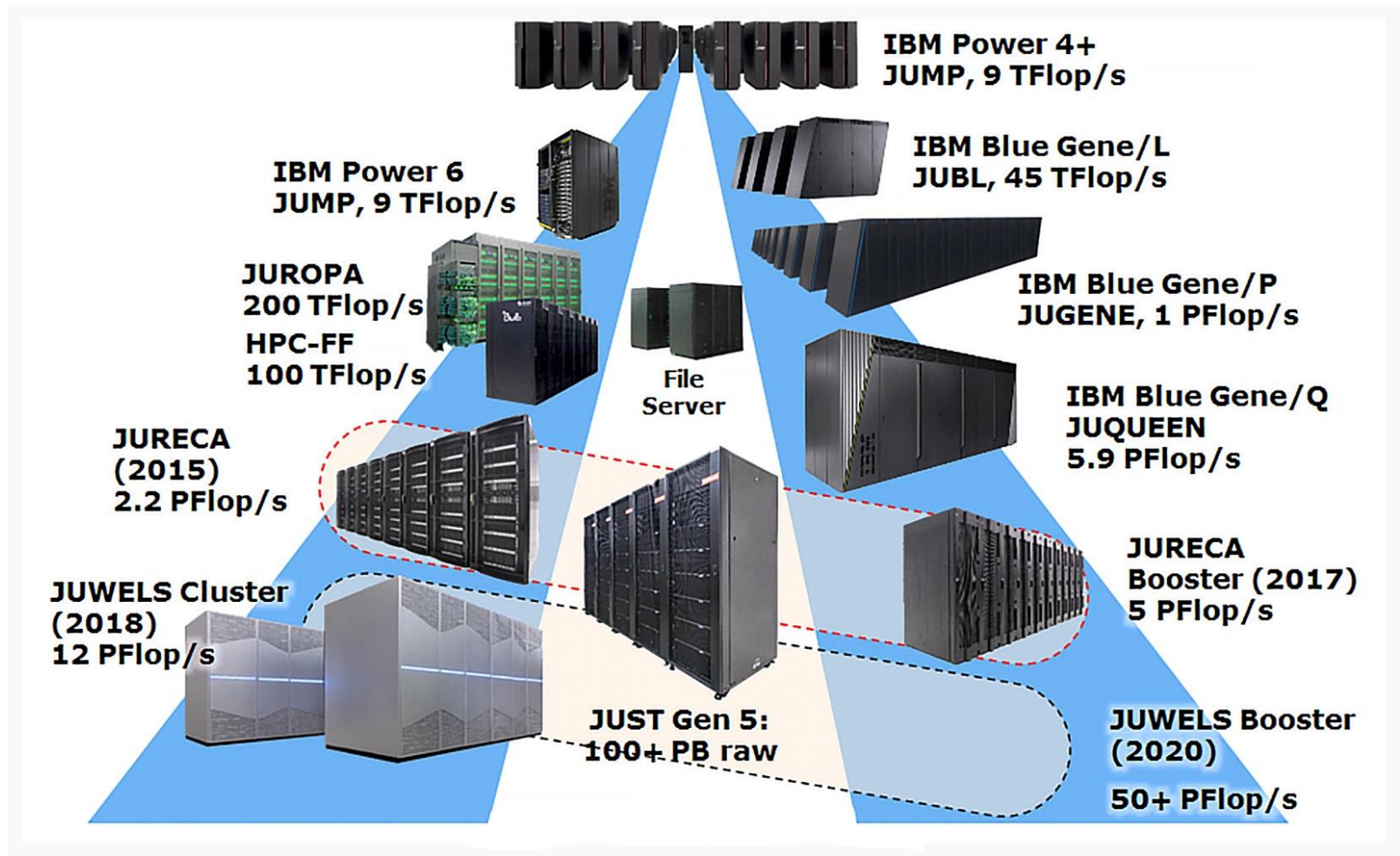


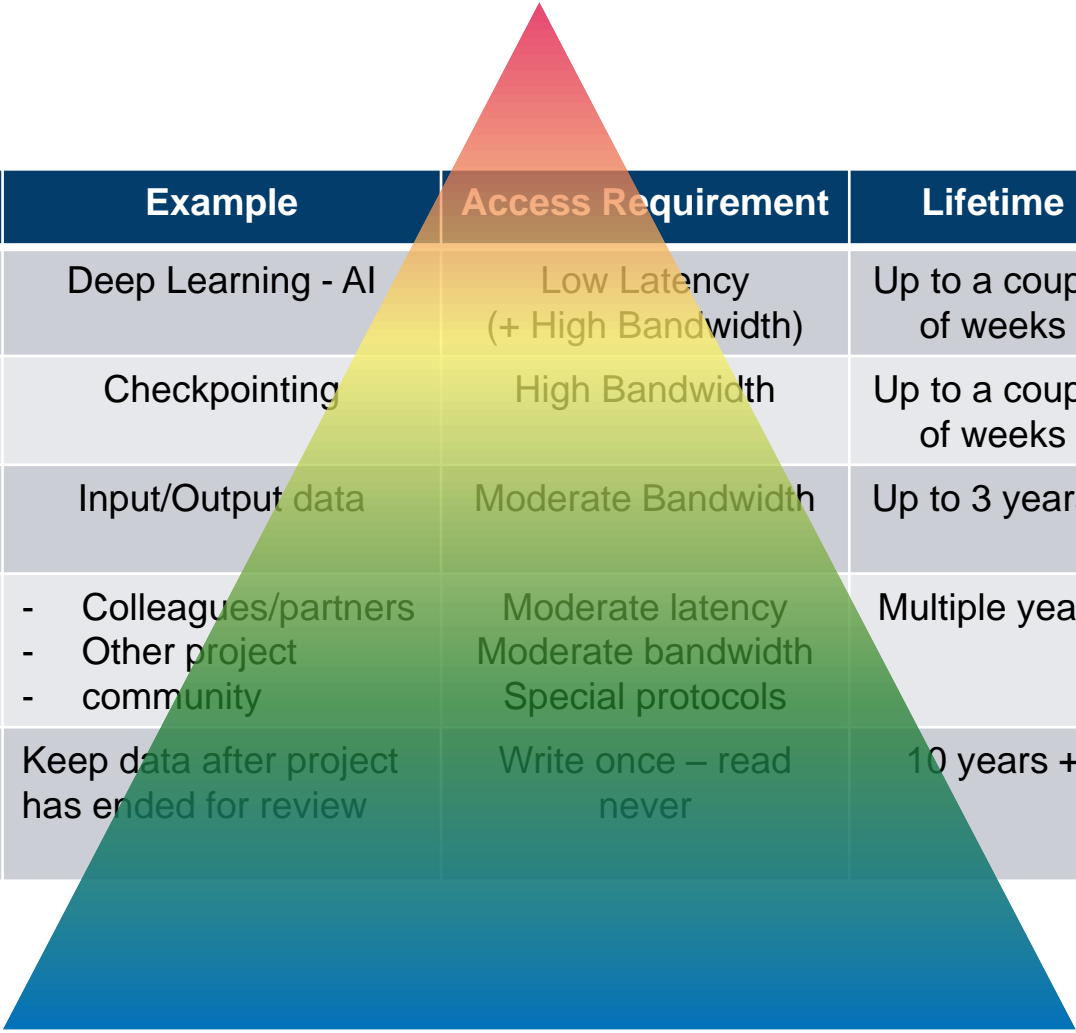
# FROM XCST TO HPST JÜLICH STORAGE CLUSTER JUST

