

Comparison between **CDFS** (Comtrade Distributed FS), **CephFS**, **HDFS** (Hadoop Distributed FS), **GPFS** (IBM Spectrum Scale)

Gregor Molan

Branko Blagojević

Ivan Arizanović

Comtrade Group / Comtrade 360

Content

- Introduction
- (Clustered) Distributed file system
- High-performance file systems:
 Selection
- High-performance file systems: Features
- High-performance file systems: **Results of comparison**
- Conclusion



Introduction

- Introduction
 - (Clustered) Distributed file system
- High-performance file systems: Selection
- High-performance file systems: Features
- High-performance file systems: **Results of comparison**
- Conclusion

(Clustered) Distributed file system





High-performance file systems: **Selection**

• Introduction

High-performance file systems: Selection

- CephFS
- HDFS (Hadoop Distributed File System)
- GPFS (IBM Spectrum Scale)
- CDFS (Comtrade Distributed FS)
- High-performance file systems: Features
- High-performance file systems: Results of comparison
- Conclusion



CephFS

- Massively scalable
- Part of the Linux kernel since 2010
- RADOS: Ceph's foundation

Interfaces:

- S3-compatible
- Swift-compatible



HDFS (Hadoop Distributed File System)

Basic features

- Distributed
- Scalable
- Portable
- Written in Java

HDFS services

- Name Node (master)
- Data Node
- Secondary Name Node
- Job tracker
- Task Tracker



GPFS (IBM Spectrum Scale)

- #1 HPC in 2019 (Oak Ridge National Laboratory)
- OS support for server:
 - AIX
 - Linux
 - Windows
- Configuration update on a mounted file system
- Data replication
- Active File Management (AFM)
 - Data share across clusters



CDFS (Comtrade Distributed FS)

CDFS = appliance of CERN EOS at Comtrade

- Comtrade: The industry partner of CERN EOS
- The initial usage of CERN EOS
 - Data collection for CERN LHC experiments
- Current usage of CERN EOS
 - Data collection for all CERN experiments
 - The primary data storage software (including CERN staff data)
- The fastest file system for parallel data collection
- Cluster size: > 700 PB
 - Node resync: < 15 minutes



High-performance file systems: Features

- Introduction
- High-performance file systems: Selection
- High-performance file systems: Features
 - Requirements
 - Fault tolerance
 - Advantages of RAIN
- High-performance file systems: Results of comparison
- Conclusion



Requirements

- High throughput
- Low latency
- Expendability



Fault tolerance

• CephFS

- Using JBOD instead of RAID
- Snapshots
- Replication
- File and directory layouts

• HDFS

- Using JBOD instead of RAID
- Stores each file as a sequence of blocks which are replicated for fault tolerance
- The block size and replication factor are configurable per file

• GPFS

- Using IBM Spectrum Scale RAID
- Snapshots
- Synchronous and asynchronous replication
- CDFS (based on EOS)
 - Uses JBOD in the form of RAIN



About the RAIN – Checksums

RAIN = Software implementation of the RAID concept across independent servers on the network





About the RAIN – File Layouts

RAIN = Checksums calculated and recorded for every file (chunk)





Advantages of RAIN

Advantages

- Scalability
- Reliability
- Cost (JBOD without RAID controller)
- Geotag policies are applied during file placement to improve data loss prevention and IO performance.

Drawbacks

- All communication is done via the network
- Increased is IO and computational effort for nonsequential writes and server draining



High-performance file systems: **Results of comparison**

- Introduction
- High-performance file systems: Selection
- High-performance file systems: Features
- High-performance file systems: Results of comparison
 - Throughput measurements
 - Different file sizes
- Conclusion



Testing environment

- DFS as a single node*
- Clusters on different HDDs
- Identical disk drives
- Identical HW components:
 - Motherboards
 - Network adapters
 - Memory

	Cores	Memory	HDD
Client: Linux	2	4 GB	80 GB
Client: Windows	4	10 GB	80 GB
Servers: Linux	8	20 GB	500 GB



Different file sizes

Small files

- Size: 1 MB
- Transfer: 100 files at once
- Potential issue:
- Authentication time overhead

Medium files

- Size: 100 MB
- Transfer: 10 files at once
- Potential issue: -

Larger files

- Size: 2 GB
- Transfer: 2 files at once
- Potential issue:
 - Time out



Testing description

Download

- Create new test files on the server space
- Clear file cache on the client and the server machines
- Download the files from the server space to the client machine
- Verifying the MD5 hash and calculating the transfer speed from execution time
- · Remove created and copied test files

