# AF update

A. Forti
WLCG workshop
9 November 2022

# Analysis Facilities



UChicago

UK

Coffea
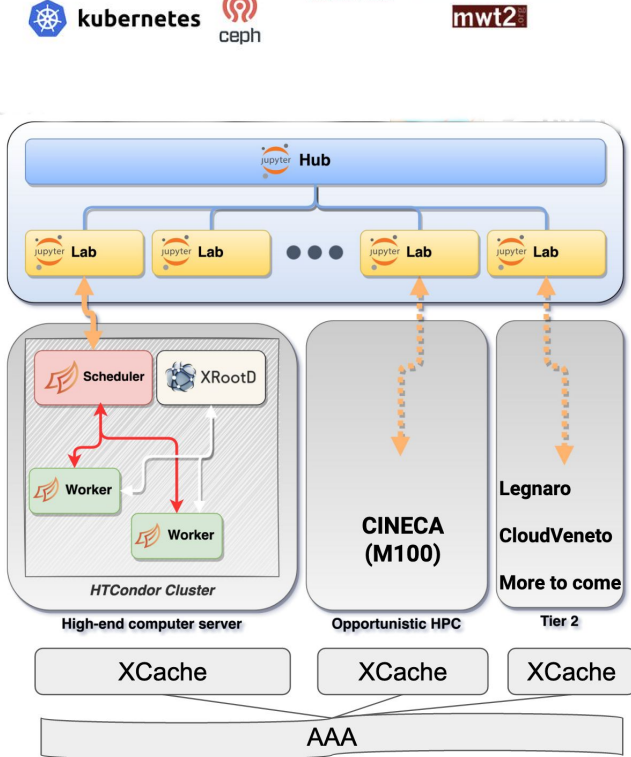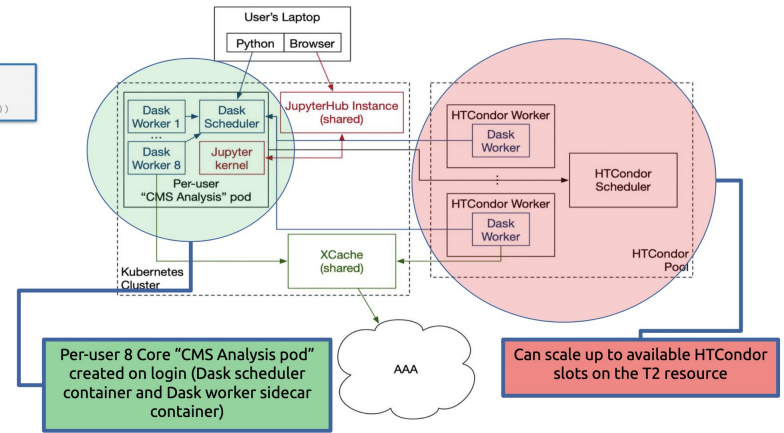
ESCAPE DLaaS

INFN

NAF

# Analysis Facilities

- AF definition in HSF is still
  - "Infrastructure and services that provide integrated data, software and computational resources to execute one or more elements of an analysis workflow. These resources are shared among members of a virtual organization and supported by that organization."
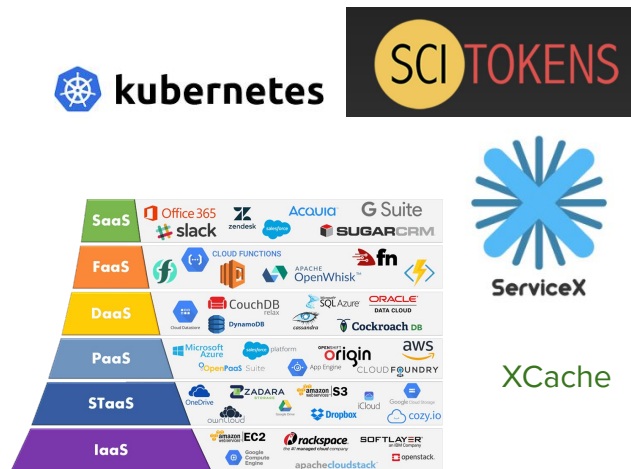- Whatever shape an AF takes the important is being able to integrate the technologies to support the analysis evolution.
  - Some of thee technologies might be adopted by T3s like CVMFS in the past

