# Data services for LHC computing

SLAC
CHEP 2016

Xavier Espinal
on behalf of IT/ST

Reliable

Fast Processing
DAQ Feedback loop

DAQ to CC
8GB/s+4xReco

Hot files

WAN aware
Tier-1/2 replica, multi-site

High throughout to tape
350+MB/s/drive - 12GB/s Pb-Pb

back-up

Filesystem 'feeling'
$HOME, SW-dist, Data

Consistent

∞

Few fast streams
CDR 2x40Gbps

Non-LHC and Local
Less structured, small communities
Unexpected usage  Catalogue=Namespace

disk and gc?

Many slow clients
Repro, reco, analysis  constant >20k

Endpoint Mounts
ie. /atlas in the WNs

CERN

IT-ST

2

cernbox

EOS
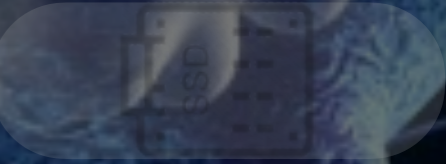
AFS

cvmfs

NFS

RBD

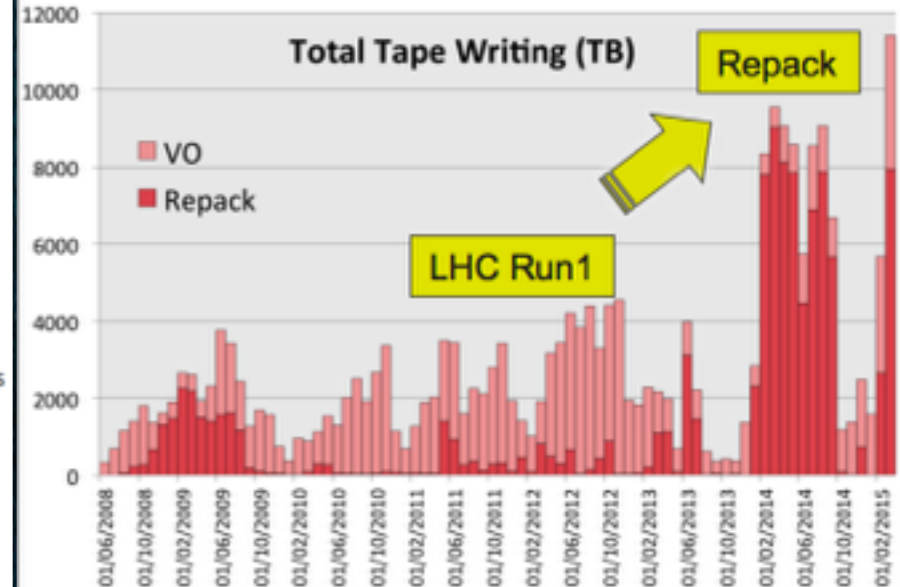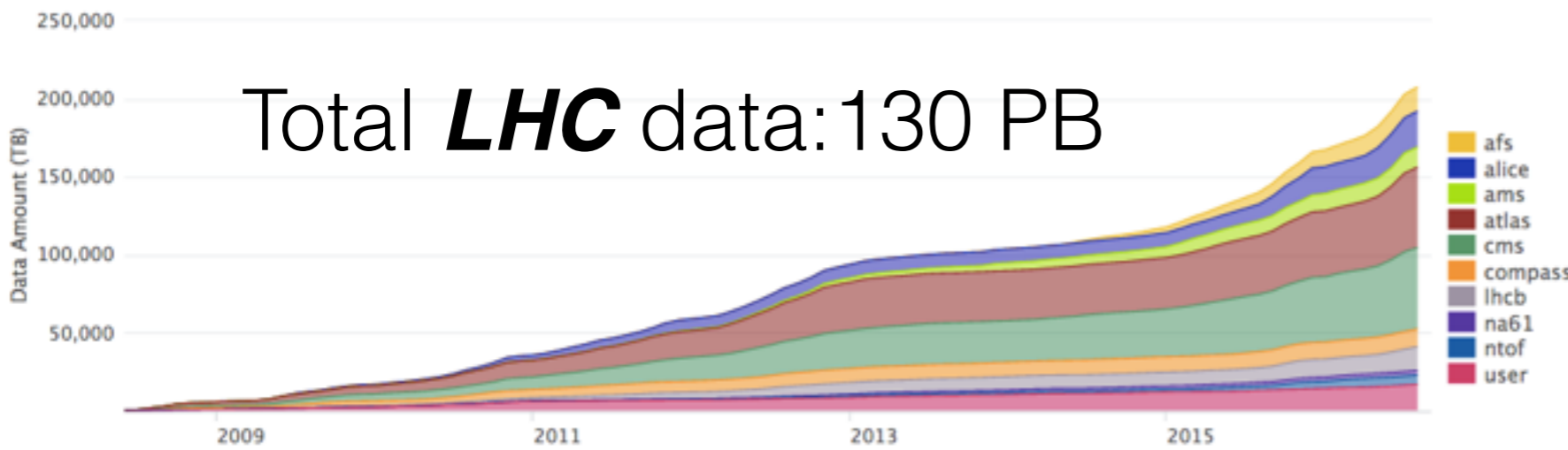S3

ceph

**CASTOR** — CERN Advanced STORage manager

Evolved to Tape oriented system
Key feature Per stream speed

Biggest physics-repo worldwide 175PB and +500M files
Towards a pluggable tape backend (EOS)
Cold by definition: hight throughput, high latency

