# Overview

- **Introduction**
  - what is EOS?
  - history
  - architectural evolution

- EOS **service** at CERN
  - dimension & challenges

- EOS in a science **ecosystem**
  - EOS, CERNBox & SWAN

- EOS as a **filesystem**

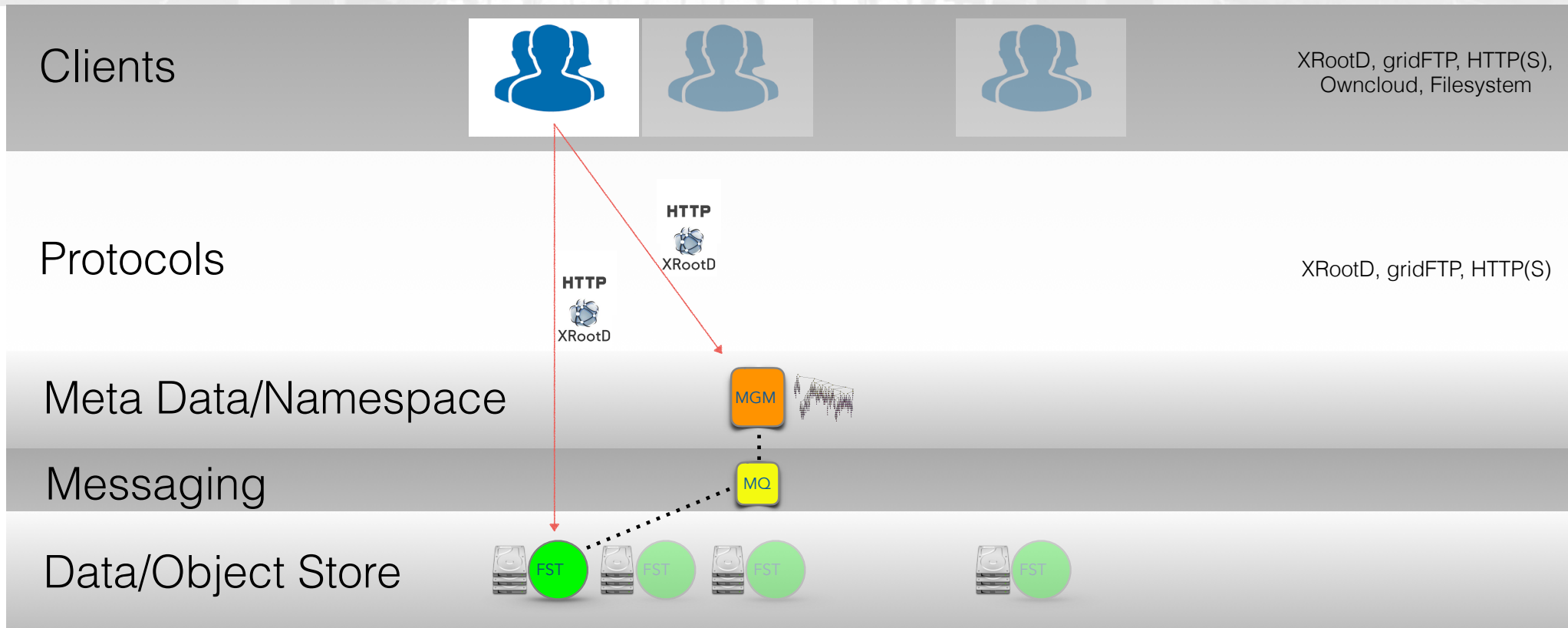- Evolution **data processing** - object storage models

# What is EOS ?

- **disk storage system** designed to serve **physics analysis** use cases
  high concurrency, pseudo-random access, LAN/WAN clients

- implemented as plug-ins into the **XRootD** framework

- native transport protocol is **XRootD** [ optimized for latency compensation ]

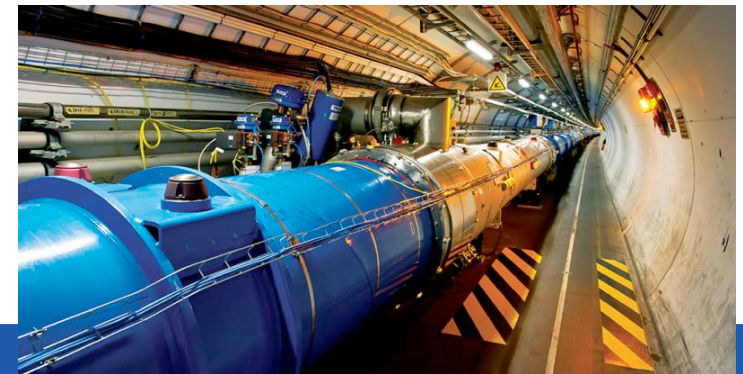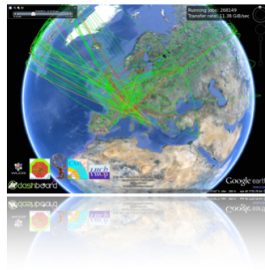- code is written in C++ in IT-ST group at CERN

