



Leibniz Supercomputing Centre
of the Bavarian Academy of Sciences and Humanities

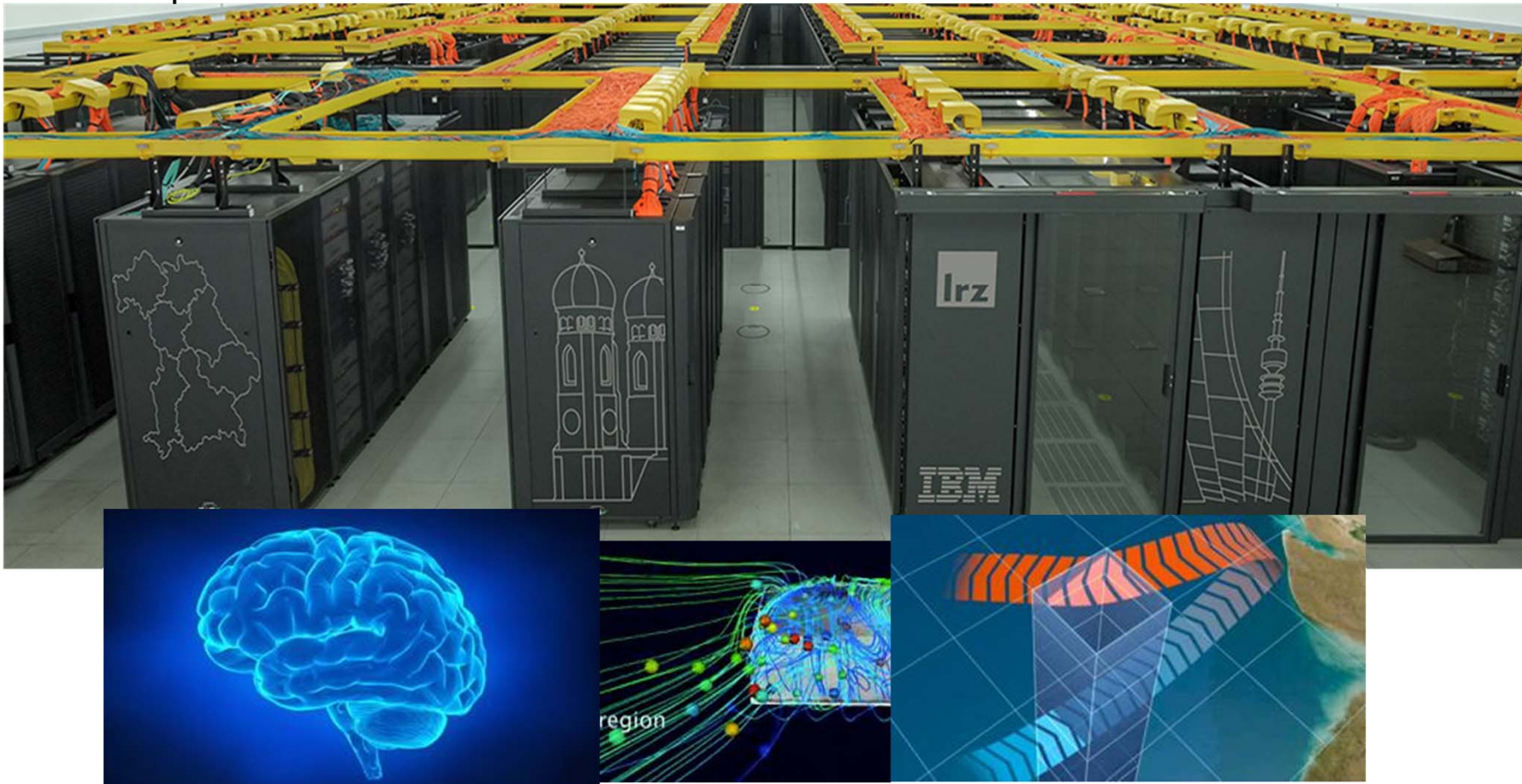


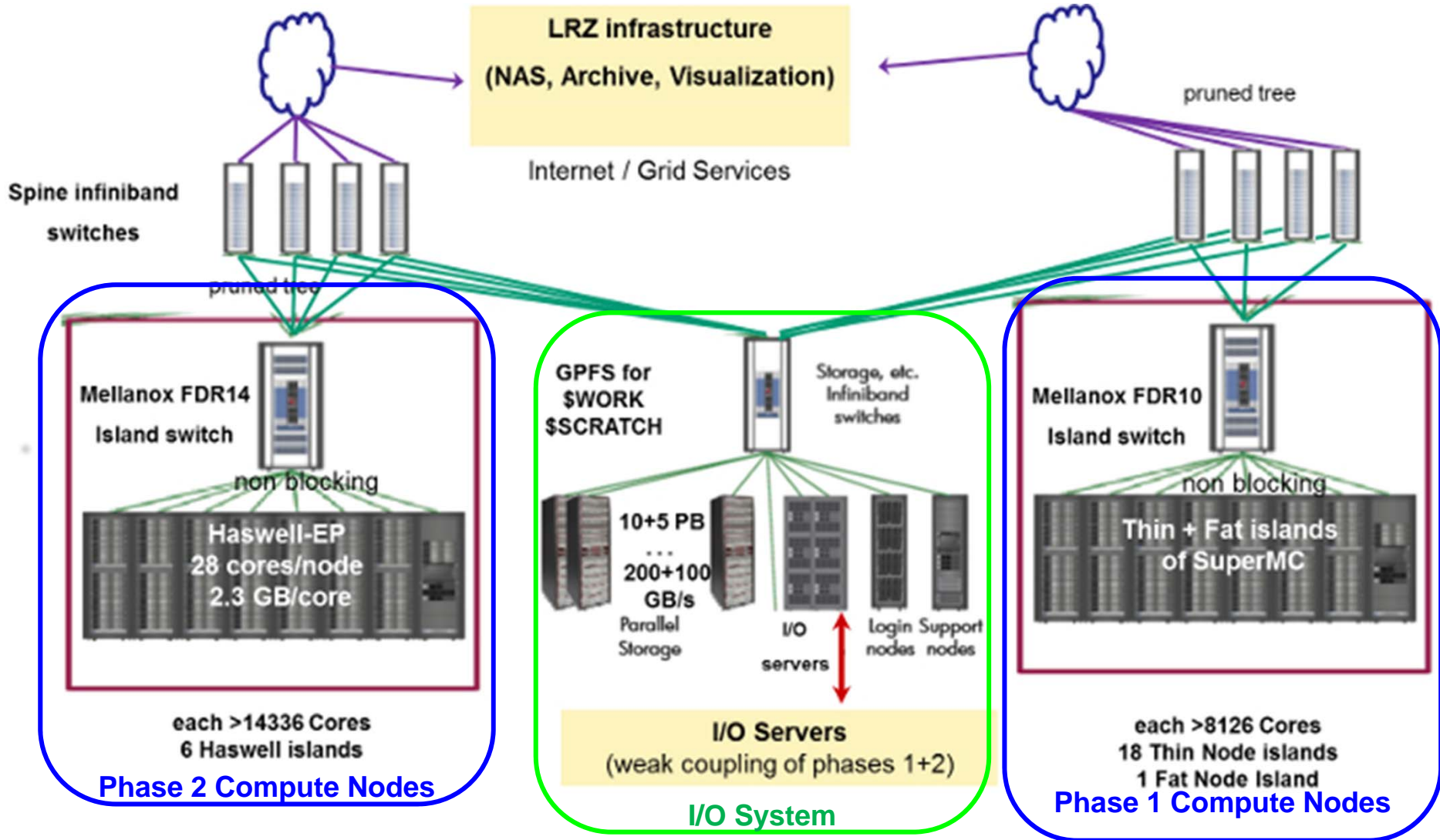
Analyzing the High Performance Parallel I/O on LRZ HPC systems

Sandra Méndez.
HPC Group, LRZ.
June 23, 2016

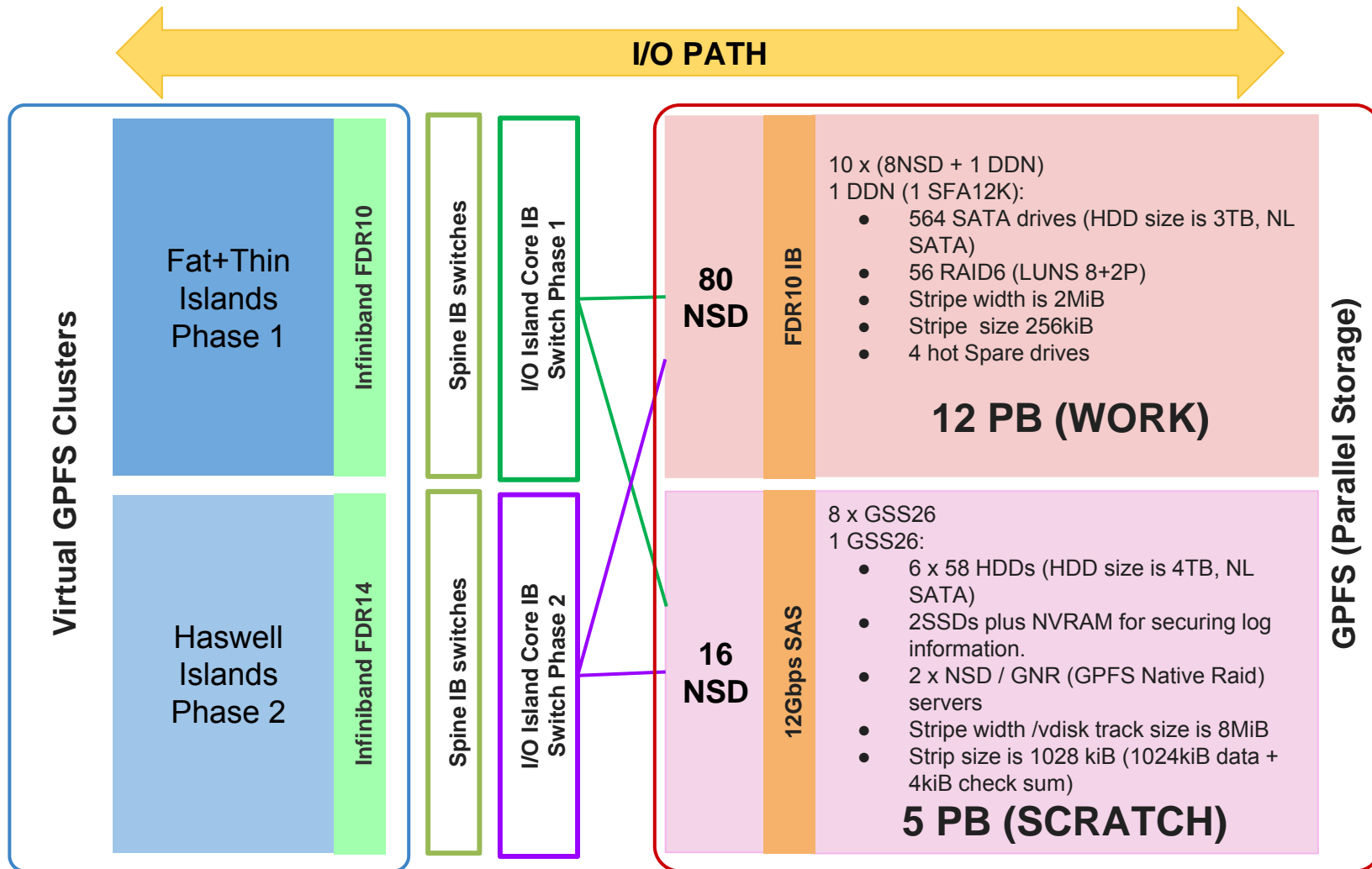
- **SuperMUC supercomputer**
- **User Projects**
- **Monitoring Tool**
- **I/O Software Stack**
- **I/O Analysis Tool**
- **Analyzing I/O Problems**
- **Conclusions**

- Member of the Gauss Centre for Supercomputing (GCS). Tier-0 centre for PRACE, the Partnership for Advanced Computing in Europe.
- 2012 SuperMUC Phase 1 and 2015 SuperMUC Phase 2. Total Peak Performance 6.4 PFlop/s.





