



# IO Performance Evaluation on Massive Data Processing at NCI High Performance Computing Platform

Rui Yang, Ben Evans

National Computational Infrastructure (NCI)  
Australia

- NCI is the peak National Computational Centre for Research in Australia
  - ANU Supercomputer Facility (1987-2011)
  - APAC (Australian Partnership for Advanced Computing) (2000-2007)
  - NCI (2007-)
- NCI: move from academic HPC centre to a full national high performance computing centre:
  - Partnership includes ANU, Bureau of Meteorology, CSIRO, Geoscience Australia
  - Particular Focus on Climate, Weather, Earth Systems, Environment, Geophysics & Water Management.
- Two key drivers around the work in this talk
  - HPC Scaling and Optimisation (supported by Fujitsu)
  - High Performance Data analysis
- Application areas:
  - Weather forecasting for the Bureau of Meteorology operations (see Tim Pugh talk)
  - Weather, Climate, Earth systems, water mgt for research (gov & unis) (See Marshall Ward talk)
  - Satellite Earth Observation data with Geoscience Australia and Bureau of Meteorology
  - Geophysics applications for Geoscience Australia , state surveys, and uni's



Raijin is a Fujitsu Primergy high-performance, distributed-memory cluster installed in 2013. It comprises:

- 57,472 cores (Intel Xeon Sandy Bridge technology, 2.6 GHz) in 3592 compute nodes
- 160 TBytes of main memory
- Infiniband FDR interconnect
- 10 PBytes of usable fast filesystem (for short-term scratch space).

- The National Computational Infrastructure has now co-located a priority set of over 10+ PetaBytes of national data collections for our priority areas.
- The facility provides an integrated high-performance computational and storage platform, or a High Performance Data (HPD) platform, to serve and analyse the massive amounts of data across the spectrum of environmental collections – in particular from the climate, environmental and geoscientific domains.
- The data is managed to support the government agencies, major academic research communities and collaborating overseas organizations.
- By co-locating the vast data collections with high performance computing environments and harmonizing these large valuable data assets, new opportunities have arisen for Data-Intensive interdisciplinary science at scales and resolutions not hitherto possible.
- Note, there are a lot of elements of this that are not touched on in this talk.
- As well as addressing management and performance issues, our work was to also help transform the communities to using better ways of processing, managing and analysing data.

## User Applications

local or remote, serial or parallel, file formats

## High-level I/O Library

**GDAL:** Geospatial Data Abstraction Library  
converting among many formats like GeoTIFF, NetCDF

**NetCDF4:** Network Common Data Form  
Simplified data abstraction, self-describing

**HDF5:** Hierarchical Data Formats  
chunking, compression

## I/O Middleware

**MPIIO:** higher level of data abstraction  
MPI hints

**POSIX I/O:** full and low-level control of serial I/O  
transfer regions of bytes

## Parallel File System

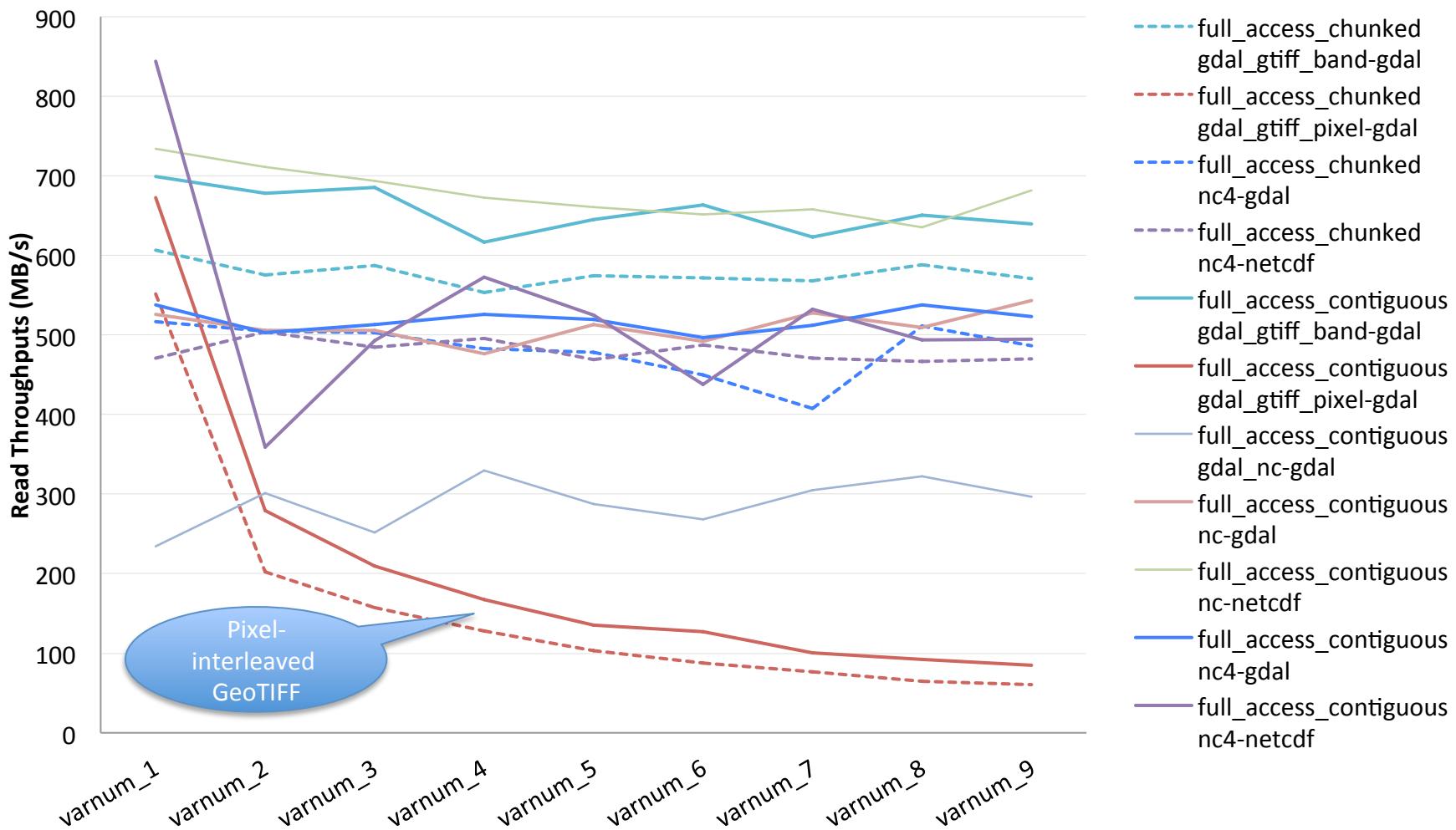
**Lustre:** Parallel access to multiple OSTs

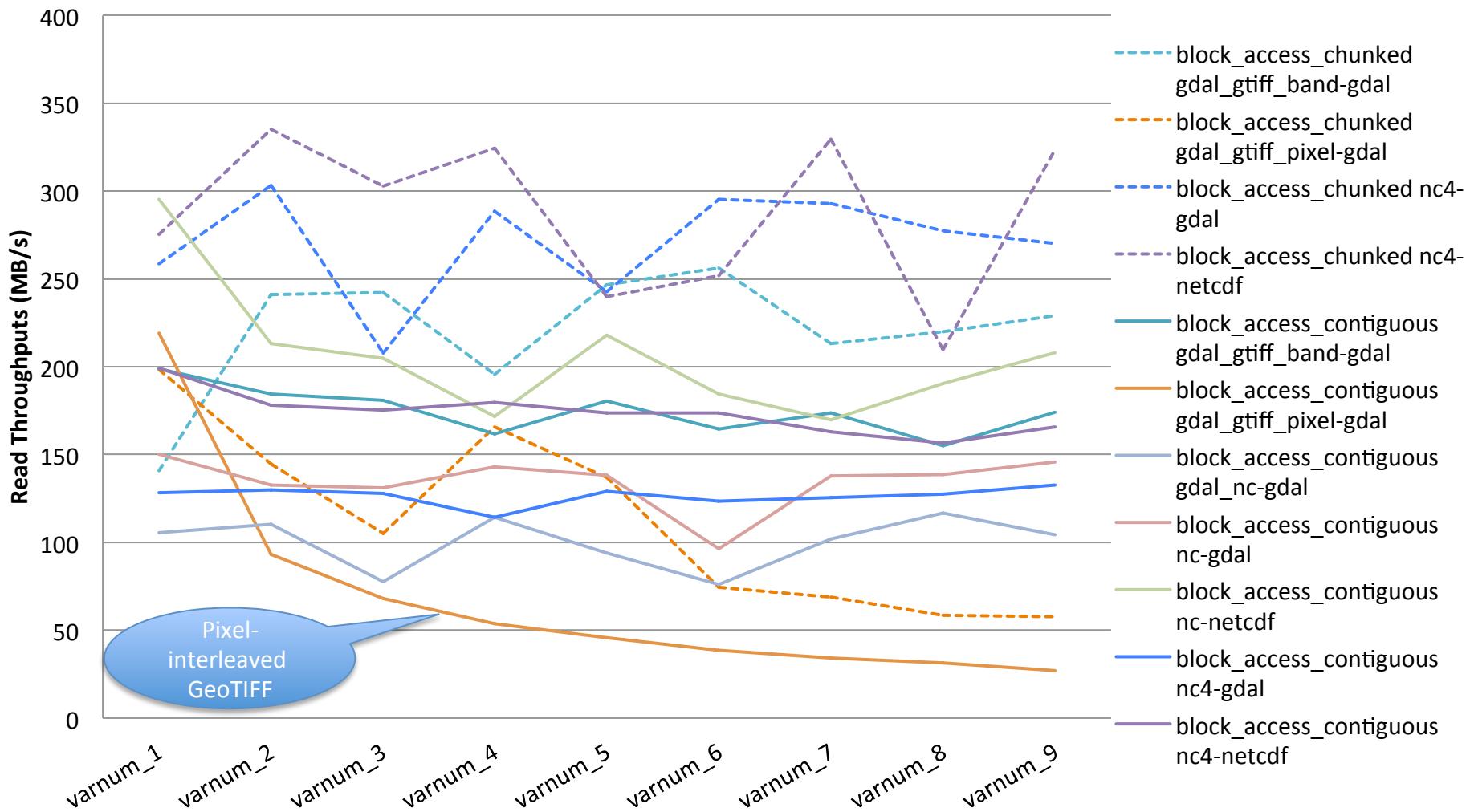
Metrics	Serial IO	Parallel IO
<b>User application</b>		
Transfer size	✓	✓
File size	✓	✓
Subset selection	✓	✓
Concurrency	N/A	✓
Access remote DAP server	✓	
<b>IO interfaces</b>		
NetCDF4/HDF5	✓	✓
MPIIO	N/A	✓
POSIX	✓	✓
GDAL/GeoTIFF	✓	
GDAL/NetCDF(4) Classic	✓	

