# Data storage at CERN

Overview:
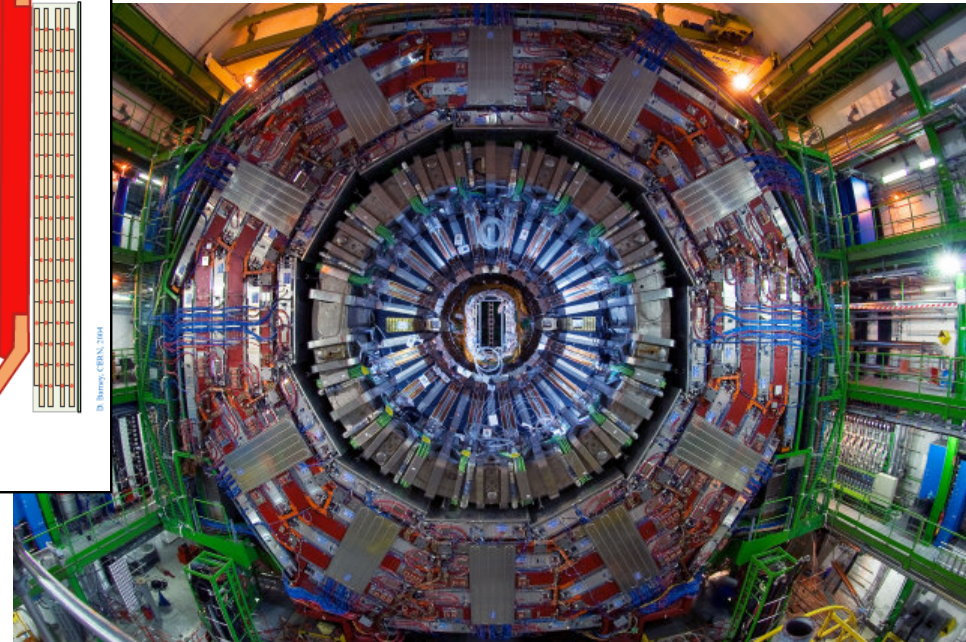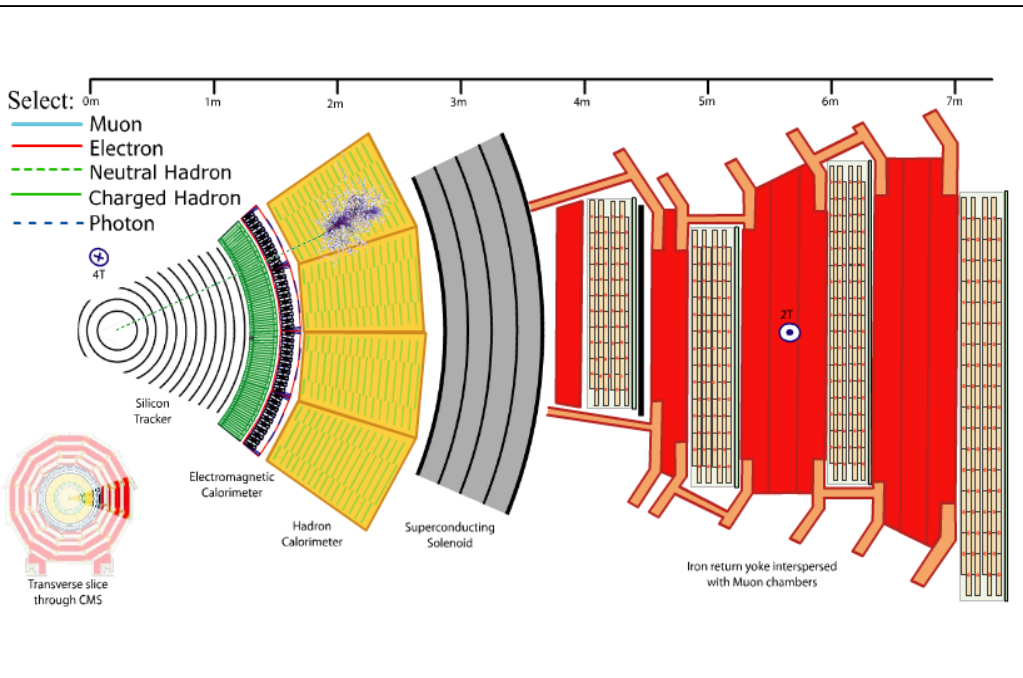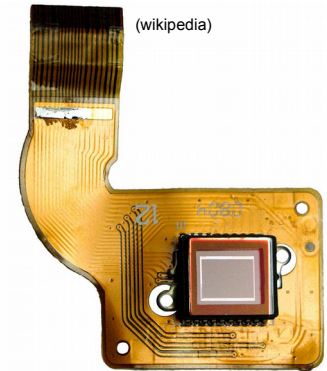
- Some CERN / HEP specifics
    - Where does the data come from, what happens to it
- General-purpose data storage @ CERN
- Outlook

