

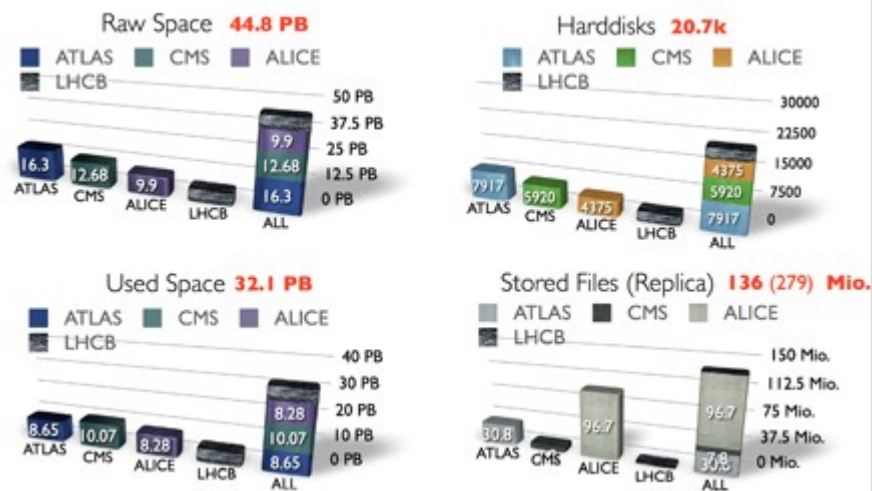


DSS

Areas of developments in storage at CERN

Massimo Lamanna

Source: A.J.Peters and FDO section



- EOS: Large disk farms for physics and beyond
 - Currently ~25 PB used quota
 - → 100 PB quota (@LHC Run2)
- Developed in CERN/IT (DSS)
- Original goal
 - Large scale (PBs for 100s/1000s independent scientists) analysis of LHC data
 - Arbitrary level of data durability via cross-node file replication or RAIN using commodity hardware
- Status
 - Open to non-physics use cases
 - NB: large number of protocols available!

EOS logo and navigation bar.

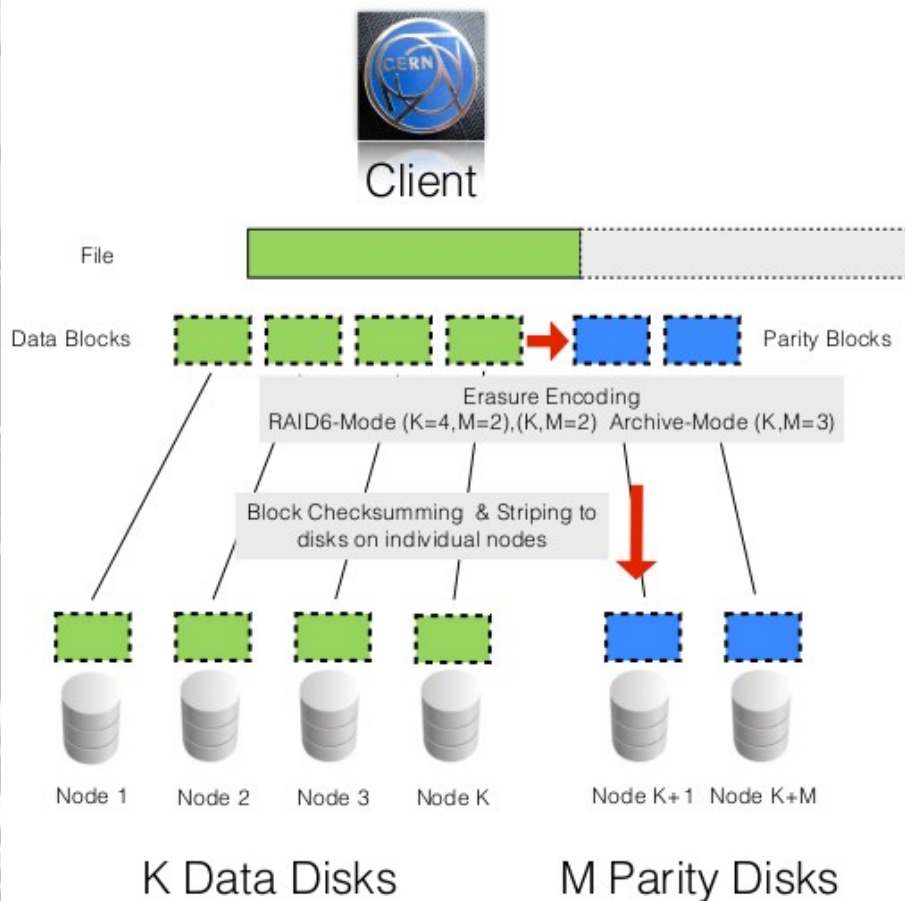
Quota Whoami Who Upload Mkdir Rmdir Info Space Nodes Groups Filesystems

