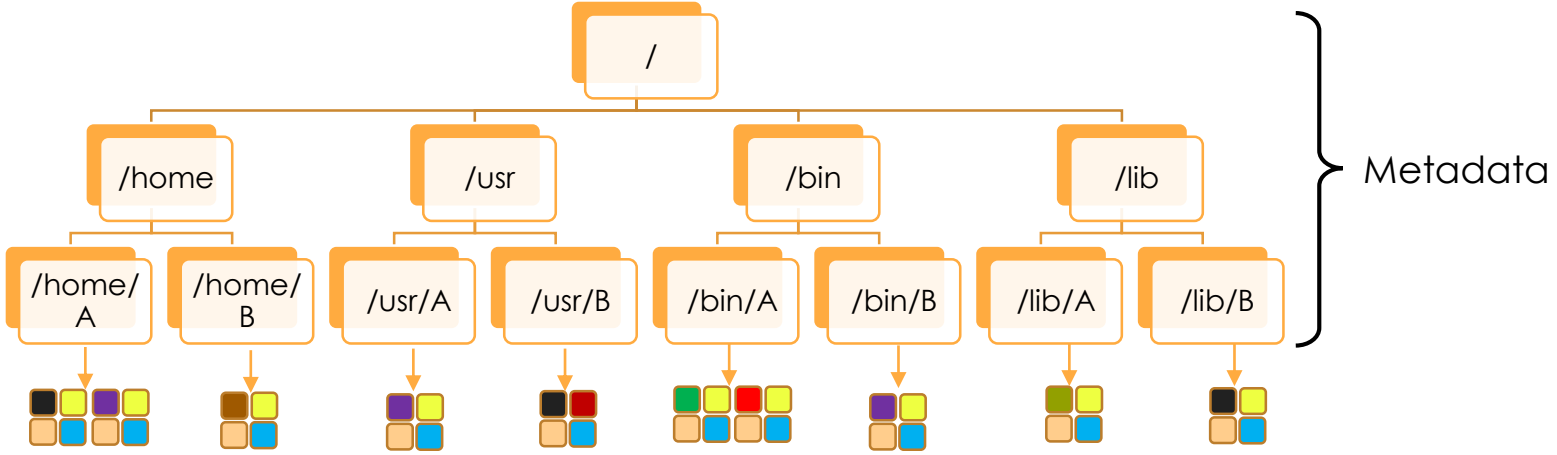# HopsFS
# Scaling Distributed Hierarchical File Systems Using NewSQL Databases
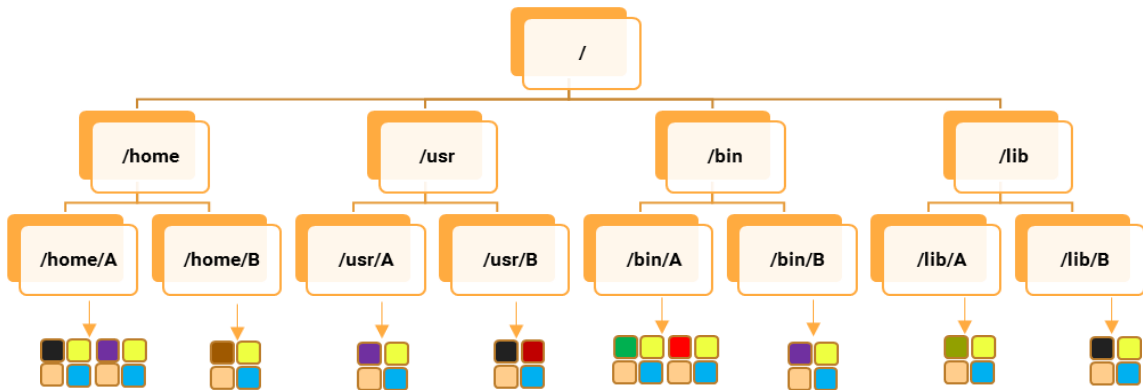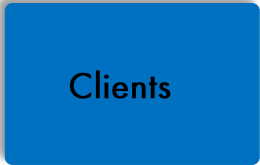
Salman Niazi

