Understanding Burst Buffers with NetCDF-Bench

Eugen Betke, Julian Kunkel

Research Group German Climate Computing Center

SIG-IO-UK Workshop 2018-06-06

Eugen Betke, Julian Kunkel Understanding Burst Buffers with NetCDF-Bench

イロト イポト イヨト イヨト

= non

Table Of Content

Introduction

2 Benchmarks

- IOR
- NetCDF-Bench

3 Evaluation

- Kove XPD
- DNN IME240

4 Summary

イロン 不得 とくほ とくほう

= nar

Introduction

Capturing current performance values

- How much performance can we achieve on our system?
- Does performance changes over time?
- Understanding I/O performance
 - How different access patterns or features influence the I/O performance?
 - What are the most efficient settings?
- HPC procurement
 - Which technologies are beneficial for our applications / next generation HPC?

Climate data

Large scale climate applications can access **terabytes of data** on parallel file systems in a single run. Mostly, **NetCDF4** file format is used to access data.



Figure: Regular grid [2] (legacy)

Figure: Icosahedral grid [2] (new)

4 B K 4 B K B

Data arrangement in a file depends on earth grid used by application. Most large experiments on Mistral (HPC of DKRZ) are using regular and icosahedral grids.

Is there any benchmark, that can **approximate the behaviour** and provide a **realistic performance** values?

Under the hood, NetCDF4 uses HDF5. Maybe, IOR with HDF5-API can do the job.

Eugen Betke, Julian Kunkel Understanding Burst Buffers with NetCDF-Bench

Table Of Content

1 Introduction

2 Benchmarks

- IOR
- NetCDF-Bench

3 Evaluation

- Kove XPD
- DNN IME240

4 Summary

イロト イポト イヨト イヨト

= nar

Running IOR with HDF5-API

mpiexec -n 2 ./ior -s 5 -t 5242880 -b 5242880 -o ./testfile.h5 -i 1 -v -k -f tests/HDF5

IOR example (-s - number of segments; -b - block size; -t - transfer size)

1 IOR START 2 reorderTasksConstant=1 3 taskPerNodeOffset=1 4 api=HDF5 5 RUN 6 IOR STOP

Configuration file

- reorderTasksConstant reorders tasks by a constant node offset for writing/reading neighbor's data from different nodes [0=FALSE]
- **taskPerNodeOffset** for read tests. Use with -C & -Z options. [1] With reorderTasks, constant N. With re-ordertasksrandom, >= N

Insight into HDF5 testfile

mpiexec —n 2 ./ior —s 5 —t 5242880 —b 5242880 —o ./testfile.h5 —i 1 —v —k —f tests/HDF5

```
~ h5dump — A testfile.h5
   HDF5 "testfile.nc" {
     GROUP "/" {
       DATASET "Dataset-0000.0000" {
         DATATYPE H5T STD I64LE
         DATASPACE SIMPLE \{(1310720)/(1310720)\}
       DATASET "Dataset-0000.0001" {
         DATATYPE H5T STD I64LE
         DATASPACE SIMPLE { ( 1310720 ) / ( 1310720 ) }
11
       DATASET "Dataset-0000.0002" {
13
         DATATYPE H5T STD I64LE
14
         DATASPACE SIMPLE { ( 1310720 ) / ( 1310720 ) }
16
       DATASET "Dataset-0000.0003" {
17
         DATATYPE H5T STD I64LE
18
         DATASPACE SIMPLE { ( 1310720 ) / ( 1310720 ) }
19
       DATASET "Dataset-0000.0004" {
         DATATYPE H5T STD I64LE
         DATASPACE SIMPLE { ( 1310720 ) / ( 1310720 ) }
24
25
26
```

IOR constraints

- Segments are stored in separate datasets
- transfer size = block size
- HDF5 features can not be tested

Conclusion

- No suitable access patterns, due to many metadata operations.
- Usually, transfer size in climate applicatons is smaller than block size. Data is written not as single block, but piecewise.

