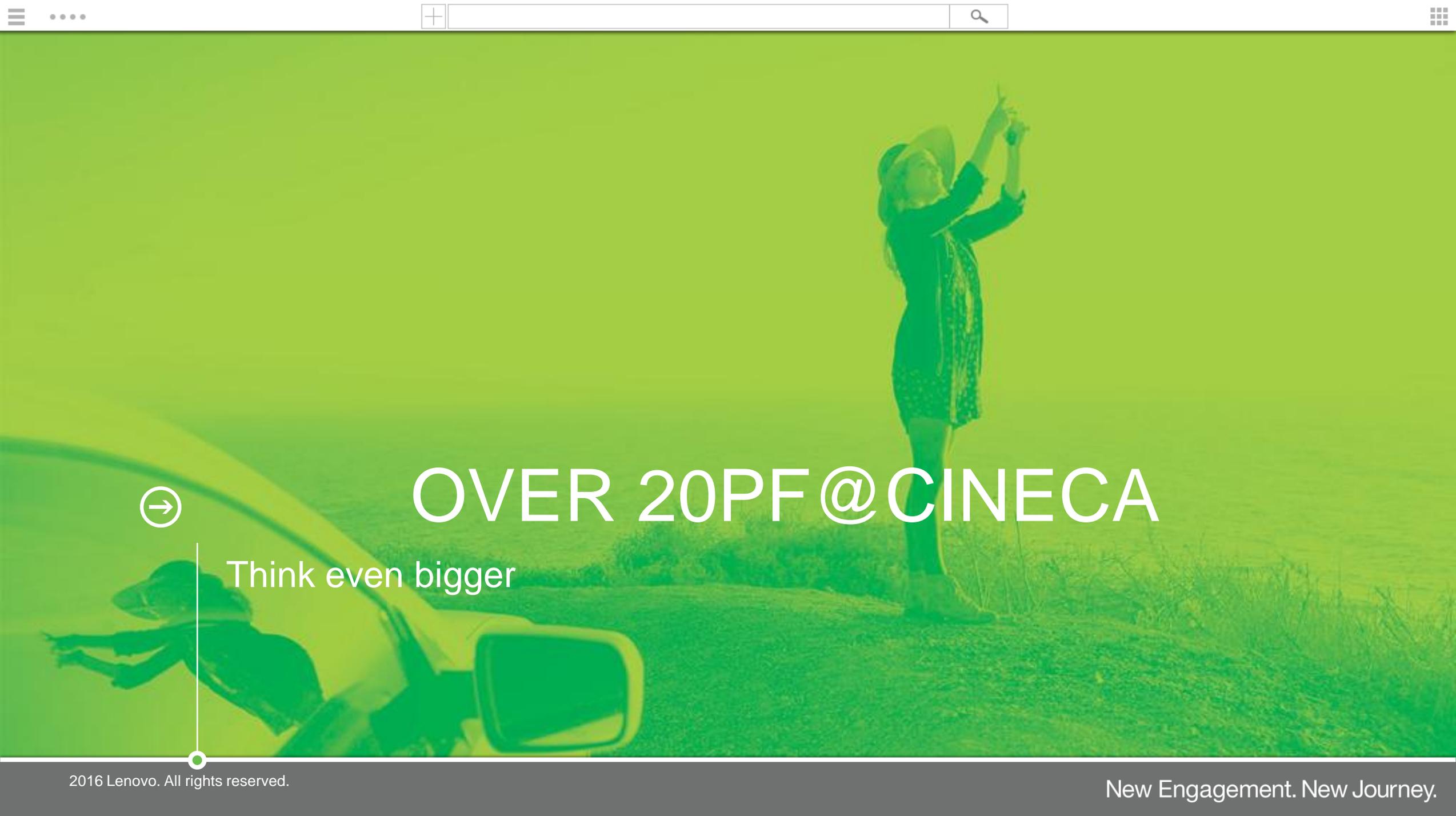Technology to Maximize 5-Year NPV for an Existing Construction

