



EOS Open Storage

eosxd & eoscfds FUSE filesystems



CernVM Workshop 2024

16–18 Sept 2024  
CERN  
Europe/Zurich timezone

---

**Elvin Alin Sindrilaru**

**& Andreas-Joachim Peters**

*for the EOS Project - CERN IT - Storage Group*

---

- EOS at CERN
  - Architecture
  - Production usage of FUSE mounts
- FUSE clients in EOS
  - eosxd
    - Evolution
    - Outlook
  - eoscfds
    - Design

## About EOS

EOS provides a service for storing large amounts of physics data and user files, with a focus on interactive and batch analysis.

### Flexible



EOS is a storage solution for central data recording, analysis and processing++

### Adaptable and Scalable



EOS supports thousands of clients with random remote I/O patterns with multi protocol support  
WebDAV, CIFS, FUSE, XRootd, GRPC.

### Over 930 PB at CERN



Designed for high capacity and low latency.



### Security

EOS offers a variety of authentication methods: KRB5, X509, OIDC, shared secret, and JWT and proprietary token authorisation.



### Sync & Share

EOS provides Sync&Share functionality for the **CERNBox** front-end services.



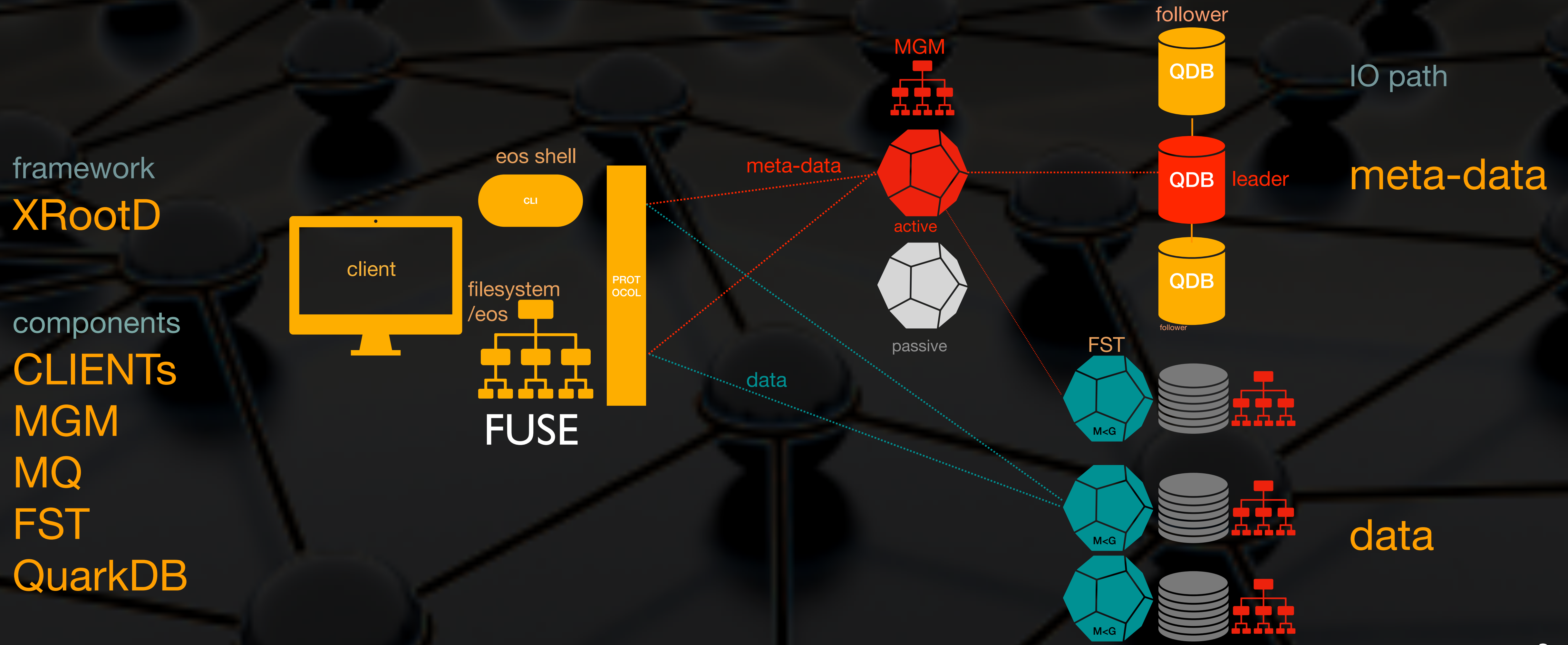
### Tape Storage

EOS includes tape storage in combination with the **CTA** Cern Tape Archive software.

### Main differences to mainstream filesystems:

- no client trust
- krb5, gsi, oidc
- clients are fully audited
- clients are regulated (rate limiting)
- subtree space accounting
- sync time aggregation
- file checksums
- user, group project quota
- rich ACL model ( e-groups )
- distributed byte-range & flock
- not 100% POSIX
- supporting thousands of clients
- remote stack trace, logging & eviction

