

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

Gregor Molan

Branko Blagojević

Ivan Arizanović

Content

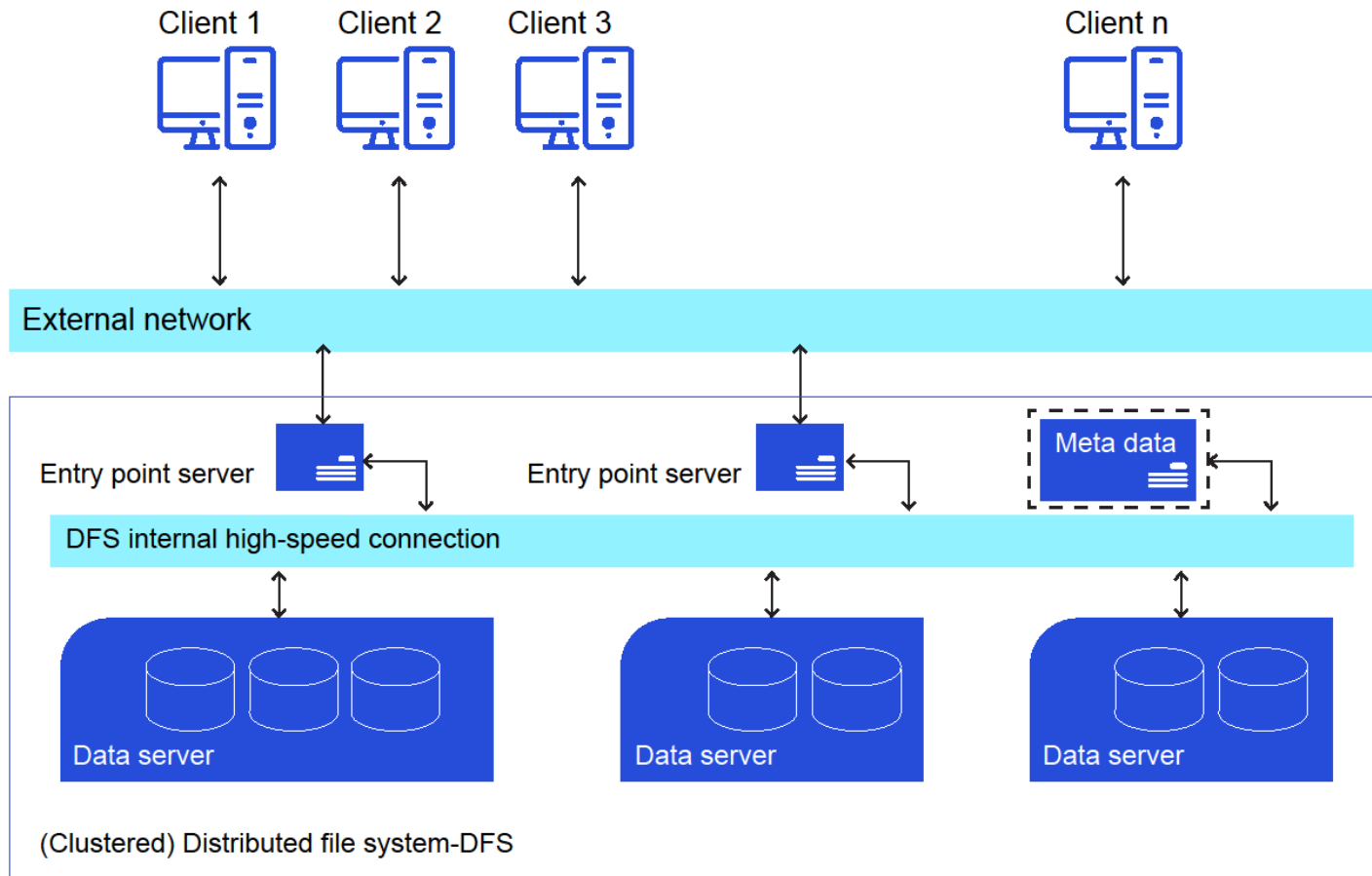
- Introduction
- (Clustered) Distributed file system
- High-performance file systems: **Selection**
- High-performance file systems: **Features**
- High-performance file systems:
Results of comparison
- Conclusion



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
- High availability
- High reliability

