



SCALING STORAGE FOR EXASCALE JUST6 & EXASTORE AT JUELICH

01. APRIL 2026 | STEPHAN GRAF, OLEG TSIGENOV, MARTIN LISCHWESKI & SALEM EL SAYED

FZJ AT FIRST GLANCE

Founded in 1956



Revenue:
€ 987 million in
2023



Research priorities:
information, energy, bioeconomy



Research campus
with 14 institutes and
18 branch offices in
Germany and abroad
7400 Employees
114 countries



Shareholders:
Federal Republic of
Germany (90%),
federal state of
North Rhine-
Westphalia (10%)

Member of the
Helmholtz
Association



~300 Employees
(~280 FTE)



210 Scientists

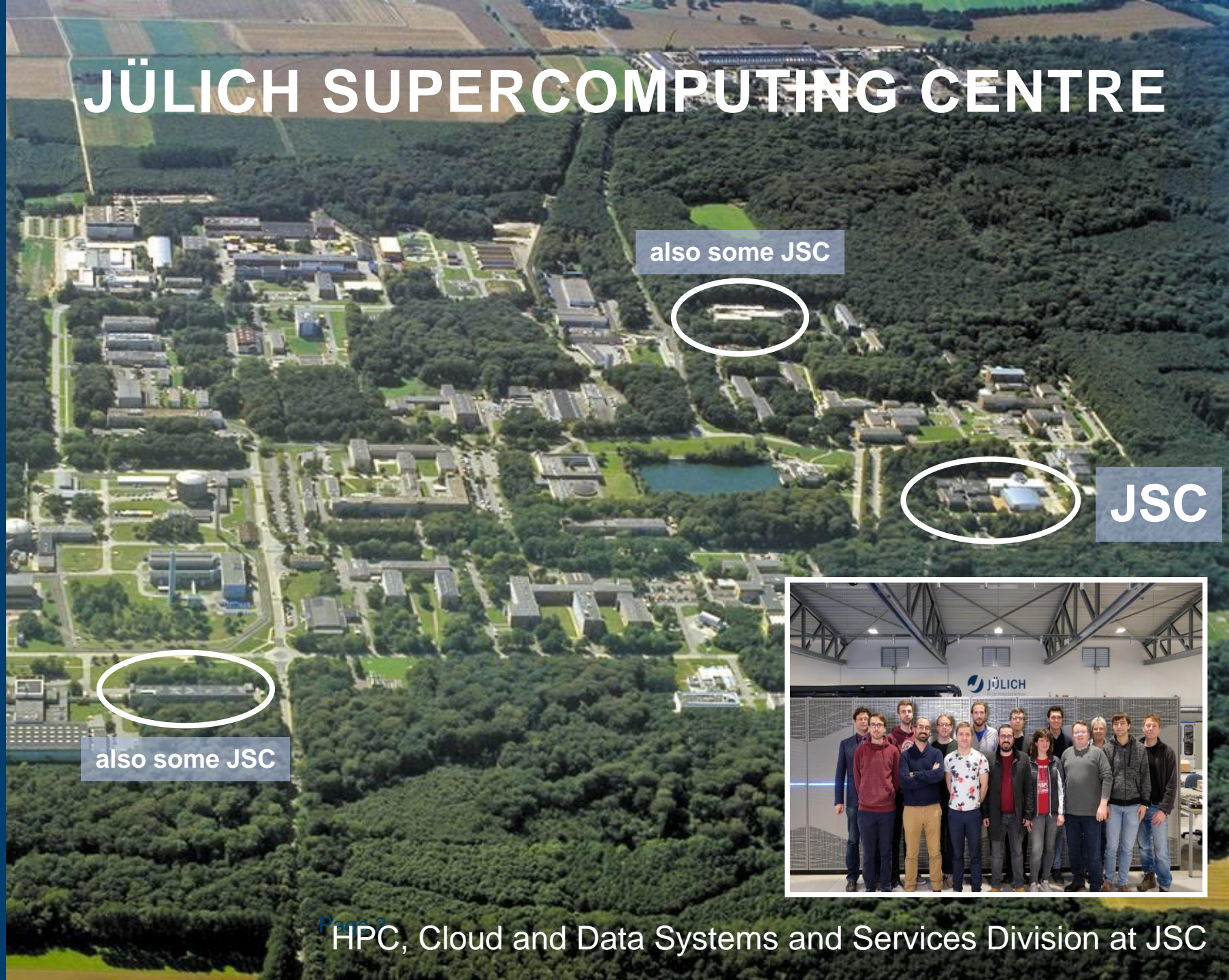


22 PhD Students (+ externals)
27 Students (Bachelor/Master)



fz-juelich.de/jsc

JÜLICH SUPERCOMPUTING CENTRE



also some JSC

JSC

also some JSC



HPC, Cloud and Data Systems and Services Division at JSC

JÜLICH SUPERCOMPUTING CENTRE

- **Supercomputer operation for**

- Centre – FZJ
- Region – RWTH Aachen University
- Germany – Gauss Centre for Supercomputing (GCS)
John von Neumann Institute for Computing (NIC)
- Europe – EuroHPC, PRACE, EU projects

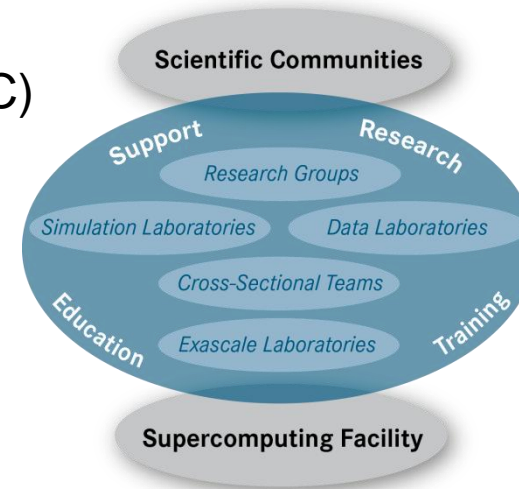
- **Application support**

- Unique support & research environment at JSC
- Peer review support and coordination

- **R&D work**

- Methods and algorithms, computational science, performance analysis and tools
- Scientific Machine Learning and AI with HPC
- Computer architectures, Co-Design, Exascale Labs together with IBM, Intel, NVIDIA

- **Education and training**



**First Exascale
supercomputer
in Europe**

DEEP



JUPITER

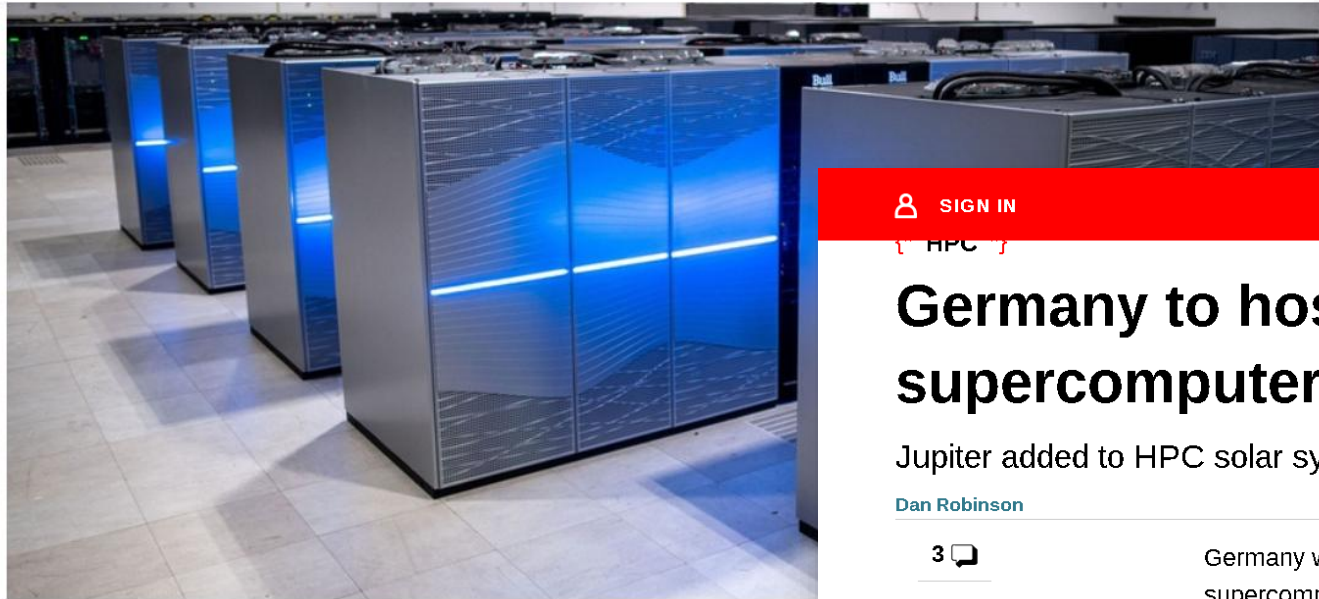
Europa's first ExaScale
Supercomputer

JUPITER – HOSTING ENTITY DECISION

15.06.2022



Startseite ▶ Wirtschaft ▶ Technologie ▶ Hochleistungs-Rechner: Supercomputer "Jupiter" kommt nach Jülich



Hochleistungs-Rechner

Supercomputer "Jupiter" k

Stand: 15.06.2022 16:43 Uhr

Das Forschungszentrum Jülich wird Standort für die ersten Exascale-Computers. "Jupiter" soll die Schal Rechenoperationen in der Sekunde durchbr

SPIEGEL Netzwelt

»Jupiter«

Jülich bekommt Europas ersten Exascale-Supercomputer

Das Forschungszentrum Jülich bekommt für eine halbe Milliarde Euro einen neuen Vorzeigerechner. Er soll helfen, Fragen zum Klimawandel zu beantworten – mit mehr als einer Trillion Rechenoperationen pro Sekunde.

15.06.2022, 16:52 Uhr

SIGN IN

The Register

HPC

Germany to host Europe's first exascale supercomputer

Jupiter added to HPC solar system

Dan Robinson

Thu 16 Jun 2022 // 07:38 UTC

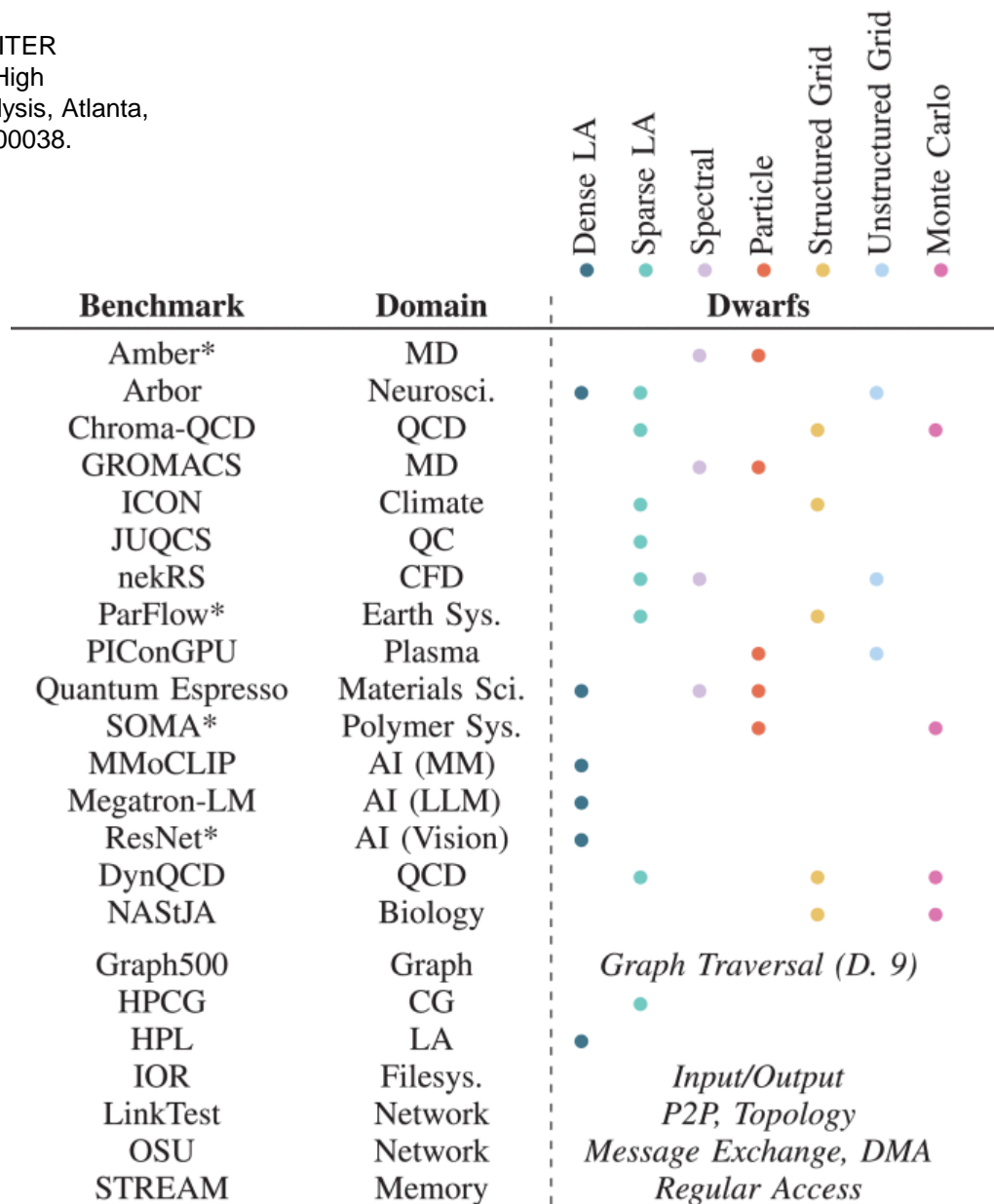
3

Germany will be the host of the first publicly known European exascale supercomputer, along with four other EU sites getting smaller but still powerful systems, the European High Performance Computing Joint Undertaking (EuroHPC JU) announced this week.

Germany **will be** the home of Jupiter, the "Joint Undertaking Pioneer for Innovative and Transformative Exascale Research." It should be switched on next year in a specially designed building on the campus of the **Forschungszentrum Jülich research centre** and operated by the Jülich Supercomputing Centre (JSC), alongside the existing Juwels and **Jureca** supercomputers.

PROCUREMENT BENCHMARK EFFORTS

A. Herten et al., "Application-Driven Exascale: The JUPITER Benchmark Suite," SC24: International Conference for High Performance Computing, Networking, Storage and Analysis, Atlanta, GA, USA, 2024, pp. 1-45, doi: 10.1109/SC41406.2024.00038.



JUPITER – THE BOOSTER

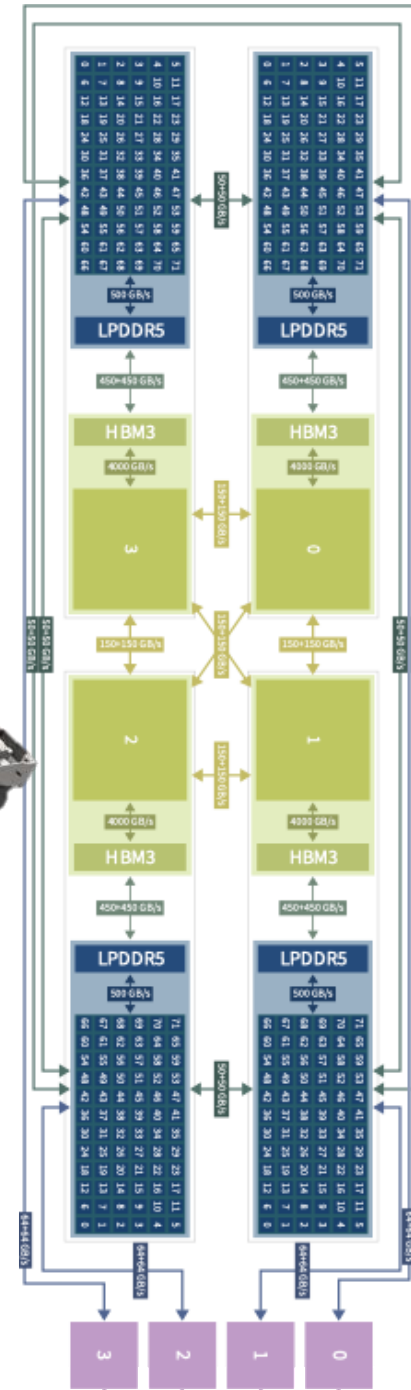
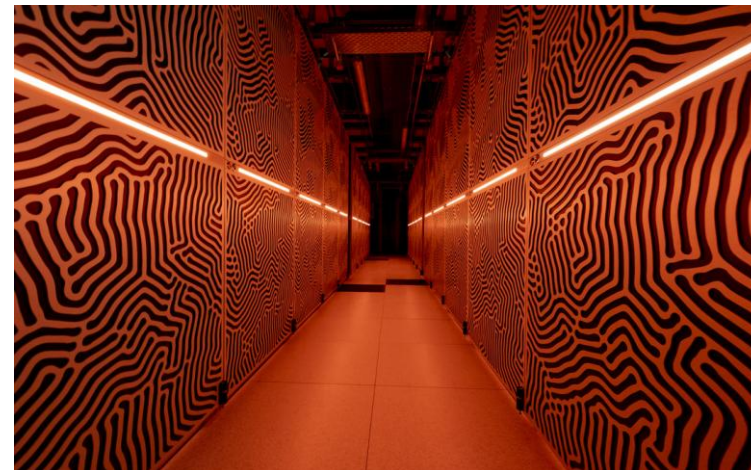
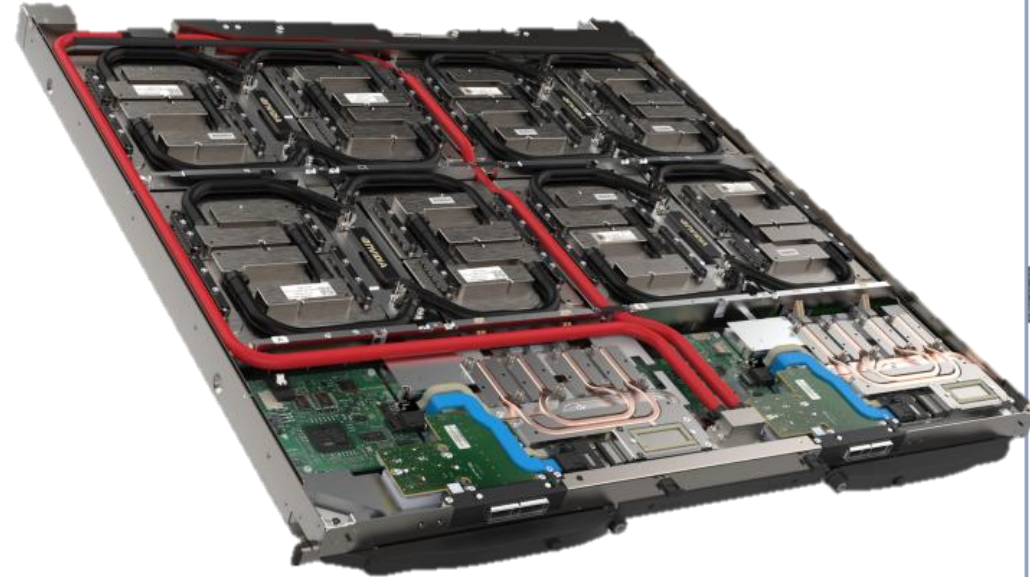
HIGHLY-SCALABLE MODULE FOR HPC AND AI WORKLOADS

EVIDEN

an atos business



- 1 EXAFLOP/S (FP64, HPL)
- NVIDIA GRACE-HOPPER CG1
 - ~5900 compute nodes
 - 4x CG1 chips per compute node
- NVIDIA MELLANOX NDR
 - 4x NDR200 NICs per compute node
- BULLSEQUANA XH3000
 - Direct Liquid Cooled blades
 - 2x compute node per blade



Backup Cold Water Cooling: 1MW

**Network:
3,2 Tbit/s**

Power (Campus): 2*60-80 MVA

Water (River Rur): up to 30 cbm/h



JUPITER ASCENDING

SINCE JANUARY 2025



JUWELS Booster:

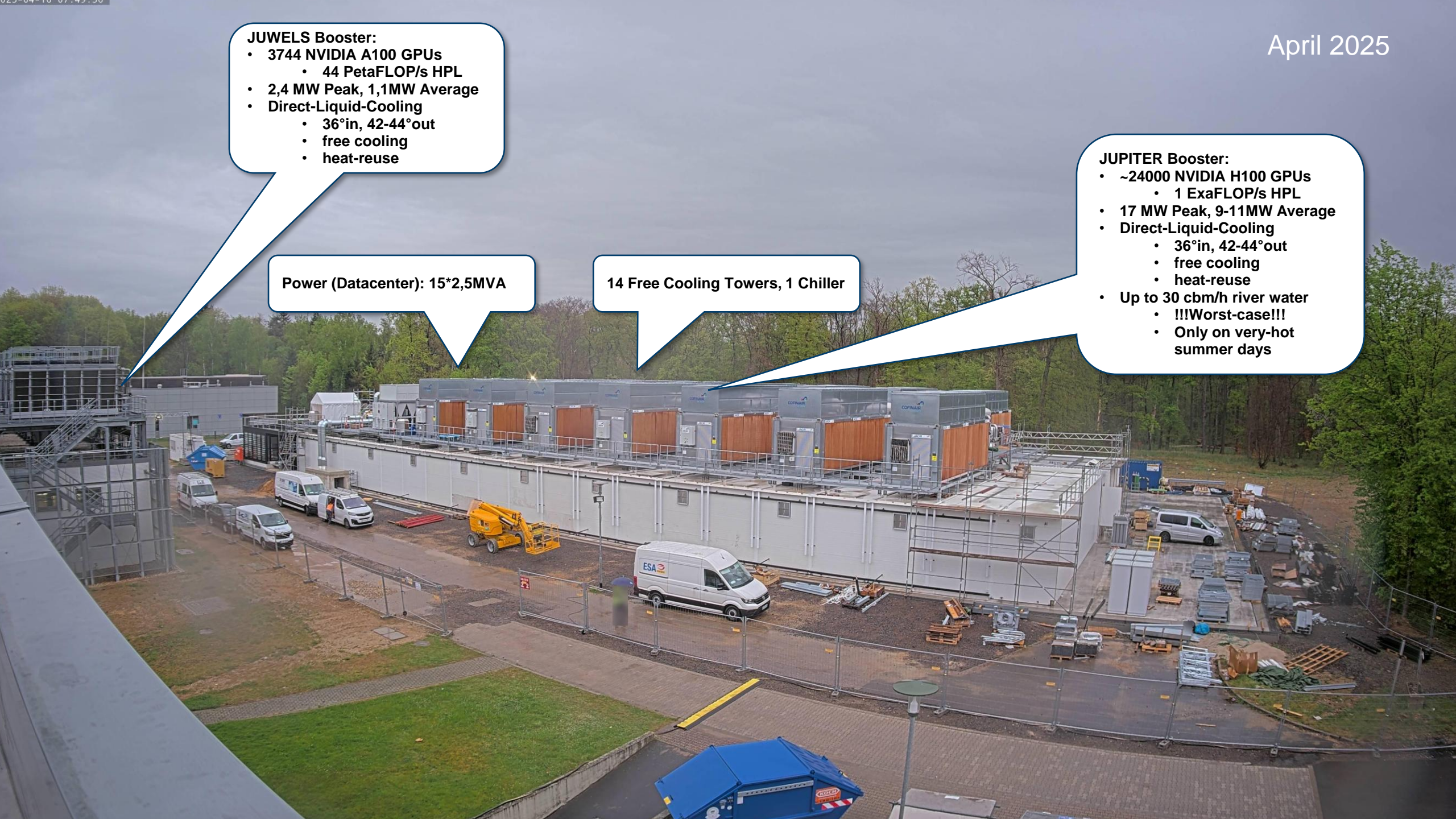
- 3744 NVIDIA A100 GPUs
 - 44 PetaFLOP/s HPL
- 2,4 MW Peak, 1,1MW Average
- Direct-Liquid-Cooling
 - 36°in, 42-44°out
 - free cooling
 - heat-reuse

