

# Standards-based, Energy-efficient HPC Systems: Trends, Implementations and Solutions.

Dr. Frank Baetke  
HP ISS/SCI Global Technology Programs

EnA-HPC Conference, Hamburg, September 16 – 17, 2010



# Purpose-built HPC Servers



# The Most Successful Architecture Ever to Enter the TOP500



# The Most Successful Architecture Ever to Enter the TOP500 – the **BL-Series (c-Class)**



# New Performance/Density for HPC: HP ProLiant BL2x220c G6



## BL2x220c G6

<b>Processor</b>	Two 80W or 60W dual- or quad-core Intel Xeon 5500 Series processors per server node*
<b>Memory</b>	Registered or Unbuffered DDR3 6 DIMM Sockets per server 96GB max per server
<b>Internal Storage</b>	1 Non-Hot Plug SFF SATA HDD per server
<b>Networking</b>	2 integrated 1GbE Ethernet ports per server
<b>Mezzanine Slots</b>	1 PCIe Gen2 x8 mezzanine expansion slot per server
<b>Additional Features</b>	Internal USB 2.0 connector Optional internal SD Card slot (consumes the USB slot)
<b>Management</b>	ProLiant Onboard Administrator powered by iLO2
<b>Density</b>	32 server nodes in 10U enclosure



# HP BladeServer c-Class 2p servers (subset)

	BL280c G6	BL460c G6	BL490c G6	BL465c G7 (June)
Processor	Up to 2P, up to 6c Intel Xeon 5500/5600 series			Up to 2, up to 12-Core AMD Opteron 6100 Series
Max Memory	12 DDR3 slots Max memory: 192GB		18 DDR3 slots Max memory: 192GB	16 DDR3 Sockets Max memory: 256GB
Storage	2 non-hot plug SFF SATA/SAS/SSD drives	Up to 2 Hot Plug SFF SAS/SATA	Up to 2 non-hot plug SSD	Up to 2 Hot Plug SFF SATA/SAS/SSD
Networking	2 integrated Multifunction GbE ports	2 integrated Multifunction 10GbE ports with Flex-10 support		
Form factor	16 per 10u enclosure 8 per 6u enclosure			
Usage	General Compute	General Compute, with hot-plug drives & 10GbE	Large Memory, with 10GbE	General Compute, with hot-plug drives & 10GbE



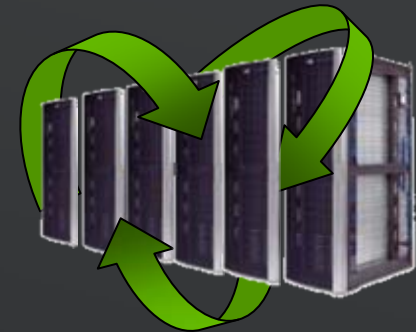
# PURPOSE DRIVEN SCALE-OUT PRODUCT LINES



Density optimized for the data center



Shared infrastructure for accelerated service delivery



Extreme scale out datacenters with lean management

	DL	BL	SL
Design center	Rack	Blade enclosure in rack	Rack
Design focus	Versatility & value	Integrated & optimized, maximum redundancy	Cost & features optimized for extreme scale out
Application	General purpose	General purpose / private cloud / scale out	Web 2.0 / cloud / scale out
Management	Essential and advanced management HP Insight Dynamics	Advanced management- accelerated service delivery & change in minutes	Home grown management Basic management via IPMI/DCMI

**Forget  
everything  
but remember  
SL**





# Purpose-built HPC Nodes

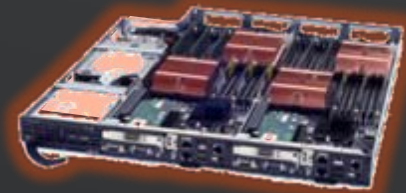
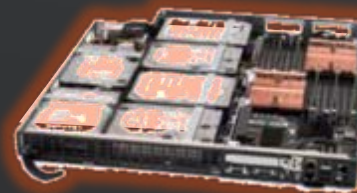
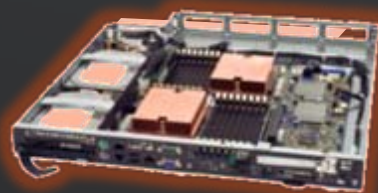
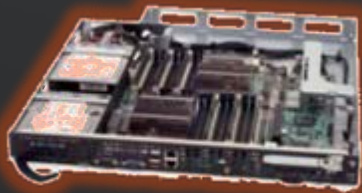