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

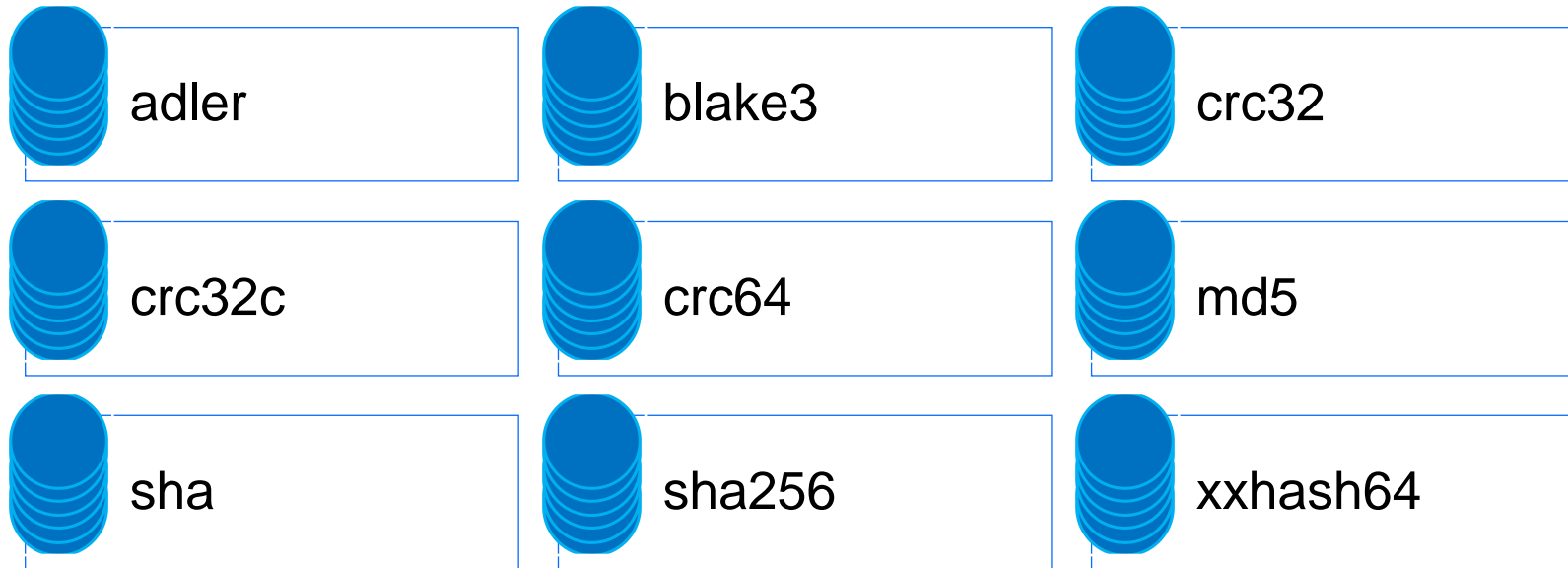
About the RAIN – File Layouts

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



About the RAIN – Checksums

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 non-sequential 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

Iterations (EOS)		21	(checksums OK)				
Iterations (IBM)		21	(checksums OK)				
Iterations (Ceph)		23	(checksums OK)		Number of files	100	
Iterations (Hadoop)		14	(checksums OK)		File size [MB]	1	
Test [MB/s]							
Upload	Linux	EOS: xrdcp command	142,86	181,82	165,24	165,87	★ 6,05
		EOS Fusex	45,81	52,03	49,27	49,32	★ 20,30
		IBM Spectrum Scale	54,08	58,82	55,99	55,92	★ 17,86
		Ceph on Linux	145,99	156,25	150,62	150,50	★ 6,64
		Hadoop on Linux	3,46	3,64	3,55	3,55	★ 281,92
		EOS-wnc	14,25	15,18	14,75	14,76	★ 67,78
Download	Linux	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
		IBM Spectrum Scale	158,73	187,27	174,76	175,16	★ 5,72
		Ceph on Linux	7,47	110,13	94,19	97,40	★ 10,62
		Hadoop on Linux	3,68	4,30	3,97	3,97	★ 251,71
		EOS-wnc	10,66	11,17	10,88	10,88	★ 91,90
Upload	Windows	EOS-drive ST	9,70	10,00	9,88	9,89	★ 101,22
		EOS: Samba	22,68	24,13	23,29	23,28	★ 42,94
		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-wnc	14,25	15,18	14,75	14,76	★ 67,78
		EOS-drive ST	9,70	10,00	9,88	9,89	★ 101,22
Download	Windows	EOS: Samba	22,68	24,13	23,29	23,28	★ 42,94
		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-wnc	10,66	11,17	10,88	10,88	★ 91,90
		EOS-drive ST	17,95	19,03	18,44	18,43	★ 54,24
		EOS: Samba	13,19	15,82	14,28	14,24	★ 70,02

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

Iterations (EOS)		28	(checksums OK)				
Iterations (IBM)		28	(checksums OK)				
Iterations (Ceph)		52	(checksums OK)		Number of files	10	
Iterations (Hadoop)		11	(checksums OK)		File size [MB]	100	
Test [MB/s]							
Upload	Linux	EOS: xrdcp command	359,71	444,44	411,14	412,67	★ 243,22
		EOS Fusex	134,72	192,64	160,16	160,07	★ 624,38
		IBM Spectrum Scale	176,62	188,71	181,08	181,01	★ 552,25
		Ceph on Linux	131,89	162,39	140,86	139,97	★ 709,95
		Hadoop on Linux	9,43	10,17	9,94	9,97	★ 10064,96
		EOS-wnc	174,21	204,60	186,28	185,66	★ 536,82
Download	Linux	EOS: xrdcp command	301,20	436,68	412,94	417,50	★ 242,16
		EOS Fusex	186,67	217,11	206,36	207,14	★ 484,59
		IBM Spectrum Scale	306,75	345,18	322,89	322,24	★ 309,70
		Ceph on Linux	20,44	183,49	31,49	28,27	★ 3175,40
		Hadoop on Linux	8,06	10,51	9,37	9,39	★ 10668,22
		EOS-wnc	128,10	177,31	151,39	151,01	★ 660,54
Upload	Windows	EOS-drive ST	148,70	185,92	157,76	156,40	★ 633,87
		EOS: Samba	72,97	97,50	81,26	80,39	★ 1230,60
		Ceph on Win	17,63	81,00	25,54	23,66	★ 3915,54
		Hadoop on Win	4,31	4,62	4,50	4,51	★ 22217,73
		EOS-wnc	128,10	177,31	151,39	151,01	★ 660,54
		EOS-drive ST	148,70	185,92	157,76	156,40	★ 633,87

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

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	
Test [MB/s]							
Upload	Linux	EOS: xrdcp command	329,49	405,27	371,03	371,17	★ 5,39
		EOS Fusex	187,92	237,63	210,76	210,51	★ 9,49
		IBM Spectrum Scale	283,61	318,22	294,47	293,28	★ 6,79
		Ceph on Linux	141,00	163,37	157,56	158,17	★ 12,69
		Hadoop on Linux	9,74	10,10	9,91	9,91	★ 201,83
		EOS-wnc	158,40	331,09	231,25	227,75	★ 8,65
Download	Linux	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
		IBM Spectrum Scale	328,95	364,96	342,54	341,65	★ 5,84
		Ceph on Linux	188,80	355,49	265,08	264,04	★ 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	★ 11,75
Upload	Windows	EOS-drive ST	179,86	210,49	190,24	189,72	★ 10,51
		EOS: Samba	17,95	35,43	25,85	25,54	★ 77,37
		Ceph on Win	105,38	141,66	122,82	122,90	★ 16,28
		Hadoop on Win	4,30	4,73	4,55	4,55	★ 440,00
		EOS-wnc	119,92	213,86	170,17	169,49	★ 11,75
		EOS-drive ST	179,86	210,49	190,24	189,72	★ 10,51