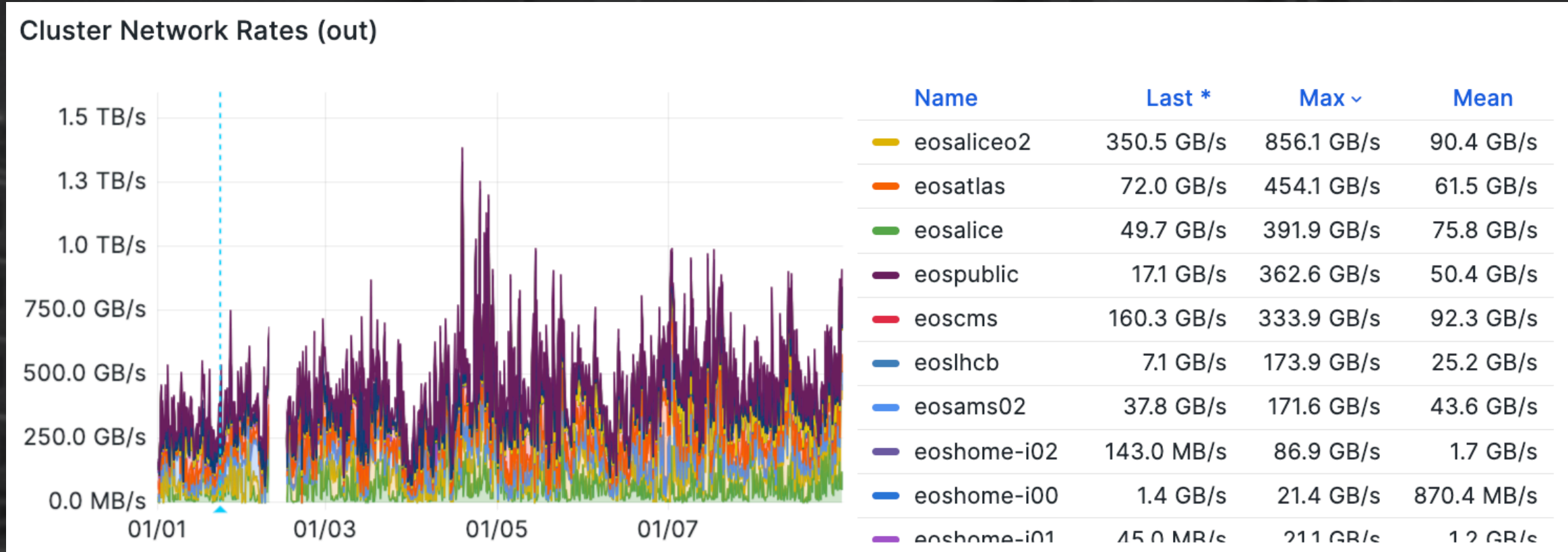
# EOS Architecture



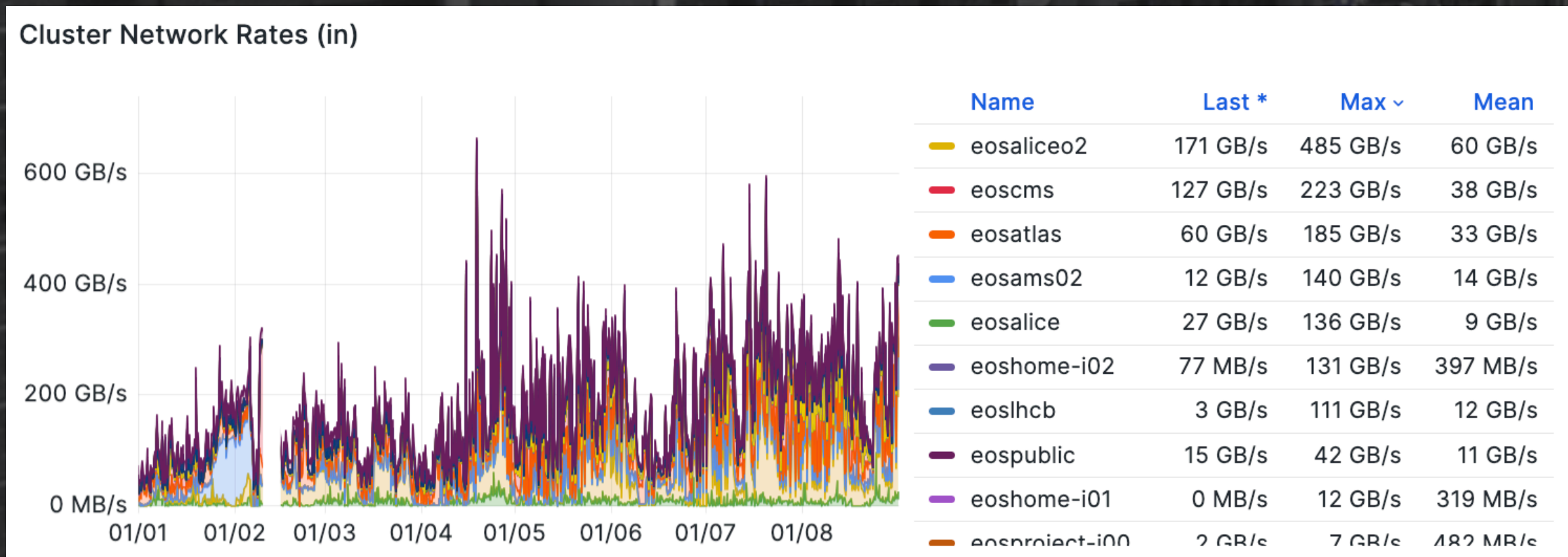
framework  
**XRootD**

components  
**CLIENTs**  
**MGM**  
**MQ**  
**FST**  
**QuarkDB**

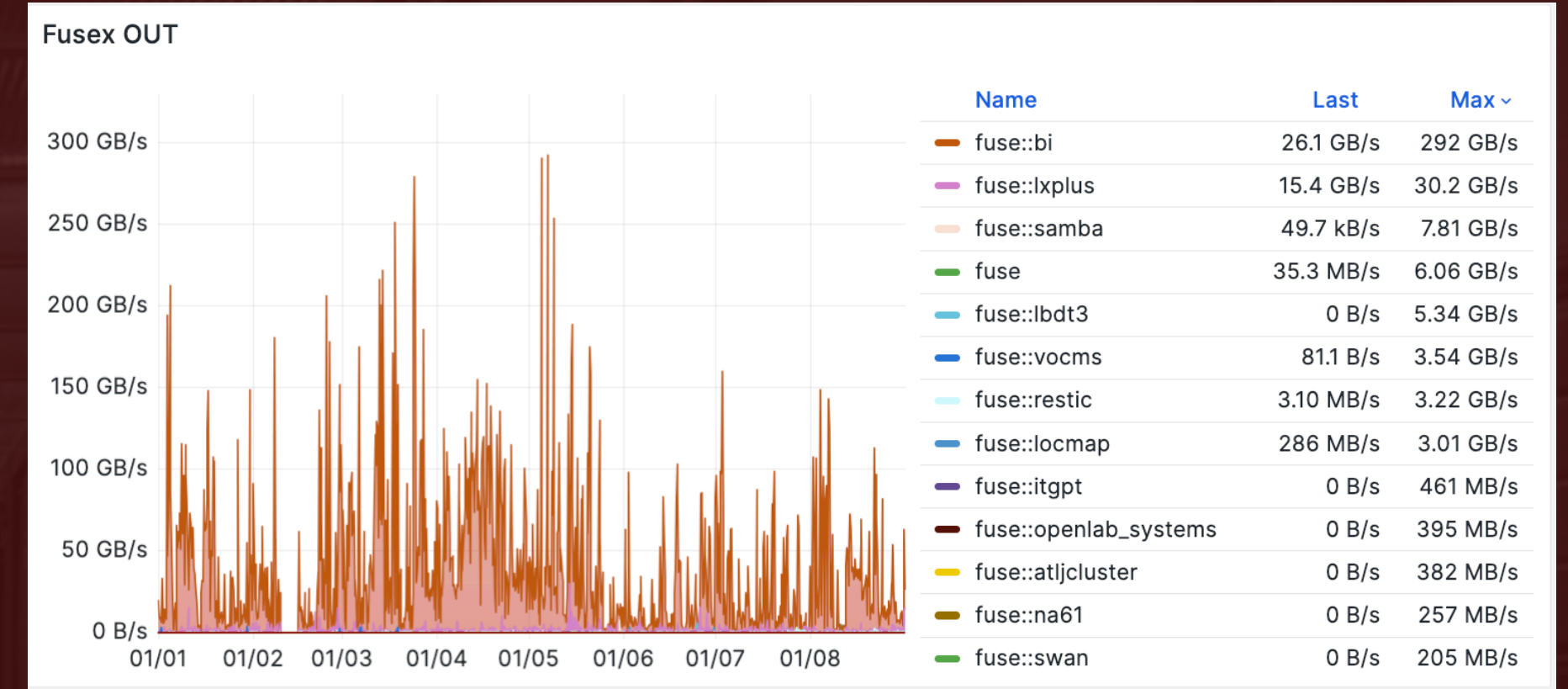
MGM meta-data server FST storage server QuarkDB meta-data persistency



