



Open Architectures for HPC and AI

Ander Ochoa Gilo – ander.ochoa.gilo@ibm.com
Cognitive Systems Technical Architect for SPGI
OpenPOWER Foundation member



<https://es.linkedin.com/in/anderotxoa>



@AnderOtxoa

Current Mayor Architectures



PROPIETARY ARCHITECTURES



- X86
- Intel
- AMD



LICENCEABLE ARCHITECTURE



- ARM
- Samsung
- Apple
- ...

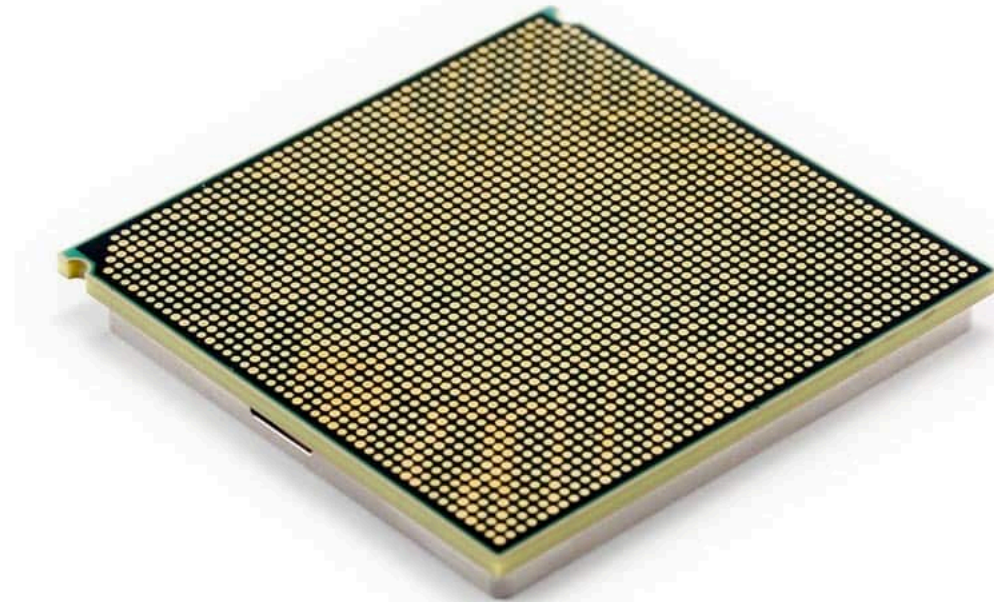
OPEN ARCHITECTURES

- RISC-V
- OpenPOWER





POWER9



OpenPOWER Foundation (members @ 09/2019)

Platinum



Gold



Silver

ZCRSI

@rchanan

ADDFOR
s.p.a.

Aerobyte

Bowmicro,
Ltd.



Shanghai Thinkforce Electronic Technology Co., Ltd.

Shenyang China-Bigdata Technology Co., Ltd.

Shenyang Yantou Cloud Computing Technology Co., Ltd.

DaoCloud

Shanghai Weida Cloud Power Co., Ltd.

Shenzhen Yu-Yuan Intelligence



Associate & Academic

A*STAR Computational Resource Centre



J.K.K. Nairra Educational Institutions



Aryant Institute of Technology

ASTRI



Barcelona Supercomputing Center

Bauman Moscow State Technical University

BV Raju Institute of Technology (BVRIT)

California Institute of Technology (Caltech)



Centre for the Development of Advanced Computing

Clemson University



College of Engineering - Anna University



Computer Network Information Center, Chinese Academy of Sciences

Computing Center FEBRAS



Custom Computing Research Group, Imperial College

Daegu Gyeongbuk Institute of Science & Technology (DGIST)

Dalian University of Technology

Delft University of Technology

Dept. of Comp. Sci. @ Univ. of Brasilia

FreeBSD

Garwood Center for Corporate Innovation

GENCI

GSIC, Titech, Global Scientific Information and Computing Center

Harbin Institute of Technology Weihai



Hasso Plattner Institute for Software Engineering

ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



ICM University of Warsaw



Indian Institute of Science (Computational Aerodynamics Lab)

Indian Institute of Science - Dept. of Computer Science & Automation



Indian Institute of Technology Roorkee

Institute of Communication and Computer Systems (ICCS)

Institute for Dev. & Research in Banking Technology (IDRB)

Institute of Physics

Interaction

International Robotics School of Myths

Istituto Nazionale di Astrofisica (INAF)



KLE Technological University

KTH Royal Institute of Technology

Kuppam Engineering College





OpenPOWER – Partner Engagement Summary (2018)

<u>New /Pending</u>	<u>Discussion</u>	<u>Concept</u>	<u>Design</u>	<u>Power ON</u>	<u>Ship Ready</u>	<u>Total</u>
1	0	2	3	7	29	39

China							
ROW							

Engagement		3Q18	4Q18	1Q19	2Q19	3Q19	4Q19
	On Track	PON					
	On Track	PON	GA				
	On Track		GA				
	On Track			GA			
	On Track						
	On Track				PON		
	On Track		GA				
	On Track		GA				
	On Track	Concept					
	On Track			GA			
	On Track					PON	

The TWO most POWERFUL HPC systems: Summit & Sierra

The United States Department of Energy together with Oak Ridge National Laboratory and Lawrence Livermore National Laboratory have contracted **IBM and Nvidia** to build two supercomputers, the Summit and the Sierra, that are **based on POWER9 processors coupled with Nvidia's Volta GPUs**. These systems went online in 2018.