# SL-Series: HP PROLIANT SL6000

Ideal environments

#1 perf/watt

SPEC<sub>power</sub>\_ssj2008  
3106\*



**HP ProLiant  
SL160z G6**

Maximum expansion  
with 18 DIMM slots  
and up to 2 PCIe slots

**HP ProLiant  
SL165z G7**

Maximum expansion  
with 12-core AMD  
processors and 24  
DIMM slots

**HP ProLiant  
SL170z G6**

Maximum storage  
capacity with up to 6 LFF  
SATA or SAS hard  
drives

**HP ProLiant  
SL2x170z G6**

Maximum compute  
density with two servers  
per tray (1U)

**Ideal Application**

HPC database tier  
Web memory-cache

**Ideal Application**

HPC database tier  
Web memory-cache

**Ideal Application**

Web Search  
Web database

**Ideal Application**

HPC compute intensive  
Web front end

\* Based on April 2010 published benchmarks. 12/11/07 SPEC announces the release of SPEC<sub>power</sub>\_ssj2008, the first industry-standard SPEC benchmark that evaluates the power and performance characteristics of volume server class computers. The competitive benchmark results stated herein reflect results published on [www.spec.org](http://www.spec.org). See [http://www.spec.org/power\\_ssj2008/results/power\\_ssj2008.html](http://www.spec.org/power_ssj2008/results/power_ssj2008.html)

SPEC®, the SPEC logo and the benchmark name SPEC<sub>power</sub>\_ssj®2008 are registered trademarks of the Standard Performance Evaluation Corporation. The SPEC logo is © 2007 Standard Performance Evaluation Corporation (SPEC), reprinted with permission.



# HP ProLiant SL Scalable System

Next generation breakthrough server family optimized for scale-out

- Affordable scale
  - Lower acquisition cost than traditional rack servers
  - Right-sized dense server solutions
  - Based on Industry Standards
- Leading performance and efficiency
  - Concentrated compute power
  - Shared high efficiency power and cooling components
  - Lower your operating costs
- Flexible and serviceable solutions
  - Modular design allows tailoring
  - Serviceability and storage capable designs
  - Works in existing data center infrastructure



# Highly Flexible SL Chassis



## *Benefits: Low cost, high efficiency chassis*

- 4U Chassis for deployment flexibility
- Standard 19" racks, with front I/O cabling
- Unrestricted airflow (no mid-plane or I/O connectors)
- Reduced weight
- Individually Serviceable Nodes
- Variety of optimized Node Modules
- Ability to mix and match nodes

## *Multi-node, Shared Power & Cooling Architecture*

- Shared Power & Fans
- Optional Hot-Plug Redundant PSU
- Energy efficient Hot Plug fans
- 3 Phase Load Balancing
- 94% Platinum Common Slot Power Supplies
- N +1 Capable Power Supplies (up to 4)



## SL Advanced Power Manager Support

- Power Monitoring
- Node Level Power Off/On

# Purpose-built HPC Storage



# Complementary Scalable Storage Solutions for High Performance Computing

## X9000 Network Storage System

- Scalable performance and capacity
  - Scalable aggregate bandwidth
  - Scalable metadata, ideal for small files
- Shared datacenter multipurpose storage
  - Linux and Windows clients
  - NFS & CIFS support
- Ideal for applications in media, FSI, bioinformatics, web/cloud

