Access Pattern Analysis in the EOS Storage System at CERN

Olga Chuchuk CERN, IT department





9 June 2020

CERN LHC Detectors & EOS













Motivation

- Detailed data about EOS File Access is collected
- Useful for
 - Understanding differences between instances (experiments)
 - Understanding popularity and lifecycle of data
 - To evaluate usefulness of a cache
 - To evaluate the impact of EC
 - To find unexpected uses





Key Metrics:

fid	osize	csize	rb	wb	ots	cts	ruid	td	
293031349	0	495399529	0	495399529	2019-01-27 11:29:43	2019-01-27 11:32:45	7947	Ihcbprod.45887:92@hlta1013.lbdaq	
23514666	1744118422	1744118422	1021229	0	2019-01-27 12:06:32	2019-01-27 12:07:30	9801	lblocal.2988:190@lbbuild39	



Analysis Workflow

1 day of logs:





Metric	LHCb	СМЅ	ATLAS
Instance Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Total Accessed	49.2 PiB	67.5 PiB	155 PiB
Total Accessed (% of Total Volume)	328 %	188 %	310 %
Writes	20.0 PiB	28.0 PiB	26.1 PiB
Writes (% of Total Volume)	134 %	77.7 %	52.3 %
Reads	28.4 PiB	39.4 PiB	129 PiB
Reads (% of Total Volume)	190 %	110 %	257 %
Repeated Reads	8.7 PiB	30.7 PiB	108 PiB
Repeated Reads (% of Read Workload)	30.7 %	77.8 %	83.7 %



Metric	LHCb	CMS	ATLAS
Instance Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Total Accessed	49.2 PiB	67.5 PiB	155 PiB
Total Accessed (% of Total Volume)	328 %	188 %	310 %
Writes	20.0 PiB	28.0 PiB	26.1 PiB
Writes (% of Total Volume)	134 %	77.7 %	52.3 %
Reads	28.4 PiB	39.4 PiB	129 PiB
Reads (% of Total Volume)	190 %	110 %	257 %
Repeated Reads	8.7 PiB	30.7 PiB	108 PiB
Repeated Reads (% of Read Workload)	30.7 %	77.8 %	83.7 %

Different in absolute numbers.



Metric	LHCb	CMS	ATLAS
Instance Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Total Accessed	49.2 PiB	67.5 PiB	155 PiB
Total Accessed (% of Total Volume)	328 %	188 %	310 %
Writes	20.0 PiB	28.0 PiB	26.1 PiB
Writes (% of Total Volume)	134 %	77.7 %	52.3 %
Reads	28.4 PiB	39.4 PiB	129 PiB
Reads (% of Total Volume)	190 %	110 %	257 %
Repeated Reads	8.7 PiB	30.7 PiB	108 PiB
Repeated Reads (% of Read Workload)	30.7 %	77.8 %	83.7 %

Different in relative numbers:

• LHCb and ATLAS have more accesses than CMS.



Metric	LHCb	CMS	ATLAS
Instance Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Total Accessed	49.2 PiB	67.5 PiB	155 PiB
Total Accessed (% of Total Volume)	328 %	188 %	310 %
Writes	20.0 PiB	28.0 PiB	26.1 PiB
Writes (% of Total Volume)	134 %	77.7 %	52.3 %
Reads	28.4 PiB	39.4 PiB	129 PiB
Reads (% of Total Volume)	190 %	110 %	257 %
Repeated Reads	8.7 PiB	30.7 PiB	108 PiB

Writes:

- LHCb produces more data than other experiments
- LHCb produces more volume than its instance's size.



Metric	LHCb	CMS	ATLAS
Instance Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Total Accessed	49.2 PiB	67.5 PiB	155 PiB
Total Accessed (% of Total Volume)	328 %	188 %	310 %
Writes	20.0 PiB	28.0 PiB	26.1 PiB
Writes (% of Total Volume)	134 %	77.7 %	52.3 %
Reads	28.4 PiB	39.4 PiB	129 PiB
Reads (% of Total Volume)	190 %	110 %	257 %
Repeated Reads	8.7 PiB	30.7 PiB	108 PiB
Repeated Reads (% of Read Workload)	30.7 %	77.8 %	83.7 %

Reads:

• ATLAS reads more data than LHCb and CMS.



