



# Analysis in the Cloud: Google Cloud Platform

Johannes Elmsheuser (BNL)  
on behalf of the ATLAS/Google R&D team

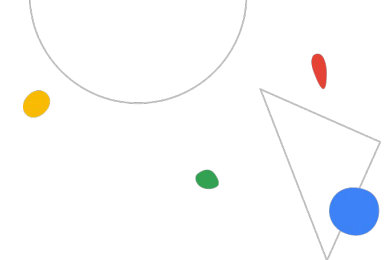
27 October 2020

Future Analysis Systems and Facilities, IRIS-HEP Blueprint Workshop



# ATLAS/Google R&D in 2020

November 1, 2019 - October 30, 2020



## Goals

- Setup US ATLAS - Google technical and management team to work together
- Test **Google Cloud Storage** as an additional component to ATLAS storage & compute on the grid
- Evaluate, test and demonstrate how Google Cloud can be **integrated with PanDA+Rucio**
- Explore ways to provide fast, reliable and easy **access to data for analysis**
- Evaluate **new data formats** for physics analysis
- Use **Google Kubernetes Engine** for HPO and GKE integration with ATLAS grid via PanDA WMS

## Tracks (as initially defined in 2019)

- Track 1: Data Management Across **Hot/Cold Storage**
- Track 2: **Machine Learning**, TPU vs. GPU for GNN training
- Track 3: **Optimized I/O** and data formats for object storage
- Track 4: **End user analysis** conducted worldwide at PB scale
- New Track: LSST/Vera C. Rubin Observatory

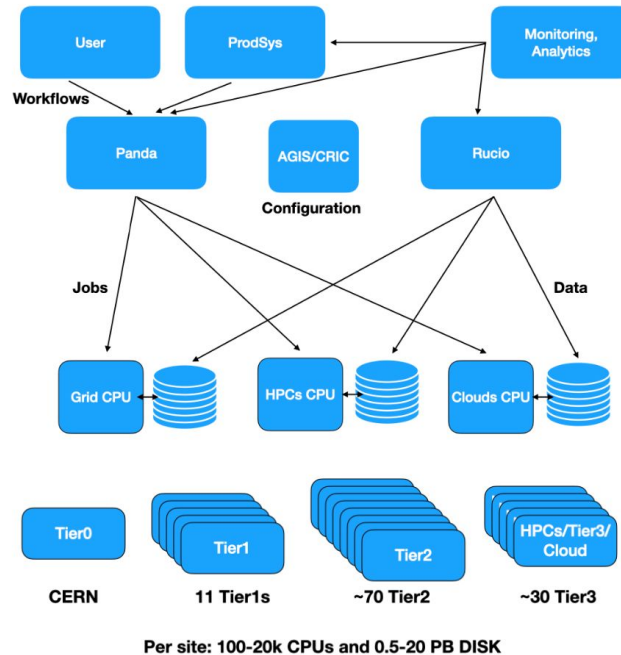
*Focus in this talk*

# ATLAS DISTRIBUTED COMPUTING OVERVIEW



The ATLAS distributed computing system is centered around:

- **Workload management system:** PanDA
- **Data management system:** Rucio
- **Many additional components:** AGIS/CRIC, ProdSys, Analytics, ...
- **Resources:** WLCG grid sites, Tier0, HPCs, Boinc, Cloud
- **Shifters:** Grid, Expert and Analysis (ADCoS, CRC, DAST)
- **Runs 24/7 all 365 days per year**

