# CHEP 2024
# Distributed computing (Track 4) Summary

● ● ●

Daniela Bauer, Fabio Hernandez,
Panos Paparrigopoulos, Gianfranco Sciacca

October 19 - 25, 2024

**CHEP 2024**

# Themes

- Tokens, Tokens, Tokens
  - Tokens being widely used and lessons are being learned
  - Current production setups
  - Development of best practice
  - Indigo-IAM development (also in response to "lessons learned")
  - Adoption of tokens outside WLCG - early stages
- Operations
  - Grid computing adapts to changing circumstances
    - Operations: Optimizing use of available resources
    - Monitoring
    - New architectures and modern resources: ARM, HPCs and Clouds!
  - Security: It's not just technology, people matter, too
- Distributed computing as part of non-WLCG computing models:
  - Gaining popularity especially in Astronomy: SKA, LSST, Einstein Telescope, CTA, HERD, but also DUNE (not astronomy)
    - Predicted SKA data volumes easily comparable to WLCG, building on WLCG experience for large scale operations!

**30 Talks 19 Posters**

Link to all track 4 talks: https://indico.cern.ch/event/1338689/sessions/553987/#all
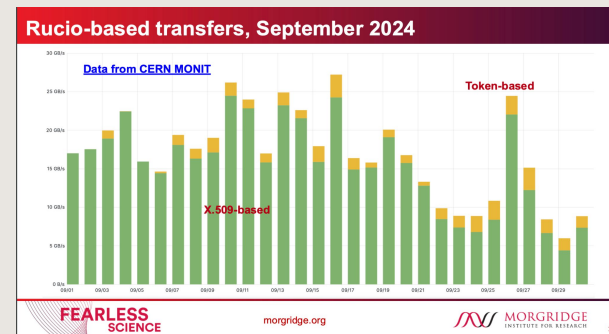
October 19 - 25, 2024
CHEP 2024

# Tokens

- Tokens are now a reality! The infrastructure is almost token ready- time to focus on the operational models!
- DC24, a major milestone: millions of transfers with tokens! Lessons learned:
  - Token implementations of middleware need to improve: FTS/Rucio/Dirac workflows/IAM: all doing a lot of work:
    - Since August ATLAS has been running tokens to 15 sites: 1-2Hz with 5Hz spikes!
  - Performance must be stress-tested
- IAM is being improved:
  - Moving from OpenShift to K8S
  - Better token lifecycle management and storage
  - OIDC/OAuth 2.0: from MitreID Connect to Spring Authorization Server
  - Open Policy Agent (OPA) to speed up policy evaluation and a move to 2FA



TOKENS

TOKENS EVERYWHERE

October 19 - 25, 2024
CHEP 2024

- CMS has made strides in token usage:
  - Every CE but 3 use tokens for pilot submissions, analysis to follow
  - CMS is using tokens in production (via Rucio) since early September (in over 30 sites)
    - 1 token per dataset, IAM can handle just fine
  - By the end of 2025 all services should be able to handle both x509 and tokens!
- Fermilabs' Vault instance is used to manage tokens:
  - All FNAL hosted experiments (minus CMS) have been using tokens with all grid jobs for over a year now, with the CILogon OIDC Provider.
  - Vault is hiding the token complexity from users:
    - Vault is paid but there is a very promising open source version.
  - CMS will do the same, possibly via a CERN instance.



Brian Bockelman - CMS Token Transition



Nick Smith - Fermilab's Transition to Token Authentication

# Tokens

- The balance between operability, security and performance needs to be found:
  - The Token Trust and Traceability WG Aims to form best practices, for users, devs, service providers + issuers.
  - Audience, lifetime, scopes are the three orthogonal parameters that one needs to tweak to meet the operational needs without compromising too much security.

- Next steps:
  - Use tokens on all grid jobs for stageout and reading
  - Users no longer need to issue x509/proxies
  - Accounting needs to be figured out



### Orthogonal(-ish) Axes: What, Where, How Long For.

Whilst not capturing everything, a useful way of visualising some of the **tunable token attributes** (but this is more than a 3-Dimensional problem).

