

Towards Disaster Recovery in the Cloud

Cliff Addison, Manhui Wang



UNIVERSITY OF
LIVERPOOL

Overview

- Liverpool in 2018 demonstrated strategic benefits of cloud for research.
- Deployment was ad hoc, but worked.
- We want to embed cloud for research and graduate teaching.



Successful cloud scenarios

- Cloud bursting – more cycles needed for a short period
 - typically for papers or presentations
- High throughput workflows
 - Current Windows Condor pool limited to circa 8 hr jobs
- Scoping studies
 - I think I need X cores and Y GB of memory for my research
- GPU nodes for Deep Learning
- **Avoiding** large data transfers in this first instance



Cloud bursting - 1

- An existing Condor pool can be extended easily to the cloud.
 - Users just request the cloud resource on local Condor server – acts as scheduler.
 - Customise a standard AWS Linux image with necessary extra software and then save this image so is ready to go.
 - Have an in-cloud manager that deploys compute images; liaison with scheduler.
 - Spot market makes the compute even more cost-effective
 - Fits perfectly with Condor cycle stealing idea



- Test that target instances are good enough
Micro instances may be too slow so more expensive for compute



How this can work...

- Researcher came to us in May 2018 with an urgent request to run 100,000 simulations related to a paper under review.
- Our AWS Condor pool ideal.
- Cost per simulation cheapest on t2.medium, but fastest on c4.large or c5.large.



- Final set up:
 - 1000 jobs with 100 simulations each, pool size of 400 (so 400 jobs at once) completed task in **7h 21m**. Serially would need about **98 days** – massive speed-up.
 - Price **\$51.16**

Paper resubmitted on time!



State of play end of 2018

- Working on AWS an eye-opening experience.
 - Some learning curve with the EC2 and with S3 storage
- Started with Alces Flight clusters and spinning specialised instances on EC2 (e.g. Condor in the cloud).
 - Just creating instances with keys for particular groups is great for small numbers of groups, but it does not scale.
- Can get major additional benefits to an on premise HPC.
- Planned use cases for 2019:
 - Seamless cloud access from local cluster
 - Replication of some cluster functionality in the cloud

Moving forward on a narrow front - HPC

- Have some basics in both AWS and Azure tenancies
 - Use Active Directory authentication
 - VPC connection so cloud resources appear on campus network
- Liverpool HPC has:
 - Defined set of users (circa 70 active every month)
 - Stable software offering (slowly growing)
 - Alces Flight provide system and software framework
- Liverpool HPC needs:
 - Better resiliency – no failover component
 - More flexible environment for new users
 - Better development / experimentation support

Resiliency

- Many core services have failover capability across two on-campus data centres.
- Hard to support two on-campus HPC systems with an active-active failover mode (active-passive is silly)
 - HPC systems often sited in a single data centre
 - HPC systems bought at different times, maybe from different vendors
 - Need to synchronise user storage, applications
- What about failing over / bursting to a cloud setting?

HPC resiliency in the cloud

- Want to have an on-demand clone available
- Compute can be brought on-line relatively quickly.
- Front-end / login node and storage need to be there through the life-time on the cluster.
- Compute node costs can be controlled via autoscaling options and by exploiting the spot-market (on AWS) – how many nodes are needed?
- Most cloud platforms want a year of always on use before offering major discounts over their on-demand price.
- How deal with storage??



HPC cloud resiliency – storage issues

- Three types of data storage to consider
 - System – node images and applications
 - Persistent, relatively stable, modestly sized
 - Probably current to within a week is fine
 - User home directories
 - Many systems keep to a small number of TBytes for local backup
 - Daily incremental back-up to the cloud with occasional full back-up should be possible.
 - User volatile / work areas
 - These can be huge.
 - If shutdown is planned, can get relevant users to pre-stage important data; typically during the local rundown before shutdown..
 - Cloud as a primary and permanent site for volatile data?
 - Will be slow and might be very expensive...