Legend:	[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

Iterations (EOS)		21	(checksums OK)		Number of files		100	
Iterations (IBM)		21	(checksums OK)		File size [MB]		1	
Iterations (Ceph)		23	(checksums OK)					
Iterations (Hadoop)		14	(checksums OK)					
		Test [MB/s]	min	max	avg	trim25%	Avg time [ms]	
Upload	Linux	EOS: xrdcp command	142,86	181,82	165,24	165,87	★	6,05
		EOS Fusex	45,81	52,03	49,27	49,32		20,30
		IBM Spectrum Scale	54,08	58,82	55,00	55,02	☆	17,86
		Ceph on Linux	145,99	156,25	150,62	150,50	★	6,64
		Hadoop on Linux	3,46	3,64	3,55	3,55		281,92
	Windows	EOS-wnc	14,25	15,18	14,75	14,76	☆	67,78
		EOS-drive ST	9,70	10,00	9,88	9,89		101,22
		EOS: Samba	22,68	24,13	23,29	23,28	★	42,94
		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
Download	Linux	EOS: xrdcp command	95,24	196,08	169,56	174,52	★	5,90
		EOS Fusex	48,42	55,02	50,87	50,85		19,74
		IBM Spectrum Scale	158,73	187,27	174,76	175,16	★	5,72
		Ceph on Linux	7,47	110,15	54,15	57,40	☆	10,62
		Hadoop on Linux	3,68	4,30	3,97	3,97		251,71
	Windows	EOS-wnc	10,66	11,17	10,88	10,88		91,90
		EOS-drive ST	17,95	19,03	18,44	18,43	★	54,24
		EOS: Samba	13,19	15,82	14,28	14,24	☆	70,02
		Ceph on Win	1,93	47,25	35,38	37,60	★	28,26

Iterations (EOS)		28	(checksums OK)		Number of files		100
Iterations (IBM)		28	(checksums OK)		File size [MB]		1
Iterations (Ceph)		52	(checksums OK)				
Iterations (Hadoop)		11	(checksums OK)				
		Test [MB/s]	min	max	avg	Avg time [ms]	
Upload	Linux	EOS: xrdcp command	359,71	444,44	410,00	410,00	★
		EOS Fusex	134,72	192,64	160,00	160,00	
		IBM Spectrum Scale	176,62	188,71	180,00	180,00	☆
		Ceph on Linux	131,89	162,39	140,00	140,00	★
		Hadoop on Linux	9,43	10,17	10,00	10,00	
	Windows	EOS-wnc	174,21	204,60	180,00	180,00	☆
		EOS-drive ST	186,54	210,24	190,00	190,00	
		EOS: Samba	165,56	231,33	190,00	190,00	☆
		Ceph on Win	102,68	141,96	130,00	130,00	★
		Hadoop on Win	4,50	5,22	5,00	5,00	
Download	Linux	EOS: xrdcp command	301,20	436,68	410,00	410,00	★
		EOS Fusex	186,67	217,11	200,00	200,00	
		IBM Spectrum Scale	306,75	345,18	320,00	320,00	☆
		Ceph on Linux	20,44	183,49	30,00	30,00	★
		Hadoop on Linux	8,06	10,51	10,00	10,00	
	Windows	EOS-wnc	128,10	177,31	150,00	150,00	☆
		EOS-drive ST	148,70	185,92	150,00	150,00	
		EOS: Samba	72,97	97,50	80,00	80,00	☆
		Ceph on Win	17,63	81,00	30,00	30,00	★

Throughput results - Medium Files

Iterations (EOS)	28	(checksums OK)
Iterations (IBM)	28	(checksums OK)
Iterations (Ceph)	52	(checksums OK)
Iterations (Hadoop)	11	(checksums OK)
Number of files	10	
File size [MB]	100	

Iterations (EOS)	28	(checksums OK)
Iterations (IBM)	28	(checksums OK)
Iterations (Ceph)	52	(checksums OK)
Iterations (Hadoop)	11	(checksums OK)
Number of files	10	
File size [MB]	100	

Iterations (EOS)	27	(checksums OK)
Iterations (IBM)	28	(checksums OK)
Iterations (Ceph)	52	(checksums OK)
Iterations (Hadoop)	11	(checksums OK)
Number of files	10	
File size [MB]	100	

