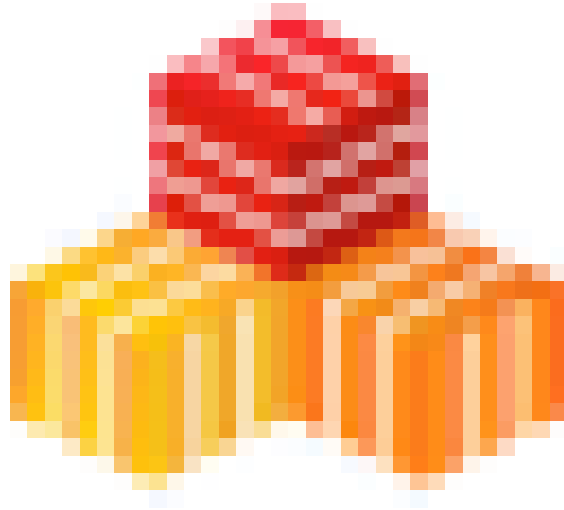


# EOS 2023 Workshop

Monday 24 April 2023 - Thursday 27 April 2023

CERN



## Book of Abstracts



# Contents

EOS for Users - how to use the CERN physics storage system most effectively . . . . .	1
EOS deployment at GRIF . . . . .	1
CTA at PIC . . . . .	2
Operation Status of CDS for ALICE experiment . . . . .	2
EOS for Users Report . . . . .	2
Prometheus Monitoring Exporter for EOS . . . . .	3
Handling failed requests . . . . .	3
Disk File Metadata for Tape Files —Migrating, Restoring, Replicating . . . . .	3
Monitoring your EOSCTA deployment - The general recipe . . . . .	4
CTA Status and Future at IHEP . . . . .	4
EOS instance at the Joint Research Centre . . . . .	4
New CTA tape lifecycle . . . . .	5
External tape readers: Integration into CTA and OSM/Enstore cases . . . . .	5
The EOS namespace locking - demystification and optimization . . . . .	6
A tool to visualize EOS data transfers over time . . . . .	6
EOS 5 client rollout . . . . .	6
EOS for ALICE O2 - Evolution and challenges (2022 - 2023) . . . . .	6
Technical challenges of tape instance consolidation at RAL . . . . .	7
Retiring LevelDBs with the new FST attr backend . . . . .	7
CTA at DESY . . . . .	7
EOS Status at IHEP . . . . .	8
EOS 5/6 Roadmap . . . . .	8
EOS 5 Developments . . . . .	8

eosxd/3 - FUSE filesystem using libfuse2/3 . . . . .	9
Local vs Remote - High Performance Benchmarking . . . . .	9
XRootD Development Update . . . . .	9
Antares: the first year in production . . . . .	9
EOS at the Fermilab LHC Physics Center . . . . .	10
EOS for a T2 Storage Element on Kubernetes . . . . .	10
How to enable EOS for tape . . . . .	10
CTA CI: Running a standalone CTA instance with latest kubernetes . . . . .	11
CTA efforts at Fermilab . . . . .	11
A computational storage plugin implemented in EOS to support in-situ data processing on storage servers . . . . .	11
Balancing and Draining Groups with the GroupBalancer and Drainer . . . . .	12
Improving File Scheduling for EOS . . . . .	12
Workshop Introduction . . . . .	13
EOS Operations at CERN . . . . .	13
Fsck to the rescue . . . . .	13
What you wish you knew about ... tokens! . . . . .	13
EOS for Administrators . . . . .	13
EOS GUI - Simple way for EOS management . . . . .	14
EOS-drive for Windows: Architecture and file transferring system . . . . .	15
EOS Windows native client: Overview . . . . .	16
Welcome and Introduction . . . . .	17
CTA Challenges and Roadmap . . . . .	17
Discussion and close-out . . . . .	17
HDFS to EOS migration - Purdue site report . . . . .	18
EOS Storage Element Status at IHEP . . . . .	18
CERNBox, the Scientific Cloud powered by EOS . . . . .	18
Restic CTA at AARNet . . . . .	19
mhVTL : Mark Harvey's Virtual Tape Library . . . . .	19

ATRESYS —Automated Tape REpacking System, a tool for managing CTA repacks and tape lifecycle . . . . . 19

Empowering CERNBox users with self-service restore functionality . . . . . 20

BoF Session - CERNBox and Sync&Share storage solutions . . . . . 20



**EOS Seminars / 1****EOS for Users - how to use the CERN physics storage system most effectively**

**Authors:** Andreas Joachim Peters<sup>1</sup>; Elvin Alin Sindrilaru<sup>1</sup>; Luca Mascetti<sup>1</sup>; Roberto Valverde Cameselle<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Authors:** luca.mascetti@cern.ch, elvin.alin.sindrilaru@cern.ch, roberto.valverde.cameselle@cern.ch, andreas.joachim.peters@cern.ch

The IT storage group provides end-user access to a 700 PB disk storage system: CERN EOS. In this seminar we will explain your possibilities to use EOS storage as a CERN user most effectively for everyone!

