# **EOS** Report, Evolution & Strategy

## **HEPiX** *online* **25.-29.4.2022**
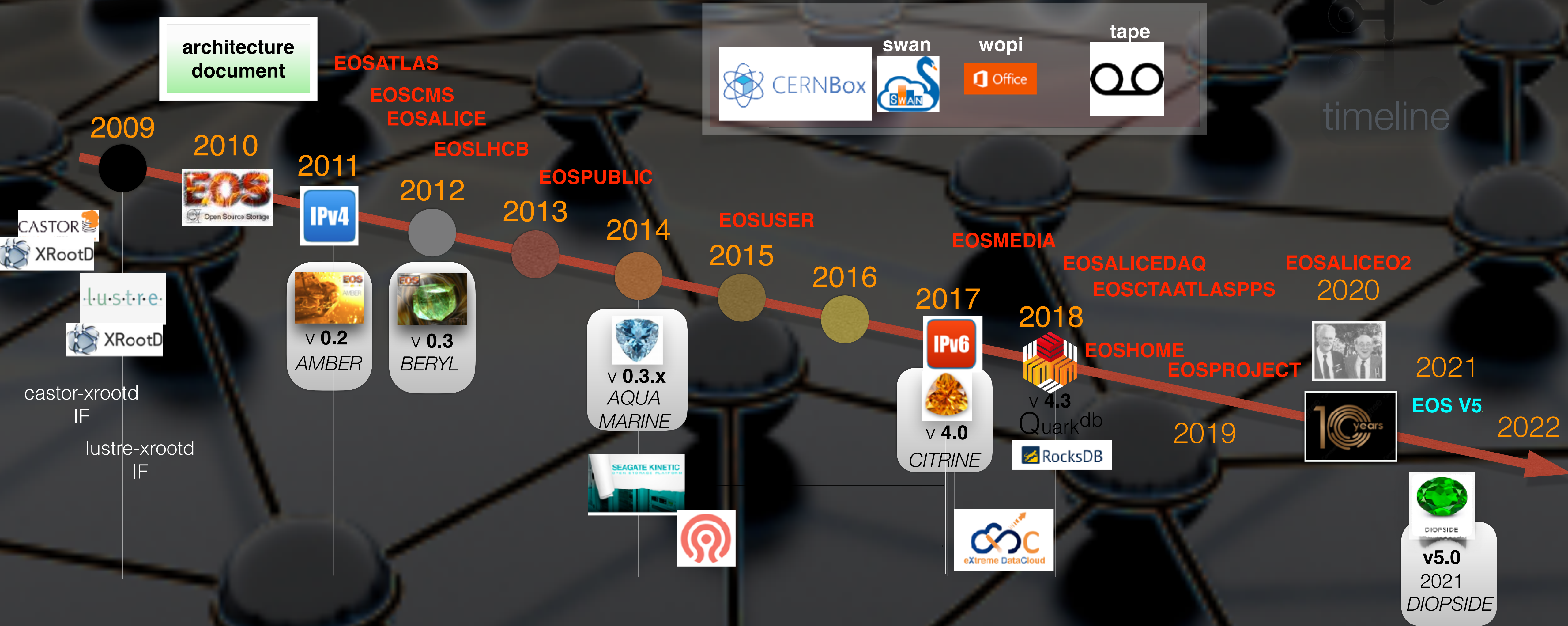
**Andreas-Joachim Peters**
CERN IT-SD for the EOS project

# Overview

- EOS Project Overview
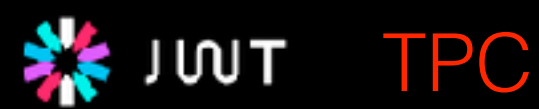- Highlights from the EOS '22 workshop
- Roadmap - Evolution - Strategy