Rank	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband , IBM DOE/SC/Oak Ridge National Laboratory United States	2,414,592	148,600.0	200,794.9	10,096
2	Sierra - IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband , IBM / NVIDIA / Mellanox DOE/NNSA/LLNL United States	1,572,480	94,640.0	125,712.0	7,438
3	Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway , NRCPC National Supercomputing Center in Wuxi China	10,649,600	93,014.6	125,435.9	15,371
4	Tianhe-2A - TH-IVB-FEP Cluster, Intel Xeon E5-2692v2 12C 2.2GHz, TH Express-2, Matrix-2000 , NUDT National Super Computer Center in Guangzhou China	4,981,760	61,444.5	100,678.7	18,482

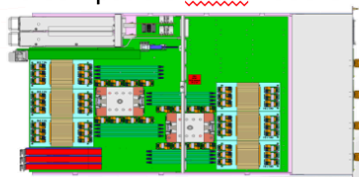


Summit Overview



Compute Node

2 x POWER9
6 x NVIDIA GV100
NVMe-compatible PCIe 1600 GB SSD



25 GB/s EDR IB- (2 ports)
512 GB DRAM- (DDR4)
96 GB HBM- (3D Stacked)
Coherent Shared Memory

Components

IBM POWER9

- 22 Cores
- 4 Threads/core
- NVLink



NVIDIA GV100

- 7 TF
- 16 GB @ 0.9 TB/s
- NVLink



Compute Rack

18 Compute Servers
Warm water (70°F direct-cooled components)
RDHX for air-cooled components



39.7 TB Memory/rack
55 KW max power/rack

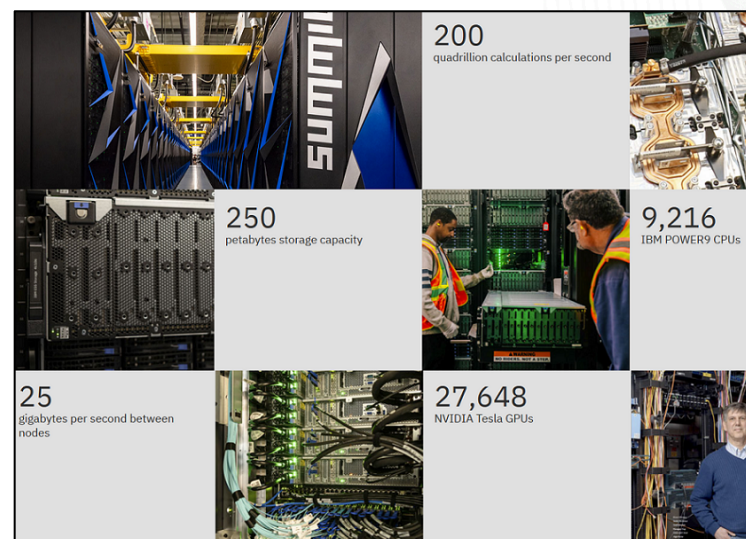
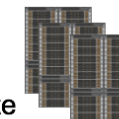
Compute System

10.2 PB Total Memory
256 compute racks
4,608 compute nodes
Mellanox EDR IB fabric
200 PFLOPS
~13 MW



GPFS File System

250 PB storage
2.5 TB/s read, 2.5 TB/s write



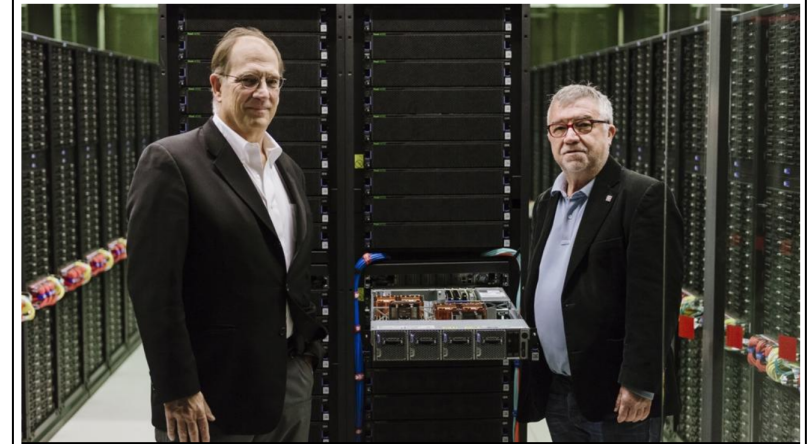
OAK RIDGE
National Laboratory
LEADERSHIP
COMPUTING
FACILITY

BSC Mare Nostrum 4



El MareNostrum 4 amplía su capacidad para investigar en IA

- El BSC se convierte en el primer centro de Europa en ofrecer acceso a las mismas tecnologías que el nuevo supercomputador Summit, de EE.UU., el más potente del mundo



David Durek, de IBM, y Mateo Valero, director del BSC (BSC-CNS)

Part of Mare Nostrum 4

- 3 Racks
- 54 Power9 Systems
 - 54 AC922 servers
 - 4x Nvidia V100
 - 512 GB RAM
- 6.4 TB NVMe storage
- 1.48 Pflops !

IBM's Pangea III is the world's most powerful commercial supercomputer



Total's Supercomputer Ranked First in Industry Worldwide



The new **IBM POWER9-based supercomputer (25 PFLOPS & 50 Pbytes)** will help Total more accurately locate new resources and better assess the potential of new opportunities.