Technology Evolution

Techniques Evolution

# Topics

- Out of the Analysis Ecosystem [workshop report](#) and other places the need to concentrate on building blocks rather than specific architectures
- Some topics
  - AAI
  - DOMA topics
    - Grid storage access and xcaches
    - Shared storage
    - Object stores
  - Environment sharing: containers
  - Tracking analysis performance
  - Resource sharing and services deployment (k8s)
  - IRIS-HEP AGC as a concrete goal to pull various threads
- Need to discuss today HOW to follow this up
  - Preference would be to have groups of people willing to put sometime for each topic to follow also in other forum and then create a coherent picture from AF point of view
  - Some topics needs tight cooperation with other projects WLCG/IRIS-HEP/EOSC/… or HSF groups

# AAI

- Technologies current analysis facilities are exploring are better suited to tokens than x509
  - AF R&D and more flexible and adaptable to change
  - AF users are testers and are more adaptable
- AF better built around this from the beginning
- Token infrastructure not production ready but things can be followed in WLCG authz and DOMA BDT
- The question of integrating AF R&D in the testbed needs proactive engagement with above people
- Separate discussion should be held about the federated identities and access policies.
  - Final decisions on this depend on funding and experiments models and should be independent from creating the token infrastructure

Priority (tokens): could help the current effort 💬‼️
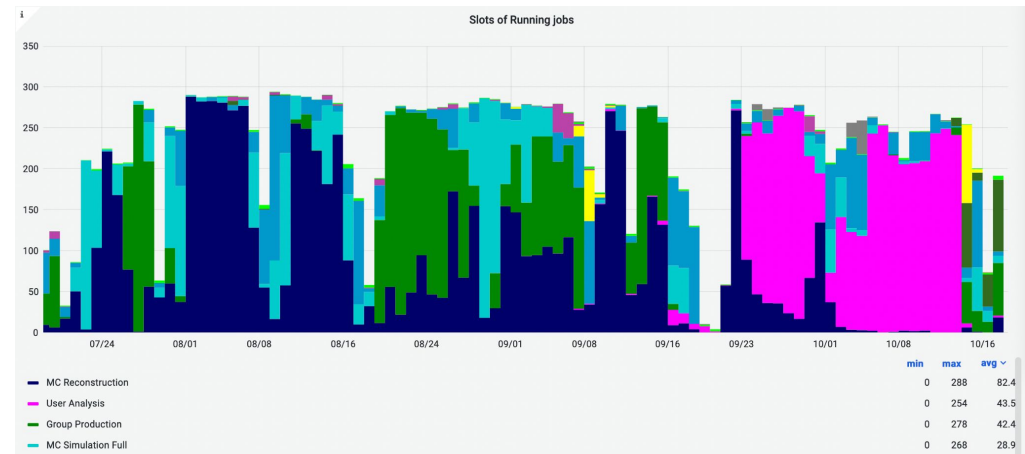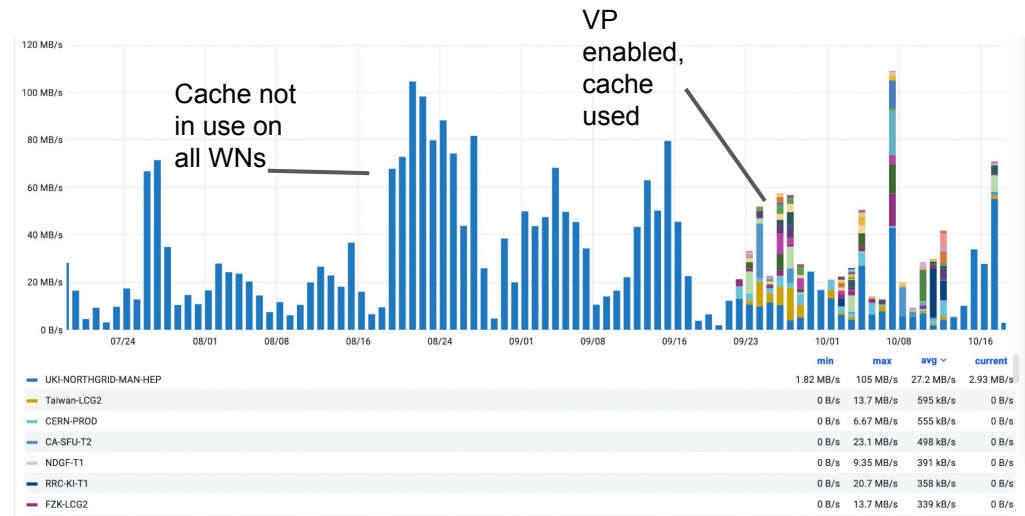
# DOMA (1)
## Bulk Storage And XCaches

