



# IT-ST: The Storage group

Alberto Pace on behalf of the IT-ST group




# The Storage group



# Mandate

- Ensure a coherent development and operation of storage services at CERN for all aspects of (physics) data
- Three goals
  - Storage
  - Distribution
  - Preservation

# Three sections

- Analytics & Developments (IT-ST-AD)
  - Design and develop central storage services and their evolution. Operate and support the Analytics and FTS services
  - Leader: Dirk Duellmann 
- File & Disk Operations (IT-ST-FDO)
  - Operate and support the storage and file system services for physics (AFS, CEPH, CERNBOX, CASTOR and EOS)
  - Leader Massimo Lamanna 
- Tape, Archives & Backups (IT-ST-TAB)
  - Design, operate and support the archive and backup services (includes robotics, drive and media, infrastructure for backup).
  - Leader: German Cancio Melia 

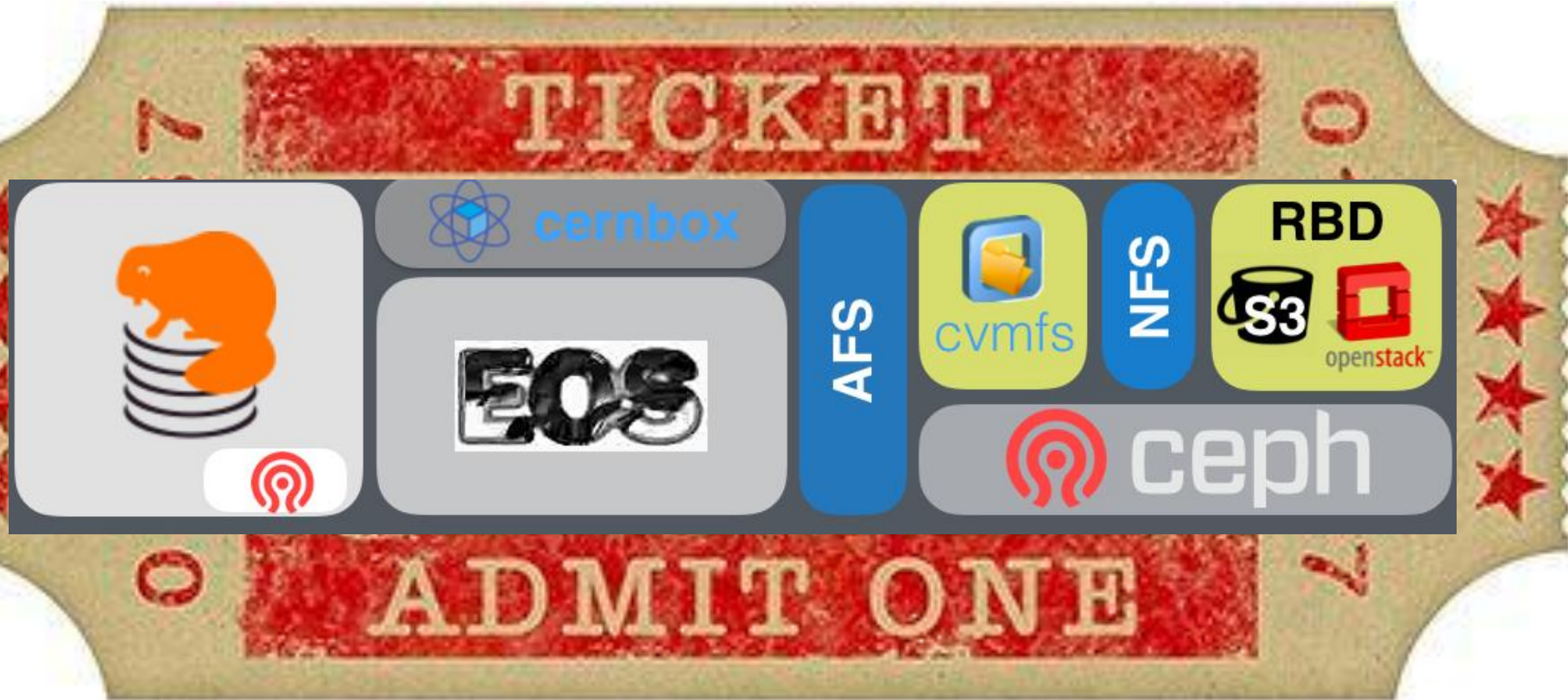
# IT-ST-FDO

Presented by Xavier Espinal

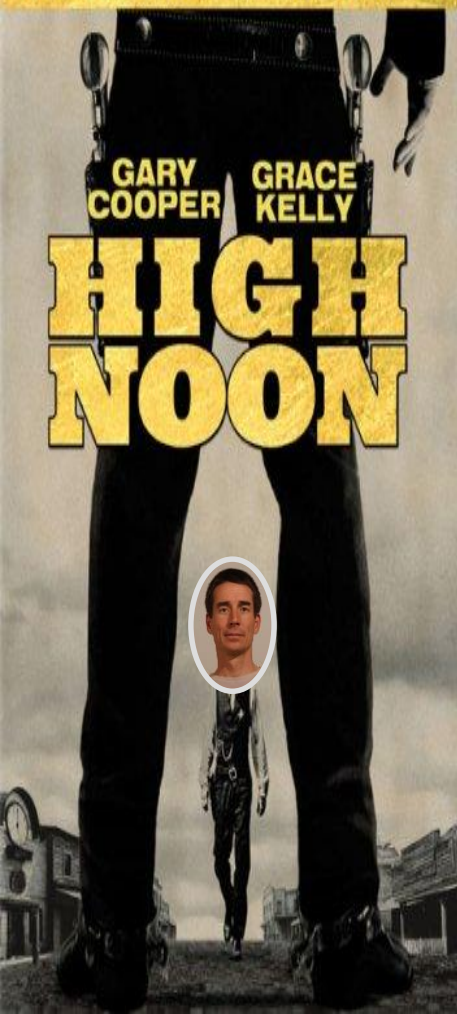
# Storage Services for Physics and more



# Storage Services for Physics and more IT-ST-FDO







# AFS (phase-out)

- Difficult task, but achievable
  - 520 TB (+16% this year vs 34% last year), 3.5 billions files (+16% vs +20% last year)
- We have solutions for nearly all current AFS usages
  - Successfully migrating **web/projects/personal files** to **EOS** and **software** to **/cvmfs**
- No hard line to “unplug” (Run3?)
- Starting a discussion with all relevant service managers concerning **home directories**:
  - Ixplus, Ixbatch, linux/win/Mac support, IT-CDA, IT-CM,...

# Filer service



- The current Filer service solves the native NFS (posix) needs
- 80TB allocated and +30 projects:
  - Twiki, Puppet, MIC, LSG, gridCE, myproxy, gitlab, boinc, openshift, Indico, Inspire, etc.
- Currently running on a single NFS server over CEPH-RBD does not provide High Availability. We are evaluating possibilities to provide HA:
  - pNFS cluster
  - CEPH-FS

## The Odd Couple



[Wiki » Ceph Advisory Board »](#)

CAB 2016-01-12 ¶

Present:

- Patrick McGarry, Red Hat (chair)
- Sage Weil, Red Hat (tech committee chair)
- Dan van der Ster, CERN (Academic Liason)



**Beesly** (5 PB + 433 TB, v0.94.9):  
Cinder (block storage volumes)  
Glance (images repo)  
Rados GW (object storage interface: S3, Swift)



**Dwight** (0.5 PB, v10.2.3):  
Preprod cluster for development (client side)  
Testing, upgrades and crazy ideas



**Erin** (4.2 PB, v10.2.3):  
New cluster for CASTOR: disk buffer/cache in front of tape drives



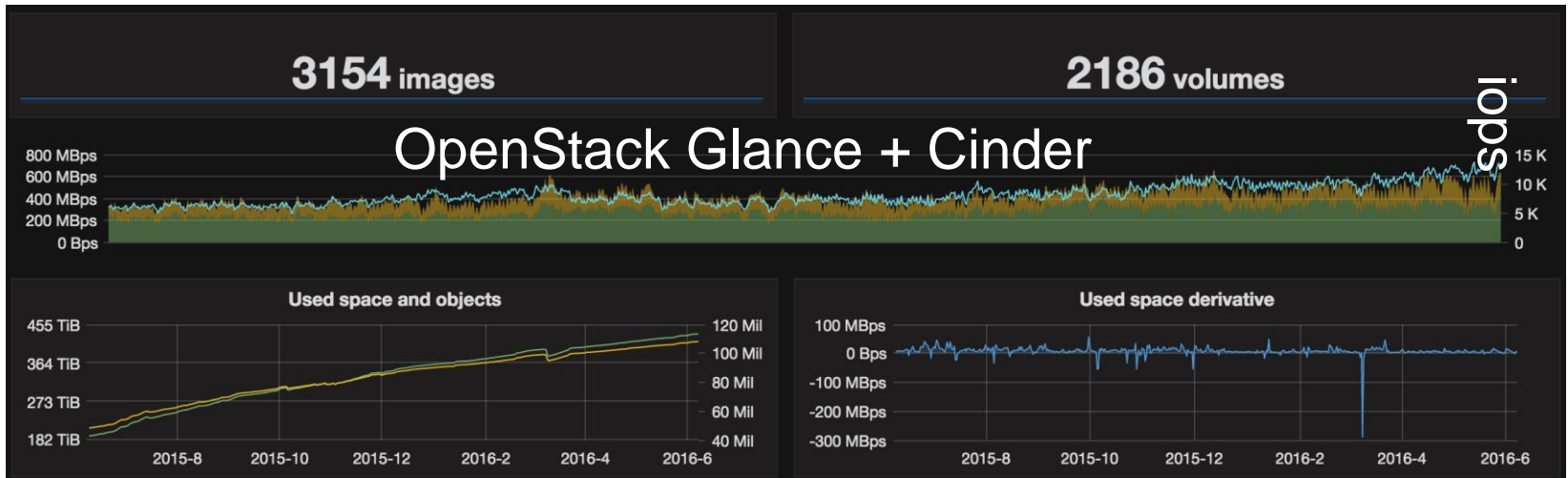
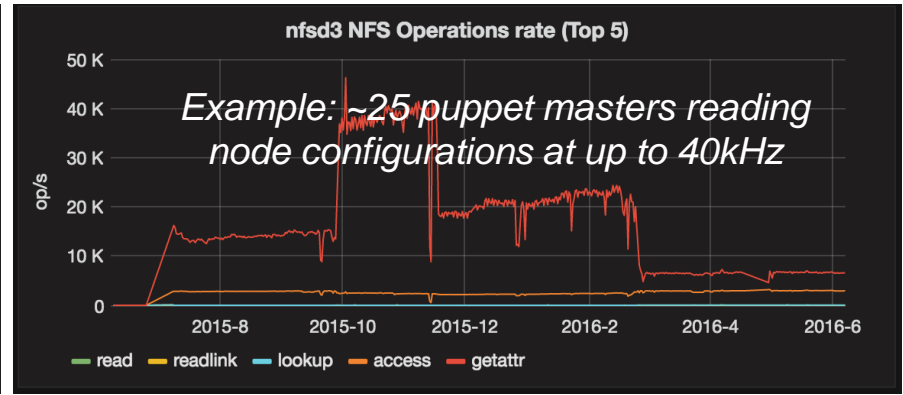
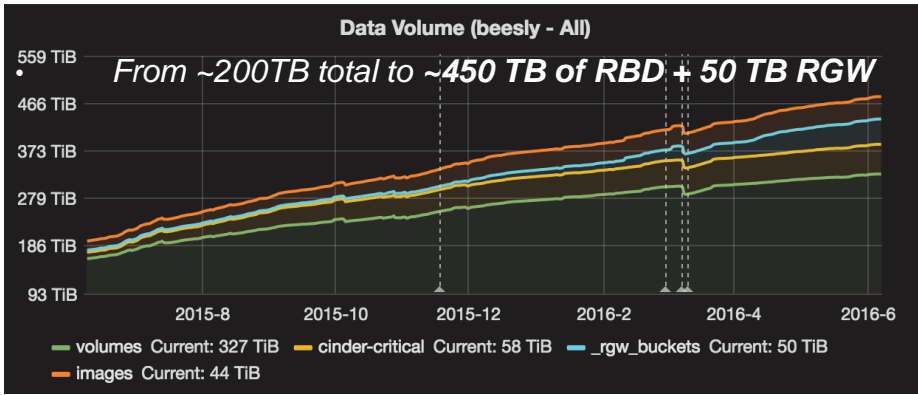
**Flax** (0.4PB, v.10.2.3 - early stage):  
Ceph-FS HPC cluster for QCD studies



**Gabe** (1PB, v.10.2.3):  
New S3 Object Store IPV6 only



**Bigbang** (~30 PB, master):  
Playground for short term scale tests  
Usually when we receives new hardware



# LES 7 SAMOURAÏS +5

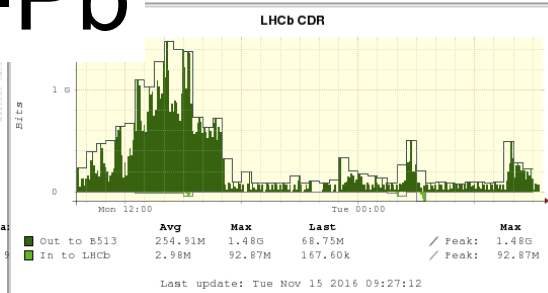
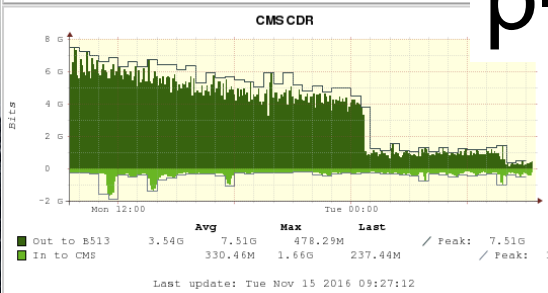
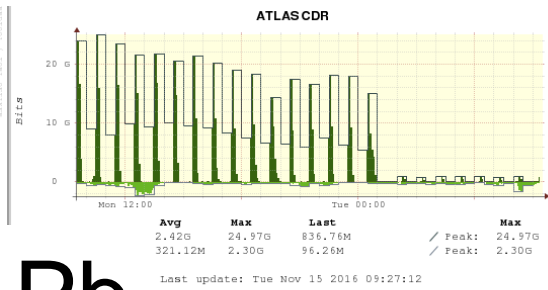
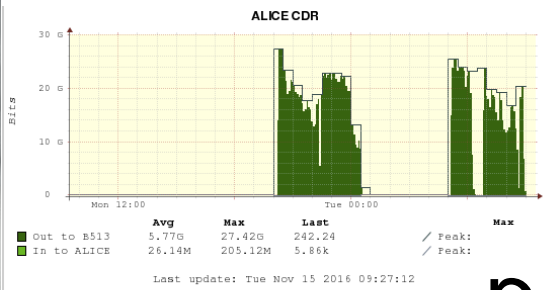
LE CHEF-D'ŒUVRE DE  
AKIRA KUROSAWA