# LHC Use Case



CERN mainstream use case

tape archive **CASTOR** CERN Advanced STORage manager

LHC Detector

local batch cluster

batch processes O(10^5)

O(GB/s)

5-10 GB/s

2017 1y averages
50k reader - 6k writer
**36 GB/s** - peak**120 GB/s**

openstack™

5-10 GB/s

Data Export to Worldwide Computing Grid

2017 1y averages
50k reader - 6k writer
**36 GB/s read** - peak **120 GB/s**

# Project History

timeline

2009
2010

architecture document

EOSATLAS
EOSCMS
EOSALICE
EOSLHCB
EOSPUBLIC

EOS Open Source Storage

CASTOR
XRootD

lustre
XRootD

castor-xrootd IF

lustre-xrootd IF

2011
IPv4

2012

2013

v **0.2**
*AMBER*

v **0.3**
*BERYL*

2014

v **0.3.x**
*AQUA MARINE*

EOSUSER
2015

2016

EOSMEDIA

2017

IPv6

v **4.0**
*CITRINE*

EOSALICEDAQ
2018 EOSCTAATLASpps
EOSHOME

v **4.3**
Quark$^{db}$
RocksDB

sync&share — swan — wopi — tape
CERNBox   SWAN   Office

SEAGATE KINETIC OPEN STORAGE PLATFORM — kinetic R&D

librados R&D

eXtreme DataCloud — QOS, caches & data lakes

CERN

EOS Disk Space at CERN

| 100 PB | | | | | |
|---|---|---|---|---|---|
| 2 PB | 12 PB | 40 PB | 150 PB | 250 PB | |
| 1 PB | | | | | |
| 2010 | 2012 | 2014 | 2016 | 2018 | |

# Introduction
# Architectural evolution

2010-2017                    2017++

Beryl Aquamarine **V 0.3**   →   Citrine **V 4**

**stateful in-memory**

META DATA

read/write | MGM Master | MGM Slave | read only   stateless   active MGM | standby MGM | standby MGM   **stateless + cache in-memory**

**scale-up**

Persistency

DATA

FST FST FST FST   **scale-out**   FST FST FST FST

KINETIC Open Storage Project   ceph

Parallel Sessions
Scaling the EOS namespace

**meta data service daemon becomes stateless**

# Introduction
# Architectural evolution

FS Clients

**external** entity
uni-directional communication

MGM

FST FST FST FST

FS Client

**internal** entity
bi-directional communication

MGM

FST FST FST FST

better support of filesystem semantics requires
FS clients receiving call-backs

# EOS service at CERN

designed as a 'lossy' service with two replica CERN/Wigner file replication

2017
statistics

| | |
|---|---|
| bytes read | 1.00 EB/a |
| bytes written | 0.25 EB/a |
| disk IO | 7.90 EB/a |
| hard disks | ~ 50k |
| streams | ~ 55k |
| fileloss rate | ~ $O(10^{-6})$/a |

Parallel Sessions
Providing large-scale disk storage at CERN

Parallel Sessions
Disk failures in the EOS setup at CERN

| Number of Files | Number of Directories | Total Space | Free Space | MGM # of open FDs |
|---|---|---|---|---|
| 2.969 Bil | 181 Mil | 244 PB | 66.2 PB | 44775 |

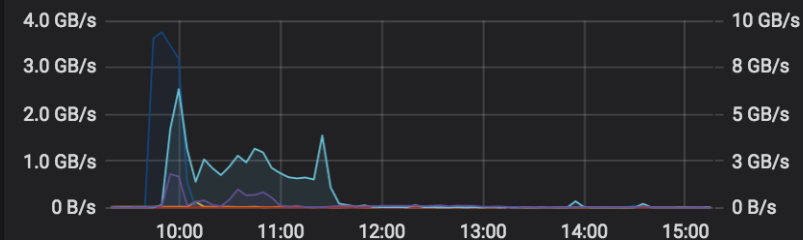| Current Writers | Current Readers | IOPS | Write Throughput | Read Throughput |
|---|---|---|---|---|
| 27.0 K | 62.5 K | 188 K | 3.94 GB/s | 35.2 GB/s |