25. MARCH 2021 | STEPHAN GRAF (JSC)

# DUAL ARCHITECTURE STRATEGY

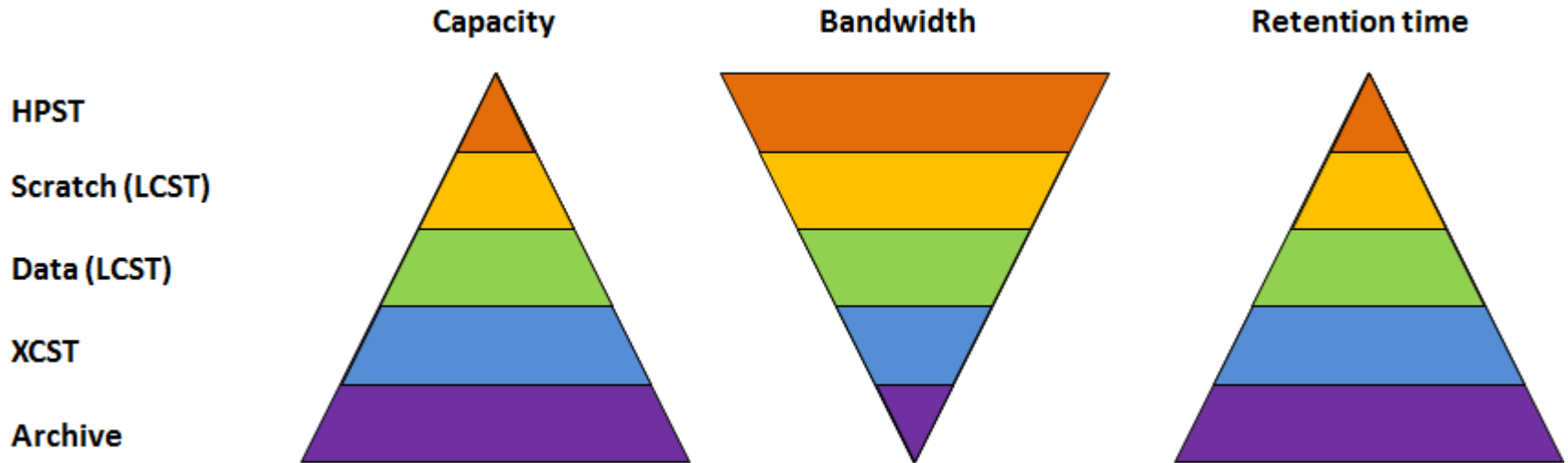


# IO USE CASES



Scenario	Example	Access Requirement	Lifetime	Data Reliability
Small random IO	Deep Learning - AI	Low Latency (+ High Bandwidth)	Up to a couple of weeks	Can be recreated easily
Streaming IO	Checkpointing	High Bandwidth	Up to a couple of weeks	Can be recreated easily
Store data for the project	Input/Output data	Moderate Bandwidth	Up to 3 years?	Backup
Share data	<ul style="list-style-type: none"> <li>- Colleagues/partners</li> <li>- Other project</li> <li>- community</li> </ul>	Moderate latency Moderate bandwidth Special protocols	Multiple years	Backup
Archive data	Keep data after project has ended for review	Write once – read never	10 years +	Backup

# TIERED STORAGE OFFERING



- High Performance Storage Tier (HPST):  
NVMe based Storage (low latency+high bandwidth) → pilot phase
- Large Capacity Storage Tier (LCST): disk based, bandwidth optimized
- Extended Capacity Storage Tier (XCST): disk based, capacity optimized
- Archive: Tape storage + disk based cache

# IBM SPECTRUM SCALE (GPFS)

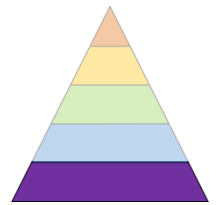
## Cluster File System

- General Parallel File System
- Optimized for HPC parallel file access
- POSIX compliant
- Client-Server architecture
- Special features included
  - RDMA
  - Hierarchical Storage Management (**HSM**, Spectrum Protect)
  - Optimized backup (mmBackup, Spectrum Protect)
  - Cluster Export Services – Protocol support (**NFS**, SMB, **S3/Swift**, HDFS)
  - AFM (Active File Management)
  - Encryption ...



# ARCHIVE

## GPFS combined with HSM



Move data to  
\$ARCHIVE

GPFS Building Block  
(500 x 10TB)



After 7 days  
moved to tape  
and replaced by  
stub file

Tape Library



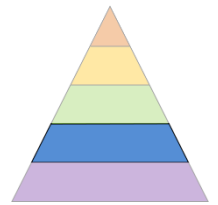
```
[zdv124@judac01:/arch/zam/zdv124> ls -liah
```

```
total 320K
```

```
 407977 128K drwx----- 2 zdv124 zam 64K May 18 10:01 .
 407555 128K drwxr-xr-x 316 root sys 64K May 24 15:00 ..
18062260 64K -rw-r--r-- 1 zdv124 zam 5 Sep 2 2011 date.txt
12920848 0 -rw-r--r-- 1 zdv124 zam 12G Jun 3 2015 Vervet_s0050_tiff.tgz
```

