

PARALLEL I/O AND PORTABLE DATA FORMATS OPTIMIZATION AND PROFILING

29.01.2020 | SEBASTIAN LÜHRS (S.LUEHRS@FZ-JUELICH.DE)



I/O patterns

continuous

Large continuous data blocks for each individual process

striped

Pattern often found while handling multi dimensional arrays



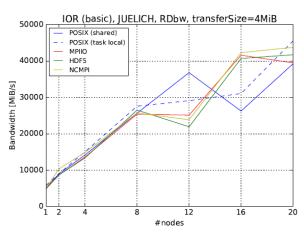


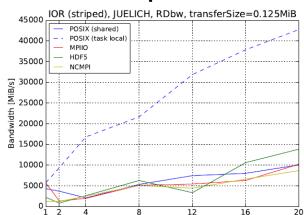
I/O pattern bandwidth

continuous

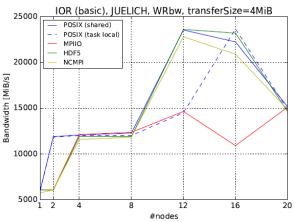
uous striped

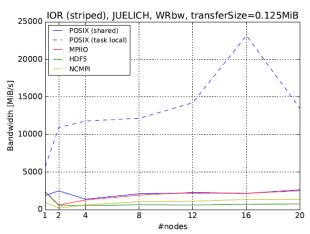
read bandwidth





write bandwidth





#nodes

Measurements on JURECA at JSC

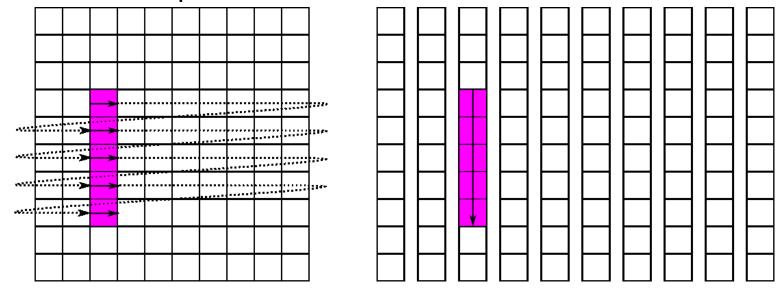
This work was supported by the Energy oriented Centre of Excellence (EoCoE), grant agreement number 676629, funded within the Horizon2020 framework of the European Union.



Performance hints

Chunking

- Contiguous datasets are stored in a single block in the file, chunked datasets are split into multiple chunks which are all stored separately in the file.
- Additional chunk cache is possible



https://www.hdfgroup.org/HDF5/doc/Advanced/Chunking/



Performance hints

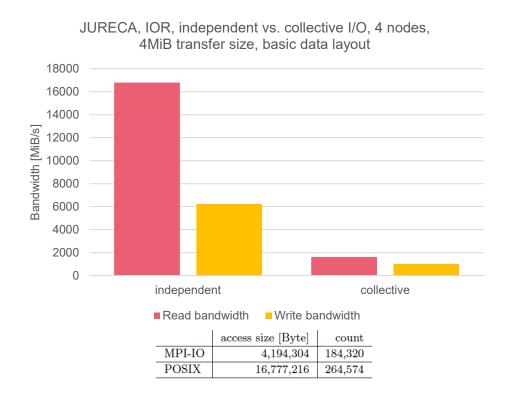
Compression

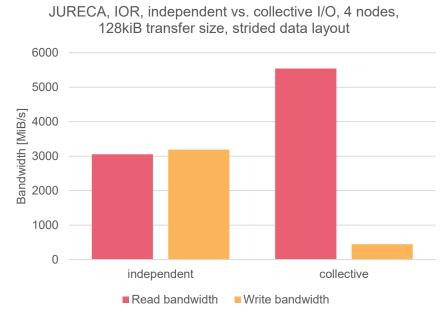
- In-transit compression can help to lower the overall datasize:
- HDF5 allows compression within a parallel, collective write commands for chunked datasets
- NetCDF4 only allows compression within serial programs (so far)
- Gzip (deflate) compression available by default (szip can be added on demand)
- Other compression techniques are available by using filters and external plugins: https://support.hdfgroup.org/services/filters.html



Collective buffering

 Collective I/O operations not always speed up the general I/O, as more data might be processed than needed





This work was supported by the Energy oriented Centre of Excellence (EoCoE), grant agreement number 676629, funded within the Horizon2020 framework of the European Union.