- Accessing, integrating and caching input from the grid bulk storage (aka Data Lake)
  - Many ideas already in DOMA Access
  - Ongoing testing in different places
  - In ATLAS use of XCaches at AF and storageless sites
    - Integration with rucio using Virtual Placement
- Hardware specs
  - Range of cache sizes and performance but there is no recommendation on the actual specs
- Xcaches monitored w/ custom scripts, there is no agreed solution
    - Like for Xrootd streamed data most problematic
    - WLCG xroot monitoring is a way forward towards consolidating also xcaches monitoring but need experiments agreement
    - Other sciences might be interested

Priority: understanding hardware and having monitoring  !!!

# XCache in UK example

- Bham ~300 cores
  - XCache points to Mcr
  - VP allows more than one source
- Before mid September XCache was not enabled on all the WNs and part of the traffic is direct access job -> Manchester SE
  - 800 Mb/s peaks are large-ish for 300 cores
- After mid September there is a clear correlation between enabling XCache+VP and the amount of Analysis jobs.
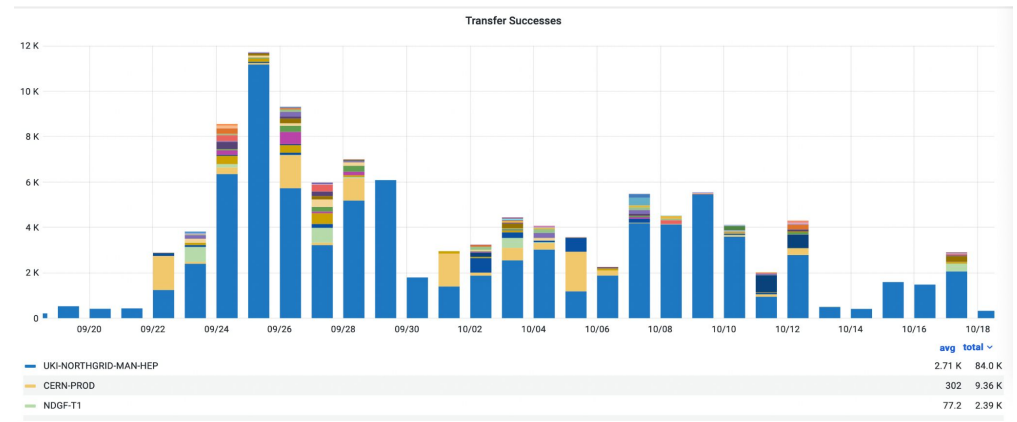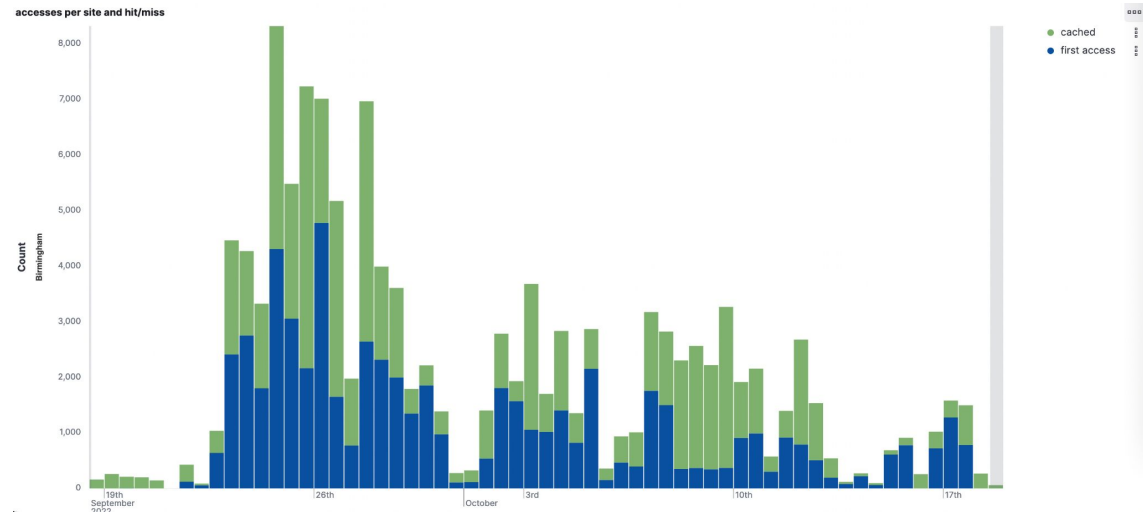