quota	gid	space	usedbytes	usedlogicalbytes	usedfiles	maxbytes	maxlogicalbytes	maxfiles	percentusedbytes	statusbytes	statusfiles
node	project	/eos/lhcb/	0	0	0	0	0	0	100.00	Ignored	Ignored

```
nobody@eospps ]: /eos/lhcb
Path      Size      Created      Mode      owner      group  Act
./        -         Nov 05 2013 10:08 drwxrwxr--+ lmascott c3
../       -         Apr 26 2011 11:03 drwxrwxr--+ root    root
grid/     -         Dec 02 2013 18:09 drwxrwxr--+ lmascott z5
opetest/  -         Nov 05 2013 11:23 drwxrwxr--+ lmascott c3
proc/     -         Nov 05 2013 11:23 drwxrwxr--+ lmascott c3
test/     -         Nov 05 2013 11:23 drwxrwxr--+ lmascott c3
testCfg/  -         Nov 05 2013 11:23 drwxrwxr--+ lmascott c3
```

- No need for a central relational database
 - In memory hashtable
- Designed for infinite scalability and arbitrary reliability
 - Replica or Erasure Code (ReedSalomon, LDPC...)
- Disantagle physical and logical view
 - Disks keep file replicas, MGM manages them (no disk – service link)
 - Easy to manage realistic hardware (heterogeneity)
- Multi site federation support
 - Wigner / Geneva > 1000 Km distance

- Cornerstone of the architecture
- Initial concern
 - Run out of memory?
 - Durability
- Where are we?
 - Low-latency file access with in-memory namespace ~ 200 M files
 - Stable and durable (Master + Slaves)
 - Demonstrated by doing it
- Investigation areas
 - Stability above the 1B-file area
 - EOS Diamond
 - EOS Ceph inbreeding

**RAIN**

- Block Striping allows parallel IO

- boost single file performance

e.g. with $K=4$ 280 MB/s streaming write & 400 MB/s streaming write

- Erasure Encoding + Block Checksumming allows
 - error correction on the fly
 - up to M concurrent disk failures without data loss

- disk space savings compared to file replication

e.g. 25% less space needed for (4+2) vs. 2 replica

Not deployed at large scale: several scenarios to be investigated:

- Interplay replica/fragments and self-healing
- New modes of operations (large n of fragments)
- ...

