



THE UNIVERSITY
of EDINBURGH



Archer Status

Andrew Washbrook

University of Edinburgh

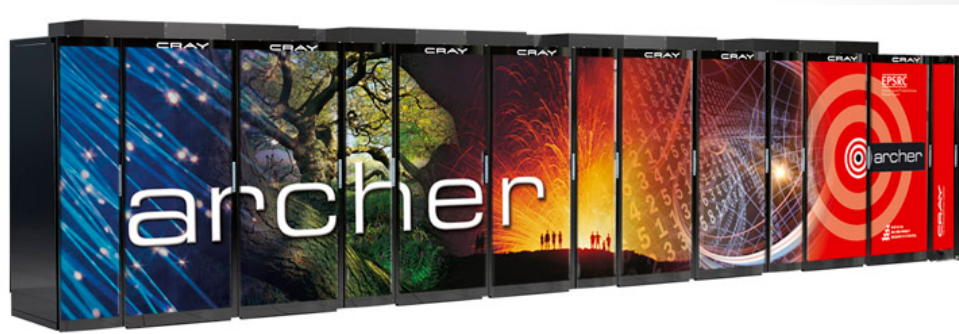
Centre of Excellence in Physics at Extreme Scales

WP5/A2: Exploitation of HPC Clusters for LHC data intensive workflows and analysis applications

19th September 2014

Archer

- **Archer** is the UK's primary academic research supercomputer
- Operational since Nov 2013



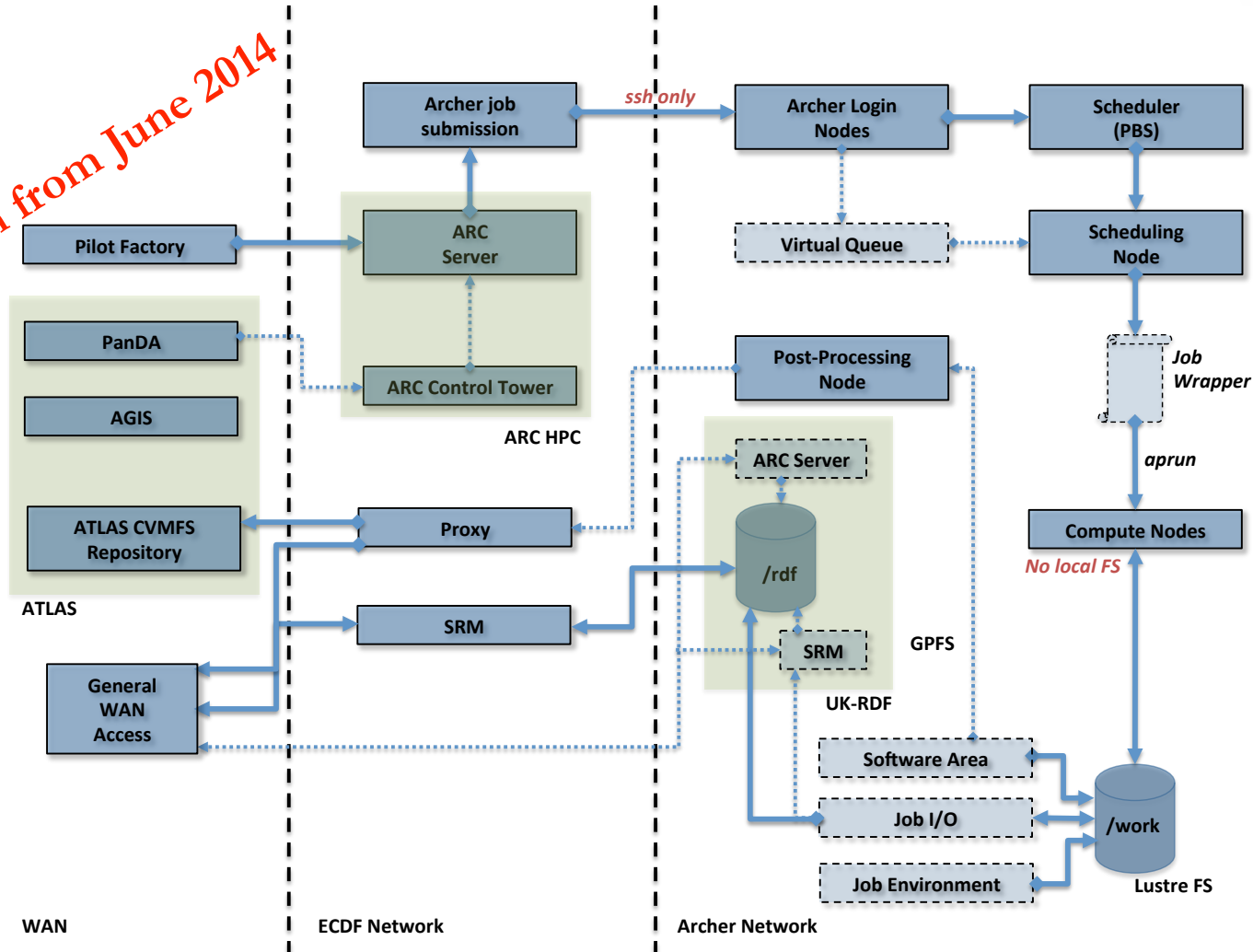
- **Cray XC30** system
- Each compute node comprises of:
 - 2 x 12-core 2.7 GHz Ivy Bridge processors
 - At least 64 GB of DDR3-1833 MHz main memory
- Cray Aries interconnect (multi-tier all-to-all connectivity)
- 4.4 PB scratch storage (Lustre)