NetCDF Performance Benchmark Tool

- Developed at DKRZ
- Written in C
- Published on GitHub under LGPL License
 - https://github.com/joobog/netcdf-bench
- Mimics scientific data
 - Data is written in time steps
- Supports different NetCDF4 features
 - Compression, chunking, independent/collective I/O, fill value, unlimited dimensions



Figure: Data geometry (t:x:y:z)=(3:6:4:3)

NetCDF-Bench output

mpiexec -n 2 benchtool ----nn=2 ----ppn=1 -b=1:1024:256:5 -d=5:2048:256:5 -r -w -f=./testfile.nc

NetCDF-Bench example (-d: data geometry; -b: block geometry)

1	Benchtool (datatype:	int)				
2	Data geometry (t:x:y	:z x sizeof(type))	5:2048:25	6:5 x 4 bytes		
3	Block geometry (t:x:	y:z x sizeof(type))	1:1024:25	6:5 x 4 bytes		
4	Datasize		5	2428800 bytes	(52.4 MB)	
5	Blocksize			5242880 bytes	(5.2 MB)	
6	I/O Access		inde	pendent		
7	Storage		cor	itiguous		
8	File length			fixed		
9	File value			no		
10				min	avg	max
11	benchmark:write	Open time		0.1929476971	0.1929488905	0.1929500839 secs
12	benchmark:write	I/O time		0.0629276498	0.0629278460	0.0629280421 secs
13	benchmark:write	Close time		0.0829069570	0.0829072705	0.0829075840 secs
14	benchmark:write	I/O Performance (w/o	open/close)	794.5583288194	794.5608056386	794.5632824579 MiB/s
15	benchmark:write	I/O Performance		147.5863619891	147.5866598368	147.5869576846 MiB/s
16				min	avg	max
17	benchmark : read	Open time		0.0039127499	0.0039127710	0.0039127921 secs
18	benchmark : read	I/O time		0.0124143478	0.0124145120	0.0124146761 secs
19	benchmark : read	Close time		0.0004221359	0.0004225459	0.0004229560 secs
20	benchmark : read	I/O Performance (w/o	open/close)	4027.4912891123	4027.5445416068	4027.5977941012 MiB/s
21	benchmark : read	I/O Performance		2985.0575402544	2985.1051135732	2985.1526868921 MiB/s

Eugen Betke, Julian Kunkel Understanding Burst Buffers with NetCDF-Bench

Insight into NetCDF-Bench test file

mpiexec -n 2 benchtool ----nn=2 ----ppn=1 -b=1:1024:256:5 -d=5:2048:256:5 -f=./testfile.nc

NetCDF-Bench example (-d: data geometry);-b: block geometry)

1	~ ncdump — h testfile.nc
2	netcdf testfile {
3	dimensions:
4	time = 5 ;
5	$\dim_1 = 2048$;
6	$\dim_2 = 256$;
7	$\dim_3 = 5$;
8	variables:
9	int data(time, dim_1, dim_2, dim_3) ;
0.	}

1 ~ h5ls testfile.nc 2 data Dataset {5, 2048, 256, 5} 3 dim_1 Dataset {2048} 4 dim_2 Dataset {26} 5 dim_3 Dataset {5} 6 time Dataset {5}

Listing content with an HDF5 tool

NetCDF header

Kove XPD DNN IME240

Table Of Content

1 Introduction

2 Benchmarks

- IOR
- NetCDF-Bench

3 Evaluation

- Kove XPD
- DNN IME240

4 Summary

イロト イポト イヨト イヨト

= nar

Kove[®] XPD[®] Specifications

- In-memory storage offers high troughput and low latency
- The Kove XPD is an example of a scalable in-memory storage system



Figure: Kove[®] XPD[®] L3 [1]

Capacity Bandwidth IOPS

Response Time Data Protection 64GB - 1.5TB per 1U unit, mesh connectable for larger capacities 27+ GB/sec (12 μ s latency and 100 ns variance), sustained 43+ million (3.0 μ s latency and 100 ns variance), sustained <2.5 μ s, random, consistent, sustainable for any time duration - Multiple modes for de-staging data from RAM onto persistent media

- Integrated instant copying
- Integrated monitoring logic for system health management and UPS

Storage Access

The XPD offers various storage access methods

- Block device
 - Can be used as ordinary block device, e.g., with EXT4
 - Needs (slow) cluster file systems on top for shared access
- Remote memory
 - L4 network cache pooled memory
 - Memory allocation functions can be replaced by using LD_PRELOAD.
- Direct access
 - Native I/O library: KDSA-API (Kove Direct Storage Access)
 - Fastest access to XPDs

But none of these access methods provides a shared storage

Our MPI-IO Driver

• The MPI-IO Driver¹

- is built on top of the KDSA-API (Kove Direct Storage Access)
- implements many MPI-IO functions
- is implemented as a shared library and usable with any MPI
- can be selected at startup of an application using LD_PRELOAD with the shared library
- checks the file name for the prefix "xpd:"
 - without the prefix, they route the accesses to the underlying MPI
 - files can be selectively stored on XPD volumes
- There is **no cache** on client side (needed)!

