# Ensuring Data Durability in EOS Systems

Maria Arsuaga-Rios

IT-ST-PDS
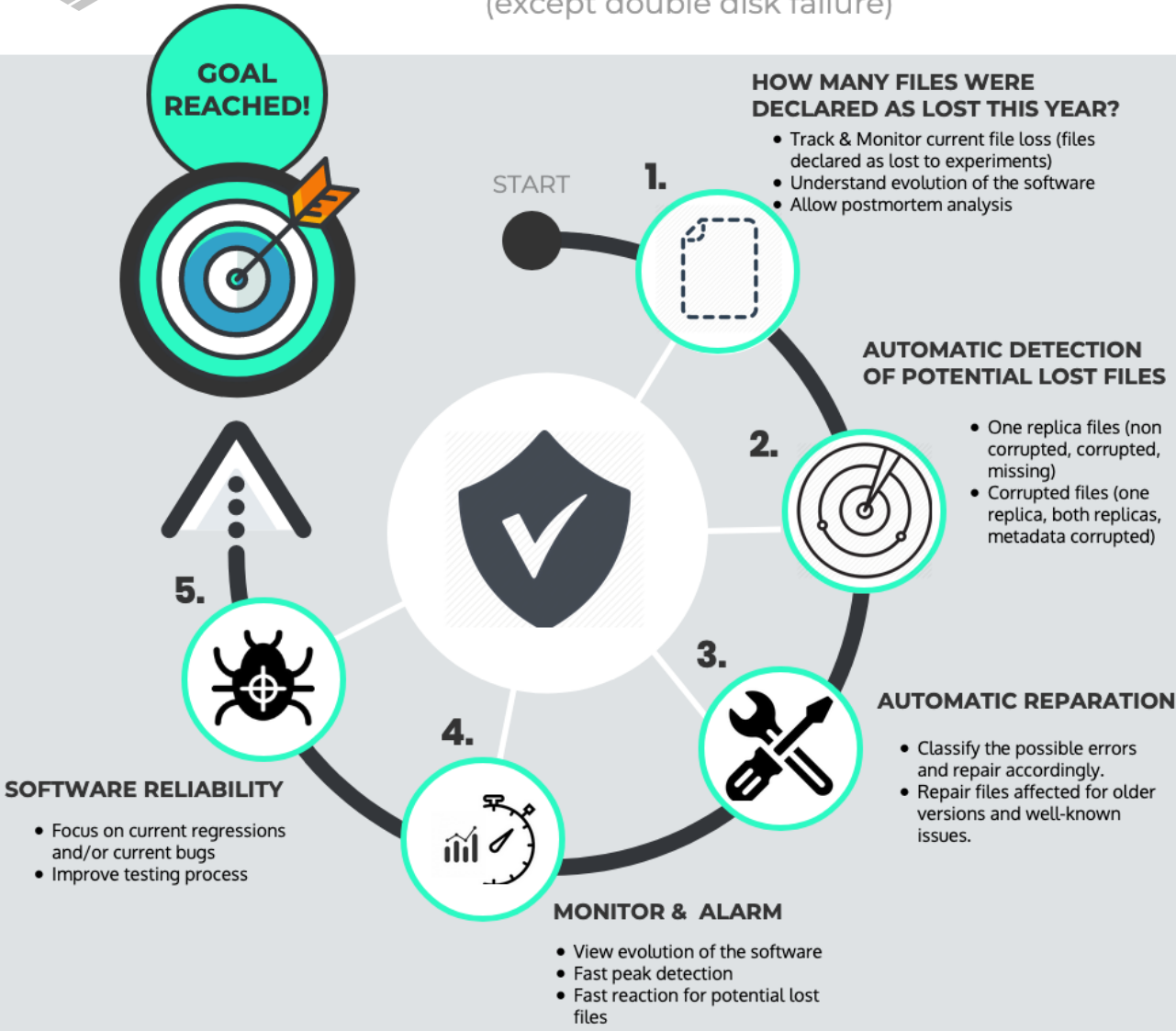
Information Technology Department

# DATA DURABILITY

Data is never lost or compromised
(except double disk failure)

**GOAL REACHED!**

**START**

**1.**
**HOW MANY FILES WERE DECLARED AS LOST THIS YEAR?**
- Track & Monitor current file loss (files declared as lost to experiments)
- Understand evolution of the software
- Allow postmortem analysis

**2.**
**AUTOMATIC DETECTION OF POTENTIAL LOST FILES**
- One replica files (non corrupted, corrupted, missing)
- Corrupted files (one replica, both replicas, metadata corrupted)

**3.**
**AUTOMATIC REPARATION**
- Classify the possible errors and repair accordingly.
- Repair files affected for older versions and well-known issues.

**4.**
**MONITOR & ALARM**
- View evolution of the software
- Fast peak detection
- Fast reaction for potential lost files

**5.**
**SOFTWARE RELIABILITY**
- Focus on current regressions and/or current bugs
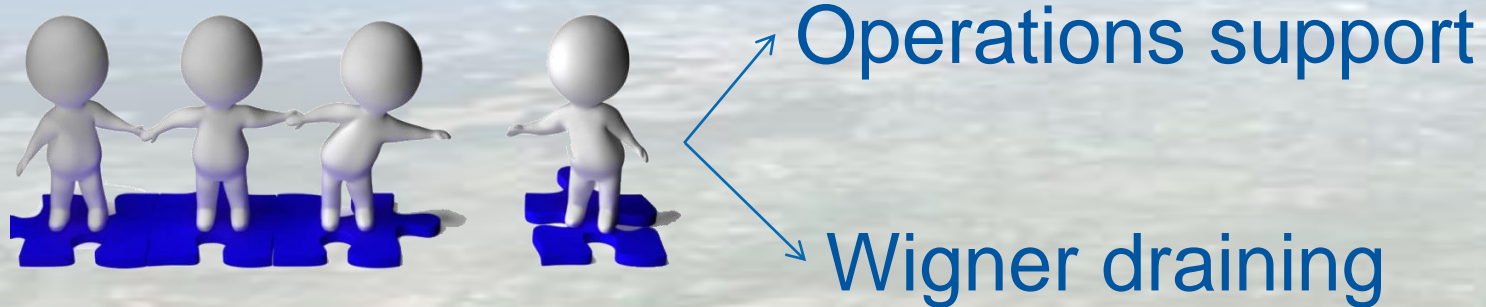- Improve testing process

EOSPHYSICS: minimal fraction of files may be lost "by design" in certain situations (e.g. double disk failure) — this is a deliberate choice due to the sheer size of the system (100s PBs) and it is trade-off between the cost of storage and required reliability.