**Part 1 : Dive into the EOS eco-system**

We will start with a brief introduction:

*How is the EOS service deployed and segmented? How do you get access to EOS storage and how you can authenticate to the service?*

We will explain the various access interfaces:

- Command line access using the shell
- Using EOS as a filesystem /eos/
- Remote access protocols root:// and https://
- Accessing EOS from applications like ROOT, C++, Python ...
- The CERNBox web interface

*You will learn, how you share access to files, folders or subtrees with your colleagues, how the permission systems of EOS and CERNBox interact, how you get an EOS drive on your Mac, Linux or Windows computer, how you can verify your quota, how you can understand where you use most of your space, how you can access EOS from outside CERN, what are the best access method for applications and many more useful hints for your daily work.*

We will finish with a short list of features, which are configured to mitigate user errors and service downtimes:

- Service & Data high-availability model
- Backup system
- File Versioning
- Undo Deletion using the EOS recycle bin

*If you deleted all your files, how can you get them back? Can you?*

**Part 2: Running workflows using EOS storage**

This part will cover best practices for running interactive and batch workflows using EOS storage with few examples on a laptop/desktop, lxplus, the batch farm etc.

*What can you do to get efficient data access and what you should never do! How can you authenticate from GITLAB to EOS?*

We will also briefly give some insights, how EOS service managers might influence or change your access to EOS.

**EOS Operation & Sites / 2****EOS deployment at GRIF**

**Author:** Emmanouil Vamvakopoulos<sup>1</sup>

<sup>1</sup> *Université Paris-Saclay (FR)*

**Corresponding Author:** emmanouil.vamvakopoulos@ijclab.in2p3.fr  
to.be.specified

**CTA / 3**

## CTA at PIC

**Author:** Elisabet Carrasco Santos<sup>1</sup>

**Co-author:** Jordi Casals Hernandez<sup>2</sup>

<sup>1</sup> *PIC-IFAE*

<sup>2</sup> *Port d'Informació Científica*

**Corresponding Authors:** jcasals@pic.es, elisabet@pic.es

At PIC, we currently have Enstore as our tape storage system, but due to the discontinuation of its support and development in the near future, we want to share our experiences and insights about the testing and implementation of Cern Tape Archive (CTA) as a potential replacement.

We have set up a CTA test instance integrated with dCache in order to evaluate its functionalities and work on how to adapt our preexisting tape infrastructure and design. Our goal is to provide some valuable information to the CTA community, including our experience, thoughts and any challenges we have encountered during this process, future steps and plans.

**EOS Operation & Sites / 4**

## Operation Status of CDS for ALICE experiment

**Author:** Sang Un Ahn<sup>1</sup>

**Co-authors:** Heejune Han<sup>1</sup>; Jeongheon Kim<sup>2</sup>; Seung Hee Lee<sup>3</sup>

<sup>1</sup> *Korea Institute of Science & Technology Information (KR)*

<sup>2</sup> *Korea Institute of Science and Technology Information*

<sup>3</sup> *KiSTi Korea Institute of Science & Technology Information (KR)*

**Corresponding Authors:** seung.hee.lee@cern.ch, sang.un.ahn@cern.ch, jh.kim@kisti.re.kr, hjhan@kisti.re.kr

We present the current operation status of CDS (the Disk-based Custodial Storage) for ALICE experiment. The CDS is based on EOS Erasure Coding implementation with four parity mode to match with Tape based archival storage in terms of data protection. We will discuss briefly the plan of CDS operation automation for hardware intervention, especially the disk replacement, and of its expansion to meet the upcoming pledges. Also we will discuss the cost analysis of CDS for long-term basis.

**EOS Ecosystem / 5**

## EOS for Users Report

**Authors:** Emmanouil Bagakis<sup>1</sup>; Roberto Valverde Cameselle<sup>1</sup>



<sup>1</sup> CERN

**Corresponding Authors:** emmanouil.bagakis@cern.ch, roberto.valverde.cameselle@cern.ch

EOS for users service, internally known as EOSHPM (EOSHOM, EOSPROJECT and EOSMEDIA) currently stores 2.8 billion files and more than 20PB of storage. We store data of more than 45,000 users and project spaces and host multimedia related use cases for the IT Department. Data is accessed via filesystem (fuse), CERNBox (Web interface, Sync/Mobile client), SAMBA and HTTP. We will be reporting on achievements and challenges in 2022 and future service Roadmap.

EOS Ecosystem / 6

## Prometheus Monitoring Exporter for EOS

**Author:** Roberto Valverde Cameselle<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** roberto.valverde.cameselle@cern.ch

Prometheus is an open-source systems monitoring and alerting toolkit originally built at SoundCloud. Since its inception in 2012, many companies and organizations have adopted Prometheus, and the project has a very active developer and user community. CERN EOS operations team have developed a Prometheus exporter for EOS, that exposes common EOS metrics in prometheus format. This presentation will give an overview of the EOS Exporter, how to set it up and what kind of information can be visualized.

CTA / 7

## Handling failed requests

**Author:** Volodymyr Yurchenko<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** volodymyr.yurchenko@cern.ch

The CERN Tape Archive (CTA) is a vital system for storing and retrieving data at CERN. However, the reliability of the CTA system can be impacted by various factors, including hardware failures, software bugs, and network connectivity issues. To ensure the continued availability of the stored data, it is critical to have robust mechanisms in place for handling failed requests.

This talk will cover the strategies employed by the CTA team at CERN for managing failed requests in the system. These include techniques such as automatic retries, requests classification and file reinjection. By implementing these measures, the CTA team ensures the majority of failed files are archived to tape without user intervention during incidents and in normal operations mode.

CTA / 8

## Disk File Metadata for Tape Files —Migrating, Restoring, Replicating

**Author:** Lasse Tjernaes Wardenauer<sup>1</sup>

<sup>1</sup> *Norwegian University of Science and Technology (NTNU) (NO)*

**Corresponding Author:** lasse.tjernaes.wardenaer@cern.ch

When updating the disk file meta data for tape files, it is necessary to do the updates in both EOS and CTA. Examples of use-cases that require these updates are migration to CTA, moving a file from one EOS instance to another, switching from single to dual copy and restoring deleted files.