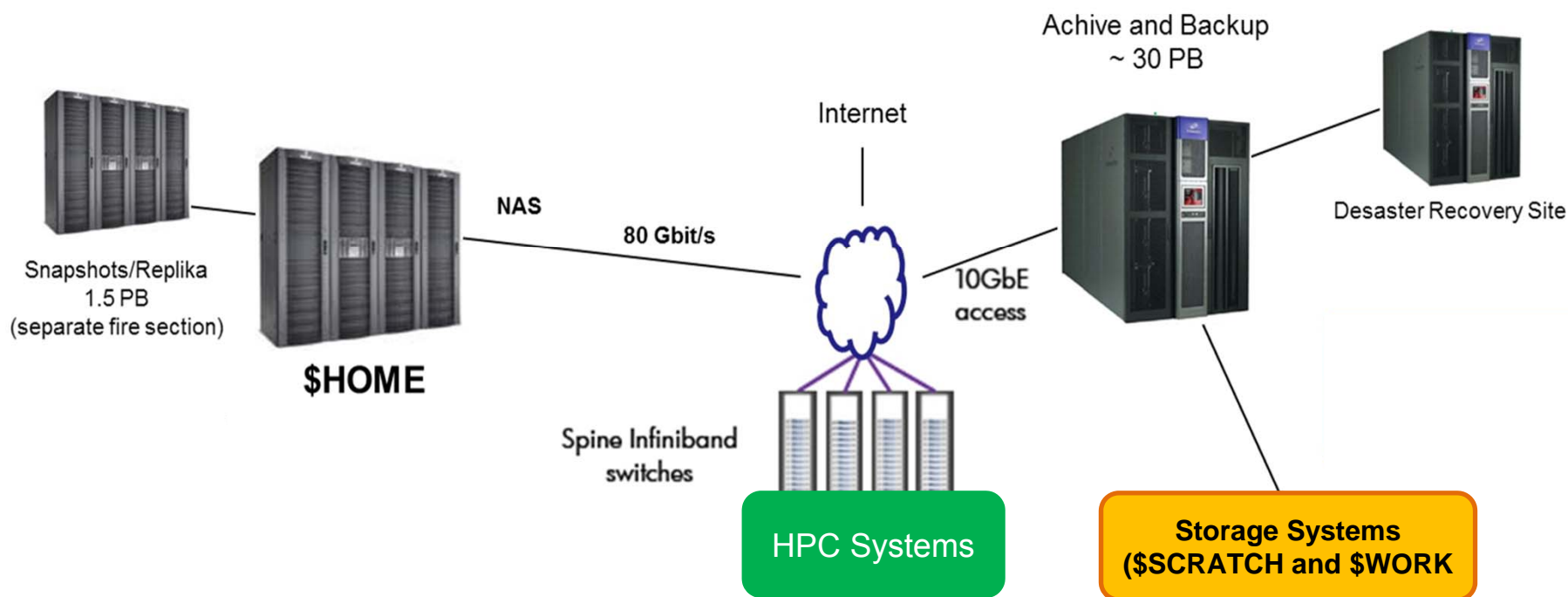
The I/O PATH on SuperMUC - Parallel Storage (WORK and SCRATCH filesystems)



- Available on all HPC cluster systems (environment variable `$HOME`)
- Shared area for all user accounts in a project

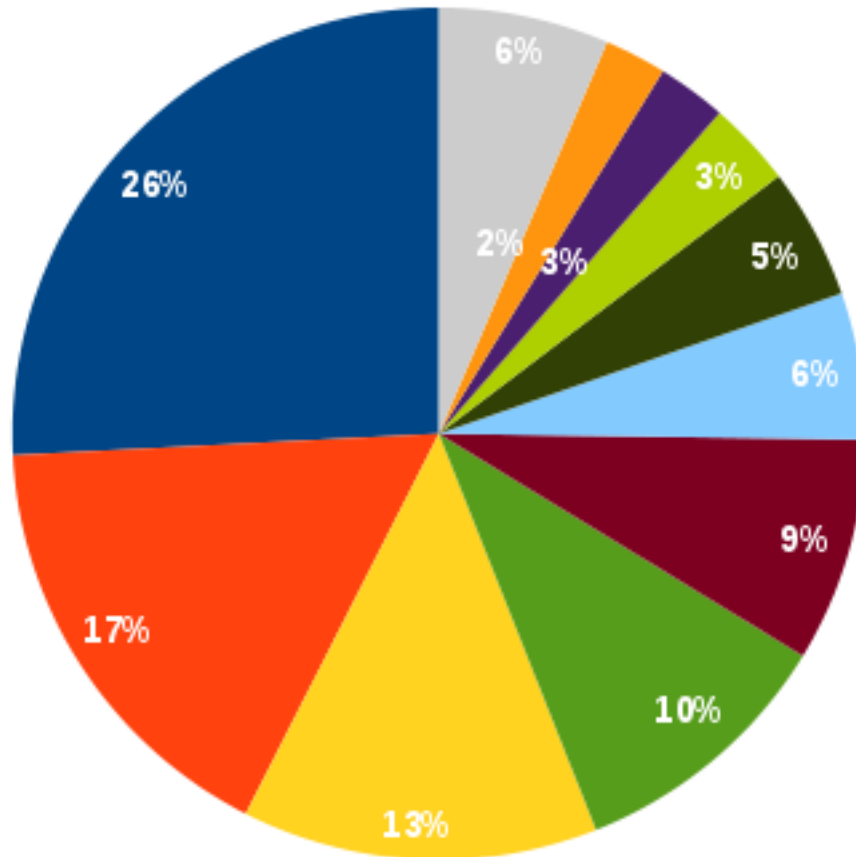
Very reliable

- user-restorable snapshots (last 10 days)
- automatic data protection by LRZ

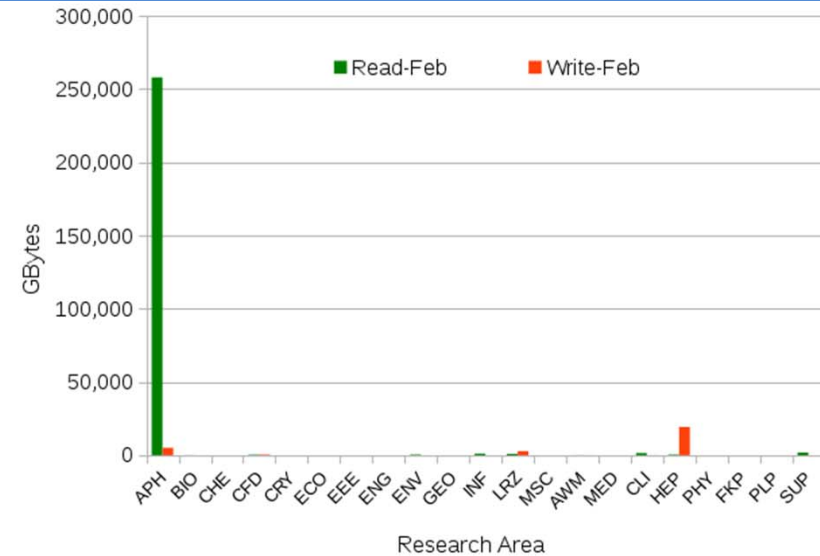
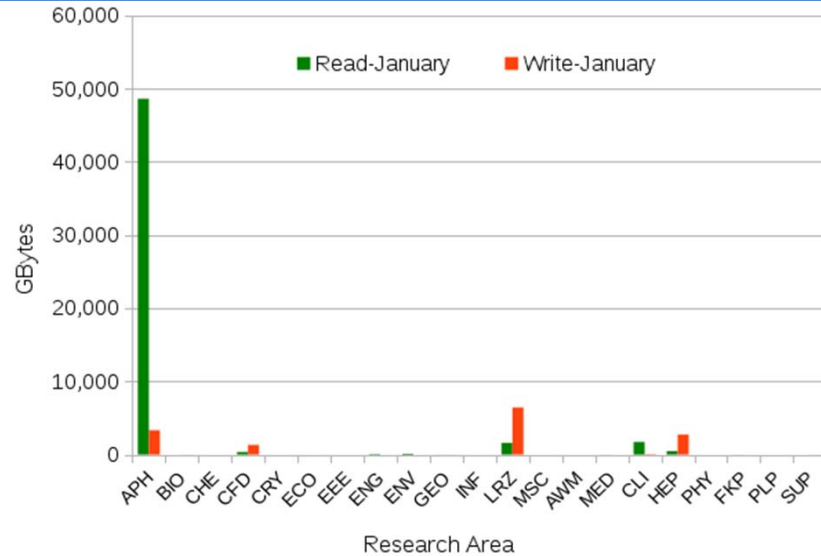


- SuperMUC supercomputer
- **User Projects**
- Monitoring Tool
- I/O Software Stack
- I/O Analysis Tool
- Analyzing I/O Problems
- Conclusions

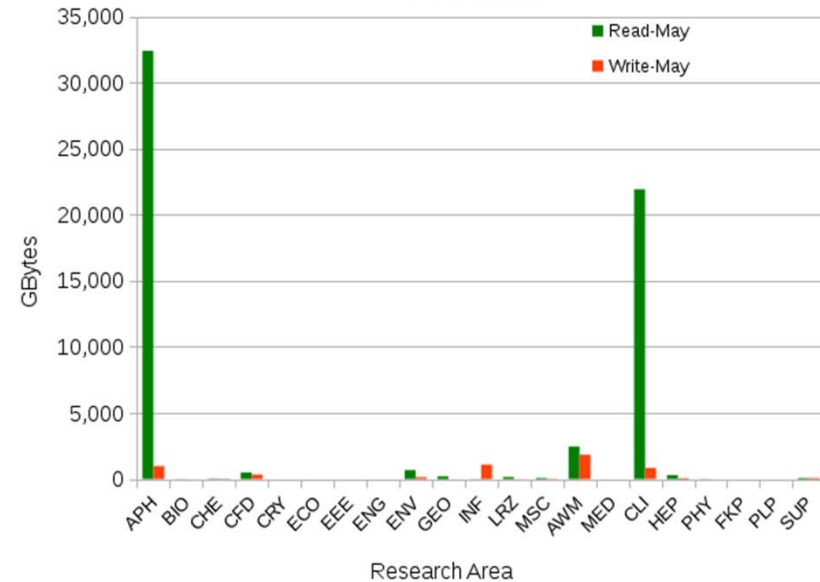
- Computational-Fluid-Dynamics (CFD)
- Astrophysics-Cosmology (APH)
- Informatics-ComputerSciences (INF)
- Chemistry (CHE)
- Biophysics-Biology-Bioinformatics (BIO)
- Physics-High-EnergyPhysics (HEP)
- Physics-Solid-State (FKP)
- Geophysics (GEO)
- Engineering-others (ENG)
- Meteorology-Climatology-Oceanography (CLI)
- Other