EOSUSER: we guarantee better reliability (and we don't expect the file loss at all) by increasing the cost of the service. This may be done in several ways (e.g. higher replication, additional backup, …).

Information Technology Department

# Flashback: March 2019

Operations support

Wigner draining

How many files are we missing?

Is there any tool to allow us post-mortem/historical analysis?

# Missing files tracking

- Track and monitor files declared as lost
- Understand evolution of the software
- Allow post-mortem analysis

**HOW MANY FILES WERE DECLARED AS LOST THIS YEAR?**

- Track & Monitor current file loss (files declared as lost to experiments)
- Understand evolution of the software
- Allow postmortem analysis

START 1.

Lost files
Metadata
(eos-ops-lostfiles)

elasticsearch

Grafana — Monitoring

kibana — Data discovery

CERNBox

eos-ops-lostfiles -f file -d 2020-02-04 -r "bootfailure_with_other_replica_0_size" -e ALICE --send

# Missing files tracking

```
send_data_to_es(data, index, _id, es):
    requestBody = ""
    requestBody+='{"index":{"_id": "'+_id+'"}}\n'
    requestBody+= json.dumps(data) +"\n"

    try:
        out = es.bulk(index='eosmon_'+index+'-
'+str(datetime.datetime.now().year), doc_type=index, body=requestBody)
        print out
    except TransportError as e:
        print e.info
    except Exception as er:
        print er

get_es_connexion():
    USER_ES = 'eosmon'
        PASSWORD_ES =  lalala
        certpath = "/etc/pki/tls/certs/ca-bundle.trust.crt"
         es = Elasticsearch(
        'https://'+USER_ES+':'+PASSWORD_ES+'@es-
eosmon.cern.ch/es',
        # turn on SSL
        use_ssl=True,
        # make sure we verify SSL certificates (off by default)
        verify_certs=True,
        ca_certs=certpath)
      return  es
```



Monitoring

Data discovery

Information Technology Department

# Missing files tracking

# Missing files tracking

Monitoring

Data discovery

# 7th May 2019 - we discovered 32K* missing files when draining ALICE

- How we communicate these information to our users?

  1. How many condition data base files are missing?

     - (5 replicas expected)

  2. Which is the distribution over mtime?

  3. Is there a correlation with an incident, ticket or known bug?

*: 0.0046% files stored

IT Information Technology Department

# 7th May 2019 - we discovered 32K missing files when draining ALICE

experiment: "ALICE"    date: "May 7th 2019, 00:00:00.000"    query: "{"wildcard":{"file":"/eos/alice/cond/*"}}"    Add a filter +

**eosmon_lost_files***

Data | Options

**Metrics**

Metric                                    Count

Add metrics

**Buckets**

Select buckets type

Split Group

Cancel

1. How many condition data base files are missing?

   (5 replicas expected)

**780**
Count

Information Technology Department

2. Which is the distribution over mtime?

2013 to 2017

# 7th May 2019 - we discovered 32K missing files when draining ALICE

3. Is there a correlation with an incident, ticket or known bug?

- 10/11/2017 Incident: "**Root cause:** Crash caused the namespace to be corrupted" OTG0040840
- 10/11/2017 to 14/12/2017 GGUS Ticket:"eosalice manual restore, missing condition files" GGUS ticket

Information Technology Department

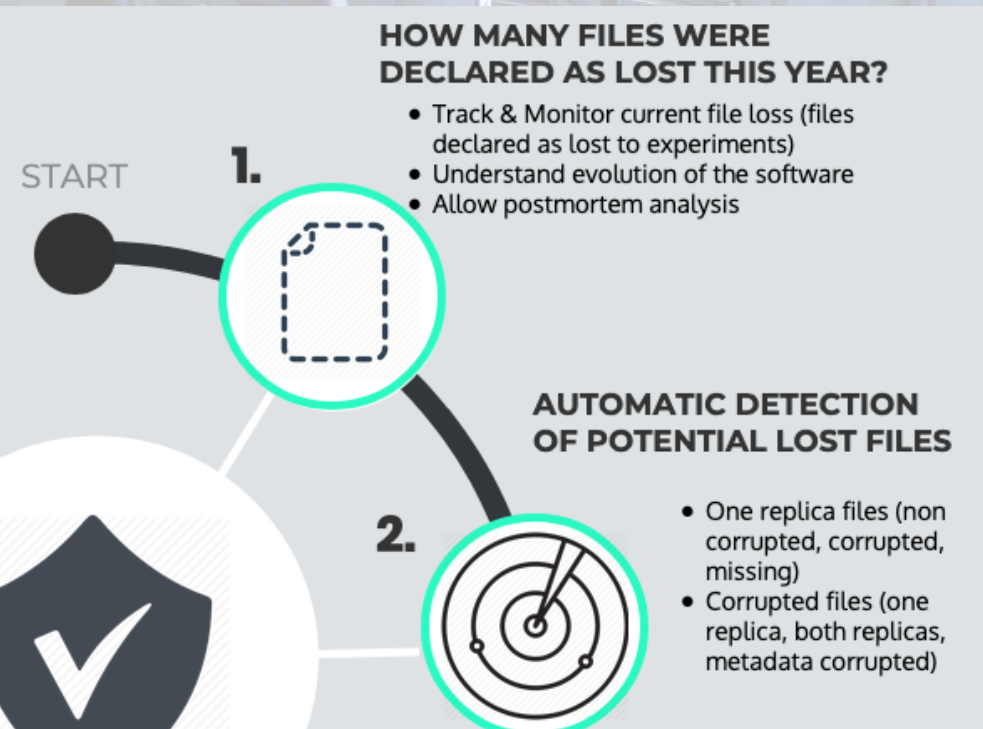# Automatic detection of *high risk* files



**HOW MANY FILES WERE DECLARED AS LOST THIS YEAR?**

START

1.

- Track & Monitor current file loss (files declared as lost to experiments)
- Understand evolution of the software
- Allow postmortem analysis

**AUTOMATIC DETECTION OF POTENTIAL LOST FILES**

2.