LHC data taking

ATLAS Point1  Max: 3.75 GB/s  Avg: 204 MB/s  Current: 14 kB/s
CMS Point5  Max: 124 MB/s  Avg: 9 MB/s  Current: 7 MB/s
ALICE Point2  Max: 0 B/s  Avg: 0 B/s  Current: 0 B/s

EOS Disk Space at CERN

| 100 PB | | | | | |
|---|---|---|---|---|---|
| 1 PB | 2 PB | 12 PB | 40 PB | 150 PB | 250 PB |
| | 2010 | 2012 | 2014 | 2016 | 2018 |

# EOS service at CERN

- *cheap* disk storage

  - > 1.300 server, 50k disks

  - JBOD with dual (geo-) replication

  - 48-192 disks per standard head-node (batch server) in production

  - new big server provide **2 PB** storage capacity

  - **BEER** containerised batch jobs on EOS disk server
    EOS storage service daemon light-weight, allows to use 90% of free CPU resources

Parallel Sessions
Sharing server nodes for
storage and compute

cheap volume storage on HDD

CERN

# A new CERNBOX backend
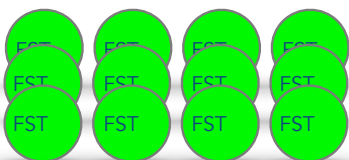## Segmented high-available service model

## CERNBOX now
EOSUSER

GW

1TB RAM    MGM    MGM    600M files

at namespace scalability limit
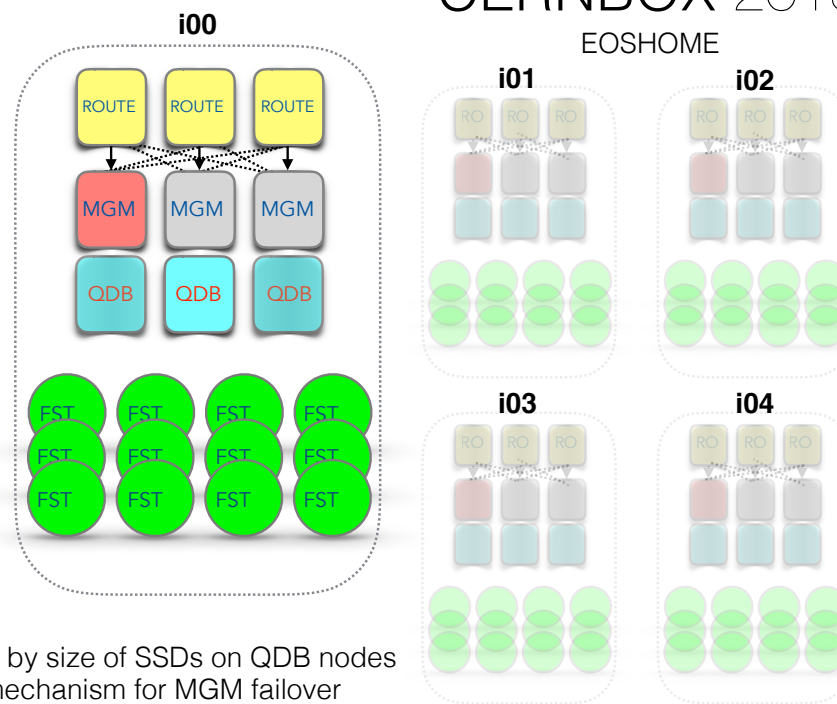availability constrained by infrequent long boot time of 2h

tested with >4B files

**i00**

ROUTE  ROUTE  ROUTE

MGM  MGM  MGM

QDB  QDB  QDB

FST  FST  FST  FST
FST  FST  FST  FST
FST  FST  FST  FST

namespace scalability limit by size of SSDs on QDB nodes
automatic built-in HA mechanism for MGM failover

## CERNBOX 2018
EOSHOME

**i01**          **i02**

RO RO RO        RO RO RO
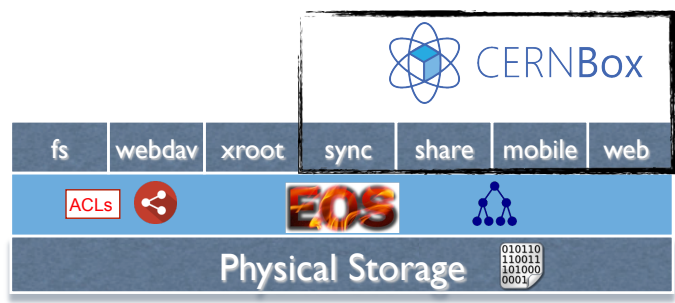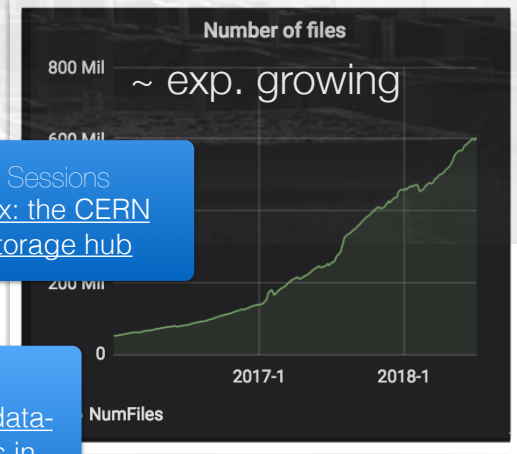
**i03**          **i04**

