# CERN/IT Storage Evolution



Xavier Espinal

on behalf of IT/ST

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
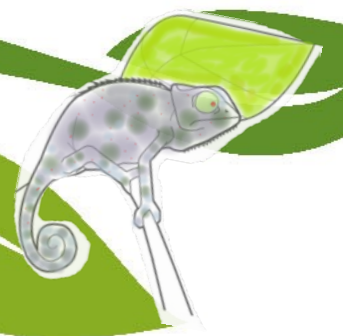350+MB/s/drive - 12GB/s Pb-Pb

back-up

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

Consistent

$\infty$

Few fast streams
CDR 2x40Gbps

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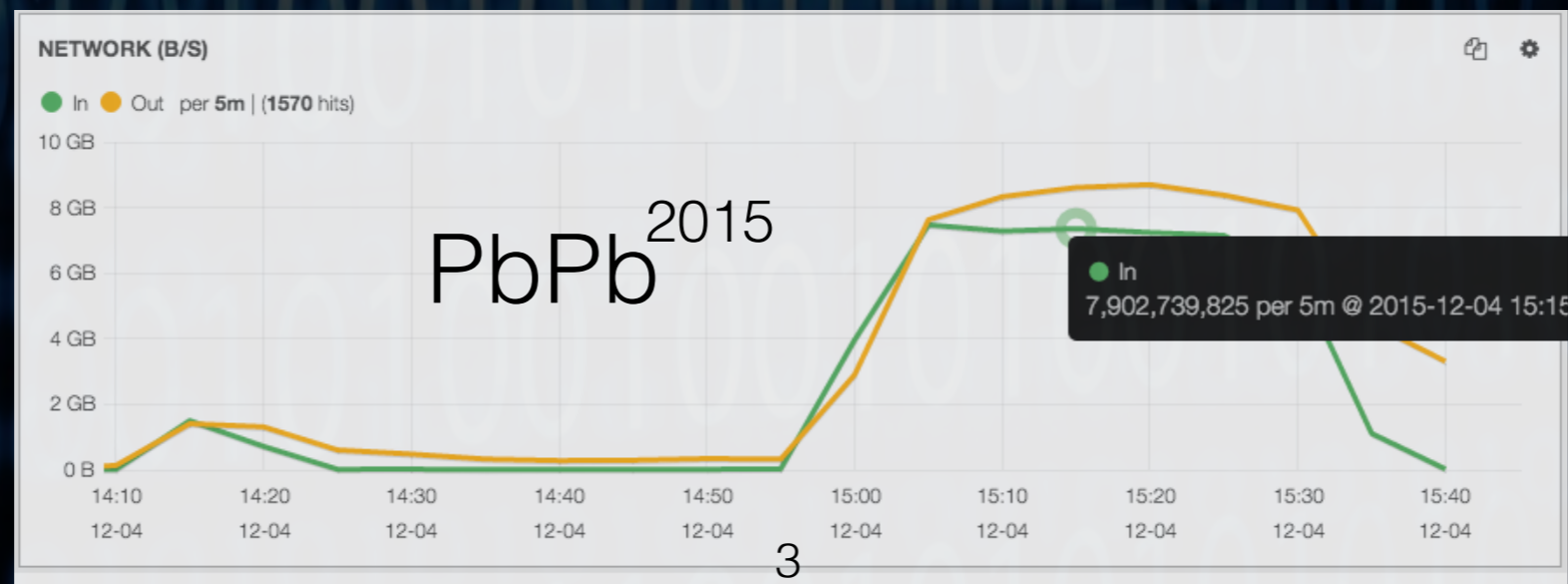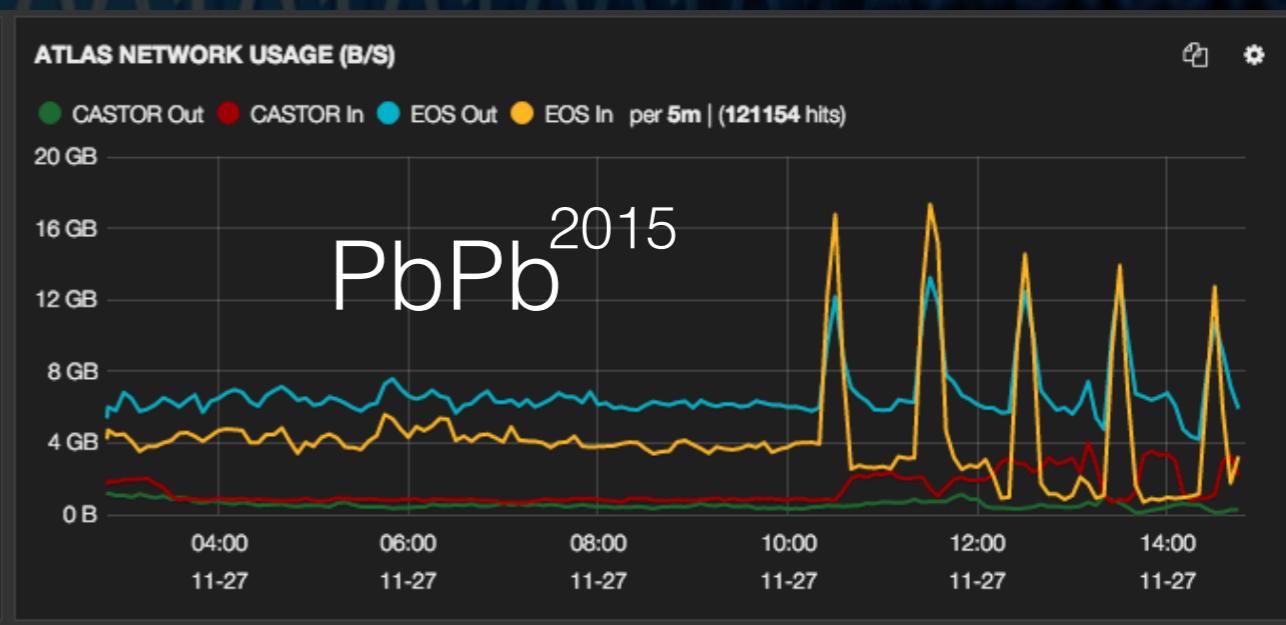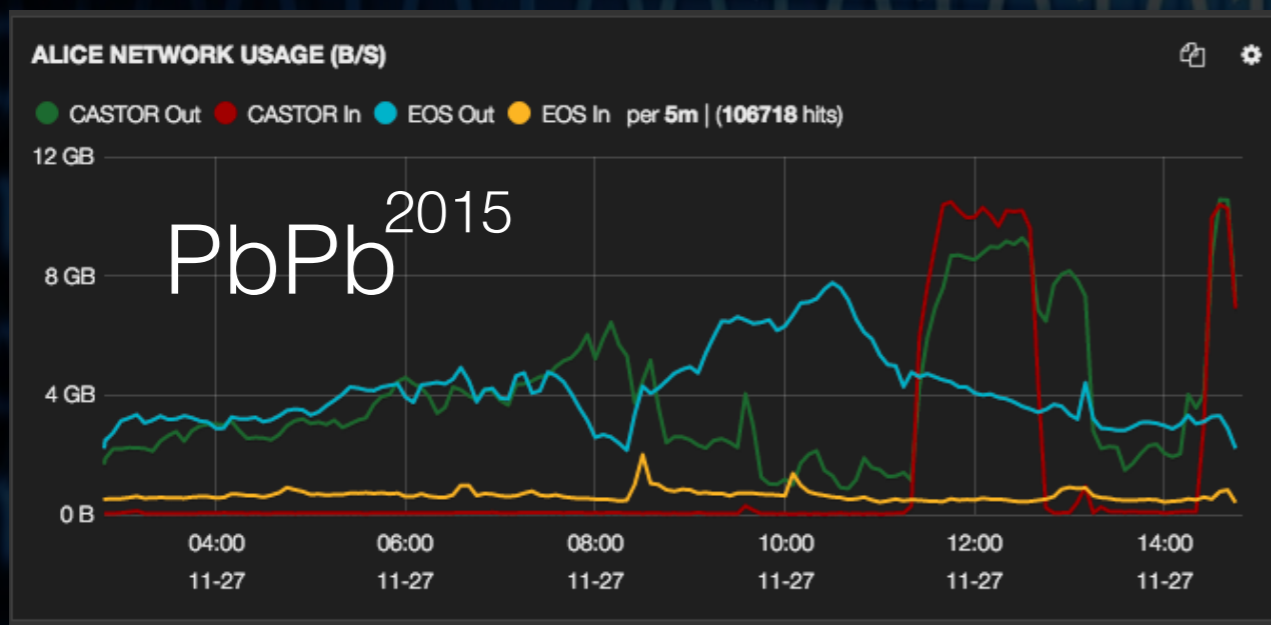Unexpected usage Catalogue=Namespace