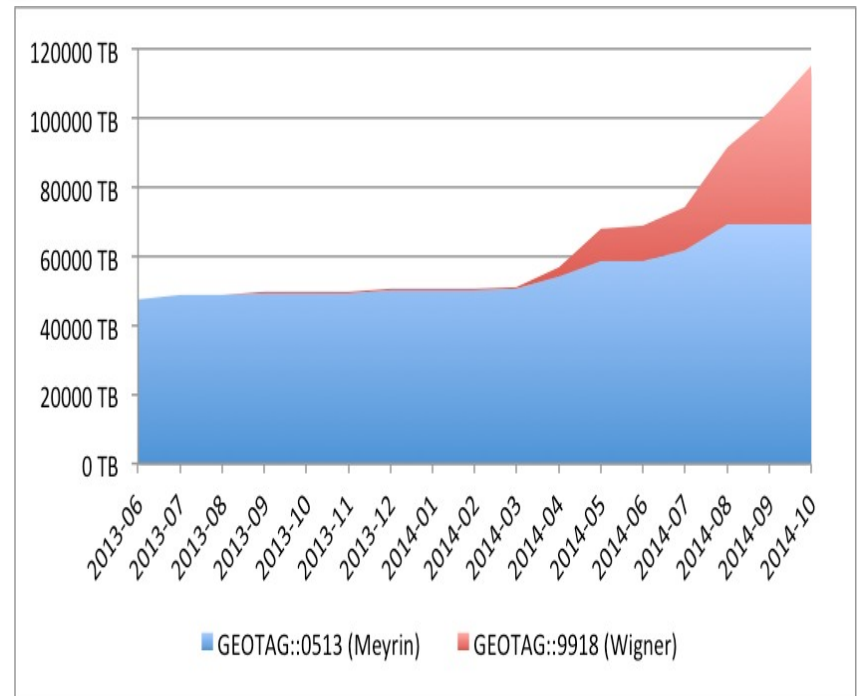
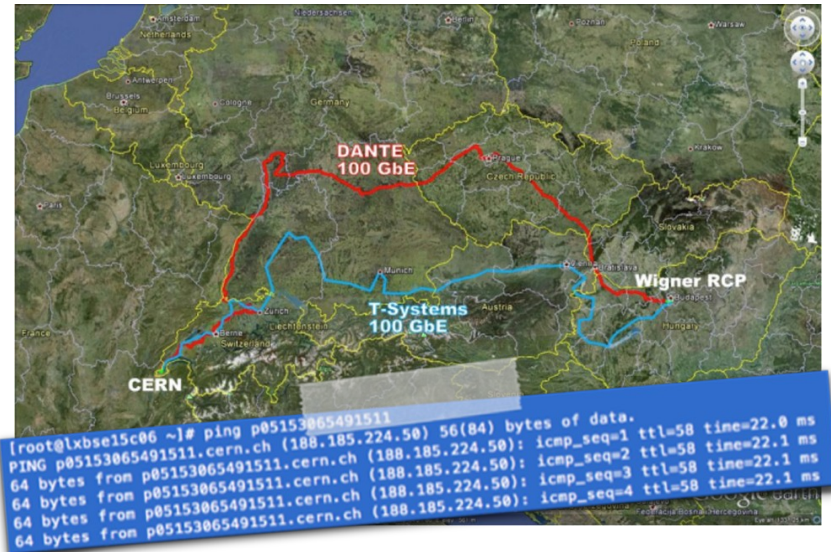
- Six multi-PB installations
 - few thousand disks per instance
 - Tape storage with CASTOR
 - New disks going mainly to EOS
 - **Operated by the same team**
- Simplified life-cycle management workflows for on-going replacement/repair of hardware
 - JBOD disks (no RAID controller) using software RAIN
 - Low-latency file access with in-memory namespace ~ 200 M files
 - Fine-grained access control and quota management
- GRID and local storage element
 - Accessed from thousands of CERN-local and remote batch nodes with Kerberos and GSI authentication
 - Full support for XRootD, GridFTP, HTTP(S) & WebDav protocol - and – mountable with FUSE as a remote file system
 - notably users love fuse...



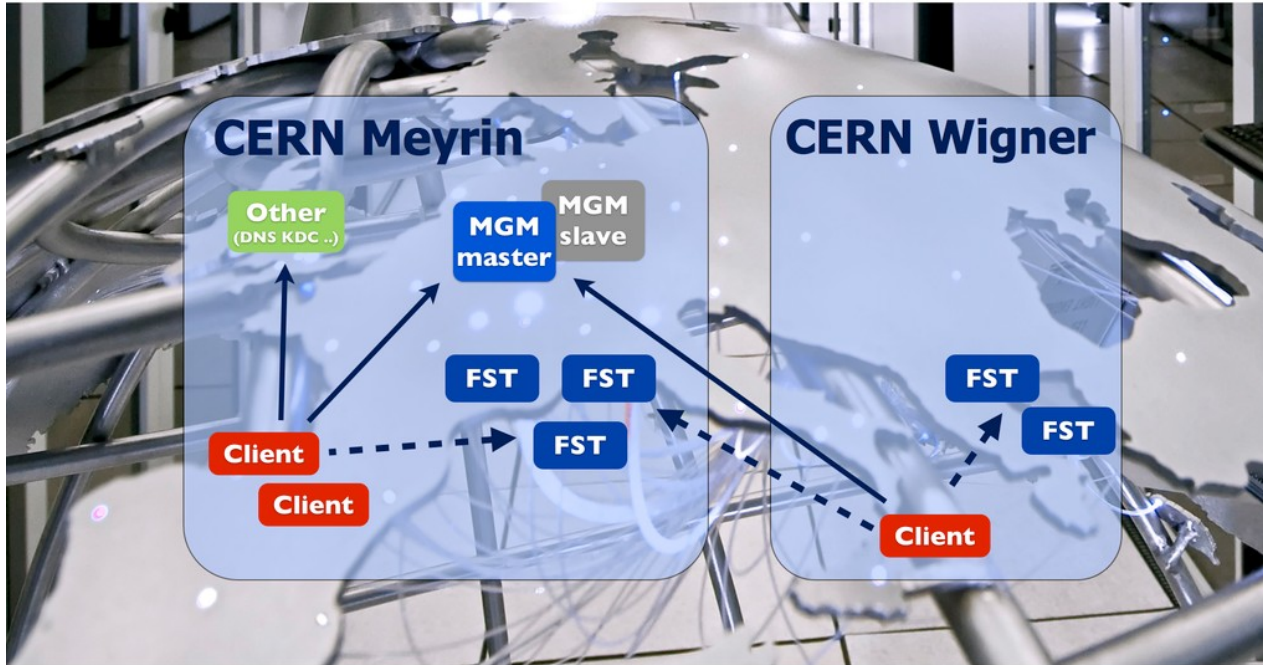
- What happens above ~ 1B files?
 - OK if static
 - What if the rate of changes goes up?
- CERNBOX
 - Marry EOS (EOSUSER) with your laptop
 - Flexibility of the sync client
 - Power of a file system (e.g. LXBATCH jobs using the EOS as a file system)
 - Both views coincide!