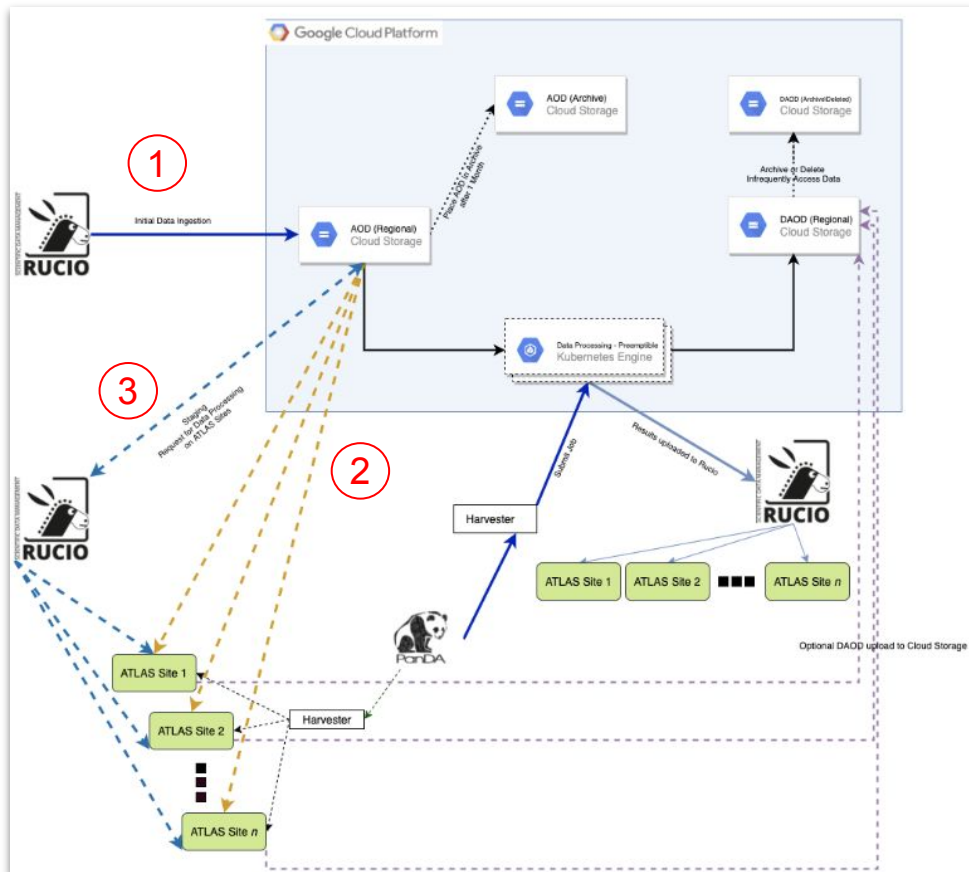


Why ?

- Run2 data15-18 are approx.  $20 \times 10^9$  data +  $40 \times 10^9$  MC events - assume 50 kB/event = 3 PB of single DAOD type
- For Run3 (202224) expect slightly more data, but for HL-LHC (2027-...) expect much more data
- Assume a data format with small event size: store and process fraction of data in GKE/GCS by individual analysers through PanDA pruned jobs or interactive

Fit Google compute and storage into Rucio and PanDA ecosystem as a Cloud site and explore interactive analysis possibilities

# Step 1: GCS storage integration in Rucio



- 1 Transfer from Grid Sites to GCS
- 2 Transfer from GCS to Grid sites
- 3 Data Carousel mode

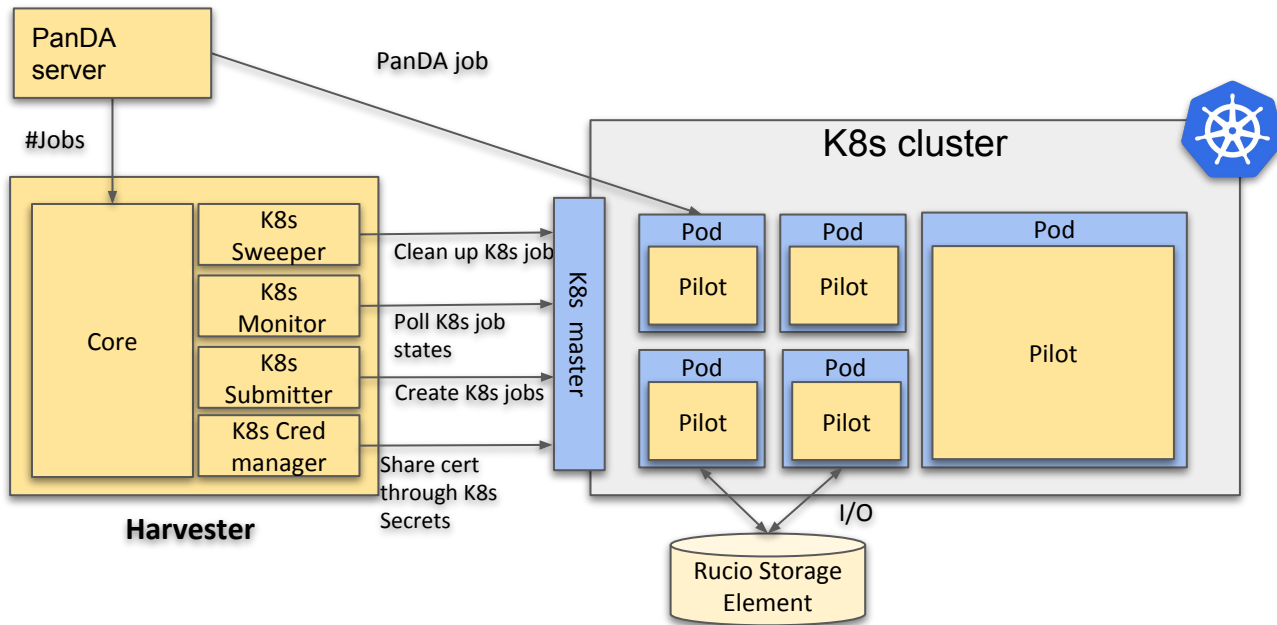
- Google Cloud storage setup as Rucio storage with **3rd-party-copy FTS transfers**
- **Fully validated** at 10 grid sites with transfers up to 15 GB/s over hours - longer term have Google CA cert in IGTF
- Direct downloads from Google to Grid worker nodes possible but blocked at a few sites
- Transfer from Google storage uses most probably GPN - might cause troubles to non-HEP activities
- Future large scale tests put on hold due to **large egress costs**