- One replica files (non corrupted, corrupted, missing)
- Corrupted files (one replica, both replicas, metadata corrupted)
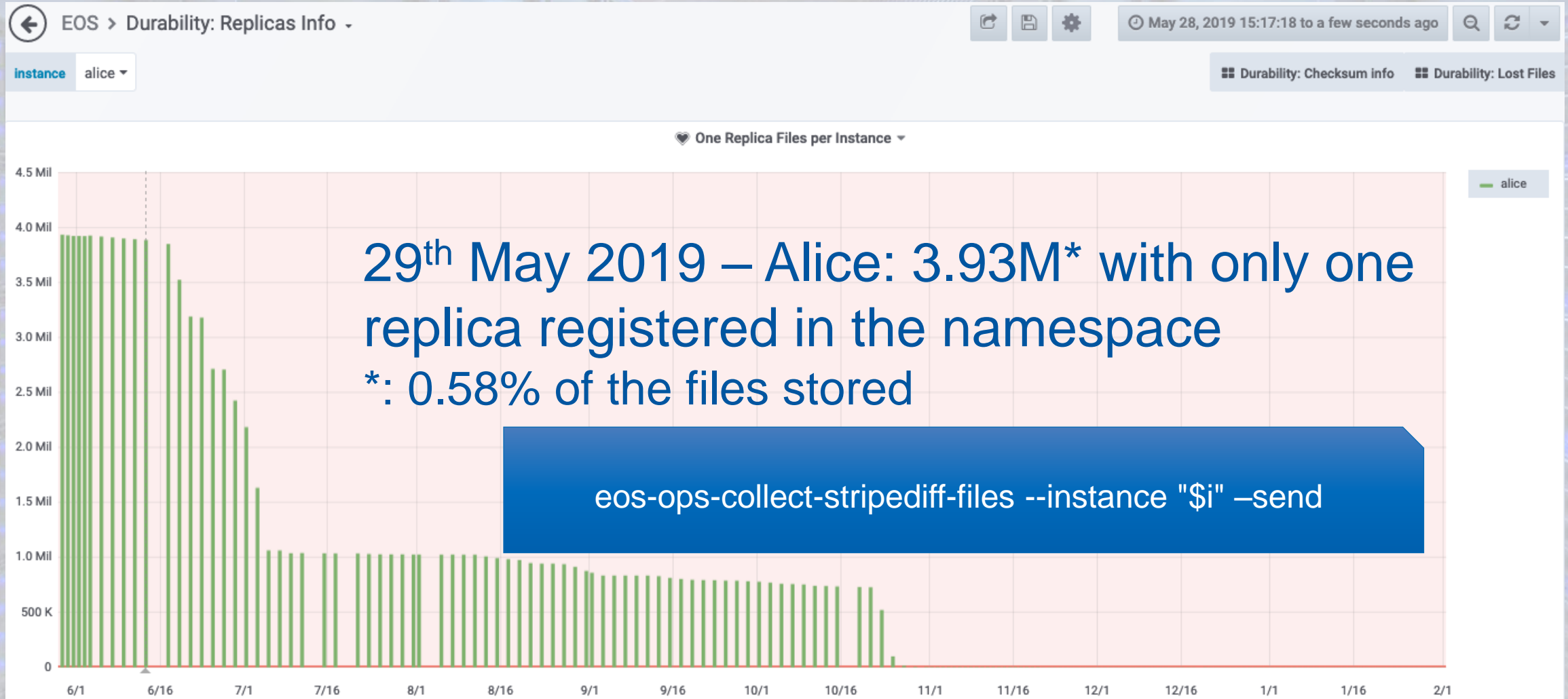
- **Detection is the first step, make it visible!**

  o Detect one replica files (non corrupted, corrupted, missing, …)

  o Detect mismatching checksums and sizes (one replica, both replicas, metadata corrupted, …)
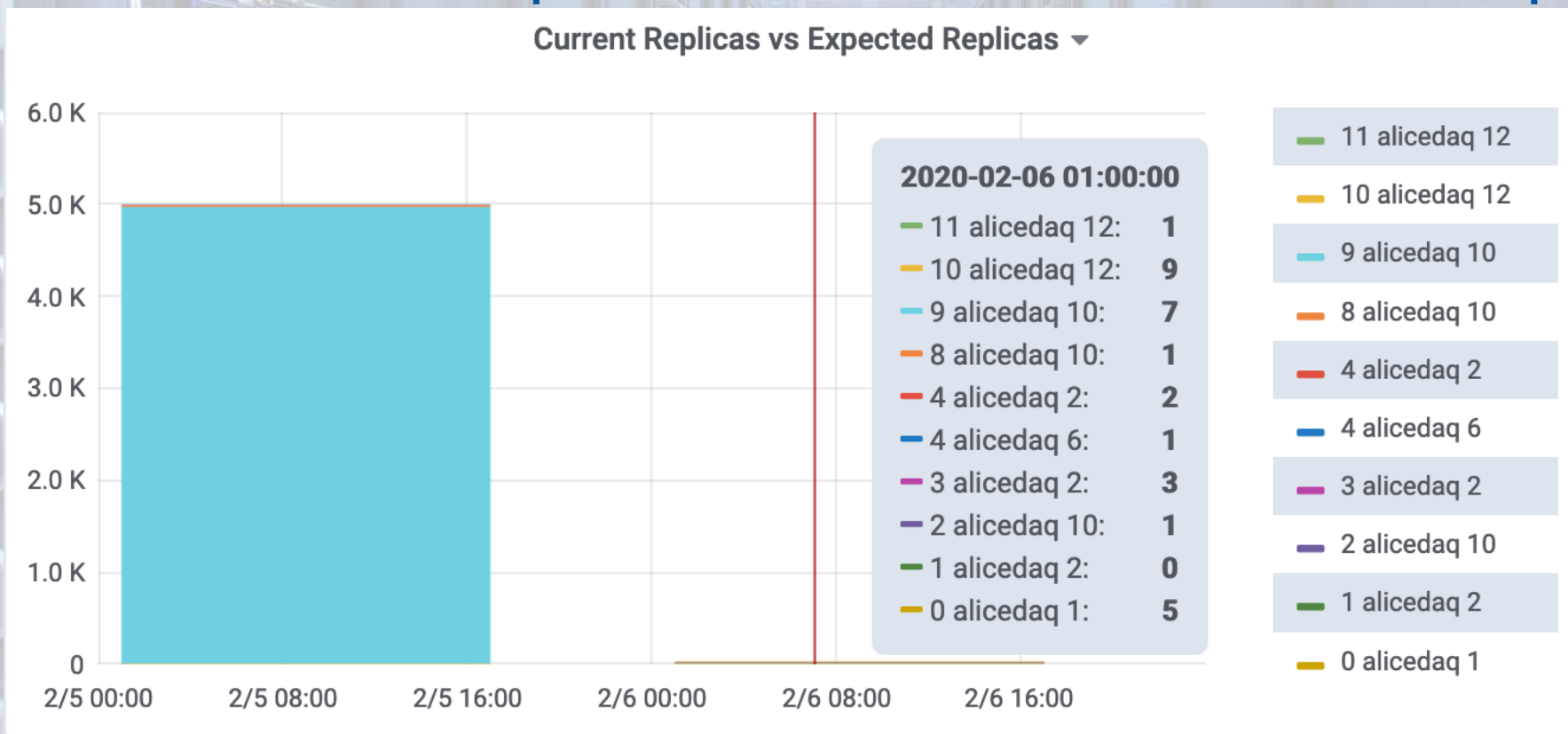
# Automatic detection of *high risk* files