- 3008 compute nodes → **72,192** cores (in 16 cabinets)
- 1.56 Petaflops of theoretical peak performance.
- Top 500 supercomputer position (June 2014): 25th
- Directly connected to the UK **Research Data Facility (RDF)**
 - 750 TB reserved for LHC Storage Activities



Archer Integration

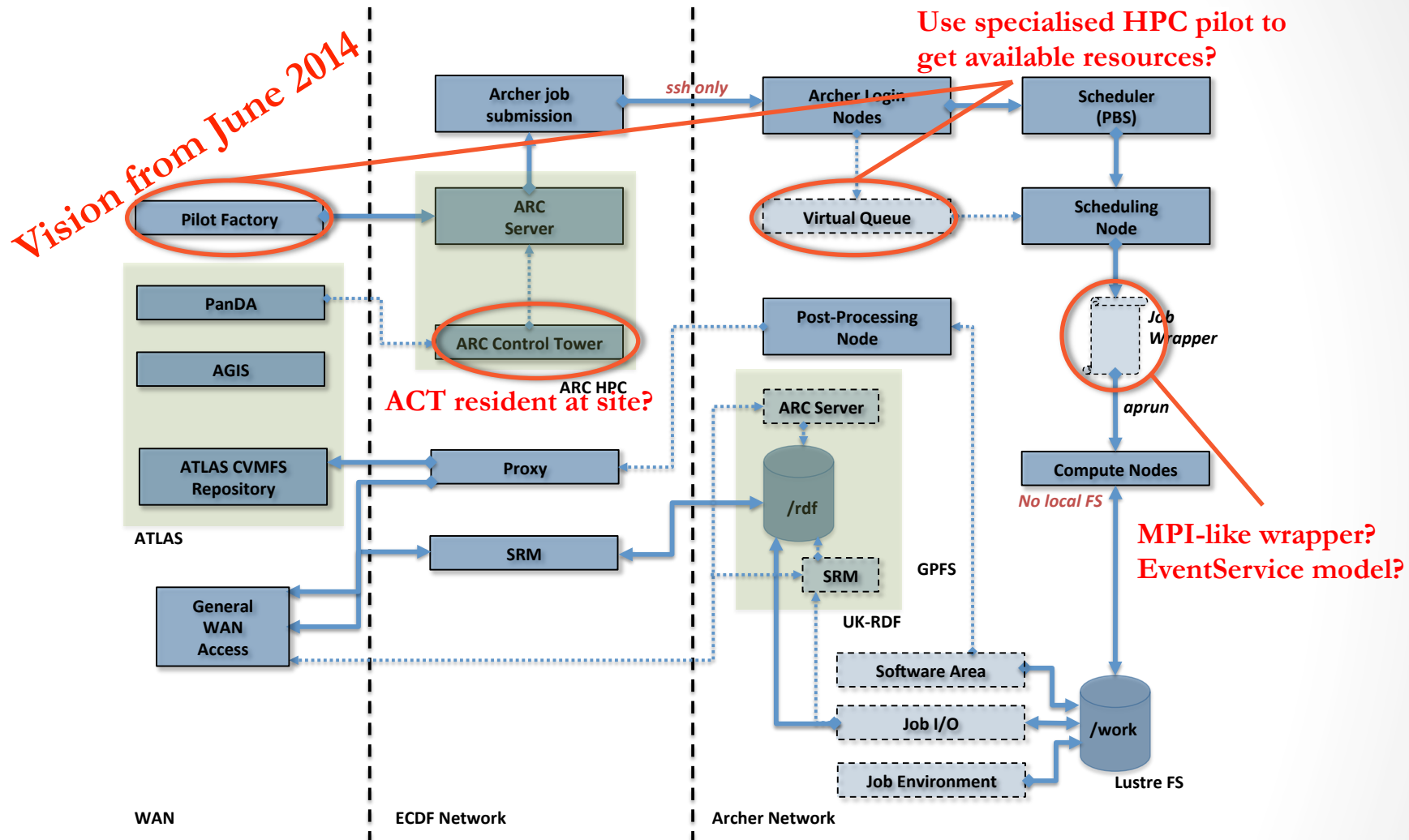
Vision from June 2014



- Tier-2 site (UKI-SCOTGRID-ECDF) and Archer reside at the same facility
- Above model is flexible - can adapt to fit successes from other HPC resources



Archer Integration Revision



- Tier-2 site (UKI-SCOTGRID-ECDF) and Archer reside at the same facility
- Above model is flexible - can adapt to fit successes from other HPC resources



Job Validation and Software Availability

- Using event generation and simulation job templates for initial validation tests
- Suitable baseline for now to flush out initial setup problems

Survey of all available ATLAS software retrieval methods to evaluate suitability for Archer:

CVMFS

- Standard CVMFS setup not possible on compute nodes
 - No WAN connectivity - possible use of port forwarding through to external proxy
 - FUSE (and CVMFS) software unavailable (for now?)
- Instead used scheduled *rsync* of selected releases
 - Modify handful of files with relative path dependencies

Pacman