# Large Scale Storage for LHC data

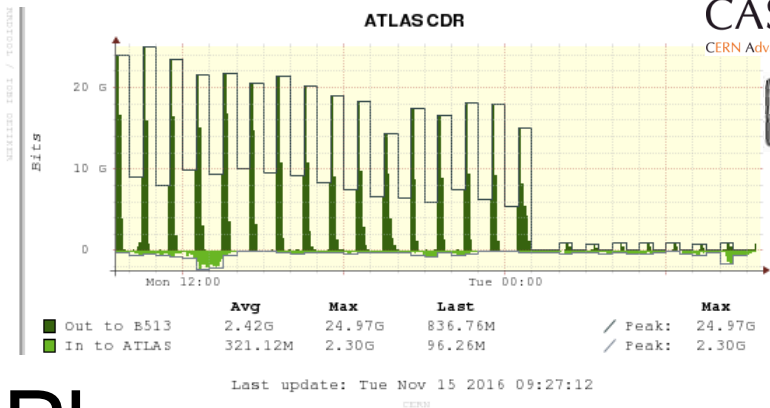
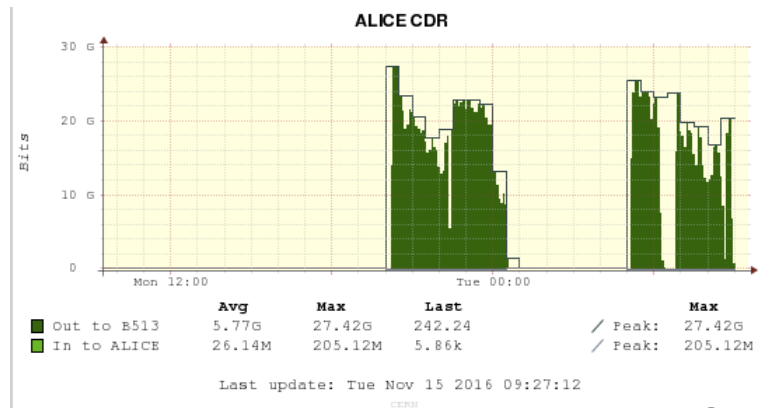
CASTOR   
CERN Advanced STORAGE manager

EOS 

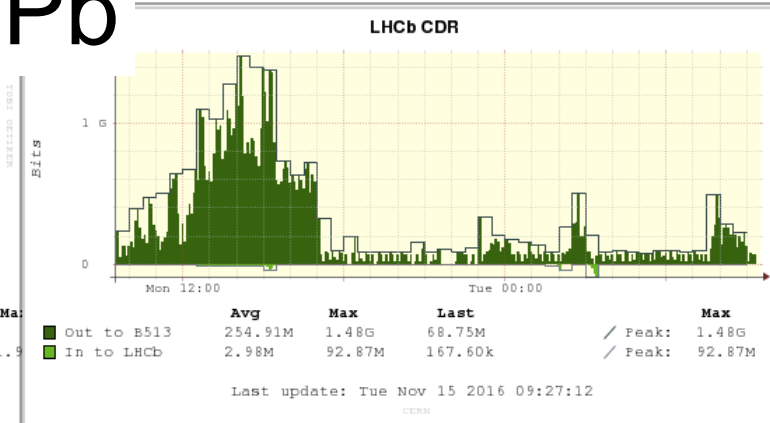
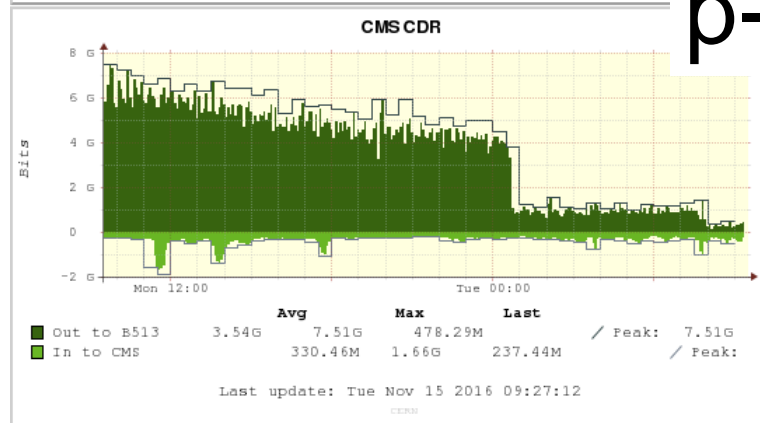


p-Pb

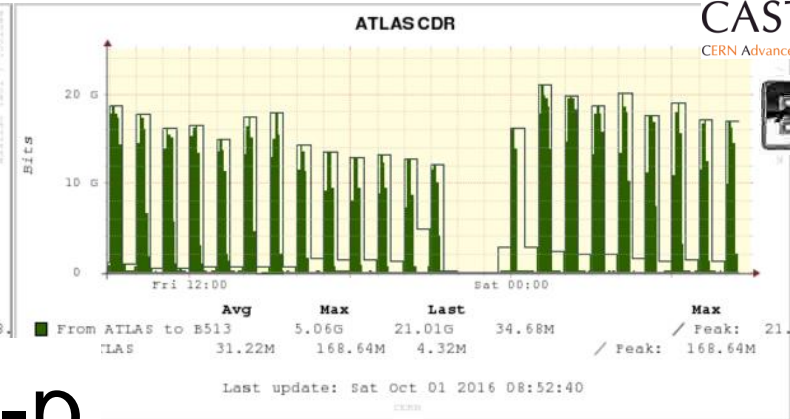
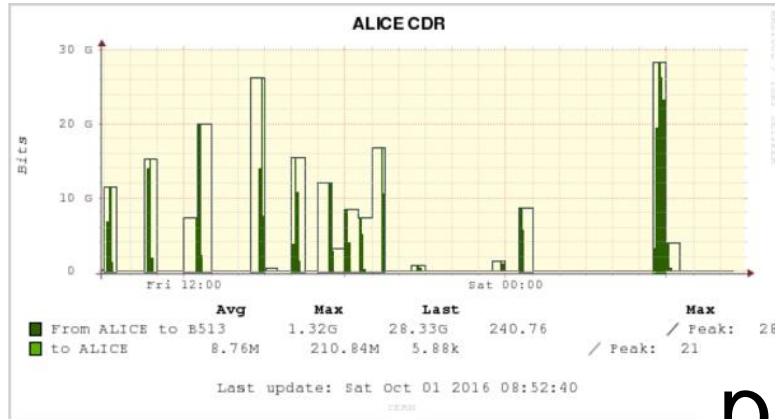
# Large Scale Storage for LHC data



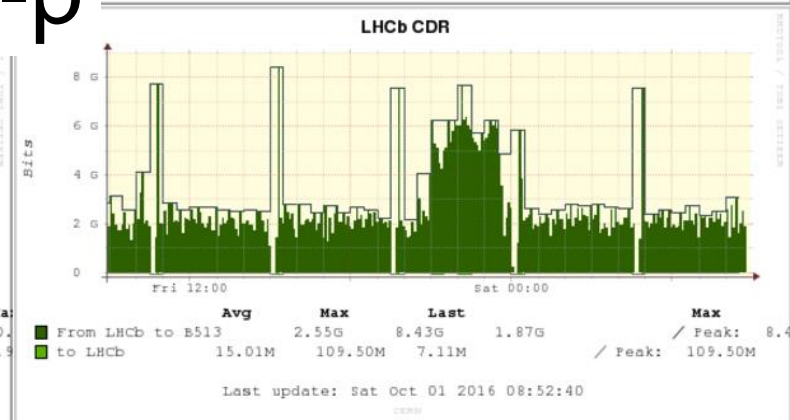
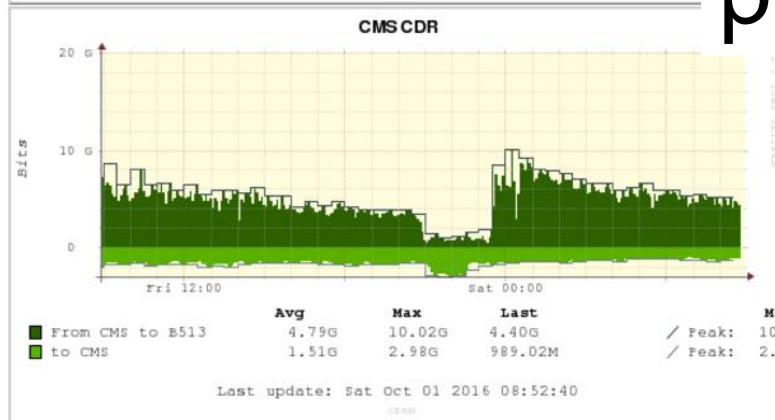
p-Pb



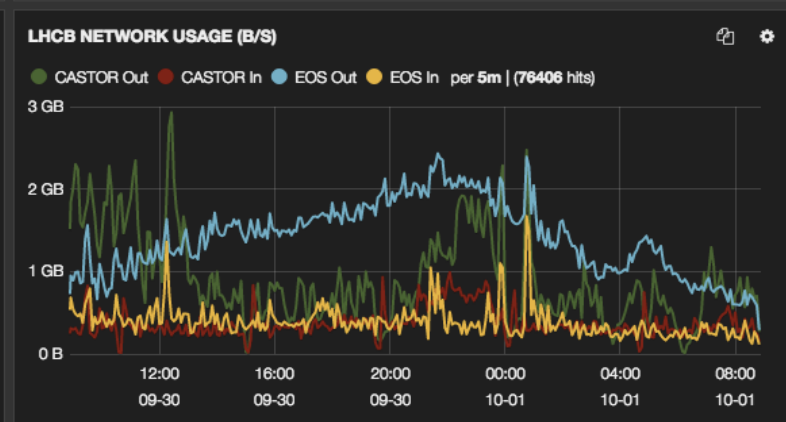
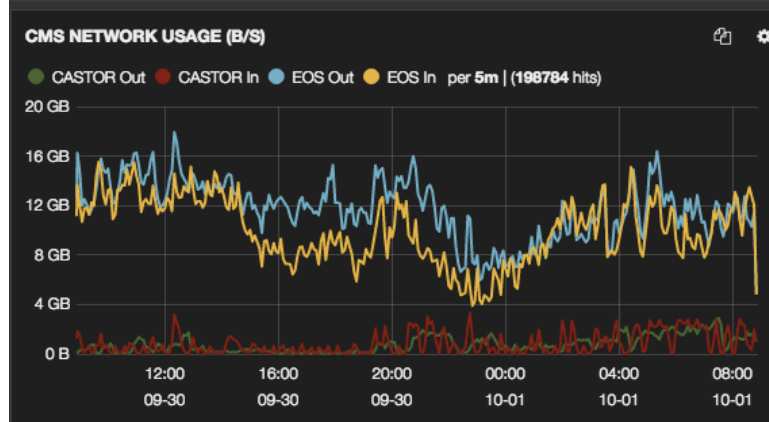
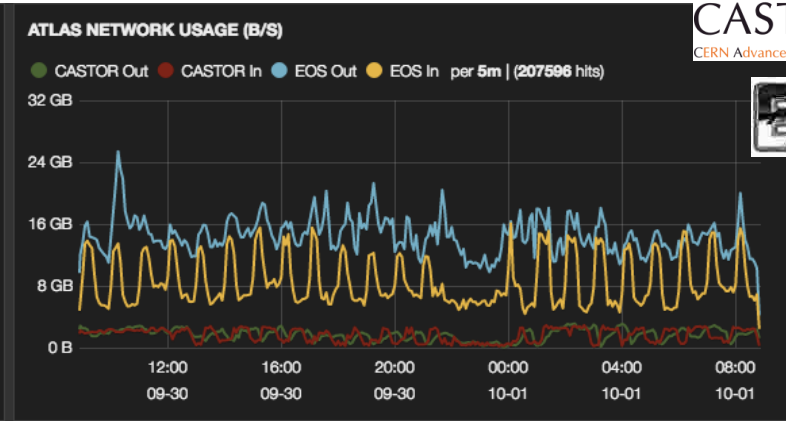
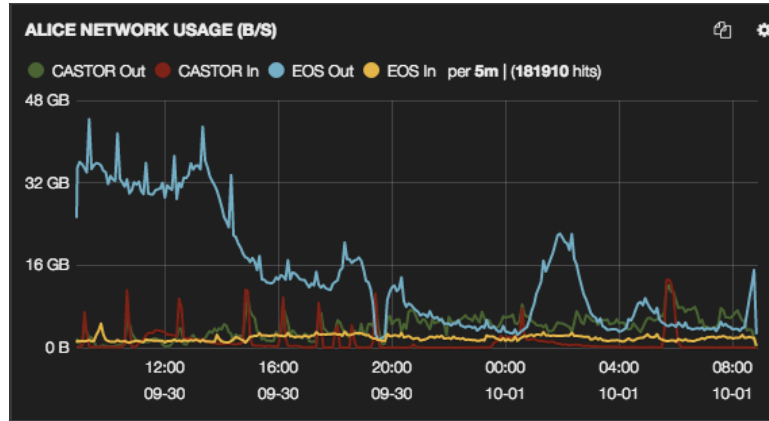
# Large Scale Storage for LHC data



p-p



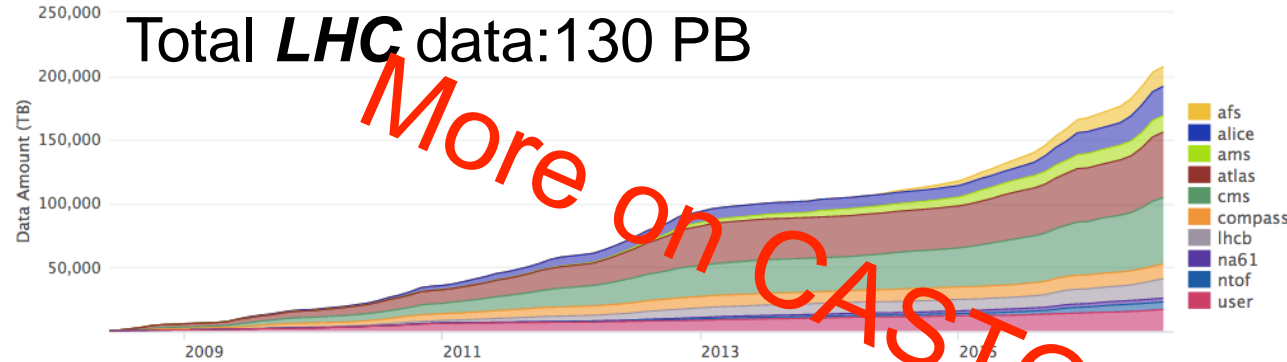
# Large Scale Storage for LHC data



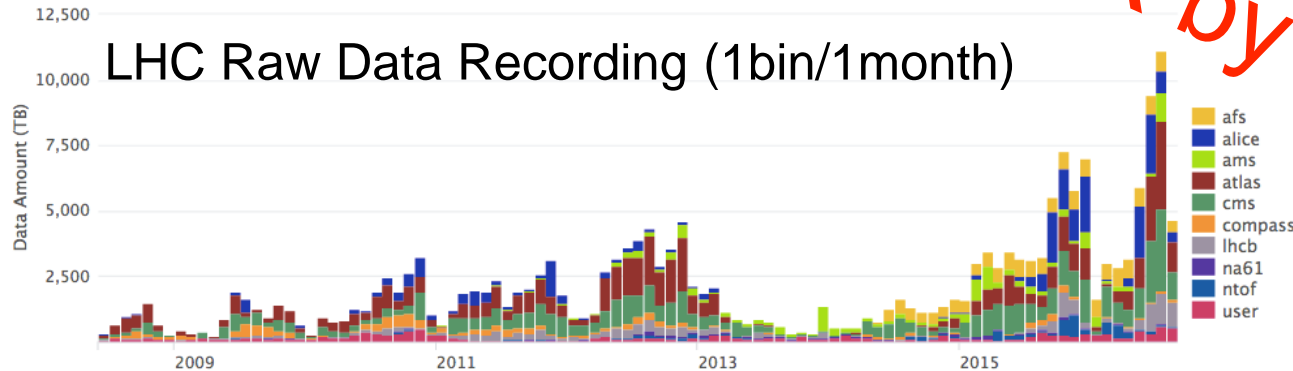


# Large Scale Storage for LHC data

Accumulative Transferred Data Amount per Virtual Organization for WRITE Requests



Transferred Data Amount per Virtual Organization for WRITE Requests



Biggest physics-repo  
worldwide :

