

Data Center and Services for the CNAF Reloaded project

Luca dell'Agnello, Gaetano Maron, Davide Salomoni

Data center: current figures

- INFN-T1 provides services and resources to more than 40 scientific collaborations (HEP, astro-particle, GW etc.)
 - 70-80% resources for WLCG experiments
- Staff composed by 22 people (including Facilities, Network management and User Support)
- 2019 pledges:
 - ~400 kHS06 on HTC farm Geographically distributed (~180 kHS06 @CINECA + 10 kHS06 @Bari-RECAS)
 - Also small HPC farm (~100 TFlops) and cloud instance (~1000 cores) available
 - ~39 PB of disk
 - ~89 PB of tapes
- Focus on efficiency and availability of services

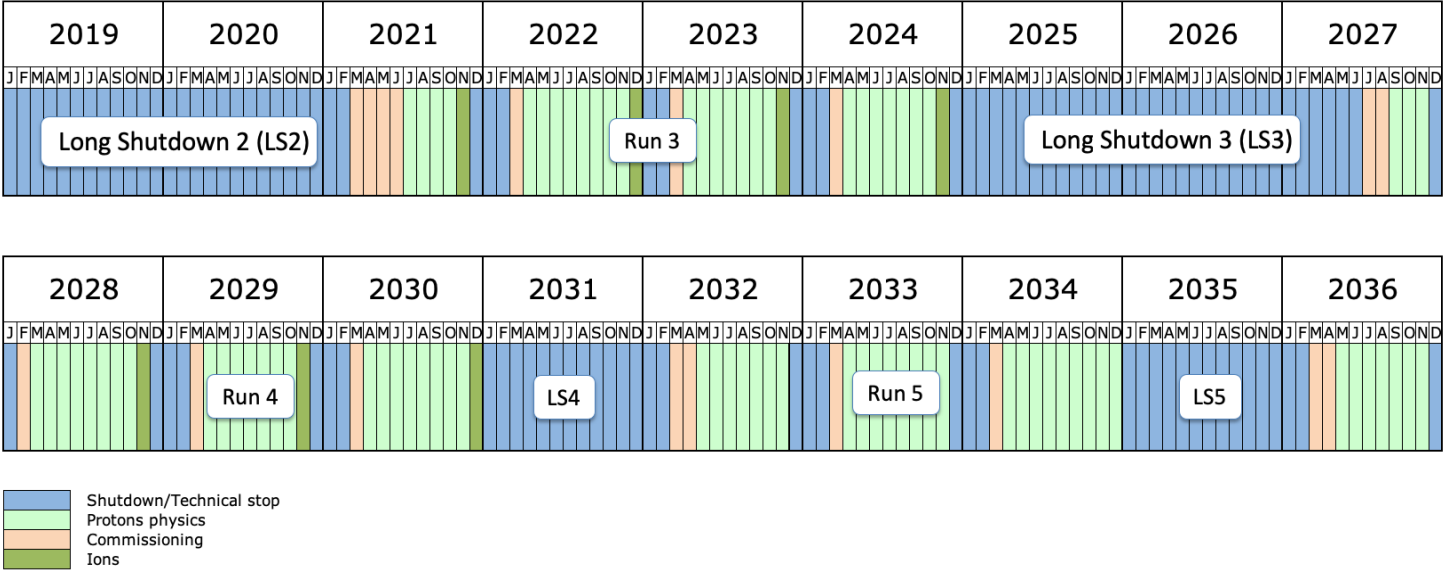
Services provided

- Services provided:
 - "Traditional" HTC farm with Grid front-end
 - Bare metal nodes managed by a batch system (LSF, now migrating to HtCondor)
 - Small HPC instance available
 - Cloud for "long tail of physics"
 - Interactive remote access currently under test
 - DM based on HSM (disk+tape) system with srm interface
 - All de-facto standard protocols provided (i.e. gridftp, xrootd, webdav)
 - LTDP
- WLCG services and [SLA](#) extended to all experiments (for relevant services)

| Service | Maximum delay in responding to operational problems | | | Average availability measured on an annual basis | |
|--|---|---|---|--|--------------------|
| | Service interruption | Degradation of the capacity of the service by more than 50% | Degradation of the capacity of the service by more than 20% | During accelerator operation | At all other times |
| Acceptance of data from the Tier-0 Centre during accelerator operation | 12 hours | 12 hours | 24 hours | 99% | n/a |
| Networking service to the Tier-0 Centre during accelerator operation | 12 hours | 24 hours | 48 hours | 98% | n/a |
| Data-intensive analysis services, including networking to Tier-0, Tier-1 Centres outwith accelerator operation | 24 hours | 48 hours | 48 hours | n/a | 98% |
| All other services – prime service hours | 2 hour | 2 hour | 4 hours | 98% | 98% |
| All other services – outwith prime service hours | 24 hours | 48 hours | 48 hours | 97% | 97% |

Data center: 2024 figures

- INFN-T1 will (likely) provide services and resources to even more scientific collaborations
 - Probably other demanding experiments besides WLCG ones
- Staff composed by 22+ people
 - Synergies with CINECA (Facilities at least)
- 2024 pledges: ~1000 kHS06 on HTC farm
 - Probably ~50% provided by Leonardo machine
 - Also hosting Cloud@INFN main instance
- ~120 PB of disk
 - Possible further increase of disk due to INFN data lake instance
- ~200 PB of tapes