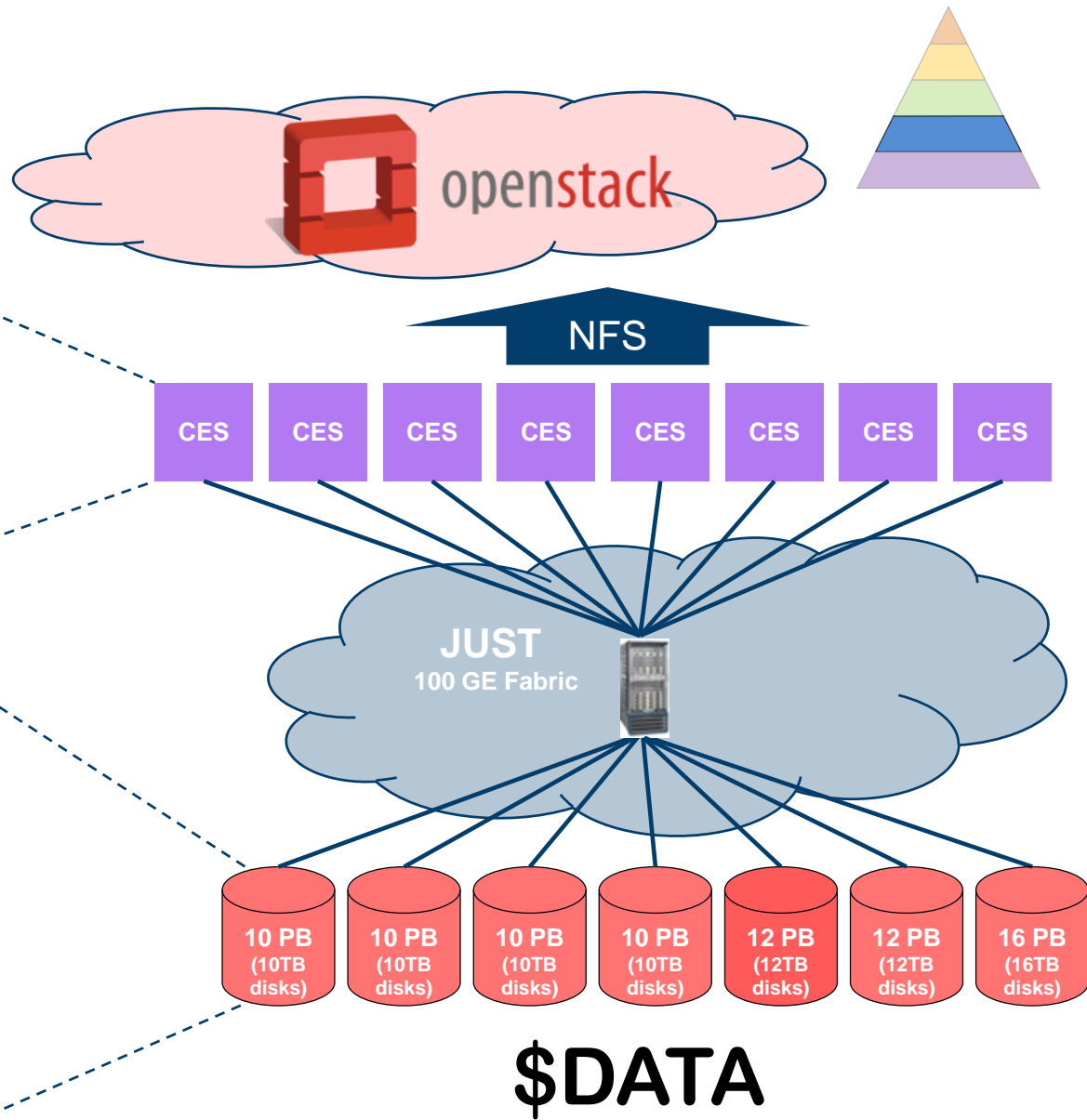
# XCST

## Extendet Capacity Storage Tier



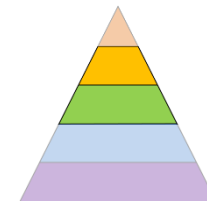
- Goal: „online“ data (disk), maximum capacity
- Campaign storage: storing and sharing your data
- Available on
  - HPC login nodes (not on the compute nodes)
  - JUDAC system (JUelich Data Access)
  - HDF Cloud (openstack)
- Phased installation: new hardware every year, but different disks (10/12/16TB)  
this year final phase 5 will be installed
- Two different repository types
  - File system access
  - Object store (in preparation)
- One global namespace