# ✚ Technology Selection for a <u>New</u> Data Center Installation



**Technology to Maximize 5-Year NPV for a New Construction**

Example TCO Results for NextScale Solution

NeXtScale WCT

NeXtScale RDHX

Electricity Price per kWh — Climate Classification

Data points: Anchorage, Boston, Denver, Chicago, Hamburg, San Francisco, Albuquerque, St. Louis, New York City, Washington DC, Atlanta, London, Philadelphia, Bordeaux, Los Angeles, Dallas

Climate Classification: Very Cold, Cold, Marine, Mixed, Dry, Mixed Humid, Hot, Dry, Hot Humid

## How to manage power

- Report

  – temperature and power consumption per node / per chassis

  – power consumption and energy per job

- Optimize

  – Reduce power of inactive nodes

  – Reduce power of active nodes

# OVER 20PF@CINECA

Think even bigger

New Engagement. New Journey.

# CINECA OBJECTIVES AND TECHNOLOGIES

## Several phases

- A1: 2 PFs peak convential architecture

- A2: >10 PFs peak non conventional architecture

- A3: >4 PF peak

- Interconnect : >40Gbs bidi between 2 nodes

- Storage :

    - S1: 10PB, >100 GB/s

- Power  < 2.0 megawatts all inclusive

## Technologies

- A1: BRDW in Lenovo NeXtScale

- A2: KNL in Intel AdamsPass and RDHX

- A3: SKL with Lenovo Stark and RDHX

- Single OPA fabric and 2:1 blocking ratio

- Storage

    - S1: 6xGSS26 with 8 TB drive

# LENOVO ECO SYSTEM FOR CINECA

**Compute Platforms**

**Software Environments**

**Interconnects**

**Storage Subsystems**

**Infrastructure**

**IBM Spectrum Scale**

## Lenovo Services

- ✓ Design
- ✓ Architecture
- ✓ Project Mgmt.
- ✓ Optimization

# BRW vs. KNL vs. SKL (based on Cineca)

| | BRW (2PFL) | KNL (11PFL) | SKL (>4,5PFL) |
|---|---|---|---|
| **Nodes** | 1512 | 3600 | 1512 |
| **CPU/node** | 2 | 1 | 2 |
| **TFlop/node** | 1.3 | 3 | 3.2 |
| **Price/node** | | | |
| **CPU** | E5-2697v4 18c@2,3GHz | 7250 68c@1.4GHz | 8160 24c@2,1GHz |
| **TFlop/Socket** | 0.65 | 3 | 1.6 |

# CINECA – OMNI-PATH FABRIC ARCHITECTURE (SINGLE FABRIC, WITH 32:15 BLOCKING)

**Total : 5 director switches in 5 racks**
**(A1 LeROM : 5 * #509)**

**SawTooth Forest 768p Director**
(max 24x 32p linecards)

**5x**

**SawTooth Forest 768p Director**
(max 24x 32p linecards)

OPA Optical Cables (50m, 30m, 20m)
3 links from each EDF to each STF

(A1 LeROM: 1 * #085)
(A2 LeROM: 1 * #142)
(A3 LeROM: 1 * #143)

**Eldorado Forest 48p Edge (~2:1)**
(15p up + 32p down)

**Eldorado Forest 48p Edge (~2:1)**
(15p up + 32p down)

EDF 48p Edge (~1:1)

EDF 48p Edge (~1:1)

EDF 48p Edge (~1:1)

EDF 48p Edge (~1:1)

EDF 48p Edge (~1:1)

8p (1p/srv)

32 NSD servers (1p/srv)

**32p** (1p per server)

**32p** (1p per server)

12p (2p/srv)

**1 rack**

**(A2 LeROM:        1 *   #508)**

**FDR**

**32 AdamsPass KNL nodes**
**(9 switches + 288 nodes in 4 racks)**

**32 NeXtScale BDW nodes**
**(9 switches + 288 nodes in 4 racks)**

**3x GSS26 @ 8TB (6 servers)**
(~6PB in 2 racks;
OPA parts shipped in #722)

**8x mgmt node**
(xCAT, IFS, misc )

**2x login node**

Total : 3600 KNL nodes in 50 racks
(A2 LeROM : 12 * #506 + 1 * #106)
Total : 1512 SKL nodes in 21 racks
(A3 LeROM : #516 placeholder)