According to Total **Pangea III requires 1.5 Megawatts, compared to 4.5 MW for its predecessor system.** Combined with the increased performance of Pangea III, Total has reported that they have observed that the new system **uses less than 10% the energy consumption per petaflop as its predecessor.**

- **Higher Resolution Seismic Imaging in exploration and development phase**
- **Reliable Development and Production Models**
- **Asset Valuation and Selectivity**





MIT SATORI “Sudden Enlightenment”

IBM POWER9 + NVIDIA

<https://researchcomputing.mit.edu/satori/home/>

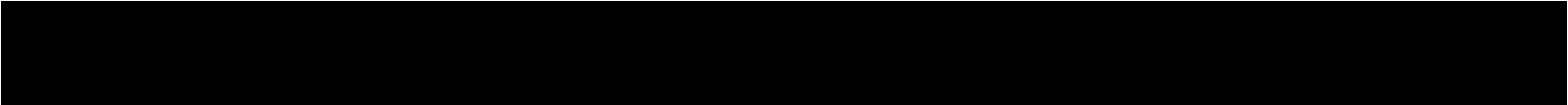


satori.mit.edu is the name of a new scalable AI oriented hardware resource for research computing at MIT. It is made possible by a donation through IBM Global Universities Program. Provided as a gift from IBM it will help further the aims of the new MIT Stephen A. Schwarzman College of Computing and other campus initiatives that are combining supercomputing power and AI algorithmic innovation.

- 64 IBM Power 9 Nodes
- 256 NVidia V100 GPUs
- EDR Infiniband
- 2PB storage
- 8TB GPU memory
- 64 TB main memory
- IBM Power AI Software
- Cloud Integration



?



IBM POWER9 + NVIDIA



?



- [redacted]
- [redacted]
- [redacted]
- [redacted]
- [redacted]
- [redacted]
- [redacted]

?

IBM® Power System™ Accelerated Compute Server (AC922)



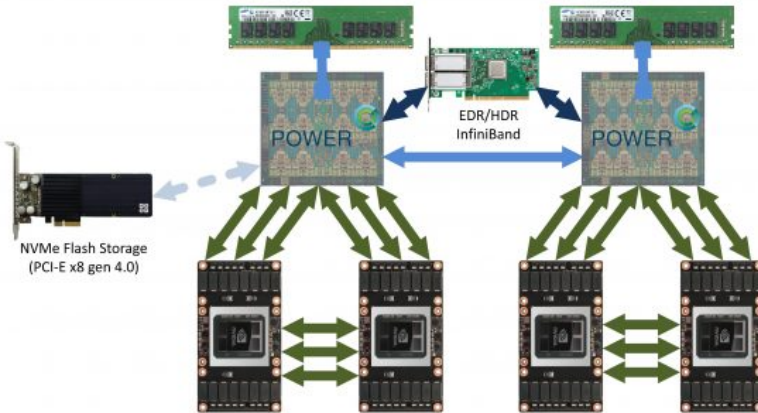
The Best Server for Enterprise AI

POWER9

An acceleration superhighway.
The only processor specifically designed for the AI era.

Server Block Diagram

Power Systems AC922 with NVIDIA Tesla V100 with Enhanced NVLink GPUs



- IBM POWER9 SMP bus
- Direct Attach DDR4 memory (~170GB/s BW per CPU)
- PCI-Express x8 (gen 4.0) bus with CAPI for IB (12.8GB/s)
- 1x PCI-E x8 4.0 from each CPU to IB (multi-socket host direct)
- PCI-Express x8 (gen 4.0) bus with CAPI (12.8GB/s)
- 25GB/s NVIDIA NVLink Interconnect (50GB/s bi-directional)
- 75GB/s of bandwidth between points (3 links)

4x

Threads per
core vs x86

9.5x

Up to 9.5x more I/O
bandwidth than x86

2.6x

More RAM
possible vs. x86

1st

CPU to deliver
PCIe gen 4

Innovation from an ecosystem of partners across the stack and open to the core



OpenPOWER™
READY

OPEN POWER

250+ OpenPOWER members co-design around the core to accelerate cognitive and general workloads



+++

OPEN SOURCE

Density, speed-up, compaction of the most ubiquitous open source engines



+++

OPEN CAPI

Laying the groundwork for faster coherent open interfaces to attach to accelerators