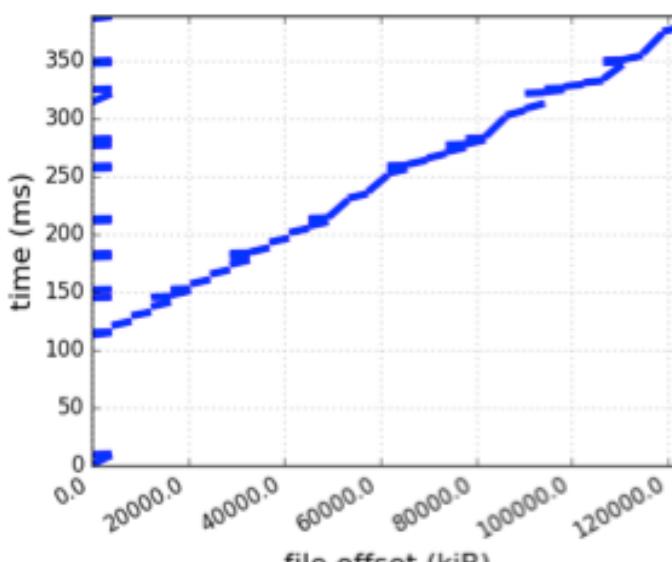
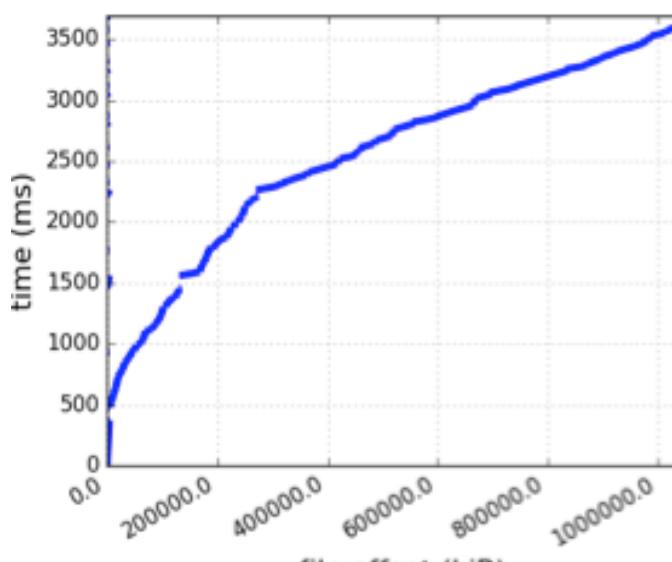
Metrics	Serial IO	Parallel IO
<b>NetCDF/HDF5</b>		
Chunk pattern	✓	✓
Chunk cache	✓	✓
Compression	✓	✓
<b>MPIIO</b>		
Independent & Collective	N/A	✓
Collective buffering	N/A	✓
Data sieving		✓
<b>Lustre file system</b>		
Stripe count	✓	✓
Stripe size	✓	✓
<b>IO profiling &amp; tracing</b>		
total	14	16

# SERIAL IO

File Formats	pixel-interleaved GeoTIFF band-interleaved GeoTIFF NetCDF Classic NetCDF4
Total Bands/Variables	9
Storage Layouts	contiguous, chunked(chunk/tiling size=640×640)
IO Libraries	GDAL, NetCDF
Access patterns	full file access, block access (hyperslab size= 2560×2560)

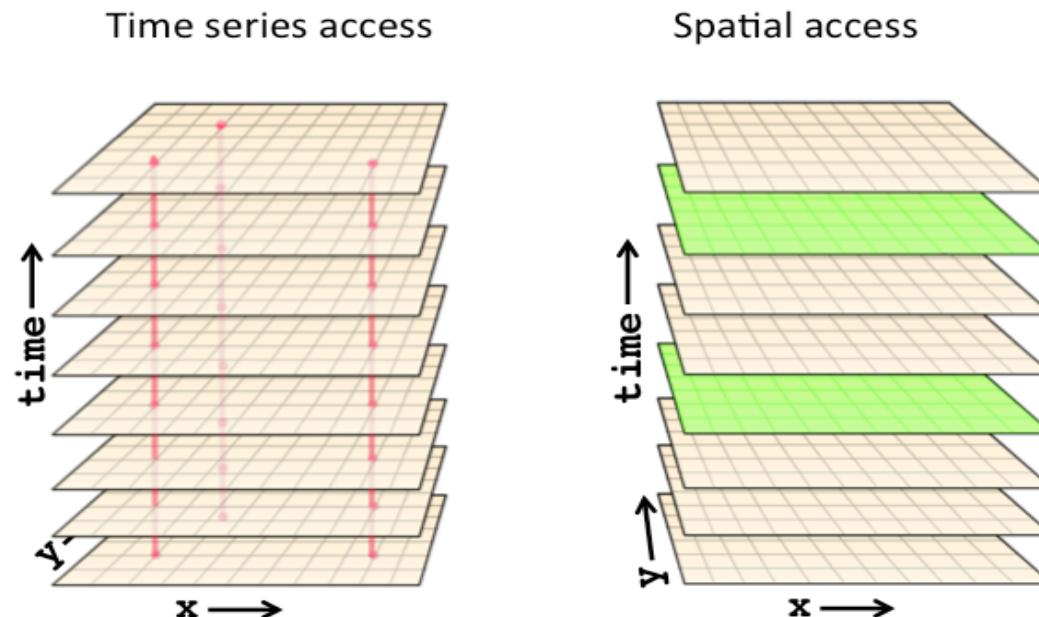


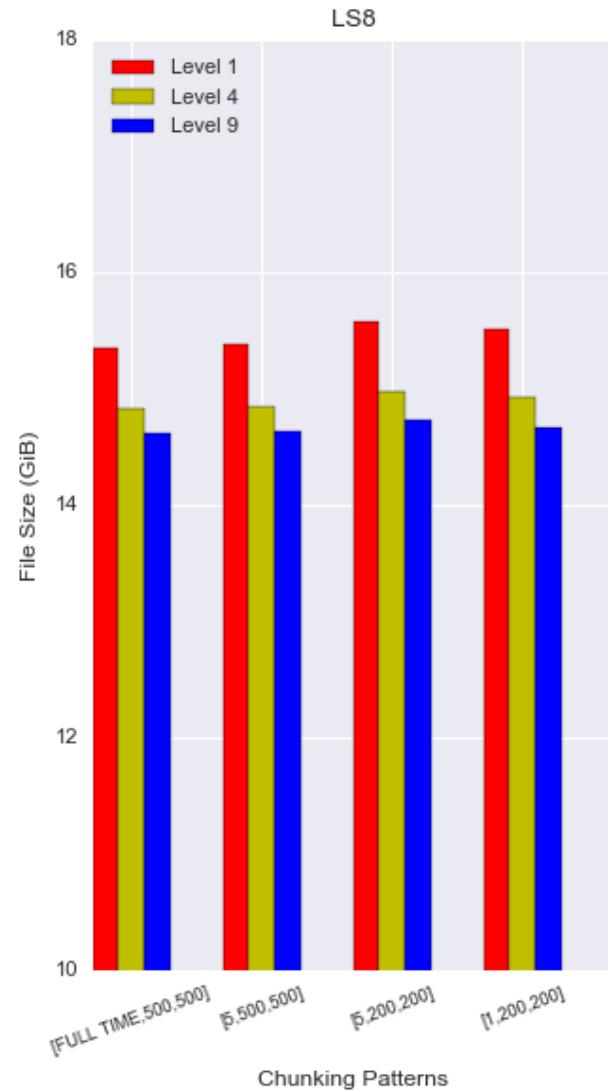
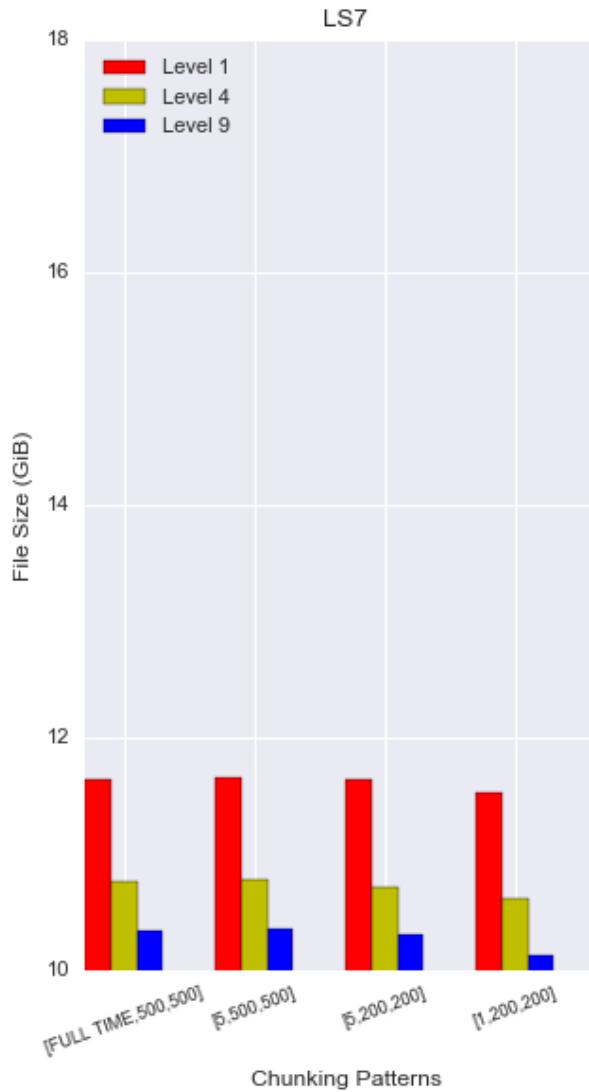
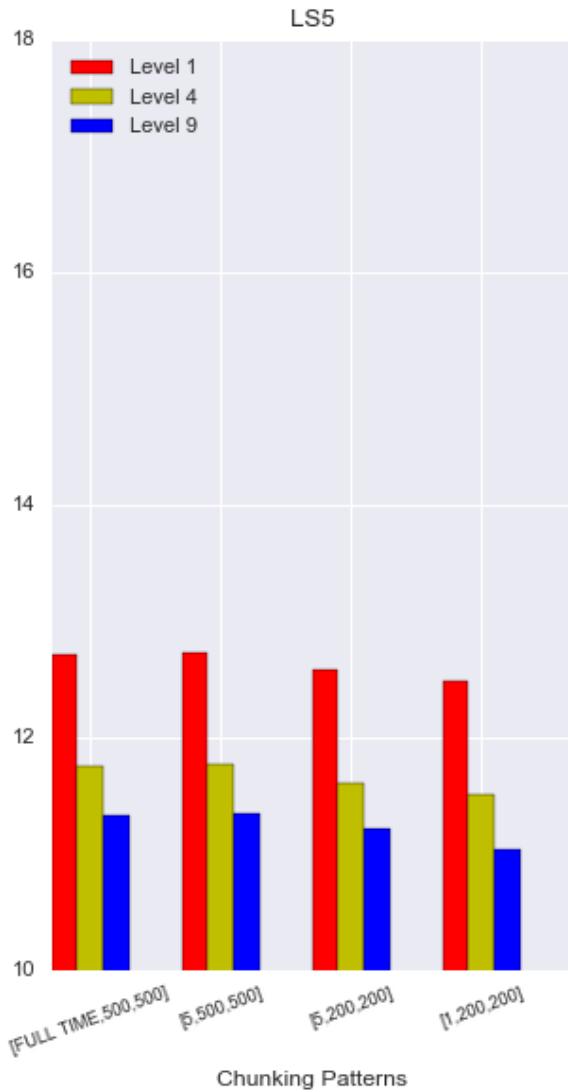