- EOS Releases
  - 29 x Citirine
  - 3 x Diopside
  - C8 Stream

- Production
  - HTTP(S) Protocol  JWT  TPC
  - Erasure Coding @ TBit/s

- Completion & Consolidation
  - FSCK - integrity
  - GRPC - EOS aaS

0064ebb1
by Elvin Sindrilaru at 2021-06-11T09:46:18+02:00
DOC: Update release notes for 5.0.0

**EOS** 5

EOS@CERNBOX
Availability **99.9999%**

# EOS Development Team

**Andreas** J. Peters
project leader & core developer

**Elvin** Alin Sindrilaru
core developer & operations

**Cedric** Caffy
core developer & operations

**Abhishek** Lekshmanan
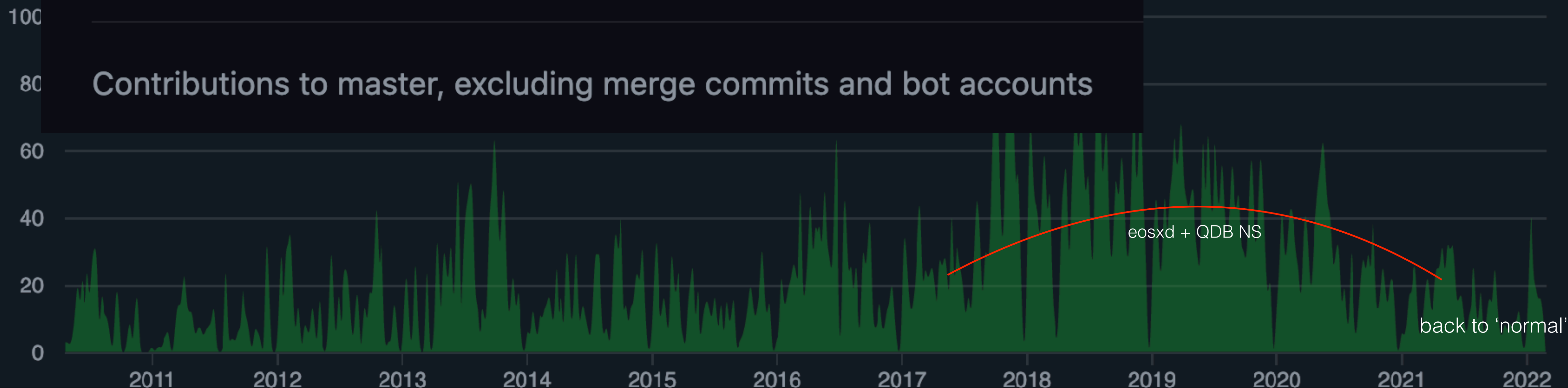core developer & operations

**Jaroslav** Guenther
development & operations

- **2018-20** were years with very **high development activity** due to the architectural changes coming with QuarkDB and eosxd FUSE implementation
- **2021** activity got **back to 'normal'**



Apr 25, 2010 – Mar 2, 2022

Contributions to master, excluding merge commits and bot accounts

eosxd + QDB NS

back to 'normal'

EOS workshop 2022

7th to 10th of march @ VIRTUAL

*platform for exchange between developers, users, sites
and people interested in storage technology*

**disk - tape - cloud - sync & share - devops**

10

https://indico.cern.ch/event/1103358/



All slides & Recordings available!

12

# Who was participating ... ?



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ■ 2017 | ■ 2018 | ■ 2019 | ■ 2020 | ■ 2021 | ■ 2021 | ■ 2022 | ■ 2022 |

*before start of the workshop

■ CERN    ■ External

up to 100 following in ZOOM +Webcast

# Scope of '22 workshop

- **Where are we today with EOS** for disks and tape storage, how does the EOS 5 version look like, what is the **direction** of the project?
- **Where and how** do people use it?
- **What** does EOS, XRootD, CTA and CERNBox offer to you**?**

**46 Presentations - 16 hours**

14

# Workshop Contents
## Monday Morning

**EOS5**, **EOS**@CERN,
**ALICEO2** Commissioning

**XRootD** Landscape, Monitoring,
Workload Replay & Benchmarking,
ScienceBox

Run-3 **Comissioning**
EOS@Vienna, Fermi, Purdue,

EOS-**PoC**@CNAF
(Kubernetes+ObjectStorage)

WLCG **Token** Support in EOS

| | | |
|---|---|---|
| **LHC Data Storage: RUN 3 Data Taking Commissioning** | | *Dr Maria Arsuaga Rios* |
| *31/3-004 - IT Amphitheatre, CERN* | | 15:45 - 16:05 |
| **EOS site report Vienna** | | *Erich Birngruber* |
| *31/3-004 - IT Amphitheatre, CERN* | | 16:05 - 16:25 |
| **EOS at the Fermilab LHC Physics Center** | | *Dan Szkola* |
| *31/3-004 - IT Amphitheatre, CERN* | | 16:25 - 16:40 |
| **EOS deployment at Purdue** | | *Stefan Piperov* |
| *31/3-004 - IT Amphitheatre, CERN* | | 16:40 - 16:55 |
| **EOS and Ceph integration with Kubernetes** | | *Federico Fornari* |
| *31/3-004 - IT Amphitheatre, CERN* | | 16:55 - 17:10 |
| **Data flowing on the Stream** | | *Cristian Contescu* |
| *31/3-004 - IT Amphitheatre, CERN* | | 17:10 - 17:30 |
| **WLCG tokens integration and support in EOS** | | *Elvin Alin Sindrilaru* |
| *31/3-004 - IT Amphitheatre, CERN* | | 17:30 - 17:40 |

16

EOS@**Kisti, GRIF, JRC**

EOS **Tools**

EOS Enhancements for **IO Shaping**
( see following presentation )

EOS **Windows** Client
EOS **Durability**

| | | |
|---|---|---|
| **Operation status of Custodial Disk Storage for the ALICE experiment** | | *Sang Un Ahn* |
| *31/3-004 - IT Amphitheatre, CERN* | | 09:00 - 09:15 |
| **EOS deployment at GRIF** | | *Dr Emmanouil Vamvakopoulos* |
| *31/3-004 - IT Amphitheatre, CERN* | | 09:15 - 09:35 |
| **EOS site report of the Joint Research Centre** | | *Armin Burger* |
| *31/3-004 - IT Amphitheatre, CERN* | | 09:35 - 09:55 |
| **EOS GroupBalancer improvements** | | *Abhishek Lekshmanan* |
| *31/3-004 - IT Amphitheatre, CERN* | | 09:55 - 10:15 |
| **EOS migration tools** | | *Dr Jaroslav Guenther* |
| *31/3-004 - IT Amphitheatre, CERN* | | 10:15 - 10:35 |
| **Coffee Break** | | |
| *31/3-004 - IT Amphitheatre, CERN* | | 10:35 - 10:55 |
| **Direct IO, IO priority and Bandwidth Policies in EOS** | | *Andreas Joachim Peters* |
| *31/3-004 - IT Amphitheatre, CERN* | | 10:55 - 11:10 |
| **Encryption and Obfuscation Support in EOS** | | *Andreas Joachim Peters* |
| *31/3-004 - IT Amphitheatre, CERN* | | 11:10 - 11:25 |
| **Taming Batch Access to EOS at CERN** | | *Andreas Joachim Peters* |
| *31/3-004 - IT Amphitheatre, CERN* | | 11:25 - 11:35 |
| **xrdcp primer** | | *Michal Kamil Simon* |
| *31/3-004 - IT Amphitheatre, CERN* | | 11:35 - 11:45 |
| **EOS Windows client productisation** | | *Gregor Molan* |
| *31/3-004 - IT Amphitheatre, CERN* | | 11:45 - 12:00 |
| **EOS Durability Summary** | | *Manuel Reis* |
| *31/3-004 - IT Amphitheatre, CERN* | | 12:00 - 12:10 |