## DDN Storage with Lustre

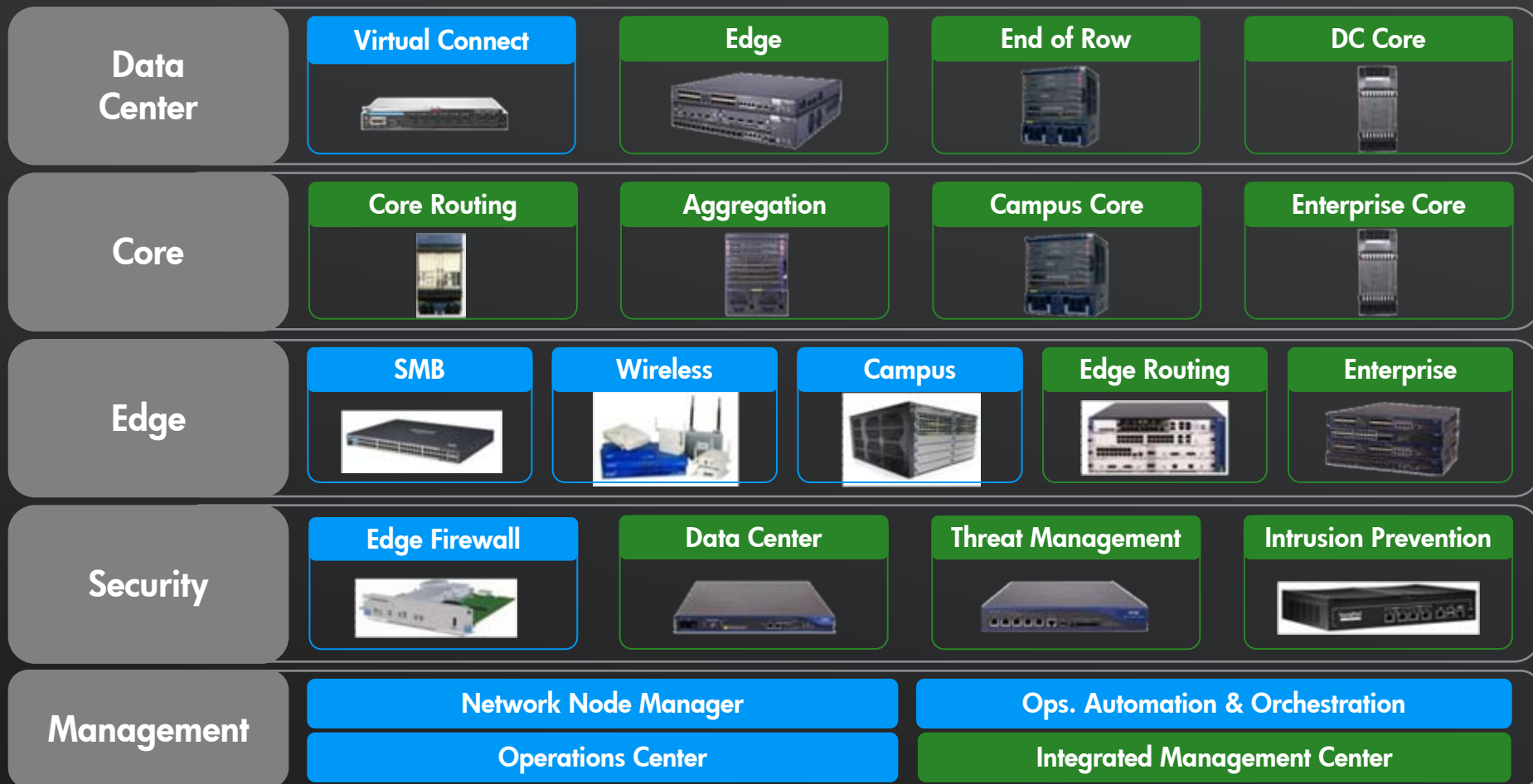
- Scalable performance and capacity
  - Scalable single-file bandwidth, with multiple writers
  - demanding bandwidth requirements
- Tightly coupled to HPC Linux clusters
- Ideal for parallel applications in traditional HPC



# Purpose-built HPC Fabrics



# HP + 3Com – Leadership from Edge to Data Center Core





# HPC Software Infrastructure



# Unified Cluster Portfolio

**HPC Technical and Enterprise services**

**HPC application, development and cloud software portfolio**

**Advanced and specialty options**

(Accelerators, Visualization, other)