+++

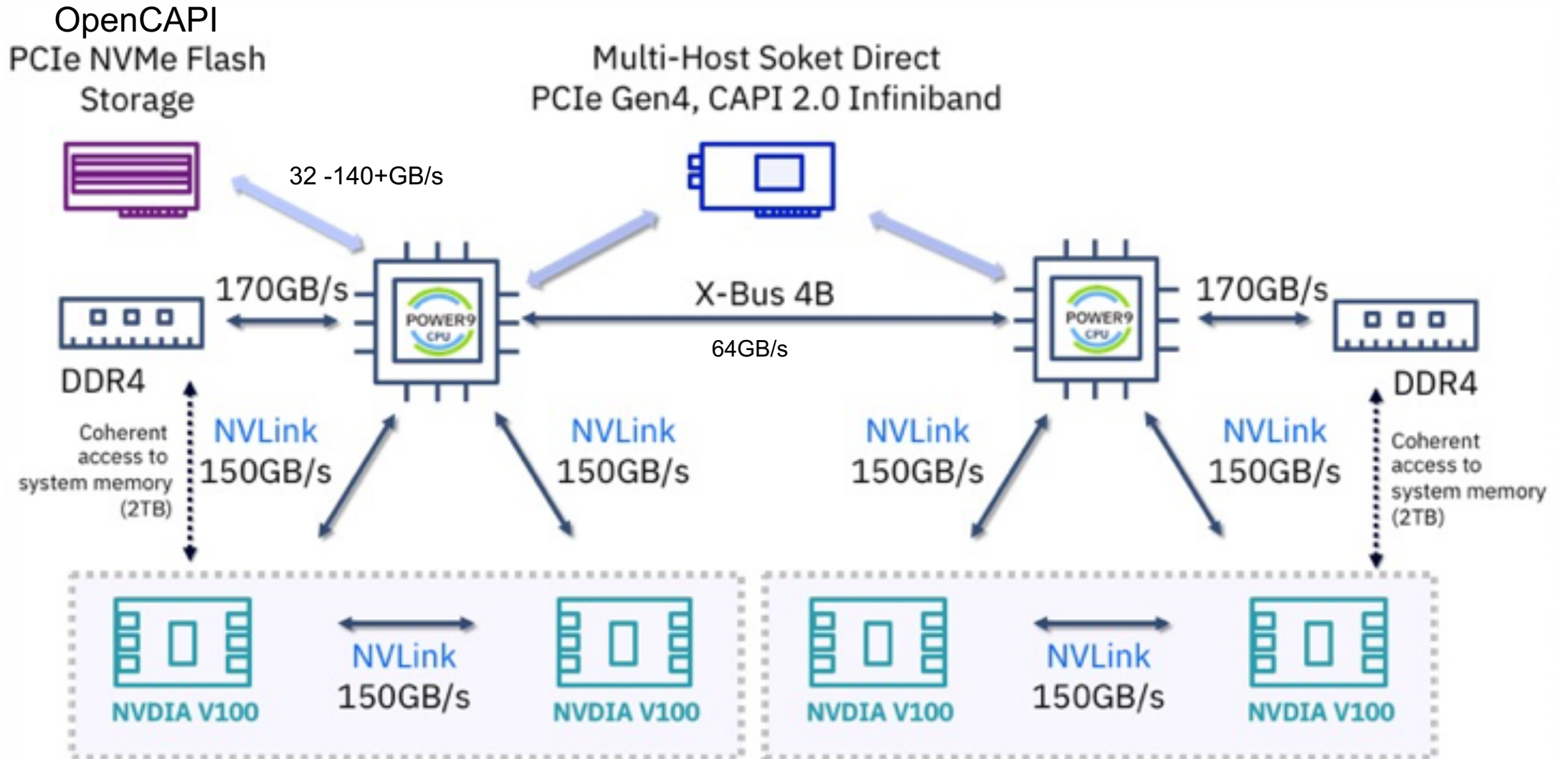
OPEN FRAMEWORKS

The industry's most ubiquitous Cognitive/AI frameworks - optimized and accelerated.



+++

AC922 System buses and components diagram





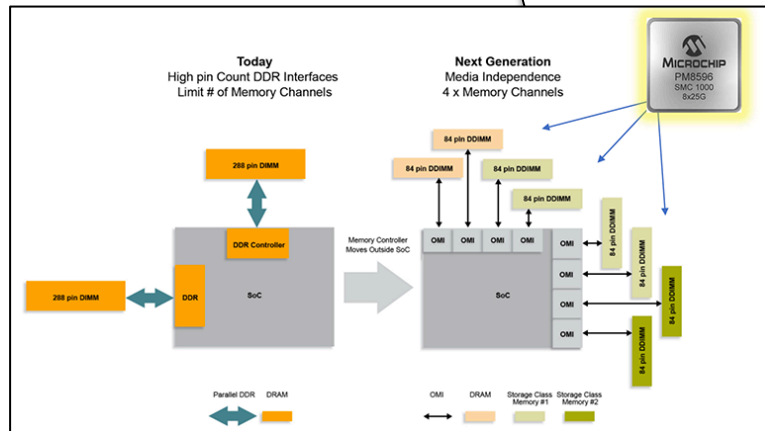
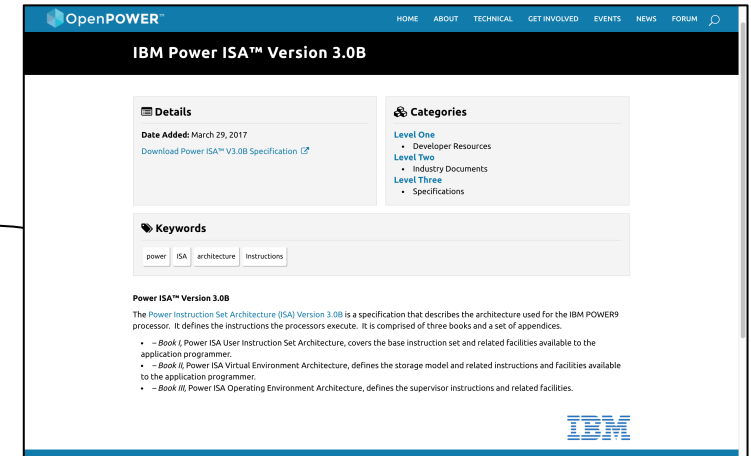
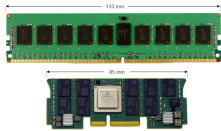
**The future of open source
hardware**

