

Linux Clusters Institute: Designing a Multipurpose Production Cluster

Irfan Elahi

National Center for Atmospheric Research

High-End Services Section Head

Program Manager (NWSC-2)

What is a Cluster

- **Compute Cluster:**
 - Simplest definition is that a **compute cluster** is a collection of network attached **computers** that can communicate with each other to execute a parallel application or workflow.
- **Parallel Application or parallel computing:**
 - The execution of processes are carried out simultaneously, so that large problems can be divided into smaller ones, which can then be solved at the same time across multiple compute nodes in the cluster.

Types of Clusters

- Dedicated Clusters
 - Commodity HPC cluster
 - Specialized supercomputer
- Shared Cluster
 - HPC Cloud Computing
 - POD, AWS/EC2, COD, Azure, JetStream, etc
 - Grid Computing
 - TeraGrid/XSEDE
 - Condominium Cluster Computing
 - University of Wyoming's Mount Moran

Cluster Building Block

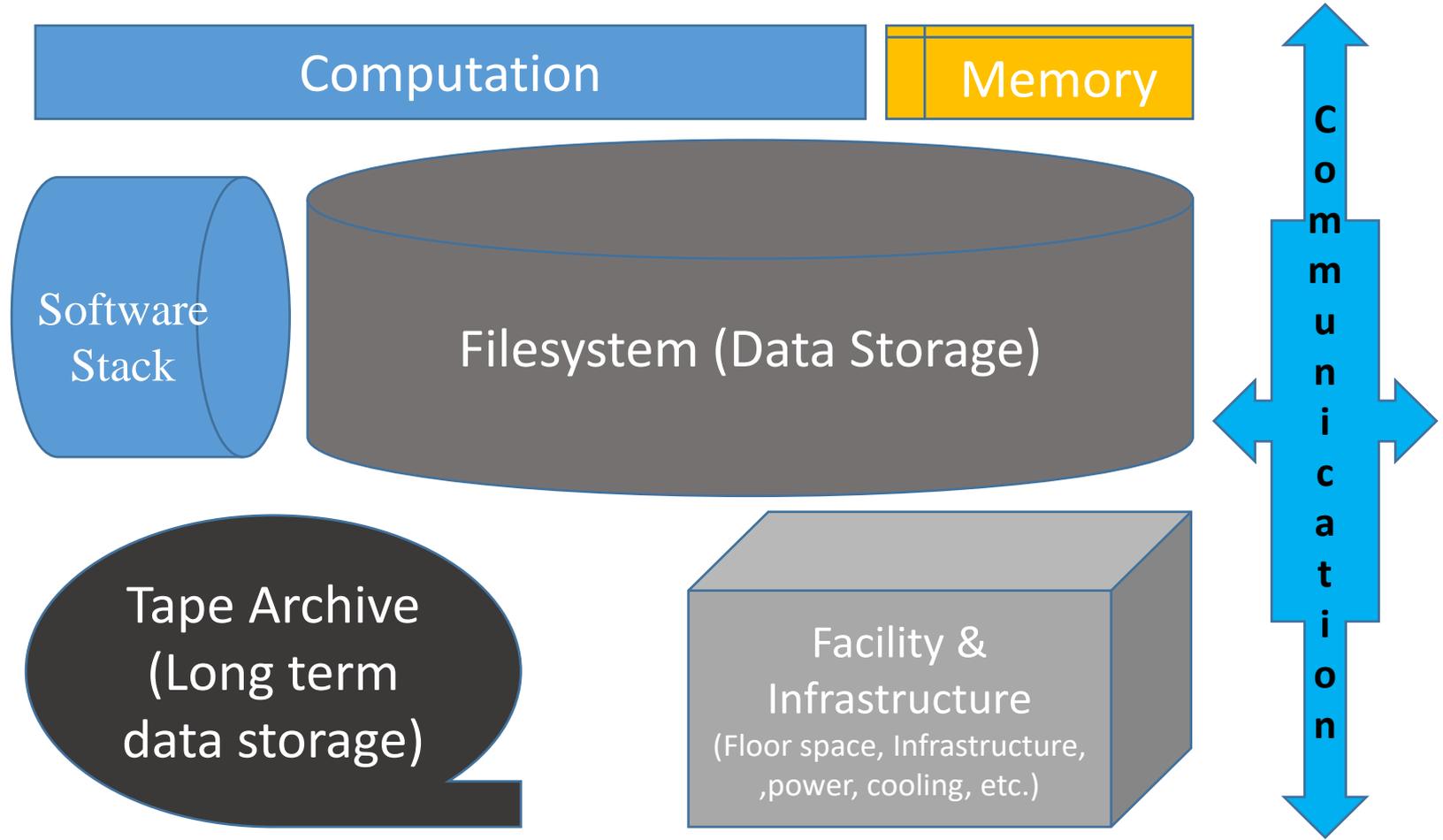
- Cluster Hardware
 - Nodes, Processors and memory
 - Interconnect
 - Storage
- Cluster Software
 - Cluster Management, Operating System
 - File System
 - Scheduler/Resource Manager
 - Development Tools
- Facility
 - Power, Cooling, Floor Space - Racking
- User/Application
- Support and Maintenance
- Cost of Ownership