Power (Datacenter): 15*2,5MVA

14 Free Cooling Towers, 1 Chiller

JUPITER Booster:

- ~24000 NVIDIA H100 GPUs
 - 1 ExaFLOP/s HPL
- 17 MW Peak, 9-11MW Average
- Direct-Liquid-Cooling
 - 36°in, 42-44°out
 - free cooling
 - heat-reuse
- Up to 30 cbm/h river water
 - !!!Worst-case!!!
 - Only on very-hot summer days





September 05, 2025: Inauguration of JUPITER

SUPERCOMPUTER

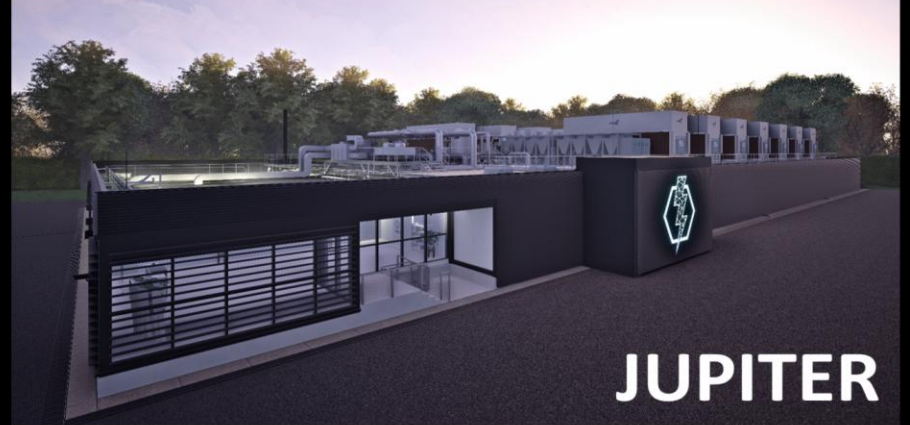
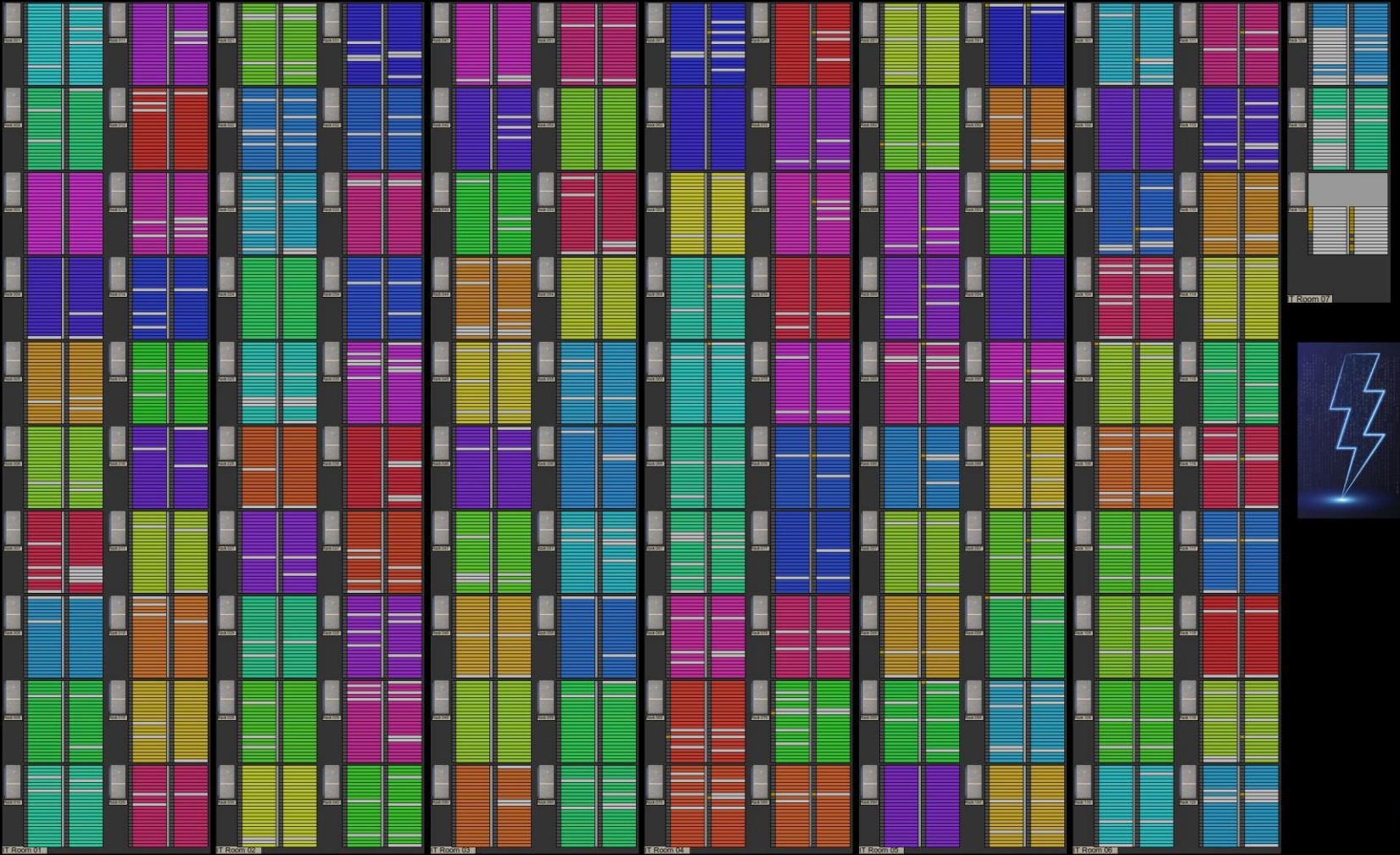
Launch of the Exascale Supercomputer JUPITER

JUPITER was unveiled to the public today at Forschungszentrum Jülich in front of high-ranking guests from politics, science, and industry. The German Chancellor Friedrich Merz and the Minister-President of North Rhine-Westphalia Hendrik Wüst attended the inauguration ceremony of Europe's fastest and most energy-efficient supercomputer.

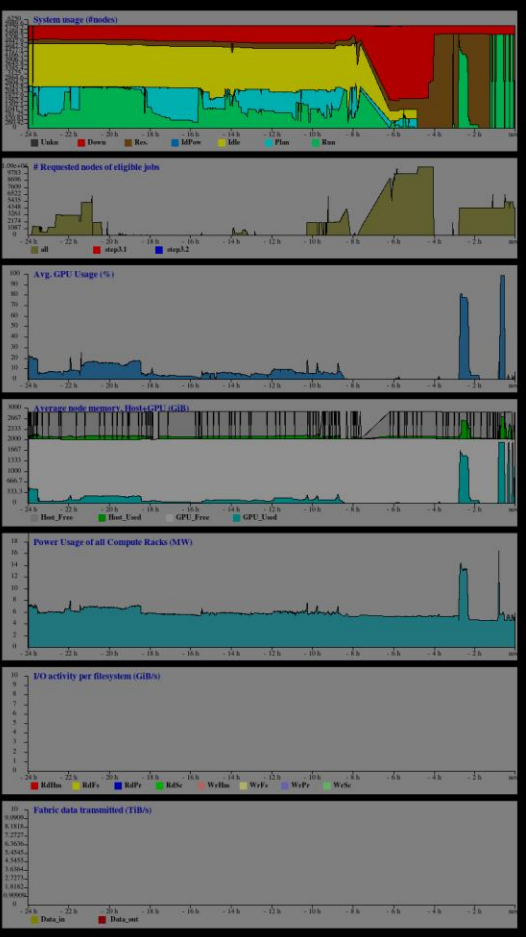


Inauguration of JUPITER: Ceremonial act with Prof. Dr. Dr. Thomas Lippert, Director of Jülich Supercomputing Centre (JSC), Prof. Dr. Kristel Michielsens, Director of Jülich Supercomputing Centre (JSC), Ekaterina Zaharieva, Commissioner Startups, Research and Innovation, European Commission, Hendrik Wüst, Minister-President of North Rhine-Westphalia, Prof. Dr. Astrid Lambrecht, Chair of the Board of Directors at FZJ, Friedrich Merz, Federal Chancellor, Dorothee Bär, Federal Minister for Research, Technology and Space, Karsten Wildberger, Federal Ministry for Digital Transformation and Government Modernisation, Ina Brandes, Minister for Culture and Science of North Rhine-Westphalia, Prof. Dr. Laurens Kuipers, Member of the Board of Directors at FZJ. Fotos: Forschungszentrum Jülich / Kurt Steinhausen.





node	job	status	start	end	cores	gpus
ppc002-47	ppc002-48	Running	17:45:00	17:50:00	1	1
ppc002-40	ppc002-40	Running	17:45:00	17:50:00	1	1
ppc002-43	ppc002-44	Running	17:45:00	17:50:00	1	1
ppc002-41	ppc002-42	Running	17:45:00	17:50:00	1	1
ppc002-38	ppc002-40	Running	17:45:00	17:50:00	1	1
ppc002-37	ppc002-38	Running	17:45:00	17:50:00	1	1
ppc002-38	ppc002-38	Running	17:45:00	17:50:00	1	1
ppc002-38	ppc002-38	Running	17:45:00	17:50:00	1	1
ppc002-31	ppc002-32	Running	17:45:00	17:50:00	1	1
ppc002-35	ppc002-36	Running	17:45:00	17:50:00	1	1
ppc002-27	ppc002-28	Running	17:45:00	17:50:00	1	1
ppc002-25	ppc002-26	Running	17:45:00	17:50:00	1	1
ppc002-23	ppc002-24	Running	17:45:00	17:50:00	1	1
ppc002-21	ppc002-22	Running	17:45:00	17:50:00	1	1
ppc002-19	ppc002-20	Running	17:45:00	17:50:00	1	1
ppc002-17	ppc002-18	Running	17:45:00	17:50:00	1	1
ppc002-15	ppc002-16	Running	17:45:00	17:50:00	1	1
ppc002-13	ppc002-14	Running	17:45:00	17:50:00	1	1
ppc002-11	ppc002-12	Running	17:45:00	17:50:00	1	1
ppc002-09	ppc002-10	Running	17:45:00	17:50:00	1	1
ppc002-07	ppc002-08	Running	17:45:00	17:50:00	1	1
ppc002-05	ppc002-06	Running	17:45:00	17:50:00	1	1
ppc002-03	ppc002-04	Running	17:45:00	17:50:00	1	1
ppc002-01	ppc002-02	Running	17:45:00	17:50:00	1	1



node	job	status	start	end	cores	gpus
ppc002-47	ppc002-48	Running	17:45:00	17:50:00	1	1
ppc002-40	ppc002-40	Running	17:45:00	17:50:00	1	1
ppc002-43	ppc002-44	Running	17:45:00	17:50:00	1	1
ppc002-41	ppc002-42	Running	17:45:00	17:50:00	1	1
ppc002-38	ppc002-40	Running	17:45:00	17:50:00	1	1
ppc002-37	ppc002-38	Running	17:45:00	17:50:00	1	1
ppc002-38	ppc002-38	Running	17:45:00	17:50:00	1	1
ppc002-38	ppc002-38	Running	17:45:00	17:50:00	1	1
ppc002-31	ppc002-32	Running	17:45:00	17:50:00	1	1
ppc002-35	ppc002-36	Running	17:45:00	17:50:00	1	1
ppc002-27	ppc002-28	Running	17:45:00	17:50:00	1	1
ppc002-25	ppc002-26	Running	17:45:00	17:50:00	1	1
ppc002-23	ppc002-24	Running	17:45:00	17:50:00	1	1
ppc002-21	ppc002-22	Running	17:45:00	17:50:00	1	1
ppc002-19	ppc002-20	Running	17:45:00	17:50:00	1	1
ppc002-17	ppc002-18	Running	17:45:00	17:50:00	1	1
ppc002-15	ppc002-16	Running	17:45:00	17:50:00	1	1
ppc002-13	ppc002-14	Running	17:45:00	17:50:00	1	1
ppc002-11	ppc002-12	Running	17:45:00	17:50:00	1	1
ppc002-09	ppc002-10	Running	17:45:00	17:50:00	1	1
ppc002-07	ppc002-08	Running	17:45:00	17:50:00	1	1
ppc002-05	ppc002-06	Running	17:45:00	17:50:00	1	1
ppc002-03	ppc002-04	Running	17:45:00	17:50:00	1	1
ppc002-01	ppc002-02	Running	17:45:00	17:50:00	1	1



CERTIFICATE

JUPITER Booster - BullSequana XH3000, GH Superchip 72C 3GHz, NVIDIA GH200 Superchip
Quad-Rail NVIDIA InfiniBand NDR200, RedHat Enterprise Linux
EuroHPC/FZJ, Germany

is ranked

No. 14

among the World's TOP500 Supercomputers
with **63.316 GFlops/watts Performance**
in the Green500 List published at the SC2025
Conference on November 18, 2025.

Wu-chun Feng
Virginia Tech

Editors

Kirk Cameron
Virginia Tech

Erich Strohmaier
NERSC/Berkeley Lab



CERTIFICATE

JUPITER Booster - BullSequana XH3000, GH Superchip 72C 3GHz, NVIDIA GH200 Superchip,
Quad-Rail NVIDIA InfiniBand NDR200, RedHat Enterprise Linux
EuroHPC/FZJ, Germany

is ranked

No. 4

among the World's TOP500 Supercomputers
with **1000.00 PFlop/s Linpack Performance**
in the 66th TOP500 List published at the SC2025
Conference on November 18, 2025.

Jack Dongarra
University of Tennessee

Horst Simon
NERSC/Berkeley Lab

Martin Meuer
Prometeu



JÜLICH
Forschungszentrum

CHOOSING A STORAGE SYSTEM

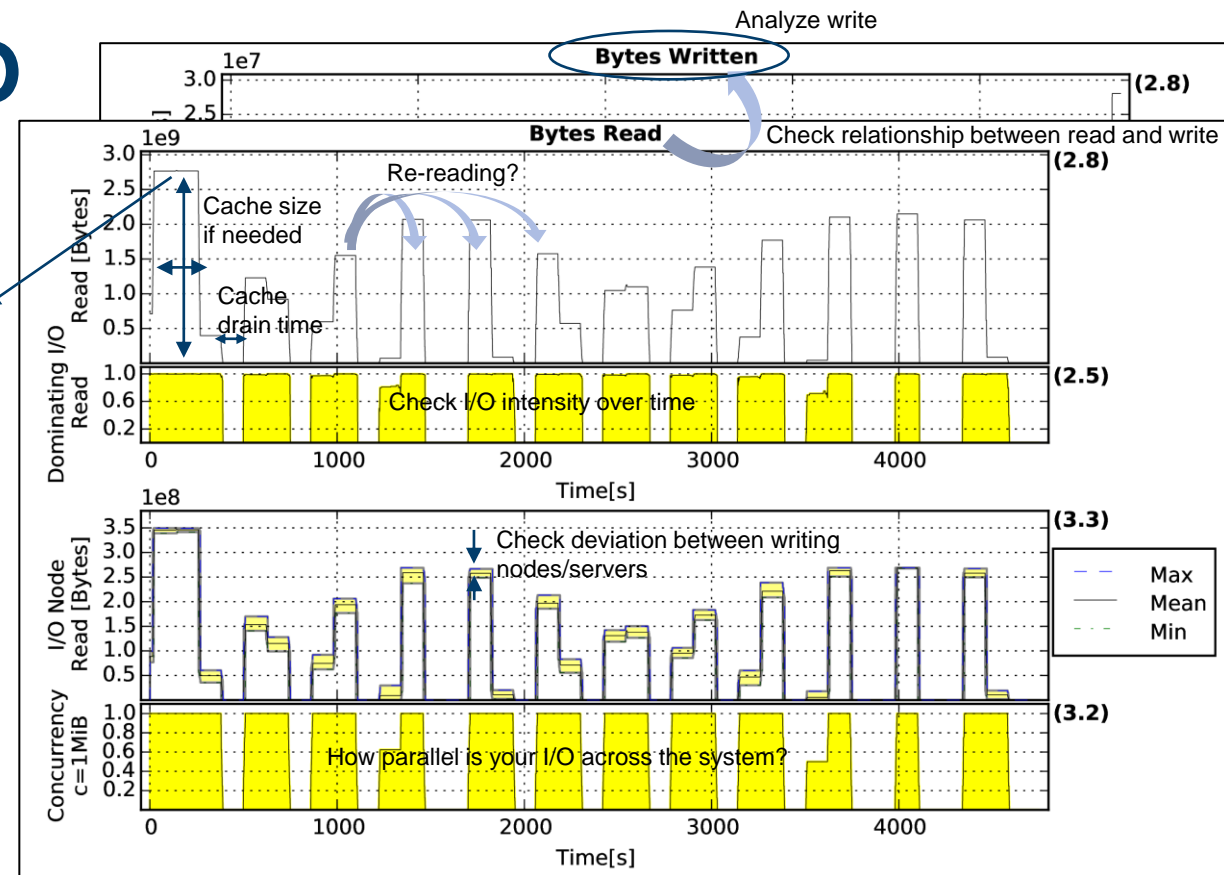
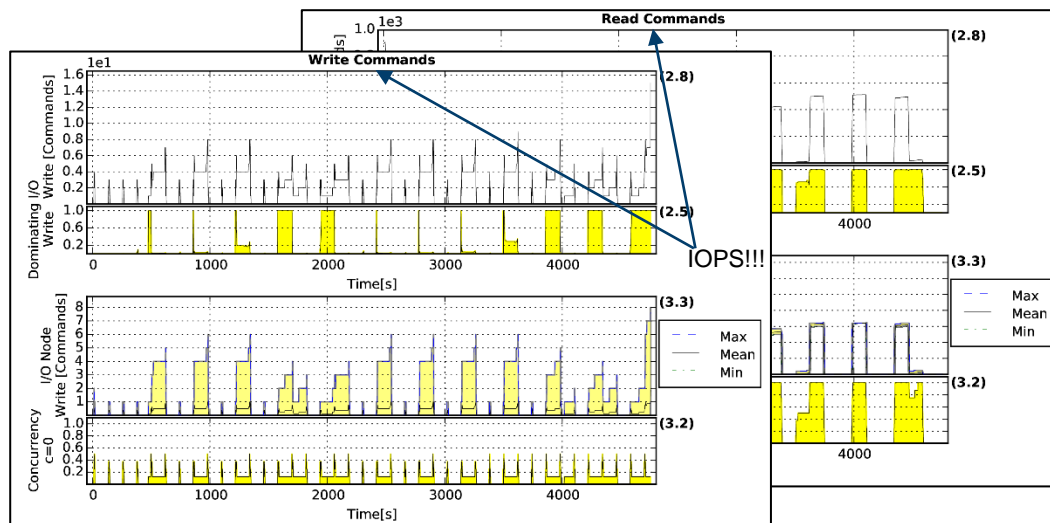
MEASURE & MATCH YOUR I/O

Without DATA it is just an opinion!

Hints:

- Do not over fit and add margins
- This was measured on your old system
- Make room for future expansion

Max Read bandwidth



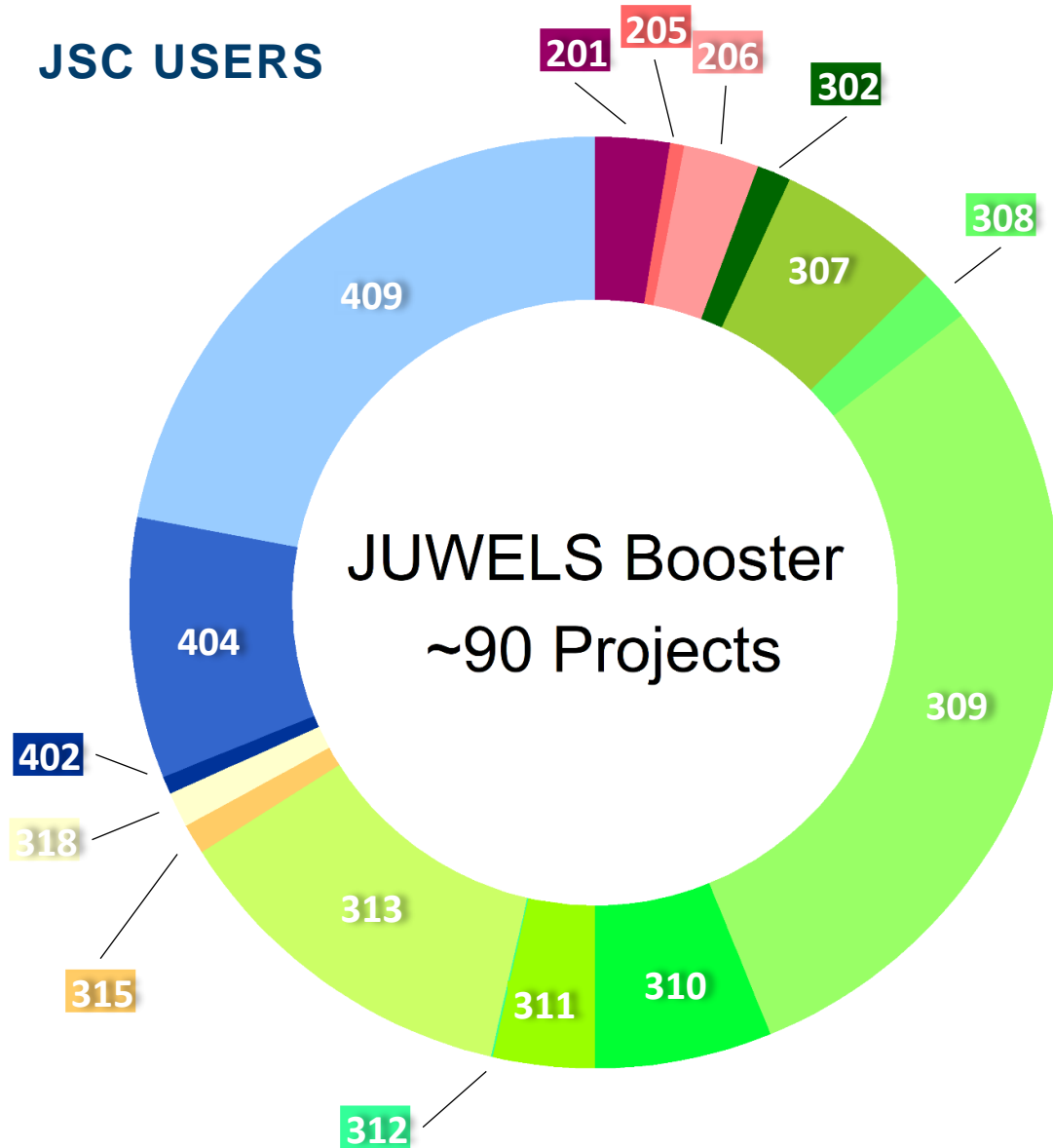
Besides performance:

- Capacity
- I/O interface
- Legacy infrastructure
- Retention time



WE DON'T HAVE THAT LUXURY

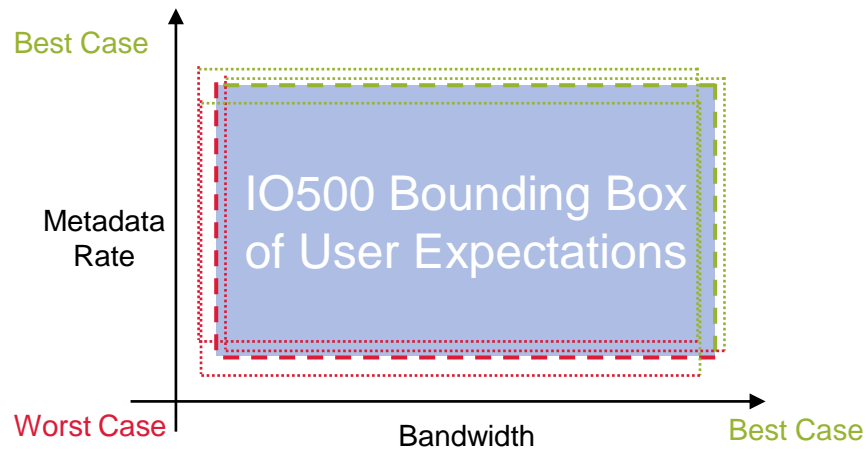
JSC USERS



Research Fields

- 201** Basic Biological and Medical Research
- 205** Medicine
- 206** Neurosciences
- 302** Chemical Solid State and Surface Research
- 307** Condensed Matter Physics
- 308** Optics, Quantum Optics and Physics of Atoms, Molecules and Plasmas
- 309** Particles, Nuclei and Fields
- 310** Statistical Physics, Soft Matter, Biological Physics, Nonlinear Dynamics
- 311** Astrophysics and Astronomy
- 312** Mathematics
- 313** Atmospheric Science, Oceanography and Climate Research
- 315** Geophysics and Geodesy
- 318** Water Research
- 402** Mechanics and Constructive Mechanical Engineering
- 404** Heat Energy Technology, Thermal Machines, Fluid Mechanics
- 409** Computer Science

BOUNDING BOX OF USER EXPECTATIONS



IO⁵⁰⁰

[\(\[HTTPS://IO500.ORG/ABOUT\]\(https://io500.org/about\)\)](https://io500.org/about)

- **IOEasy**: Continues read/write of file per process + 4k random read
- **IOHard**: Interleaved read/write of 47008-byte blocks
- **MDEasy**: Write/stat/delete zero-byte files in a unique directory per process
- **MDHard**: Write/stat/read/delete small files (3901 bytes) in a shared directory
- Find: Finding relevant objects based on patterns

R. Liem et al. "User-Centric System Fault Identification Using IO500 Benchmark," 2021 IEEE/ACM Sixth International Parallel Data Systems Workshop (PDSW), St. Louis, MO, USA, 2021, pp. 35-40, doi: 10.1109/PDSW54622.2021.00011.

