# **Discovering Structure in Unstructured I/O**

Jun He<sup>1,2</sup>, John Bent<sup>3</sup>, Aaron Torres<sup>4</sup>, Gary Grider<sup>4</sup>, Garth Gibson<sup>5</sup>, Carlos Maltzahn<sup>6</sup>, Xian-He Sun<sup>1</sup> <sup>1</sup>Illinois Institute of Technology, <sup>2</sup>New Mexico Consortium, <sup>3</sup>EMC, <sup>4</sup>Los Alamos National Laboratory, <sup>5</sup>Carnegie Mellon University, <sup>6</sup>University of California Santa Cruz

# Motivation

Extremely challenging storage requirements for exaflop supercomputers:

Checkpoints of 32 petabytes in size should complete in 300 seconds (DOE projection)

**D**PLFS (Parallel Log-structured File System) reduces checkpoint time by up to several orders of magnitude

□By transparently transforming N-1 writes (N processes write to 1 file) to N-N.

**However**, PLFS metadata (index) grows as the application size and number of data writes increase.

Considerable overhead of meatadata fetching

- Considerable memory occupation by metadata Example: an anonymous application from LANL Data: **4 GB** 
  - □ PLFS metadata (on disk): 192 MB
  - □ PLFS metadata (in memory): 12 GB

# PLFS Index

PLFS transparently transforms shared-file writing into log-structured file-per-process writing. Index has a mapping between the bytes within the logical file and their location within the physical files.

The figure below shows an example of two processes writing to a traditional PLFS file. Indices can become very large as the number of writes increases.



To shrink index size, we discover patterns in index entries and represent them in a compact way.

simplified For example, sequence (8,8,8,8,8,8,8,8,8,8) can be represented as 8^10 (ten eights). So the size is compressed.

### **Notation for Complex Patterns**

 $\Box$  *i* is the first element of the original sequence.  $\Box d[]$  (delta) is the repeating part of an array containing the distances of any two consecutive elements in the original sequence.  $\Box$  *r* is the number of repetitions. For example, (5, 7, 10, 12, 15) can be represented as  $[5, (2, 3) \land 2].$ 







<b>Time</b>
$\Box_{W}$ :
$\Box n$ :

•	
	PL
2.	Red
3.	Loo
	ima
	[0,(
	0
	14
	28
4.	fine
	offs
5.	Cal

# **Discovering Pattern Structure**

 $[i, (d[0], d[1], ...)^{\wedge}r]$ 

### **Discovering Pattern**

An example to find the logical offset patterns in index.0 of the previous figure:

0 3 7 14 17 21 28 31 35 42 46 50 54 58 3473473474444 Pattern Stack .....[3] 3 4 7 3 4 7 3 4 7 4 4 4 4 ······[3][4] 3473473474444 ······[3][4][7] 3473473474444 3473473474444 ······[(3,4,7)^2] ·····[(3,4,7)<sup>3</sup>] 3473473474444 .....[(3,4,7)^3],[4] 3473473474444 3473473474444 ····[(3,4,7)^3],[(4)^2] ..[(3,4,7)^3],[(4)^4] 3 4 7 3 4 7 3 4 7 4 4 4 4 Search window Look ahead window [0,(3,4,7)^3], [42,(4)^4]

### □ New Index.0 using pattern structures:

Physical Offset Chunk ID Logical Offset Length [2,(0,-2,2)^2] [0,(2,2,4)^3] [0,(3,4,7)^3] [24,3^4] [42,(4)^4] [3,0^4]

complexity: **O(wn)**. window size length of input sequence

### Look Up a Request in a Pattern

Assemble all indices at the time of opening FS file for read.

quest comes: read(offset=22, len=1). ook up to what position offset=22 falls in the aginary pattern matrix. For example, (3,4,7)^3]



d the corresponding 6th length and physical Eset of logical offset 22

lculate and return the corresponding physical offset and length of the request.

# Merging Isolated Patterns into a **Global One**

In order to further compress metadata, isolated perprocess index patterns can be merged into a global one.

with a fixed stride of 10 as shown.



*id*[]: [4,7,6,2,8,9,4,7,6,2,8,9] *logical:* 1000,(30)^4 *length:* 10,(0)^4 physical:  $0,(10)^{4}$ 







For example, in the figure below, all processes write

- 4.7.6.2 and so on are PIDs. Blocks of same texture represent data area that is shared by a group of processes, e.g. global strides. The writes from different processes in the figure above can be described by one global pattern:

# **Evaluation**





Write patterns of MILC. In-memory index compression rates by Pattern PLFS (higher is better): (A):37.0; (B):3.0;(C):3.6

### Others



# Conclusions

- This paper proposes efficient and practical techniques to discover structures from unstructured I/O operations, thereby enabling powerful I/O optimizations.
- Several orders of magnitude improvement in the size of the PLFS metadata.
- Up to 40% of write improvement
- Up to 480% of read improvement
- Other Possible Use Cases:
- Pre-fetching, block pre-allocation, and index compression in SciHadoop.

# Acknowledgement

The authors are thankful to Michael Lang (LANL) and Adam Manzanares (California State University) for their help toward this study. This work was performed at the Ultrascale Systems Research Center (USRC) at Los Alamos National Laboratory, supported by the U.S. Department of Energy DE-FC02-06ER25750. The publication has been assigned the LANL identifier LA-UR-12-25954.