RO RO RO        RO RO RO

## Since 2017 service running over scalability limit - new service architecture

**2015++**

**EOS**



**CERNBox**

building block of an ecosystem for scientific data repositories

Number of files

~ exp. growing

800 Mil
600 Mil
200 Mil
0

2017-1    2018-1

NumFiles

**CERNBox**

| fs | webdav | xroot | sync | share | mobile | web |

ACLs   EOS

Physical Storage

010110
110011
101000
0001

Poster Sessions
The EU Up to University Project

Parallel Sessions
CERNBox: the CERN cloud storage hub

Parallel Sessions
Cloud Storage for data-intensive sciences in science and industry

Office

since mid 2017 support for collaborative editing

- "**dropbox**" for science
- **cloud storage**, **synchronisation** and **file sharing** service
- implemented as **web services** in front of EOS backend

Daily users & IP addresses

10 K
8 K
5 K
3 K
0

8/1  9/1  10/1  11/1  12/1  1/1  2/1  3/1  4/1  5/1  6/1  7/1

Linux 21%
MacOSX 49%
Windows 30%

**connected client platforms**

**~3.000 daily users, 9k connected devices**

CERN

CERNBox Evolution

2014 | 2015-2017 | 2018 - future
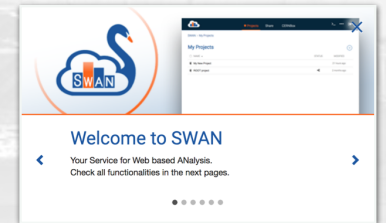
analysis  shares  auth  storage  office  webUI  roles

**CERNBox refactored using micro service approach - boost performance & functionality**

# SWAN service

## web-based analysis

swan.web.cern.ch

2016++

Parallel Sessions
Facilitating collaborative
analysis in SWAN

Welcome to SWAN
Your Service for Web based ANalysis.
Check all functionalities in the next pages.

PS Triple Bunch Splitting

RF Bucket Matching

Transverse Tracking

Detuners

SSO

jupyterhub    Web portal

Container Scheduler

User 1    User 2    ...    User n

EOS
(Data)

CVMFS
(Software)

CERNBox
(User Files)

CERN Resources

service architecture

**SWAN provides interactive analysis front-end using JUPYTER notebooks**

# SWAN & compute

SWAN interfacing to Spark Cluster

# SWAN & compute



SSO

jupyterhub   Web portal

Container Scheduler

docker User 1   docker User 2   ...   docker User n

EOS (Data)   CVMFS (Software)   CERNBox (User Files)

CERN Resources

Notebook Job

EOS

HT Condor

Batch Farm   HTCondor High Throughput Computing

Google Summer of Code

**next: SWAN interfacing to Batch Cluster**

2017++

**EOS**

packaged eco-system

# Science Box

**EOS** + **CERNBox** + **SWAN**

docker-compose

docker

kubernetes

**One-Click Demo Deployment**

- Single-box installation via **docker-compose**
- No configuration required
- Download and run services in 15 minutes

  https://github.com/cernbox/uboxed

**Production-oriented Deployment**

- Container orchestration with **Kubernetes**
- Scale-out storage and computing
- Tolerant to node failure for high-availability

  https://github.com/cernbox/kuboxed

WLCG

**SWAN**
CVMFS Client
EOS Fuse Mount
JupyterHub
Interactive Notebooks

**CERNBox**
CERNBox Backend
CERNBox Gateway
Synchronization and Sharing

**EOS**
File Storage Servers
Management Node
Disk-based Storage

CERN

# Science Box provides an easy demo & production platform

# 2018++  EOS + Tape = EOSCTA

## integrated support for tape into EOS  file on tape=offline replica

- loose service coupling between EOS and CTA via protocol buffer interface & notification events
- no SRM, using XRootD protocol only for now - integrated with FTS
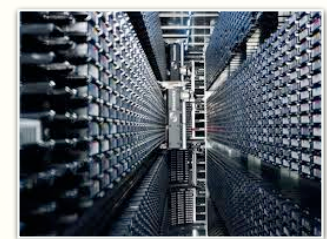- pre-production service for ATLAS available

high disk capacity

low disk capacity

Operation Model



**TPC**

EOSATLAS

EOSATLASCTA

short file lifetime

Cern Tape Archive

## Mid-term plan to migrate CASTOR data to CTA

# **eos**xd - a filesystem client for EOS

Why this is important but difficult …

- mounted filesystem access is required to enable storage access to any software *out of the box*