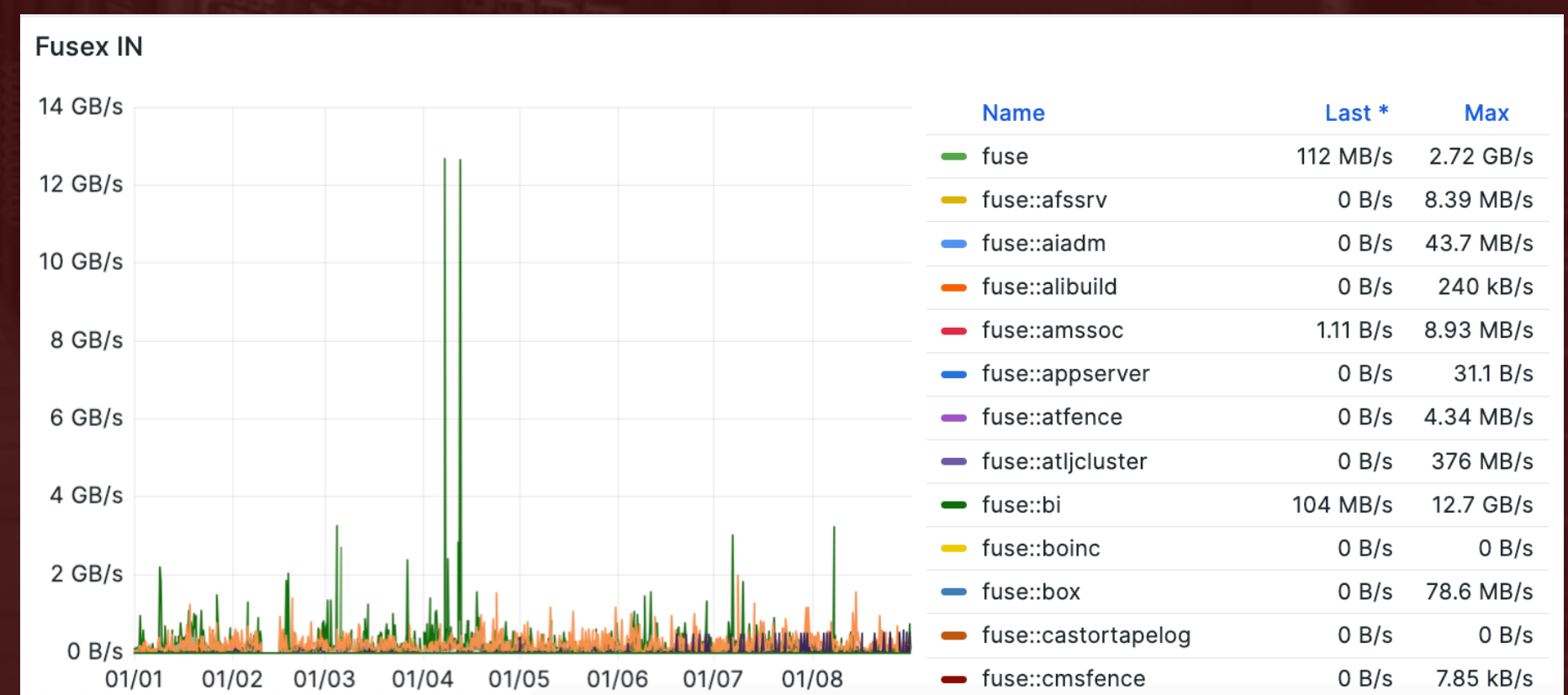
read all protocols > 1 TB/s



write all protocols > 300 GB/s

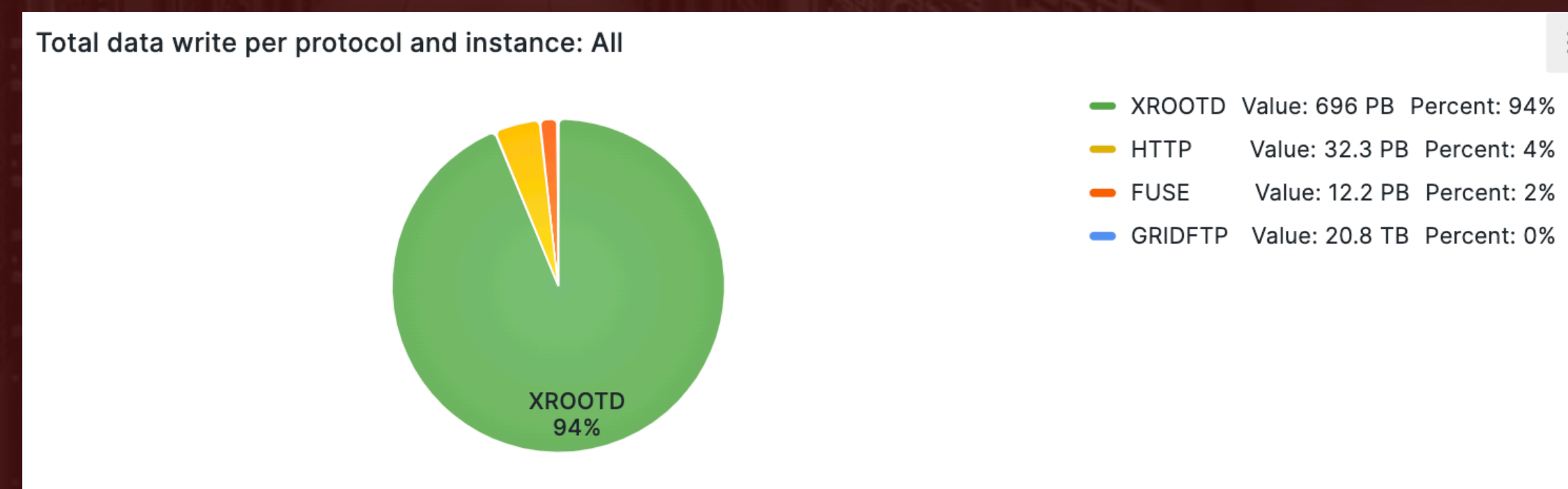
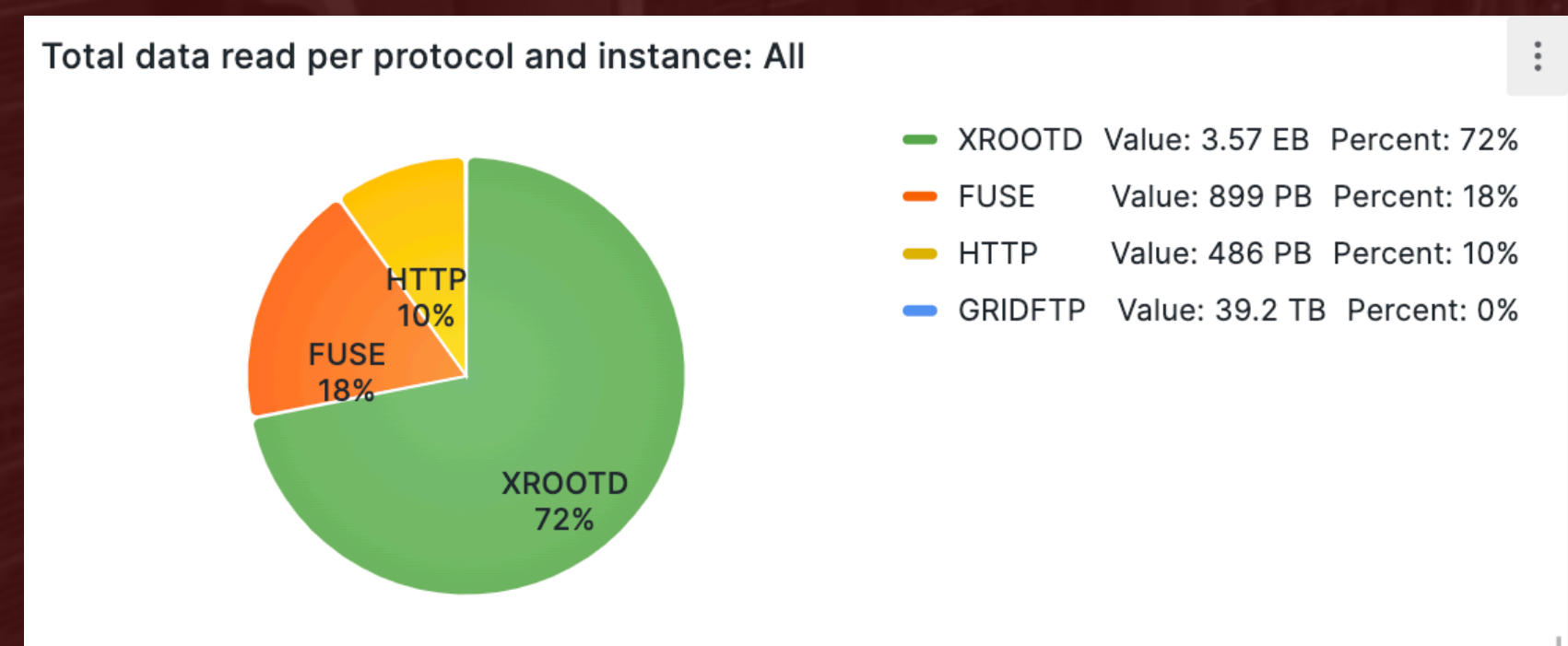
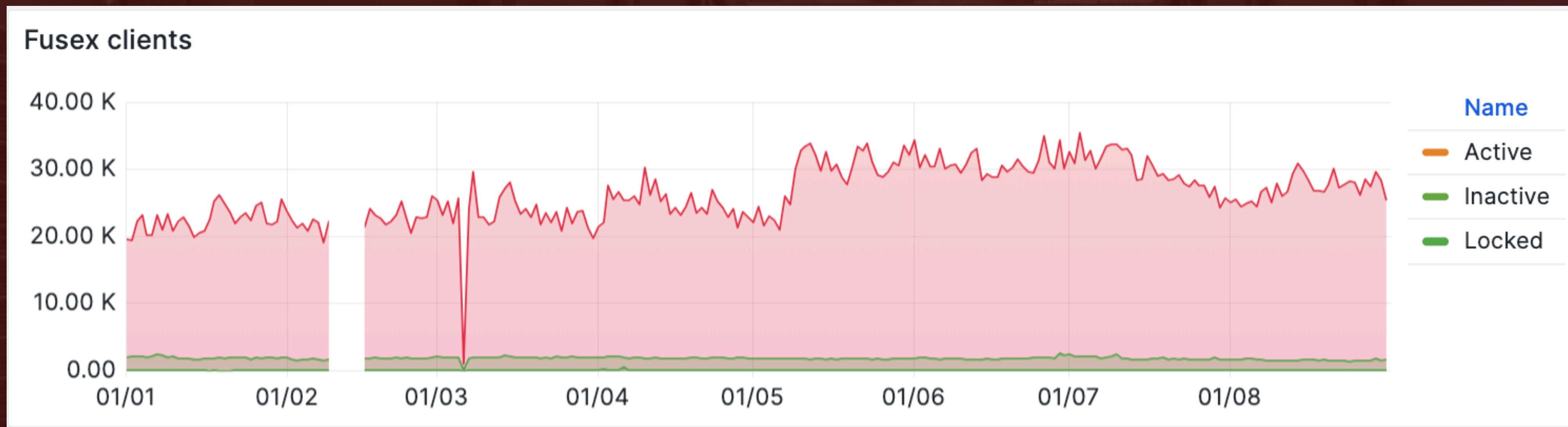