- Cumbersome for maintenance of all ATLAS releases but works in disconnected setup

Parrot

- Initial attempts shown to work, but scalability has been mentioned as an issue at other HPC sites
- Used successfully at Piz Daint



Software Dependencies

- Cray Linux Environment (CLE) OS based on SUSE Linux – assumption/hope that SL(C) libraries are compatible, but may have to recompile some libraries as we go along

Environment Setup

- Non-system paths for binary and library dependencies have to be defined before the job is executed
- Alternative paths inherited from default Archer login environment have to be removed (e.g Python)
- Archer uses *TCL modules* to define library and application setup
- We can apply the same method to define ATLAS job environment setup
- Invoke module install at start of each job

```
conflict gcc
prepend-path PATH /opt/gcc/4.8.2/bin
prepend-path LD_LIBRARY_PATH /opt/gcc/4.8.2/snos/lib64
setenv GCC_PATH /opt/gcc/4.8.2
setenv GCC_VERSION 4.8.2
setenv GNU_VERSION 4.8.2
```

**Default module
installation for gcc**

ARC CE and Job Submission

- Installed latest version of ARC CE on dedicated server
- Located at ECDF (Tier-2 UK site)
- An ARC CE hosted on the Archer subnet will be difficult if not impossible
- We will need to pursue a disconnected setup at this stage
- Password-less ssh possible from ARC server to Archer gateway nodes