OUR BOUNDING BOX OF EXPECTATIONS

Free to choose:

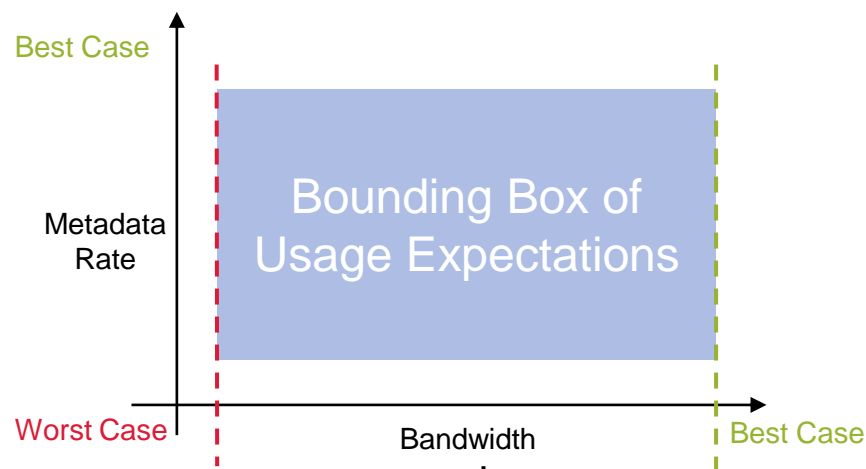
- Number of nodes
- Number of processes per node

Restrictions:

- Run at least 1min
- Avoid caching

File system block size:

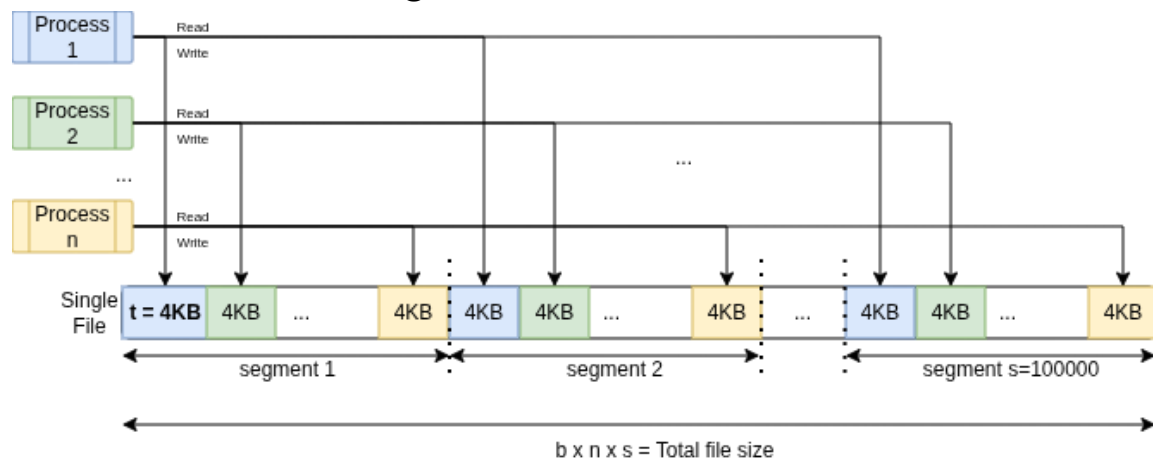
- 16MB
- 512KB



Common IOR Parameters	
-C	reorderTasksConstant
-Q	taskPerNodeOffset (choose in relation to --distribution from slurm)
-g	intraTestBarriers
-a	API
-i	Repetitions

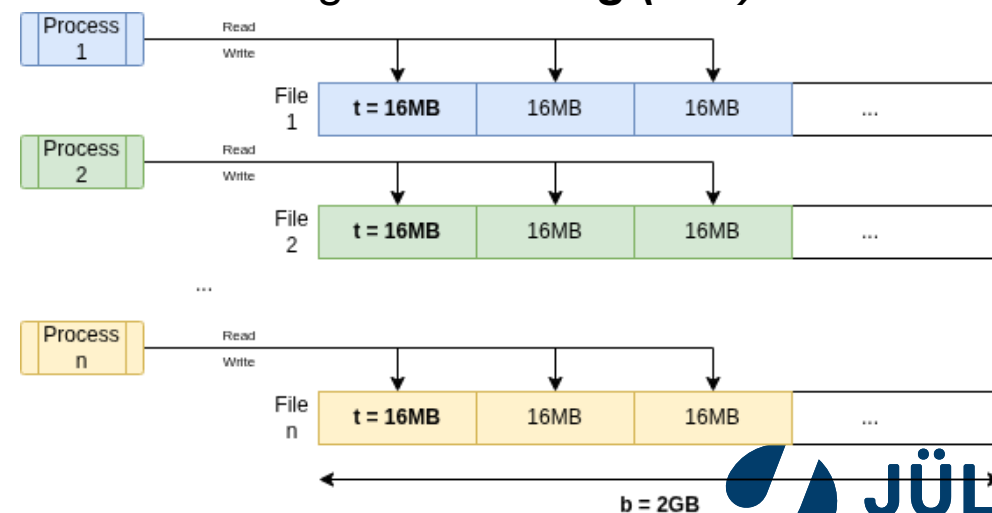
IOR Hard

`ior -C -Q 1 -g -t 4k -b 4k -s 100000 -a POSIX -i 3`



IOR Easy

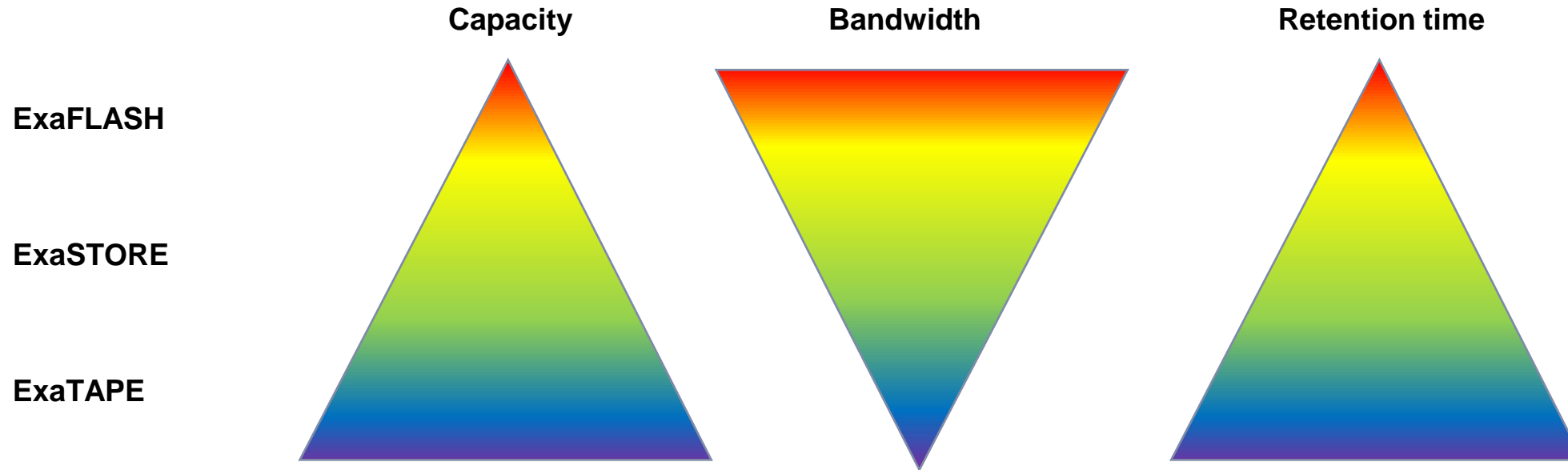
`ior -C -Q 1 -g -t 16m -b 2g (-s 1) -a POSIX -i 3`



JUST

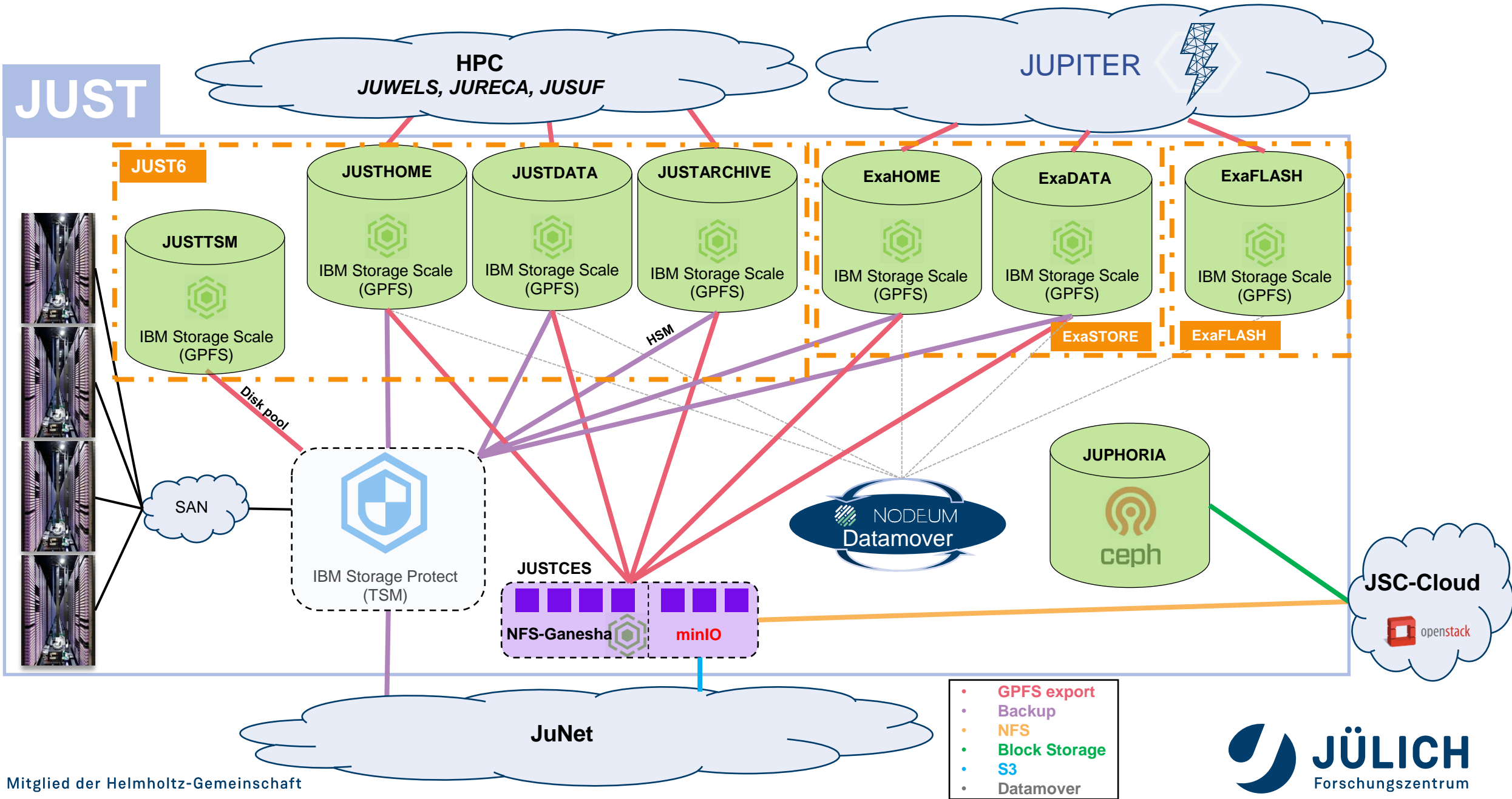
JUelich STorage

TIERED STORAGE OFFERING FOR JUPITER



- ExaFLASH: IBM SSS6K cluster (29 PB raw capacity on NVMe)
- ExaSTORE: IBM SSS6K + ESS3500 cluster (308 PB raw capacity on 22 TB HDD)
- ExaTAPE: Tape storage (Backup + GPFS&TSM-HSM as part of JUST6, 2x TS4500, ~370 PB)

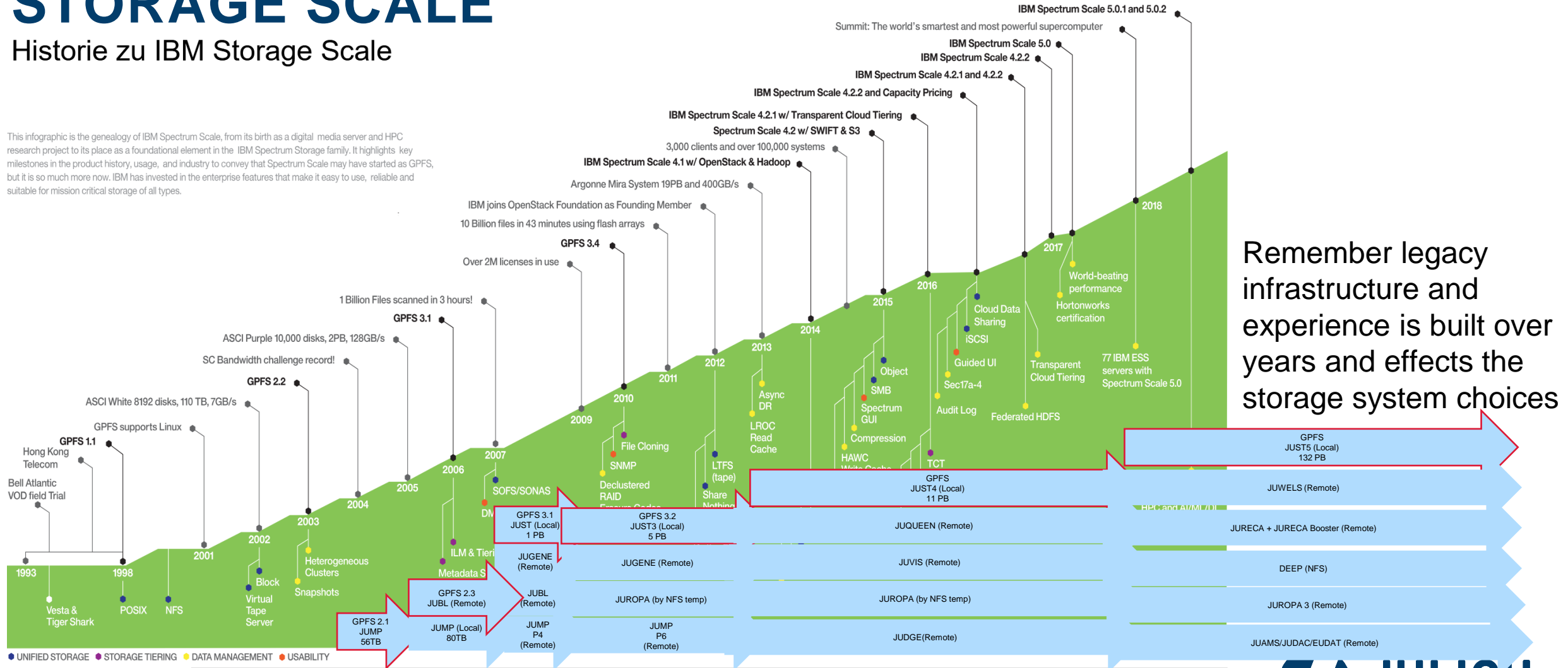
JUST



HISTORY OF GPFS / SPECTRUM SCALE / STORAGE SCALE

Historie zu IBM Storage Scale

This infographic is the genealogy of IBM Spectrum Scale, from its birth as a digital media server and HPC research project to its place as a foundational element in the IBM Spectrum Storage family. It highlights key milestones in the product history, usage, and industry to convey that Spectrum Scale may have started as GPFS, but it is so much more now. IBM has invested in the enterprise features that make it easy to use, reliable and suitable for mission critical storage of all types.



Remember legacy infrastructure and experience is built over years and effects the storage system choices

STORAGE SCALE 6000

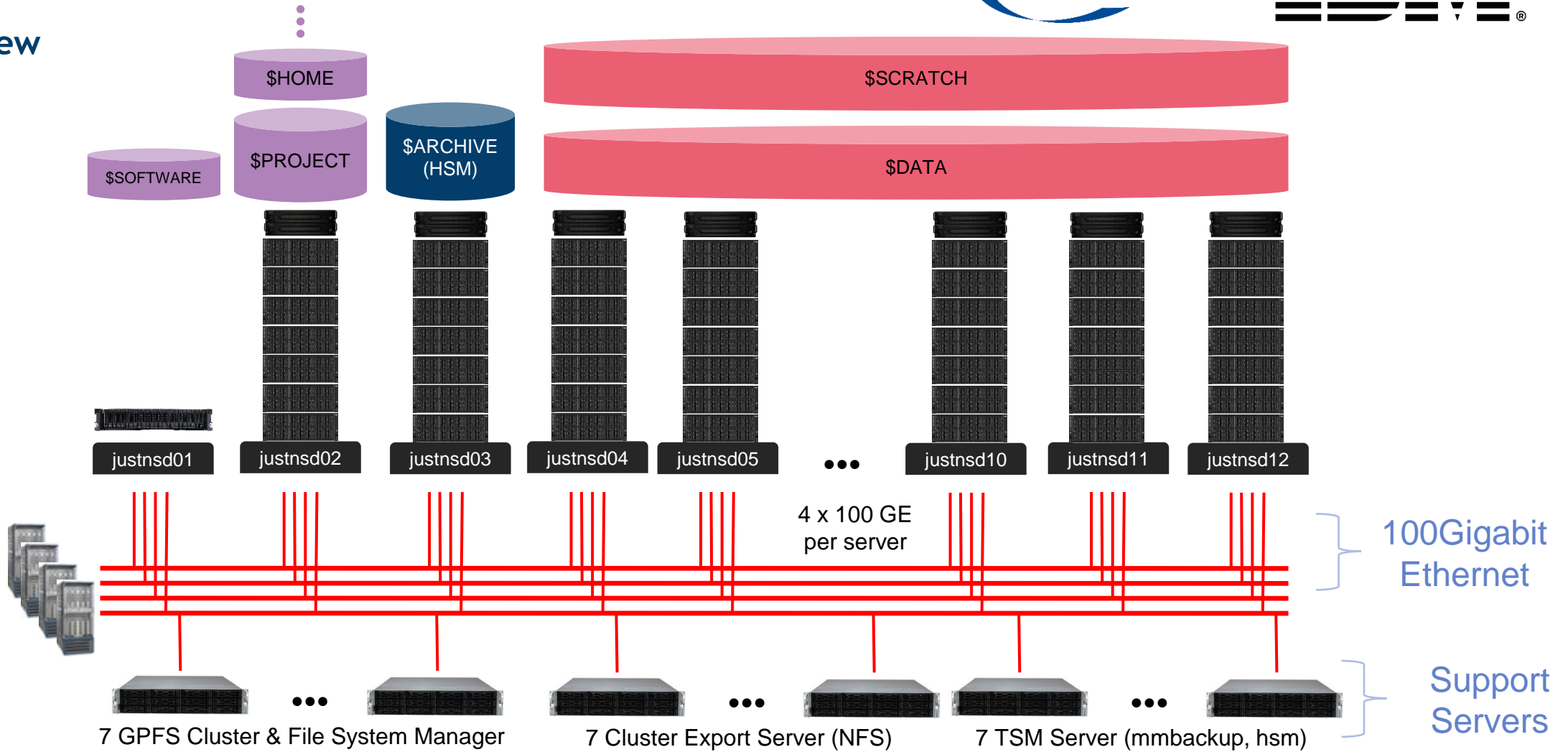
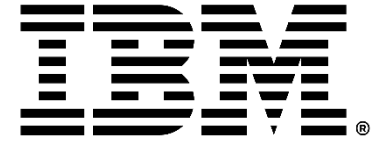


Storage Scale System 6000

- 2 x AMD EPYC Genoa 9454 CPU 48c, 3,8GHz
- 1024GB DDR5 memory in 16 x 64
- 2 x CX-7 NDR200 Dual Port
- 4 x SAS4 x8 G4 Adapter
- 7 enclosures (91 x 22TB HDD drives) and a raw capacity of 14 PB

