# Storage for High Energy Physics
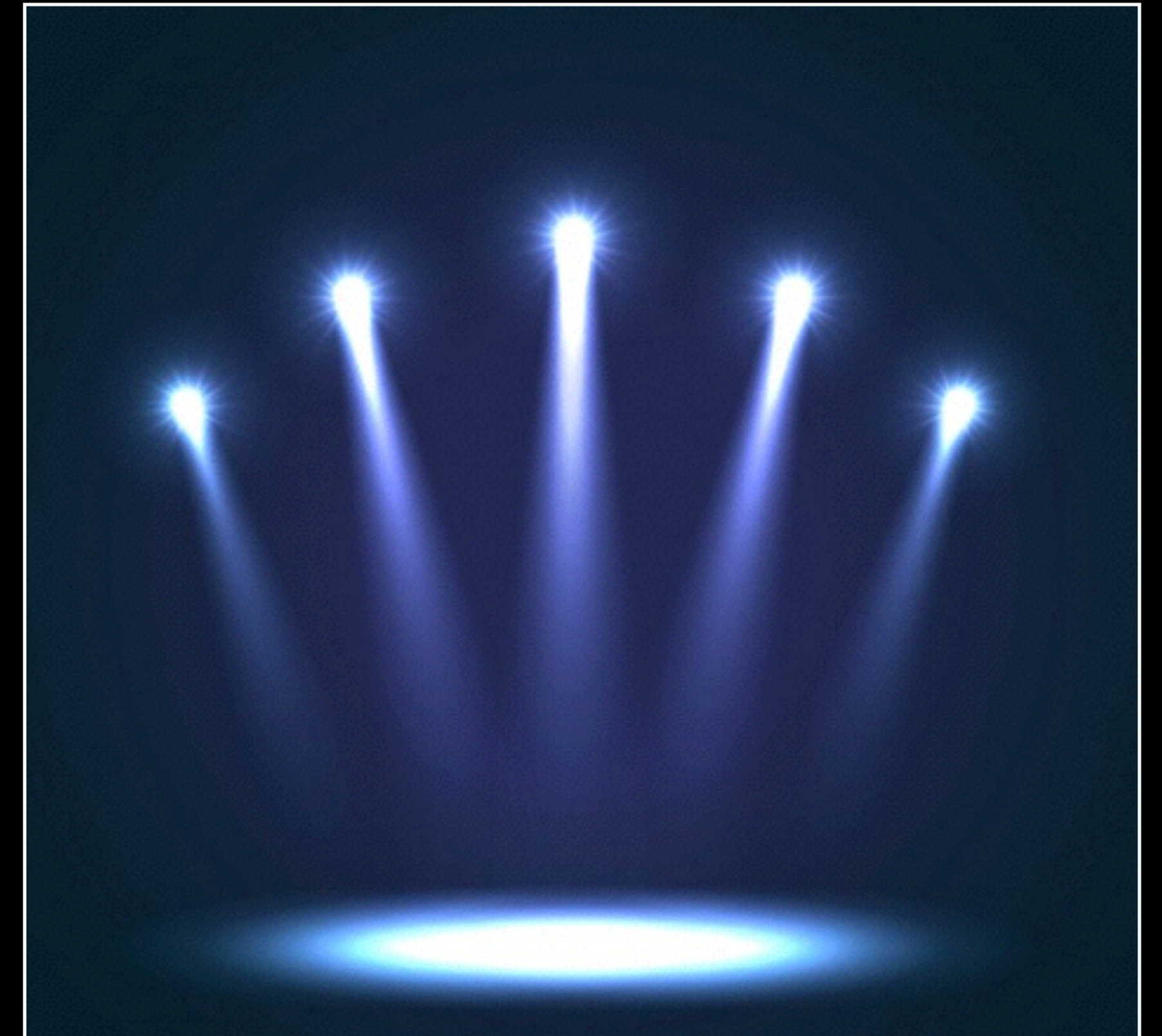
## CEPH DAY CERN 2019

**Andreas-Joachim Peters**

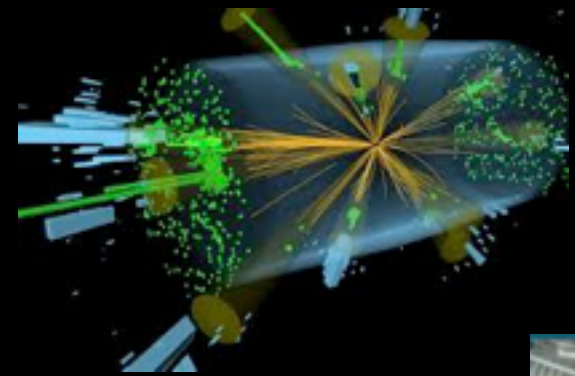CERN IT Storage Group

# Overview

- Storage in High Energy Physics

  - Past, Present & Future Challenges

  - Data Formats & Access Patterns

  - Storage Software Eco System

  - Remote Access Protocols & Security Mechanism

- Dedicated Storage Systems & Hardware

- Inventory, Summary & Vision

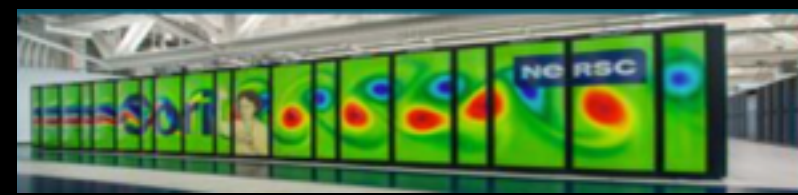# Storage in High Energy Physics


Archival & Backup Storage

Storage for Data Acquisition
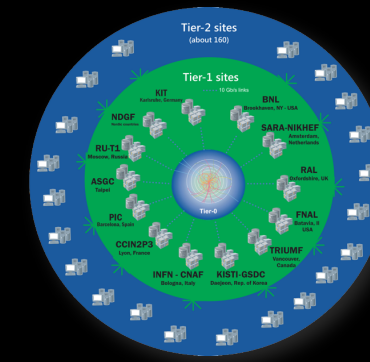

Storage for Home Directories


Storage for HPC


Storage for Applications

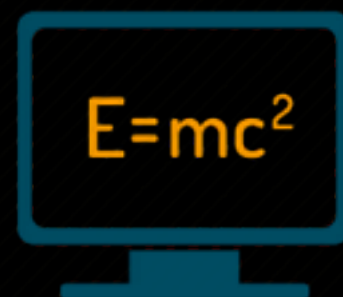Private Cloud Storage