disk and gc?

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

Endpoint Mounts
ie. /atlas in the WNs

CERN
IT-ST

2

**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$^{\text{per stream}}$)
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
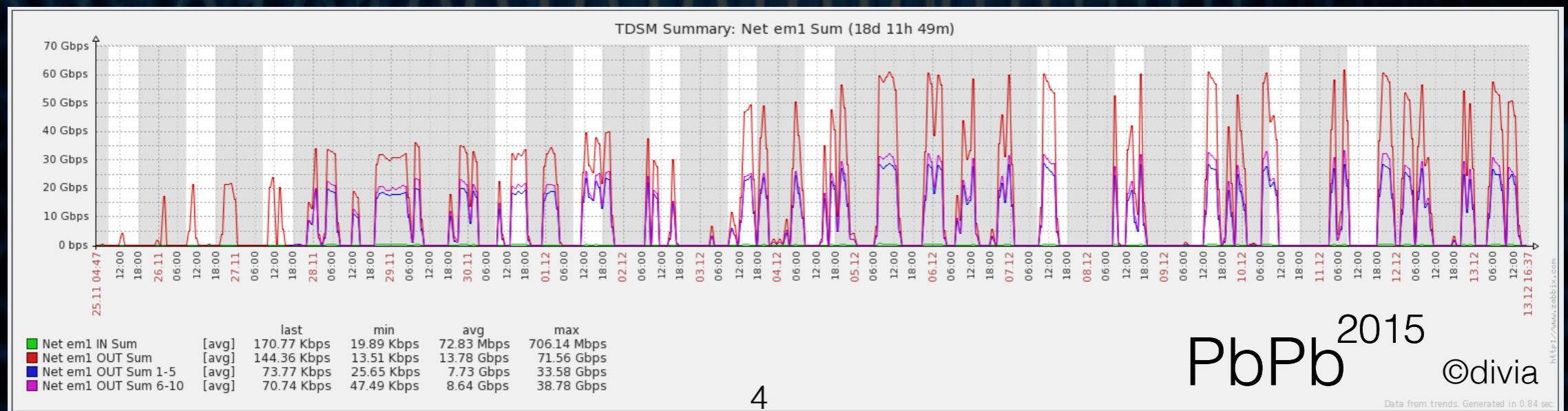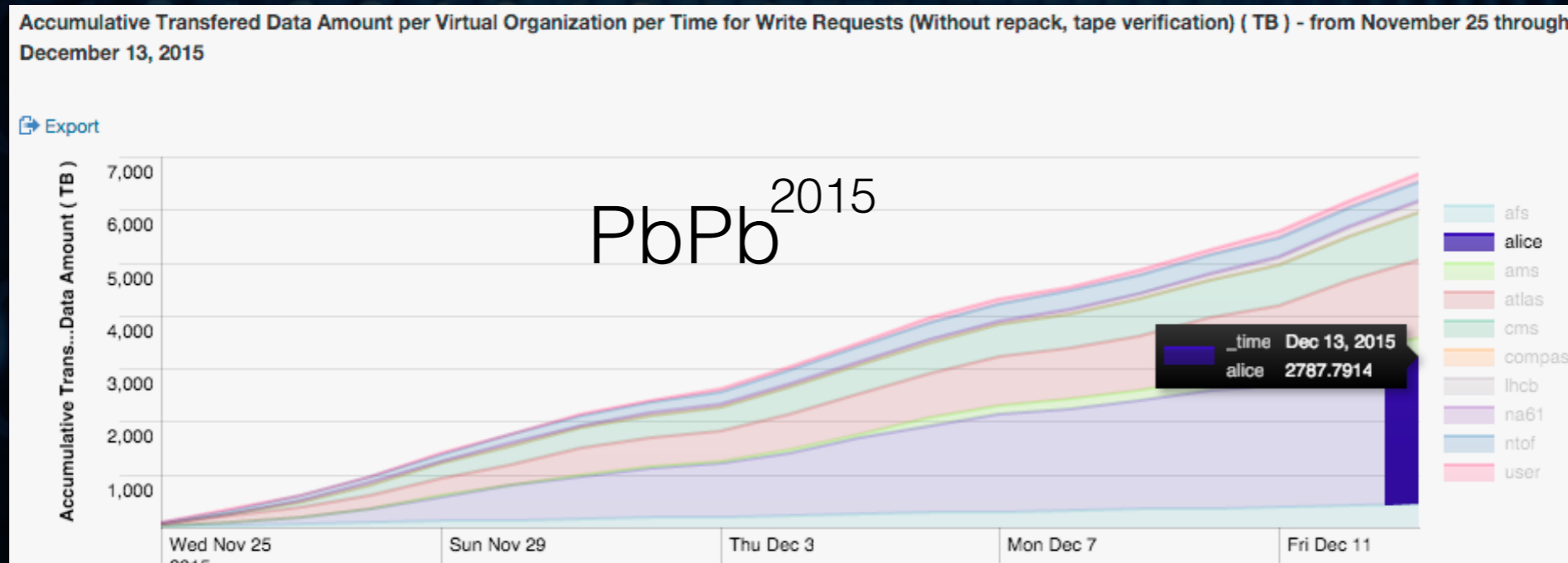Tape policies, per experiment/user/group resources