Tape best technology for data repositories: TCO media power density and resilient/reliable very large disk caches nowadays
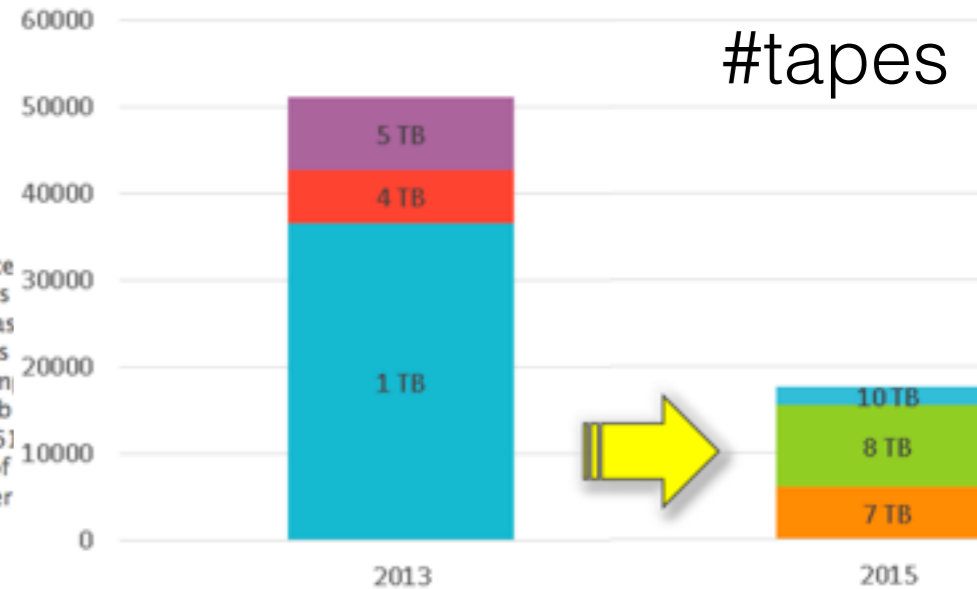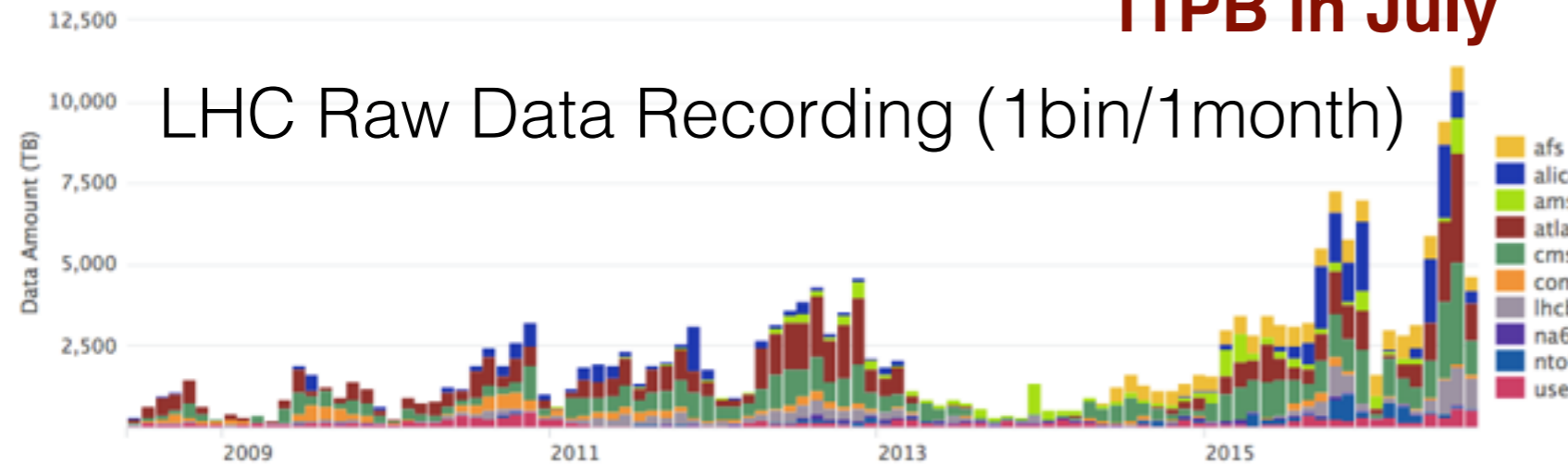
# CERN Tape Archive

Biggest physics-repo worldwide 175PB and +500M files

**Towards a pluggable tape backend (EOS)**

Cold by definition:  hight throughput, high latency

Tape best technology for data repositories: TCO media power density and resilient/reliable very large disk caches nowadays



Total *LHC* data:130 PB



Total Tape Writing (TB)

**11PB in July**

LHC Raw Data Recording (1bin/1month)

#tapes

IT-ST

# EOS now

October 2016

+1200 🖥
+45000 💾

850M files
150PB 💾

EB era

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



ALICE CDR

| | Avg | Max | Last | | Max |
|---|---|---|---|---|---|
| From ALICE to B513 | 1.32G | 28.33G | 240.76 | / Peak: | 28. |
| to ALICE | 8.76M | 210.84M | 5.88k | / Peak: | 210.84M |

Last update: Sat Oct 01 2016 08:52:40
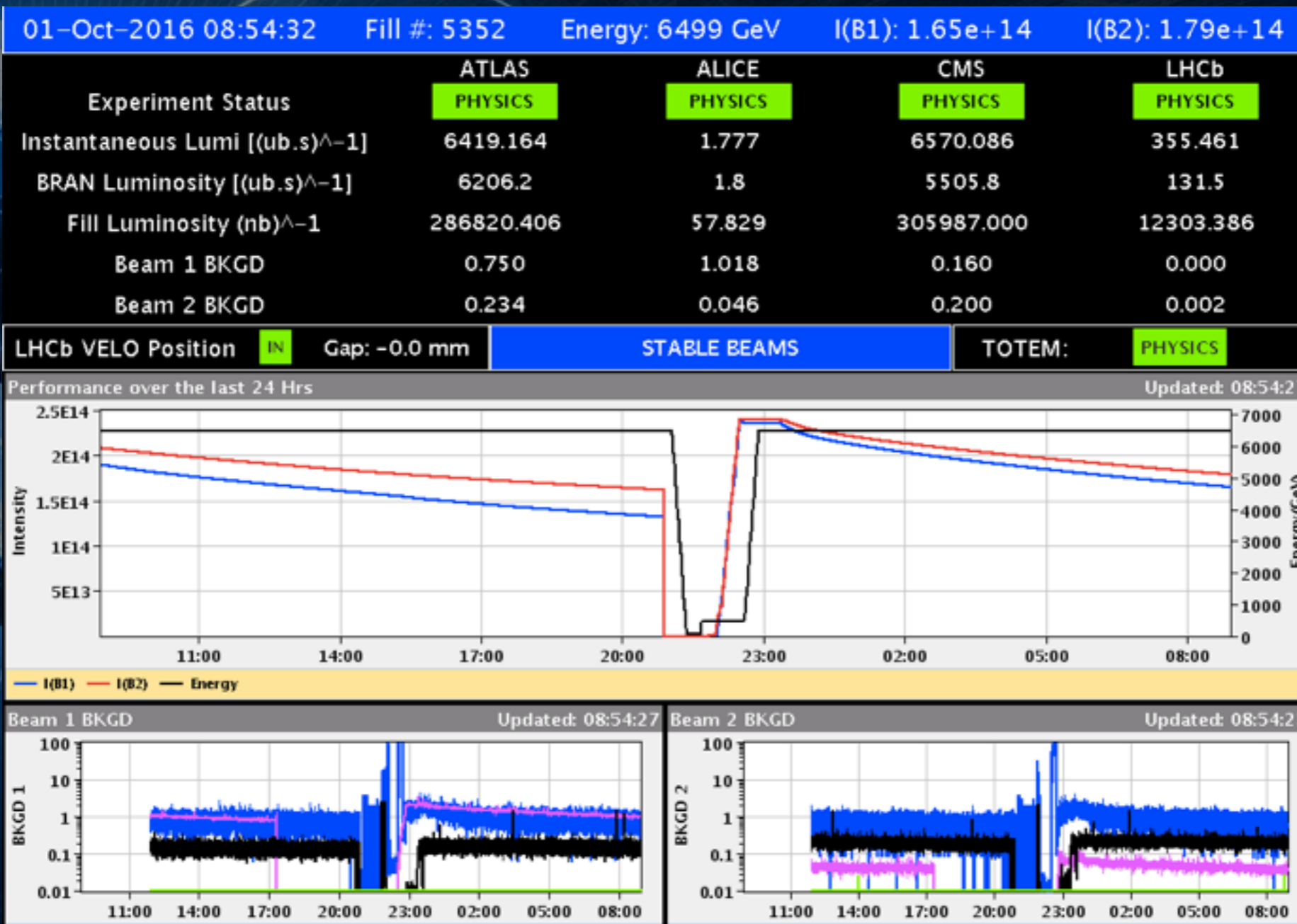
ATLAS CDR