User communities for NCAR HPC

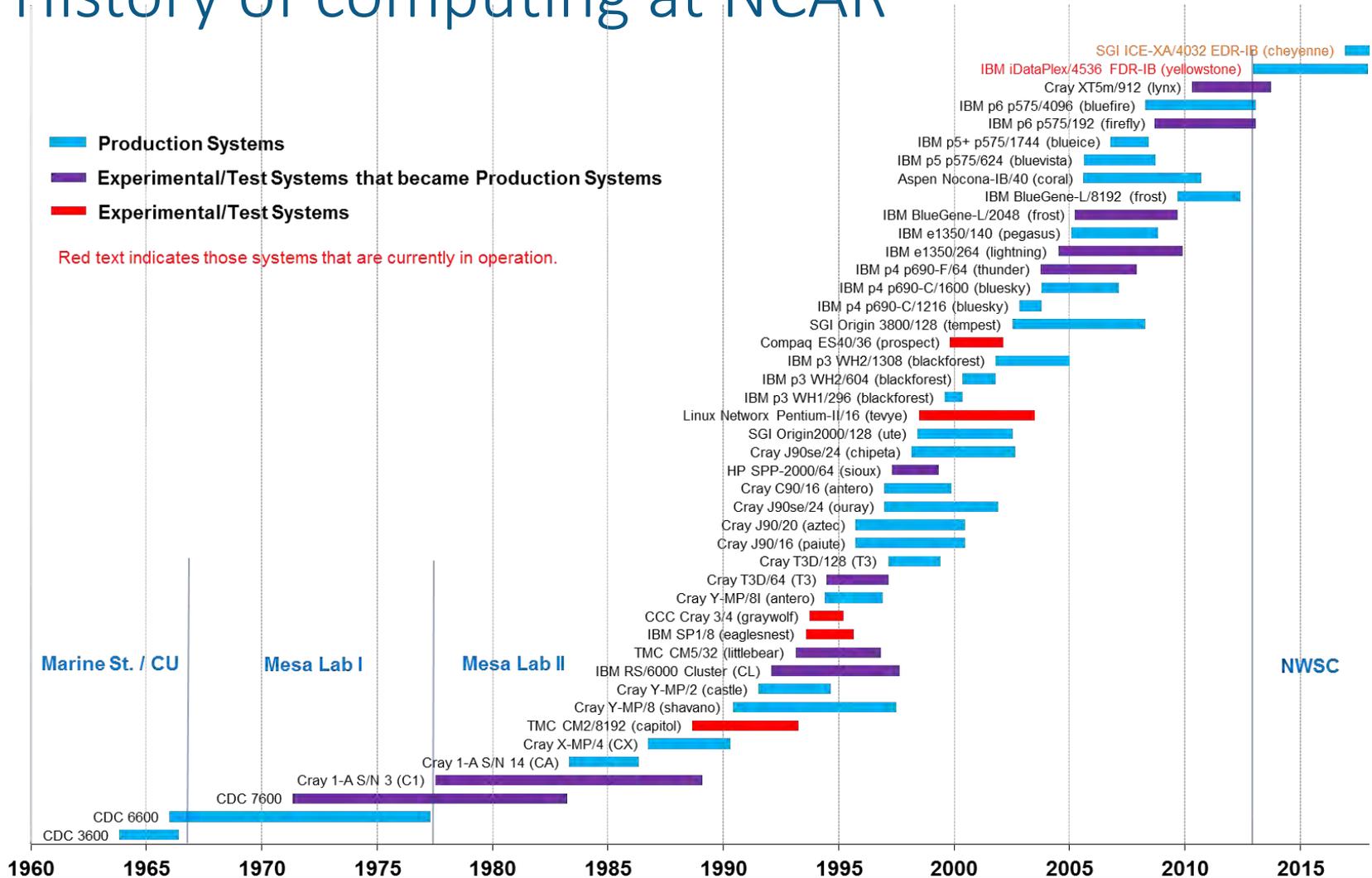
NCAR supports four user communities through policies established via various agreements with NSF or approved by NSF.

- University research
 - Eligibility: “In general, any U.S.-based researcher with an NSF award in the atmospheric, geospace or closely related sciences is eligible to apply for a University Community allocation.”
 - Large requests reviewed by CISL HPC Allocation Panel (CHAP)
 - Small allocations for projects with NSF awards usually processed in ~1 day
 - Small allocations, without NSF award, also available to graduate students, post-docs, or new faculty; also, instructional allocations (classroom, tutorial) in eligible field.
- Climate Simulation Laboratory
 - Supports large-scale, long-running climate simulations
 - Eligibility otherwise similar to University allocations
 - Also supports large annual allocation to CESM Community
 - Requests reviewed by CHAP
- NCAR Lab and NCAR Strategic Capability (NSC) activities
 - NCAR staff may engage collaborators from outside of NCAR
 - Large requests reviewed by internal NCAR Allocation Review Panel, approved by NCAR Executive Committee
- Wyoming-NCAR Alliance
 - Must be led by U Wyoming researcher
 - Must be in the Geosciences or related fields (including solid Earth geophysics)
 - Any source of funding support
 - Large requests reviewed by Wyoming Resource Allocation Panel (WRAP)

NCAR's Supercomputer Architecture

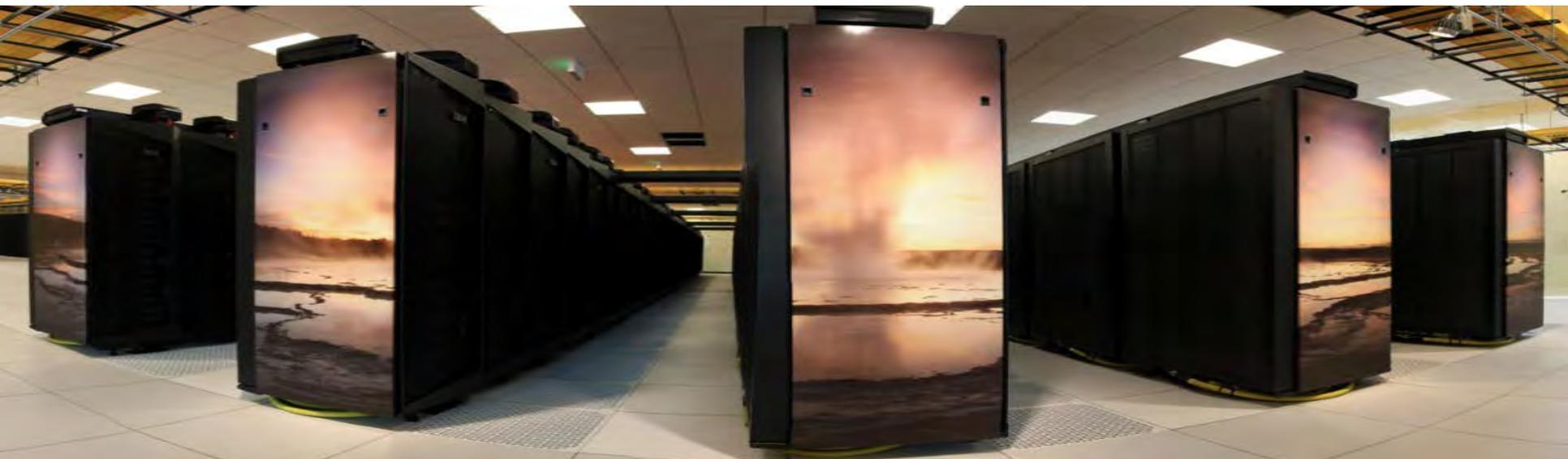


History of computing at NCAR





NCAR's Data-Centric
Supercomputing Environment.



Yellowstone (2012 -2017)

• Yellowstone Hardware

- 4,536 nodes, 72,576 cores total, 145.2 TB total memory & 1.51 PFLOPs Peak
- High-Performance Interconnect (Mellanox FDR InfiniBand full fat-tree)
- 6 Login/Interactive Nodes - Each 16 cores & 128 GB Memory
- 6 Service Nodes, 6 Boot Nodes, 3 Management Servers and 2 UFM Servers

• Yellowstone Software Stack

- Compilers, Libraries, Debugger & Performance Tools
 - Intel Cluster Studio (Fortran, C++, MPI libraries, trace collector & analyzer) - 50 users
 - Intel VTune Amplifier XE performance optimizer - 2 users
 - PGI CDK (Fortran, C, C++, pgdbg debugger, pgprof) - 50 users
 - PGI CDK GPU Version (Fortran, C, C++, pgdbg debugger, pgprof) - 2 users
 - PathScale EckoPath (Fortran C, C++, PathDB debugger) - 20 users
 - Rogue Wave TotalView debugger - 8192 floating tokens
 - Rogue Wave ThreadSpotter - 10 seats
 - Allinea DDT (debugger and profiler) - 1024 seats
- System Software
 - IBM Parallel Environment (POE) & IBM HPC Toolkit
 - IBM's LSF-HPC Batch Subsystem / Resource Manager
 - Red Hat Enterprise Linux (RHEL 6.4)
 - IBM General Parallel Filesystem (GPFS)
 - Mellanox Unified Fabric Manager
 - IBM xCAT cluster administration toolkit
 - XDMoD



