# Storage Chamäleons

2016 Hepix Spring
DESY, Zeuthen

CERN

Xavier Espinal
on behalf of IT/ST

25th Anniversary

HEPiX

Zeuthen Spring 2016

Reliable

Fast Processing
DAQ Feedback loop

DAQ to CC
8GB/s+4xReco [ALICE]

Hot files

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

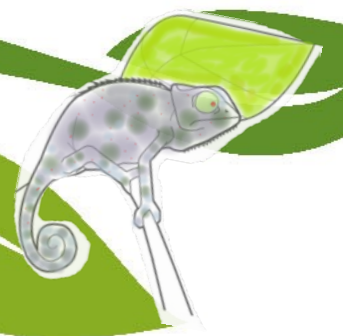High throughout to tape
350+MB/s/drive - 12GB/s [Pb-Pb]

back-up

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

Consistent

∞

Few fast streams
CDR 2x40Gbps

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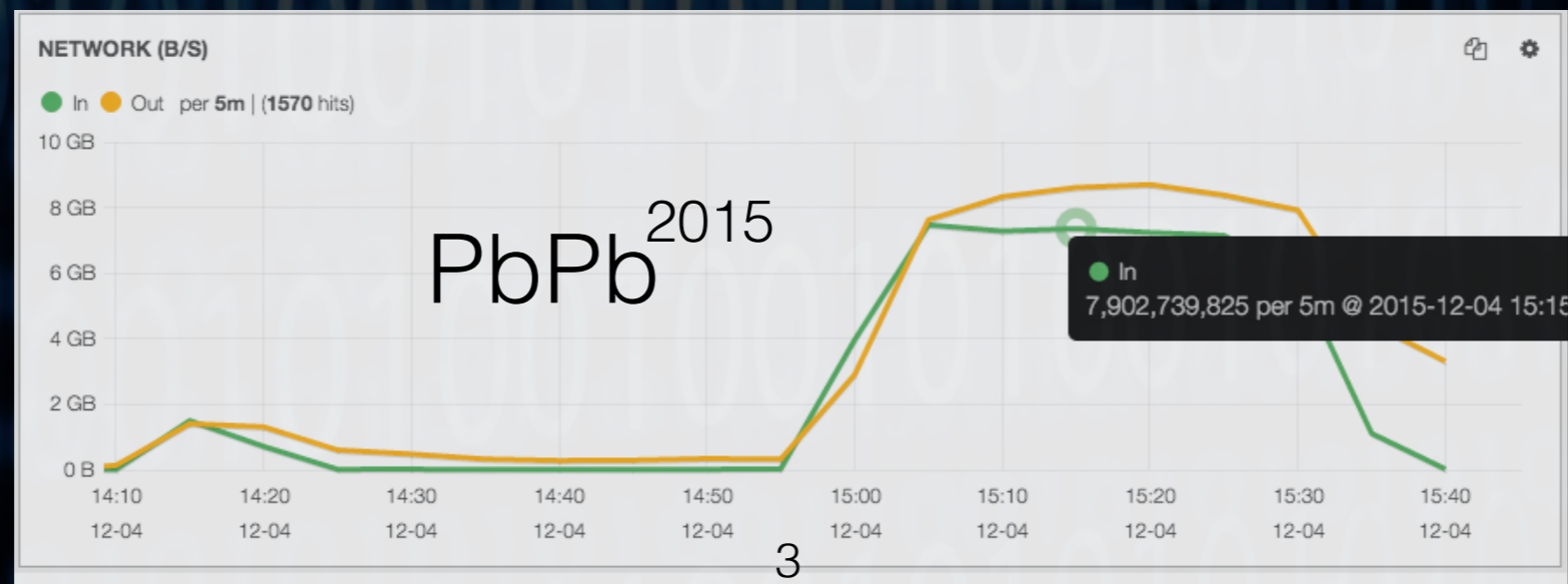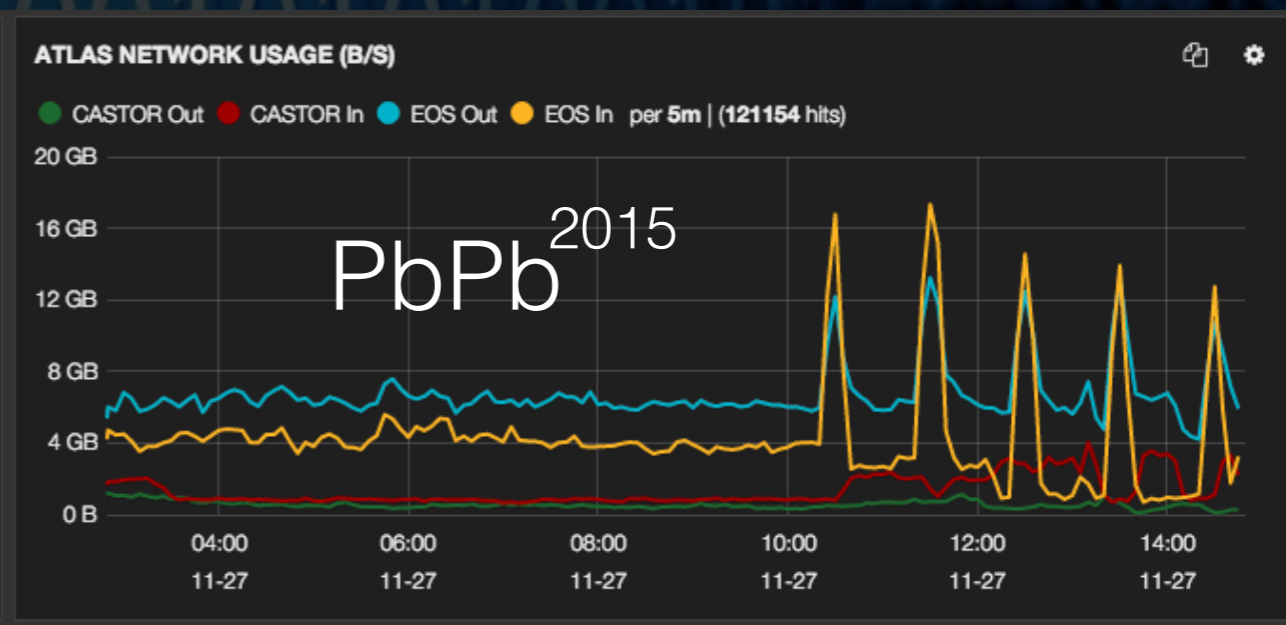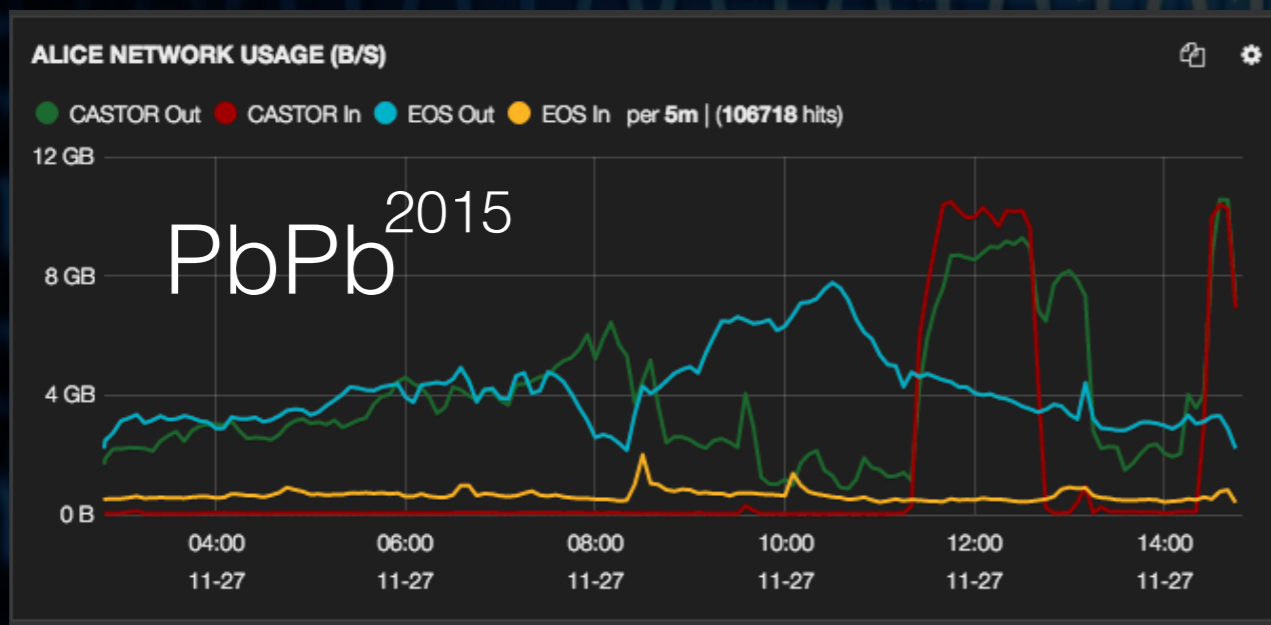Unexpected usage Catalogue=Namespace

disk and gc?