Upload

- Create new test files on the client machine
- Clear file cache on the client and the server machines
- Upload the files from the client machine to the server space
- Verifying the MD5 hash and calculating the transfer speed from execution time
- Remove created and copied test files



Throughput results

ľ	Iterations (EOS) Iterations (IBM)		21 21	(checksi (checksi	ums OK) ums OK)				
	ter	atio	ons (Ceph)	23	(checks)	(checksums OK) Number of fil		files	100
I	ter	atio	ons (Hadoop)	14	(checks	ums OK)	File size [M	B]	1
Γ	Test [MB/s]		min	max	avg	trim25%	Avgtin	ne [ms]	
			EOS: xrdcp command	142,86	181,82	165,24	165,87	23	6,05
		×	EOS Fusex	45,81	52,03	49,27	49,32		20,30
		inu	IBM Spectrum Scale	54,08	58,82	55,99	55,92	25	17,86
	_	-	Ceph on Linux	145,99	156,25	150,62	150,50	2	6,64
	oad		Hadoop on Linux	3,46	3,64	3,55	3,55		281,92
	9		EOS-wnc	14,25	15,18	14,75	14,76	1	67,78
	swopu	Mo	EOS-drive ST	9,70	10,00	9,88	9,89		101,22
		bd	EOS: Samba	22,68	24,13	23,29	23,28	23	42,94
		Ň	Ceph on Win	50,28	56,50	53,52	53,51	\$	18,68
			Hadoop on Win	3,13	3,21	3,18	3,18		314,61
			EOS: xrdcp command	95,24	196,08	169,56	174,52		5,90
		×	EOS Fusex	48,45	53,05	50,67	50,63		19,74
		inu	IBM Spectrum Scale	158,73	187,27	174,76	175,16	2	5,72
-	8	_	Ceph on Linux	7,47	110,13	94,19	97,40	23	10,62
-	ë		Hadoop on Linux	3,68	4,30	3,97	3,97		251,71
	No.		EOS-wnc	10,66	11,17	10,88	10,88		91,90
4	ŏ	W	EOS-drive ST	17,95	19,03	18,44	18,43	2	54,24
		bd	EOS: Samba	13,19	15,82	14,28	14,24	23	70,02
		Ň	Ceph on Win	1,93	47,25	35,38	37,60	23	28,26
			Hadoop on Win	1,69	2,60	2,20	2,21		454,84

Legend:	[MB/s]						
Red	0		10				
Orange	10	-	20				
Yellow	20	-	30				
Light green	30	-	40				
Green	40	-	~				

lter Iter	ratio ratio	ons (EOS) ons (IBM)	28 28	(checks) (checks)	umsOK) umsOK)			
Iter	ratio	ons (Ceph)	52	(checks	ums OK)	Number of	files	10
Iter	ratio	ons (Hadoop)	11	(checks	ums OK)	File size [M	B]	100
	_	T + (h m ()	.	-				
	_	lest [MB/s]	min	max	avg	trim25%	Avgt	ime [ms]
		EOS: xrdcp command	359,71	444,44	411,14	412,67	23	243,22
	×	EOS Fusex	134,72	192,64	160,16	160,07	23	624,38
	Ĕ.	IBM Spectrum Scale	176,62	188,71	181,08	181,01	23	552,25
-		Ceph on Linux	131,89	162,39	140,86	139,97		709,95
ga		Hadoop on Linux	9,43	10,17	9,94	9,97		10064,96
<u>a</u>	s	EOS-wnc	174,21	204,60	186,28	185,66	1	536,82
_	WS	EOS-drive ST	186,54	210,24	197,67	197,40	\$	505,89
	P P	EOS: Samba	165,56	231,33	196,82	196,65	2	508,08
	Vir	Ceph on Win	102,68	141,96	136,28	137,07		733,77
	ſ-	Hadoop on Win	4,50	5,22	4,61	4,56		21670,61
		EOS: xrdcp command	301,20	436,68	412,94	417,50	4	242,16
	×	EOS Fusex	186,67	217,11	206,36	207,14	23	484,59
	nu	IBM Spectrum Scale	306,75	345,18	322,89	322,24	23	309,70
Ð		Ceph on Linux	20,44	183,49	31,49	28,27		3175,40
loa		Hadoop on Linux	8,06	10,51	9,37	9,39		10668,22
N.		EOS-wnc	128,10	177,31	151,39	151,01	\$	660,54
Å	WS	EOS-drive ST	148,70	185,92	157,76	156,40	2	633,87
	^b	EOS: Samba	72,97	97,50	81,26	80,39	1	1230,60
	Vir	Ceph on Win	17,63	81,00	25,54	23,66		3915,54
	1	Hadoop on Win	4.31	4.62	4.50	4.51		22217.73

Legend:		[MB/s]
Red	0	-	100
Orange	100	-	150
Yellow	150	-	200
Light green	200	-	250
Green	250		00

lter Iter	ratio ratio	ons (EOS) ons (IBM)	27 28	{checksu {checksu	ıms OK) ıms OK)			
Iter	ratio	ons (Ceph)	52	(checks)	ims OK)	Number of files		2
Iter	ratio	ons (Hadoop)	11	(checksu	ims OK)	File size [M	B]	2000
		,		-				
	Test [MB/s]		min	max	avg	trim25%	Avg ti	me [s]
		EOS: xrdcp command	329,49	405,27	371,03	371,17	2	5,39
	×	EOS Fusex	187,92	237,63	210,76	210,51	25	9,49
	2	IBM Spectrum Scale	283,61	318,22	294,47	293,28	23	6,79
_		Ceph on Linux	141,00	163,37	157,56	158,17		12,69
oad		Hadoop on Linux	9,74	10,10	9,91	9,91		201,83
ā		EOS-wnc	158,40	331,09	231,25	227,75	*	8,65
_	swopu	EOS-drive ST	212,22	294,47	237,44	234,72	2	8,42
		EOS: Samba	164,11	229,82	181,25	178,59	25	11,03
	Š.	Ceph on Win	128.18	158.19	153.32	154.04		13.04
	-	Hadoop on Win	4,61	4,72	4,66	4,66		428,85
		EOS: xrdcp command	328,68	365,97	353,00	354,47	\$	5,67
	×	EOS Fusex	218,66	233,36	227,13	227,15		8,81
	2	IBM Spectrum Scale	328,95	364,96	342,54	341,65	2	5,84
5	⊡	Ceph on Linux	188,80	355,49	265.08	264.04	27	7.54
loa		Hadoop on Linux	9,28	10,63	10,12	10,15		197,66
M.		EOS-wnc	119,92	213,86	170,17	169,49	\$	11,75
õ	WS	EOS-drive ST	179,86	210,49	190,24	189,72	2	10,51
	Pg.	EOS: Samba	17,95	35,43	25,85	25,54		77,37
	Vir	Ceph on Win	105,38	141,66	122,82	122,90	\$3	16,28
	-	Hadoop on Win	4.30	4.73	4.55	4.55		440.00

Legend:	[N	[MB/s]				
Red	0 - 100					
Orange	100	-	150			
Yellow	150	-	200			
Light green	200	-	250			
Green	250	-	~			

^ST - Single-thread ^MT - Multi-thread



Throughput results - Small Files

lte	ratio	ons (EOS)	21	(checksu	ums OK)			
lte	ratio	ons (IBM)	21	(checksเ	ıms OK)			
lte	ratio	ons (Ceph)	23	(checksı	ıms OK)	Number of	files	100
lte	ratio	ons (Hadoop)	14	(checksเ	ıms OK)	File size [MB	3]	1
	Tect [MR/c]		min	may	avg	trim25%	Avgtin	ne [ms]
		EOS: xrdcp command	142,86	181,82	165,24	165,87	<u> </u>	6,05
			15,81	57,02	/10/07/	10 27		20,30
	2	IBM Spectrum Scale	54,08	58,82	55,99	55,92	12	17,86
-		Ceph on Linux	145,99	156,25	150,62	150,50	1	6,64
oac		Hadoop on Linux	3,46	3,64	3,55	3,55		281,92
ld		EOS-wnc	14,25	15,18	14,75	14,76	13	67,78
	NS NO	EOS-drive ST	9,70	10,00	9,88	9,89		101,22
	opu	EOS: Samba	22,68	24,13	23,29	23,28	25	42,94
	Vi1	Ceph on Win	50,28	56,50	53,52	53,51	1	18,68
		Hadoon on Win	3 13	3 21	3 18	3 18		314,61
		EOS: xrdcp command	95,24	196,08	169,56	174,52	%	5,90
		I FUS FUSEX	48.45	55.05	າ ບ.ບ	50.05		19,74
	<u>.</u>	IBM Spectrum Scale	158,73	187,27	174,76	175,16	1	5,72
pe		Cepir on Linux	/ ,4/	110,15	94,19	57,40	13	10,62
ol		Hadoop on Linux	3,68	4,30	3,97	3,97		251,71
W L		EOS-wnc	10,66	11,17	10,88	10,88		91,90
D	WS	EOS-drive ST	17,95	19,03	18,44	18,43	1	54,24
	opu	EOS: Samba	13,19	15,82	14,28	14,24	25	70,02
	-i-	Ceph on Win	1 93	17 25	35 38	37.60	<u> </u>	20.26