| | Avg | Max | Last | | Max |
|---|---|---|---|---|---|
| From ATLAS to B513 | 5.06G | 21.01G | 34.68M | / Peak: | 21. |
| to ATLAS | 31.22M | 168.64M | 4.32M | / Peak: | 168.64M |

Last update: Sat Oct 01 2016 08:52:40

CMS CDR

| | Avg | Max | Last | | Ma |
|---|---|---|---|---|---|
| From CMS to B513 | 4.79G | 10.02G | 4.40G | / Peak: | 10. |
| to CMS | 1.51G | 2.98G | 989.02M | / Peak: | 2.9 |

Last update: Sat Oct 01 2016 08:52:40

LHCb CDR

| | Avg | Max | Last | | Max |
|---|---|---|---|---|---|
| From LHCb to B513 | 2.55G | 8.43G | 1.87G | / Peak: | 8.4 |
| to LHCb | 15.01M | 109.50M | 7.11M | / Peak: | 109.50M |

Last update: Sat Oct 01 2016 08:52:40
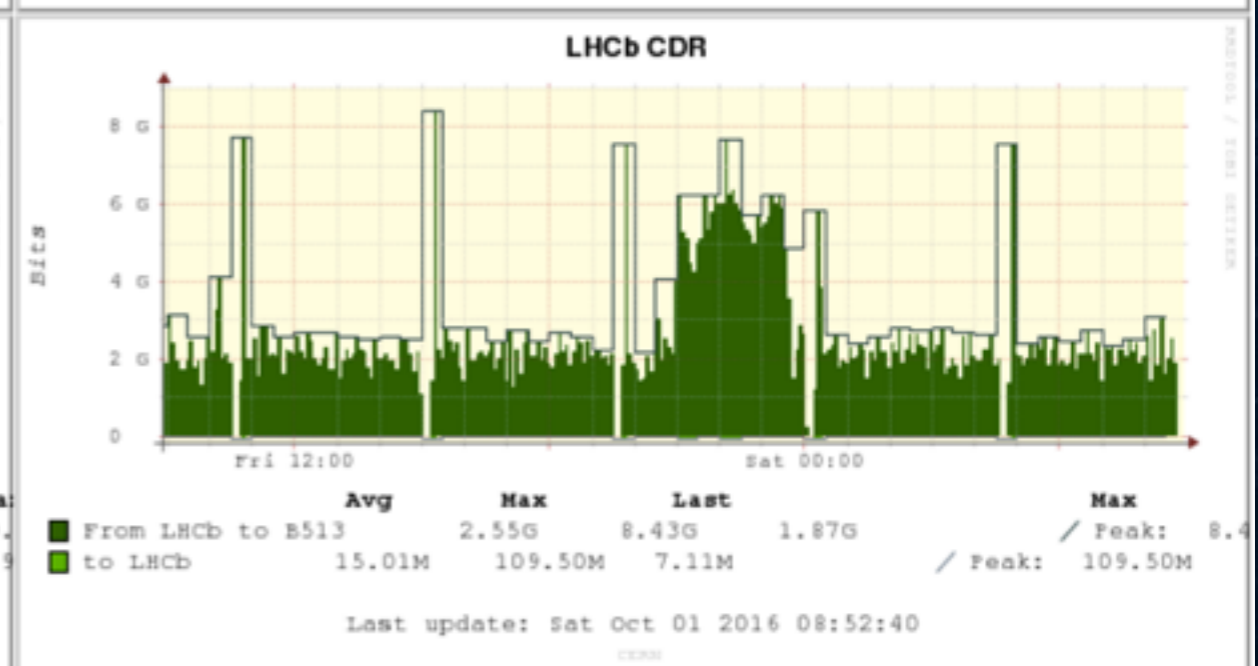
CERN
IT-ST

