



HopsFS & ePipe

Mahmoud Ismail <maism@kth.se> Gautier Berthou <gautier@sics.se>







- From HDFS to HopsFS
- ePipe to enable a searchable HopsFS
- Tutorial



A File System with a million ops/sec and searchable in sub seconds?!



3

Bill Gates' biggest product regret?



WinFS

- "WinFS was an attempt to bring the benefits of schema and relational databases to the Windows file system. ...The WinFS effort was started around 1999 as the successor to the planned storage layer of Cairo and died in 2006 after consuming many thousands of hours of efforts from really smart engineers."
 - [Brian Welcker]*



$NTFS \leftrightarrow WinFS$



*http://www.zdnet.com/article/bill-gates-biggest-microsoft-product-regret-winfs/

















• Storing NameNode metadata in JVM Heap

- Storing NameNode metadata in JVM Heap
 - Very efficient, yet



- Storing NameNode metadata in JVM Heap
 - Very efficient, yet
 - Number of files/directories are limited

- Storing NameNode metadata in JVM Heap
 - Very efficient, yet
 - Number of files/directories are limited
 - Garbage collection pause times

One Lock to rule them all





One Lock to rule them all



multi-reader, single writer concurrency semantics





• Stateless NameNode

- Stateless NameNode
- Multiple NameNodes to increase Throughput



- Stateless NameNode
- Multiple NameNodes to increase Throughput
- Throughput dependent on our chosen data store



- Stateless NameNode
- Multiple NameNodes to increase Throughput
- Throughput dependent on our chosen data store
- Choosing a data store?

- Stateless NameNode
- Multiple NameNodes to increase Throughput
- Throughput dependent on our chosen data store
- Choosing a data store?
 - An in-memory storage system that can be efficiently queried and managed. Preferably Open-Source.

- Stateless NameNode
- Multiple NameNodes to increase Throughput
- Throughput dependent on our chosen data store
- Choosing a data store?
 - An in-memory storage system that can be efficiently queried and managed. Preferably Open-Source.
 - Row-level locking



- Stateless NameNode
- Multiple NameNodes to increase Throughput
- Throughput dependent on our chosen data store
- Choosing a data store?
 - An in-memory storage system that can be efficiently queried and managed. Preferably Open-Source.
 - Row-level locking
 - Efficient Cross partition transaction





• NewSQL (Relational) DB

- User-defined partitioning
- Row-level Locking
- Distribution-aware transactions
- Partition-pruned index scans
- Real-time, 2-Phase Commit
 - 1.2 sec TransactionInactive timeouts



• NewSQL (Relational) DB

- User-defined partitioning
- Row-level Locking
- Distribution-aware transactions
- Partition-pruned index scans
- Real-time, 2-Phase Commit
 - 1.2 sec TransactionInactive timeouts

- Commodity Hardware
 - Scales to 48 nodes
 - Supports on-disk columns
- SQL API
- C++/Java Native API
- C++ Event API



• NewSQL (Relational) DB

- User-defined partitioning
- Row-level Locking
- Distribution-aware transactions
- Partition-pruned index scans
- Real-time, 2-Phase Commit
 - 1.2 sec TransactionInactive timeouts

200 Million NoSQL Ops/sec

- Commodity Hardware
 - Scales to 48 nodes
 - Supports on-disk columns
- SQL API
- C++/Java Native API
- C++ Event API



HopsFS





From Memory to Database





Apache NameNode Internals

Client: mkdir, getblocklocations, createFile,	Client		
NameNode			

Journal Nodes

HOPSFS & EPIPE



Apache NameNode Internals



Journal Nodes

HOPSFS & EPIPE



Apache NameNode Internals



Journal Nodes




Journal Nodes





Journal Nodes





Journal Nodes















Client: mkdir, getblocklocations, createFile,	Client	
NameNode		

DAL-Impl





DAL-Impl





DAL-Impl





DAL-Impl









14





Fine grain Locking



< 15 >

HOPSFS & EPIPE

Fine grain Locking

• Hierarchical Locking (Implicit Locking)