Storage for GRID Computing

Public Cloud Storage

Storage for Software Distribution

Storage for Data Analytics

Storage for Physics Analysis

Storage for Sync&Share

# What is HEP Storage?

- **modular stack of layered services** on top of **disk & tape** based persistent storage

- very **heterogeneous resource** composed of s.c. GRID & non-GRID resources

  - various storage systems in computer centres
    **filesystems, file stores, object stores, tape storage**

  - temporary / time-limited

    - HLT/HPC facilities burst buffers/scratch space

    - public cloud R&D CERN Openlab …

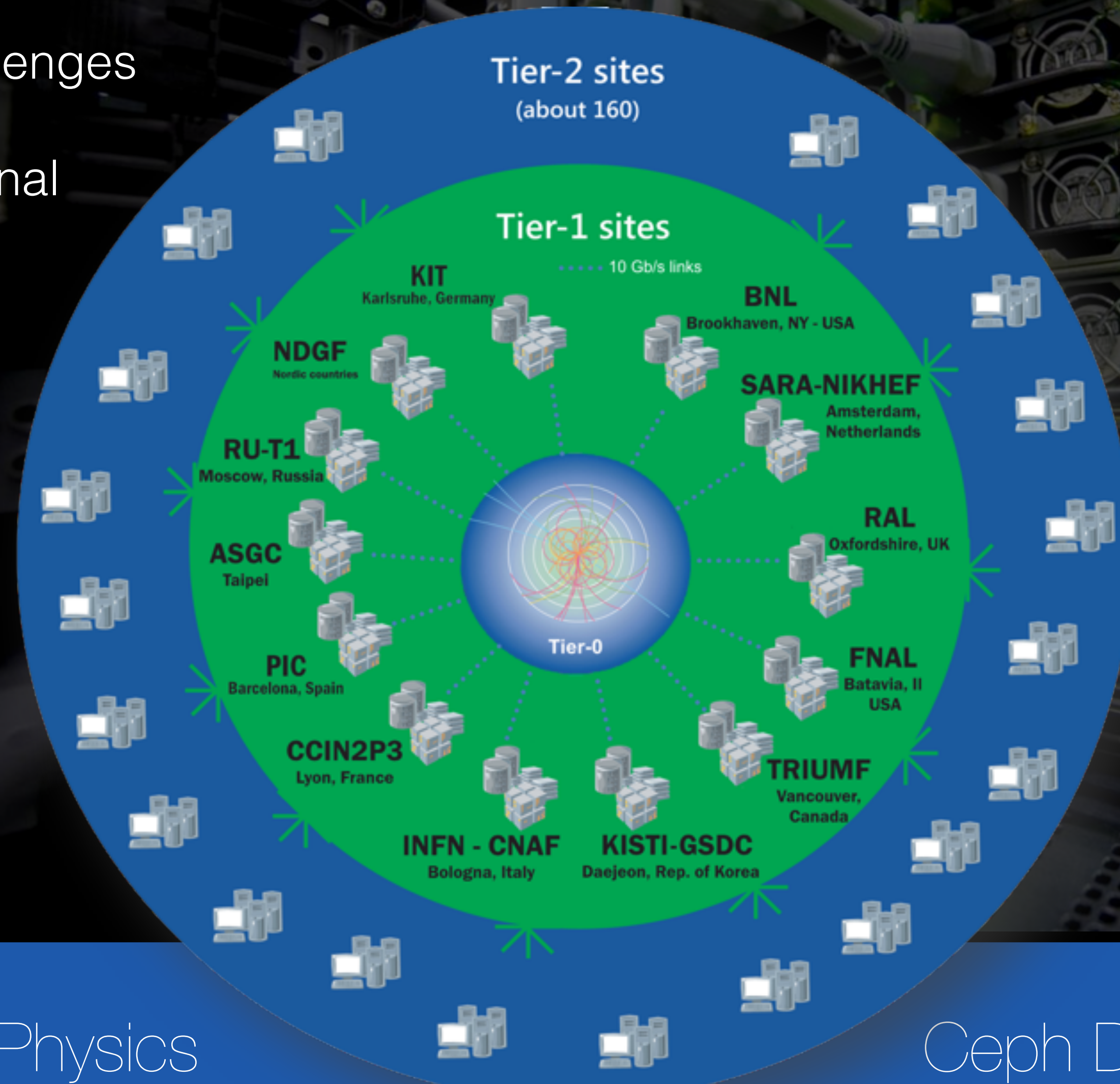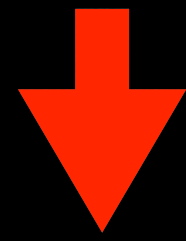# Origins of HEP GRID storage architecture