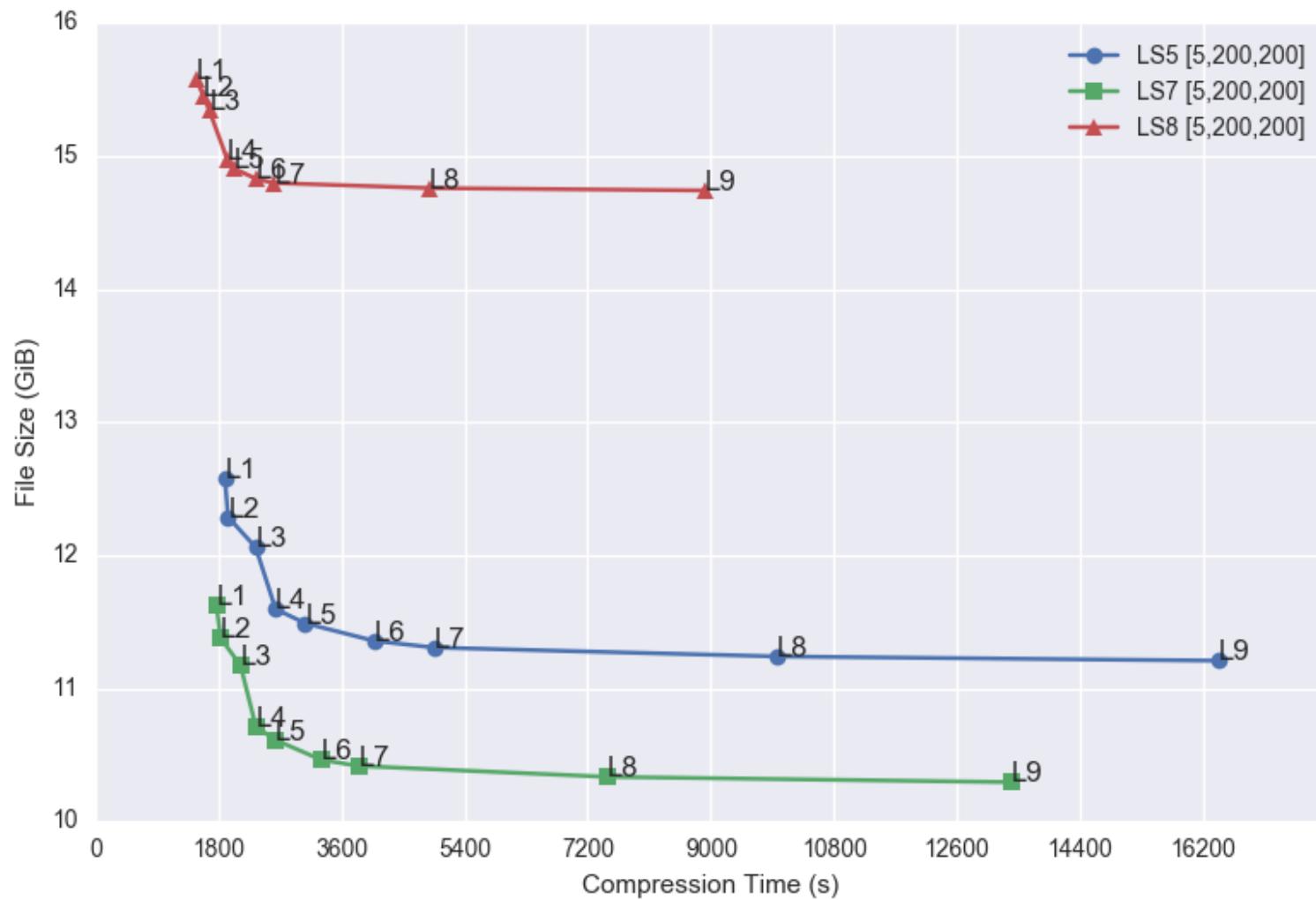
	Band-interleaved	Pixel-interleaved																														
Read Total	1 variables 9 variables	1 variables 9 variables																														
Read Time Throughputs	0.273s 446.3MB/s	3.208s 38MB/s																														
Touch range	122 MB	1071 MB																														
Touch Size	258 MB	1214 MB																														
Access Patterns	 <p>Graph showing time (ms) vs file offset (kiB) for band-interleaved access pattern. The x-axis ranges from 0.0 to 1200000.0 kiB, and the y-axis ranges from 0 to 350 ms. The data shows a linear increase in time as file offset increases.</p> <table border="1"> <caption>Data for Band-interleaved Access Pattern</caption> <thead> <tr> <th>file offset (kiB)</th> <th>time (ms)</th> </tr> </thead> <tbody> <tr><td>0.0</td><td>0</td></tr> <tr><td>200000.0</td><td>120</td></tr> <tr><td>400000.0</td><td>240</td></tr> <tr><td>600000.0</td><td>360</td></tr> <tr><td>800000.0</td><td>480</td></tr> <tr><td>1000000.0</td><td>600</td></tr> <tr><td>1200000.0</td><td>720</td></tr> </tbody> </table>	file offset (kiB)	time (ms)	0.0	0	200000.0	120	400000.0	240	600000.0	360	800000.0	480	1000000.0	600	1200000.0	720	 <p>Graph showing time (ms) vs file offset (kiB) for pixel-interleaved access pattern. The x-axis ranges from 0.0 to 1000000.0 kiB, and the y-axis ranges from 0 to 3500 ms. The data shows a non-linear, increasing trend in time as file offset increases.</p> <table border="1"> <caption>Data for Pixel-interleaved Access Pattern</caption> <thead> <tr> <th>file offset (kiB)</th> <th>time (ms)</th> </tr> </thead> <tbody> <tr><td>0.0</td><td>0</td></tr> <tr><td>200000.0</td><td>1500</td></tr> <tr><td>400000.0</td><td>2200</td></tr> <tr><td>600000.0</td><td>2800</td></tr> <tr><td>800000.0</td><td>3200</td></tr> <tr><td>1000000.0</td><td>3500</td></tr> </tbody> </table>	file offset (kiB)	time (ms)	0.0	0	200000.0	1500	400000.0	2200	600000.0	2800	800000.0	3200	1000000.0	3500
file offset (kiB)	time (ms)																															
0.0	0																															
200000.0	120																															
400000.0	240																															
600000.0	360																															
800000.0	480																															
1000000.0	600																															
1200000.0	720																															
file offset (kiB)	time (ms)																															
0.0	0																															
200000.0	1500																															
400000.0	2200																															
600000.0	2800																															
800000.0	3200																															
1000000.0	3500																															

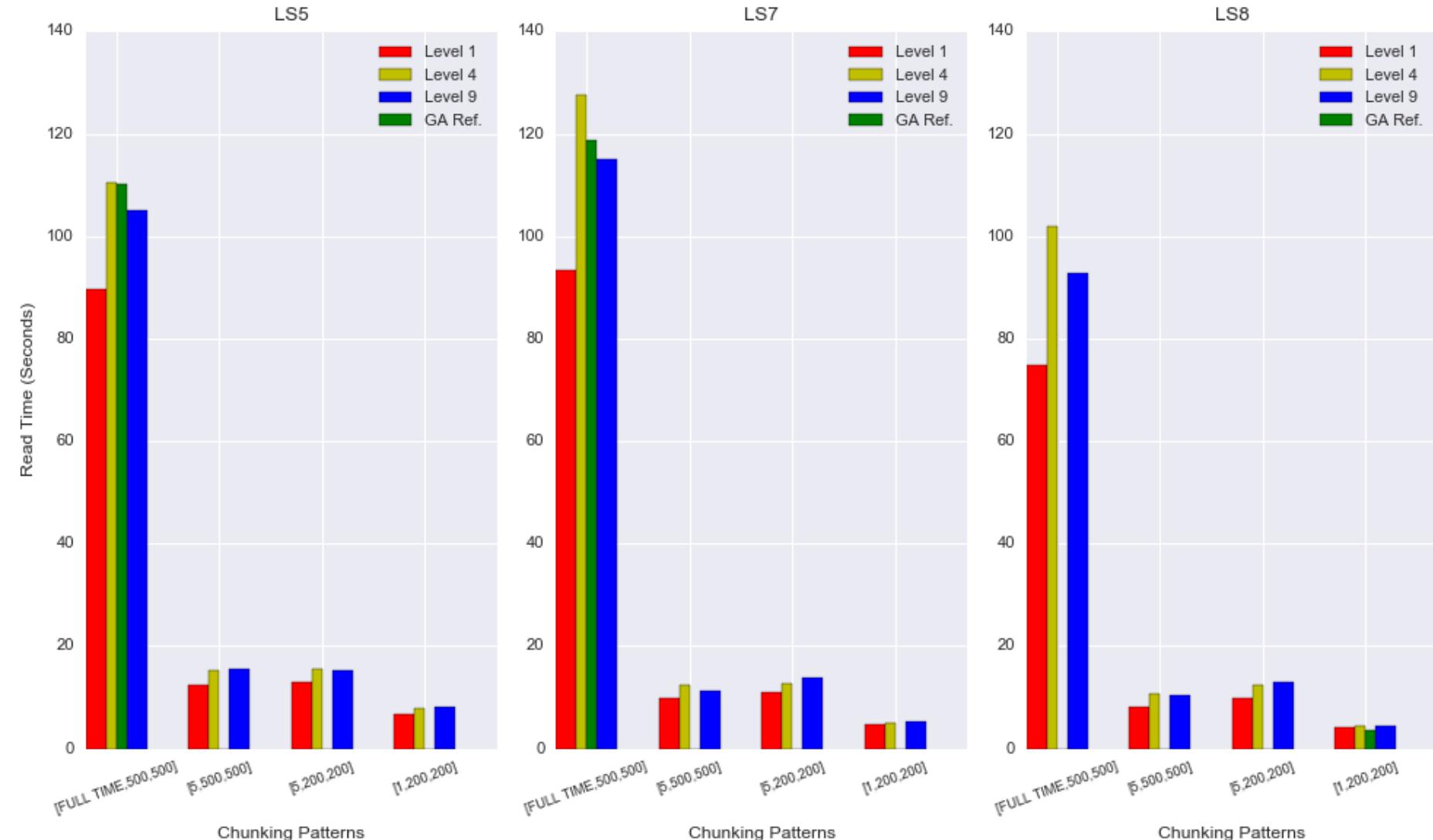
Raw GA NCF	LS5_TM_NBAR, LS7_ETM_NBAR, LS8_OLI_TIRS_NBAR
Variable Size/Band/file	[Full Time, 4000 pixel, 4000 pixel] e.g. 1 year x 100KM x 100KM
Chunking Shapes	[Full Time, 500 pixel, 500 pixel ], [5, 500, 500], [5, 200, 200], [1, 200, 200]
Deflation Levels	1, 4, 9
Data Type	Short
	Spatial access: 13 days x 400 KM x 400KM (16 files) Time series access: 49 years x 1 pixel x 1 pixel (49 files)

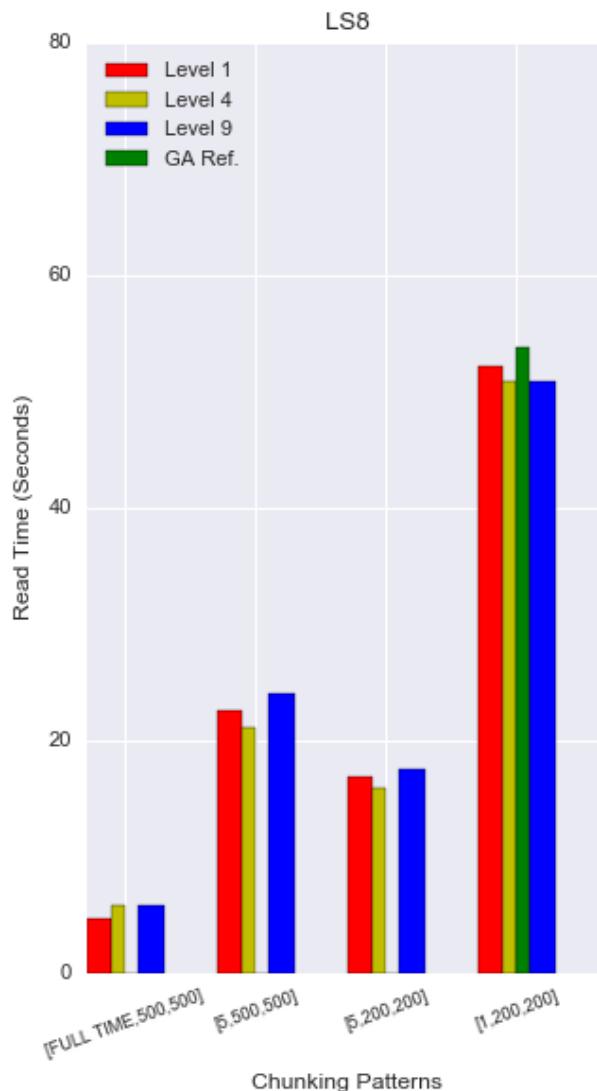
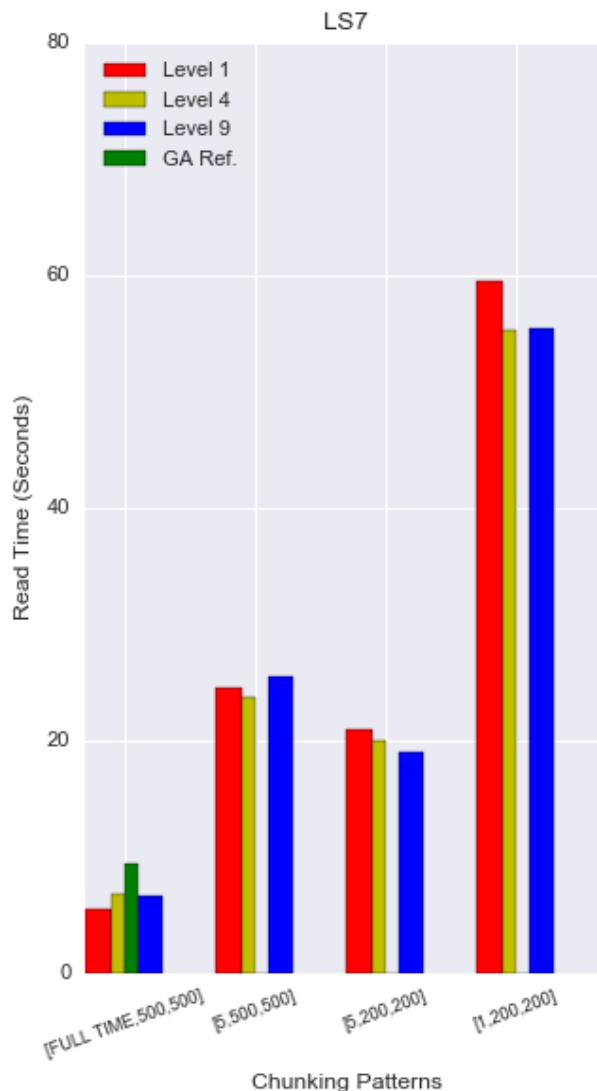
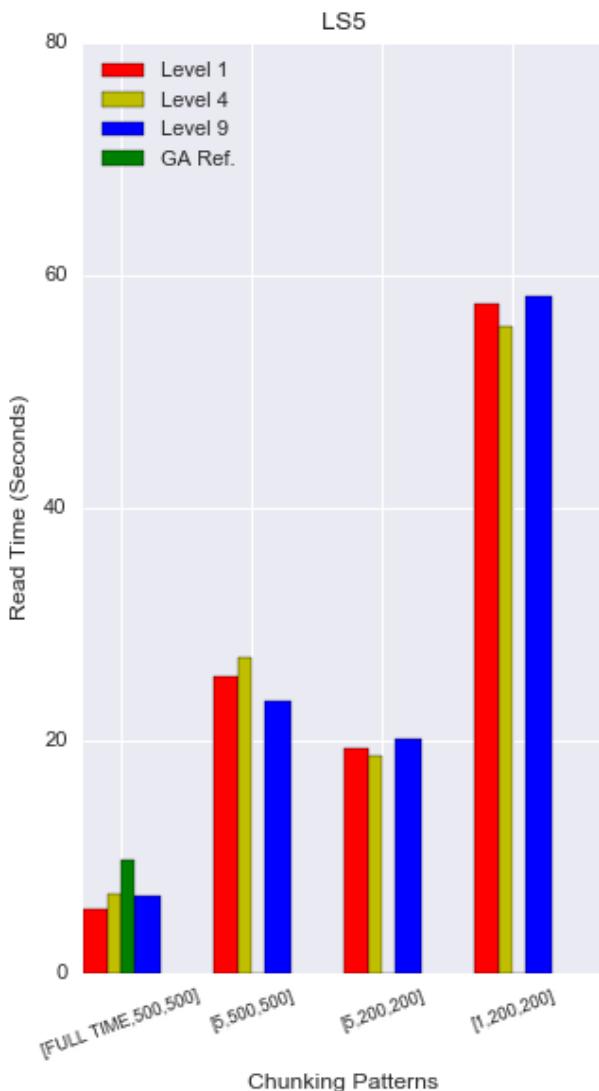
### Access patterns