Open POWER moves to Linux Foundation

IBM Demonstrates Commitment to Open Hardware Movement

https://openpowerfoundation.org/?resource_lib=power-isa-version-3-0

- **Open POWER ISA** (instruction set architecture)
 - **Microwatt** (Power ISA soft core)
- **Open CAPI** (coherent accelerator processor interface)
- **OMI** (Open Memory Interface)



Power Systems **OpenCAPI 4.0: Asymmetric Open Accelerator Attach** IBM

Roadmap of Capabilities and Host Silicon Delivery

Accelerator Protocol	CAPI 1.0	CAPI 2.0	OpenCAPI 3.0	OpenCAPI 4.0	OpenCAPI 5.0
First Host Silicon	POWER8 (GA 2014)	POWER9 SO (GA 2017)	POWER9 SO (GA 2020)	POWER9 AIO (GA 2021)	POWER10 (GA 2021)
Functional Partitioning	Asymmetric	Asymmetric	Asymmetric	Asymmetric	Asymmetric
Host Architecture	POWER	POWER	Any	Any	Any
Cache Line Size Supported	128B	128B	64/128/256B	64/128/256B	64/128/256B
Attach Vehicle	PCIe Gen 3 Tunneled	PCIe Gen 4 Tunneled	25 G (open) Native DL/TL	25 G (open) Native DL/TL	32/50 G (open) Native DL/TL
Address Translation	On Accelerator	Host	Host (secure)	Host (secure)	Host (secure)
Native DMA to Host Mem	No	Yes	Yes	Yes	Yes
Atomics to Host Mem	No	Yes	Yes	Yes	Yes
Host Thread Wake-up	No	Yes	Yes	Yes	Yes
Host Memory Attach Agent	No	No	Yes	Yes	Yes
Low Latency Short Msg	4B/8B MMIO	4B/8B MMIO	4B/8B MMIO	128B push	128B push
Posted Writes to Host Mem	No	No	No	Yes	Yes
Caching of Host Mem	RA Cache	RA Cache	No	VA Cache	VA Cache

© 2019 IBM Corporation 14

<https://openpowerfoundation.org/the-next-step-in-the-openpower-foundation-journey/>
<https://www.talospace.com/2019/09/a-beginners-guide-to-hacking-microwatt.html>

OpenCAPI

- **What is OpenCAPI?**

- OpenCAPI is an **Open Interface Architecture** that allows any microprocessor to attach to
 - Coherent user-level accelerators and I/O devices
 - Advanced memories accessible via read/write or user-level DMA semantics
 - Agnostic to processor architecture

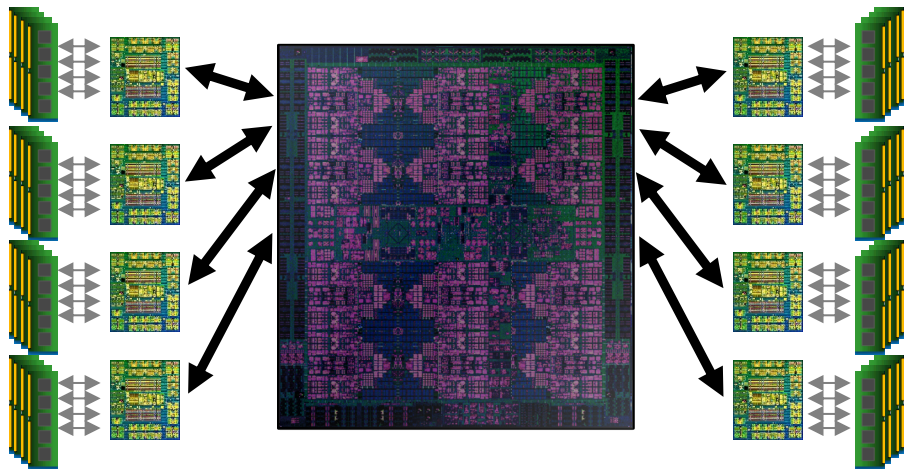
- **Key Attributes of OpenCAPI**

- High-bandwidth, low latency interface optimized to enable streamlined implementation of attached devices
 - 25Gbit/sec signaling and protocol built to enable very low latency interface on CPU and attached device
 - Complexities of coherence and virtual addressing implemented on host microprocessor to simplify attached devices and facilitate interoperability across multiple CPU architectures
- Attached devices operate natively within an application's user space and coherently with processors
 - Allows attached device to fully participate in application without kernel involvement/overhead
- Supports a wide range of use cases and access semantics
 - Hardware accelerators
 - High-performance I/O devices
 - Advanced memories
- 100% Open Consortium / All company participants welcome / All ISA participants welcome