## Distributed Computing Model

Since early 2000 HEP embraced
'**the GRID**' model to tackle LHC challenges

- federation of national and international GRID initiatives

- in line with funding structure



WLCG
Worldwide LHC Computing Grid



Tier-2 sites
(about 160)

Tier-1 sites
10 Gb/s links

KIT
Karlsruhe, Germany

BNL
Brookhaven, NY - USA

NDGF
Nordic countries

SARA-NIKHEF
Amsterdam, Netherlands

RU-T1
Moscow, Russia

RAL
Oxfordshire, UK

ASGC
Taipei

Tier-0

FNAL
Batavia, Il USA

PIC
Barcelona, Spain

CCIN2P3
Lyon, France

TRIUMF
Vancouver, Canada

INFN - CNAF
Bologna, Italy

KISTI-GSDC
Daejeon, Rep. of Korea

# HEP Computing Community

- **GRID** resources are shared among many experiments and sciences : Network - CPU - Storage

- **LHC** experiments **consume >> 90%** of the accounted computing capacity

- The remaining part is consumed by 155 other experiments/sciences



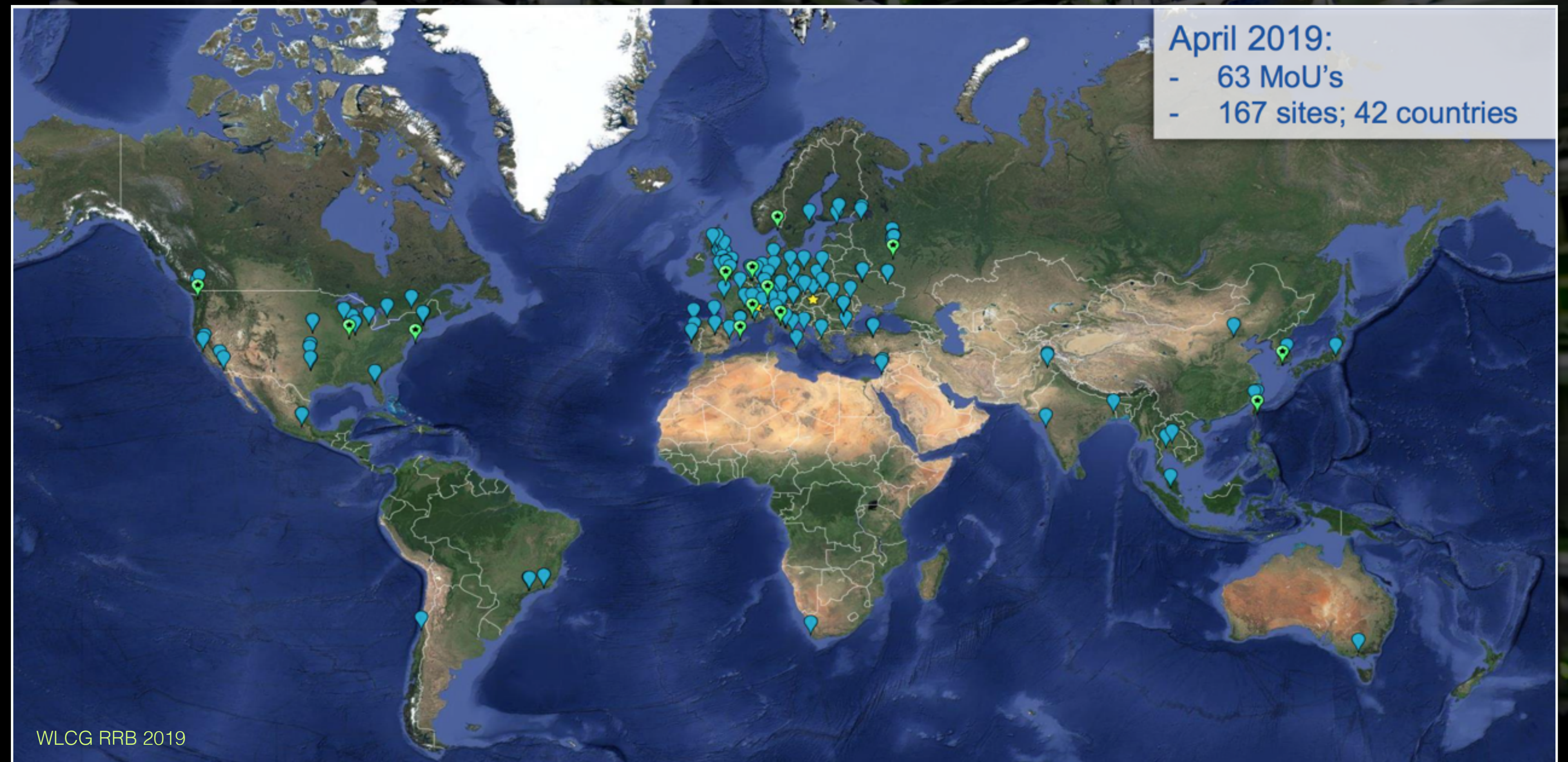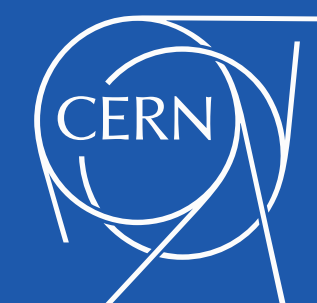Today **LHC** is in a leading position in steering evolution of infrastructure