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

Endpoint Mounts
ie. /atlas in the WNs

CERN
IT-ST

2

# Core Systems

**Evolved to**
## Tape oriented system

**Key feature**
## Per stream speed

Biggest scientific-repo worldwide 138PB and +500M files
High throughput from DAQ, high throughput to tape

Moved from Raid1 to Raid60 (100MB/s to >350MB/s$^{per\ stream}$)
Evaluating common disk layer
Tape policies, per experiment/user/group resources

3

**CASTOR**
CERN Advanced STORage manager

**Evolved to**
## Tape oriented system

Biggest scientific-repo worldwide 138PB and +500M files
High throughput from DAQ, high throughput to tape

**Key feature**
## Per stream speed

Moved from Raid1 to Raid60 (100MB/s to >350MB/s$^{per\ stream}$)
Evaluating common disk layer
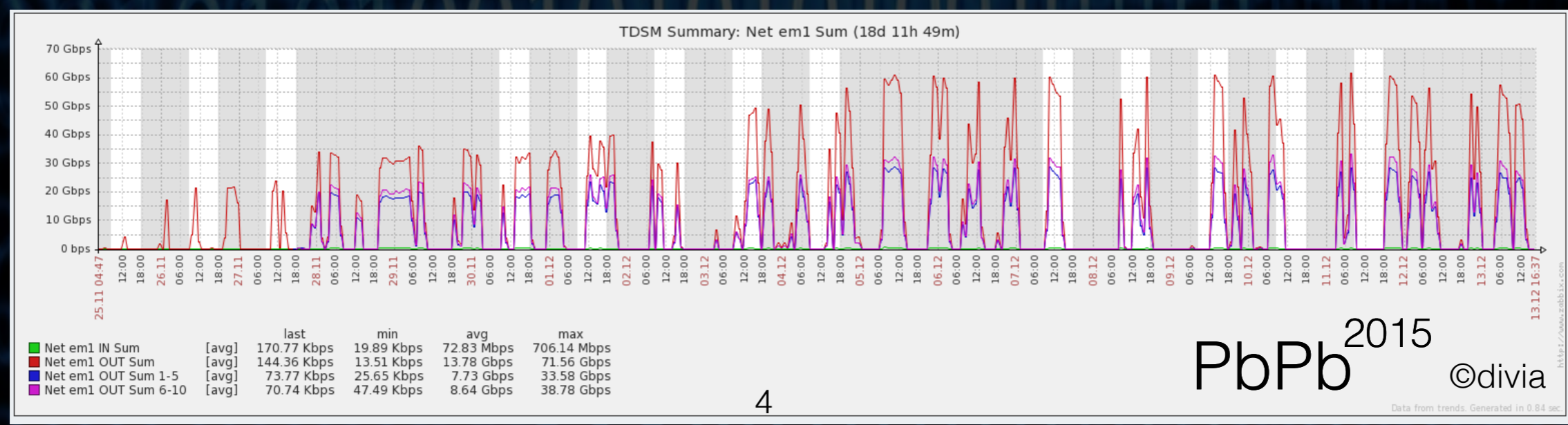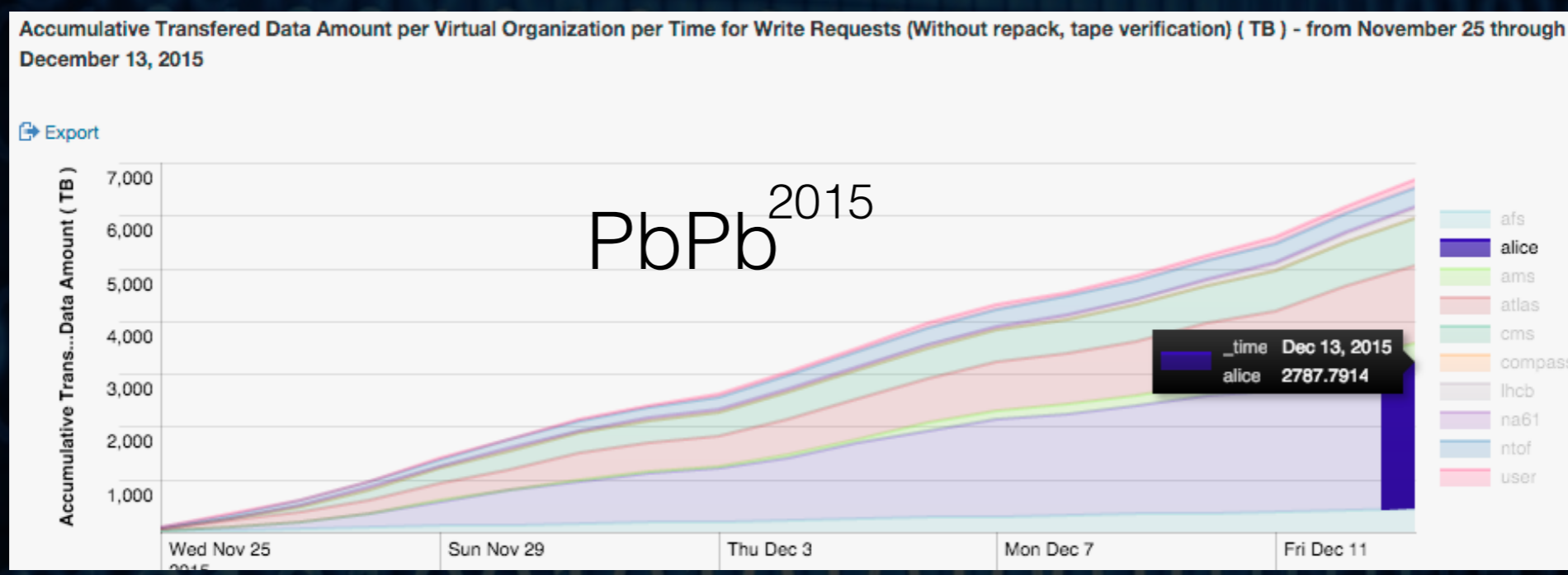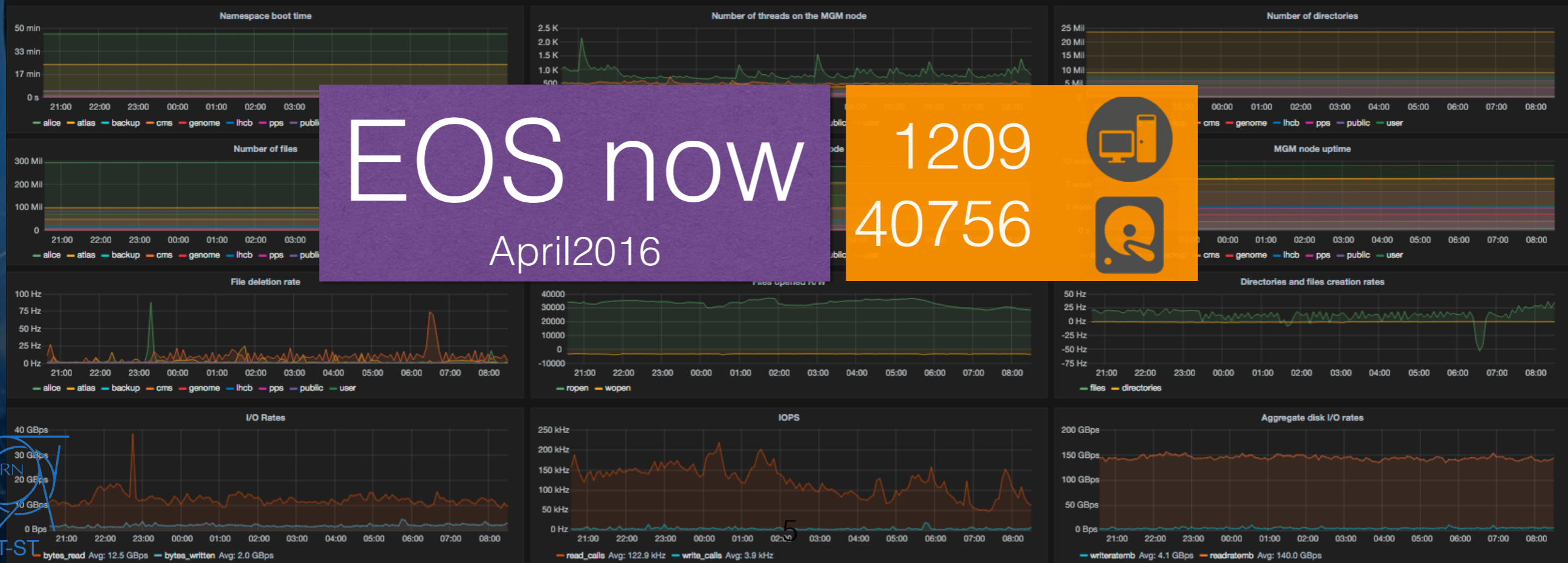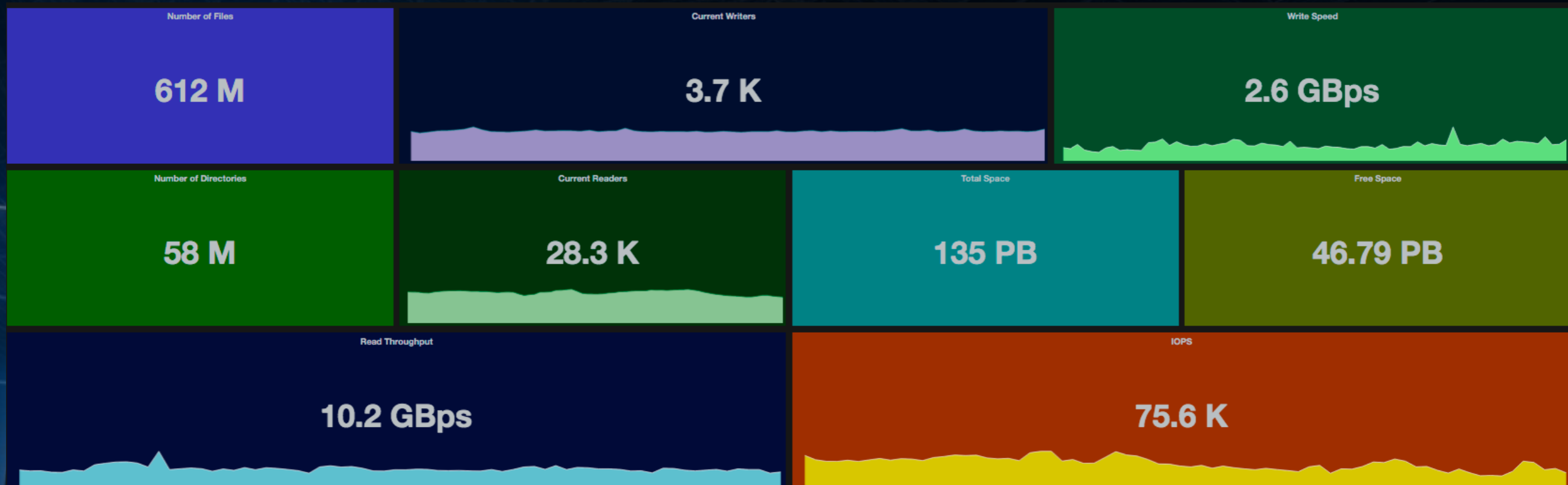Tape policies, per experiment/user/group resources