MPI-IO hints

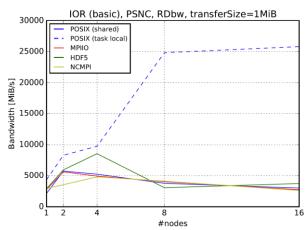
- romio cb read: Enable collective buffering (reading)
- romio cb write: Enable collective buffering (writing)
- cb buffer size: Collective buffering, buffer size
- cb nodes: Aggregator nodes
- romio ds read: Enable data sieving (reading)
- romio ds write: Enable collective buffering (writing)

```
export ROMIO HINTS=romio hints file
```

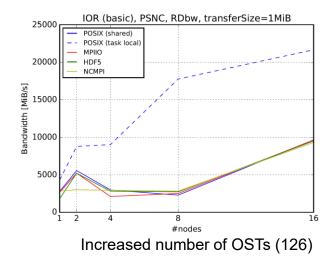


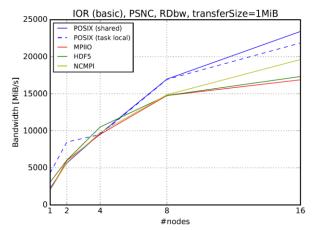
Filesystem specific options

 On Lustre filesystems the user can influence the striping size and the number of involved object storage targets



Default number of OSTs (12) and default strip-size setting (1MiB)





Increased stripe size to align with the individual amount of data per process (256MiB)

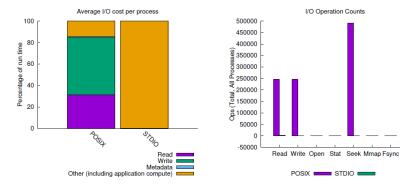
Measurements on Eagle at PSNC



Profiling with Darshan

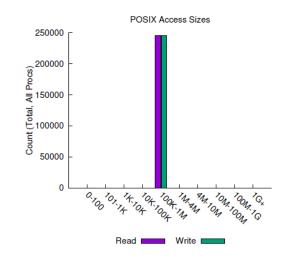
- I/O profiling tool for parallel applications
 - http://www.mcs.anl.gov/research/projects/darshan/
- Integration by using LD_PRELOAD:
 - LD PRELOAD=.../lib/libdarshan.so

I/O performance *estimate* (at the POSIX layer): transferred 37431 MiB at 6692.22 MiB/s I/O performance *estimate* (at the STDIO layer): transferred 0.0 MiB at 5.27 MiB/s





Profiling with Darshan



Most Common Access Sizes (POSIX or MPI-IO)

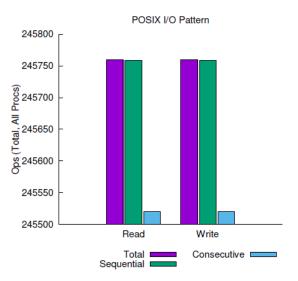
(1 com of mil 1 co)					
	access size	count			
POSIX	131072	491520			

File Count Summary (estimated by POSIX I/O access offsets)

•	,		
type	number of files	avg. size	max size
total opened	4	7.6G	30G
read-only files	1	711	711
write-only files	2	1.7K	3.2K
read/write files	1	30G	30G
created files	3	11G	30G



Profiling with Darshan



sequential: An I/O op issued at an offset greater than where the previous I/O op ended. *consecutive*: An I/O op issued at the offset immediately following the end of the previous I/O op.

Variance in Shared Files (POSIX and STDIO)

File	Processes	Fastest		Slowest		σ			
Suffix		Rank	Time	Bytes	Rank	Time	Bytes	Time	Bytes
ehrs/IOR/2_1	48	35	7.507493	1.3G	33	9.180811	1.3G	0.397	0
or_input.cfg	48	32	0.003404	711	2	0.006366	711	0	0
<stdout></stdout>	48	1	0.000000	0	0	0.000392	3.2K	0	455
<stderr></stderr>	48	1	0.000000	0	0	0.000014	119	0	17



Darshan: Usage example on JURECA

- Load module
 - module load darshan-runtime
- Tell srun to use Darshan (in submit script)
 - LD_PRELOAD=\$EBROOTDARSHANMINRUNTIME/lib/libdarshan.so \ DARSHAN_LOG_PATH=/path/to/your/logdir \ srun/executable
- Analyse output
 - module load darshan-util
 - darshan-job-summary.pl <logfile>.darshan
 - xpdf <logfile>.darshan.pdf



Interested in individual HPC support? Apply for PRACE support activities!

Preparatory Access

- Prepare, scale and optimise your application on the European high class HPC systems
- Receive individual PRACE expert support
- http://www.prace-ri.eu/prace-preparatoryaccess/

SHAPE

- PRACE's SME HPC Adoption
 Programme in Europe
- Free support for European businesses to adopt highperformance computing
- http://www.prace-ri.eu/prace-shapeprogramme/

Interested in PRACE regular calls, white papers, best practise guides and much more:

>http://www.prace-ri.eu