

GridPP

UK Computing for Particle Physics



THE UNIVERSITY
of EDINBURGH

Experiences with Tier-2 operations on shared university resources

Andrew Washbrook

University of Edinburgh

GridPP 34

29th April 2015

ECDF and Eddie

Eddie Mark 2

- Phase 1 - 130 x IBM dx360M3 iDataPlex servers (2 x Xeon E5620 quad-core)
- Phase 2 - 156 x IBM dx360M3 iDataPlex servers (2 x Xeon E5645 six-core)
 - ~3,000 cores
- GPU and large memory systems
- Single queue for single core and multi-core workloads

Eddie Mark 3 - **Available from August 2015**

- Now tendering for £1M of new equipment
 - We had early input into machine specification
- Expected to be similar in scope to Eddie Mk2
- Similar operational model
 - "Free" at point of use
 - Paid-for jobs have higher priority
 - Opportunistic use encouraged
- Hosting service for additional compute purchased by university research groups
 - Spare rack capacity for bespoke equipment

ACF Hosting

- Equipment hosting provided by Advanced Computing Facility (ACF)
- ACF provide:
 - Infrastructure management
 - Power
 - Cooling
 - Security
 - Routine system tasks (e.g. disk replacement)



- All our Grid middleware servers are located in the same machine room as the Eddie cluster
- Rely on remote server management tools (e.g. idrac)
- Occasional site access needed for server maintenance

ECDF Customer Base

- ECDF provides computing resources across the university for:
 - Physics
 - Geoscience
 - Engineering
 - Life sciences
 - Veterinary medicine
 - Informatics
 - Biology

Support and feedback

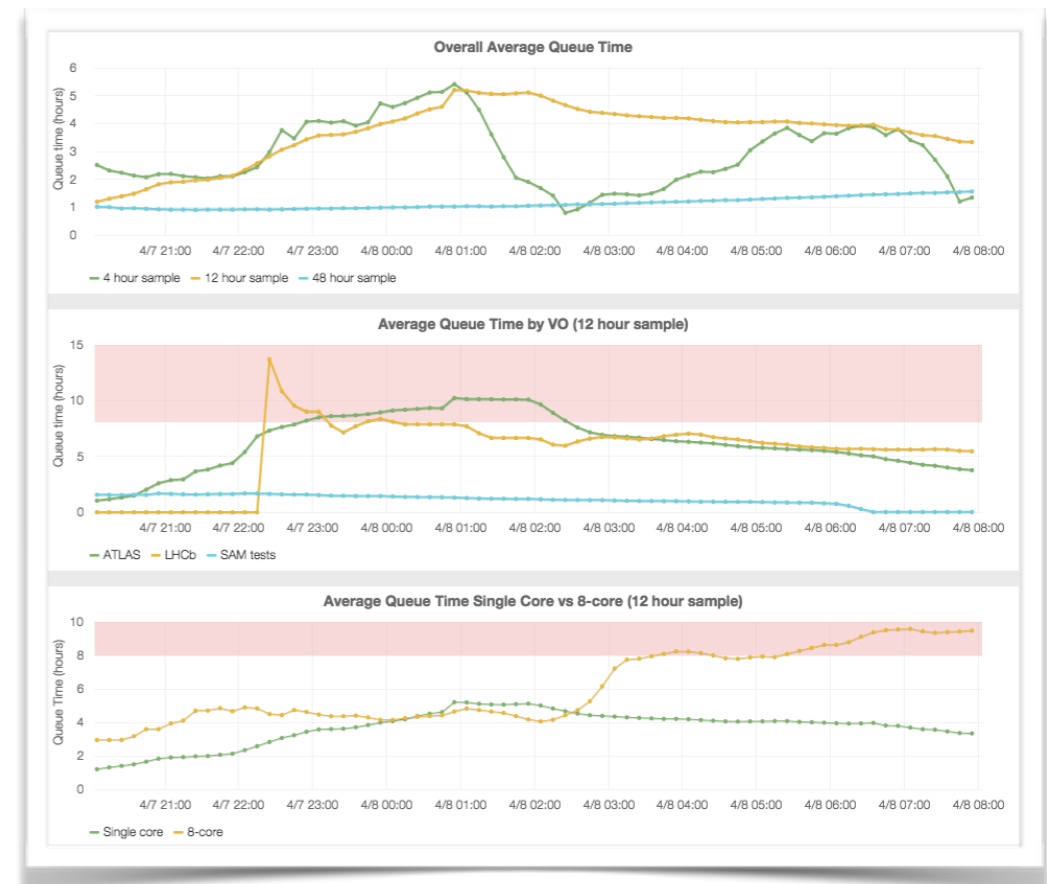
- Ticketing system manages incident response calls and simple change requests
 - Ongoing collaboration with groups that have bespoke requirements (i.e. GridPP)
 - Software troubleshooting and optimisation
 - Regular drop-in sessions
-

Shared Facility Benefits

- We don't have to care about:
 - Cluster and batch system setup and configuration
 - Continual equipment maintenance
 - System wide troubleshooting
- Leveraging of opportunistic resources
 - This benefits GridPP - in principle we always have work to process 24/7/365
- Better resource size to FTE ratio

