# The second generation of the ATLAS Production System: expertise and future evolution

...

M. Borodin (U of Iowa)
On behalf of F. Barreiro, M. Golosova, D. Golubkov, A. Klimentov, T. Maeno, P. Nilsson, M. Titov and ATLAS collaboration

# Introduction

- PanDA -  **P**roduction **an**d **D**istributed **A**nalysis System
  - Designed to meet ATLAS production/analysis requirements for a data-driven workload management system capable of operating at LHC data processing scale
- New generation of ATLAS production system was developed for Run 2 and beyond – **ProdSys2**
  - Improved resource utilization
  - New types of computing resources: HPC, Clouds
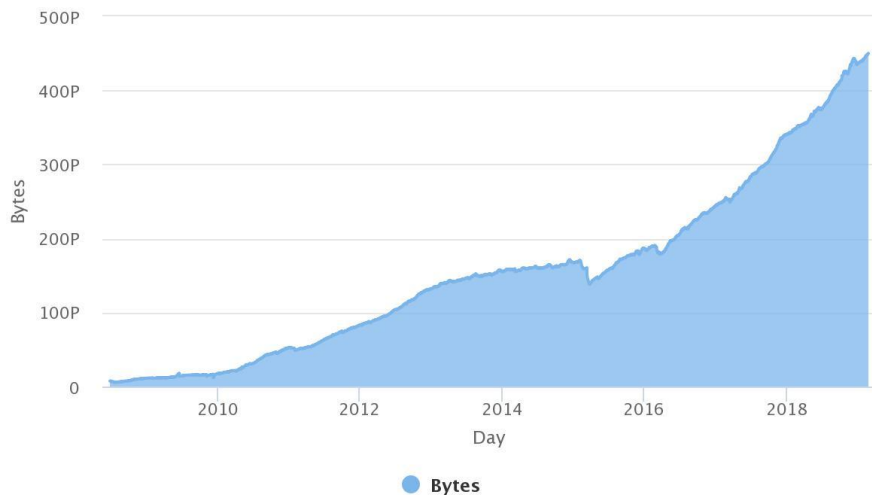  - Improved usability and robustness

# ATLAS production system design goals

- Deliver transparency of data processing in a distributed computing environment
- Achieve high level of automation to reduce operational effort
- Flexibility in adapting to evolving hardware, computing technologies and network configurations
- Scalable to the experiment requirements
- Support diverse and changing middleware
- Insulate user from hardware, middleware, and all other complexities of the underlying system
- Support custom workflow of individual physics groups
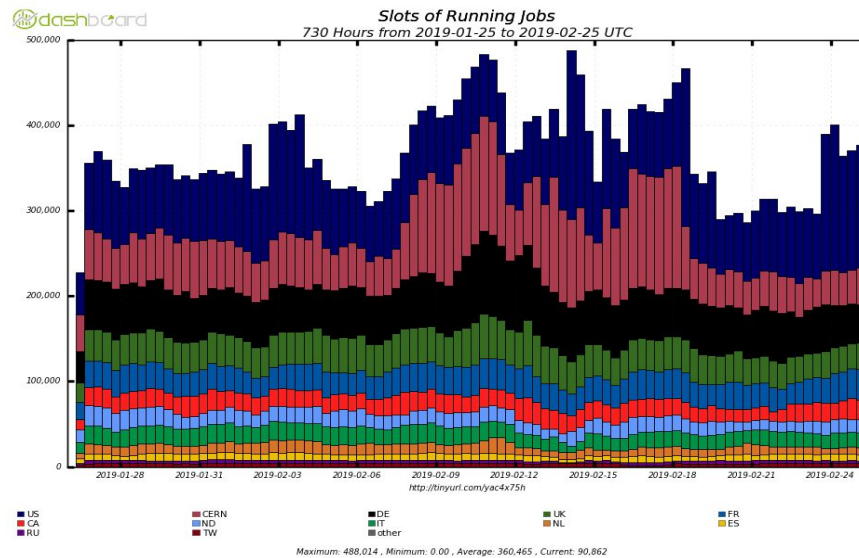- Incremental and adaptive software development