# Worldwide LHC Computing GRID

- **WLCG** is a shared resource for ~12.000 physicist

  - 1M cores

  - 1EB storage



April 2019:
- 63 MoU's
- 167 sites; 42 countries

WLCG RRB 2019

S.Campana Granada

# Data Distribution Model
## Usage of the **LHC** **O**ptical **P**rivate **N**etwork

**Data Distribution in WLCG**

- global transfer rates exceeding regulary **60 GB/s**

- **830 PB** and 1.1B files transferred until end of Run 2

- main **challenge** is to have the **useful data close** to available computing resources match storage/compute/network

Running jobs: 365644
Active CPU cores: 807139
Transfer rate: 21.54 GiB/sec

# The HEP Data Volume Challenge

## LHC future: **ten times more data** … but there is more than just LHC …



Several experiment will require relatively large amount of resources in the future.
Several factors less then HL-LHC

Comparable data volume to LHC

FAIR Disk Storage Requirements

DUNE foresees to produce ~70PB/year in the mid 2020s

*J. Eschke @ ESCAPE kick-off*

**Storage**

*Y. Kato @ HOW 2019*

S.Campana Granada

# The LHC Storage Cost Challenge



**ATLAS** Preliminary
Disk resource needs
- 2017 Computing model
- 2018 estimates:
  - Baseline model
  - Reduced storage model
- Flat budget model (+15%/year)

## LHC Challenge of High Luminosity Run 3/4

Assuming a flat budget, the storage requirements of the future cannot be satisfied anymore by technology evolution

Possible improvements
- **changes** to **computing models**
  - **media shift** towards cheapest media combined with experiment driven storage tiering (data carousel)

- **capacity** increase using erasure encoding and optimised replication strategies
  ( but often EC ~ RAID volume )

- **Data lake/cloud** model
  - composed of fewer HA storages and satellite caches
  - more efficient storage management and optimised redundancy using QOS interfaces

# "Data Lake" exploring geographical distributed storage systems

n sites
k replicas
k<<n
latency > 1ms
caches + persistent storage

13

More background to Data Lake / DOMA activities

WLCG Report 2019

# HEP Data Volume in perspective



50 — LHC Raw 2016
200 — LHC Science
98 — Google Searches
180 — Facebook Upload
300 — SKA Phase 1 data/year - 2023
600 — LHC HL-Run Raw 2026
1'000 — LHC HL-Run Physics Data 2026
1'000 — SKA Phase 2 Science Data
15'000 — Google Internet Archive

10000 PB
100 PB
1 PB

S.Campana Granada

# Tape Archive at CERN

prediction

- 5000 PB
- 3750 PB
- 2500 PB
- 1250 PB
- 0 PB

2018

2021-22

2024

2027

2030

4'300 PB

G. Cancio - IT-ST

Storage for High Energy Physics

Ceph Day CERN 2019

# Storage Hardware

# Physics Data Storage Hardware CERN

- Profiting from **_economy of scale_**
  - minimise price per GB
- System Unit:
  - 8 physical cores (16 virtual) 64-128GB RAM
  - disk-tray of 24x 4-6-10-12TB HDDs

- Running different generations
  - 2 trays per system unit - 48 disks
  - next gen.: **2 dense trays** per system unit - 120 disks
  - **4 trays per system unit - 96 disks**
  - 8 trays per system unit - 192 disks
    up to 2.4 PB per disk server

# Tapeless Archive Project

@ KISTI / S. Korea

**Goal**:
low-cost **disk-only erasure encoded archival storage** requiring deletion/integrity safety features

10 nodes - 15 PB usable space

# Tape: Lower Cost & Data Safety

B. Panzer-Steindel CERN IT CTO

## 100 PB storage example I

**Processing Cluster 5000 nodes**

~50000 streams

**1 PB Storage Cache (SSD)**

~40 streams

**IBM tape library: 40 LTO tape drives, 100 PB tape media**

Tape infrastructure requires disk front-end storage
Requires 'impedance' matching between clients and tape drives → SSDs needed
Total cost estimate : ~1.7 MCHF
Performance : ~ 10 GB/s        116 days to read 100 PB
One drive for ~200 tapes

~50000 streams

100 PB disk storage distributed across 44 disk server
Total cost estimate : ~3.5 MCHF
Performance : ~440 GB/s
3 days to read 100 PB
One server for ~200 disks

**100 PB Storage Cluster (HDD)**

20. September 2018                    Bernd Panzer-Steindel, CTO CERN/IT

- still (50%) cheaper
- physical deletion is slow
- however
  - single vendor problem (enterprise)
  - media shipments shrinking

**The Register**
*Biting the hand that feeds IT*

Data Centre ▸ Storage

**Did Oracle just sign tape's death warrant? Depends what 'no comment' means**

Big Red keeps schtum over the status of StreamLine

By Chris Mellor 17 Feb 2017 at 10:44    29 ☐    SHARE ▼

Oracle's StorageTek (StreamLine) tape library product range will be end-of-lifed, El Reg has learned.