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

⤷ Export

PbPb$^{2015}$

| _time | Dec 13, 2015 |
| alice | 2787.7914 |

afs
alice
ams
atlas
cms
compass
lhcb
na61
ntof
user

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 |

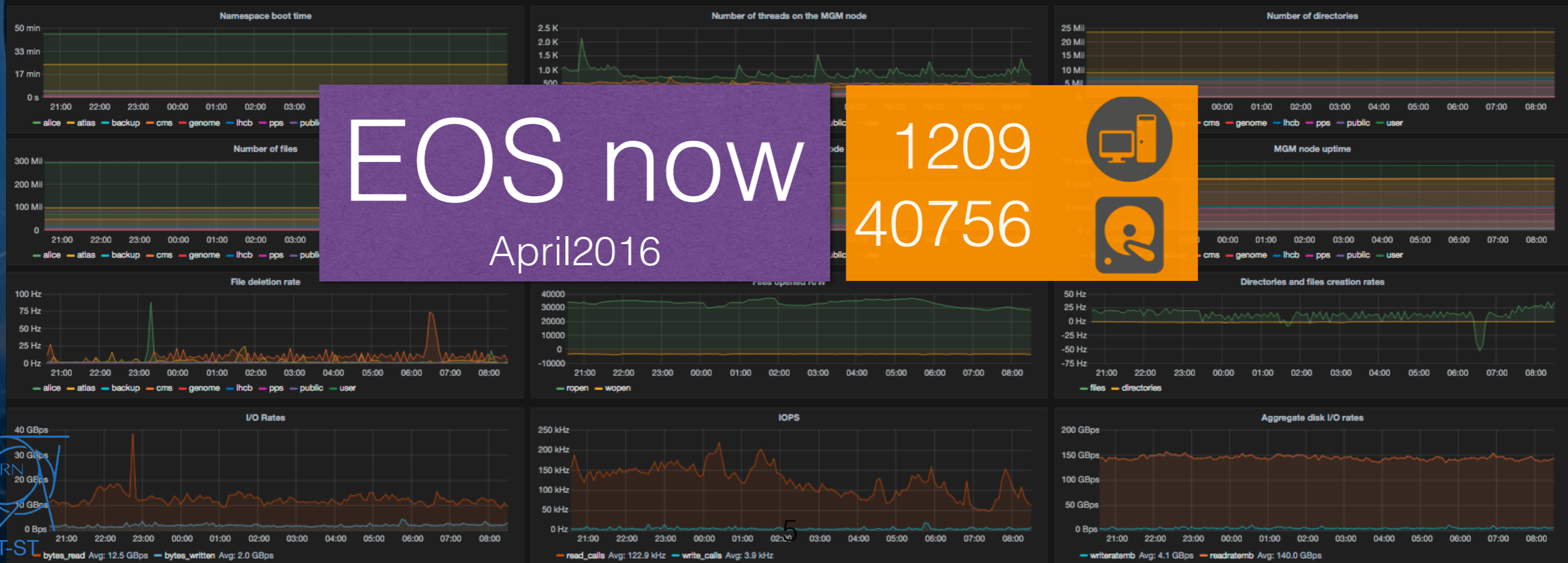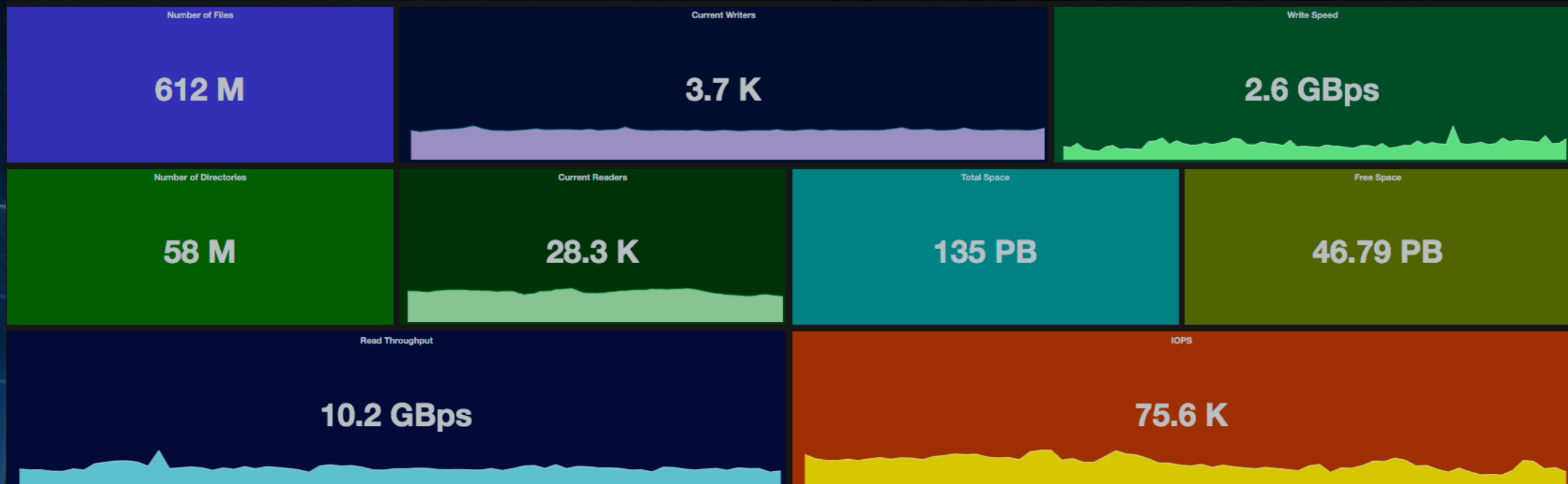PbPb$^{2015}$

©divia

4

CERN
IT-ST

EB era
- Scaling well on #disks
- Performant and manageable
- Main storage platform

NS future
- Fast and consistent
- Horizontally scalable (no single box limitiation)
- zero boot time

**EOS now**
April 2016

| | |
|---|---|
| Number of Files | 612 M |
| Current Writers | 3.7 K |
| Write Speed | 2.6 GBps |
| Number of Directories | 58 M |
| Current Readers | 28.3 K |
| Total Space | 135 PB |
| Free Space | 46.79 PB |
| Read Throughput | 10.2 GBps |
| IOPS | 75.6 K |

1209
40756

# 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,…
Openlab/COMTRADE JRC, Univ. Vienna, Univ. Trieste
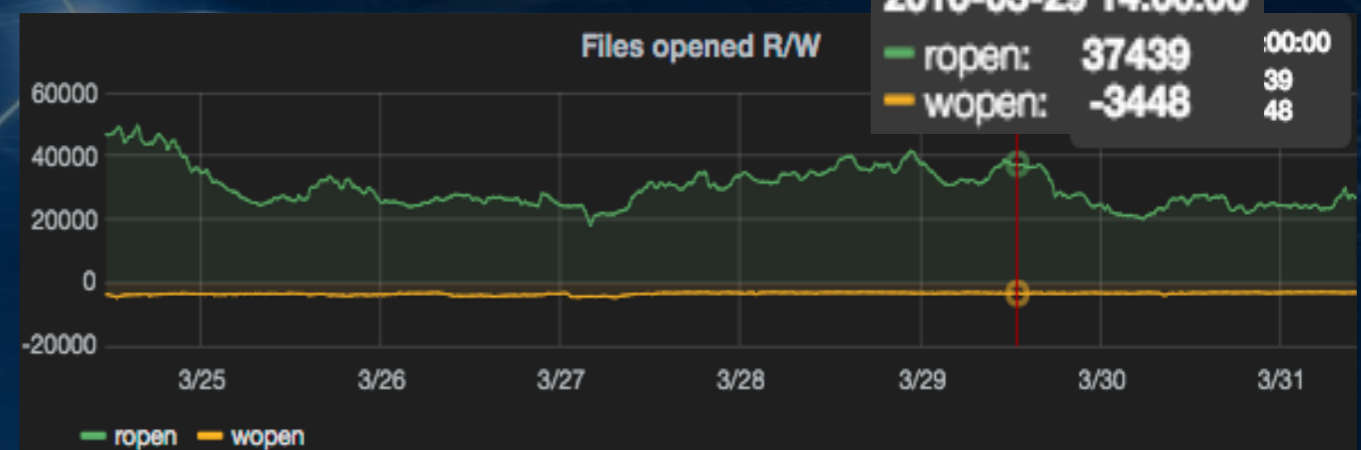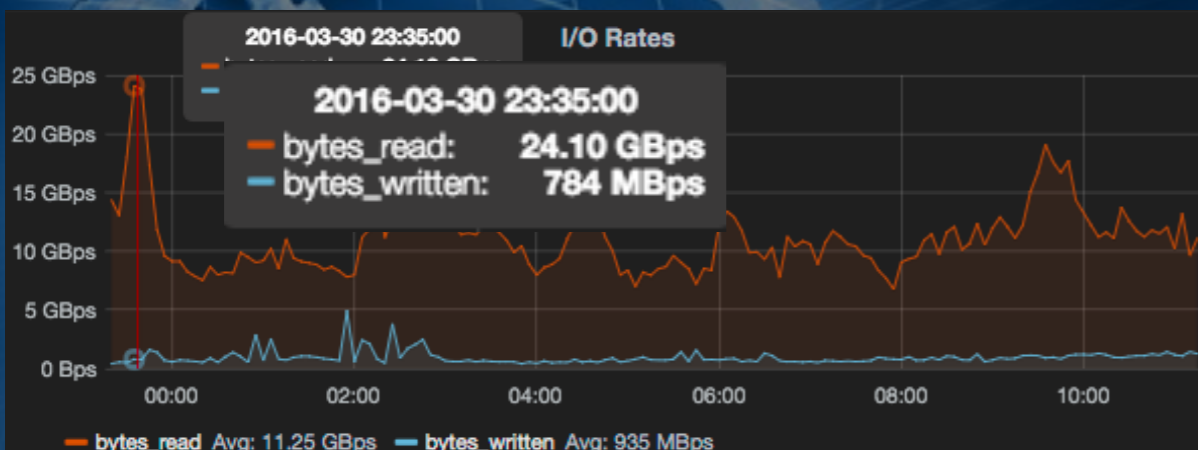
CDR
Data processing
Adaptable User Analysis
Catering with different uses CernBOX

Community data
DmaaS (iJupyter)
Sync share

CERNBox

I/O Rates
2016-03-30 23:35:00
2016-03-30 23:35:00
— bytes_read:      24.10 GBps
— bytes_written:    784 MBps

25 GBps
20 GBps
15 GBps
10 GBps
5 GBps
0 Bps
00:00   02:00   04:00   06:00   08:00   10:00
— bytes_read Avg: 11.25 GBps — bytes_written Avg: 935 MBps

Files opened R/W
2016-03-29 14:00:00
:00:00
— ropen:   37439     39
— wopen:   -3448     48

60000
40000
20000
0
-20000
3/25   3/26   3/27   3/28   3/29   3/30   3/31
— ropen — wopen

# Future Shared FS ? more later…

IOPS
2016-03-28 23:15:00
2016-03-28 23:15:00
— read_calls:   320 kHz
— write_calls:    3 kHz

400 kHz
300 kHz
200 kHz
100 kHz
0 Hz
3/25   3/26   3/27   3/28   3/29   3/30   3/31
— read_calls Avg: 142 kHz — write_calls Avg: 5 kHz