**Scalable data management**

(HP x9000 NSS, Lustre Cluster FS)

**Cluster management layer**

**HP CMU**

**Partner and Open  
Source choice**

**Microsoft Windows  
HPC Server 2008**

**Operating environment and OS extensions**

**Linux**

**Windows**

**HP cluster platforms**

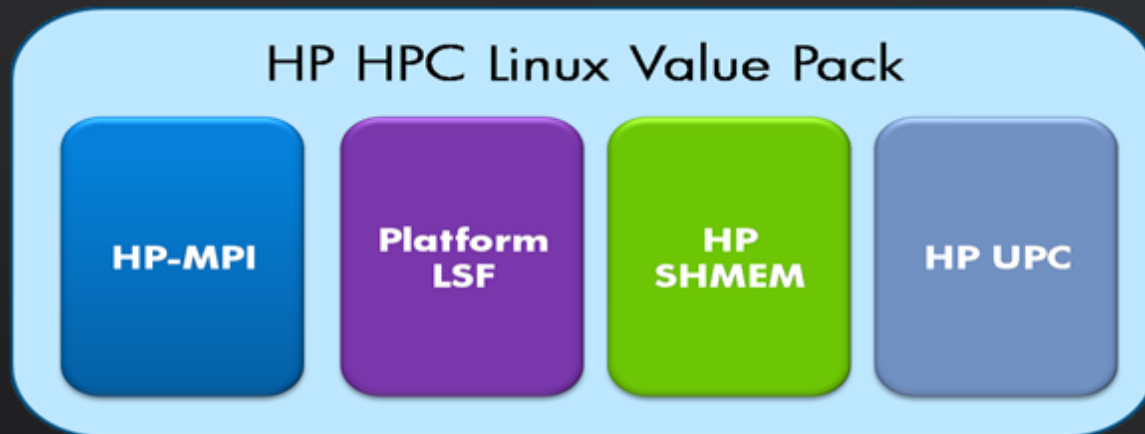
HP ProLiant servers, HP BladeSystem, multiple interconnects

**HP Datacenter Products & Services**



# A la Carte cluster options for HP Clusters

- Operating systems: RHEL, SLES, or customer-supported community distributions; Microsoft Windows HPC Server 2008
- Cluster Management: HP CMU, or third party, via SLMS or customer installed (e.g., ROCKS, Platform Cluster Manager)
- MPI: HP-MPI, or third party/open source; Windows MPI
- Workload manager: Platform LSF (via SLMS now), Altair PBS Pro (HP SKU), Adaptive Computing Moab (via SLMS)



# Datacenters – Good looking

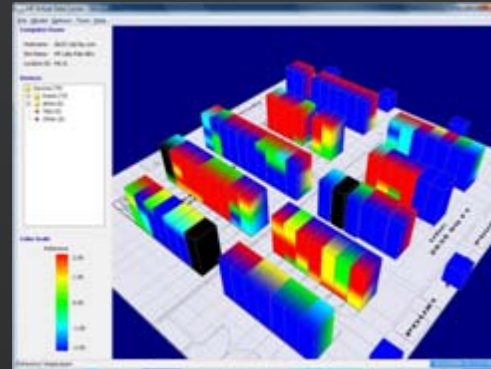


# ADAPTIVE AND SCALABLE SOFTWARE FOR HPC Datacenters

Edge 3D Visualization

## Edge Futures

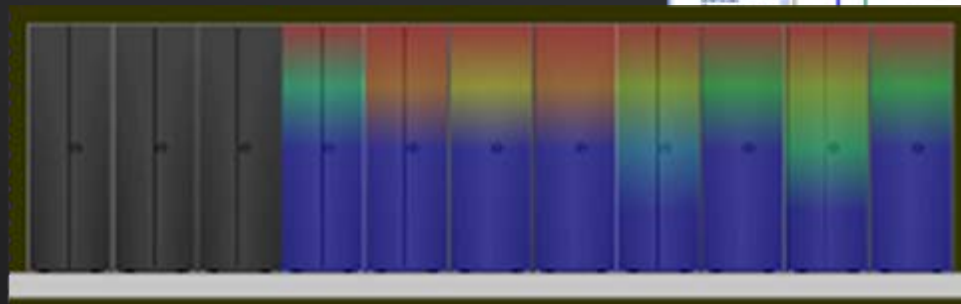
- Integration with Insight Control
- One pane of glass power and cooling visualization
- True 3D visualization
- Macro Data Center view
- Micro rack level view



