

SPARTAN

Performance & Flexibility

An HPC-Cloud Chimera

OpenStack Summit, Barcelona, October 27, 2016

Lev Lafayette, University of Melbourne

Greg Sauter, University of Melbourne

Linh Vu, University of Melbourne

Bernard Meade, University of Melbourne

HPC, Cloud, and Edward

- HPC systems provide excellent performance, especially for multinode computational tasks.
- Cloud architectures provide flexibility and cost-efficiency - however they don't have the focus of HPC.
- University of Melbourne's general purpose HPC system (Edward) was due for retirement.
- Edward supported 886 users and 371 projects.
- Also provided specialist queues for particular departments.

Why change?

- Debate about the value proposition of HPC
- Profiling of Edward users indicated that usage was very heavily biased towards lots of single node and low memory tasks.
- Demand for more cores, more RAM, faster storage and interconnect is **actually** a demand for shorter queues
- Access to national peak facilities is not simple and many researchers choose to DIY with grant money
- Substantial investment in the Research Cloud

HPC Planning

- Two ways to design HPC:
 - Meet high-profile use cases-> suits few cases very well, but 90% cases only ok
 - Meet general case-> suits 90% cases very well, but not special cases, which are often pet projects of Profs, who have the money -> private clusters, often abandoned

The Spartan Vision

- A mix of bare-metal and virtualized resources
- Operated as a service in the Research Cloud
- All clusters have base requirements (management, login, etc.)
- Dynamically scalable depending on demand and availability