# EOS now

October 2016

+1200 🖥
+45000 💾

850M files
150PB 💾

EB era

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

---

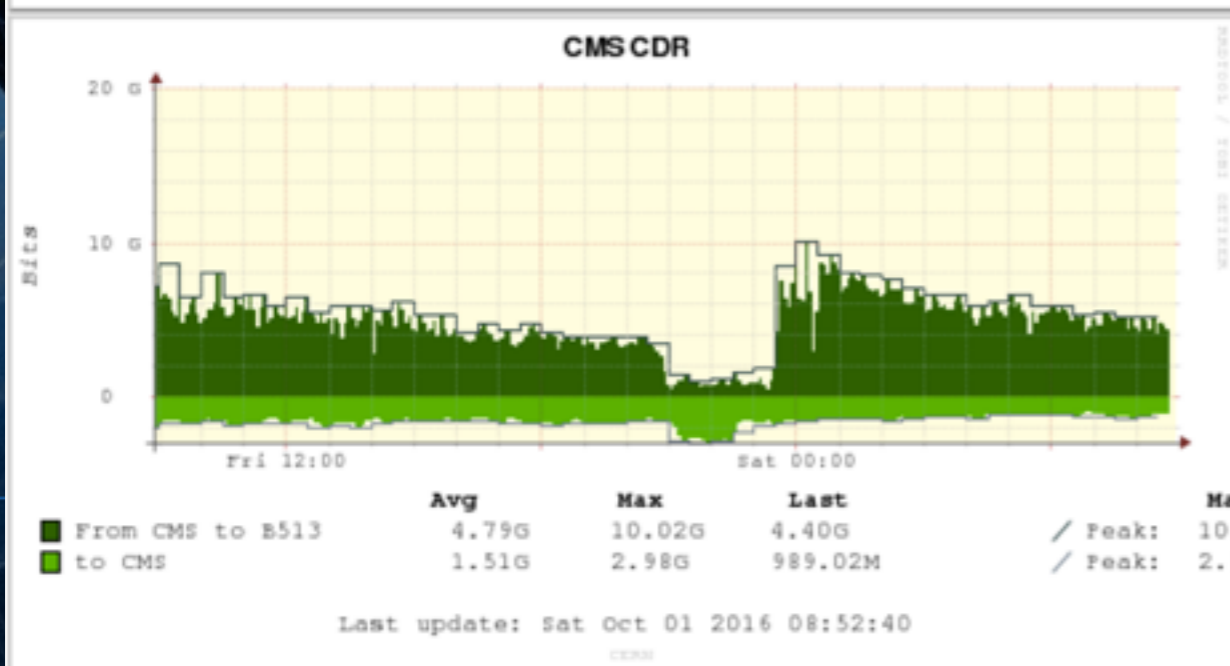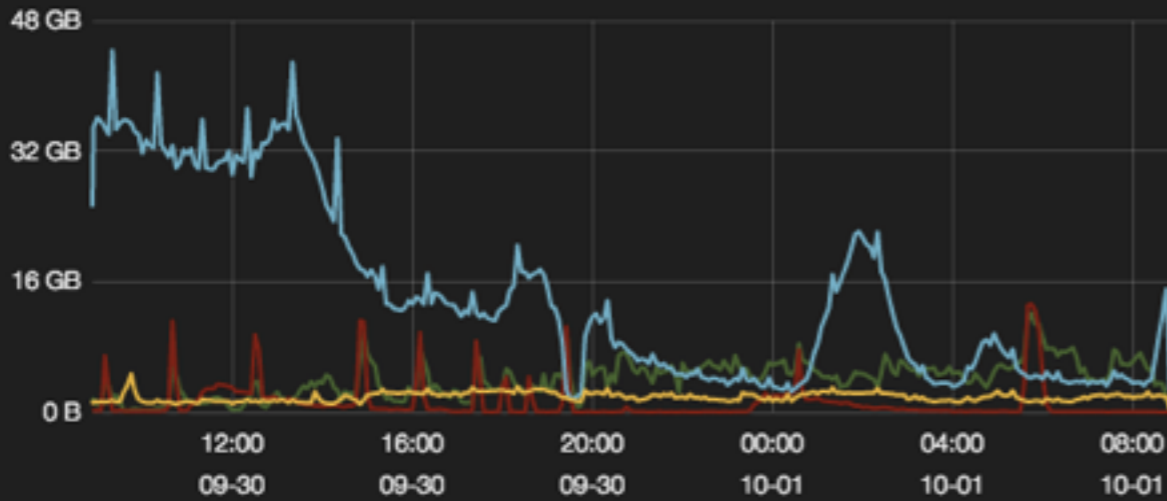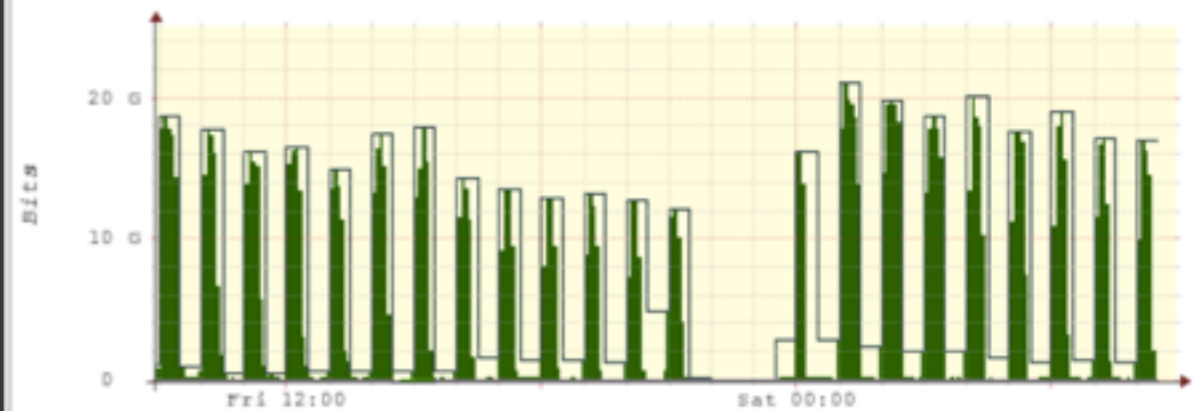**ALICE NETWORK USAGE (B/S)**