- One replica files (one replica layouts included):
  - Automatic and daily full scan in all EOS quarkDB instances
    - ✓ stripediff option requested to Georgios for the eos-ns-inpect-tool

- Draining failures
  - Automatic and daily full scan in all file systems marked as drained failure (which have problematic files that prevent the completion of the draining process)

- Backup errors
  - Automatic and daily detection of files that couldn't be backup

# Automatic detection of *high risk* files



29th May 2019 – Alice: 3.93M* with only one replica registered in the namespace
*: 0.58% of the files stored

eos-ops-collect-stripediff-files --instance "$i" –send

# Automatic detection of *high risk* files

## Under-replication of Rain files in Alicedaq



eos-ops-collect-stripediff-files --instance "$i" –send

# Automatic detection of *high risk* files
## Draining failures – last month (January 2020)



`eos-ops-collect-drain-failed -i $instance --repair --send`

Information Technology Department

# Automatic reparation



**HOW MANY FILES WERE DECLARED AS LOST THIS YEAR?**

1.
- Track & Monitor current file loss (files declared as lost to experiments)
- Understand evolution of the software
- Allow postmortem analysis

**AUTOMATIC DETECTION OF POTENTIAL LOST FILES**

2.
- One replica files (non corrupted, corrupted, missing)
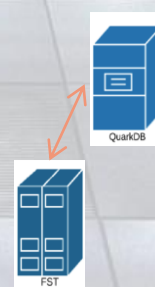- Corrupted files (one replica, both replicas, metadata corrupted)

**AUTOMATIC REPARATION**

3.
- Classify the possible errors and repair accordingly.
- Repair files affected for older versions and well-known issues.

- Namespace full scan is not enough
  - ○ We need to go deeper and get the storage nodes information

- Classify the possible errors and repair accordingly
  - ○ **Divide & Conquer:** 14 categories for one replica files out of 21 categories in total

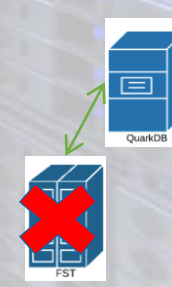- Repair files affected for older versions and well-known issues/cases

One replica and size and checksum match the namespace

One replica with checksum 0 in namespace

One replica with size and checksum mismatch with namespace

Missing replica

Information Technology Department

# Automatic reparation

**HOW MANY FILES WERE DECLARED AS LOST THIS YEAR?**

1.
- Track & Monitor current file loss (files declared as lost to experiments)
- Understand evolution of the software
- Allow postmortem analysis

**AUTOMATIC DETECTION OF POTENTIAL LOST FILES**

2.
- One replica files (non corrupted, corrupted, missing)
- Corrupted files (one replica, both replicas, metadata corrupted)
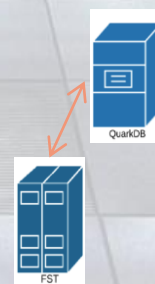
**AUTOMATIC REPARATION**

3.
- Classify the possible errors and repair accordingly.
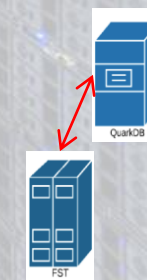- Repair files affected for older versions and well-known issues.

- Namespace full scan is not enough
  - We need to go deeper and get the storage nodes information

- Classify the possible errors and repair accordingly
  - **Divide & Conquer:** 14 categories for one replica files out of 21 categories in total

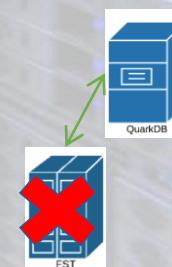- Repair files affected for older versions and well-known issues/cases