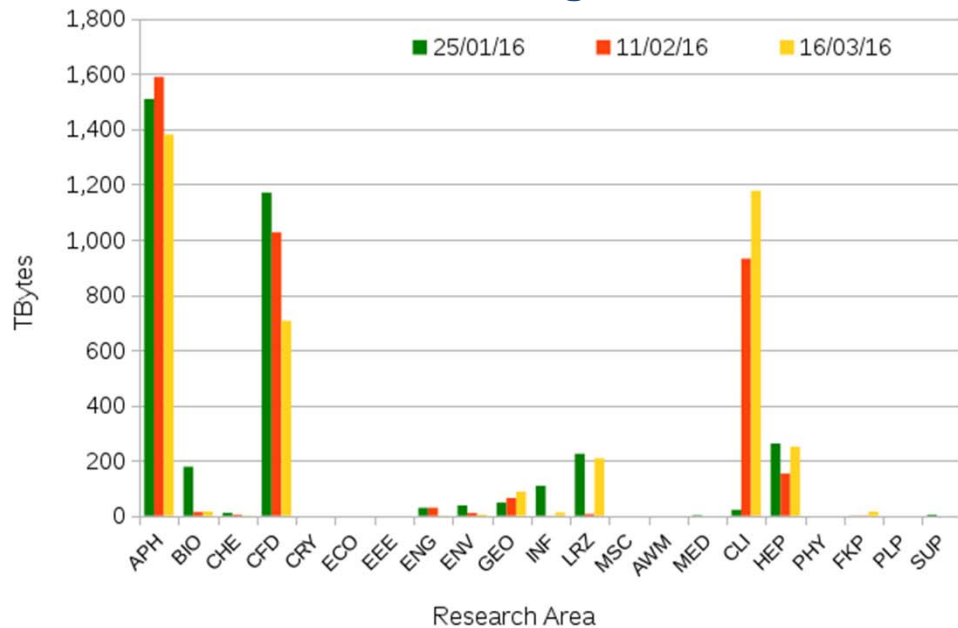
Data Transferred by Research Area



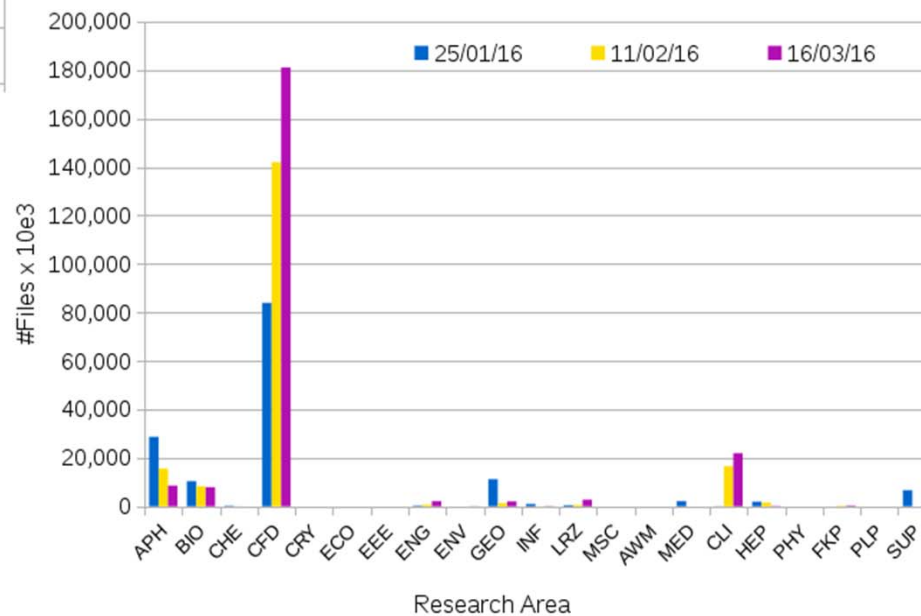
APH	Astrophysics-Cosmology
BIO	Biophysics-Biology-Bioinformatics
CHE	Chemistry
CFD	Computational-Fluid-Dynamics
CRY	Crystallography
ECO	Economics
EEE	Engineering-ElectricalEngineering
ENG	Engineering-others
ENV	Enviromental-Sciences
GEO	Geophysics
INF	Informatics-ComputerSciences
LRZ	LRZ
MSC	Material-Science
AWM	Mathematics
MED	Medicine
CLI	Meteorology-Climatology-Oceanography
HEP	Physics-High-EnergyPhysics
PHY	Physics-others
FKP	Physics-Solid-State
PLP	Plasma-Physics
SUP	Support-Benchmarking



Storage



Count of Files



I/O Libraries

- HDF5 15%, NetCDF or PnetCDF 10%; POSIX, MPI-IO, or an I/O library locally installed 75%.

Storage Parallel

- WORK (70% Capacity) -> 5 fold increase
- SCRATCH (80% Capacity) -> 8 fold increase

Checkpointing and large scale output with a connection to a visualization cluster.

Checkpointing (for the Large-Scale Projects):

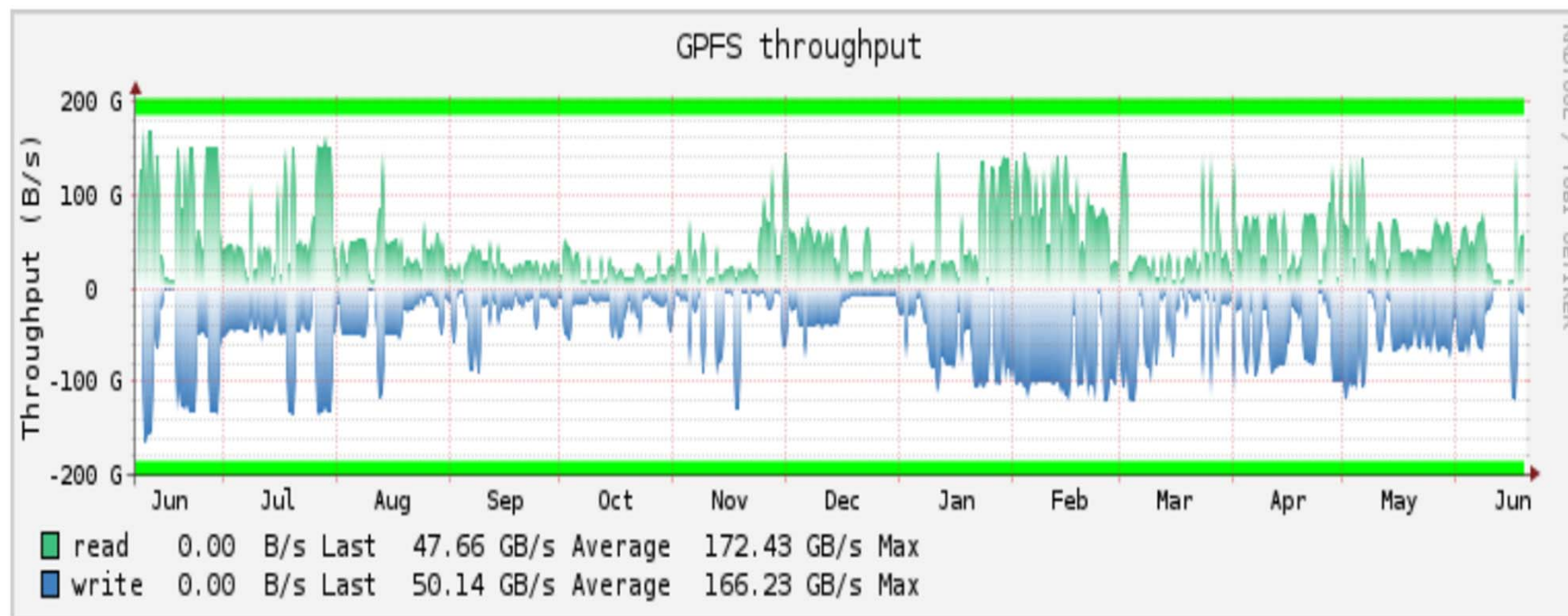
Periods: 5 min to 8 hours

Size: 100 GB -> 38%
1TB -> 10%
5TB -> 7%
10TB -> 1%
35TB -> 2%
70TB -> 1%
< 100GB -> 41%

- SuperMUC supercomputer
- User Projects
- **Monitoring Tool**
- I/O Software Stack
- I/O Analysis Tool
- Analyzing I/O Problems
- Conclusions

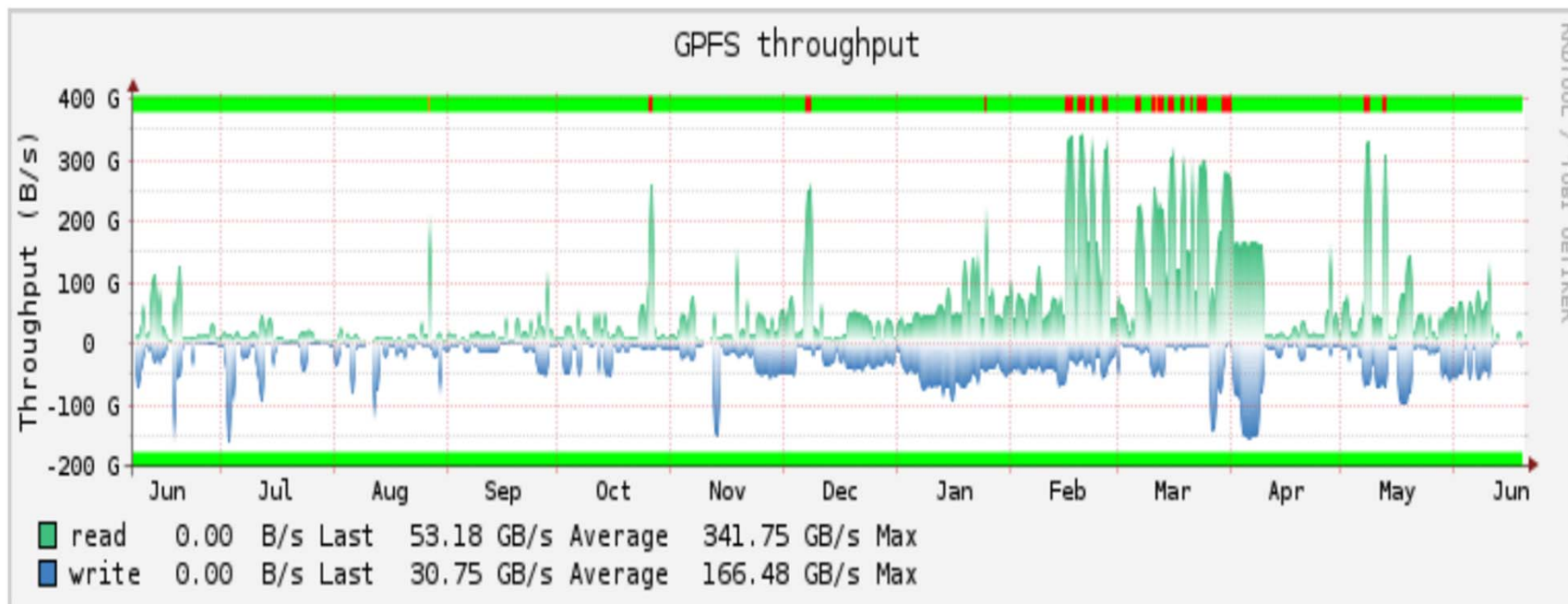
One Year (06.06.15 19:01 - 20.06.16 19:01)