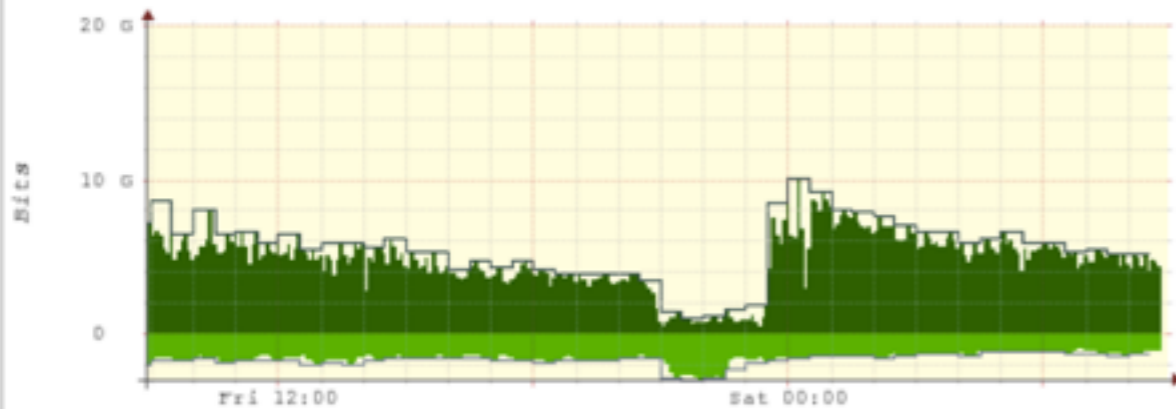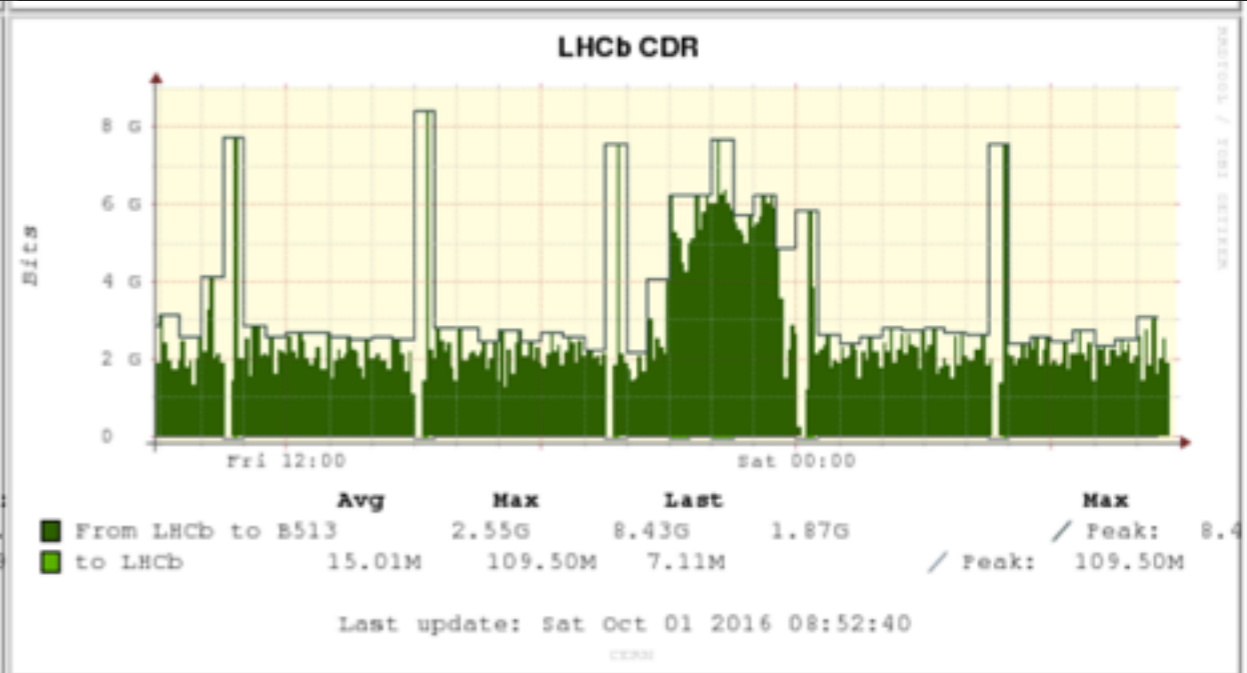● CASTOR Out ● CASTOR In ● EOS Out ● EOS In  per **5m** | (**181910** hits)

48 GB
32 GB
16 GB
0 B

12:00 / 09-30   16:00 / 09-30   20:00 / 09-30   00:00 / 10-01   04:00 / 10-01   08:00 / 10-01

---

**ATLAS NETWORK USAGE (B/S)**

● CASTOR Out ● CASTOR In ● EOS Out ● EOS In  per **5m** | (**207596** hits)

32 GB
24 GB
16 GB
8 GB
0 B

12:00 / 09-30   16:00 / 09-30   20:00 / 09-30   00:00 / 10-01   04:00 / 10-01   08:00 / 10-01

---

**CMS NETWORK USAGE (B/S)**

● CASTOR Out ● CASTOR In ● EOS Out ● EOS In  per **5m** | (**198784** hits)

20 GB
16 GB
12 GB
8 GB
4 GB
0 B

12:00 / 09-30   16:00 / 09-30   20:00 / 09-30   00:00 / 10-01   04:00 / 10-01   08:00 / 10-01

---

**LHCB NETWORK USAGE (B/S)**

● CASTOR Out ● CASTOR In ● EOS Out ● EOS In  per **5m** | (**76406** hits)

3 GB
2 GB
1 GB
0 B

12:00 / 09-30   16:00 / 09-30   20:00 / 09-30   00:00 / 10-01   04:00 / 10-01   08:00 / 10-01

CERN
IT-ST

# CERN made for LHC experiments needs, but…

EOS

Data processing

**Adaptable**
Catering with different uses

User Analysis

LHC Data Recording

Sync&Share

**Community** storage

Collaborate
CERNBOX Share
Offline work
Sync

SWAN (Jupyter)

CERNBox

| | | |
|---|---|---|
| **Users** | 6000 ($\Delta^{7d}$=60) | |
| **#files** | 110M ($\Delta^{7d}$=250K) | |
| **#dirs** | 14M | |
| **Quota** | 1TB/user | |
| **Used Space** | 240TB | |
| **Deployed Space** | 1.5PB | |

IT-ST

# Goals

Make data access easy
Make analysis simple
Facilitate Science

**My Laptop**
Small scale analysis
Test jobs

**AFS**
$home

/cvmfs

**batch/interactive services**
Large scale experiment processing
User extensive analysis

protocols
(xrdcp,rfio,*)

## Data Access
Main experiment data repositories

IT-ST

# 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

EOS CERNBOX does *"your files"* /cernbox/jdoe
EOS "e*xperiment"* does *"big data"* /eos/lhcb
Different QoS, different patterns, overlaps

IT-ST

# 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**

goldenrun052014

---

Submit jobs to lxbatch/wlcg 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 **synced** on his laptop as the
↳ if desired
jobs are being completed

---

Working on **final plots** on his **laptop** and Latex-ing the paper directly on **/mycernnbox/topmass/paper**

---

**Share** on-the-fly:
**n-tuples**
**Final plots**
**Publication**

via **/mycernbox**

---

EOS is the enabling technology binding all this
Multi QoS   Access patterns   Protocols   Redundancy

IT-ST

# Goals

Facilitate Science

**BLUE WATERS**
SUSTAINED PETASCALE COMPUTING

SIGN IN   NCSA

YOUR BLUE WATERS   ABOUT   SCIENCE AT BLUE WATERS   USING BLUE WATERS   EDUCATION & TRAINING   NEWS & EVENTS   HELP

Mapping Proton Quark Structure in Momentum and Coordinate Space using PetaByte Data-Sets from the COMPASS Experiment at CERN.