# Hierarchical File System



Metadata

- Strongly consistent metadata.
  - Atomic file system operations, such as, move and create

2

# Distributed Hierarchical File System



3

# Typical Hierarchical File System Operation

- `{operation} [flags] {path(s)}`
  - cat /home/F1

# Data Blocks Storage Layer

- Thousands of servers

- Uses data replication and erasure coding for high availability



Block Storage

# Metadata Service Layer

- Atomic File System Operations
- Due to complexity of metadata service monolithic architecture is the most popular solution
  - HDFS, GFS, AFS



Metadata Data

# Why not use Databases?

- Metadata consists of lots of very small data
- Databases specialize is storing and manipulating large amounts of small data



- Traditional databases do not provide high throughput required by distributed file systems
- High operational latencies for resolving file paths

# Using Databases

- WinFS by Microsoft
- GiraffaFS, CassandaraFS, CalvinFS



| inode | ... |
|-------|-----|
| / | |
| /usr | |
| /usr/F1 | |
| /usr/F2 | |
| /usr/F3 | |

Denormalized

| ID | name | PID | |
|----|------|-----|---|
| 0 | / | -1 | |
| 1 | usr | 0 | |
| 2 | F1 | 1 | |
| 3 | F2 | 1 | |
| 4 | F3 | 1 | |

Normalized

# HopsFS



Hadoop Software Stack

# HDFS Architecture

HDFS Client

**File1**

Where can I save the file?

DataNodes Addresses

NameNode

File System Metadata

| File | Blocks Mappings |
|------|-----------------|
| File1 | Blk1 → DN1, Blk2 → DN5, Blk3 → DN3 |
| File2 | Blk1 → DN1, Blk2 → DN4 |
| File3 | Blk1 → DN1, Blk2 → DN2, Blk3 → DN3 |
| File4 | Blk1 → DN100 |
| File5 | Blk1 → DN4, Blk2 → DN2, Blk3 → DN9 |
| | ... ... ... ... |
| FileN | Blk1 → DN2, Blk2 → DN8 |

DataNodes

# HDFS 2.0 High Availability

Standby
NameNode

NameNode

FAIL

Journal Nodes

DataNodes

# HDFS Limitations

- HDFS has been scaled to store 100 PB – 200 PB on 4000 – 5000 datanodes

- Namespace size upper bound: ~ 500 million files

- At most 70 – 80 thousand file system operations / sec

# HopsFS Architecture



Distributed Database

File System Metadata

NameNode

DataNodes

| File | Blocks Mappings | | |
|------|-----------------|---|---|
| File1 | Metadata | 2 | DN5 |
| File2 | Metadata | 2 | |
| File3 | Metadata | 2 | DN3 |
| File4 | Metadata | | |
| File5 | | 2 | DN9 |
| FileN | Blk1 → DN2, Blk2 → DN8 | | |

| File | Blocks Mappings |
|------|-----------------|
| File3 | Metadata |
| File4 | Metadata |
| File7 | Metadata |
| File8 | Metadata |

| File | Blocks Mappings |
|------|-----------------|
| File6 | Metadata |

# NewSQL DB

# MySQL Cluster: **Network Database Engine (NDB)**

- **Open Source**
- **Commodity Hardware**
    - Scales to 48 database nodes
        - 200 Million Read Ops/Sec*  using NDB native API
    - Read Committed Transaction Isolation
        - Row-level Locking
    - User-defined partitioning