Datasource Throughput



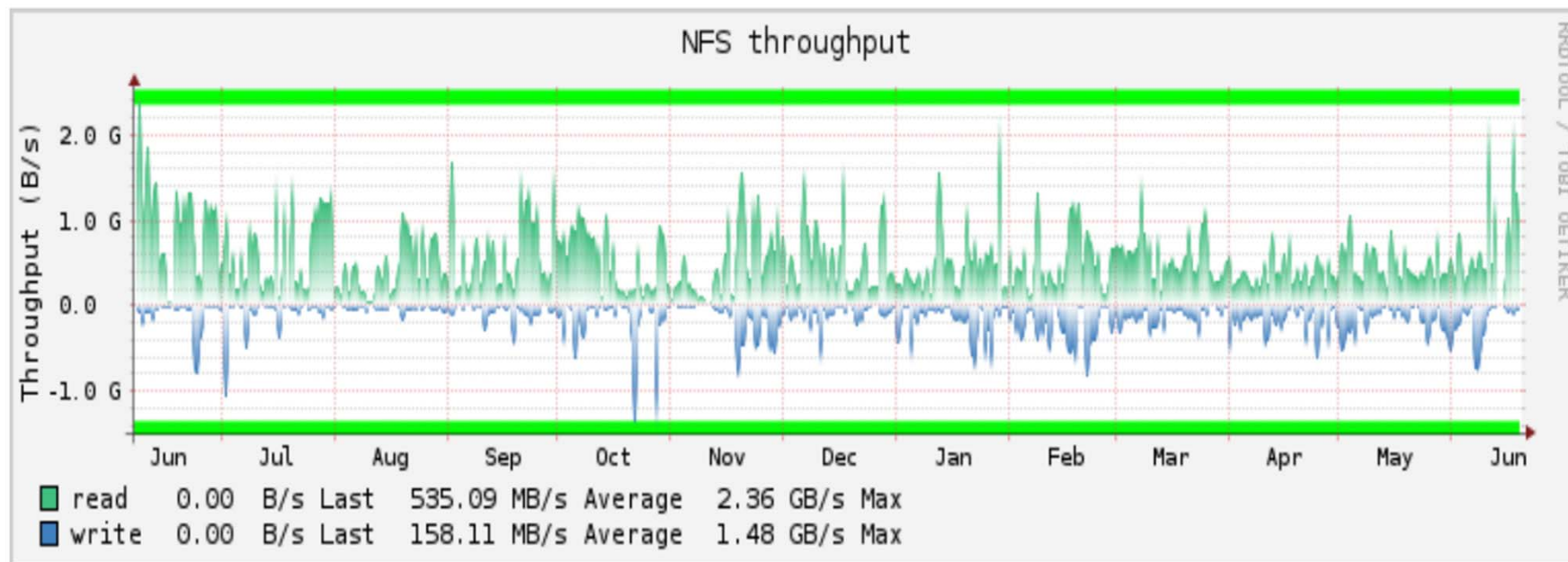
One Year (06.06.15 19:00 - 20.06.16 19:00)

Datasource Throughput



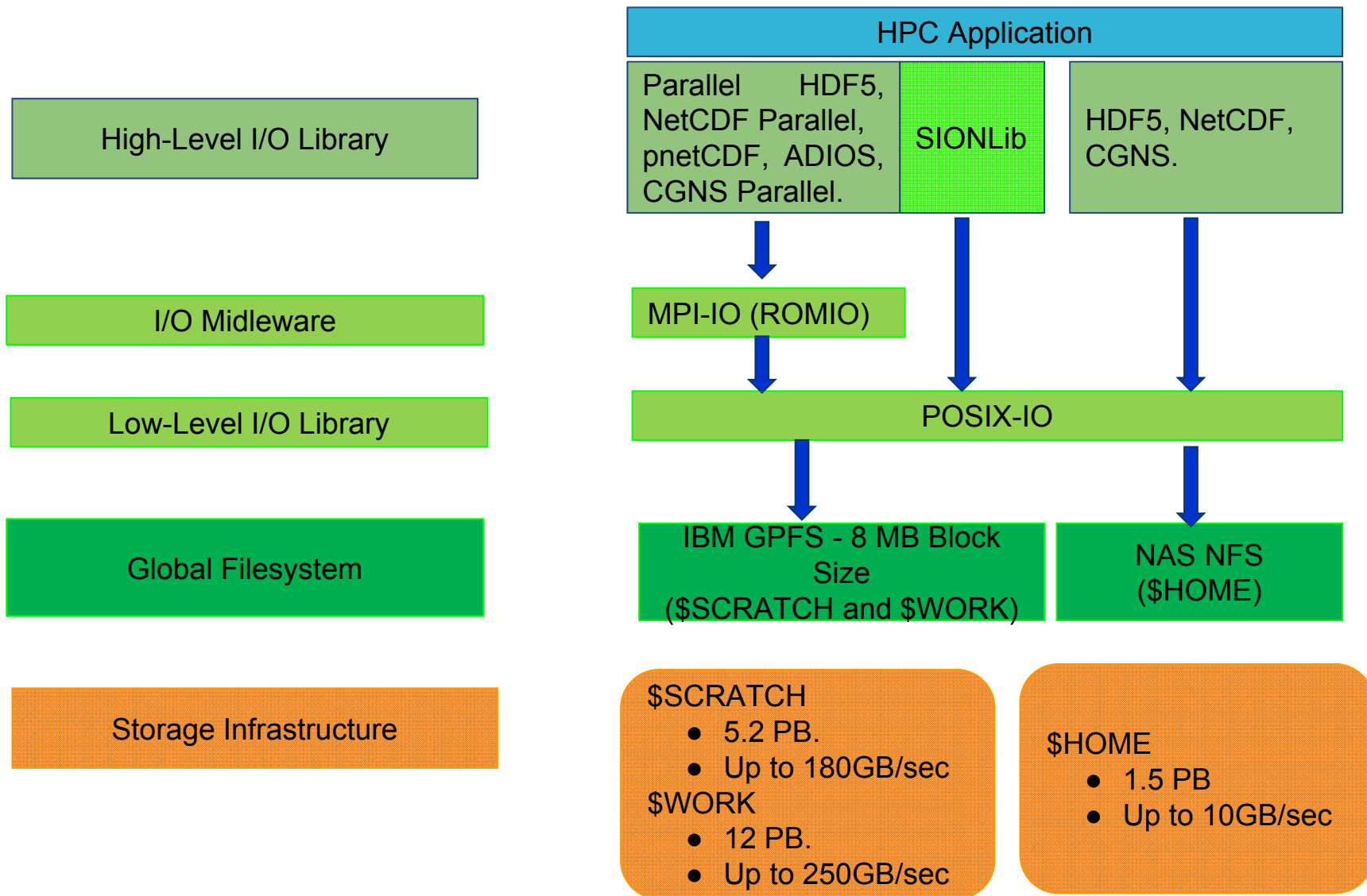
One Year (06.06.15 18:59 - 20.06.16 18:59)

Datasource Throughput





- SuperMUC supercomputer
- User Projects
- Monitoring Tool
- **I/O Software Stack**
- I/O Analysis Tool
- Analyzing I/O Problems
- Conclusions



- **IBM Parallel Environment** (MPI-IO implementation a ROMIO version)
- **Intel MPI** (MPI-IO implementation a ROMIO version)
 - Set the `I_MPI_EXTRA_FILESYSTEM` environment variable to `on` to enable parallel file system support.
 - Set the `I_MPI_EXTRA_FILESYSTEM_LIST` environment variable to request native support for the specific file system.

ROMIO Hints

Data Sieving:

- `ind_rd_buffer_size`
- `ind_wr_buffer_size`
- `romio_ds_read`
- `romio_ds_write`

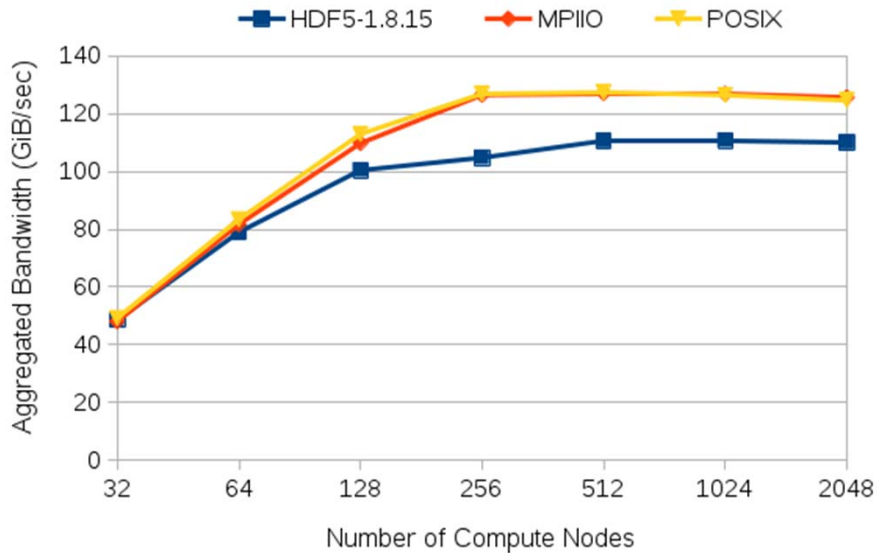
Collective buffering (Two-Phase I/O)