The tools for handling these use-cases are not atomic, but they are idempotent and consistency is monitored. As a result, multiple executions by the operator might be necessary to ensure that EOS and CTA agree on the metadata. In the presentation we will show the steps for setting up these tools, as well as the workflow when using them.

CTA / 9

## Monitoring your EOSCTA deployment - The general recipe

**Author:** Richard Bachmann<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** richard.bachmann@cern.ch

Reliable and effective monitoring is essential for smooth operations and for tailoring an EOSCTA deployment to users' needs. Short-term monitoring provides alerting for abnormal system states, and long-term monitoring allows us to track system usage and performance over time.

In this presentation we walk you through the general setup we use for Tier-0 storage at CERN, which allows us to monitor multiple large EOSCTA MGM instances, more than 200 tape servers, and various other machines.

The tech stack is easy to access and based on the open source technologies, such as Fluentd, InfluxDB, Rundeck, and Grafana. We will give some examples on how to ingest CTA log files, aggregate monitoring data, and how turn these data points into useful metrics.

CTA / 10

## CTA Status and Future at IHEP

**Authors:** Qiuling Yao<sup>1</sup>; Yaodong CHENG<sup>2</sup>; Yaosong Cheng<sup>3</sup>; Yujiang BI<sup>4</sup>

<sup>1</sup> *Institute of High Energy Physics Chinese Academy*

<sup>2</sup> *IHEP, Beijing*

<sup>3</sup> *Institute of High Energy Physics Chinese Academy of Sciences, IHEP*

<sup>4</sup> *Institute of High Energy Physics, Chinese Academy of Sciences*

**Corresponding Authors:** biyujiang@ihep.ac.cn, ycheng@cern.ch, yaoql@ihep.ac.cn, chengys@ihep.ac.cn

We will share our experiences on EOS CTA and talk about our plan for the future of CTA. All experiments of IHEP have adopted CTA as the main tape storage management system, and preparing a new tape library for TIER1 of LHCb. We've test the tape restful API with X509 and token auth to access EOS & CTA via HTTP as well as XRootD. In the future, we shall upgrade our production instances to EOS & CTA 5, and deploy FTS services for CTA data transmission.

**EOS Operation & Sites / 11**

## EOS instance at the Joint Research Centre

**Author:** Armin Burger<sup>1</sup>

**Co-author:** Franck Eyraud<sup>1</sup>

<sup>1</sup> *JRC*

**Corresponding Authors:** armin.burger@ec.europa.eu, franck.eyraud@ext.ec.europa.eu

The Joint Research Centre (JRC) of the European Commission is running the Big Data Analytics Platform (BDAP) to enable the JRC projects to store, process, and analyze a wide range of data. The platform evolved as a core service for JRC scientists to produce knowledge and insights in support of EU policy making.

EOS is the main storage system of the BDAP for scientific data. It is in increasing use at JRC since 2016. The Big Data Analytics Platform is actively used by more than 70 JRC projects, covering a large variety of data analytics activities. The EOS instance at JRC has currently a gross capacity of 30 PB with an additional increase planned throughout 2023.

The presentation will give an overview about EOS as storage back-end of the Big Data Analytics Platform. It covers the general set-up and current status, experiences made, issues discovered, and an outlook of planned activities and changes in 2023.

CTA / 12

## New CTA tape lifecycle

**Author:** Joao Afonso<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** joao.afonso@cern.ch

As the amount of stored data, scope and operation complexity of CTA grew, it became necessary to improve the level of control over each tape lifecycle and to provide mechanisms that allow for an improved automation of CTA operations, such as repacking.

In this presentation we will talk about the new CTA tape states, which expand the behaviour of the disabled state. This includes support for repacking tapes as a separated state. In addition, we created a new mechanism that automatically reschedules user requests after a tape state change that no longer allows it to be used for user retrieves.

CTA / 13

## External tape readers: Integration into CTA and OSM/Enstore cases

**Author:** Jorge Camarero Vera<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** jorge.camarero.vera@cern.ch

At the EOS Workshop 2022 BoF, it was decided that CTA should add support for reading OSM/dCache and Enstore tape formats. To make this feature work seamlessly within CTA, we refactored our codebase to accommodate different tape file readers.

In this presentation, we will discuss the design and implementation of external tape format readers into CTA. We will also cover the unit and functional tests that were implemented, including testing via CTA's CI system using an image of an OSM tape.

**EOS Development / 14**

## **The EOS namespace locking - demystification and optimization**

**Author:** Cedric Caffy<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** cedric.caffy@cern.ch

[To be filled]

**EOS Ecosystem / 15**

## **A tool to visualize EOS data transfers over time**

**Author:** Cedric Caffy<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** cedric.caffy@cern.ch

[To be filled]

**EOS Ecosystem / 16**

## **EOS 5 client rollout**

**Author:** Manuel Reis<sup>1</sup>

<sup>1</sup> *Universidade de Lisboa (PT)*

**Corresponding Author:** manuel.b.reis@cern.ch

A report about the deployment of the major version update of the eos client stack. From which XRootD v5 and Fuse v3 upgrades stand out.

**EOS Operation & Sites / 17**

## **EOS for ALICE O2 - Evolution and challenges (2022 - 2023)**

**Authors:** Andreas Joachim Peters<sup>1</sup>; Cristian Contescu<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Authors:** cristian.contescu@cern.ch, andreas.joachim.peters@cern.ch

2022 was a critical year which had some operational impact on all systems and services. In this talk we are going to present a few operational decisions, their implication on the EOS storage for the ALICE data taking and how we implemented mitigations by bringing improvements to the operations model as well as to the software stack.