17

# Workshop Contents
## Wednesday CTA Day

**CTA** Project, Status & Community
CTA **Operation**
Tape **Rest API**

CTA@**AARNet, IHEP, dCache, FNAL, RAL**
**Tapeformat**

CERN
Tape Archive

| | |
|---|---|
| The CTA project, team and community | Oliver Keeble |
| 31/3-004 - IT Amphitheatre, CERN | 08:55 - 09:05 |
| **CTA at AARNet** | Mr Denis Lujanski Not Supplied |
| 31/3-004 - IT Amphitheatre, CERN | 09:05 - 09:20 |
| **EOS and CTA Status at IHEP** | Yujiang Bi |
| 31/3-004 - IT Amphitheatre, CERN | 09:20 - 09:35 |
| **CTA Status and Roadmap** | Michael Davis |
| 31/3-004 - IT Amphitheatre, CERN | 09:35 - 09:55 |
| **How to enable EOS for tape** | Julien Leduc |
| 31/3-004 - IT Amphitheatre, CERN | 09:55 - 10:15 |
| **Break** | |
| 31/3-004 - IT Amphitheatre, CERN | 10:15 - 10:35 |
| **Configuring user access control in CTA** | Volodymyr Yurchenko |
| 31/3-004 - IT Amphitheatre, CERN | 10:35 - 10:50 |
| **Tape Drive Status Lifecycle** | Jorge Camarero Vera |
| 31/3-004 - IT Amphitheatre, CERN | 10:50 - 11:05 |
| **EOSCTA file restoring** | Miguel Barros |
| 31/3-004 - IT Amphitheatre, CERN | 11:05 - 11:20 |
| **Maintaining consistency in an EOSCTA system** | Richard Bachmann |
| 31/3-004 - IT Amphitheatre, CERN | 11:20 - 11:40 |

| | |
|---|---|
| **Evaluation of CTA for use at Fermilab** | Ren Bauer |
| 31/3-004 - IT Amphitheatre, CERN | 16:00 - 16:20 |
| **An HTTP Rest API as SRM replacement for tape access** | Cedric Caffy |
| 31/3-004 - IT Amphitheatre, CERN | 16:20 - 16:40 |
| **CTA at RAL** | Dr George Patargias |
| 31/3-004 - IT Amphitheatre, CERN | 16:40 - 17:00 |
| **dCache integration with CTA** | Mr Tigran Mkrtchyan |
| 31/3-004 - IT Amphitheatre, CERN | 17:00 - 17:20 |
| **CTA tape format support : BoF discussion** | Michael Davis |
| 31/3-004 - IT Amphitheatre, CERN | 17:20 - 18:00 |

# Workshop Contents
## Thursday Morning

**XRootd** Erasure Encoding
EOS & XCache **Access Analytics** at CERN
EOS Run-3 **Roadmap**

| | |
|---|---|
| **Native XRootD EC @ SLAC** | *Michal Kamil Simon* 📎 |
| *31/3-004 - IT Amphitheatre, CERN* | 09:00 - 09:20 |
| **EOS and XCache data access performance for LHC analysis at CERN** | *Dr Andrea Sciabà* 📎 |
| *31/3-004 - IT Amphitheatre, CERN* | 09:20 - 09:45 |
| **EOS 5 during Run-3 Roadmap** | *Andreas Joachim Peters* 📎 |
| *31/3-004 - IT Amphitheatre, CERN* | 09:45 - 10:05 |
| **Community Feedback & Open Discussion** | 📎 |
| *31/3-004 - IT Amphitheatre, CERN* | 10:05 - 11:00 |

**EOS5** & **XRootD5** is in production (1st instance EOSAMS )
- not a completely smooth ride so far, but we are almost there - bugs, race conditions++
- main interest for EOS5 encrypted wire protocol in XRootD5, few new client features