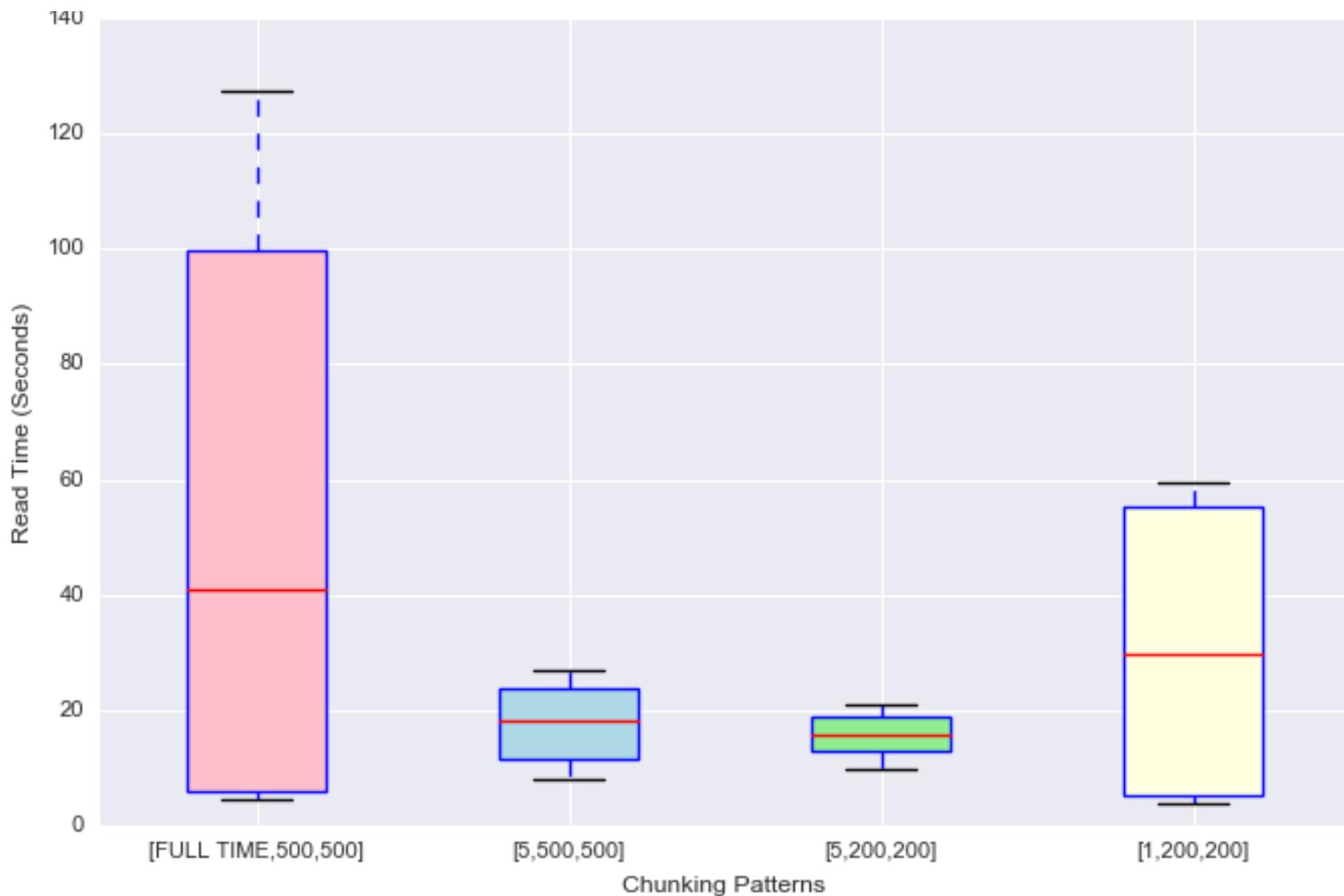






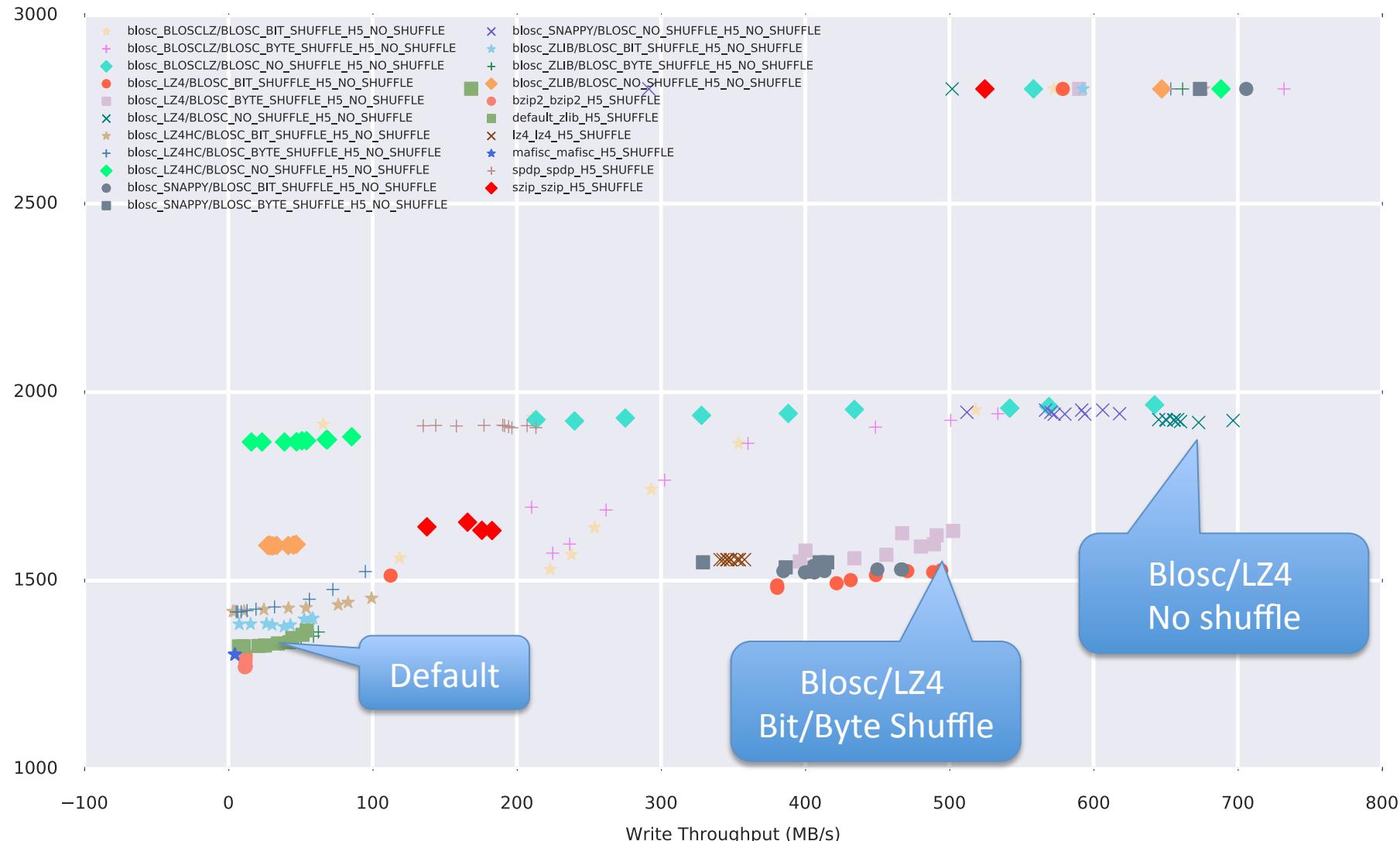


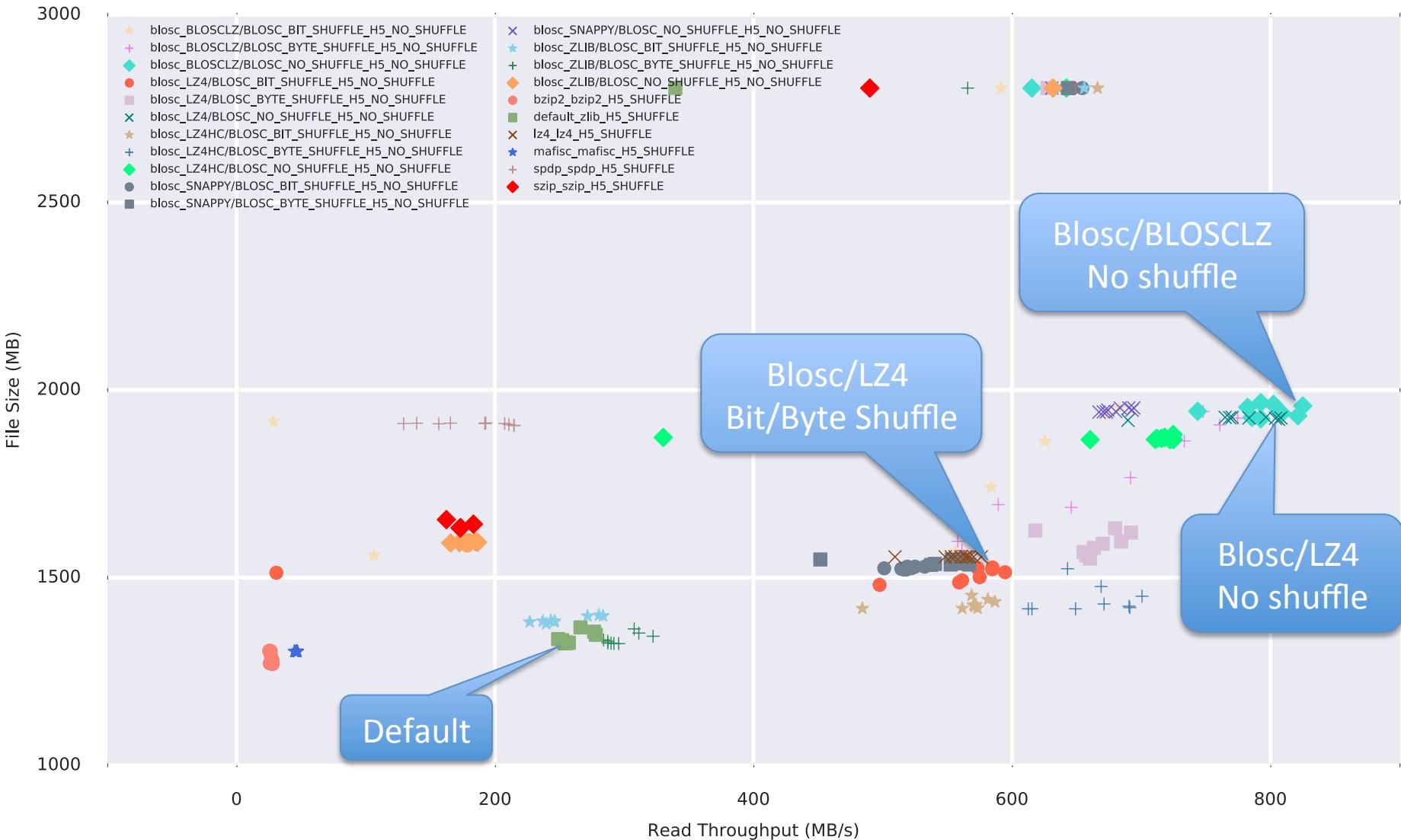


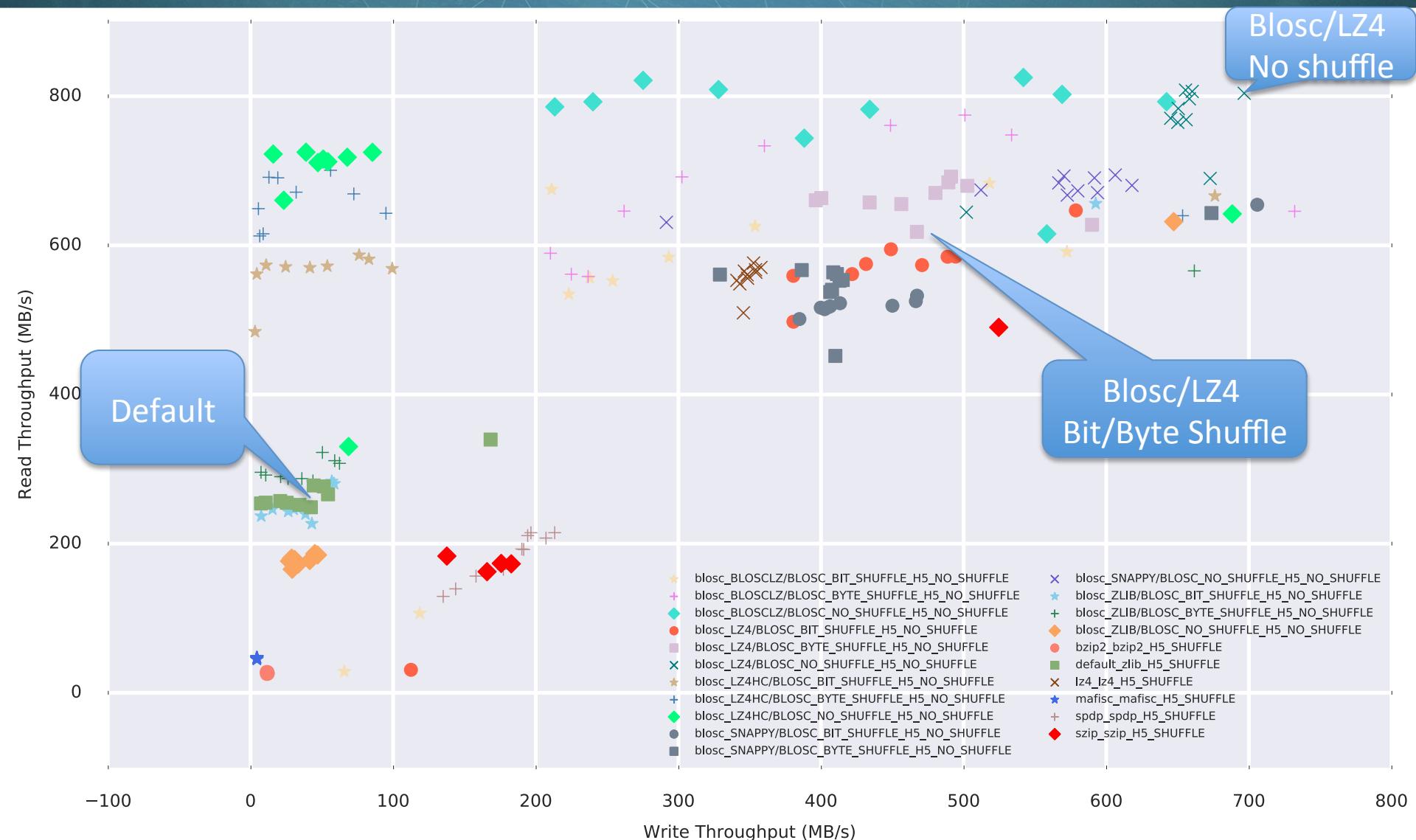


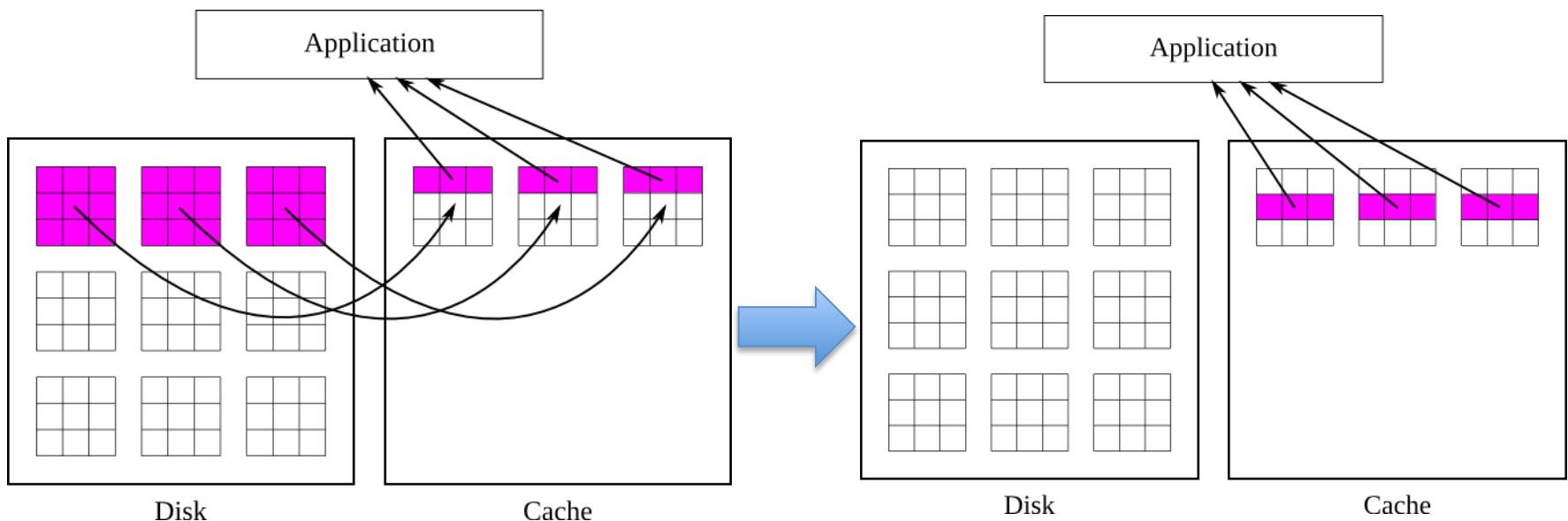
Library	Default	Dynamic Filter