VP enabled, cache used

Cache not in use on all WNs

| | min | max | avg ∨ | current |
|---|---|---|---|---|
| UKI-NORTHGRID-MAN-HEP | 1.82 MB/s | 105 MB/s | 27.2 MB/s | 2.93 MB/s |
| Taiwan-LCG2 | 0 B/s | 13.7 MB/s | 595 kB/s | 0 B/s |
| CERN-PROD | 0 B/s | 6.67 MB/s | 555 kB/s | 0 B/s |
| CA-SFU-T2 | 0 B/s | 23.1 MB/s | 498 kB/s | 0 B/s |
| NDGF-T1 | 0 B/s | 9.35 MB/s | 391 kB/s | 0 B/s |
| RRC-KI-T1 | 0 B/s | 20.7 MB/s | 358 kB/s | 0 B/s |
| FZK-LCG2 | 0 B/s | 13.7 MB/s | 339 kB/s | 0 B/s |

Slots of Running jobs

| | min | max | avg ∨ |
|---|---|---|---|
| MC Reconstruction | 0 | 288 | 82.4 |
| User Analysis | 0 | 254 | 43.5 |
| Group Production | 0 | 278 | 42.4 |
| MC Simulation Full | 0 | 268 | 28.9 |

Last 3 months

7

# XCache in UK example

- A good 35-40% of the data is accessed from the cache
  - The rest is still accessed via XCache (buffer cache in = out)
- Monitoring might not be universal needs to be looked at.
  - For streamed analysis differences with monit.



accesses per site and hit/miss



Transfer Successes

| | avg | total |
|---|---|---|
| UKI-NORTHGRID-MAN-HEP | 2.71 K | 84.0 K |
| CERN-PROD | 302 | 9.36 K |
| NDGF-T1 | 77.2 | 2.39 K |

Atlas data Last 1 month

# Doma(2)
# Shared Storage

- Recurring topic: Local shared storage for people to seamlessly run from different resources and share with colleagues
    - Discussed during the SWAN/ScienceBox, EOSC and the INFN multi-site AF presentations.
        - First two use EOS at CERN, the latter has a local cephFS storage
    - Users repeatedly report it as a main feature at CERN
- As soon as the AF is distributed this becomes problematic
- If access is from remote or with xrootd gateways it's still not as straightforward and can cause problems to the site
    - Typical example users want to access/mount EOS
- For different solutions monitoring of traffic and access needs to be harmonised
- Needs discussion and recommendations

