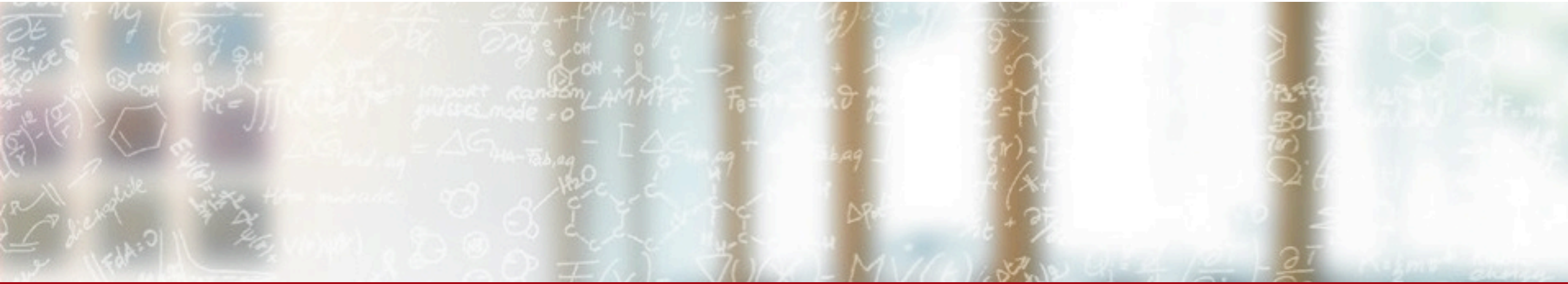




CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich



Architecture and Design of CSCS experimental testbed (Ault)