# CERN vs Experiments

- CERN = (Conseil* Européen* pour la Recherche Nucléaire*) →
  European Organization for Nuclear Research
  - Est.1954, international treaty, 21 states
  - Provides lab facilities: water, electricity, cooling, offices, network, computing, various flavours of particle beams,..
  - ~2300 staff – clerical, engineers, firemen, ..

- Experiments: international scientific collaborations, own funding
  - "HEP": high-energy physics
  - Build & install detectors
  - Use lab facilities ("MoU")
    - Computing: yearly review
  - Generate & use & manage data
  - 10...3000 physicists each
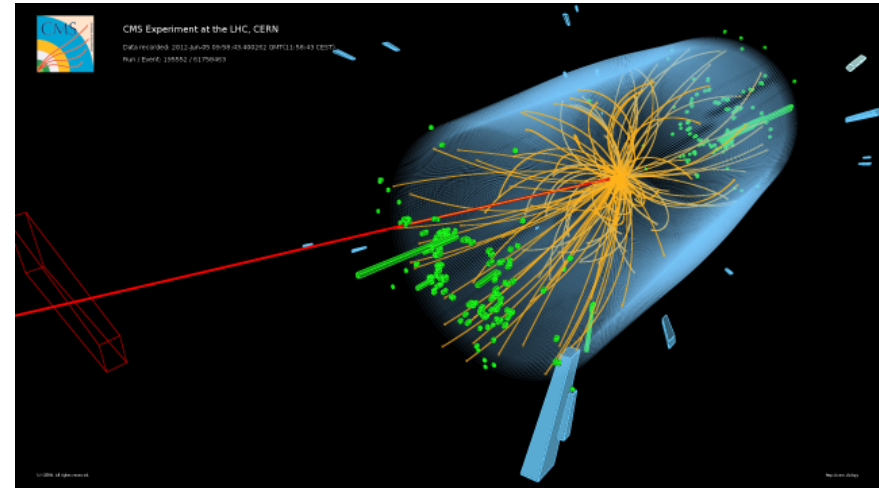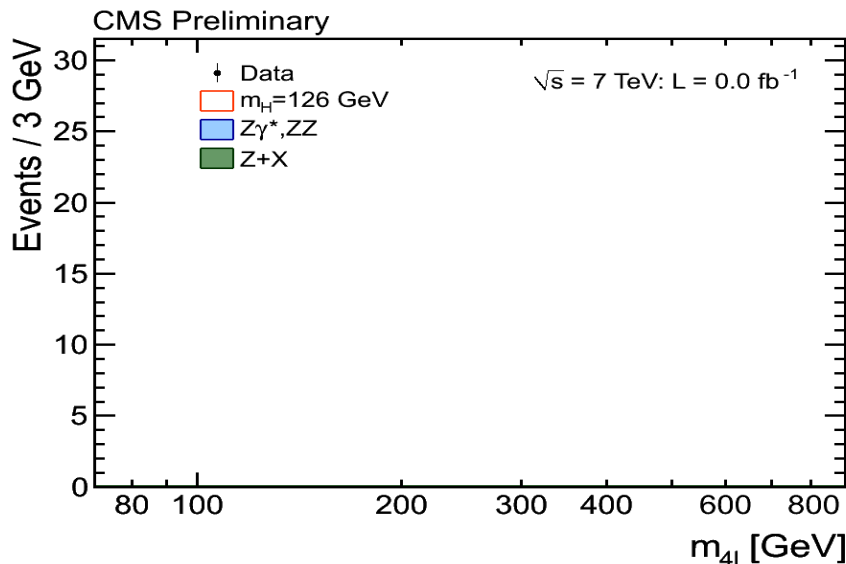- WLCG: computing grid for LHC



*: wrong, nowadays

# "Data taking" in HEP

- ## Think "digital video camera"..
  - – Unwieldy & complicated & expensive

(wikipedia)



