# Operations and comparison of a commercial datacentre

Peter Love
Lancaster University
HEPiX Spring 2015 - Oxford

# Introduction

- The deployment of academic and commercial facilities is becoming more widespread

- What does it take to use these opportunistic resources?

- HEP experiments are in a good position to take advantage of these resources, a variety of tools are being used

- Described here are some experiences when commissioning the site

- Conclude with a comparison with an existing grid facility using a few key metrics

## Limit Summary

**Instances**
Used 28 of 120

**VCPUs**
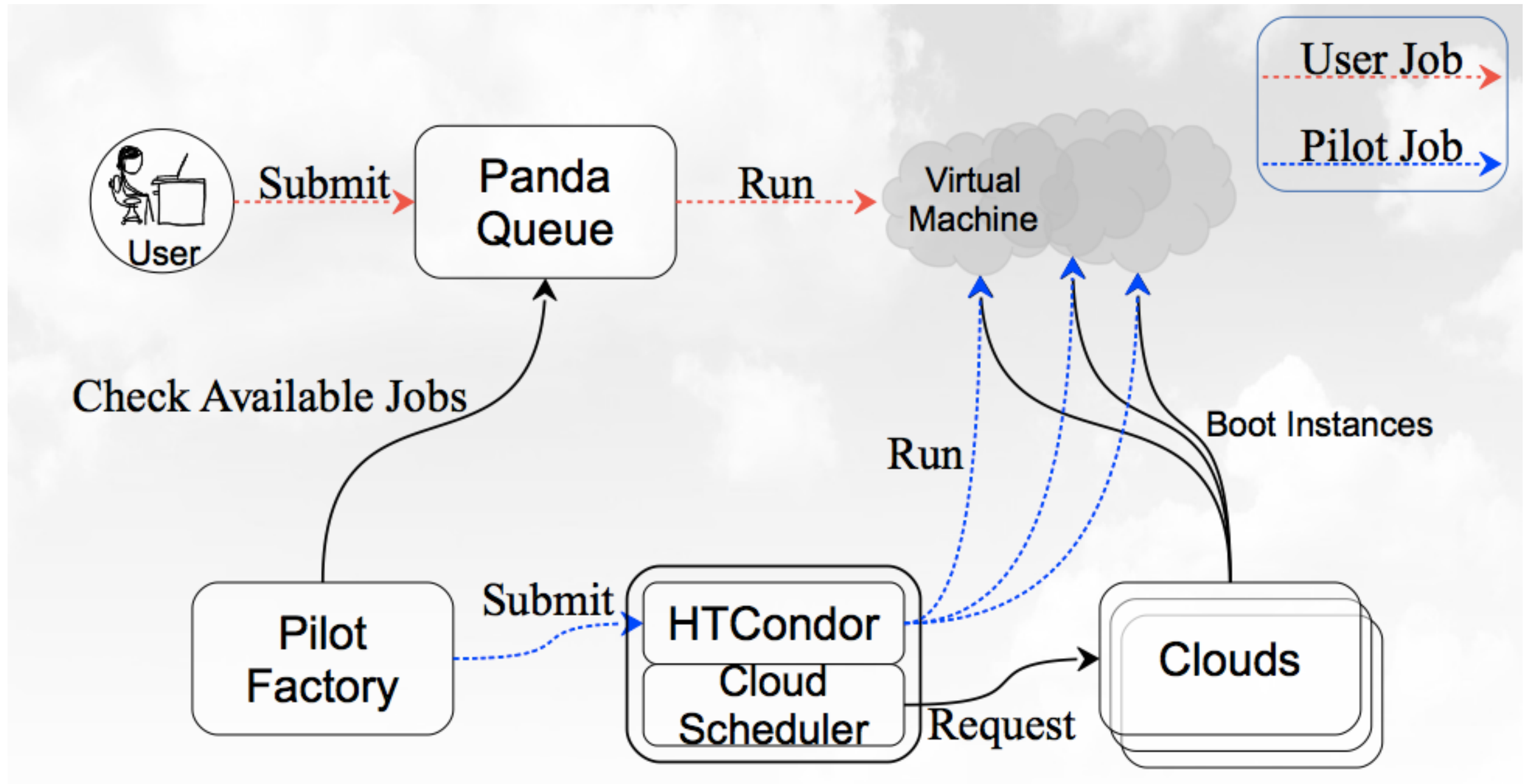Used 218 of 220

**RAM**
Used 434.0GB of 500.0GB

**Volumes**
Used 1 of 50

**Volume Storage**
Used 100.0GB of 700.0GB

- Relatively new commercial operation providing co-lo facilities and a cloud hosting service