LRMS Communication via ssh-based methods

- ARC development already for SLURM but less community interest in PBS
 - Volunteered to develop similar functionality for PBS
-
- Could deploy an ARC Control Tower (ACT) at the site to pull in Panda workload and input data before job submission
 - Other alternative: ARC CE located at RDF Data Analytic cluster



Job Steering

- Some form of HPC job wrapper envisaged to handle HPC (and site) specific job environment setup
- Could be incorporated into a HPC pilot similar to US HPC sites
 - Network connectivity
 - WAN connectivity setup if available
 - Uninstall default modules that clash with ATLAS job setup
 - ATLAS software location and dependencies
 - Scheduling logic
 - Pilot execution
 - Parallel job submission to compute node (e.g. `aprun -n 24`)
 - Output staging
 - Move job output from Lustre FS to externally addressable storage

**HPC job wrapper
template**

- Can we just run standard pilots instead of defining a new job wrapper?
 - **Yes** - but very inefficient for us
- What resources should a pilot reserve?
 - Site policy is for each user to run a maximum of **16** queued jobs, **8** running jobs

Job Scheduling

- Job scheduling remains the biggest unknown in the proposed setup
- Different scheduling environment compared to HTC, even for shared facilities such as ECDF
- No clear solutions exist at present - the simpler the better

In-flight Heuristic Approach

- One possibility is to look at the number of resources available at the time of the pilot execution
- Titan using the showbf command (PBS/Maui/Moab scheduler) to determine what resources are available for immediate use
- Derive the backfill opportunity given the job runtime and resource requirements
- Request resource for (pilot) job based on transient queue conditions, project allocation and site policy
- Archer uses the Cray ALPS scheduler so the showbf command not available
- Investigating alternative approaches with Archer support



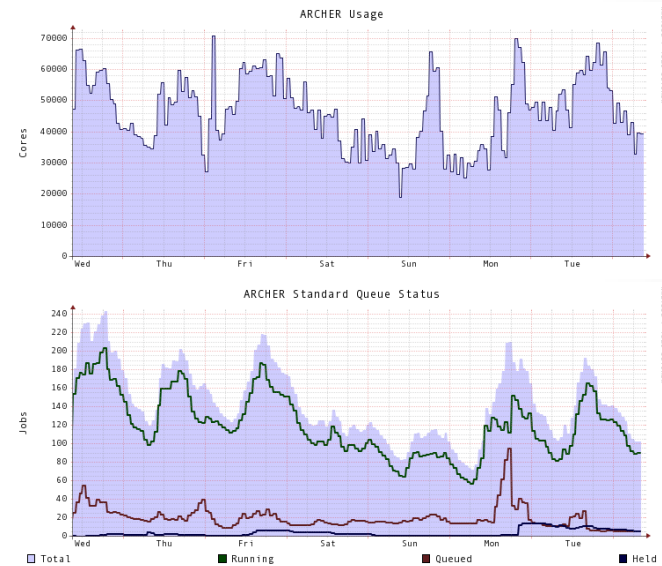
Job Scheduling

Job Efficiency

- Multi-core - but not multi-node - workloads could lead to risk of wasting paid-for resources if job termination depends on “speed of the slowest”
- Is there a tolerable failure rate?
- Event Service like model?
- Pre-emption / Application checkpointing?

Resource Allocation

- Pay for CPU quota and burn through allocation over time
 - Easier to manage and scale back usage when necessary
 - Need to use allocation as efficiently as possible
- Use low priority queue to leverage free cycles
 - Sluice gate during periods of low utilisation
 - Difficult to anticipate when this will be available at any given time



Archer Utilisation (1 week)



Outlook

- Continuing to work with Archer support to design best method to handle general ATLAS production activities
- Current focus is on persistent solution based on feedback from Archer and from exploratory work by other HPC sites
- Site restrictions are not static and changing due to demands of the “mainstay” users
 - More general outbound access being explored without need for port forwarding tricks
- VM instances close to RDF storage will enable feasibility of HPC-based SRM and CE
- Working with Archer support to develop advanced queue monitoring tools
 - Will benefit other HPC users and aid make decision on reservation requests for ATLAS
- Discussion with with ARC developers on ssh-based job submission and management backend for PBS