**Yearly Cartridge Shipments**

# Storage Software Development Projects in HEP

https://rucio.web.cern.ch

**RUCIO**

SCIENTIFIC DATA MANAGEMENT

LEARN MORE

**Data Management**

Storage for High Energy Physics      Ceph Day CERN 2019

CS³

Home | Programme | Community | Abstracts | Committees

www.cs3community.org

Cloud Services for Synchronisation and Sharing

28 - 30 January 2019, Roma

Previous Workshops

Krakow 2018 - Amsterdam 2017 - Zurich 2016 - Geneva 2014

CERN  Storage for High Energy Physics

Ceph Day CERN 2019

# HEP Storage Software Development

## Why is there so much HEP Storage related Software Development?

HEP

| Amga |
| AliEn |
| Castor |
| CERNBOX |
| CTA |
| DAVIX |
| dCache |
| Dirac |
| Dynafed |
| DPM |
| EOS |
| FTS |
| GFAL |
| Phedex |
| ROOT |
| Rucio |
| VOMS |
| XCache |
| XRootD |
| … |

BigData

| AI Store |
| AlluxIO |
| AWS |
| CEPH |
| GCS |
| Hadoop |
| SkyllaDB |
| Spark |
| … |

minimal overlap in
development projects

many projects date back to the GRID area before BigData skyrocket

# HEP GRID Storage Ecosystem

**Physics Applications / Storage Clients**

ROOT | GFAL | DAVIX | XRootD

**Storage Applications**

Sync & Share | Jupyter Notebooks

CERNBox | Nextcloud | SWAN

Software

CernVM File system

**Data Management Services & Global Namespaces**

Alien | Rucio | Dirac | Indigo | Phedex

**File Transfer Service**

FTS

**Auth/Authz / Authz Translation**

DYNAFED | Authz Token | Macaroons | VOMS | OAUTH2

**Remote Access Protocols**

S3/GCS | DAV(S) | HTTP(S) | gridFTP | NFS4 | XRootD

**File Storage Services**

XrdCeph | DPM | XRootD/EOS | dCache

**Cloud Storage**

**File Systems**

CTA/Castor

HPSS

# HEP & BigData Technology

- If we would use BigData Analytics in physics, we could profit from all the existing BigData storage technologies, protocols & analytics frameworks

- Why is that not yet mainstream?

# Physics Data Formats

**unstructured raw data** - each physics event is stored in a compound block - events are assembled during data taking from many detector systems



| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |

File

N

compressed raw data with individual events

**structured data** - data is stored optimised for volume and access patterns

1..8

| z1 z2 z3 z4 z5 z6 z7 z8 | y1 y2 y3 y4 y5 y6 y7 y8 | x1 x2 x3 x4 x5 x6 x7 x8 | . . . . . | a1 a2 a3 a4 a5 a6 a7 a8 |

File

C(Z) C(Y) C(X) . . . . . C(a)

File

√ROOT
Data Analysis Framework

# Data Formats & Storage Access Patterns
## in selective analysis use cases

read pattern (read) in a selective physics analysis workflow

| | |
|---|---|
| **ROOT (inflated)** | 119 MB (7.99 %) |
| **ROOT (zlib)** | 67 MB (6.36 %) |
| **ROOT (LZ4)** | 77 MB (6.48 %) |
| **Protobuf (inflated)** | 1740 MB (100.00 %) |
| **Protobuf (gzip)** | 1177 MB (100.00 %) |
| **SQlite** | 1675 MB (100.00 %) |
| **HDF5 (row-wise)** | 1501 MB (100.00 %) |
| **HDF5 (column-wise)** | 98 MB (6.55 %) |
| **Parquet (inflated)** | 1502 MB (99.99 %) |
| **Parquet (zlib)** | 1322 MB (99.99 %) |
| **Avro (inflated)** | 1368 MB (100.00 %) |
| **Avro (zlib)** | 1058 MB (100.00 %) |

Jakob Blomer / CERN

- **sparse access** pattern
  *cry for* certain access
  protocol capabilities in
  LAN & WAN environments

- predictable read patterns
  allows to use **asynchronous
  multi-byte-range read** requests to
  to compensate latencies

- good news: most of traffic in
  HEP is still mainly sequential
  forward-seeking IO
  jobs@CERN like 100.000 people watching all a different movie
  with 1 MB/s streaming average

physics analysis uses high parallelism with relatively slow streams  ( tens of MB/s ) - no need for high throughput clients in the GRID

# Data Formats & Compression Algorithms



Compression speed vs Compression Ratio for compression algorithms

Shadura / Bockelmann

Test node: Haswell+    SSD

- **compression** done on **application side**

- **LZMA** cheapest for storage, most expensive for CPU

- best **algorithm** has to be selected **per use case** (de-/compression speed)

- compression inside storage systems rarely a benefit for physics data

- de-duplication marginal

# HEP Data Access Protocols

| prot/usage | LAN | WAN | WAN Transfer third party transfers |
|------------|-----|-----|------------------------------------|
| Mounted FS | high | - | - |
| XRootD | high | high | medium |
| HTTP(S) | low | low | comissioning |
| S3 | low | low | - |
| gridFTP | - | - |  |



## http://xrootd.org

- data **client/server framework**
  think of NFS or HTTP server written in C++ with own protocol

- optimised for remote access in **LAN & WAN**
  - arbitrary request redirection
  - third-party transfer between XRootD server with credential delegation

- front-end **protocol plugins** XrootD & HTTPs

- **storage** back-end **plugins** XRootD & HTTPs & S3

- **authentication** plugins (krb5, x509, sss, unix)

- **authorization** plugins (rule-based, tokens, macaroons)

- **proxy** & **cache** plug-ins, **clustering** support

# HEP Authentication & Authorization

CLIENT          SERVICE                                              STORAGE



Digital Certificate     Time Limited          Delegated Time Limited          Mapping          id
                        Proxy Certificate     Proxy Certificate                                uid/gid

                        VOMS    role                                                           local account in storage system

                                                                                               current way

who    what
for whom

                                                                                               id
                                                                                               uid/gid

user/pwd                                                                                       local account in storage system
social logins
GITHUB                                                                                         future way
...

