



Pattern-Aware File Reorganization in MPI-IO

Jun He¹, Huaiming Song¹, Xian-He Sun¹, Yanlong Yin¹, Rajeev Thakur² 1: Illinois Institute of Technology, Chicago, Illinois 2: Argonne National Laboratory, Argonne, Illinois



Outline

- Motivation
 - Examples
 - o Basic idea

Design

- System Overview
- Trace collecting
- Pattern classification
- o I/O Trace analyzer
- Remapping table
- MPI-IO remapping layer

Evaluation

- Remapping overhead
- Pattern variation
- o Benchmarks
- Conclusion & Future Work





Motivation

• • •







Parallel File Systems



- Important Factors
 - Number of requests
 - Contiguousness of accesses

Network overhead IOPS Locality

. . .





Mismatch

Logical data

- Developer's understanding, for programmability and runtime performance
- -> Logical organization -> Access pattern

Physical data

- Where the data blocks are stored
- -> Physical data organization

Good logical organization

!=

Good physical organization for better I/O performance



Potential benefit:

Better spatial locality Easier for some optimization to take effect Less disk head movements

An Example for Regular 2-d Array Default Organization



A 2-D array 4 5 6 7

 4
 3
 6
 7

 8
 9
 10
 11

 12
 13
 14
 15





Read a Subarray

			I/O Server #0				I/O Server #1				I/O Server #2				I/O Server #3								
						1	2	3			1	2	3			1	2	3			1	2	3
						1	2	3			1	2	3			1	2	3		j	1	2	3
A 2-D array				:					:					:				:					
	1	2	3		4	5	6	7		4	5	6	7		4	5	6	7	4		5	6	7
					4	5	6	7		4	5	6	7		4	5	6	7	4		5	6	7
	5	6	7		:					÷					:				:				
:	9	10	11		8	9	10	11		8	9	10	11		8	9	10	11	8		9	10	11
					8	9	10	11		8	9	10	11		8	9	10	11	8		9	10	11
2	13	14	15		:			:					:				:						
					12	13	14	15		12	13	14	15		12	13	14	15	12	2	13	14	15
					12	13	14	15		12	13	14	15		12	13	14	15	12	2	13	14	15
				:					÷					:				:					

PDSW'11





After Re-organizing







A Messier One

- Irregular data
- Very complex data model
- Computation which involves multiple data fields





Pattern-Aware Reorganization

- Be aware of repeating non-contiguous access patterns
 n-d strided and irregular
- Try to reorganize the data so that data is contiguous.
 - o Less network overhead
 - Less IO operations
 - Better locality
 - Beneficial for other optimizations, e.g. data sieving...
- Motivating Scenarios
 - Application start-up
 - Data analysis, visualization
 - 0 ...
- Where it does not apply
 - Patterns do not repeat from run to run.



















Trace Collecting

- Wrap the original function call
 - Add recording function
 - Call original function inside
- Process ID, MPI rank, file path, type of operation, offset, length, data type, time stamp, and file view







Pattern Classification







I/O Trace Analyzer

- Pattern matching
 - Sort Traces by time
 - Separate by process
 - Find out patterns
- I/O Signature

{I/O operation, initial position, dimension, ([{offset Pattern}, {request size pattern}, {pattern of number of repetitions}, {temporal pattern}], [...]), # of repetitions}







I/O-signature-based Remapping Table







MPI-IO Remapping Layer

- Convert old offsets to new ones
 Example:
- Read m bytes data from offset f.
- Whether this access falls in a 1-d strided pattern ?
 - o starting offset off
 - o read size **rsz**
 - o hole size **hsz**
 - $\circ\,$ number of accesses of this pattern ${\it n}$
- (f-off)/(rsz+hsz) <n
- (f-off)%(rsz+hsz) = 0
- m = rsz (3) newoff = off+rsz*(f-off)/(rsz+hsz)







Evaluation

 $\bullet \quad \bullet \quad \bullet$





Experiment Environment

- Dual 2.3GHz Opteron quad-core processors
- 8G memory
- 250GB 7200RPM SATA hard drive
- 100GB PCI-E OCZ Revodrive X2 SSD (read: up to 740 MB/s, write: up to 690 MB/s).
- Ethernet/Infiniband
- Ubuntu 9.04 (Linux kernel 2.6.28-11-server)
- **PVFS2** 2.8.1: stripe size 64 KB
- MPICH2 1.3.1





Remapping Overhead

1-D Strided Remapping Table Performance (1,000,000 accesses)

Table Type	Size (bytes)	Building time (sec)	Time of 1,000,000 lookups (sec)
1-to-1	64,000,000	0.780287	0.489902
I/O Signature	28	0.00000269	0.024771

Who use 1-to-1: PLFS uses 1-to-1 mapping table in index file. Most OS file systems also use similar table to store free blocks in disk.





Request Size Variation

• X: different of request size. For example, 5% means the actual request size is 5% less than the one assumed.





Variation of Starting Offset

• X: difference of starting offsets. 5% means that the starting offset moved to the 5%th of the whole access.



R/W Performance – on IOR

• 4 I/O clients, 4 I/O servers. 64 processes with HDD and Infiniband



• PDSW'11

Performance on MPI-TILE-IO

• 4 I/O clients, 4 I/O servers. 64 processes with HDD and Infiniband. Elements in a tile: 1024x1024.





• 4 I/O clients, 4 I/O servers. 64 processes with SSD and Infiniband. Elements in a tile: 1024x1024.







Conclusion & Future Work

Conclusion

- Different file organizations lead to very different performance.
- Bridging logical data and physical data Access pattern

 better organization

 better performance

Future Work

- Multiple replicas with different organizations.
- More complicated access patterns, patterns with hints
- File reorganization for emerging storage medias, such as SSD





Acknowledgement

- Hui Jin and Spenser Gilliland (Illinois Institute of Technology)
- Ce Yu (Tianjin University, China)
- Samuel Lang (Argonne National Laboratory)
- NSF grant CCF-0621435, CCF-0937877
- Office of Advanced Scientific Computing Research, Office of Science, U.S. DOE, under Contract DEAC02-06CH11357.

Thanks!