6

©lmascetti

# Can go distributed, can be shared and synced

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 ~]#
```
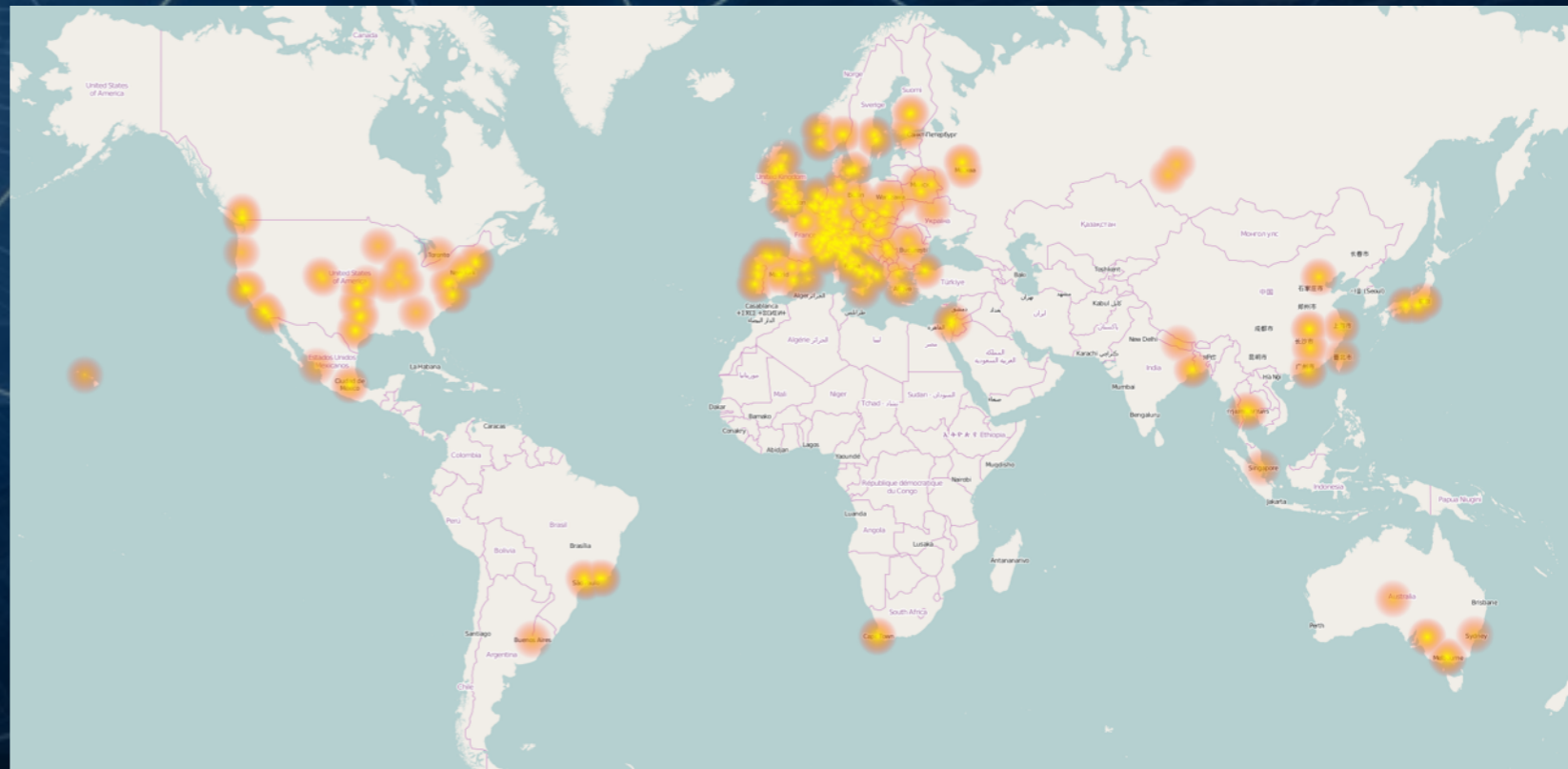
Community data DmaaS (iJupyter) share Sync

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


©lmascetti

Clients delocalization   Used from 22 ms to 300 ms

Multi-site deployment

```
[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 ~]#
```
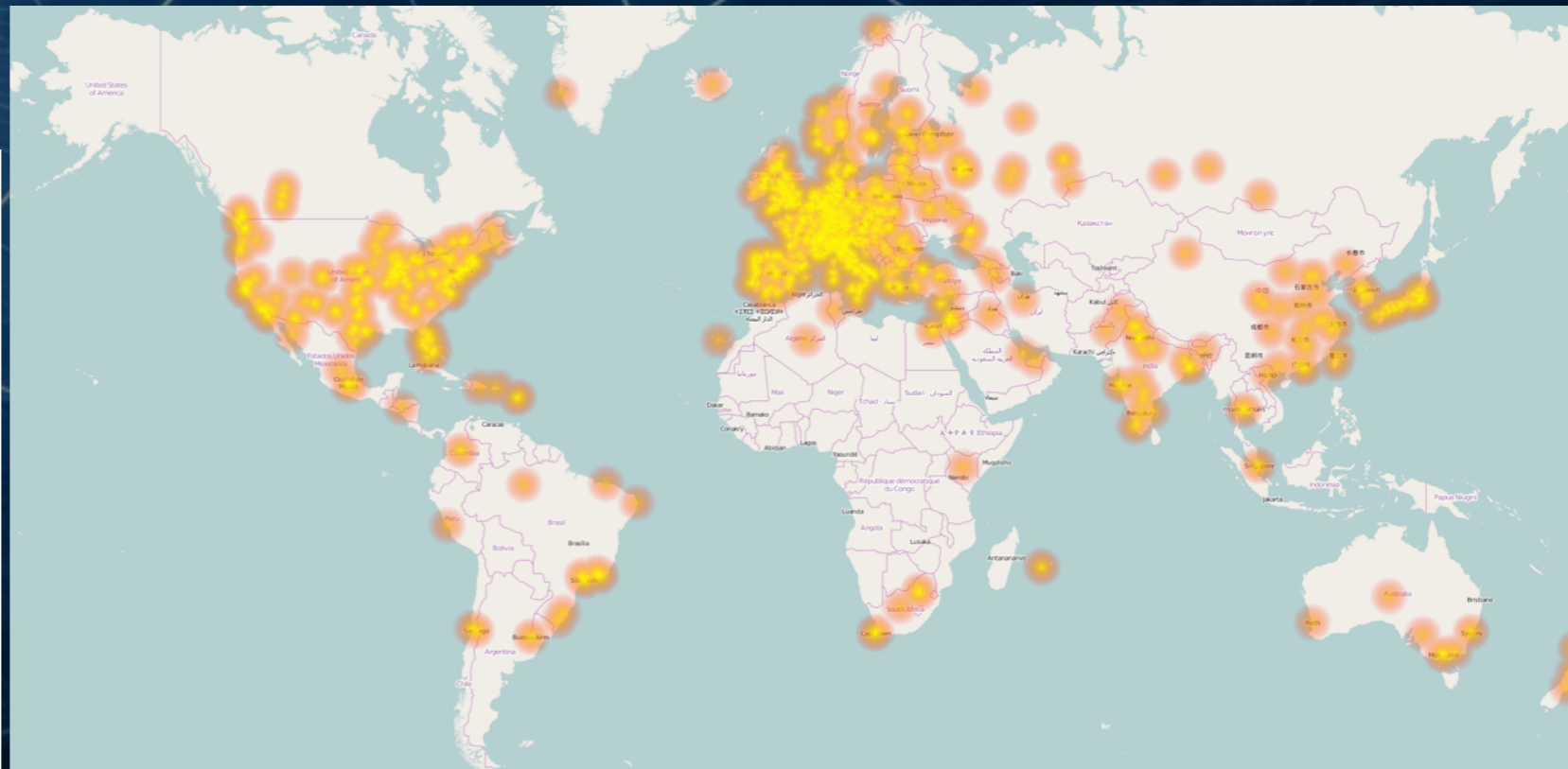
Community data share
DmaaS (iJupyter) Sync

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



©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 ~]#
```
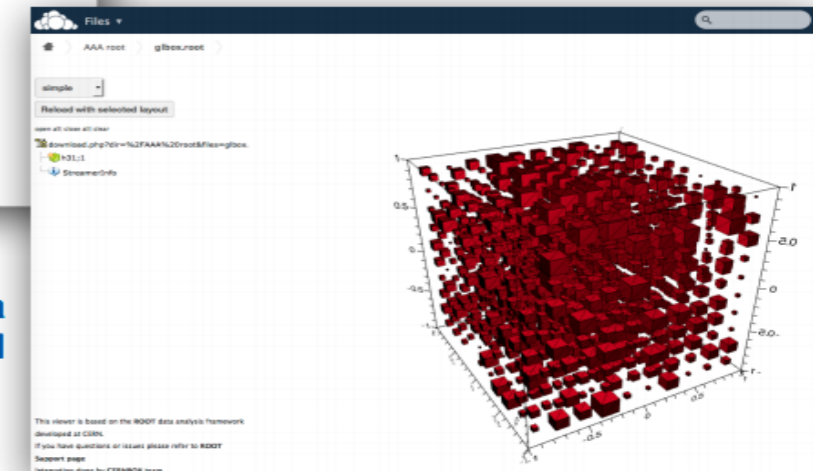
Community data
DmaaS (iJupyter) Sync
share