**180PB** and **500M**  
files

Evolved towards a **Tape**  
**oriented** system during LS1


**CEPH** backend in  
production (Alice, Repack)

**Cold** by definition: high  
throughput, high latency

Future evolution for a  
pluggable tape backend,  
see CTA (Julien)

# Large Scale Storage for LHC data



+1200  nodes  
+45000  disks

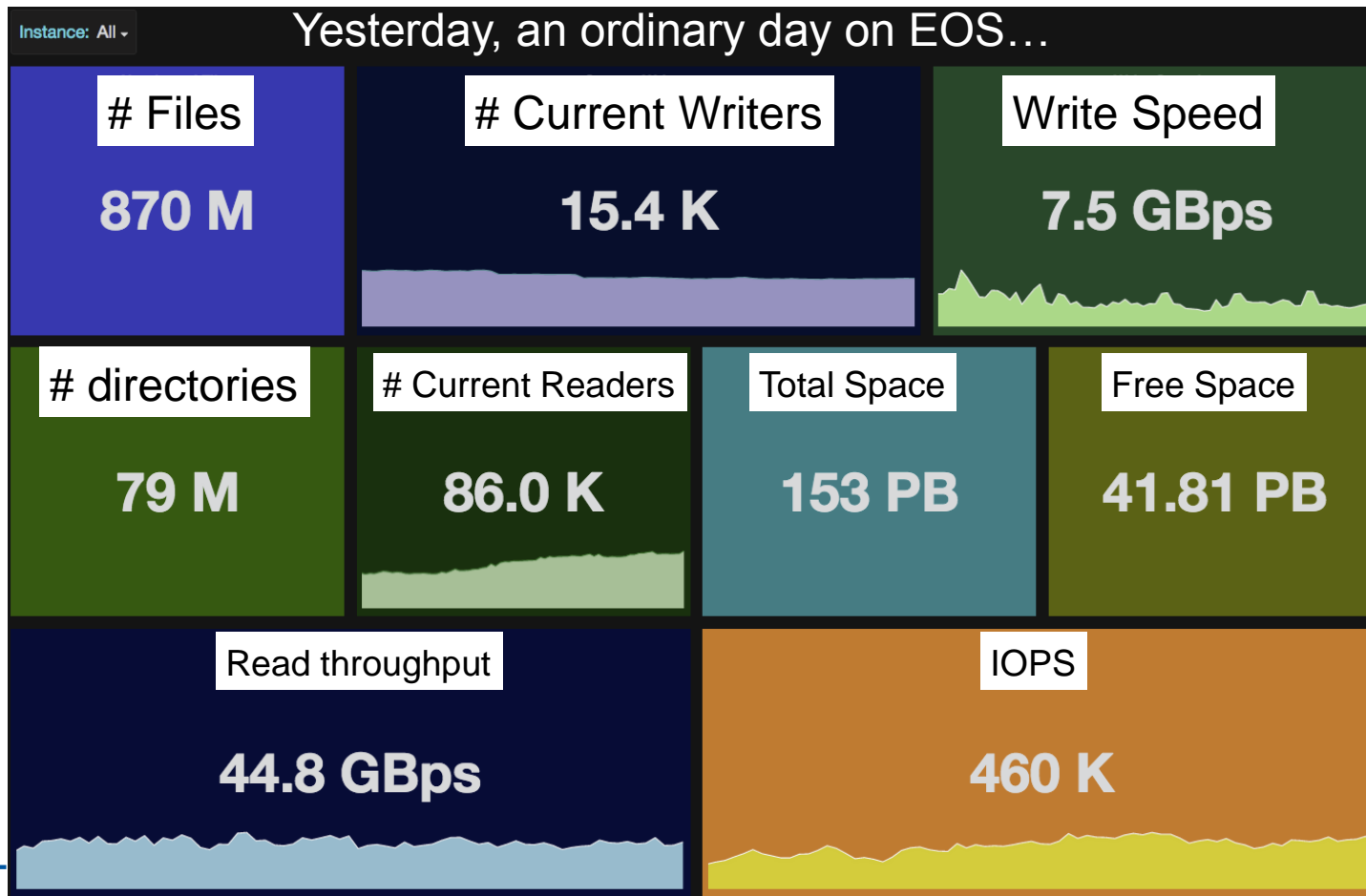
**EB**  
era

Easily scalable (#disk #servers)  
Performant and manageable  
LHC main storage platform

850M  files  
150PB  volume usage

Concurrency peaks observed at  
+25000  $W_{rites}$   
+100000  $R_{eads}$

Throughput peaks at  
+80GB/s  
 $R_{eads}$





# Large Scale Storage for (not only) LHC data



**Users** 6200  
( $\Delta^{7d} + 60$ )

**#files** 125M  
( $\Delta^{7d} + 250K$ )

**#dirs** 15M

**Quota** 1TB/user

**Used Space** 310TB

**Deployed Space** 1.5PB

Collaborate

Share

Offline work

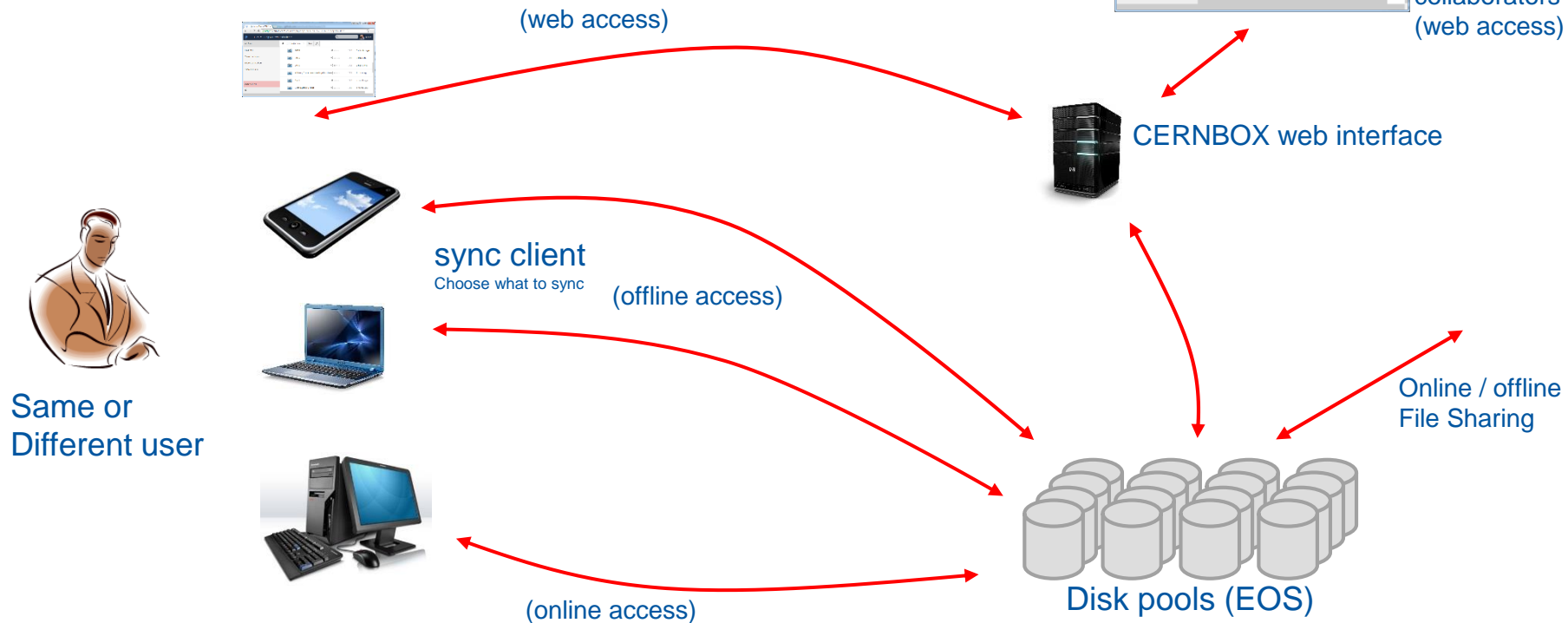
**Community storage**

Sync

SWAN



# Multiplatform / Multiuser





# Goals

Make data access easy  
Make Analysis simple  
Facilitate Science

# Goals

Make data access easy

Make Analysis simple

Facilitate Science

## My Laptop

Small scale analysis

Test jobs

## batch/interactive services

Large scale experiment processing

User extensive analysis

### Mounts

squids  
/cvmfs/athena

fuse  
/mycernbox

fuse  
/eos/atlas

## Data Access

Main experiment data repositories



# Goals

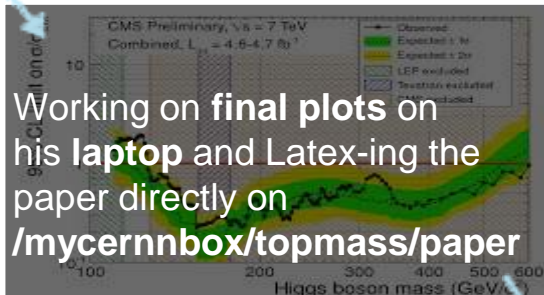
Make data access easy  
Make Analysis simple  
Facilitate Science

Physicist code: `topmass.kumac`  
on his laptop on `/mycernbox`  
and sync'd via `cernbox` client

Physicist identify an  
interesting dataset  
`/eos/atlas/phys-top`

Submit jobs to `lxbatch/wlbg` to  
process the data  
EOS Fuse: `/eos/atlas/phys-top`  
EOS Fuse: `/mycernbox/topmass.kumac`  
Experiment SW: `/cvmfs/athena`

Results (ntuples) aggregated on  
`/mycernbox/topmass` are  
sync'd on his laptop as the  
if desired  
jobs are being completed



Share on-the-fly:  
n-tuples  
Final plots  
Publication  
via `/mycernbox`

GLOBAL CONSERVATION LAWS AND MASSLESS PARTICLES\*  
by R. Utiyama, R. M. Wald, and S. Deser  
(Received 12 October 1964)

In all of the fairly numerous attempts to date to formulate a consistent field theory possible in the presence of broken symmetry, Goldstone's remarkable theorem<sup>1</sup> has played an important role. This theorem, briefly stated, asserts that if a Lagrangian is invariant under a continuous group of transformations which in no way reflects on the conserved operator  $Q_i$  such that

$$[Q_i, A_j(x)] = \sum_k t_{ijk} A_k(x)$$

and if it is possible consistently to take  $t_{ijk} \neq 0$ , then  $A_j(x)$  has a zero-mass particle in its spectrum. It has more recently been observed that the assumed Lorentz invariance essential to the proof<sup>2</sup> may allow one the hope of avoiding such massless particles through the introduction of a nonlocal conservation law usually

\* This work is based on a paper by R. Utiyama, R. M. Wald, and S. Deser, *Journal of Mathematical Physics*, **5**, 913 (1964).  
585

# Goals

Make data access easy

Make Analysis simple

Facilitate Science

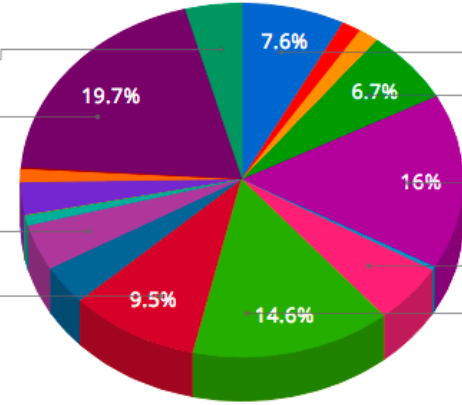
The image shows a screenshot of the Blue Waters website. The header is dark blue with the text "BLUE WATERS" in large white letters and "SUSTAINED PETASCALE COMPUTING" below it. To the right of the header are the "SIGN IN" link, the "NCSA" logo, and a logo for the University of Illinois. Below the header is a navigation bar with links: "YOUR BLUE WATERS", "ABOUT", "SCIENCE AT BLUE WATERS", "USING BLUE WATERS", "EDUCATION & TRAINING", "NEWS & EVENTS", and "HELP" with a search icon. The main content area is white and features a large heading: "Mapping Proton Quark Structure in Momentum and Coordinate Space using PetaByte Data-Sets from the COMPASS Experiment at CERN."