- filesystem development is difficult and lengthy

  AFS V1,2,3 - **35 years**

  NFS V1,2,3,4 - **34 years**

  cephfs - **12 years** - production version announced after 10 years!

- EOS filesystem client rewrite started Q4 2016: **eos**d => **eos**xd

Question: **how far can you get with a user-space filesystem implementation?**

# **eos**xd

## filesystem daemon

### Architecture



- enough **POSIX**ness
- file **locks**, byte-range locks
- hard **links** within directories
- rich **ACL** client support
- local **caching**
- **bulk deletion**/protection
- strong **security** & mount-by-key
- user, group & project **quota**
- implemented using **libfuse**

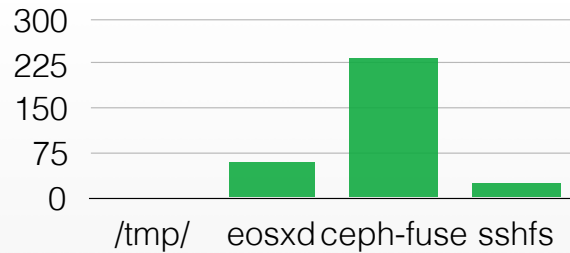**eos**xd **provides POSIXness very similar to AFS**
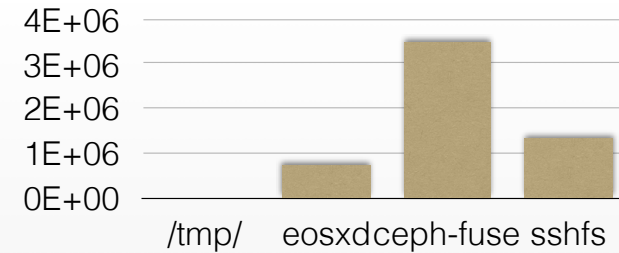
# **eos**xd
## Strong Security Model

application runtime: **export KRB5CCNAME or X509xxx**

kerberos or X509 authentication

ACL per directory by mapped uid/gid

**sys.acl=u:foo:rwx**

before mount: **export XrdSecsssENDORSEMENT=<secret>**

shared secret authentication

ACL per directly by exported secret

**sys.acl=k:B8E776C5-F5B2-4EF1-B2C3-64CB7C158FF3:rwx**

clients **exports environment variables** in application context
to configure strong authentication - *root* role on client is **unavailable**

# eos*xd*

## sub-mount feature

glue external filesystems

- **Question**: can we integrate seemingly external filesystems into an EOS mount keeping their full performance?

- automount is a proven solution, but it has a static configuration and can not be configured by a user on the fly

```
/eos/user/f/foo/                    → EOS area
/eos/user/f/foo/software/root6      → software image
/eos/user/f/foo/hpc                 → manila share
/eos/user/f/foo/s3                  → S3 bucket
/eos/user/f/foo/backups             → backup snapshots
```

**Short answer: yes we can!**

# eosxd

## sub-mount feature

glue external filesystems

- allows **eos**xd to mount on-the fly any kind of filesystem described by a symbolic link in the EOS namespace

  - implemted: **squashfs** images with e.g. software distributions …

    - extremely space efficient file distribution with **zstd** compression, export millions of small files as a single image file

    - high-performance kernel module or FUSE module available

```
-rw-r--r--   1 nobody    nobody          256622592 Jun 29 18:04 .gcc-4.9.3.sqsh
lrwxrwxrwx   1 nobody    nobody                  1 Jun 29 18:04 gcc-4.9.3 -> squashfuse:
```

  - envisaged: **external filesystem** areas e.g. high-performance *manila* shares, *s3* buckets etc. …

    - store cephx or s3 key as private extended attribute in EOS

  - envisaged: *restic* **backup snapshots** of user areas with restore password in extended attributes in EOS