Total : 1512 BDW nodes in 21 racks
(A1 LeROM : 5 * #512 + 1 * #084] )

Total : ~12 PByte in 4 racks
(A1 LeROM : 1 * #515 + 1 * #517)

Management Rack
(A1 LeROM: 1 * #722)
(A2 LeROM: 1 * #122)
(A3 LeROM: 1 * #146)

38

# HPC Cloud Low Latency Networking

Spectrum Scale building blocks attached to 40Gbe

CES servers exporting Spectrum Scale user FS (user-specific subtrees)

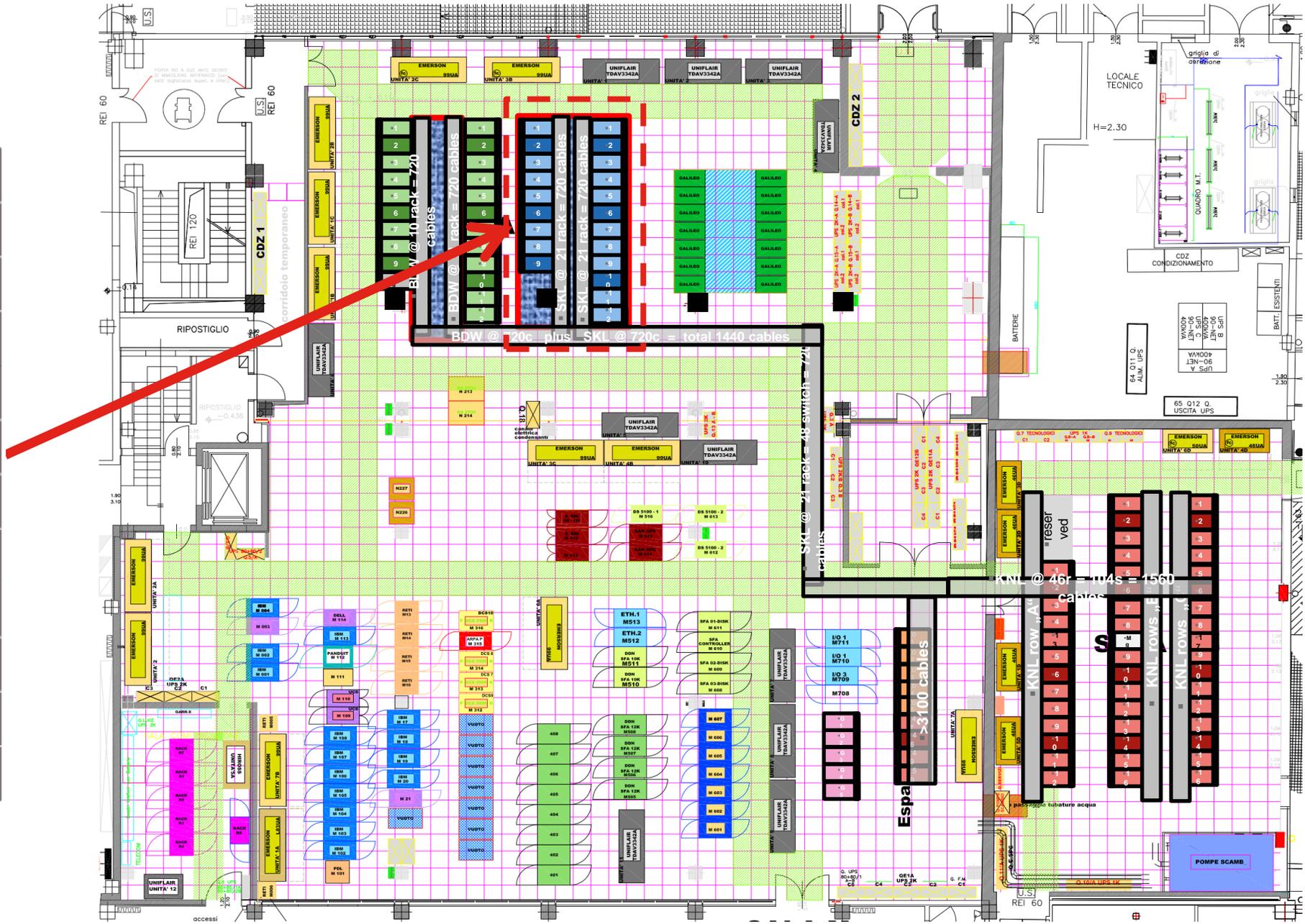CES server (NFS, SWIFT) If not iSCSI

Spectrum Scale Building Blocks

Optionally, for using existing Spectrum Scale storage as Cloud backend storage for Block (Cinder), Image (Glance), Compute (Nova), Cloud controllers as well as compute nodes have to be part of a Spectrum Scale cluster and mount Specific file systems / file Sets for this purpose

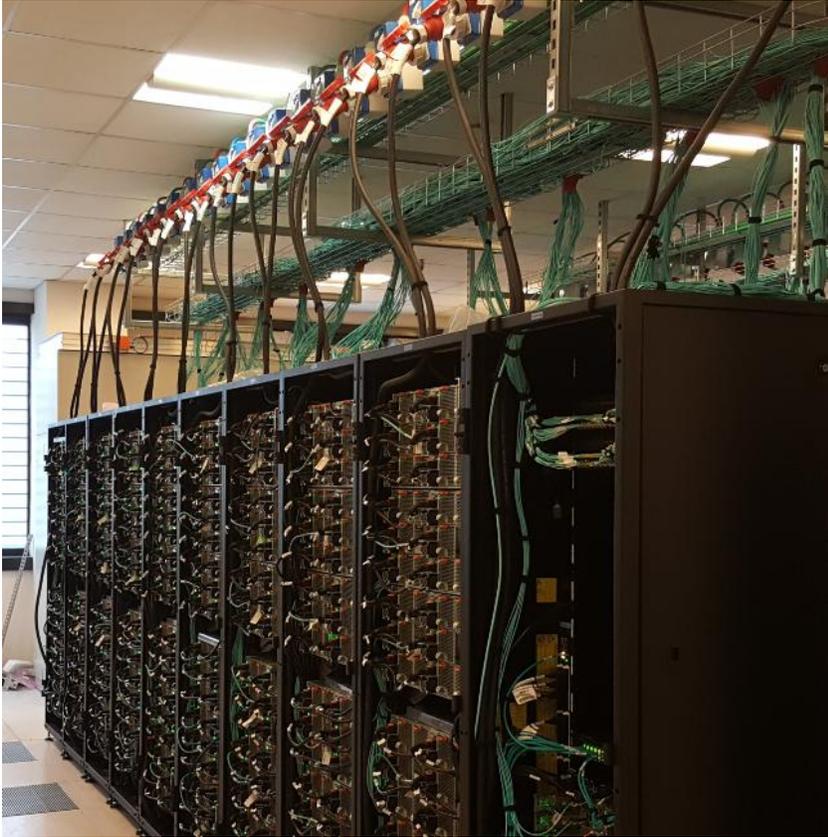HPC cloud using Spectrum Scale File Systems through CES servers or iSCSI block storage

vLAG 25 GbE

Cloud Backbone Netw. 25 GbE
3:1 oversubscription

Management Netw. 1 GbE

vLAGs 25 GbE

XX * iSCSI block storage if not CES server

XX * Cloud Deployment from Galileo

XX x Cloud Controllers from Galileo

XXX Cloud Compute Nodes

XX GPU Cloud Compute Nodes

# Cineca A3 Floor Plan

## A3

| | | |
|---|---|---|
| skl-r129 (N300) | | skl-r138 (N400) |
| skl-r130 (N301) | | skl-r139 (N401) |
| skl-r131 (N302) | | skl-r140 (N402) |
| skl-r132 (N303) | | skl-r141 (N403) |
| skl-r133 (N304) | | skl-r142 (N404) |
| skl-r134 (N305) | | skl-r142 (N405) |
| skl-r135 (N306) | | skl-r144 (N406) |
| skl-r136 (N307) | | skl-r145 (N407) |
| skl-r137 (N308) | | skl-r146 (N408) |
| | | skl-r147 (N409) |
| pillar | | skl-r148 (N410) |
| | | skl-r149 (N411) |