Klein, CSCS

March 25, 2019

High Level Overview

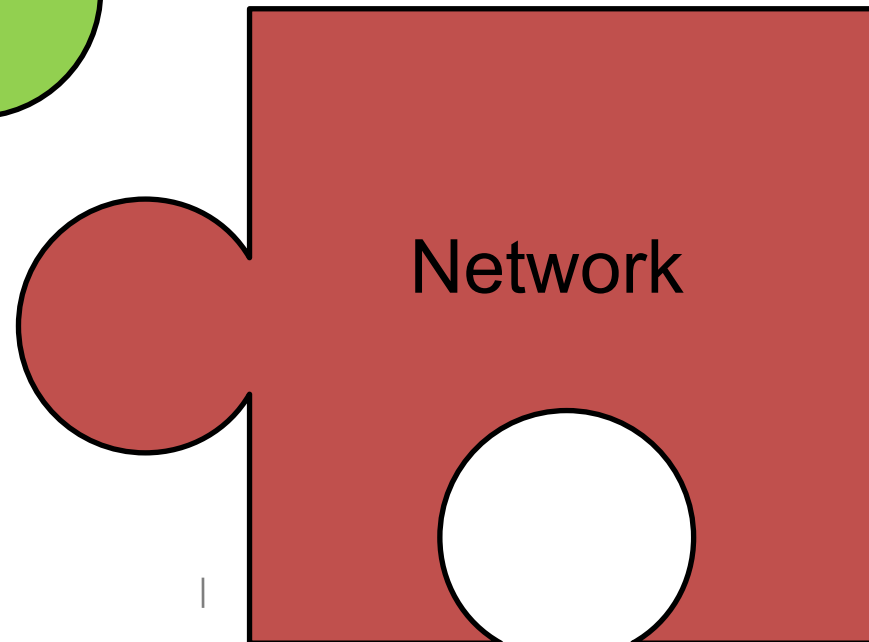
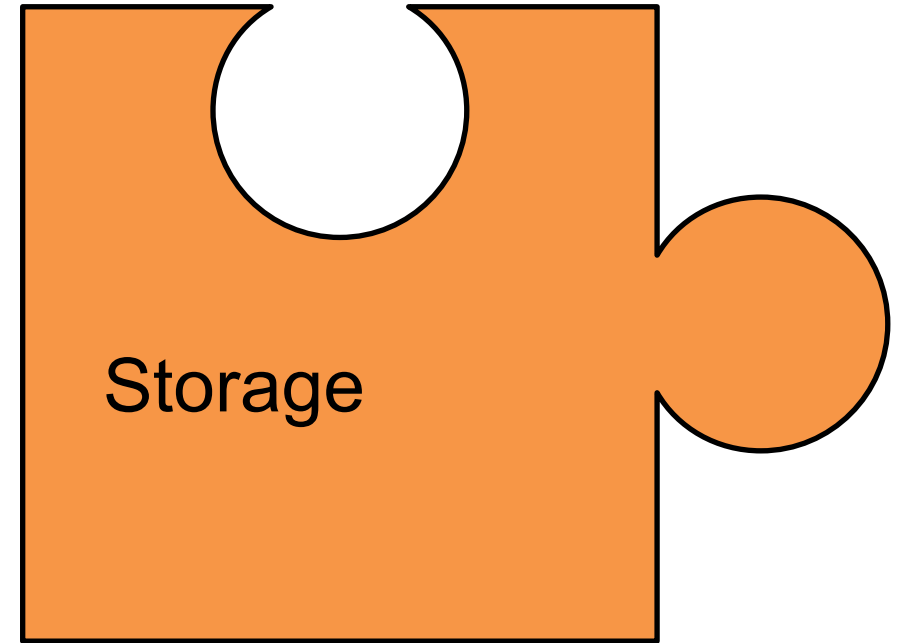
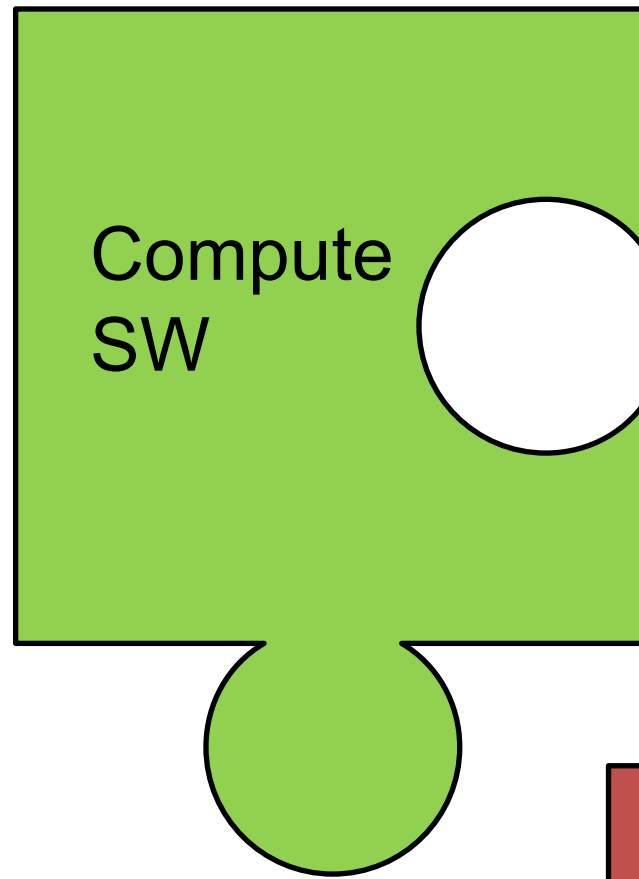
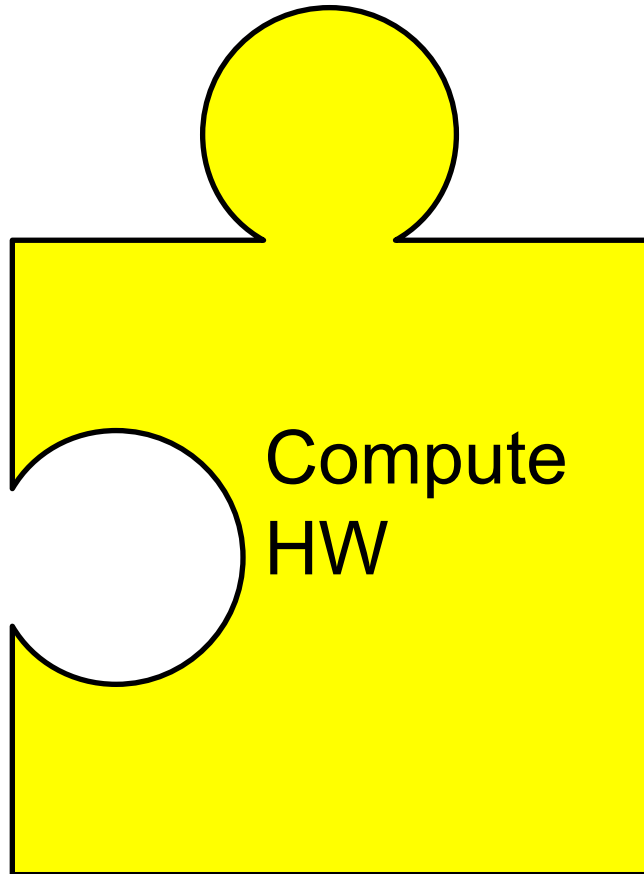
Ault Goals