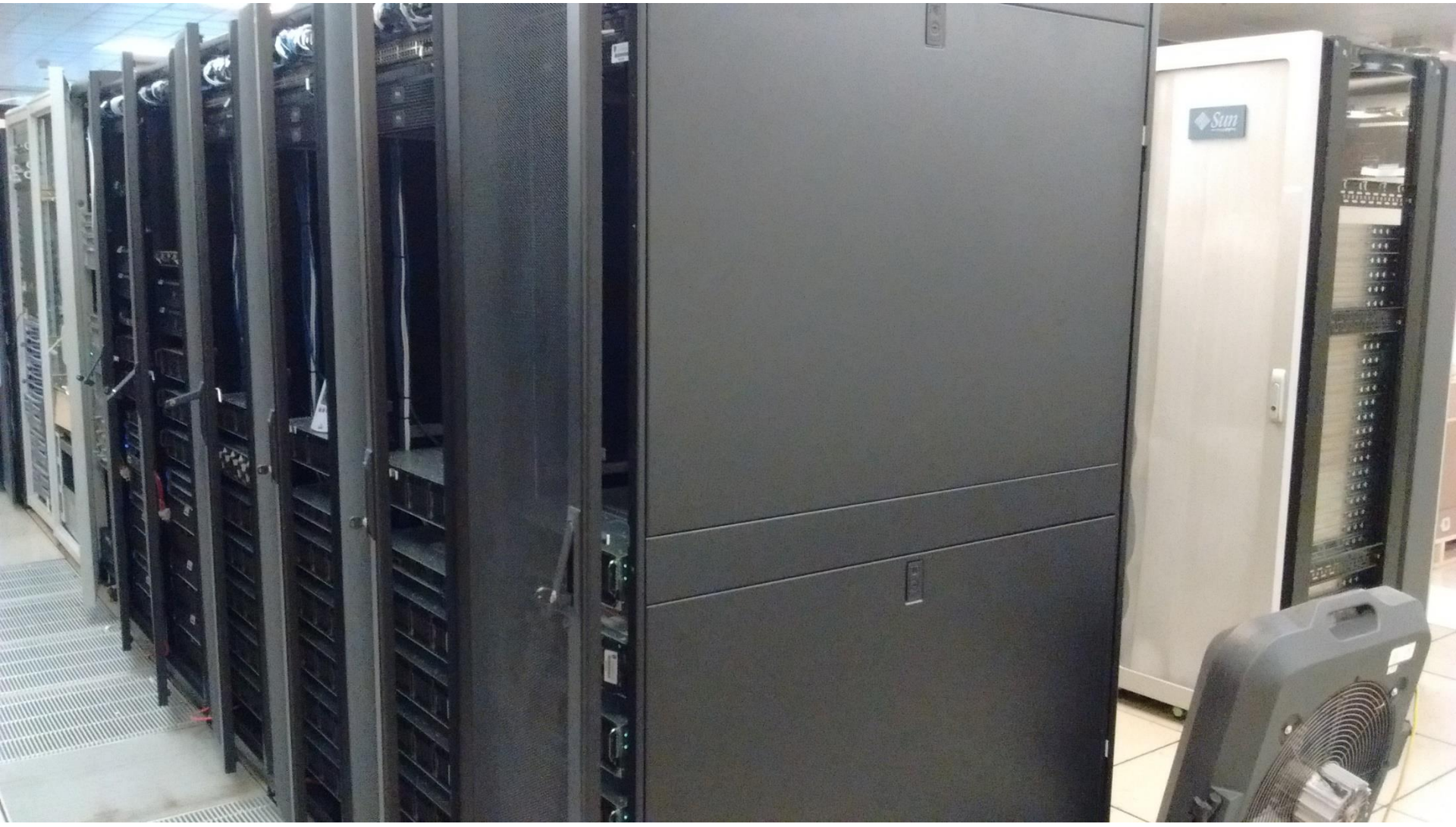
My current understanding of AWS storage

- There are the traditional 3 storage layers:
 - Elastic Block Storage – fast access for active storage tied to hardware instances. Always need some of this on cluster. [100 GB costs about \$8.10 per month]
 - Standard S3 Object Storage – slower but accessible from anywhere in the AWS cloud (and elsewhere with S3 supported logical devices) [100 GB costs about \$2.30 per month]
 - S3 Intelligent Tiering, S3 Standard Infrequent access – slowly changing; not often accessed [100 GB/month \$2.40, \$1.31 resp.]
- Also there are the archival options, not for HPC(?)
 - S3 Glacier and S3 Glacier Deep Archive – 6 month no-change?
- Other cloud vendors have similar arrangements.