IPM Rack View



Edge 2D Front View



# Datacenters – PODs



# Trend? Are Next Generation Data Centers Ugly?



# Advantages: Fast Deployment and Time to Operation: Efficient to Build and Rebuild



- Container backed into truck bay on mfg floor
- Racks assembled and then put into containers
- Truck pulls out with fully-configured container to the customer site...



# HP POD products and concepts

- 22 50U racks 40ft
- 600kW power capacity
- Designed for high density deployments – max 34kW per rack
- Flexible for redundant or non-redundant deployments



- 10 50U racks 20ft
- Modular design for better supply chain efficiency
- Flexibility to customize



- Rugged exterior
- EMI shielding
- Designed for portability



# Future Peta-scale Centers



# **Peta-scale Implementation Example: TITECH Tsubame 2.0**



# TSUBAME 2.0 Overview

- Compute nodes: 2.4PFlops (CPU+GPU)
  - **New SL-node** >>1408 thin nodes, each with 2 Westmere-EP and 3 NVIDIA M2050
    - 1347 with 54GB and SSD 60GB, 41 with 96GB and SSD 120GB
    - **Suse Linux Enterprise Server or Windows HPC Server**
  - **DL580 G7** Medium (24) and Fat (10) nodes, with 2 NVIDIA S1070
    - Medium: 128GB plus SSD 120GB x4
    - Fat: 256GB plus SSD 120GB x4
- QDR InfiniBand, full bisection, non-blocking
  - Spine: **Voltaire Grid Director 4700** 12 x 324port
  - Edge: **Voltaire Grid Director 4036** 179 x 36 port and **4036E** 6 x 34port/10GbE 2 port
- Storage: 5.93PB
  - Lustre file system 5.93PB: **DDN SFA 10000** (10/rack, 5 racks) and DL360 G6 (30)
  - Home file system: 1.2PB: **DDN SFA 10000** (10/rack, 1 racks), BlueArc Mercury 100 (2) and DL360 G6 (30)
- Press release (Japanese):
  - <http://www.gsic.titech.ac.jp/sites/default/files/pdf/TSUBAME/press.pdf>



# TSUBAME 2.0 System Overview

## Storage system $\square$ Total 7.13PB (Lustre+ home)

Lustre file system (DDN SFA10K) 5.93PB



x5

MDS,OSS  
 HP DL360 G6 30nodes  
 Storage  
 DDN SFA10000 x5  
 ( 10 enclosure x5)  
 Lustre  $\square$  5File System  $\square$   
 OSS: 20 OST: 5.9PB  
 MDS: 10 MDT: 30TB

OSS x20 MDS x10

Home directory region 1.2PB



Storage Server  
 HP DL380 G6 4nodes  
 BlueArc Mercury 100 x2  
 Storage  
 DDN SFA10000 x1  
 $\square$  10 enclosure x1  $\square$

NFS,CIFS  $\square$  x4 NFS,CIFS,iSCSI  $\square$  x2

Existing system

Tapa System

SupreTitenet

SupreSinet3

## Interconnect: Full bi-section non-blocking

Core Switch



12switches

Voltaire Grid Director 4700 12switches  
 IB QDR: 324port

Edge Switch



179switches

Voltaire Grid Director 4036  
 179switches IB QDR : 36 port

Edge Switch(10GbE port  $\square$   $\square$ )



6switches

Voltaire Grid Director 4036E 6  
 switches  
 IB QDR:34port  
 10GbE: 2port

Management nodes

## Compute nodes $\square$ 2.4PFlops(CPU+GPU)

"THIN" nodes



1408nodes (32node x44 Rack)