Services to be provided @Tecnopolo

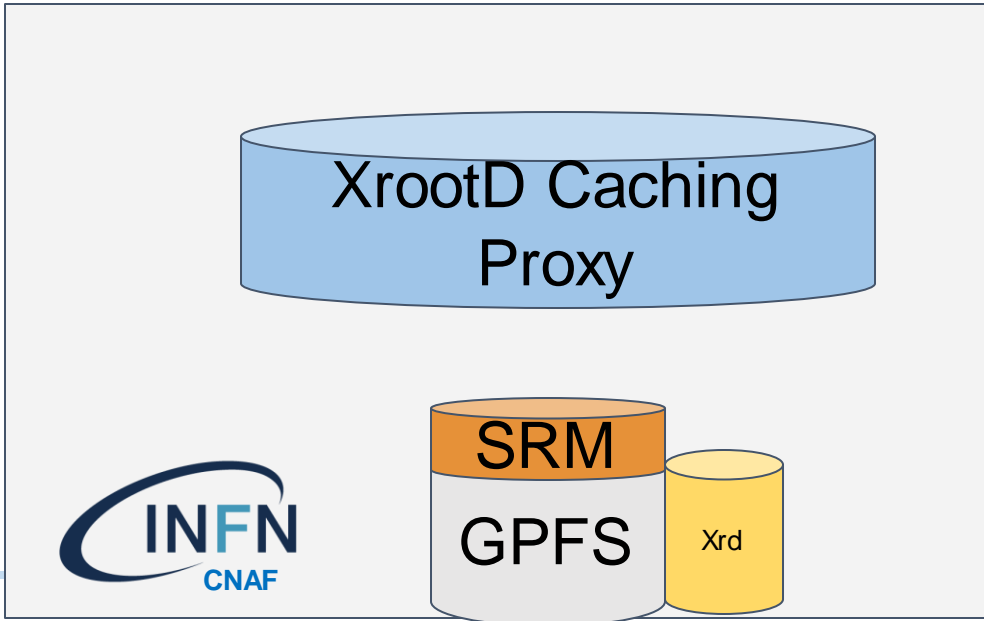
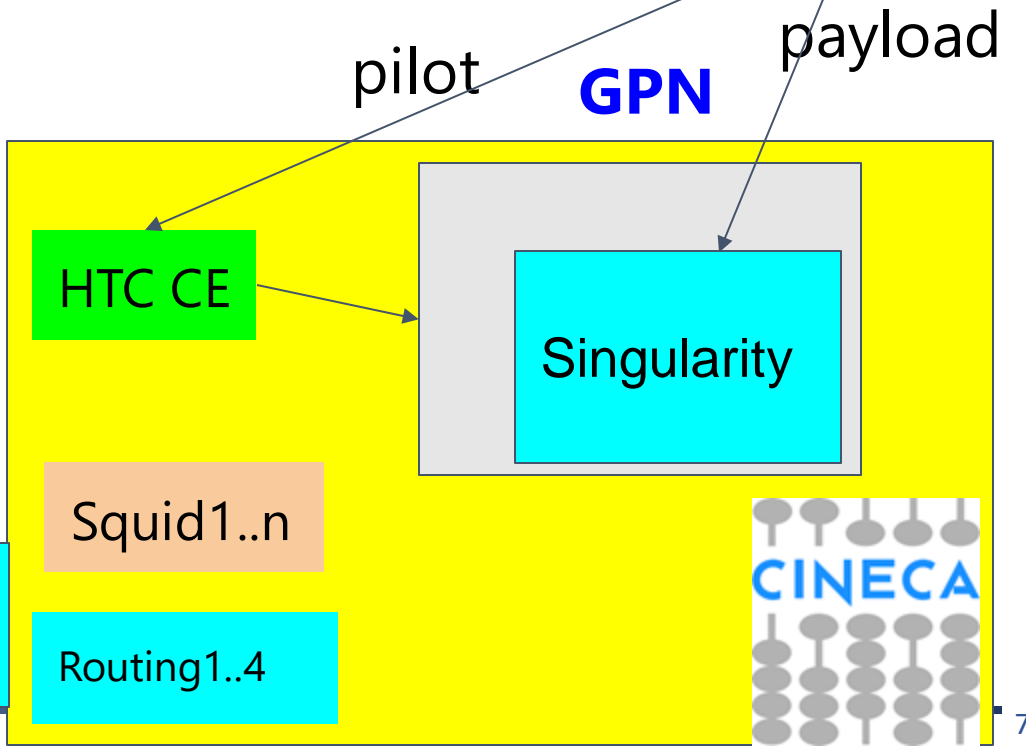
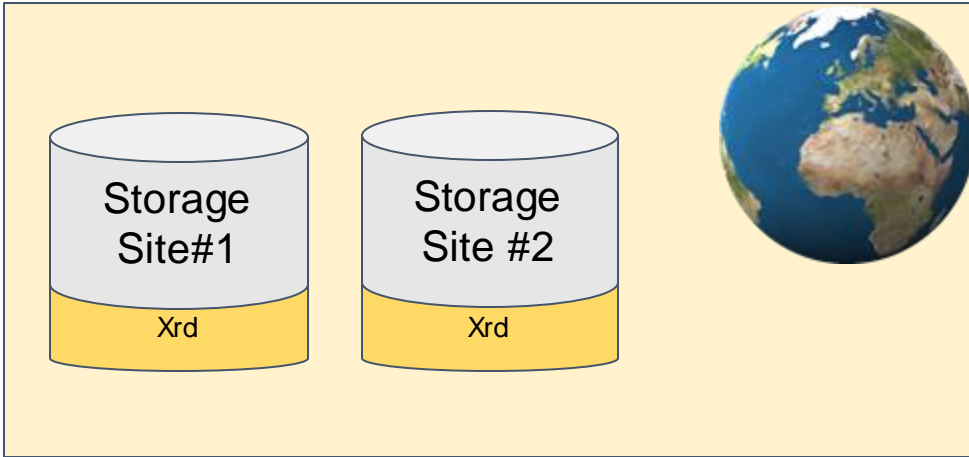
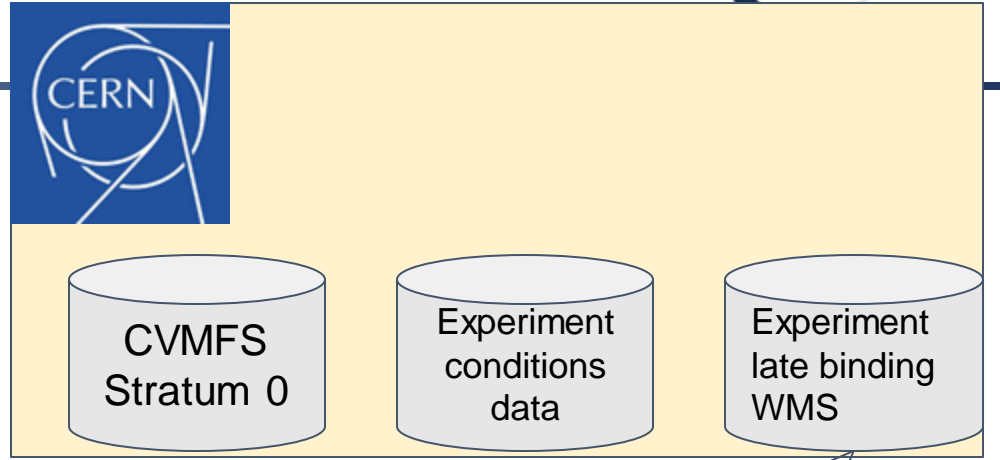
- At least the same set of services provided today
- But support also for:
 - INFN-cloud
 - INFN data lake
- More experiments: increased quantity of resources
 - (Probably) increased number of objects to manage and maintain
- Doing more with (most probably) the same number of FTEs
- ***We need to reduce even more human intervention and increase the overall efficiency***



HTC farm configuration and first evolution

- “Traditional” configuration of the HTC farm
 - Bare metal WNs managed by LSF (migrating to HtCondor)
 - Skipped the OpenStack phase on farm
 - Investigating K8s
 - Cream CE as interface (being replaced by CondorCE)
- Providing feedback to CINECA for defining the architecture of the pre-exascale machine to be installed at Tecnopolo....
 - and to clarify the operational model
 - See later for HPC/HTC integration activities.

