# Consolidation of Cloud Computing in ATLAS

Ryan P Taylor,[1] Cristovao Jose Domingues Cordeiro,[2] Domenico Giordano,[2]
Alessandro Di Girolamo,[2] John Hover,[3] Tomas Kouba,[4] Peter Love,[5]
Andrew McNab,[6] Jaroslava Schovancova[2] and Randall Sobie,[1]
on behalf of the ATLAS Collaboration

[1] University of Victoria
[2] CERN
[3] Brookhaven National Laboratory
[4] Czech Academy of Sciences
[5] Lancaster University
[6] University of Manchester

## Sim@P1 Improvements

The Sim@P1 project enables the ATLAS High-Level Trigger farm to be used for offline production activities. Using Openstack, up to 70,000 cores in 2,300 compute nodes are exploited for running primarily CPU-intensive jobs, such as event generation and Monte Carlo production, when not needed for TDAQ purposes.
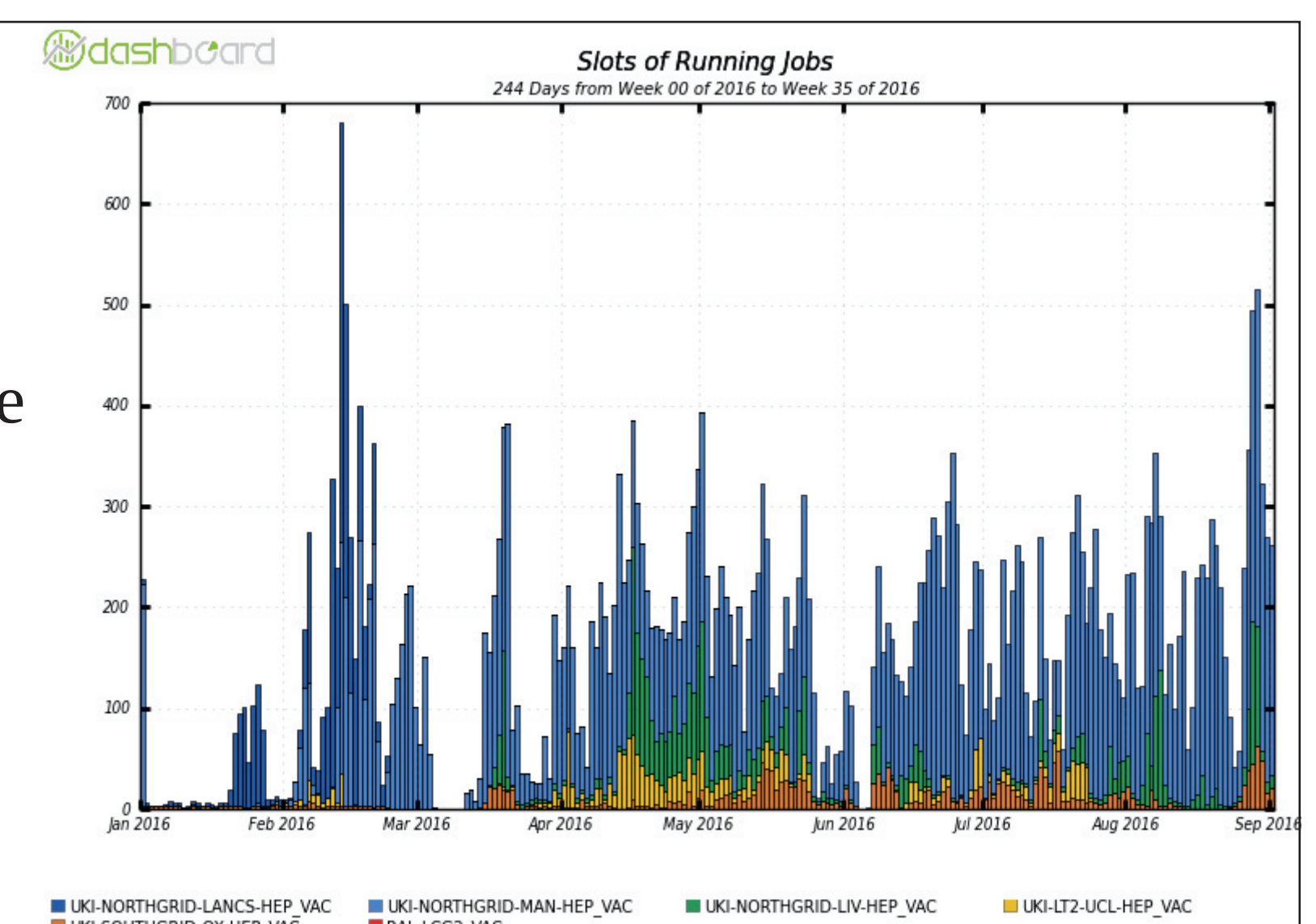