Metric	LHCb	CMS	ATLAS
Instance Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Total Accessed	49.2 PiB	67.5 PiB	155 PiB
Total Accessed (% of Total Volume)	328 %	188 %	310 %
Writes	20.0 PiB	28.0 PiB	26.1 PiB
Writes (% of Total Volume)	134 %	77.7 %	52.3 %
Reads	28.4 PiB	39.4 PiB	129 PiB
Reads (% of Total Volume)	190 %	110 %	257 %
Repeated Reads	8.7 PiB	30.7 PiB	108 PiB
Repeated Reads (% of Read Workload)	30.7 %	77.8 %	83.7 %

Repeated Reads:

- Potential for caching
- ~80% for CMS and ATLAS
- Only 30% for LHCb would profit less from caching.



Operations Classification

- (osize == 0 and csize > 0 and wb > 0 and rb == 0)
- (osize > 0 and csize == osize and wb == 0 and rb > 0)

(! create **and** ! read)



Metric	LHCb	CMS	ATLAS
Other Operations (% of related files)	0.06 %	0.26 %	0.61 %
Other Operations (% of Total Volume)	0.89 %	0.05 %	0.14 %

- Not much influence from abnormal operations.
- Assumption: Data is immutable.



Created/Read Volume Unique data files

Metric	LHCb	CMS	ATLAS
Total Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Created Volume	20.3 PiB	28.0 PiB	26.1 PiB
Created Volume (% of Total Volume)	136 %	77.7 %	52.2 %
Read Volume	22.7 PiB	22.0 PiB	26.6 PiB
Read Volume (% of Total Volume)	151 %	61.1 %	53.2 %
Repeated Read Volume (% of Read Volume)	20.7 %	54.7 %	55.1 %
Average Fraction of File Read	88.8 %	55.1 %	81.6 %

LHCb:

- Produce and read a lot of data
- Well-organized workflow
- Create \rightarrow Read \rightarrow Delete
- Don't have enough space.



Created/Read Volume Unique data files

Metric	LHCb	CMS	ATLAS
Total Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Created Volume	20.3 PiB	28.0 PiB	26.1 PiB
Created Volume (% of Total Volume)	136 %	77.7 %	52.2 %
Read Volume	22.7 PiB	22.0 PiB	26.6 PiB
Read Volume (% of Total Volume)	151 %	61.1 %	53.2 %
Repeated Read Volume (% of Read Volume)	20.7 %	54.7 %	55.1 %
Average Fraction of File Read	88.8 %	55.1 %	81.6 %

CMS & ATLAS:

- Produce less data
- Use only part of the space
- Higher chance of reuses.



Created/Read Volume Unique data files

Metric	LHCb	CMS	ATLAS
Total Volume (EOS Control Tower)	15.0 PiB	36.0 PiB	50.0 PiB
Created Volume	20.3 PiB	28.0 PiB	26.1 PiB
Created Volume (% of Total Volume)	136 %	77.7 %	52.2 %
Read Volume	22.7 PiB	22.0 PiB	26.6 PiB
Read Volume (% of Total Volume)	151 %	61.1 %	53.2 %
Repeated Read Volume (% of Read Volume)	20.7 %	54.7 %	55.1 %
Average Fraction of File Read	88.8 %	55.1 %	81.6 %

Fraction of Files Read:

- 80-90% for ATLAS and LHCb
- Only 55% for CMS
- Bad for caching.



Normalized probability distribution of file sizes:





New/Old Files Relevance

Metric	LHCb	CMS	ATLAS
Created and Read	19.2 PiB	16.2 PiB	21.0 PiB
Created and Not Read	1.1 PiB	11.7 PiB	5.1 PiB
Created and Read (% of Created Volume)	94.4 %	58.0 %	80.4 %
Created and Not Read (% of Created Volume)	5.62 %	42.0 %	19.6 %
Old and Read	3.5 PiB	5.8 PiB	5.6 PiB
Old and Not Read	11.5 PiB	31.2 PiB	40.4 PiB
Old Before and Read (% of Old Volume)	23.5 %	16.0 %	11.2 %
Old Before and Not Read (% of Old Volume)	76.6 %	84.0 %	88.8 %

Usage of new data:

- LHCb uses almost all the data they produce
- For ATLAS and CMS the numbers are lower.



New/Old Files Relevance

Metric	LHCb	CMS	ATLAS	
Created and Read	19.2 PiB	16.2 PiB	21.0 PiB	
Created and Not Read	1.1 PiB	11.7 PiB	5.1 PiB	
Created and Read (% of Created Volume)	94.4 %	58.0 %	80.4 %	
Created and Not Read (% of Created Volume)	5.62 %	5.62 % 42.0 %		
Old and Read	3.5 PiB	5.8 PiB	5.6 PiB	
Old and Not Read	11.5 PiB	31.2 PiB	40.4 PiB	
Old Before and Read (% of Old Volume)	23.5 %	16.0 %	11.2 %	
Old Before and Not Read (% of Old Volume)	76.6 %	84.0 %	88.8 %	

Usage of old data:

• Rarely used in future.



