

---

# Hadoop

Hardware sizing – lessons learnt

'The Quest for the ~~Holy Grail~~ Ideal Hadoop Server



Marcel van Drunen  
Senior Enterprise Technologist HPC & Cloud  
Dell EMEA

---



**3B** People will connect electronically via mobile or Internet technology by 2014



**20B** pieces of content shared on Facebook every month

**\$440B** 2011 IT spending in emerging markets - an increase of 10.4% over 2010



**2x** IDC estimates that the Digital Universe will double every 18 months

**35 ZB** By 2020, the Digital Universe will be 44 times as big as it was in 2009

# Outline

- Intro
- Hardware scaling for Hadoop clusters
- Big Data influences Big Dell

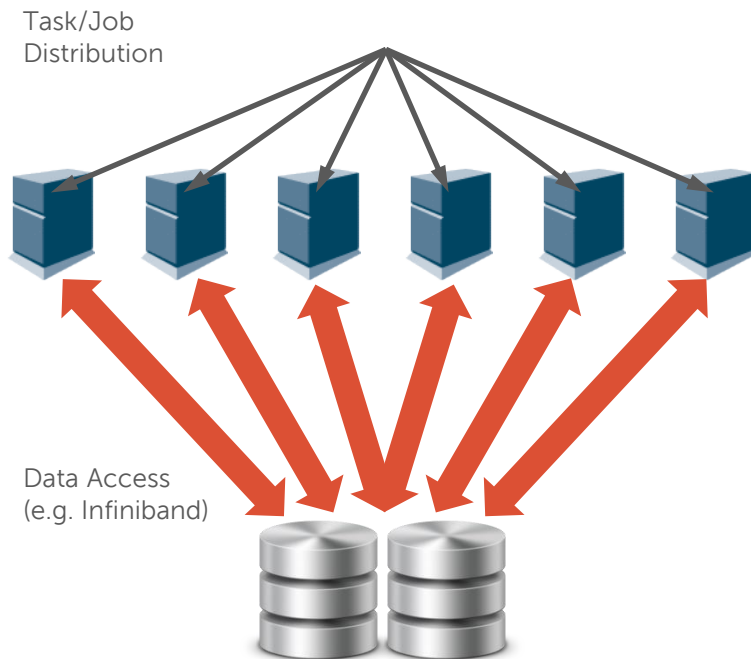


# Big Data Infrastructure vs. HPC Clusters

## Two Forms of Distributed Computing

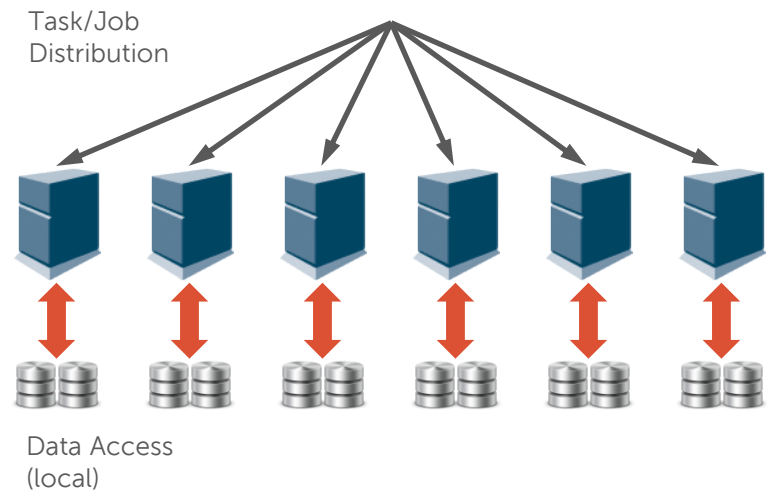
### High Performance Compute Clusters

- Parallel File System
  - High throughput
  - All nodes can access all data
  - Compute-centric workloads



### Hadoop Cluster Architecture

- Distributed File System
  - Global namespace (ingest!)
  - Nodes just work on local data
  - Data/IO-centric workloads



# Hadoop & HPC scaling

- Combining Hadoop and regular HPC
  - not a good idea
- Rules of thumb: Data nodes
  - minimum 1 core per disk
  - most workloads HyperThreading counts as second core
  - min. 4 GB per disk, 8 GB is better
- More smaller nodes (12 disks per node maximum)
  - Rebuild times
  - Performance loss during rebuild
  - Network load during rebuild

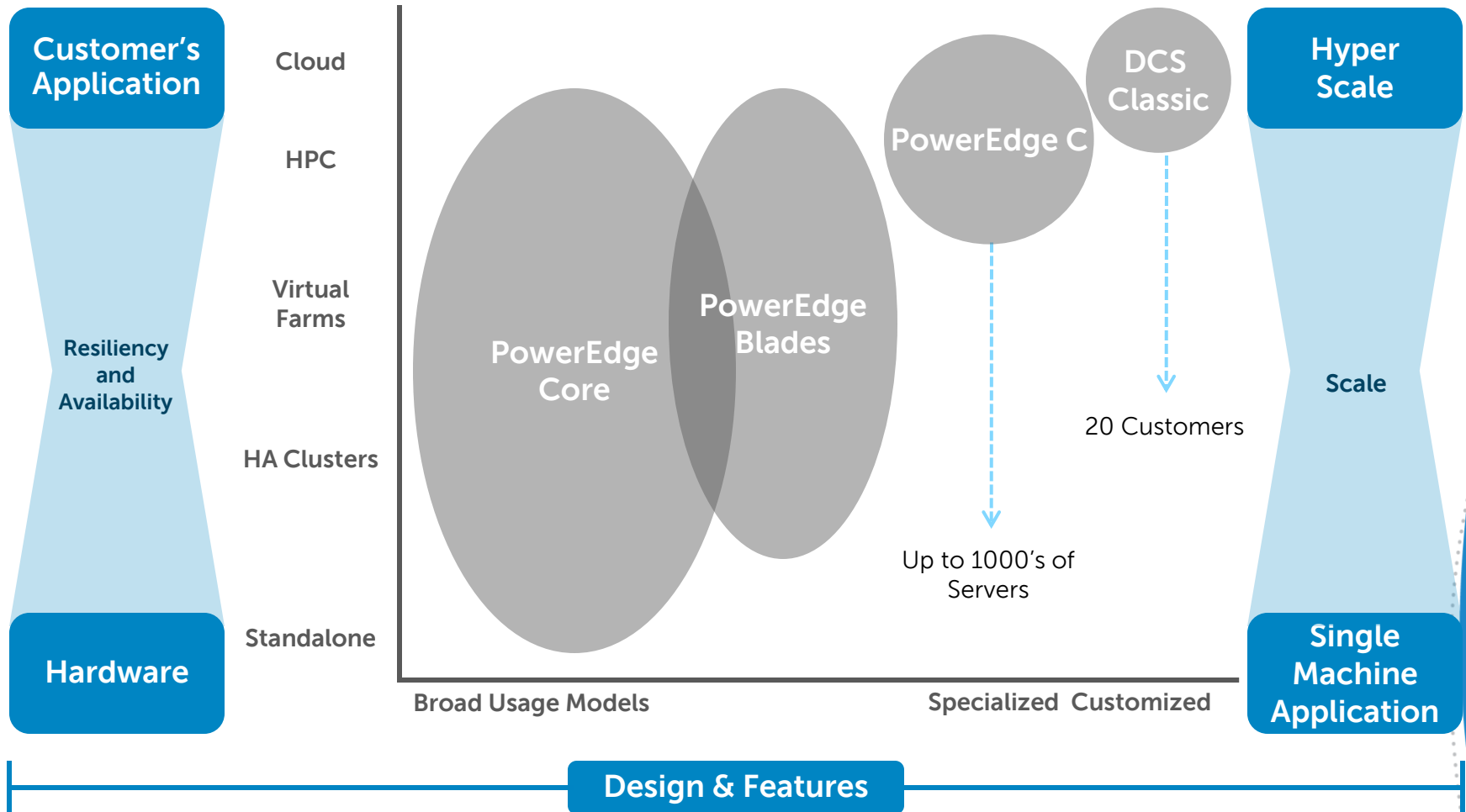


# Hadoop & HPC scaling

- Rules of Thumb: Name Node:
  - 64GB RAM
    - › 1GB files → ~13PB of data
    - › 128MB files → ~1.3PB of data
    - › 1MB files → ~82TB of data
  - Plus OS & JobTracker daemon requirements
  - Plenty of network bandwidth
  - Less nodes, so often overspecced



# Big Data influence on Dell portfolio



Please note the diagram is conceptual and not to scale



# Dell | Cloudera Solution for Apache Hadoop 2.x

HW + SW + Services		
Hardware	HW Reference Architecture	<ul style="list-style-type: none"> <li>• PE-R &amp; PE-C Servers*</li> <li>• Storage and compute</li> <li>• PowerConnect switches</li> </ul>
	Configuration	<ul style="list-style-type: none"> <li>• Min of 5 nodes</li> <li>• Deployment guide</li> </ul>
Software	Software	<ul style="list-style-type: none"> <li>• Hadoop Installer (Crowbar)</li> <li>• CDH Hadoop Enterprise</li> <li>• CDH mgmt. applications</li> <li>• Other SW elements installed by Crowbar</li> </ul>
	Operating System	<ul style="list-style-type: none"> <li>• RHEL (if Dell deploys)</li> <li>• Cent OS (if customer deploys)</li> </ul>
Services	Deployment	<ul style="list-style-type: none"> <li>• Onsite HW Install</li> <li>• Onsite SW Install</li> <li>• Whiteboard session &amp; training (via Cloudera)</li> </ul>
	Support	<ul style="list-style-type: none"> <li>• HW: Dell ProSupport</li> <li>• SW: Hadoop support (via Cloudera)</li> </ul>



## Example: PowerEdge R720xd

- Designed with big data in mind
- Compact 2U form factor
- Capacity, performance, flexibility
- Expansive disk storage
- Quantity:
  - Master nodes: 2 for redundancy
  - Edge node: 1 minimum
  - Data nodes: as many as needed
  - Admin node: 1 minimum

\* supported today: C6100, C6105, C2100, C8000, R720, R720XD





# Some servers in Dell reference design

- C2100: the server many (DCS) customers ran their first few PB on five years ago