The "goal" is to get a vector in "token-trust space" with as small a magnitude as you can and still meet your **operational needs**. The closer to the "origin" the better.

These considerations are made **per workflow**, and are ultimately a form of **risk analysis**.

X.509 proxies would exist almost "off the charts" on all 3 axes!

The "units" of the axes correspond roughly to:
- "Broadness" and "Power" of Scopes.
- Number and "Sensitivity" of Audiences
- Time

Matt Doidge - Early recommendations from the Token Trust and Traceability WG



### Next milestones

- **M.9 (Mar 2025): Grid jobs** use tokens for reading and stageout.
  - Implies **significant changes** in workload management systems
    - Tokens to be provided just in time?
    - Scopes? Audiences? Lifetimes?
    - Scalability concerns?
    - Fallback on X509 + VOMS during transition period?

- **M.10 (Mar 2026): Users** no longer need X.509 certificates
  - **Tools** should be sufficiently smart to obtain the correct tokens for specific operations
  - **Auxiliary** services such as **Vault + htgettoken** or **MyToken** may be needed to simplify the user experience, used under the hood by tools for job and/or data management
  - Investigations in this space are already underway within some experiments

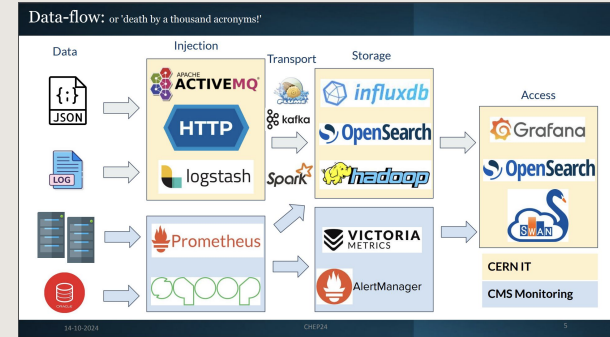Tom Dack - WLCG transition from X.509 to Tokens: Progress and Outlook

- WLCG central ops:
  - Size and complexity of infrastructure grows: new architectures, non-grid sites person-power doesn't.
  - Focus is on integrating those heterogeneous resources, while keeping the grid operating (while pushing for common tools and approaches!)
- ATLAS Hammercloud
  - Automatic exclusion/recovery of sites
- Monitoring:
  - Unified Experiment Monitoring
  - Experiment specific: CMS
    - Common tools and technologies to minimize maintenance and operations



Brij Kishor Jashal - Advanced monitoring capabilities of the CMS Experiment for LHC Run3 and beyond



Ewoud Ketele - Unified Experiment Monitoring

Fri 25th Oct - Track 4 summary: Distributed Computing - 6

# Operations - Optimisations

- Optimising the use of available resources:
  - ATLAS: HEP benchmark
    - Distinguishing fact/reality from fiction/ideal case
  - ATLAS: Results of the review of the ATLAS workflow management system (PanDa)
  - ALICE: Job Optimizers
    - Submit jobs faster and to the correct sites
  - ALICE: Whole node scheduling
    - Better exploit node resources with tuned payloads
  - ALICE: Unprivileged subdivision of job resources within the ALICE grid
  - CMS pilot overloading



Natalia Szczepanek - Optimization of ATLAS computing resource usage through a modern HEP Benchmark Suite via HammerCloud and PanDA



Marta Bertran Ferrer- Whole-node scheduling in the ALICE Grid: Initial experiences and evolution opportunities

# Operations - Alternative architectures



### Meet Perlmutter

- Runs HPE Cray OS
- Only whole-node scheduling
- 3072 nodes running 2 x AMD EPYC 7763, 64 core, 512 GB memory
- Jobs run in Shifter containers
- Outside connection from nodes
- Full CVMFS support