Time distribution of file accesses:







99% of volume covered by converting <3% of files - 2M files



20

Input to the future plans





Would we benefit from having several QoS levels?

How do we decide which level to use for each piece of data?



Popularity Prediction



Thank you Any questions?

Olga Chuchuk

olga.chuchuk@cern.ch

CERN, IT-SC-RD, Geneva, Switzerland Taras Shevchenko KNU, Kyiv, Ukraine

Dirk Duellmann

dirk.duellmann@cern.ch CERN, IT-SC, Geneva, Switzerland





9 June 2020

Backup slides



Analysis Environment

D	File Edit Code	View Plots Session	n Build Debug Profile	Tools Help					ocloca	al 🕞 🕘		
1 • 🕲 • 🕲 🖆 • 🔒 🖨 • 🕞 📥 🖉 • Go to file/function												
🗣 Presentation.Rmd x 🗣 Operations.Rmd x 🗣 Files.Rmd x 🗣 FilesStats.Rmd x												
	$(\Box) = ABC $					n 🗸 💁 🖌 🖹	🚭 🕞 🖙 Import Dataset 🗸 🔮 List 🗸 🎯					
64	64 # convert bytes to human readable					~~ _	🔒 Global Environment 🗸 🛛 🔍					
65	65 all_volume_stats[, total_volume := hsize(total_volume)]						Data					
66 67	66 all_volume_stats[, total_turnover := hsize(total_turnover)] 67 all_volume_stats[write_workload := hsize(write_workload)]						💽 all_files… 400000 obs. of 2 varia… 🔲					
68	<pre>68 all_volume_stats[, read_workload := hsize(read_workload)]</pre>						Values					
69 70	<pre>69 all_volume_stats[, repeated_read_workload := hsize(repeated_read_workload)] 70 </pre>						end_date 2019-06-30					
71							experiment	t "alice"				
72	<pre>'``{r} all volume state</pre>					\$ 🔺 🕨	experimen.	cnr 11:41	"Incb" "cm	s		
74							Files Plots	Packages	Help Viewe	r 🗖		
	experiment <chr></chr>	total_volume	total_turnover		read_workload	<i>≈</i> × × ►	ONew Folder	0 Upload	🕴 Delete 📮	Rename		
				write_workload <chr></chr>			ome > olga >	logs_analysis	_remote > note	ebooks		
	lhcb	15.0 PiB	28.5 PiB	11.6 PiB	16.3 PiB		▲ Nar	ne		Size		
	cms	37.0 PiB	35.4 PiB	14.7 PiB	20.6 PiB		Т					
	alice	43.0 PiB	14.9 PiB	4.7 PiB	10.2 PiB		Chang					
	atlas	46.0 PiB	69.8 PiB	10.6 PiB	58.9 PiB		🗌 🧰 Explo	ring				
	4 rows 1-5 of 10 columns						🗌 🔍 🍯 Files.F	Rmd		3 KB		
							🗌 🔍 FilesS	tats.Rmd		4.8 KB		
75	<pre>75 * ```{r} 76 for_save <- dcast(melt(all_volume_stats, id.vars = "experiment"), variable ~ experiment) 77 for save</pre>						🗌 🔍 Opera	tions.Rmd		5.2 KB		
76							🗌 📄 OperS	stats.Rmd		4.5 KB		
78	···-					🗌 🧰 Parsin	Iq					
1:1	🗰 R Notebook 💲					R Markdown 💲	Result	ts				
Console	2					a n		al for ALICE				



CMS Files Volume Cumulative Distribution:

CMS Files Size Cumulative Distribution:



87% of volume covered by converting <3% of files - 17M files



1.00 -1.00 -Files Volume Cumulative Density Files Count Cumulative Density 0.20 -0.22 -0.75 -0.50 -0.25 -0.00 -0.00 1 kB 1 MB 1 ĠB 1 kB 1 TB 1 MB 1 ĠB 1 TB Files Size Files Size Colour Current Namespace (all files including old, but excluding deleted) ____ EOS Logs (only files that were created or accessed) _

ATLAS Files Volume Cumulative Distribution:

ATLAS Files Size Cumulative Distribution:

70% of volume covered by converting <3% of files - 6M files





65% of volume covered by converting <3% of files - 5M files



ALICE Files Volume Cumulative Distribution:

ALICE Files Size Cumulative Distribution: