



ECFS

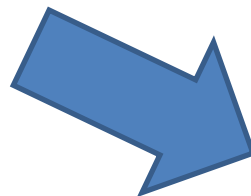
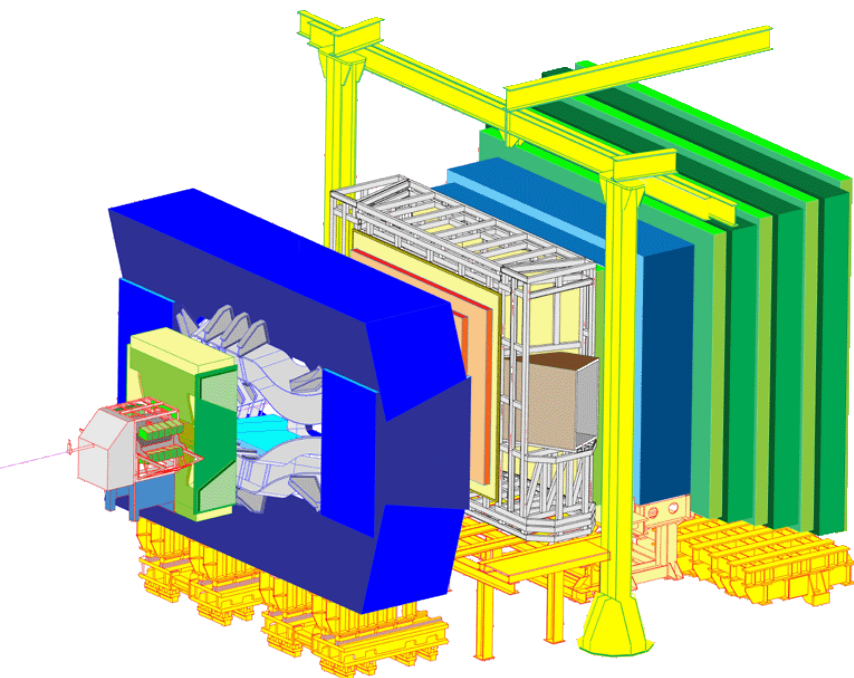
A decentralized, distributed and fault-tolerant FUSE filesystem for the LHCb online farm