## 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
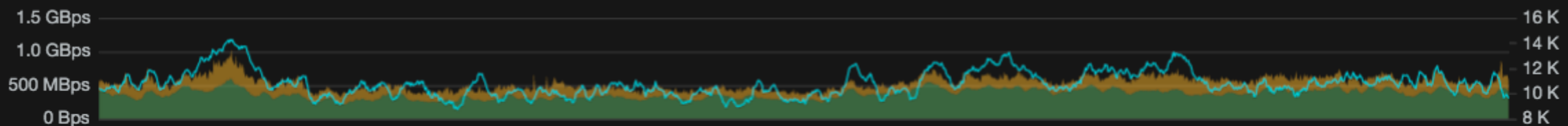Under evaluation

Large contribution
Community

Code development
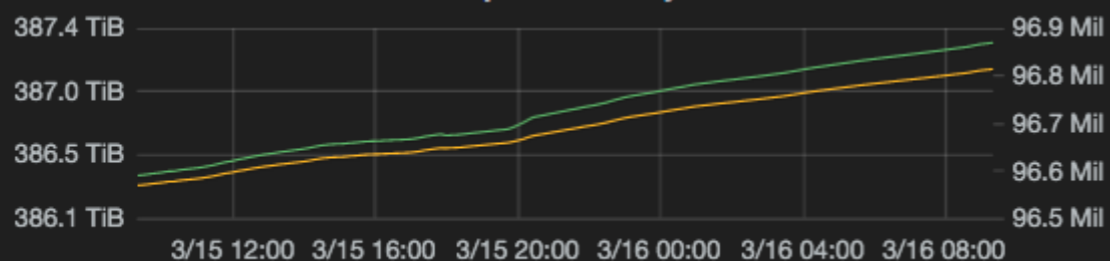CERN-IT/ST

Largest Cluster 30PB
Deployed to date
40k OSDs

Multi-site
In production

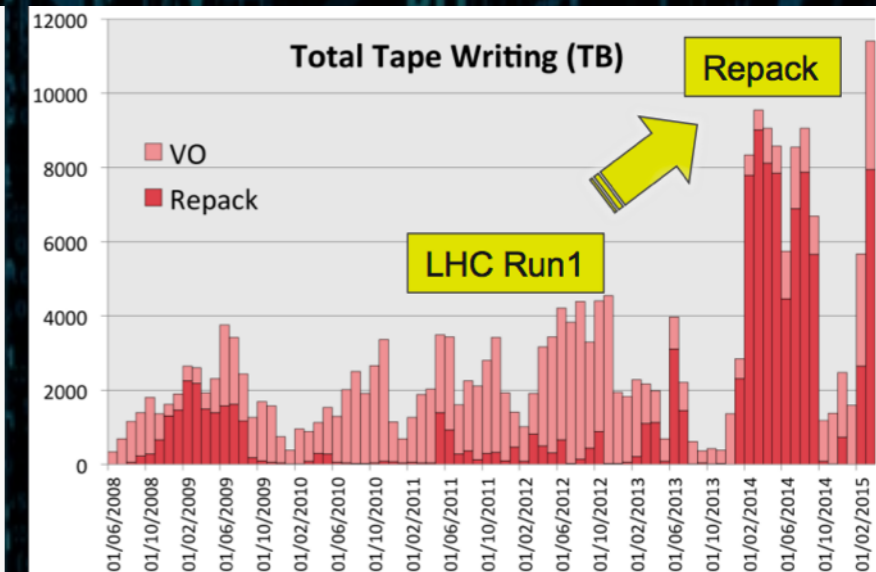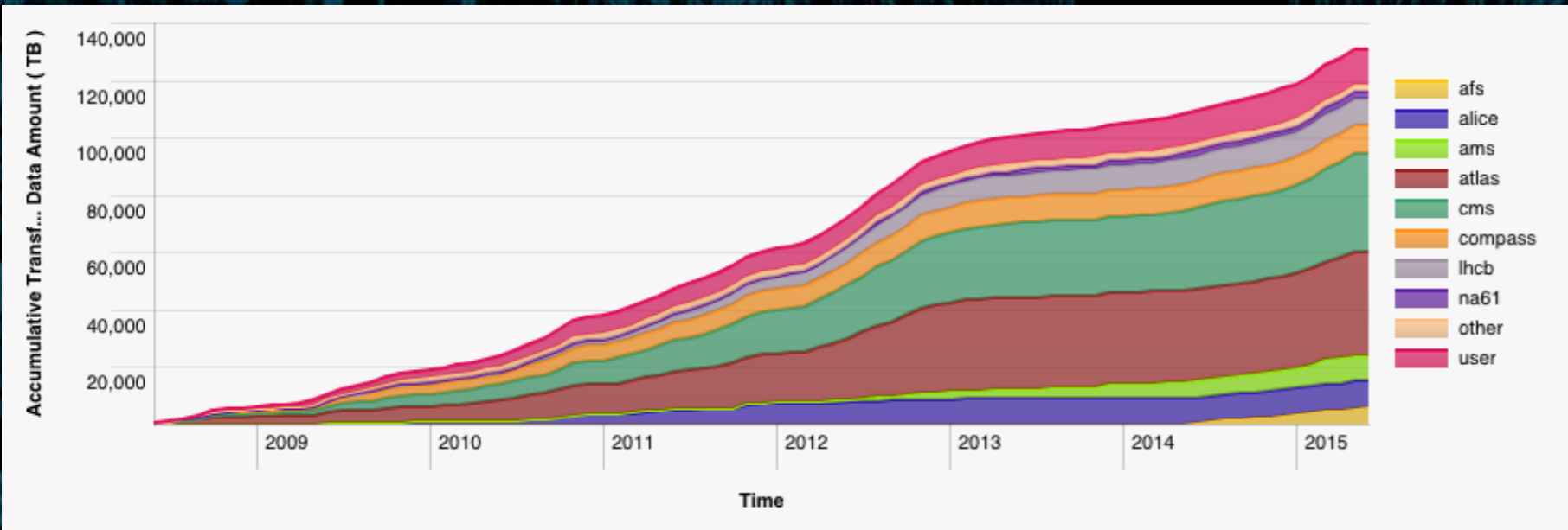3PB@wigner
1PB@meyrin