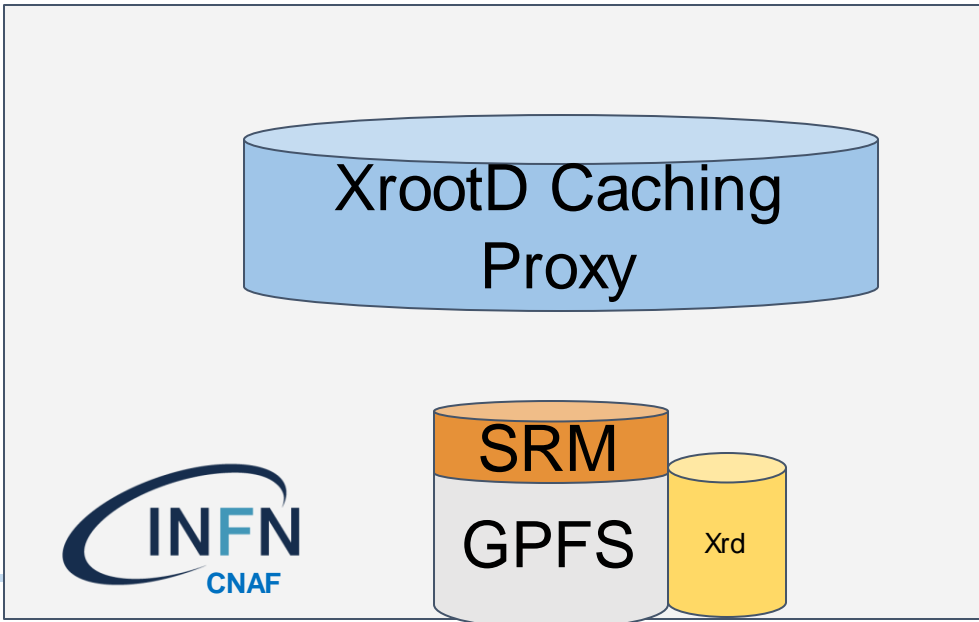
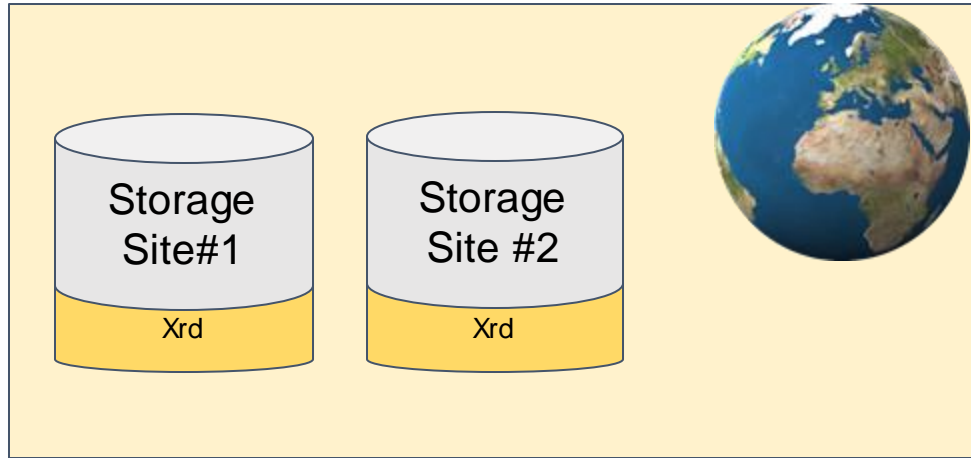
Technical setup #1: jobs