		Test [MB/s]	min	max	avg	trim25%	Avg time [ms]	
Upload	Linux	EOS: xrdcp command	359,71	444,44	411,14	412,67	★	243,22
		EOS Fusex	134,72	182,64	160,16	160,07	☆	624,38
		IBM Spectrum Scale	176,62	188,71	181,08	181,01	★	552,25
		Ceph on Linux	131,89	162,39	140,86	139,97		709,95
		Hadoop on Linux	9,43	10,17	9,94	9,97		10064,96
	Windows	EOS-wnc	174,21	204,60	186,28	185,66	☆	536,82
		EOS-drive ST	186,54	210,24	197,67	197,40	★	505,89
		EOS: Samba	165,56	231,33	196,82	196,65	★	508,08
		Ceph on Win	102,68	141,96	136,28	137,07		733,77
		Hadoop on Win	4,50	5,22	4,61	4,56		21670,61
Download	Linux	EOS: xrdcp command	301,20	436,68	412,94	417,50	★	242,16
		EOS Fusex	188,87	217,11	208,38	207,14	☆	484,59
		IBM Spectrum Scale	306,75	345,18	322,89	322,24	★	309,70
		Ceph on Linux	20,44	183,49	31,49	28,27		3175,40
		Hadoop on Linux	8,06	10,51	9,37	9,39		10668,22
	Windows	EOS-wnc	128,10	177,31	151,39	151,01	★	660,54
		EOS-drive ST	148,70	185,92	157,76	156,40	★	633,87
		EOS: Samba	72,97	97,50	81,26	80,39	☆	1230,60
		Ceph on Win	17,63	81,00	25,54	23,66		3915,54

		Test [MB/s]	min	max	avg	trim25%	Avg time [ms]	
Upload	Linux	EOS: xrdcp command	329,49	405,27	372,38	372,38	★	243,22
		EOS Fusex	187,92	237,63	212,77	212,77	☆	624,38
		IBM Spectrum Scale	283,61	318,22	300,91	300,91	★	552,25
		Ceph on Linux	141,00	163,37	152,18	152,18		709,95
		Hadoop on Linux	9,74	10,10	9,92	9,92		10064,96
	Windows	EOS-wnc	158,40	331,09	244,74	244,74	☆	536,82
		EOS-drive ST	212,22	294,47	253,34	253,34	★	505,89
		EOS: Samba	164,11	229,82	196,96	196,96	★	508,08
		Ceph on Win	128,18	158,19	143,18	143,18		733,77
		Hadoop on Win	4,61	4,72	4,66	4,66		21670,61
Download	Linux	EOS: xrdcp command	328,68	365,97	347,32	347,32	★	242,16
		EOS Fusex	218,66	233,36	226,01	226,01	☆	484,59
		IBM Spectrum Scale	328,95	364,96	346,95	346,95	★	309,70
		Ceph on Linux	188,80	355,49	272,14	272,14		3175,40
		Hadoop on Linux	9,28	10,63	9,95	9,95		10668,22
	Windows	EOS-wnc	119,92	213,86	166,89	166,89	★	660,54
		EOS-drive ST	179,86	210,49	195,17	195,17	★	633,87
		EOS: Samba	17,95	35,43	26,70	26,70	☆	1230,60
		Ceph on Win	105,38	141,66	123,52	123,52		3915,54

Throughput results - Large Files

s	10
	100
Avg time [ms]	
★	243,22
☆	624,38
★	552,25
	709,95
	10064,96
☆	536,82
★	505,89
★	508,08
	733,77
	21670,61
★	242,16
☆	484,59
★	309,70
	3175,40
	10668,22
★	660,54
★	633,87
☆	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	
	Test [MB/s]	min	max	avg	trim25%	Avg time [s]	
Upload	Linux	EOS: xrdcp command	329,49	405,27	371,03	371,17	★ 5,39
		EOS Fusex	187,97	237,63	210,76	210,51	☆ 9,49
		IBM Spectrum Scale	283,61	318,22	294,47	293,28	★ 6,79
		Ceph on Linux	141,00	163,37	157,56	158,17	12,69
		Hadoop on Linux	9,74	10,10	9,91	9,91	201,83
	Windows	EOS-wnc	158,40	331,09	231,25	227,75	★ 8,65
		EOS-drive ST	212,22	294,47	237,44	234,72	★ 8,42
		EOS: Samba	164,11	229,82	181,25	178,59	☆ 11,03
		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
Download	Linux	EOS: xrdcp command	328,68	365,97	353,00	354,47	★ 5,67
		EOS Fusex	218,99	233,39	227,13	227,13	8,81
		IBM Spectrum Scale	328,95	364,96	342,54	341,65	★ 5,84
		Ceph on Linux	188,80	203,43	203,08	204,04	☆ 7,54
		Hadoop on Linux	9,28	10,63	10,12	10,15	197,66
	Windows	EOS-wnc	119,92	213,86	170,17	169,49	★ 11,75
		EOS-drive ST	179,86	210,49	190,24	189,72	★ 10,51
		EOS: Samba	17,95	35,43	25,85	25,54	77,37
		Ceph on Win	105,38	141,66	122,82	122,90	☆ 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