- `cb_buffer_size`
- `romio_cb_read`
- `romio_cb_write`

Exp.	MPI Processes	Compute Nodes	Access Pattern	Request Size	Hints	Transfer Rate(GiB/sec)
1	512	32	Sequential	1 MiB	romio_cb_read = automatic romio_cb_write = automatic	write = 25.92 read = 23.80
					romio_cb_read = disable romio_cb_write = disable	write = 75.34 read = 67.58
					independent I/O	write = 80.39 read = 69.62
2	512	32	Strided	1 MiB	romio_cb_read = automatic romio_cb_write = automatic	write = 1.63 read = 17.74
					romio_cb_read = enable romio_cb_write = enable	write = 25.49 read = 26.10
					independent I/O	write = 5.15 read = 12.60
3	512	32	Sequential	256 MiB	romio_cb_read = automatic romio_cb_write = automatic	write = 74.48 read = 46.67
					romio_cb_read = disable romio_cb_write = disable	write = 83.37 read = 65.88
					independent I/O	write = 82.29 read = 64.22
4	512	32	Strided	256 MiB	romio_cb_read = automatic romio_cb_write = automatic	write = 71.12 read = 41.80
					romio_cb_read = disable romio_cb_write = disable	write = 77.22 read = 70.21
					romio_cb_read = enable romio_cb_write = enable	write = 24.81 read = 24.90
					independent I/O	write = 71.25 read = 67.71

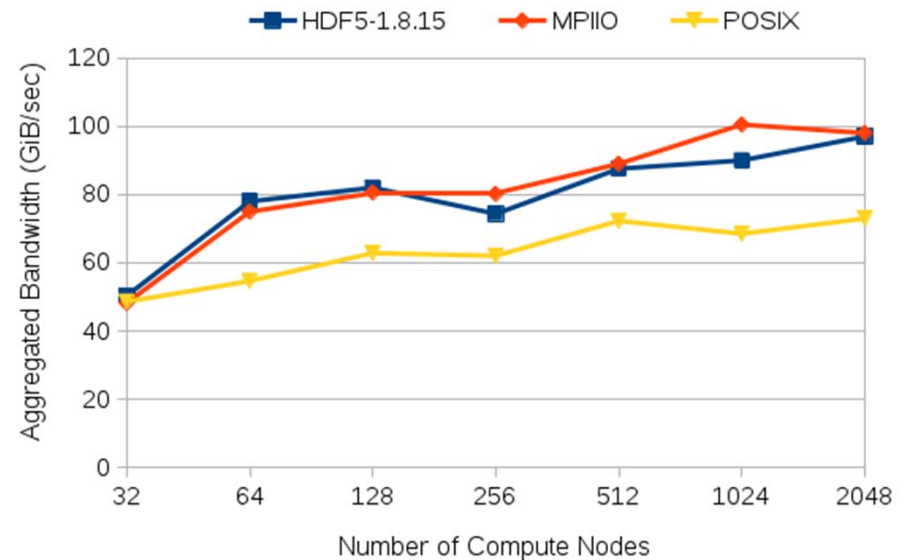
GPFS Scalability on SuperMUC: Weak scaling up to 4 Islands on SCRATCH

IOR Data Transfer Rate - SuperMUC - \$SCRATCH - GPFS
 Read - 64 GB per MPI Proc - 1 MPI proc per Compute Node



**Islands of Thin Nodes.
 Peak Performance 130 GB/sec
 (From Phase 1 to SCRATCH).**

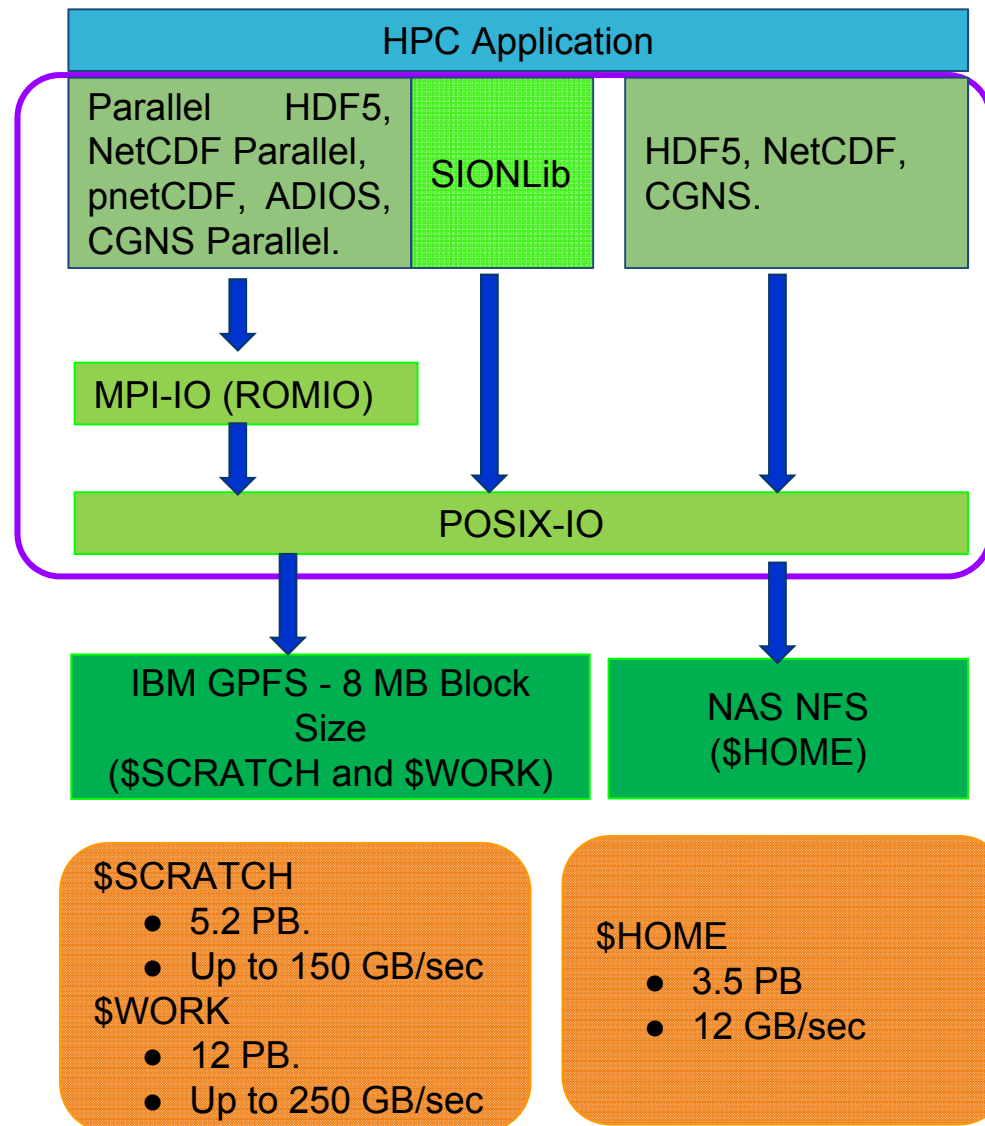
IOR Data Transfer Rate - SuperMUC - \$SCRATCH - GPFS
 Write - 64 GB per MPI Proc - 1 MPI proc per Compute Node

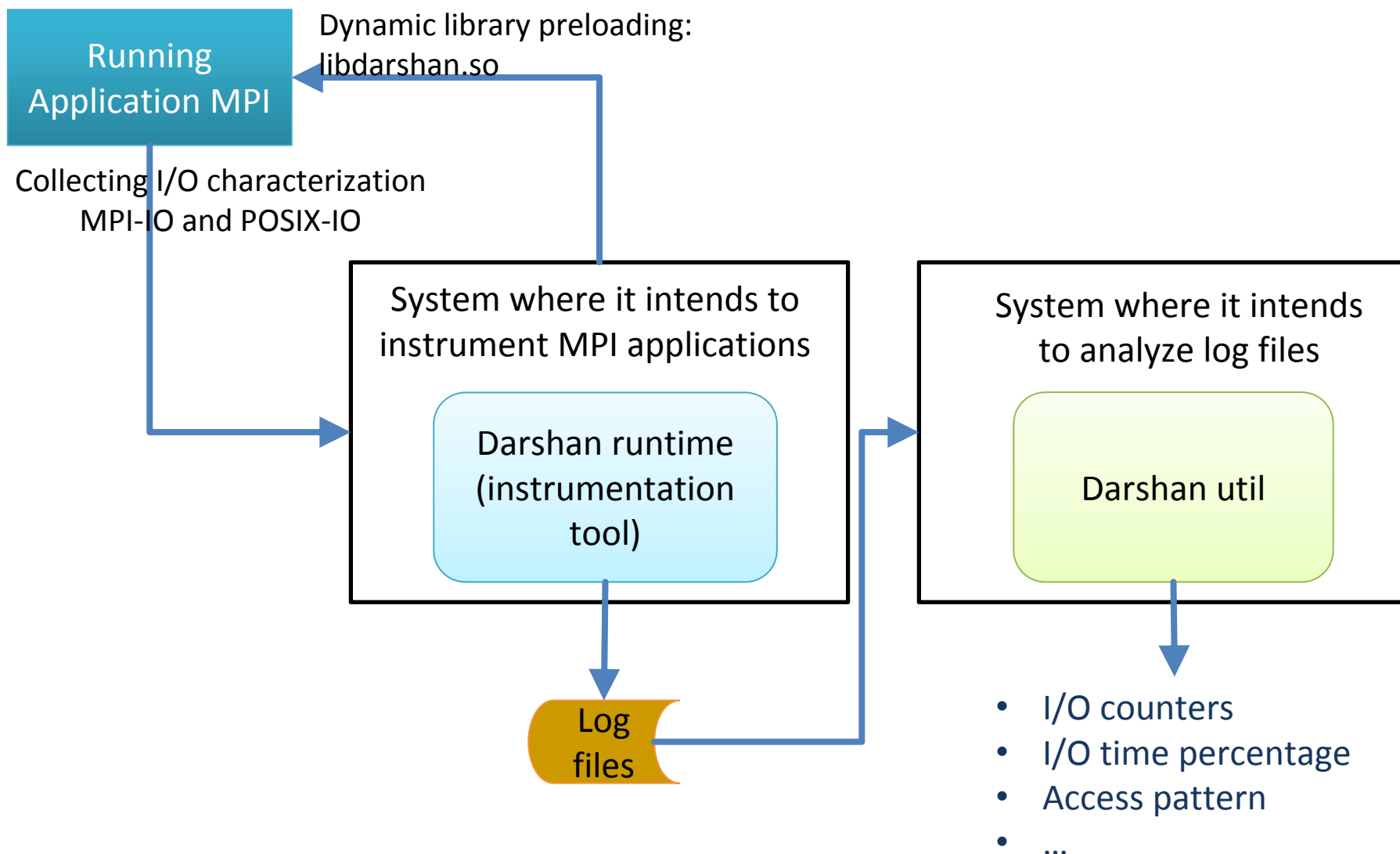


- SuperMUC supercomputer
- User Projects
- Monitoring Tool
- I/O Software Stack
- **I/O Analysis Tool**
- Analyzing I/O Problems
- Conclusions

I/O Profiling:

- DARSHAN (POSIX-IO, MPI-IO)





- **To make use of Darshan in its version 2.3 and 3.0, the module appropriate must be loaded.**

```
module load darshan
```

- **Set up the variable FORTRAN_PROG in “true” if the program is a Fortran program and false if it's not.**

```
FORTRAN_PROG=true
```