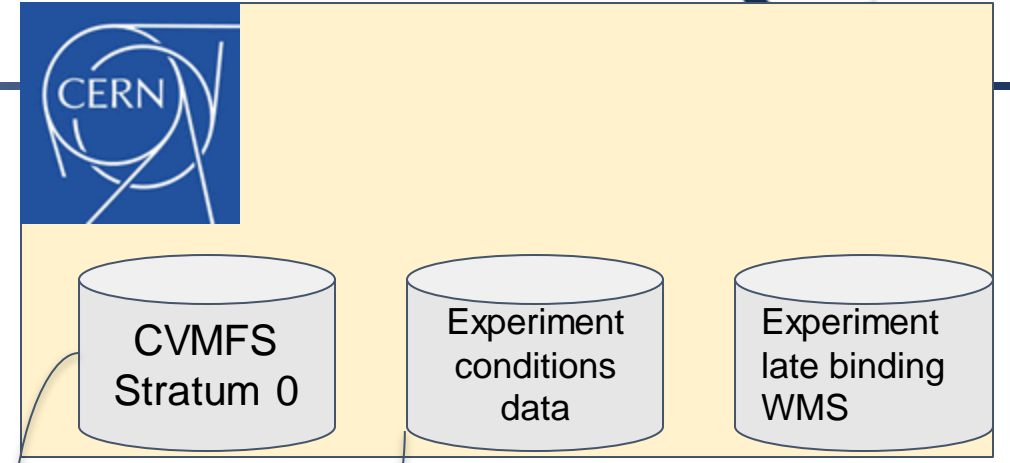
by T. Boccali – CHEP 2019

Technical setup #2: sw and conditions

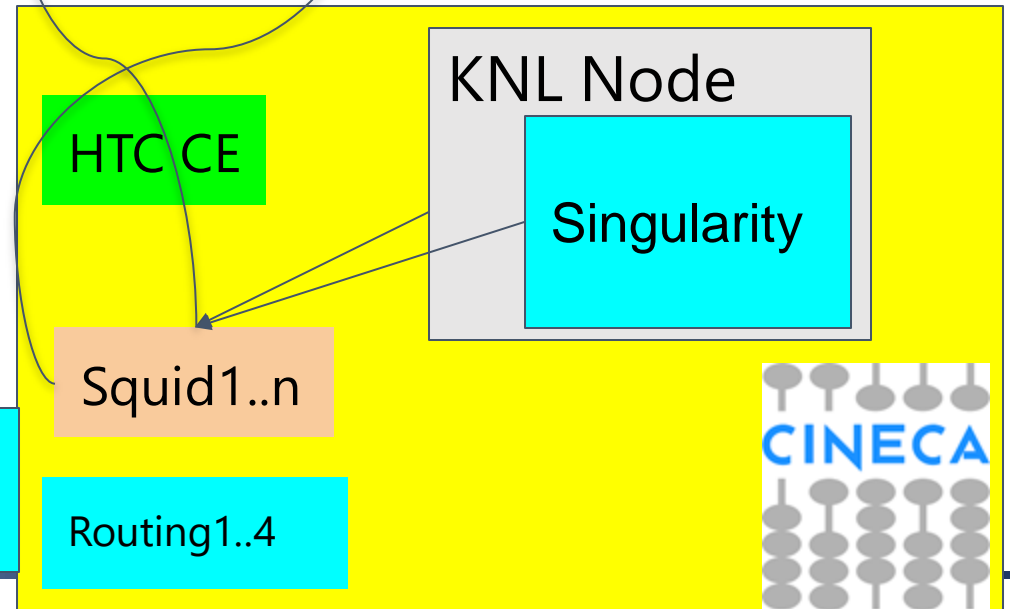
by T. Boccali – CHEP 2019



Direct link

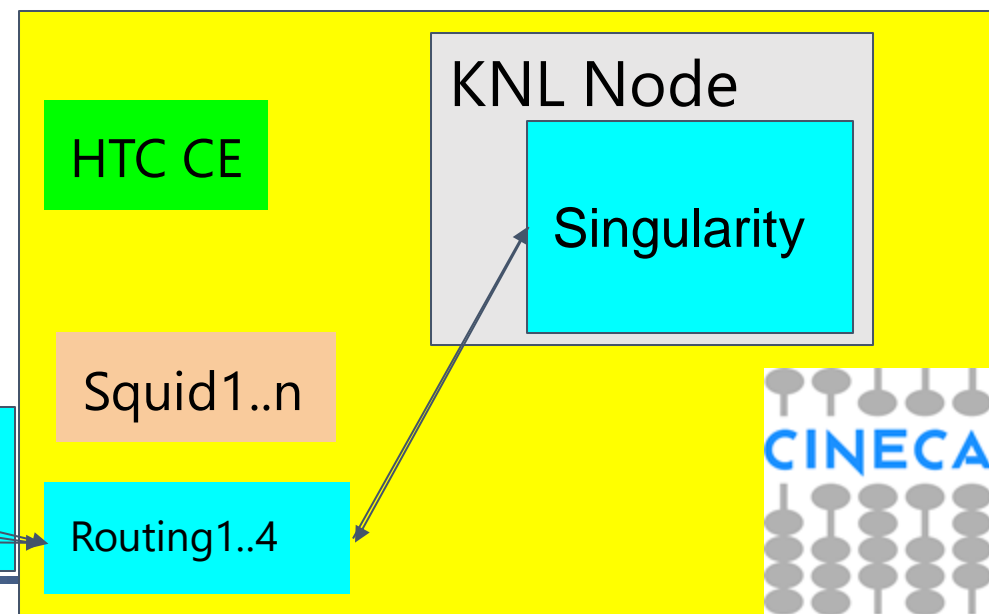
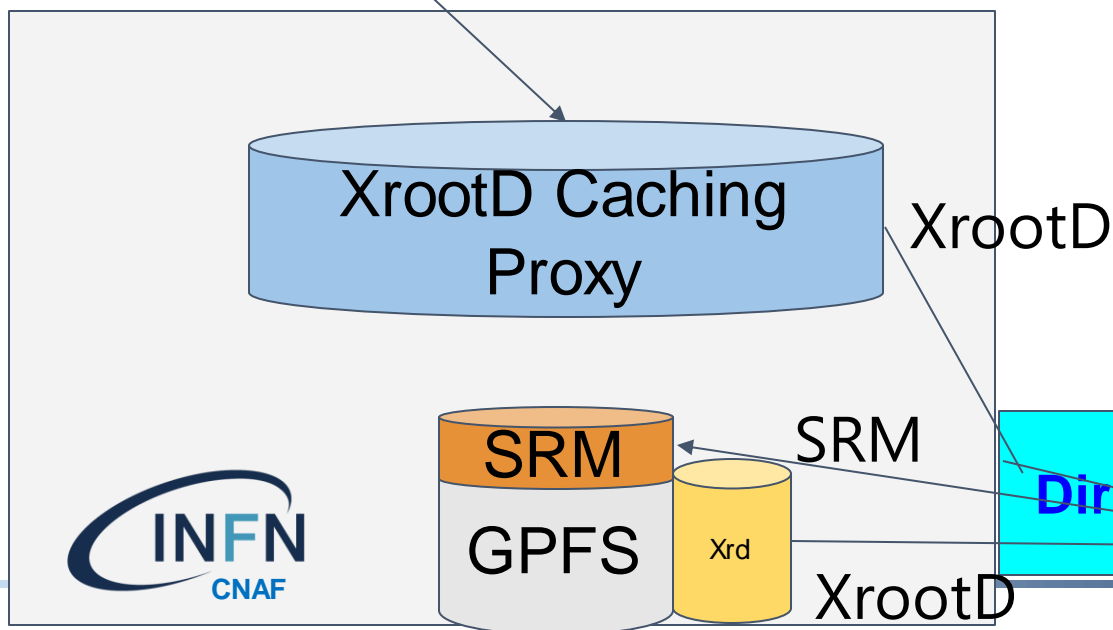
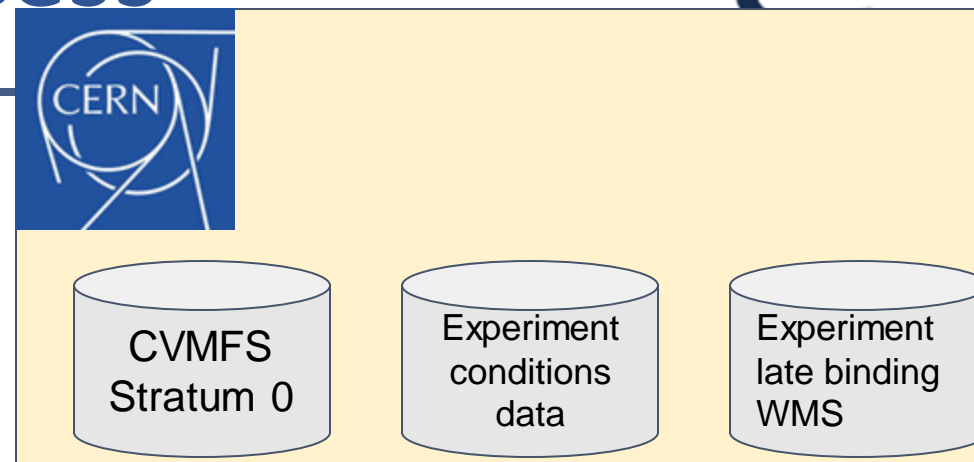
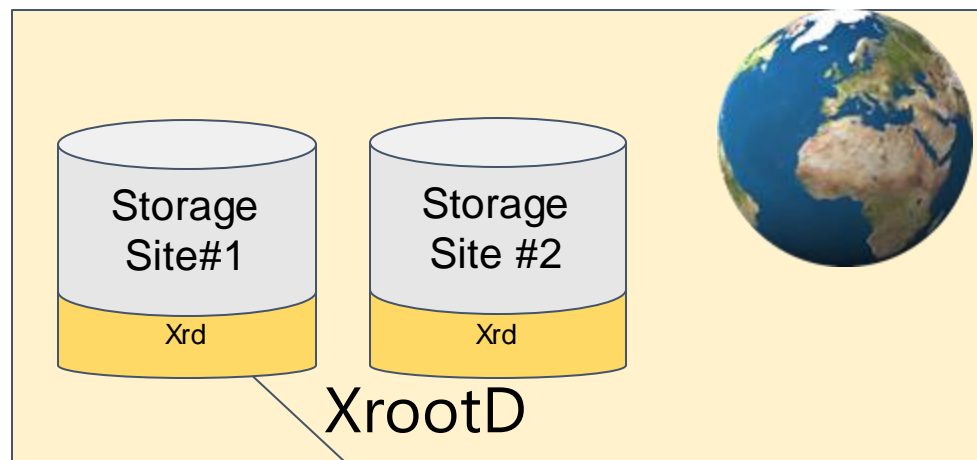


GPN



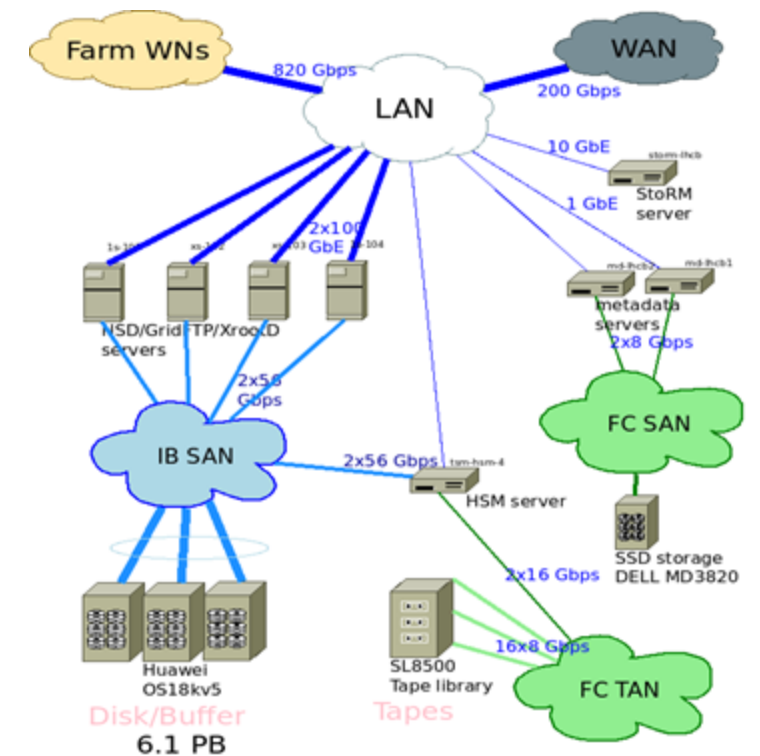
Technical setup #3: data access

by T. Boccali – CHEP 2019



Storage and data management (1/3)

- Custodial site for WLCG and most of the other experiments
- Also supporting Long Term Data Preservation for CDF
- All storage managed by GEMSS (GPFS+TSM+StoRM+“glue”) and based on FC/IB SAN
 - Large (~3 PB) building blocks for storage
 - Storage servers connected with 2x100 Gbps to LAN
 - Oldest installation with 4x10 Gbps servers
 - Storage density ranging 0.8-3 PB/server
 - Two libraries (35 drives) available
- Supported protocols: xrootd, srm, gridftp, webdav/http



LHCb example

- 4 as GridFTP, XrootD and NSD
 - 2 as metadata servers
 - 1 (VM) as StoRM FE/BE
 - 1 as HSM
- 4 I/O servers and 4 service nodes for 6PB of data!