Select:
— Muon
— Electron
- - - Neutral Hadron
— Charged Hadron
- - - Photon

0m    1m    2m    3m    4m    5m    6m    7m

4T

2T

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

Iron return yoke interspersed with Muon chambers

Transverse slice through CMS

D. Barney, CERN, 2014
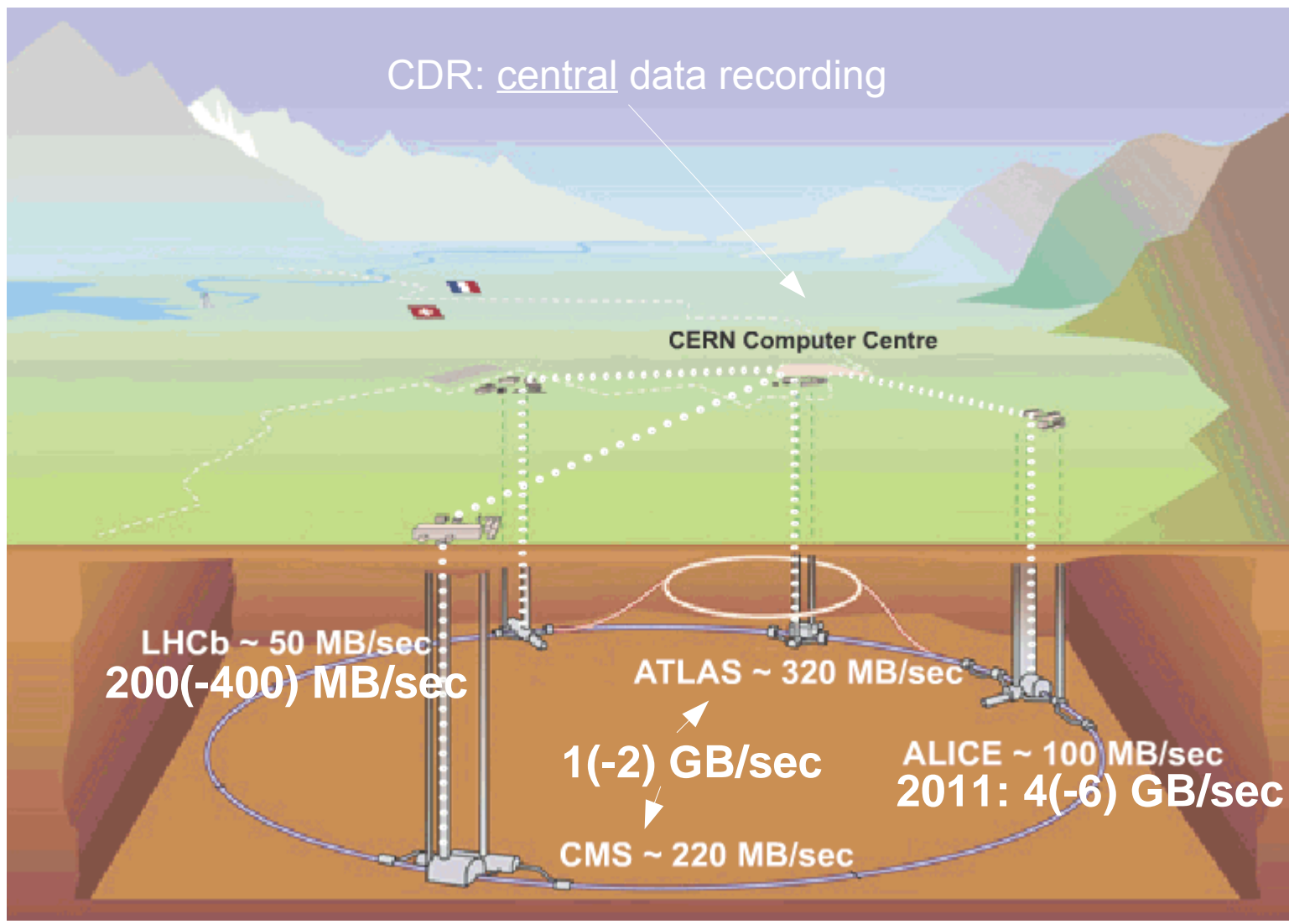
# "Data taking"

- .. but gives 4d "pictures"



- Significant postprocessing required: calibration, track reconstruction



- .. and subsequent analysis

- Result: good statistics on very rare events ≈

scientific papers.