Accumulative Transfered Data Amount per Virtual Organization per Time for Write Requests (Without repack, tape verification) ( TB ) - from November 25 through December 13, 2015

PbPb$^{2015}$

| | | |
| afs |
| alice |
| ams |
| atlas |
| cms |
| compass |
| lhcb |
| na61 |
| ntof |
| user |

_time  Dec 13, 2015
alice  2787.7914



TDSM Summary: Net em1 Sum (18d 11h 49m)

| | | last | min | avg | max |
|---|---|---|---|---|---|
| Net em1 IN Sum | [avg] | 170.77 Kbps | 19.89 Kbps | 72.83 Mbps | 706.14 Mbps |
| Net em1 OUT Sum | [avg] | 144.36 Kbps | 13.51 Kbps | 13.78 Gbps | 71.56 Gbps |
| Net em1 OUT Sum 1-5 | [avg] | 73.77 Kbps | 25.65 Kbps | 7.73 Gbps | 33.58 Gbps |
| Net em1 OUT Sum 6-10 | [avg] | 70.74 Kbps | 47.49 Kbps | 8.64 Gbps | 38.78 Gbps |

PbPb$^{2015}$

©divia

CERN
IT-ST

4

# made@CERN

Designed and tailored for experiments needs

Experts in-house: Adapt when required
Re-design if needed

Being used outside: Fermilab, Russia-T1, EsNET, Aarnet
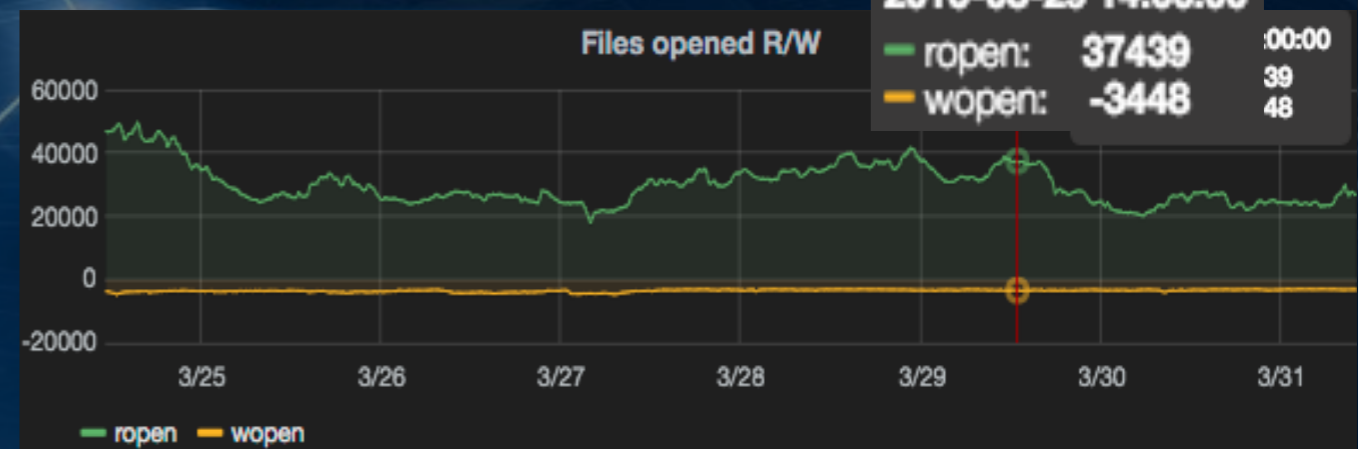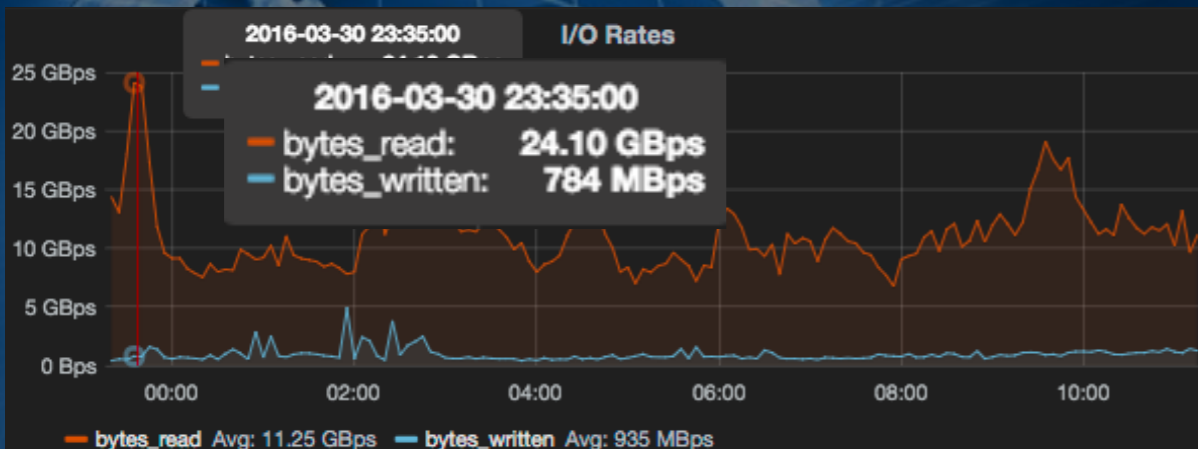Openlab/COMTRADE JRC, Univ. Vienna, INFN-Trieste

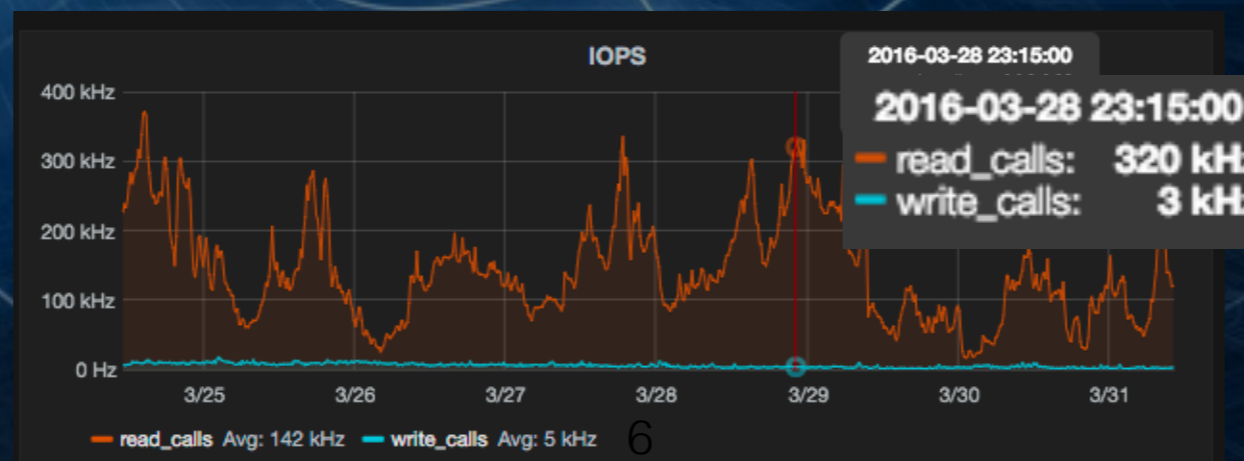## Adaptable
Catering with different uses