One replica and size and checksum match the namespace
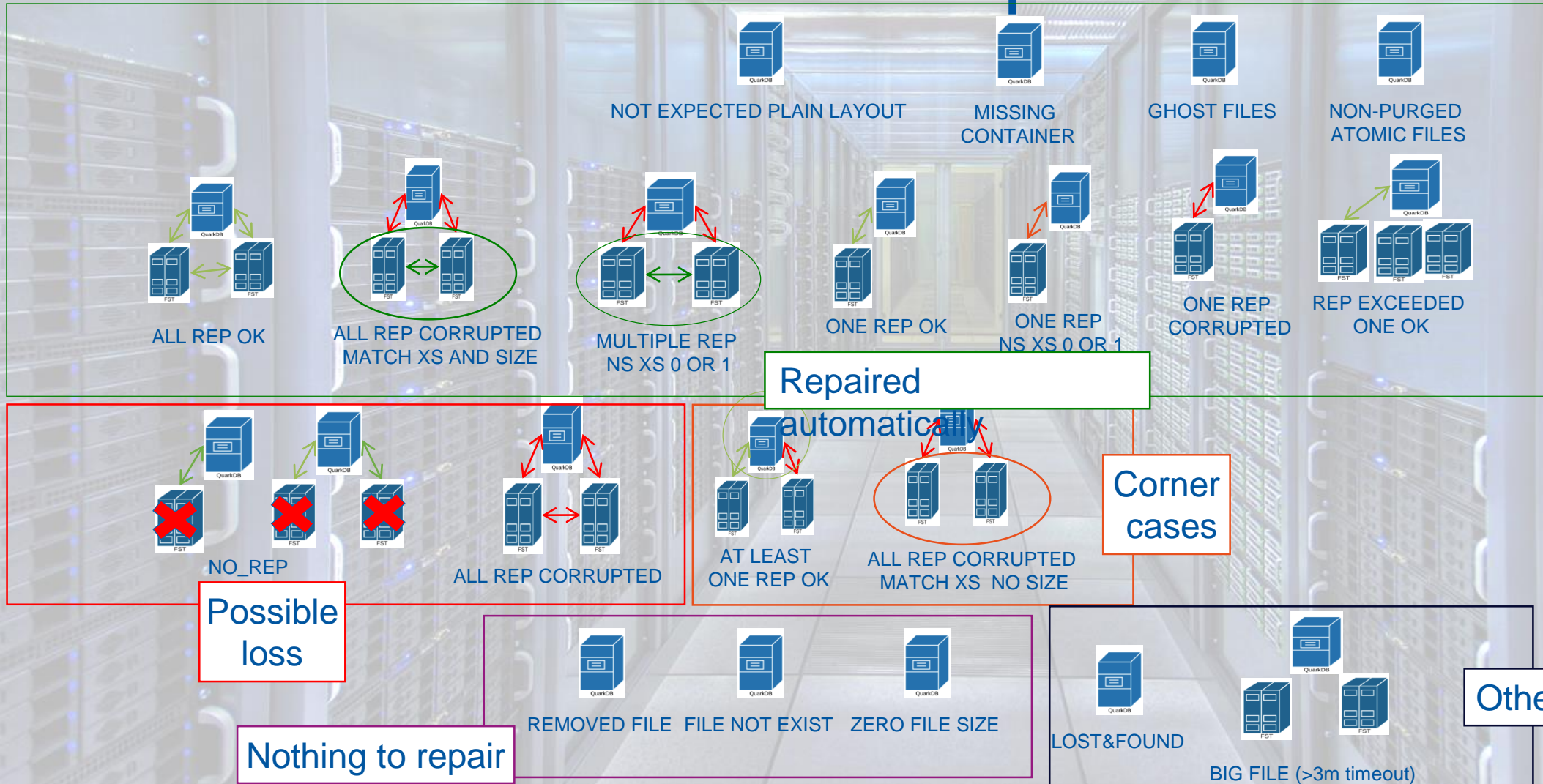
One replica with checksum 0 in namespace

One replica with size and checksum mismatch with namespace

Missing replica

Information Technology Department

# Automatic reparation

21 + 1 categories

NOT EXPECTED PLAIN LAYOUT

MISSING CONTAINER

GHOST FILES

NON-PURGED ATOMIC FILES

ALL REP OK

ALL REP CORRUPTED MATCH XS AND SIZE

MULTIPLE REP NS XS 0 OR 1

ONE REP OK

ONE REP NS XS 0 OR 1

ONE REP CORRUPTED

REP EXCEEDED ONE OK

Repaired automatically

Corner cases

NO_REP

ALL REP CORRUPTED

AT LEAST ONE REP OK

ALL REP CORRUPTED MATCH XS NO SIZE

Possible loss

REMOVED FILE    FILE NOT EXIST    ZERO FILE SIZE

LOST&FOUND

Other filters

Nothing to repair

BIG FILE (>3m timeout)

# Automatic reparation

21 + 1 categories

NOT EXPECTED PLAIN LAYOUT

MISSING CONTAINER

GHOST FILES

NON-PURGED ATOMIC FILES

ALL REP OK

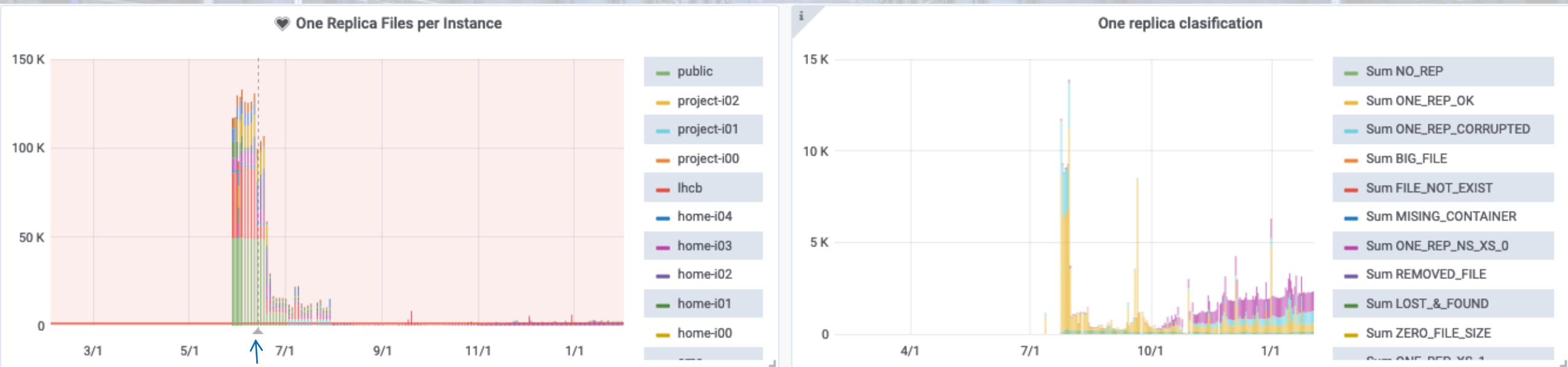ALL REP CORRUPTED MATCH XS AND SIZE

MULTIPLE
NS XS 0 C

Possible loss categories can have a chance!:
1. Retrieve finding replicas in other storage nodes (recycle bin, orphan files, lost&found)
2. Retrieve from versions
3. Retrieve from backup (eosbackup)
4. Retrieve from restic (independent system backup)
5. Retrieve from sync client

NO_REP

ALL REP CORRUPTED

Possible loss

Nothing to repair

REMOVED FILE    FILE NOT EXIST    ZERO FILE SIZE

LOST&FOUND

Other filters

BIG FILE (>3m timeout)