CTA / 18

## Technical challenges of tape instance consolidation at RAL

**Author:** Thomas Byrne<sup>None</sup>

**Co-authors:** Alison Packer ; George Patargias ; Mahalakshmi Agilandamurthy ; Tim Folkes

**Corresponding Authors:** george.patargias@stfc.ac.uk, mahalakshmi.agilandamurthy@stfc.ac.uk, tim.folkes@stfc.ac.uk, alison.packer@stfc.ac.uk, tom.byrne@stfc.ac.uk

At RAL, we intend to consolidate the two CASTOR instances into a single CTA instance with multiple EOS disk buffers, similar to the CERN architecture. The 'WLCGTape' Castor instance has been fully migrated to CTA at RAL, and has been running in production for LHC run 3 data taking.

In preparation for the migration of the 'Facilities' CASTOR instance at RAL onto our CTA instance, there have been various technical hurdles to overcome. The analysis of namespace (archive file ID) clashes between CTA and the remaining CASTOR instance, and the process of resolving clashes in a repeatable and safe manner were of particular note. This talk describes the details of the analysis, and the tooling developed to enable the consolidation of our two tape instances.

EOS Development / 19

## Retiring LevelDBs with the new FST attr backend

**Author:** Abhishek Lekshmanan<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** abhishek.lekshmanan@cern.ch

Until EOS version 5.1.8, FST metadata on FSTs were stored in a leveldb, which was often heavily contended during writes. We added a feature to move the metadata to attributes. With a minimal configuration, we should be able to switch to the new backend and FSTs automatically move from one backend to another at startup. Additionally there is some tooling to inspect all this. We briefly explain all of this so that sites can get ready to use the new backend.

CTA / 20

## CTA at DESY

**Author:** Mwai Karimi<sup>1</sup>

<sup>1</sup> DESY

**Corresponding Author:** mwai.karimi@desy.de

Since early 2021 CTA has been on a test bed at DESY. Having observed no flies in the ointment and seamless integration with dCache in-place, CTA advances to production in 2023. This presentation will give an overview of the current migration and deployment status as well as future plans at DESY.

## EOS Operation & Sites / 21

### EOS Status at IHEP

**Author:** LI Haibo lihaibo<sup>None</sup>

**Co-authors:** Xuantong Zhang<sup>1</sup>; Yujiang BI<sup>2</sup>; 程耀东 chyd

<sup>1</sup> *Chinese Academy of Sciences (CN)*

<sup>2</sup> *Institute of High Energy Physics, Chinese Academy of Sciences*

**Corresponding Authors:** xuantong.zhang@cern.ch, lihaibo@ihep.ac.cn, biyujiang@ihep.ac.cn, chyd@ihep.ac.cn

The Institute of High Energy Physics undertakes many large scientific engineering projects in China. These large scientific projects generate a large amount of data every year and require a computing platform for analysis and processing.

EOS is one of the main storage system at IHEP since 2016. The EOS instance at IHEP has currently a gross capacity of 50 PB. Currently we have deployed 6 instances and will add new experimental instances in the future, such as HERD experiment.

The presentation will give an overview about the deployment status at IHEP, the work we are doing around EOS and the development plans for 2023.

## EOS Ecosystem / 22

### EOS 5/6 Roadmap

**Author:** Andreas Joachim Peters<sup>1</sup>

**Co-author:** Elvin Alin Sindrilaru<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Authors:** elvin.alin.sindrilaru@cern.ch, andreas.joachim.peters@cern.ch

We will give an overview about the development roadmap for 2023 and beyond.

## EOS Development / 23

### EOS 5 Developments

**Author:** Elvin Alin Sindrilaru<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** elvin.alin.sindrilaru@cern.ch

An overview about the developments since the last workshop.

**EOS Development / 24****eosxd/3 - FUSE filesystem using libfuse2/3**

**Author:** Andreas Joachim Peters<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** andreas.joachim.peters@cern.ch

This presentation will highlight the changes and improvements for the EOS filesystem access using libfuse2/3.

**EOS Development / 25****Local vs Remote - High Performance Benchmarking**

**Author:** Andreas Joachim Peters<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** andreas.joachim.peters@cern.ch

This presentation includes performance benchmarks comparing local and remote IO for various use-cases, storage stacks and protocols.

**EOS Ecosystem / 26****XRootD Development Update**

**Author:** Guilherme Amadio<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** amadio@cern.ch

Latest updates about XRootD development and the March XRootD Workshop.

**CTA / 27****Antares: the first year in production**

**Author:** George Patargias<sup>None</sup>

**Co-authors:** Tom Byrne<sup>1</sup>; Mahalakshmi Agilandamurthy<sup>1</sup>; Alison Packer<sup>1</sup>; Alastair Dewhurst<sup>1</sup>

<sup>1</sup> *STFC - UKRI*

**Corresponding Authors:** george.patargias@stfc.ac.uk, tom.byrne@stfc.ac.uk, alastair.dewhurst@stfc.ac.uk, mahalakshmi.agilandamurthy@stfc.ac.uk, alison.packer@stfc.ac.uk

Antares is the new tape archive service at RAL Tier-1 that went into production on 4th of March 2022. The service is built around the EOS/CTA technologies developed at CERN. EOS is the user

facing service that manages the incoming namespace requests and a thin SSD buffer, and CTA is deployed as the tape back-end system. In this talk, we describe the setup of ANTARES and discuss the service's performance as well as the main operational issues since the beginning of LHC Run-3. Finally, we provide an overview of the future plans for the expansion of the service to cover the whole of the high energy physics and astronomy community in the UK.

