

Exceptional service in the national interest



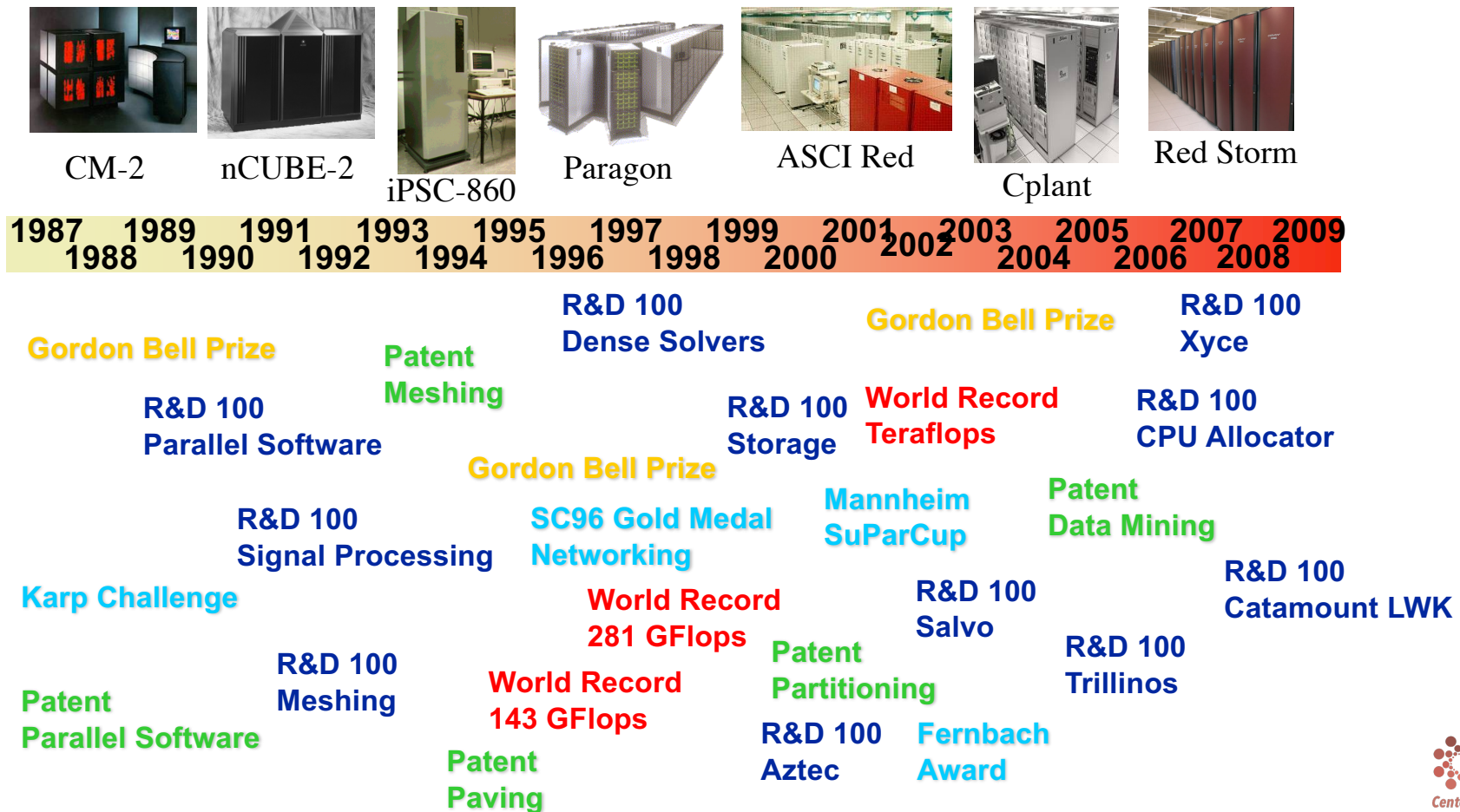
HPC Co-Design

Ron Brightwell, R&D Manager



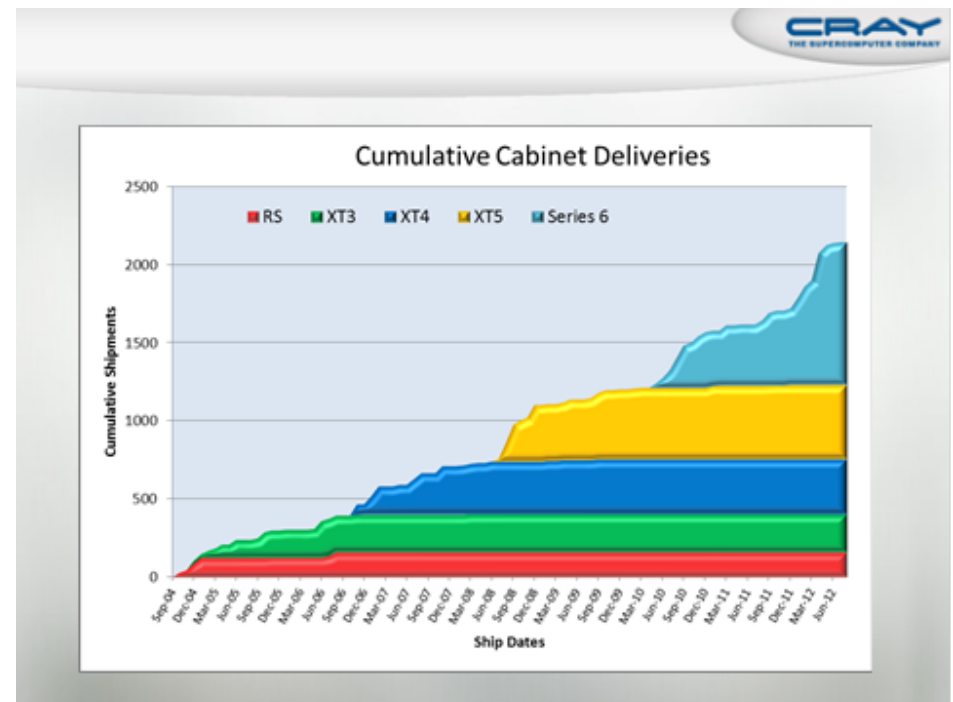
Sandia National Laboratories is a multi-mission laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000. SAND NO. 2011-XXXXP

Track Record of HPC Success



Red Storm – Prototype for Cray XT Series

- Architected by Sandia, engineered jointly with Cray
 - Sandia contributed to the design of the SeaStar network interface and router
- Sandia also developed
 - Lightweight kernel compute node OS
 - Scalable parallel job launching system
 - Portals high-performance interconnect programming interface
 - SeaStar firmware
- 140+ systems to 80 different customers worldwide
 - Including ORNL, NERSC, and LANL
- Following Red Storm, Cray's market share rose from 6% in 2002 to 21% in 2007*
- Revenue of \$1B +
- Basis of Cray's business today

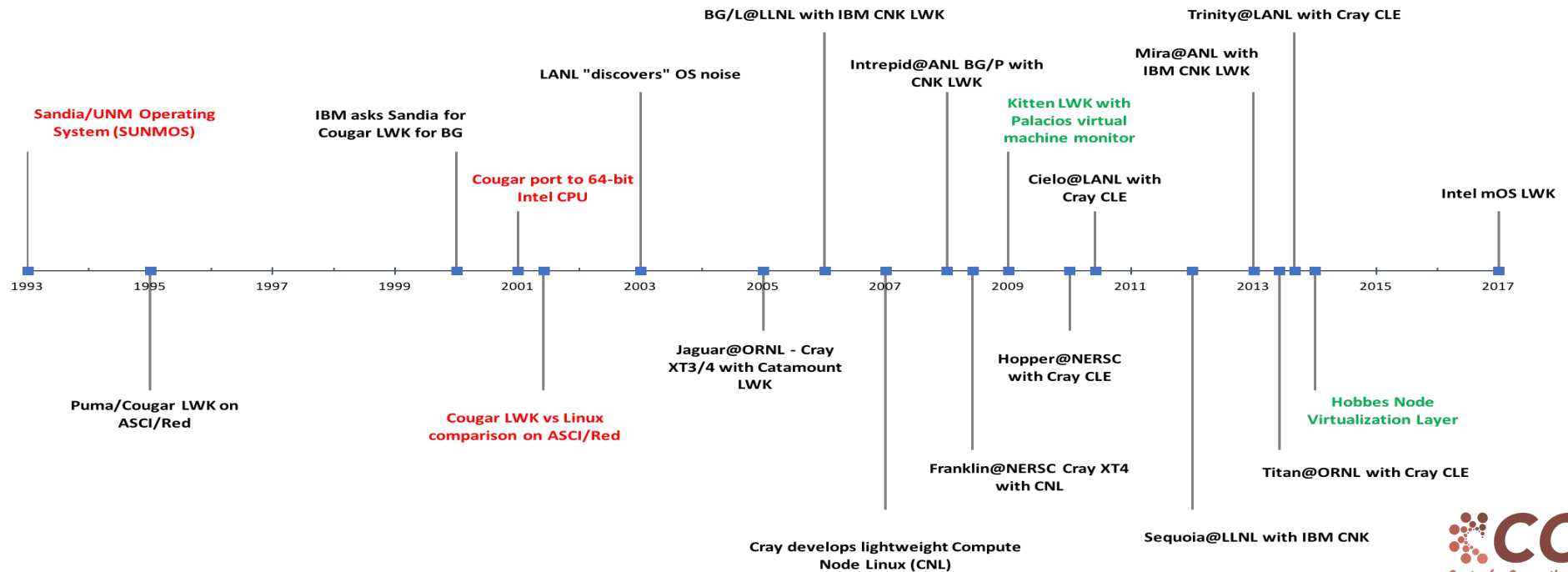


*Source: IDC #209251 *Technical Computing Systems: Competitive Analysis*, November 2007

Sandia's LWK Approach Has Had Broad Impact

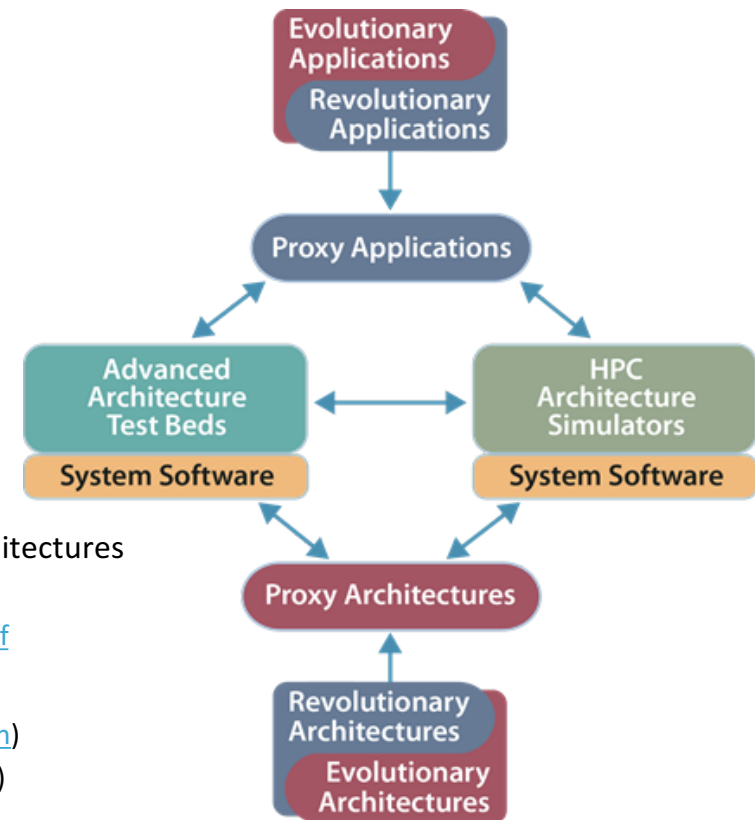
Sandia is the only DOE laboratory to partner with vendors to deploy a custom OS in production

- SUNMOS LWK on Intel Paragon; Cougar LWK on ASCI/Red; Catamount on Cray Red Storm
- Other vendors have followed the LWK model: IBM CNK for BG/{L,P,Q}; Cray's Linux Environment
- Every large-scale DOE distributed memory machine in the past 25 years has deployed a lightweight OS



Sandia HPC Co-design Capabilities

- Proxy Applications (Mantevo):
 - Application source for architecture-centric optimization and analysis
 - <http://mantevo.org>
- HPC Architectural Simulation Framework:
 - Structural Simulation Toolkit (SST)
 - Flexible to accommodate fidelity/speed tradeoffs
 - <http://SST-simulator.org>
- ASC Advanced Architecture Testbeds:
 - Evolving examples of COTS “state-of-the-art”
 - http://www.sandia.gov/asc/computational_systems/HAAPS.html
- Abstract Machine Model (AMM) Definitions and associated Proxy Architectures
 - Developed by SC/ASCR Computer Architecture Lab (SNL & LBL)
 - http://crd.lbl.gov/assets/pubs_presos/CALAbstractMachineModelsv1.1.pdf
- System software components
 - Kitten lightweight operating system (<https://github.com/HobbesOSR/kitten>)
 - Qthreads user-level thread library (<https://github.com/Qthreads/qthreads>)
 - Portals interconnect API (<https://github.com/Portals4/portals4>)



Portals Interconnect Programming Interface

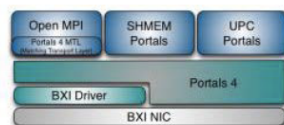
- Developed by Sandia, U. New Mexico, Intel
- Previous generations of Portals deployed on several production massively parallel systems
 - 1993: 1800-node Intel Paragon (SUNMOS)
 - 1997: 10,000-node Intel ASCI Red (Puma/Cougar)
 - 1999: 1800-node Cplant cluster (Linux)
 - 2005: 10,000-node Cray Sandia Red Storm (Catamount)
 - 2009: 18,688-node Cray XT5 – ORNL Jaguar (Linux)
- Focused on providing
 - Lightweight “connectionless” model for massively parallel systems
 - Low latency, high bandwidth
 - Independent progress
 - Overlap of computation and communication
 - Scalable buffering semantics
 - Protocol building blocks to support higher-level application protocols and libraries and system services
- At least three hardware implementations currently under development
- Portals influence can be seen in IB Verbs (Mellanox) and libfabric (Intel & others)



Portals Hardware and Co-Design Activities

BXI application environment

BXI comes with a complete software stack to provide optimal performance and reliability to all traditional HPC components.



BXI Computing stack

- ▶ Parallel applications can take full advantage of the capabilities of the Bxi network using MPI, SHMEM or UPC communication libraries.
 - ▶ All components are implemented directly using the Portals 4 API.
 - ▶ Kernel services are also implemented using the kernel Portals 4 implementation.
 - ▶ A Portals 4 LND (Lustre Network Driver) provides the Lustre parallel filesystem with a direct / native access to Portals 4.
 - ▶ The IPoPtI (IP over Portals) component makes it possible to have large scale, efficient and robust IP communication for legacy software.
-
- The diagram illustrates the software stack for high-performance computing. At the top, 'Lustre' (containing 'Portals 4 LND' and 'Lustre Network Driver') and 'IPoPtI (IP over Portals)' are shown. Below them, 'Portals 4 (Kernel)' is connected to the 'Bxi Driver'. The 'Bxi Driver' is then connected to the 'Bxi NIC' at the bottom. Arrows indicate the flow of data and communication between these components.



BXI Kernel services

For more information: Please contact hpc@atos.net

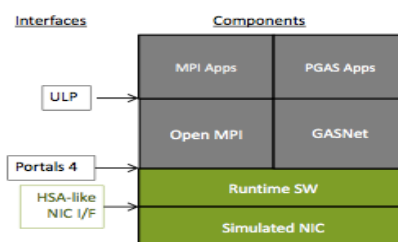
Your business technologists. Powering progress

AtoS



FASTFORWARD NIC SOFTWARE STACK

- ▲ Portals 4 API chosen for initial investigation
 - Supports multiple programming models: PGAS, MPI
- ▲ Implemented in thin software layer over hardware interface
- ▲ Leverage existing ULPs that have Portals 4 implementations
 - GASNet
 - Open MPI



EXPERIMENTAL FRAMEWORK

RESULTS

AMD

- ▲ All data collected in gem5^[6]
 - System call emulation mode (no OS)
 - AMD GPU model^[7]
 - Full Support for HSA
 - Tightly coupled system

- Portals 4-based NIC mode^[8]
 - Low-level RDMA network programming API currently supported by:
 - MPICH, Open MPI, GASNet, Berkeley UPC, GNU UPC, and others
 - XTQ implemented as an extension of the Portals 4 remote Put operation

CPU and Memory Configuration	
CPU Type	8-wide OOO, 4GHz, 8 cores
L1-Cache	64K, 2-way, 1 cycle
L2-Cache	2MB, 8-way, 4 cycles
L3-Cache	16MB, 16-way, 20 cycles
DRAM	DDR3, 4 Channels, 800MHz

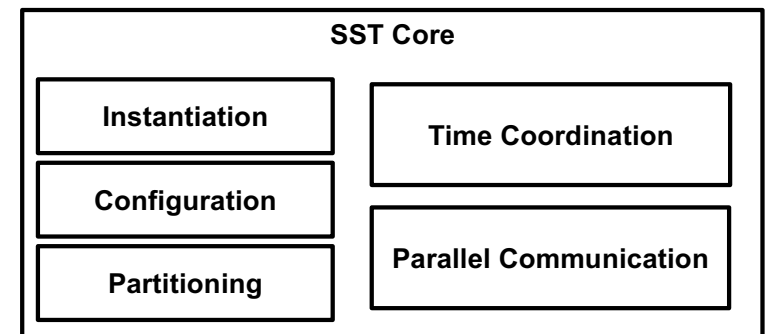
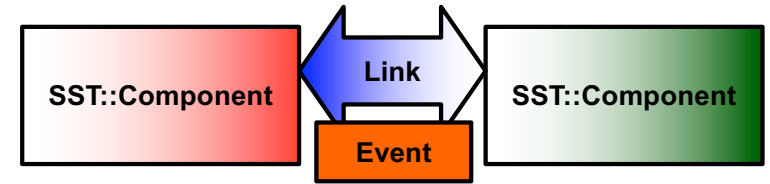
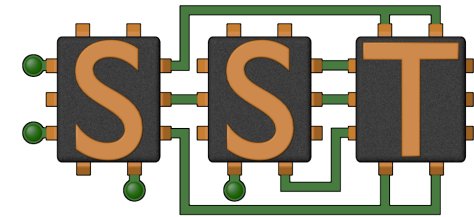
GPU Configuration	
GPU Type	1 GHz, 24 Compute Units
D-Cache	16KB, 64B line, 16-way, 4 cycles
I-Cache	32KB, 64B line, 8-way, 4 cycles
L2-Cache	768KB, 64B line, 16-way, 24 cycles

NIC Configuration	
Link Speed	100ns/ 100Gbps
Network API	Portals 4
Topology	Star

[19] EXTENDED TASK QUEUING: ACTIVE MESSAGES FOR HETEROGENEOUS SYSTEMS AUGUST 10, 2017

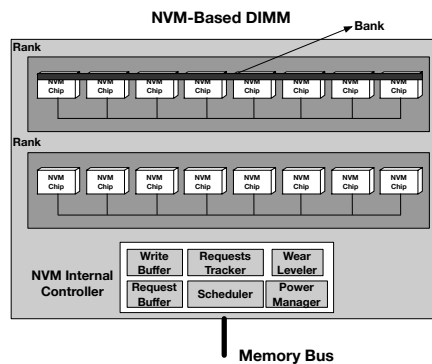
SST in a Nutshell

- Use Supercomputers to Design Supercomputers
- Parallel Discrete-Event Simulator Framework
 - Flexible framework allows multitude of custom simulators
 - Demonstrated scaling to over 512 processors
- Comes with many built-in simulation models
 - Processors, Memory, Network
- Open API
 - Easily extensible with new models
 - Modular framework
 - (Non-Viral) Open-source core
- Time-scale independent core
 - Handles Micro-, Meso-, Macro-scale simulations
- “Best of Breed” – Bring together work from Labs, Industry, Academia

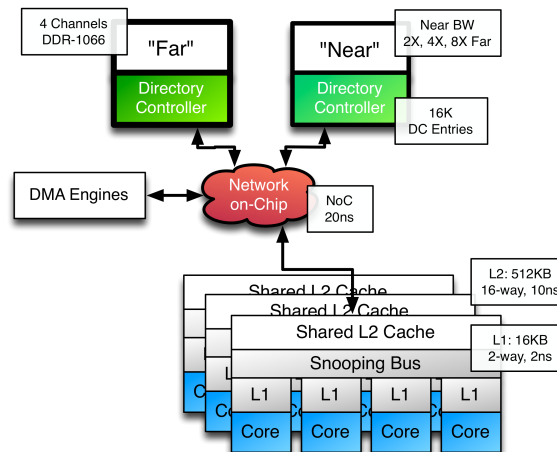


SST Use Cases

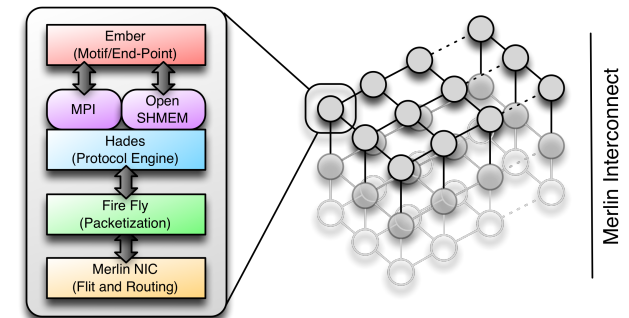
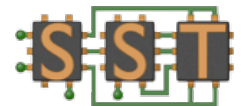
Emerging NV Memory Technologies



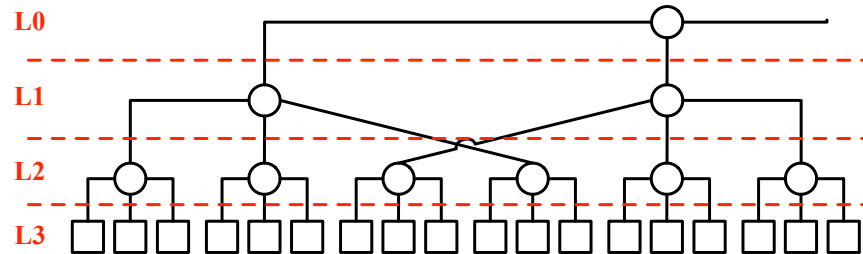
Multi-Level Memory (HBM+DDR+NV)



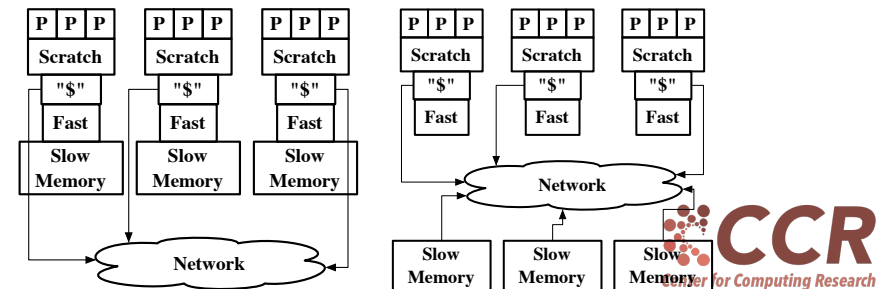
Communication/ Interconnect Modeling



Photonic Network Topology & Routing



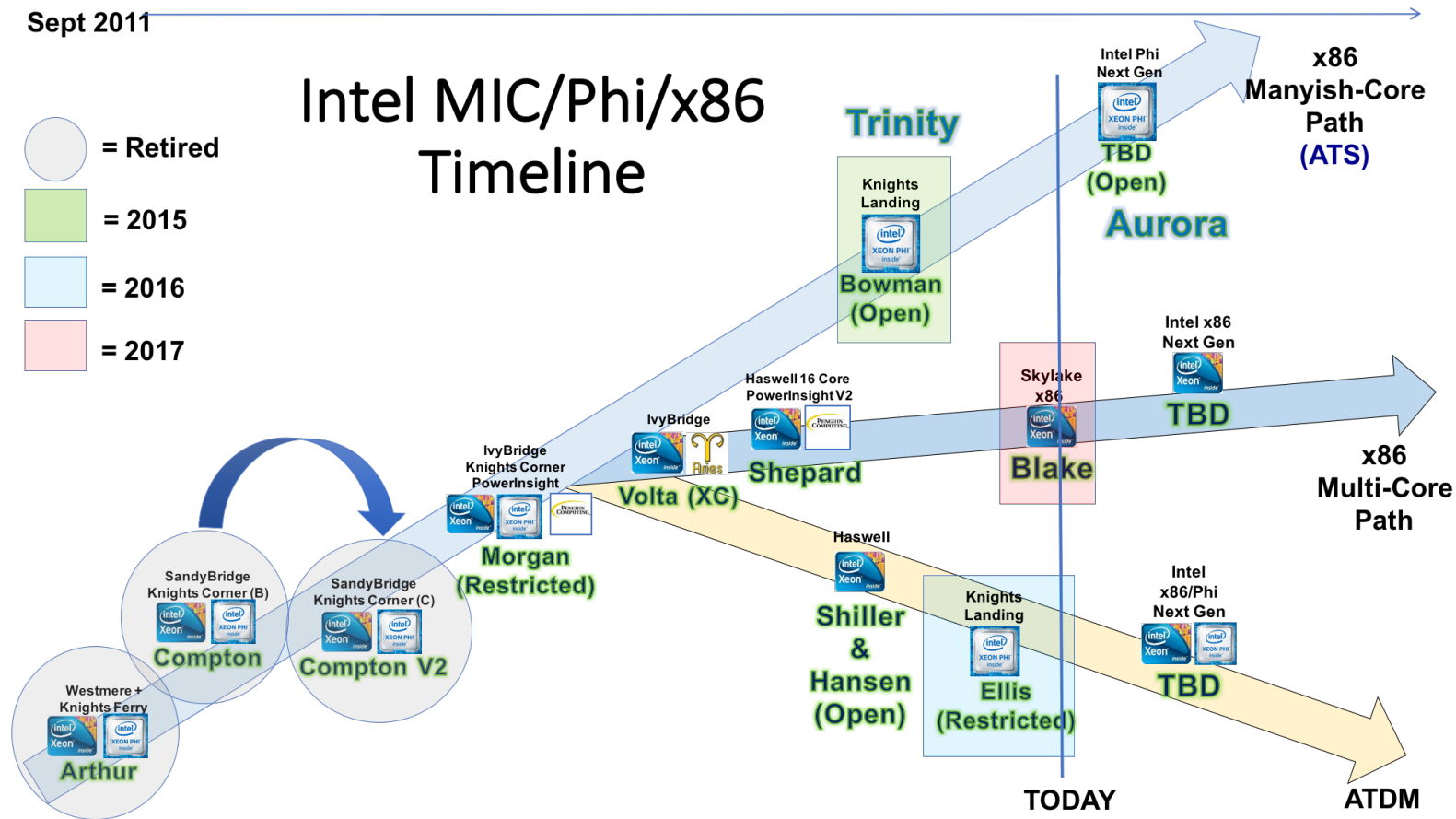
Disaggregated Memory



Advanced Architecture Testbeds

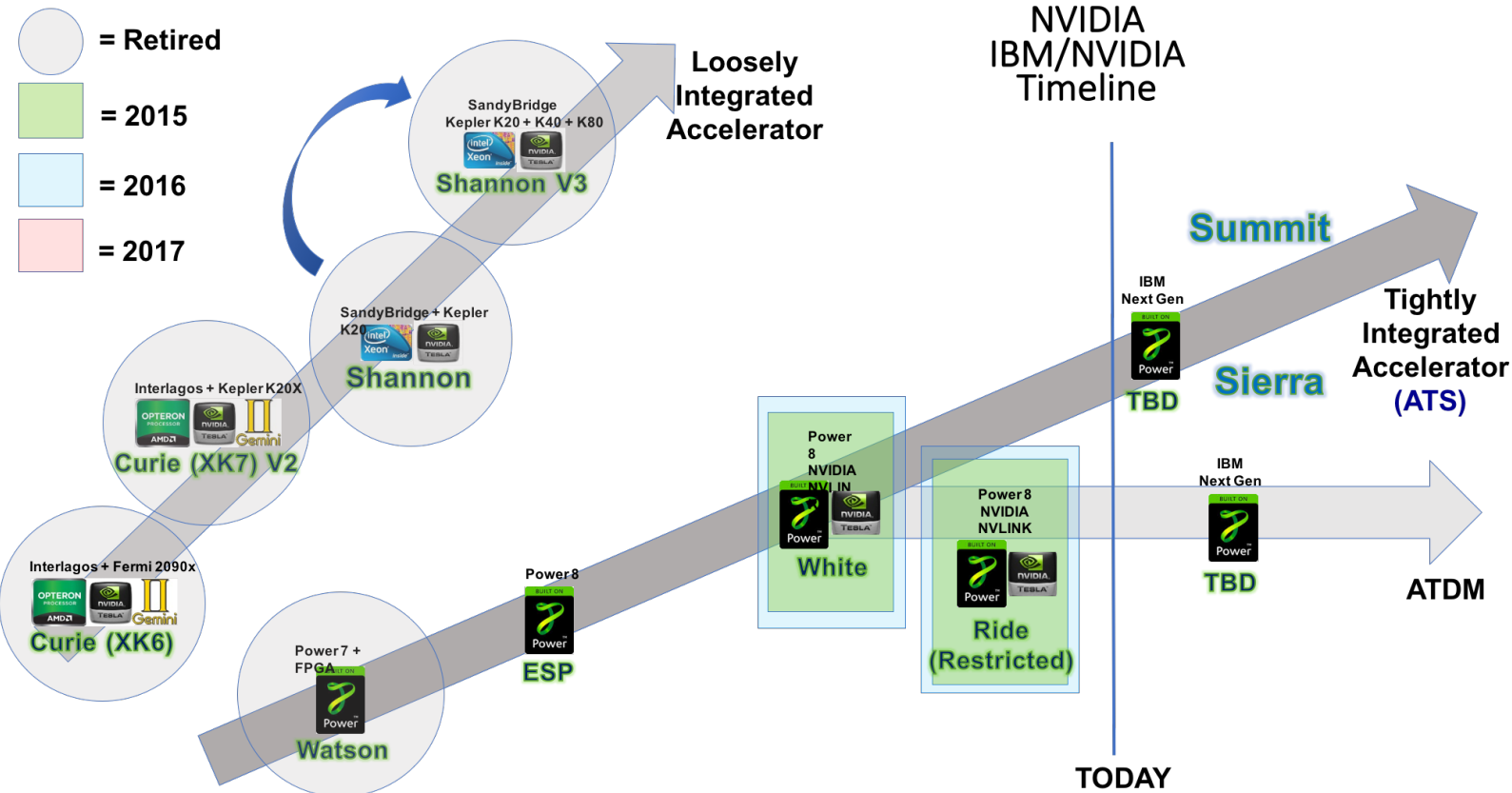
- Pre-production vendor systems
 - Intentionally closer to prototypes than production
 - Not intended for production application cycles
 - Priority is to explore wide and diverse set of emerging architectural alternatives
 - Rack-scale
- Facilitates co-design
 - Can potentially identify enhancements or influence configurations of follow-on components
- Helps prioritize vendor efforts
- Reduces application porting time on production systems