read **FUSE** traffic peak 292 GB/s



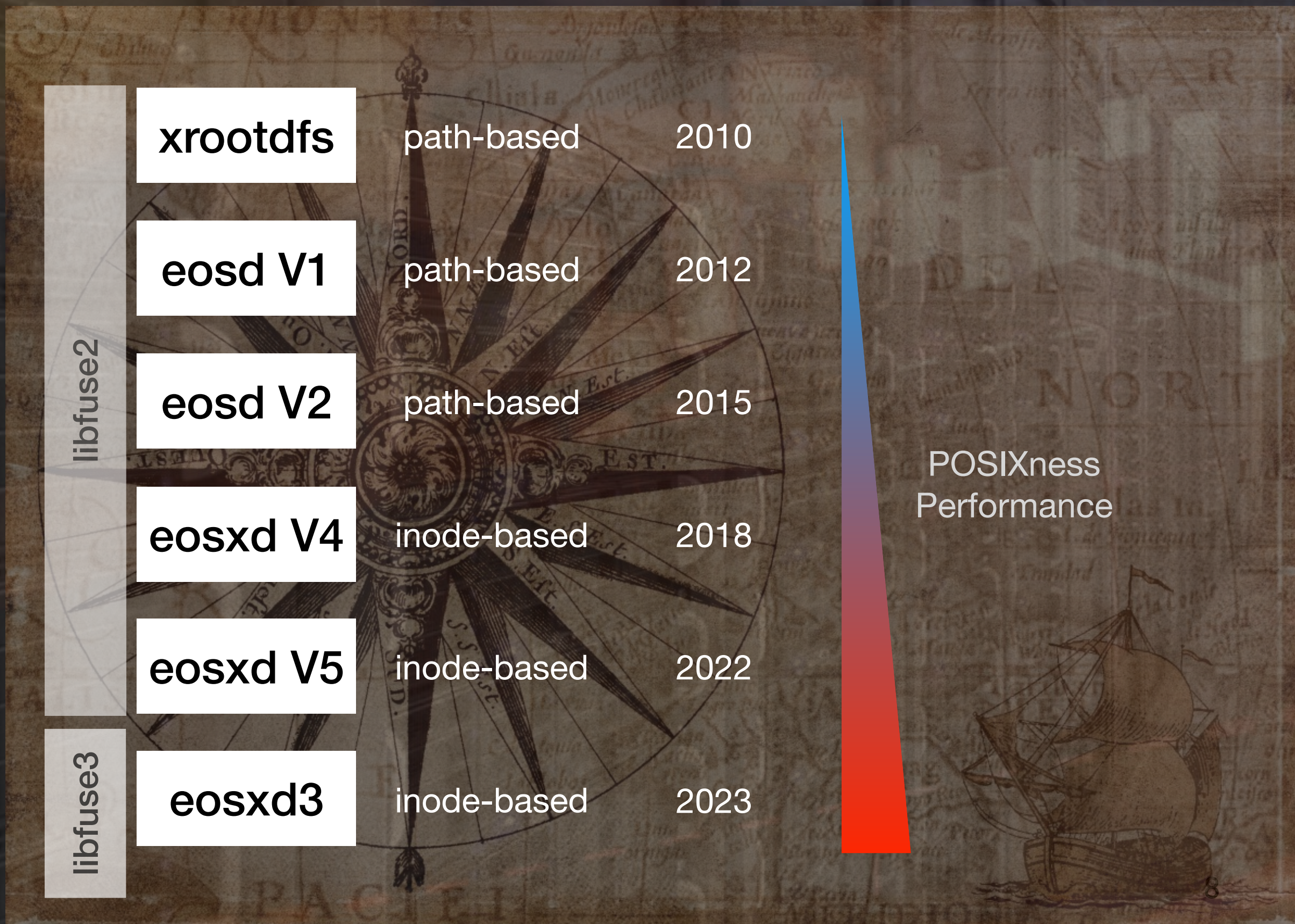
write **FUSE** traffic peak 6 GB/s

Over 30k mount clients - 900 PB read 32 PB written via FUSE this year so far

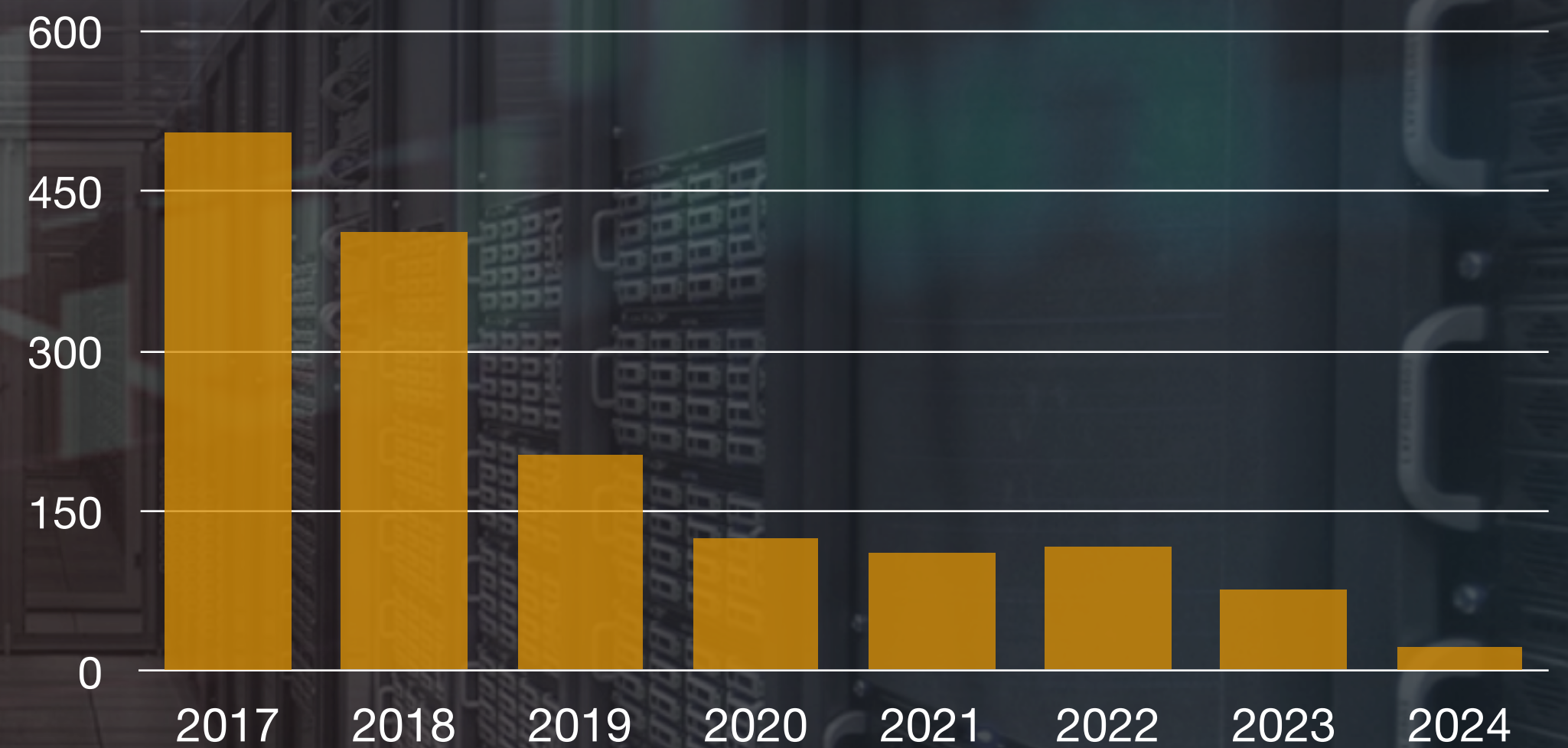


- until 7-2024 mainly **CentOS7** clients **libfuse2**
- since 7-2024 mainly **ALMA9** clients **libfuse2**
- never seen crashes due to FUSE kernel module or **libfuse** itself 🙌🙌🙌
- **libfuse3** clients only used in SAMBA gateways for now
- today we have **few crashes per month** in population of 30k client
  - in the past often due to unsafe concurrent access on meta-data objects - never from **libfuse**
- in the past many issues due to the notification callback mechanisms of FUSE (inval, lookup, forget)
  - typical dead-lock scenario due to locked objects resulting in D state processes