EOS installed across the CERN computer centres

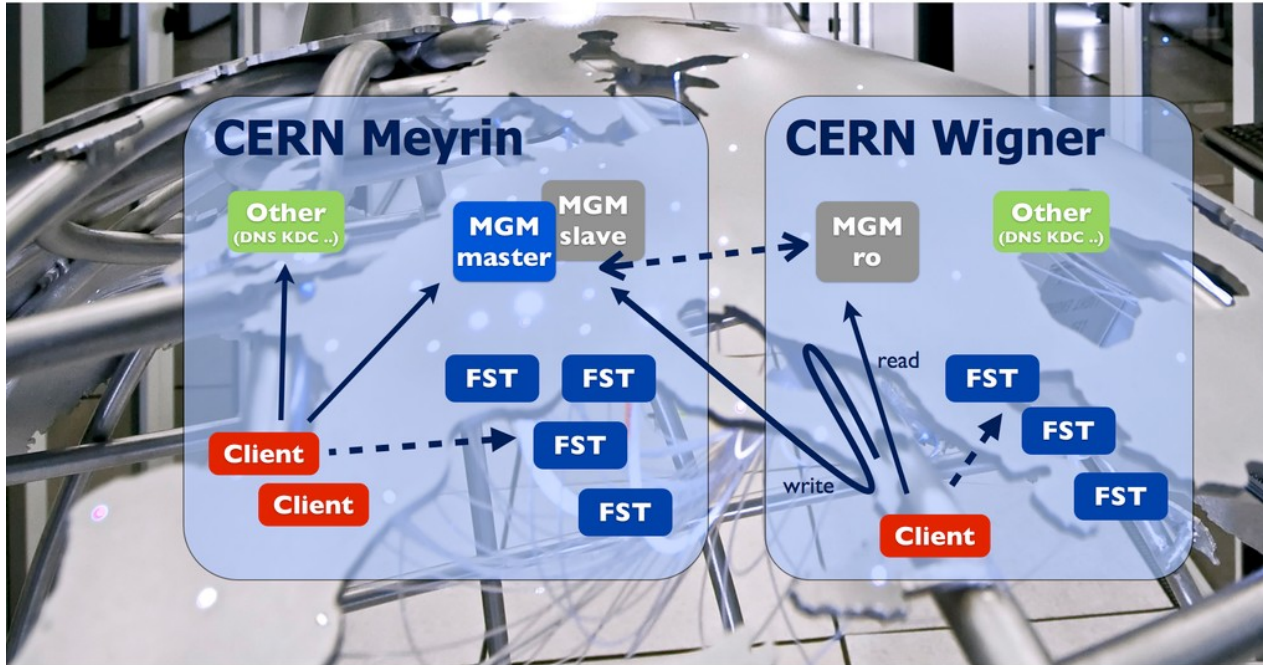
- EOS takes advantages of the two CERN computer centres
 - Coping with ~20-ms latency
 - Distributing copies across the two sites for dependability and performance
- Status
 - As today we are crossing the 40% mark



EOS 2014 Deployment



EOS 2015 Deployment





DSS

Challenges in distributed installations

- 2 sites across 1000 km: OK
- Can we do more?
- Locality in a several-site installation
- EOS hit-the-road-jack kit :)
- ...

Workshop on Cloud Services for File Synchronisation and Sharing

All material is online (slides and presentations recording)

>80 participants

Broad interest (large participation outside our “traditional” community)