BDW @ 720c plus SKL @ 720c = total 1440 cables

KNL @ 46r = 104s = 1560 cables

>3100 cables

# Installation Pictures – A1 Broadwell
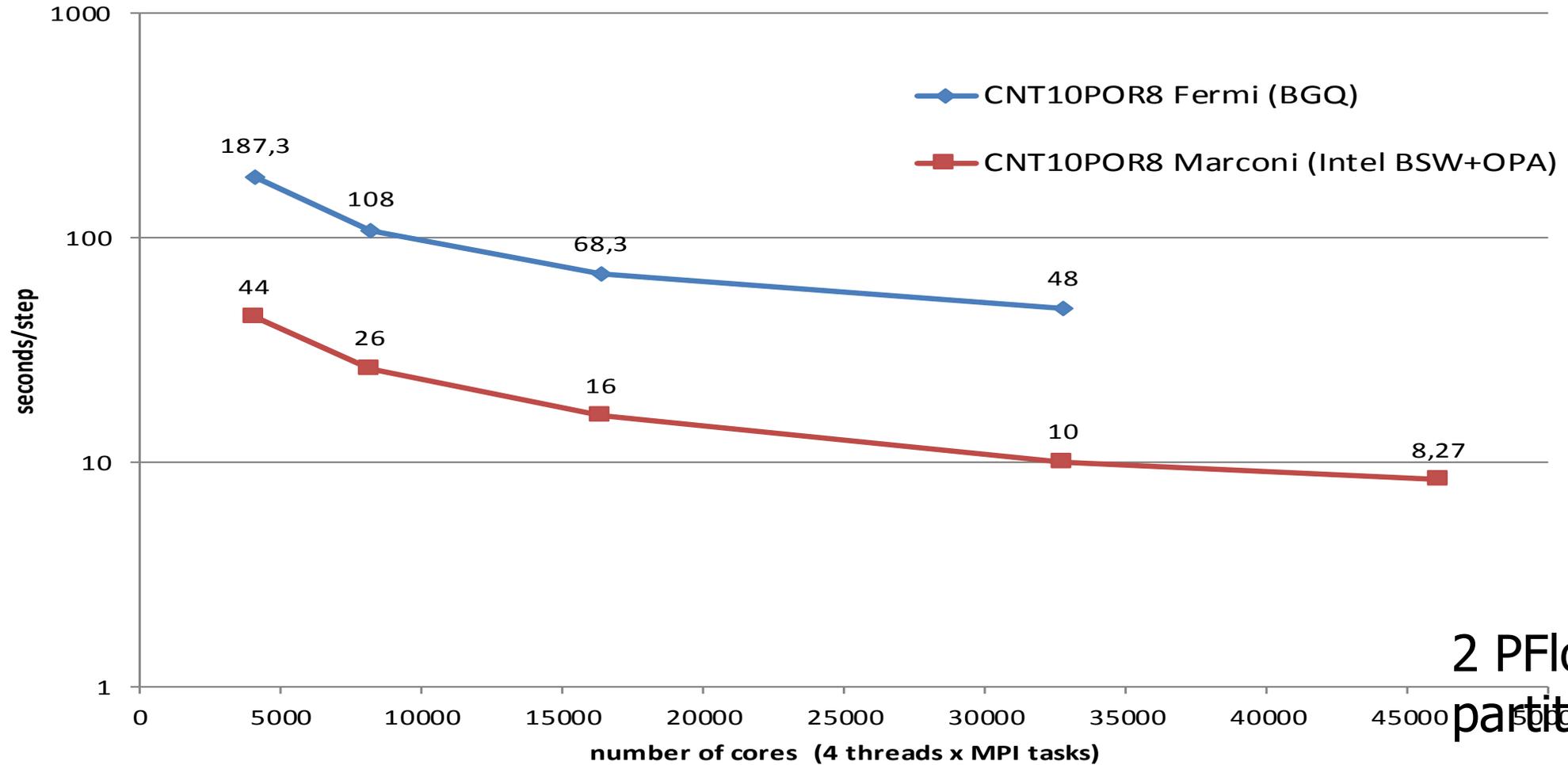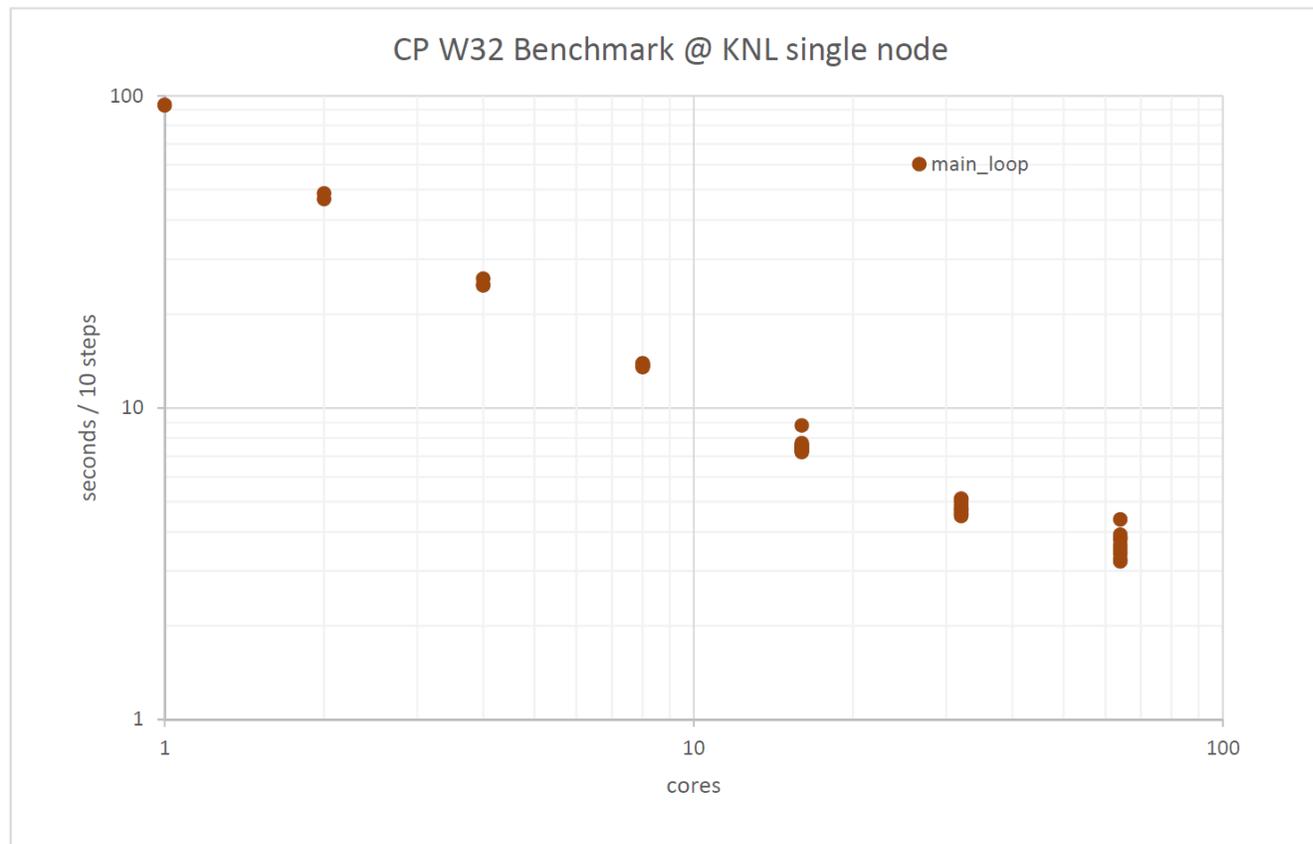


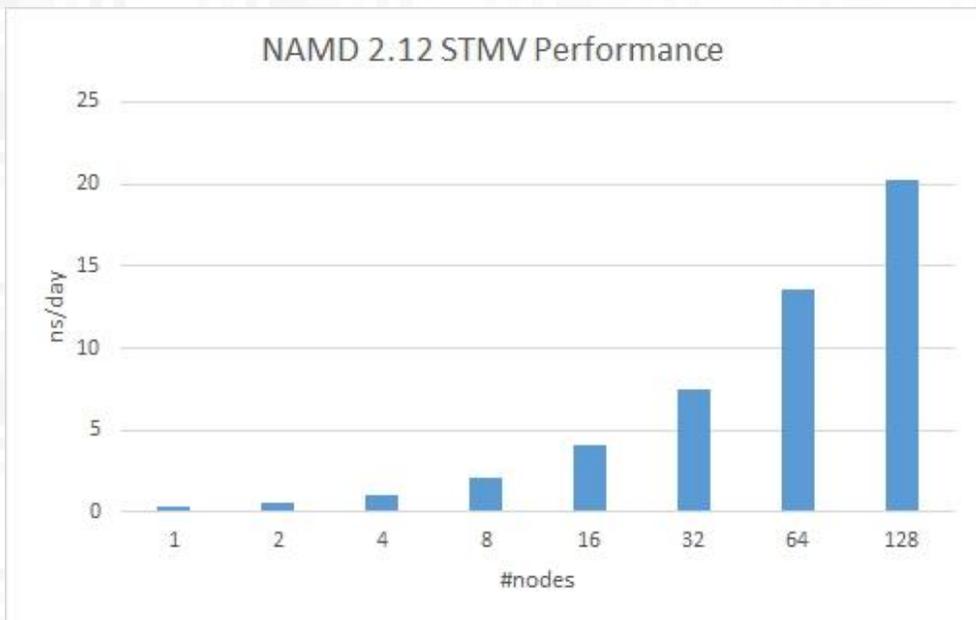*Mgmt & Compute Racks (hot aisle)*

# Installation Pictures – A2 KNL

# QE scaling benchmark (cp.x)

**Input dataset: http://www.qe-forge.org/gf/download/frsrelease/49/63/CNT10POR8.tgz**

Legend:
- CNT10POR8 Fermi (BGQ)
- CNT10POR8 Marconi (Intel BSW+OPA)

Fermi (BGQ) data points: 187,3 · 108 · 68,3 · 48

Marconi (Intel BSW+OPA) data points: 44 · 26 · 16 · 10 · 8,27

Y-axis: seconds/step (1, 10, 100, 1000)

X-axis: number of cores (4 threads x MPI tasks) (0, 5000, 10000, 15000, 20000, 25000, 30000, 35000, 40000, 45000)