    - browse/recover existing backups stored in an external instance without help from a service manager
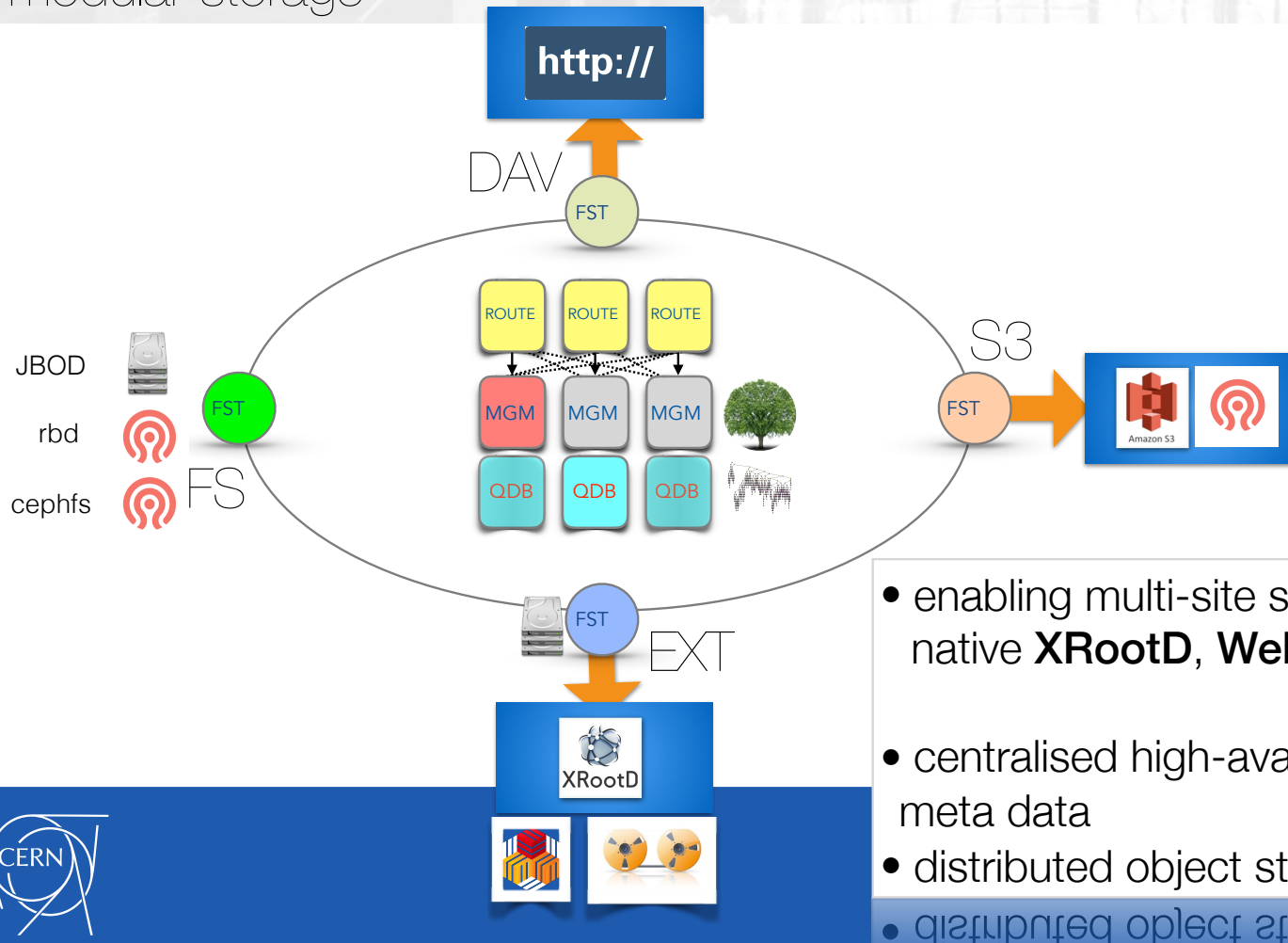
WORK
IN
PROGRESS

```
EOS Console [root://localhost] |/eos/wfe/submount/> squash
usage: squash new <path>                              : create a new squashfs under <path>
       squash pack <path>                             : pack a squashfs image
       squash unpack <path>                           : unpack a squashfs image for modification
       squash info <path>                             : squashfs information about <path>
       squash rm <path>                               : delete a squashfs attached image and its smart link
```

**eos**xd **leverages performance of external optimised filesystems**

# Distributed Storage Architecture

2018++

**EOS** modular storage

**http://**

DAV

FST

ROUTE ROUTE ROUTE

JBOD

rbd

cephfs

FS

FST

MGM MGM MGM

QDB QDB QDB

S3

FST

Amazon S3

FST

EXT

XRootD

- enabling multi-site storage supporting native **XRootD**, **WebDav**, AWS/CEPH **S3** or **FS** storage