# Goals

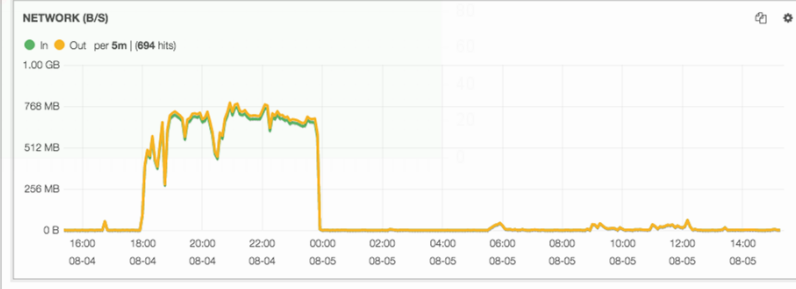
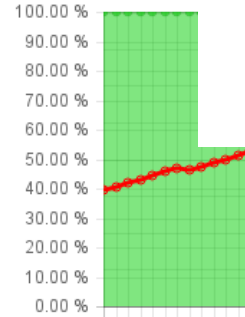
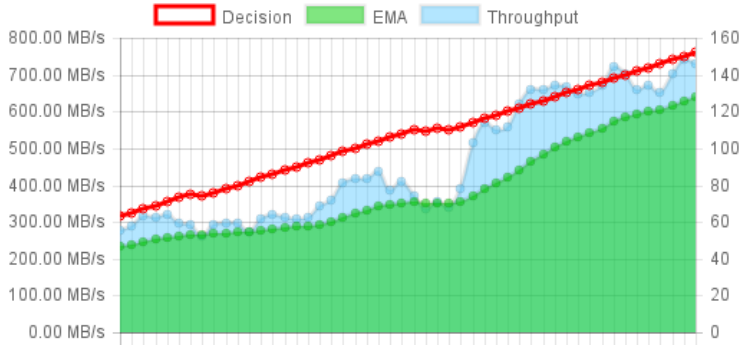
Make data access easy  
 Make Analysis simple  
 Facilitate Science

## CURRENT RUNNING JOBS BY SCIENCE AREA

Stellar  
 Astronomy and  
 Astrophysics  
 4.2%  
 Physics  
 19.7%  
 Fluid, Particulate,  
 and Hydraulic  
 Systems  
 4.2%  
 Earth Sciences  
 9.5%



Details for <srm://castorpublic.cern.ch> → <gsiftp://ie15.ncsa.ill>



First Previous 1 2 Next Last

Timestamp	Decision	Running	Queue	Success rate (last 1min)	Throughput	EMA
2016-08-05T13:57:24	154	152	1898	100.00%	735.688 MB/s	648.032 MB/s

Make data access easy  
Make Analysis simple  
Facilitate Science

# Goals

- Scale-out filesystem underneath the ownCloud app, using the eosd fuse interface for file IO

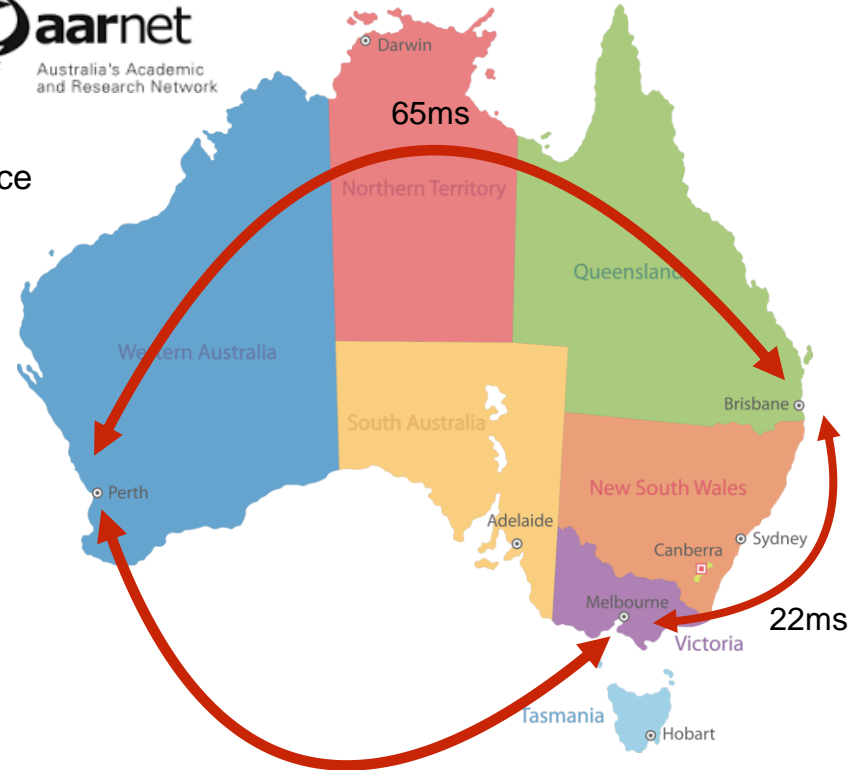


- Geo-distributed setup: Brisbane, Melbourne, Perth
  - ~1PB (scale to ~20PB next year)



Australian National University

- Australian National University, in Acton Canberra: mirror archives of both genome sequences and open or freely available software distributed among three sites



Make data access easy  
 Make Analysis simple  
 Facilitate Science

# Goals

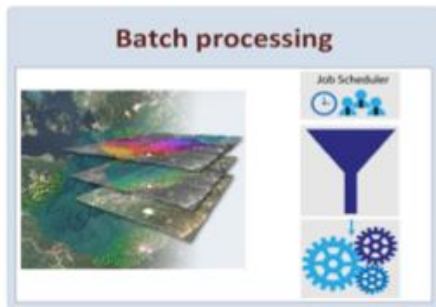
## Components of the JRC Earth Observation Data Processing Platform (JEODPP)



Joint Research Centre  
 @V.Vasilev,F.Eyraud  
 (JRC)

**Batch processing interface**  
 direct access interface through HTCondor scheduler for experiment/projects in cases

**Interactive Processing Interface**  
 web access interface for end users based on Jupyter, Leaflet and custom built image processing libraries



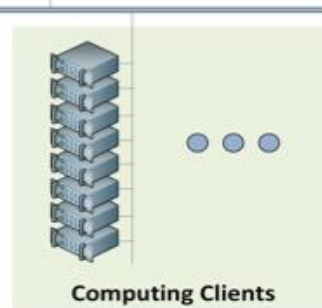
High performance network



**FILESYSTEM**  
 EOS  
 2 MGM servers  
 10 FST servers  
 240 6 TB disks

**TOTAL SPACE:**  
 1.44 PB

Current usage:  
 24.5M files  
 336TB



**PROCESSING**  
 9 2U processing servers  
 16 1U processing Servers

**TOTAL:**  
 600 cores  
 7+TB memory

# IT-ST-TAB

Presented by Julien Leduc





# IT-ST-TAB

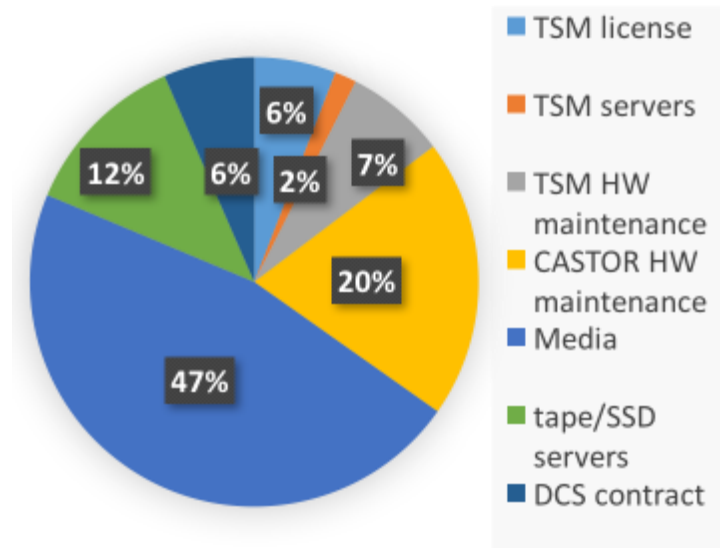
# Storage Services for Physics and more





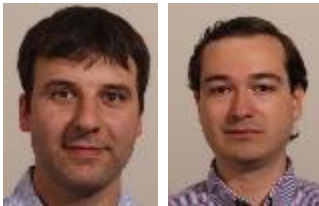
# Tape, Archives and Backups (TAB)

- Mandate
  - Design, operate and support the archive and backup services.
  - This includes the tape-based software backend for CASTOR, tape robotics, drive and media for physics, infrastructure for backup and restore of file servers and databases.
  - "Tape from A to Z"
- People
  - 7 staff + visitors + students
  - 1 external contract
- Budget
  - 1.7MCHF in 2016



# Backup Service

- Operate, maintain and monitor CERN Backup infrastructure:
  - Critical service backups
  - High reliability, sustained daily traffic: 70TB/day
- Infrastructure:
  - 2 tape libraries (9k tapes)
  - 55 tape drives
  - 20 TSM servers

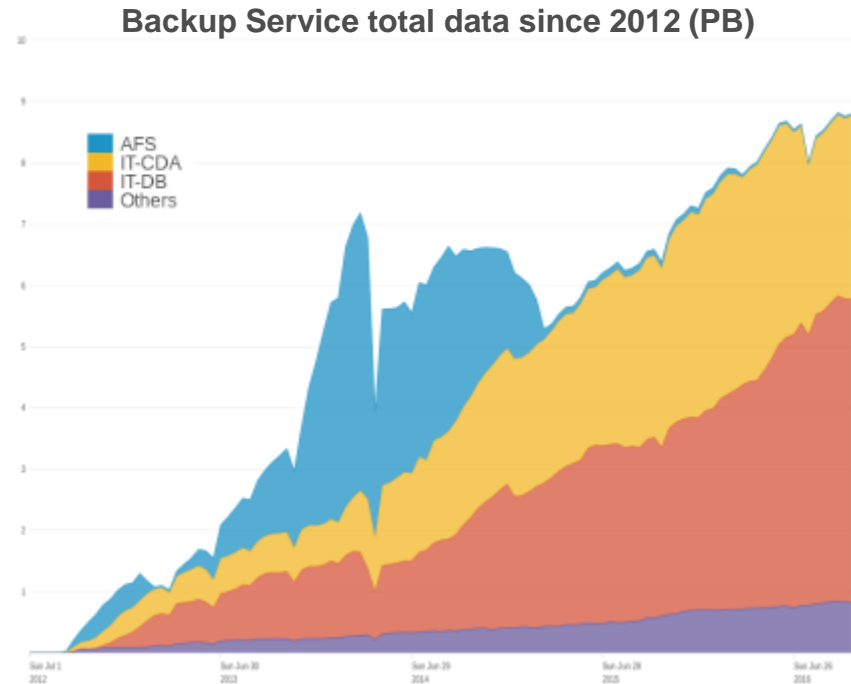


Backup Service total data since 2012 (PB)



# Backup Service

- TSM cost is proportional to backed up volume
  - Minimize backup volume in view of upcoming license retendering
- Ongoing efforts to limit annual growth rate: 24% (50% in the past):
  - Moved AFS backups to CASTOR
  - Worked with IT-CDA to remove redundant copies
  - Working with IT-DB to move Oracle backups out of TSM

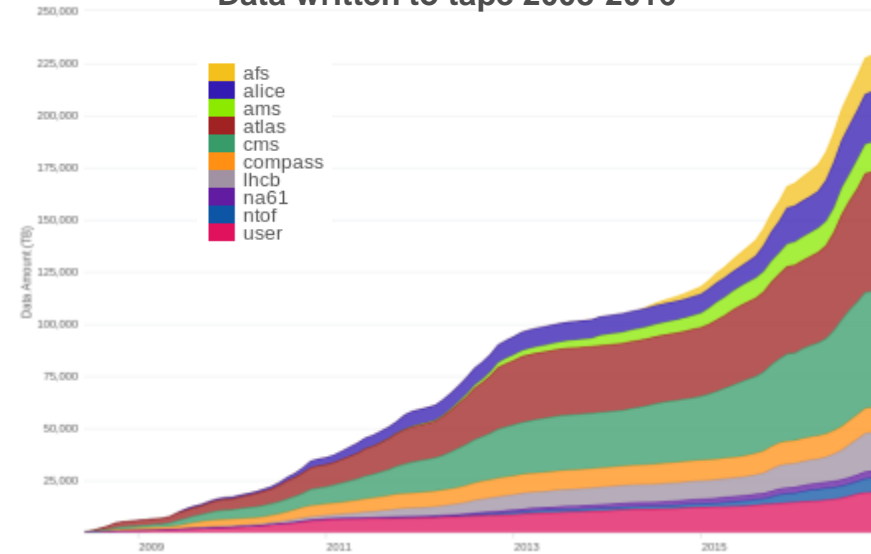


# Data Archiving

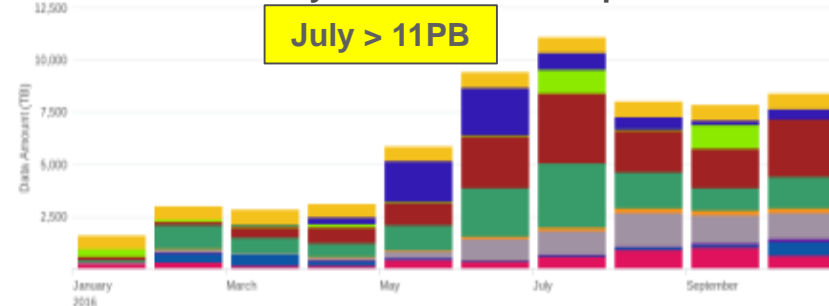
- Development, operation and monitoring of CERN's tape data archive:
  - Storage and long term (ad aeternum) archiving of CERN physics data
  - High throughput, high reliability, cost efficient
- Infrastructure:
  - 7 tape libraries, 83 tape drives / tape servers, 23k tapes
  - Current use: 180PB
  - Current capacity: 0.6 EB
  - Max throughput: 25GB/s