POWER9 Family Memory Architecture

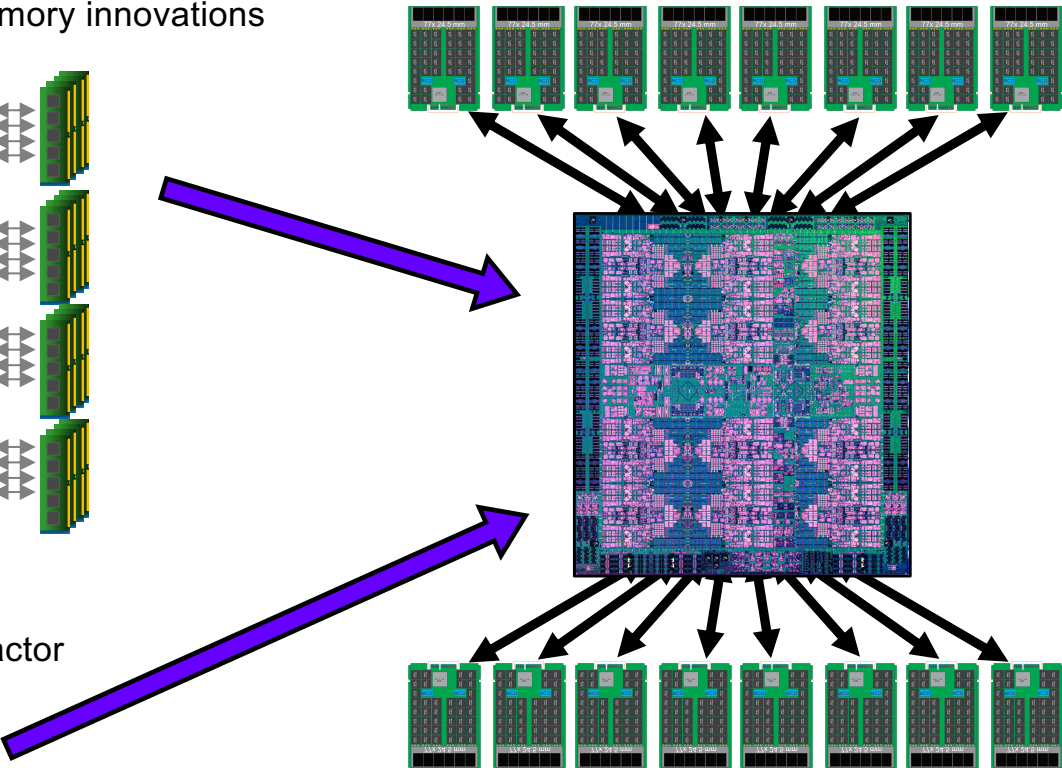
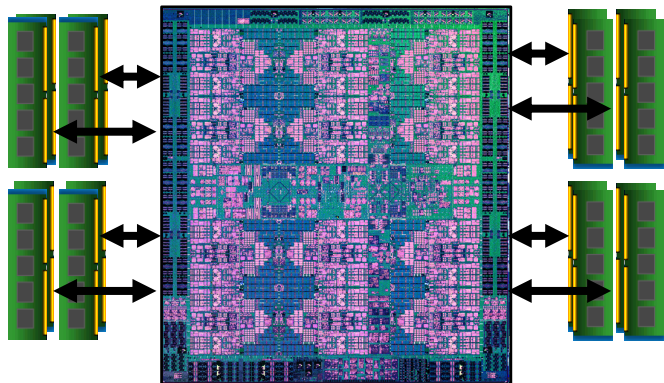
Scale Up
Buffered Memory

Superior RAS, High bandwidth, High Capacity
Agnostic interface for alternate memory innovations