## EOS Operation & Sites / 28

### EOS at the Fermilab LHC Physics Center

**Author:** Dan Szkola<sup>1</sup>

<sup>1</sup> *Fermi National Accelerator Lab. (US)*

**Corresponding Author:** dszkola@fnal.gov

Fermilab has been running an EOS instance since testing began in June 2012. By May 2013, before becoming production storage, there was 600TB allocated for EOS. Today, there is approximately 13PB of storage available in the EOS instance.

The LPC cluster is a 4500-core user analysis cluster with 13 PB of EOS storage. The LPC cluster supports several hundred active CMS users at any given time.

An update of our current experiences and challenges running an EOS instance for use by the Fermilab LHC Physics Center (LPC) computing cluster. Planning the upgrade to EOS 5 and moving to Almalinux before EOL for SL7.

## EOS Operation & Sites / 29

### EOS for a T2 Storage Element on Kubernetes

**Author:** Ryan Taylor<sup>1</sup>

<sup>1</sup> *University of Victoria (CA)*

**Corresponding Author:** rptaylor@uvic.ca

I will discuss our efforts to deploy EOS on Kubernetes at the University of Victoria T2 site for ATLAS, using a Helm chart and CephFS storage.

## CTA / 30

### How to enable EOS for tape

**Author:** Julien Leduc<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** julien.leduc@cern.ch

An EOSCTA instance is an EOS instance - commonly called a tape buffer - configured with a CERN Tape Archive (CTA) back-end.

This EOS instance is entirely bandwidth oriented: it offers an SSD based tape interconnection, it can



contain spinning disks if needed and it is optimized for the various tape workflows. This talk will present how to enable EOS for tape using CTA and the Swiss horology gears in place to maximize tape hardware usage while meeting experiment workflow requirements for xrootd and HTTP protocols.

CTA / 31

## CTA CI: Running a standalone CTA instance with latest kubernetes

**Author:** Julien Leduc<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** julien.leduc@cern.ch

This hands-on session will focus on installing and configuring a standalone CTA CI runner:

- single host kubernetes cluster in Alma9
- 1 Virtual tape library with CTA CI requirements

At the end of this sessions the participants should be able to run CTA Continuous Integration test on their box.

CTA / 32

## CTA efforts at Fermilab

**Authors:** Eric Vaandering<sup>1</sup>; Ren Bauer<sup>1</sup>; Scarlet Norberg<sup>None</sup>

<sup>1</sup> *Fermi National Accelerator Lab. (US)*

**Corresponding Authors:** renbauer@fnal.gov, ewv@fnal.gov, norberg@ou.edu

Fermilab has decided to replace Enstore, its locally developed tape management system, with CTA. Fermilab runs two Enstore instances: CMS with a small, dedicated tape buffer and Public operating like an HSM with tight integration between dCache and Enstore

This talk will cover:

- Metadata migration from Enstore to CTA
- Results of dCache integration with CTA at FNAL
- Performance testing and monitoring
- Thoughts on integrating Enstore's "small file aggregation" (SFA) feature (Public instance only)
- Timeframe for migrations

Reading of physical Enstore tapes is discussed in Jorge's OSM talk

EOS Operation & Sites / 33

## A computational storage plugin implemented in EOS to support in-situ data processing on storage servers

**Author:** Yaodong Cheng<sup>1</sup>

**Co-author:** Haibo LI

<sup>1</sup> *Institute of High Energy Physics, Chinese Academy of Sciences*

**Corresponding Authors:** lihaibo@ihep.ac.cn, chyd@ihep.ac.cn

Computational storage involves integrating compute resources with storage devices or systems to enable data processing within the storage device. This approach reduces data movement, enhances processing efficiency, and reduces costs. To facilitate in-situ data processing on storage servers, we developed a computational storage plugin that can be added to EOS FST. This plugin enables users to deploy compute resources directly within the storage servers, allowing them to perform data processing operations on the data stored in the FST nodes without having to move the data to a separate computing system. This can reduce latency and improve overall performance, especially when processing large volumes of data.

The plugin can be extended to support a variety of data processing tasks, including data filtering, compression, encryption, and machine learning. The computational storage function is defined in a configuration that can be implemented in scripting languages or evolved independently of the storage system in the form of containers.

When an FST node receives a request to open a file, the plugin is executed first. It then calls the target program on the storage server by parsing the parameters of the command to open the file. At this time, the input file must be on the FTS storage server, and the plugin also writes the output file to the node. At the end of the task execution, the output file is automatically registered into the MGM server.

Client access is fully compatible with XRootD's API and EOS commands. Users can add tasks and parameters to be performed in the open option. The plugin has been tested and applied in the data processing of the Large High Altitude Air Shower Observatory (LHAASO), and the results show that the efficiency of data decoding is more than 5 times higher than the original method.

**EOS Development / 34**

## Balancing and Draining Groups with the GroupBalancer and Drainer

**Author:** Abhishek Lekshmanan<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** abhishek.lekshmanan@cern.ch

New improvements in the existing GroupBalancer and introduction of functionality to drain whole groups. We look at the various configuration options to run these and how these work under the hood.

**EOS Development / 35**

## Improving File Scheduling for EOS

**Author:** Abhishek Lekshmanan<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** abhishek.lekshmanan@cern.ch

We take a look at the geoscheduler and see how we can introduce a new lock-free scheduling algorithm

**EOS Development / 36**

## Workshop Introduction

**Corresponding Author:** andreas.joachim.peters@cern.ch

Logistics & Information.