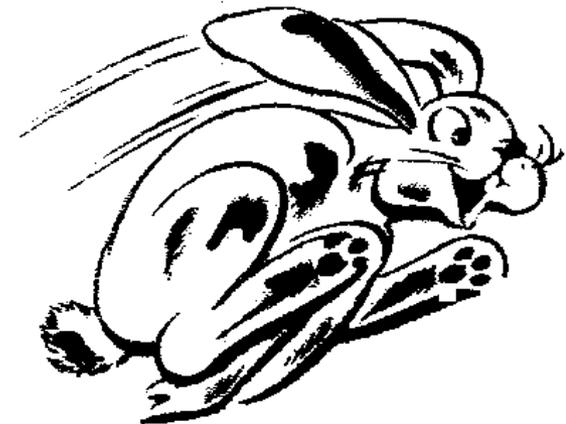
Cheyenne (2017-2021)

- Cheyenne Hardware

- 4,536 nodes, 72,576 cores total, 145.2 TB total memory & 1.51 PFLOPs Peak
- High-Performance Interconnect (Mellanox FDR InfiniBand full fat-tree)
- 6 Login/Interactive Nodes - Each 16 cores & 128 GB Memory
- 6 Service Nodes, 6 Boot Nodes, 3 Management Servers and 2 UFM Servers

- Cheyenne Software

- Compilers, Libraries, Debugger & Performance Tools
 - Intel Parallel Studio XE Cluster
 - Fortran, C++, performance & MPI libraries, trace collector & analyzer
 - Intel VTune Amplifier XE performance optimizer
 - PGI CDK (Fortran, C, C++, pgdbg debugger, pgprof)
 - Allinea Forge and Performance Reports
 - SGI Message Passing Toolkit (MPT)
- System Software
 - Altair PBS Pro Batch Subsystem / Resource Manager
 - SuSE Linux (Operating System)
 - IBM Spectrum Scale Parallel File System software (GPFS)
 - Mellanox Unified Fabric Manager
 - SGI Management Center (Cluster administration)
 - SGI Foundation Software (tools/utilities)



DAV Systems: Geyser, Caldera & Pronghorn

- **Geyser: Large Memory System**
 - 16 IBM x3850 X5 nodes; Intel Westmere-EX processors
 - 40 2.4 GHz cores, 1 TB memory per node
 - 1 nVIDIA Quadro K5000 GPU per node
 - Mellanox FDR full fat-tree interconnect
- **Caldera: GPU-Computation/Visualization System**
 - 16 IBM dx360 M4 nodes; Intel Xeon E5-2670
 - 16 2.6 GHz cores, 64 GB memory per node
 - 2 x nVIDIA Tesla K20X GPUs per node
 - Mellanox FDR full fat-tree interconnect
- **Pronghorn: Started as a KNC Cluster now used for share queue**
 - Similar configuration to Caldera
 - 2 x Intel KNC 6110 Xeon Phi per node (Decommisioned)

Erebus: Severity 1 Production Cluster

- Erebus (AMPS):
 - Dedicated to running operational weather forecasts for Antarctica to support the National Science Foundation's South Pole Station and its operations."
- IBM iDataPlex Compute Cluster
 - 84 IBM dx360 M4 Nodes; Intel Sandy Bridge EP processors
 - 16 cores, 32 GB memory per node; 2.6 GHz clock
 - 1,344 cores total - 28 TFLOPs peak
 - Mellanox FDR-10 InfiniBand full fat-tree
- Login Nodes
 - 2 IBM x3650 M4 Nodes; Intel Sandy Bridge EP processors
 - 16 cores & 128 GB memory per node
- Dedicated GPFS filesystem
 - 2 IBM x3650 M4 GPFS NSD servers
 - 57.6 TB usable disk storage

GLADE—GLOBally Accessible Data Environment

Glade Tier 1 (2012-2018)

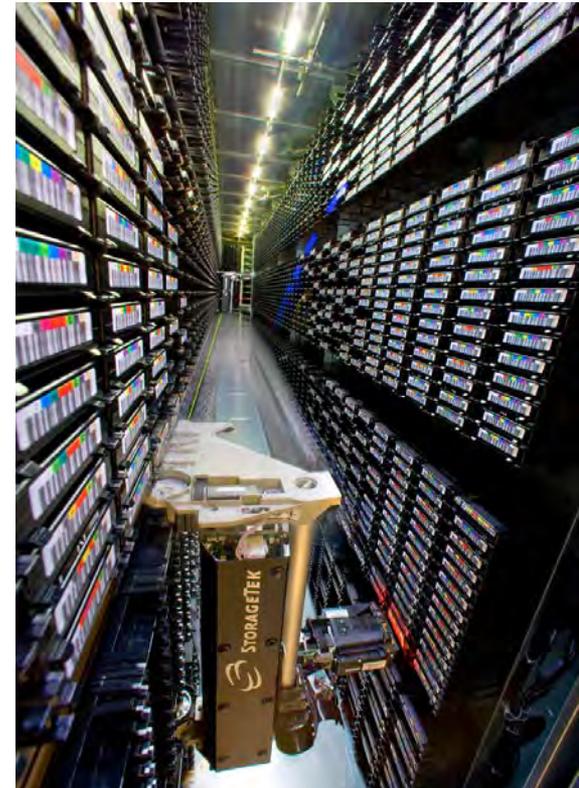
- 90 GB/s bandwidth
- 16 PB useable capacity
- 76 IBM DCS3700
- FDR & 10 GbE
- 6840 3TB drives
 - shared data + metadata
- 20 NSD servers
- 6 management nodes

Glade Tier 2 (2017-2023)

- 220 (312) GB/s bandwidth
- 40 PB useable capacity
- 8 DDN SFA14KXE
- EDR & 40 GbE
- 6720 x 8TB drives
 - data only
- 96 x 800GB SSD
 - Metadata

NCAR HPSS archive resource

- Added two SL8500 libraries
- 2 week ATP completed
- HPSS/Archive resources
 - Four SL8500 robotic libraries
 - 40,000 cartridge/slot capacity
 - 46 T10000C tape drives
 - 5-TB “C” cartridges
 - 46 T10000D tape drives
 - 8-TB “D” cartridges
 - 320 PB capacity
 - Current total holdings: +65 PB
 - Yellowstone growth rate: 12-14 PB/yr
 - Projected Cheyenne growth rate: +30 PB/yr