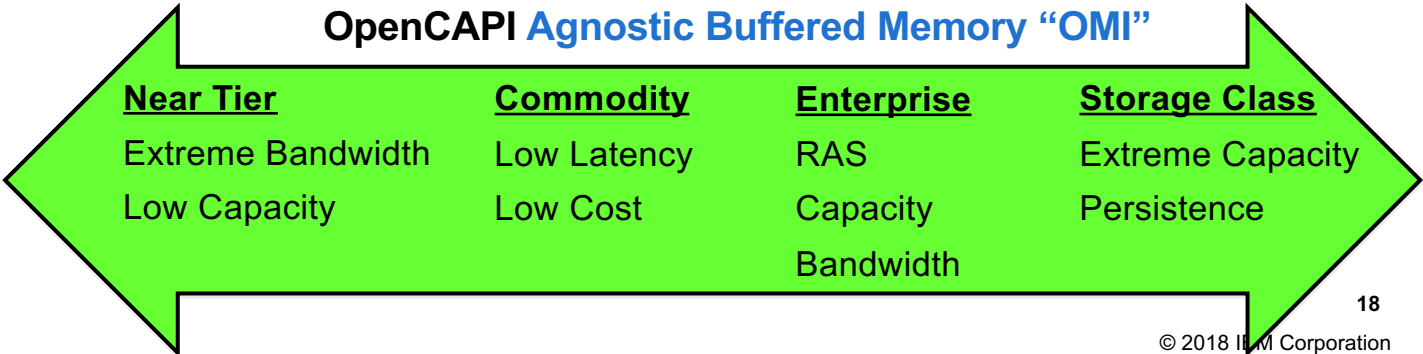


Scale Out
Direct Attach Memory

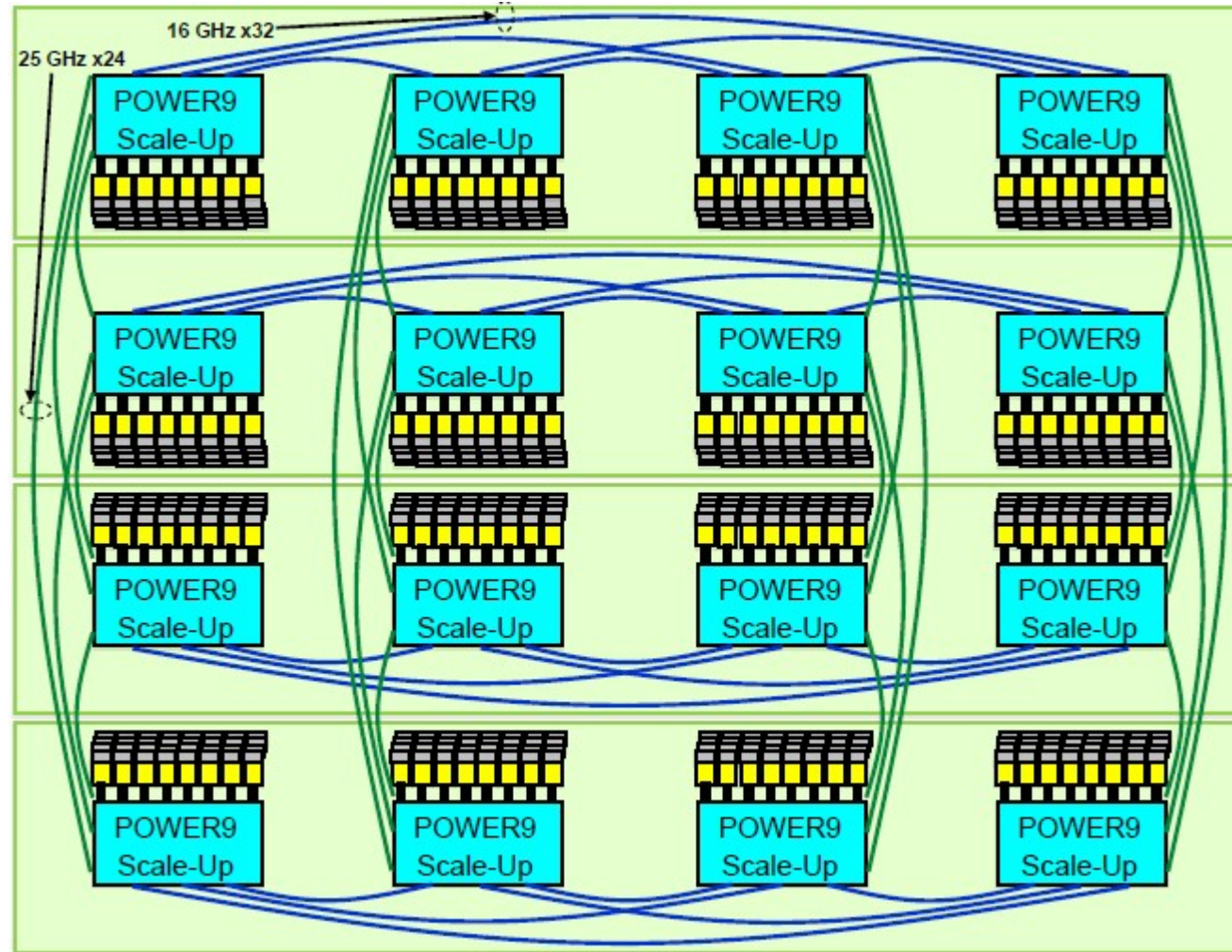
Low latency access
Commodity packaging form factor



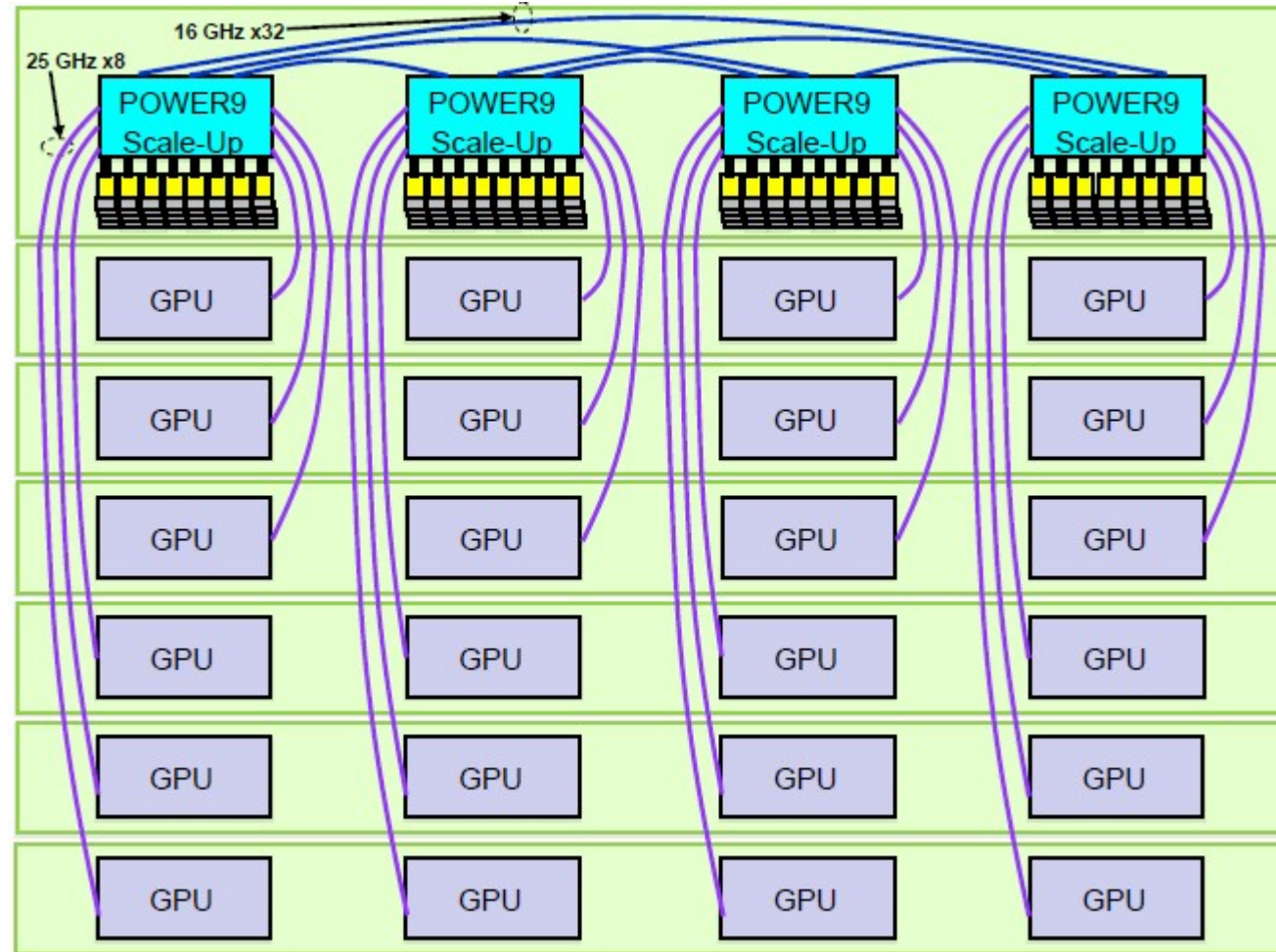
OpenCAPI Agnostic Buffered Memory “OMI”