CDR
Data processing
User Analysis
CernBOX

## Community data
DmaaS (iJupyter)
Sync share

CERNBox



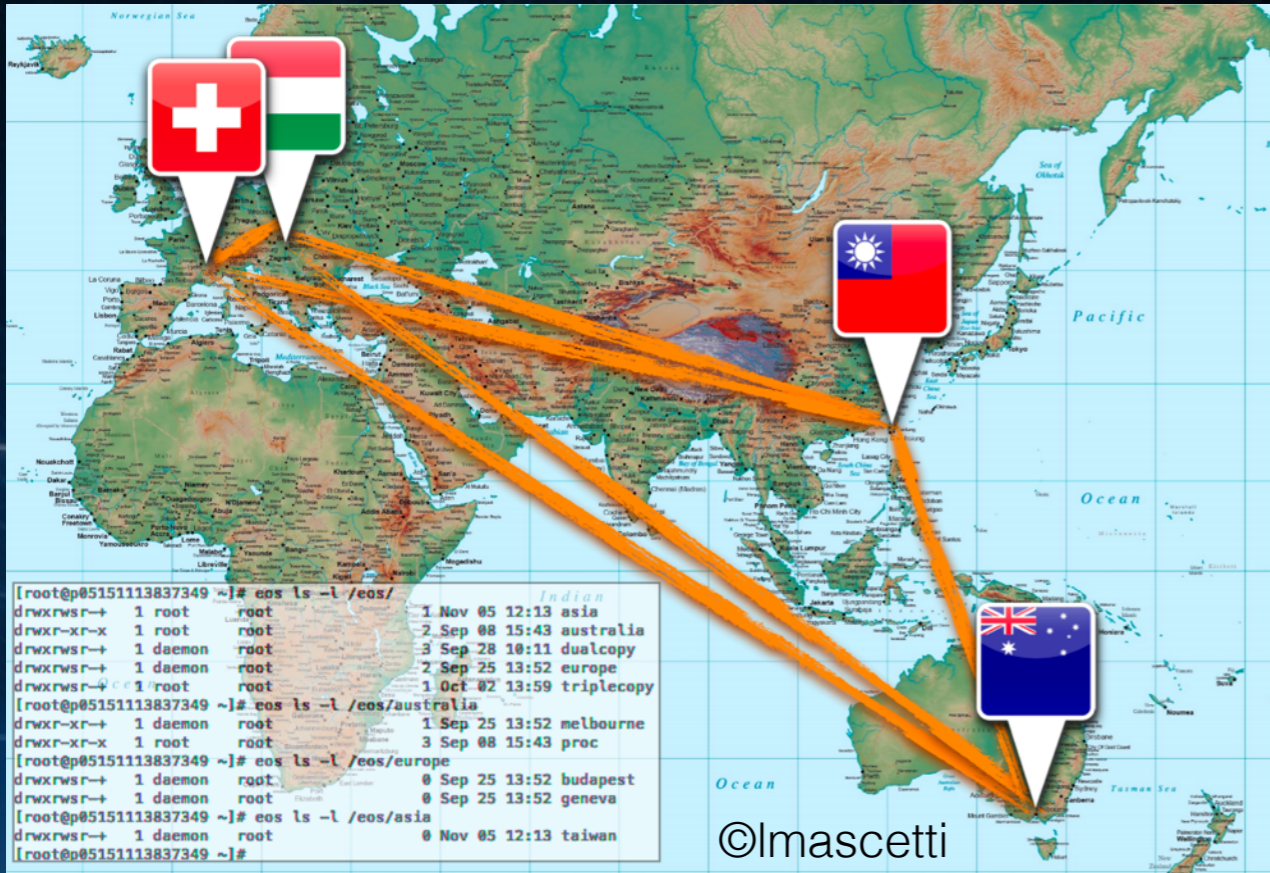## Future Shared FS ?more later…



6

©lmascetti

# Can go distributed, can be shared and synced

©lmascetti

Clients delocalization
Multi-site deployment

Used from 22 ms to 300 ms

```
[root@p05151113837349 ~]# eos ls -l /eos/
drwxrwsr-+   1 root      root              1 Nov 05 12:13 asia
drwxr-xr-x   1 root      root              2 Sep 08 15:43 australia
drwxrwsr-+   1 daemon    root              3 Sep 28 10:11 dualcopy
drwxrwsr-+   1 daemon    root              2 Sep 25 13:52 europe
drwxrwsr-+   1 root      root              1 Oct 02 13:59 triplecopy
[root@p05151113837349 ~]# eos ls -l /eos/australia
drwxr-xr-+   1 daemon    root              1 Sep 25 13:52 melbourne
drwxr-xr-x   1 root      root              3 Sep 08 15:43 proc
[root@p05151113837349 ~]# eos ls -l /eos/europe
drwxrwsr-+   1 daemon    root              0 Sep 25 13:52 budapest
drwxrwsr-+   1 daemon    root              0 Sep 25 13:52 geneva
[root@p05151113837349 ~]# eos ls -l /eos/asia
drwxrwsr-+   1 daemon    root              0 Nov 05 12:13 taiwan
[root@p05151113837349 ~]#
```

## Community data
### DmaaS (iJupyter) share Sync

CERN IT-ST

**CERNBox**

| Users | 4719 |
|---|---|
| # files | 70 Million |
| # dirs | 9 Million |
| Quota | 1TB/user |
| Used Space | 125 TB |
| Deployed Space | 1.5 PB |

20%  60%  20%

©lmascetti

# Can go distributed, can be shared and synced

Clients delocalization

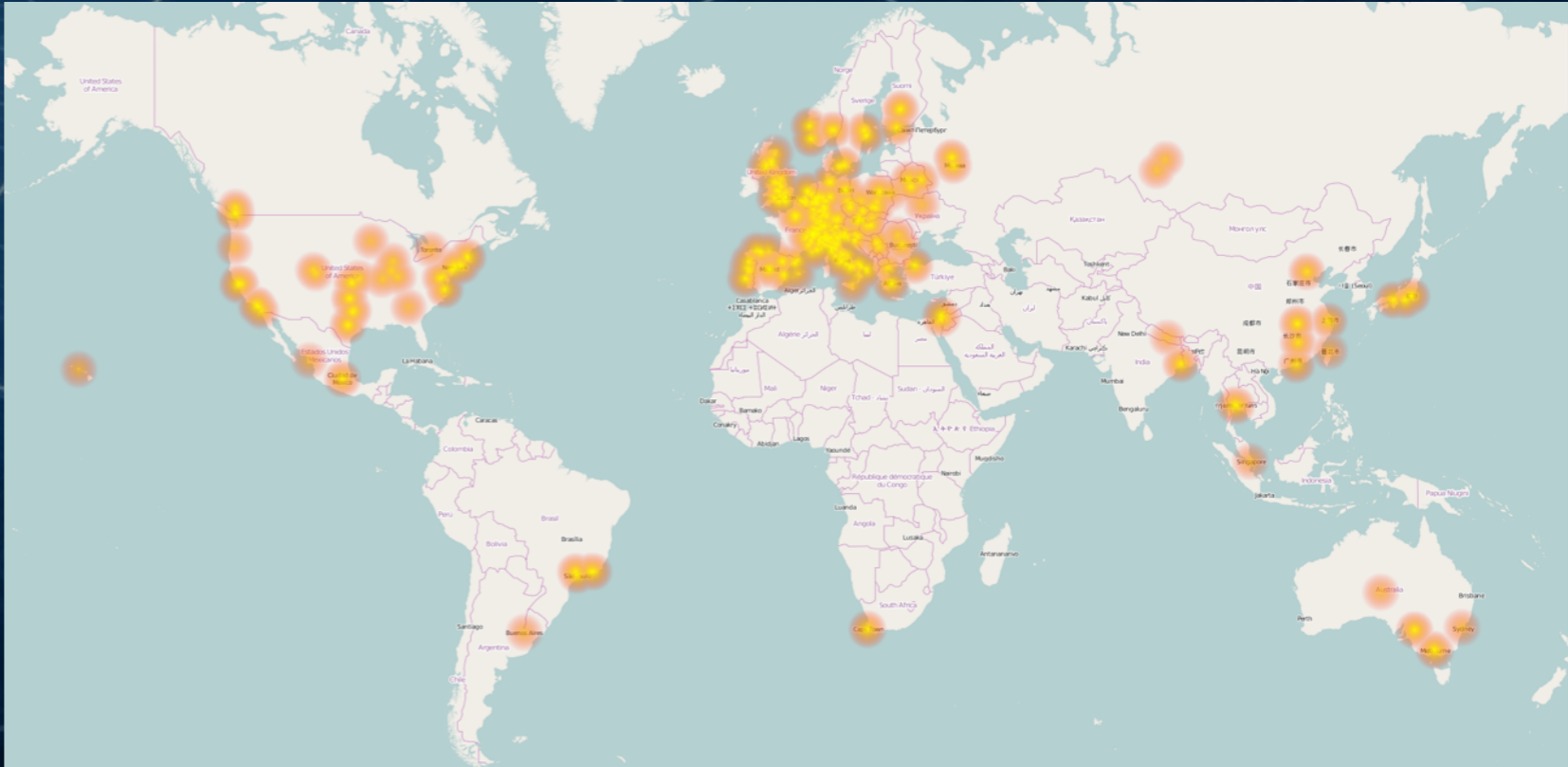Multi-site deployment

Used from
22 ms
to
300 ms

```
[root@p05151113837349 ~]# eos ls -l /eos/
drwxrwsr-+   1 root     root          1 Nov 05 12:13 asia
drwxr-xr-x   1 root     root          2 Sep 08 15:43 australia
drwxrwsr-+   1 daemon   root          3 Sep 28 10:11 dualcopy
drwxrwsr-+   1 daemon   root          2 Sep 25 13:52 europe
drwxrwsr-+   1 root     root          1 Oct 02 13:59 triplecopy
[root@p05151113837349 ~]# eos ls -l /eos/australia
drwxr-xr-+   1 daemon   root          1 Sep 25 13:52 melbourne
drwxr-xr-x   1 root     root          3 Sep 08 15:43 proc
[root@p05151113837349 ~]# eos ls -l /eos/europe
drwxrwsr-+   1 daemon   root          0 Sep 25 13:52 budapest
drwxrwsr-+   1 daemon   root          0 Sep 25 13:52 geneva
[root@p05151113837349 ~]# eos ls -l /eos/asia
drwxrwsr-+   1 daemon   root          0 Nov 05 12:13 taiwan
[root@p05151113837349 ~]#
```
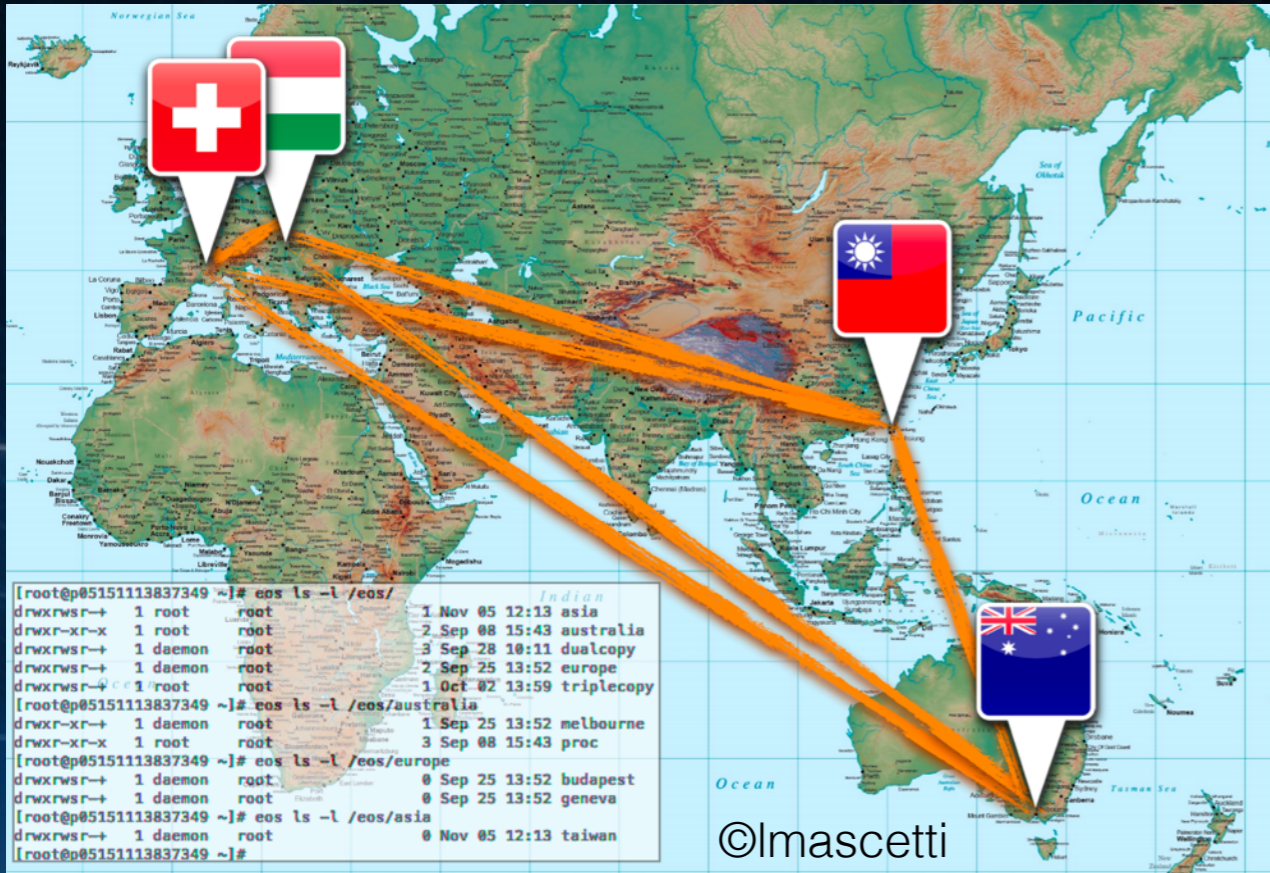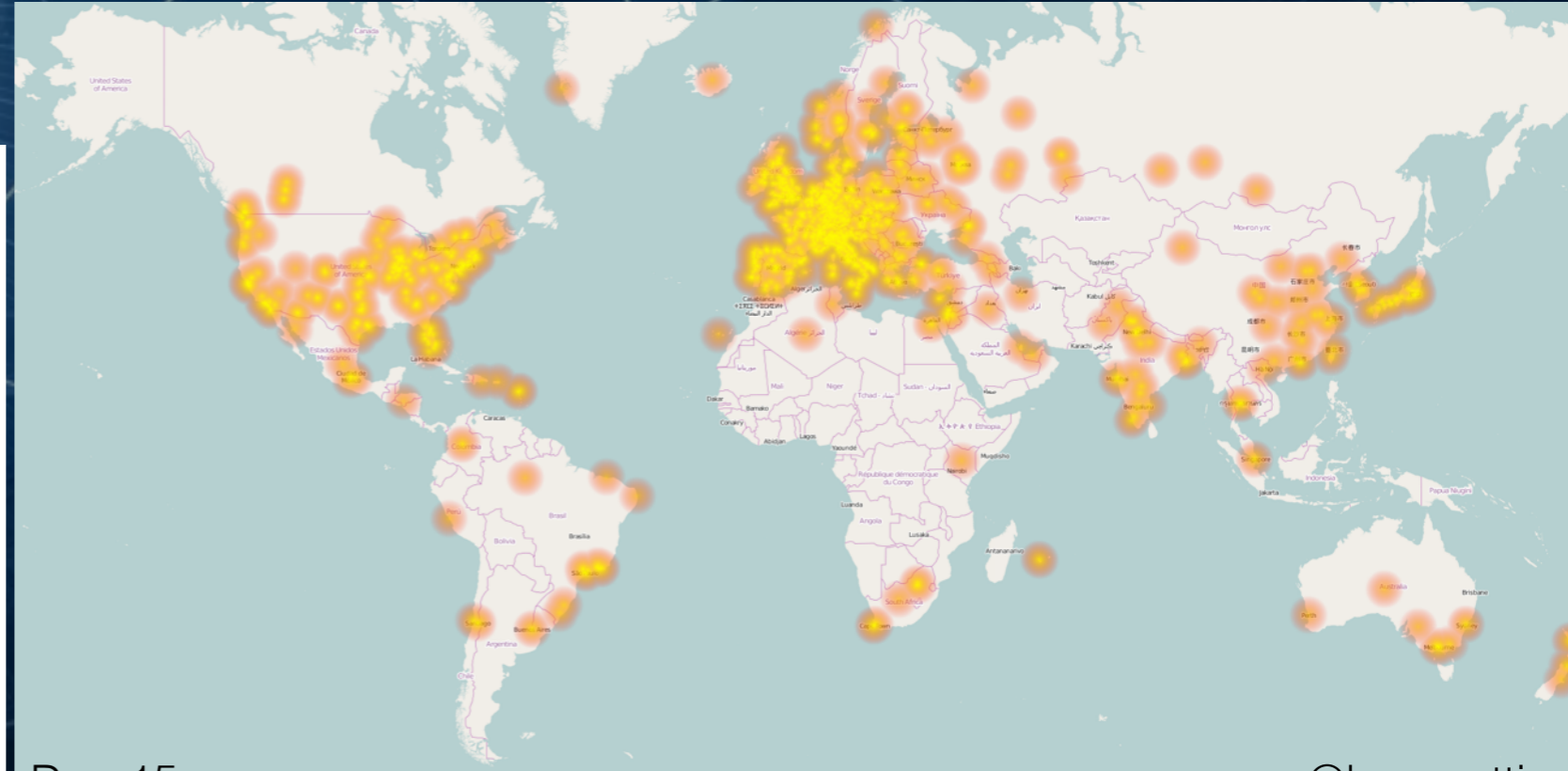
©lmascetti

Community data
DmaaS (iJupyter) share Sync

CERN IT-ST



## CERNBox

| Users | 4719 |
|---|---|
| # files | 70 Million |
| # dirs | 9 Million |
| Quota | 1TB/user |
| Used Space | 125 TB |
| Deployed Space | 1.5 PB |

20%  60%  20%

Dec-15    8    ©lmascetti

Clients delocalization

Multi-site deployment

Used from 22 ms to 300 ms

©lmascetti

```
[root@p05151113837349 ~]# eos ls -l /eos/
drwxrwsr-+  1 root    root        1 Nov 05 12:13 asia
drwxr-xr-x  1 root    root        2 Sep 08 15:43 australia
drwxrwsr-+  1 daemon  root        3 Sep 28 10:11 dualcopy
drwxrwsr-+  1 daemon  root        2 Sep 25 13:52 europe
drwxrwsr-+  1 root    root        1 Oct 02 13:59 triplecopy
[root@p05151113837349 ~]# eos ls -l /eos/australia
drwxr-xr-+  1 daemon  root        1 Sep 25 13:52 melbourne
drwxr-xr-x  1 root    root        3 Sep 08 15:43 proc
[root@p05151113837349 ~]# eos ls -l /eos/europe
drwxrwsr-+  1 daemon  root        0 Sep 25 13:52 budapest
drwxrwsr-+  1 daemon  root        0 Sep 25 13:52 geneva
[root@p05151113837349 ~]# eos ls -l /eos/asia
drwxrwsr-+  1 daemon  root        0 Nov 05 12:13 taiwan
[root@p05151113837349 ~]#
```
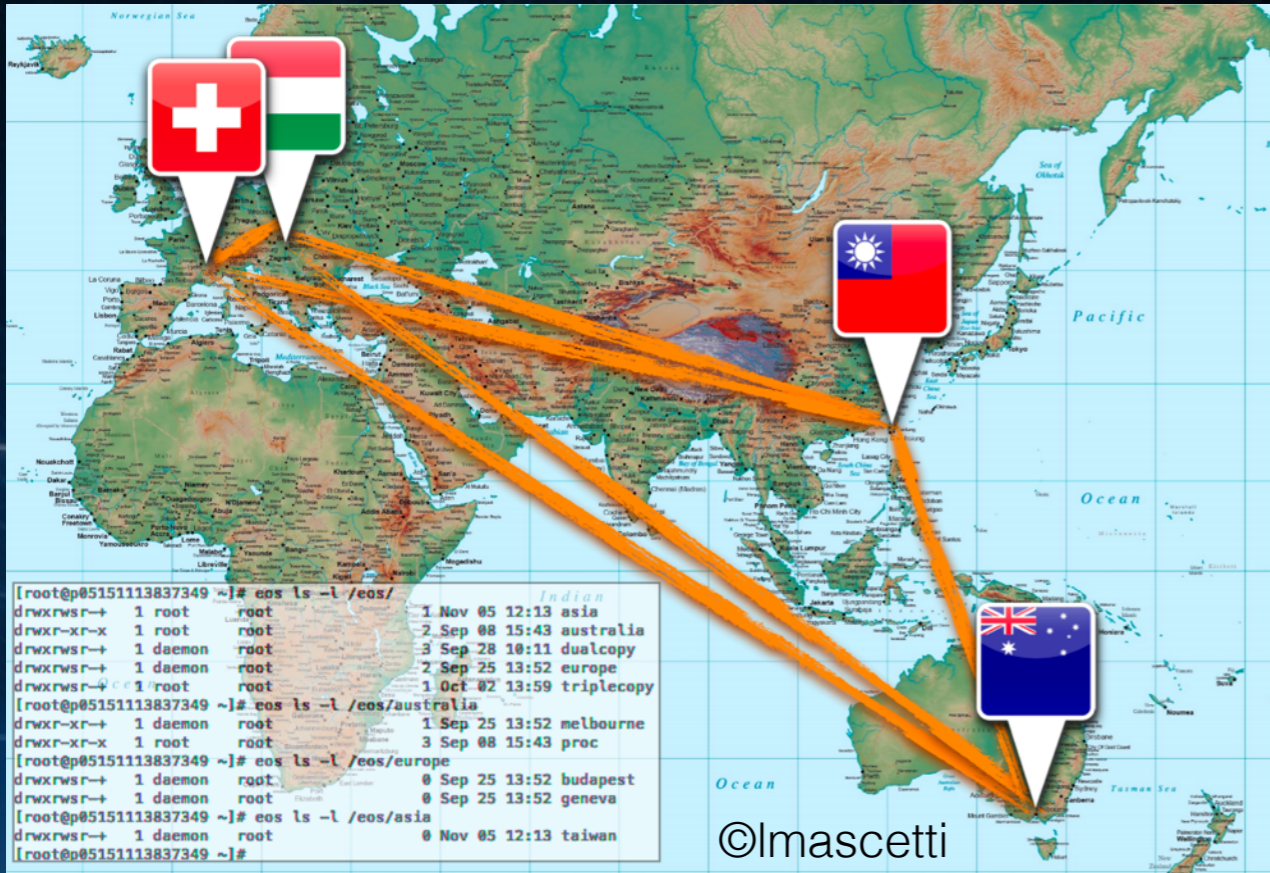
Community data

DmaaS (iJupyter) share Sync

CERN

IT-ST

## CERNBox

| | |
|---|---|
| Users | 4719 |
| # files | 70 Million |
| # dirs | 9 Million |
| Quota | 1TB/user |
| Used Space | 125 TB |
| Deployed Space | 1.5 PB |

20%  60%  20%

# Embedded ROOT Viewer

©dpiparo

ROOT
Data Analysis Framework

The viewer is based on the ROOT data analysis framework developed at CERN by PH-SFT.

Integration done by CERNBox team.

9

# BLOCK STORAGE

ceph

Openstack VM
Cinder Volumes
S3

RADOS FS

File stripper
CASTOR backend
Under evaluation
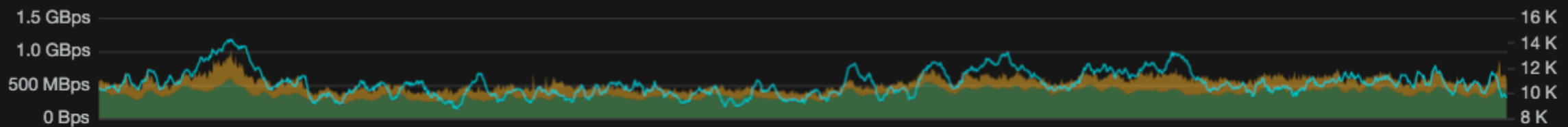
Large contribution

