

Using Property Graphs for Rich Metadata Management in HPC Systems

Dong Dai, Robert B. Ross, Philip Carns, Dries Kimpe, and Yong Chen

- The data used to describe other data
 - Simple Metadata

Attribute	Description	Attribute	Description
ino	inode number	ctime	change time
pino	parent inode number	atime	access time
name	file name	owner	file owner
type	file or directory	group	file group
size	file size	mode	file mode
mtime	modification time		

- A. Leung, et. al. Magellan: A Searchable Metadata Architecture for Large-Scale File Systems

- Rich Metadata
- HPC systems heavily rely on these metadata
 - *inode* attributes for file management
 - *location* information for directories and files stored across metadata server
 - *provenance* information partially collected and stored



Busy directory hashed across many MDS's

- S. A. Weil, et. al. Ceph: A Scalable, High-Performance Distributed File System



- Wf4Ever Research Object Model 1.0, <u>http://wf4ever.github.io/ro/</u>





2





- 1. Diverse metadata need to be managed;
- 2. Relationships need to be captured





3







- 1. Diverse metadata need to be managed;
- 2. Relationships need to be captured





3



- 1. Diverse metadata need to be managed;
- 2. Relationships need to be captured





3

Rich Metadata Challenges

Metadata Integration

- diverse metadata should be collected from different components
- diverse metadata should be managed in a unified way

• Storage System Pressure

- large volume of metadata generated from different components
- high concurrent insert rates from parallel applications

Efficient Processing and Querying

- some operations exist in the critical execution path of applications
- some operations require complex query and searching











• Based on Property Graph Model

























Based on Property Graph Model •













Based on Property Graph Model •





















- Entity => Vertex
 - Data Object: represents the basic data unit in storage
 - Executions: represents applications including Jobs, Processes, Threads
 - User: represents real end user of a system
 - Users allowed to define their own entities









- Entity => Vertex
 - Data Object: represents the basic data unit in storage
 - Executions: represents applications including Jobs, Processes, Threads
 - User: represents real end user of a system
 - Users allowed to define their own entities
- Relationship => Edge
 - Relationships between different entities are mapped as edges
 - User *runs* Executions. An edge with type '*Run*' is created between them
 - Reversed relationships also are defined
 - Users allowed to define their own relationships









- Entity => Vertex
 - Data Object: represents the basic data unit in storage
 - Executions: represents applications including Jobs, Processes, Threads
 - User: represents real end user of a system
 - Users allowed to define their own entities
- Relationship => Edge
 - Relationships between different entities are mapped as edges
 - User *runs* Executions. An edge with type '*Run*' is created between them
 - Reversed relationships also are defined
 - Users allowed to define their own relationships
- Attributes => Property
 - On both Entity and Relationship
 - Stored as Key-Value pairs attached on vertices and edges







lsw.org





Sample Graph: Size



Number	Basic	With	With	With		
		I/O Ranks	Full Ranks	Directory		
Vertices	34.6 M	41.7 M	147.8 M	147.9 M		
Edges	126.5 M	133.6 M	239.8 M	366.3 M		
		-	-			
	Road Graph	Twitter [2]	Facebook [3]	Web Page		
	Road Graph USA [1]	Twitter [2]	Facebook [3]	Web Page (2002) [4]		
Vertices	Road Graph USA [1] 24 M	Twitter [2] 645 M	Facebook [3] 1.28 B	Web Page (2002) [4] 2.1 B		
Vertices Edges	Road GraphUSA [1]24 M29 M	Twitter [2] 645 M 81.4 B	Facebook [3] 1.28 B 256 B	Web Page (2002) [4] 2.1 B 15 B		

nne

DNAL RATORY

Sample Graph: Structure



(c) Process node degree (d

(d) File node degree



law distribution



9



Sample Graph: Structure



(d) File node degree



law distribution

•

•



9



Sample Graph: Structure

Common Attribute •

most entities have small



over, $\mathbf{r}(a) \propto a$

Further investigation also • confirm they fit the powerlaw distribution



pdsw.org

(c) Process node degree

(d) File node degree

 $\alpha = 3$ $\alpha = 1.7$

 $\alpha = 0.5$





9

Operations on the Graph: Namespace Traversal

- Hierarchical Namespace Traversal
 - Present logical layout of data sets to users
 - traditional POSIX-style tree-structure directory
 - The metadata graph already contains
 - belongs/contains relationships between Data Objects vertices
 - directory can be considered as Data Object entity too
 - locate files by given path
 - 1. locate the root directory in the graph
 - 2. repeatedly travel through *contains* edges from *directory* vertices to directory or files vertices