<

15

(>

Fine grain Locking

- Hierarchical Locking (Implicit Locking)
- Subtree Locking



<

15

(>

Implicit Locking





Implicit Locking





Implicit Locking















Spotify Workload

Op Name	Percentage	Op Name	Percentage
append file	0.0%	content summary	0.01%
mkdirs	0.02%	set permissions	0.03% [26.3%*]
set replication	0.14%	set owner	0.32 % [100%*]
delete	0.75% [3.5%*]	create file	1.2%
rename	1.3% [0.03%*]	add blocks	1.5%
list (<i>listStatus</i>)	9% [94.5%*]	stat (<i>fileInfo</i>)	17% [23.3%*]
read (getBlkLoc)	68.73%		

Relative frequency of operations on a Spotify HDFS cluster. List, read, and stat operations account for \approx 95% of the metadata operations in the cluster.

*Of which, the relative percentage is on directories

[Niazi, Salman, et al. "HopsFS: Scaling Hierarchical File System Metadata Using NewSQL Databases." arXiv preprint (2016)]











Bigger Clusters

	Number of Files		
Memory	HDFS	HopsFS	
1 GB	2.3 million	0.44 million	
50 GB	115 million	22 million	
100 GB	230 million	44 million	
200 GB	460 million	88 million	
500 GB	Does Not Scale	220 million	
1 TB	Does Not Scale	440 million	
24 TB	Does Not Scale	10.8 billion	





• The Database (NDB) is the single source of truth



- The Database (NDB) is the single source of truth
- Extending INodes (files/directories)



- The Database (NDB) is the single source of truth
- Extending INodes (files/directories)
 - Adding a new table with a foreign key to the nodes table



- The Database (NDB) is the single source of truth
- Extending INodes (files/directories)
 - Adding a new table with a foreign key to the nodes table
- Attaching metadata to a file/directory

21

- The Database (NDB) is the single source of truth
- Extending INodes (files/directories)
 - Adding a new table with a foreign key to the nodes table

21

- Attaching metadata to a file/directory
 - Schema-less

- The Database (NDB) is the single source of truth
- Extending INodes (files/directories)
 - Adding a new table with a foreign key to the nodes table
- Attaching metadata to a file/directory
 - Schema-less
 - Schema-based



Schema-less Metadata



Schema-less Metadata

inodelD	Name	parentId	
1	/	0	
2	Users	1	
3	alice.txt	2	


Schema-less Metadata

inodelD	Name	parentId	
1	/	0	
2	Users	1	
3	alice.txt	2	

attach /Users/alice.txt '{"age" : 20, "gender" : "female", "about": "I am alice"}'



Schema-less Metadata

inodelD	Name	parentId	inodelD	metadata
1	/	0		{"age" : 20, "gender" : "female", "about": "I am alice"}
2	Users	1		
3	alice.txt	2		

attach /Users/alice.txt '{"age" : 20, "gender" : "female", "about": "I am alice"}'

22

<

 $\langle \rangle$

Schema-less Metadata

metadata	inodelD	parentId	Name	nodelD
{"age" : 20, "gender" : "female", "about": "I am alice"}		0	/	1
	3	1	Users	2
		2	alice.txt	3

attach /Users/alice.txt '{"age" : 20, "gender" : "female", "about": "I am alice"}'

22

<

(>)



< 23 >





< 23 >























HOPSFS & EPIPE







23 >

 $\langle \rangle$

HOPSFS & EPIPE















hdfs -find



- hdfs -find
 - limited



- hdfs -find
 - limited
 - inefficient by design



- hdfs -find
 - limited
 - inefficient by design
- Pig/MapReduce



- hdfs -find
 - limited
 - inefficient by design
- Pig/MapReduce
- NameNode is a critical point, avoid overloading



- hdfs -find
 - limited
 - inefficient by design
- Pig/MapReduce
- NameNode is a critical point, avoid overloading

24

• SQL Query on NDB

- hdfs -find
 - limited
 - inefficient by design
- Pig/MapReduce
- NameNode is a critical point, avoid overloading
- SQL Query on NDB
- One Size does not fit all

Polyglot Persistence: the Right Tool for the Job





Monolithic vs Polyglot Persistence



[http://martinfowler.com/]





Eventual Consistency for Metadata. Metadata Integrity maintained by Asynchronous Replication and Metadata Immutability.