Community

Code
development
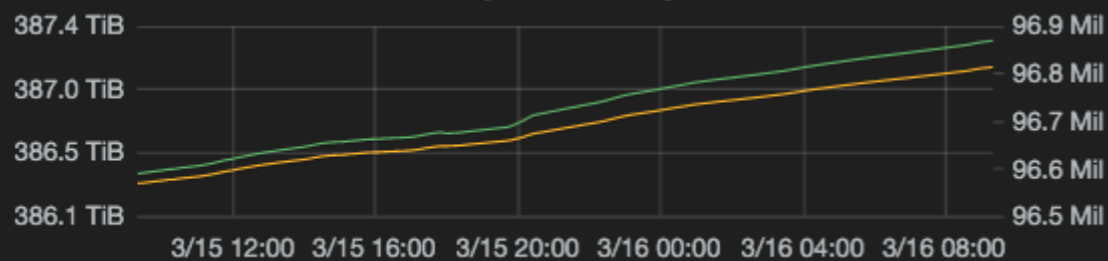CERN-IT/ST

Largest Cluster 30PB

Deployed to date

40k OSDs

Multi-site

In production

3PB@wigner
1PB@meyrin



**2870** images

**2037** volumes

1.5 GBps
1.0 GBps
500 MBps
0 Bps

16 K
14 K
12 K
10 K
8 K

3/15 12:00   3/15 16:00   3/15 20:00   3/16 00:00   3/16 04:00   3/16 08:00

**Used space and objects**

387.4 TiB
387.0 TiB
386.5 TiB
386.1 TiB

96.9 Mil
96.8 Mil
96.7 Mil
96.6 Mil
96.5 Mil

3/15 12:00   3/15 16:00   3/15 20:00   3/16 00:00   3/16 04:00   3/16 08:00

**Used space derivative**

60 MBps
40 MBps
20 MBps
0 Bps
-20 MBps
-40 MBps

3/15 12:00   3/15 16:00   3/15 20:00   3/16 00:00   3/16 04:00   3/16 08:00

CERN
IT-ST
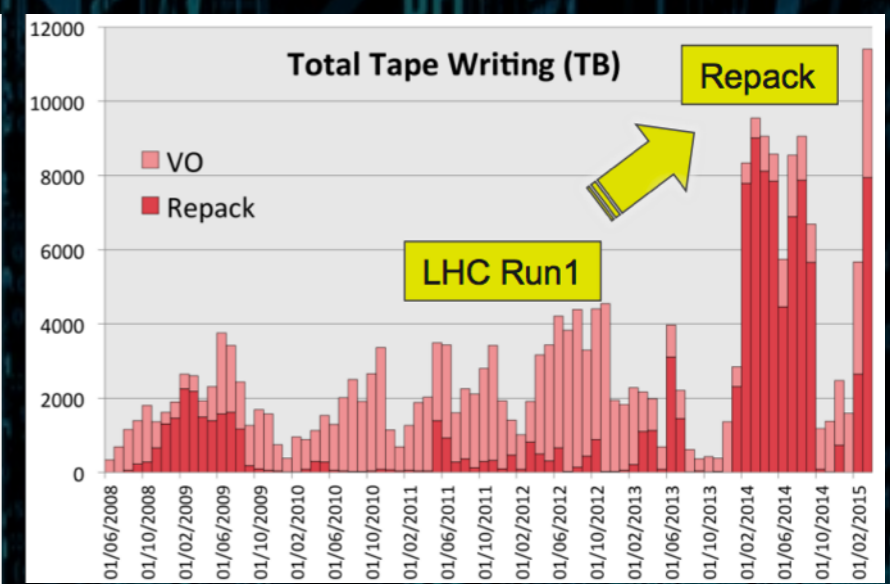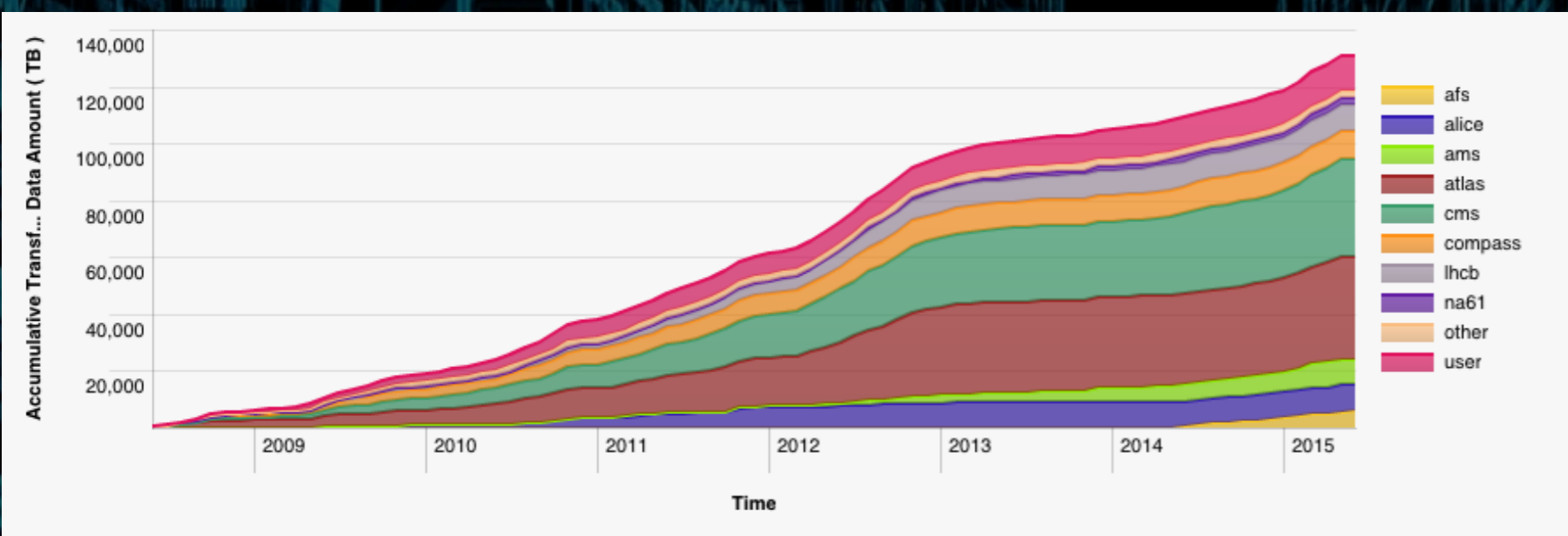
10

# CERN Tape Archive

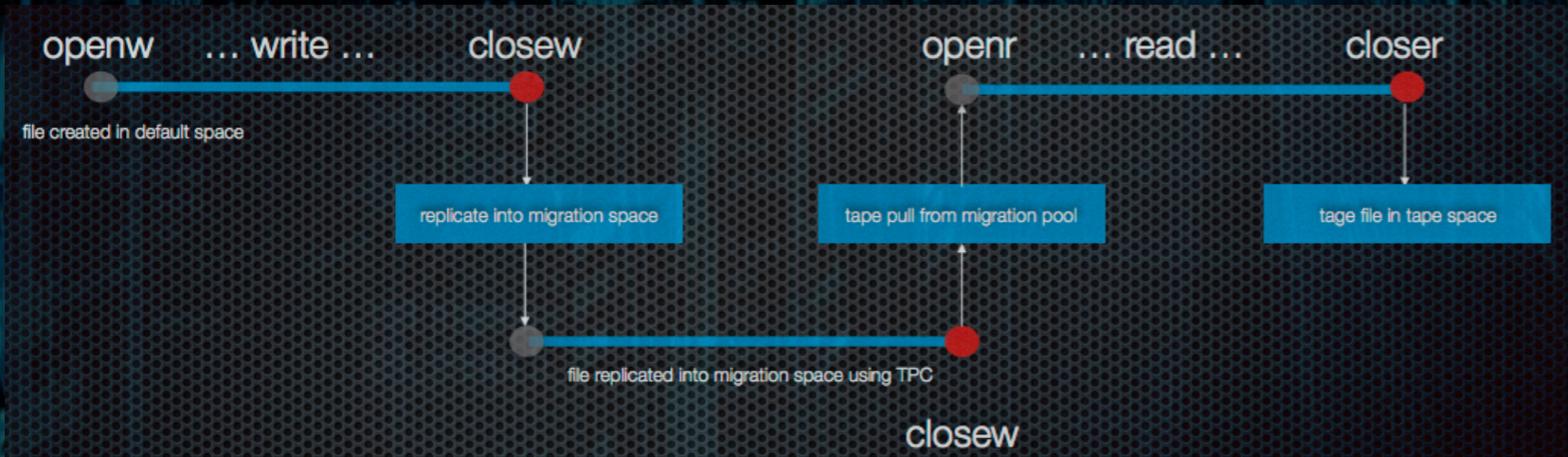Technology driven: new medias brings ↑density ↑speed