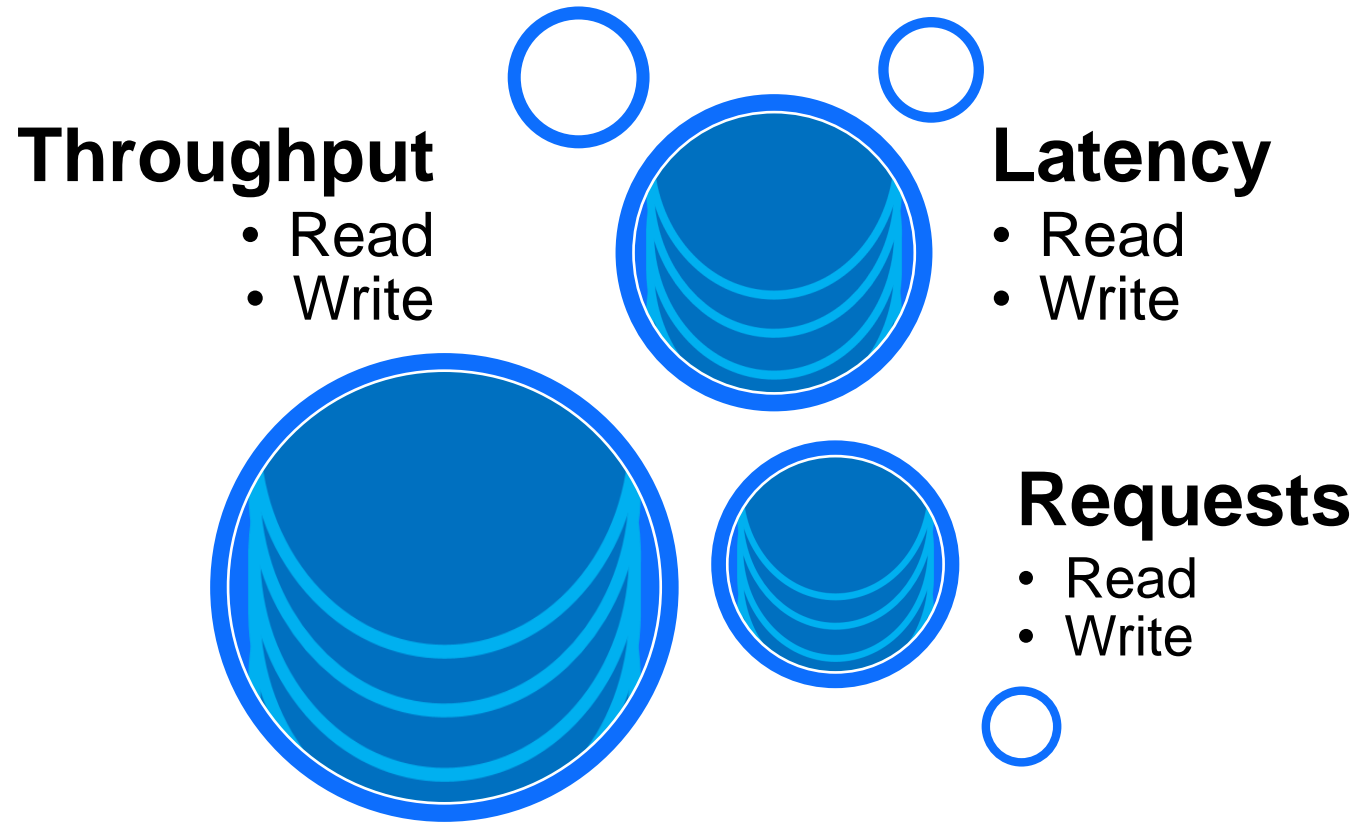
The best

- Small files
 1. EOS on Linux
 1. Ceph on Linux
 1. GPFS on Linux
- Medium files
 1. EOS on Linux
 2. GPFS on Linux
- Large files
 1. EOS on Linux
 2. 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
 - Mean time between failures
- 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