# XCST SETUP

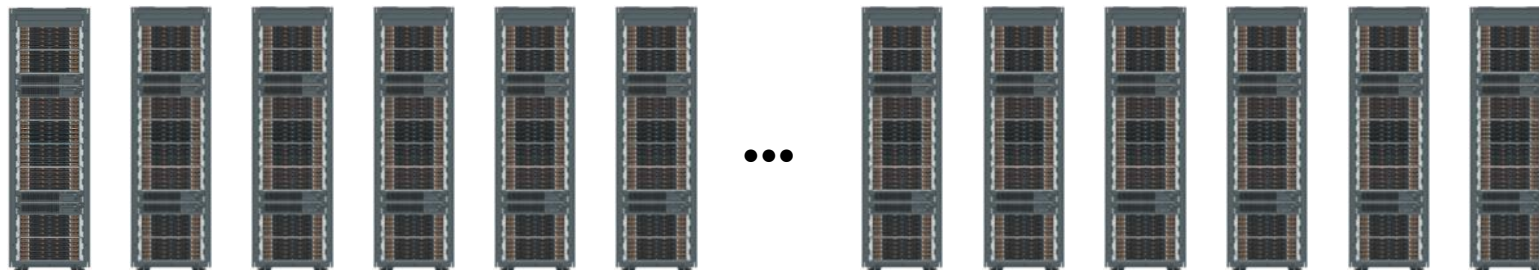


# LCST

## Large Capacity Storage Tier



21 x GPFS Building Blocks → 42 x NSD Server, 84 x Enclosure → +7.000 x 10TB disks



3 x 100 GE



2 Cluster Export  
Server (NFS)

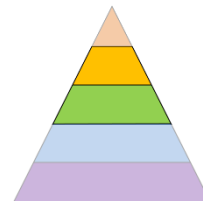
### Characteristics:



- Spectrum Scale (GPFS) + GNR (GPFS Native RAID)
  - Declustered RAID technology
  - End-to-End data integrity
- Spectrum Protect (TSM) for Backup
- Hardware:
  - x86 based server + RHEL 7
  - IBM Power 8 + AIX 7.2
  - 100GE network fabric
- 70 PB gross capacity
- Bandwidth: **400 GB/s**



# LCST

## Multiple file systems in one cluster

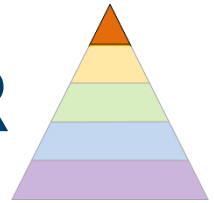


- **/p/scratch**   
→ temporary file system for compute jobs, high bandwidth
  - No backup, data are deleted after 90 days
- **/p/fastdata**   
→ high bandwidth file system for persistent data to compute (e.g. training data)
  - Data are in backup
- 3 building blocks for the other file systems
  - **/p/project** → compute project repository
  - **/p/software** → data repository for sharing software between projects
  - ...