	Iter	ratio	ons (EOS)	28	(checks	ums
	Iter	atio	ons (IBM)	28	(checks	ums
	Iter	atio	ons <mark>(</mark> Ceph)	52	(checks	ums
	lter	atio	ons (Hadoop)	11	(checks	ums
			Test [MB/s]	min	max	
			EOS: xrdcp command	359,71	444,44	4
		Linux	EOS Fusex	134,72	192,64	1
			IBM Spectrum Scale	176,62	188,71	1
	T		Ceph on Linux	131,89	162,39	1
	oac		Hadoop on Linux	9,43	10,17	
	١pl		EOS-wnc	174,21	204,60	1
	_	S No	EOS-drive ST	186,54	210,24	1
		pdq	EOS: Samba	165,56	231,33	1
		Wi	Ceph on Win	102,68	141,96	1
			Hadoop on Win	4,50	5,22	
			EOS: xrdcp command	301,20	436,68	4
		×	EOS Fusex	186,67	217,11	2
		inu	IBM Spectrum Scale	306,75	345,18	3
	ad		Ceph on Linux	20,44	183,49	
	nlo		Hadoop on Linux	8,06	10,51	
	IMO	10	EOS-wnc	128,10	177,31	1
	Õ	Mo	EOS-drive ST	148,70	185,92	1
		nde	EOS: Samba	72,97	97,50	
			Cenh on Win	17.63	81.00	

Throughput results - Medium Files

5	100
	1
gtin	ne [ms]
3	6,05
	20,30
3	17,86
5	6,64
	281,92
S.	67,78
	101,22
3	42,94
3	18,68
	314,61
a l	5,90
	19,74
3	5,72
5	10,62
	251,71
	91,90
a la	54,24
3	70,02
<u>.</u>	28.26

lte Ite	Iterations (EOS) Iterations (IBM)		28 28	(checksı (checksı	ums OK) ums OK)			
lte Ite	Iterations (Ceph) Iterations (Hadoop)		52 11	(checksı (checksı	(checksums OK) (checksums OK)		Number of files File size [MB]	
		Tect [MB/c]	min	may	avg	trim25%	Avgt	ime [ms]
	Г	EOS: xrdcp command	359,71	444,44	411,14	412,67		243,22
	×		12/1 77	107,61	160,16	160,07	1	624,38
	inu	IBM Spectrum Scale	176,62	188,71	181,08	181,01	1	552,25
		Ceph on Linux	131,89	162,39	140,86	139,97		709,95
oac		Hadoop on Linux	9,43	10,17	9,94	9,97		10064,96
lqL		EOS-wnc	174,21	204,60	186,28	185,66	1	536,82
	No.	EOS-drive ST	186,54	210,24	197,67	197,40	2	505,89
	bpu	EOS: Samba	165,56	231,33	196,82	196,65	2	508,08
	N.	Ceph on Win	102,68	141,96	136,28	137,07		733,77
		Hadoon on Win	4 50	5 22	4 61	4 56		21670,61
	Π	EOS: xrdcp command	301,20	436,68	412,94	417,50	*	242,16
	×			/ / / / /	70b 3b	707-14	1	484,59
	inu	IBM Spectrum Scale	306,75	345,18	322,89	322,24	12	309,70
ad		Серноп спих	20,44	103,49	51,49	20,27		3175,40
iol	Downlos Downlos	Hadoop on Linux	8,06	10,51	9,37	9,39		10668,22
N L		EOS-wnc	128,10	177,31	151,39	151,01	2	660,54
ă		EOS-drive ST	148,70	185,92	157,76	156,40	2	633,87
		EOS: Samba	72,97	97,50	81,26	80,39	25	1230,60
	-i-	Ceph on Win	17.63	91 00	25 54	23.66		2015 5/