- Testbed platform arose from ADAC discussion
 - Allow collaborators access to hardware they might not have available
 - Shared between ADAC sites.
 - Reciprocal: CSCS access to ORNL ExCL 2.0 for example
- Difference between a testbed platform (Ault) and a testbed service (Greina).
 - Create a testbed platform that can also support a traditional TDS environment
- Ault is designed to allow for testing beyond what is possible easily on a shared resource.

Ault Goals

- Different levels of access for different research requirements.
 - Benchmarking/TDS
 - Greina-Like Service (Historical CSCS TDS testbed Service)
 - Privileged access (Low-level debugging/tweaking/etc)
 - Easily allow for approved users to modify OS, Kernel, BIOS settings on bare-metal servers
 - CLI access to BIOS configs
 - Can prototype and test designs of systems quickly and easily.
- Local storage provisioned on demand
- Provisioning/commissioning system verifies hardware and cleans up when done.

Assemble your test



Ault Testbed (Metal as a a Service)

- This Use-case mode still being iterated
 - not ready for complete self-directed usage yet
 - Interested early testing users please contact me
- Fully Reserve Hardware
 - Removed from shared service
 - Loses access to shared /scratch storage
 - Isolated from other resources
 - Can build up own network of multiple nodes
 - Cloud-init compatible images
 - CentOS, Ubuntu, SLES, and Windows of various versions already in library

Ault Testbed (Metal as a Service)

- On release of reservation node is:
 - cleaned up
 - Bios reset, firmware re-flashed.
 - re-provisioned back to the shared service environment
 - Automated as much as possible
 - Care has been taken to acquire hardware supporting this mode.
 - Those nodes that can't be automatically cleaned and validated, will be unable to used for dedicated testing.
- MaaS.io for documentation.