Several companies

CERNBox presented

CERN 17-18 November 2014

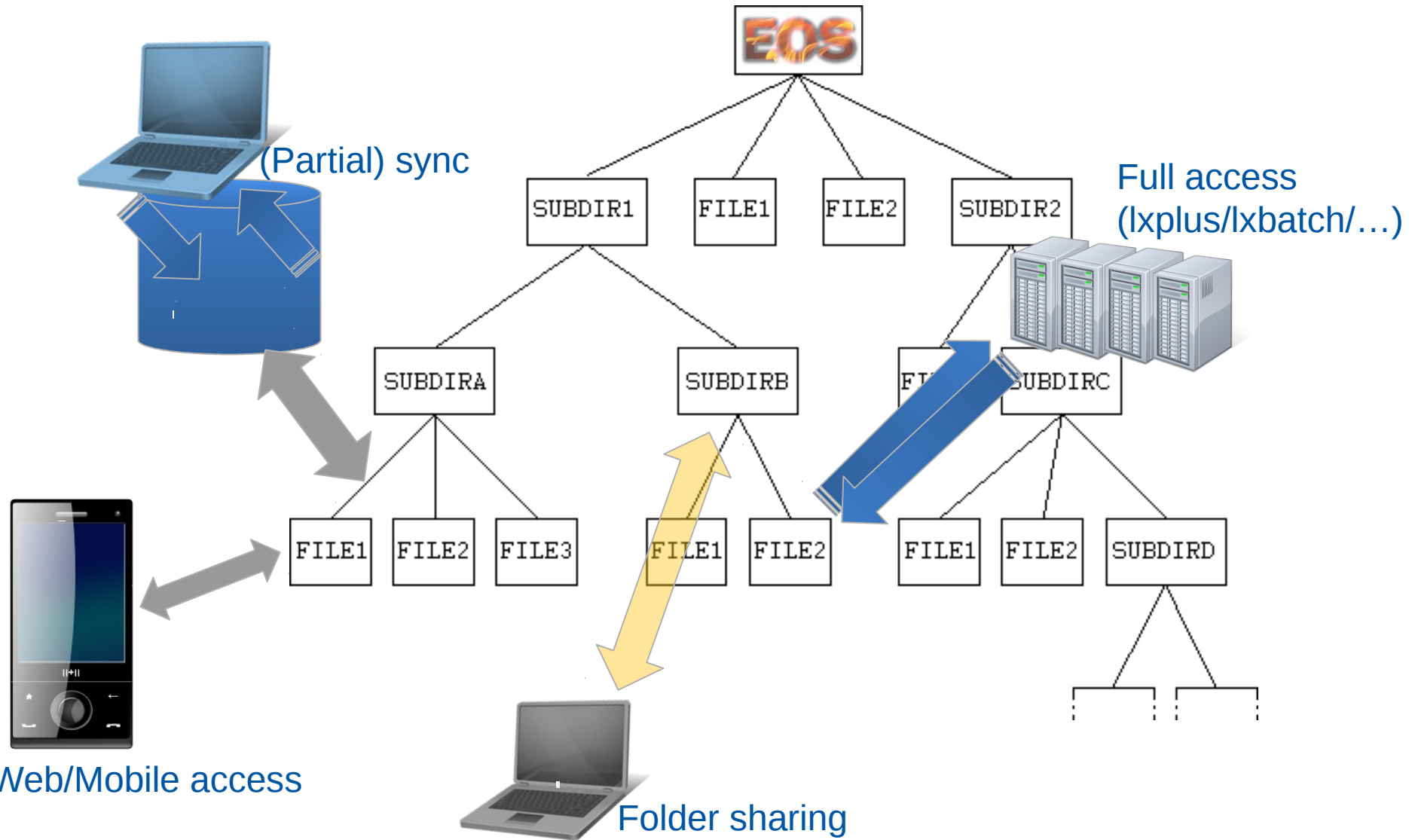
<https://indico.cern.ch/event/336753/>
Abstract submission: 30 September
Registration: 31 October

ORGANISING COMMITTEE
Miguel Branco
Massimo Lamanna
Jakub T. Moscicki

** We have modify ObjectStoreStorage to avoid de-frees and create our dist Cache, because the not an interface implement.*

Diagram description: A hand-drawn UML class diagram on a chalkboard. It shows an 'Interface IFilesStorage' extending from 'Interface IObjectStore'. 'IFilesStorage' is implemented by 'Abstract class IObjectStorageCommon', which is further implemented by 'class IObjectStorage' and 'class IObjectStorageStorage'. 'IObjectStorage' is implemented by 'class IObjectStorageStorage' and 'class EOS'. 'IObjectStorageStorage' is implemented by 'class DropBox, Sugar, S3'. 'DropBox, Sugar, S3' is labeled as 'external fs'. 'EOS' is marked with an asterisk. A note explains that 'ObjectStoreStorage' was modified to avoid de-frees and create a dist cache because it's not an interface implementation.

CERNBox



Conclusion

- EOS is our flagship project
- Essential for the LHC Run 2
 - Revised utilisation of all DSS projects
- EOS evolution
 - Number of challenging questions
 - Important for LHC and non-LHC usage