Information Technology Department
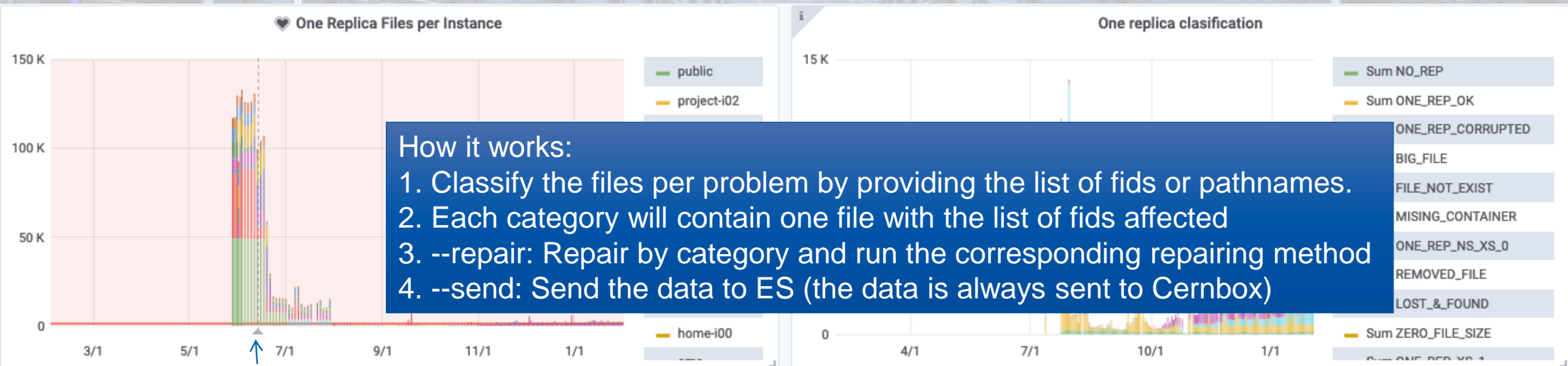
# Automatic reparation

- All instances (Alice excluded): One replica files classification and reparation



```
eos-ops-repair -f "$i"_one_replica_$(date +%Y-%m-%d -d "yesterday").txt --id_type dec -i "$i" -l DEBUG --rundeck
--one_rep --repair --send
```

Information Technology Department

# Automatic reparation

- All instances (Alice excluded): One replica files classification and reparation



How it works:
1. Classify the files per problem by providing the list of fids or pathnames.
2. Each category will contain one file with the list of fids affected
3. --repair: Repair by category and run the corresponding repairing method
4. --send: Send the data to ES (the data is always sent to Cernbox)
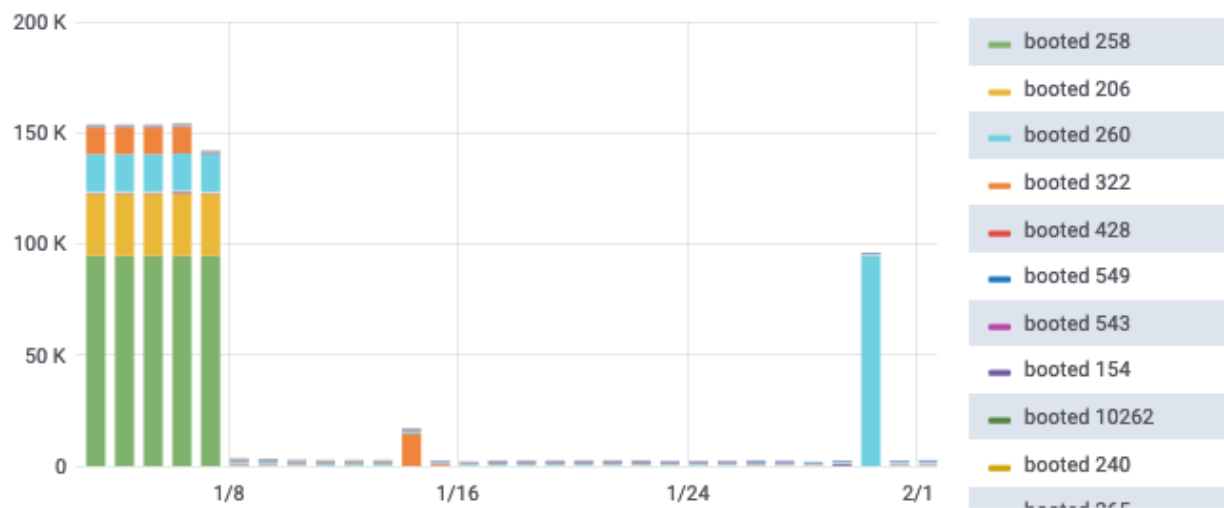
```
eos-ops-repair -f "$i"_one_replica_$(date +%Y-%m-%d -d "yesterday").txt --id_type dec -i "$i" -l DEBUG --rundeck
--one_rep --repair --send
```
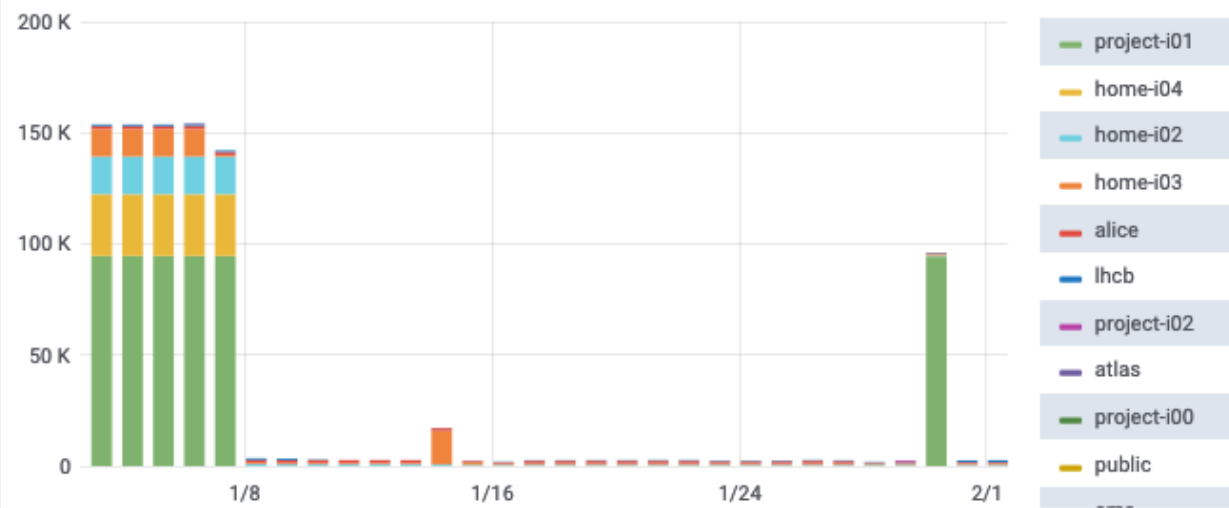
# Automatic reparation