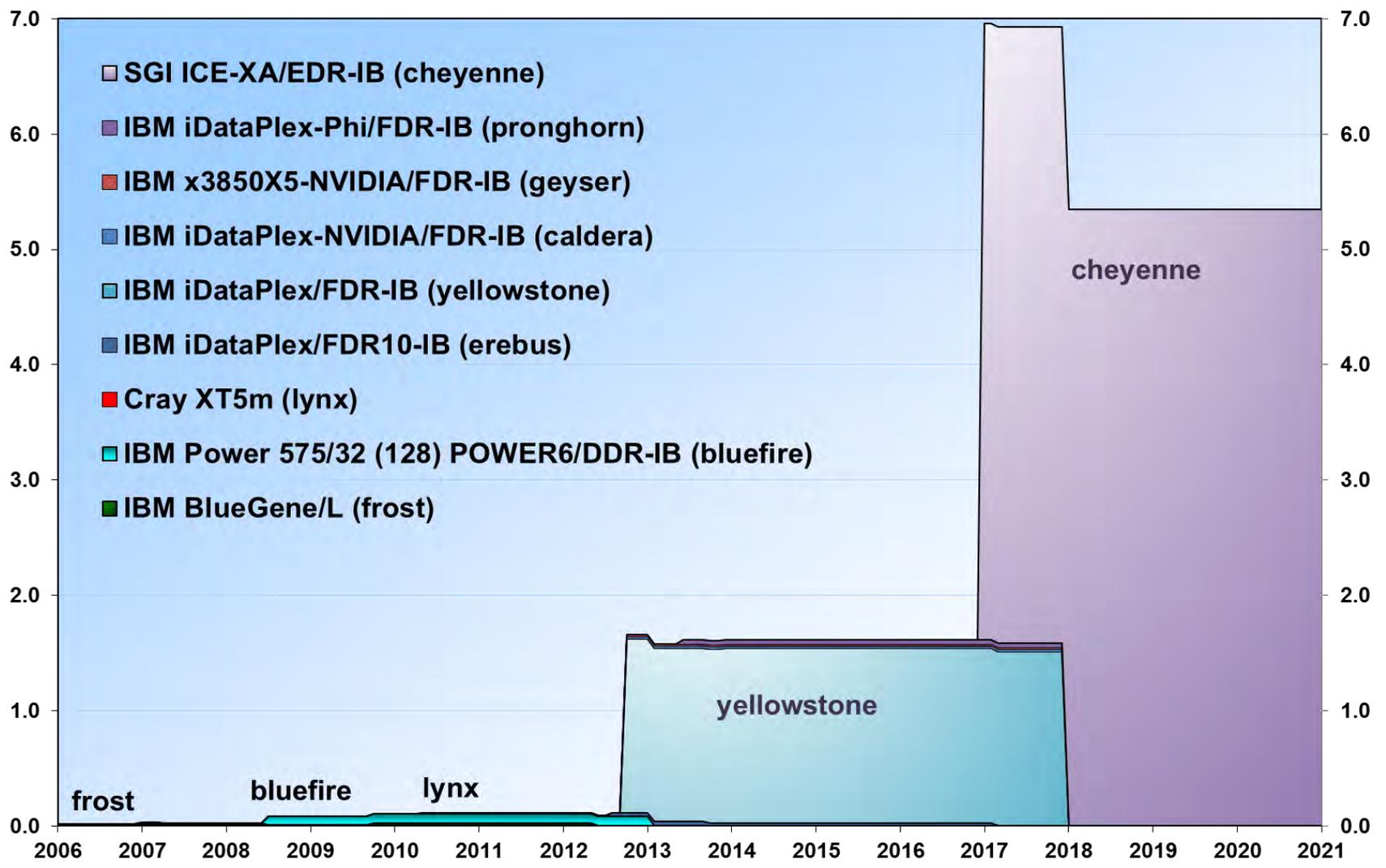


NCAR: Physical Infrastructure

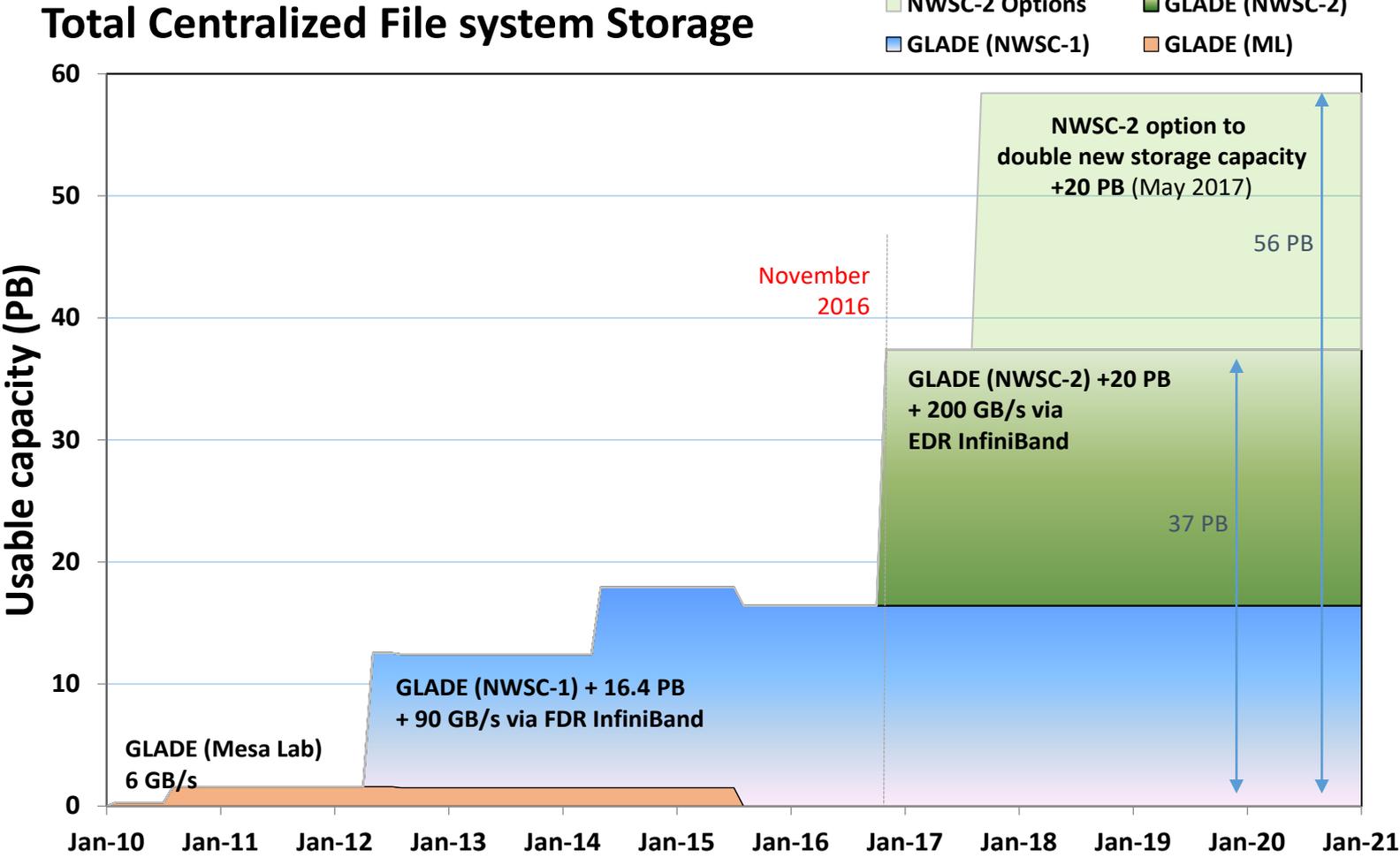
Resource	# Racks	
Cheyenne	14	SGI ICE XA E-Cells each containing 2 water-cooled E-Racks & heat exchanger, and 16 Mellanox 36-port EDR InfiniBand switches
	2	air-cooled storage & service racks including login nodes
GLADE+	8	DDN SFA14KX racks containing 24 NSD servers and storage
GLADE	19	NSD server, controller, and storage racks
	1	19" Rack (I/O aggregator nodes, management, switches)
DAV	1	IBM iDataPlex Rack (Caldera, Pronghorn)
	2	19" Racks (Geyser, management, IB switch)