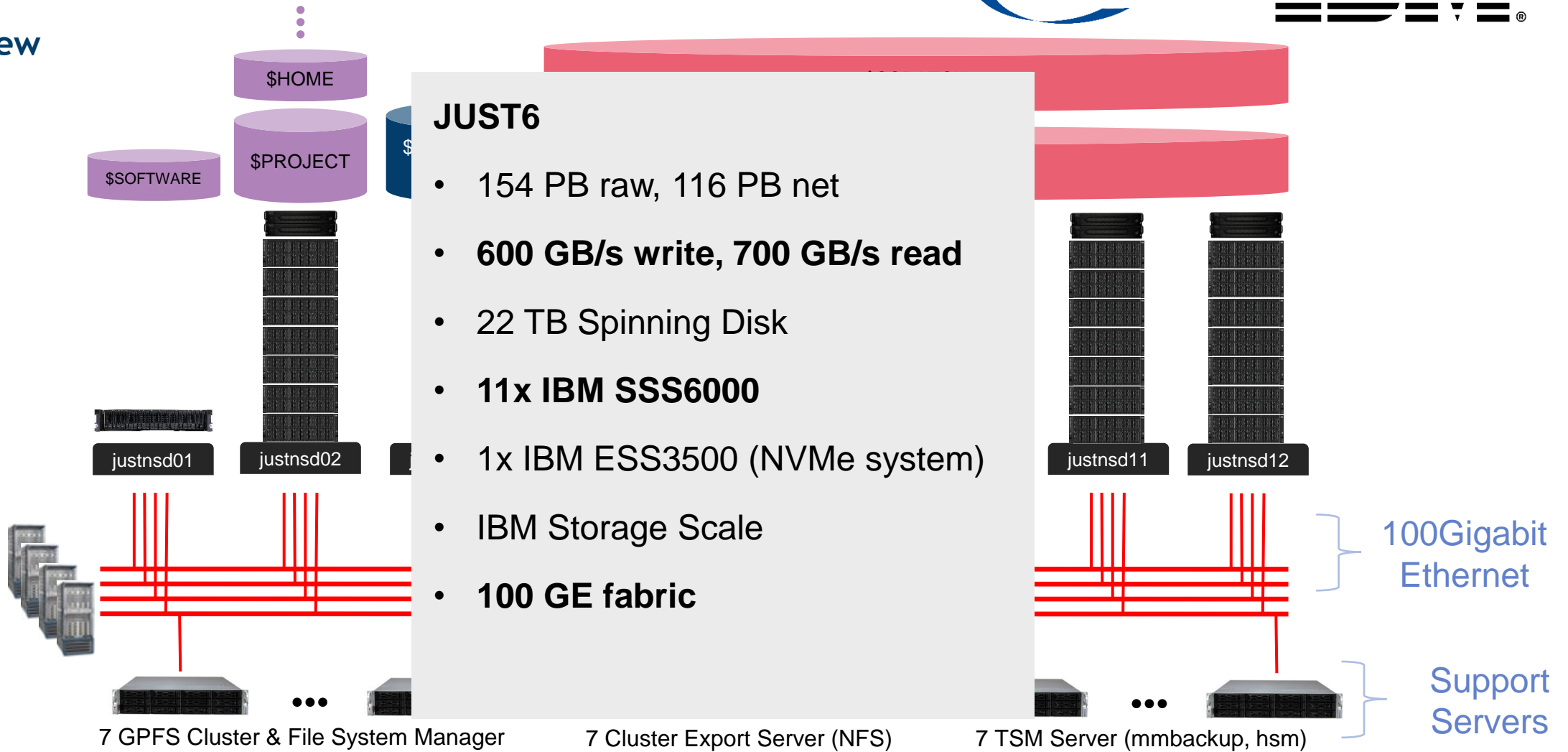
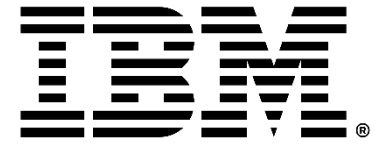
JUST6

Physical view



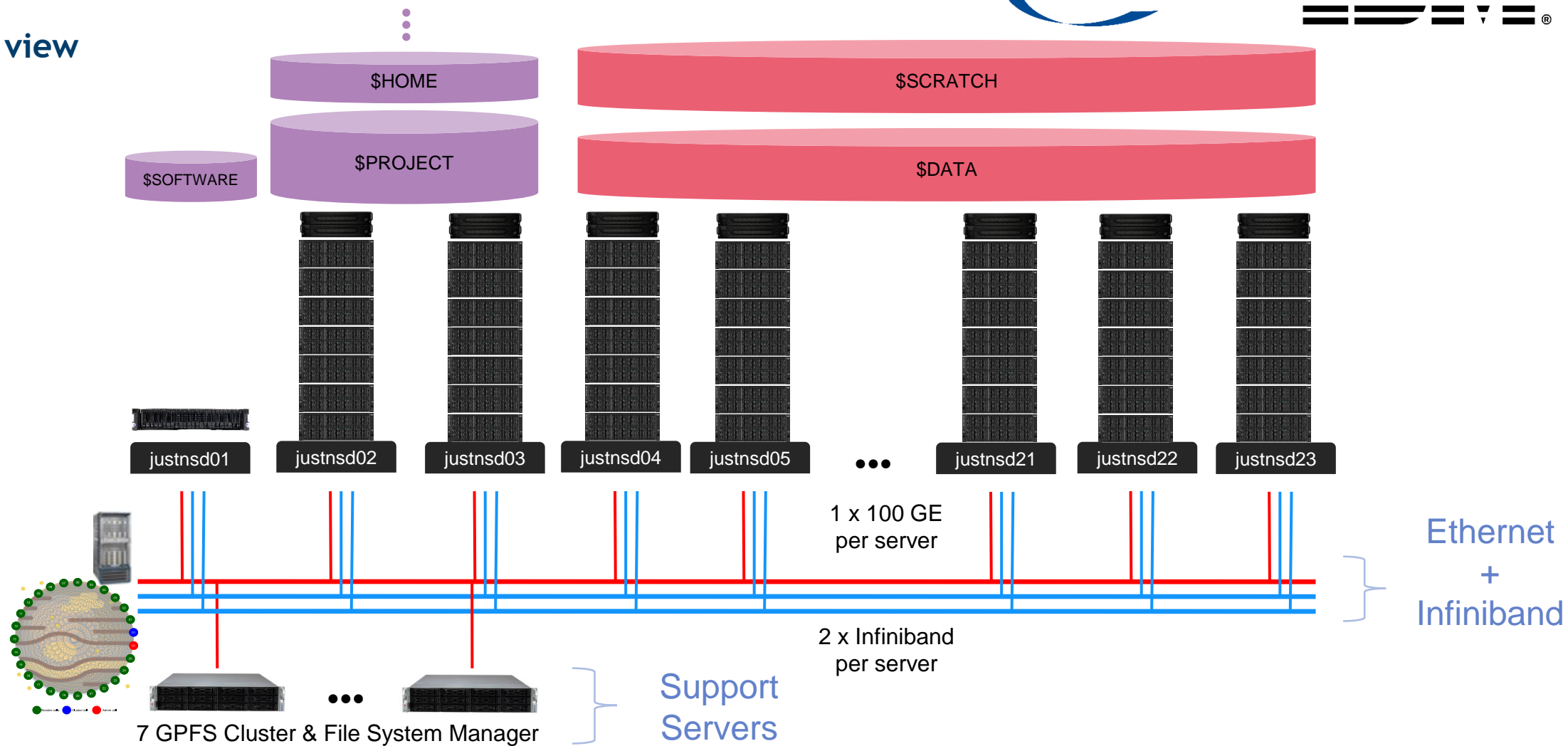
JUST6

Physical view



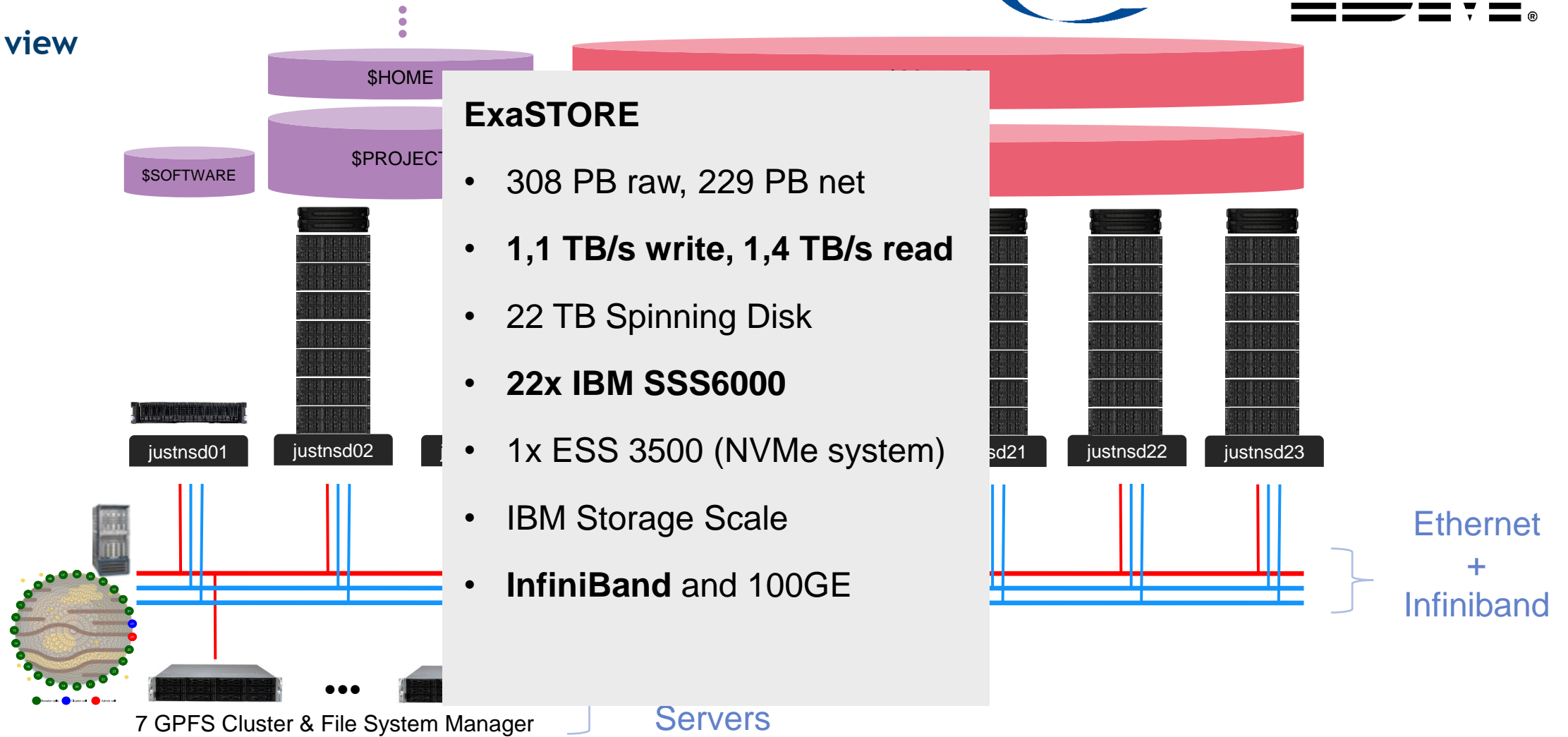
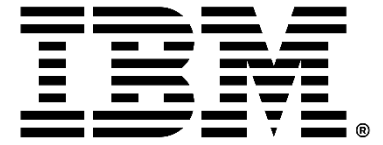
EXASTORE

Physical view



EXASTORE

Physical view

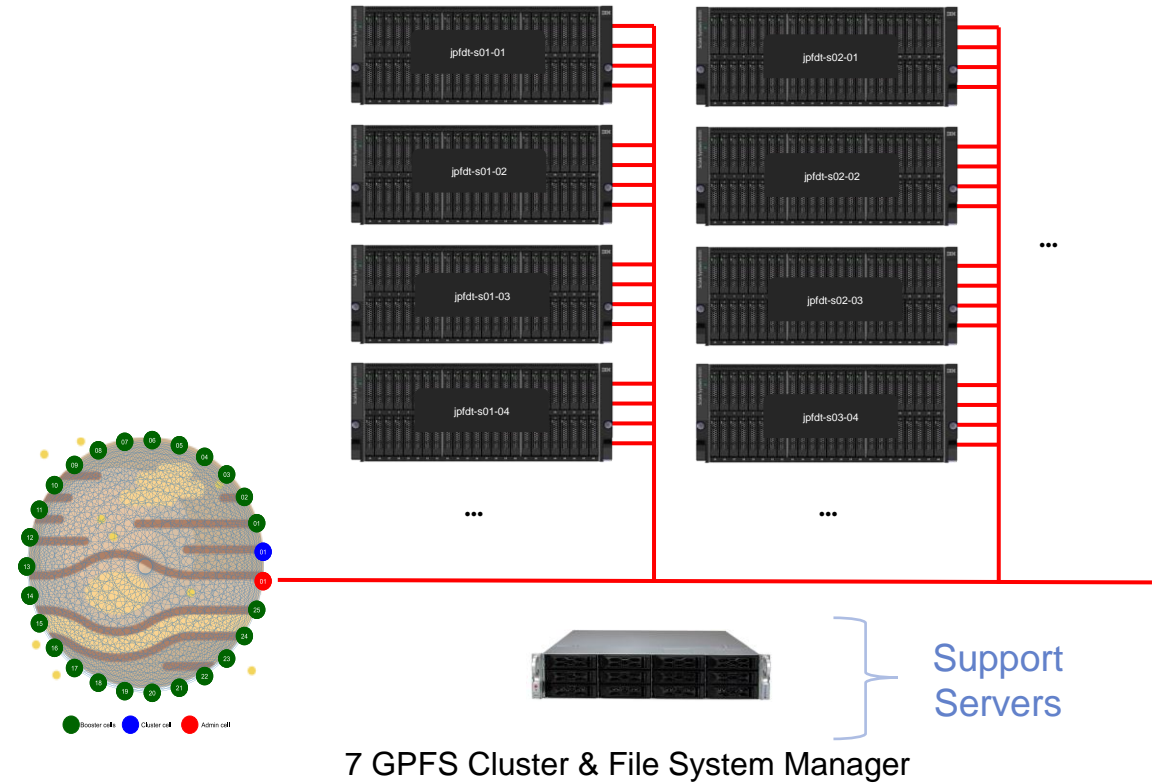


EXAFLASH

Physical View

ExaFLASH

- 29 PB raw, 21PB net
- **2,1 TB/s write, 3,1 TB/s read**
- 30 TB NVME
- **20x IBM SSS6000**
- IBM Storage Scale
- **InfiniBand**



ACCEPTANCE !



ACCEPTANCE TEST LIST

IOR Easy

Each Building Block

IOR Easy 16MB blocksize file system

IOR Easy 512KB blocksize file system

All Building Blocks

IOR Easy 16MB blocksize file system

IOR Easy 512KB blocksize file system

IOR Hard

Each Building Block

IOR Hard 16MB blocksize file system

IOR Hard 512KB blocksize file system

All Building Blocks

IOR Hard 16MB blocksize file system

IOR Hard 512KB blocksize file system

MDTest

Each Building Block

MDTest 16MB block size file system

MDTest 512KB block size file system

All Building Block

MDTest 16MB block size file system

MDTest 512KB block size file system

DISCLAIMERS

- **This is a summary:** Not all acceptance events can/should be listed here. The following points are a subset of the hundreds of tests and benchmarks. The examples have been extracted from hundreds of emails with long subject lines containing many `Re:` and `Aw:`. Additionally, any meeting minutes and acceptance documents available have been used.
- **Based on real events:** However some events have been dramatized for effect, others we wish were not true. The persons involved are real and their pain has been scientifically proven to exist.
- **Depends on remembering:** This one too was in a galaxy far far away! At least it feels that way. Events are partially pulled out of our Random Access Memory. Expect events not to be in the correct order or missing curtail details. Ask us or better yet, one of our trusted IBM colleagues.

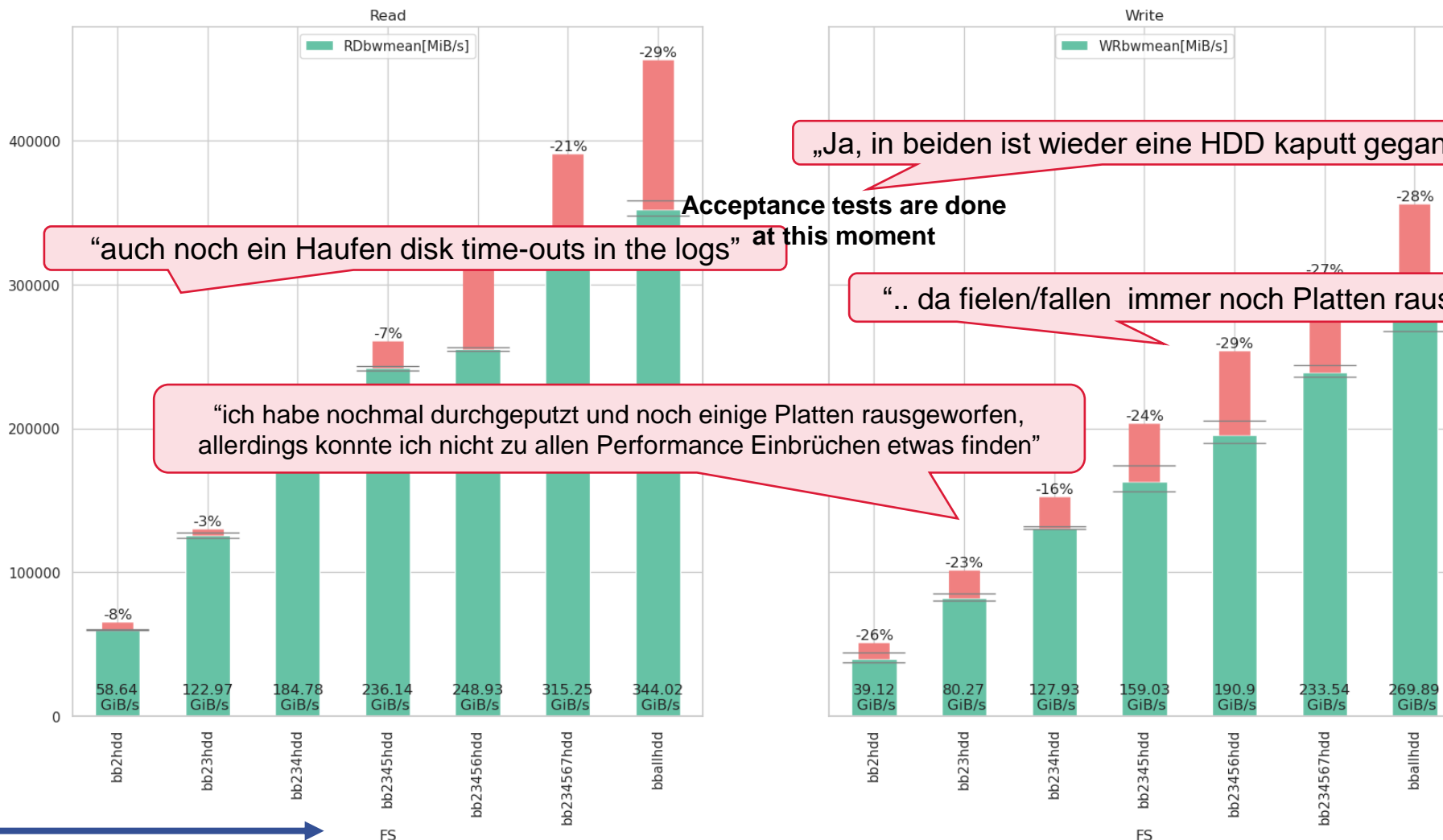
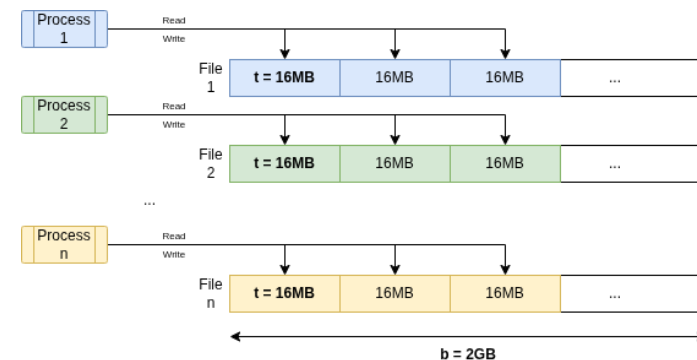
Warning:

You might not be able to keep track of all the details in following tests. We couldn't either and we had more than 20min

IOR EASY - JUST6 PHASE 1

7 x SSS6000 Building Blocks

IOR easy scale Building Blocks



„Ja, in beiden ist wieder eine HDD kaputt gegangen.“

“auch noch ein Haufen disk time-outs in the logs” at this moment

“.. da fielen/fallen immer noch Platten raus”

“ich habe nochmal durchgeputzt und noch einige Platten rausgeworfen, allerdings konnte ich nicht zu allen Performance Einbrüchen etwas finden”

Acceptance tests are done

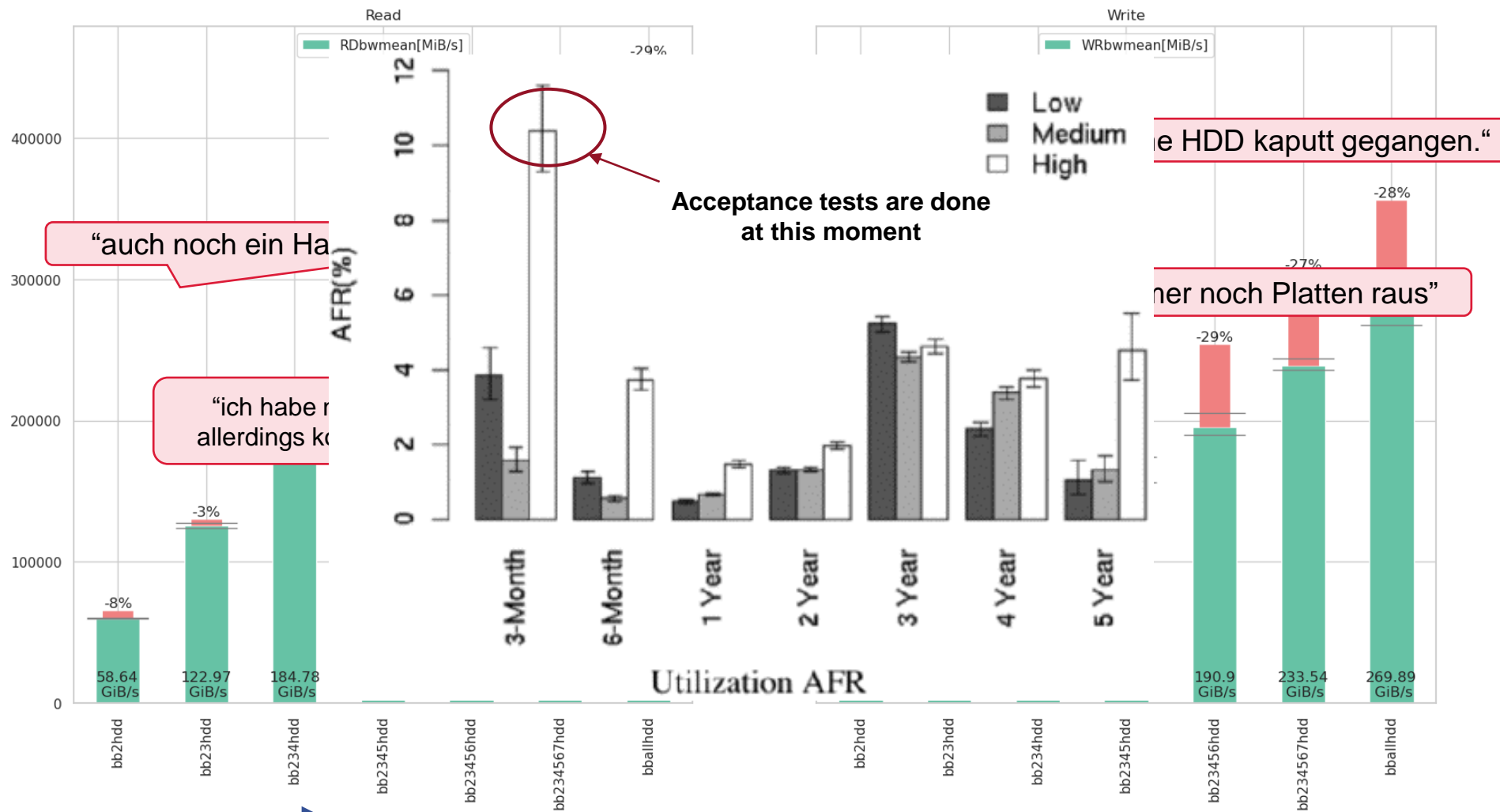
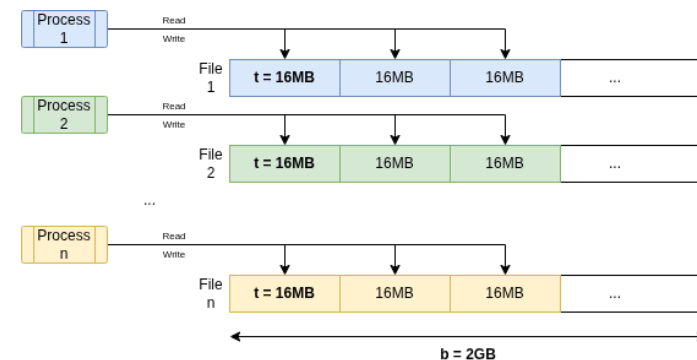
Burn IN