- **Load the appropriate library.**

```
export LD_PRELOAD=`darwin-user.sh $FORTRAN_PROG`
```

- **Set up Darshan job identifier with loadleveler job identifier.**

```
export JOBID_LL=`darwin-JOBID.sh $LOADL_STEP_ID`
```

- **Set up environment variable DARSHAN_JOBID to environment variable name that contain the job identifier of loadleveler.**

```
export DARSHAN_JOBID=JOBID_LL
```

- **Set up Darshan log path**

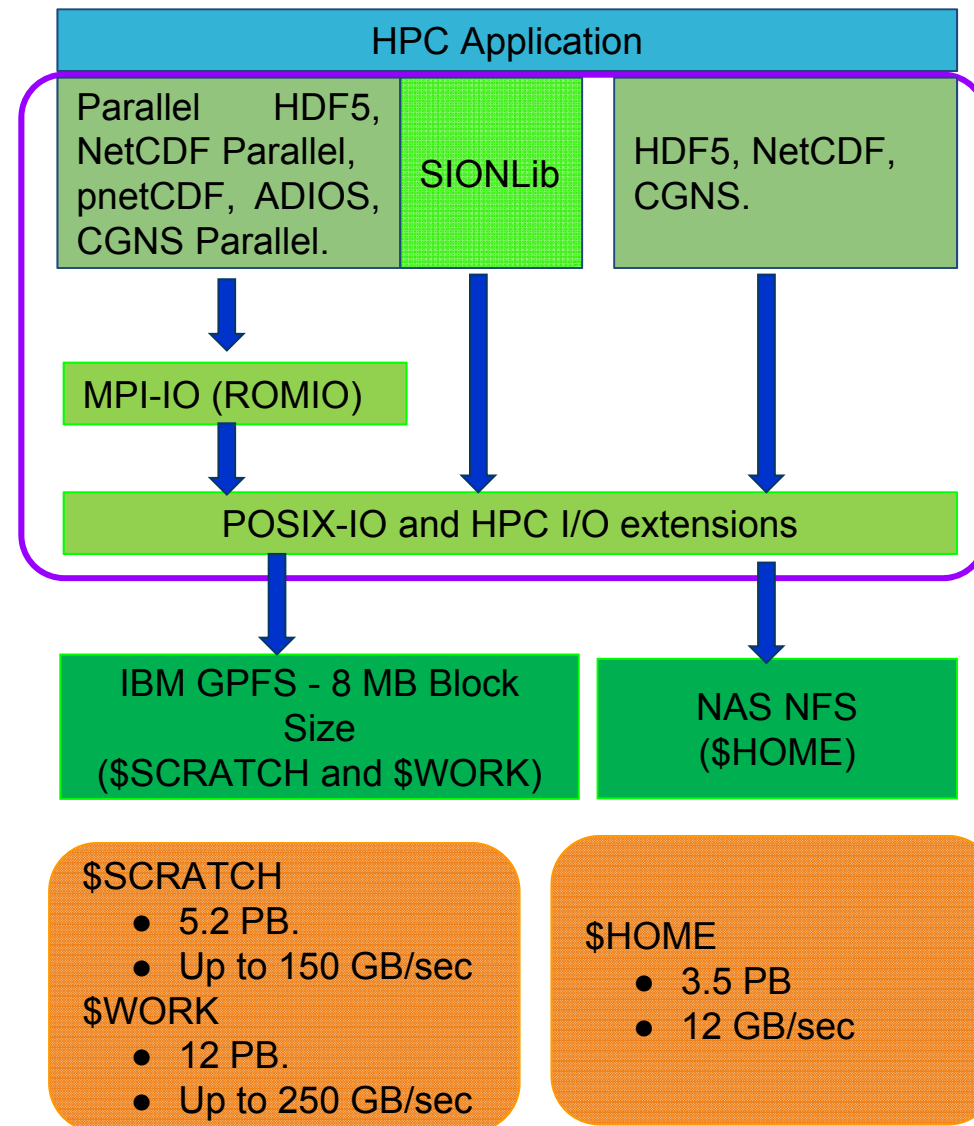
```
export LOGPATH_DARSHAN_LRZ=`darwin-logpath.sh`
```

- **Darshan - Splunk format**

```
darwin-splunk.sh $JOBID_LL $LOGPATH_DARSHAN_LRZ $LOADL_STEP_ID
```

Performance Analysis Tools:

- VAMPIRTrace (POSIX-IO, MPI-IO)
- Scalasca (MPI-IO)



- SuperMUC supercomputer
- User Projects
- Monitoring Tool
- I/O Software Stack
- I/O Analysis Tool
- **Analyzing I/O Problems**
- Conclusions

1. I/O Patterns

- Time-out error
- Slow I/O performance

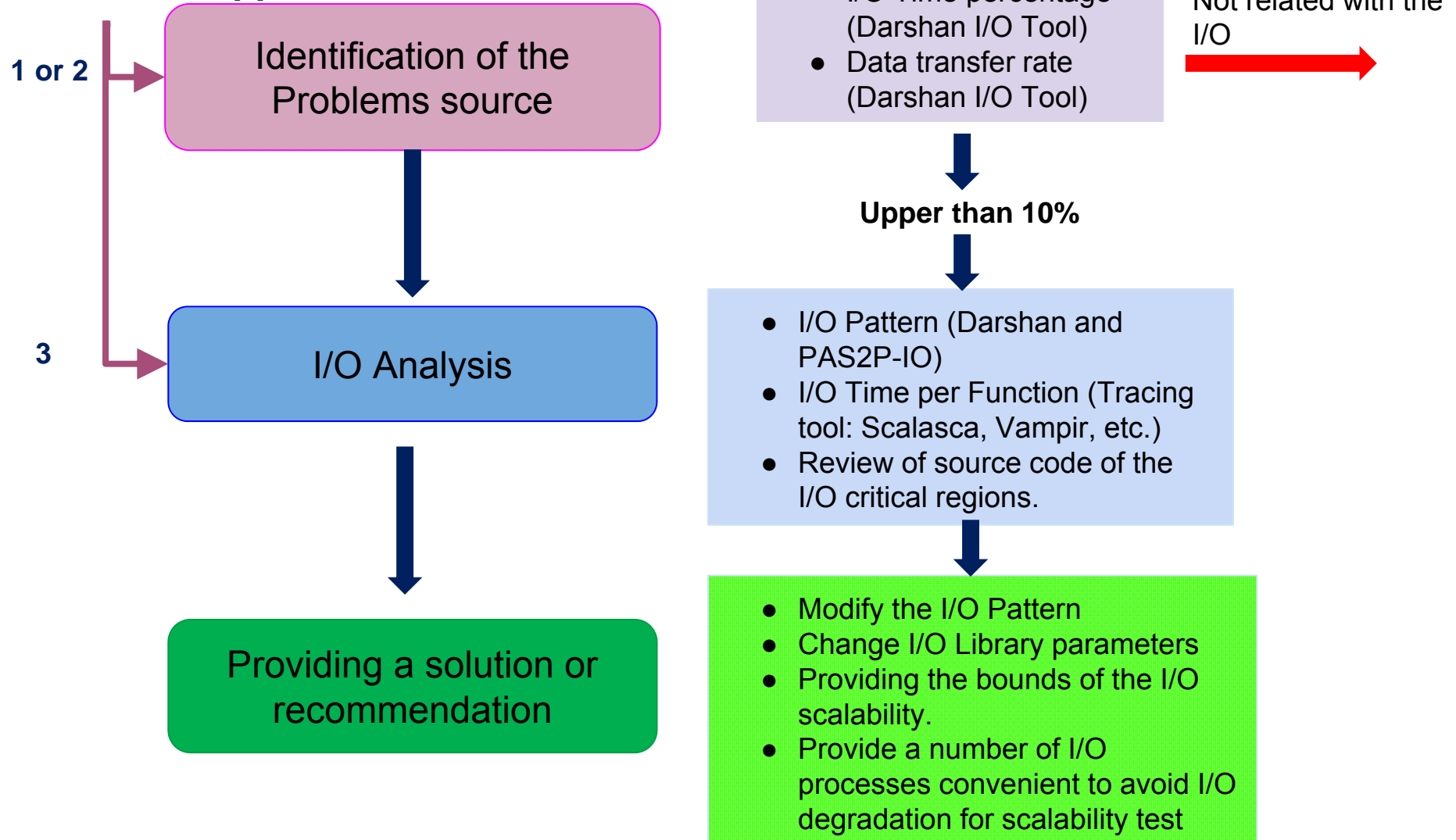
2. MPI-IO Hints configuration

- Slow I/O performance for Collective operations

3. I/O Scalability

- Requirement for large scale parallel applications

Parallel Application



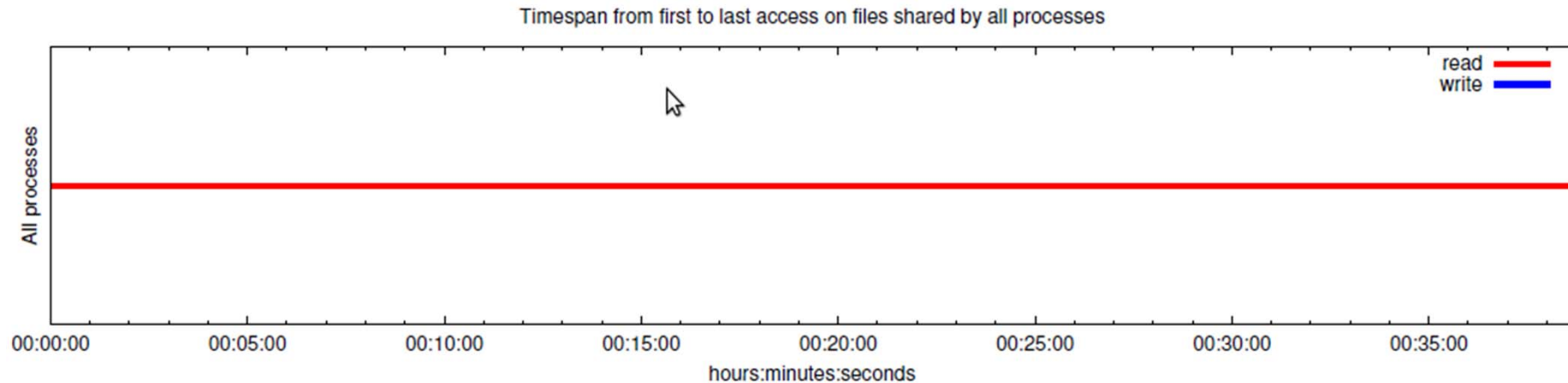
Experimental Environment:

- **An Island of thin nodes (512 compute nodes)**
- **2 processors per compute node and 8 cores per processor.**
- **Size of shared memory per node is 32 GBytes (2 GBytes per core).**
- **\$WORK filesystem (GPFS version 3.5). Up to 180 GB/s.**
- **Block size equal to 8,388,608 bytes (8MiB) and a minimum fragment size of 262,144 bytes.**
- **IBM MPI 1.3, NetCDF 4.3.3, Scalasca 2.2.2 and Darshan 2.3.1.**

- Error Report: Time-out or Slow performance.
- Master-Worker Application (Workers perform the I/O).
- NetCDF Serial (POSIX-IO).
- Darshan report the I/O time of 90%

File with more impact (Read Only)	Size
=====	
1. AK135f_5s_fwd_8MB/MZZ/Data/ordered_output.nc4	70GB
2. AK135f_5s_fwd_8MB/MXX_P_MYY/Data/ordered_output.nc4	70GB
3. AK135f_5s_fwd_8MB/MXZ_MYZ/Data/ordered_output.nc4	104GB
4. AK135f_5s_fwd_8MB/MXY_MXX_M_MYY/Data/ordered_output.nc4	104GB
5. AK135f_5s_bwd_8MB/PZ/Data/ordered_output.nc4	70GB
6. AK135f_5s_bwd_8MB/PX/Data/ordered_output.nc4	104GB
=====	
Total	522 GB

- Scalasca profile shows that the problem is related with the function `load_strain_point_interp` (I/O time of 70%) in the call `nc_getvar()` for FWD in the `readfields` routine.