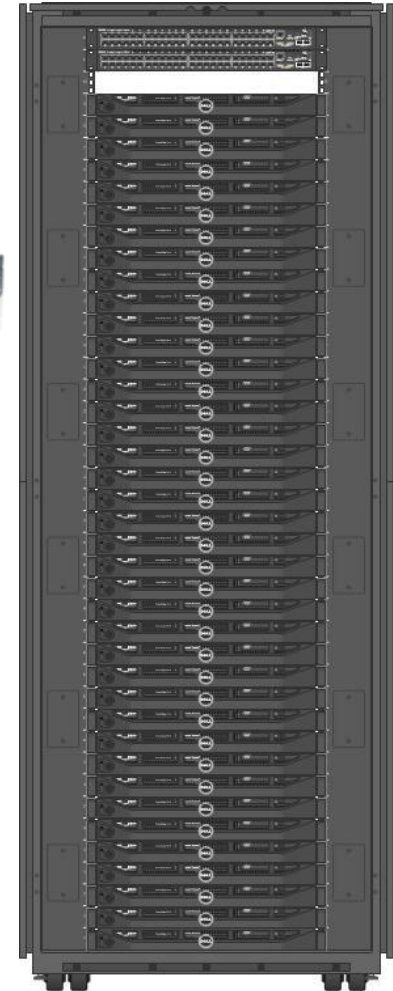


- R720xd: the server they are buying now



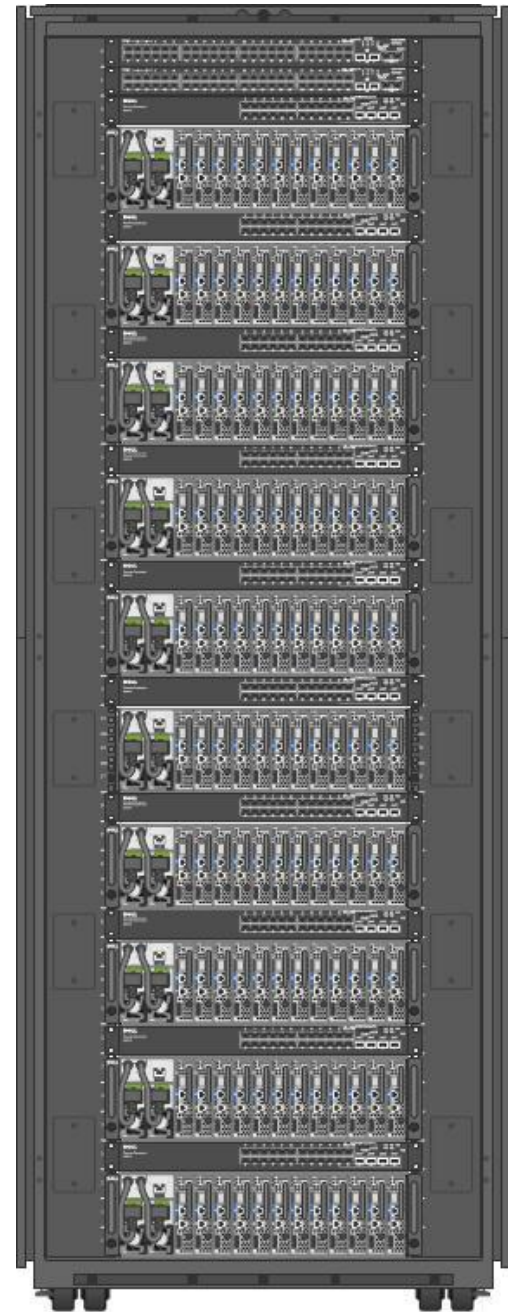
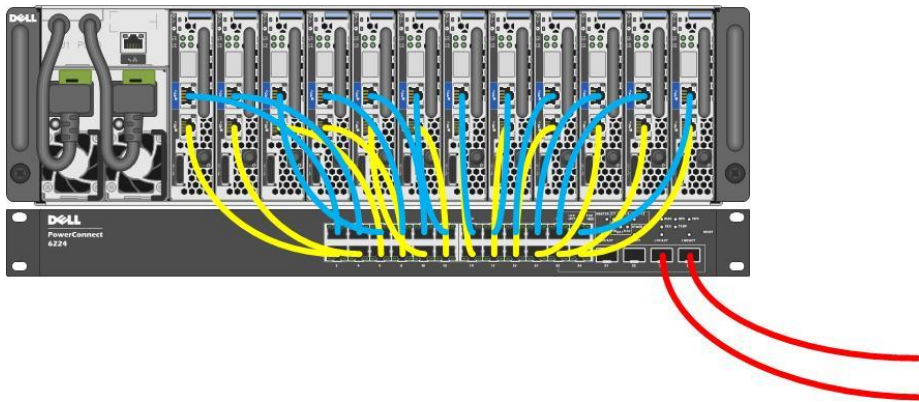
# Scenario at an unnamed customer

- Wanted single socket nodes
- Four drives per node
- 16 Gb per node
- -> Dell PowerEdge R320
  - › - 39 servers per rack (3 x 13)
  - › - 156 TB per rack raw, 52 TB usable
  - › So what happened when
  - › they wanted a PetaByte?