Interfaced **CERN storage services** with **Blue Waters NCSA** using **WLCG's FTS3** to manage the data workflow
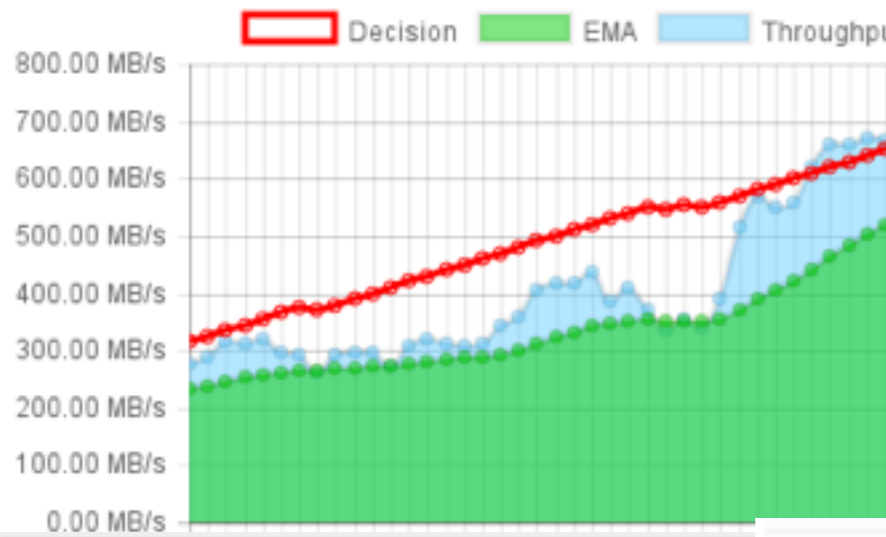
Open door for HPC environments to link with our HTC and Distributed Computing expertise

IT-ST

# Goals
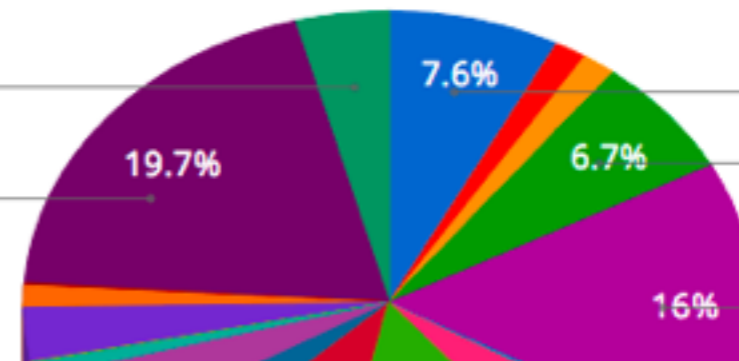Make data access easy
Make analysis simple
Facilitate Science

COMPASS · WLCG · BLUE WATERS — SUSTAINED PETASCALE COMPUTING

## Details for srm://castorpublic.cern.ch → gsiftp://ie15.ncsa.illinois.edu

Legend: Decision · EMA · Throughput

### CURRENT RUNNING JOBS BY SCIENCE AREA

Stellar Astronomy and Astrophysics — 19.7%
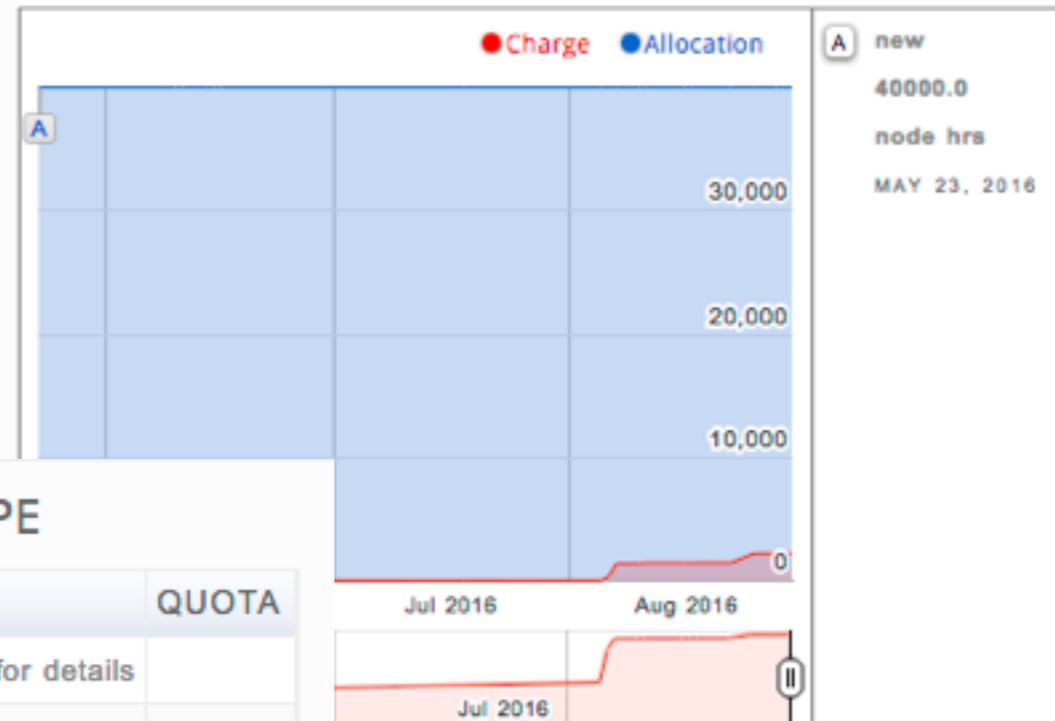Physics — 4.2%
Fluid, Particulate, and Hydraulic Systems