- Built on Openstack and Ceph

- Large facility 1800sq.m2, enough for 850 30kW racks with UPS backup

- Involvment in Helix Nebula becoming member in early 2015. Engaging with big science projects.

- Lancaster collaboration at early stages of commissioning the facility, we have a small tenancy

- Objective from our side is to exploit opportunistic resources in as simple way as possible. This is not AWS or GCE.

# What approach do we use to make use of this opportunity?



Slide: Frank Berghaus

# VMs are uCernVM using Shoal for squid discovery

## List of Active Squids

**12 active in the last 180 seconds**

| # | Hostname | Public IP | Private IP | Bytes Out | City | Region | Country | Latitude | Longitude | Last Received | Alive | Verified | Access Level |
|---|----------|-----------|------------|-----------|------|--------|---------|----------|-----------|---------------|-------|----------|--------------|
| 1 | pygrid-kraken.hec.lancs.ac.uk | 194.80.35.16 | 10.41.52.16 | 677 kB/s | Lancaster | | United Kingdom | 54.0667 | -2.8333 | 1s | 54h59m28s | ✗ | Global |
| 2 | ca17.cern.ch | 128.142.163.110 | | 86840 kB/s | Geneva | | Switzerland | 46.1956 | 6.1481 | 1s | 54h59m17s | ✗ | Same Domain Only |
| 3 | atlascaq3.triumf.ca | 142.90.110.68 | | 0 kB/s | Vancouver | | Canada | 49.2765 | -123.2177 | 4s | 54h59m15s | ✔ | Global |
| 4 | t2software03.physics.ox.ac.uk | 163.1.5.175 | | 3991 kB/s | Oxford | | United Kingdom | 51.75 | -1.25 | 5s | 10h17m17s | ✔ | Global |
| 5 | ca02.cern.ch | 188.185.165.173 | | 86864 kB/s | Cern | | Switzerland | 46.2324 | 6.0502 | 6s | 54h59m49s | ✗ | Same Domain Only |
| 6 | ca19.cern.ch | 188.185.163.104 | | 29127 kB/s | Cern | | Switzerland | 46.2324 | 6.0502 | 6s | 54h59m30s | ✗ | Same Domain Only |
| 7 | ca06.cern.ch | 188.184.148.164 | | 96207 kB/s | Cern | | Switzerland | 46.2324 | 6.0502 | 9s | 54h59m13s | ✗ | Same Domain Only |
| 8 | squid-test01.gridpp.rl.ac.uk | 130.246.183.249 | | 0 kB/s | Appleton | | United Kingdom | 51.7 | -1.35 | 14s | 54h59m11s | ✗ | Global |
| 9 | ca05.cern.ch | 128.142.152.230 | | 52777 kB/s | Geneva | | Switzerland | 46.1956 | 6.1481 | 18s | 54h59m32s | ✗ | Same Domain Only |
| 10 | ca16.cern.ch | 188.184.135.75 | | 51580 kB/s | Cern | | Switzerland | 46.2324 | 6.0502 | 22s | 54h59m28s | ✗ | Same Domain Only |
| 11 | ip-172-31-36-99.us-west-2.compute.internal | 52.11.50.231 | 172.31.36.99 | 1 kB/s | Boardman | | United States | 45.8399 | -119.7006 | 25s | 54h59m24s | ✔ | Same Domain Only |
| 12 | kraken01.westgrid.ca | 206.12.48.249 | 172.22.2.25 | 800 kB/s | Vancouver | | Canada | 49.2836 | -123.1041 | 30s | 54h59m26s | ✔ | Global |

# Early days - 2014

- First incarnation was Havana with nova networking

- Hardware was testbed quality

- Difficult to find stability, metadata service had hardware limits and was generally unreliable

- At this point adding workarounds was hard due to contextualization being buried in puppet modules located on a private repo - an organisational issue

- Debug cycle was slow

# Later on

- Later upgraded to Icehouse and Neutron with production quality hardware, HA etc.