Data written to tape 2008-2016

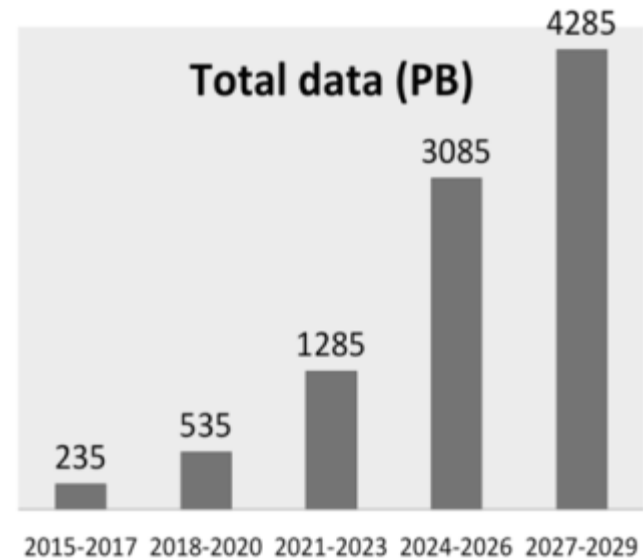


Monthly data written on tape in 2016



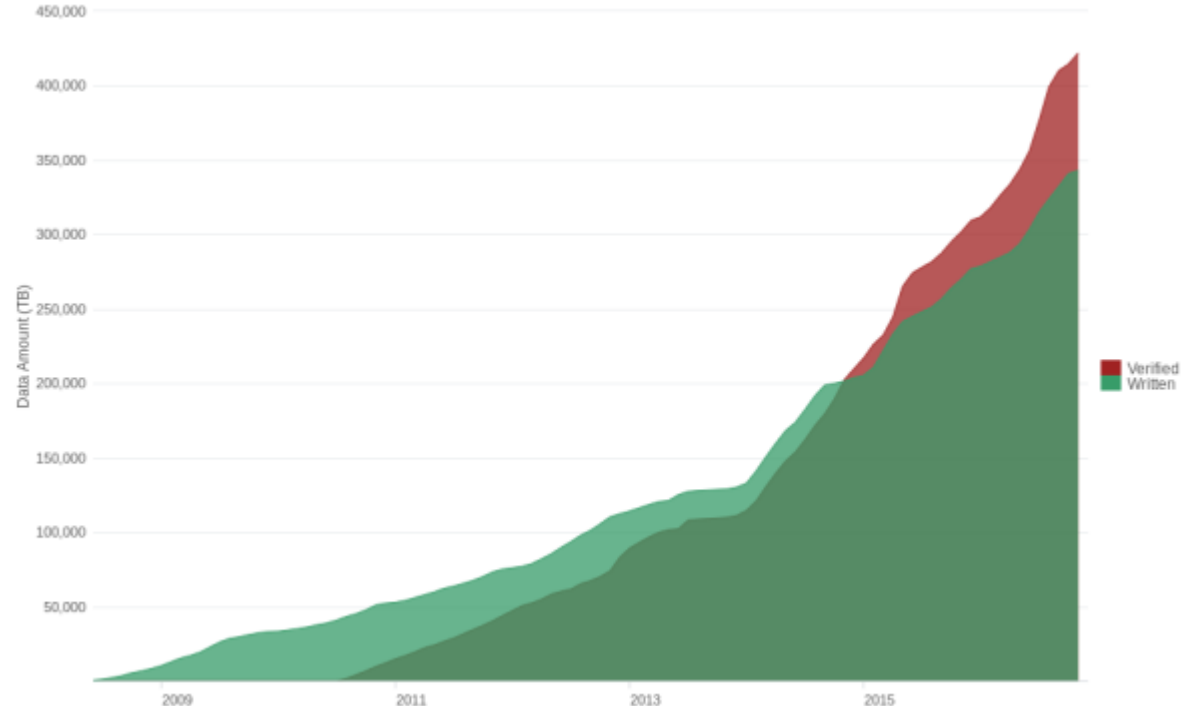
# Data Archiving

- Anticipating the future ...
  - 2017-2018: 100PB/year
  - 2021++ : 150PB/year
- ... and preserving the past
  - Protect against bit rot (data corruption, bit flips, environmental elements, media wear out and breakage ...)
  - Migrate data across technology generations, avoiding obsolescence
    - Some of our data is 40 years old



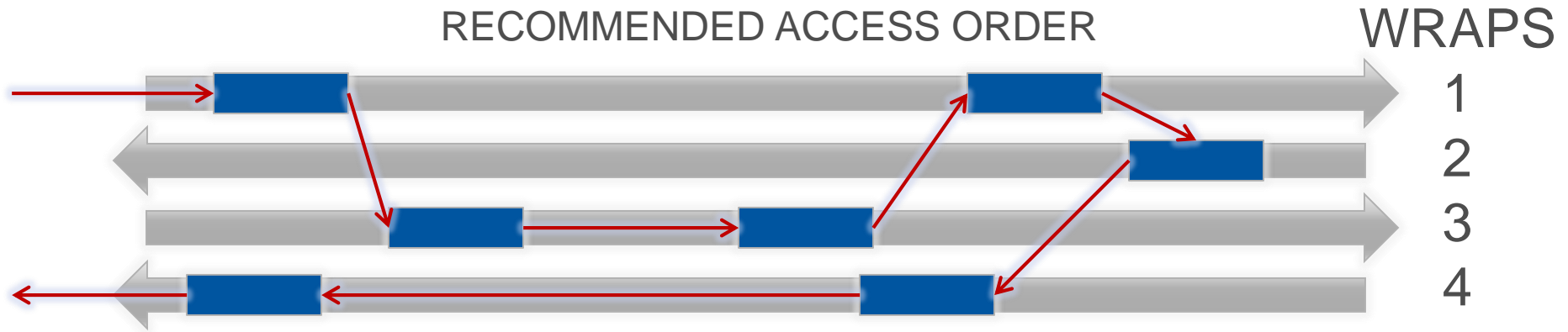
# Data Archiving

- Systematic media verification



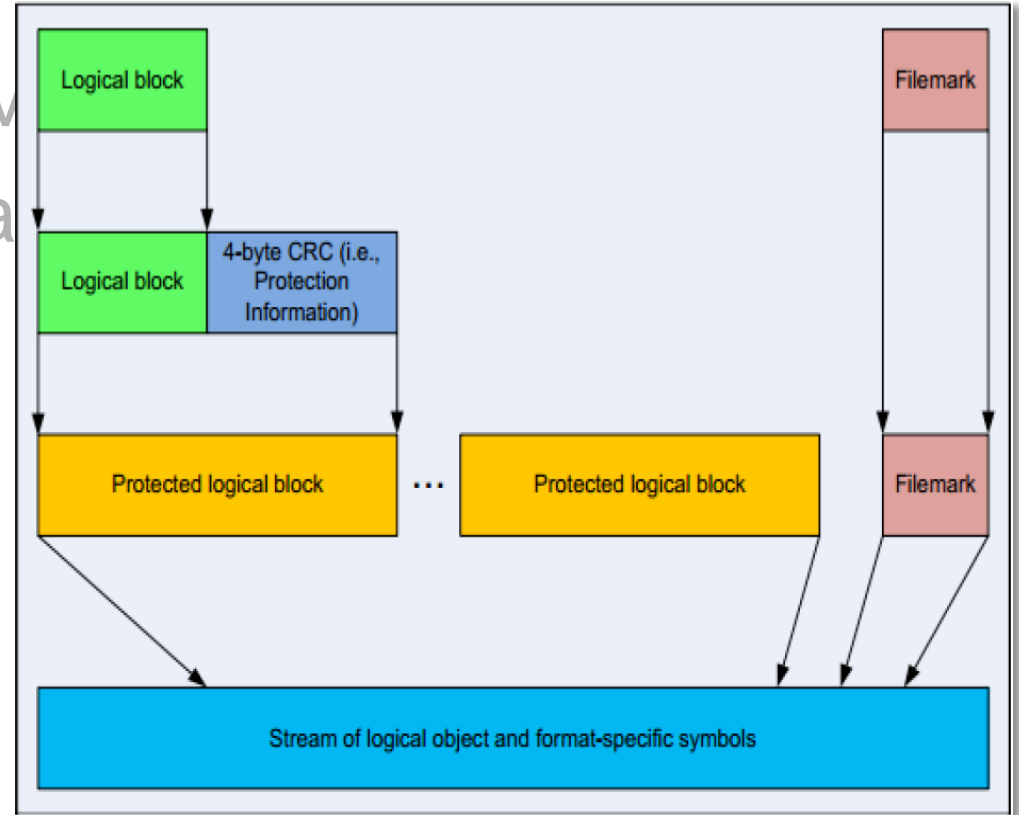
# Data Archiving

- Systematic media verification
- Optimising media access



# Data Archiving

- Systematic media verification
- Optimizing media allocation
- Improve reliability





# Data Archiving

- Systematic media verification
- Optimizing media access

• Im  
• In

```
2016-01-26T05:46:03.594249+01:00 tpsrv220 tapeserverd[3335]:  
LVL=Info TID=3350 MSG="Logging volume statistics"  
firmwareVersion="460E" lifetimeBOTPasses="1486"  
lifetimeMOTPasses="1556" lifetimeVolumeMounts="202"  
lifetimeVolumeRecoveredReadErrors="167"  
lifetimeVolumeRecoveredWriteErrors="30"  
lifetimeVolumeUnrecoveredReadErrors="4"  
lifetimeVolumeUnrecoveredWriteErrors="2"  
validity="1" volumeManufacturingDate="20110603"
```

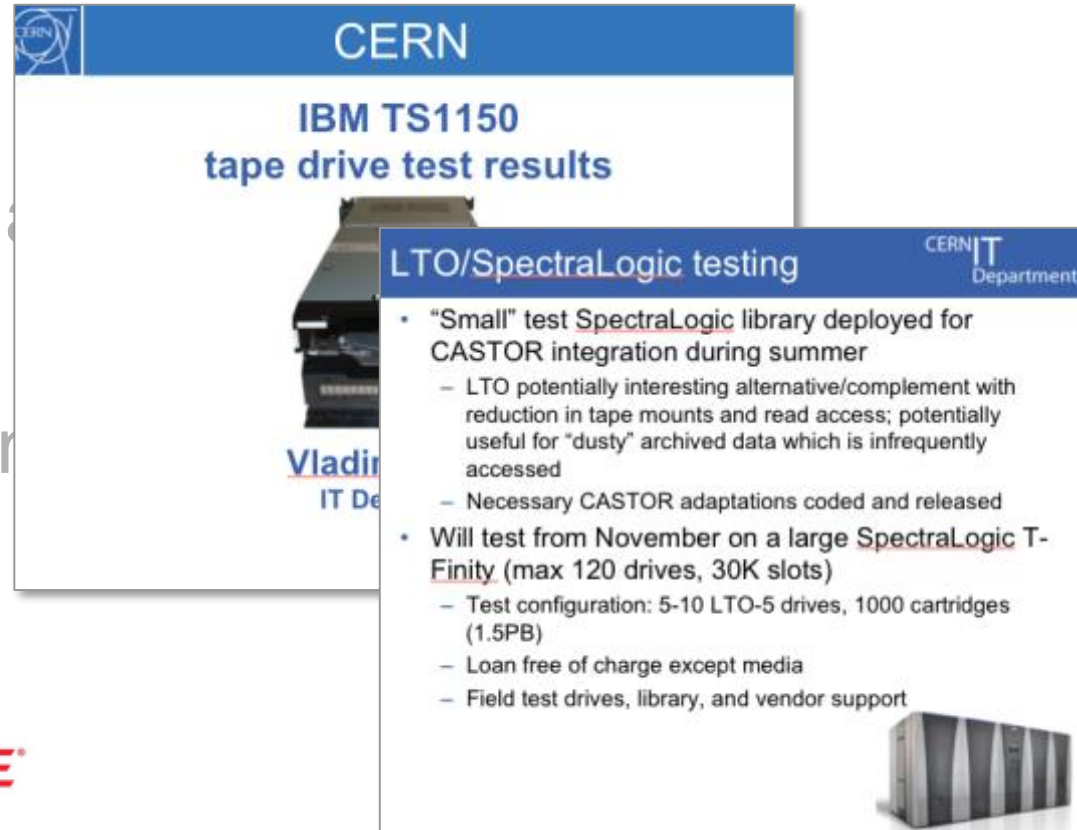
bring

→ Not good..

→ Really bad!

# Data Archiving

- Systematic media
- Optimizing media
- Improve reliability
- In depth performance
- Collaborations




**CERN**

**IBM TS1150**  
tape drive test results

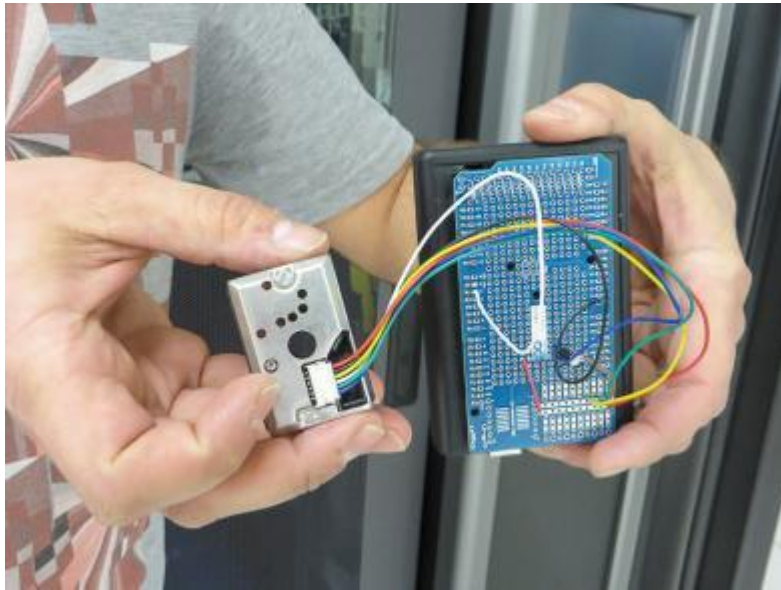
**LTO/SpectraLogic testing** CERN IT Department

- “Small” test SpectraLogic library deployed for CASTOR integration during summer
  - LTO potentially interesting alternative/complement with reduction in tape mounts and read access; potentially useful for “dusty” archived data which is infrequently accessed
  - Necessary CASTOR adaptations coded and released
- Will test from November on a large SpectraLogic T-Finity (max 120 drives, 30K slots)
  - Test configuration: 5-10 LTO-5 drives, 1000 cartridges (1.5PB)
  - Loan free of charge except media
  - Field test drives, library, and vendor support

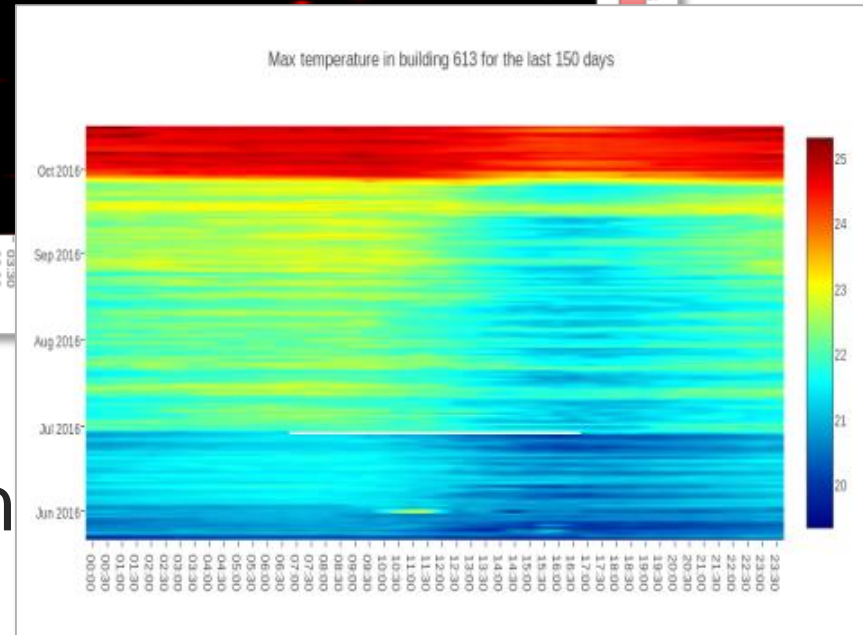
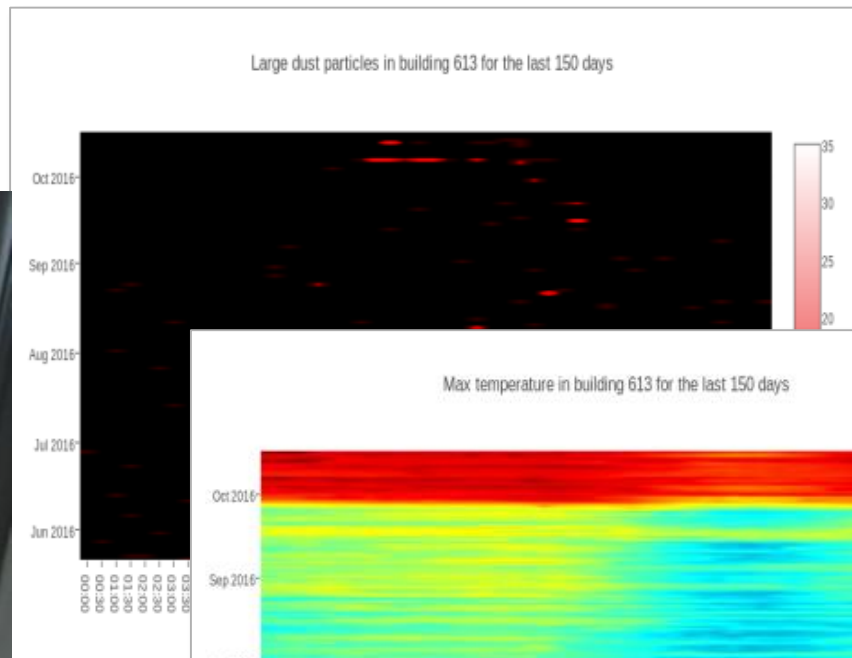
Vladimir  
IT De