Astronomical Sciences — 7.6%
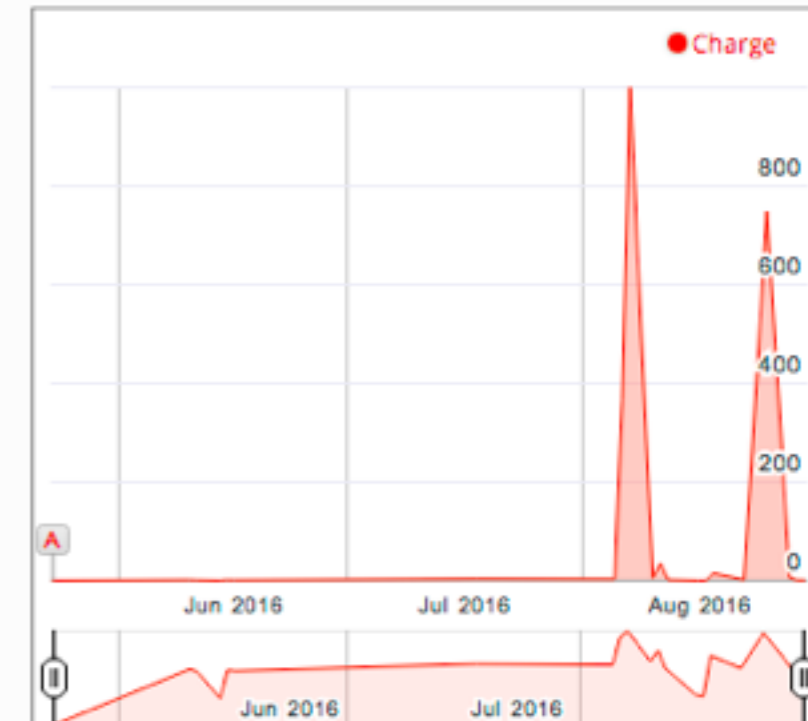Biological Sciences — 6.7%
Biophysics — 16%

### NETWORK (B/S)
● In ● Out per 5m | (694 hits)

### HISTORICAL CHARGED USAGE OVER TIME
(81 DAYS UNTIL EXPIRATION)
● Charge ● Allocation

A new 40000.0 node hrs MAY 23, 2016

### HISTORICAL DAILY CHARGED USAGE
● Charge

### PROJECT STORAGE USAGE BY TYPE

| TYPE ▼ | FILE COUNT | USAGE | QUOTA |
|---|---|---|---|
| Project Online | | N/A - See MOTD for details | |
| Project Nearline | 323,559 | 11.8 TiB | 50.0 TiB |
| Online Scratch | | 76.2 TiB | 500 TiB |

# Pushing boundaries

Raw data recording for Jan – August

# Pushing boundaries
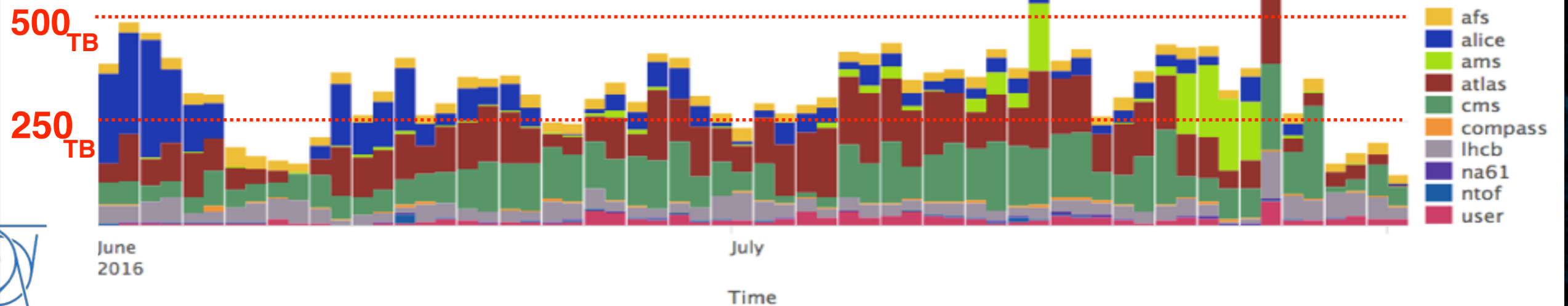
Transfered Data Amount per Virtual Organization for WRITE Requests