Ite	ratio	ons (EOS)	27	(checksı	ım
Ite	ratio	ons (IBM)	28	(checksu	ım
lte	ratio	ons (Ceph)	52	(checksu	ım
Ite	ratio	ons (Hadoop)	11	, (checksu	ım
				-	
		Test [MB/s]	min	max	
		EOS: xrdcp command	329,49	405,27	3
	Linux	EOS Fusex	187,92	237,63	2
		IBM Spectrum Scale	283,61	318,22	2
-		Ceph on Linux	141,00	163 <i>,</i> 37	1
oac		Hadoop on Linux	9,74	10,10	
lq		EOS-wnc	158,40	331,09	2
	N S	EOS-drive ST	212,22	294,47	2
	Windo	EOS: Samba	164,11	229,82	1
		Ceph on Win	128,18	158,19	1
		Hadoop on Win	4,61	4,72	
		EOS: xrdcp command	328,68	365,97	3
	×	EOS Fusex	218,66	233,36	2
	inu	IBM Spectrum Scale	328,95	364,96	3
ad		Ceph on Linux	188,80	355,49	2
		Hadoop on Linux	9,28	10,63	
Ň		EOS-wnc	119,92	213,86	1
ă	No.	EOS-drive ST	179,86	210,49	1
	opc	EOS: Samba	17,95	35,43	
	15	Cenh on Win	105 38	141 66	1

Throughput results - Large Files

S	10	
	100	
Avg ti		
<u>i</u> z	243,22	
(r	624,38	
(s	552,25	
	709,95	
	10064,96	
(r	536,82	
<u>i</u>	505,89	
Cr.	508,08	
	733,77	
	21670,61	
2	242,16	
63	484,59	
3	309,70	
	3175,40	
	10668,22	
3	660,54	
<u>a</u>	633,87	
3	1230,60	
	3915 54	

						-		
Iterations (EOS)		27	(checksums OK)					
Iterations (IBM)		28	(checksums OK)					
Iterations (Ceph)		52	(checksums OK)		Number of files		2	
Iterations (Hadoop)		11	(checksums OK)		File size [MB]		2000	
	Tect [MR/c]		min	may avg		trim 25% Avg ti		me [s]
		EOS: xrdcn command	329.49	405.27	371.03	371 17		5 39
			12/07	72762	210.76	210 51	್ಷ	9,00
	Ž	IBM Spectrum Scale	283.61	318.22	294.47	293.28	845 	6.79
			141.00	163.37	15/56	158.1/		12 69
bad		Hadoop on Linux	9.74	10.10	9.91	9.91		201.83
plq		EOS-wnc	158,40	331,09	231,25	227,75	1	8,65
	ŴS	EOS-drive ST	212,22	294,47	237,44	234,72	1	8,42
	opc	EOS: Samba	164,11	229,82	181,25	178,59	15	11,03
	Vir	Ceph on Win	128,18	158,19	153,32	154,04		13,04
		Hadoon on Win	4.61	4 72	4 66	4 66		428,85
		EOS: xrdcp command	328,68	365,97	353,00	354,47	1	5,67
Download Linux	×	I FUS FUSEX	1218.00	255.50	///.15			8,81
	ī.	IBM Spectrum Scale	328,95	364,96	342,54	341,65	1	5,84
		Ceph on Linux	100,00	<u>,45</u>	205,00	204,04	1	7,54
		Hadoop on Linux	9,28	10,63	10,12	10,15		197,66
		EOS-wnc	119,92	213,86	170,17	169,49	1	11,75
	N S	EOS-drive ST	179,86	210,49	190,24	189,72	1	10,51
	ndc	EOS: Samba	17,95	35,43	25,85	25,54		77,37
	15	Cenh on Win	105 38	141 66	177 87	122 90	1	16.28

Conclusion

- Introduction
- High-performance file systems: Selection
- High-performance file systems: Features
- High-performance file systems: Results of comparison
- Conclusion
 - Interpretation of results
 - Metrics used
 - Need to compare in future



Interpretation of results

The best

- Small files
 - GPFS on Linux
 - EOS on Linux
- Medium files
 - EOS on Linux
 - GPFS on Linux
- Large files
 - EOS on Linux
 - GPFS on Linux

Not the best

- All file sizes
 - Hadoop on Win
 - Hadoop on Linux
 - Samba



Metrics used (data and metadata)





Need to compare in future

High Availability metrics

- MTBF
- TBW
- Failover resync time
- Resync of replaced disk

High Availability requirements

- Load balancing
- Data scalability
- Geographical diversity
- Backup to tape





Thank you

Comparison between **CDFS** (Comtrade Distributed FS), **CephFS**, **HDFS** (Hadoop Distributed FS), **GPFS** (IBM Spectrum Scale)

Gregor Molan gregor@comtrade.com Branko Blagojević

Ivan Arizanović

Comtrade Group / Comtrade 360