Above: cores used by Sim@P1 over the last year. About 5,000 cores are generally available for use, with higher opportunistic usage during technical stops, etc.

Recent developments have been made on a toolkit for fast switching between the TDAQ and Sim@P1 modes of operation, as well as improved resource partitioning and monitoring. Enhancements in resource provisioning now allow the entire farm to be ramped up for running jobs in as little as two hours. Testing of Event Service pilots is also underway.[1]

## Vac Integration

In 2016, the virtual machine (VM) definition for the Vacuum platform[2] was rewritten from scratch to use the industry-standard cloud-init format for contextualization. This allowed common elements of the contextualization recipe used across cloud platforms in ATLAS to be consolidated.
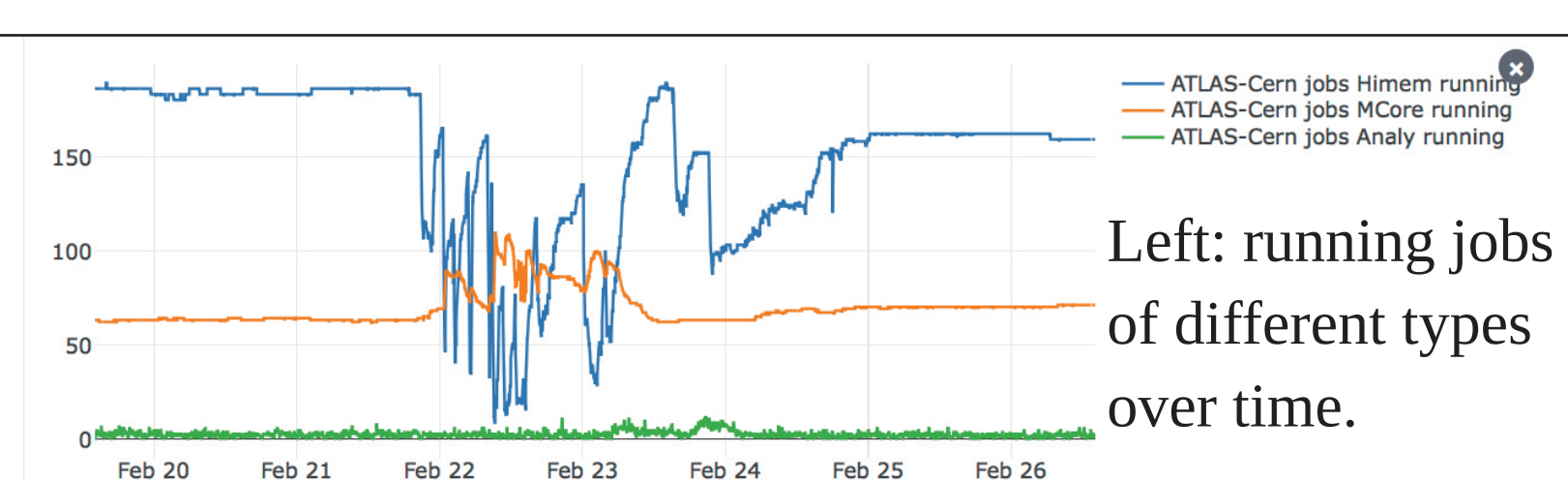


Above: Number of ATLAS jobs running on six Vac sites in the UK over the last year, peaking around 700.

Moreover, VM instances now run as HTCondor workers rather than directly retrieving payloads from PanDA. This aligns the workflow with that of other cloud and batch systems used in ATLAS, while retaining the feedback-based philosophy of the vacuum model. Future work is planned to continue converging the various ATLAS virtual machine definitions.

