

ATLAS Analysis on Cloud Facilities

Fernando Barreiro Megino (University of Texas at Arlington)

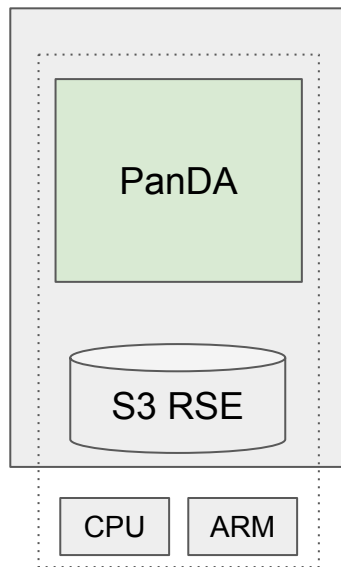
Analysis Facilities Forum Kick Off Meeting, 25 March 2022



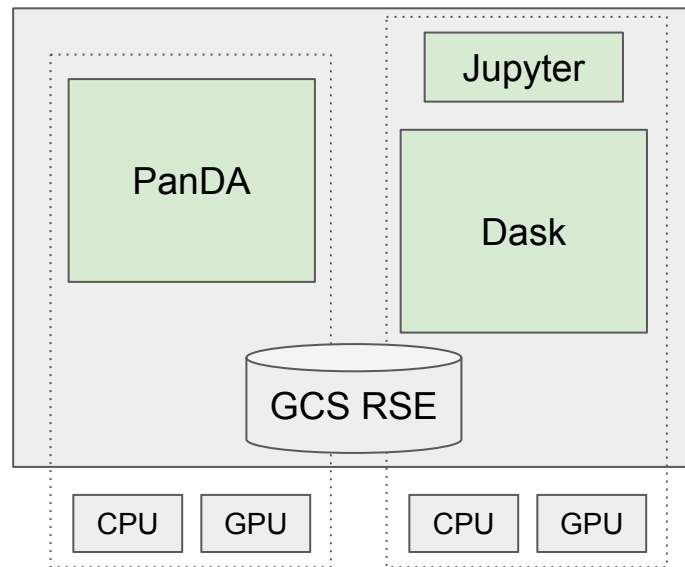
UNIVERSITY OF
TEXAS
ARLINGTON

What's available?

Amazon (CSU Fresno grant)



Google (USATLAS grant)

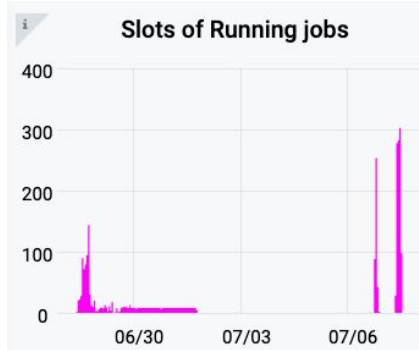


- + Users can explore Amazon/Google infrastructure and services on their own: FPGA/GPU/ARM/XXL nodes, cloud AI platforms...

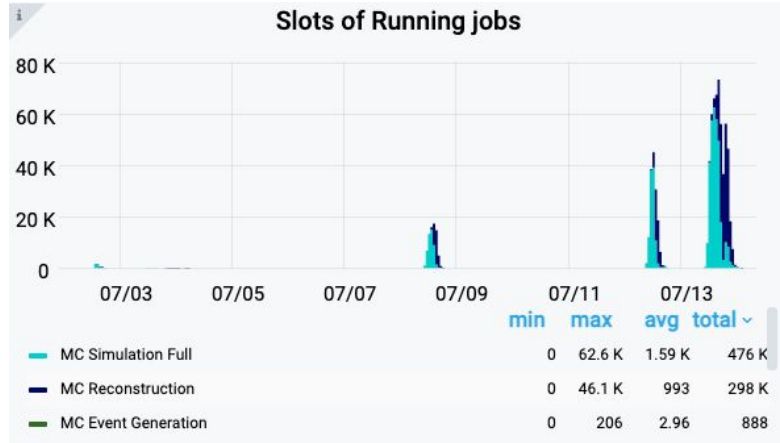
What's behind the setups?