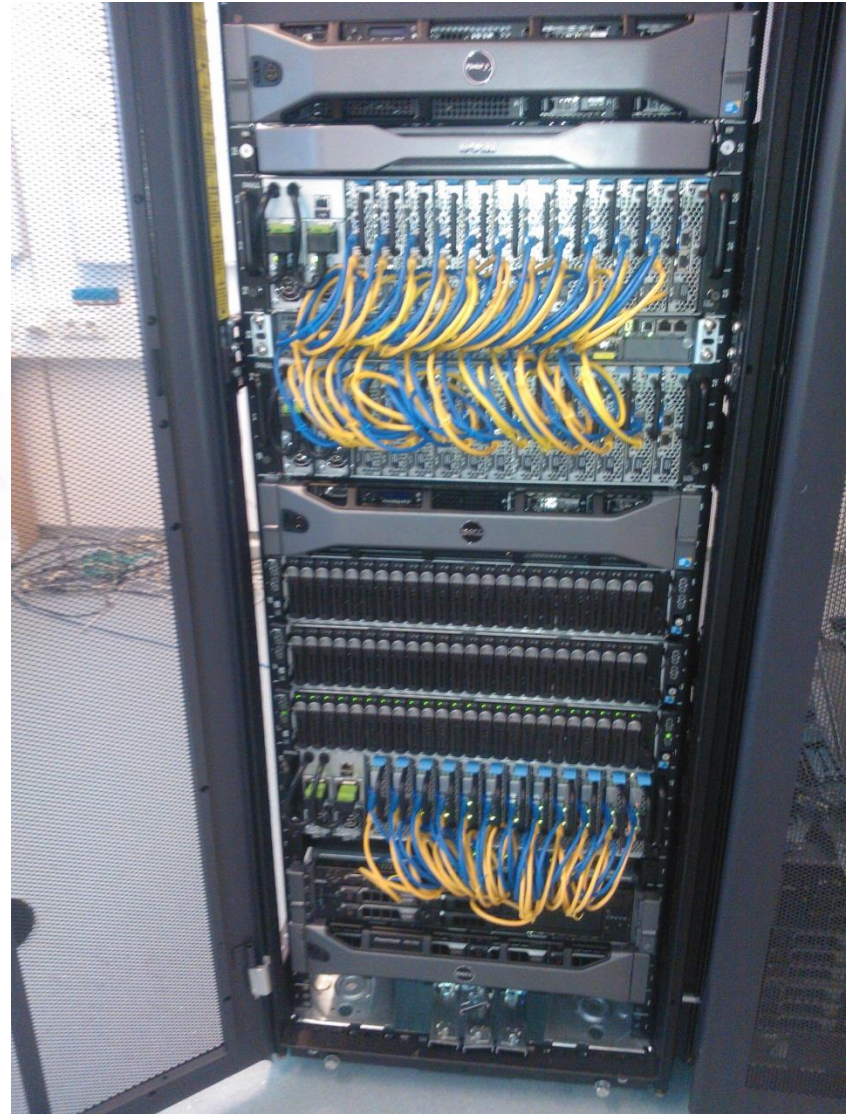


# Where density makes sense

- Dell PowerEdge C5220
  - 12 1S servers per 3U
  - Below design with switch per 4U
- 120 nodes per rack, for 480 drives
- That's 3 x 160 TB
- One PB requires 6 racks, plus one for other nodes
- And humongous network bandwidth



I eat my own dog food!