## Cloud Monitoring

When operating several HTCondor/Cloud Scheduler instances, each with several



Left: running jobs of different types over time.

clouds, job types, and VM types, it is essential to have a unified monitoring interface to present a coherent view of all VMs and jobs.[3]

The Cloud Monitor stack accomplishes this using standard tools such as Ganglia, ELK, Graphite, Grafana, Sensu, Plotly.js and Flask. All monitored metrics can be visualized in time-series graphs in a flexible manner. One can select any VM to see the jobs which ran on it, and examine the job logs as well.
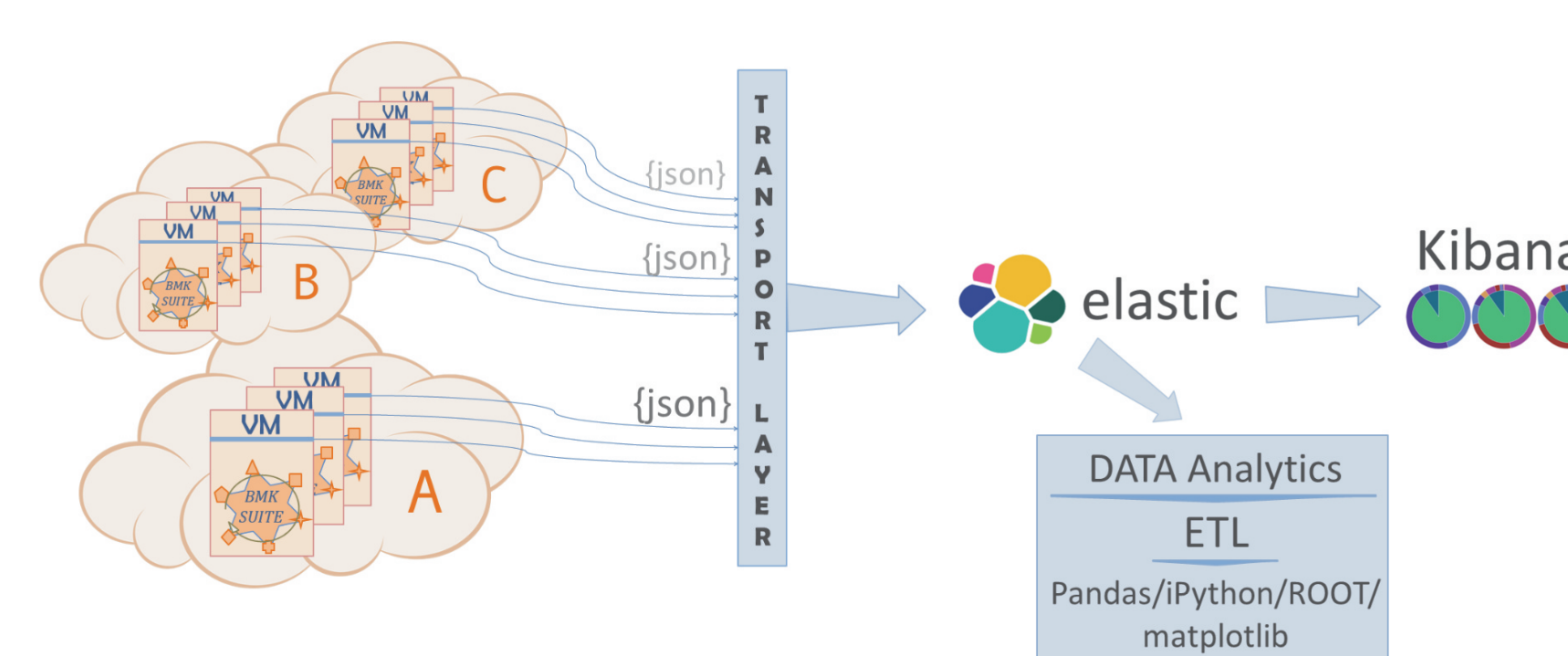


Above: the current quantity of VMs and jobs in different states on one cloud.