Data &
Storage
Services

Tier 0 at CERN: Acquisition, First pass
reconstruction,  Storage & Distribution

CERN
IT
Department

CDR: central data recording

CERN Computer Centre

LHCb ~ 50 MB/sec
**200(-400) MB/sec**

ATLAS ~ 320 MB/sec

**1(-2) GB/sec**

ALICE ~ 100 MB/sec
**2011: 4(-6) GB/sec**

CMS ~ 220 MB/sec

# Some history of scale…

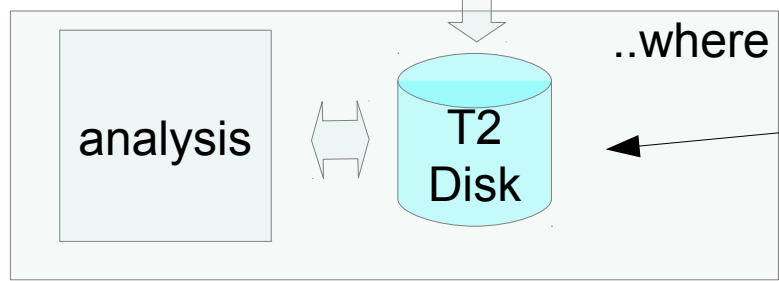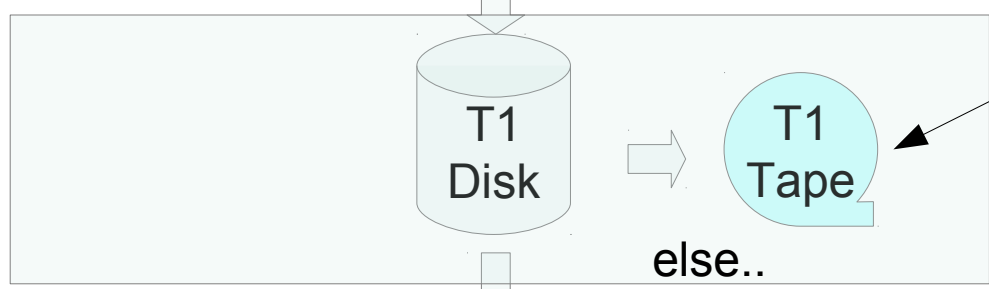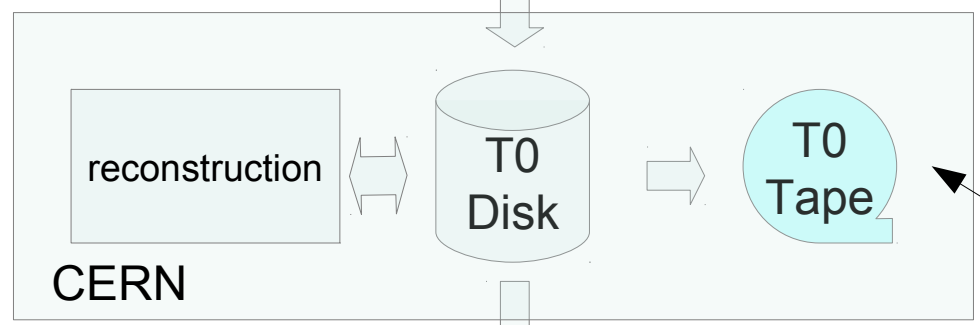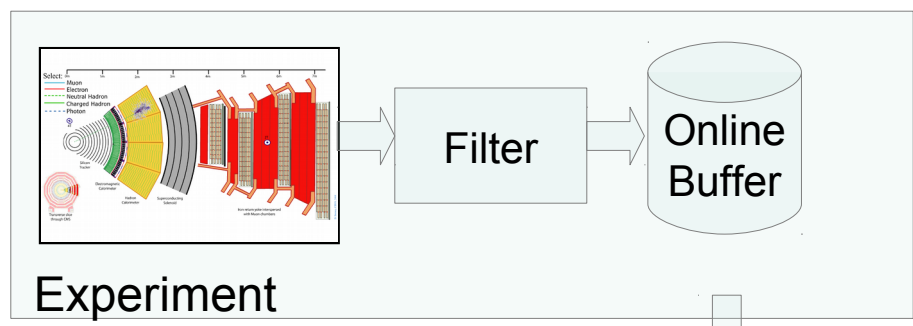| Date | Collaboration sizes | Data volume, archive technology |
|---|---|---|
| Late 1950's | 2-3 | Kilobits, notebooks |
| 1960's | 10-15 | kB, punchcards |
| 1970's | ~35 | MB, tape |
| 1980's | ~100 | GB, tape, disk |
| 1990's | 700-800 | TB, tape, disk |
| 2010's | ~3000 | PB, tape, disk |

For comparison:
1990's: Total LEP data set ~few TB
Would fit on 1 tape today
Today: 1 year of LHC data ~25 PB

Experiment

Filter

Online Buffer

reconstruction

T0 Disk

T0 Tape

CERN

T1 Disk

T1 Tape

else..