Immediate problem to hand

- New data centre with better cooling and generator-backed power for all systems coming on stream now.
- Need to move Dell / Alces system to new home.
- Old Bull (SandyBridge) cluster available during this time, but lose 4000 cores for circa 10 days.
- Idea – augment SandyBridge cluster with some cloud-based Cascade Lake (AVX-512 support) and AMD nodes – great general purpose + GPU

Today – 6 racks air-cooled



**Today –
Racks
generally
look like
this –
cabling
and
storage
challenge**



Going to these water cooled racks



Trial run – June 2019

- Data centre 24 hr outage to swap power supplies
- Use outage to test cloud cluster.
- Huge advantage with Alces Flight.
 - Environment largely cloud ready
- Deals vary with provider and time needed etc.
 - AWS better this test period.
 - Sacrificed faster interconnect for more nodes
- Plan – basic system login node, storage, small test node available several days before and after outage.

System configuration

- Login node – Skylake 24 cores, 350 GB memory
- 10 TB shared storage for all nodes
- 2 x 2C/16GB small compute nodes for testing
- 1 x Single Nvidia V100 GPU compute node
- 20 x 36C/128GB Skylake compute nodes
 - Only available just before poweroff and for 3 days after
- 100 GB of data / day down load from cluster
- Whilst local system up, easy to copy files.
- Used existing usernames with ssh keys for access.
- Home filestore **NOT** copied over.



Important considerations

- Our tailored Alces Flight Gridware preinstalled, we needed to copy over local application files.
 - Local module files completely replicated on cloud system
- Emphasis was on SMP parallel or coarse grain parallel plus Deep Learning on GPU.
- Nodes were hyperthreaded - users told to ask for exclusive access to avoid overloading.
- ssh key access – users told where to grab this from and the name of the cluster to ssh to.
- Cluster was only accessible from on-campus or via VPN to local system and then to cloud.

Lessons learned

- Once access obtained, people had no problem editing slurm scripts to run jobs.
- ssh keys, off-campus access slight niggles.
- Fewer people than expected used the system (only over a weekend).
 - Next time – check on how many likely users there will be.
 - Had more compute nodes than necessary; gpu node was used
- Cost of main compute nodes ~ 75% of overall cost
 - GPU node ~ 20%
- Similar look and feel to local system big help.
- Test nodes not helpful because lacked AVX-512



Next steps – the major outage

- Had hoped cloud Barkla would be happening already
 - Move delayed until 21/10
- Cloud cluster will be needed for 10 days
- Going for 100 gbps InfiniBand interconnect
- Bulk Cascade Lake with 4 AMD nodes for Gaussian and codes without AVX-512 builds.
- Now helping people with file transfer and getting started tests (both easy from Barkla).
- Trying to understand user demand for next week
 - Can limit number of nodes on offer to match need



Plans after data centre move

- Work out cost of storing user home directories and system images for resilient deployment.
- Integrate cloud cluster with university network
 - Active directory for authentication
 - VPC so cluster appears as on the University network
- Try bringing up cloud cluster alongside Barkla to test “easy access” and cloud bursting potential.
- Experiment with spot market on AWS
 - Massive savings but small number of nodes
- Get firm University budget to sustain resiliency – storage plus occasional compute
 - Compute costs for full cluster mount very quickly!!



Summary – general issues

- Need cloud cluster to have a similar look and feel to the local cluster.
- Integrate the cloud cluster into your local environment
 - Active Directory / VPC so appears on campus network
- What storage is put where?
 - Maybe local back-up that can push onto cloud as needed?
 - Local storage that is pushed to the cloud avoids lock-in – flexibility is good!
- Compute and login nodes created on-demand
 - How many compute nodes makes sense?? Interconnect??
 - Spot-market for some / all nodes
 - Auto-scaling is a definite must.