# PowerEdge C Microserver Advantage

- **Rare occurrence: a really true marketing slide**

Cooling



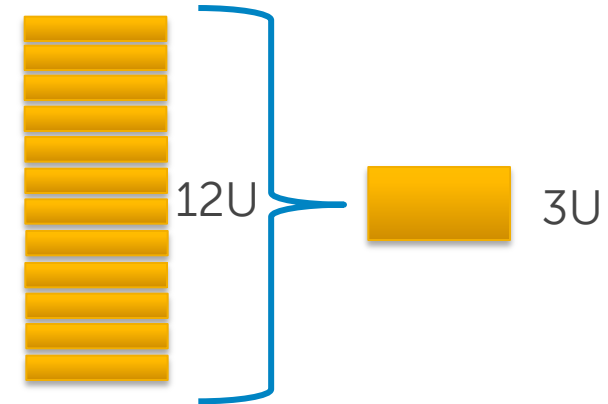
1/6<sup>th</sup> less amount of fans

Power

C5220 sled	Comparable 1U
41W	71W

40% less power per node

Mechanicals



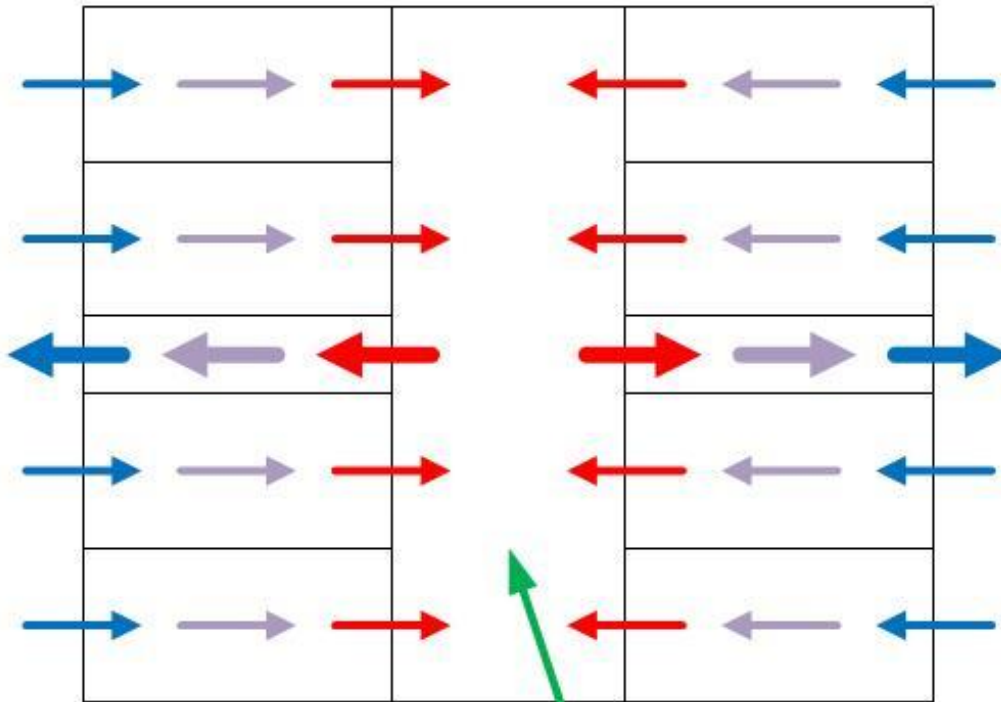
4x the density

# The Advantage of Front side serviceability

Up to 40° C

Delta can be over 20° C

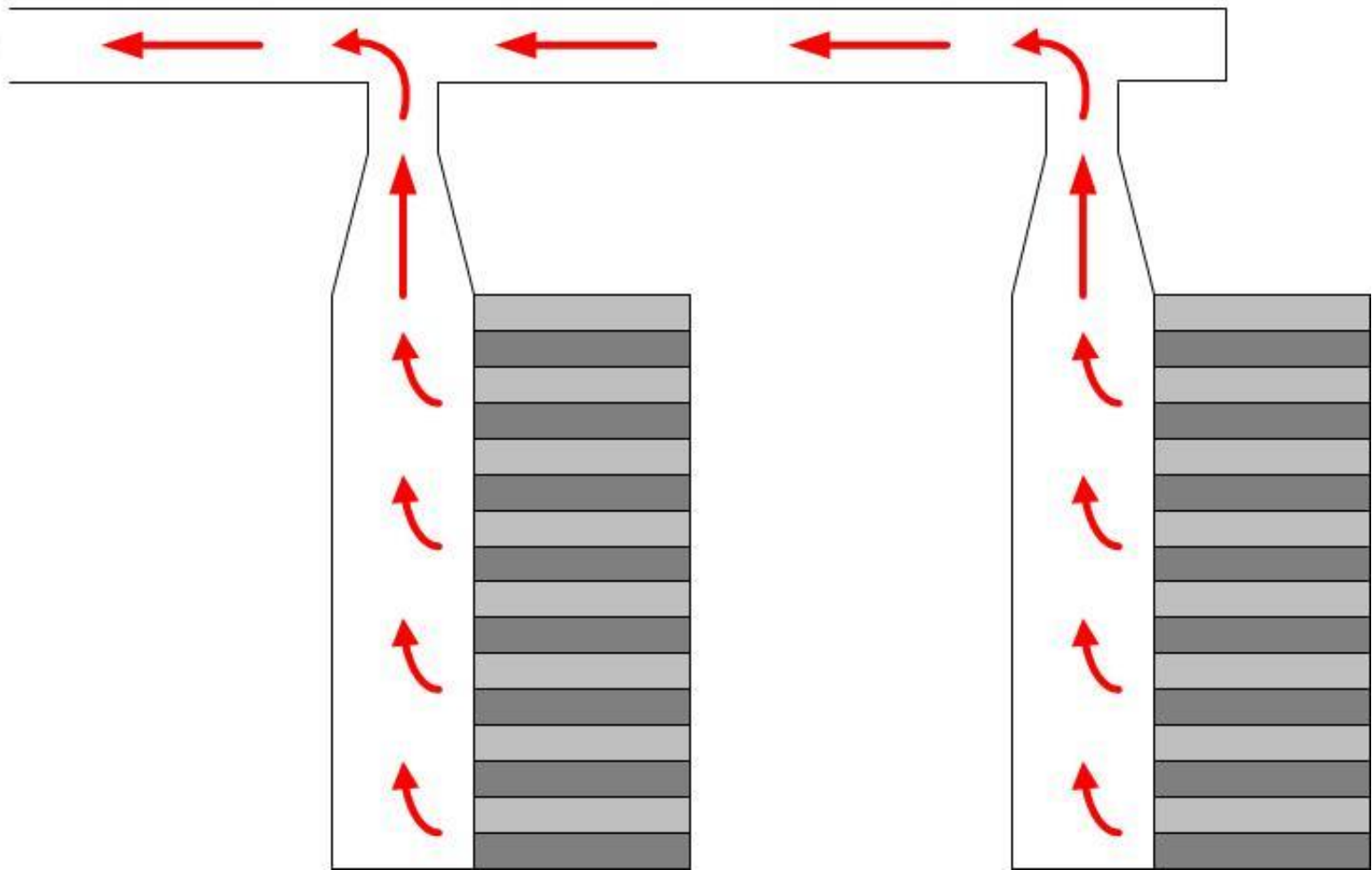
Over 60° C



Cables nor people like to be here



# Cooling with style (at an Italian customer)



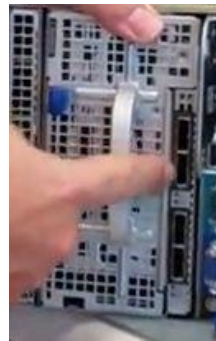
# PowerEdge C8000 Details



**C8000 Chassis  
(Internal Power Supplies)**



**C8000XD Storage Sled**





---

Thank You!



# PE-C8220 Single-Width Compute Sled features



**Architecture** 2S Intel Xeon E5-2600 Series CPUs  
Intel C600 Chipset

**Memory** 16 DIMM Slots DDR3 ECC RDIMM  
Max 256GB  
1600MHz