Total Power	~2.0 MW
HPC	1.75 MW
GLADE	0.21 MW
DAV	0.063 MW

NCAR HPC profile, 2006–2021



NCAR disk storage capacity profile



Cheyenne: Power Efficiency

	Cheyenne 2017–2021	Yellowstone 2012–2017	Bluefire 2008–2013
Processor	Xeon E5-2697v4 2.3 GHz	Xeon E5-2670 2.6 GHz	POWER6 4.7 GHz
Total Batch Processor Cores	145,152	72,288	3,744
HPC portion peak PFLOPs	5.342	1.510	0.072
Power consumption	1.75 MW	1.4 MW	0.54 MW
Watts/peak GFLOPS	0.33	0.93	7.5
Peak GFLOPS/watt	3.05	1.08	0.133
Average workload floating point efficiency	1.1% (estimate)	1.56% (measured)	3.9% (measured)
Sustained MFLOPS/watt (on NCAR workload)	~34	~17	5.2
Bluefire-equivalents	70.8	28.9	1
Yellowstone-equivalents	>2.45	1	0.035

For 3.2x more power, Cheyenne delivers 71x more computational performance than Bluefire.

Let's Design a Cluster

What are your Goals?

- What is the allocated budget?
 - TCO analysis is your friend
- What are your user & application requirements?
 - Computational, post-processing, analysis, I/O (data), workflow, etc.
- What are the science drivers that will impact the requirements?
 - What are your new science drivers or new applications
 - Changes or optimization of existing applications
 - New campaigns or workflows in the horizon
- What is the projected workload and characteristics?
- What applications are ready, or expected to be able to take advantage of the emerging technologies
 - GPGPU, Many Core, ARM, NVMe over Fabric, Burst Buffer, ZCA, etc)?

What are your Goals & Capabilities?

- When do you want to deploy your cluster?
 - How long do you plan to keep the cluster?
- Dedicated vs Shared node resources.
- What are the anticipated service levels for these systems? (production, experimental, or experimental-transition-to-production)
- What are the overarching software, maintenance and support expectations? (vendor-supplied, mixture, rolls-your-owns)
- What racking/stacking & power/cooling capabilities does your data center has and/or can provide?
 - Make sure you include the cost in your TCO analysis.

What are your specs?

- What is your CPU/Memory requirements?
 - Good balance of CPU & RAM, etc
- What are the specific storage requirements?
 - What type of I/O: Sequential, random/bursty, MPI/IO, small/large files.
 - Spinning Storage vs SSD.
 - What is the I/O requirement of your applications: size/speed?
- What are the specific network requirements?
 - InfiniBand/OPA vs Ethernet
 - I/O vs MPI traffic, capability vs capacity jobs (average job sizes).
- What is your software stack requirement
 - Are you tied to a specific OS, Compiler, MPI, Development Tools, etc.
- What is your power/cooling and data center space requirements?
- What is your support and maintenance requirements?
 - What is the service Level Agreement you have with your user community

Other things to consider

- Cluster should be manageable by your technical and operations team.
- Availability & MTBSF requirements.
 - Minimizing single points of failure for robustness will cost you.
- Do you have any performance expectations from the new machine?
 - What benchmarks do you plan to use?
 - Benchmarking/profiling statistics on critical applications and workflows
 - Vendor optimization of your benchmarks
 - ATP (Acceptance Test Plan)
- Do you have staff who can maintain the hardware
 - Commodity replacement parts can do the job in some instances.
- Do you really want to spend money on purchasing proprietary software
 - Open source vs vendor supplied software
- Can you independently do the cluster software/firmware stack upgrade
 - Some vendors may want to charge you for this activity.
- Unintended consequences
 - Specifications/requirements, language, decisions, etc.