# Step 2: Google Kubernetes Engine + PanDA

Run ATLAS G4/Fast simulation with storage at CERN

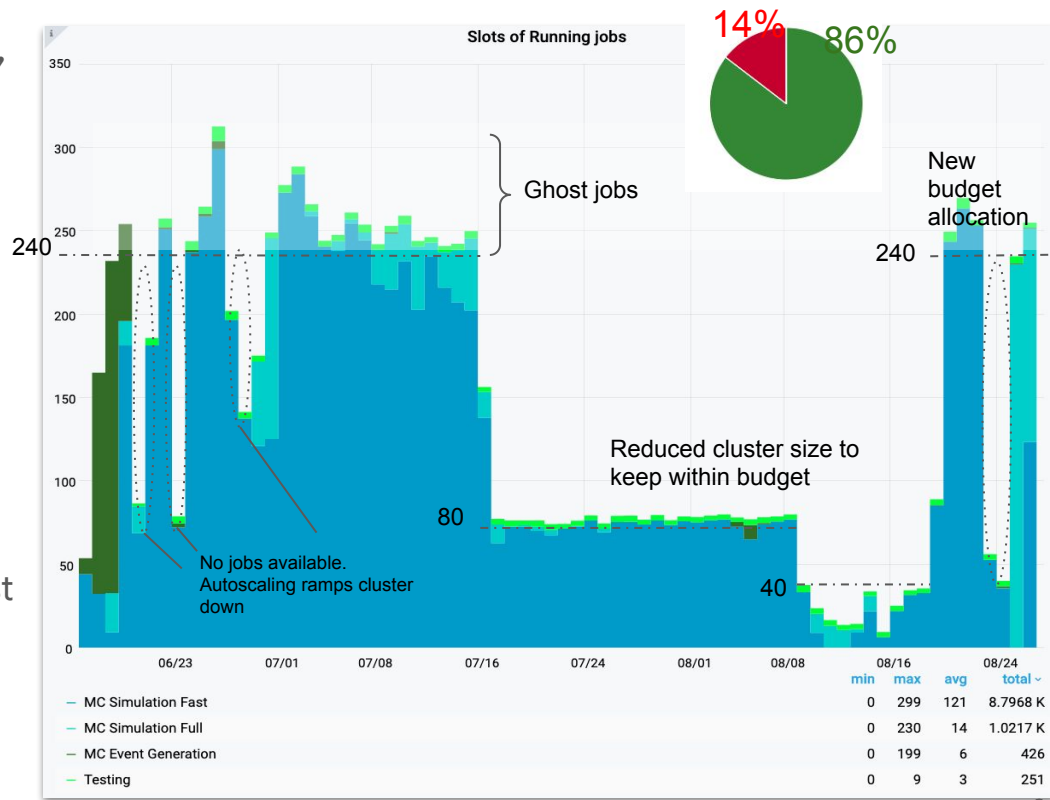
- GKE setup for compute
- Very light I/O jobs