## Cloud Benchmarking Suite

The Cloud Benchmarking Suite[4] is a configurable framework for running a selection of benchmarks, and transporting, storing and analysing the measured results. It enables users to profile the performance of computing resources and quickly identify possible issues, in a unified way.
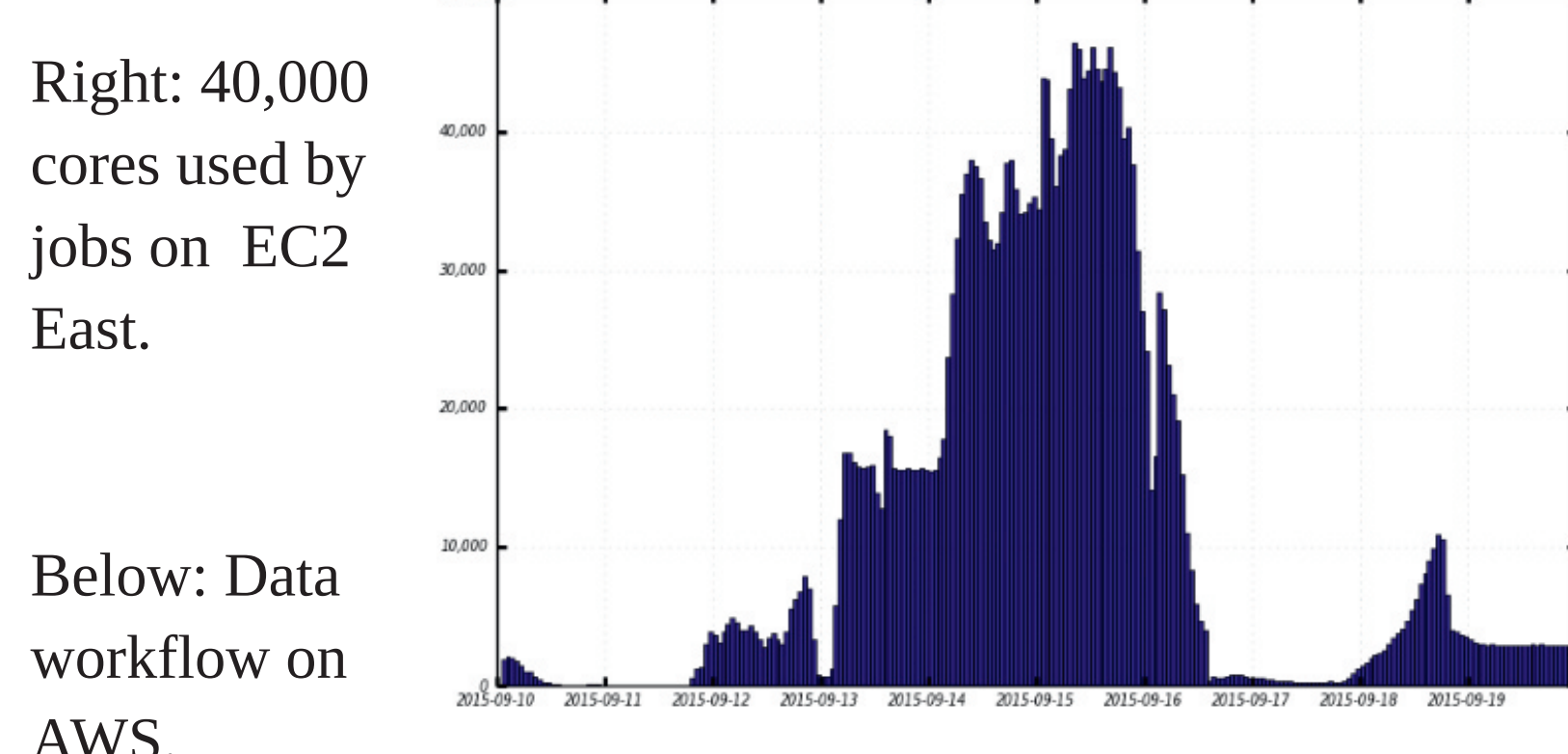
The benchmarking suite is being used by ATLAS to quantify the performance of every VM on several clouds, in HS06-equivalent units. A fast benchmark runs at bootup, and the result is joined with the delivered walltime and CPU time for each VM, to generate HS06-normalized accounting data.