- The performance and stability was fixed but tweaking things was still cumbersome.

- Workaround was to spin-up persistent VMs via 'nova boot'.

- Eventually moved to specific cloud-init yaml contextualization, <u>hosted on github</u>. Things were much more tranparent.

- The puppet approach was scrapped and these yaml scripts are now the way we provision the production clouds.

- Flexibility was needed to workaround issues with uCernVM ganglia cloudinit module and also ganglia app version to deal with override_hostname.

Hardware and Icehouse upgrade
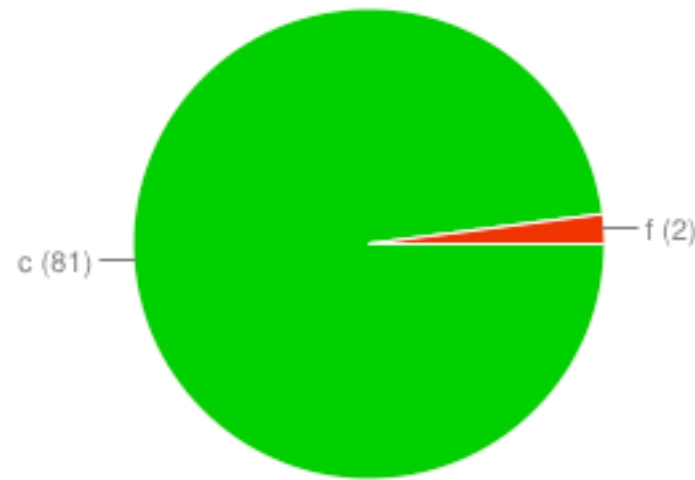Nova API response time

# Comparison of a Grid and Openstack site for ATLAS production

- UKI-NORTHGRID-LANCS-HEP_SL6 (~2000 cores on Lancs grid site)

- UKI-NORTHGRID-LANCS-HEP_CLOUD (~200 cores on commercial Openstack)

- Several ATLAS Hammercloud Stress tests were run on both sites and metrics compared

- Each stress test ran for 24 hours and consisted of a continuous stream of jobs

- Jobs were mc12 AtlasG4_trf 17.2.2.2 using a single input dataset located on the grid storage ~100MB

- Metrics are compared on following slides with results as one may expect
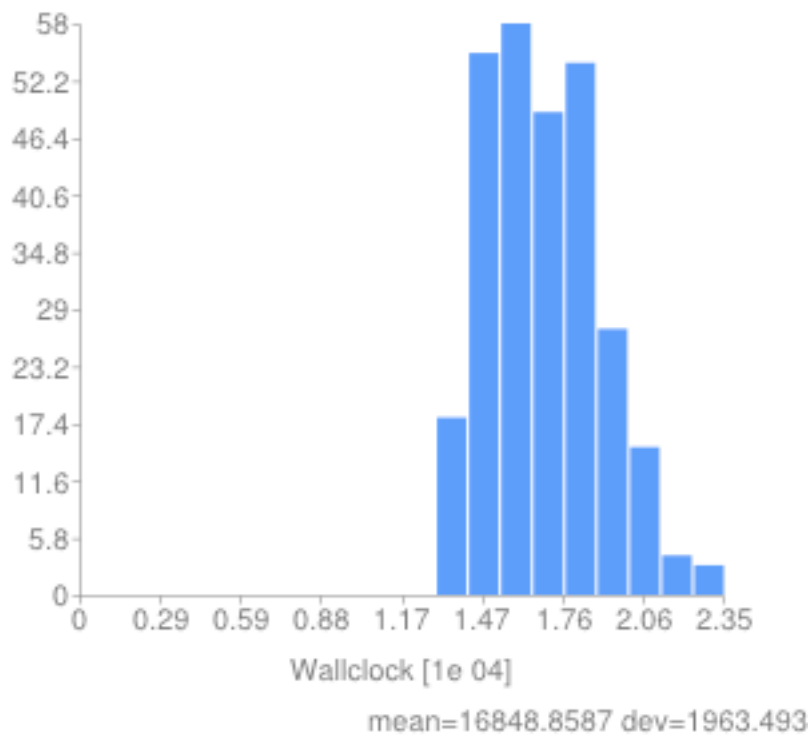
UKI-NORTHGRID-LANCS-HEP_SL6