Session: Data Stores, Data Bases,
and Storage Systems

tomasz.rybczynski@cern.ch*
enrico.bonaccorsi@cern.ch
niko.neufeld@cern.ch

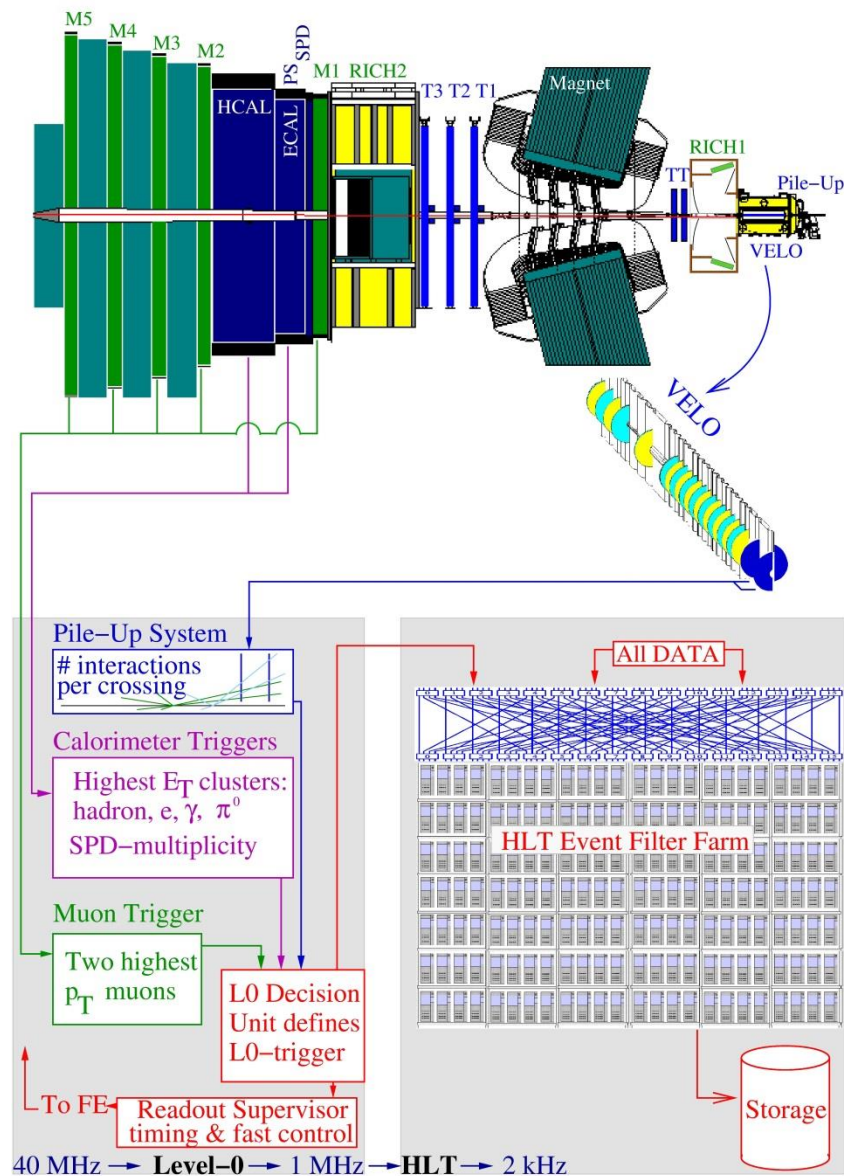
Introduction



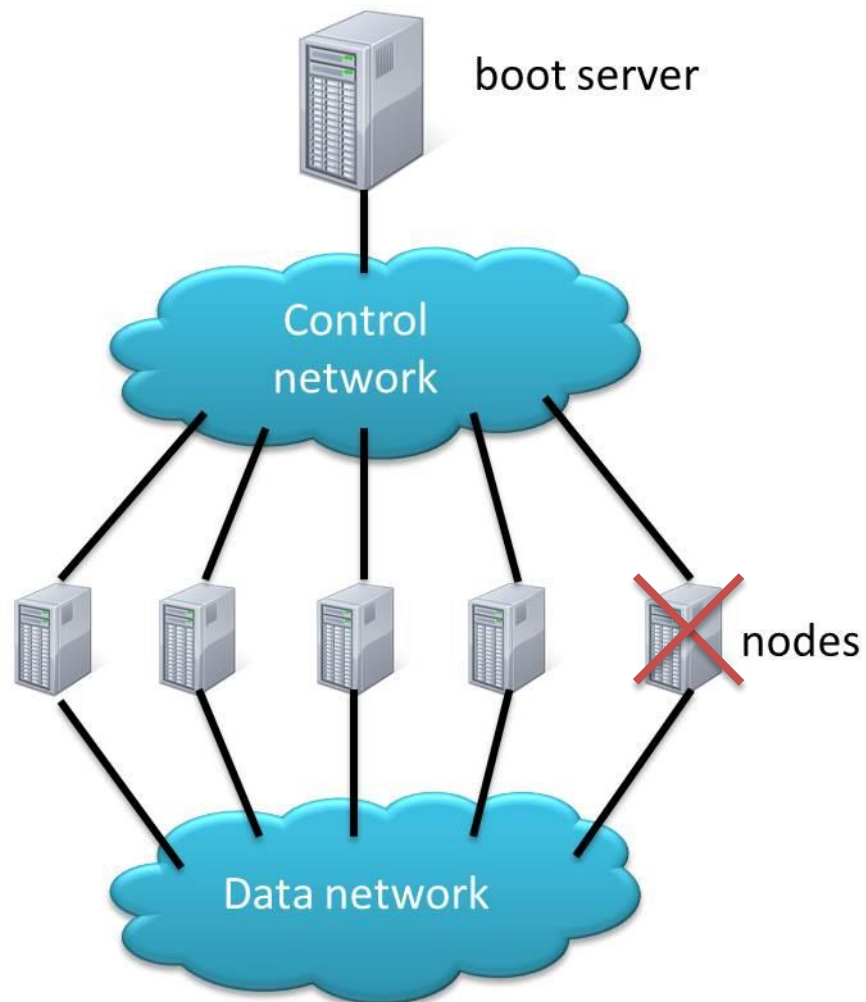
Event Filter Farm
 ~2000 nodes



High Level Trigger



Event Filter Farm



20-25% of the time

Requirements

- Fault tolerance – data redundancy
- Single namespace
- POSIX semantics
- Write once read many
- Sufficient performance

Searching for a solution

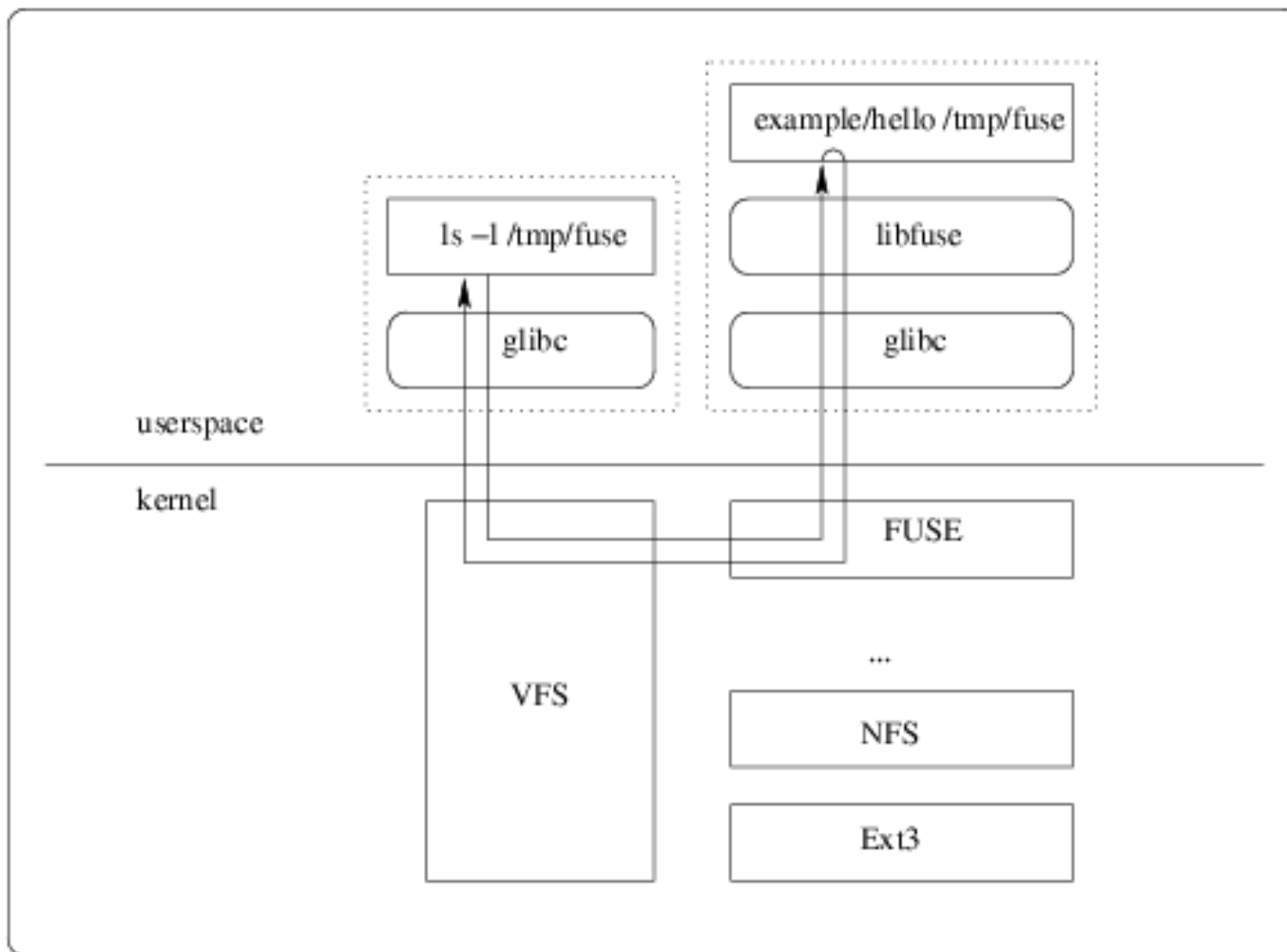
- Linux RAID over HA cluster
- GlusterFS
- Tahoe-LAFS tahoe-lafs.org
- ...