- centralised high-available namespace in KV store for meta data
- distributed object store for data
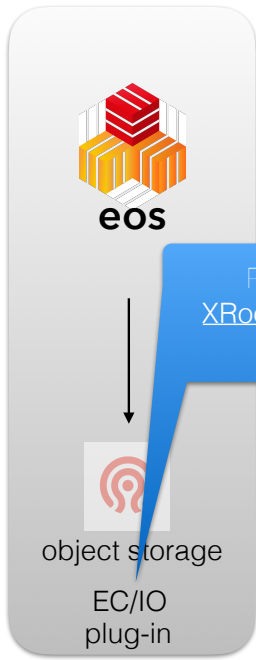
CERN

eXtreme DataCloud

# Modular Storage
## XRootD http ecosystem

**http access**

the end of FTP/SRM

XRootD

| HTTP TPC |
| Token Auth |
| Macaroons |

Parallel Sessions
Capability-Based
Authorization for HEP

eos

HTTP
Protocol Support

protocol gateways

| gridFTP | gridFTP | gridFTP | SRM |

eos

protocol gateways    XrdHttp

| XRootD | XRootD | XRootD | XRootD | XRootD |
| HTTP TPC | HTTP TPC | HTTP TPC | HTTP TPC | HTTP TPC |
| Token Auth | Token Auth | Token Auth | Token Auth | Token Auth |
| Macaroons | Macaroons | Macaroons | Macaroons | Macaroons |

eos

**XRootD is growing a complete set of plug-ins for HTTP enabled storage allowing decommissioning of gridFTP/SRM soon(ish)**
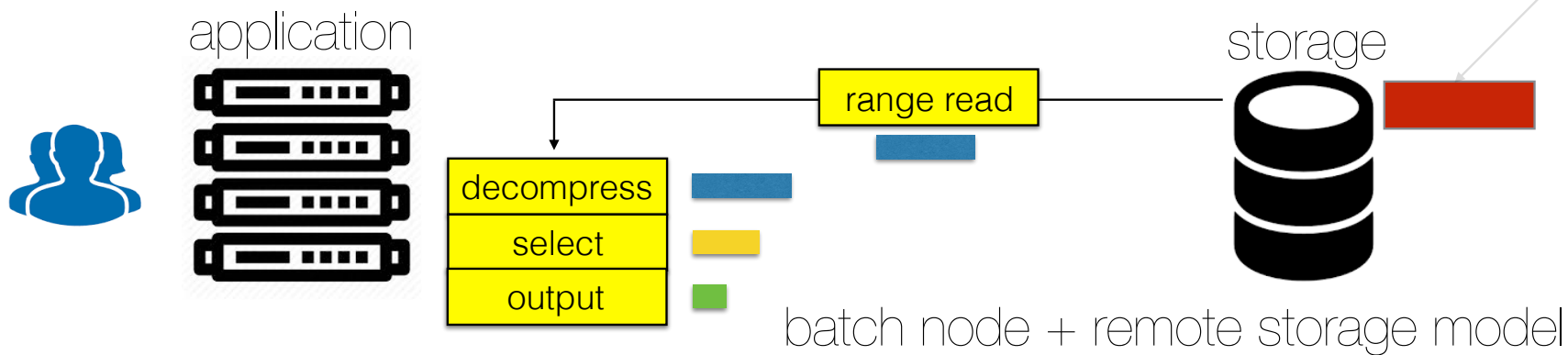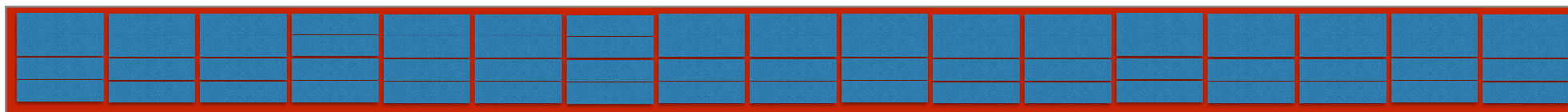
# Our conventional file processing model

Parallel processing of a large file by e.g. 10k subtasks is not very scalable/efficient when using POSIX I/O.
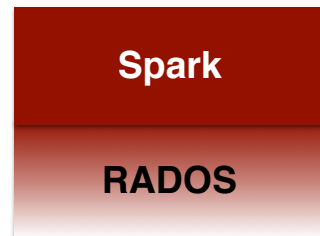
25 GB file



application

storage

range read

decompress
select
output

batch node + remote storage model

**Do we need to change this simple model?**

# Why Spark on Ceph? (Part 1 of 3)

*Posted on: June 25, 2018*

https://redhatstorage.redhat.com/2018/06/25/why-spark-on-ceph-part-1-of-3/

**Sounds HADOOP-like**

| Spark |
|---|
| **RADOS** |

but **means** only **S3 remote reading**

| Spark |
|---|
| **S3A** |

↓

| **S3** |
|---|
| **RADOS** |

**Conclusion in this article:**

Not highest possible performance when **storage and compute** are **separated** but the most flexible model when you have many people sharing infrastructure.