c (283) — f (11)

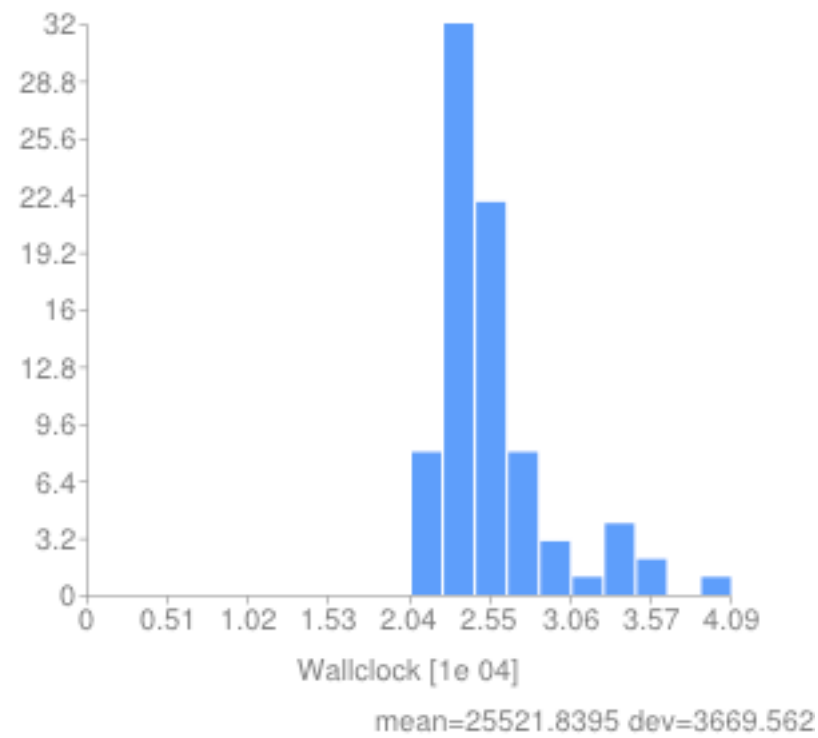UKI-NORTHGRID-LANCS-HEP_CLOUD

c (81) — f (2)