Left: architecture of the cloud benchmark system. Results are transported using ActiveMQ, stored in ElasticSearch, and visualized with Kibana.
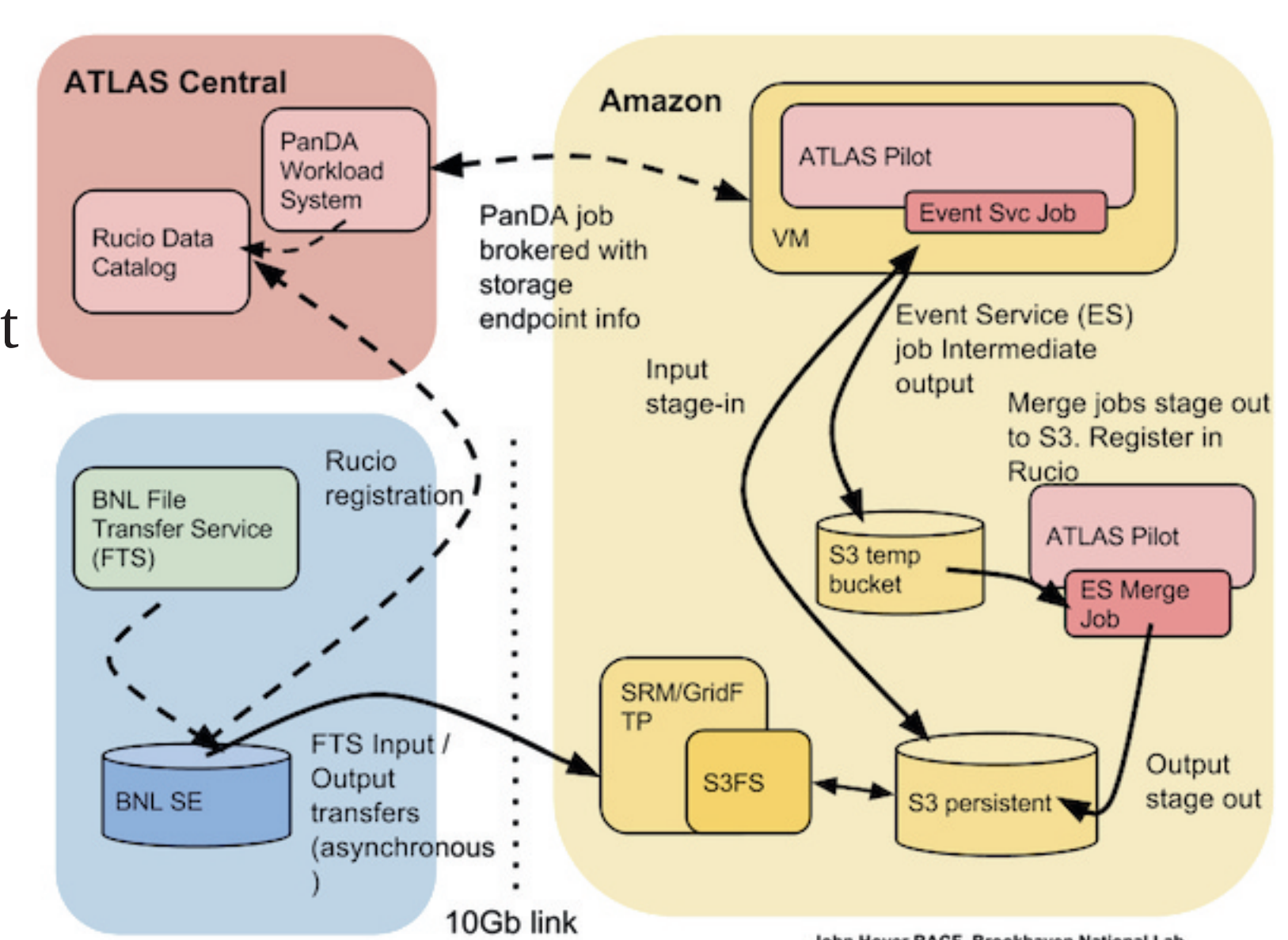
## Amazon Scale Testing

Building on successful scaling tests in 2014 of Amazon EC2 at the level of 20,000 cores,[5] a 40,000 core run was conducted in 2015. High-bandwidth direct network peering and S3 integration have been established to enable running at this scale.