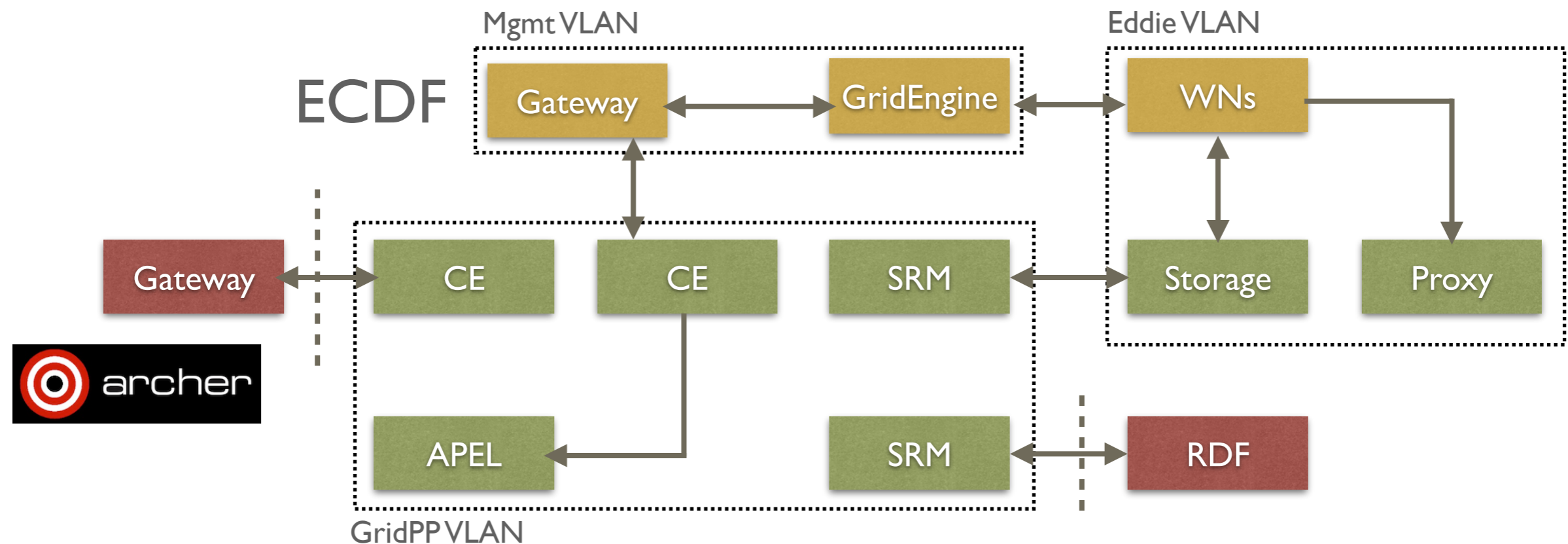
Mutual Benefits

- Collaborate on troubleshooting
 - We are more sensitive (and responsive) to job failures
 - Issue resolution generally benefits other projects
- Increase overall resource utilisation by continual submission of a steady stream of jobs
- Development opportunities
 - Recent CHEP work on grid site-oriented analytics being extended to provide coverage for all ECDF customers
 - Many-core devices - provided Xeon Phi experience
 - Cloud and Virtualisation experience from the GridPP and LHC communities



Middleware Deployment

- Grid middleware services are decoupled from the ECDF internal network
- Use passwordless ssh to interact with GridEngine (qsub, qstat) via Gateway servers
- Experience using this model has enabled us to branch out to other resources such as RDF and Archer



Operational Conflicts

- We underuse the native high performance storage (GPFS)
- Relatively small worker node local storage given ATLAS requirements
- Hyper-threading was not palatable for all ECDF projects

Middleware Compatibility

- A shim of configuration and hackery is needed for most services
 - Persistent definition issue for site information publishing
 - "How many cores do you have?"
 - Do not have full control of worker node configuration
 - Use "tarball" method where applicable
 - Have to (gently) push back on services requiring worker node admin access
-

Operational Challenges

Worker Nodes

- Change management is non-trivial compared to other sites
- Not as agile due to the diverse customer base
 - Must not impact on jobs from other non-Grid users
 - All packages have to be evaluated by site administrators before deployment
- Some steps are unavoidable for us
 - CVMFS package and configuration updates
 - Security updates

Fairshare scheduling “burstiness”

- Occasionally completely throttled by fairshare model
- Fairshare tree needs continual tuning and pruning

Incident Response

- Rely on issue resolution by ECDF systems team in some circumstances
 - They are generally quick to respond - but it is out of our hands
-

Cloud Provisioning

- Production cloud service pilot planned for new Eddie Mark 3 cluster
 - Proposed bare-metal/cloud hybrid model
 - Ability to rapidly switch worker nodes between cloud and traditional batch system based on user demand
 - Primary motivation from Biology - software pipeline exclusively in Biolinux OS

 - How does this fit in with future GridPP operations?
 - Some of our operational challenges could be alleviated with a cloud-based solution
 - Greater control over worker node environment
 - Is there a risk of over-engineering a solution to match the proposed hybrid model?
 - I am signed up as an early adopter
-

Reflections

- Running a Tier-2 Grid site on a shared university facility has been shown to work
- Continual effort required from **both** parties to keep the site up and running

Working Relationship

- A lot of our requests are unique compared to other projects running on the facility
 - Evolving practices - we have no script to work from
 - Other shared facilities have different MoUs
 - Not the role of the facility to understand the idiosyncrasies of Grid computing
 - Correspondence through ticketing system only works for incident response
 - Essential to have a direct line to the system administrators
 - In contact at least once a week, often daily
 - Regular face-to-face reviews very useful

 - Possible move to a hybrid bare metal/cloud model may help to harmonise our site operations with other Tier-2 sites
-