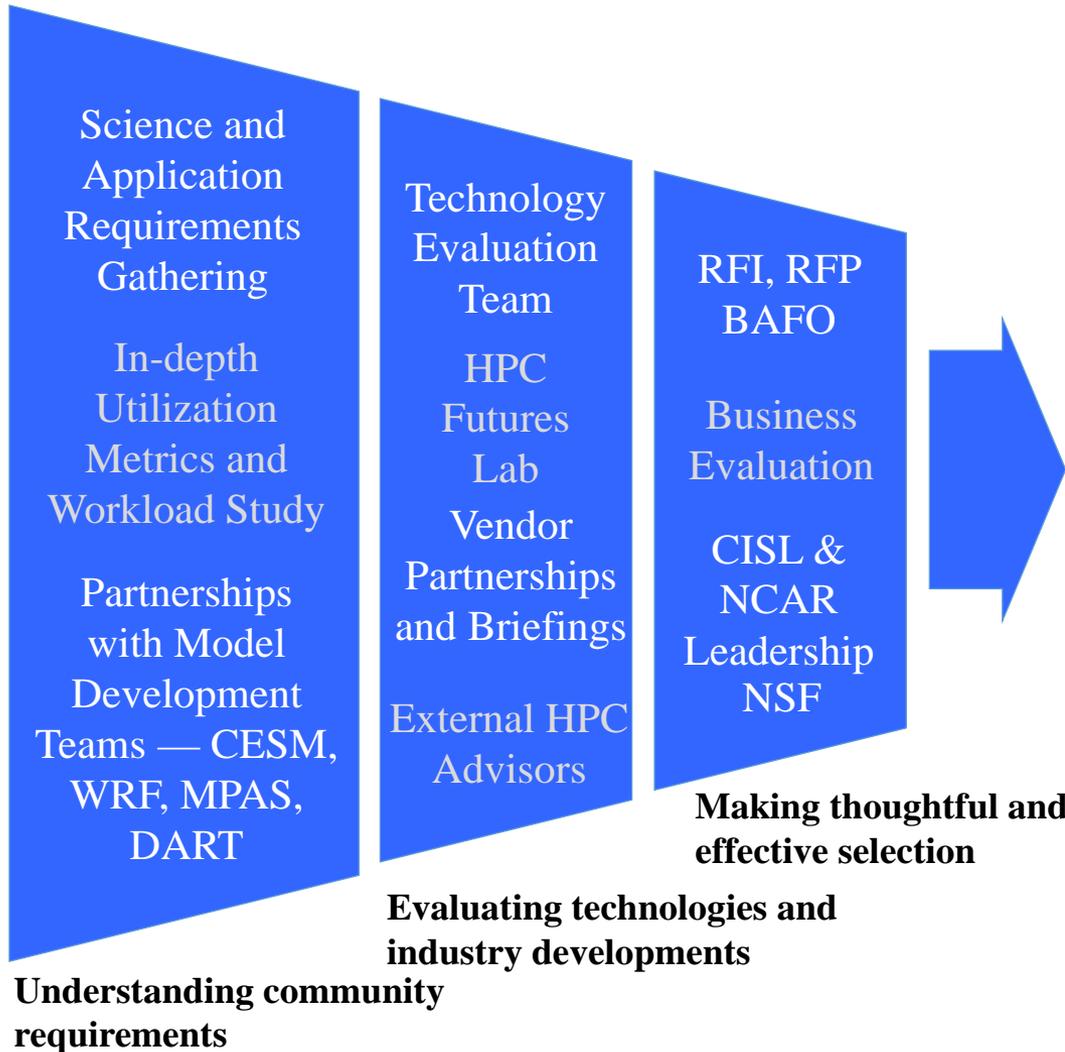
Utilization Metrics and Statistics

- Do you have utilization data from current/past Clusters
 - CPU Load
 - Memory Utilization
 - Network utilization
 - GPU utilization (GPGPU rather than just graphics)
 - Scheduler Statistics
 - Storage and file system utilization
- Storage and File System
 - Data produced by your applications. Types/size of files produced by your applications. Network bandwidth used by data traffic
 - Projected data size/type, latency & bandwidth requirements (Interpolation/extrapolation)
- Availability Statistics
 - Failure root causes (How to avoid them in future)
 - Gathas, issues, lessons learned from past systems

Where to begin?

- Do you know the current HPC trends and technology?
- Have you held vendor briefings, attended workshops, conferences and built a relationship with prospective vendors?
- Gather user requirements. Meet with top/heavy users, solicit feedback and information. Ask questions and be open.
- Know the boundaries and organizations procurement rules.
 - What/who will be providing the oversight?
 - Who will be supporting you on contracts and legal work?
 - Who is the approving/decision making authority?
- Have you thought of releasing an Request For Information (RFI) to the prospective vendors?
 - You can learn what is available in your time frame and what your money can buy. Budgetary estimates can help with sizing & requirements.
- Create a core “Technology Team”, a core “Business Team” and a core “Science/User Team” before you launch an RFP.

Yellowstone/Cheyenne Procurement



Yellowstone
NWSC-1 — 2013-2017



Cheyenne
NWSC-2 — 2017-2021

Sample User Survey Questions

Q1. The “top 20” users could be interviewed:

Q2. What do you currently use Cluster for or plan to use the new cluster for?

Q3. Please describe your current workflow patterns:

Which filesystems do you use?

Do you Archive your? If so...

- How much data do you plan to create? How much of that will be archived to tape?
- How often and how much data do you retrieve/read-back from Archive?
- Do you foresee any change in your workflow that will increase the data archive/read-backs on the new system?
- Which applications do you have that can run on a cluster? Characterize how you use these applications:
- What I/O patterns do these applications have?

What are your files like?

Number and size of application’s input files

Number and size of application’s output files

Intermediate and scratch file usage

What are the I/O characteristics of these applications?

many small I/O requests?

large-block I/O requests?

Fortran and/or C I/O?

MPI I/O?

POSIX (cp, cat, tar, etc.)?

Sample User Survey Questions . . .

Do you run these applications in serial or parallel?

Serial

Multi-threaded (e.g. OpenMP)? How many threads?

Multi-process (e.g., MPI)? How many processes?

Would you prefer the new system to have more analysis nodes with smaller memory or prefer having fewer analysis nodes with larger memory?

Do you use GPUs on the analysis nodes? If so...

Do you use GPUs for visualization? What applications?

Do you use GPUs for computation?

If "yes"...

Is the application single- or multi-node?

Do you use more than one GPU per node?

If "no"...

Have you considered optimizing your code for GPU/GPU Direct?

Would you consider it if the option was available?

Q4. How do you expect your workflow to change in the future? If so...

Do you expect just the same workflow, only more of it?

Will the character of your workflow change?

Do you expect to use different applications? Which ones?

Do you expect to use more nodes per application? How many?

Do you expect to use more GPUs per application? How many?

Sample User Survey Questions ...

Q5. What have you found frustrating or difficult in using a cluster?

For data processing and analysis?

For visualization?

For GPGPU computation?

Q6. Can your use of cluster be best characterized as interactive, batch, or a combination of both?

Q7. How would you best describe your I/O workload?

Is your workload “read-heavy” or “write-heavy”?

Does your workload exhibit substantial random I/O or sequential I/O?

Does your workload exhibit parallel or multi-threaded I/O?

What is the average I/O Request Size of your workload?

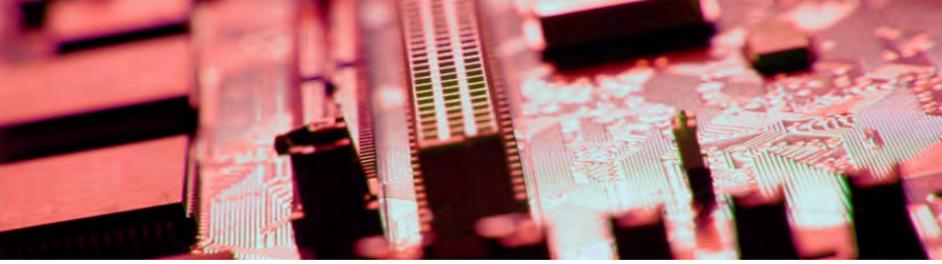
Is your workload I/O bound?

Do your applications and workflow demand higher I/O performance?

Q8. Do you believe that using Flash SSD could help you in your workflow?

Would node-local Flash SSD or node-shared flash/SSD better suit your workflow requirements? Or could you take advantage of both?

Would your workflow require data on flash/SSD to persist between jobs? If so, characterize how long such persistence would be required/desired by your workflow.



Cluster Hardware

Nodes/ Servers & Accelerators, Interconnect & Storage

Servers & Accelerators

- Compute Nodes
 - Accelerator Nodes
 - Accelerators
- Data Analysis, Post processing & visualizations Nodes
 - Fat Nodes (Large Memory)
 - GPGPU/GPU
- Support/Service nodes
 - Management Servers
 - Scheduler, Fabric Managers, Rack Leaders/Boot Nodes, etc
 - Login/Head Nodes
 - Storage Servers
 - File Transfer Gateways
 - Infrastructure (DNS, SSO, License Servers, etc)

Yellowstone/Cheyenne - Compute Nodes

Yellowstone

- Scientific computation nodes
 - IBM iDataPlex cluster
 - 4,536 dual-socket nodes
 - 16 cores, 2.6 GHz clock, Intel Xeon E5-2670 processors
 - 333 GFLOPs per node
 - 72,576 cores total
 - 145.2 TB total memory (32 GB per node)
 - 1.51 PFLOPs Peak