Success rate similar, grid site
processed four times more jobs
283 vs. 81 jobs

UKI-NORTHGRID-LANCS-HEP_SL6
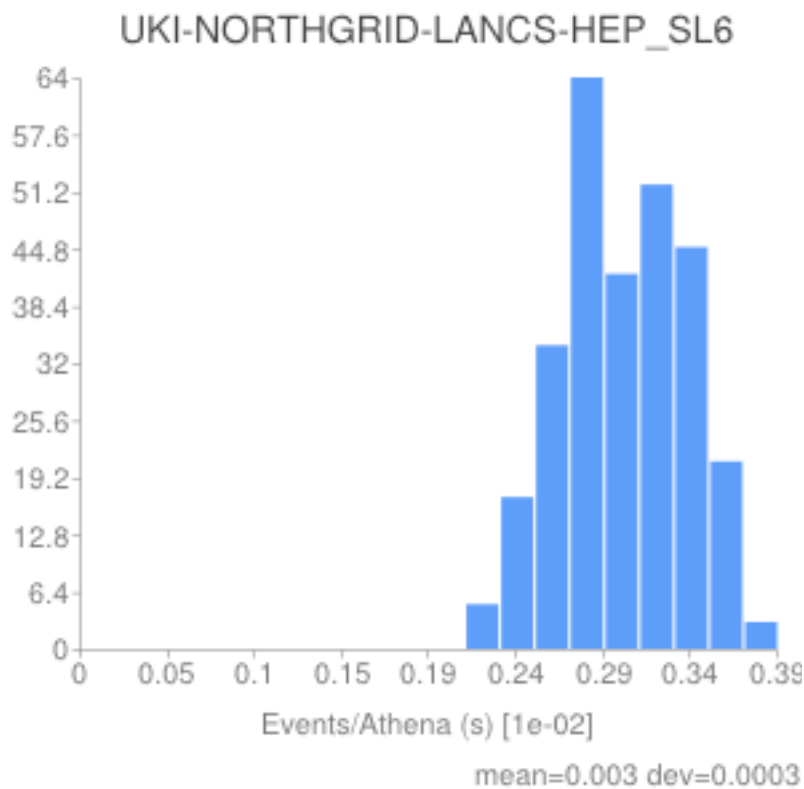
Wallclock [1e 04]

mean=16848.8587 dev=1963.493

UKI-NORTHGRID-LANCS-HEP_CLOUD

Wallclock [1e 04]

mean=25521.8395 dev=3669.562
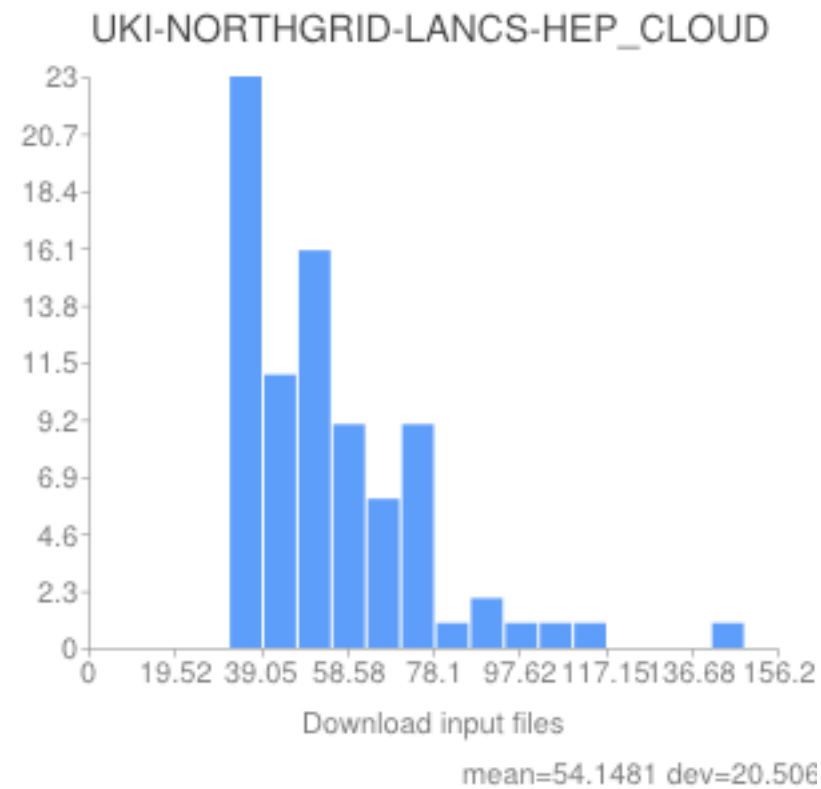
Wallclock twice as long on cloud site
with greater spread in runtimes

UKI-NORTHGRID-LANCS-HEP_SL6

mean=98.1 dev=1.0

UKI-NORTHGRID-LANCS-HEP_CLOUD

mean=92.8 dev=3.6

CPU percentage slightly down on the cloud site with a greater spread in efficiency

UKI-NORTHGRID-LANCS-HEP_SL6

mean=0.003 dev=0.0003

UKI-NORTHGRID-LANCS-HEP_CLOUD

mean=0.002 dev=0.0002

Events/Athena (s) although slower on cloud resource, the spread is similar to the grid site