        OIDC                    Storage Service                        • positive evolution:
        OAUTH2                                                           **adaption of industry standards** tokens/macaroons

---

Storage for future LHC Online Systems

# Storage for future LHC Online Systems

- high **capacity** & high **IO** requirements - confined environment

- wide range of solutions possible: from distributed high performance **parallel filesystems** to **object storage** - key issue **cost**

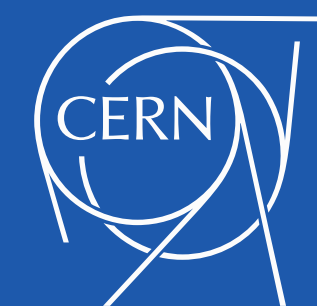| TDR | ALICE | ATLAS | CMS | LHCB |
|---|---|---|---|---|
| IO Rate | 200 GB/s | 60 GB/s | 61 GB/s | 100 GB/s |
| Capacity | 60 PB | 36 PB | 5.7 PB | 100 PB |

Storage for HPC

# Storage for HPC

- typical use case: **MPI applications** requiring low latency access & high stability
  - most of LHC related computing done with **HTC** (trivial parallelism via batch jobs)

- playground of high performance filesystems Lustre, Spectrum Scale and others
  - e.g. NERSC Lustre 700 GB/s
  - CERN 'exotic' pioneering with CephFS



**SCs** significant resource for opportunistic computing
a common problem is the availability of storage clients for these platforms (e.g. FUSE based filesystems, services for data injection and extraction) and the external connectivity of HPC facilities

# Storage in public clouds

# Storage in Public Clouds

- Public Cloud Services like AWS or GCS allow **time limited access to CPU resources**
  in times of high computing demands

  - **simplest use case simulations** mainly producing data

- Public Cloud Storage S3-like **easily integrated** as GRID resource
  **pricing** for storage and data access **not competitive** to replace HEP storage systems

- CERN successfully demonstrated **feasibility** of physics workflows in public clouds    CERN openlab collaborations

CERN scientists "rediscover" the Higgs boson live on stage at KubeCon using Google Cloud. Solution used Google Kubernetes Engine, Memorystore, and Storage (with network traffic peaking at 175G/s)! #k8s5    CERN R. Rocha et al



**CERN openlab**    **Google Cloud**

Google Cloud Storage → Cluster on GKE → Job Results → jupyter

**70 TB Dataset**    Cluster on GKE    Job Results    Interactive Visualization

Max **25000 Cores**

Single Region, 3 Zones

**25000 Kubernetes Jobs**    Aggregation

# Storage for HOME directories

# Storage for HOME directories

- several centres in HEP use commercial solutions like NetApp, Spectra Scale, DFS
  *hit by unexpected increase in license costs*

- CERN started replacing **DFS**

- **MALT** project: CERN strategy to decrease risk of vendor lock-in

- CERN also looking into long-term alternative for **AFS** *future unclear*



## The MALT Project
Re-assessing the IT provisioning Strategy for Core Services at CERN

### Increasing our technology, data and vendor independence

Our strategy, enacted through the MALT project, seeks open software solutions and products with simple exit strategies and low switching costs.

The project aims to deliver services inclusive of all the CERN community. The project's principles of engagement are to deliver the same service to every category of CERN user, to avoid vendor lock-in so as to decrease risk and dependency, to keep hands on the data and to serve the common use-cases.

https://malt.web.cern.ch/malt/

# Storage for Software Distribution

CernVM File System

**CernVM File System** is a network file system based on HTTP and optimized to deliver experiment software in a fast, scalable, and reliable way
- typical use case: need to start any kind of software in 100k batch jobs at the same time