## eosxd (EOS FUSE daemon) *genealogy*



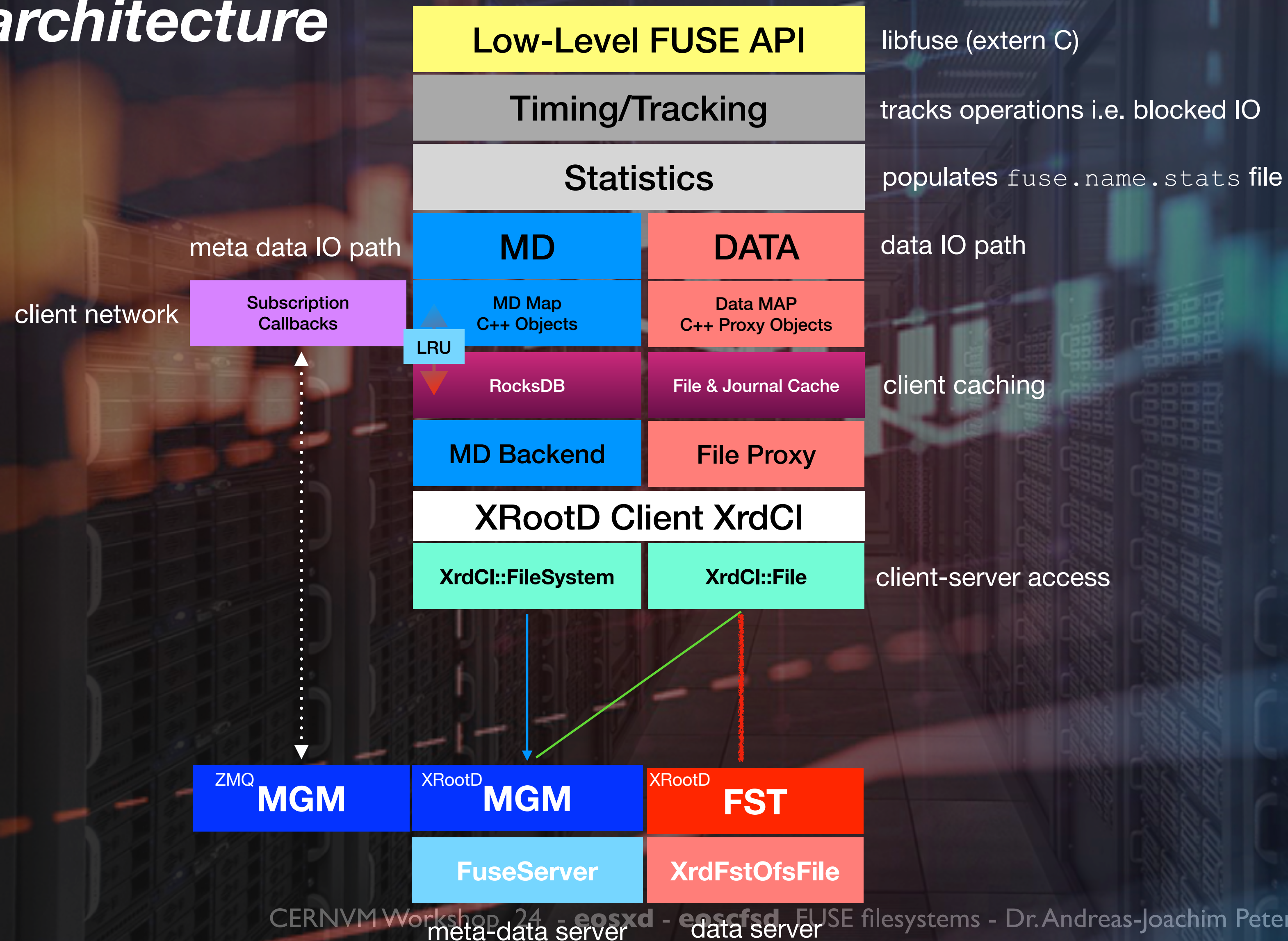
### EOS Fuse Dev - Commits per Year



XRootD5



## eosxd architecture



- **EOS FUSE** supports **various authentication mechanisms**
  - Kerberos 5 ( FILE, KEYRING )
  - GSI Proxy Certificates
  - Shared **Secrets**
  - **UNIX** (uid,gid) based authentication
  - **OIDC** ( OAuth2 token )
  - **JWT** Token (SciTokens)
  - **EOS** Token
- The available **AUTHs** methods are **pre-set by configuration**
  - run-time AUTH configuration is **read from process environment** of the calling process or its parent process e.g. **KRB5CCNAME**
  - Heuristic to resolve **dead-lock occurring** when reading environment of a **forking process**

- **EOS FUSE** uses available **credentials in configured order** to create authenticated connections to the namespace server
  - authenticated connections are bound to process IDs
- **Permissions are always delegated from server-side**
  - either via **authenticated call** to namespace or **delegated capability** valid by default for 5 minutes until revoked
    - **every file open** is authorised locally and remote by the namespace
    - **every rm, mkdir, rmdir** is authorised locally and remote by the namespace
    - only **stat & ls** are authorised locally using a delegated capability

# eosxd

## Namespace FuseServer - central service for FUSE clients

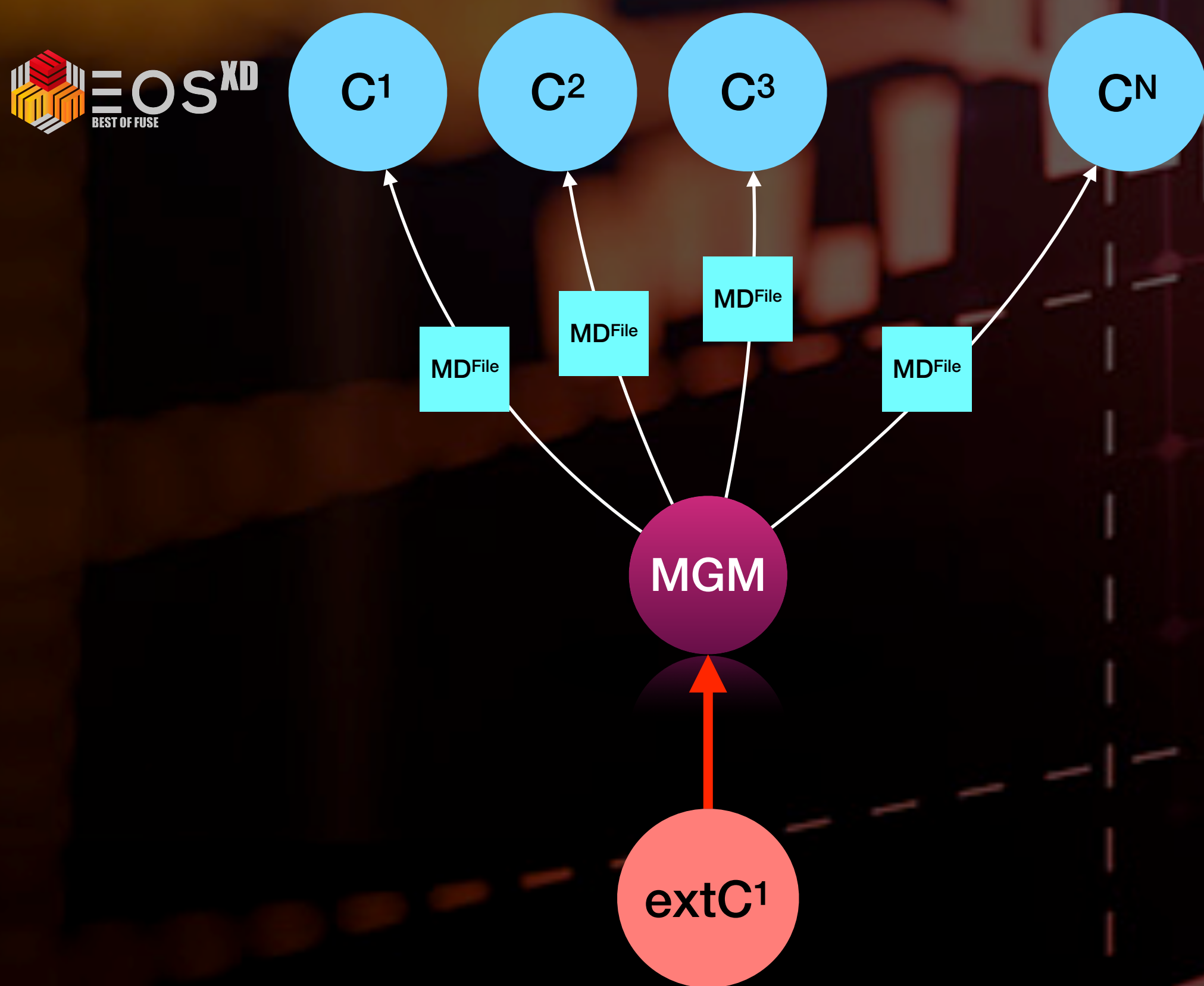
- **The Fuse Server server-side CLI is used to ...**
  - **show connected clients**
    - and most important state information (stuck IO etc.)
  - **evict unwanted clients**
  - **show notification subscriptions and permission delegation (capabilities)**
  - **drop state information**
    - capabilities (directory subscriptions)
    - locks (delegated locks given to a client)
  - **configure heartbeat and broadcast configuration**
    - how often each client has to send a heartbeat
    - audience suppression to avoid too high callback rates