- **CVMFS:** Installed through daemonset + k8s volumes
- **Frontier Squid:**
  - Installed on dedicated VM for this exercise
  - Now also possible to install internally in the K8s cluster



# Step 2: GKE + PanDA running simulation

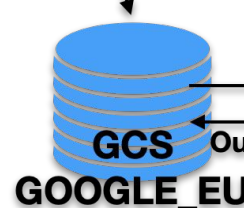
- Limited to Simulation (low I/O) jobs, since storage at CERN
- Preemptible nodes
  - Causing most of the failures
  - Limiting job duration to <5 hours
  - Attractive deal: 80% cost reduction, slightly higher failure rate
  - Good failure rate considering pre-emption effect. Limiting jobs to <5h duration, short jobs not always available
- Autoscaled cluster
  - Cluster ramps down and lowers the cost when no jobs queued
- Costs (remote storage, 120-160 cores)
  - July: 2.3k USD/month (76.6 USD/day)
  - Aug: 1.67k USD/month (54.4 USD/day)



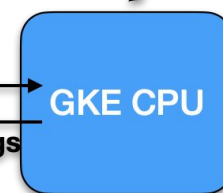
# Step 3: User Analysis with GKE and GCS



DAOD dataset



DAOD input file

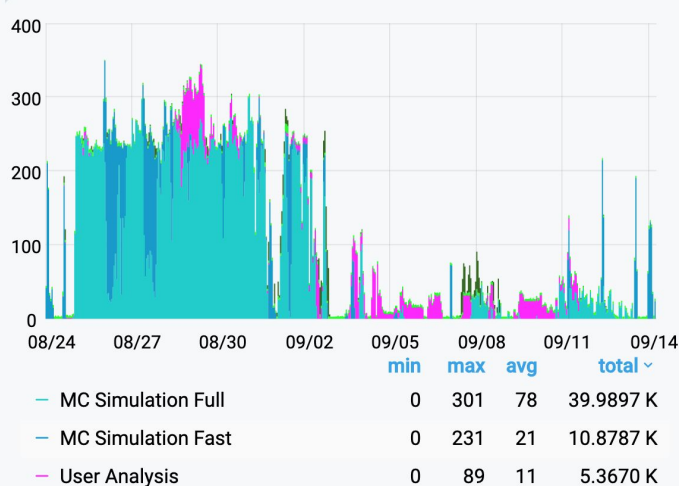


Output files/logs

How ?

- Replicate 1 TB DAOD\_SUSY5 dataset to GCS Rucio storage GOOGLE\_EU
- Run regular ATLAS analysis submitted pruned to PanDA with/without systematics (30 min/20h)
- Store outputs back to GOOGLE\_EU

Slots of Running jobs



# Step 3: GKE/GCS user analysis

What works:

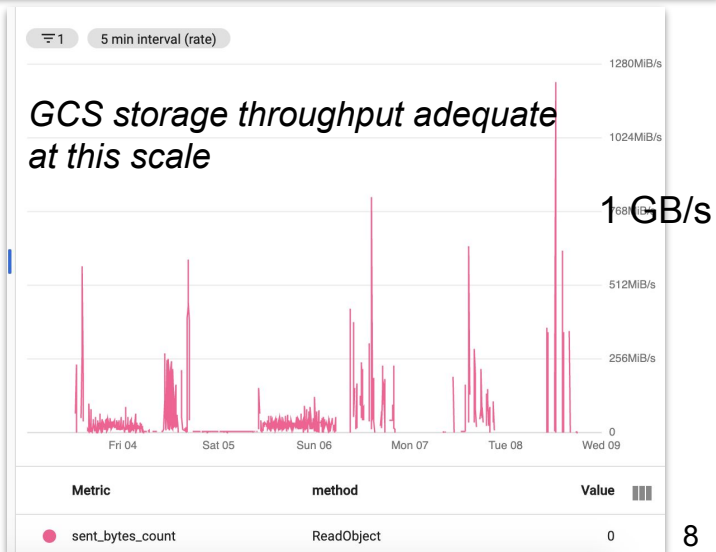
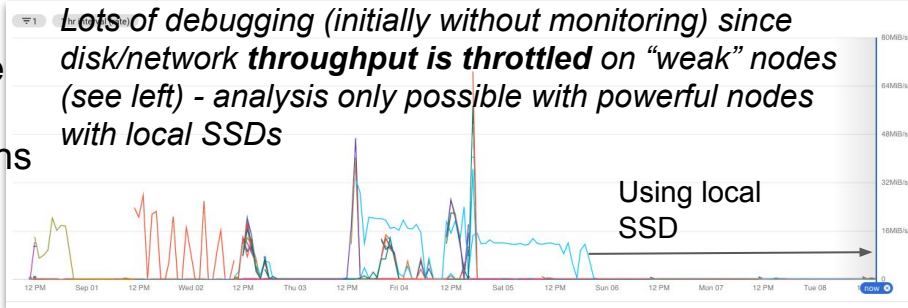
- **Successful GKE/GCS integration for the first time with full Rucio/PanDA workflow**
- PoC for analysis works stable after extensive iterations of GKE node setup with copy-to-scratch input
- Essential to use powerful well connected GKE nodes
- Usage of preemptible nodes seems only useful for workloads under a few hours