¹ Repository: http://github.com/JulianKunkel/XPD-MPIIO-driver

Test-Setup

Cooley

- Visualization cluster of Mira on ALCF
- 3x XPD with 14 FDR connections in total
- Connected to a GPFS system

Mistral

- HPC deployed at DKRZ
- 3000 compute nodes
- FDR interconnect
- Lustre storage system with 54 PByte
 - Peak transfer rate 450 GiB/s



Figure: Cooley XPDs (Kove)

イロン 人間 とくほ とくほとう

= nar

Collective vs. Indepenent vs. Ind.-Chunked (NetCDF-Bench)





ъ.

4 4 5 4 5 5 4 5 5

Observations:

XPDs

- show always best performance
- seem to be insensitive to collective, indendent and independent-chunked I/O

DDN IME240 Specifications



CPUDual Intel E5 Series ProcessorMemory8x 16GB, DDR4-2400 RAMFabric Ports2x InfiniBandTM EDR/FDRCapacityup to 23x 1.2TB NVMe SSD drivesData ProtectionErasure coding

Access modes

- IME-native
 - Native I/O library provides the fastest and direct access to IMEs
- IME-FUSE
 - Stores data on IME storage and metadata on Lustre metadata servers
- IME-Lustre
 - A Lustre-backed IME storage

4 D F 4 B F 4 B F 4 B F

ъ.

PMPI-IO with IME support (in development)

Test cluster

Compute Nodes (1-10x)

CPU 2x Intel(R) Xeon(R) CPU E5-2650 v2 @ 2.60GHz RAM DDR3 8x8GB Network 1x EDR mlx 5.0 single ported

Storage (4x DDN IME240)

CPU Intel(R) Xeon(R) CPU E5-2650 v4 @ 2.20GHz RAM DDR4 8x16GB

Network 2x EDR mlx 5.0 single ported



Eugen Betke, Julian Kunkel Understanding Burst Buffers with NetCDF-Bench

Testing NetCDF4 features: Test setup

- Each test configuration is repeated 10 times.
- All experiments are conducted with block sizes 16, 100, 1024, and 10240 KiB.
- **IME-Lustre**, **IME-FUSE** cache data on IME. The data is synchronized with Lustre storage.

Comparison: IOR-MPIIO vs. NetCDF-Bench on IME-Lustre









Blocksize in KiB: -- 16 -- 100 -- 1024 -- 10240

Figure: NetCDF-Bench

Figure: IOR sequential

Figure: IOR random

Eugen Betke, Julian Kunkel

Understanding Burst Buffers with NetCDF-Bench

Comparison: IOR-MPIIO vs. NetCDF-Bench on IME-FUSE





Blocksize in KiB: -- 16 -- 100 -- 1024 -- 10240



MPIIO

MPIIO

Blocksize in KiB: -- 16 -- 100 -- 1024 -- 10240

Figure: NetCDF-Bench

Figure: IOR sequential

Figure: IOR random

Eugen Betke, Julian Kunkel

Understanding Burst Buffers with NetCDF-Bench

Testing NetCDF4 features with NN=10 and PPN=8 [1]



Blocksize in KiB: 📋 16 🖨 100 🚔 1024 🛱 10240

Observations:

- Large variability for large block sizes
- Independent I/O with chunking works best for large block sizes

ъ.

b) A (B) b)

Figure: Linear scale

Eugen Betke, Julian Kunkel Understanding Burst Buffers with NetCDF-Bench

Testing NetCDF4 features with NN=10 and PPN=8 [2]



Observations:

 Collective I/O without chunking or independent I/O with chunking works best for small block sizes

4.35

ъ

Figure: Logarithmic scale

Table Of Content

Introduction

2 Benchmarks

- IOR
- NetCDF-Bench

3 Evaluation

- Kove XPD
- DNN IME240



イロト イポト イヨト イヨト

= nar

Scalability (direct access over native interfaces)







Figure: XPD read

Figure: XPD write

Figure: IME read/write

Eugen Betke, Julian Kunkel Understanding Burst Buffers with NetCDF-Bench

Summary

IOR is

- good for pushing systems to the limits
- bad for
 - generating access patterns of scientific applications
 - testing special HDF5 features
- NetCDF-Bench mimics data access of scientific applications
- Burst Buffers
 - Kove XPD
 - MPI-IO driver with XPD support makes NetCDF4 performance insensitive to collective/independent I/O and chunking
 - DDN IME
 - Large performance gap between accessing small and large block sizes
 - No information if data is located in IME or on PFS
 - Read performance and variability depends on data location
 - Applications require tuning for better performance