# eosxd

## Broadcasts & Callbacks

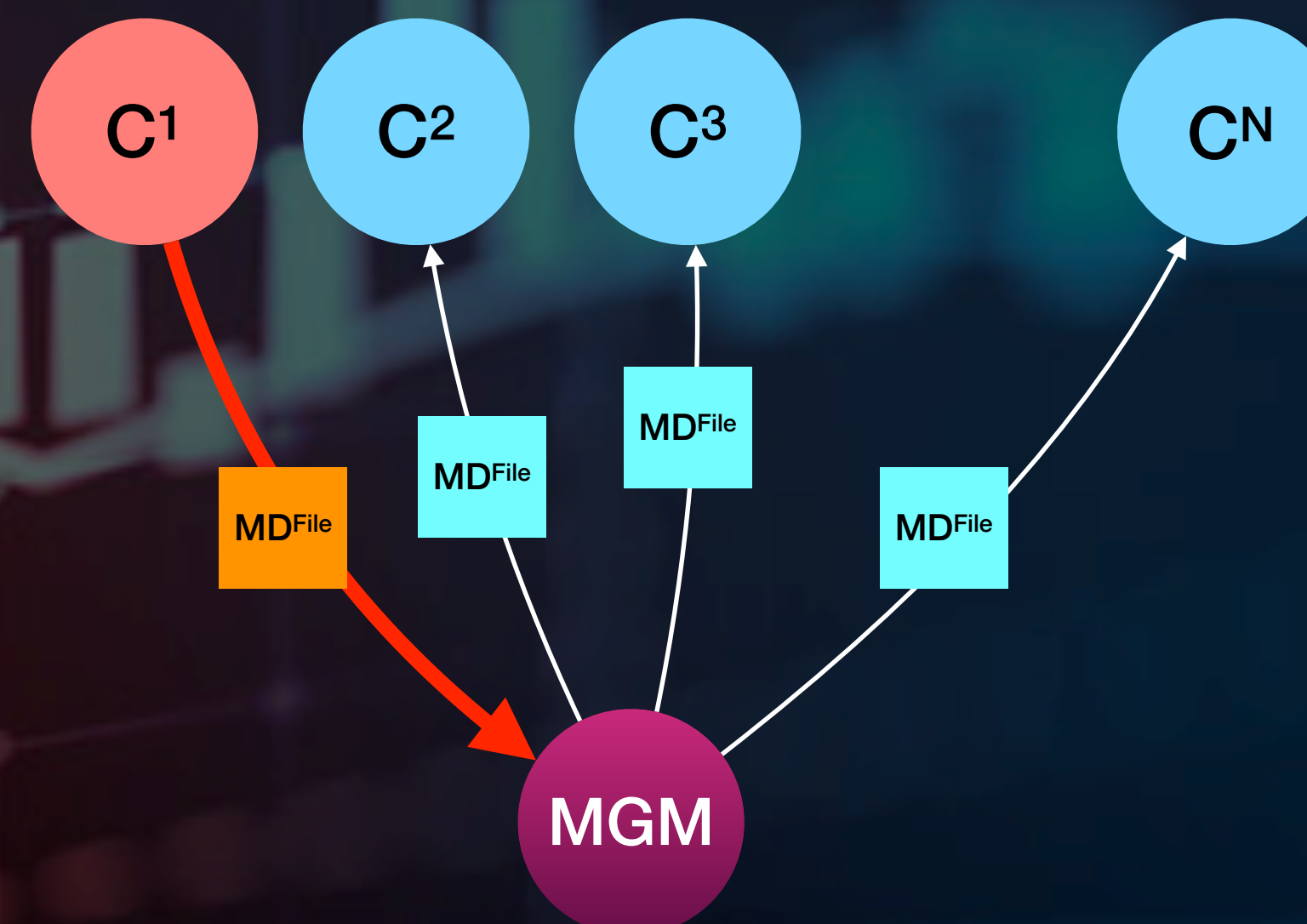
### external broadcasts

an external Client (not FUSE) triggers a broadcast of file meta-data



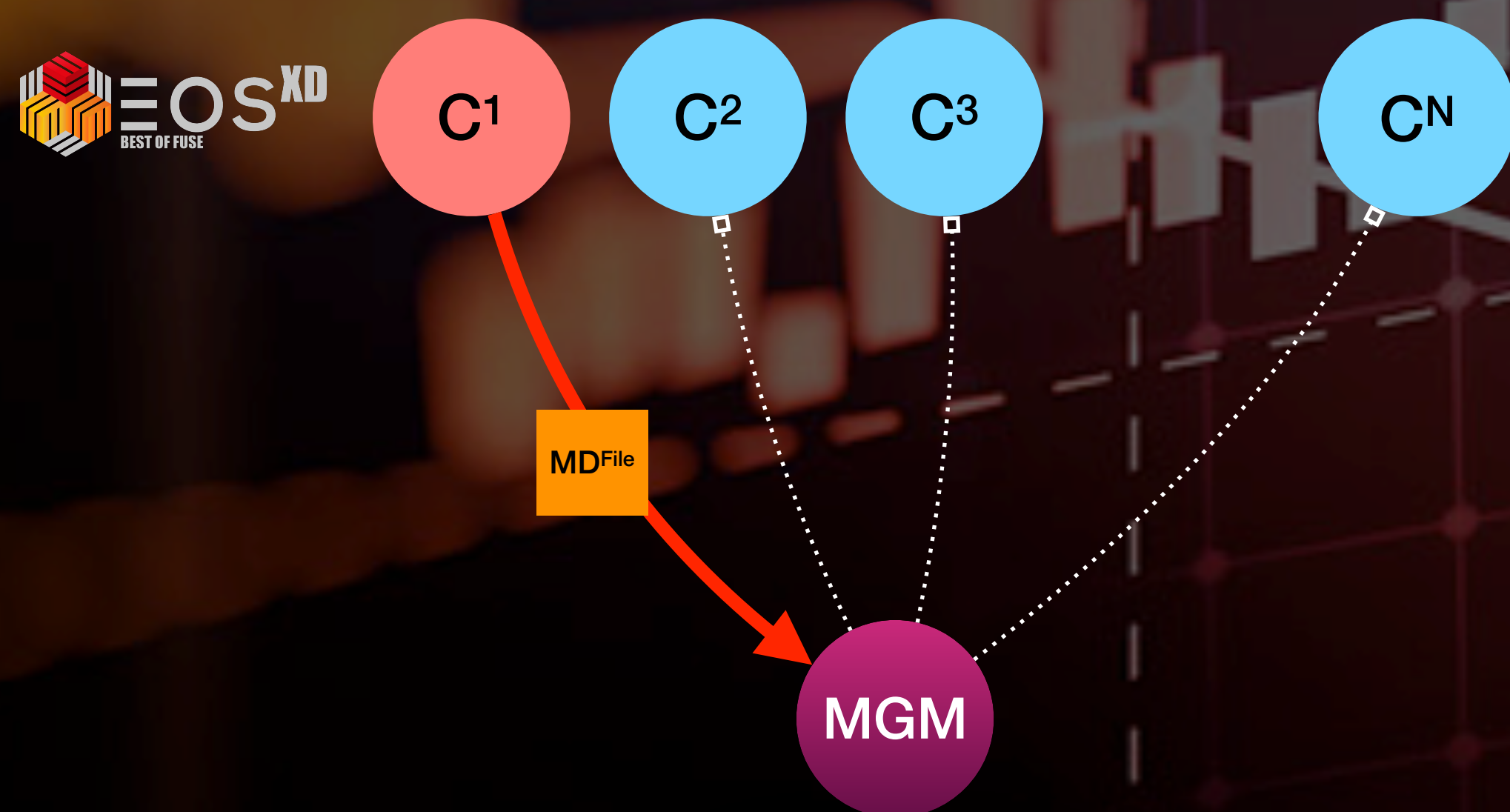
### internal broadcasts

an internal client triggers a broadcast of file meta-data



# eosxd

## Broadcast Suppression



- an internal client would trigger a **broad cast** of file meta-data
- if the receiving **audience is too large** and the broadcast gets suppressed
- the not updating clients keep with an **unsynchronised state** until their subscriptions expire and a directory refresh is done
  - **up to 300s by default**

# eosxd

## Namespace FuseServer

- admin CLI for managing client mounts

<b>eos fusex ls</b>	list client nodes
<b>eos fusex ls -l</b>	list client nodes and some statistics per client
<b>eos fusex evict</b>	force unmount, trigger a stack trace, ask for logfile, truncate logfiles
<b>eos fusex conf</b>	configure heartbeat interval configure audience suppression
<b>eos fusex caps</b>	show directory subscriptions
<b>eos fusex dropcap(s)</b>	drop one or all subscriptions from a client
<b>eos fusex droplocks</b>	drop all locks for given inode and process

```

eos fusex ls -l
client : eosxd                               lxpplus700.cern.ch 5.1.14  online
.....  ino                               : 40868
.....  ino-to-del                          : 0
.....  ino-backlog                         : 0
.....  ino-ever                            : 511893
.....  ino-ever-del                        : 18426
.....  threads                             : 60
.....  total-ram                           : 29.988 GB
.....  free-ram                            : 0.280 GB
.....  vsize                               : 2.016 GB
.....  rsize                               : 0.193 GB
.....  wr-buf-mb                           : 0 MB
.....  ra-buf-mb                           : 24 MB
.....  load1                               : 19.33
.....  leasetime                           : 300 s
.....  open-files                          : 22
.....  logfile-size                        : 18576
.....  rbytes                              : 20397308261
.....  wbytes                              : 24307570949
.....  n-op                                : 4656820
.....  rd60                               : 0.00 MB/s
.....  wr60                               : 0.33 MB/s
.....  iops60                              : 303.83
.....  xoff                                : 830
.....  ra-xoff                             : 0
.....  ra-nobuf                            : 9489
.....  wr-nobuf                            : 0
.....  idle                                : 24
.....  blockedms                           : 0.00 [i]

```

# eosxd

## Client side Disk Caching

### File & Journal Cache

```
"cache" : {  
  "type" : "disk",  
  "size-mb" : 512,  
  "size-ino" : 65536,  
  "file-cache-max-kb", 256,  
  "file-journal-max-kb", 131072  
  "journal-mb" : 2048,  
  "journal-ino" : 65536,  
  "clean-threshold" : 85.0,  
  "location" : "/var/cache/eos/fusex/cache/",  
  "journal" : "/var/cache/eos/fusex/journal/",  
  "read-ahead-strategy" : "static",  
  "read-ahead-bytes-nominal" : 262144,  
  "read-ahead-bytes-max" : 2097152,  
  "read-ahead-blocks-max" : 16,  
  "max-read-ahead-buffer" : 134217728,  
  "max-write-buffer" : 134217728  
}
```

- There are **two data caches** which are persisted while a mount is still up
  - **journal cache**
    - journalling writes
      - (and also reads if a file has `attr:sys.file.cache=1`)
      - with 128MB max. size per file
        - is used to recover writes
  - **file (start) cache**
    - caches the first 256kb of each file
      - triggered when a file is opened



# eosxd

## Virtual/Functional Attributes

- client CLI to interact with the mount

Information/Getter

Virtual Ext. Attribute	Meaning
eos.btime	file birth time
eos.ttime	latest modification time in subtree
eos.tsize	total size of directory subtree here
eos.dsize	<b>total size of files in this directory</b>
eos.dcount	number of subdirectories here
eos.fcount	number of files here
eos.name	EOS instance name
eos.hostport	EOS hostname:portname
eos.mgmurl	EOS instance URL root://...
eos.url.xroot	EOS XRootD TURL of the given file
eos.stats	get statistics output file
eos.quota	get quota output for the calling user

Authentication IF/Getter

Virtual Ext. Attribute	Meaning
eos.identity	get identity of the calling client
eos.identityparent	get identity of the calling client for the parent directory
eos.reconnect	force to drop current connection and create a new one
eos.reconnectparent	as above but for the parent

Modifications / Setter

Virtual Ext. Attribute	ARG	Meaning
system.eos.debug	<level>	Set debug level
system.eos.dropcap	-	Drop subscription here
system.eos.dropcaps	-	Drop all subscriptions
system.eos.resetstat	-	Reset Stat Counter
system.eos.log	private public	Change Logfile Perm.
system.eos.fuzz	all/config	Toggle Fuzzing mode

# eosxd vs eosxd3 changes

- **readdirplus**

- the function returning a directory listing attaches the stat information for each listed entry
  - avoids additional `getattr` calls for each listed child!