UKI-NORTHGRID-LANCS-HEP_SL6

Download input files
mean=10.8516 dev=3.8615

UKI-NORTHGRID-LANCS-HEP_CLOUD

Download input files
mean=54.1481 dev=20.506

Input data stage-in time.
Remote vs. local storage.
(5x slower)

UKI-NORTHGRID-LANCS-HEP_SL6

Setup Software Time
mean=15.4134 dev=6.7218

UKI-NORTHGRID-LANCS-HEP_CLOUD

Setup Software Time
mean=45.0247 dev=14.9757

Software setup time. Via cvmfs
and influenced by squid.
(3x slower)

# Summary of metric comparisons

- These are a few metrics of interest showing the differences in performance.

- The result are to be expected given the architectural differences in hardware and network between the grid and cloud site.

- No real worries although clear where more work is needed. Immediately:

  1. persistent local squid

  2. persistent local ARC CE
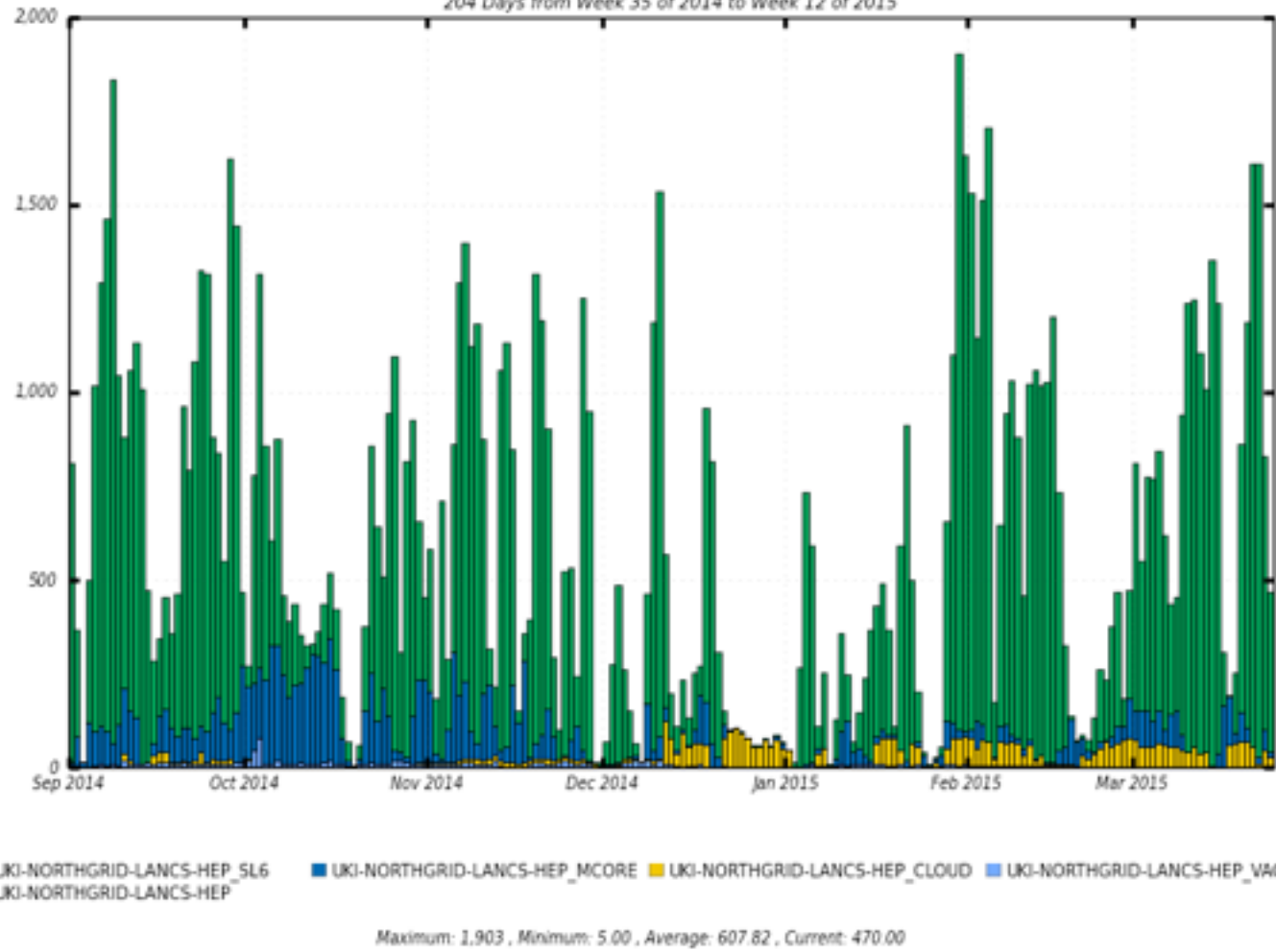
# Cloud local object store