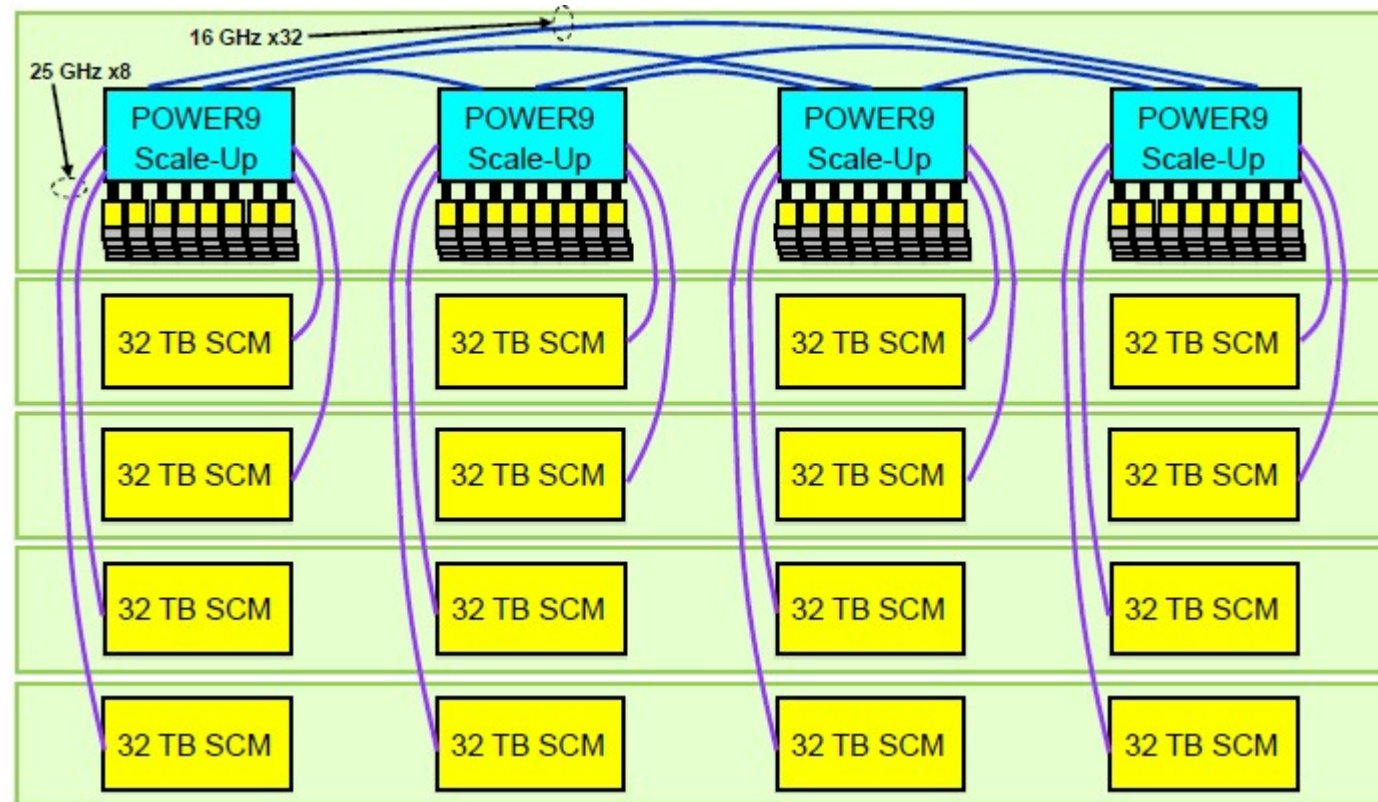
Somewhere in a possible future... (16 ways POWER9 server using OpenCAPI?)



Somewhere in a possible future... (24 GPU server using OpenCAPI/NVLINK?)



Somewhere in a possible future... (512TB Storage Class Memory using OpenCAPI?)



Somewhere in a possible future...

OpenPOWER X



“The Only Stupid Question is
the One You Don’t Ask”