- Automatic reparation for draining failures:

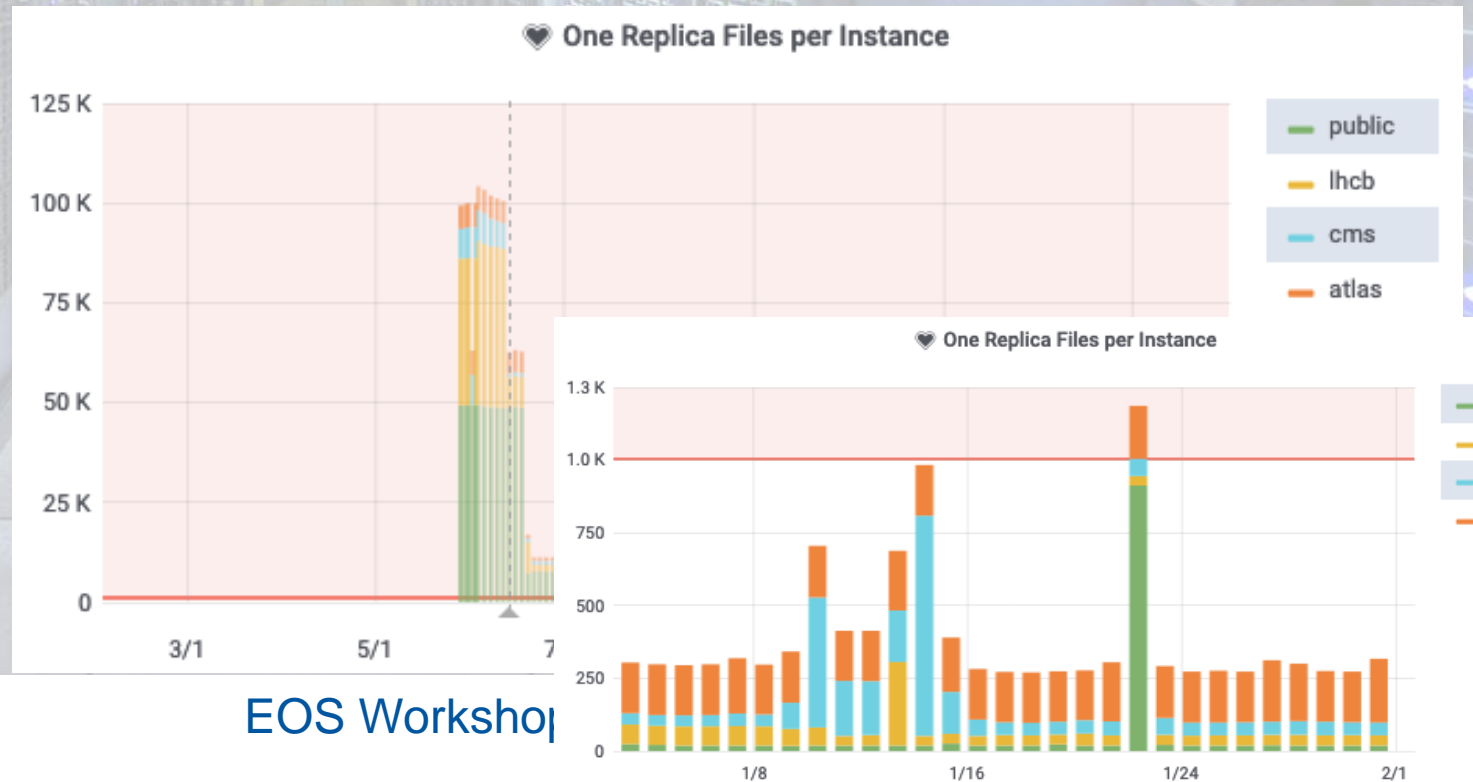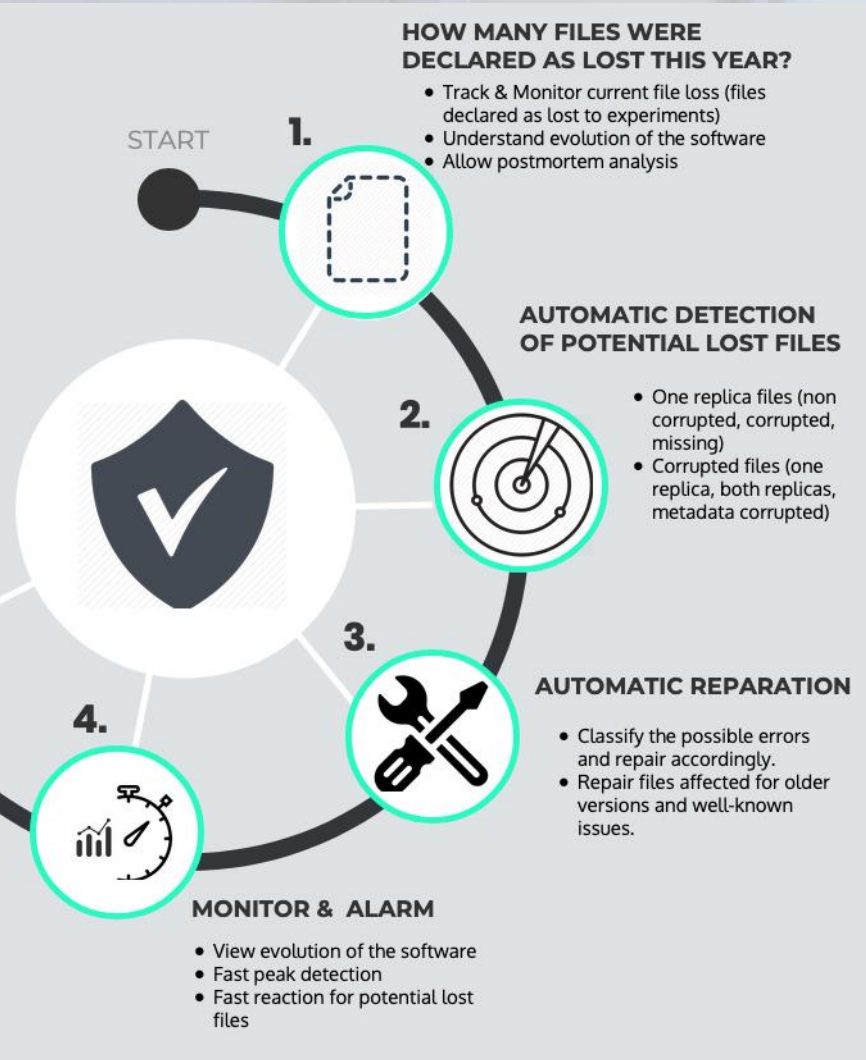  - Every day at 3pm: *"Detect + Classify + Repair + Drain" = Less human effort*



eos-ops-collect-drain-failed -i $instance --repair --send

# eos-ops-durability toolkit

- Code & automatic rpm creation:
  - ✓ Gitlab: https://gitlab.cern.ch/eos/eos-ops-durability
  - ✓ EOS repo: http://storage-ci.web.cern.ch/storage-ci/eos-ops-durability