**2870** images

**2037** volumes

1.5 GBps
1.0 GBps
500 MBps
0 Bps

16 K
14 K
12 K
10 K
8 K

**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

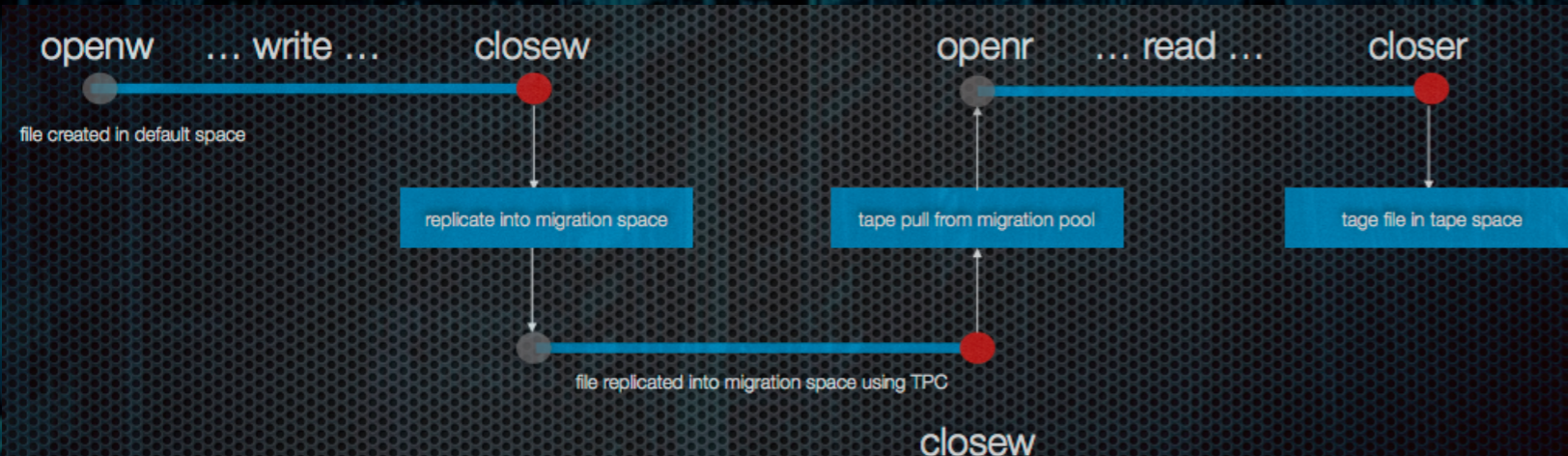3/15 12:00  3/15 16:00  3/15 20:00  3/16 00:00  3/16 04:00  3/16 08:00

CERN
IT-ST

# CERN Tape Archive

Technology driven: new medias brings ↑density ↑speed

Towards a pluggable tape backend (EOS)

Cold by definition: hight throughput, high latency





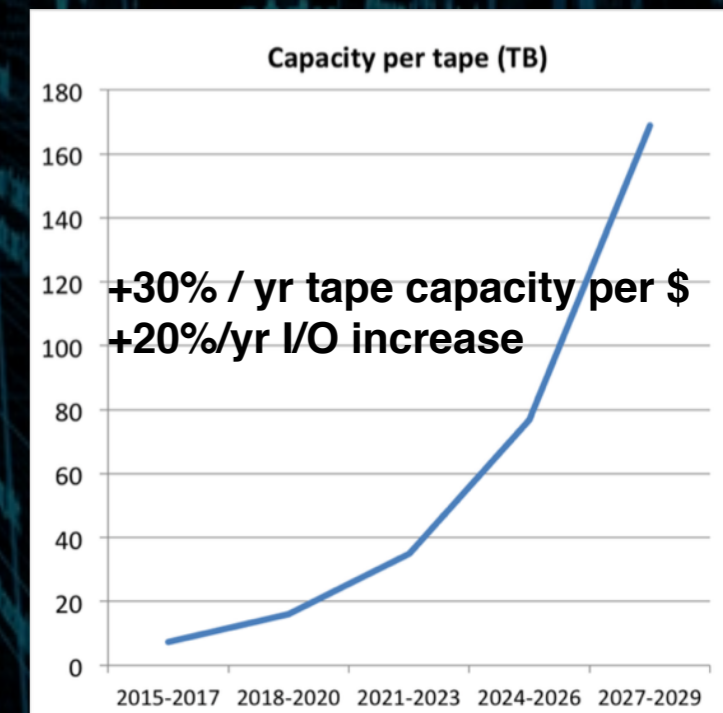Tape best technology for data repositories: TCO media power density and resilient/reliable
very large disk caches nowadays