EOS@CERN growing and
growing and…
 ~**660-680 PB** in '22

**CERN Services**



EOS for Physics: Numbers

| 6 + 3 Instances | 2.67 Bill Number of Files +8% | 215 Mil Number of Directories +13% | 514.84 PB Total Space +75% |

- Protocol Usage at CERN
  (#file accessed not volume!)
  - **XRootD** is most versatile
  - **FUSE** is for convenience
  - **HTTP(S)** mainly TPC
  - **gridFTP** disappearing



Total reads per protocol

Total writes per protocol

- XROOTD
- FUSE
- GRIDFTP
- HTTP

**EOS**

Storage@CERN **Scale Change** = **TBit/s** per Instance **GB/s** per client

- 200 GB/s R/W in **EOSALICEO2**

- 8 GB/s for a single multithreaded analysis application

- 100GE rocks! … but 12 GB/s Disks + 100GE Network != 100GE IO over network - 50%

- **O2 Benchmarking**
  - with new 100GE disk server we push the maximum performance with erasure coding to 6-7 GB/s per disk server

  - although machines have a theoretical performance of 10 GB/s we cannot exploit this for the time being

- **IO bound analysis** on 100GE with EC demonstrated …
  - … that we can reach **8 GB/s** data INGRES on a single client using parallel IO and EC files using `root` protocol
  - … that xrdcl-record/replay is able to sample IO of a complex multithreaded application and replay with identical timing - easing future benchmarking

**Site Reports** : various reports about existing & new deployments
with RAID, Replication & Erasure Coding

**Storage Virtualisation** : topic for PoC **CNAF**, PoC **CERNBOX**,

but saw also interest from other external sites

- new feature in latest EOS version:
  - *available*: **local redirects** to shared filesystem
    requires read-only access on shared filesystem for all clients
  - *prototype*: **file registration/adoption** from
    local filesystem
    - interesting for people who want to put a WLCG
      stack on top of shared filesystem



25

**EOS**

**CTA Day**

- CTA in full production at CERN, AARNet, IHEP, RAL
- modular design, envisaged for dCache

- RAL has migrated **CASTOR** to **ANTARES** [CTA]
  ECHO = disk       ANTARES = tape

- FNAL to decide about CTA options

- HTTP Tape API
  - agreement between CTA, dCache, Storm (& FTS)
    - hopefully: use and forget

    - write a JSON file and a three-line CURL command
      and enjoy - it is not for end-users!

26

**EOS**

- •**XRootD** brings client-side erasure coding based on INTEL ISA-L library
  - • XRootD native EC - see presentation this afternoon Michal Simon
  - • until now only PoC prototype in EOS - goal is to coalesce XRootD & native EOS EC

- • A study at CERN evaluated the need for **dedicated** high-performance **storage for analysis for Run-3** and xCache as cache-frontend to EOS

EOS and XCache data access performance for LHC analysis at CERN

Dirk Duellmann, Bernd Panzer-Steindel, Markus Schulz, Andrea Sciabà, David Smith (CERN IT-SC)

- • results: see presentation by Andrea Sciabà **this afternoon**

27

- **EOS5** roadmap during Run-3
  - Simplification / Deprecation
  - Faster Namespace ( Locking Model )
  - Storage Virtualisation ( Stateles Storage Server )
  - High-Availability automatisation
  - Erasure Coding Extension
  - Filesystem Latency **/eos**
  - More supported platforms: ARM, RH8, Ubuntu
  - Documentation Trilogy



Latest EOS5 testing release **5.0.18**
Latest EOS4 testing release **4.8.83**

- support for EOS **token with HTTP & GRPC** using authorisation header/authz field

- fine-grained **IO policies** - configure IO by user/group, directory or application

- **FSCK** improvements for EC files

- new implementation of **monitoring** finished transfers

- per file optional **File encryption/obfuscation** of on-disk format

  - **codec** client-side for FUSE clients, server-side for remote-access protocols

- prototype: **Share ACLs** - ACL syntax extension tailored to share permissions

- merged: **TAPE REST API**

- EOS5 **EGI** release - C8S + RockyLinux8 available, C9S in CI (not KOJI)