Locate -> Traversal -> Filter -> Traversal







N.Org





Operations on the Graph: Data Audit

- Data Audit
 - The metadata graph already contains
 - run relationships between Users and Executions
 - read/write relationships between Executions and Data Objects
 - additional *attributes* are also recorded with these relationships
 - locate files accessed by a specific user in a given time frame
 - 1. locate the given user in the graph
 - 2. travel through run edges from User to Execution
 - 3. filter execution based on the time frame
 - 4. travel through read edges from Executions to Data Objects

Locate -> Traversal -> Filter -> Traversal











Operations on the Graph: Provenance Search

Provenance Support

- Wide range of use cases
 - data sharing, reproducibility, work-flow
- The metadata graph already contains
 - *Relationships* between different entities
 - User-defined attributes and relationships
- #8 in the first Provenance Challenge
 - 1. Use graph to abstract the workflow executions
 - 2. Search all Executions with model "AlignWarp"
 - 3. Travel through *read* edges to Data Objects entities
 - 4. Filter based on property 'center' ('UChicago')

Search Attributes -> Traversal -> Filter -> Traversal











#8 Problems: Given a *fMRI* workflow with multiple stages processing.

Try to find the **Execution whose model** is 'AlignWarp' and inputs have annotation ['center':'UChicago']











Operations on the Graph: Provenance Search

Provenance Support

- Wide range of use cases
 - data sharing, reproducibility, work-flow
- The metadata graph already contains
 - *Relationships* between different entities
 - User-defined attributes and relationships
- #8 in the first Provenance Challenge
 - 1. Use graph to abstract the workflow executions
 - 2. Search all Executions with model "AlignWarp"
 - 3. Travel through *read* edges to Data Objects entities
 - 4. Filter based on property 'center' ('UChicago')

Search Attributes -> Traversal -> Filter -> Traversal











#8 Problems: Given a *fMRI* workflow with multiple stages processing.

Try to find the **Execution whose model** is 'AlignWarp' and inputs have annotation ['center':'UChicago']



Requirements

- Read Search/Locate, Travel, Filter Pattern
 - Search graph vertices and edges by their attributes => **Indexing**
 - Fast locate vertices and edges by global ID => Partitioning
 - Efficient multi-step traversal in large graph => Traversal Speed
 - Customized filter function during traversal => **Filtering**

Write - HPC Environment

- High Volume: rich metadata are actually 'big data'
- Lots of Clients: millions of cores generate metadata concurrently
- High Contention: clients modify the same vertex or edge at the same time
 - Creating files under the same directory
 - All applications read/write the same file







Existing Graph Infrastructure







14



Proposed Solutions

Basic Needs

- Property Graph Model
- Distributed Writes/Reads
- User-defined Indexing
- Graph Traversal

Performance Requirements

High Contention Writes
Efficient Graph Traversal
Consider the Graph Structure

Proposed Solution

╈

Existing Graph Infrastructure (Titan + Cassandra)



Burst Write Partition Strategy Fast Server-side Traversal Caching strategy for power-law like graphs

pdsw.org





15

Prototyped Graph Infrastructure









16



Conclusion

- We observed that a property graph representation seems to be a good match for rich metadata in HPC storage
- We generated an example metadata property graph using access data from a real, large-scale system over a year period
- We observed properties of this graph, compared it to graphs in other contexts, and identified some challenges for processing these graphs for HPC metadata storage











References

- [1] C. Demetrescu, A. V. Goldberg, and D. S. Johnson, The Shortest Path Problem: Ninth DIMACS Implementation Challenge. American Mathematical Soc., 2009, vol. 74.
- [2] "Twitter Statistics," <u>http://www.statisticbrain.com/twitter-statistics/</u>.
- [3] A. Ching, "Giraph: Production-Grade Graph ProcessingInfrastructure for Trillion Edge Graphs," in ATPESC, ser.ATPESC '14, 2014.
- [4] J.-L. Guillaume, M. Latapy et al., "The Web Graph: anOverview," in Actes d'ALGOTEL'02, 2002.
- [5] A. Leung, I. Adams, and E. L. Miller, "Magellan: A Searchable Metadata Architecture for Large-Scale File Systems," University of California, Santa Cruz, Tech. Rep. UCSC-SSRC-09-07, 2009.
- [6] L. Moreau, B. Clifford, J. Freire, J. Futrelle, Y. Gil, P. Groth, N. Kwasnikowska, S. Miles, P. Missier, J. My- ers *et al.*, "The Open Provenance Model Core Specifi- cation (v1. 1)," *Future Generation Computer Systems*, vol. 27, no. 6, pp. 743–756, 2011.
- [7] S. A. Weil, S. A. Brandt, E. L. Miller, D. D. Long, and C. Maltzahn, "Ceph: A Ccalable, High-Performance Distributed File System," in Proceedings of the 7th symposium on Operating Systems Design and Implemen- tation. USENIX Association, 2006, pp. 307–320.









Thanks & Questions