### The truth

Both uses the same disks  
(18 building blocks), so  
performance is the same

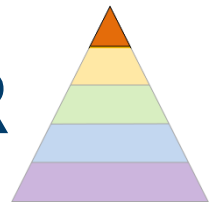
# HIGH PERFORMANCE STORAGE TIER



## NVMe based cache to underlying file system

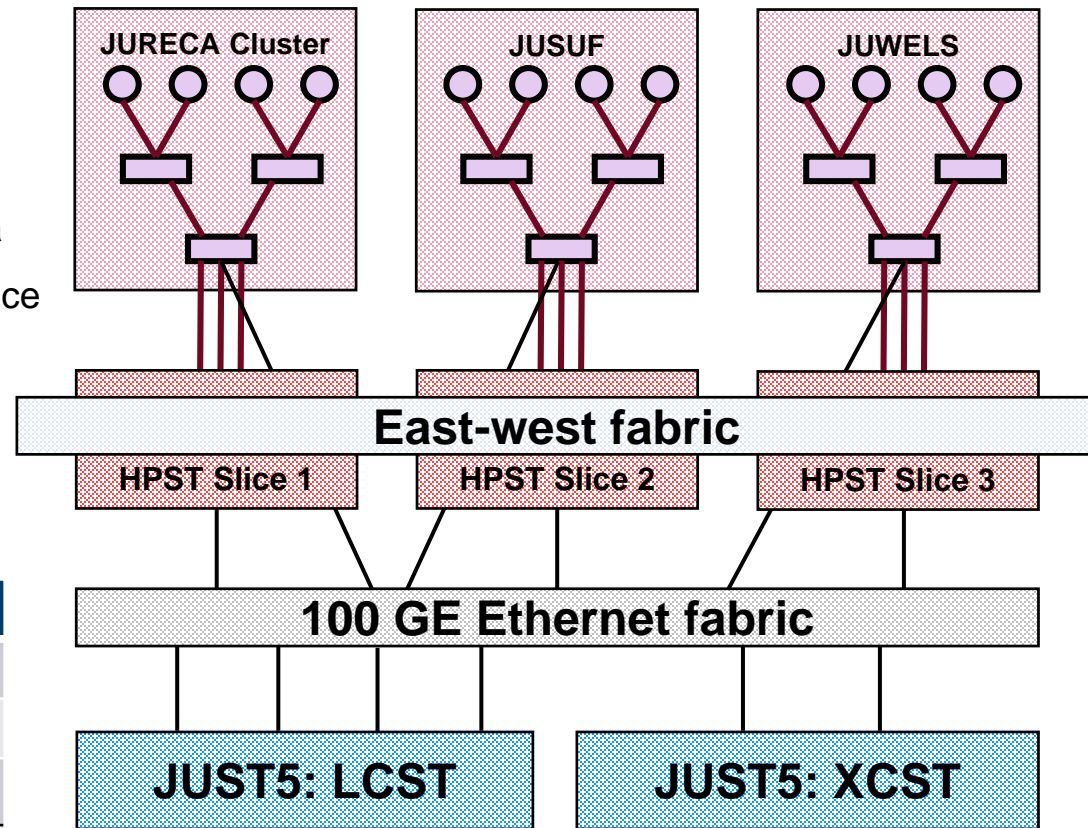
- Goal: High Bandwidth + low Latency for compute jobs
- Solution: NVMe Storage cluster + Infinite Memory Engine (IME) from DDN
- Design decision:
  - Each cluster has it's own “slice” of the HPST, but
  - One global namespace (each cluster has access to data on “foreign slice”)
  - Direct integrated in cluster fabric (IB)
- Client access by mount point (FUSE) or IME native interface
- Cache of the \$SCRATCH file system
- CLI + API for data migration & data staging (e.g. flush, sync, pre-stage, pin)
- Path: `/p/ime-scratch/fs/<group>` → is a mapping to `/p/scratch/<group>`
- Access is managed by special HPST group ID (GID)

# HIGH PERFORMANCE STORAGE TIER



## Infinite Memory Engine (IME) from DDN

- 110 Server IME-140, each
  - 2 x HDR 100 (North connection)
  - 1 x HDR 100 (East-West connection)
  - 1 x 100 GE (South connection)
  - 10 x NVMe drives (2TB) for HPST data
  - 1 x internal NVMe IME commit log device



Slice	Size	Bandwidth*
JUWELS	1036 TB	1024/600 GB/s
JURECA	844 TB	800/600 GB/s
JUSUF	230 TB	240/220 GB/s
<b>Total</b>	<b>2110 TB</b>	<b>2064/1420 GB/s</b>

\*committed values

# TIERED STORAGE OFFERING

## Summary

Layer	ENV Name	characteristics	IO type	Write	Read
HPST	-	<ul style="list-style-type: none"><li>- NVMe</li><li>- Cache</li><li>- Compute with hard IO pattern</li></ul>	Single client Multiple client Multiple client (small IO)	3,3 GB/s 600 GB/s 550 GB/s	7,4 GB/s 730 GB/s 320 GB/s
Scratch (LCST)	<b>\$SCRATCH</b>	<ul style="list-style-type: none"><li>- Temporary files of compute jobs</li><li>- streaming</li></ul>	Single client Multiple client	8GB/s >300 GB/s	9 GB/S 400 GB/s
Data (LCST)	<b>\$FASTDATA</b>	<ul style="list-style-type: none"><li>- Persistent files for compute jobs</li><li>- streaming</li></ul>	Single client Multiple client	8GB/s >300 GB/s	9 GB/S 400 GB/s
XCST	<b>\$DATA</b>	<ul style="list-style-type: none"><li>- Persistent files</li><li>- Share data with external users (Web services)</li></ul>	Single client Multiple client	8 GB/s >50 GB/s	9 GB/s > 70 GB/s
Archive	<b>\$ARCHIVE</b>	<ul style="list-style-type: none"><li>- Archive data</li><li>- Read data seldom</li></ul>	Single client	~8 GB/s	~200 MB/s

# DATA SHARING INSIDE JÜLICH HPC

Different use cases and solutions for sharing data between users:

1. **Use compute project repository (\$PROJECT)**

Any user can be joined to project without access to project's compute resources

2. **Use data project**

Members of different compute projects can join a common data project

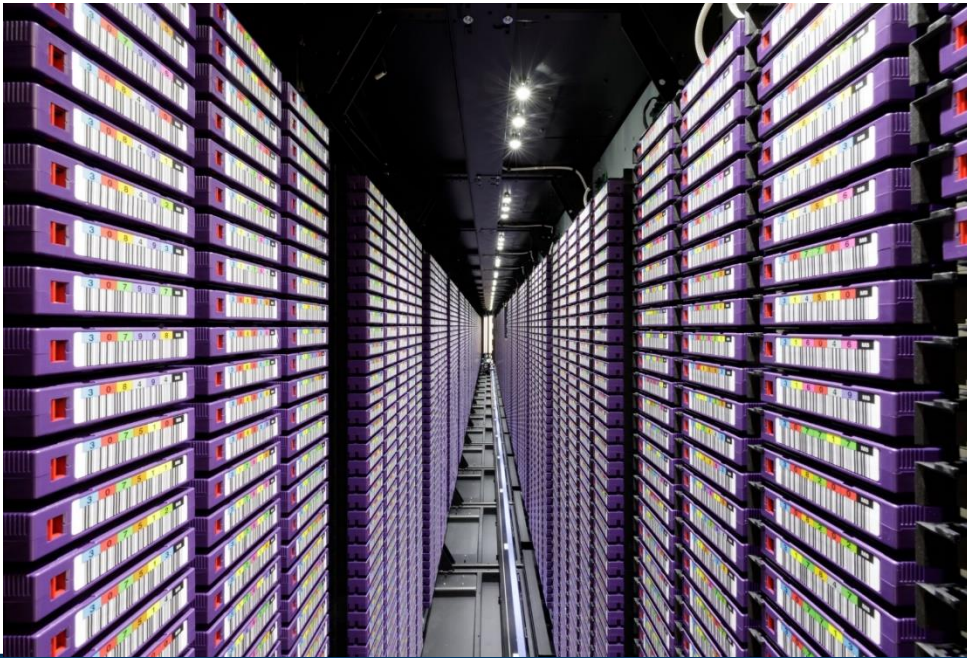
3. **Single files**

All users can access common directory „**\$SCRATCH**/.../share“.

Remember the automatic file deletion after 90 days!