**EOS Ecosystem / 37**

## EOS Operations at CERN

**Corresponding Author:** maria.arsuaga.rios@cern.ch

**EOS Development / 39**

## Fsck to the rescue

**Author:** Elvin Alin Sindrilaru<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** elvin.alin.sindrilaru@cern.ch

This talk will describe the fsck mechanism, the various options when it comes to controlling the repair process and the internal process of deciding whether a file can be fixed or not.

**EOS Development / 40**

## What you wish you knew about ... tokens!

**Author:** Elvin Alin Sindrilaru<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** elvin.alin.sindrilaru@cern.ch

This presentation gives an overview of the token support in EOS. We'll discuss the configuration options, what plugins need to be enabled for the various protocols and how to configure them. Besides this, we'll trace one particular request using tokens to see how it interacts with the existing authentication/authorization features that already exists in EOS and provide some helpful examples.

**EOS Seminars / 41****EOS for Administrators**

**Corresponding Authors:** elvin.alin.sindrilaru@cern.ch, michael.davis@cern.ch, andreas.joachim.peters@cern.ch

In this seminar we will go through the architecture of EOS, showcase some EOS instance configuration and follow with an introduction to CTA!

We will explain some generic concepts, deployment models, hardware requirements, redundancy models, storage layout, scheduling & file placement - CPU, Storage and Network requirements and few tricks to optimize these. We highlight in detail the XRootD framework and how EOS is implemented using it. We will go through all components and sub services of the namespace service MGM, what their role is and point you to the relevant documentation.

The second part will showcase some EOS instances and some of their configuration at CERN.

In the third part of the presentation we will introduce the CERN Tape Archive CTA.

We hope there is enough time available to keep the seminar interactive and allow for questions during the presentation.

**EOS Ecosystem / 42****EOS GUI - Simple way for EOS management**

**Author:** Gregor Molan<sup>1</sup>

**Co-authors:** Branko Blagojevic<sup>2</sup>; Ivan Arizanovic<sup>1</sup>

<sup>1</sup> *Comtrade 360's AI Lab*

<sup>2</sup> *Comtrade*

**Corresponding Authors:** ivan.arizanovic@cern.ch, branko.blagojevic@cern.ch, gregor.molan@cern.ch

The Graphical User Interface (GUI) for CERN EOS could be crucial in the interaction between potential users and the EOS storage technology. The GUI could serve as an interface between a user and the complex EOS infrastructure, enabling non-experts to learn and discover EOS features seamlessly and effectively. This would help users interact with the storage infrastructure without needing to delve into technical complexities, making it easier to decide on the architecture and proposal for large storage organizations. Additionally, the GUI can provide a visually appealing and user-friendly interface that can enable users to carry out informative tasks such as monitoring data storage usage.

The first proposal for CERN EOS GUI is designed for the Microsoft Windows platform. The proposal specifies technologies for GUI that allow extensions to all major operating systems, resulting in an interface accessible to all users regardless of their preferred platform, including Linux and MacOS. The GUI focuses on usability and functionality, featuring intuitive navigation, clear labeling of buttons and controls, and informative feedback mechanisms. The goal is to cover three functionalities related to EOS Windows Native Client:

- a) Interface to EOS commands (EOS-shell).
- b) Interface to EOS cluster mounted as Windows drive letter (EOS-drive).
- c) Interface to all functionalities covered by EOS commands.

EOS GUI is intended to provide two ways of starting:

- 1) Start a program executable from the start menu or command line.
- 2) Start from a popup list from the system tray icon.

Thus, starting the same program can be made available from multiple entry points, making it easier for users to access the GUI.

The EOS GUI should be organized into three window forms:

- a) Main EOS window.
- b) Popup list from the system tray icon.
- c) Other popup windows.

The main window would cover all EOS features available with EOS commands. The format of the main EOS window is proposed as tabs, icons, or “office style”. For each of these proposed formats, all EOS features should be grouped in three to five groups represented by separate tabs, icons, or “office-style” buttons. The popup list from the system tray icon should provide shortcuts to frequently used functions, such as connecting/disconnecting EOS drives and opening the main EOS window. Other popup windows should cater to specific functionalities, such as showing detailed storage usage.

The GUI for EOS should be technically implemented as a program/application or as a Web GUI. Overall, the proposed EOS GUI should provide a user-friendly and accessible interface that allows users to carry out tasks related to EOS commands and data storage. The proposed EOS GUI is designed to provide an efficient and accessible interface for users of all levels to perform tasks related to EOS commands and data manipulation. It is not designed to replace EOS CLI but to complement it while providing a more user-friendly alternative for those less comfortable with command-line interfaces.

In summary, the proposed EOS GUI aims to provide a user-friendly interface for managing CERN EOS storage. The GUI is proposed for the Microsoft Windows platform; however, the technologies used would allow extensions to all major operating systems.

**EOS Ecosystem / 43**

## **EOS-drive for Windows: Architecture and file transferring system**

**Author:** Ivan Arizanovic<sup>1</sup>

**Co-authors:** Branko Blagojevic<sup>2</sup>; Gregor Molan<sup>1</sup>

<sup>1</sup> *Comtrade 360's AI Lab*

<sup>2</sup> *Comtrade*

**Corresponding Authors:** [ivan.arizanovic@cern.ch](mailto:ivan.arizanovic@cern.ch), [gregor.molan@cern.ch](mailto:gregor.molan@cern.ch), [branko.blagojevic@cern.ch](mailto:branko.blagojevic@cern.ch)

EOS-drive is part of the EOS Windows Native Client package, it mounts the EOS filesystem as a Windows disk drive by which Windows applications interact with the EOS filesystem.