1408 SL nodes  
 CPU Intel Westmere-EP 2.93GHz  
 Turbo boost 3.196GHz  $\square$  12Core/node  
 Mem: 54GB (1347 nodes)  
 96GB (41 nodes)  
 GPU NVIDIA M2050 515GFlops,3GPU/node  
 SSD 60GB x 2 120GB (54GB nodes)  
 120GB x 2 240GB (96GB nodes)  
 OS: Suse Linux Enterprise Server  
 Windows HPC Server

CPU Total: 215.99TFLOPS(Turbo boost 3.196GHz)  
 CPU+GPU: 2391.35TFlops  
 Memory Total  $\square$  80.55TB  
 SSD Total  $\square$  173.88TB

"Med" nodes



DL580 G7 24nodes  
 CPU Intel Nehalem-EX 2.0GHz  
 32Core/node  
 Mem:137GB(=128GiB)  
 SSD 120GB x 4 480GB  
 OS: Suse Linux Enterprise Server

CPU Total: 6.14TFLOPS

"Fat" nodes



DL580 G7 10nodes  
 CPU Intel Nehalem-EX 2.0GHz  
 32Core/node  
 Mem:274GB(=256GiB)  $\square$  8nodes  
 549GB(=512GiB)  $\square$  2nodes  
 SSD 120GB x 4 480GB  
 OS: Suse Linux Enterprise Server

CPU Total: 2.56TFLOPS

PCI -E gen2 x16 x2slot/node

GSIC:NVIDIA Tesla S1070GPU



# Trends in Efficiency

	5 yrs. ago	2010	2015
PUE	2, 3, Higher	1.1 Great	?
UPS Efficiency (Part of PUE)	94%	98%+	?
Power Supply Efficiency	75%	94%+	?
Fan Power per 2s Node	60+ W	2-10 W ( $< 1\%$ ) (some think 0)	?