NetCDF4	Deflate(Zlib)	N/A
HDF5	Deflate(Zlib)	Bzip2,mafisc,spdp,szip,LZ4 Blosc(blosclz,lz4hc,lz4,SNAPPY,ZLIB)

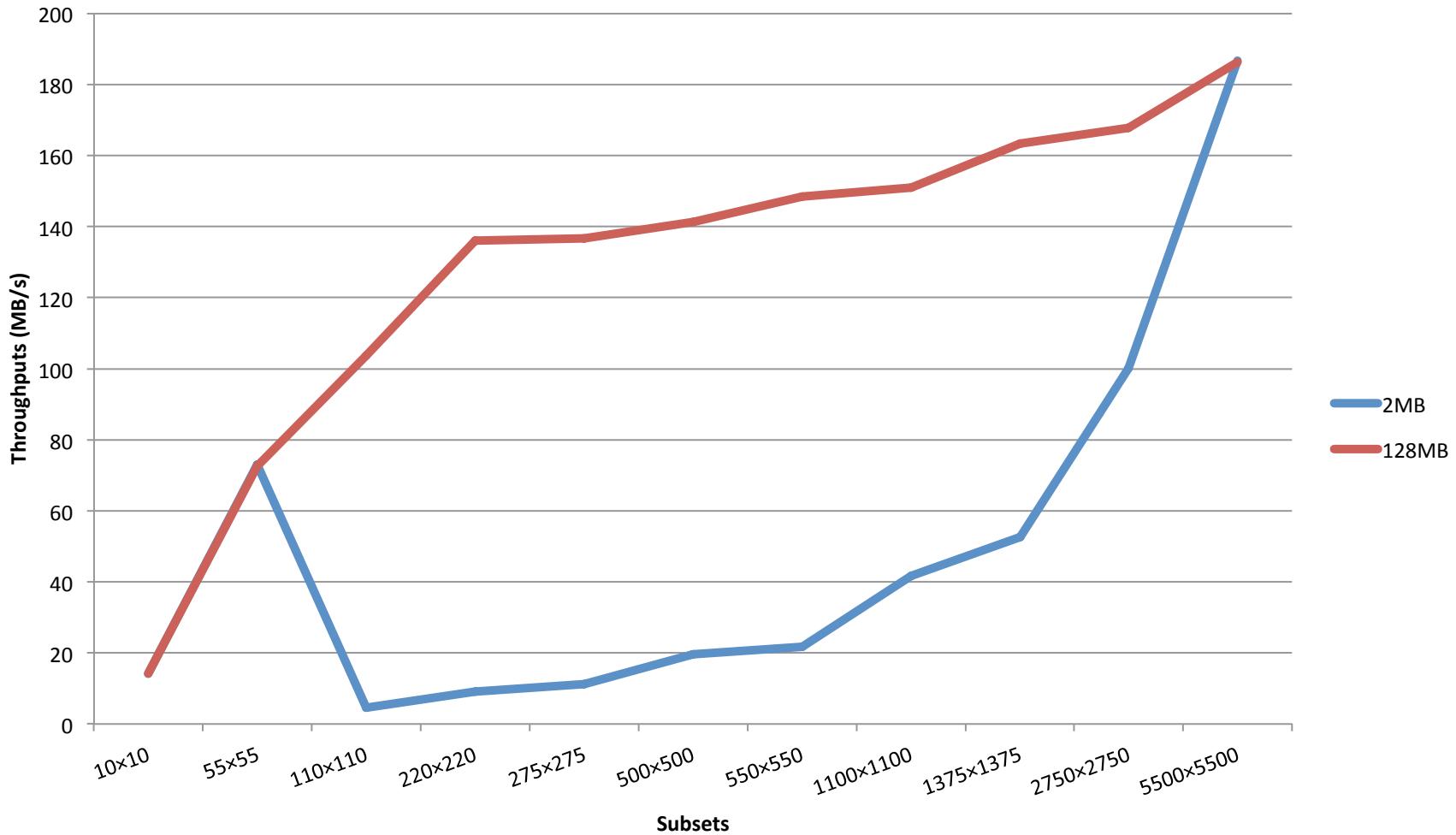
The NetCDF4 library is revised to enable different compression algorithms via dynamic filter.