Increasing number of building blocks

IOR EASY - JUST6 PHASE 1

7 x SSS6000 Building Blocks

IOR easy scale Building Blocks



Burn
IN

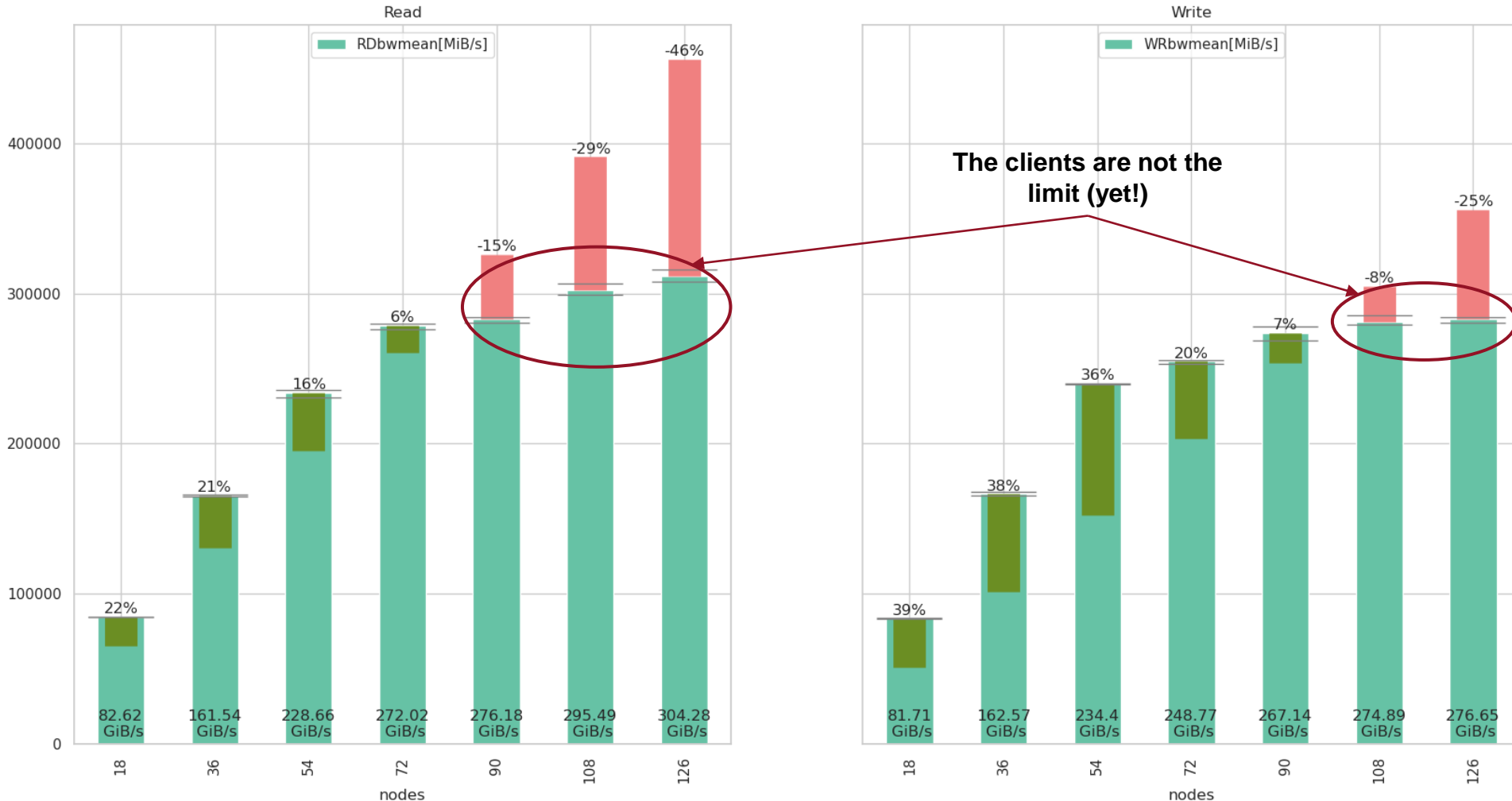
E. Pinheiro, et al. "Failure trends in a large disk drive population". 2007. In Proceedings of the 5th USENIX conference on File and Storage Technologies (FAST '07). USENIX Association, USA, 2.

IOR EASY – JUST6 PHASE 1

Don't forget the clients - 7 x SSS6000 Building Blocks

IOR easy scale Clients

Originally planned for 18 nodes per Building Block (126 for 7)

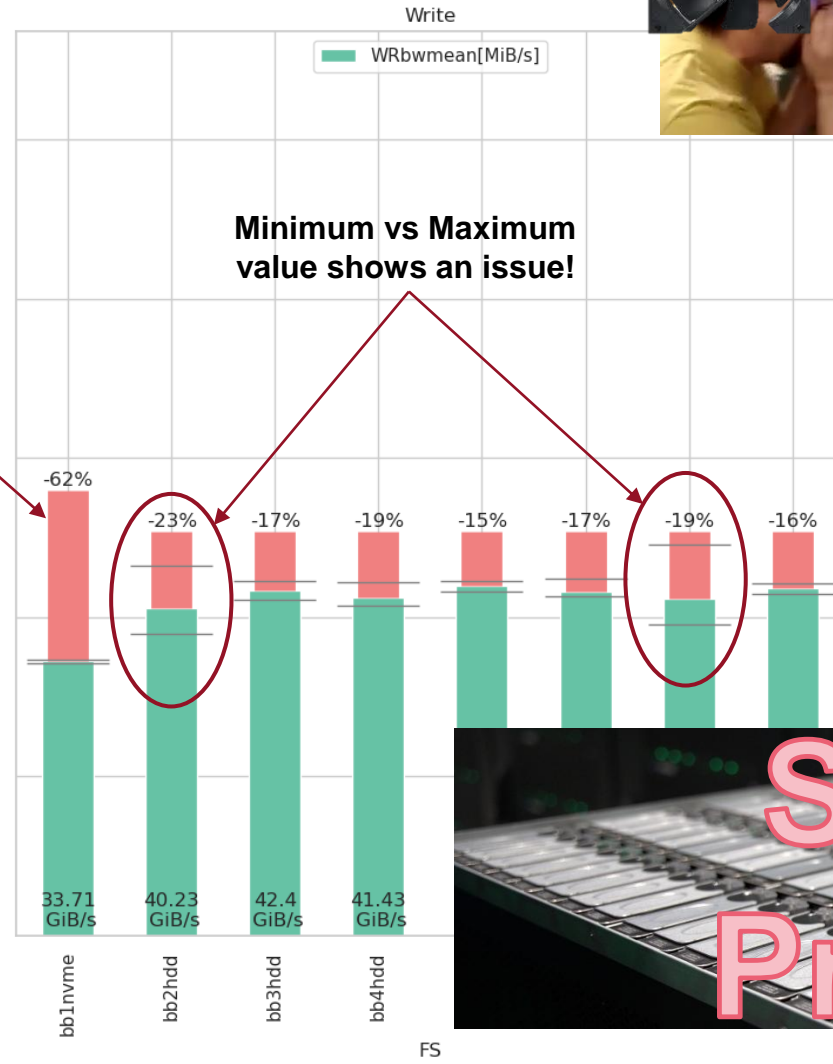
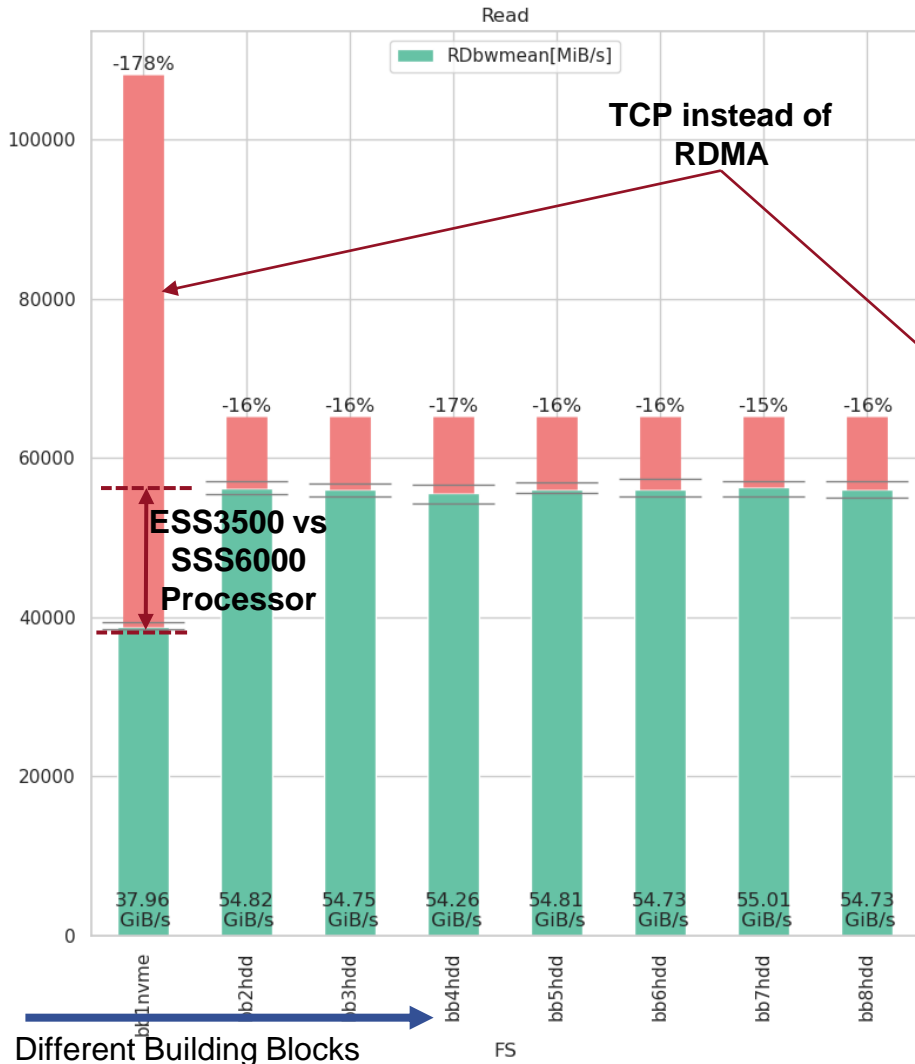


Increasing number of client nodes

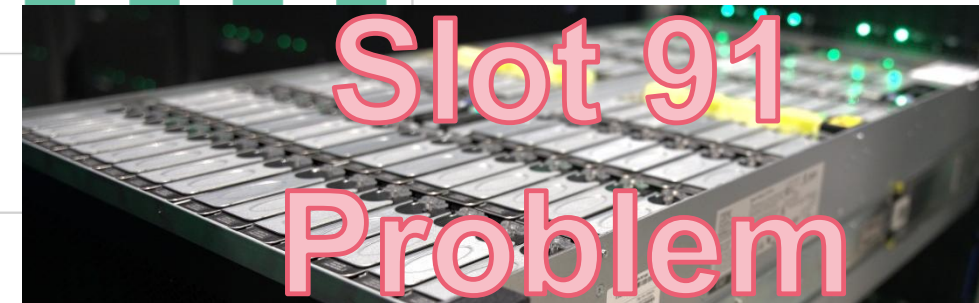
IOR EASY – JUST6 PHASE 1

Let's test individual BB - 7 x SSS6000 Building Blocks

IOR easy per Building Block



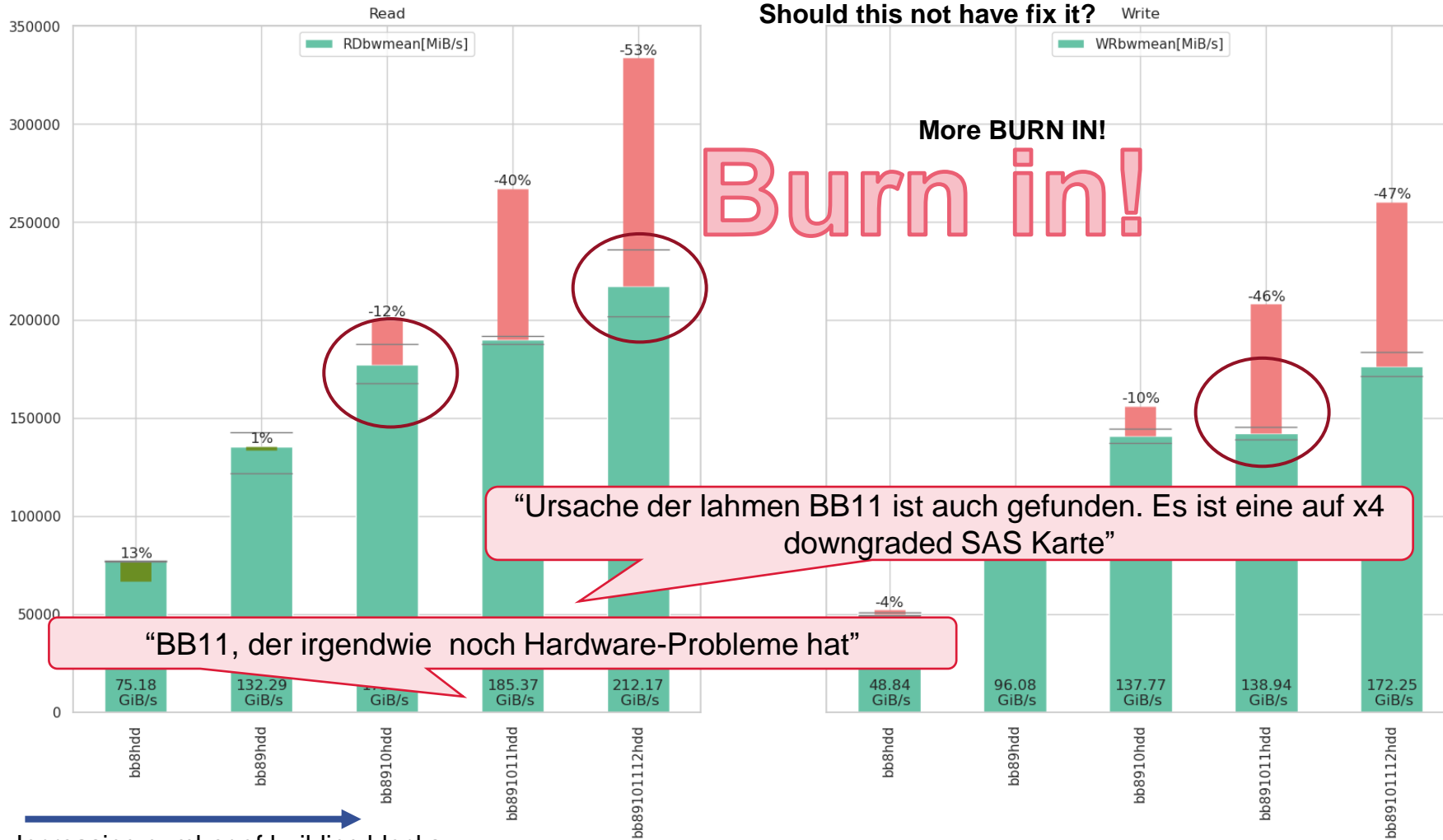
Fix with Firmware
Phase 1 goes into production before fix can be deployed!



IOR EASY – JUST6 PHASE 2

5 x SSS6000 Building Blocks

IOR easy scale Building Blocks PHASE2 Firmware update

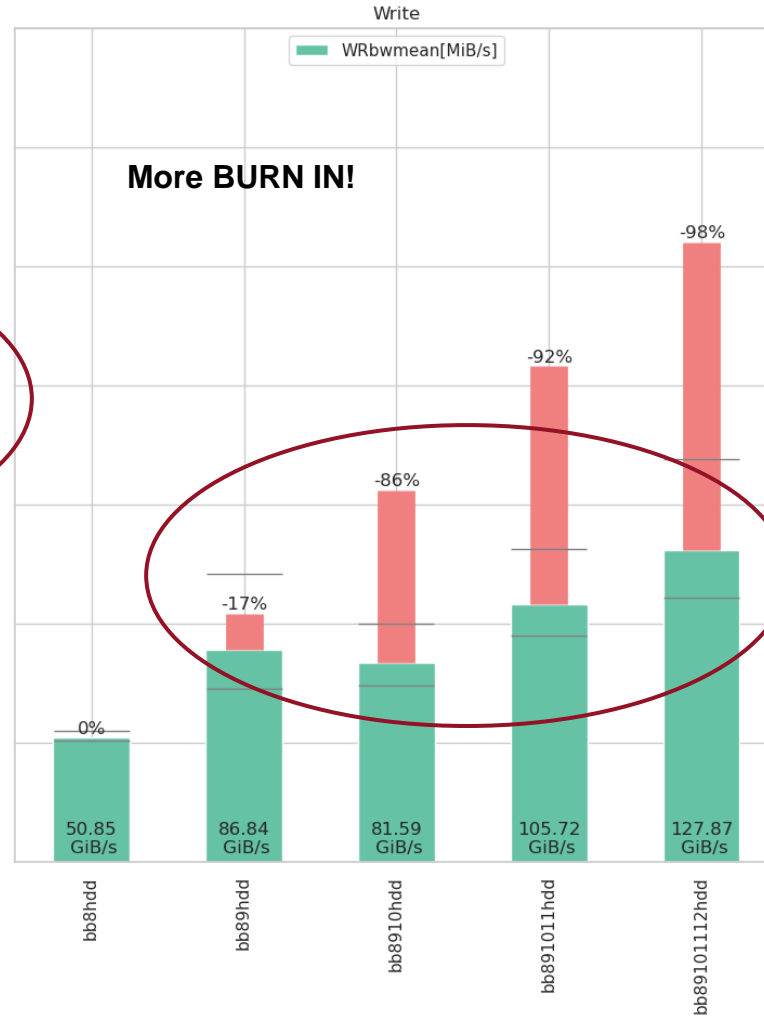
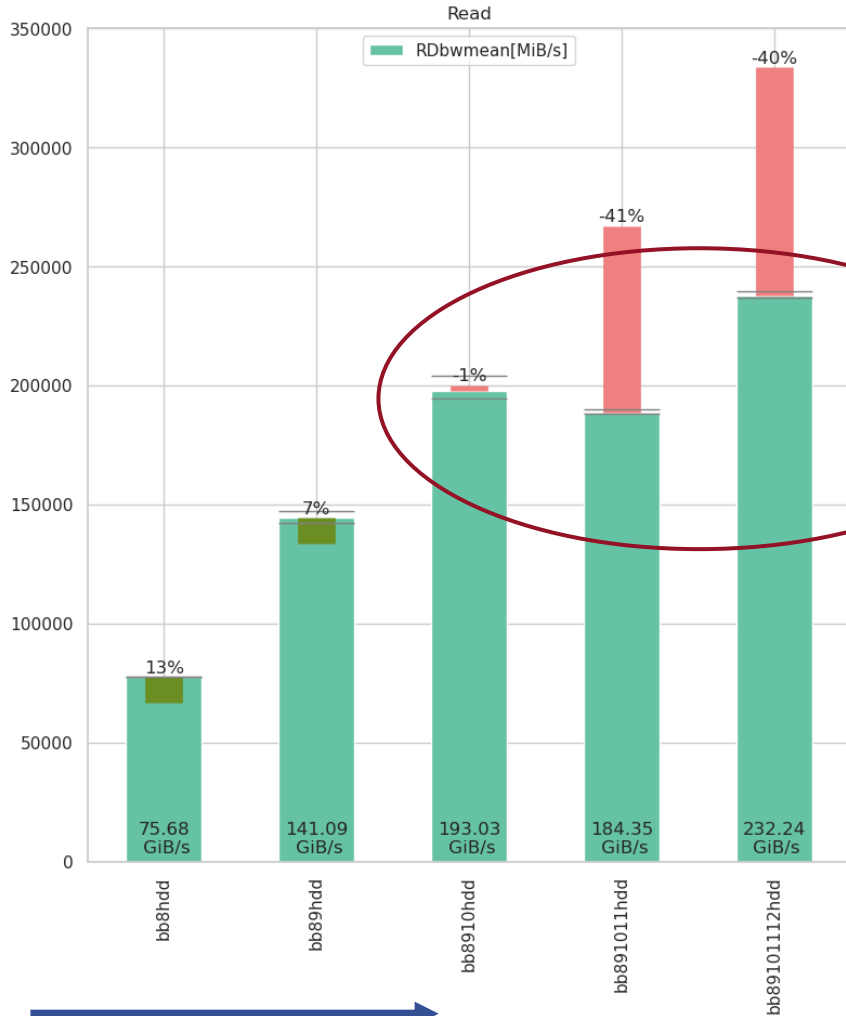


IOR EASY – JUST6 PHASE 2

5 x SSS6000 Building Blocks

IOR easy scale Building Blocks PHASE2 PagePool 32GiB

Let's see if the pagepool can help
(Only for debugging)



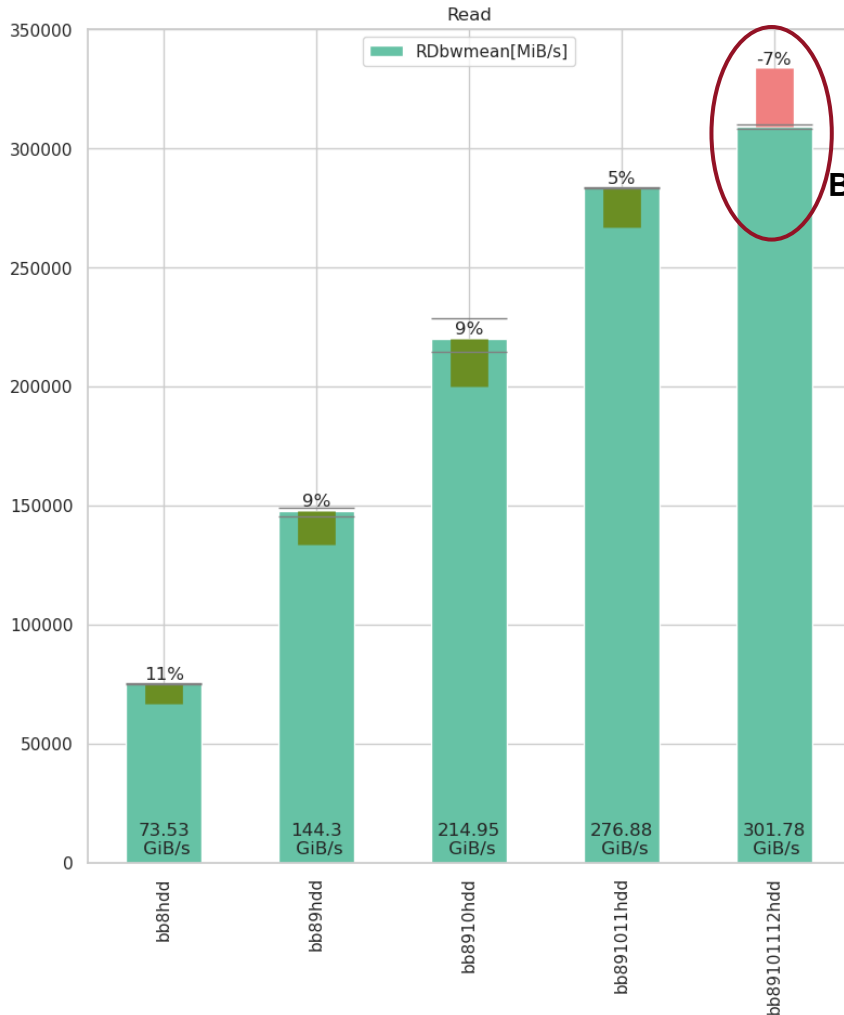
Increasing number of building blocks →

IOR EASY – JUST6 PHASE 2

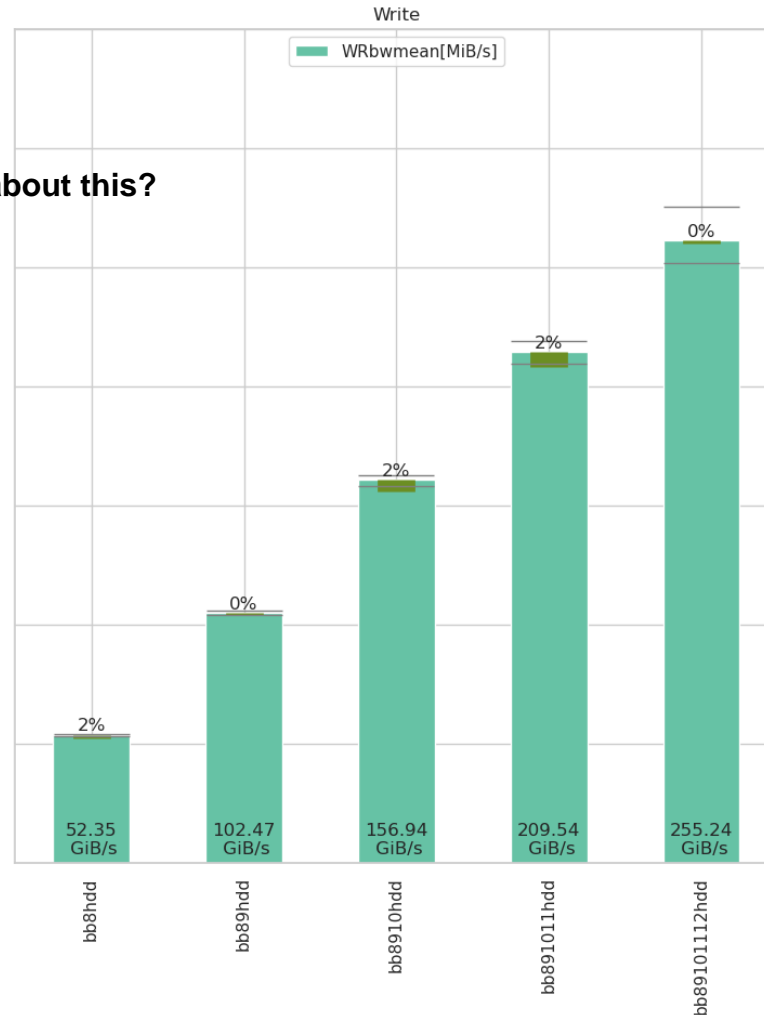
5 x SSS6000 Building Blocks

IOR easy scale Building Blocks PHASE2 126 clients

Now we are getting somewhere



But what about this?