https://cernvm.cern.ch/portal/filesystem



Jakob Blomer CERN

- client is implemented as a FUSE based filesystem

- works like a content delivery network

- the central repository is published in Ceph S3 at CERN

- very popular and widely adapted

- already more than a read-only filesystem for software

# CERN Storage Software for Tape

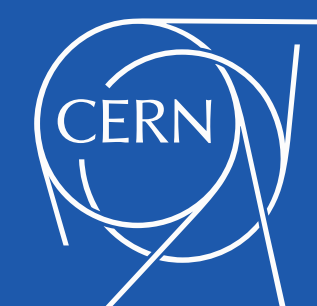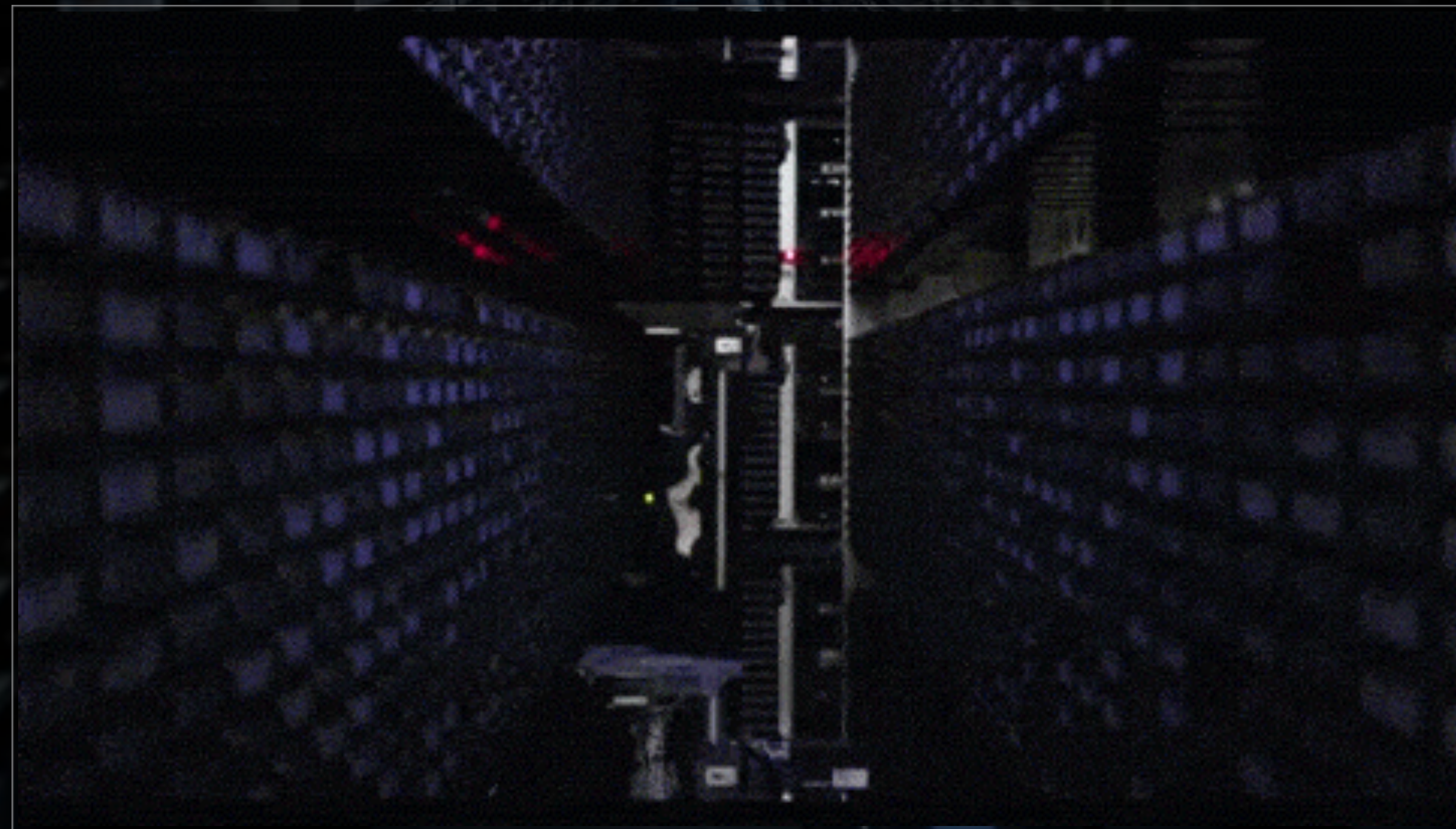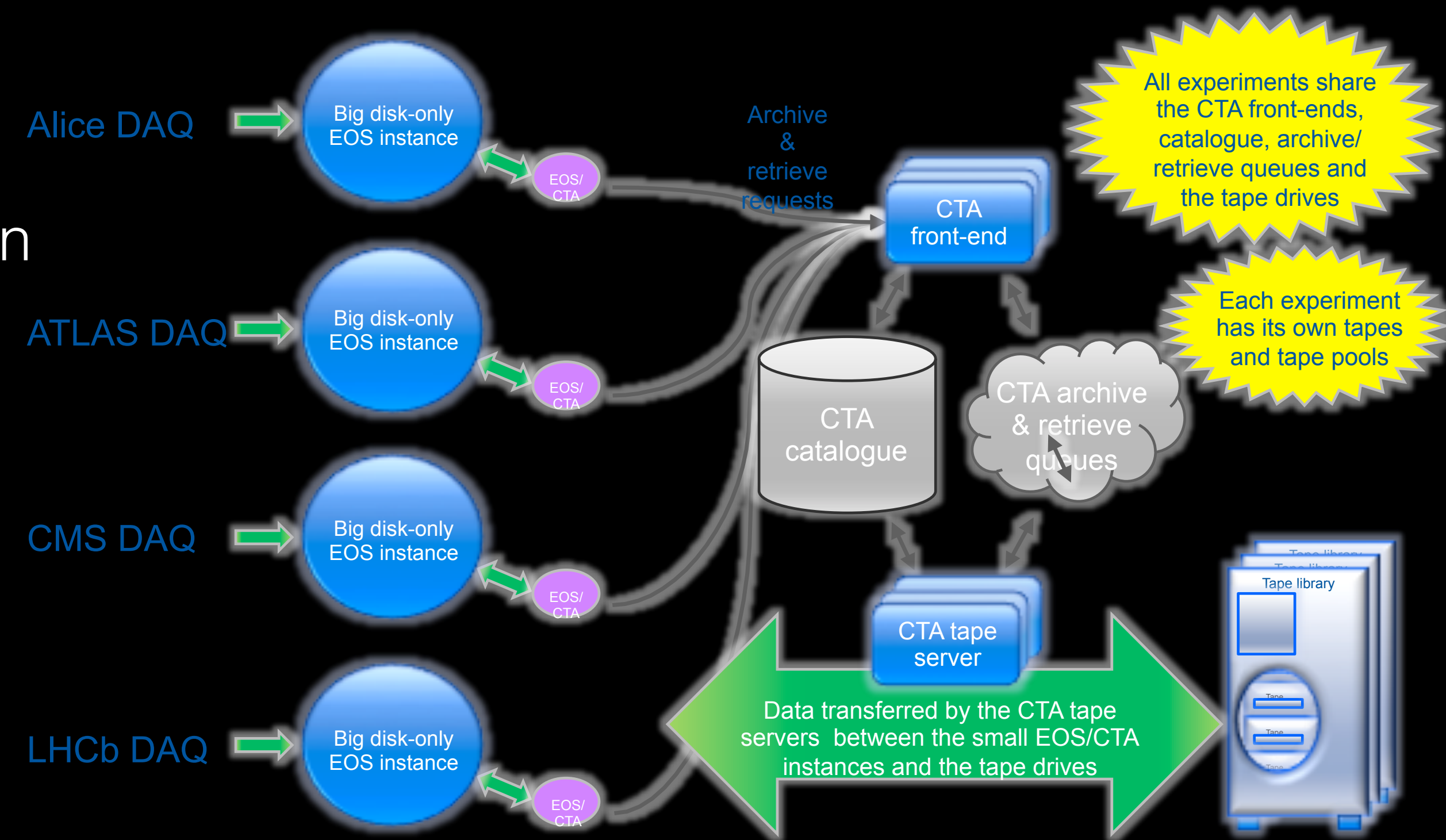# CERN Storage Software for Tape