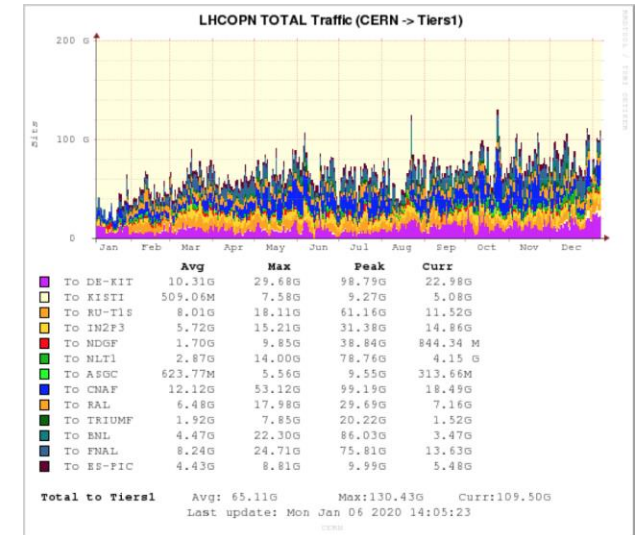
Storage and data management (2/3)

- Due to uncertainty in future licensing costs for GPFS, exploring alternatives (at least) for disk-only part
 - Licence cost affordable for Run 3 but unclear for HL-LHC
- CEPH and EOS as possible GPFS replacement
 - First step: finalize present test-bed of CEPH
- Also the hardware infrastructure depends on the storage management model
 - Few servers with big buckets of storage vs. many storage bricks
 - Which in turn impacts the network level (add an aggregation layer?)
- A clear understanding of TCO of storage model is mandatory
- Evaluation of object storage for users/small experiments

Storage and data management (3/3)

- In view of building the INFN data lake, need to investigate services of replication and synchronization among sites
 - First target are small experiments and scientific communities not able to invest a very sophisticated data management system
- FTS is the natural candidate (activity covered in ESCAPE and DOMA)
- The proposal discussed with CERN is its integration with Sync&Share services
 - CNAF would therefore collect requirements, deploy a prototype FTS service, work on its certification involving the user communities and close the feedback loop with the developers at CERN.

- Current LAN based on a star-centered topology (Edge-CORE)
 - 100 Gbps LAN
 - Core composed of 2 interconnected Cisco NEXUS 9516 switches (VPC)
 - Disk-servers (2x10 Gbps -- 2x100 Gbps) and ToR switches (2x10 Gbps) for CPU nodes and services directly connected to both core switches
 - Complete separation between storage and CPU islands
- What topology to adopt for the scale of HL-LHC not clear
 - Introduce another level of switches for disk-servers?
 - Mixed storage-CPU island?
 - Different Topology? Spine Leaf?
 - Scale to nx100 Gbps / 400Gbps?
 - Virtualization of the network?

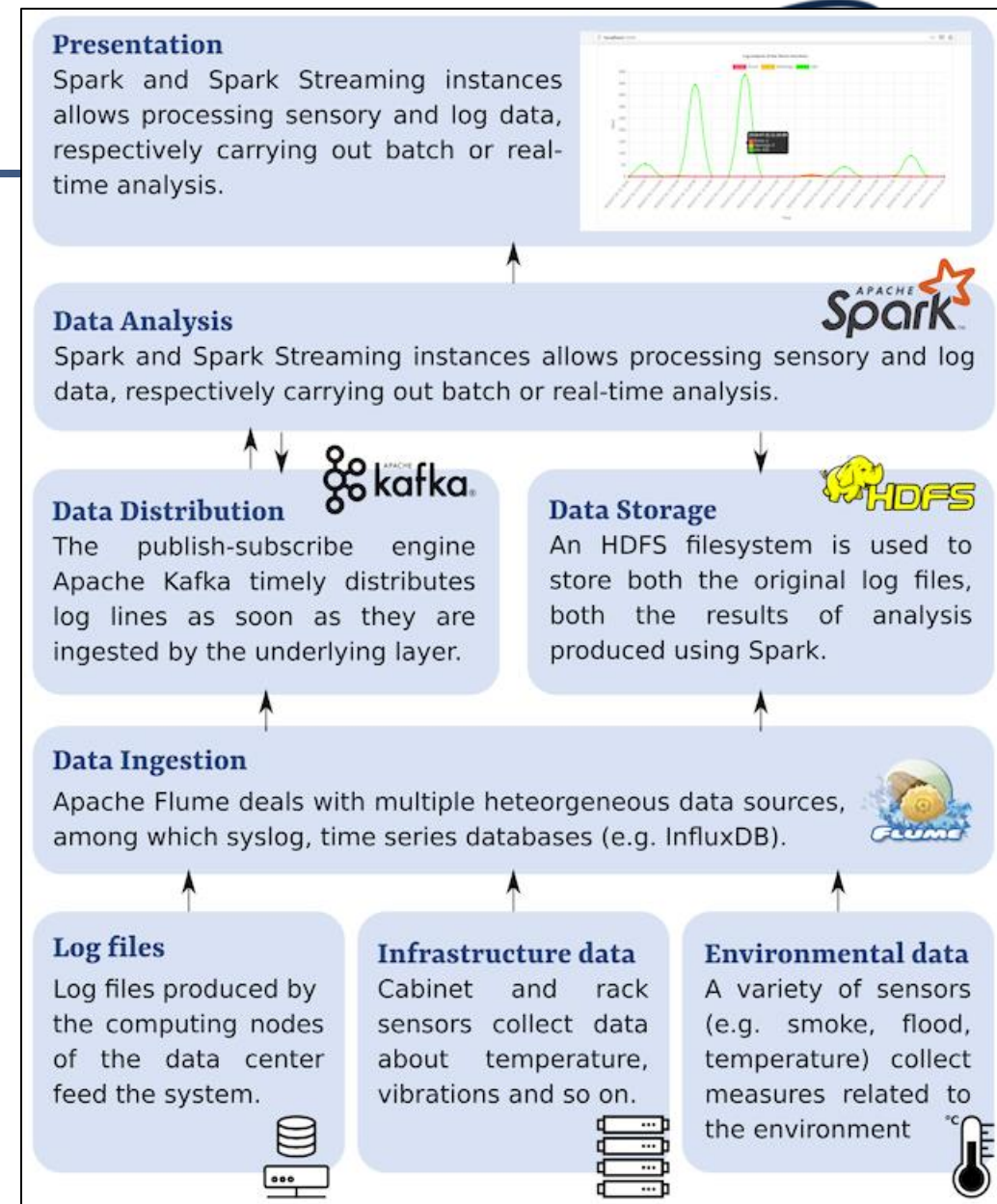
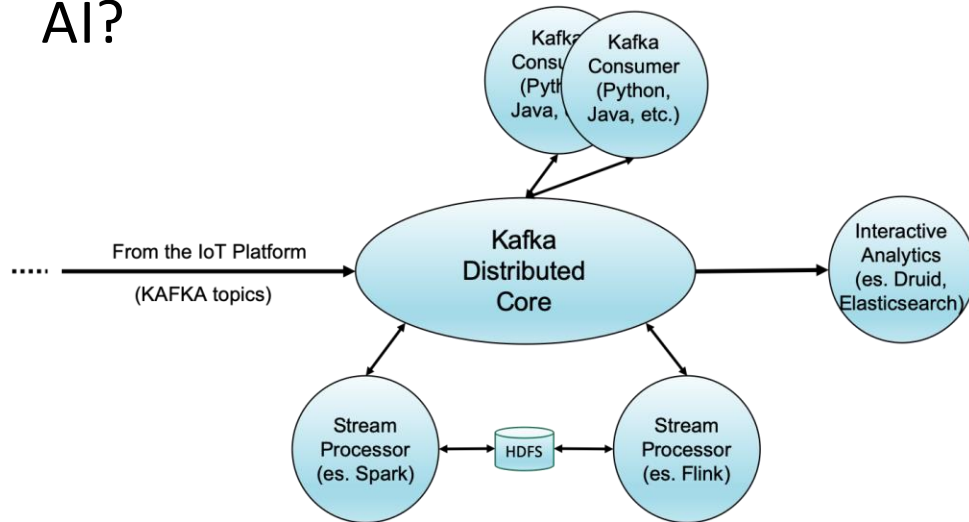