..where

analysis

T2 Disk

Custodial copy

Working copies

Experiment

merge

compress

reconstruction

fast analysis

T0 Disk

T0 Disk

T0 Tape

CERN

analysis

T1 Disk

T1 Tape

analysis

T2 Disk

Side note: CERN is small..

WLCG Disk Growth

Tier2
Tier1
CERN
15% Growth
2008-12 linear

# Physics data - CERN IT part

- Physics Storage systems in CERN-IT:
  - **CASTOR**: HSM
  - **EOS**: diskonly low-latency access, recent
- Both:
  - Homegrown
  - [Non|HEP]-standard protocols (**XrootD**, RFIO, SRM, gridftp)

# CASTOR HSM

CASTOR
CERN Advanced STORage manager

92PB    316M    350k
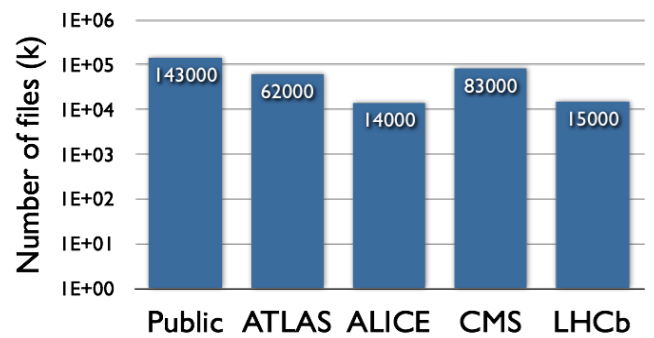
- Born in 1999
- Common Namespace
- Main Role: **data recording**,
    Tier-1 data export, production activities
- Mainly **tape-backed data**
- Focus on tape performance
    (latency can be high)
- **Database** centric
- Not optimized for concurrent access:
  - Currently Raid-1 configuration
- **Aimed** at DAQ activities: limited transfer slots = QoS
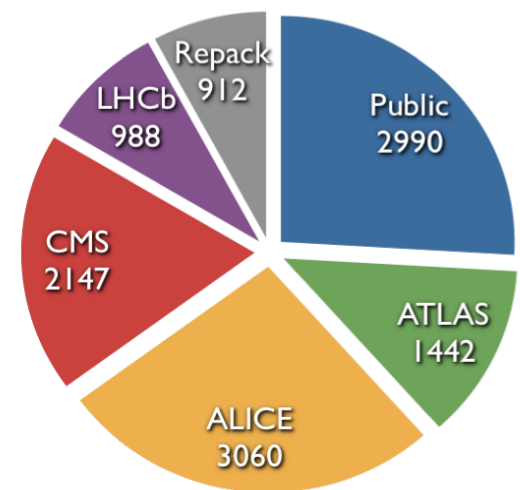- No (real) quotas**..**

# CASTOR: current setup

- 7 instances:
  - LHCs: ATLAS, ALICE, CMS and LHCb
  - PUBLIC: users, non-LHC experiments and specific pool for DAQ activities for AMS, COMPASS, NA61, NA62
  - Repack (tape media migration and compacting) and PPS (pre-production)
- Totals:
  - 92PB, 316M files, ~650 diskservers
- Mature release cycle:
  - ~1 major release per year
- In production also at RAL and ASGC