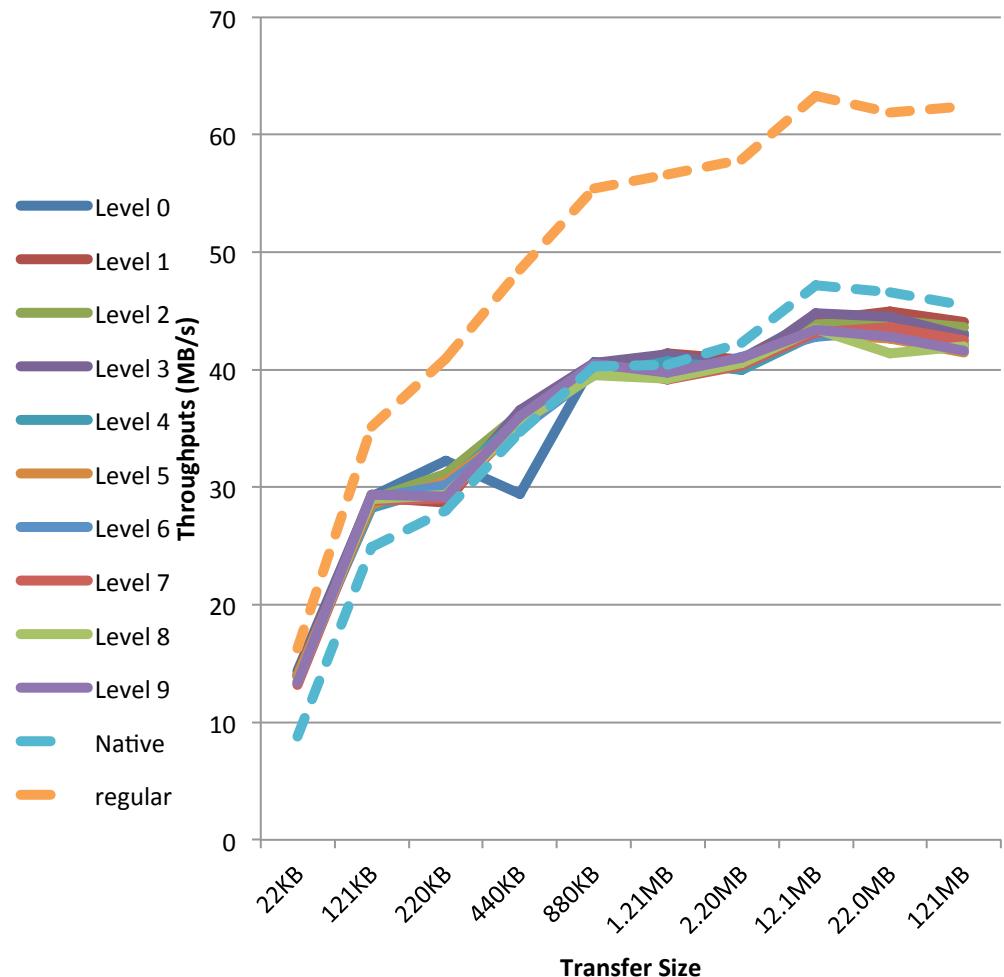
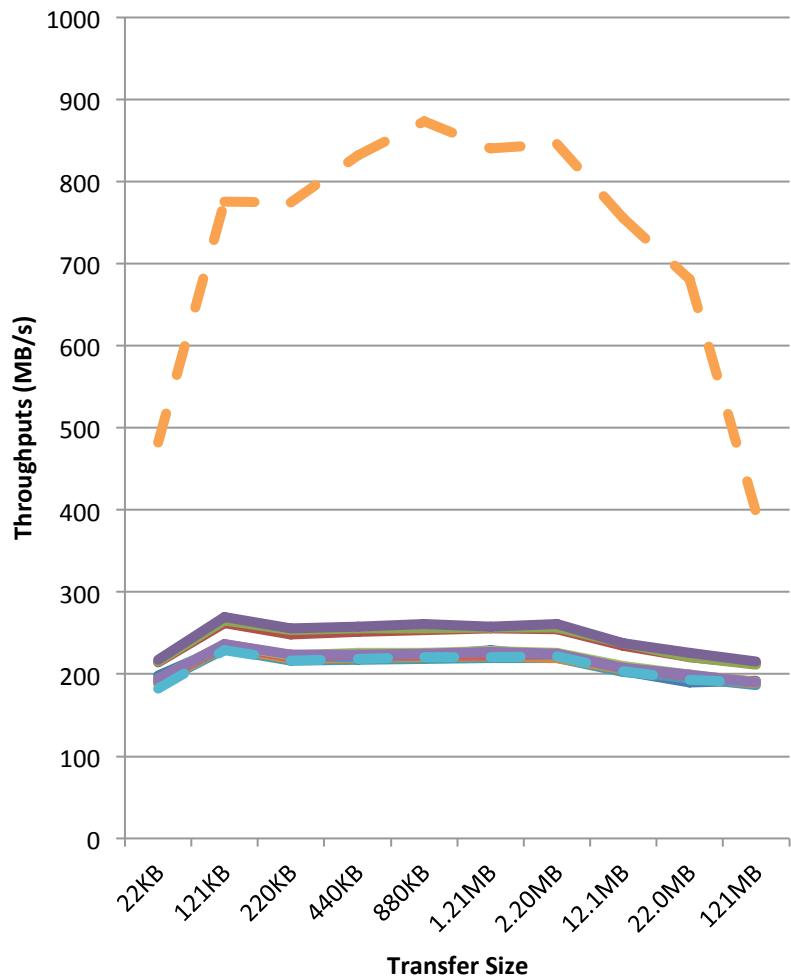






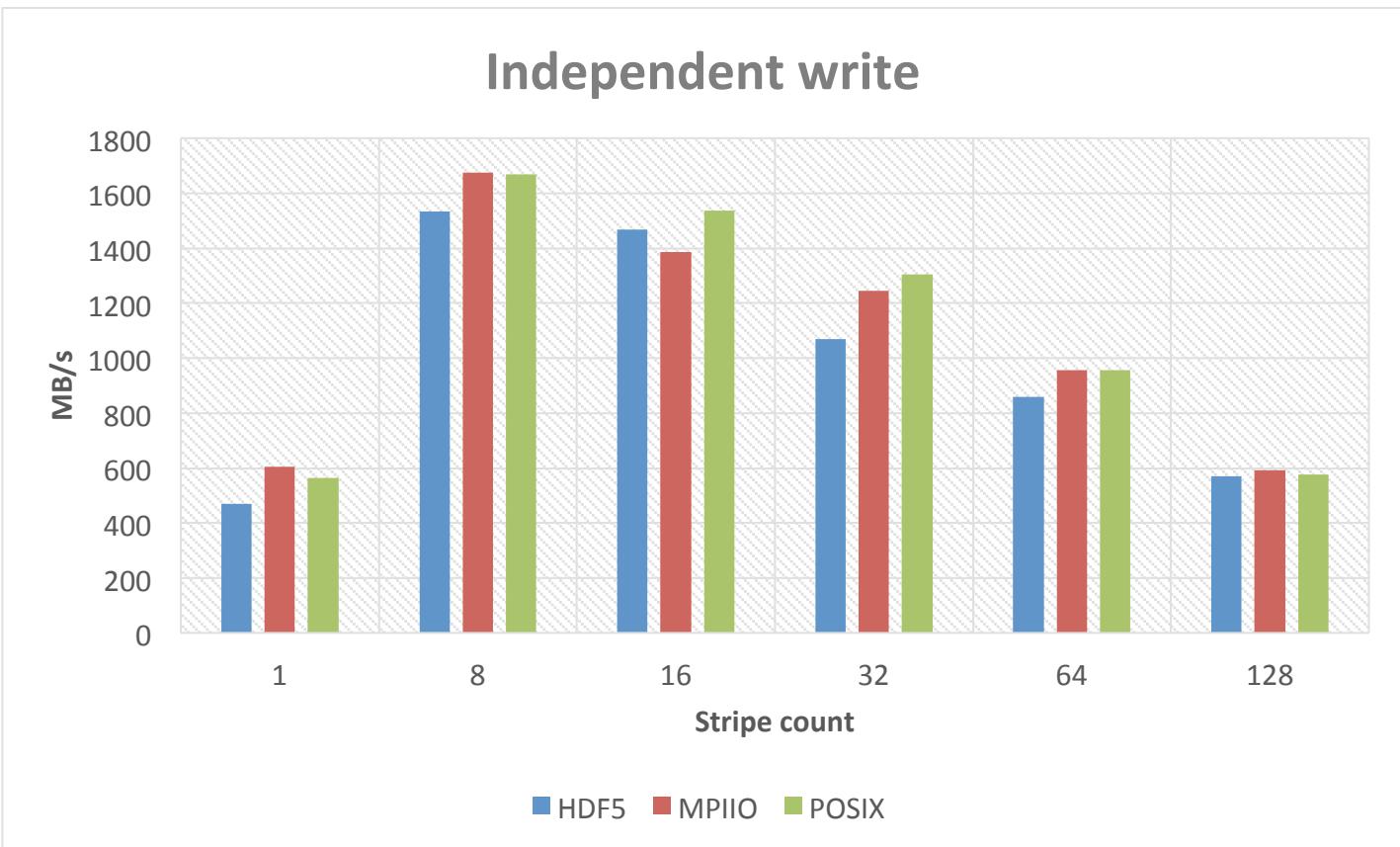



 Chunk size= $1 \times 5500$



# PARALLEL IO

IOR Benchmark  
MPI size = 16  
Stripe size = 1M  
Block size = 8G  
Transfer size = 32M



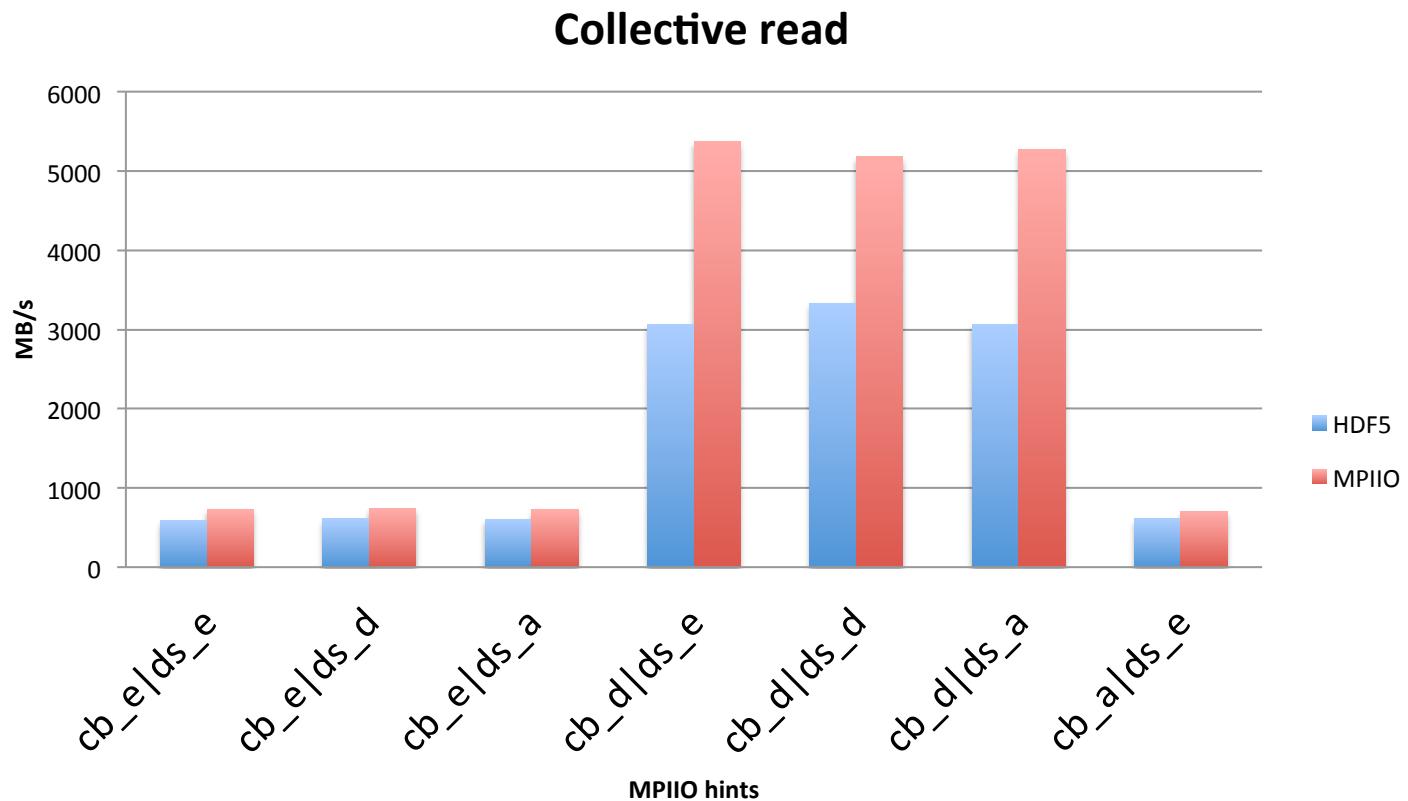
IOR Benchmark  
MPI size = 16  
Stripe size = 1M  
Block size = 8G  
Transfer size = 32M