## Cross-system workflows

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

©apeters

CERN
IT-ST

# 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
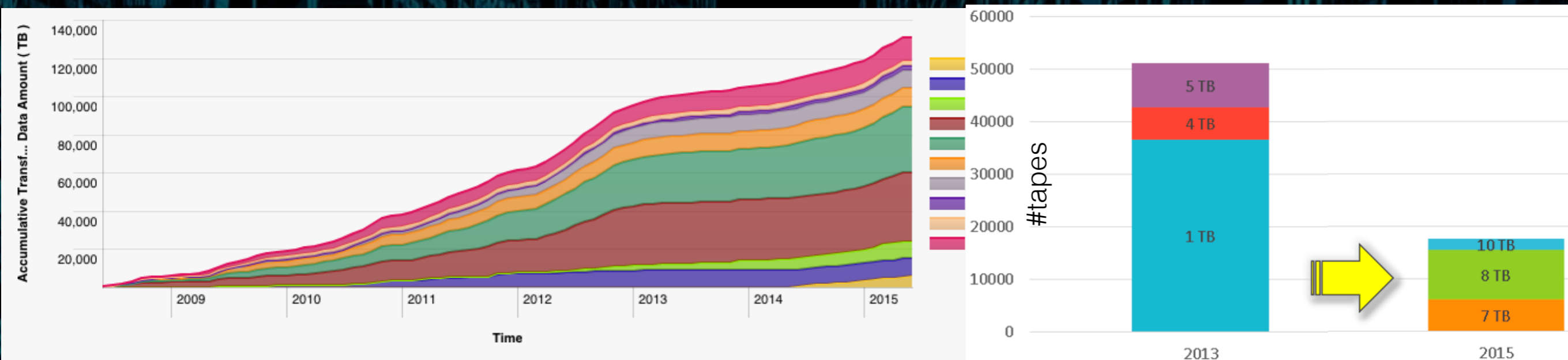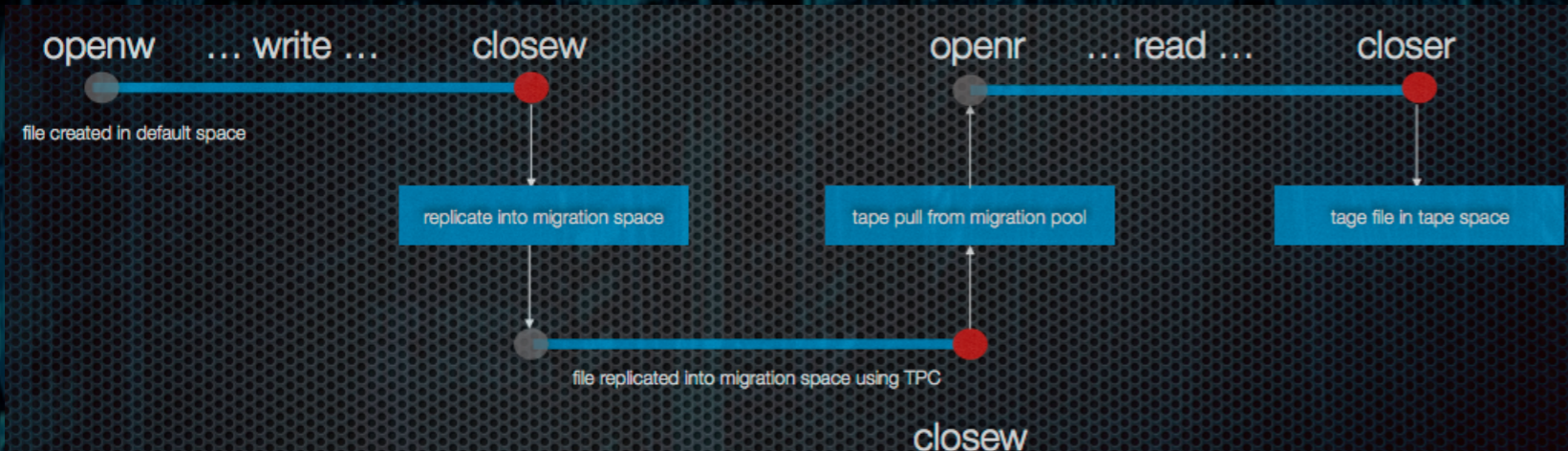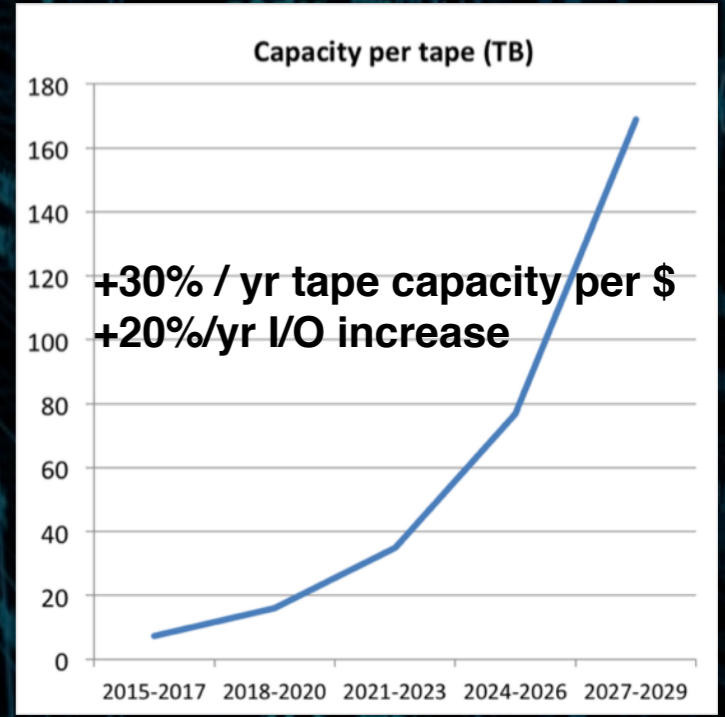


Cross-system workflows

CLIENT    WF 1    WF 2    ©apeters

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



12

# Goals ¹Make data access easy

$HOME — Batch /home, Laptop use ➡

**?**

DATA ACCESS — Protocol based xroot, rfio, etc. ➡

AFS being (slowly) ramped down

**My Laptop**
Small scale analysis
Test jobs

AFS $home

**lxbatch/interactive**
Large scale experiment processing
User extensive analysis

AFS $home
/cvmfs

xroot, *

**Data Access**

**/eos/atlas/topphys**, /mycernbox, /cvmfs/athena

# Goals

## ¹Make data access easy

$HOME    Batch /home
Laptop use    →    **SharedFS**    /cernbox
Syncing option

DATA ACCESS    Protocol based
xroot, rfio, etc.    →    **Large Scale Storage Access**    MountPoints
/eos, /cernbox
/cvmfs

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

## My Laptop
Small scale analysis
Test jobs

## lxbatch/interactive
Large scale experiment processing
User extensive analysis

## Mounts

/cvmfs/athena

/mycernbox

/eos/atlas

## Data Access

**/eos/atlas/topphys**, /mycernbox, /cvmfs/athena

# Goals

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.root**
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

CERN
IT-ST

# Goals
summary

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

# 3 We are here for you

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