**Number of files in CASTOR Namespace**

Number of files (k)

| | 143000 | 62000 | 14000 | 83000 | 15000 |

Public  ATLAS  ALICE  CMS  LHCb

*Disk Space Installed (TB)*

Repack 912
LHCb 988
CMS 2147
Public 2990
ATLAS 1442
ALICE 3060

# CASTOR historical data



Experiments Production Data in CASTOR

LHC stop,
EOS

LHC first
beam

EOS
migration

Generated on Mar 25, 2014



Tape "cold" data verification

Repacking tapes is major activity



Tape data written, 2011-2018

25PB   158M      73k

- Born in 2010
- **in-memory namespace** – split per instance
- Main role: end-user analysis
- **Disk-only** storage
- Focus on **low latency**
- Optimized for **concurrency**
- **Multi-replica** on different diskservers
- No limit on transfer slots – throttle via overload
- **Quota** system: users&groups (for volume and files)
- **Strong** authentication: krb5, X509
- Diskserver: **JBOD** configuration

# EOS: current setup

- 6 instances:
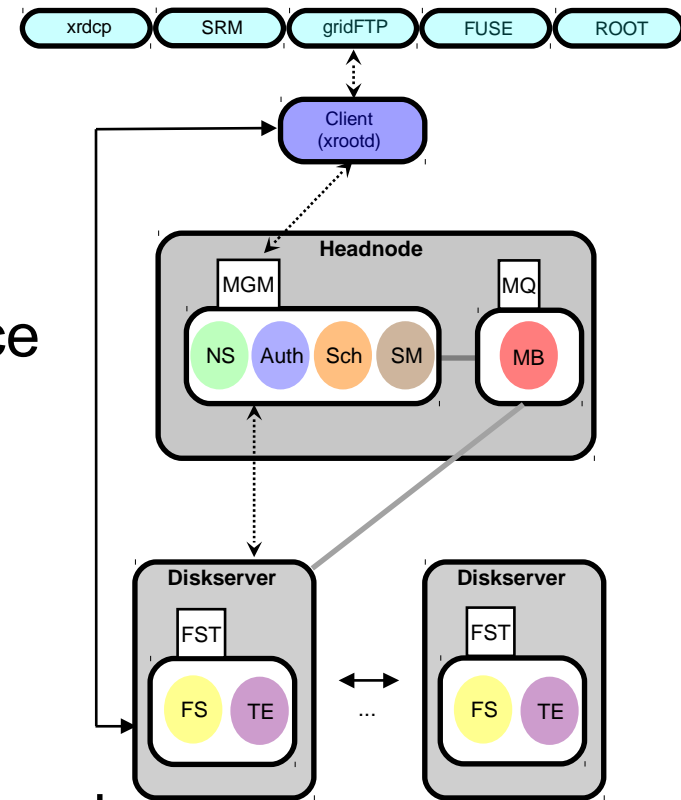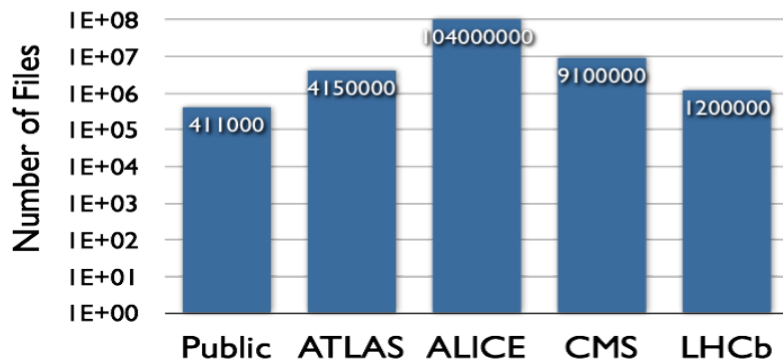  - LHCs: ATLAS, ALICE, CMS and LHCb
  - PUBLIC: recently deployed - AMS and COMPASS experiments
- Totals:
  - 20PB, 158M files, ~1100 diskservers
- Release lifecycle driven by functionality
  - ~2 major release per year, constant updates
- Used at Fermilab (Tier-3 functionality)

Number of files on EOS Namespace

Public 411000, ATLAS 4150000, ALICE 104000000, CMS 9100000, LHCb 1200000

Disk Space Installed (PB)

LHCb 3.3, Public 1.0, ATLAS 9.4, CMS 6.5, ALICE 5.4

Data &
Storage
Services

# EOS  deployment

CERN**IT**
Department

EOS Disk Space deployment

**Data & Storage Services**

(Data storage – everywhere at CERN...)

- Experiments:
    - ex. DAQ 'disk buffers' – up to several days of data taking
    - Often: prototype solution, but trouble with running long-term

- Department & group-scale independent solutions
    - CERN IT services ought to cover these..
    - .. but not always do.
    - ("NIH"-syndrome?)