# Orders of magnitude



ATLAS Data Overview — Worldwide



Slots of Running Jobs
730 Hours from 2019-01-25 to 2019-02-25 UTC

https://bigpanda.cern.ch

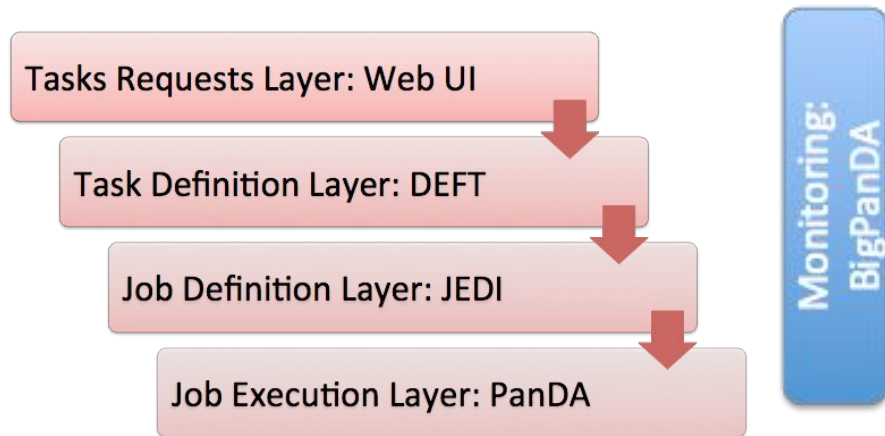400 PB of data is managed by ATLAS DDM system (Rucio)

More than 300K cores used by simultaneously running jobs in the system

# ATLAS production system components

- **Web UI** for Managers and Users provides the interface for task* and production request managing and monitoring at the higher level
- Database Engine for Tasks (**DEFT**): is responsible for formulating the tasks, chains of tasks and also task groups (production request), complete with all necessary parameters
  - It also keeps track of the state of production requests, chains and their constituent tasks

Tasks Requests Layer: Web UI

Task Definition Layer: DEFT

Job Definition Layer: JEDI

Job Execution Layer: PanDA

Monitoring: BigPanDA

*Task consists of jobs that all run the same program.

# ATLAS production system components (cont.)

- Job Execution and Definition Interface (**JEDI**): is an intelligent component in the **PanDA** server to have capability for **task-level** workload management.
  - Key part of it is **'Dynamic'** job definition, which highly optimizes resources usage compared to 'Static' model used in ProdSys1.
    - Dynamic job definition in JEDI is also crucial for multi-core, HPCs and other new requirements
- Monitoring (**BigPanDA**): progress, status and error diagnostics for all components.
- The PanDA **pilot** is an execution environment used to prepare the computing element, request the actual payload (a production or user analysis job), execute it, and clean up when the payload has finished. Input and output are transferred from/to storage elements, including object stores.

# Harvester

Harvester is a resource-facing service between the PanDA server and collection of pilots for resource provisioning and workload shaping. It is a lightweight stateless service running on a VObox or an edge node of HPC centers to provide a uniform view for various resources. The following picture shows how harvester interacts with PanDA and resources.

# DEFT data model

- Model is represented by multilevel relational instances:
  - **Request** -> **Slice**(chain of steps) -> **Step** -> **Task**
  - Depending on workflow each instance could play a role of a template
  - Tasks are created by initiating a step instance.
  - **Hashtags** are used to union an arbitrary number of tasks

# DEFT workflows

- ATLAS production workflows were implemented in chosen model
  - **MC simulation** is composed of many steps: generate hard-processes, hadronize signal and minimum-bias events, simulate energy deposition in the ATLAS detector, digitize electronics response, simulate triggers, reconstruct data, transform the reconstructed data into reduced forms for physics analysis

| Hard-scattering or min-bias | Event generation | Detector simulation | Digitization and pileup events | Trigger simulation | Reconstruction | Group production | Analysis |
|---|---|---|---|---|---|---|---|

- **Data Reprocessing** workflow has a tree structure, where output of one task can be an input for several more tasks
- **Derivation** is using so called "train" model, there each input runs on some of many predefined outputs
- **Tier-0** workflow
- **HLT, EventIndex,** ...