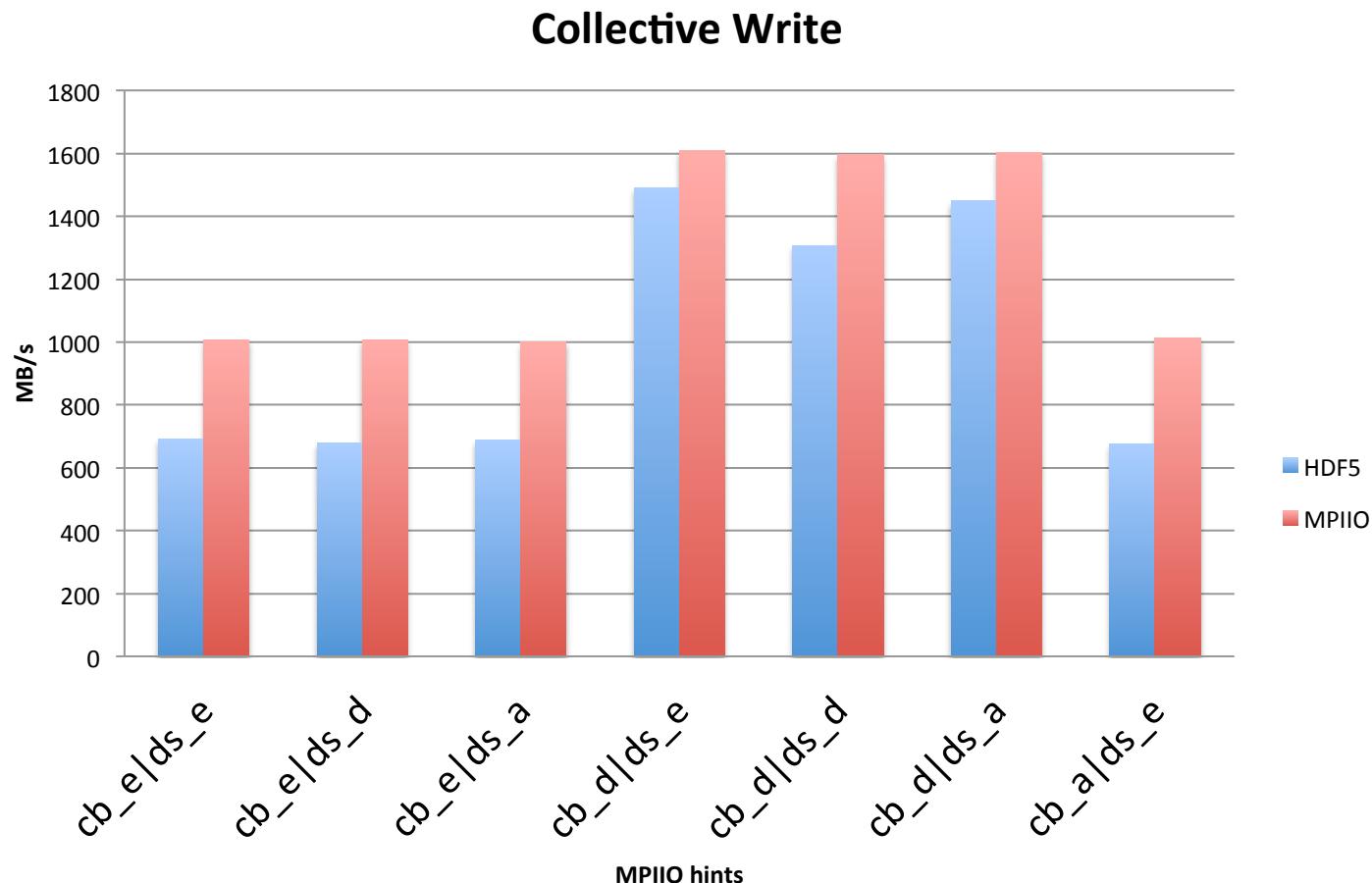
- Default value for MPI ranks using 8 cores/node \* 2 nodes:
  - *direct\_read = false*
  - *direct\_write = false*
  - *romio\_lustre\_co\_ratio = 1*
  - *romio\_lustre\_coll\_threshold = 0*
  - *romio\_lustre\_ds\_in\_coll = enable*
  - *cb\_buffer\_size = 16777216 (16M)*
  - *romio\_cb\_read = automatic*
  - *romio\_cb\_write = automatic*
  - *cb\_nodes = 2*
  - *romio\_no\_indep\_rw = false*
  - *romio\_cb\_pfr = disable*
  - *romio\_cb\_fr\_types = aar*
  - *romio\_cb\_fr\_alignment = 1*
  - *romio\_cb\_ds\_threshold = 0*
  - *romio\_cb\_alltoall = automatic*
  - *ind\_rd\_buffer\_size = 4194304*
  - *ind\_wr\_buffer\_size = 524288*
  - *romio\_ds\_read = automatic*
  - *romio\_ds\_write = automatic*
  - *cb\_config\_list = \*:1*
  - *striping\_unit = 1048576*
  - *striping\_factor = 1*
  - *romio\_lustre\_start\_iodevice = 0*
- Collective buffering
  - Romio\_cb\_read/write
    - auto
    - enable
    - disable
  - Romio\_config\_list (aggregators)
  - Cb\_buffer\_size
- Data sieving
  - Romio\_ds\_read/write
    - auto
    - enable
    - disable

- Data sieving
  - I/O performance suffers considerably when making many small I/O requests
  - Access on small, non-contiguous regions of data can be optimized by grouping requests and using temporary buffers
  - This optimization is local to each process (non-collective operation)
- Collective buffering
  - processes to match data layout in file.
  - Mix of I/O and MPI communications to read or write data
  - Communication phase to merge data from different processes into large chunks
  - File accesses are done only by selected processes (called aggregators), the others communicates with them
  - Large operations are split into multiple phases (to limit the size of the buffers and to overlap communications and I/O)

MPI size = 16  
 Stripe count = 8  
 Stripe size: 1M  
 Block size:  
 1G (cfg 1,2,3,7)  
 4G (cfg 4,5,6)  
 Transfer size = 8M

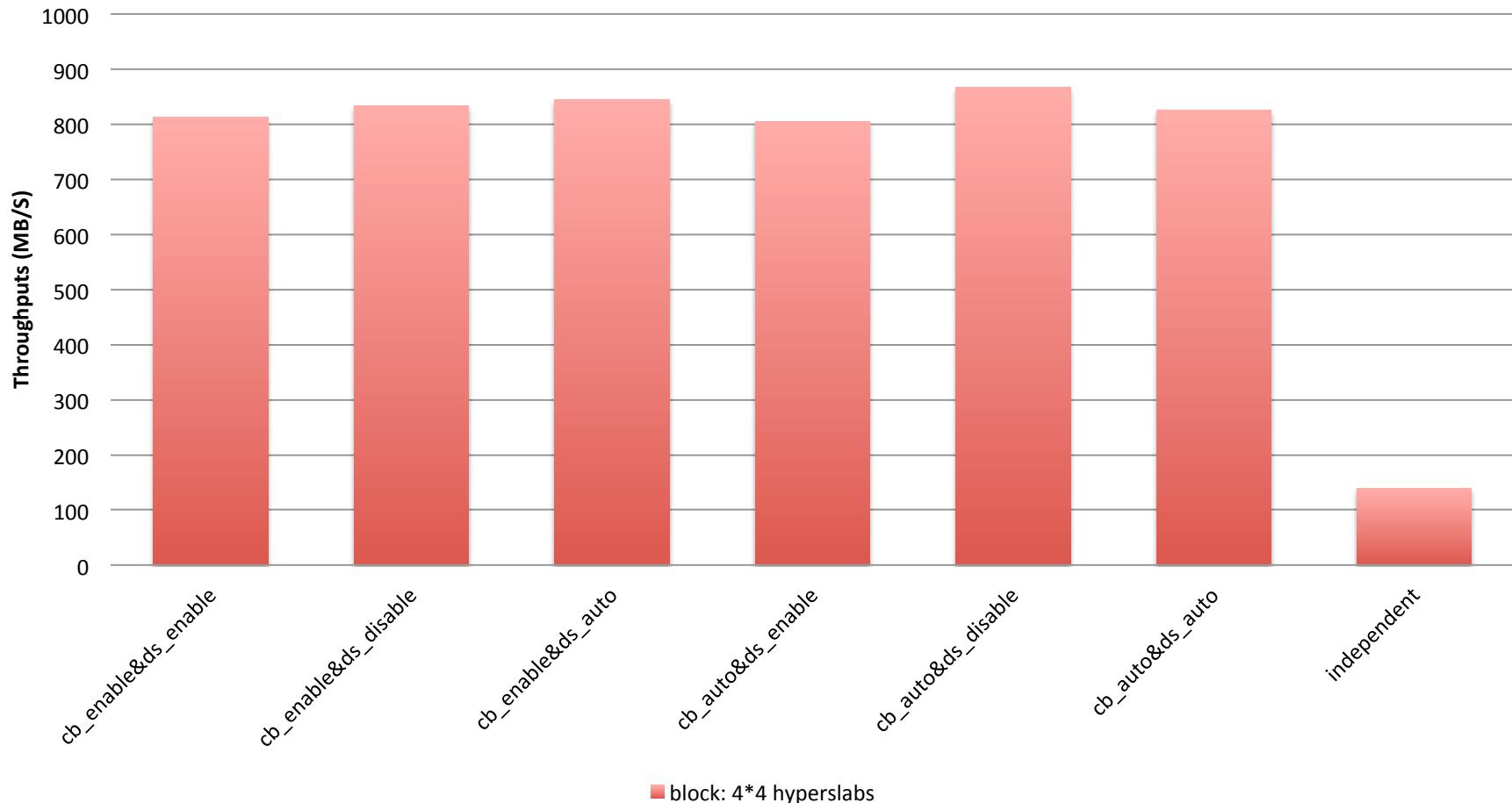


Disable collective buffering to get better performance for contiguous access.