Database







Eventual Consistency for Metadata. Metadata Integrity maintained by Asynchronous Replication and Metadata Immutability.

Elasticsearch



























HopsFS | ElasticSearch



28 >

<)



Supporting Project-Level Multi-Tenancy



How can we introduce GitHub-style projects to Hadoop?

Tutorial


• Karamel Automated Installation



- Karamel Automated Installation
- <u>http://www.hops.io/?q=boss</u>



- Karamel Automated Installation
- <u>http://www.hops.io/?q=boss</u>
- HopsWorks Clusters



- Karamel Automated Installation
- <u>http://www.hops.io/?q=boss</u>
- HopsWorks Clusters
 - <u>http://52.210.205.125:8080/hopsworks/</u>



- Karamel Automated Installation
- <u>http://www.hops.io/?q=boss</u>
- HopsWorks Clusters
 - <u>http://52.210.205.125:8080/hopsworks/</u>
 - <u>http://52.209.143.68:8080/hopsworks/</u>



- Karamel Automated Installation
- <u>http://www.hops.io/?q=boss</u>
- HopsWorks Clusters
 - <u>http://52.210.205.125:8080/hopsworks/</u>
 - <u>http://52.209.143.68:8080/hopsworks/</u>
 - <u>http://52.18.186.135:8080/hopsworks/</u>

30

Conclusions

- Hops is a next-generation distribution of Hadoop.
- HopsWorks is a frontend to Hops that supports multitenancy, free-text search, interactive analytics with Zeppelin/Flink/Spark, and batch jobs.
- Looking for contributors/committers
 - Check out our github (<u>github.com/hopshadoop</u>)

31















Automated Installation

- Vagrant/Chef to spin up on a single host
- Karamel/Chef to deploy on AWS/GCE/OpenStack or on-premises

name: HopsWorks ec2: type: m3.medium cookbooks: hadoop: github: "hopshadoop/hopsworks-chef" version: "v0.1" groups: ui: size: 1 recipes: - hopsworks metadata: size: 2 recipes: - hops::nn - hops::rm datanodes: size: 50 recipes: - hops::dn - hops::nm



HopsWorks



< 37 >













Alice has only one Kerberos Identity.

Neither attribute-based access control nor dynamic roles supported in Hadoop.



Solution: Project-Specific UserIDs



Member of



Users__Alice



Solution: Project-Specific UserIDs



Member of

Project Users

Users__Alice



Solution: Project-Specific UserIDs



How can we share DataSets between Projects?



Sharing DataSets between Projects

(>)

39

<





Sharing DataSets between Projects





39

<)

(>)

Sharing DataSets between Projects



< 39 >

HopsWorks









Authenticate











HopsWorks









Authenticate









(>)

40

<



(>)

40

<

X.509 Certificate per project specific user



<

Project







- A project is a collection of
 - Members
 - HDFS DataSets
 - Kafka Topics
 - Notebooks, Jobs
- A project has an owner

42

(>

• A project has quotas



Project Roles

Members	×
Find member	Rew members
Members to be added Add members	
Member	Role
test2@kth.se	Data scientist 🕤 💼
	Save
Members	Role Action
Admin Admin (me) admin@kth.se	Data owner 🔹 💼
Test1 Test1 test1@kth.se	Data scientist 🚽 💼
Test3 Test3 test3@kth.se	Data scientist 🔹 💼

- Data Scientist Privileges
 - Write and Run code



Project Roles

Members	×
Find member	e +New members
Members to be added Add members	
Member	Role
test2@kth.se	Data scientist 🕒 📋
	Save
Members	Role Action
Admin Admin (me) admin@kth.se	Data owner 🕤 📋
Test1 Test1 test1@kth.se	Data scientist 🔹 💼
Test3 Test3 test3@kth.se	Data scientist 🕤 🗊

Data Scientist Privileges

- Write and Run code

We delegate administration of privileges to users