# DEFT and web UI development and deployment

- Key development points
  - Agile methodology: continuous meetings with the main users and often releases
  - Using open source
    - Django, Celery, AngularJS
  - «Model View ViewModel» approach



- Using CERN SSO for authentication and authorization

# Web UI

Request management

Request creation interface



Tasks management

# Production request processing

- Task request Web UI provides many general and experiment specific features:
  - **Bookkeeping**. Storing metadata, including arbitrary hashtags, allows to provide fine tuning statistics for running and historical tasks.
  - **Approval management.** E.g. MC production request required several levels of approval.
  - **Monitoring**. User can easily follow progress of a running tasks.
  - **Error Handling.** Task could fail because of many permanent (e.g. bug in software) and temporal (storage is down) reasons. To be able to quickly understand the root of the problem and fix it by redefining the task is one of the major features of the production system.
  - **Chaining** one production to the other. E.g. derivation production could be chained to MC or reprocessing task, that significantly speeds them up.
  - **Automation** of task submission. User can define a pattern and when new data appears tasks are started automatically.
  - ...

# DKB - Data Knowledge Base



DKB is ElasticSearch based system. It is being developed to consolidate different metadata which are related to the ATLAS production system. It's useful for troubleshooting, statistics, workflows optimization.

# Addressing future challenges

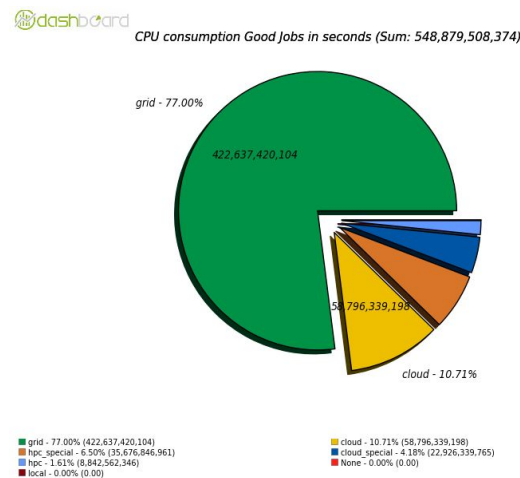- ATLAS Distributed Computing was very
  successful in the last years with clouds, HPC
  and HTC integration and using opportunistic
  **computing resources** for the Monte Carlo
  production



CPU consumption Good Jobs in seconds (Sum: 548,879,508,374)

grid - 77.00%
422,637,420,104

58,796,339,198

cloud - 10.71%

- grid - 77.00% (422,637,420,104)
- hpc_special - 6.50% (35,676,846,961)
- hpc - 1.61% (8,842,562,346)
- local - 0.00% (0.00)
- cloud - 10.71% (58,796,339,198)
- cloud_special - 4.18% (22,926,339,765)
- None - 0.00% (0.00)

- The HL-LHC era **data storage** estimated requirements are several times bigger
  than the present forecast of available resources, based on flat budget assumption
  - "Data Carousel" is a new project, which should allow orchestration between
    workload management, data management and storage services whereby a
    bulk production campaign with its inputs resident on a cheaper storage(e.g.
    tape), is executed by staging and promptly processing a sliding window of
    inputs

# Conclusion

- Constantly increasing luminosity and always limited computing budget require to find ways for further efficient and economical use of traditional and new computing resources

- The ATLAS production system development and operation experience gained during LHC Run 2 creates an excellent base to face upcoming challenges