CTA is the next generation storage software for tape archiving and backup use cases developed by the CERN storage group

- replacing CASTOR software after 20 years

- design is decoupling disk pool implementation from tape storage

- flat namespace

- no HSM model, no complex GC



All experiments share the CTA front-ends, catalogue, archive/ retrieve queues and the tape drives

Each experiment has its own tapes and tape pools

Alice DAQ → Big disk-only EOS instance → EOS/CTA

ATLAS DAQ → Big disk-only EOS instance → EOS/CTA

CMS DAQ → Big disk-only EOS instance → EOS/CTA

LHCb DAQ → Big disk-only EOS instance → EOS/CTA

Archive & retrieve requests

CTA front-end

CTA catalogue

CTA archive & retrieve queues

CTA tape server

Tape library

Data transferred by the CTA tape servers between the small EOS/CTA instances and the tape drives

# Inventory

- **HEP community** has built a very **modular stack of storage related software components**
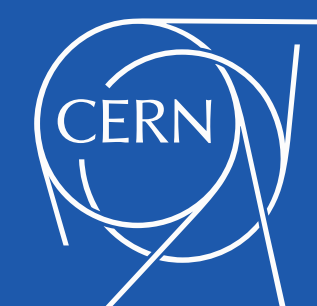  => allowed to **integrate** more or less **any storage solution** into the GRID

  - in the context of the **HEP Software Foundation** and Initiatives like the **XDC, IRIS** and **ESCAPE** projects many of these components are made **available to other sciences**



- **integrate-everything approach** is **not** the **most efficient**
  => inline with the community diversity

- **changes** in storage service implementations/technology **are slow** but happening
  => stateful nature of the service and limited resources to adapt new technologies

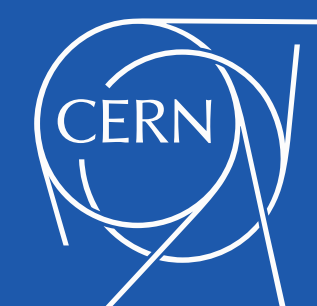  - lifespan of storage software in HEP is decades

# Inventory

- **file storage** is common denominator for most HEP storage systems

  - **site storage** systems have to fit also **requirements** of **local community**

- direct use of **object storage** appealing for physics analysis
  - diversity of the infrastructure is not compatible with global enforcement of object storage
  - **object storage** good match as a generic backend-solution - hide behind files
    find the economic/performant approach to gateways/security mechanisms

- Ancient community driven products/protocols **dying off**
  - **SRM** storage resource manager to handle tape storage access -  lack of success - decommissioning ongoing
  - **gridFTP** - globus filte transfer protocol - replacement in testing
  - **rfio**- remote access protocols initially used in **Castor** - decommissioned

# Inventory

- Storage Tiering: **HSM\*** model for GRID DM has died  *hierarchical storage management

  - no storage system can predict better access to data than the community who wants to use the data

    - manual(user driven) storage tier migrations have proven very successful and are part of data management frameworks used in HEP

- **Modern authentication & authorisation mechanisms** are slowly adapted by the community - they are not usable directly in filesystems - only via gateways

  - OIDC/OAUTH2

  - Macaroons

  - Tokens

# Summary & Outlook

- **HEP storage** is a **diverse universe** of commercial & open source storage components
  - diversity is a **blessing & a curse** at the same time

- **HEP storage delivered** required functionality & infrastructure for LHC & other experiments

- In the future **cost** is becoming a **hard limitation** on *what is doable* and at the same time more competition on resources

- **Tape** is the strategic media *still*

- **Network** is a strategic resource to enable remote access for storage & caching

- **HEP** community has ability to **shape the future** and aim for simpler & more efficient storage within budgets
  - 😉 open source technology helps in achieving this goal
  - 😨 amount of data produced is not influenced by storage technology but physics & computing models …

*Thanks for the attention!*
*Questions?*