EOS-drive communicates with Windows applications through the user-mode Dokan library and kernel-mode Dokan driver. File operation requests from applications (e.g., CreateFile, ReadFile, WriteFile...) are sent to the Dokan driver, which is then forwarded to the Dokan library and subsequently to the EOS-drive. The results of this routine are sent back to the Windows application as a response to the operation request. The Dokan file system driver and library therefore acts as a proxy between Windows applications and the EOS-drive.

To collect data that should be sent to the Windows application as a response to an operation request, EOS-drive communicates with the EOS cluster through HTTPS protocol. EOS-drive uses the cURL library for file transferring and gRPC for metadata requests.

WriteFile and ReadFile operations are used for file transferring. Windows applications requests a process over chunks (data fragments) by providing offset and chunk size to the EOS-drive. Then, EOS-drive is made to serve Windows applications by performing WriteFile or ReadFile operations over specified chunks.

There are three mechanisms of file transferring in EOS-drive, these are: transfer chunk by chunk, transfer with the EOS-drive buffer, and single-session transfer.

Transferring chunk by chunk means that each requested chunk by the Windows applications will be uploaded/downloaded in a separate session. This method is the slowest but can be realized in any case.

Transferring with an EOS-drive buffer can be done in three different scenarios. In case of upload, EOS-drive buffers some of the chunks received from the Windows driver, sorts them if needed, and uploads them together per second. With regards to downloading “small” files, instead of downloading the requested chunk, EOS-drive will download the whole file and then use downloaded bytes for chunk requests. There is also a buffer between the upload and download process, i.e., the first part of an uploaded file will be buffered and then used for subsequent download requests.

Single-session transfer means that one session will be used for the whole file, not only for one chunk. Afterwards, if chunk requests from the driver match the incoming chunk from the EOS cluster, the chunk from the EOS cluster will then only be forwarded to the driver. Or in the opposite direction for the upload process. This method is the fastest, however, it has several restrictions.

Functionality and transfer speed are tested locally at Comtrade, as well as at CERN, therefore the best solution for file transfer in EOS-drive is single-session transfer, although this method is impossible in all situations due to many restrictions. The next best solution as per speed, is to implement transfer using the EOS-drive buffer, it is however limited per number of files to save memory resources. Subsequently, if none of the previous mechanisms can be realized, transfer chunk by chunk will be used.

**EOS Ecosystem / 44**

## **EOS Windows native client: Overview**

**Author:** Branko Blagojevic<sup>1</sup>

**Co-authors:** Gregor Molan<sup>2</sup>; Ivan Arizanovic<sup>2</sup>

<sup>1</sup> *Comtrade*

<sup>2</sup> *Comtrade 360's AI Lab*

**Corresponding Authors:** gregor.molan@cern.ch, branko.blagojevic@cern.ch, ivan.arizanovic@cern.ch

### **Context**

An overview of the EOS Windows native client for EOS users on Windows operating systems

### **Objectives**

EOS Windows native client should provide Windows platform users with native access to EOS cluster for both file transferring and command requests, giving them improved user experience compared to EOS Linux client.

### **Method**

EOS Windows native client comes with two interfaces - EOS-shell (EOS-wnc command line interface) which can be run inside Command Prompt or PowerShell terminal, and EOS-drive which mounts EOS file system as Windows drive. These interfaces bring EOS file system to Windows users with the user-friendly experience they expect and require.

EOS-shell supports all the commands as EOS Linux client, together with some improvements in terms of uniform manuals for all commands, additional checks for specified parameters and auto-complete functionality for commands, command arguments and paths inside EOS space. This means that a user, specifically used to EOS Linux client can easily make the switch to EOS Windows native client seamlessly.

EOS-drive provides classic Windows experience by representing EOS file system as a drive on Windows, meaning that any user can easily manage files and directories inside EOS space, as if they

are on Windows machine itself, with no need for knowledge of how to use command line interface.

By constantly testing EOS Windows native client, in Comtrade's testing environment and against large EOS instances at CERN, alongside comparing it with other ways to access EOS file system and even with different distributed file systems, we are ensuring that EOS users on Windows will not lag behind regarding performances.

### **Conclusion**

Developing and maintaining high-performance and user-friendly EOS client for Windows platforms should always be a priority while seeking to constantly provide possible improvements that ensure better user experience. Additionally, Windows users should be given the opportunity to benefit from EOS file system within environments they are familiar with.

CTA / 45

## **Welcome and Introduction**

**Author:** Michael Davis<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** michael.davis@cern.ch

Welcome to the second annual CTA Day at the EOS Workshop!

In mid-2022, the LHC awoke from its second Long Shutdown and restarted physics operations. Run-3 data-taking rates are several times higher than during Run-2; already CTA hit a new record of 26 PB archived in one month in November 2022.

This presentation will introduce the CTA Project, Team and Community, as well as an overview of the challenges and achievements during the first year of Run-3.

CTA / 46

## **CTA Challenges and Roadmap**

**Author:** Michael Davis<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** michael.davis@cern.ch

CTA software development is primarily driven by the needs of the CERN experimental programme. Looking beyond Run-3, data rates are set to continue to rise exponentially into Run-4 and beyond. The CTA team are planning how to scale the software and service to meet these new challenges.

CTA is also driven by the needs of the community outside CERN. The landscape of tape archival for scientific data is consolidating, and CTA is constantly adapting to a wider range of use cases.

This talk will present the short-term and medium-term roadmap for CTA development and new features.