Raw data recording **per month** in 2016



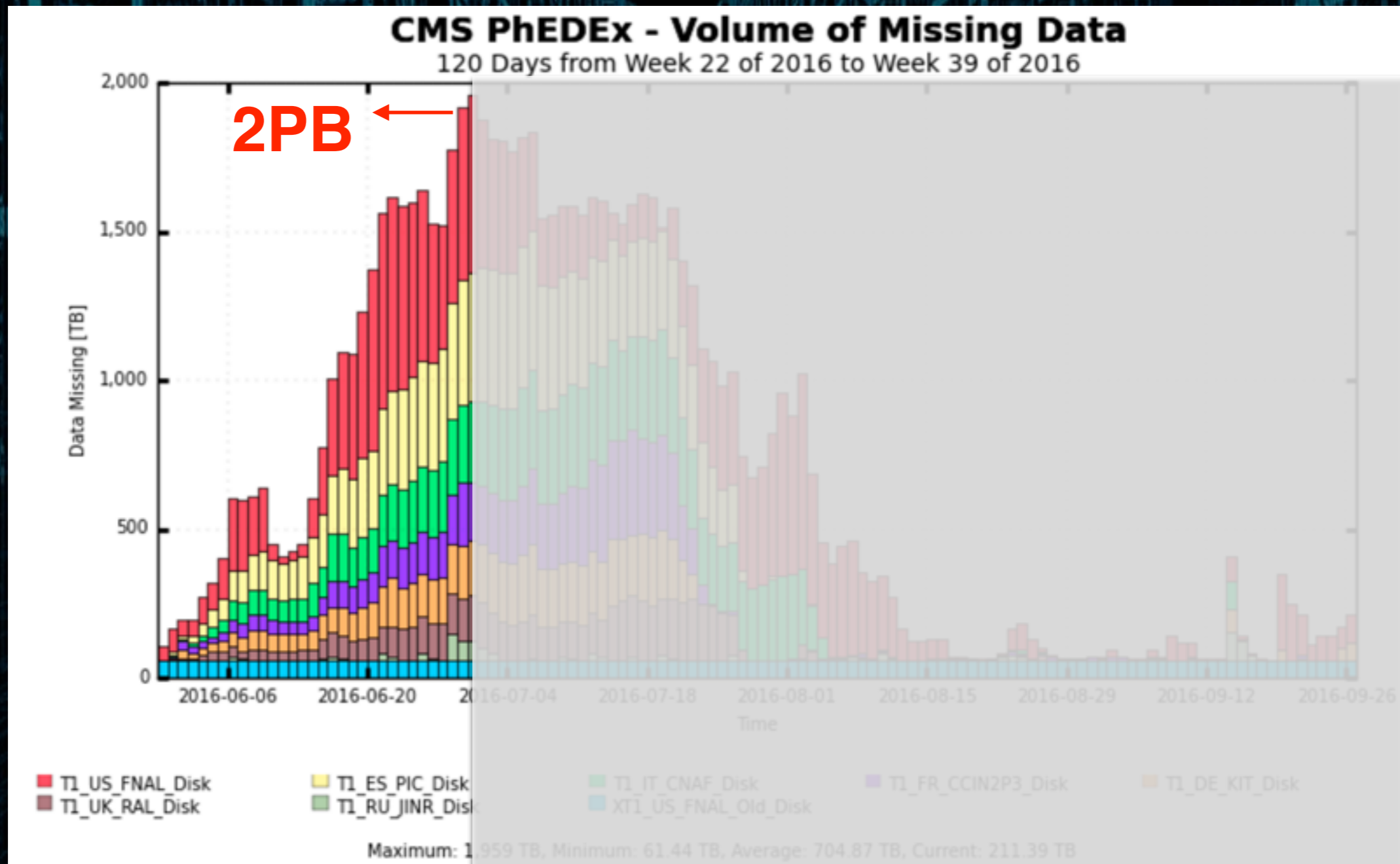Transfered Data Amount per Virtual Organization for WRITE Requests

Raw data recording for June and July, **per day**

21

# Pushing boundaries

LHC running at **full speed** before ICHEP 2016, **unprecedent** amounts of data

Systems **reaching limits** and not exporting data fast enough from Tier-0 to WLCG

REGI SERA **WOLFGANG PETERSEN**
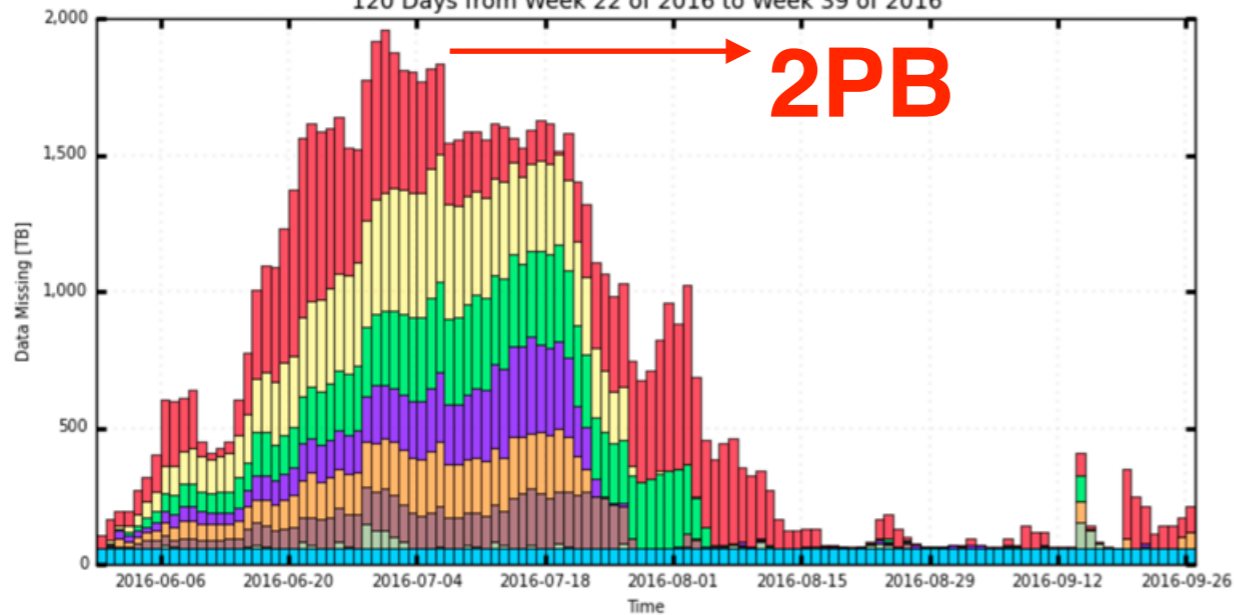
Das
**Boot**

RÉTISERA WOLFGANG PETERSEN

Das
Boot

# Pushing boundaries

CMS accumulated backlog    CMS
Tier-0 to T* data export speed



**2PB**

**7GB/s**

Issue **solved after one week**. Many investigations, many actions and many experts from Storage, Network, Experiments and FTS worked together to identify the issues and solve them.

# Goals
summary

Ensure a coherent development and operation of storage services at CERN for all aspects of physics data

Keep developing and operating Storage Services for Physics at the highest level

Communicating
Understanding
Delivering

Keep the ability to adapt and react fast

Problem/solution
Ask/Implement
In-house knowhow

Evaluate and investigate evolutions in technologies for better service/$

More for less
Operational costs
New applications

Envision new models on data mananagement and analysis

Sync&Share
SWAN
LHC@myPC

27

CERN

IT-ST

More on **CERNBOX**:

**CERNBox: the data hub for data analysis** (J.Moscicki) - Poster session

More on **SWAN**:

**SWAN: a Service for Web-Based Data Analysis in the Cloud**(D.Piparo/E.Tejedor)12th/Oct-11:45 (SierraB)

More on distributed EOS distributed:

**Global EOS: exploring the 300-ms-latency region** (L.Mascetti) - Poster session

More on **Cern Tape Archive (CTA)**:

**An efficient, modular and simple tape archiving solution for LHC Run-3** (S. Murray) - Poster Session

**From Physics to industry: EOS Outside HEP** (XE)

IT-ST