Cheyenne

- Scientific computation nodes
 - SGI ICE XA cluster
 - 4,032 dual-socket nodes
 - 18-core, 2.3-GHz Intel Xeon E5-2697v4 processors
 - 145,152 “Broadwell” cores total
 - 313 TB total memory (64-GB & 128-GB nodes)
 - 5.34 PFLOPs peak (1.325 TFLOPs per node)
 - *>3 Yellowstone equivalents on NCAR Benchmark Suite*

Processor Types

- Intel Xeon – X86 Architecture
 - **Xeon:** x86 microprocessors from Intel with multi-socket capabilities, higher core counts, larger cache memory, and support for ECC memory. Large number of the Top500 systems have Xeons.
 - **Intel Xeon Phi:** coprocessors, based on Intel® Many Integrated Core Architecture. Initially it was introduced as an add-on card (KNC), now it is a standalone CPU (KNL/KNH).
- IBM POWER (Also used by Bull, Hitachi, etc)
 - IBM has a series of microprocessors called **POWER** followed by a number designating generation.
- AMD
 - **Zen** is the codename for a computer processor from AMD. AMD launched it with their Ryzen series of CPUs in February 2017. Zen is based on a SoC design/implementation. Servers are expected to ship later this year.
- ARM
 - **ARM**, originally **Acorn RISC Machine**, later **Advanced RISC Machine**, is a reduced instruction set computing (RISC) architecture based processors.
- GPGPU
 - **GPGPU** is the use of a graphics processing unit or GPUs, to perform computational work that was traditionally handled by CPUs.
 - NVIDIA, AMD, etc

GPGPU/Many-Core/Accelerators

- Similar to Compute Nodes
- May have a different form factor
- Designed TDP and slots for accelerator cards
- Geyser has 16 Quadro 5000 NVIDIA cards
- Caldera has 32 NVIDIA K20x cards (2 per node)
- Rey has 8 K20x NVIDIA cards per node
 - PCI Switch allows GPU direct communication
- IBM POWER 8 system uses NVLINK for GPU direct
- Intel's Xeon Phi:
 - started as a card/adaptor (KNC) but now is a bootable host (KNL/KNH).

Miscellaneous Nodes

- Fat Node
 - Large Memory
 - Mostly used for data analysis and post processing
 - Shared queue nodes
 - Machine Learning/Deep Learning
 - Other uses
- Service/Support Nodes
 - Management Node(s)
 - Boot Nodes/Rack leaders
 - Scheduler/Resource Manager
 - Fabric Manager
 - DHCP, DNS, License servers, etc
- Login/Head nodes
 - Dedicate compile nodes
 - Nodes for interactive work

Diskless or not?

- Diskful/Stateful Nodes
 - Data Analysis and post processing may require local disk and larger system image, additional RPM, software stack, etc.
 - Nodes used for compilation
 - Login/Head nodes
 - Scheduler/RM
- Diskless/Stateless Nodes
 - Compute node images are small so you may save on cost of local disk on each compute node.
 - Cost of memory vs disk used for image
 - Provisioning stateless nodes vs stateful node

Virtualization or not?

- It depends
 - Application performance & resource requirement
 - Dedicated vs shared resource use
 - Network congestion
 - Software licensing
 - Private or Hybrid Cloud
 - Resource containment
- Container or VM
 - Container may be good to deliver as a series of stateless microservices while virtual machines may be good for traditional monolithic models.
 - Containers: Software Virtualization vs VMs: Hardware Virtualization
 - Containers started as Jails, chroot to segregate namespaces in a Linux operating system for security purposes.
 - Containers are isolated but can share OS, Libraries, Binaries where appropriate.

Storage

- How much Storage do you need
- What type of data protections do you require
- What is the I/O performance requirements Size of your users/applications
- What should the layout of the storage (User Visible)
- Home Directories
- Scratch/Work Space
 - Slow scratch
 - Fast scratch
- Project/Campaign Storage
- High IOPS Storage, Burst Buffer, etc
- Backups &/or Archive

Home vs Scratch

- Home:
 - Intended for application binaries, development space, small input files, etc – not where the bulk data should reside.
 - Governed by quota limits
 - User expectation is that it must be backed up regularly
- Scratch:
 - Intended for large data files, typically outputs of computational runs, etc.
 - Could be tiered into slow & fast disk pools based on requirements.
 - Could also have quota limits, and may be oversubscribed.
 - Generally users do not expect the contents to be backed up.
 - Could be under a purge policy (60, 90, 120, nnn days residency)

Storage Availability & Reliability

- If the storage/file system is offline, the cluster is not usable.
- Storage reliability is crucial to productivity in HPC.
- For enhanced availability consider:
 - Ask vendor for 99% Availability as a requirement
 - Eliminate single point of failures
 - Require high availability, redundancy &/or failover capability for drives, controllers, and other hardware components.
 - Full data path protection
 - What type of RAID do you require & can afford.
- For enhanced reliability consider:
 - Protection against disk drive data loss.
 - Hot spares, sector-based checksum, sniffing of disk media, error correction capabilities, ability to save write-cached data upon power failure, etc.
- Consider a robust and proven file system.
 - Some file system have added availability and reliability features in addition to what the storage subsystem provides.

Storage Performance

- You get what you pay for.
- What is the I/O requirements
 - What I/O patterns do your applications have?
 - many small I/O requests?
 - large-block I/O requests?
 - Fortran and/or C I/O?
 - MPI I/O?
 - POSIX (cp, cat, tar, etc.)?
 - What are your files like?
 - Number and size of application's input files
 - Number and size of application's output files
 - Intermediate and scratch file usage

Backups/Archive

- Do you have a requirement to backup and/or archive user data?
 - Backup/Archive can get out of hand so be careful of what you wish for.
 - If you cannot recreate the data than must backup/archive it.
 - Does your organization have a mandate to backup/archive research data?
- NCAR has Oracle SL8500 tape libraries
 - 46 T10000D Drives. 320 PB capacity & current holding of over 65PB
 - Yellowstone produces 12-15 PB a year for archival
 - Cheyenne is projected to produce 30PB a year for archival
- Software: NCAR uses HPSS. You can use open source SW.
- Do not try to backup or archive scratch.
- Users should be involved and know what/why they are saving.
 - Make sure you have a lifecycle management policy for all your data.

The Purge: Scratch

- Usually scratch space is oversubscribed.
 - Some clusters administrators prefer to have quota for scratch while other may not.
- Any file in a scratch partition over 60 days old is “deleted.”
- Send periodic and timely notices, to all users who have scratch files, about what those files’ fate will be, including the timeline.
- This should applies to all of scratch.
- Without “The Purge”, scratch would become unmanageable.