CTA / 47

## Discussion and close-out

**Author:** Michael Davis<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** michael.davis@cern.ch

Final comments, questions and discussion. Segue into the apéro in R2 where we can continue talking.

**EOS Operation & Sites / 48**

## HDFS to EOS migration - Purdue site report

**Author:** Stefan Piperov<sup>1</sup>

<sup>1</sup> *Purdue University (US)*

**Corresponding Author:** stefan.piperov@cern.ch

In 2022 the CMS Tier-2 at Purdue University migrated its 10PB storage system from HDFS to EOS. Here we report on the details of the process, the difficulties we encountered and the ways in which we solved them. We also report on the current status of the storage system, and our future plans.

**EOS Operation & Sites / 49**

## EOS Storage Element Status at IHEP

**Author:** Xuantong Zhang<sup>1</sup>

**Co-authors:** LI Haibo lihaibo ; Yaodong CHENG<sup>2</sup>; Yujiang BI<sup>3</sup>

<sup>1</sup> *Chinese Academy of Sciences (CN)*

<sup>2</sup> *IHEP, Beijing*

<sup>3</sup> *Institute of High Energy Physics, Chinese Academy of Sciences*

**Corresponding Authors:** xuantong.zhang@cern.ch, ycheng@cern.ch, lihaibo@ihep.ac.cn, biyujiang@ihep.ac.cn

The EOS system serving as a grid storage element at IHEP, CAS started since 2021, working for JUNO experiment. A CTA with its EOS SE buffer also started its service for JUNO since 2023. In this talk, we would like to share our experiences and thoughts about the SE operations, including deployment, monitoring, data transfer performance, authentication management with VOMS and Sci-token, etc. Meanwhile, as EOS SE will replace the DPM as our new Beijing-T2 site storage system, this talk will also share our plan and status about EOS upgrading.

The Beijing LHCb T1 site storage construction status will also be included in this talk.

**EOS Ecosystem / 50**

## CERNBox, the Scientific Cloud powered by EOS

**Authors:** Diogo Castro<sup>1</sup>; Hugo Gonzalez Labrador<sup>1</sup>



<sup>1</sup> CERN

**Corresponding Authors:** hugo.gonzalez.labrador@cern.ch, diogo.castro@cern.ch

CERNBox combines the ease of use of a file sync and share service with the power of the scientific data processing infrastructure at CERN. Built on top of EOS and ownCloud, it provides a simple and uniform way to access over 15PB of research, administrative and engineering data across more than 2 billion files.

In this talk we will go through the latest advances made possible with the new version, released in 2022, and the new functionalities planned for this year and the future. From notifications to search or from better sync client integration to federation of heterogeneous storages, both CERNBox and EOS are evolving together to provide a more powerful and user friendly system for our community.

CTA / 51

## Restic CTA at AARNet

**Author:** Denis Sergeevich Lujanski<sup>None</sup>

**Corresponding Author:** denis.lujanski@aarnet.edu.au

Since 2021, AARNet have been using restic backup software in conjunction with CTA to backup user data in production EOS clusters. The road to production has not been without its challenges, requiring us to modify restic and create a custom backup scheduler and client workflows. This presentation will aim to cover the architecture, with a focus on restic: the basics, the customisations and integration with CTA.

CTA / 53

## mhVTL : Mark Harvey's Virtual Tape Library

**Author:** Mark Harvey<sup>None</sup>

**Corresponding Author:** markh794@gmail.com

Virtual Tape Libraries (VTLs) have been widely used in data storage environments. mhVTL design goals are unique in it is not attempting to be a better tape library, simply emulate real hardware—for those situations where it is impractical to carry a physical tape library with you.

CTA / 54

## ATRESYS — Automated Tape REpacking System, a tool for managing CTA repacks and tape lifecycle

**Author:** Vladimir Bahyl<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** vladimir.bahyl@cern.ch

A 'Repack' is the process of moving or copying data from one tape cartridge to one or multiple others. Such a process may be needed for various reasons, such as transferring data to more compact media,

creating additional copies, and recovering data from faulty media. At CERN the CTA software manages the data transfer itself, but more steps are needed in order to complete the full repack process for each tape cartridge.

In this presentation we will give an overview of the repack process at CERN and present ATRESYS, the tool developed to automate most of it. By using ATRESYS we are able to queue batches of hundreds of repacks and have them run their course without micro-management by operators. The tool will soon be available as free software as part of the CTA operator tools.

**EOS Ecosystem / 55**

## **Empowering CERNBox users with self-service restore functionality**

**Corresponding Author:** gianmaria.del.monte@cern.ch

The IT storage group at CERN is responsible to ensure integrity and security of all the stored data for physics and general computing services. In the last years a backup orchestrator, cback, has been developed based on the open source backup software restic. Cback is able to backup EOS, CephFS and any local mountable file system, like NFS or DFS. cback is currently used to daily backup CERNBox data (2.5 billion of files and 18PB), including experiment project spaces and user home directories.

The data copy is stored in a disk-based S3 cluster in another geographical location in the CERN campus 4km away from the main data center (protecting against natural disasters). The usage of restic allows us to reduce the storage costs thanks to the deduplication of the data. In the last months, the cback portal server has been implemented, exposing a set of REST APIs to allow the integration with end-user backup utilities to navigate snapshots and restore data.

In this presentation, we will describe the architecture and the implementation of cback, the integration with CERNBox and the future integration with tape archive (CTA) for long term data preservation.

**BoF Session - CERNBox and Sync&Share storage solutions / 56**

## **BoF Session - CERNBox and Sync&Share storage solutions**

**Corresponding Author:** hugo.gonzalez.labrador@cern.ch