# Fundamental Research



# The Prediction of a New Circuit Element: the Memristor

Ohm  
1827

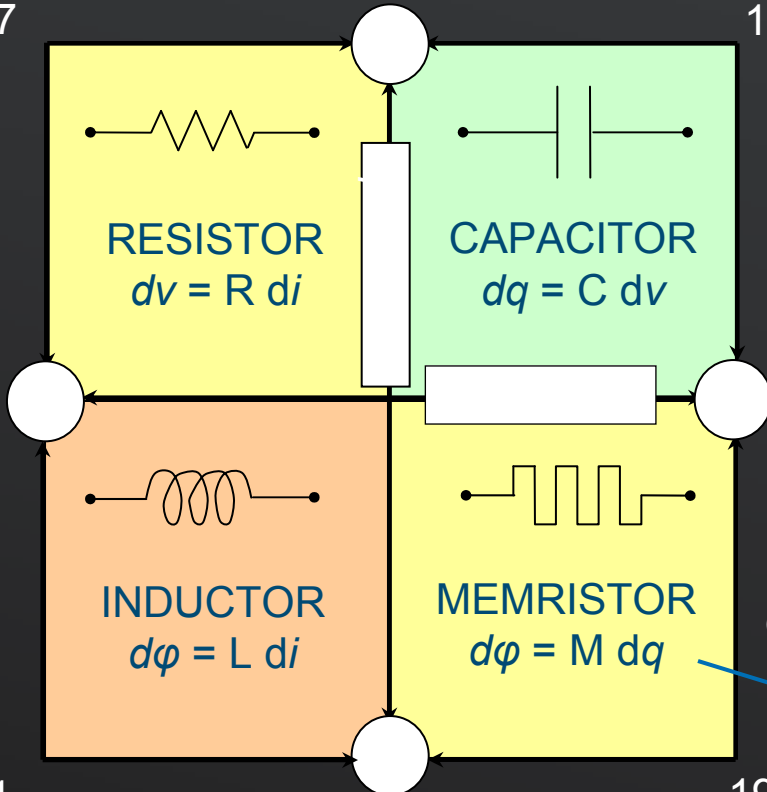
Von Kleist  
1745



L. O. Chua, *IEEE Trans. Circuit Theory*  
18, 507 (1971)

1831  
Faraday

1971  
Chua



rigorous  
definition

$$v(t) = R[w, i(t)]i(t)$$

Quasi-static conduction eq.-  
 $R$  depends on state variable  $w$

$$\frac{dw(t)}{dt} = f[w, i(t)]$$

Dynamical equation –  
Evolution of state in time





# First Hybrid CMOS-Memristor Chip

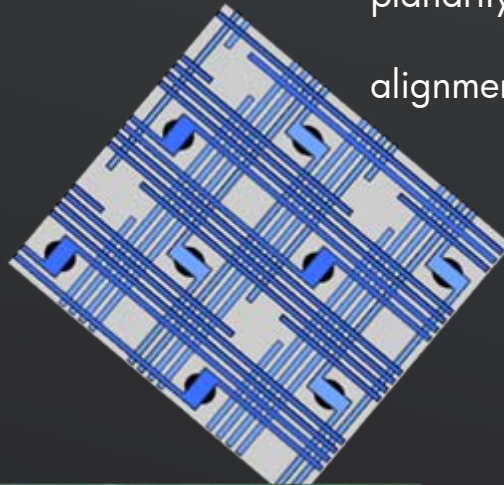
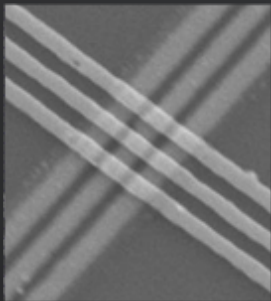
Issues that had to be overcome:

planarity

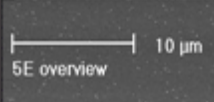
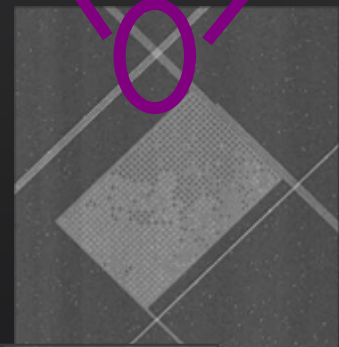
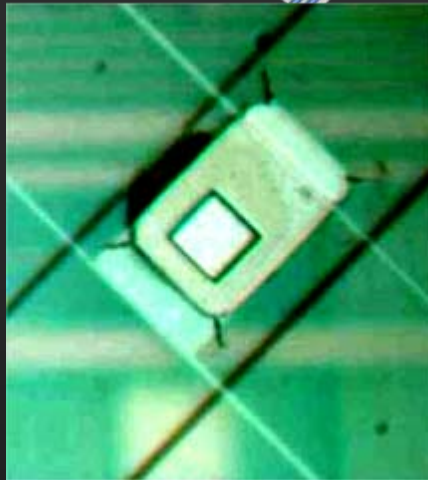
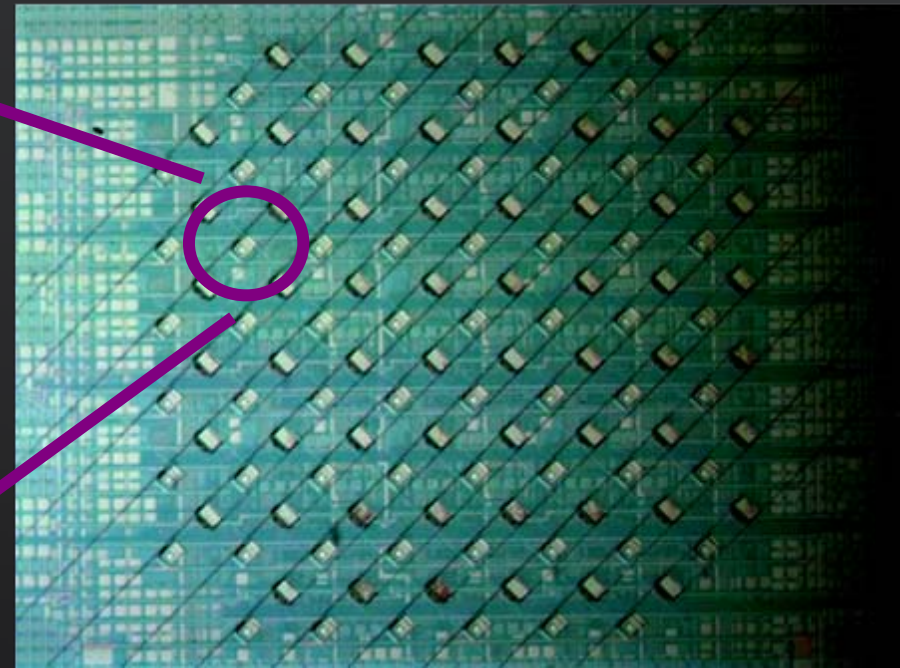
alignment of fine features