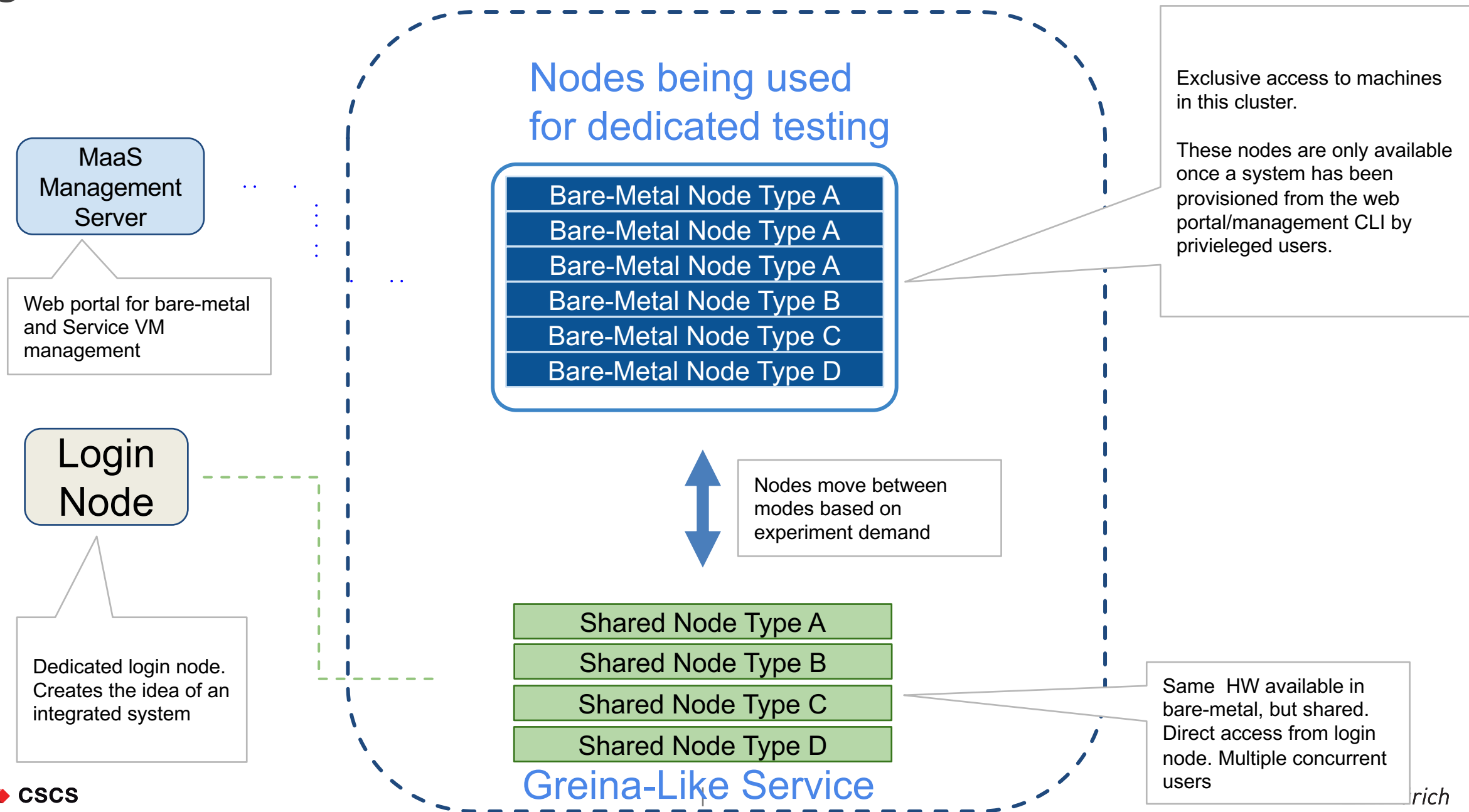
Early Success tested already

- MCH Cosmo Benchmarks
 - Frequency Scaling policies across different CPU SKUs
 - Comparing OS governors and BIOS settings
- Compatibility of OS/Application Library (SLES15/Old-Cuda Apps)
 - And Cuda 10+ compat layers going forward
- On-Demand BeeGFS Container Prototyping using Intel Rulers

Early Success: Greina as a Service!

- A shared TDS built on top of testbed platform
- Default state of hardware
 - SSH access to login node
 - submit jobs through Slurm to run development and benchmark runs
 - Container Workloads through Singularity
- Useful to test new HW arch, develop and benchmark codes, and do things that users can do.
- Basically for most users, no change to workflows, should support everything previously by the greina testbed service
 - Early testing has been very helpful, but still some issues to be worked out





High Level Overview



Access to Resources

- When systems are reserved for dedicated testing, they will appear Down in Slurm of the shared resource testbed.
- Currently dedicated system allocation is manual process.
 - Request opened
 - Operator assign a resource to a user
 - Operator redeploys shared environment when done
- Investigating better methods (using APIs)