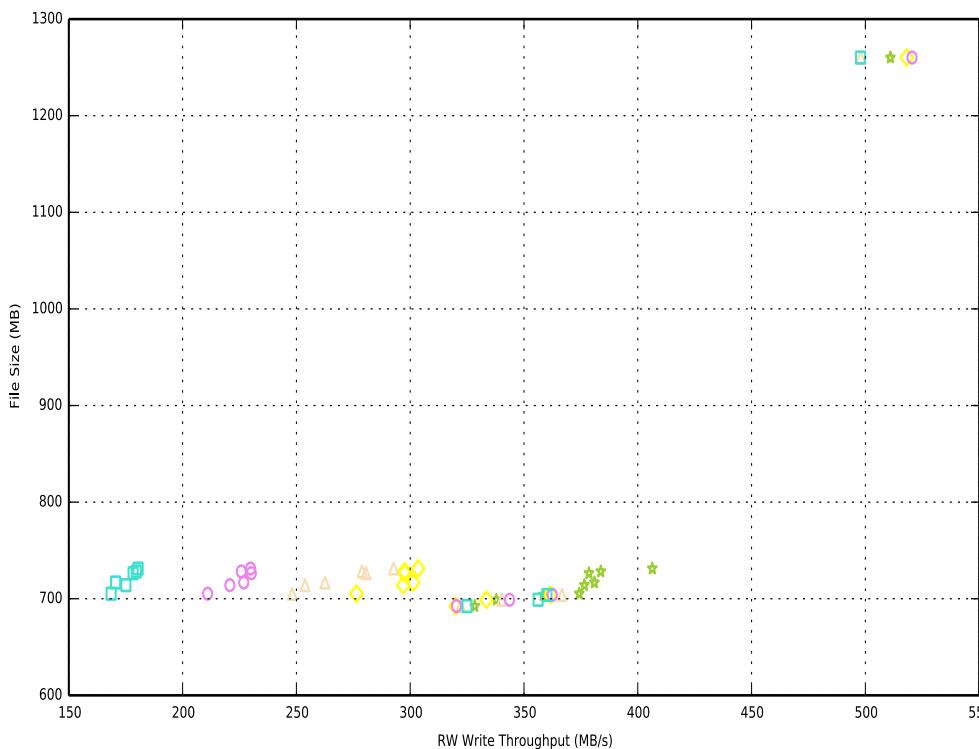


Disable collective buffer to get better performance for contiguous access.

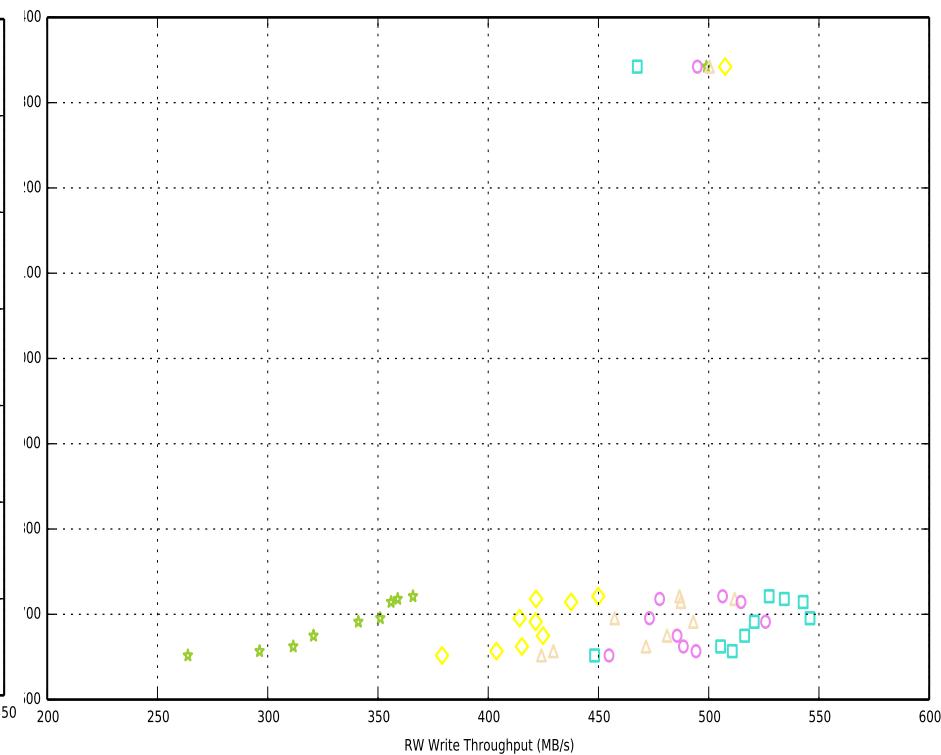
## Read with hyperslab [1:1:512:256]

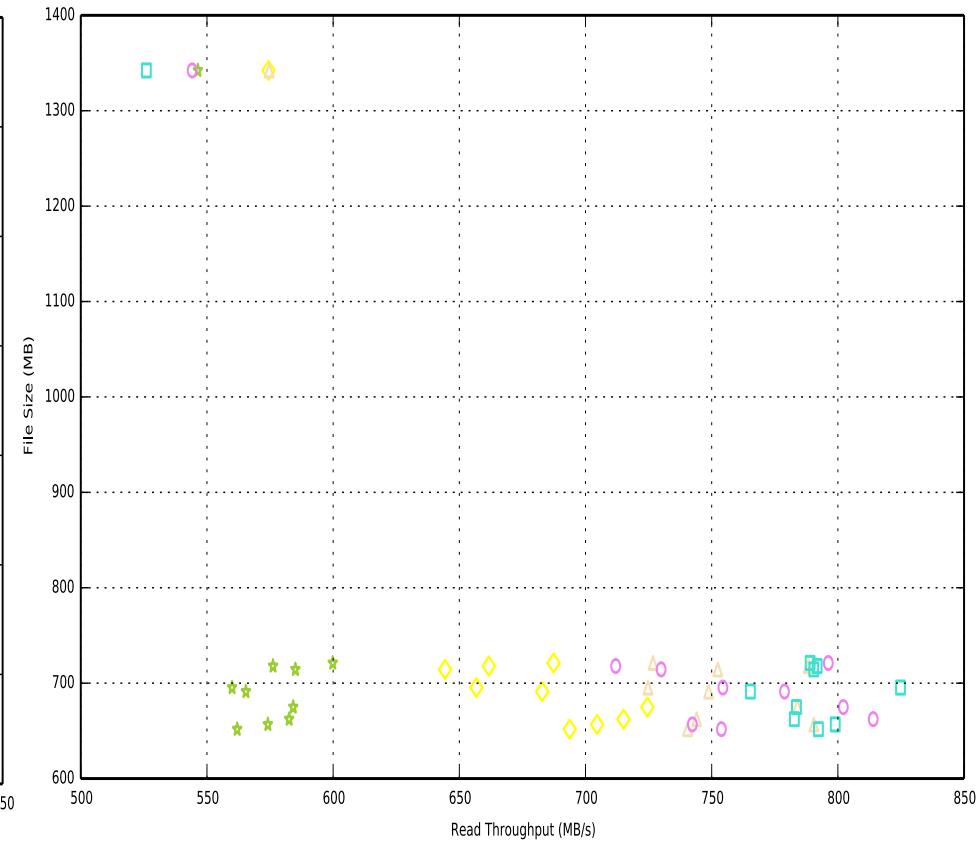
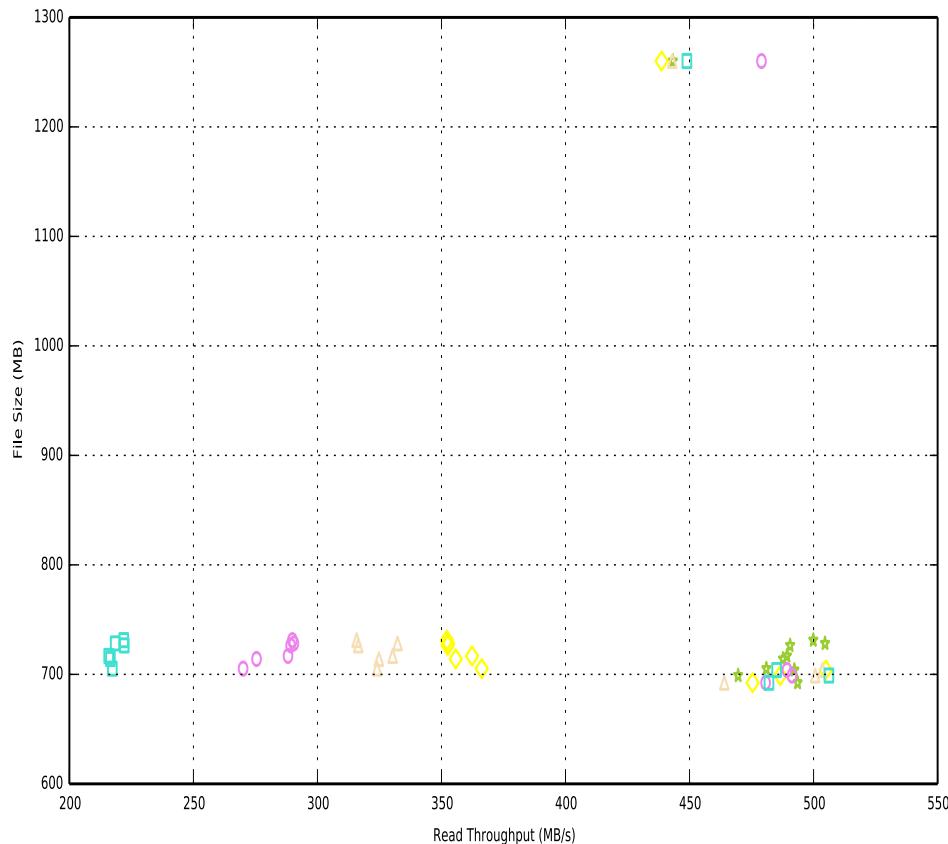
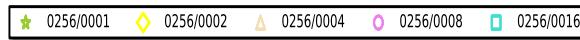


★ 0256/0001    ◆ 0256/0002    ▲ 0256/0004    ○ 0256/0008    □ 0256/0016



★ 4096/0001    ◆ 4096/0002    ▲ 4096/0004    ○ 4096/0008    □ 4096/0016

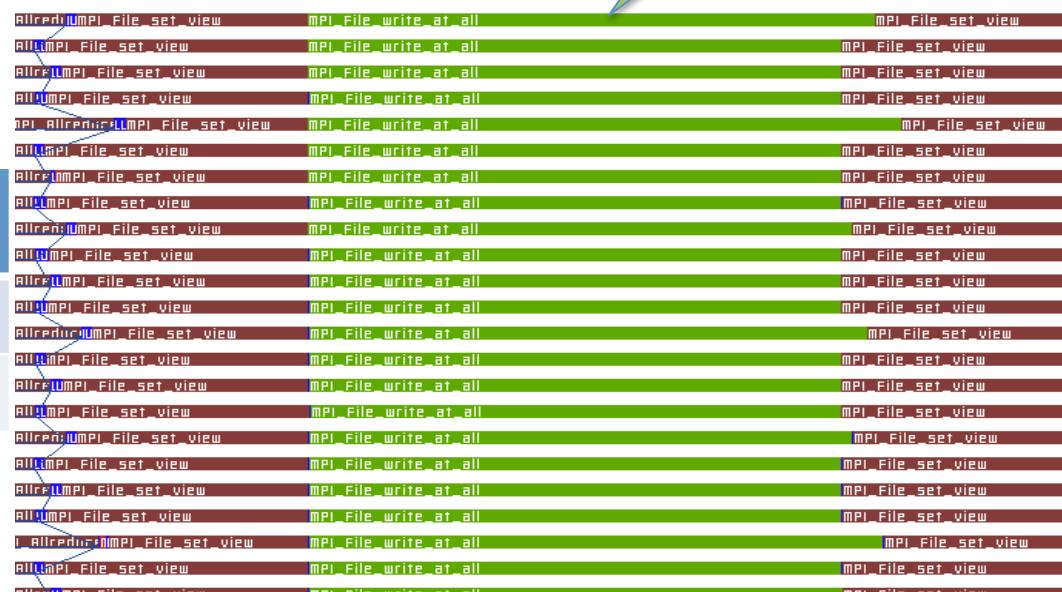






Independent  
metadata  
operations

Collective  
metadata  
operations



Elapsed Time(s)	Serial IO	Official HDF5/NetCDF4	Patched HDF5/NetCDF4
Case 1	4909	10173	3766
Case 2	1625	8543	1341

- IO performance tuning involves many parameters of interest
  - File System (Lustre): stripe count, stripe unit
  - MPI-I/O: collective buffer size, collective nodes etc.
  - NetCDF/HDF5: alignment, sieve buffer size, data layout, chunking, compression etc.
  - User application: access patterns,
  - Access to data: network (remote) vs in-situ (local)

Thank you!