3x3 100nm nanowire

Crossbar junctions



CMOS chip with memristive devices

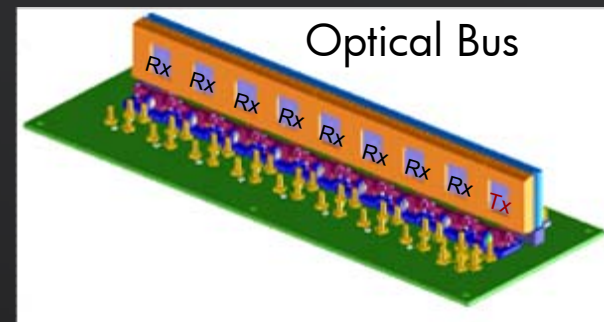
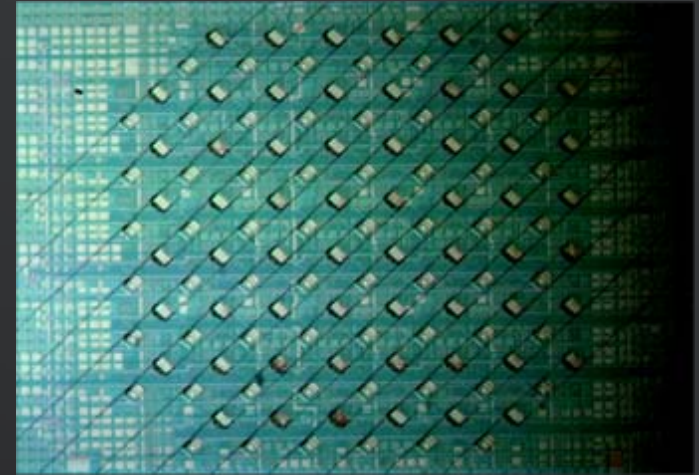


Connecting the CMOS layer with the nanowire crossbar junctions

# Long-term Trends in HPC: Examples of HP Labs Innovation

- Capacity - Memristor (short for **memory resistor**)
  - Scales to extremely high density (many terabits/sq cm)
  - Non-volatile – essentially infinite data retention time
  - Reasonably fast (ns) and low energy (pJ)
- Bandwidth – Photonics
  - High bandwidth, and highly energy efficient
  - Photonic interconnects between systems available now
  - Long term research leading to photonic interconnects within systems and chips

CMOS chip with memristive devices



# Invitation



Attend HP-CAST in New Orleans, right before SC10, November 12 – 13 !!  
Worldwide User Group Conference  
**Focus Session: Energy Efficient Peta-scale Computing** - see [www.hp-cast.org](http://www.hp-cast.org)

# HP-CAST

**HP Consortium for Advanced Scientific and Technical Computing  
Word-Wide User Group Meeting  
Scalable Computing Infrastructure (ISS/SCI) Organization  
InterContinental Hotel, Fontenay 10, 20354 Hamburg, Germany  
May 28<sup>th</sup> – 29<sup>th</sup>, 2010**

## HP-CAST 14

**World-wide User Group Conference with Participation of  
NTIG (Nordic Technical Interest Group) & HP-CAST IBÉRICA  
Draft Agenda V2.1p**

**Thursday, May 27<sup>th</sup> – Registration & Get-Together**

<b>17:00 – 22:00</b>	<b>Registration</b>	
<b>19:00 – 22:00</b>	<b>HP-CAST Welcome Reception</b>	<b>All Attendees</b>

**Friday, May 28<sup>th</sup> – Conference**



# Thank You