- This hosting service provides a Ceph object store

- Used as a backend for both Swift and S3 interfaces

- How can we (as users) exploit this facility?

- Various approaches are in development

  - ATLAS Event Service (John's talk on Thursday)

  - Via FTS3 and special pilot settings (John's talk on Thursday)

  - ARC-CE as a gateway to pre-staging data

# Provisioning technologies - continue to assess merits

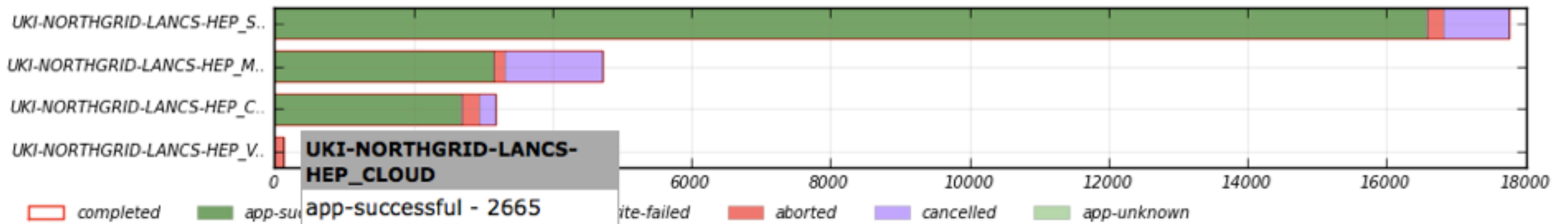- We have experince (via ATLAS) of using the following approaches to provisioning VMs on Openstack

  - cloudscheduler - UVic project very successful gathering resources around the globe

  - htcondor - mainly via BNL expertise described in John's talk on Thursday

  - vcycle - using the VAC model as decribed in Andrew's talk on Thusday

Running jobs
204 Days from Week 35 of 2014 to Week 12 of 2015

■ UKI-NORTHGRID-LANCS-HEP_SL6   ■ UKI-NORTHGRID-LANCS-HEP_MCORE   ■ UKI-NORTHGRID-LANCS-HEP_CLOUD   ■ UKI-NORTHGRID-LANCS-HEP_VAC
■ UKI-NORTHGRID-LANCS-HEP

Maximum: 1,903 , Minimum: 5.00 , Average: 607.82 , Current: 470.00

Completed Jobs per site

UKI-NORTHGRID-LANCS-HEP_S..
UKI-NORTHGRID-LANCS-HEP_M..
UKI-NORTHGRID-LANCS-HEP_C..
UKI-NORTHGRID-LANCS-HEP_V..

□ completed   ■ app-suc...   ...ite-failed   ■ aborted   ■ cancelled   ■ app-unknown

**UKI-NORTHGRID-LANCS-HEP_CLOUD**

app-successful - 2665
app-failed - 28
site-failed - 0
aborted - 250
cancelled - 229
app-unknown - 0
completed - 3172

# Summary

- New cloud resources need commissioning and we have procedures to do this quickly.

- Tools are flexible out of necessity. Once stable the site can be quickly integrated into production machinery.

- Development continues to optimize the performance and also create a recipe for a self-contained facility, relying less on outside services.