- Kubernetes: no intermediate layers and installing only necessary software
- PanDA/Harvester submits jobs directly to Kubernetes
- Jupyter & Dask installed via Helm chart
- CVMFS plugin for Kubernetes is the weakest part in the setup
- Autoscaling: clusters scale according to usage
 - Sometimes even Google or Amazon might be out of a specific resource
- Clusters optimized for cost: standard vs preemptible nodes
- Rucio is integrated with cloud storage through signed URLs.
 - Users have to request permission to be able to use cloud storage through Rucio (cost protection)

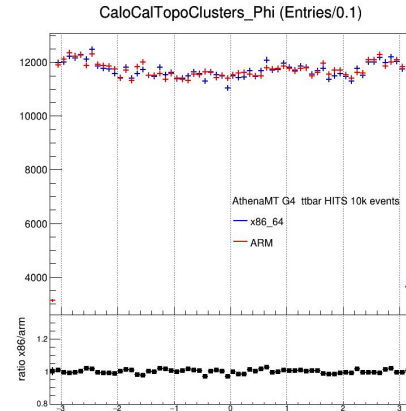
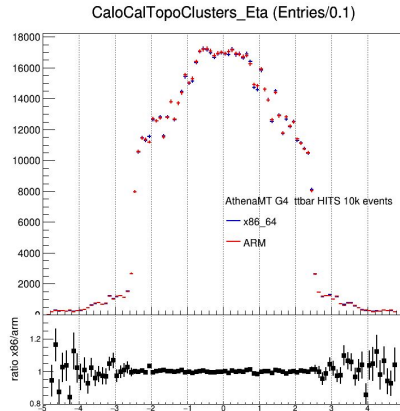
PanDA examples



Analysis of 100% of PHYSLITE dataset
(with Nikolai Hartmann)



Simulation chains for Active Learning
(with Lukas Heinrich)



First 10k AthSimul HITS events: ARM (Amazon) vs x86 (Prague)
(with Johannes Elmsheuser)

Now: Jupyter/Dask walkthrough

(See Nikolai Hartmann's demo for more advanced Dask usage: <https://indico.cern.ch/event/1135251/>)

Backup

Future thoughts: from R&D to production

- Very promising results in batch and interactive
- Improve CVMFS plugin
- Cluster settings and setup
 - Determine final cluster settings
 - Programmatic creation of clusters and simplify installation
- Controls for interactive clusters
 - Usage accounting, monitoring and limits for Dask/Jupyter
 - Housekeeping of Dask clusters
 - Customize Spawner
- Improve images (CVMFS, RAPIDS) and image management
- Interactive/local data management and storage
 - Compare to SWAN's EOS integration
 - Are independent user disks the best option? How can they be migrated?
 - How to access data from GCS? Are Rucio signed URLs user-friendly enough?
- What are the distributed frameworks to support?

Notes on DDM-Cloud integration

- Rucio and FTS are able to manage 3rd party transfers and direct download/upload through signed URLs
 - ⇒ Having an S3 or GCS RSE is possible
- Main challenge: WLCG relies IGTF CA certificates
 - Google and Amazon are not part of IGTF, i.e. their CAs are not trusted by WLCG sites
 - Workarounds based on load balancers or “friendly” sites installing cloud certificates
- Egress costs

PanDA GKE setup

GOOGLE100
Optimized for
PHYSLITE exercise
and IO intensive
jobs. Moderate scale
expected



GOOGLE_BULK
Optimized for Active
Learning exercise
and simulation jobs
at large scale



