# Use of Archer for Particle Physics

**Andrew Washbrook**
University of Edinburgh

Physics at Extreme Scales, Edinburgh
15th April 2014

# Motivation

- The Edinburgh Particle Physics group had access to a share of resources on HECToR facility - and now on Archer

**CHEP 2013 conference note:**
**Leveraging HPC resources for High Energy Physics**
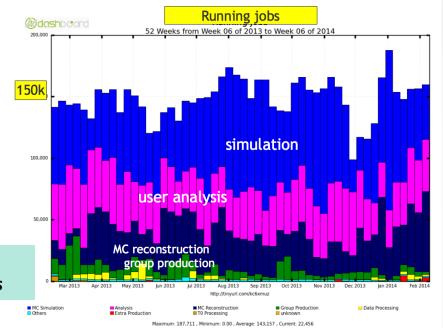http://indico.cern.ch/event/214784/session/9/contribution/438

- A feasibility study was performed to determine how an example HPC resource could be incorporated into a WLCG Tier-2 Grid site hosted at the same facility

- Now would like to move from HECToR feasibility studies a production-level service using Archer

- We have recently been looking at the issues and challenges in using the new Archer HPC facility for HEP-EX

# ATLAS Computing Context

- The ATLAS experiment at the LHC:
  - Processes and manages more than 130 PB of data
  - Uses more than 150k CPUs distributed across 100 computing centres managed by central workload management system (PanDA)

**ATLAS Dashboard 2013/14 Running Jobs**



- Run 2 (2015-2018) data processing will require a **lot** more computing and storage resources
- Can HPC and Leadership Class Facilities help with the increased demand?

- Several HPC sites in Europe and the US are working with ATLAS including:
  - Mira
  - Titan
  - Stampede
  - Hydra, RZG Munich
  - *Archer*

3

# HPC vs. High Throughput Computing

The following restrictions apply for HPC usage compared with traditional high throughput computing methods:
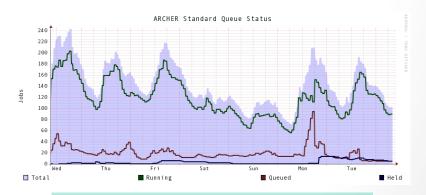
- Network access to and from HPC compute nodes more restrictive than Grid worker nodes
  - No WAN connectivity available
- OS deployed on compute nodes (CLE) is more lightweight than WN OS
  - Optimise code execution by limiting the number of interruptions to compute processes
  - Standard software libraries and packages not available
- No local disk on compute node
  - All job data is expected to reside on the shared filesystem
  - Not designed to cater for applications handling large input data sets and sustained I/O calls during job execution
- Separate identity management policy cannot be coupled to the federated systems we use on the Grid
  - All jobs submitted through my local account (for now)


- Restrictions are mostly driven by HPC user expectations rather than by strict technical barriers
- Exploring where adjustments to system configuration can be potentially adapted to accommodate ATLAS workloads

# Job Submission and Scheduling

- Deployed Grid Middleware services at our existing Tier-2 site (ECDF) to enable jobs from the Grid to be routed to Archer

- ATLAS software not currently suited for MPI-type jobs **but** can efficiently process *multi-core* workloads
- Submit single HPC job can steer hundreds of *wholenode* jobs
- Other options explored include offloading critical sections of workload well suited for HPC resources
- Job resource request size can be adapted to queue conditions
- Backfilling could generate slots for HEP-EX use without loss of service to other HPC users

**Archer Utilisation (1 week view)**

**Archer Queue Status (1 week view)**

# Outlook

- Aim to provide a production level HPC service in concert with local HEP-EX Tier-2 operations at ECDF in Edinburgh

- Currently investing time and effort into a building a robust setup and to resolve compatibility issues

- Previous experience with running LHC software at a shared cluster facility (ECDF) is proving useful

- Novel solutions will be required to fit the computing environment expectations from ATLAS and other HEP experiments

- Incorporating ideas and solutions from other HPC facilities in the US and Europe rather than working in isolation

- Edinburgh and Archer are well placed to contribute in this area