**PCI Expansion** 1x PCI-E x16 Slot (Single-Width)

**Mezzanine Slot** PCI-E x8 Mezz

**Drive Controller** Intel C600 (Patsburg)

**Drive Bays** 2 x 2.5" Internal

**HDDs** SAS/SSD/SATA/NLSAS

**NIC** Dual Intel 825xx 10 GbE NIC

**Management** IPMI 2.0  
DCMI 1.0  
AST2300 (iKVM)  
Dedicated Management



# PE-C8220X Double-Width Compute Sled features

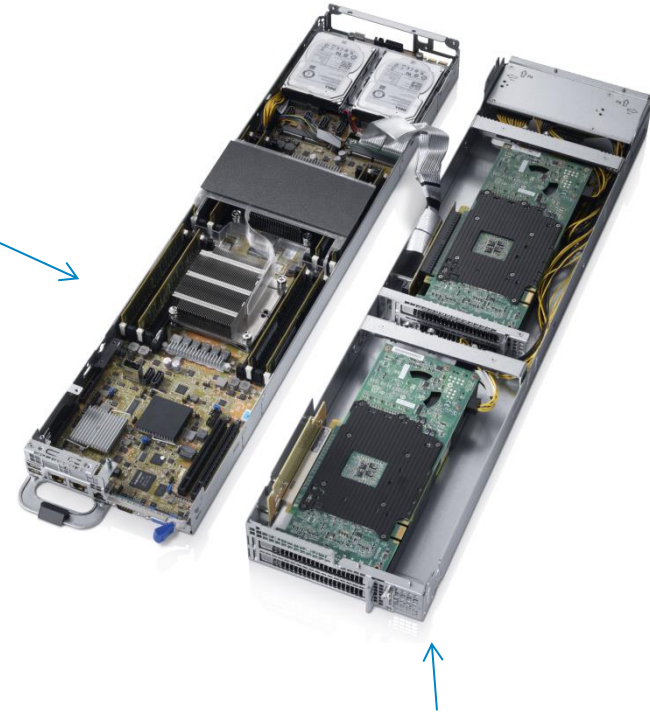
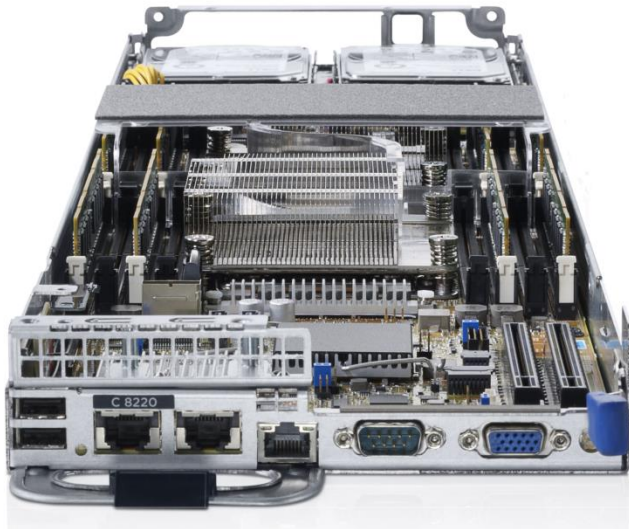
<b>Architecture</b>	2S Intel Xeon E5-2600 Series CPUs Intel C600 Chipset
<b>Memory</b>	16 DIMM Slots DDR3 ECC RDIMM Max 256GB 1600MHz
<b>PCI Expansion</b>	2x PCI-E x16 Slots (Stacked or In-line for GPU Support)
<b>Mezzanine Slot</b>	PCI-E x8 Mezz
<b>Drive Controller</b>	Intel C600 (Patsburg)
<b>Drive Bays</b>	Up to 8 x 2.5" or 4x 3.5" Internal 2x 2.5" External (Full-Width sled only)
<b>HDD's</b>	SAS/SSD/SATA/NLSAS
<b>NIC</b>	Dual Intel 825xx 10 GbE NIC
<b>Management</b>	IPMI 2.0 DCMI 1.0 AST2300 (iKVM) Dedicated Management
<b>GPGPU*</b>	Up to 2 GPGPU controllers