What does not work (so far):

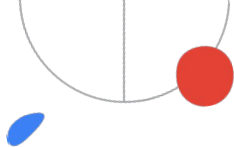
- “weak” GKE nodes
- ROOT direct I/O via DAVIX access of inputs to GCS broken - could avoid parts of local storage troubles

ToDo:

- Detailed cost estimation per wall clock hour or processed TB
- Scale to really large datasets, many users and long payloads







# Further tracks

## Track 3: **Optimized I/O and data formats** for object storage

- Focus on python-based analysis (batch or interactive)
- DAOD\_PHYSLITE input (calibrated physics objects) on GCS storage
- Data conversion for efficient use (e.g. ROOT to parquet)
- Conversion of simple branches (leaf) is straight forward C++ object conversion requires more work and thought.

## Track 2: **Machine Learning**, TPU vs. GPU for GNN training

- Detailed comparison of GPU V100 at NERSC Cori vs. Google TPUs - see CERN IML workshop ([link](#))
- Key metrics: Accuracy, Latency, Cost, Heat dissipation

## Track: **LSST PanDA+Rucio & Google Cloud Storage**

- Reuse ATLAS+Google PoC tools: Harvester plugin for landing payload to GKE clusters
- Reuse Cloud LSST features: Data accessible globally with S3 protocol and direct IO, Metadata accessible globally using Postgres DB
- *The Vera C. Rubin Observatory signed a contract with Google Cloud for the Interim Data Facility that provides production data management and processing. It is a multi-year contract to support their data processing needs*

# ATLAS/Google R&D summary and plans



## Summary

- Integrated Google Cloud Storage as Rucio storage element - offers automatic way to transferring large datasets
- Integrated Google Kubernetes Engine as PanDA queue - offers scalable way to process large datasets
- PoC of full ATLAS analysis with PanDA/Rucio within GCS/GKE
- Other means of processing data through GCS/GKE underway (and use e.g. GPU/TPU)
- Easy(er) integration with other Clouds or HEP communities
- Many easy to use Google services: Bigquery, AutoML, ...  
-> added value for an Analysis Facility in the Cloud

## Next R&D plan:

- Focus on analysis facilities and interactive analysis using compact data formats
- Transfer experience to other Clouds

# ATLAS & Google Cloud Team

Alexei Klimentov, Brookhaven National Laboratory

Kaushik De, University of Texas Arlington

Fernando Barreiro Megino, University of Texas at Arlington

Johannes Elmsheuser, Brookhaven National Laboratory

Mario Lassnig, CERN

Cedric Serfon, BNL

Misha Borodin, University Iowa

Tobias Wegner, University of Wuppertal

Xianyang Ju, Lawrence Berkeley National Laboratory

Paolo Calafiura, Lawrence Berkeley National Laboratory

Andy Hanushevsky, SLAC National Accelerator Laboratory

Ricardo Rocha, CERN

Siarhei Padolski, Brookhaven National Laboratory

Doug Benjamin, Argonne National Laboratory,

Karan Bhatia, Google Cloud

Ema Kaminskaya, Google Cloud

Miles Euell, Google Cloud

Usman Qureshi, Google Cloud

Ross Thomson, Google Cloud

Kevin Jameson, Google Cloud

Dom Zippilli, Google Cloud

Kevin Kissel, Google Cloud