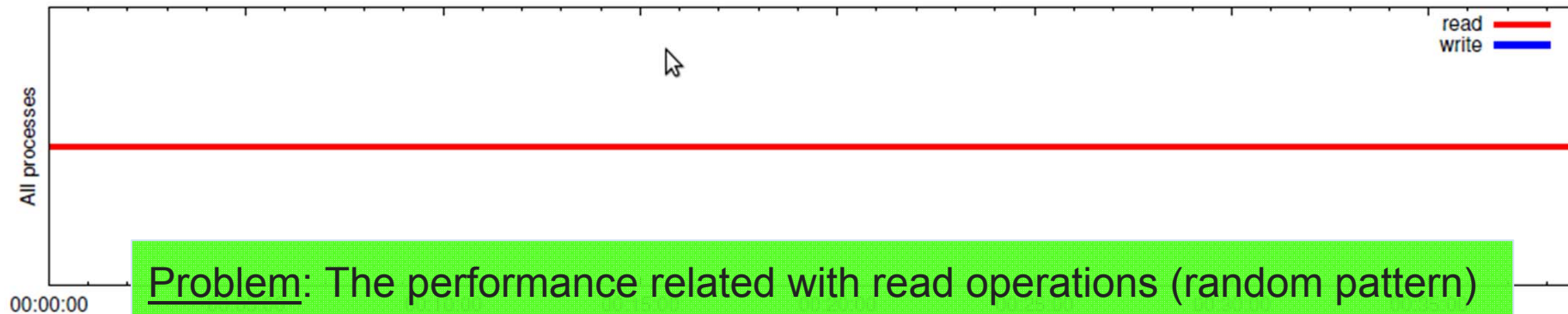


Data Transfer Per Filesystem

File System	Write		Read	
	MiB	Ratio	MiB	Ratio
/gss/scratch	0.00000	0.00000	1400386.82951	1.00000
/home	18.55555	1.00000	2.07610	0.00000

Darshan reports that the application is moving 1.5 TB or 2 TB, usually each process is moving 6GB per file for file sizes of 104 GB and 4GB for the file sizes of 70GB.

Timespan from first to last access on files shared by all processes



Problem: The performance related with read operations (random pattern) which impacts on memory utilized. The data must be read during all the execution. Each worker loads the files in memory because each worker is an I/O process.

Recommendation: The user must be avoid the random pattern (that is possible) and reduce the number of I/O processes to avoid loads several times the same file.

Darshan reports that the application is moving 1.5 TB or 2 TB, usually each process is moving 6GB per file for file sizes of 104 GB and 4GB for the file sizes of 70GB.

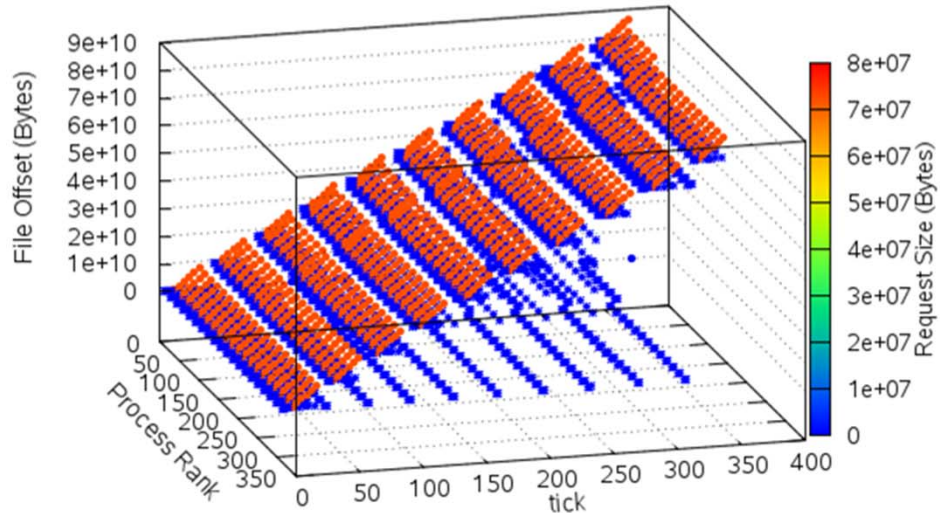
Experimental Environment:

- **An Island of thin nodes (512 compute nodes)**
- **2 processors per compute node and 8 cores per processor.**
- **Size of shared memory per node is 32 GBytes (2 GBytes per core).**
- **\$WORK filesystem (GPFS version 3.5). Up to 180 GB/s.**
- **Block size equal to 8,388,608 bytes (8MiB) and a minimum fragment size of 262,144 bytes.**
- **IBM MPI 1.3, Parallel HDF5 1.8.14, and Darshan 2.3.1.**

- A Particle-in-Cell Application.
- I/O Problem: I/O Scalability Analysis
- Parallel HDF5.

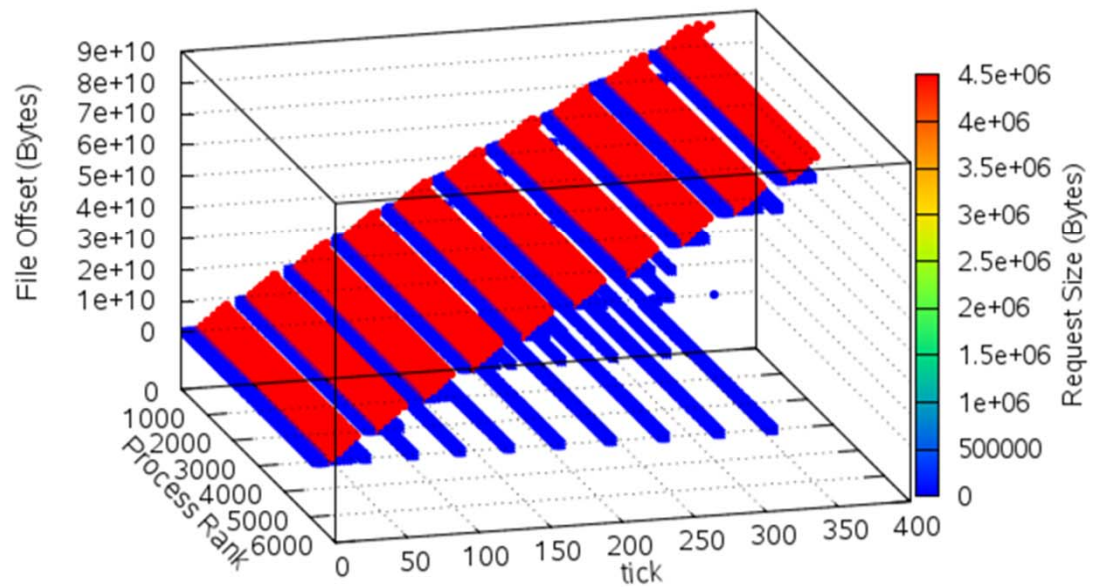
I/O Call tree

Order	MPI-I/O Operation	Data Access Aspect
1	<code>MPI_File_open</code>	
2	Once only by rank 0. <code>MPI_File_get_size</code>	
3	From seven to twelve times <code>MPI_File_read_at</code>	blocking, noncollective, explicit offset
4	Six times (once for each field) <code>MPI_File_set_view</code> <code>MPI_File_write_at_all</code> <code>MPI_File_set_view</code> <code>MPI_File_read_at</code>	blocking, collective, explicit offset blocking, noncollective, explicit offset
5	Only for the first seven I/O processes <code>MPI_File_write_at</code>	blocking, noncollective, explicit offset
6	<code>MPI_File_set_size</code>	
7	<code>MPI_File_close</code>	



Blue = read operations
Red = write operations

**I/O Phases Identified
(red with more data to transfer)**



Application I/O Parameters

I/O Parameter	Values
Global Simulation Size	(x, y, z)
Local Simulation Size	$(x_{loc} = x, y_{loc} = y, z_{loc} = \frac{z}{np})$
Compute Nodes	cn
Simulation step	st
fields	fi
writer processes	$wp = cn$
Data Size (Bytes)	ds
RequestSize(Bytes)	$rs = x_{loc} \times y_{loc} \times z_{loc} \times ds$
FileSize(Bytes)	$fz = cn \times rs \times st \times fi$
Data per st (Bytes)	$D_{st} = cn \times rs \times fi$
Data per 1 cn per st (Bytes)	$D_{cn \times st} = rs \times fi$

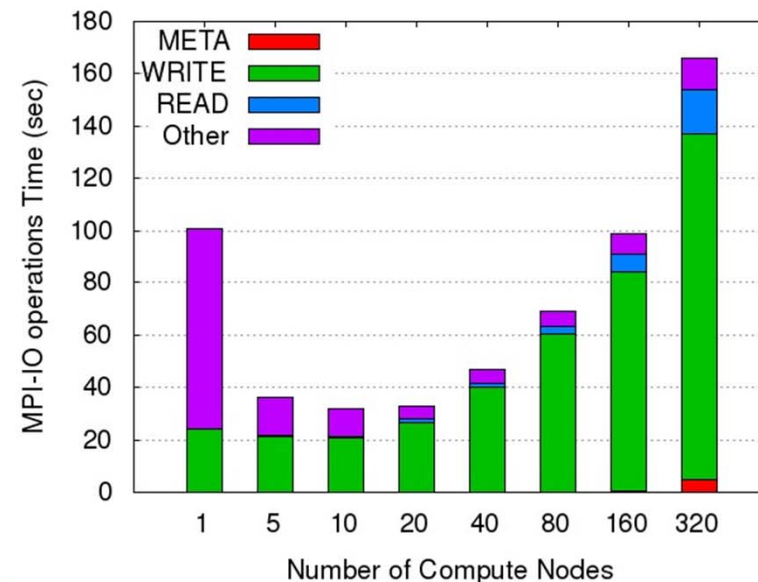
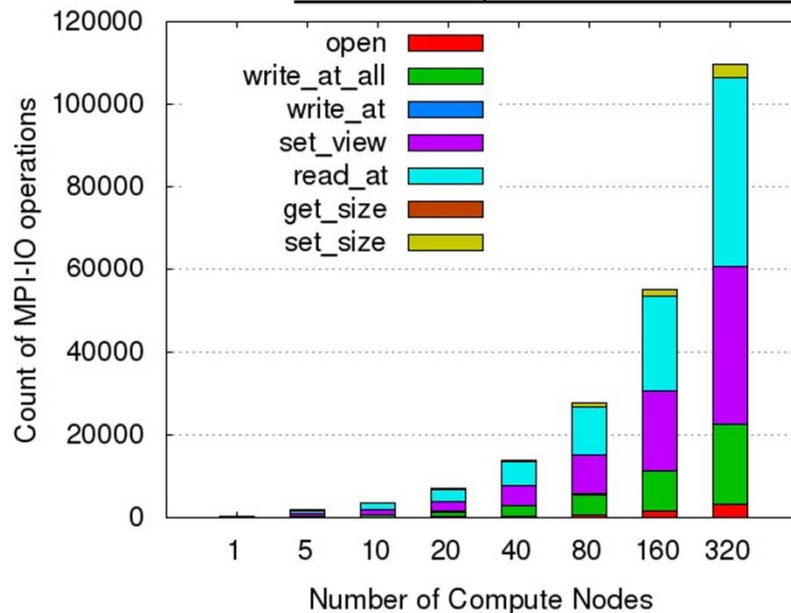
MPI-IO operations considering the I/O Parameters