# Developments / News
## Finished Transfer Monitoring

```
╭─> Transfer (tf) sample info every 5 min: tf time for 90/95/99% of data, max tf and report times, average tf size, tf count.
```

| io  | application    | 90% [s] | 95% [s] | 99% [s]  | max [s]  | max report [s] | avg tf size | tf #    | sample end time           |
|-----|----------------|---------|---------|----------|----------|----------------|-------------|---------|---------------------------|
| out | eoscp          | 3       | 3       | 4        | 4        | 2              | 104.66 M    | 61      | Tue Apr 26 11:25:18 2022  |
| out | eos/gridftp    | 679     | 717     | 747      | 754      | 2              | 1.07 G      | 10      | Tue Apr 26 11:24:17 2022  |
| out | eos/converter  | 0       | 0       | 0        | 0        | 0              | 0           | 0       | Tue Apr 26 11:23:49 2022  |
| out | eos/replication| 0       | 0       | 0        | 0        | 0              | 0           | 0       | Tue Apr 26 11:24:49 2022  |
| out | fuse           | 0       | 0       | 0        | 0        | 0              | 0           | 0       | Tue Apr 26 11:23:15 2022  |
| out | other          | 475     | 4.53 K  | 17.62 K  | 26.19 K  | 2              | 96.19 M     | 1.04 K  | Tue Apr 26 11:22:51 2022  |
| out | fuse::lxplus   | 0       | 0       | 0        | 0        | 0              | 0           | 0       | Tue Apr 26 11:25:39 2022  |
| out | fuse::bi       | 11      | 12      | 16       | 752      | 2              | 3.87 M      | 92.69 K | Tue Apr 26 11:22:34 2022  |
| out | fuse::amssoc   | 23      | 29      | 38       | 44       | 2              | 1.54 K      | 125     | Tue Apr 26 11:24:12 2022  |
| out | tpc            | 0       | 0       | 0        | 0        | 0              | 0           | 0       | Tue Apr 26 11:23:04 2022  |
| in  | eoscp          | 17      | 18      | 19       | 20       | 2              | 2.40 G      | 25      | Tue Apr 26 11:23:09 2022  |
| in  | eos/gridftp    | 78      | 88      | 102      | 116      | 2              | 324.98 M    | 49      | Tue Apr 26 11:23:07 2022  |
| in  | eos/converter  | 0       | 0       | 0        | 0        | 0              | 0           | 0       | Tue Apr 26 11:23:54 2022  |
| in  | eos/replication| 0       | 0       | 0        | 0        | 0              | 0           | 0       | Tue Apr 26 11:24:49 2022  |
| in  | fuse           | 0       | 0       | 0        | 0        | 0              | 0           | 0       | Tue Apr 26 11:25:13 2022  |
| in  | other          | 30      | 32      | 33       | 33       | 2              | 345.02 M    | 15      | Tue Apr 26 11:23:06 2022  |
| in  | fuse::lxplus   | 15      | 16      | 16       | 16       | 1              | 365         | 1       | Tue Apr 26 11:21:03 2022  |
| in  | fuse::bi       | 2       | 3       | 3        | 3        | 2              | 7.26 M      | 38      | Tue Apr 26 11:23:11 2022  |
| in  | fuse::amssoc   | 30      | 35      | 40       | 41       | 2              | 299         | 52      | Tue Apr 26 11:22:22 2022  |

```
[root@host ~]# export EOS_FUSE_SECRET=8ae6e775-300a-440d-9555-f05786622bdd
[root@host ~]# eos cp /tmp/hw root://localhost//eos/encryption/encrypted
[root@host ~]# eos file info /eos/encryption/encrypted --fullpath
  File: '/eos/encryption/encrypted'  Flags: 0640
  Size: 13
…
  #Rep: 1
 Crypt: encrypted
```