13 Machines 1 Resource pool

Filters		Search										
<input type="checkbox"/>	FQDN ▼ MAC IP	POWER	STATUS	OWNER, TAGS	POOL	ZONE	FABRIC, VLAN	CORES	RAM	DISKS	STORAGE	
<input type="checkbox"/>	 ault.cscs.ch 148.187.104.126 (PXE)	 On Virsh	CentOS 7	kleinm virtual	default	default	fabric-0 Default VLAN	16	64 GiB	1	40 GB	
<input type="checkbox"/>	ault01.cscs.ch 148.187.104.72 (PXE)	 On Ipmitool	ault	kleinm skylake, intel	default	default	fabric-0 Default VLAN	72	384 GiB	1	240 GB	
<input type="checkbox"/>	ault02.cscs.ch 148.187.104.73 (PXE)	 On Ipmitool	ault	kleinm skylake, intel	default	default	fabric-0 Default VLAN	72	384 GiB	1	240 GB	
<input type="checkbox"/>	ault03.cscs.ch 148.187.104.74 (PXE)	 On Ipmitool	ault	kleinm skylake, intel	default	default	fabric-0 Default VLAN	72	384 GiB	1	240 GB	
<input type="checkbox"/>	ault04.cscs.ch 148.187.104.75 (PXE)	 On Ipmitool	ault	kleinm skylake, intel	default	default	fabric-0 Default VLAN	72	384 GiB	1	240 GB	
<input type="checkbox"/>	ault05.cscs.ch 148.187.104.76 (PXE)	 On Ipmitool	ault	kleinm skylake, gpu, nvi...	default	default	fabric-0 Default VLAN	36	768 GiB	2	480 GB	
<input type="checkbox"/>	ault06.cscs.ch 148.187.104.77 (PXE)	 On Ipmitool	ault	kleinm skylake, nvidia, g...	default	default	fabric-0 Default VLAN	36	768 GiB	2	480 GB	
<input type="checkbox"/>	ault07.cscs.ch 148.187.104.78 (PXE)	 On Ipmitool	ault	kleinm amd, gpu, vega, ...	default	default	fabric-0 Default VLAN	64	512 GiB	2	480 GB	
<input type="checkbox"/>	ault08.cscs.ch 148.187.104.79 (PXE)	 On Ipmitool	ault	kleinm amd, naples, gp...	default	default	fabric-0 Default VLAN	64	512 GiB	2	480 GB	
<input type="checkbox"/>	ault09.cscs.ch 148.187.104.80 (PXE)	 On Ipmitool	ault	kleinm amd, naples, epy...	default	default	fabric-0 Default VLAN	64	512 GiB	2	480 GB	
<input type="checkbox"/>	ault10.cscs.ch 148.187.104.81 (PXE)	 On Ipmitool	ault	kleinm nvidia, amd, nap...	default	default	fabric-0 Default VLAN	64	512 GiB	2	480 GB	

Initial Hardware

Hardware Available: AMD New

- AMD Epyc Naples 2S Vega
 - 2 socket = 64core/128thread
 - 3xVega 10
 - 16GB HBM2 (half-width memory bus)
 - 512GB DDR4
 - EDR Infiniband
- AMD Epyc Naples 2S NV
 - 2 socket = 64core/128thread
 - 2xNVidia V100
 - 32GB HBM2
 - 512GB DDR4
 - EDR Infiniband

Hardware Available: Intel New

- Intel Skylake CPU Only 2S
 - 2 socket = 36core/72thread @ 3GHz
 - 384GB DDR4
 - EDR Infiniband
- Intel Skylake 2S GPU
 - 2 socket = 36core/72thread @ 2.3GHz
 - 768GB DDR4
 - 4xV100 Nvidia GPU
 - 32GB HBM2
 - EDR Infiniband

HW Moving from Greina (Post-Decommision)

- Persistent Shared FS
 - Only available to shared-service environment
- Power8 Nodes
- Thunder X2 Node
- FPGA Nodes

Near Future: HW Additions

- Monitored PDU
 - Power monitoring for researchers external to BMC
- Vega 20 GPUs
 - ETA: Late April
- AMD Rome CPUs
 - 2H 2019
- Cascade Lake CPUs
- Additional Storage options
 - Rulers/SFF storage options promising
- 100G Ethernet switch (In Addition to EDR IB)

Arch Supported by MaaS (and current choices)

Architectures

☒ amd64

☒ arm64

☐ armhf

☐ i386

☒ ppc64el

☐ s390x

Next Steps

- Still work in Progress
- Shared Service Operational Level
 - PE/Libs “/apps” environment
 - Mixed architecture issues with shared-users compiling on head node
 - Performant Shared-FS
- Dedicated Access
 - Address open security concerns (Should be good to go now)
 - Finish testing of cleanup procedures to support more automated dedicated workload testing
 - Federated Auth to maas UI for dedicated mode
 - Github login Preferred (currently tested using Google)
 - Any interested users, contact me to help coordinate some test plans

Thank you for your attention.