- (Structured data / databases – not considered here)

Here: looking at CERN IT(-DSS) services.

# General-Purpose storages

- **General-purpose shared filesystems**

  Home directories = Untrusted clients = strong authentication

  – **AFS – Linux/Mac -** (see CERN site report), 950TB, 3G files
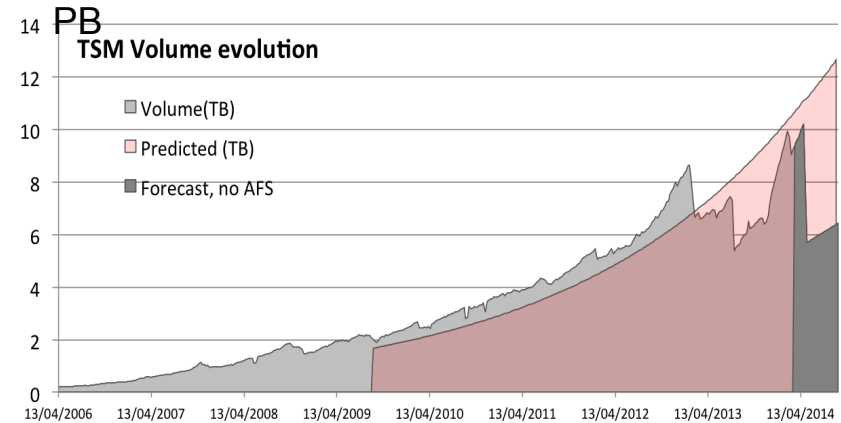
  – DFS - Windows-only

  – Future: (NFSv4), (FUSE-mounted EOS), (local FS+OwnCloud)

- **"cluster filesystems"**

  Typically: weak authentication

  – Not used for "computing" at CERN

    • general problem: "open" network + (too) many machines

    • CERN computing is "embarrassingly parallel"

    • CERN computing is worldwide

  – (experiments have own networks)

  – **NetApp** Filers: NFSv3 (190TB, 210M files)

    • Re-export: **CVMFS**, higher-level services (TWiki, VCs, ..) , DB

    • "standard", "will not void warranties"

- # Archive/tape
  - ## **TSM** (9PB, 2.1G files)
    - (See CERN AFS Backup talk)
    - Most machines are *not* backed up
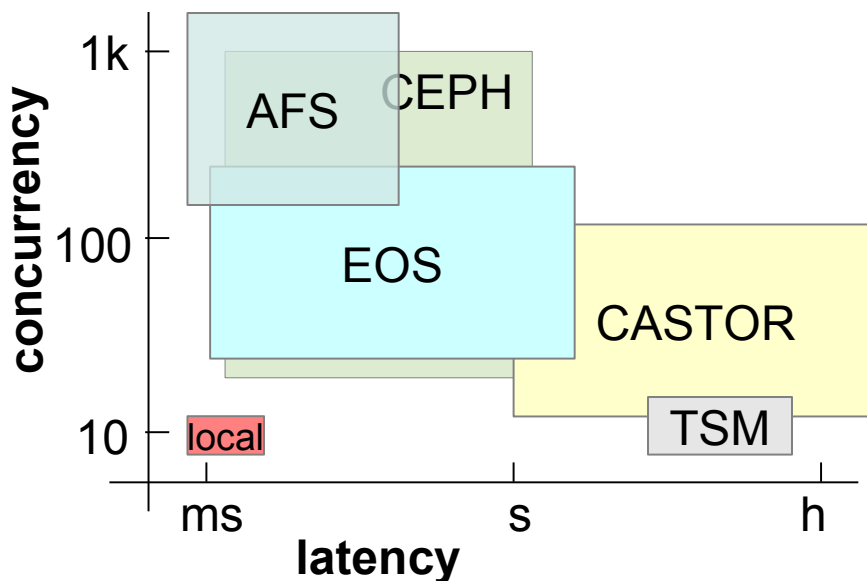  - ## (CASTOR)
  - (EOS reed-solomon)

- # Block storage: **CEPH**

- R&D: Hadoop/HDFS, Huawei S3, SWIFT, ..

- # lots of auxiliary storage services
  - protocol gateways (SRM, gridftp, http/webdav,..)
  - File transfer engines (FTS), **OwnCloud**
  - Bookkeeping, accounting, monitoring,...