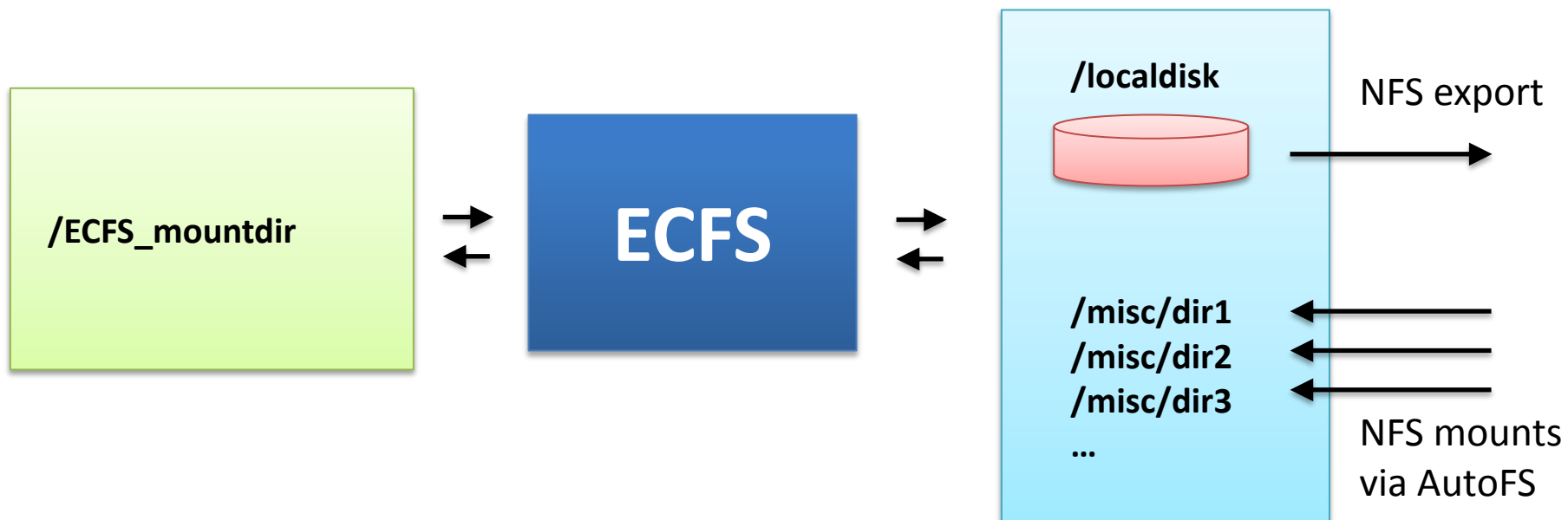
ECFS - distributed, fault-tolerant file system

Project based on:

- FUSE
- NFS
- Erasure Codes

FUSE





Coding parameters:

$$k = 6$$

$$n = 2$$

Original message of a size of k symbols



Encoding



Code word of a size of $k + n$ symbols

Successful decoding possible if the number of missing symbols is less or equal n

Code word with some symbols missing



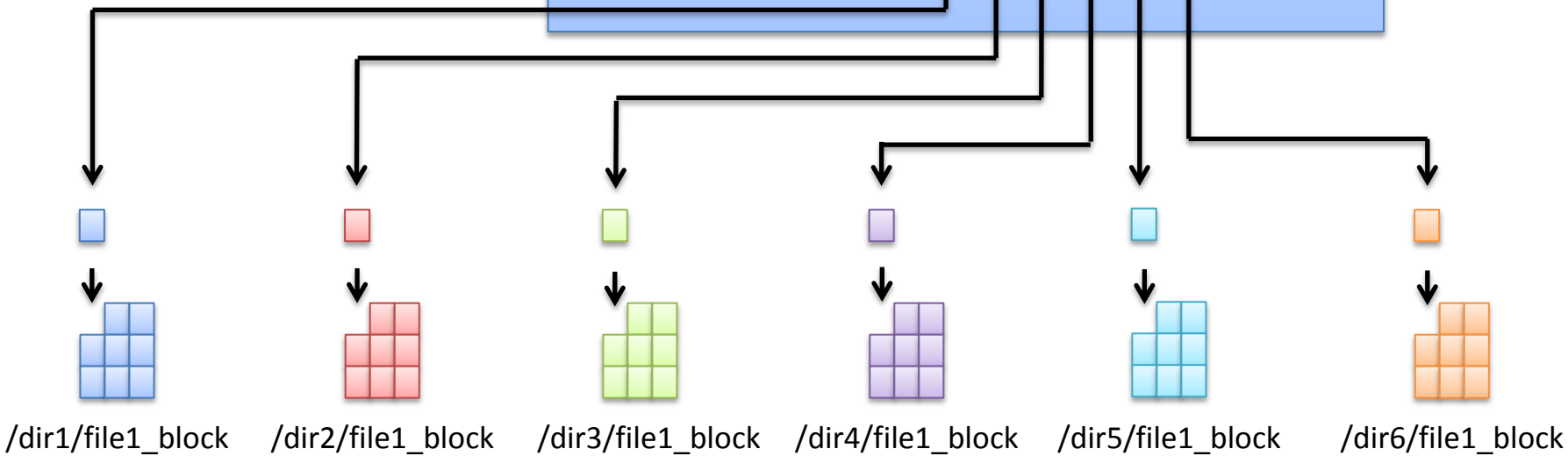
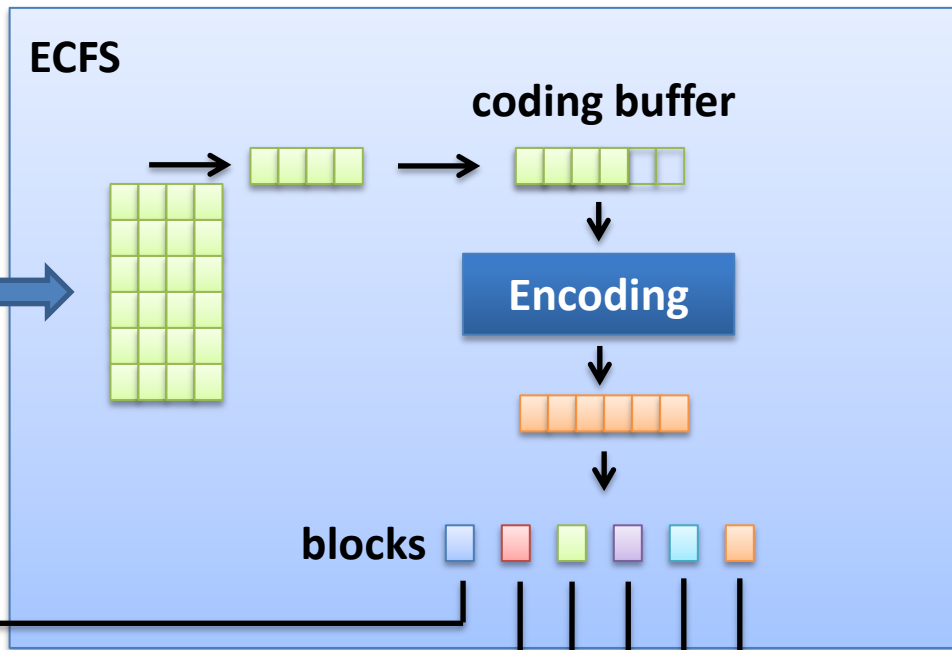
Decoding