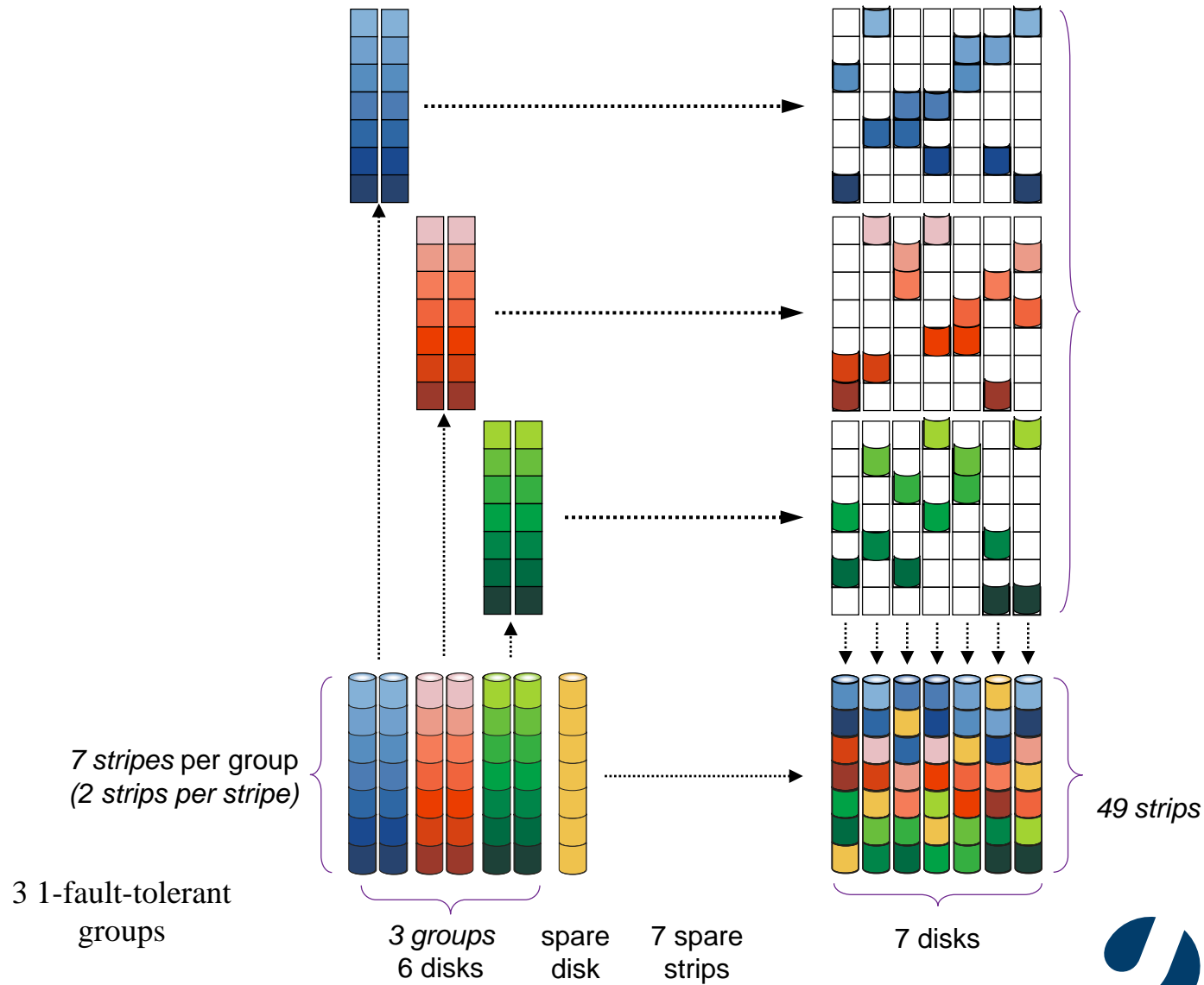
4. **Software project**

Special data project which is mounted on compute nodes

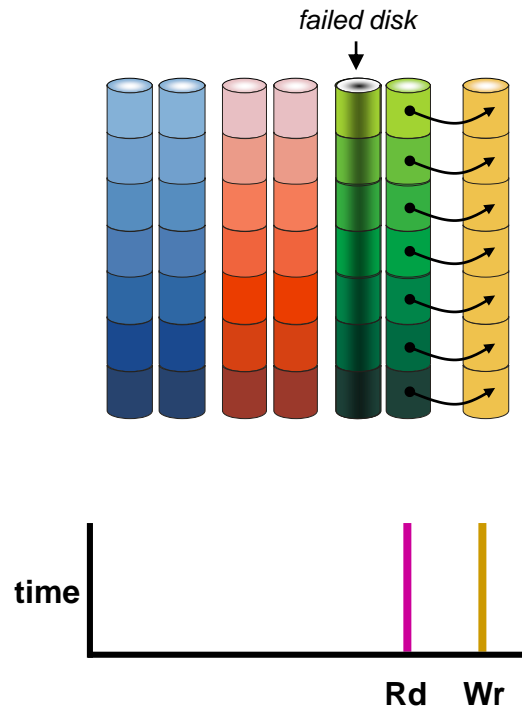


# THANK YOU FOR YOUR ATTENTION

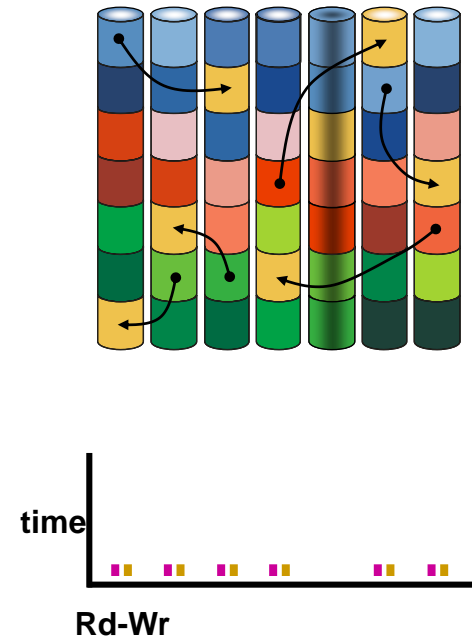
# DECLUSTERED RAID



# DECLUSTERED RAID REBUILD EXAMPLE



- Disk failure causes disk rebuild
- Volume degraded for a long time
- performance impact for file system



- Disk failure causes strips rebuild
- all disc involved
- Volume degraded for a short time
- minimized performance impact