Current AAI infrastructure

- Three sources of authority
 - Data center AAI based on Kerberos (Authn) and LDAP (authz)
 - Grid access via “standard” VOMS
 - Cloud access granted via IDEM + IAM
- Aim: consolidation/unification
 - Studying convergence to consistent AAI system (possibly integrated with IAM)
- Focus on keycloak/FreeIPA+INDIGO-IAM
 - Also add the possibility to use INFN IDP for authentication

Predictive maintenance

- Great interest in predictive maintenance
- First steps :
 - Following the CERN line
 - Log collection
 - Tests on correlation of logs
 - AI?



Strategy toward Tecnopolo

- Probably the strategy adopted so far will be not enough
 - Offer solutions common to all users (simplification)
 - Adoption of industry standards (robustness)
 - Implement redundancy for all services and systems (resiliency)
- Study and test cloud/virtualization paradigm
 - Es. synergy with CERN-IT and WLCG for K8s



Tecnopolo Services

Overall Service Deployment

- **Guideline:** starting from the current Tier-1 architecture, we want to **move toward “Cloud” interfaces** (Cloud: remote access, virtual everything, simpler and autonomous resource provisioning) and **simplify service deployment**.
 - We have been using Kubernetes for several years for both dev and ops in many services (e.g. StoRM, VOMS, Argus, INDIGO-IAM)
- **Activity:** evaluate if and how to move the entire data center services (farm and storage) toward a **k8s-based structure**.
 - Configuring Kubernetes (or any other container orchestrator) “right” for large scale installations is not trivial. We need to gather experience with this and find proper automatization / deployment strategies for k8s.
 - Main goals: resiliency and automated CI/CD pipelines (e.g. system and service upgrades, rollbacks, etc.).

Computing and Storage Resources

- **Guideline:** with or without k8s, we want to **move as much as possible to containers** (both services and computing resources).
- **Activity:** work out **allocation / reclaim policies** for resources.
 - How to connect this with accounting tools measuring pledged / opportunistic resources? (for WLCG or else)
 - GPUs are a precious and scarce asset – we need to better manage their allocation policies (take ML-INFN as an application example).
 - Understand interworking between OpenStack and k8s.
 - Can K8s be “the next batch system”?

Data Management

- **Guideline:** there are diverse user requirements (see next slide), where POSIX is still fairly requested. However, the move toward a **“Cloud-based” approach to storage** is very important.
- **Activity:**
 - POSIX is still requested but see the box on the right **e.g. wrt CEPH.**
 - Consider **multi-VO RUCIO.**
 - Enhance the current usage of remote access protocols (http, xrootd) **with dynamic instantiation of http and xrootd caches** to optimize access, and with progressive adoption of **object-based storage services.** This requires also actions at the **user support** level.

Ceph Object Store

- Question: Are Object stores fundamentally different from traditional Grid storage?
 - Answer: With WLCG use cases not really.
- Ceph uses algorithmic data placement:
 - No central catalogue for meta data queries.
 - Vector reads not supported.
- Very few WLCG use cases need a file system and it is mostly about educate user / fixing bugs due to invalid assumptions.

RAL QoS @ DOMA, 26/2/2020

CephFS, local disk throughput

- Current CephFS: 280 HDD, 750 TB, 7 servers
 - Metadata on 20 small SSDs
- Bottlenecks:
 - Can reach up to 20k IOPS, ~4GB/s (current LAN on nodes is the limit)
 - Before ceph wpg, frequent problem with slow requests
 - Currently: 2 OSD HDD/batch node - faster for input than single local HDD
 - Too slow for workdir (large mds stress, frequent small iops)
- Node size “problem” with upcoming hw
 - 128C/256HT Rome, more in the future - ~4000 hs06/node
 - Local HDDs out of question
 - ATLAS heavy jobs use 2-3Gb/s LAN
 - Local disk: 5TB with the WLCG recommendations, expensive for SSD/NVMe - fast shared FS might be cheaper and more performant