- Installed via puppet in all MGMs (AliceDaq excluded)

- Running in MGMs / rundeck

- Output sent to:
  - Cernbox: https://cernbox.cern.ch/index.php/apps/files/?dir=/__myprojects/eos/Durability&
  - Data source – Elasticsearch + Data Discovery - Kibana: https://es-eosmon.cern.ch/kibana/app/kibana#/discover
  - Monitoring - Grafana: https://filer-carbon.cern.ch/grafana/d/JzDQWU7Zz/durability-classification

Information Technology Department

# Monitor & Alarm



**HOW MANY FILES WERE DECLARED AS LOST THIS YEAR?**

- Track & Monitor current file loss (files declared as lost to experiments)
- Understand evolution of the software
- Allow postmortem analysis

**AUTOMATIC DETECTION OF POTENTIAL LOST FILES**

- One replica files (non corrupted, corrupted, missing)
- Corrupted files (one replica, both replicas, metadata corrupted)

**AUTOMATIC REPARATION**

- Classify the possible errors and repair accordingly.
- Repair files affected for older versions and well-known issues.

**MONITOR & ALARM**

- View evolution of the software
- Fast peak detection
- Fast reaction for potential lost files

- Allow us to monitor the software evolution
- Fast peak detection and fast reaction



One Replica Files per Instance
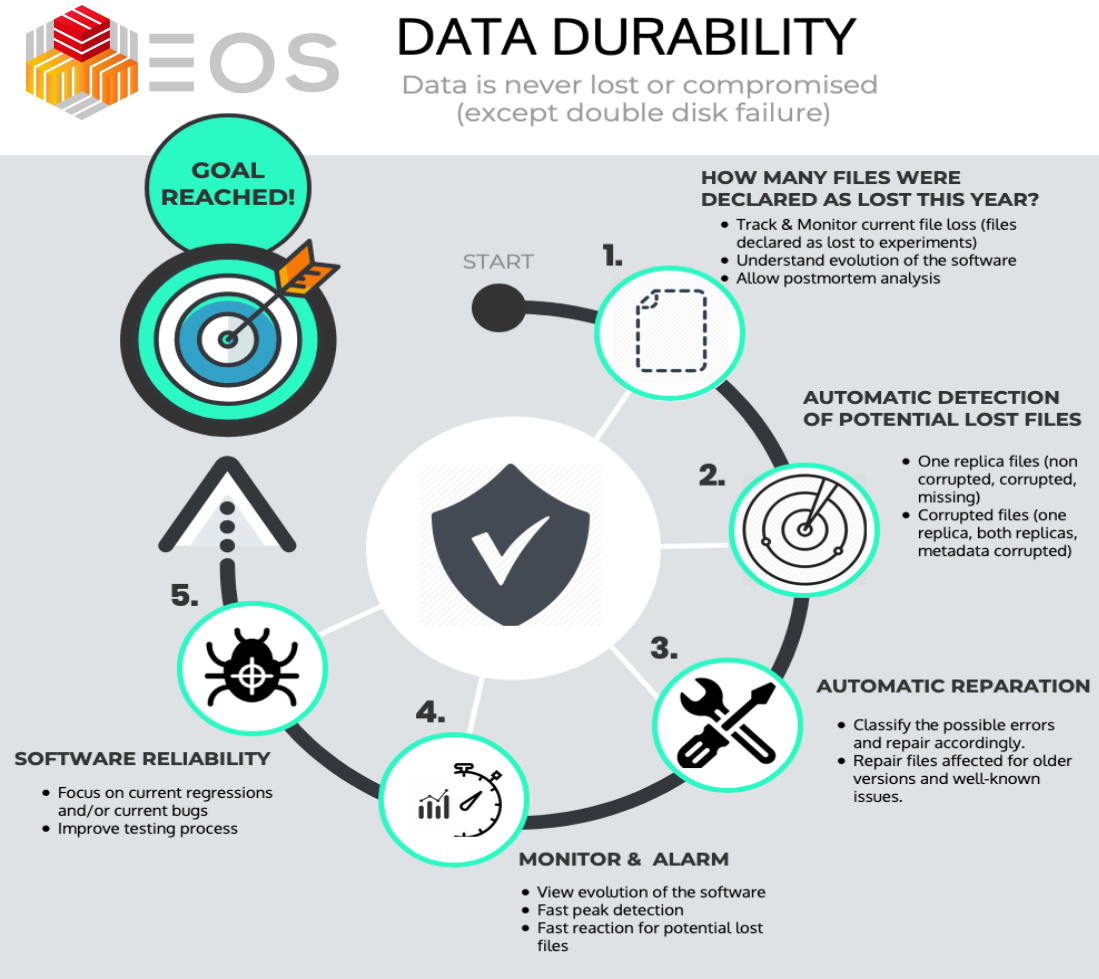
EOS Workshop

16

# Software reliability

- Focused on current regressions and/or bugs

- Improve testing process

  - Eg. Helping the testing of the new generation of fsck

# Future work
## Top 4 objectives for next round

- Automatic missing files retrieval from backup and restic in home instances

- Evaluate the new generation of fsck and complement its actions

- Provide external configuration for elasticsearch data sources and eos directories

- Include data durability checks for erasure coding – AliceDaq (complementing fsck)

# Conclusions



- Software evolution monitoring
- Better communication and better incidents understanding
- Less human effort in operations support (rota) and draining processes
- Faster peak detection and reaction
- Focus on current regressions/bugs, avoiding noise from the past
- Testing processes improvement