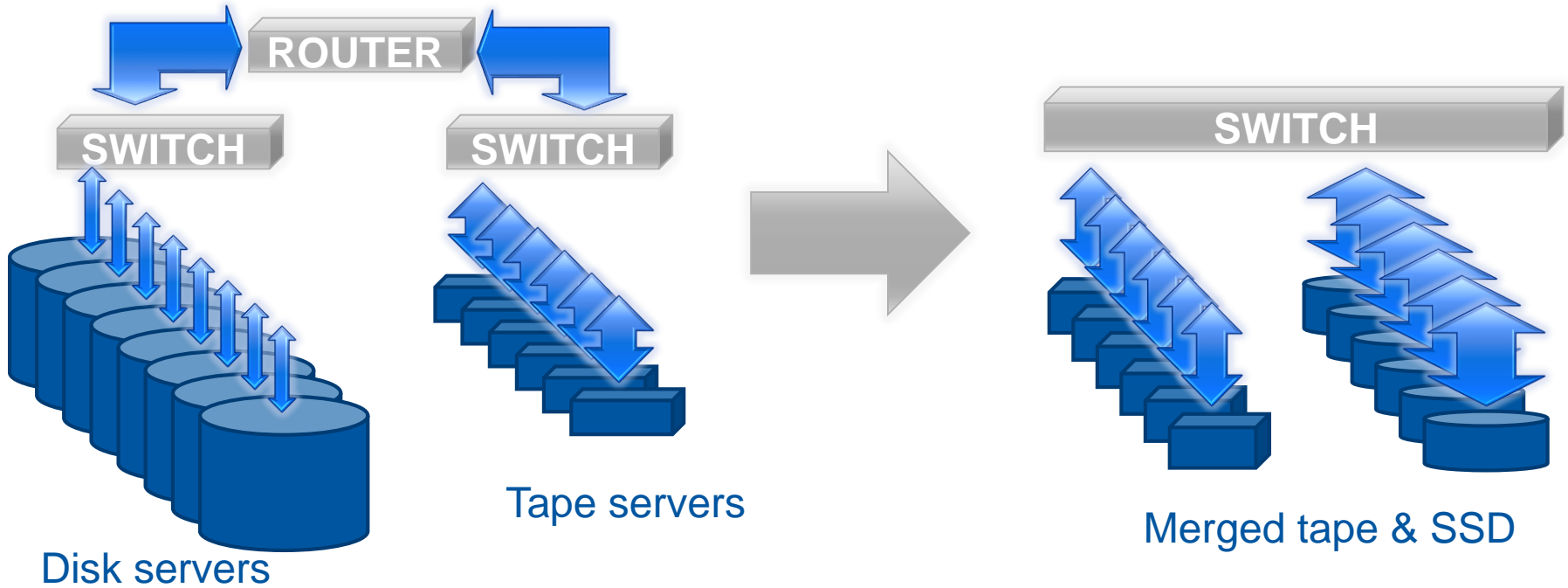


# Data Archiving



- Environmental protection



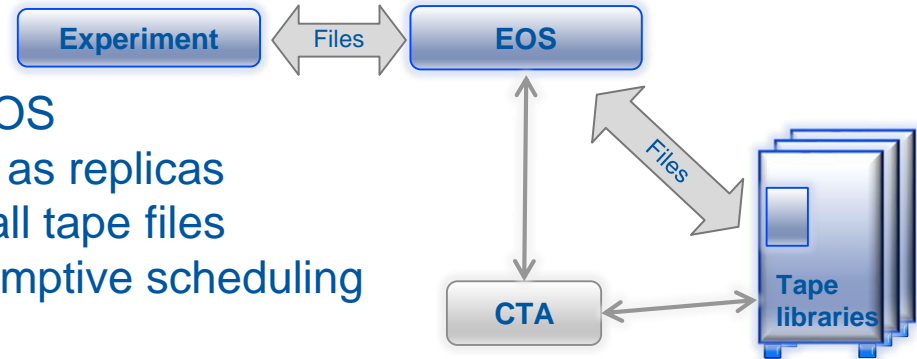


- Large scale media migration

# Data Archiving: Evolution

- EOS + tape ...
  - EOS is the strategic storage platform
  - Tape is the strategic long term archive medium
- EOS + tape = ❤️
  - Meet CTA : the CERN Tape Archive
  - Streamline data paths, software and infrastructure

- CTA is glued to the back of EOS
- EOS manages CTA tape files as replicas
- CTA contains a catalogue of all tape files
- CTA provides optimised, preemptive scheduling



# Data Archiving: Evolution

- CTA Timeline
  - End 2016: First functional prototype release
  - April 2017: First release for additional copy use cases
  - 2018: Production-ready version
- Easy migration path from CASTOR to EOS+CTA
  - Only metadata needs to be migrated
  - CASTOR tape format will be reused

# IT-ST-AD

Presented by Andreas Peters

# IT-ST-AD

## Storage - Analytics & Developments



DPM  
DynaFED



XRootD

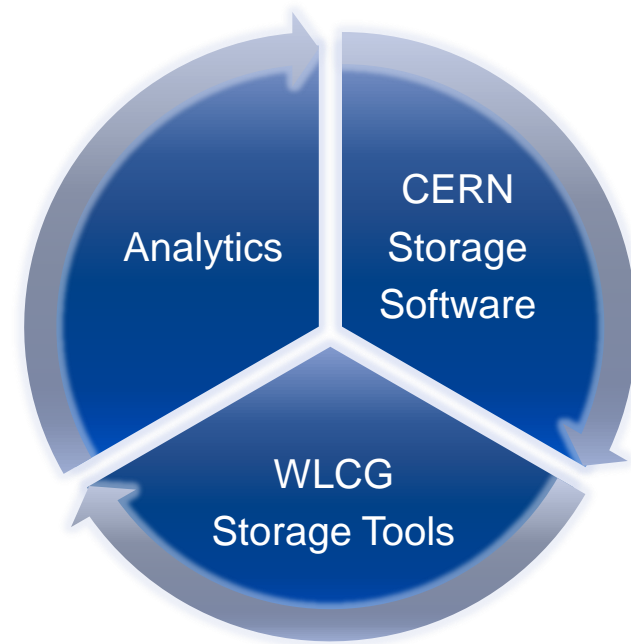


# Storage Services for Physics and more



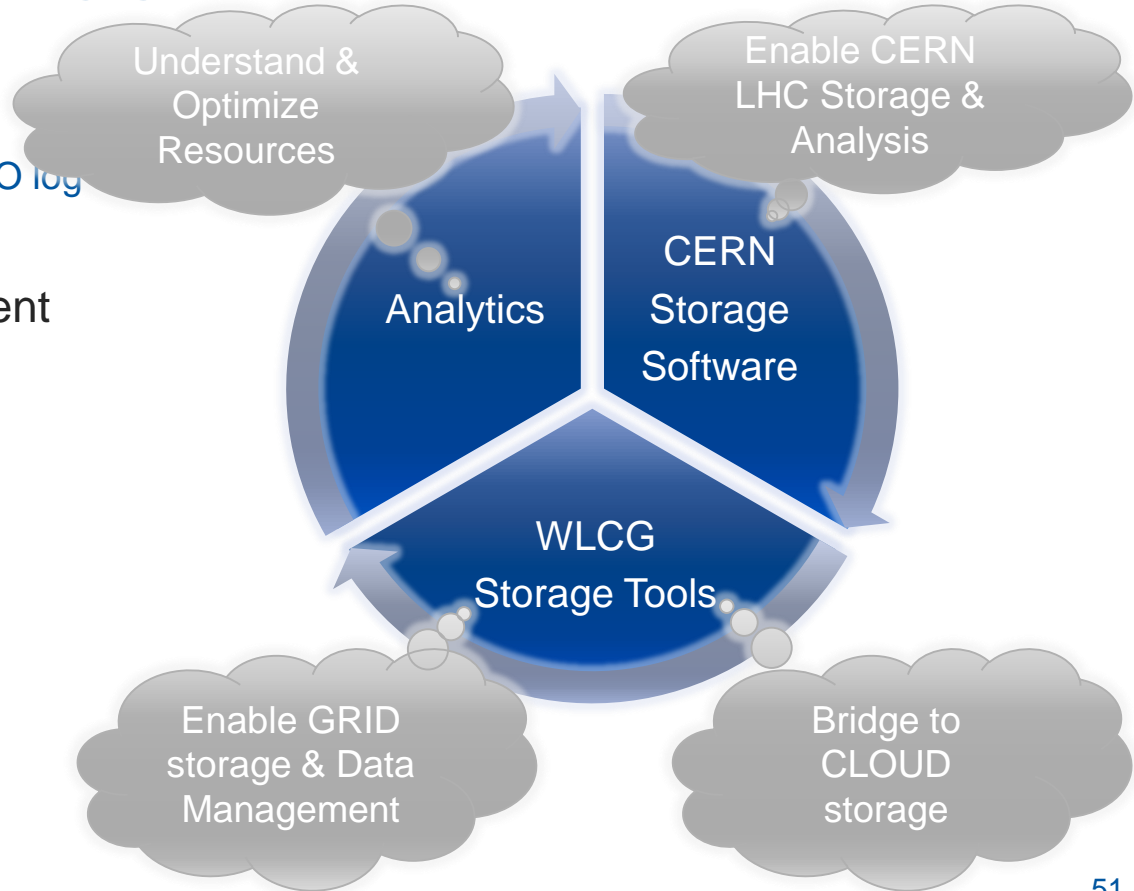
# Section Activities

- Analytics
  - Activities/Working Group - CPU/IO log analysis
  - Hadoop Service
- Storage Software Development
  - **CERN**
    - CASTOR<sup>XRootD</sup>
    - EOS / CernBOX
    - XRootD Client & Release Management
  - **WLCG**
    - DPM
    - Dynafed
    - GFAL & DAVIX
    - FTS3<sup>Dev&Ops</sup>



# Section Activities

- Analytics
  - Activities/Working Group - CPU/IO log
  - Hadoop Service
- Storage Software Development
  - **CERN**
    - CASTOR<sup>XRootD</sup>
    - EOS / CernBOX
    - XRootD Client & Release Management
  - **WLCG**
    - DPM
    - Dynafed
    - GFAL & DAVIX
    - FTS3<sup>Dev&Ops</sup>



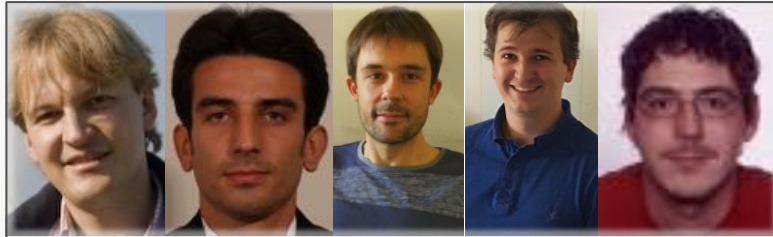
# Section Members

Dirk Düllmann Christian Nieke Rainer Többicke Luca Menicchetti



Analytics

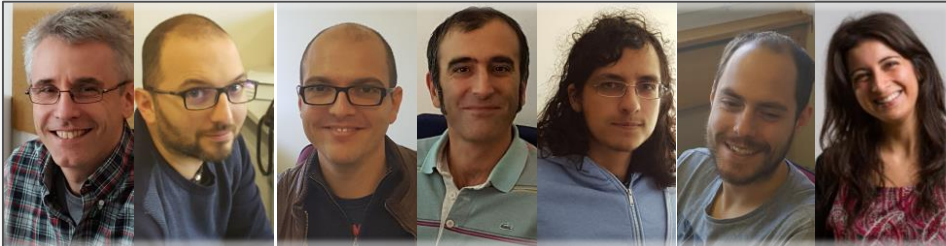
Andreas J. Peters Elvin A. Sindrilaru Geoffray Adde Michal K. Simon Paul Lensing



CERN Storage  
Development

WLCG Storage  
Development

Oliver Keeble Andrea Manzi Alejandro Alvarez Fabrizio Furano Georgios Bitzes Gerhard Rzehorz Maria Arsuaga Rios

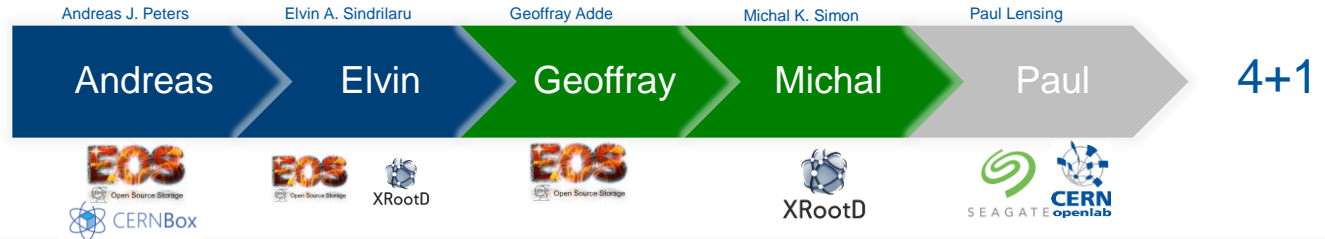


# Section Structure

## Analytics



## CERN Storage Development



## WLCG Storage Development



PhD Christian Nieke

# ANALYTICS

# Analytics – Hadoop Service



3 production cluster – joint activity of IT-ST & IT-DB

	lxhadoop	analytix	hadalytic
Purpose	general, de facto Atlas	monitoring, log analysis	beams dep, new development
Nodes	20	38	14
Size (used) PB (%)	1.5 (60)	2.5 (66)	2 (55)
Files	35 M	8 M	6M
Jobs/week	2500	6600	2500