PanDA K8S cluster

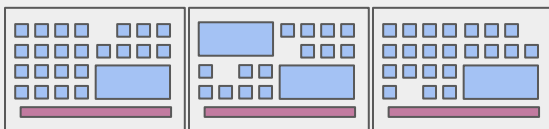
SSD static, preemptible node pool

- 1 node: 8 vCPU, 32 GB RAM, 375 GB SSD



SSD autoscaled, preemptible node pool

- Autoscaled nodes: 32 vCPU, 128 GB RAM, 1.5 TB SSD



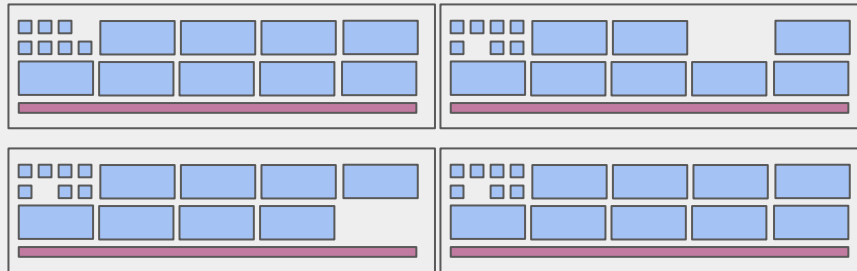
HDD static, preemptible node pool

- 1 node: 8vCPU, 32 GB RAM, 160 GB standard pdd



HDD autoscaled, preemptible node pool

- Autoscaled nodes: 80 vCPU, 320 GB RAM, 1.6 TB standard pdd



4 Frontier squids with 2 vCPU,
16GB RAM, 200GB disk

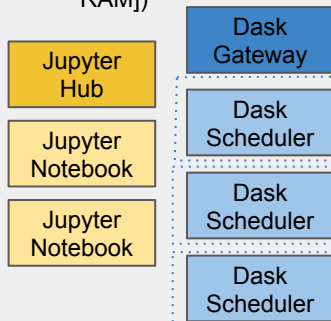


Jupyter/Dask setup on Google

Jupyter/Dask K8S cluster

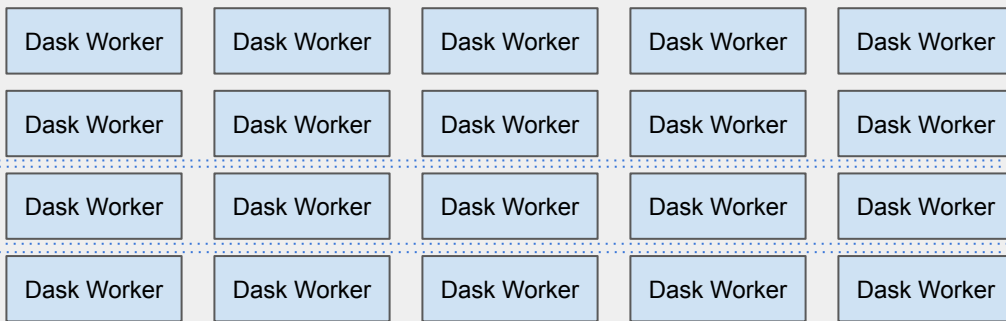
Critical node pool

- Non-preemptible nodes
- Static size (curr. 2 nodes x [16 vCPU, 64GB RAM])



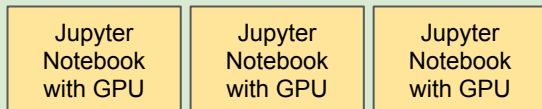
Worker node pool

- Preemptible nodes
- Autoscaled (up to hundreds of nodes x [8 vCPU, 52 GB RAM])



Critical GPU node pool

- Non-preemptible nodes
- Autoscaled (curr. up to 3 nodes x [8 vCPU, 52GB RAM, 1 GPU])
- Node is fully dedicated to one user