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



Cross-system workflows

**+30% / yr tape capacity per $**
**+20%/yr I/O increase**

©apeters

# CERN Tape Archive

Technology driven: new medias brings ↑density ↑speed

Towards a pluggable tape backend (EOS)
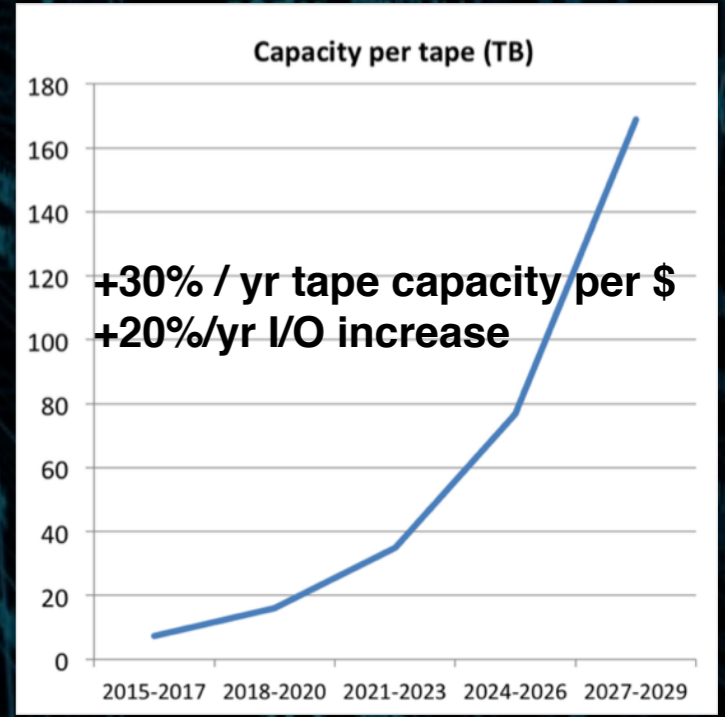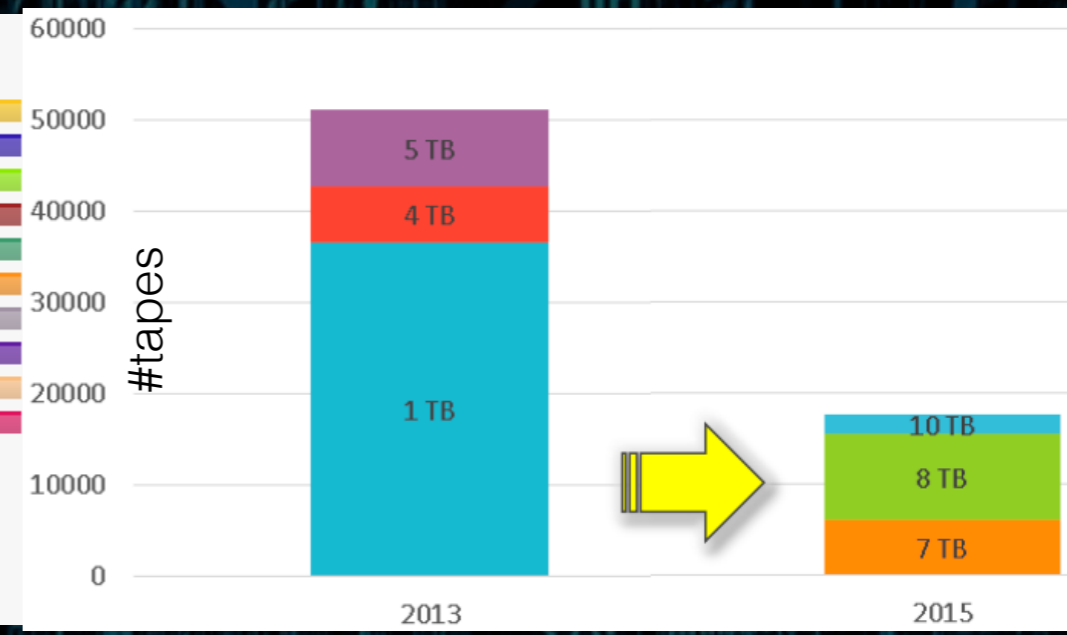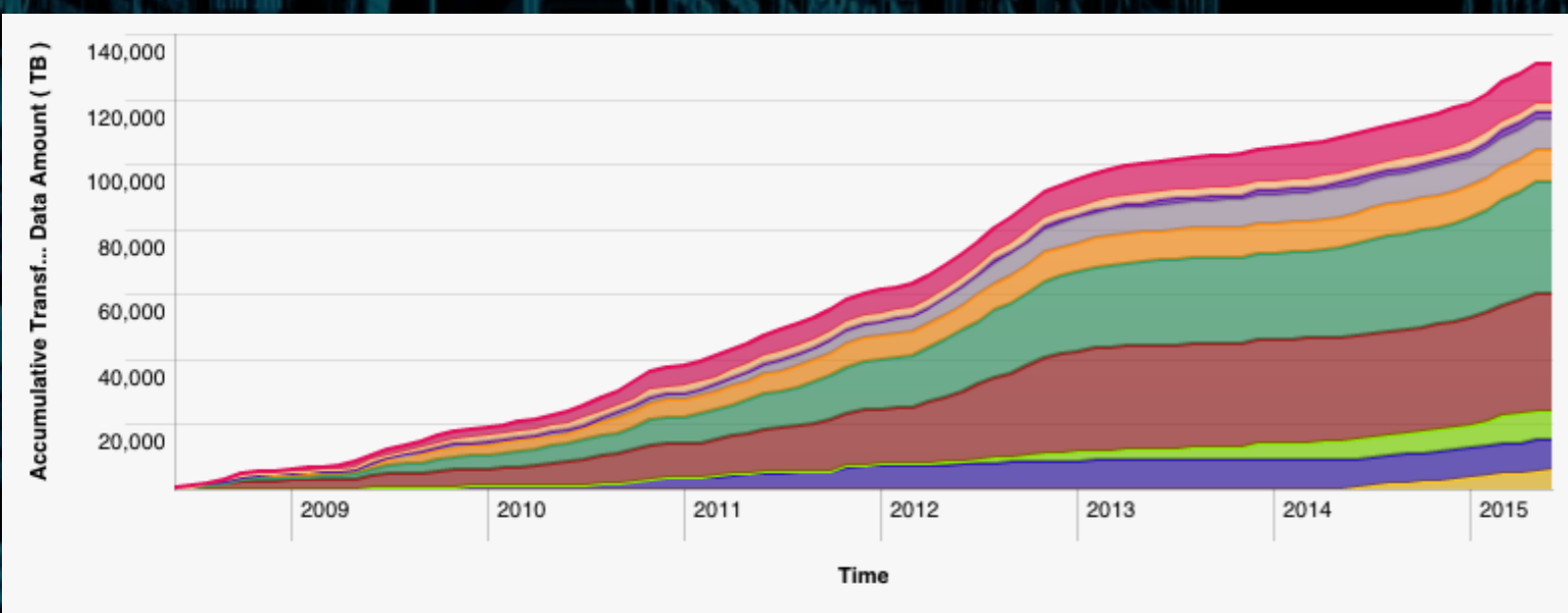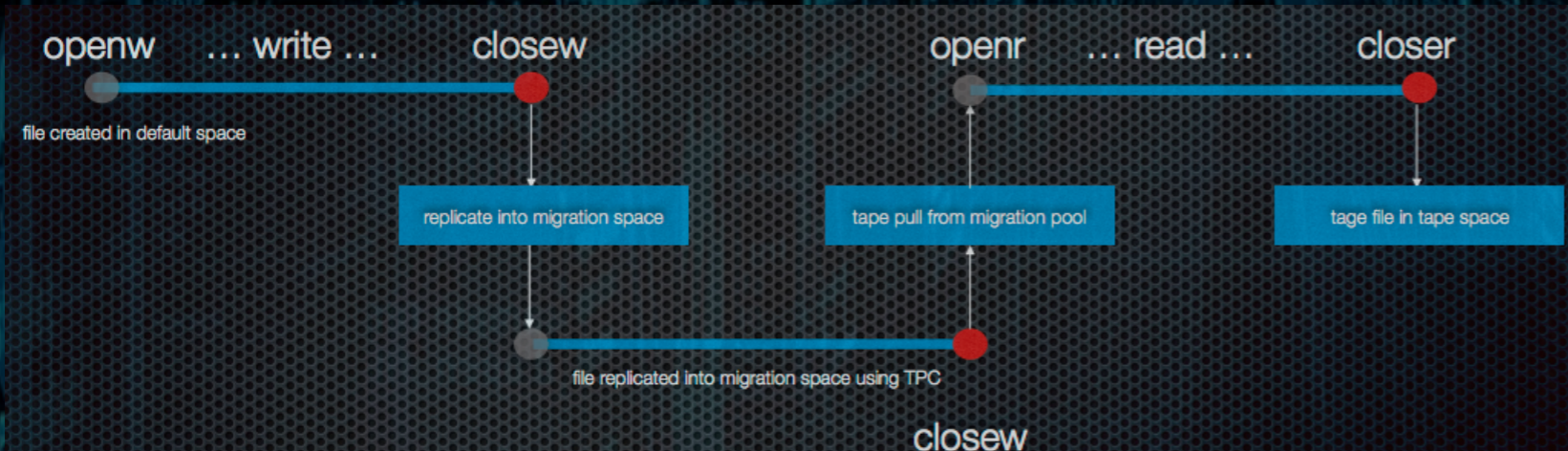
Cold by definition: hight throughput, high latency

Core Systems



Tape best technology for data repositories: TCO media power density and resilient/reliable

very large disk caches nowadays



## Cross-system workflows

©apeters

12

+30% / yr tape capacity per $
+20%/yr I/O increase

CERN IT-ST

Goals    Make data access easy

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

# Goals  Make data access easy

## My Laptop
Small scale analysis
Test jobs

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

## Data Access
Main experiment data repositories

### Mounts

AFS  $home

squids
/cvmfs/athena

fuse
/mycernbox

fuse
/eos/atlas

CERN

IT-ST

# Goals

# Make data access easy

**My Laptop**
Small scale analysis
Test jobs

**Mounts**

squids
/cvmfs/athena

fuse
/mycernbox

fuse
/eos/atlas

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

## Data Access
Main experiment data repositories

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

Backup
15

# Goals

# Make Analysis Simple

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

He/she submits 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**

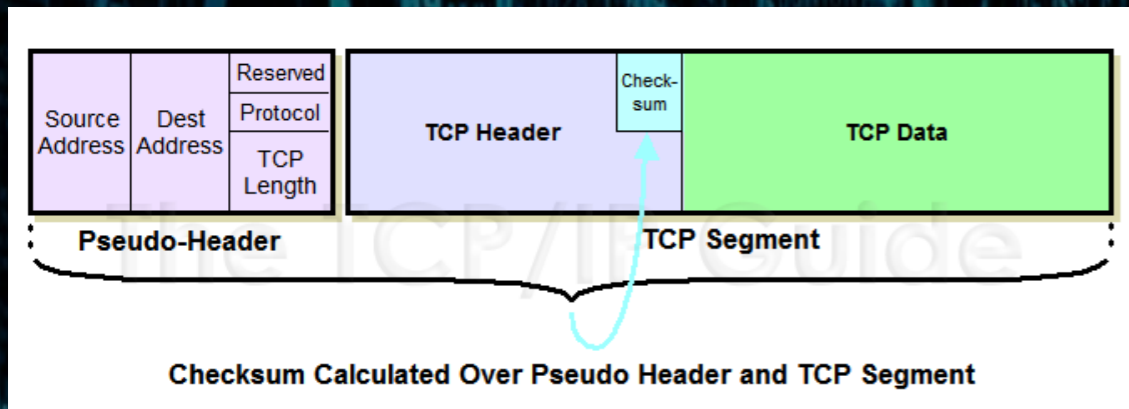is the enabling technology binding all this
Multi QoS   Access patterns   Protocols   Redundancy

IT-ST

# Before concluding... D4T4 C0RRUP710N   Adler32?

Found corrupted files after an incident caused by a a faulty router



Checksum Calculated Over Pseudo Header and TCP Segment

Brute force retransmits eventually went through due to TCP checksum collision

adler32 effectively gives less than 32-bits of verificaton power on file transfers affected by TCP checksoum collissions

(atlas) 7.8K files, 11 TB were staged in from T1 tapes for a md5 comparison, 17 corruptions with a colliding adler32 and a different md5 checksum —> 0.22% silent corruption rate

Observed courrption probability of 0.22% —> $2^{-9}$ (wrt. $2^{-32}$)

http://cern.ch/go/wr8j

# 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

LHC@myPC
Sync&Share
DmaaS

CERN

IT-ST