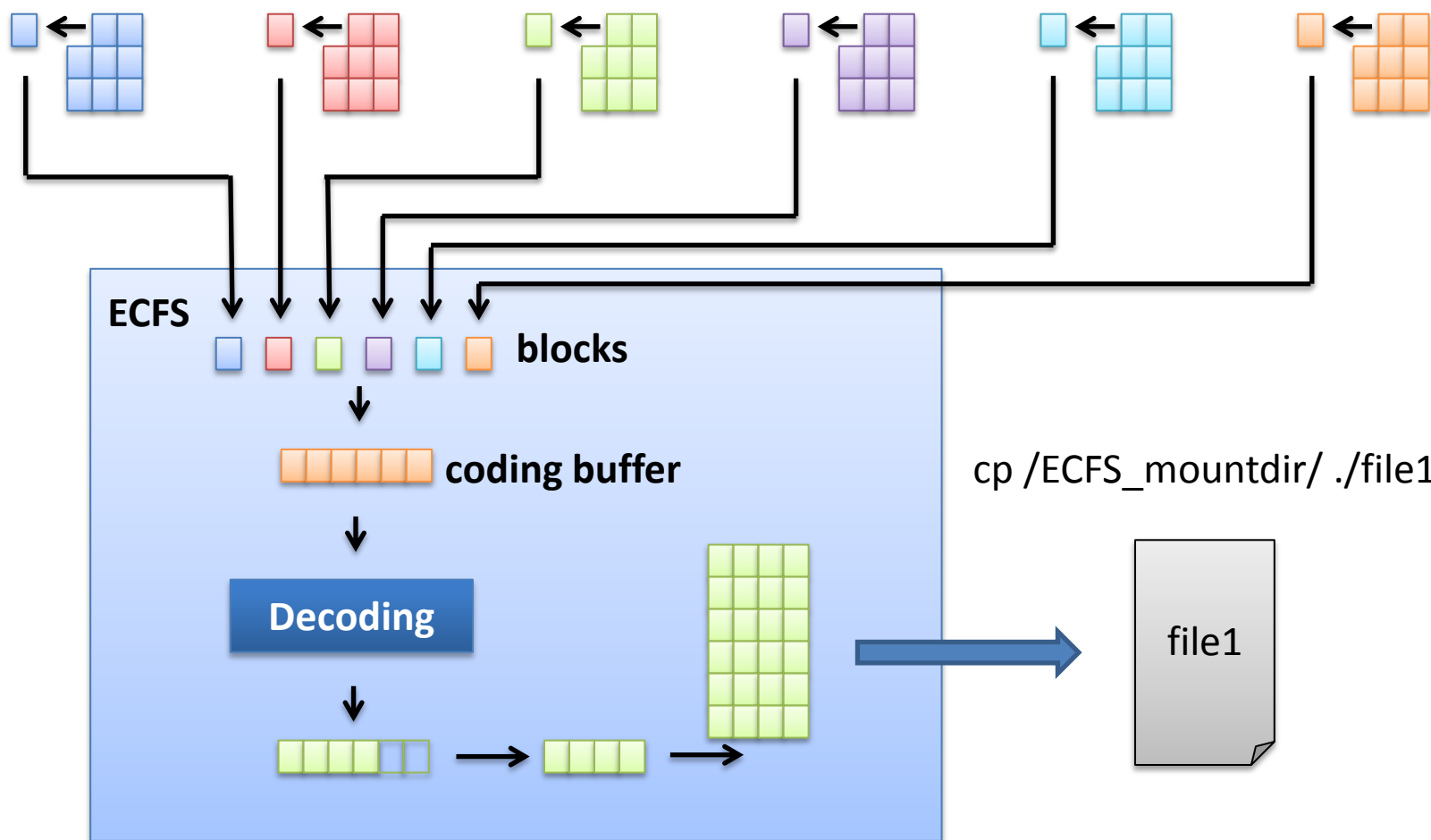
Original message

File encoding

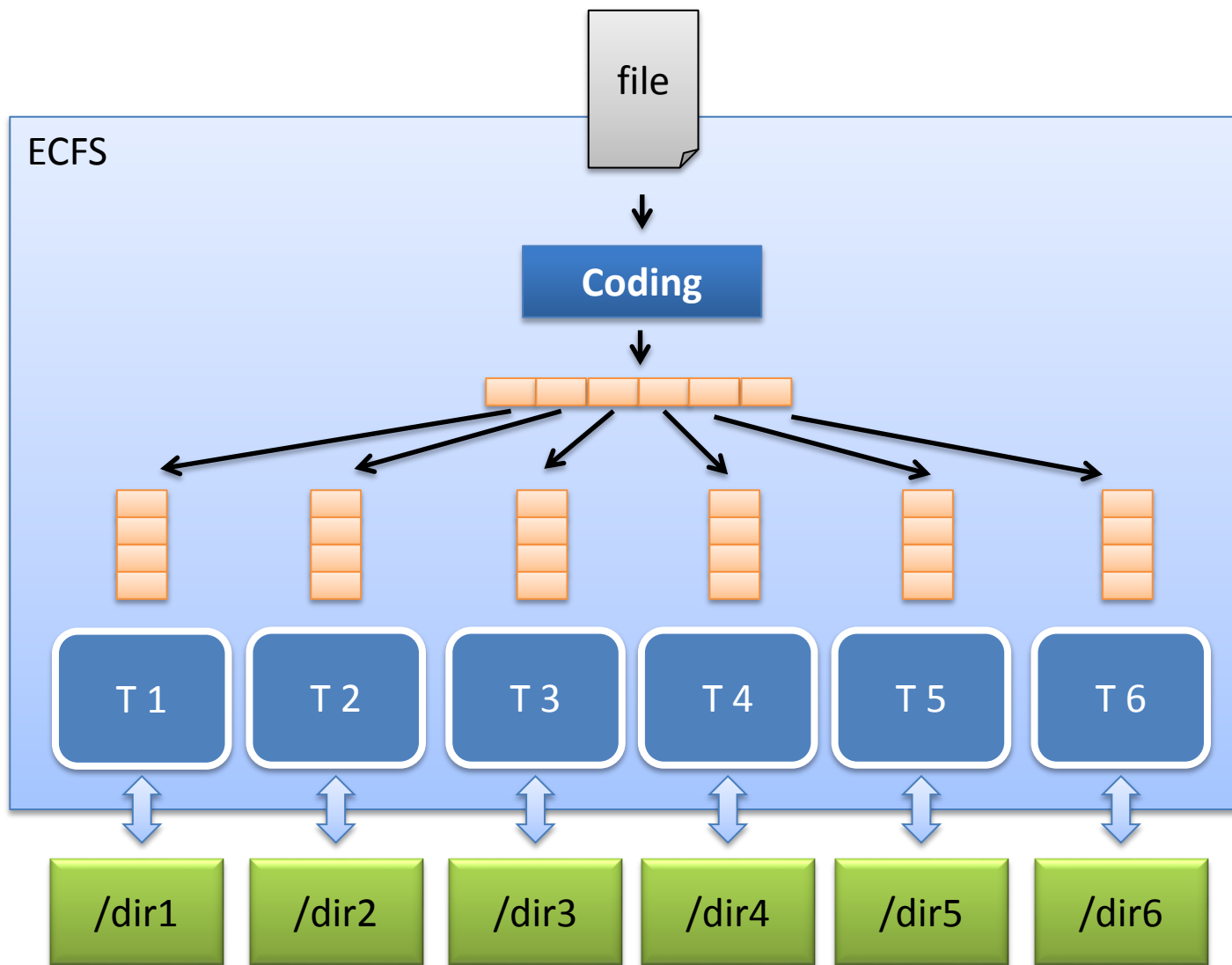
cp file_1 /ECFS_mountdir/