Priority: any user strong requirement

# Doma(3) Object Stores

- Another topic covered is the possibility to use object stores
  - Several grid sites installed ceph with cephFS for POSIX access
  - Object stores advantage is scalability because they don't need any metadata db (true or false?)
- AF workflows may have different requirements: which?
  - IRIS-HEP is introducing object stores in their testing
  - ServiceX backend storage is an object store recently moved to use S3
    - Users don't access S3 directly
- Evaluating object stores without posix access for standard experiment software is a large body of work
  - There is a CMS R&D proposal
  - Experiments requirements on this need to reviewed and initial recommendations if it is worth pursuing written

Priority: !

# Environment sharing

- Users want to share with colleagues their setup,code, configuration, small amount of input data….
  - Shared storage is the traditional way but not the only way
    - Conda (LHCb), Containers (CMS, ATLAS)
- Solutions easy to setup and should help also preservation
- Containers
  - We still don't have an official way to distribute images or a supported registry
  - Need to look at building and supporting base images
  - CVMFS unpacked.cern.ch has now 3000 images
    - Pulls images from any public register
    - Uses harbor.cern.ch as a proxy to sanitize images
  - Users using harbor directly could open various development roads
    - Images life cycle management
    - Possibly solving private images problem at least for containerd

Priority: demand is growing but scale not yet clear  !

# Tracking Analysis Performance

- Need an agreed upon list of metrics
  - Workflow ID,
  - CPU, RAM, swap,
  - I/O (local storage and network),
  - Software stack,
  - Job failure rate,
  - Time To Completion (TTC),
  - Data source local or cached from a Data Lake,
  - Formats used on input (PHYS, PHYSLITE, DAOD, NTuple,etc..),
  - Formats written (columns), ratios
  - ……
- Need also an infrastructure where to make them accessible
  - Centralised monitoring may take a long time to develop
  - Some request to instrument jobs like ATLAS does on the grid
- Other types of site monitoring in the same situation (networking, tape, xrootd, benchmarking…)
  - May draw from that experience
  - Some of the metrics may as well be the same

Priority: monitoring has always high priority    !!!

# Shared Resources and kubernetes

- Recommendation is to colocate AFs with existing sites to maximise resource sharing
- k8s is the proposed method of deployment for services
  - As well as an alternative backend to run users jobs
- It has been embraced by a number of large sites
  - But not everyone agrees it will solve problems
- In WLCG it is recurringly talked about
  - Last time at June pre-GDB (agreed another pre-GDB mid next year)
  - Ryan's talk later
- There aren't recommendations on deployment or for k8s service developers and users
  - Short document for app devs to get things to work both on okd and vanilla k8s
- CERN has a really robust documentation but it looks CERN specific
- So there is still no coordinated effort for this
  - Ricardo Rocha gave a comprehensive list of forums people can attend also report from kubecon on batch system features development

Priority: in US it is a requirement in Europe up to the sites 💬!!! 💬!! 💬!

# IRIS-HEP AGC

- [AGC status](#) at the last week
- AGC designed to measure new techniques and new services
- Official challenge with Dask and OpenData
- RDataFrame and non public data are
- Non IRIS-HEP sites can also participate
  - INFN and UK plan to participate with their infrastructures and workflows (even non-LHC)
  - Monthly ops meetings (next [tomorrow](#))
  - analysis-grand-challenge [google group](#) for announcements and discussion
- Need to couple sites and users to be more productive

Priority:

# HSF AF Forum


Kick-off meeting


EOSC update


Kubernetes


UX SWAN&coffea


XCache


Escape


Multi site AF model at INFN


Containers in CVMFS


UX DESY NAF

- HSF AF page
- Indico Category

# WLCG Discourse

- Proposal to use [wlcg-discourse.web.cern.ch](wlcg-discourse.web.cern.ch) for discussion since there is no other general forum for sites starting or needing an answer.
- Pros/cons
  - Pros: All in one place and easy to search.
  - Cons: Not well known enough to generate useful discussions