2 PFlops partition

# NAMD on A1 Broadwell and CP on A2 KNL single node



NAMD 2.12 STMV Performance



CP W32 Benchmark @ KNL single node

Courtesy by Carlo Cavazzoni - CINECA

# QE-CP: A1 Broadwell vs A2 KNL



Courtesy by Carlo Cavazzoni - CINECA

## MARCONI-A1 (physical view)
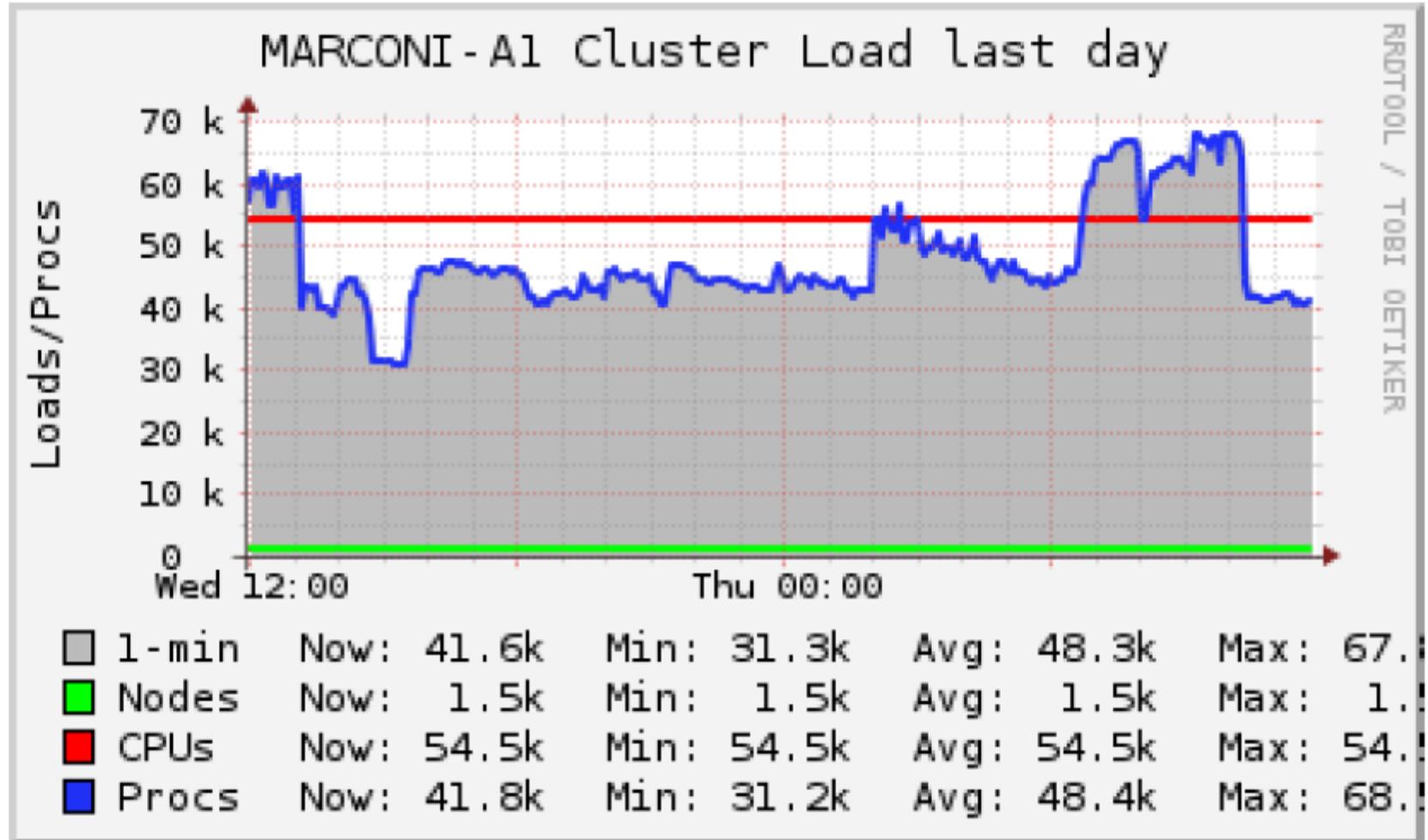
CPUs Total: **54524**

Hosts up: **1520**

Hosts down: **0**

Current Load Avg (15, 5, 1m):
  **76%, 76%, 76%**

Avg Utilization (last day):
  **89%**

Localtime:
  2017-01-19 10:50

MARCONI-A1 Cluster Load last day

| | | Now: | Min: | Avg: | Max: |
|---|---|---|---|---|---|
| ☐ | 1-min | 41.6k | 31.3k | 48.3k | 67. |
| 🟩 | Nodes | 1.5k | 1.5k | 1.5k | 1. |
| 🟥 | CPUs | 54.5k | 54.5k | 54.5k | 54. |
| 🟦 | Procs | 41.8k | 31.2k | 48.4k | 68. |