Cluster Networks

- Infiniband
 - Very high throughput and very low latency.
 - FDR (4x) : 56Gb/s & 0.7 μ s, EDR (4x): 100Gb/s & 0.5 μ s
 - HDR (4x) is projected to be 200Gb/s & 90ns
- Ethernet
 - 10GE, 10GigE, 10GbE or 10 Gb/s (Latency is in range 4-5 μ s)
 - 40 GigE (& 100 GigE)
 - Switch port-to-port latency, 230ns. RDMA application latency of 1.3 μ s.
Sockets application latency of 4 μ s.
- OPA - Omni-Path Architecture (Latency 100-110 ns)
 - High-performance communication architecture from Intel delivering speeds up to 100 Gbit/s.
- Cray's XC Series Network (Aries)
 - Only used in Cray's XC Series system
 - Bandwidth: 4.7-5.25 GB/s, Latency: 2 us)
 - Uses Dragonfly Topology

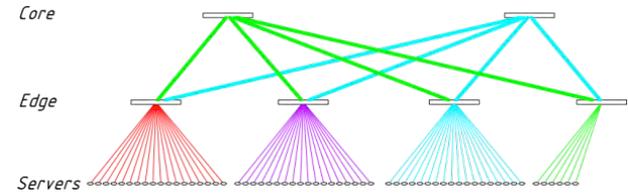
Why Both Infiniband & Ethernet?

- InfiniBand: very low latency (~.5 microseconds)
 - Good for many very small messages, which is how most parallel computing works.
- Ethernet: higher latency
 - Cost: Much cheaper than InfiniBand.
 - Separation: Keeps the I/O traffic separate from the message passing traffic – it's hard to optimize a network for both. You can opt for dual rail InfiniBand but that will be very costly.
 - Failover: If the InfiniBand fails, then the message passing can go over Ethernet, though substantially slower – a slowdown is better than not getting any work done at all.
 - Cluster management is mostly over Ethernet. Though most clusters have a dedicated management network.

High-Performance Interconnect

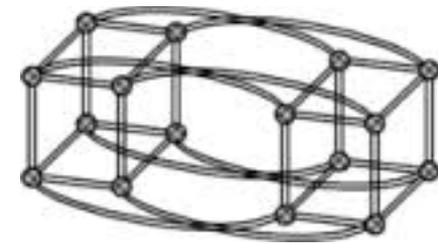
- Fat Tree

- Full Fat Tree (Yellowstone at NCAR)
- Island Topology (SuperMUC at LRZ)
- Oversubscribed (Summit Supercomputer at UCB/CSU 2:1)
- Clos
- (<http://www.mellanox.com/clusterconfig/>)



- HyperCube (Connection Machines)

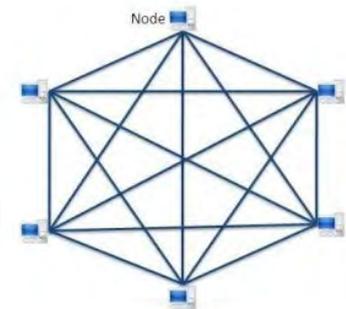
- SGI ICE XA used HyperCube (Cheyenne at NCAR)



- "Dragonfly" topology

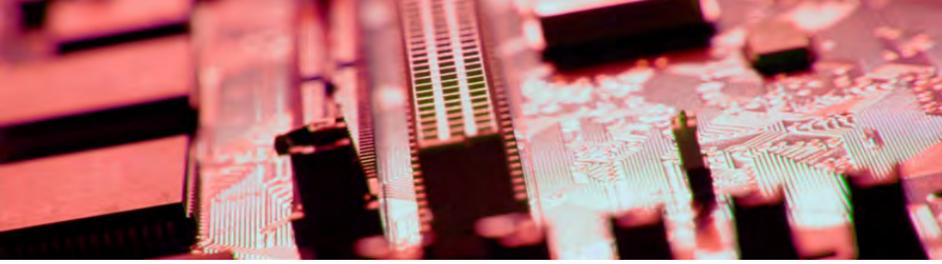
- Mesh Topology

- Every node relays data for the network.
- Cray Gemini network is a good example of a mesh (3D Torus)
- Fujitsu has a 6D Torus



Why Oversubscribe?

- InfiniBand is expensive.
- Each component is expensive:
 - Core switch
 - Leaf switches
 - Cables (copper is expensive, fiber is very expensive)
 - Cards
- Issues governed by oversubscription rate:
 - Bandwidth (may not be important to your applications)
 - Injection rate or Latency (messages per second – may be very important)
- In some cases oversubscription provides a good balance of cost vs performance.
- Know your application and user requirements before you spend the money.
- Yellowstone is Full Fat-Tree.
 - IB utilization is below what was projected.
 - We ran benchmarks by converting Yellowstone to a 2:1 oversubscribed fat tree and saw no degradation in application performance. This changed our requirement for the next generation and that is why Cheyenne is a hypercube.



Cluster Software

Operating System, development tools, cluster and network management software, etc

Software Stack: OS & Cluster Mgmt

- Operating System
 - RHEL, SuSE, AIX, CLE (Cray Linux Environment),
 - Centos, OpenSUSE, Debian, etc
- Cluster/System Management
 - Cray: System Management Support, SGI: SMC, HPE: CMU, xCAT, SaltStack, etc
 - xCAT, Rocks Cluster Distribution, Kubernetes, Stacki, etc
- Scheduler/Resource Manager
 - IBM LSF, Altair PBS, Adaptive MOAB, etc
 - OpenPBS, SLURM, Cobalt, etc
- Compiler
 - Intel Compilers, PGI, PathScale, CAPS, OpenACC, etc
 - GNU, Watcom, Oracle Developer Studio, etc
- Development Tools/Libraries
 - Intel MPI, MPT, PE, etc
 - OpenMPI, MVAPICH, etc

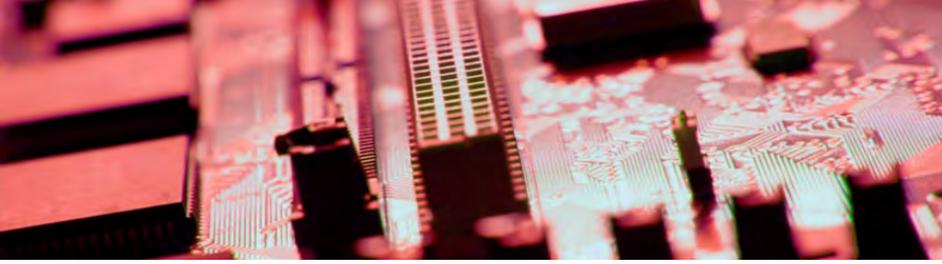
To Learn More About Yellowstone,
Cheyenne and Glade

<http://www.ucar.edu/>

<https://www2.cisl.ucar.edu/resources/resources-overview>

Acknowledgements

- Henry Neeman, University of Oklahoma for his 2015 LCI Slides.
- Erik Scott (Harris) , Jared David Baker (University of Wyoming) , Pamela Hill (NCAR), Shilo Hall (NCAR), Nathan Rini (NCAR), Ben Matthews (NCAR), Jon Roberts (NCAR), Thomas Engel (NCAR), , Jeffrey R. Lang (University of Wyoming) , Jonathan Anderson (University Colorado, Boulder), Brian Dale Haymore (University of Utah), Leslie Ann Froeschl (University of Illinois) ,Tim Brewer (University of Wyoming), Stormy Knight (NCAR), Robert McLay (TACC) for reviewing and providing feedback on the slides.
- LCI for the opportunity.



Thanks for your attention!

Questions?