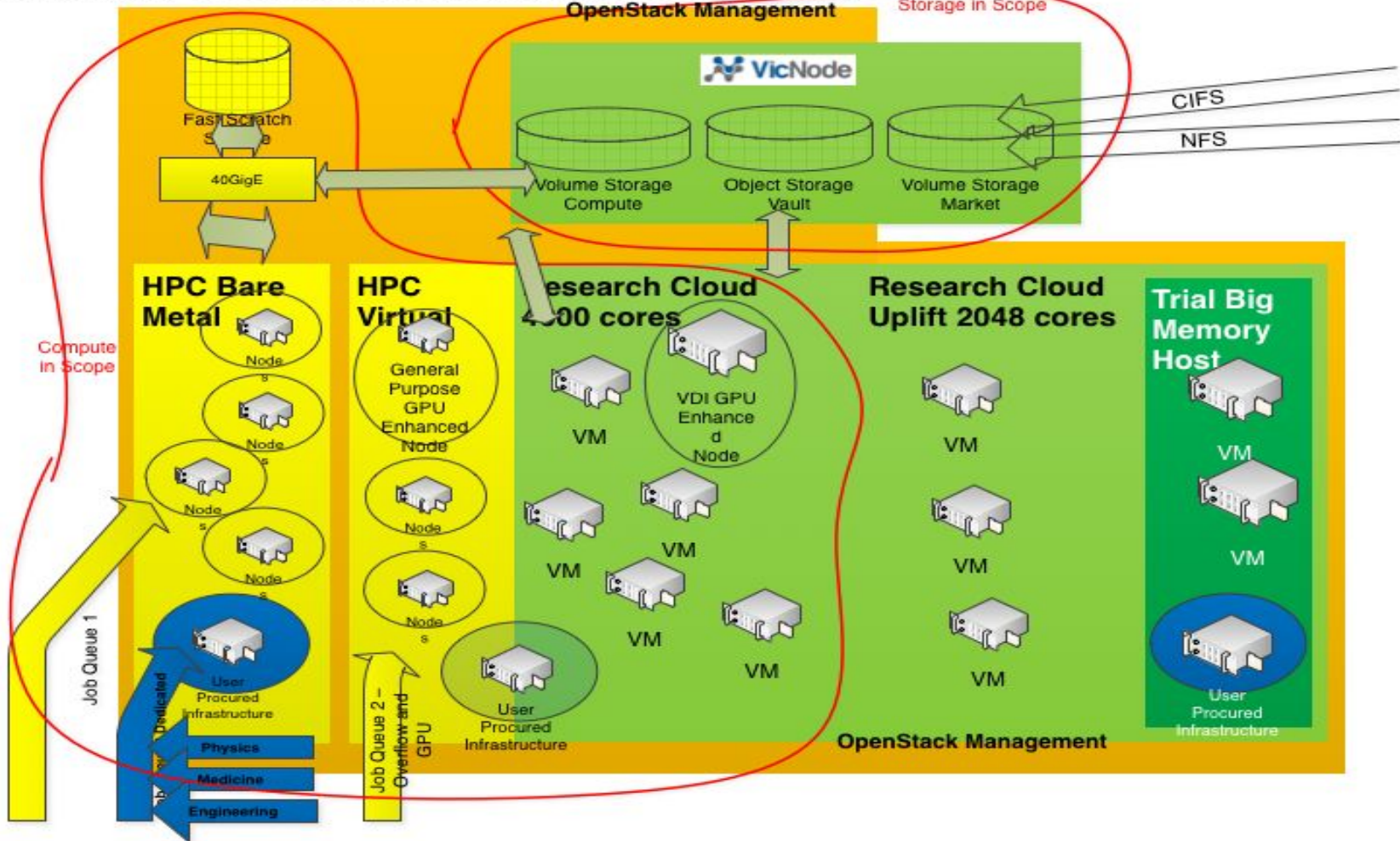
The Spartan Vision

- Investment can match the usage demand
- Greatly improves the Research Cloud vCPU utilisation
- No longer need forklift upgrades
- Grant money/departmental investment can be used to expand the required aspects of the service

Concept Diagram - Research Platform Services Compute and Storage Services

OpenStack Management

Storage in Scope



The Plan

".. no plan of operations extends with any certainty beyond the first contact with the main hostile force." (Helmuth von Moltke the Elder, 'On Strategy', 1871)

- Requirements gathering - 2014
- Partnering with Victorian Life Sciences Computation Initiative (VLSCI)
- Input from other research institutions
- External architect contracted for solution design

The Implementation

- Extensive training and transition programme.
- Integration with campus Active Directory was changed to LDAP instead.
- Noisy neighbour effects caused issues. Turn off overcommit on cloud nodes.
- All vHPC converted to cloud-burst using Slurm's power-saving feature.
- Connections to VicNode Aspera Shares and specific I/O nodes for data transfer.

Spartan Infrastructure

Breaking the HPC Stereotype Mould

- There is an entrenched mindset on how HPC systems should be built and operated:
 - You need dedicated hardware
 - You need Infiniband networking
 - You need a parallel file system such as Lustre/GPFS
 - You must have a rigid upgrade/patching policy
 - You need lots of money!

Spartan Infrastructure

- Faced with a very small budget and lots of unused cloud compute capacity, our solution:
 - Implement the HPC control plane on the cloud
 - Use cloud compute for low resource intensive jobs
 - Use bare metal for high resource intensive jobs
 - Use commodity SSD's to provide high performance shared storage for MPI jobs
 - Use the cloud CI/CD environment to manage configurations, upgrades and patching

Spartan Hardware

Compute

- Bare Metal Compute
 - High performance, Dell R630 with dual CPU 3.4Ghz 6 cores, 256GB RAM, 21GB per core, 40Gb network
 - Multi purpose, Dell R630 with dual CPU 2.3GHz 18 cores and 1.5TB RAM, 43GB per core, 25Gb network
 - High CPU cores, Dell 6320p Xeon phi 64 core 1.3Ghz 16GB on die RAM, 384GB RAM, 25Gb network
 - GPU, Dual Nvidia Tesla K80, rest specs as per high performance, 25Gb network

Spartan Hardware

Compute

- Cloud Compute
 - Dell C6300 chassis and 4 x 6320 modules
 - Each module has dual CPU E5-2683 v4 2.1GHz
16 cores and 768GB RAM, 24GB per core
 - 10Gb network connections

Spartan Hardware

Networking

- Bare Metal
 - Mellanox 2100 with 16 x 100Gb ports with a mixture of 25/50/100Gb, maximum of 64x25Gb connections
 - RDMA over ethernet
 - Cumulus Linux OS
 - Automation with Puppet
 - Full virtual network test environment with 1 real switch to test upgrades integrated into CI/CD system
 - Routing from the compute nodes

Network Latency Performance

- MPI ping pong test between two compute nodes on separate infrastructure
- Network Environments:
 - UoM large HPC service using Mellanox FDR14 56Gb Infiniband
 - Legacy Edward service using 10Gbe (Cisco Nexus)
 - Spartan HPC cloud service using 10Gbe (Cisco Nexus)
 - Spartan HPC bare metal service using Mellanox ConnectX4 cards and SN2100 ethernet switch