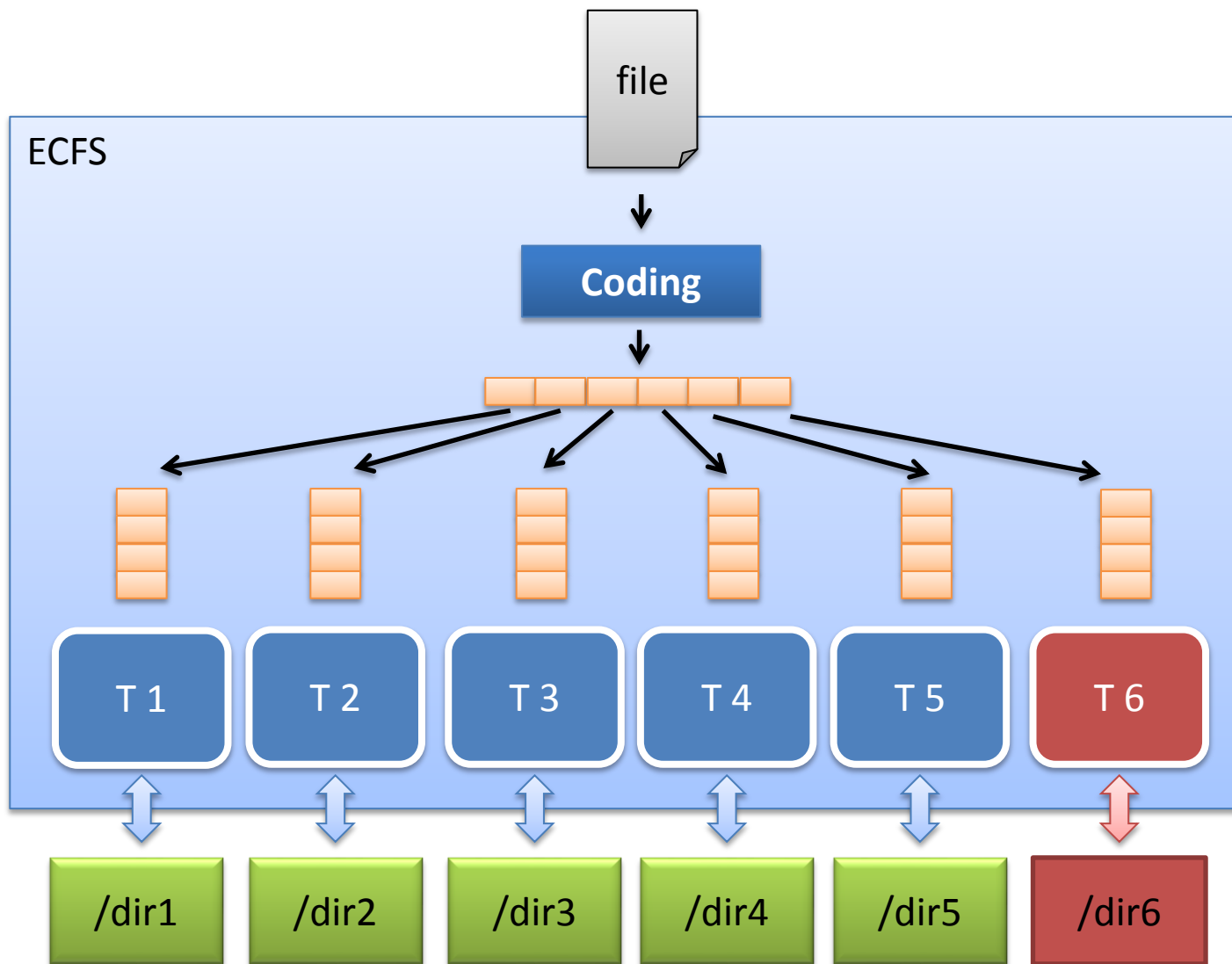
/dir1/file1_block /dir2/file1_block /dir3/file1_block /dir4/file1_block /dir5/file1_block /dir6/file1_block