Worker GPU node pool

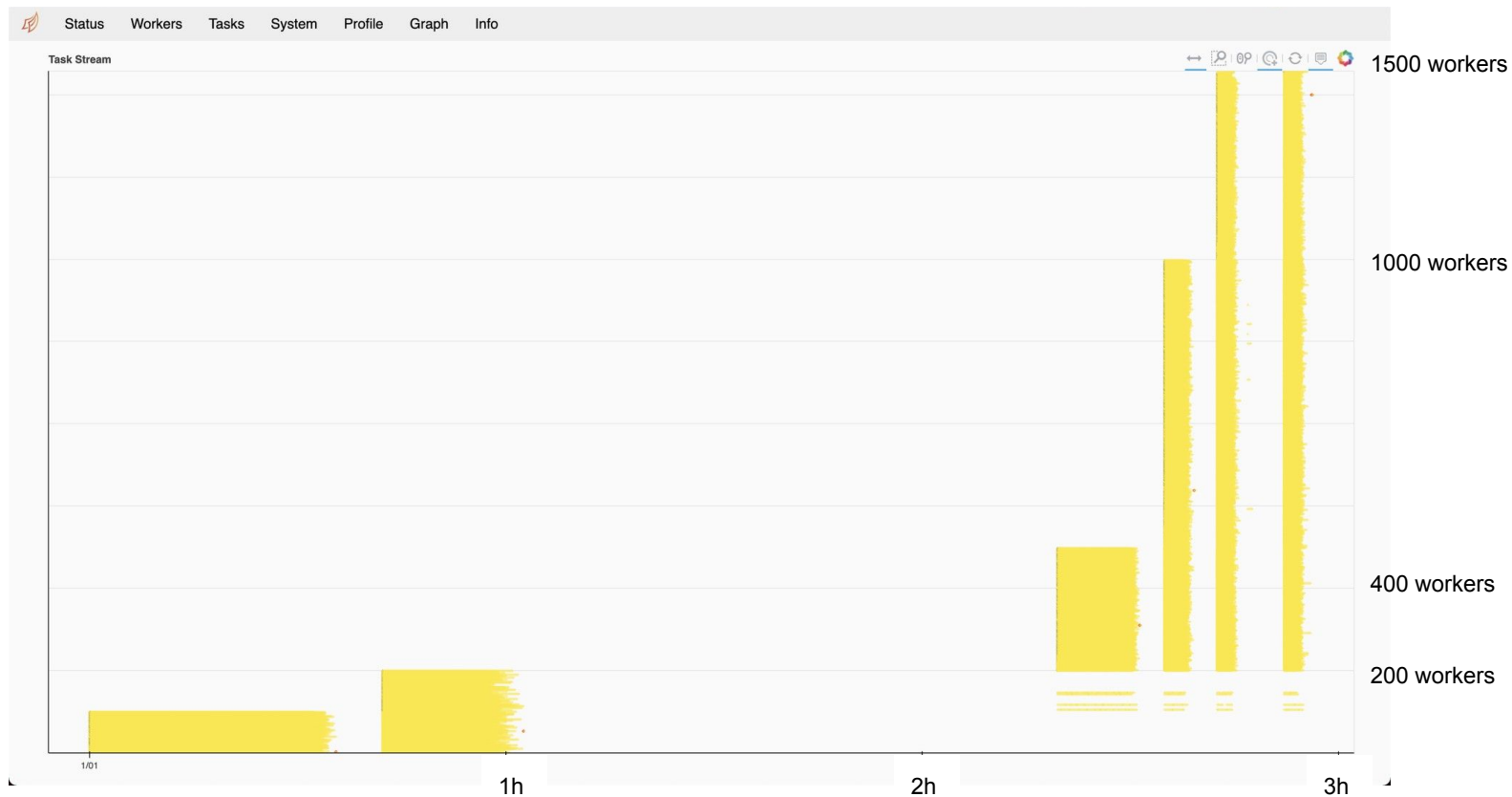
- Preemptible nodes
- Autoscaled (curr. up to couple of nodes)
- Has not been used yet



Scaling Dask

(Lukas Heinrich)

Task stream profile at various cluster sizes



See <https://indico.cern.ch/event/1131909/> for a Dask demo