- **forgetmulti**

- when the kernel wants to forget inodes, instead of recalling one by one a bulk request reduces the number of `forget` calls substantially

- **clone\_fd** option

- each FUSE thread creates a separate device file descriptor for each processing thread, which might improve performance.

- **write\_back\_cache** option

- small writes are not hitting directly eosxd, but they are written into the buffer cache and then flushed as larger operations
- increases single byte write performance from 16kHz to 0.5MHz
- but does not allow remote invalidation

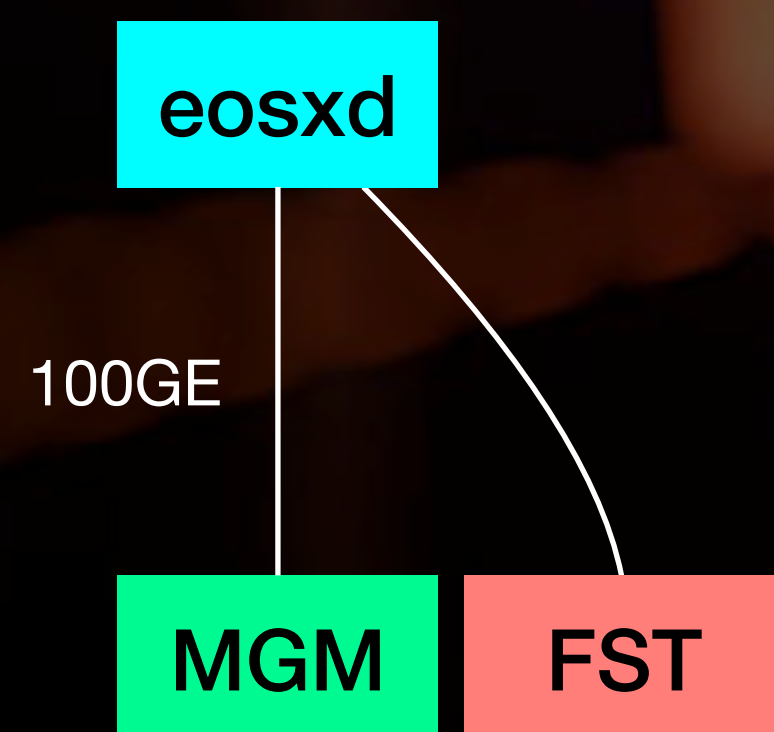
- **FUSE\_SET\_ATTR\_CTIME**

- allows to set the change time with the `utimes` function

# eosxd3

## Performance

- 100 GE client/server idle instance test



Single Client	eosxd	eosxd3
Seq. Creation		700 Hz
Par. Creation		1000 Hz
w IOPS	16 kHz	500 kHz
Seq. Read		1.6 GB/s
Seq. Write		900 MB/s

- FUSE3: **write-back cache** cannot be used because callbacks are suppressed for known inodes - see [here](#)



apeters1971 on Nov 22, 2023

When I enable write-back cache with a 4.18 or 5.14 kernel, callbacks invalidating attributes of an inode e.g. a remote size change are not applied and the current size is stale. Is that a known issue or a limitation?

Is that due to this comment in the kernel source (in this case 5.14):

```
/*  
 * In case of writeback_cache enabled, the cached writes beyond EOF  
 * extend local i_size without keeping userspace server in sync. So,  
 * attr->size coming from server can be stale. We cannot trust it.  
 */
```

↑ 1

- **double-mounting** on ALMA9 when using autofs
  - when moving from CENTOS7 to ALMA9 we had to add protection to avoid mounting a filesystem twice

20

- most relevant **performance limits** in **eosxd** are **addressable in MGM + FST implementation** - FUSE bottlenecks play only minor role in our environment
- **eosxd** would benefit from **XATTR caching** in kernel - which does not exist
- An architecture where **meta-data** caches can be directly **shared between user-space and kernel** would be desirable and simplify the user-space implementation - had a look at extFUSE



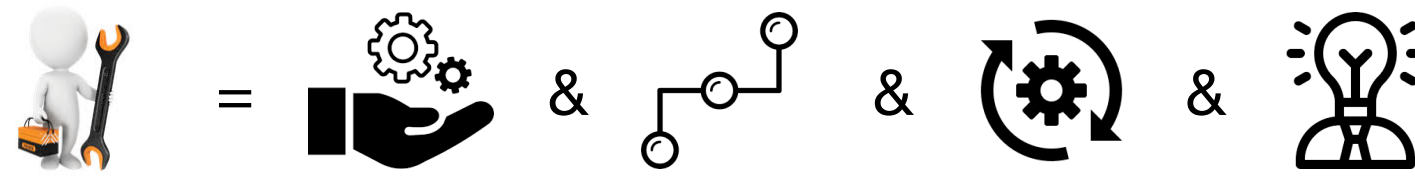
**eoscfsd**

# CFS

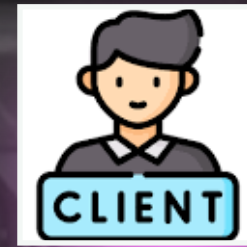
**CLIENT FILE SYSTEM**

# What is **CFS** ?

CLIENT FILE SYSTEM



- **CFS** is a FS abstraction layer written as a FUSE passthrough filesystem
  - provides **homogeneous POSIX API** connecting to storage systems
    - as much POSIX as the layer it abstracts + tiny FUSE POSIX violations
  - allows to support **arbitrary authentication methods** for any back-end
  - allows shared **extensions to standard POSIX API**
  - provides an **exit strategy** from any back-end
- a tiny **R&D project** - which is now part of EOS releases and usable as **eoscfds**



core comes from example implementation of libfuse-3

low-level FUSE API

**FUSE Passthrough Filesystem**

convenience functionality

**Mountable Filesystem**

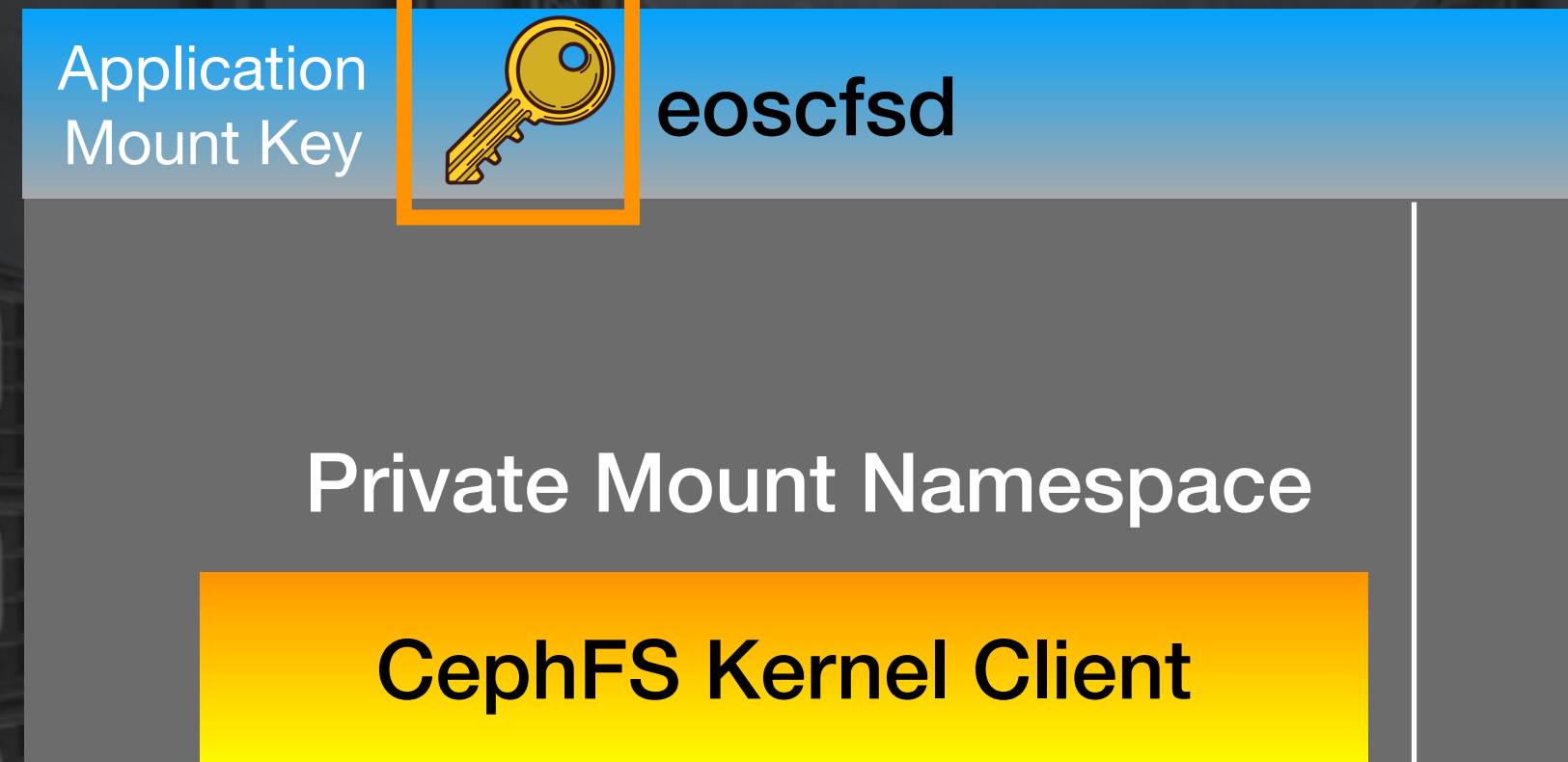
anything which can be mounted





**eoscfsd** provides a high-performance pass-through implementation for POSIX filesystems. It adds kerberos authentication, remote configuration and mount key obfuscation. **eoscfsd** fetches the mount instructions for a named mount from a configurable HTTPS server.

```
[root~]# df /cern/home/
Filesystem      1K-blocks      Used Available Use% Mounted on
cernhome        104857600  41705472  63152128   40% /cern/home
```