Sergiu Weisz — 6

Sergiu Weisz - Integrating the Perlmutter HPC system in the ALICE Grid

- LHAASO trying to integrate Chengdu Supercomputing Centre
  - Dedicated link to avoid firewall issues and SLURM/HTCondor/XrootD to the rescue
- ALICE integrating the Perlmutter HPC
  - Integrated successfully getting the resources equivalent of a T2
- JAliEn (ALICE's Grid Framework) evolved to support ARM!
  - Also riscv64 architecture, as a proof of concept
- CVMFS performance upgrades:
  - New cache manager to open fewer files and improvements on parallel downloads!
- HEPCloud, after 6 years and a lot of problem solving is now a mature provisioning system which provides access to compute resources (HPC + clouds) similar to the size of the US CMS Tier-1 facility at Fermilab!
- SPECTRUM
  - EU funded project: focus on strategy, but also technical blueprint for data-intensive science and infrastructures

# Security

- Threat from cyber attacks is persistent: strategy and a plan are needed

- People are the key: Communicate, collaborate, share

- Security operations centre (SOC) fits with an overall cybersecurity plan such as the Trusted CI Framework.
  - Be proactive to prevent cybersecurity incidents: monitor, detect, respond

- The pDNSSOC package was suggested as a lightweight way for smaller sites to get the benefits of a SOC; more volunteers/testers of this would be welcome.



David Crooks - Designing Operational Security systems: People, Processes and Technology

# Life outside the WLCG - 1

- SKA (all purpose radio telescope):
    - At full operations expects rates up to 400 PB/year by 2030 - easily comparable to LHC experiments
    - Construction is planned in stages and data from the very first stage is available
    - Computing organized in ~9 SRCNet resource centres, using common tools like IAM
- CTA (gamma ray astronomy):
    - Observations planned to start at 2030. Simulations are running since 2011
    - Production system centred around DIRAC (to move to DiracX)* with a CTA specific extension and soon Rucio
- Einstein Telescope (gravitational wave observatory):
    - (Data) challenge driven iterative development for computing model
    - Still multiple iterations expected

* see plenary: https://indi.to/DHhVG



**Square Kilometer Array: Transforming radio astronomy**

The Square Kilometer Array (SKA) Observatory (SKAO) is a next-generation radio astronomy facility which will cover the frequency range from 50 MHz to 15 GHz.

A mosaic illustrating the main science drivers for the SKA

Composite image of the SKA telescopes, blending real hardware already on site with artist's impressions. credit: SKA Observatory

Credit: SKA Observatory

Ian Collier - The path to exabyte astronomy: SRCNet v0.1 for the Square Kilometre Array



**The Cherenkov Telescope Array Observatory**

- The next generation **ground-based** observatory for **gamma-ray astronomy** at **very high energy**
- **64 telescopes** located on **two sites**
- Operations expected to start in **2030**
- Observe **extreme cosmic events**:
    - supernovae, neutron stars, black holes ...
    - Transient phenomena: Gamma Ray Burst

La Palma, Spain (North Site)

Paranal (ESO), Chile (South Site)

CHEP – October 2024 - Natthan Pigoux

Natthan Pigoux - The Cherenkov Telescope Array Observatory Production System Status and Development

# Life outside the WLCG - 2

- <u>DUNE</u> (neutrinos)
  - Worldwide distribution of data and compute
  - FNAL legacy systems were used in the first prototypes: these have now been replaced
  - Uses HTCondor + GlideinWMS (from CMS) for job submission, Rucio for data management.
  - Participated in DC24, validating the new stack
- <u>Vera C Rubin</u> (sky survey)
  - Final phases of construction. Planned to start at 2025
  - Dealing with monitoring and log keeping / analysis challenges.
  - Uses PanDa (from ATLAS)/RUCIO/FTS/IAM
- <u>HERD</u> (High Energy cosmic-Radiation Detection facility - in space !)
  - > 90 PB, ~16000 CPU cores, in 10 years
  - Uses DIRAC/Rucio from WLCG
  - T0/T1/T2 distributed model, with T2 doing simulations only



Jacob Michael Calcutt - <u>Evolution of DUNE's Production System</u>



Fabio Hernandez - <u>Preparation of Multi-Site Data Processing at the Vera C. Rubin Observatory</u>

Fri 25th Oct - Track 4 summary: Distributed Computing - 11

# Thank you!

- Thanks to all the Track 4 speakers and poster presenters!
- To the organization: the program committee and the team of track convenors for their outstanding work in making this a success!
- To all of you, making this conference a wonderful experience!
- Looking forward to seeing you all again soon!