| no. | fs-id | host | schedgroup | path | boot | configstatus | drain | active | geotag | physical location |
|-----|-------|------|------------|------|------|--------------|-------|--------|--------|-------------------|
| 0 | 2 | ajp.cern.ch | default.0 | /data/02 | booted | rw | nodrain | online | ajp | /data/06/000000b5/001ba29c |

```
[root@host ~]# cat /data/06/000000b5/001ba29c
#u˥o4c
[root@host ~]# eoscp -s -n root://localhost//eos/encryption/encrypted -
Hello World!
[root@host ~]# cat /eos/ajp/encryption/encrypted
Hello World!
```

```
[root@host ~]# env EOS_FUSE_SECRET=8ae6e775-300a-440d-9555-f05786622bdd eos cp -S root://localhost//
eos/ajp/encryption/encrypted.1 /var/tmp/128M.1
[eoscp] encrypted.1                 Total 122.07 MB|====================| 100.00 % [423.8 MB/s]
[eoscp] ##############################################################################
[eoscp] # Date                          : ( 1646311354 ) Thu Mar  3 13:42:34 2022
[eoscp] # auth forced=<none> krb5=FILE:/tmp/krb5cc_0_mN5SPEj0n9 gsi=<none>
[eoscp] # Source Name [00]              : root://localhost//eos/ajp/encryption/encrypted.1
[eoscp] # Destination Name [00]         : /var/tmp/128M.1
[eoscp] # Data Copied [bytes]           : 128000000
[eoscp] # Realtime [s]                  : 0.302000
[eoscp] # Eff.Copy. Rate[MB/s]          : 423.841063
[eoscp] # INGRESS [MB/s]                : 694.433142
[eoscp] # EGRESS [MB/s]                 : 1478.948098
[eoscp] # Write Start Position          : 0
[eoscp] # Write Stop  Position          : 128000000
[eos-cp] copied 1/1 files and 128.00 MB in 0.36 seconds with 354.46 MB/s
```

# Software & Service Strategy

- **EOS4** will reach EOL in 2022

  - **move all LHC instances** to **EOS5** when possible

    - EOSAMS running with EOS5

    - EOS5 clients in testing at CERN

  - recommendation to **move external deployments to EOS5** when possible

    - new functionality/performance NS improvements limited to EOS5

- **EC** Erasure Coding **in production**

  - EC in ALICEO2 EC(10+2) for all files

  - EC in CMS via policy in selected subtrees converting large files >1GB to EC (10+2)