no caching in kernel FUSE layer  
read/write redirect to file descriptor

back-end invisible/inaccessible for root user

Local Mount Key /etc/eoscfsd.key



These keys don't remove the need to trust *root* on a host!

mount (2)

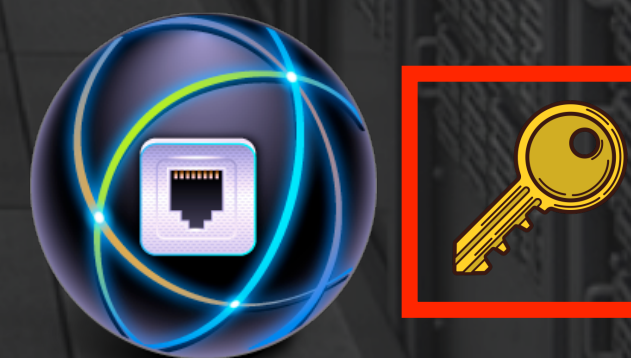


FS1 .. N



https GET (1)

Encrypted Mountpass



- RPMs available - package **eos-cfsd**
  - **kerberos** authentication
  - virtual **/.proc/** interface
    - kerberos ID to name **translation**
    - **quota** bool per user
    - **enabled** bool per user
    - **recycle** bin (enable/disable during mount)
    - **bulk deletion** wrapper via shell alias  
(remark: a shell wrapper opens up the possibility to use the recycle bin as free storage resource .. maybe we don't want to allow that)
  - virtual **xattr** '**cfs.id**' = 'who am i'
  - **statistics file** in JSON format with operations/s etc.
  - **autofs** support

[CFS Documentation](#)

26

- **unpacking the Linux kernel** on CentOS9 office desktop with CephFS disk-based backend /cern/home/
  - cephfs native backend: untars ~ **1100** files/s
  - eos`cf`sd: untars ~ **1000** files/s [90% of back-end]
  - afs: untars ~ **250** files/s
- IO **bottleneck** introduced by FUSE is **almost invisible** (reading a file inside the CephFS kernel cache via eos`cf`sd)

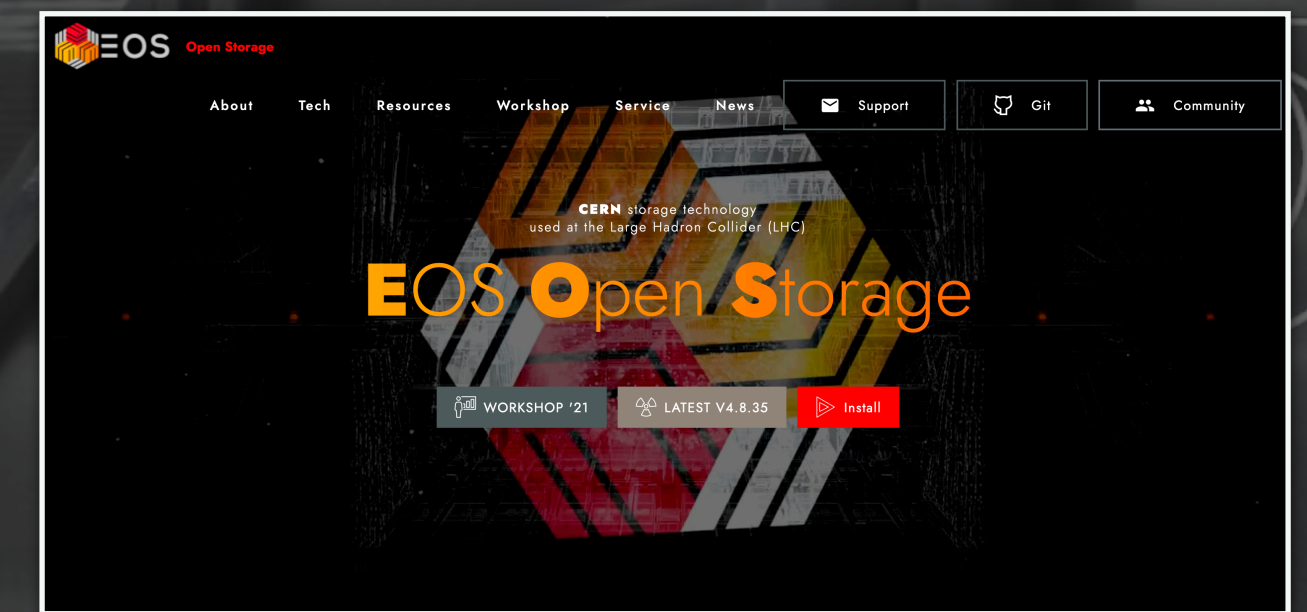
```
[apeters@engine apeters]$ dd if=1GB of=/dev/null bs=1M count=1000  
1000+0 records in  
1000+0 records out  
1048576000 bytes (1.0 GB, 1000 MiB) copied, 0.311914 s, 3.4 GB/s
```

27

- In kernel version < 6.9. FUSE **passthrough** layer receives **all** read/write requests and redirects them to the back-end filesystem file descriptor
- In kernel version 6.9 FUSE **passthrough** layer can be **bypassed** for all read/write operations completely - which means you experience the **native IOPS of the backend** ( IOPS limit we see in FUSE 16-20kHz )
- **eoscfsd** is a promising platform to create a **back-end agnostic filesystem abstraction** with customisable functionality add-ons
- If you are interested - don't hesitate to reach-out, test and/or contribute!

# Useful Links

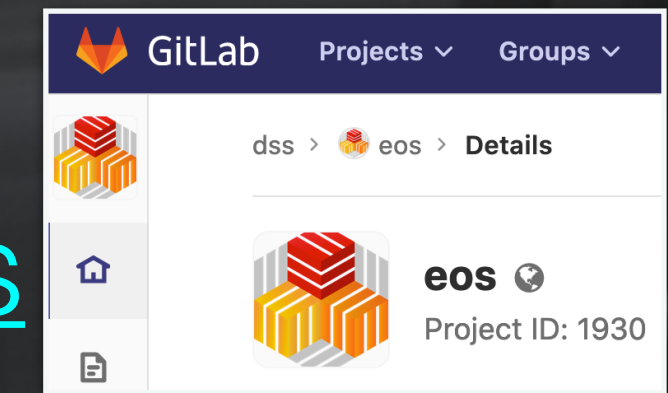
Web Page <https://eos.cern.ch>



29

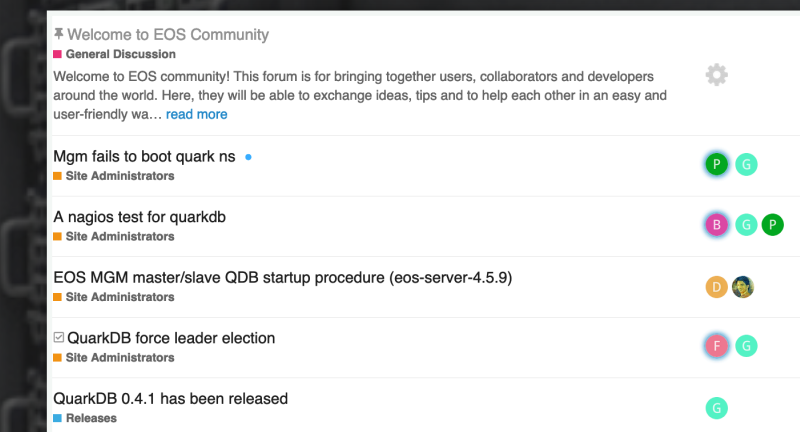
GITLAB Repository <https://gitlab.cern.ch/dss/eos>

GITHUB Mirror <https://github.com/cern-eos/eos>



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

email: [eos-community@cern.ch](mailto:eos-community@cern.ch)



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

Support email: [eos-support@cern.ch](mailto:eos-support@cern.ch)



29

***Thank you for your attention!***  
***Questions?***