Intel Testbed Timeline



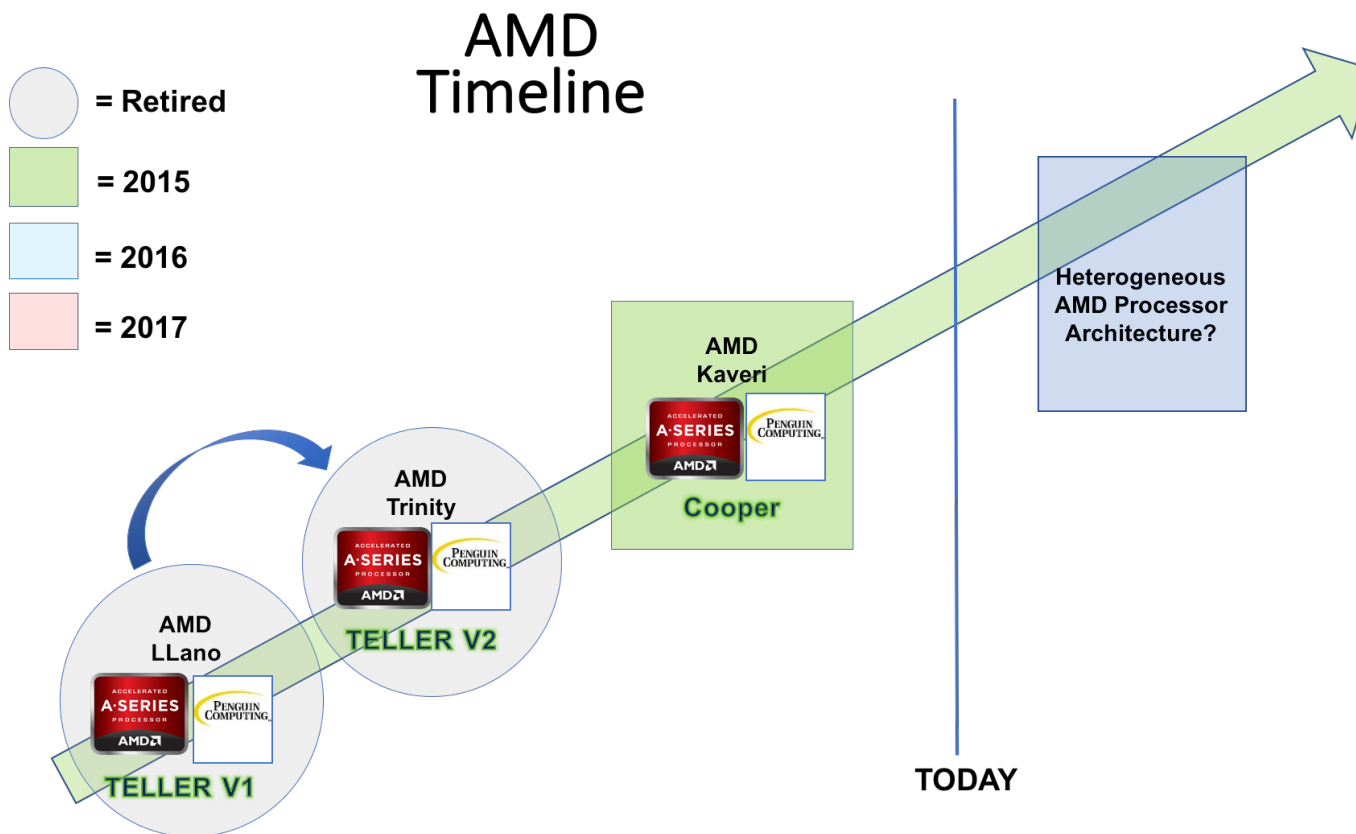
IBM Testbed Timeline

Sept 2011



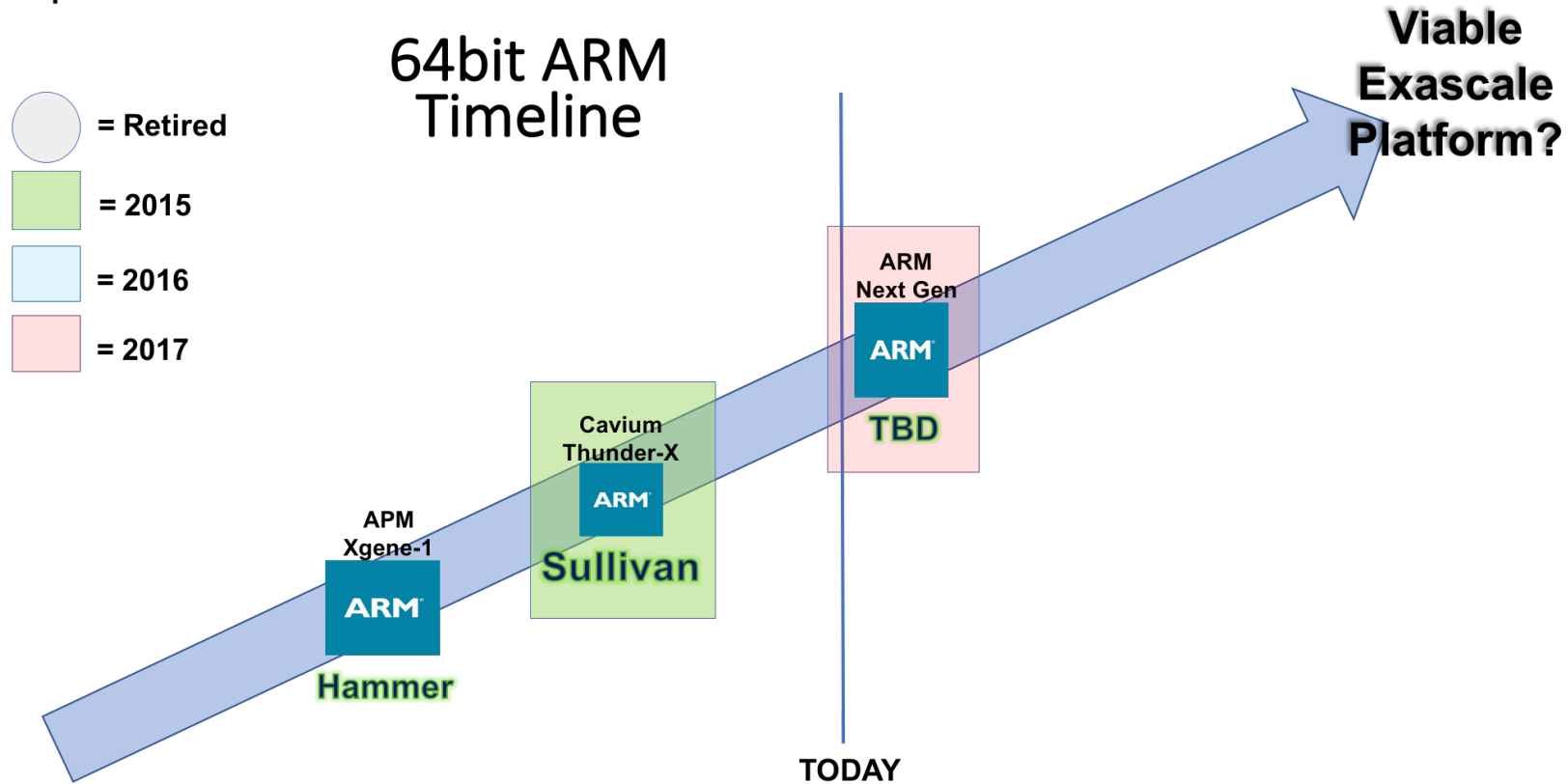
AMD Testbed Timeline

Sept 2011



ARM Testbed Timeline

Sept 2011



Large-Scale ARM Platform at Sandia

- 10-15 PFLOPS
- Key Requirements for this 2019 Platform
 - Proof of Concept for future leadership class DOE platform
 - Competitive HPC 64 Bit ARM processor technology
 - Integrated On-package Memory
 - First of a kind, Advanced Interconnect technology
 - Large focus on maturing the ARM software ecosystem
- Opportunities to *prototype* innovative test hardware
 - Logic in NIC to improve interconnection network performance electrical or optical
 - Logic in/near memory to support sparse linear algebra acceleration