Right: 40,000 cores used by jobs on EC2 East.
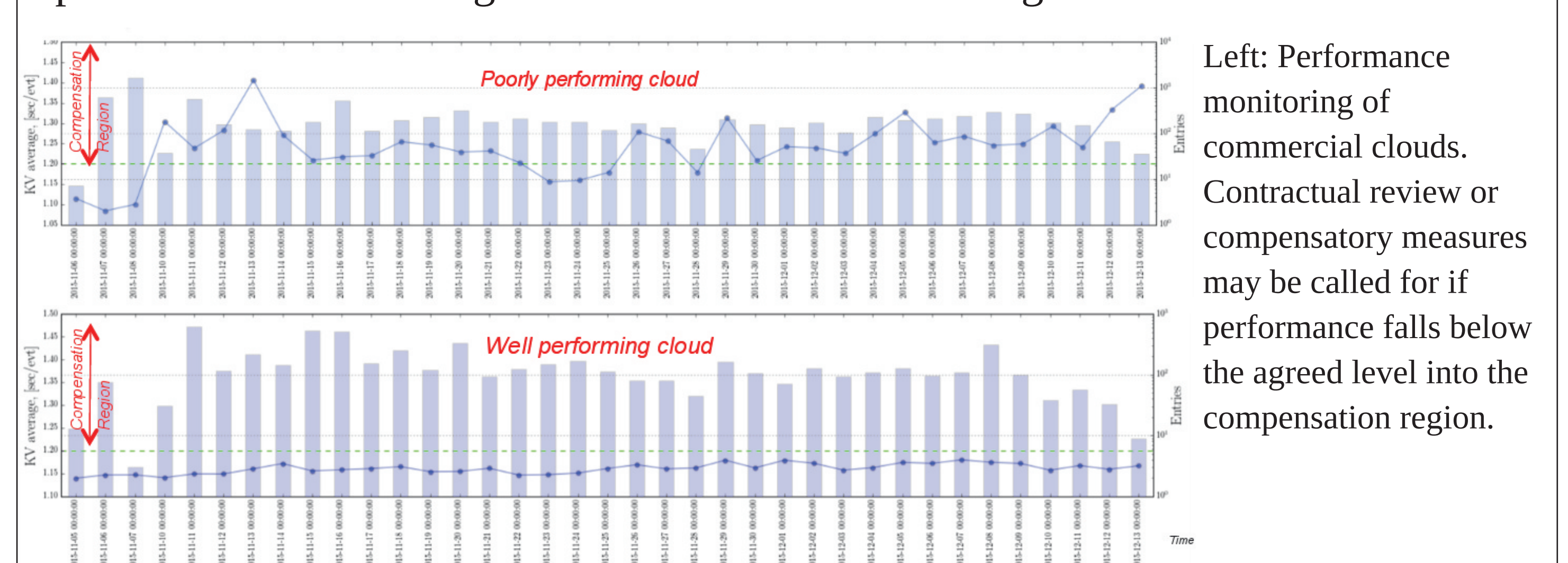
Below: Data workflow on AWS.

To efficiently exploit the spot market, an interruptible workload is needed. Event Service jobs meet this need, and natively support data stage-out to S3. The S3 storage is integrated with the grid using FTS and a SRM and GridFTP server on S3FS.



## Commercial Cloud Procurement

CERN has conducted several exercises in commercial cloud resource procurement, resulting in valuable experience in integrating these resources into the ATLAS distributed computing framework, which will lead to increased adoption in the future.[6]

In the context of commercial providers, accounting is of particular importance for understanding performance, usage and efficiency. Data from the existing monitoring infrastructure was leveraged to provide accounting information, and performance monitoring is based on the benchmarking suite.



Left: Performance monitoring of commercial clouds. Contractual review or compensatory measures may be called for if performance falls below the agreed level into the compensation region.

1. T. Kouba, "Evolution and Experience With the ATLAS Simulation at Point1 Project", presented at CHEP 2016 poster session
2. A. McNab, "The Vacuum Platform", presented at CHEP 2016 poster session
3. R. Sobie, "Context-Aware Cloud Computing for HEP", proceedings of ISGC 2016
4. D. Giordano et al, "Benchmarking Cloud Resources", presented at CHEP 2016 poster session
5. R. Taylor et al, "The Evolution of Cloud Computing in ATLAS", 2015 J. Phys.: Conf. Ser. 664 022038
6. D. Giordano et al, "CERN Computing in Commercial Clouds", presented at CHEP 2016 lecture session