

# Parallel HDF5 and compression filters with synchrotron scattering data

Zdenek Matej, Andreas Mattsson MAX IV Laboratory, Lund University



#### Ingredients

- parallel HDF5 with MPI
- image like synchrotron data
- compression





### Parallel HDF5 with MPI

- message passing interface (MPI) is a standard for parallel computing architecture; invented in 1991-1992; popular at high-performancecomputing (HPC) clusters
- Sophie Servan (DESY) et al., Technical-Workshop-Survey (2020): 9/10 facilities have SLURM HPC cluster <a href="https://github.com/ExPaNDS-eu/ExPaNDS/blob/master/WP4/20201009-Technical-Workshop-Survey.pdf">https://github.com/ExPaNDS-eu/ExPaNDS/blob/master/WP4/20201009-Technical-Workshop-Survey.pdf</a>
- "reference" HDF5 library implementation
  - a complex sw and have numerous technical limitations one should be aware of
  - MPI enables distributed multi-task application can effectively read/write data to a single HDF5 file and even dataset
  - well defined parallel HDF5 library and compatible software (h5py/mpi4py) are available off-the-shelf in HPC sw distributions (e.g. <u>EasyBuild</u>)
  - parallel HDF5 is an "traditional" feature, i.e. 1.8
  - **sw**: savu (Diamond), PtyPy both are mpi4py & h5py
- alternatives with serial HDF5
  - virtual datasets distributed over multiple files
  - direct chunk write/read
  - disadvantage: 1.10 required, some sw cannot read such data



#### Image like synchrotron data

- general synchrotron data have complex structure/format
- "detector" data (99% volume of all data): usually 2D or 3D datasets
- here focus on performance
  - example JungFRAU detector (PSI)
    - 4M pixels (2 bytes depth) x 2kHz
    - 16 GB/s uncompressed, bslz compression factor ~ 4x, finally 4 GB/s compressed



# pHDF5, image like data, no compression

- tools: h5perf\_parallel, ior (-a HDF5), hdf5\_pwrite3dc, h5py-snippets/gist
- hdf5 pwrite3dc is deeply inspired by code from *Timoty Brown* on <u>Stackoverflow</u>



European HDF Users Group Summer 2021

#### pHDF5, image like data, no compression

typical figures

	compression	write [GB/s]	read [GB/s]	ntasks	comment
ior	no	5.7	4.1	8	1 x FDR
pwrite	no	5.8	4.8	8	1 x FDR
pwrite	no	10.9	8.2	8	2 x FDR
h5py-parallel	no	4.5	-	16	1 x FDR
h5py-serial (dask)	no	-	5.4	16	1 x FDR

• Nice ! It works. Write/read rates limited mainly by bandwidth to storage.



## pHDF5, image like data, compression filter

typical figures

	compression	write [GB/s]	read [GB/s]	ntasks (omp)	comment
pwrite	no	5.8	4.8	8	1 x FDR
pread-eiger	bslz4	-	<b>23.8</b> √	20(2)	compression rate: 8.1
pwrite-bslz4	bslz4	<b>4.9</b> (?)	-	24(2)	compression rate: 8.1

- Table: raw i.e. uncopressed data rates
- 4.9 / 8.1 ~ 0.6 GB/s storage is likely not a bottleneck
- bslz4 compression is very effective, in particular at least 2x faster
- uncompressed write faster than compressed, read OK
- with multiple files, VDS, in-memory bslz4 compression and direct chunk write better performance can be achived
- parallel HDF5 is great but writing data with compression-filters is not so shiny (?)