JSI QoS @ DOMA, 7/2/2020

Exp vs protocol vs AAI matrix @CNAF

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S |
|----|------------------|---------------------------|------|------------|------------|-------------------------|---------|------------|------------|--------------|-------|------------|---------|------------|----------|------|------------|-----------------------------------|--|
| 1 | Experiment | Data Management protocols | | | | Data Transfer protocols | | | | Technologies | | | | AAI | | | | Specific needs/problems | |
| 2 | | POSIX | SRM | WebDAV | XrootD | POSIX | GridFTP | HTTP | XrootD | StoRM | GEMMS | XrootD | GridFTP | dataclient | VOMS | GRID | dataclient | token-based | |
| 3 | ALICE | no | no | no | yes | no | no | no | yes | no | no | yes | no | no | no | no | no | no | |
| 4 | ATLAS | no | yes | yes | yes | yes | yes | soon | yes | yes | yes | yes | yes | yes | no | yes | no | no | interested |
| 5 | CMS | no | yes | no | yes | yes | yes | no | yes | yes | yes | yes | yes | yes | no | yes | no | no | interested |
| 6 | LHCb | yes | yes | no | yes | no | yes | no | soon | yes | yes | yes | yes | yes | no | yes | no | no | have to be interested |
| 7 | | | | | | | | | | | | | | | | | | | |
| 8 | white: no answer | | | | | | | | | | | | | | | | | | |
| 9 | AGATA | yes | tape | no | no | yes | tape | no | no | tape | yes | no | no | no | yes | no | no | no | |
| 10 | AMS | no | yes | no | no | no | yes | no | yes | yes | yes | yes | no | no | yes | no | no | interested (see DODAS) | |
| 11 | ARGO | yes | yes | no | no | yes | yes | no | no | tape | yes | no | no | no | yes | no | no | no | closing |
| 12 | AUGER | maybe local users | yes | no | no | no | yes | no | no | yes | no | no | no | no | yes | no | no | no | no tape, use Dirac |
| 13 | BELLE | no | yes | yes | no | no | yes | yes | no | yes | no | no | no | no | yes | no | no | interested | tape not used |
| 14 | BOREXINO | yes | no | no | no | yes | no | no | no | no | no | no | no | no more | no | no | no more | no | tape not used, closing |
| 15 | COMPASS | no | yes | interested | interested | no | yes | interested | interested | yes | no | no | no | no | yes | no | no | interested | no tape |
| 16 | CORELIB | yes | no | no | no | no | yes | no | no | no | no | no | yes | no | no | yes | no | no | no tape |
| 17 | COSMO_WNEXT | yes | no | no | no | no | yes | no | no | no | no | no | yes | no | no | yes | no | no | tape not used |
| 18 | CTA | no | yes | no | future | no | yes | future | no | yes | yes | no | no | no | yes | no | no | interested | valido per MC e DIRAC, non per utenti locali |
| 19 | CUORE | yes | no | no | no | yes | no | no | no | no | no | no | no | no | no | no | no | no | no tape, rsync (100G/day) |
| 20 | CUPID | yes | no | no | no | no | yes | no | no | no | no | no | yes | no more | no | yes | no more | no | tape not used |
| 21 | DAMPE | yes | tape | no | no | no | yes | no | no | tape | yes | no | yes | no | yes | yes | no | interested | |
| 22 | DARKSIDE | yes | no | no | no | no | yes | no | no | no | no | no | yes | no | not used | yes | no | | tape? |
| 23 | ENUBET | yes | no | no | no | yes | no | no | no | no | no | no | no | no | no | no | no | interested | no tape |
| 24 | FAMU | yes | no | no | no | yes | yes | no | no | no | no | no | yes | no | no | yes | no | interested | no tape |
| 25 | GERDA | yes | tape | no | no | no | tape | no | no | tape | yes | no | no | no | yes | no | no | interested | |
| 26 | GLAST | no | yes | no | interested | no | yes | no | interested | yes | no | no | no | no | yes | no | no | | tape not used |
| 27 | ILDG | no | yes | no | no | no | yes | no | no | yes | no | no | no | no | yes | no | no | no | no tape |
| 28 | ICARUS | no | yes | no | no | no | yes | no | no | yes | yes | no | no | no | yes | no | no | | |
| 29 | JUNO | yes | yes | yes | no | yes | yes | yes | no | yes | yes | no | yes | no | yes | no | no | maybe in the future | no tape |
| 30 | KM3NET | yes | yes | no | interested | no | yes | no | no | yes | no | interested | yes | yes | yes | yes | yes | interested | tape not used |
| 31 | LHAASO | yes | no | no | no | yes | no | no | no | no | no | no | no | no | no | no | no | no | no tape, closing |
| 32 | LHCf | yes | no | no | no | yes | no | no | no | no | no | no | no | no | no | no | no | no | no tape |
| 33 | LIMADOU | yes | no | interested | no | no | yes | interested | no | no | no | no | yes | no | no | yes | no | interested | no tape |
| 34 | MAGIC | no | tape | no | no | yes | tape | no | no | tape | yes | no | no | no | yes | no | no | no | convergono con CTA |
| 35 | NA62 | no | yes | no | no | no | yes | no | no | yes | no | no | no | no | yes | no | no | only if it becomes DIRAC standard | tape not used; use xrootd (in other sites) |
| 36 | NEWCHIM | yes | no | no | no | no | yes | no | no | no | no | no | yes | no more | no | yes | no more | | tape: write to buffer with guc. Recall? |
| 37 | PADME | no | yes | no | interested | no | yes | no | interested | yes | yes | no | no | no | yes | no | no | interested | |
| 38 | PAMELA | no | tape | no | no | yes | tape | no | no | tape | yes | no | no | no | yes | no | no | no | |
| 39 | THEOPHYS | no | tape | no | no | yes | tape | no | no | tape | yes | no | no | no | yes | no | no | no | |
| 40 | VIRGO | yes | yes | yes | no | yes | yes | yes | no | yes | yes | no | yes | no | yes | yes | no | interested | also a storage area for storm no voms |
| 41 | XENON | no | yes | no | no | no | yes | no | no | yes | yes | no | no | no | yes | no | no | | |

Network and Security

- **Guideline:** the current Tier-1 Data Center is built around a traditional star-shaped, VLAN-based topology (manually maintained). However, the shared nature of the Tier1 and the diverse requirements of multiple tenants could be simplified through **a more virtualized vision of the network.**
- **Activity:**
 - Integrate or extend the **new INFN AUP for IaaS and other Cloud layers.**
 - Should we consider **other networking topologies**? A key point: integrate at the scale of a very large data center **automated, virtual networking, implemented at the “middleware layer”** (e.g. OpenStack, Kubernetes).
 - Deploy **scalable network monitoring / security facilities** and work out the implications of providing root (or “Cloud-native”) access to resources (IaaS, PaaS, SaaS, FaaS).

- **Guideline:** we want to **evolve the current LDAP+Kerberos –based AAI system** in use at CNAF, as well as the X.509-based approach for job AuthN/Z.
- **Activity:**
 - We already have some experience, but with a limited set-up, with **FreeIPA + Keycloak + INDIGO-IAM**. We now need to work out if and how to scale this to the whole data center.
 - Finalize **token-based authentication and authorization with HTCondor**.

Data Lake & National Services

- **Guideline:** we want to **integrate two complementary initiatives within INFN:**
 - **IDDLS**, i.e. a “proto-data-lake” based on a nation-wide photonic layer, in collaboration with GARR (with provider-based quasi-dynamic reconfiguration of links).
 - **INFN Cloud**, an initiative linking the larger INFN data centers with a uniform set of high-level services offered to our user base. Note: this is NOT “just another INFN-CC”.
 - Key points: simplify resource usage, service offering, user support and allow user-level service customization (see next slide).
- **Activity:**
 - Key **CNAF participation to the 5 INFN Cloud WPs** (Infrastructure and Architecture, Support and Documentation, Monitoring and Accounting, Security and Rules of Participation, Maintenance and Evolution).
 - Exploit both **IDDLS and INFN Cloud for early adopters.**

INFN Cloud




Transparent, multi-site federation for users of Cloud resources belonging to INFN and/or to other Cloud providers (private or public)

Authentication *can* be enabled for::

- Local username/password
- Google accounts
- EduGAIN (e.g. University, research centers, etc.)
- Other OIDC providers

Access to the Cloud services through a common dashboard, with different views depending on the users / user groups.

Composed, high-level services easily customizable a configurable directly by users



INDIGO - DataCloud

Welcome to **dodas**

Sign in with your dodas credentials

Username

Password

Sign in

[Forgot your password?](#)