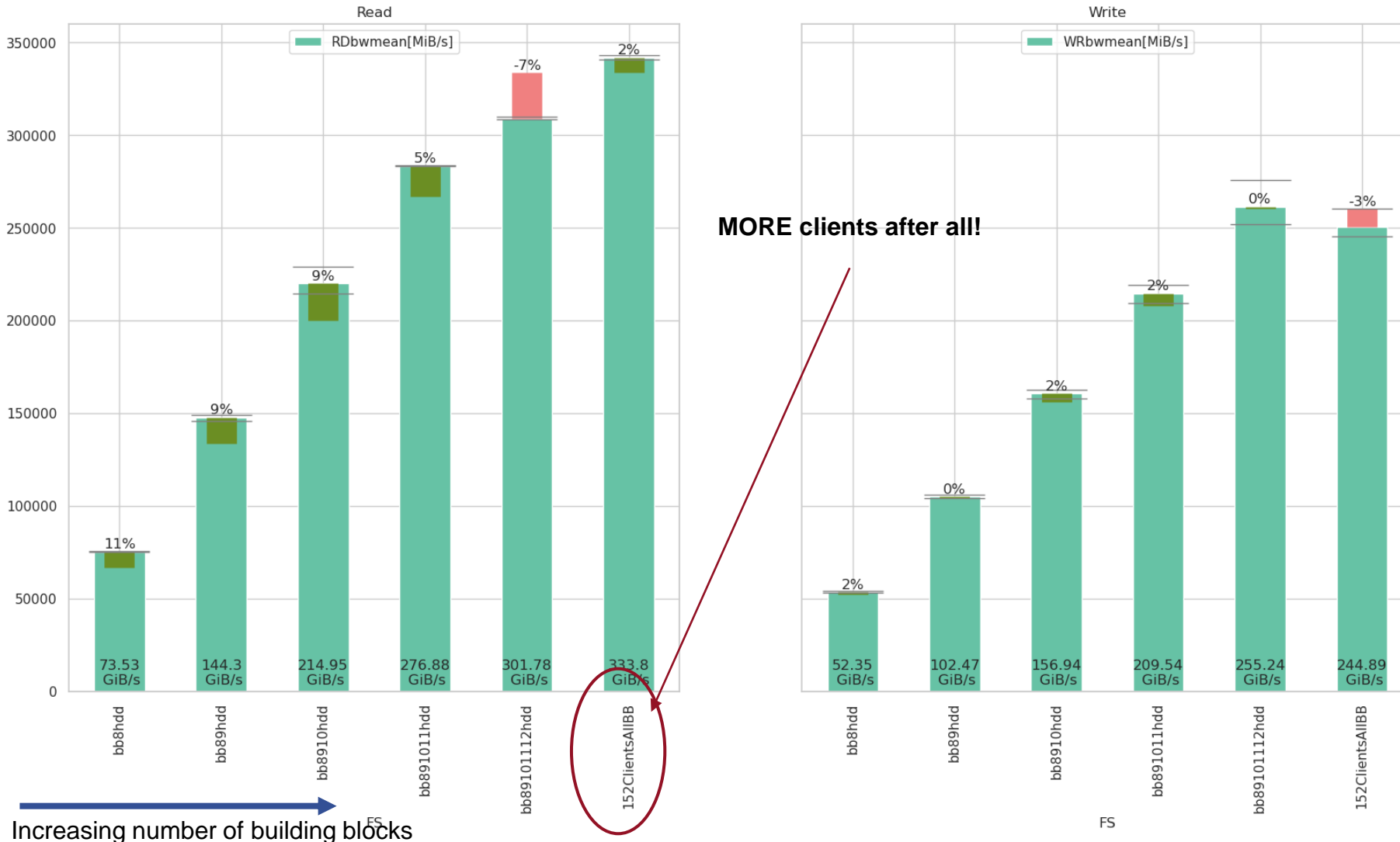


Increasing number of building blocks

IOR EASY – JUST6 PHASE 2

5 x SSS6000 Building Blocks

IOR easy scale Building Blocks PHASE2 126-152 clients

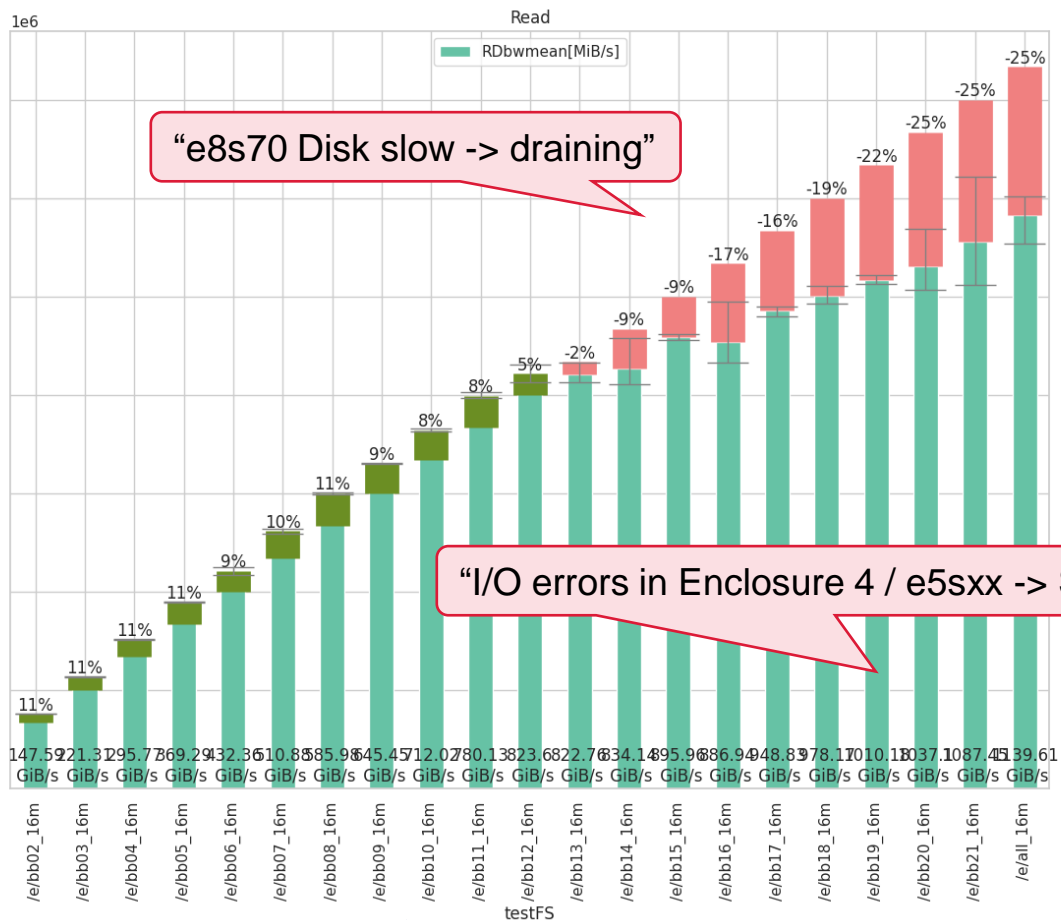


IOR EASY – EXASTORE

22 x SSS6000 Building Blocks

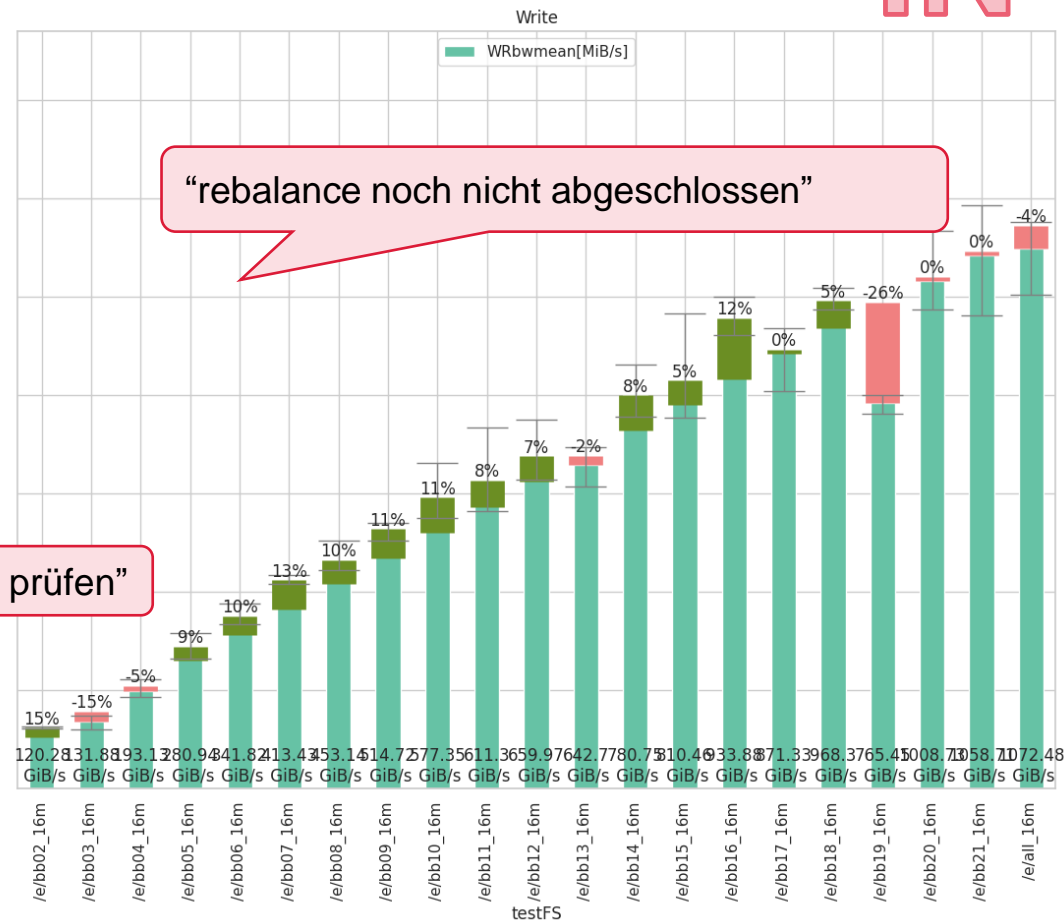
Burn
IN

20250627_ExaSTORE_easy_BB_scale



“e8s70 Disk slow -> draining”

“I/O errors in Enclosure 4 / e5sxx -> SAS Kabel prüfen”



“rebalance noch nicht abgeschlossen”

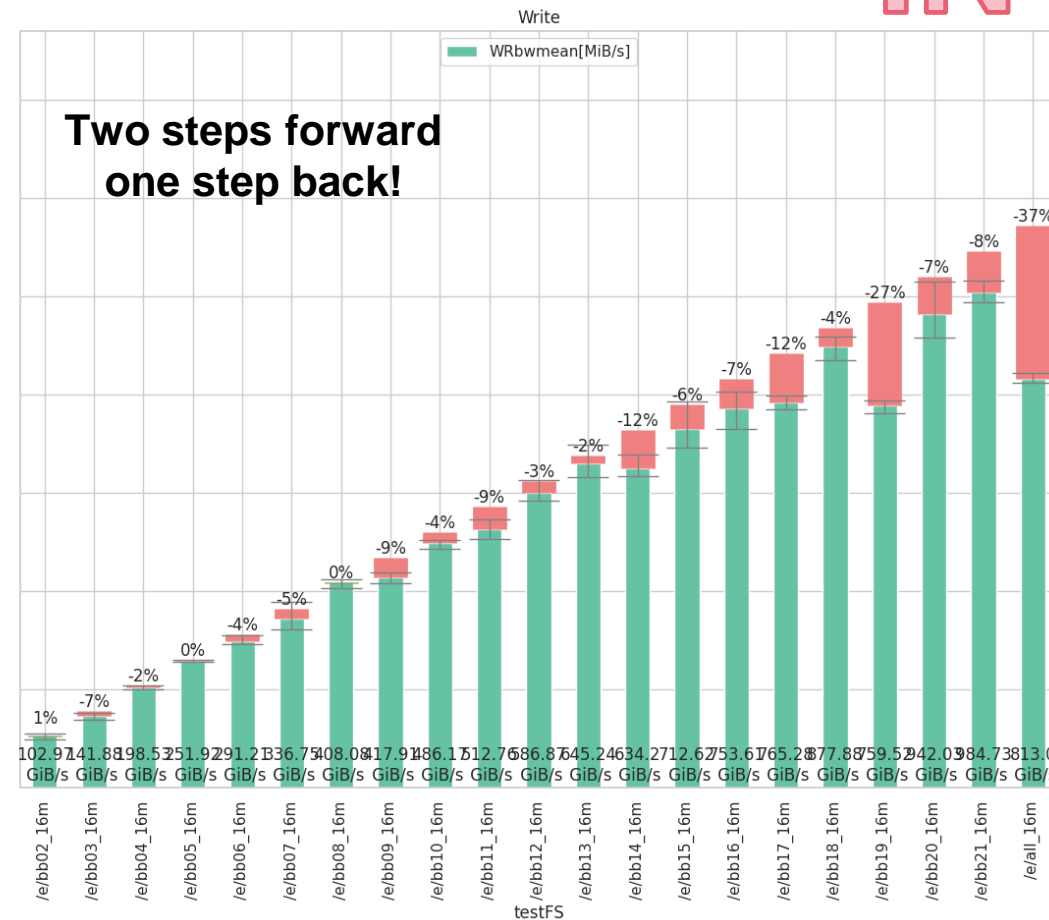
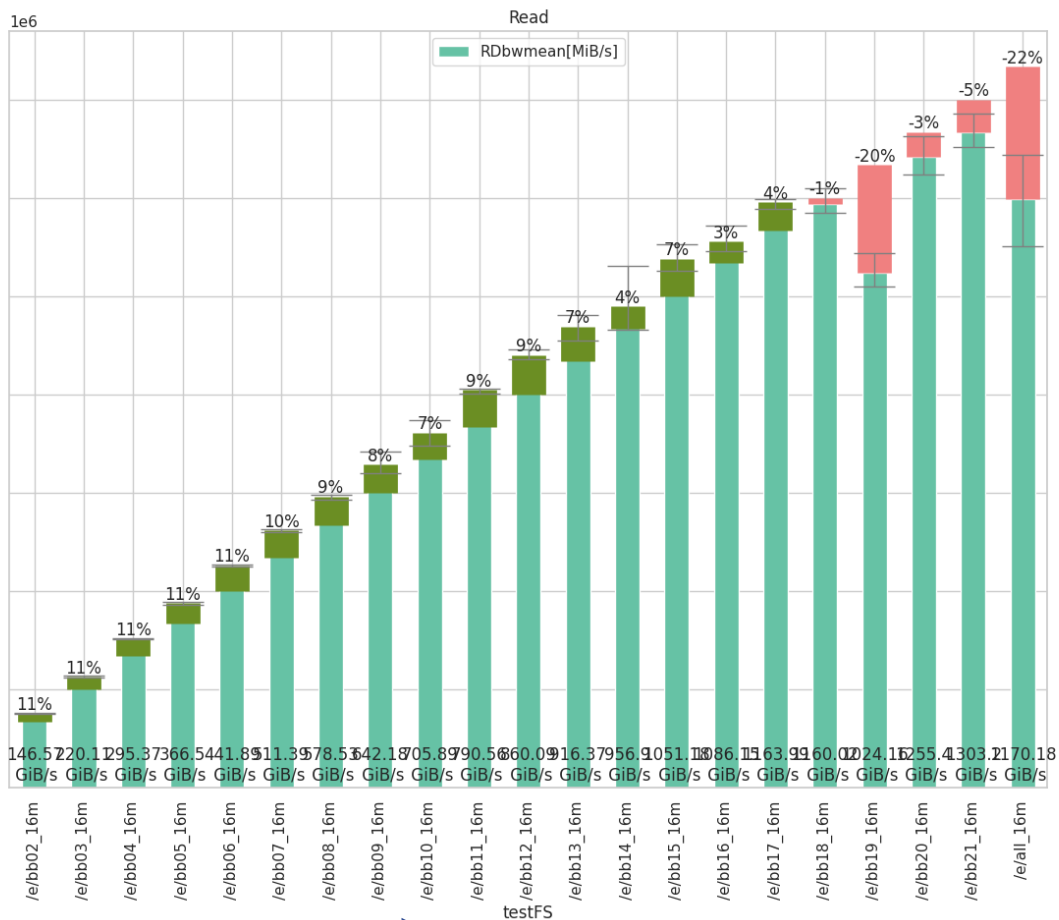
Increasing number of Building Blocks

IOR EASY – EXASTORE

22 x SSS6000 Building Blocks

20250703_ExaSTORE_easy_BB_scale

Burn
IN



Increasing number of Building Blocks

Mitglied der Helmholtz-Gemeinschaft

IOR EASY – EXASTORE

Better clients, more problems - 22 x SSS6000 Building Blocks

„ ... irgendein HPC node hat da seinen mount verloren“

„Wieder mounts verloren“

„ bb04 sah sehr gut aus, zumindest bis ein Knoten die mounts verloren hat, `“

• IPoIB Problems:

1. Compute expelling Storage nodes (and vice versa):

Temp solutions:

Ping all compute and storage nodes before mounting

Increase ARP cache and extend Garbage collector thresholds

Increase timeout till nodes are expelled (`failureDetectionTime`)

Final solution: Hope the stabilization of the IB network goes well

2. Compute expelling compute:

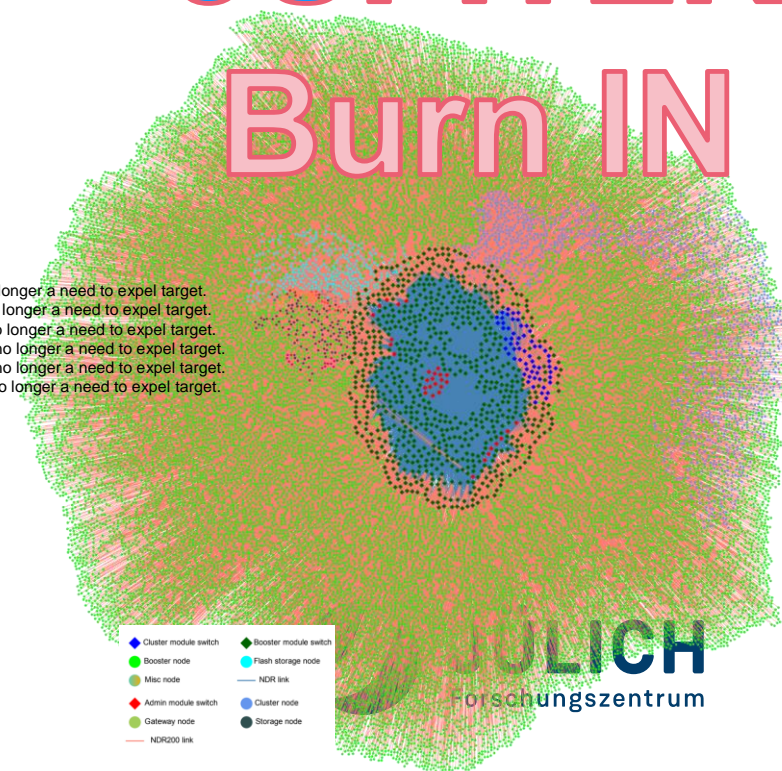
[N] Expel 10.128.16.8 (jdbo-001-08-interconnect-1.jupiter.internal in j request from 10.128.19.137 (jdbo-019-41-interconnect-1.jupiter.internal in j. Expelling 10.128.16.8 (jdbo-001-08-interconnect-1.jupiter.internal in j based on no longer a need to expel target.
[N] Expel 10.128.18.84 (jdbo-013-20-interconnect-1.jupiter.internal in j request from 10.128.16.91 (jdbo-002-43-interconnect-1.jupiter.internal in j. Expelling 10.128.18.84 (jdbo-013-20-interconnect-1.jupiter.internal in j based on no longer a need to expel target.
[N] Expel 10.128.18.90 (jdbo-013-26-interconnect-1.jupiter.internal in j request from 10.128.17.196 (jdbo-010-20-interconnect-1.jupiter.internal in j. Expelling 10.128.18.90 (jdbo-013-26-interconnect-1.jupiter.internal in j based on no longer a need to expel target.
[N] Expel 10.128.16.116 (jdbo-003-20-interconnect-1.jupiter.internal in j request from 10.128.26.162 (jdbo-035-02-interconnect-1.jupiter.internal in j. Expelling 10.128.16.116 (jdbo-003-20-interconnect-1.jupiter.internal in j based on no longer a need to expel target.
[N] Expel 10.128.18.104 (jdbo-013-40-interconnect-1.jupiter.internal in j request from 10.128.16.230 (jdbo-005-38-interconnect-1.jupiter.internal in j. Expelling 10.128.18.104 (jdbo-013-40-interconnect-1.jupiter.internal in j based on no longer a need to expel target.
[N] Expel 10.128.18.106 (jdbo-013-42-interconnect-1.jupiter.internal in j request from 10.128.18.59 (jdbo-012-43-interconnect-1.jupiter.internal in j. Expelling 10.128.18.106 (jdbo-013-42-interconnect-1.jupiter.internal in j based on no longer a need to expel target.
...