## Highlights

- Analytix repository for EOS log files
- **SWAN** hadoop integration
- CERNbox/EOS access for map/reduce and **Spark**
- HDFS **backup** to Castor

## Future Challenge

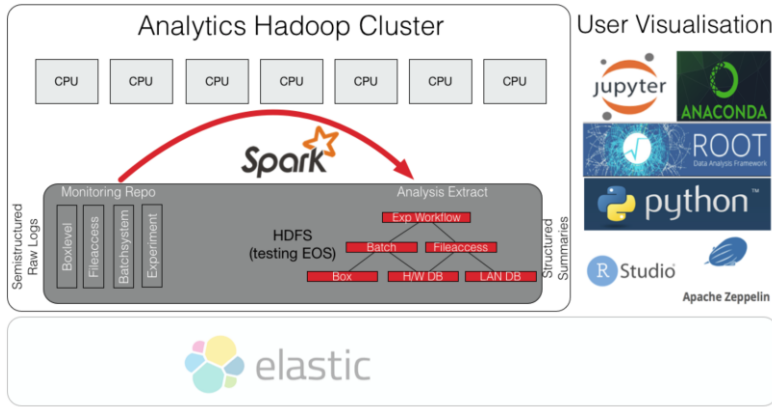
- Hadoop on EOS disk server

# Analytics – Log Storage & Analysis

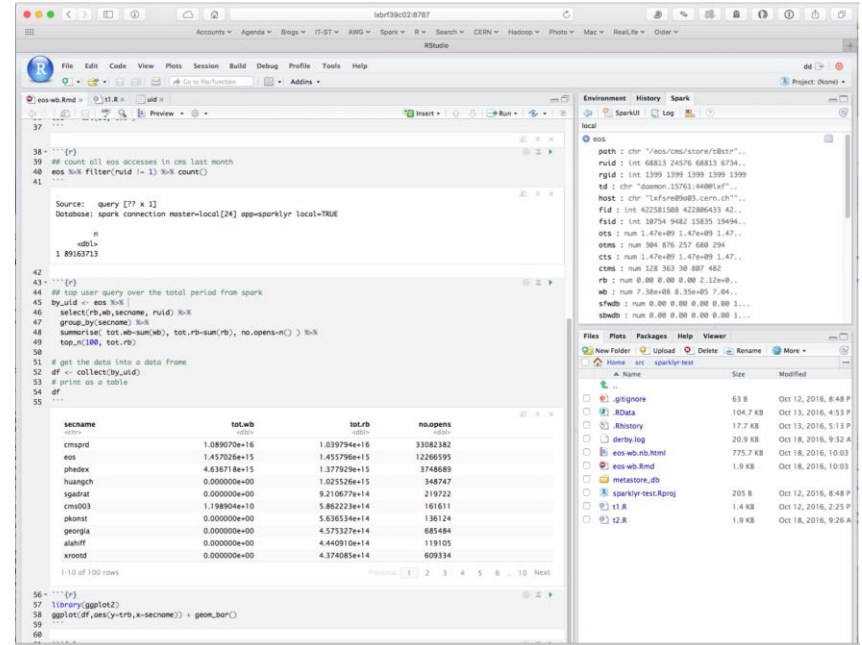


... feed log records via IT monitoring chain

... fast interactive log data access



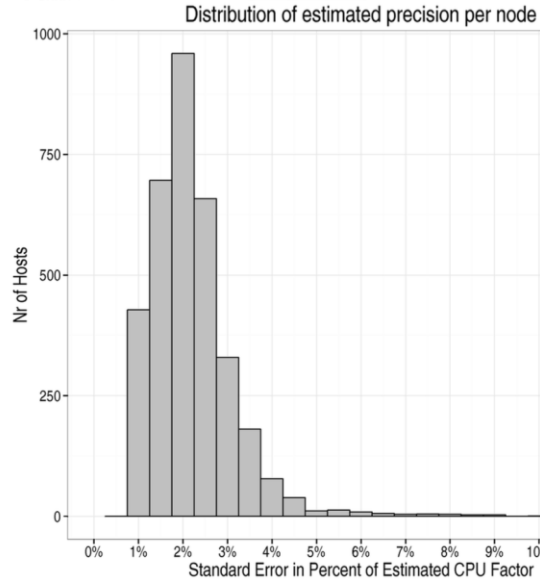
Hadoop & Friends – CHEP 2016





# Analytics – Performance Projects

PhD Christian Nieke



Study proposed by Dirk Düllmann and Darrell Long (University of California)

**BACKBLAZE** Personal Backup Business Backup B2 Cloud Storage Blog

Follow us: [f](#) [t](#) [s](#)

Cloud backup. Mac or PC. Unlimited data. \$5/month. And you can try it for free today.

## Hard Drive SMART Stats

November 12th, 2014

I've shared a lot of Backblaze data about hard drive failure statistics. While our system handles a drive failing, we prefer to predict drive failures, and use the hard drives' built-in SMART metrics to help. The dirty industry secret? SMART stats are inconsistent from hard drive to hard drive.

## Passive CPU benchmark

performance prediction error better than 5%

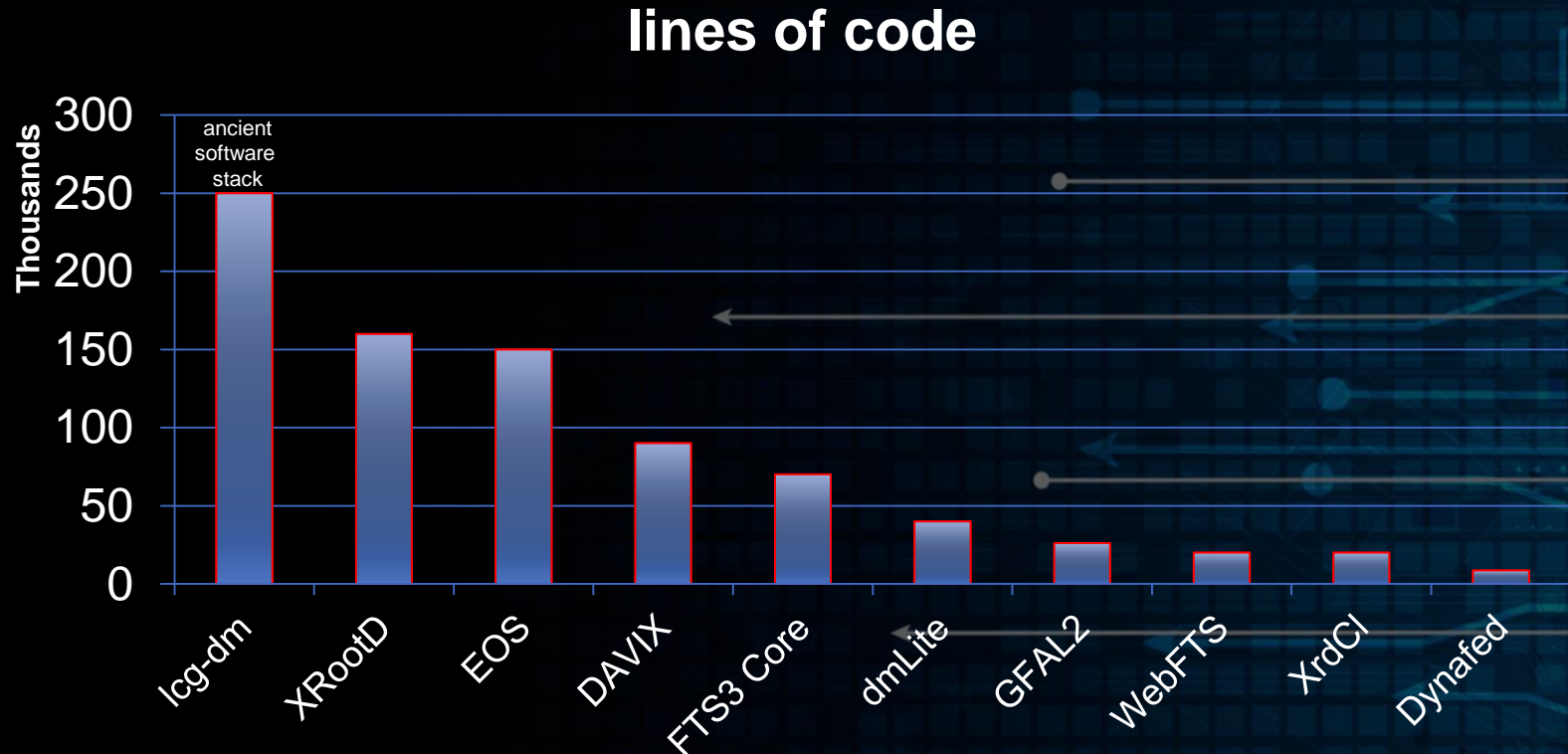
## Harddrive failure analysis

Evaluate enterprise, consumer, network drives  
For service planning & predictive maintenance



**CERN**  
STORAGE SOFTWARE DEVELOPMENT

# Open Source Software Projects

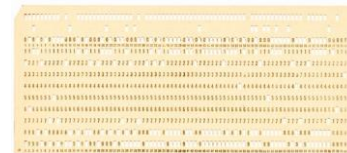


# XRootD



## What is that? What is that for?

- provides high performance, scalable fault tolerant access to data repositories
- main protocol for WLCG WAN access – latency optimized
  - >>99% of LAN data access at CERN via XRootD protocol
  - tight integration in ROOT framework
- provides a data oriented client-server C++ framework(plugin structure) - used by
  - EOS
  - LSST
- Participation via collaboration - core development SLAC



## Our Responsibilities

- **client development** since 2010
- **release management & distribution**

## Good to know ...

- since 2010 XRootD **V3** – production (e.g. EOS)
- since 2014 XRootD **V4** **IPV6, FileCache, HTTP(S) bridge, CEPH bridge** production (e.g. CASTOR, Ixplus, LHC exp.)  
wire protocol security added **V4.5**
- **V5** expected in 2017

## Future Challenges...

- multi-source support, TPC<sup>2</sup>



**GITHUB** <https://github.com/xrootd/xrootd.git>  
1000 source files – 160k lines of code



# EOS



## What is that?

- Large scale data storage system (150 PB@CERN)
  - Multi protocol XRootD HTTP(S) WebDAV, FUSE
  - Strong Authentication KRB5 X509
  - Flexible provisioning and life-cycle management
  - Large user community features (limits, quotas)
- since 2010 developed in-house in storage group
- since 2011 production service at CERN
  - **+11 HEP EOS sites not at CERN**  
LBL, ORNL, SUBATECH, IHEP, SASKE, FNAL, UNAM, MEFH, SINICA, INFN, UNIVIE
- simple scalable architecture
  - **Meta data scale-up 0.5b files**
  - **Data scale-out 15k disks**
  - **Messaging 1.5Mhz pipelined**

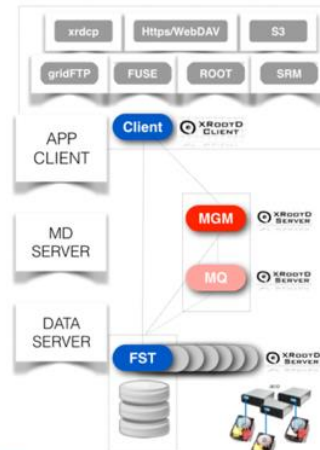


October 2016

**150 PB** disk space > **110 PB** used

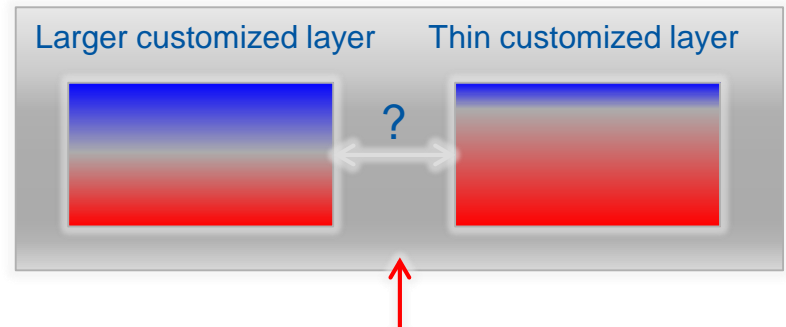
• **1143** storage nodes

1y growth rates: **+85%** more files **+50%** more space used

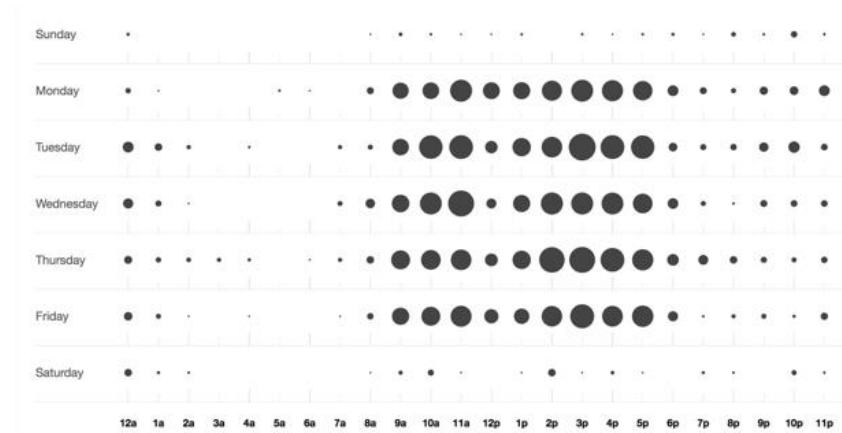
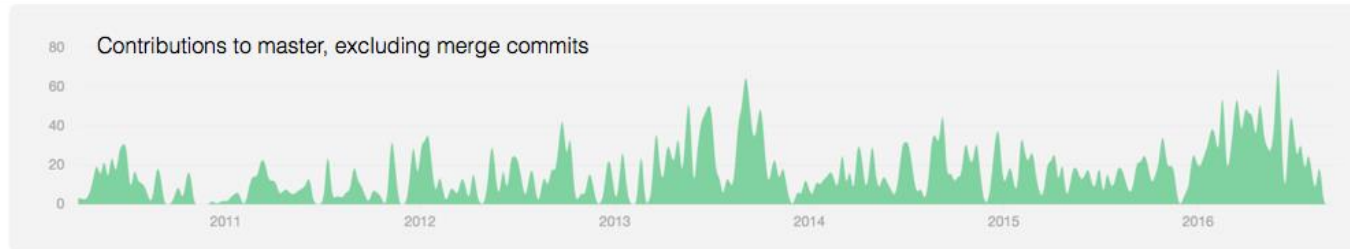