  - expect more and more usage during Run-3

CERN storage technology
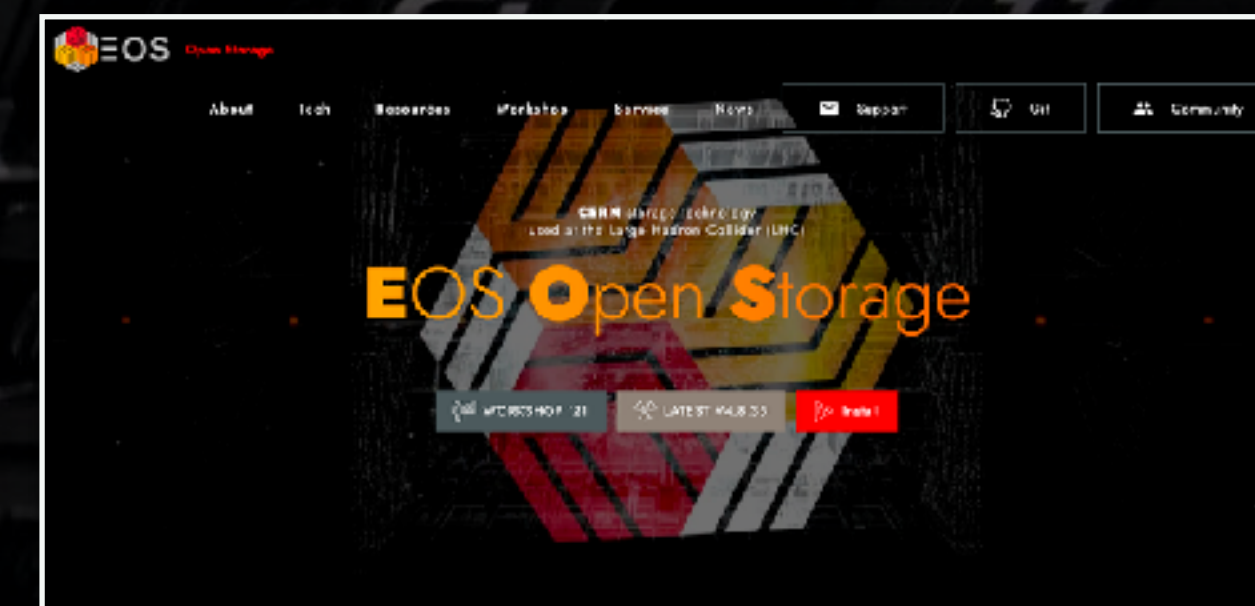used at the Large Hadron Collider (LHC)
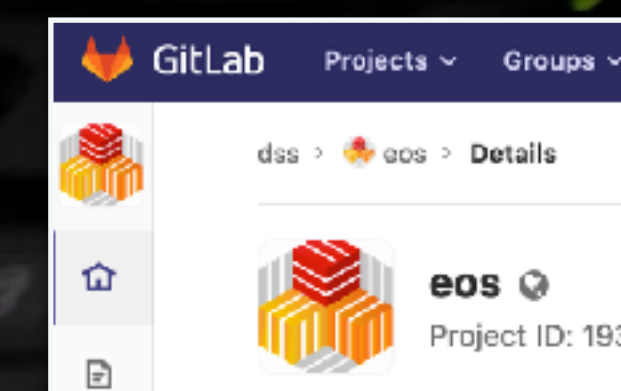
# EOS Open Storage

**Thank you!**

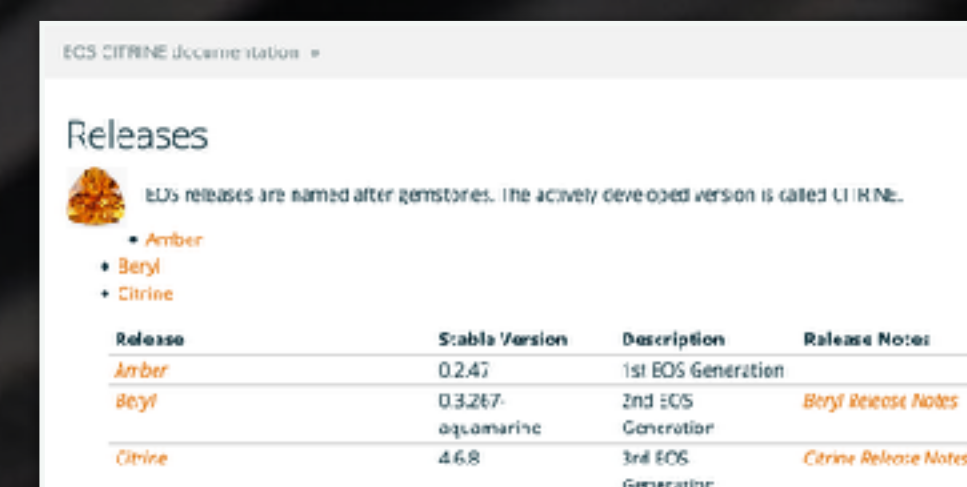**Question or Comments?**

eos.web.cern.ch

Web Page    https://eos.cern.ch

GIT Repository   https://gitlab.cern.ch/dss/eos

Community Forum   https://eos-community.web.cern.ch/

email: eos-community@cern.ch

Documentation    http://eos-docs.web.cern.ch/eos-docs/

Support   email: eos-support@cern.ch