Temp solution:

Avoid inter client node communication (`preferDesignatedMnode=yes`)

Final solution: Hope the stabilization of the IB network goes well

JUPITER Burn IN



Cluster module switch
Booster node
Misc node
Admin module switch
Gateway node
NDR200 link
Booster module switch
Flash storage node
NDR link
Cluster node
Storage node

IOR EASY - EXASTORE

Too many clients for Scale? - 22 x SSS6000 Building Blocks

Remember:

5884 Nodes with NVIDIA Mellanox NDR 4x NDR200 NICs

- **Problem:** At 3-4k nodes client manager hits (ENOMEM – not enough memory)
Managers start to cycle

- **Temporary Solution:**

- `maxReceiverThreads=128` [quorum]
- `maxTcpConnsPerNodeConn=1`
- `proactiveReconnect=no`
- `verbsPorts=mlx5_0/1 ~~mlx5_1/1~~ ~~mlx5_2/1~~ ~~mlx5_3/1`~~

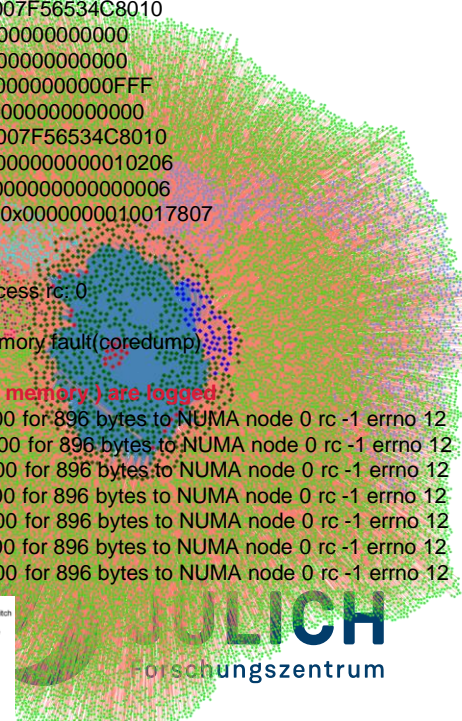
- **Final Solution (Untested):**

- `PF_LOG_BAR_SIZE=10` (Recommended by Nvidia)

```
2025-07-21_13:11:38.476+0200: [E] Signal 11 at location 0x55E345A1F3BB in process 1125817, link reg 0xFFFFFFFF
2025-07-21_13:11:38.476+0200: [I] rax 0x0000000000000000 rbx 0x0000000000000000
2025-07-21_13:11:38.476+0200: [I] rcx 0x0000000000000000 rdx 0x0000000000000000
2025-07-21_13:11:38.476+0200: [I] rsp 0x00007F54D0ACFAE0 rbp 0x00007F56534C8010
2025-07-21_13:11:38.476+0200: [I] rsi 0x0000000000100000 rdi 0x0000000000000000
2025-07-21_13:11:38.476+0200: [I] r8 0x00000000FFFFFFFF r9 0x0000000000000000
2025-07-21_13:11:38.476+0200: [I] r10 0x000055E349B1F000 r11 0x00000000000000FF
2025-07-21_13:11:38.476+0200: [I] r12 0x00007F5653568448 r13 0x0000000000000000
2025-07-21_13:11:38.476+0200: [I] r14 0x00007F54D0ACFBBC r15 0x00007F56534C8010
2025-07-21_13:11:38.476+0200: [I] rip 0x000055E345A1F3BB eflags 0x0000000000010206
2025-07-21_13:11:38.476+0200: [I] cs:gsfs 0x002B000000000033 err 0x0000000000000006
2025-07-21_13:11:38.476+0200: [I] trapno 0x000000000000000E oldmsk 0x0000000010017807
2025-07-21_13:11:38.476+0200: [I] cr2 0x0000000000000020
2025-07-21_13:11:38.480+0200: [N] Starting mmsdrserv: enter
2025-07-21_13:11:38.480+0200: [N] Starting mmsdrserv: started child process rc: 0
malloc(131072) for operator new() failed.
/usr/lpp/mmf/bin/runmmfs[494]: numactlWrapper: line 9824: 1125817: Memory fault(coredump)
no further stack information.
```

before that a lot of NUMA warning with err_12 (ENOMEM - not enough memory) are logged:

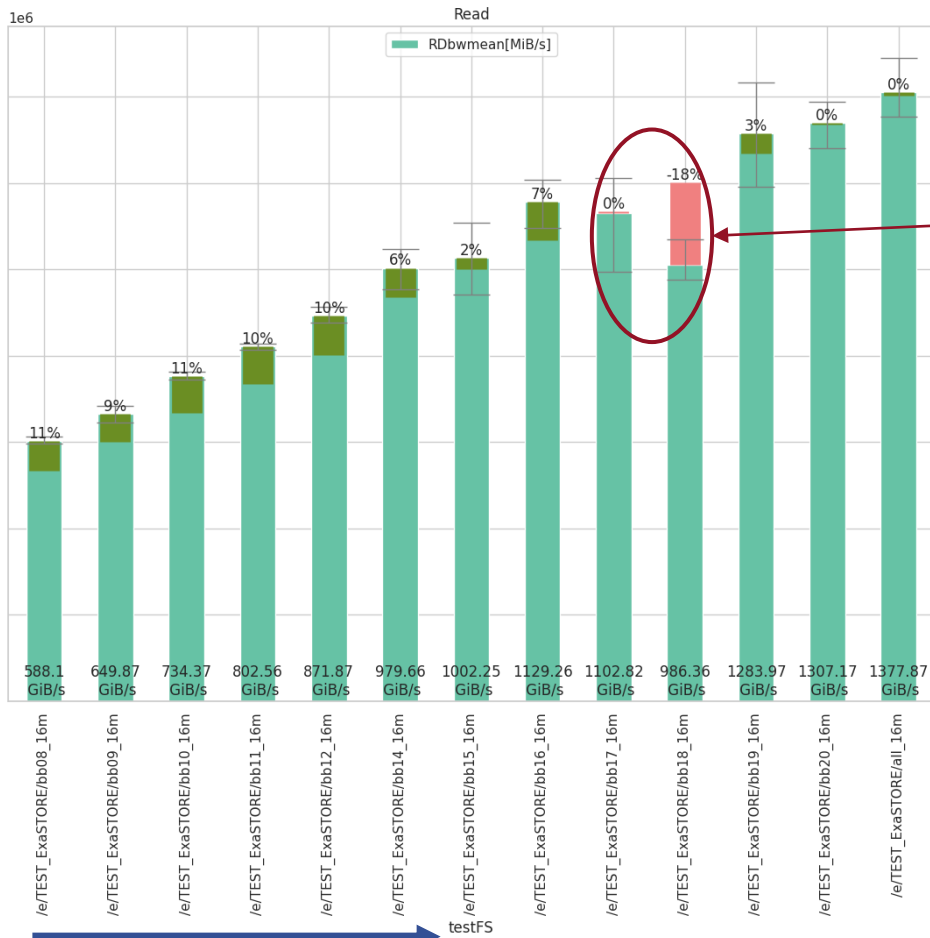
```
2025-07-21_13:11:26.734+0200: [W] NUMA failed to bind 0x7F48A1477000 for 896 bytes to NUMA node 0 rc -1 errno 12
2025-07-21_13:11:26.742+0200: [W] NUMA failed to bind 0x7F48A14AC000 for 896 bytes to NUMA node 0 rc -1 errno 12
2025-07-21_13:11:26.750+0200: [W] NUMA failed to bind 0x7F48A14E1000 for 896 bytes to NUMA node 0 rc -1 errno 12
2025-07-21_13:11:26.759+0200: [W] NUMA failed to bind 0x7F48A1516000 for 896 bytes to NUMA node 0 rc -1 errno 12
2025-07-21_13:11:27.288+0200: [W] NUMA failed to bind 0x7F48A154B000 for 896 bytes to NUMA node 0 rc -1 errno 12
2025-07-21_13:11:27.289+0200: [W] NUMA failed to bind 0x7F48A1580000 for 896 bytes to NUMA node 0 rc -1 errno 12
2025-07-21_13:11:27.291+0200: [W] NUMA failed to bind 0x7F48A15B5000 for 896 bytes to NUMA node 0 rc -1 errno 12
...
```



IOR EASY - EXASTORE

21 x SSS6000 Building Blocks (1 BB already in production)

Rev_202508291218_ExaSTORE_easy_BB_scale_INCOMPLETE



Despite this, we made it

Increasing number of Building Blocks
Mitglied der Helmholtz-Gemeinschaft

ACCEPTANCE TEST LIST – HOW FAR ARE WE?

✓ IOR Easy

✓ Each Building Block

✓ IOR Easy 16MB blocksize file system

IOR Easy 512KB blocksize file system

✓ All Building Blocks

✓ IOR Easy 16MB blocksize file system

IOR Easy 512KB blocksize file system

IOR Hard

Each Building Block

IOR Hard 16MB blocksize file system

IOR Hard 512KB blocksize file system

All Building Blocks

IOR Hard 16MB blocksize file system

IOR Hard 512KB blocksize file system

MDTTest

Each Building Block

MDTTest 16MB block size file system

MDTTest 512KB block size file system

All Building Block

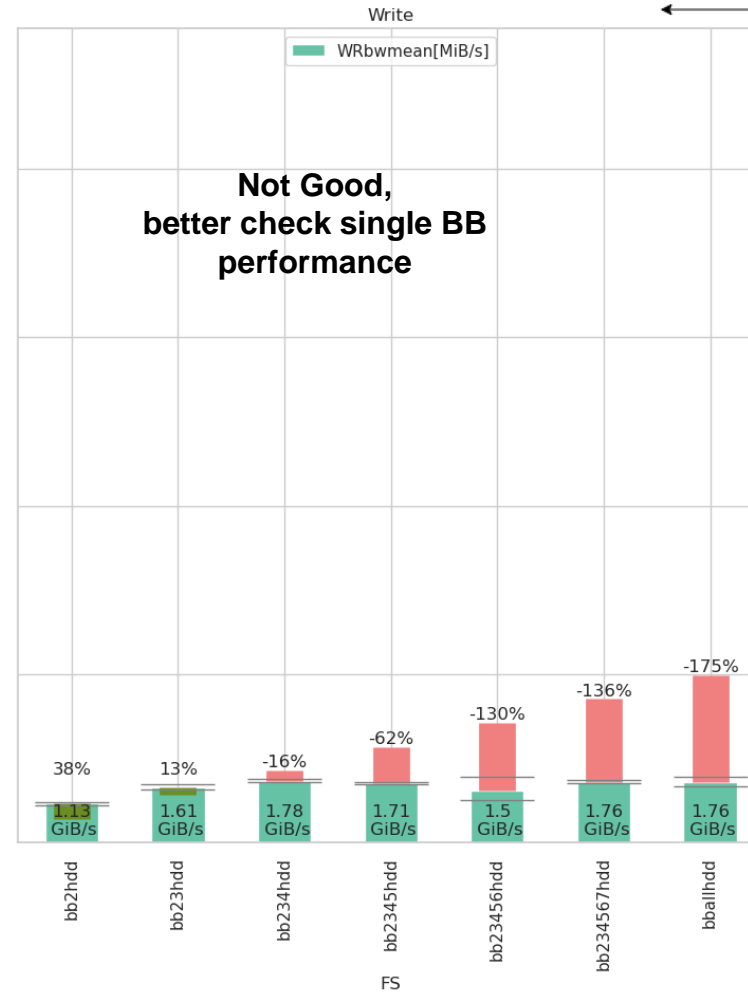
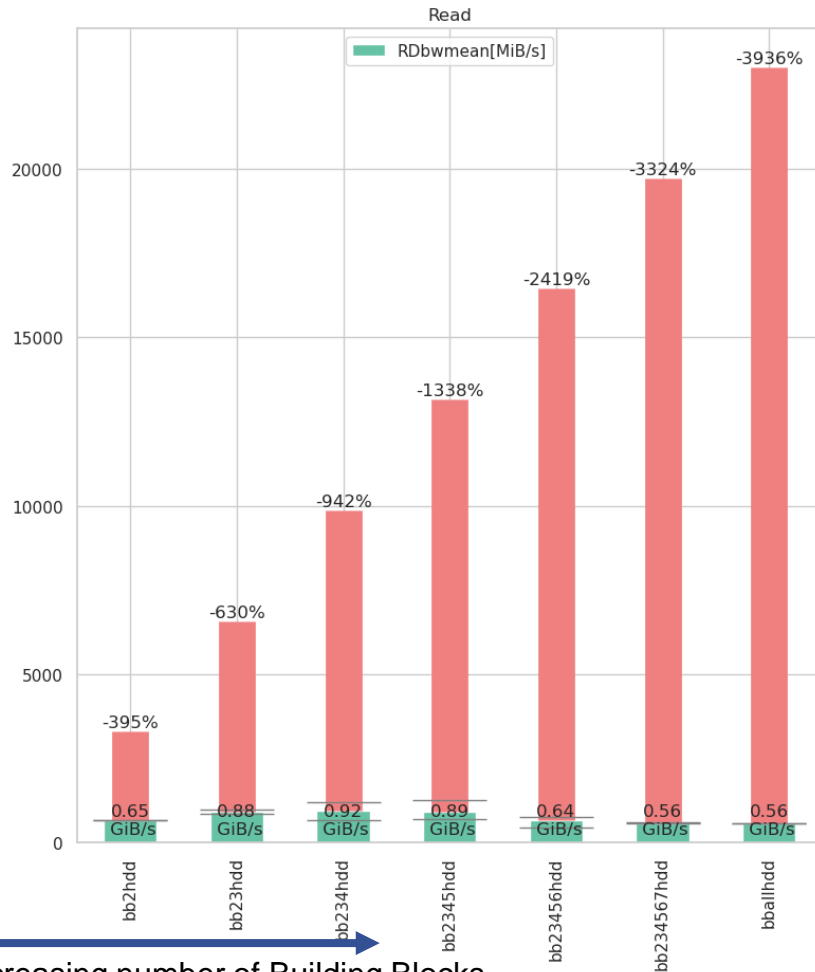
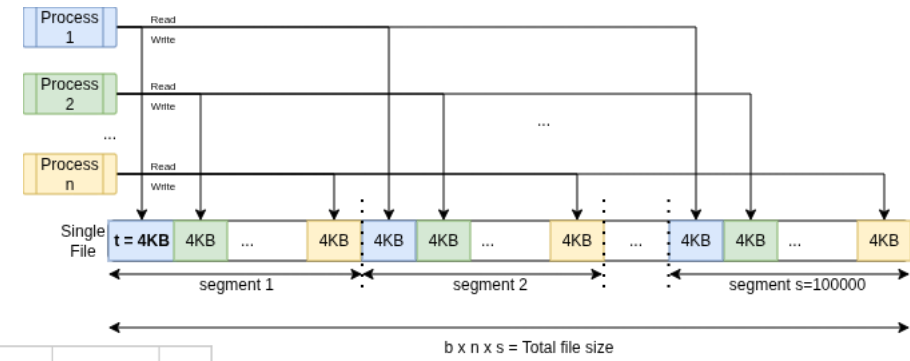
MDTTest 16MB block size file system

MDTTest 512KB block size file system

IOR HARD - JUST6 PHASE 1

7 x SSS6000 Building Blocks

IOR hard scale Building Blocks

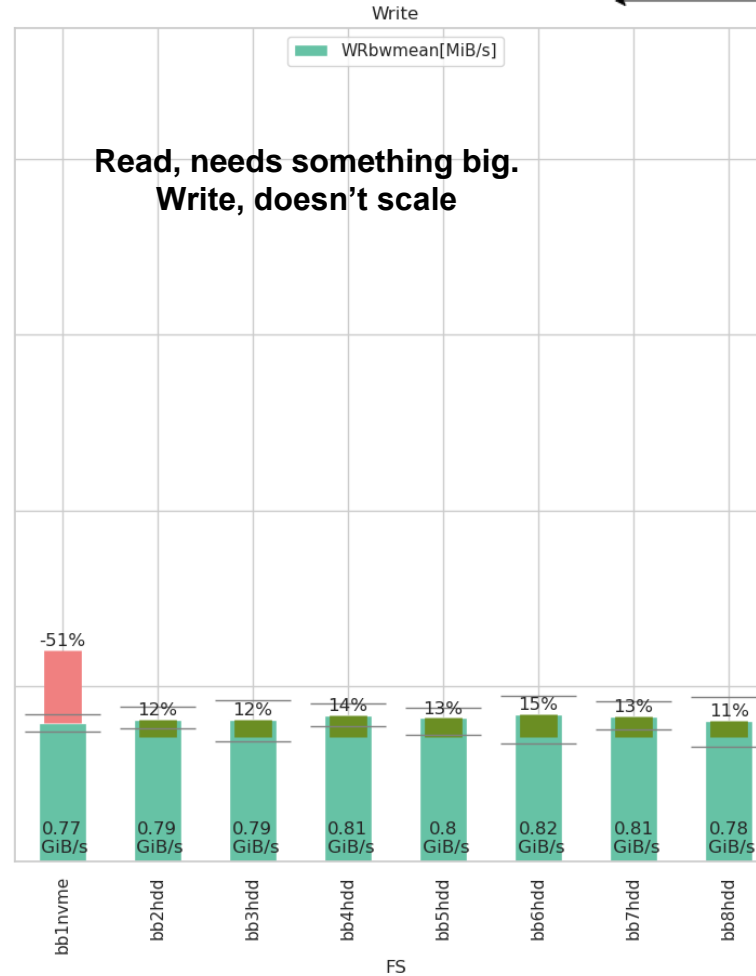
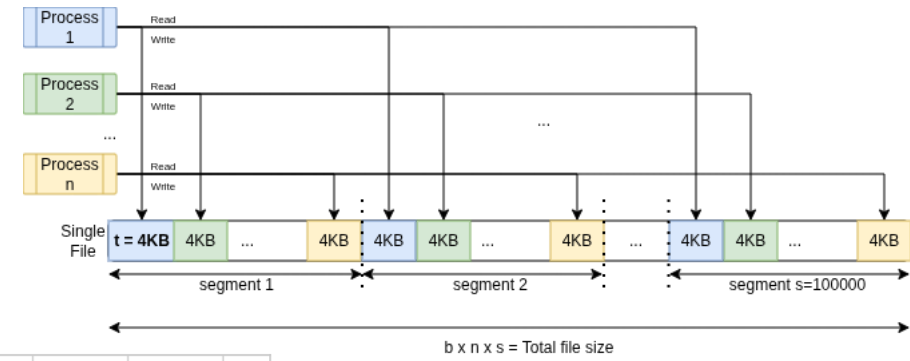


Increasing number of Building Blocks
 Mitglied der Heimloitz-Gemeinschaft

IOR HARD - JUST6 PHASE 1

Single BB Tests - 7 x SSS6000 Building Blocks

IOR hard per Building Block



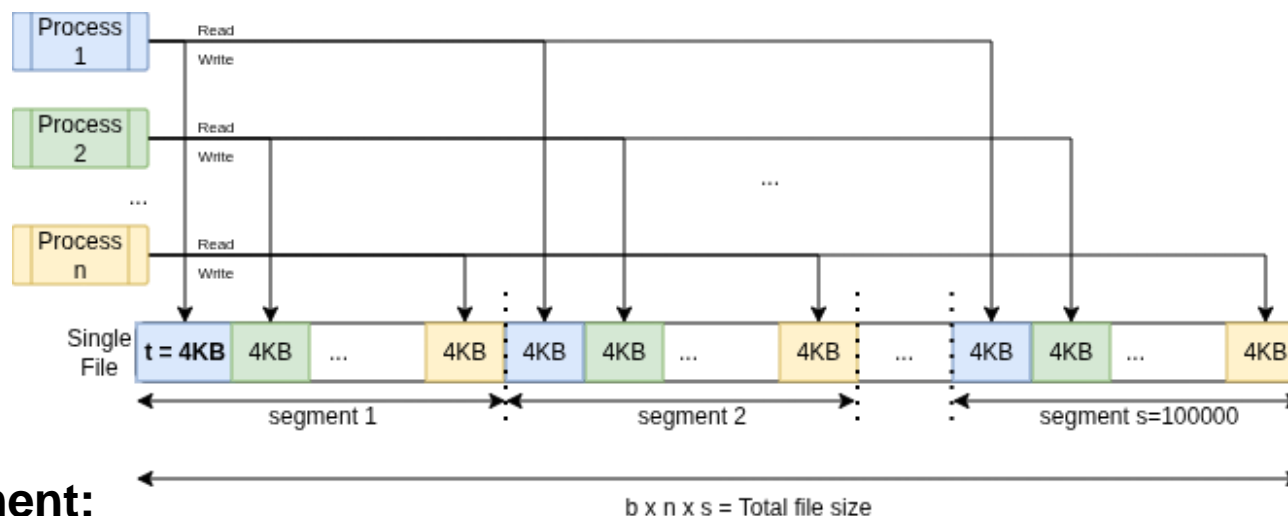
Read, needs something big.
Write, doesn't scale

Software issue that requires:

- `finegrainwritesharing` and `finegrainreadshraing`.
- Many many many parameter changes
- Additional efix and updates to GPFS

IOR HARD - JUST6 PHASE 1

Why it is harder than you think? - GPFS Tokens



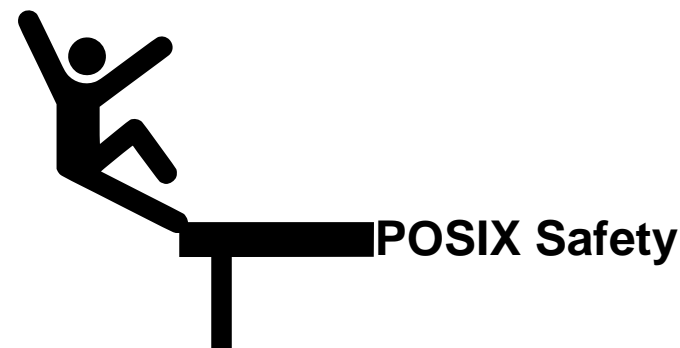
GPFS implements a distributed token management:

- Each Node (GPFS daemon) has to acquire a token for its regions
- The acquisition requires contacting token server or a client node that holds the token for this region

REQUIRED TO BE POSIX COMPLIANT!