Network Latency Performance

Service	Network Device	Network Speed	Protocol	Latency (micro secs)
UoM HPC Traditional	Mellanox	56Gb	Infiniband FDR14	1.17
Legacy Edward HPC	Cisco Nexus	10Gbe	TCP/IP	19
Spartan Cloud nodes	Cisco Nexus	10Gbe	TCP/IP	60
Spartan Cloud nodes	Cisco Nexus	10Gbe	TCP/IP and SRIOV on hypervisor	????
Spartan Bare Metal	Mellanox	40Gbe	TCP/IP	6.85
Spartan Bare Metal	Mellanox	25Gbe	RDMA over ethernet	1.84
Spartan Bare Metal	Mellanox	40Gbe	RDMA over ethernet	1.15
Spartan Bare Metal	Mellanox	56Gbe	RDMA over ethernet	1.68
Spartan Bare Metal	Mellanox	100Gbe	RDMA over ethernet	1.3

Spartan Hardware

Storage

- Bare metal nodes:
 - /scratch shared storage for MPI jobs, Dell R730 with 14 x 800GB mixed use SSD's providing 8TB of usable storage, NFS over RDMA
 - /var/local/tmp for single node jobs, pcie SSD 1.6TB
- Root and volume storage for cloud VM's, CEPH RBD
- /project and /home for user data & scripts, NetApp SAS aggregate 40TB usable, NFS (pNFS planned)
- Trialing a SanDisk 128TB SSD storage presented through CEPH (CephFS)

Spartan Hardware

Storage

- /scratch performance figures using “flexible IO test”, running 4 processes and transferring 85GB of data:

Action	Transport	Throughput	Latency 750 us
Random Read	TCP	2.84GB/s	36%
Random Write	TCP	2.71GB/s	69%
Random Read	RDMA	4.62GB/s	99.96%
Random Write	RDMA	3.08GB/s	97%

Spartan Run Time Results

Job	Description	Resources	Node Type	Data in/out	Runtime		
					Super computer (2.3 GHz)	Spartan Bare Metal (3.4 GHz)	Spartan Cloud (2.1 GHz)
		Required					
BWA	Short read DNA Sequencing	Disk intensive	8 core Single Node	15/11 GB	1:18:49	1:02:56	1:40:21
GROMACS	Molecular dynamics, simulation of large biomolecular systems	Compute intensive	128 core Multinode	2/0.12 MB	0:30:02	0:30:10	0:30:32
NAMD	Molecular dynamics, simulation of large biomolecular systems	Compute & I/O intensive	128 core Multinode	175/53 MB	1:11:41	1:00:46	1:55:54

Spartan Configuration with Openstack

- Uses SLURM power saving feature to do cloud bursting
- Pre determined custom cloud flavour 8C 64GB RAM
- RHEL 7 Glance image created from puppetised golden node
- Build time 30 seconds, 5 min (configurable) idle-to-suspend
- Bare metal currently custom PXE+kickstart+puppet system but investigating Ironic
- Consistency in OS and libraries in the management, login, i/o, physical, and cloud compute nodes maintained through Puppet reviewed through Gerrit. Upgrade image released 2 week cycle

API Interfaces:

Nova - compute VM management

Neutron - networking

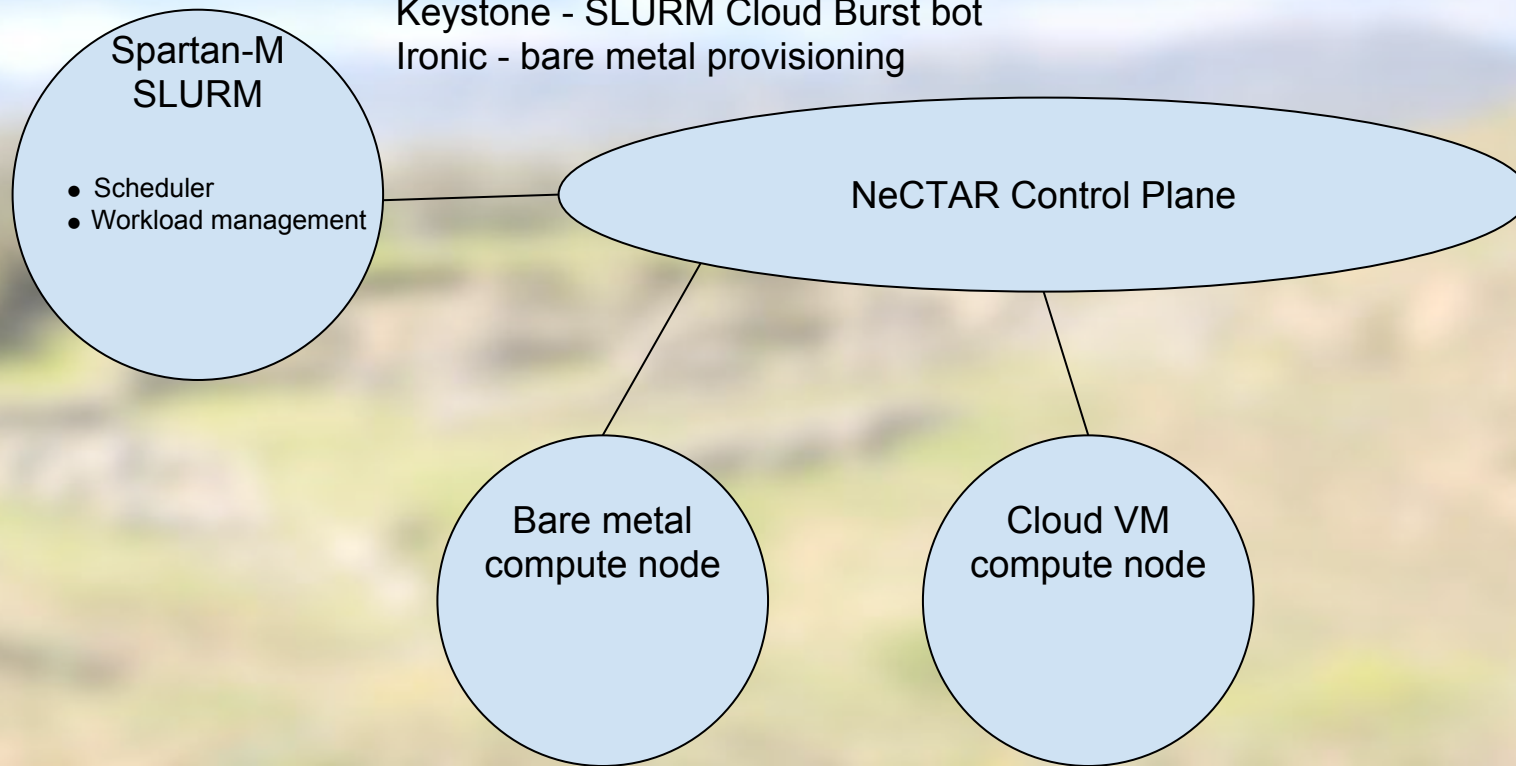
Cinder - VM root and volume storage

Glance - golden image

Designate - DNS

Keystone - SLURM Cloud Burst bot

Ironic - bare metal provisioning



Scheduler and Applications

- SLURM Workload Manager for resource manager and job scheduler.
- EasyBuild is used for software installations; builds from source codes with specified toolchains and dependencies.
- Automatically generated LMod environment module files.

Training Programme

- Extensive training programme to assist researchers; every week a day is set aside.
- Training days include: Introduction to Linux and HPC, Edward to Spartan HPC Workshop, Shell Scripting for HPC, Introduction to Parallel Programming.
- Application and department specific courses being developed for Economics (R, Stata), and Mechanical Engineering (OpenFOAM), plus profiling courses.

The Future

- Large separate partition for proteomics project.
- Closure and transition of Edward users (Dec 2016 deadline).
- Separate partitions for different hardware and separate login nodes for special groups (e.g., classes).

Prediction: Spartan's architecture will be the model for high throughput general purpose computing in the future.

The Europe 2016 Tour

Prior to OpenStack visits were carried out at several HPC centres in Europe including Center for Scientific Computing (CSC), Goethe University Frankfurt; High Performance Computing Center (HLRS), University of Stuttgart; High Performance Computing, Albert-Ludwigs-University Freiburg,; European Organization for Nuclear Research (CERN), Centre Informatique National de l'Enseignement Supérieur (CINES) Montpellier, and next week to the Centro Nacional de Supercomputación, Barcelona

Chimeras and Cyborgs

- The HPC centre in Freiburg also has an HPC-cloud hybrid.
- However the UniMelb system is a chimera (a multiheaded beast) the Freiburg system is a cyborg (an admixture of the two technologies).
- An OpenStack Nova Client runs on the HPC nodes, and OpenStack management services run on a cluster. There is essentially one big queue, which can run either traditional HPC compute or a cloud instance.