Or sign in with

Google

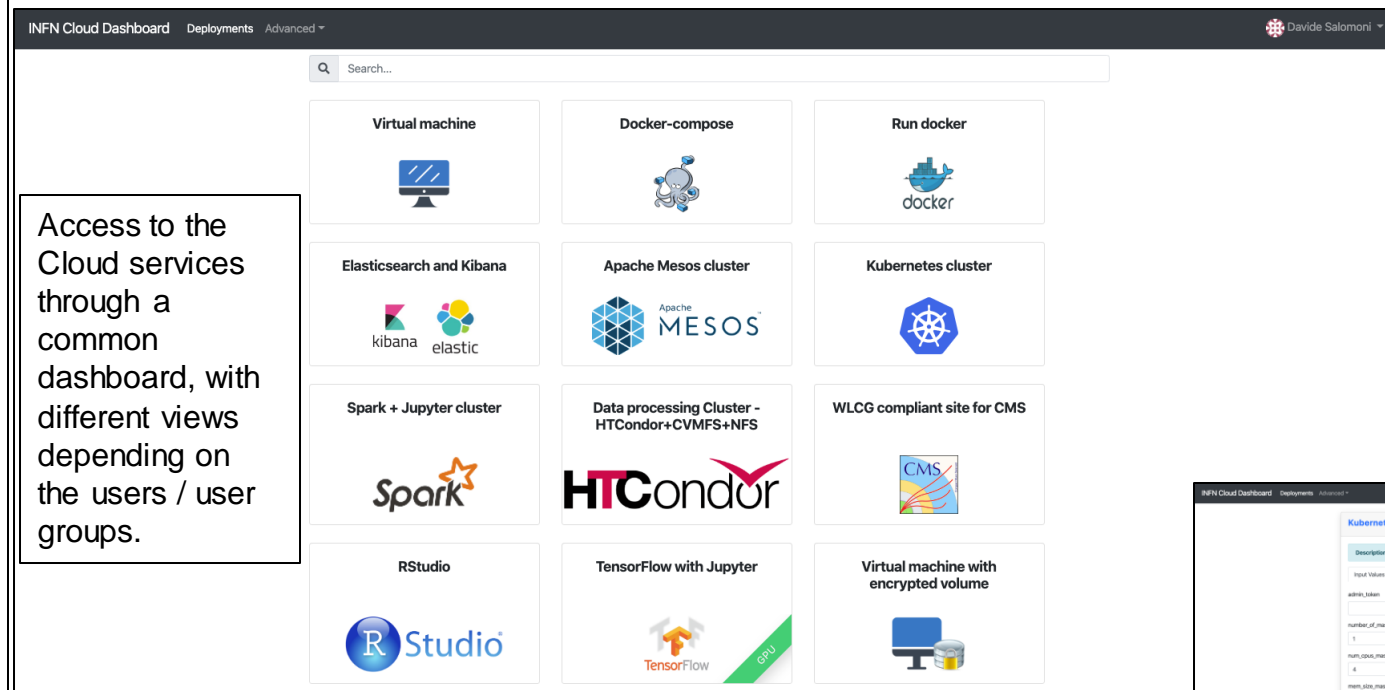
eduGAIN

EGI

Not a member?

Register a new account

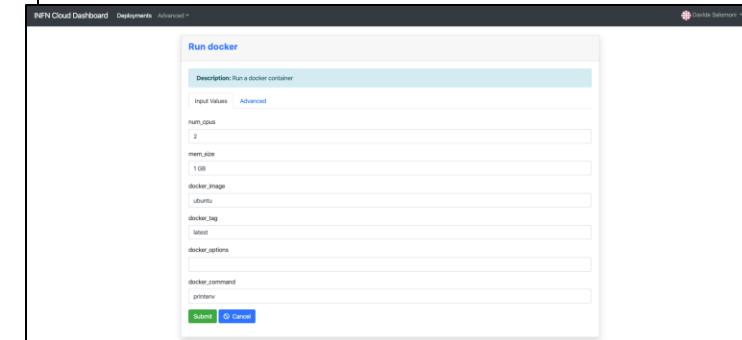
[Privacy policy](#)



INFN Cloud Dashboard Deployments Advanced ▾ Davide Salomoni ▾

Search...

- Virtual machine
- Docker-compose
- Run docker
- Elasticsearch and Kibana
- Apache Mesos cluster
- Kubernetes cluster
- Spark + Jupyter cluster
- Data processing Cluster - HTCondor+CVMFS+NFS
- WLCG compliant site for CMS
- RStudio
- TensorFlow with Jupyter
- Virtual machine with encrypted volume



INFN Cloud Dashboard Deployments Advanced ▾ Davide Salomoni ▾

Run docker

Description: Run a docker container

Input Values Advanced

num_cpus: 2

mem_size: 1 GB

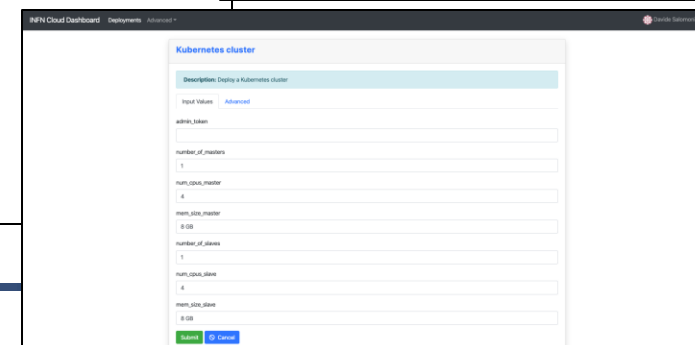
docker_image: ubuntu

docker_tag: latest

docker_options:

docker_command: python

Submit Cancel



INFN Cloud Dashboard Deployments Advanced ▾ Davide Salomoni ▾

Kubernetes cluster

Description: Deploy a Kubernetes cluster

Input Values Advanced

address:

number_of_masters: 1

num_cpus_master: 4

mem_size_master: 8 GB

number_of_slaves: 1

num_cpus_slave: 4

mem_size_slave: 8 GB

Submit Cancel

HPC-HTC Integration

- **Guideline:** we want to **more easily exploit CINECA-based HPC resources.**
- **Activity:**
 - Finalize and expand the several scenarios **making HPC and HTC integration as transparent as possible** for users.
 - Should we consider how to **facilitate the integration of HPC resources at the application level?** (how?)

ISO Certifications (for sensitive data)

- **Guideline:** so far, we have certified part of the CNAF Data Center (ISO 27001). We now want to **extend the certification toward other parts of the data center** (linking it also e.g. to ISO 31000 on Risk Management).
- **Activity:**
 - Rework the **network structure of the ISO area** so that scalability and proper isolation of tenants is realized.
 - Consider a **“Data lake ISO 27001 certification” for the INFN Cloud backbone**. This would allow us to support INFN-wide use cases related to increasingly important physics in medicine activities. This is relevant also for CNAF-specific activities.

Projects

- **Guideline:** we want to **continuously integrate the results of the many projects** running at CNAF.
- **Activity examples:** (see the Team “Progetti @ CNAF” for more info):
 - Enhance the **data and compute management platform** with **XDC** and **DEEP** results.
 - Introduce the **IoT Platform** from **We Light** and **IoTwin**s for e.g. environmental sensor monitoring.
 - Evolve the **Big Data Platform** from **ACC** and **ML-INFN**.
 - Profit from the **data lake experience (HPC, AAI, Rucio)** from **ESCAPE** and **IDDLS**.
 - Introduce high-level **TOSCA templates and user interfaces** from **INFN Cloud**.
- **Important note:** we are seeing **many requests to participate to projects**. We want to perform internal preliminary assessments to select the most promising ones and align scope and effort.

Collaboration with CERN: topics

1. **Large scale infrastructures management, including KPIs and monitoring: A. Chierici, C. Duma**
 - a) Kubernetes and virtualization: D. Michelotto, A. Ceccanti
 - b) GPUs: S. Dal Pra, T. Boccali
 - c) Network: S. Zani, P. Veronesi
 2. **Authentication and Authorization: A. Ceccanti, V. Ciaschini**
 3. **Storage and data management: V. Sapunenko, A. Costantini**
 - a) CEPH and EOS as possible GPFS replacement: Storage Team, SDDS
 - b) Remote access in the context of DOMA and ESCAPE: D. Cesini, E. Vianello
 - c) FTS as general purpose for general science: L. Morganti
 - d) CTA (only interesting if EOS is chosen): E. Fattibene
 4. **Integrating HPC in a WLCG site in a transparent way: T. Boccali, D. Spiga, S. Dal Pra**
 5. **Quantum simulation: F. Giacomini**
 6. **Data center certification: B. Martelli**
- **Activity:** complete the **description of work for these topics** before the next meeting between CERN and CNAF (March 18, 2020, videoconference).

Next steps

- Given the topics presented here, **this week we want to focus on:**
 - **Identifying chairs and participants** for the various Tecnopolo areas.
 - **Writing the first implementation work plans**, corresponding to the “Activity” bullets of the previous slides.

- **This work should be presented at the next Tecnopolo meeting, on 9/3/2020.**