Solution, let's not be:

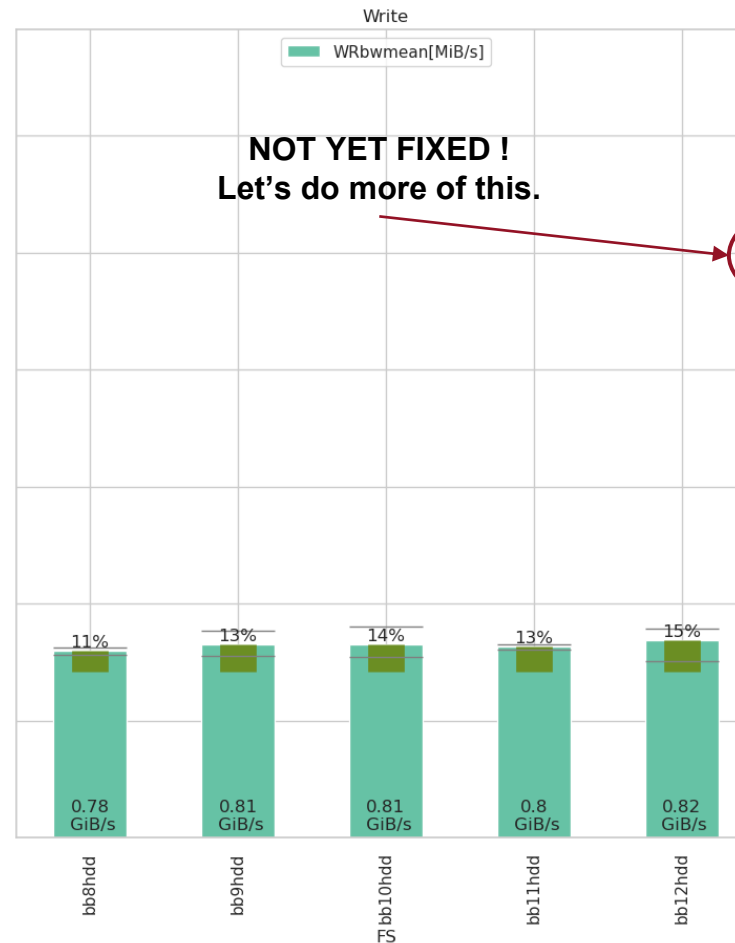
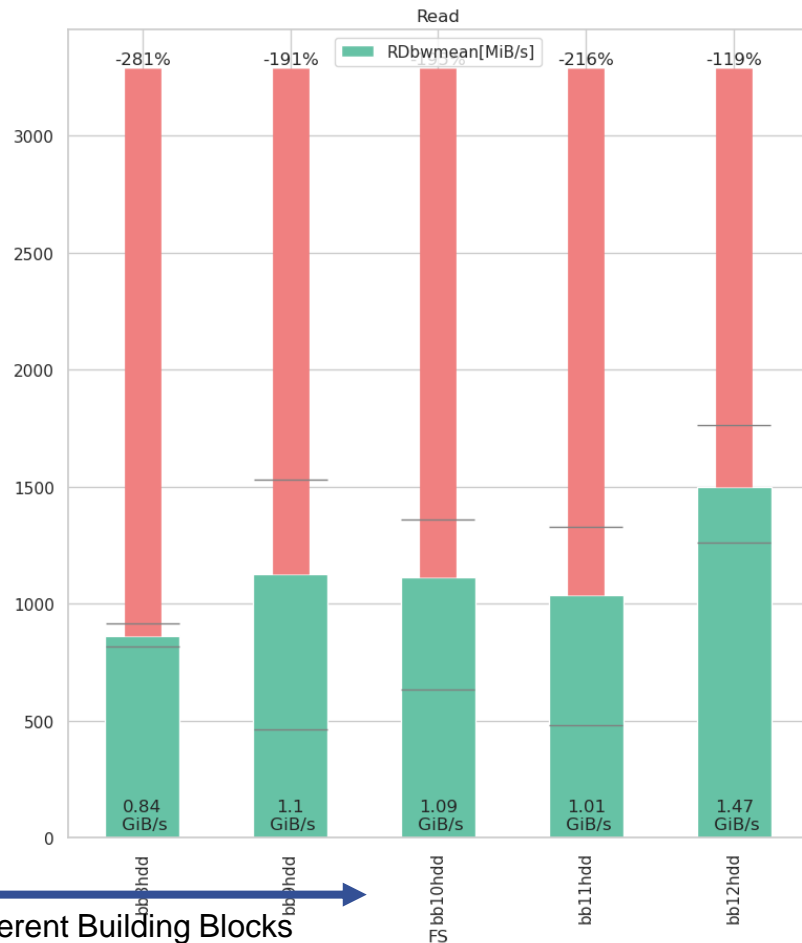
- ``finegrainwritesharing`` and ``finegrainreadshraing`` IOR
= No token management for that operation
(I know what I'm doing)



IOR HARD - JUST6 PHASE 2

Single BB Tests - 5 x SSS6000 Building Blocks

IOR hard per Building Block Phase2



Software issue that requires:

- `finegrainwritesharing` and `finegrainreadshraing`.
- Many many many parameter changes
- Additional efix and updates to GPFS

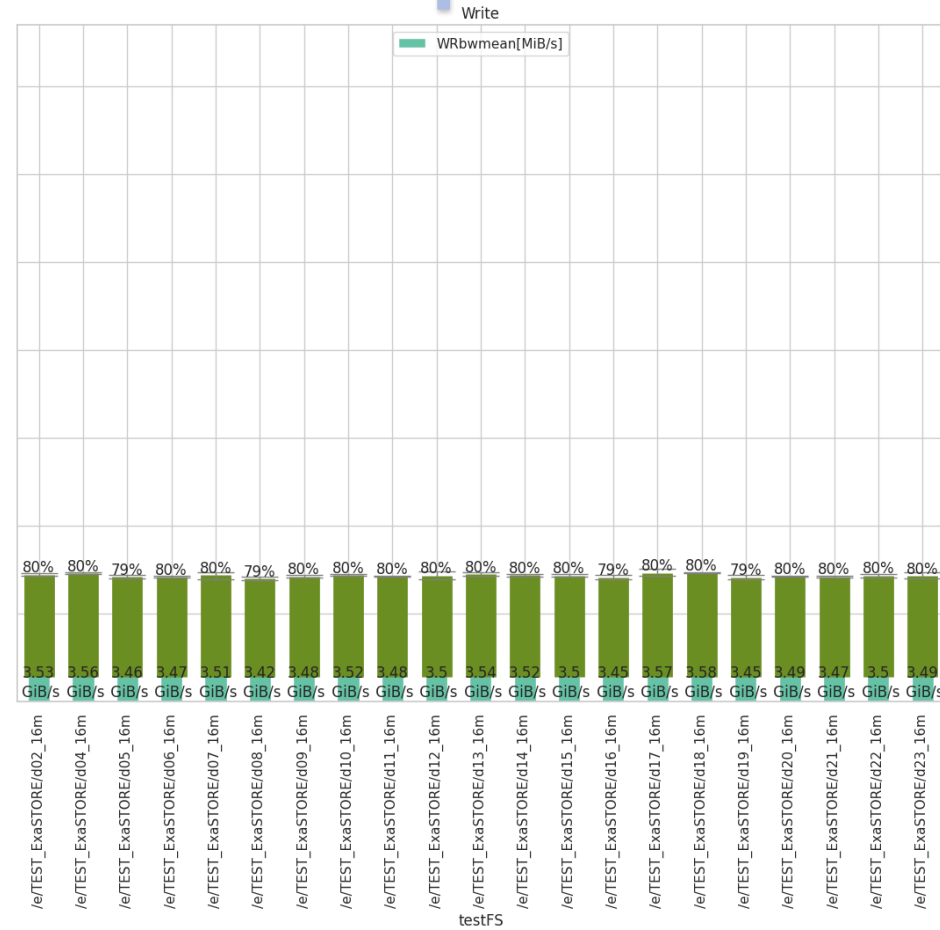
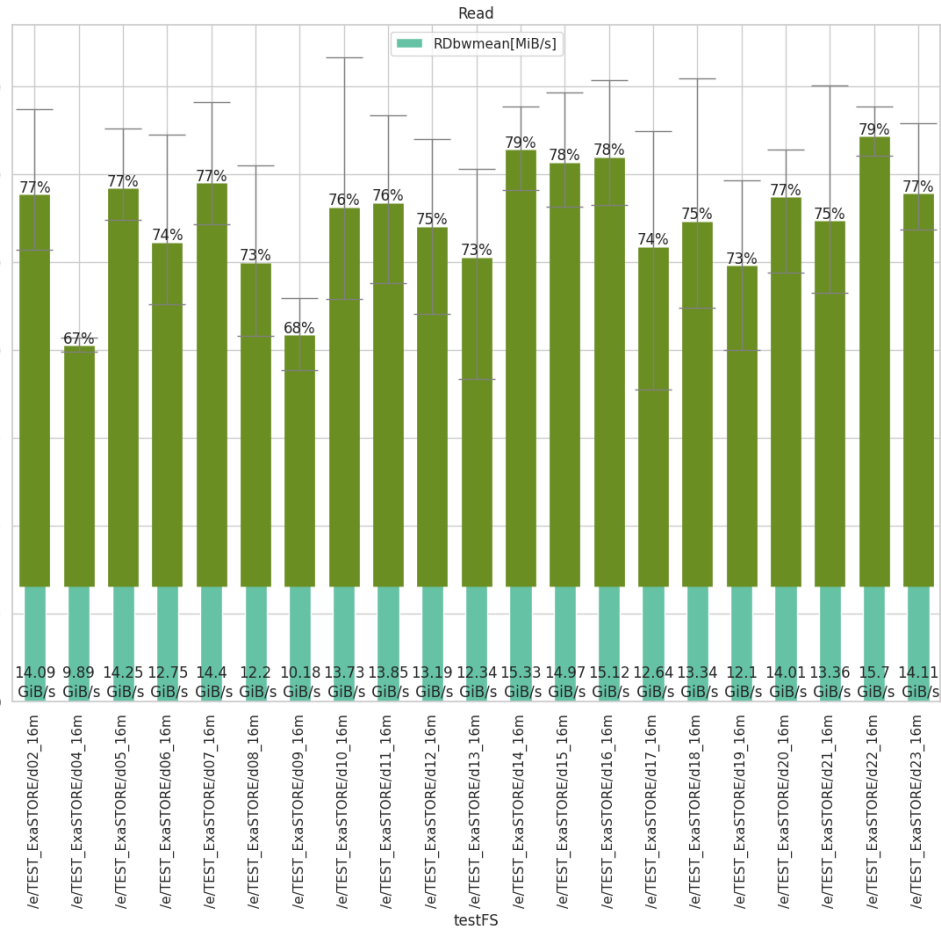
Let's revisit in **ExaSTORE**

IOR HARD - EXASTORE

22 x SSS6000 Building Blocks

20250625 IOR hard per Building Block

Impressive!





 Different Building Blocks

 Mitglied der Helmholtz-Gemeinschaft

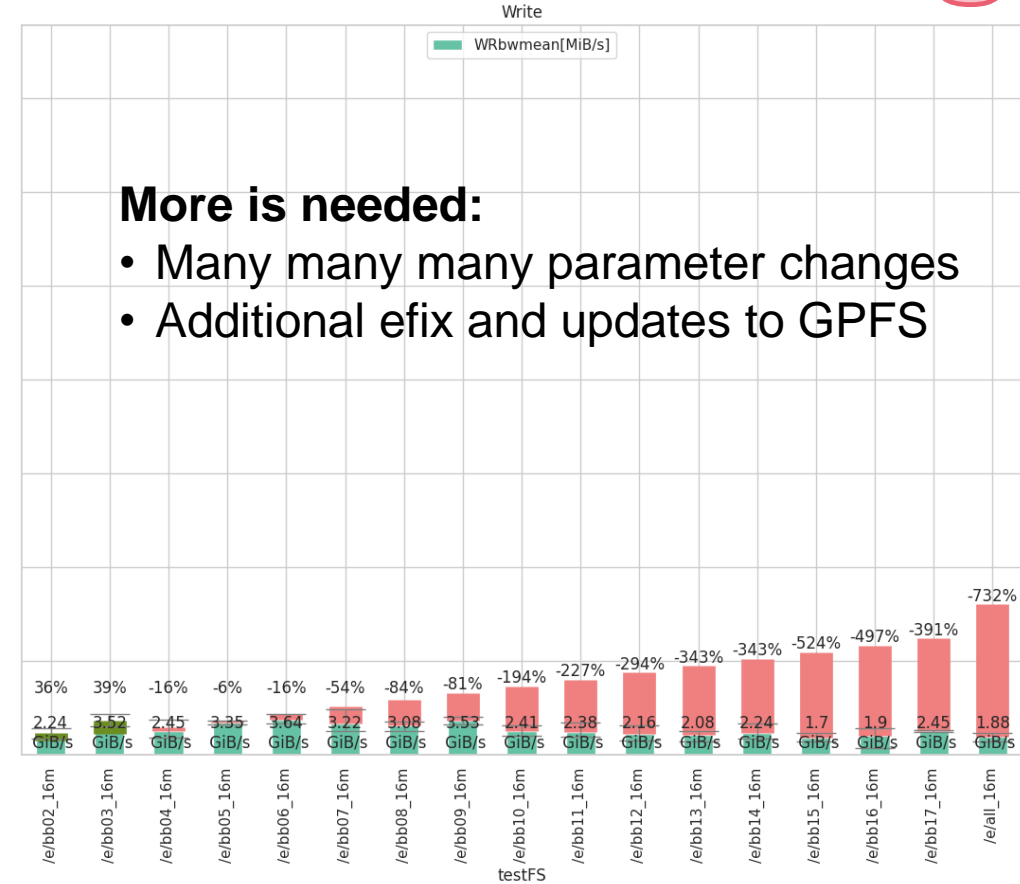
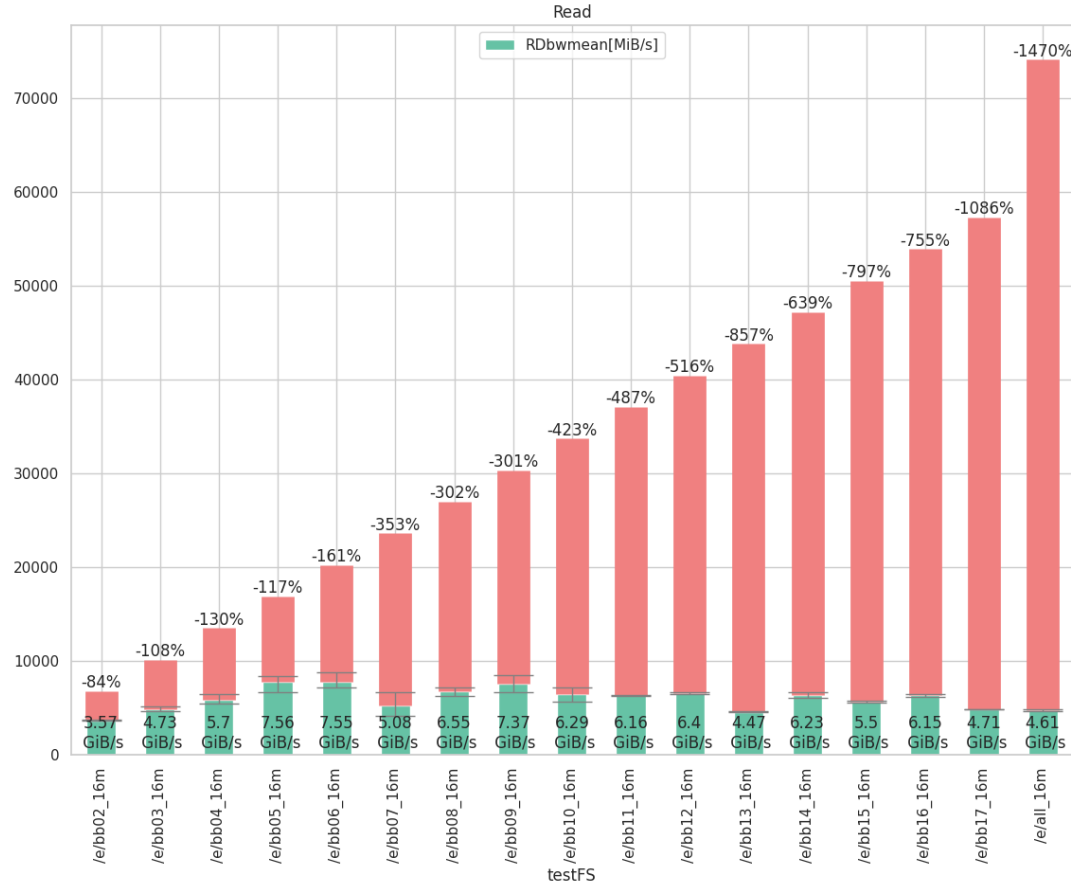
But does it scale?

IOR HARD - EXASTORE

Scale test - 22 x SSS6000 Building Blocks

20250629_IOR hard scale Building Blocks INCOMPLETE

Not enough



More is needed:

- Many many many parameter changes
- Additional efix and updates to GPFS

➔ Different Building Blocks

Mitglied der Helmholtz-Gemeinschaft

Did it work?

IOR HARD – EXASTORE

22 x SSS6000 Building Blocks

Still working on it

Cliff hanger!

PRESENT ENDEAVORS

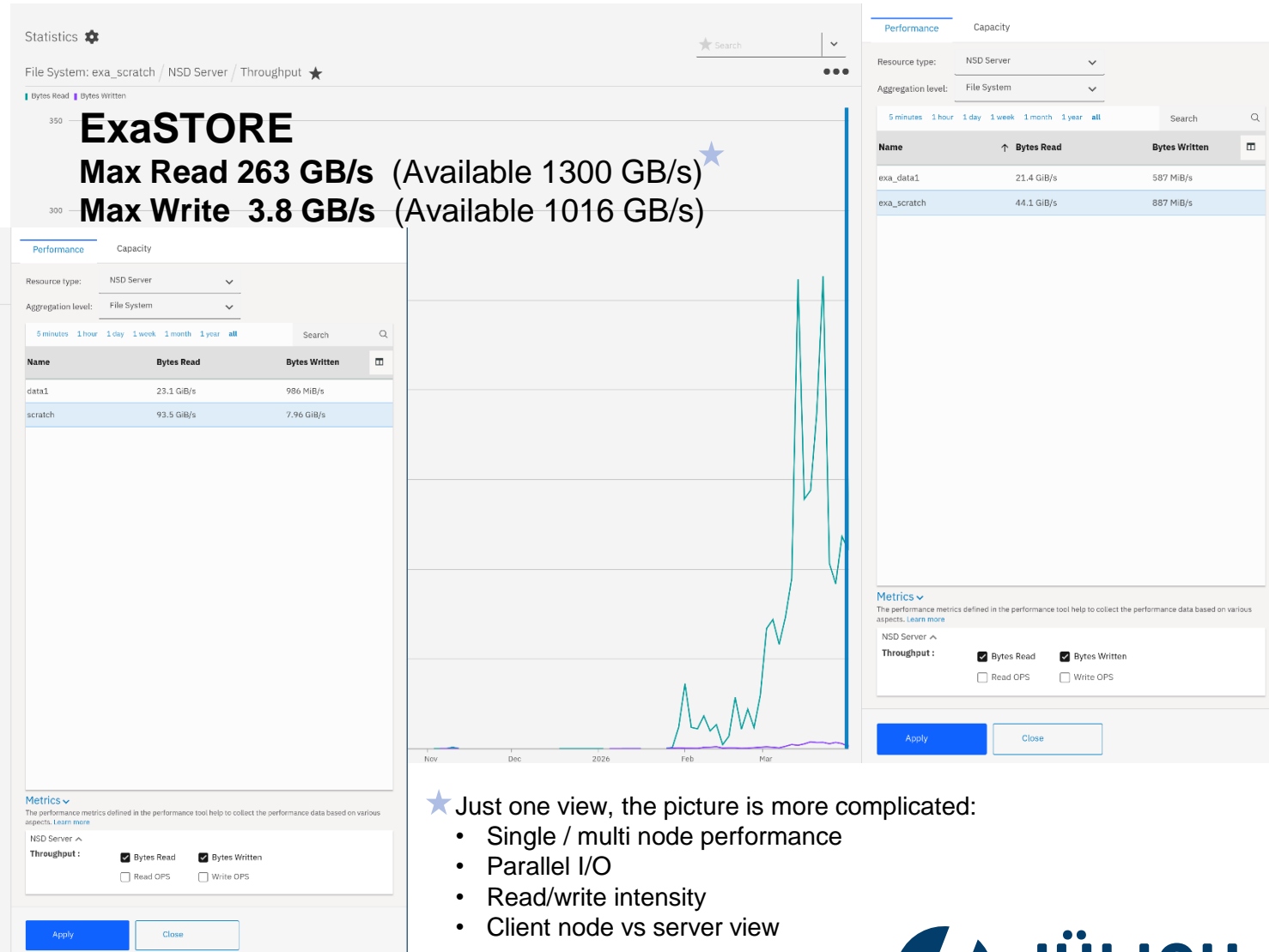
- IOR Hard scaling tests
- IO500 on ExaSTORE
- ExaFLASH Acceptance tests
- IO500 on ExaFLASH
- Continuous Benchmarking of I/O performance
- ...

BUT !

Who cares about I/O?

ARE USERS DOING I/O?

I/O reported by GPFS



★ Just one view, the picture is more complicated:

- Single / multi node performance
- Parallel I/O
- Read/write intensity
- Client node vs server view
- ...

AN OBSERVED MISMATCH

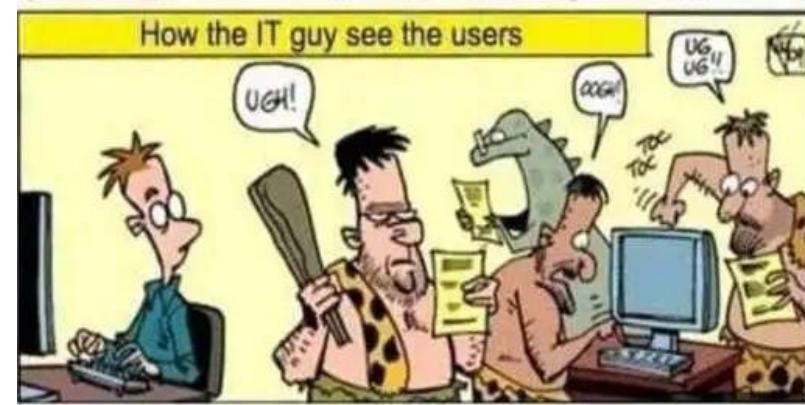
USERS

- Regularly claim they need more I/O performance



Admins

- Monitoring shows little I/O performed



WHICH IS IT?

**Or better,
how do we find the truth?**

REMEMBER - PROCUREMENT BENCHMARK

A. Herten et al., "Application-Driven Exascale: The JUPITER Benchmark Suite," SC24: International Conference for High Performance Computing, Networking, Storage and Analysis, Atlanta, GA, USA, 2024, pp. 1-45, doi: 10.1109/SC41406.2024.00038.

Benchmark	Domain	Dense LA	Sparse LA	Spectral	Particle	Structured Grid	Unstructured Grid	Monte Carlo
Amber*	MD			●	●			
Arbor	Neurosci.	●	●				●	
Chroma-QCD	QCD		●			●		●
GROMACS	MD			●	●			
ICON	Climate		●			●		
JUQCS	QC		●					
nekRS	CFD		●	●			●	
ParFlow*	Earth Sys.		●			●		
PICongPU	Plasma				●		●	
Quantum Espresso	Materials Sci.	●		●	●			
SOMA*	Polymer Sys.				●			●
MMoCLIP	AI (MM)	●						
Megatron-LM	AI (LLM)	●						
ResNet*	AI (Vision)	●						
DynQCD	QCD		●			●		●
NAStJA	Biology					●		●
Graph500	Graph							
HPCG	CG		●					
HPL	LA	●						
IOR	Filesys.							
LinkTest	Network							
OSU	Network							
STREAM	Memory							

Just one side of the story.

The only I/O benchmark and it's synthetic!



IT'S A DISCUSSION !

Do we need to understand the users requirement?

Another,
Cliff hanger!

THANKS !

GREETINGS FROM THE JSC STORAGE TEAM