# PE-C8220 and C8220X Comparison

The PE-C8220 and C8220X share the same base chassis. The C8220X expansion holds additional PCI cards, GPU and Hard Drives

C8220 Base



C8220X Expansion sled



# PE-C8000 Sled Overview

C8000 Sleds	Specification
SWC Sled C8220	<ul style="list-style-type: none"> <li>•One Radon MB</li> <li>•Up to 2x 2.5" SATA HDDs</li> <li>•1x PCIE LPX16 add-in controller</li> <li>•1x MEZZ controller</li> </ul>
DWC Sled (Non-Hot plug) C8220X	<ul style="list-style-type: none"> <li>•One Radon MB</li> <li>•Up to 2x2.5" SATA HDDs</li> <li>•Up to 8x2.5" SATA/SAS HDDS(non-HotPlug), or Up to 4x3.5" SATA/SAS HDDs(non-hotplug)</li> <li>•1x LP PCIEX16 + 1xFH PCIEx16 add-in controllers</li> <li>•1x MEZZ controller</li> </ul>
DWC Sled (Front Hot plug) C8220X	<ul style="list-style-type: none"> <li>•One Radon MB</li> <li>•Up to 2x2.5" SAS/SATA Front Hot plug HDDs</li> <li>•Up to 8x2.5" SATA/SAS HDDS(non-HotPlug), or Up to 4x3.5" SATA/SAS HDDs(non-hotplug)</li> <li>•2x LP PCIEX8 add-in controllers</li> <li>•1x MEZZ controller</li> </ul>
DWC Sled (GP GPU) C8220x + GPU	<ul style="list-style-type: none"> <li>•One Radon MB</li> <li>•Up to 2x GPGPU controllers</li> <li>•1x MEZZ controller</li> </ul>
DWS Sled C8000XD	<ul style="list-style-type: none"> <li>•LSI SAS Expander with 2 x LSISAS2x28 chips</li> <li>•Up to 12 x 3.5" SATA/SAS HDDs hotplug</li> <li>•Up to 12 x 2.5" 15mm SATA/SAS HDDs hotplug</li> <li>•Up to 24 x 2.5" 9.5mm SATA/SAS HDDs non-hotplug</li> </ul>
Power Sled	<ul style="list-style-type: none"> <li>•2 x 1400w Delta Platinum PSU</li> <li>•Up to 2 PSU sleds in system</li> </ul>

PE-C8220



PE-C8220X



PE-C8000XD