## What is the current development strategy?

- Evolve **flexible** storage infrastructure as **scale-out solution** towards exabyte scale
  - one storage – many views on data – analysis/analytics service integration
  - suitable for LAN & WAN deployments
  - enable cost optimized deployments
- Evolve **file-system** interface /eos
  - multi-platform access using standard protocol bridges CIFS, NFS, HTTP
- Evolve scale-out **storage back-end** support
  - object disks, object stores, public cloud, cold storage
- Evolve **namespace scalability**
  - scale-out implementation
- **Review** regularly break-point between in-house and community open-source layers



# EOS Development



GITHUB punch card

master: 181k LOC  
production: 150k LOC  
excluding external projects  
jerasure/sqlite/gf-complete

70 tags in production branch  
34 tags in master branch

8 contributors

+now several external contributors  
AARNET + COMTRADE



# EOS Challenges



- 1 Meta data *scale-up* ⇒ **scale-out**
  - *in-memory namespace* ⇒ *in-memory scale-up namespace cache + persistency in scale-out KV store* → 1<sup>st</sup> key project: EOS Citrine & QuarkDB
- 2 *Remote access APIs* ⇒ **filesystem API**  
→ 2<sup>nd</sup> key project: EOS FUSE rewrite



# EOS Scale-Out Namespace

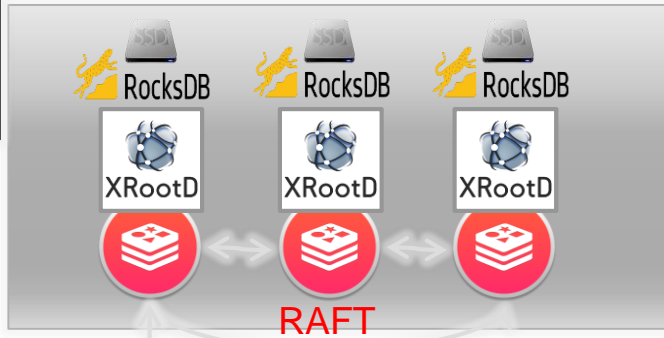


Open Source Storage

1st Key Project

## QuarkDB

- built on open source software/ideas

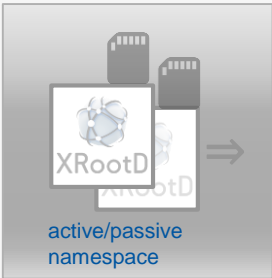


Facebook KV library for SSD persistency

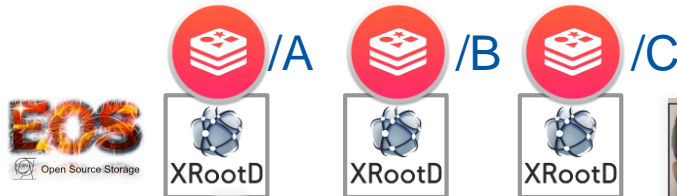
Redis XRootD plug-in

Redis protocol KV,SET,HASH

current EOS



REDIS



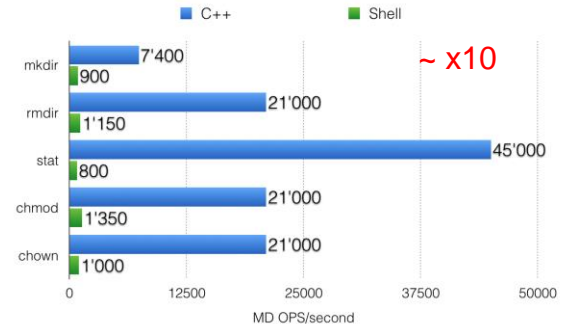
meta data server with namespace cache & plug-in for KV backend  
⇒ stateless front-end architecture



# EOS FUSE improvements

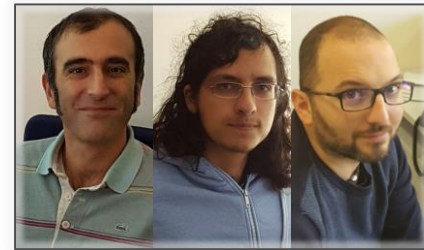


- FUSE was added as an external layer to an existing server API [path-inode translations]
  - 1<sup>st</sup> implementation 2011 clone of XRootD FUSE xrootdfs C-interface
  - 2<sup>nd</sup> implementation 2016 introducing many latency optimizations and pure C++  
**gained 10x + more POSIX**
    - meta-data consistency very difficult/race conditions without server side changes – MD updates still synchronous
  - 3<sup>rd</sup> implementation on the way **gain 10x + POSIX**
    - client leases, async meta-data pipelining
    - optional persistent meta-data & data cache
    - **libfuse2** ⇒ **libfuse3** 21k seq IOPS ⇒ 500k seq IOPS [wb cache]





# DPM



## What is that?

- CERN **Disk Pool Manager** used in many GRID sites (outside CERN)
- originally storage software fork from **CASTOR1** code developed since 16 years
- *the* Tier-2 GRID storage solution – deployed at 160 sites world-wide providing 70 PB of storage space
- First refactoring included dmLite stack as plug-in architecture in 2011
- Latest release **1.9** includes now DOME component 2016
  - allows to to remove old legacy code & daemons ⇒ simplify and consolidate the architecture
    - better modularization
    - single internal communication protocol
  - SRM less operation
  - quota, space reporting, caching hooks
- Caching laboratory for ‘unmanaged’ Tier-2 storage, low-effort deployment, long-term consolidation
- DPM Workshop Paris <https://indico.cern.ch/event/559673/>

# Key Components to Cloud Storage

- **DAVIX**



- WebDAV/S3 C++ client – integrated in ROOT
- developed since 2012 – focus on long term stability
  - latest features
    - compatible with Azure & AWS v4 authentication
    - performance optimizations by range coalescing and request parallelization

- **DynaFED**



- Provides in-memory cached namespace federation of multiple HTTP enabled storage systems
- Provides a Grid⇒Cloud bridging functionality translating X509 identities to signed URLs using DAVIX

- **GFAL**



- Grid file access library
  - provides protocol abstraction layer [⇒ FTS]

<https://gitlab.cern.ch/dmc/gfal2>



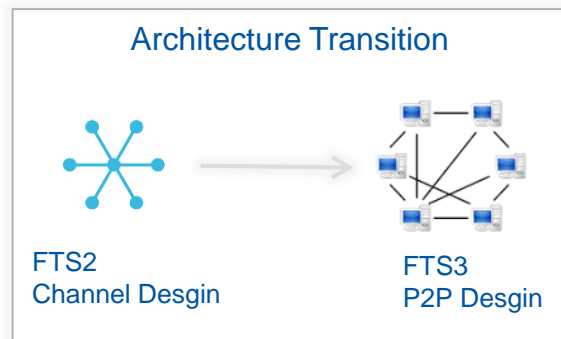


## What is that?

- 3<sup>rd</sup> generation of file transfer service building core platform for scheduled LHC data transfers
  - developed since a 2001 at CERN
  - 1<sup>st</sup> production release 2013
  - key component in WLCG
  - transfers hundreds of peta bytes per year worldwide
  - Latest release 3.5
    - Transfer Optimization based on network performance analysis
    - Decision visualization

## Future Challenges

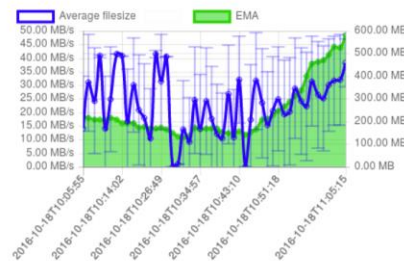
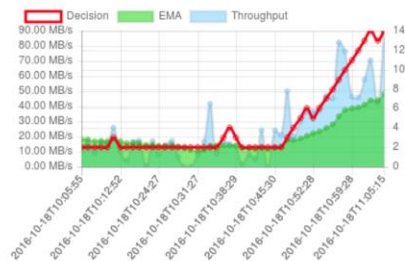
- Boost Scalability of database centric design
  - Initiated R&D project **Flutter** trying more radical redesign solutions
- HPC integration (dual credential support)
- Alarms & Containerization



# FTS Optimizer Visualization



Details for srm://ccsrm.in2p3.fr → srm://atlassrm-fzk.gridka.de



First Previous 1 Next Last

Timestamp	Decision	Running	Queue	Success rate (last 1min)	Throughput	EMA	Diff	Explanation
2016-10-18T11:06:2'	13	21	3940	100.00%	48.359 MB/s	48.426 MB/s	-1	Good link efficiency, throughput deterioration
2016-10-18T11:05:1'	14	14	3948	100.00%	89.522 MB/s	48.434 MB/s	1	Good link efficiency, current average throughput is larger than the preceding average
2016-10-18T11:04:0'	13	14	3968	100.00%	42.308 MB/s	43.869 MB/s	-1	Good link efficiency, throughput deterioration
2016-10-18T11:02:5'	14	13	3986	100.00%	70.434 MB/s	44.042 MB/s	1	Good link efficiency, current average throughput is larger than the preceding average
2016-10-18T11:01:4'	13	14	4006	99.00%	57.463 MB/s	41.109 MB/s	1	Good link efficiency, current average throughput is larger than the preceding average

## Optimizer Decision Visualization

## Optimizer Status Explanation



# FTS Operations



- AD responsible for operations since January '16
  - 2 cluster production & pilot
  - Moved to CC7
  - IPV6 enabled
  - Service Monitoring  
[https://monit.cern.ch/app/kibana#/dashboard/\\_project-FTS-Service-Level](https://monit.cern.ch/app/kibana#/dashboard/_project-FTS-Service-Level)
  - Transfer Log archival in EOS on the way



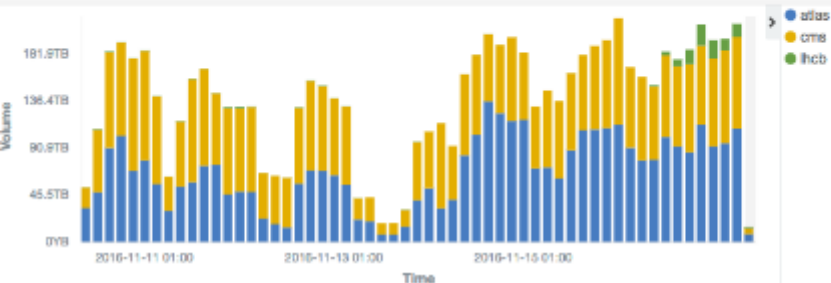


# FTS @ CERN last week stats

MONIT FTS Dashboard Navigation

Home - Overview - Transfer Plots - Matrix View - Failures - Custom Views - Servers Configuration

MONIT FTS Volume Transferred (per VO)



MONIT FTS Volume Transferred

9.3PB

Data Transferred

MONIT FTS Efficiency

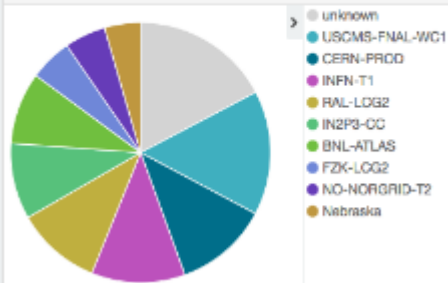
0.905

Transfer Efficiency

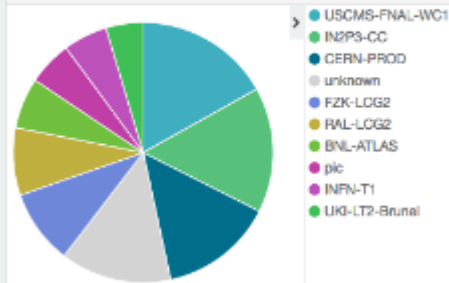
MONIT FTS Transfer Efficiency (per VO)



MONIT FTS Top 10 Source Sites



MONIT FTS Top 10 Destination Sites



# WebFTS & LastMile

<https://webfts.cern.ch>

- WebFTS



- end-user transfer portal
- targets small VOs

- LastMile

- end-point less FTS transfers
  - transfer data to/from workstations/laptops
  - run a client – not a server



**What is WebFTS?**

WebFTS is a file transfer and management solution which allows users to invoke reliable, managed data transfers on distributed infrastructures.

Created following simplicity and efficiency criteria, WebFTS allows the user to access and interact with multiple storage elements. Their content becomes browsable and different filters can be applied to get a set of files to be transferred. Transfers can be invoked and capabilities are provided for checking the detailed status of the different transfers and resubmitting any of them with only one click.

The "transfer engine" used is FTS3, the service responsible for distributing the majority of LHC data across WLCG infrastructure. This provides WebFTS with reliable, multi-protocol (gridftp, srm, http, xrootd), adaptively optimised data transfers

Continue...

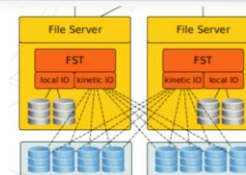
# IT-ST-AD Collaborations



EOS white-paper – now core software contributor



EOS OpenKinetic interface



Joint Research Centre

EOS collaboration

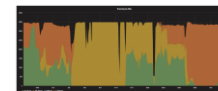


Australia's Academic and Research Network

EOS/CERNbox collaboration



DPM Grid storage in the cloud



Dynafed demonstrator

University of Victoria (CA), INFN and Belle-II

<https://indico.cern.ch/event/394788>



Data Intensive Cloud Research PhD focused on modeling ATLAS performance

IT-ST



[www.cern.ch](http://www.cern.ch)

# IT-Storage Plans for 2017

- EOS as the strategic storage platform
  - New, **scalable name server** (new REDIS-based implementation)
  - **Improved FUSE** to satisfy more requirements
- CERNBOX as the strategic service to access EOS storage
  - Offer **migration paths out of AFS** (target = EOS)
  - Start a discussion with relevant service managers on the **home directory service**
    - Define an architecture to host Ixplus home directory
    - Improve Windows support, Samba service for Windows Clients
- CERN Tape Archive
  - First internal release interfacing EOS planned end 2016 using IBMLIB1 (test library)
  - Apr 2017: D1T1 and D0T1 workflows available in EOS
- TSM
  - New license to be negotiated, investigate alternative options for Oracle backups
- FTS: improvements on manageability and scalability
  - Files to/from laptops
- DPM: Consolidation
  - Going SRM-less with DPM
- CEPH the backend storage solution
  - Understand if **CEPHFS** can replace the Filer service

# Strategic Vision for IT-ST

- Build a flexible storage infrastructure
  - **Unlimited storage for LHC experiments, at exabyte scale**
    - Disk pools (EOS on-demand reliability, on-demand performance)
    - CERN Tape Archive (high reliability, low cost, data preservation)
    - Backbone of the CERN cloud (CEPH block storage for Openstack and Filers)
    - Cluster(s) for analytics (provided with Hadoop)
  - **A generic home directory infrastructure**
    - Fuse mounts, NFS and SMB exports, (and HADOOP).
    - CERNBOX (Sync client and Web/HTTP/DAV Access)
- Maintain and enhance (grid) data management tools
  - **To empower global scientific data workflows**
    - Data transfers (FTS), data caches (DPM, CVMFS), ...



[www.cern.ch](http://www.cern.ch)