I/O Operation	Count
open	$st \times cn$
write_at_all	$st \times fi \times cn$
write_at	$7 \times st$
set_view	$st \times fi \times cn \times 2$
read_at	$2 \times fi \times st \times cn + 23 \times cn$
get_size	st
set_size	$st \times cn$
close	$st \times cn$

Example: I/O Scalability (4)

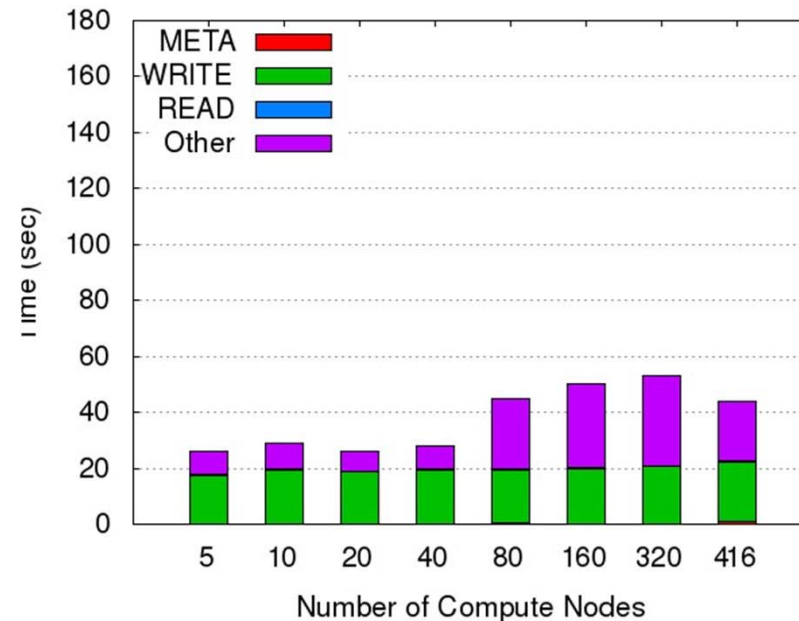
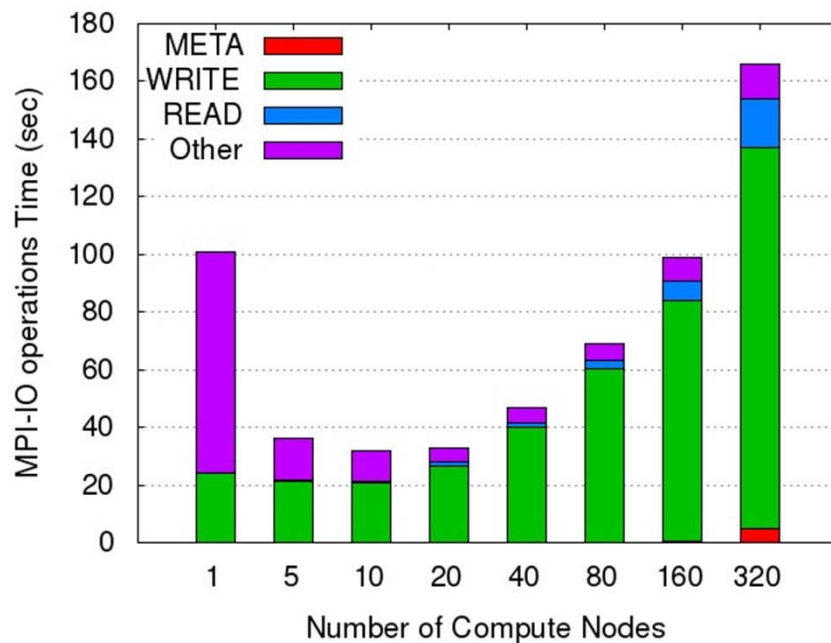
Global Simulation Size is (52,52,66560), File size = 82 GiB, 16 processes per compute node, 8.05 GiB per Simulation Step, 10 step simulation, 6 fields, 128 data size.

Compute Nodes (<i>cn</i>)	Number of Processes <i>np</i>	Local Simulation Size	Request Size <i>rs</i> (MiB)	Data per 1 <i>cn</i> per <i>st</i> $D_{cn \times st}$ (MiB)
1	16	(52,52,4160)	1373.13	8238.75
5	80	(52,52,832)	274.63	1647.75
10	160	(52,52,416)	137.31	823.88
20	320	(52,52,208)	68.66	411.94
40	640	(52,52,104)	34.33	205.97
80	1280	(52,52,52)	17.16	102.98
160	2560	(52,52,26)	8.58	51.49
320	5120	(52,52,13)	4.29	25.75



Problem: A scalability problem is produced for the strong scaling. The user writes the same amount of data and only increases the compute workload. If the number of I/O processes grows as increases the number of compute nodes then the I/O will impact in the run time.

Recommendation: reduce the number of I/O processes. As consequence the I/O Time remains constant (Right Figure)



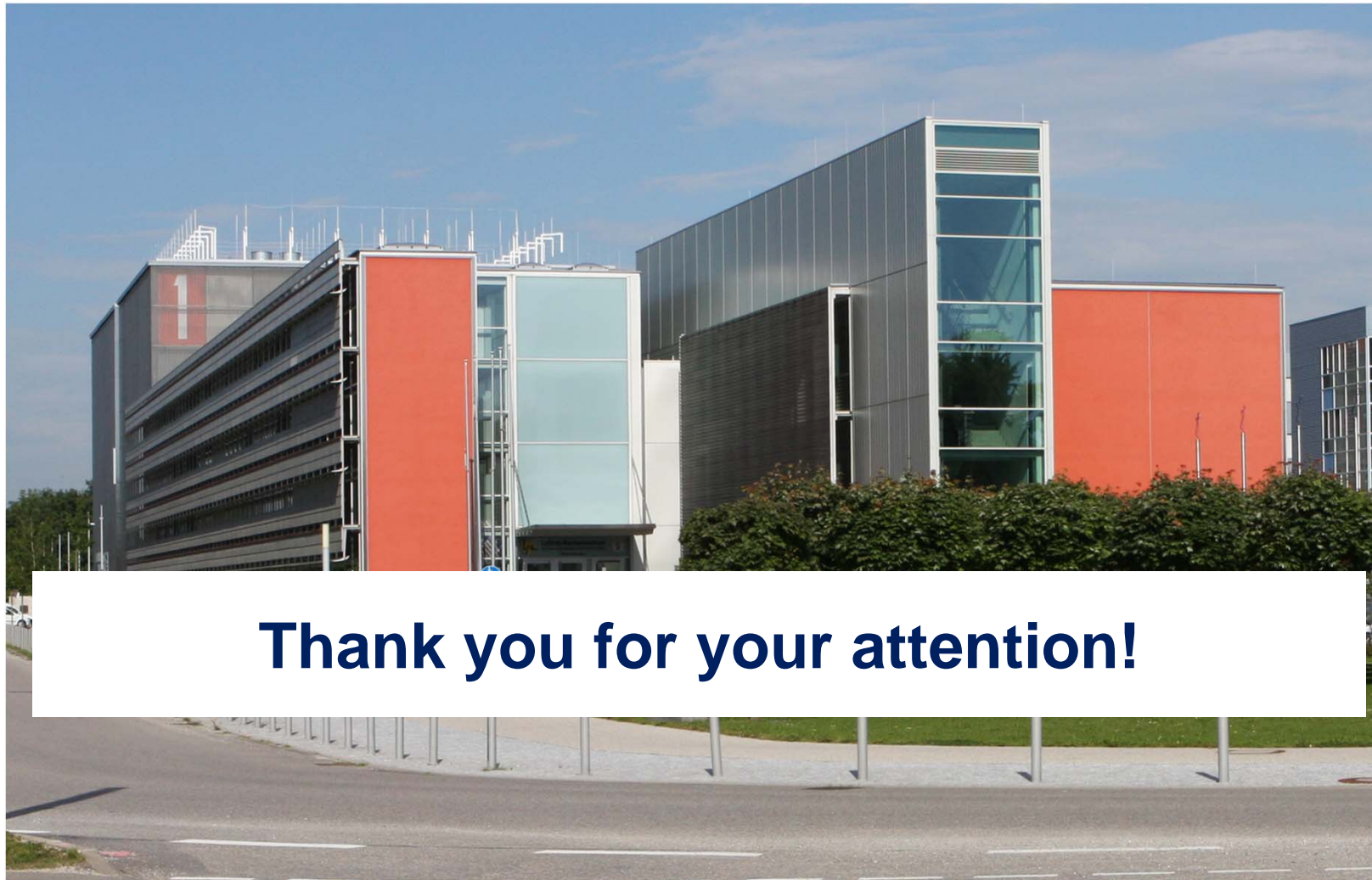
- SuperMUC supercomputer
- User Projects
- Monitoring Tool
- I/O Software Stack
- I/O Analysis Tool
- Analyzing I/O Problems
- **Conclusions**

Conclusions:

- The I/O Pattern is usually the source of slow performance.
- POSIX-IO is the library more used in SuperMUC for small and medium jobs (1 Island)
- Parallel I/O is being including at large scale (more than 2 Islands).
- I/O Aggregation has more impact on the scalability (Number of I/O processes per compute nodes).

Future Work:

- I/O Pattern and Performance Analysis at compute node level with Persyst Tool.
- Automatic I/O profiling with Darshan Tool on SuperMUC.
- I/O Scalability analysis using a formal method for detection of I/O phases and I/O operations counters.
- Integration of Darshan logs into Splunk monitoring for identifying data patterns and diagnosing problems at system level.



Thank you for your attention!