TSM Volume evolution
- Volume(TB)
- Predicted (TB)
- Forecast, no AFS

# Looks like a zoo?

- Yes. But: different Dimensions:
  - Usage
  - Size/#Files
  - I/O pattern
  - Layering
  - Lifecycle (eval/prototype/production/legacy)

- "tool for the job"-approach
  - (and a bit of history. And Co-Evolution at work..)
  - Historically, lots of "use cases" started on AFS...
    - .. and once "big enough", moved somewhere else. Mostly.

CERN IT Department

(Orders-of-magnitude only. No sweat.)



- Map new use cases to "our" toolbox
  - Size, performance
  - Acceptable limitations
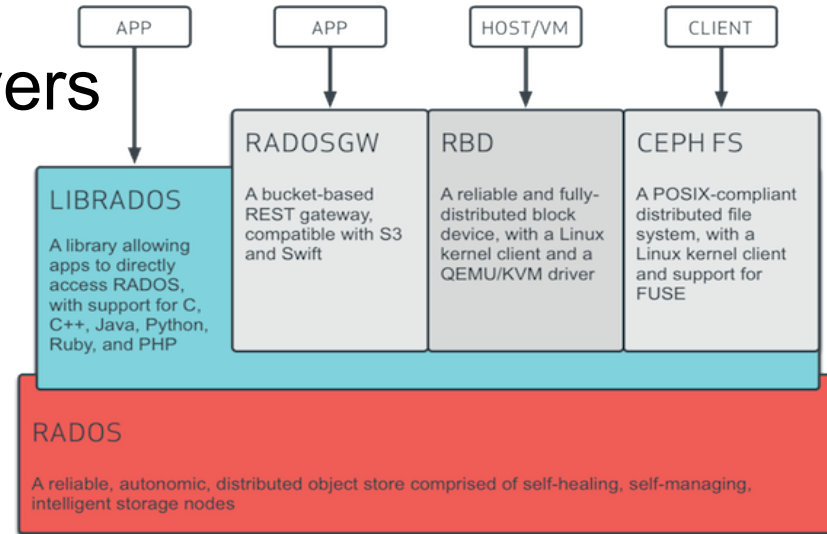
# Our building blocks

- *batchusserverus communalus*
  - 32..64GB RAM; 2..3 HD; Intel or AMD; ScientificLinuxCERN
  - Used for (replicated) headnodes, metadata servers

  + SATA tray(s) + 10GbE →

- *discusserverus vulgaris*
  - 24x 3TB + 3x2TB, 1x 10GbE,
  - used by EOS, CASTOR, AFS, CEPH, TSM, HDFS
  - $$$/TBmonth varies by replication factor (+manpower)

- Build up higher-value services from more basic services:
  - OwnCloud = EOS+FUSE+HTTP or NetApp+HTTP
  - CVMFS = NetApp+HTTP
  - AFS backups = (TSM or CASTOR) + scripts

- Redo when needed

# All is well?

- Problems / inefficiencies:
  - Manpower issues: training, split
    - But: No formal standby for any data service)
  - Overprovisioning, per-service safety margins
    - But: tend to split too-big services anyway into "instances"
  - "One size fits all" rarely does
    - c.f AFS servers are half-empty – not enough cold data
    - idle CPU capacity on diskservers
    - Spare disks on other machines (CC has 107PB raw disk)
  - Allocation & procurement cycle is sloow
    - Formal tendering, anti-corruption safeguards, national interests ..

Data &
Storage
Services

CERN IT Department

- Nice 'blocks' -fit at several layers
- Have 3PB 'prototype'
  - Used for OpenStack
    (VM images+volumes)



- Ogled by CASTOR, EOS, (AFS), (NFS)
  - To replace current talk-to-disk-and-handle-errors layer

- (see D. Van der Steer "Ceph – storage for the cloud", http://indico.cern.ch/event/300076/ )

- 2nd computer center in Wigner Institute/Budapest
  - "LAN" access now means 23ms..
- Currently LHC is not running (LS1)
  - Restart for "run 2" early 2015 – expect higher data rates (~2x) and different data flows
  - Other experiments will start earlier
  - Run 3: 75GB/s from ALICE?
- RAID works less and less (disk size vs reconstruction speed) : forced to RAIN

# Summary

- Some CERN data storage use cases are "special", some are not

- Toolbox / building block approach

- Common HW reduces cost
  - But manpower / know-how is an issue

- Nothing cast in stone
  - but some upheavals take a long time..