

### Ensuring Data Durability in EOS Systems

Maria Arsuaga-Rios IT-ST-PDS





Information Technology Department

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, ...).

#### EOS Workshop 2020 - Maria Arsuaga-Rios

1



### Flashback: March 2019



How many files are we missing?

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







T Information Technology Department

### 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*	0
Data Options	
Metrics Metric Count Add metrics	1. How many condition data base files are missing?
Buckets Select buckets type	(5 replicas expected)
Cancel	



Information Technology Department

EOS Workshop 2020 - Maria Arsuaga-Rios

Count



Information Technology Department 

#### 6

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

		cond/*")}* experiment "alice" Add a filter +	
eosmon_lost_files*	O 260		
Metrics Y-Axis Count Add metrics	3.	Is there a correlation with an incident,	
Buckets           X-Axis         mtime per month         I         X		ticket or known bug?	
Split Ser Obscending C I X Add sub-buckets		<ul> <li>10/11/2017 Incident: "Root cause: Crash caused the</li> </ul>	
		namespace to be corrupted" OTG0040840	
		<ul> <li>10/11/2017 to 14/12/2017 GGUS Ticket:"eosalice</li> </ul>	
		manual restore, missing condition files" GGUS ticket	



T Information Technology Department

### Automatic detection of high risk files

#### START • Track & Monitor current file loss (files declared as lost to experiments) • Understand evolution of the software • Allow postmortem analysis • Automatic Detection orrupted, corrupted, corrupted, missing) • One replica files (non corrupted, corrupted, nissing) • Corrupted files (one replica, both replicas, metadata corrupted)

HOW MANY FILES WERE

#### • Detection is the first step, make it visible!

• Detect one replica files (non corrupted, corrupted, missing, ...)

 Detect mismatching checksums and sizes (one replica, both replicas, metadata corrupted, ...)



IT Information Technology Department

### Automatic detection of high risk files

- One replica files (one replica layouts included):
  - o 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





T Information Technology Department

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



#### Filesystems affected by boot status





Files not drained - with drain failed by instance -



Filesystems affected by boot status







Information Technology Department

### Automatic reparation

#### 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.

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



EOS Workshop 2020 - Maria Arsuaga-Rios



2.

Information Technology Department

### Automatic reparation

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





First big reparation because of a massive adjustreplica performed without classification

### Automatic reparation

Automatic reparation for draining failures:

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





IT Information Technology Department

### eos-ops-durability toolkit

- Code & automatic rpm creation:
  - ✓ Gitlab: <u>https://gitlab.cern.ch/eos/eos-ops-durability</u>
  - ✓ 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: <u>https://cernbox.cern.ch/index.php/apps/files/?dir=/\_myprojects/eos/Durability&</u>
  - Data source Elasticsearch + Data Discovery Kibana: <u>https://es-eosmon.cern.ch/kibana/app/kibana#/discover</u>
  - Monitoring Grafana: <u>https://filer-</u> carbon.cern.ch/grafana/d/JzDQWU7Zz/durability-classification



### **Monitor & Alarm**

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



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







 Focused on current regressions and/or bugs
 Improve testing process

Software

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