We figured that out already. **That is what we did and do**!

positive**+** CEPH S3 buckets can be configured to be index-less removing a scalability limitation [sacrifying listings & accounting]
negative**-**  CEPH S3 for HEP analysis misses multi byte-range request and data flows via gateways. Good news: that could be fixed!

## Most people mean S3 when they talk about Object Storage
### In fact applications know nothing about objects

# Object Storage Usage Models
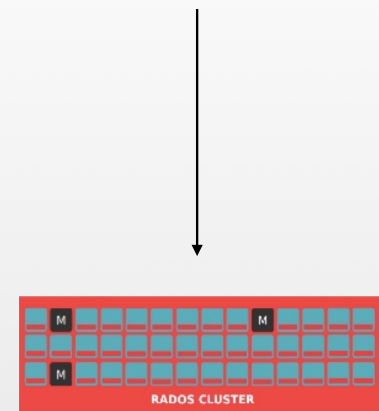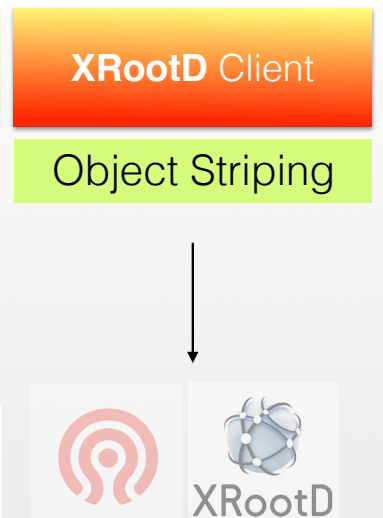
... mainly about **Parallel IO**

authenticated + universal trust

authenticated +universal trust

**Object Striping**

authenticated +universal trust

finegrained access control

finegrained access control

**POSIX**
cephfs or libcephfs

**XRootD** Client

**S3** GW

**XRootD** Server

Object Striping

Object Striping

Object Striping

Object Striping

MDS

RADOS CLUSTER

XRootD

RADOS CLUSTER

XRootD

RADOS CLUSTER

standard + secure

secure

transparent/portable

simple

most powerful

**There are many ways to do the same thing with subtle differences in complexity & functionality. Which one is the best? ... depends ...**

# Data processing
## with application object awareness

application

single roundtrip
client & server *share*
- selection
- decompression
- output processing

decompress 2
select
output

pre-selection criteria

preselect
decompress 1
output

object store

**Allows to move some IO processing inside the object storage**
**non-generic but use-case optimised approach - nice R&D**

# Summary & Outlook

- EOS has been under **steady evolution** since 8 years.

  - major promoter of XRootD as a framework and remote access protocol in HEP

  - CERN service had overrun **design limitation** in meta-data & data size during 2017 with visible impact

  - this year marks a major architectural change for scalability, availability & usability

- EOS converges towards an **integrative platform** of external storage components and services **for scientific data processing**

  - it leaves flexibility to integrate new ideas & requirements easily e.g. CERNBOX/SWAN/EOS eco-system

  - open to paradigm shift: leverage low-level components and implement high-level storage functionality

- Exabyte-scale Object Storage is an interesting technology to consider for LHC Run3

  - requires a detailed evaluation of the performance/cost model for storage and possible application benefit. Simplest approach is to build storage tiers and hide objects completely from applications. In this case: nothing visible will change for applications!

THANK YOU

QUESTIONS ?

Parallel Sessions
CERNBox: the CERN cloud storage hub

Parallel Sessions
Cloud Storage for data-intensive sciences in science and industry

Parallel Sessions
CERN Tape Archive: From Development to Production Deployment

Parallel Sessions
Sharing server nodes for storage and compute

Parallel Sessions
Facilitating collaborative analysis in SWAN

Poster Sessions
The EU Up to University Project

Parallel Sessions
Scaling the EOS namespace

Parallel Sessions
Providing large-scale disk storage at CERN

Parallel Sessions
Testing of complex, large-scale distributed storage systems

Parallel Sessions
Disk failures in the EOS setup at CERN

Parallel Sessions
Ceph File System for the CERN HPC infrastructure

Parallel Sessions
Capability-Based Authorization for HEP