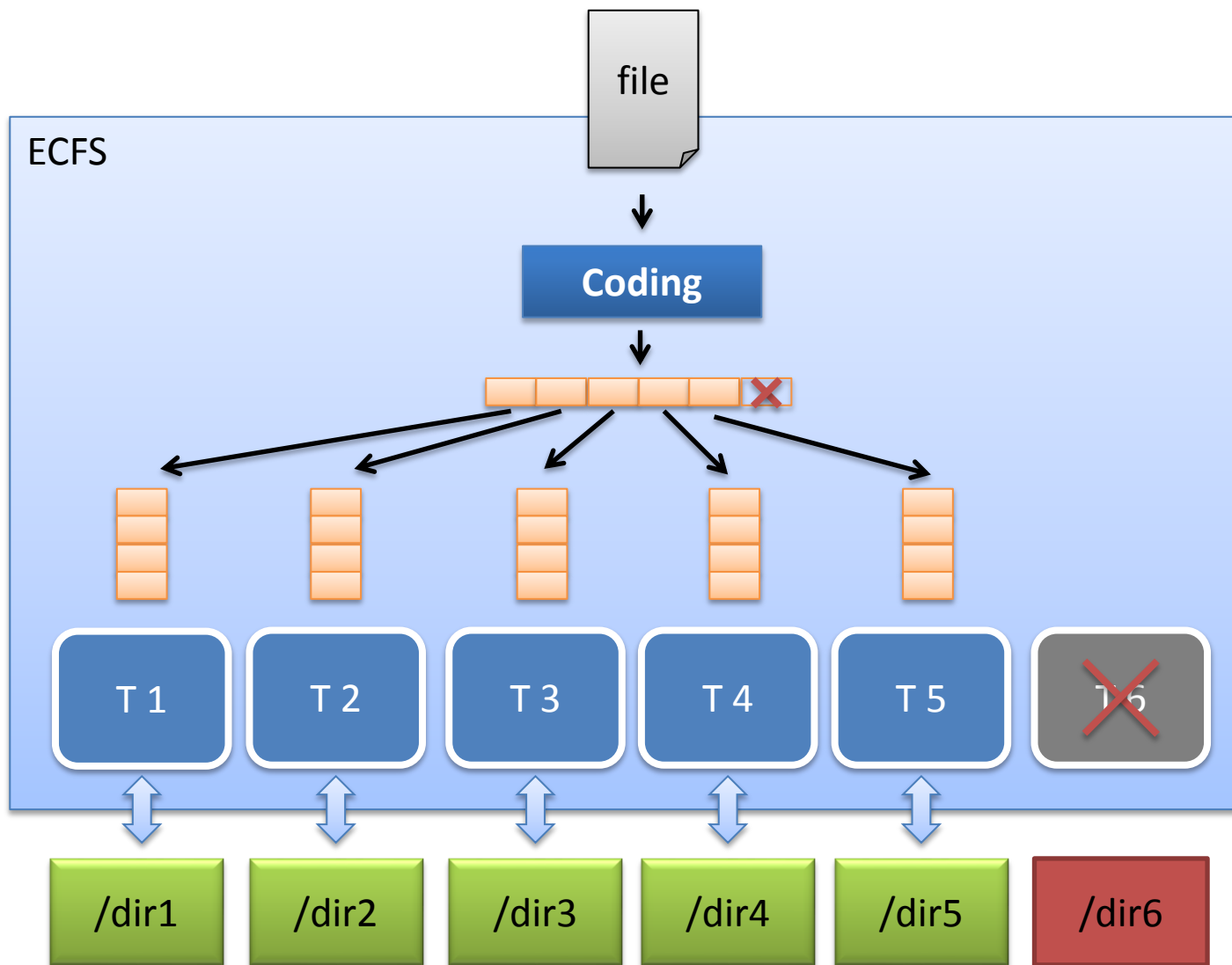
Parallel IO



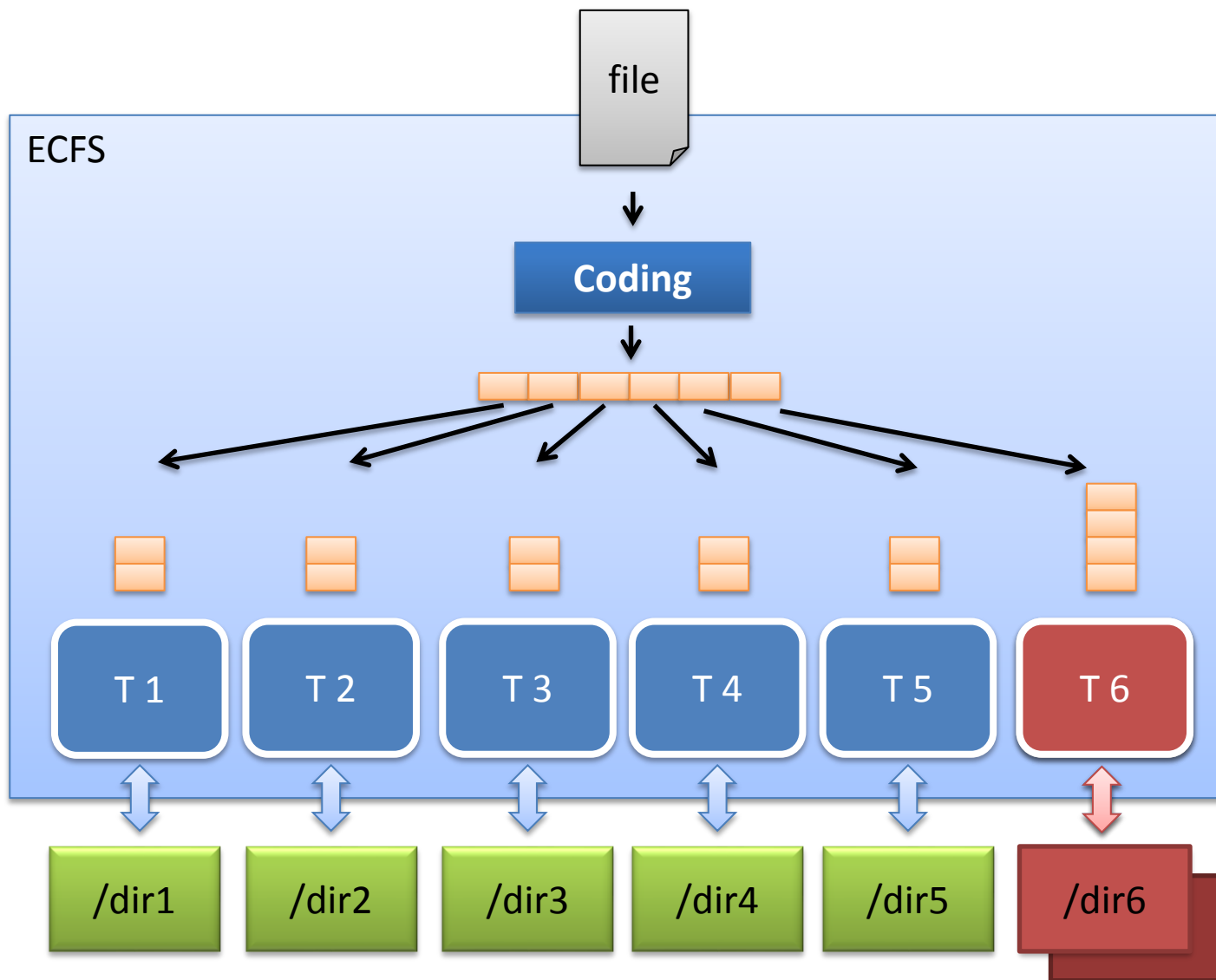
Fault detection

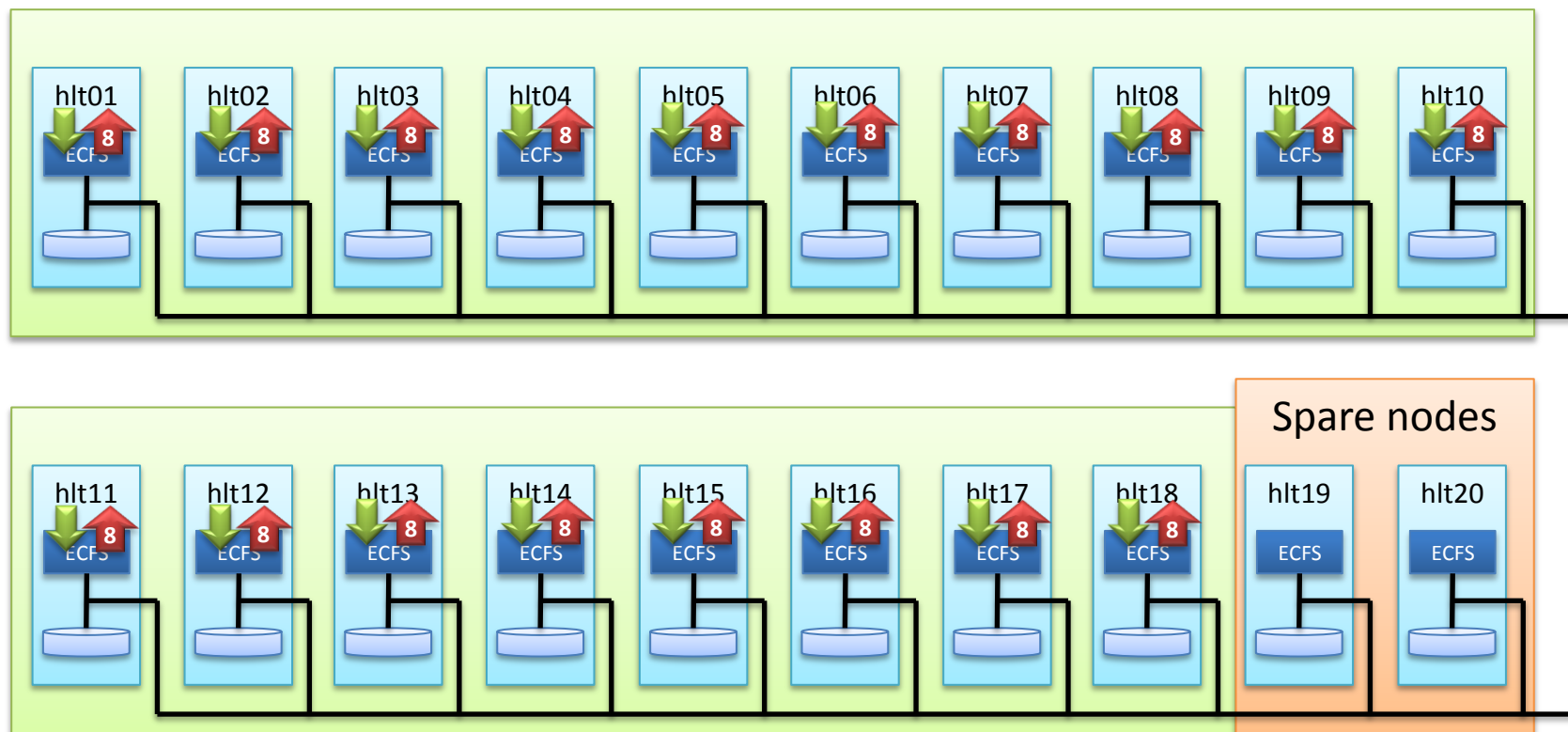


Degraded write



Write to a spare

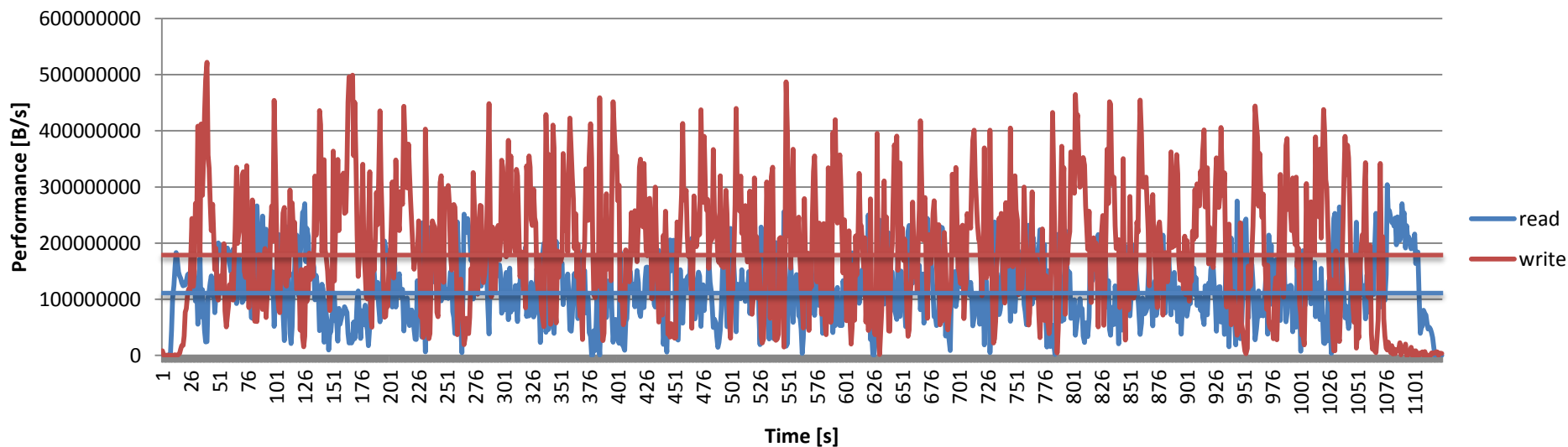




Each of the nodes:

- Writes one file
- Reads eight files

Aggregated HDD performance



Average write: 188 MB/s ~ 10 MB/s per node

Average read: 119 MB/s ~ 7 MB/s per node

Conclusion

- ECFS meets technical requirements of the project
- Performance lower than expected
- Much improvement has to be done