*https://www.mysql.com/why-mysql/benchmarks/mysql-cluster/

# Transaction Isolation

# Transaction Isolation Levels

MySQL Cluster Network Database Engine only supports Read-Committed Transaction
Isolation Level

| Isolation level | Dirty reads | Non-repeatable reads | Phantoms |
|---|---|---|---|
| Read Uncommitted | may occur | may occur | may occur |
| **Read Committed** | - | **may occur** | **may occur** |
| Repeatable Read | - | - | may occur |
| Serializable | - | - | - |

# Read Committed Transaction Isolation Level

Transaction 1

Transaction 2

Start Tx
Read  x ( x = 1 )

Start Tx
Read x ( x = 1 )
Update x = 10

Read x ( x = 1 )

Commit

Read x ( x = 10 )

**The value of X has changed before the transaction is committed**

# Read Committed Transaction Isolation

create file /user/F4

create file /user/F4

F1  F2  F3

# Read Committed Transaction Isolation



Client 1

**Start Transaction**
  Metadata Read
  • /
  • User
  • F4 (Does not exists)
  New Metadata
  • F4
**Commit Transaction**

Client 2

**Start Transaction**
  Metadata Read
  • /
  • User
  • F4 (Does not exists)
  New Metadata
  • F4
**Commit Transaction**

Read root dir
Read user dir
Read F4 file

Read root dir
Read user dir
Read F4 file

/
user
F4  F1  F2  F3  F4

# Read Committed Transaction Isolation

Use row level locking to serialize conflicting file  operation

# Read Committed Transaction Isolation With Locking



Client 1

**Start Transaction**
  Metadata Read
  ● /
  ● User
  ● F4
  throw FileAlreadyExists Exception
**Abort Transaction**

Client 2

**Start Transaction**
  Metadata Read
  ● /
  ● User
  ● F4 (Does not exists)
  New Metadata
  ● F4
**Commit Transaction**

Read root dir

Read user dir using exclusive locks

Read F4 file

Read root dir

Read user dir using exclusive locks

Read F4 file

/

User

F1  F2  F3  F4

# Database Operations
# &
# Data Partitioning

# Distributed Metadata Design

# Distributed Metadata Design

| ID | name | PID |
|----|------|-----|
| 1 | home | 0 |
| 7 | C | 2 |
| 10 | F | 3 |
| | ... | |

DB Node 0

| ID | name | PID |
|----|------|-----|
| 5 | A | 1 |
| 8 | D | 2 |
| 4 | lib | 0 |
| | ... | |

DB Node 1

| ID | name | PID |
|----|------|-----|
| 6 | B | 1 |
| 3 | bin | 0 |
| 11 | G | 4 |
| | ... | |

DB Node 2

| ID | name | PID |
|----|------|-----|
| 2 | usr | 0 |
| 9 | E | 3 |
| 12 | H | 4 |
| | ... | |

DB Node 3

TC    TC    TC    TC

ls  /home

/

home    usr    bin    lib

A    B    C    D    E    F    G    H

25

# Distributed Full Table Scans

drwxrwx--- /home/A
dr_____-- /home/B

| TC | ID | name | PID |
|----|----|------|-----|
| | 1 | home | 0 |
| | 7 | C | 2 |
| | 10 | F | 3 |
| | | ... | |

DB Node 0

| TC | ID | name | PID |
|----|----|------|-----|
| | 5 | A | 1 |
| | 8 | D | |
| | 4 | | |
| | | ... | |

| TC | ID | | |
|----|----|--|--|
| | usr | 0 | |
| | 9 | E | 3 |
| | 12 | H | 4 |
| | | ... | |

DB Node 2    DB Node 3

**Full Table Scan Operations
Do not Scale
Max 5.5 K ops/ sec using 4 node NDB Cluster**

**Select * from Inodes where PID = 1**

ls /home

CPU    Network    Latency

# Distributed Index Scan Operations

drwxrwx--- /home/A
drwxrwx--- /hom



| TC | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ID | name | PID | ID | name | PID | ID | name | PID | | | |
| 1 | home | | 5 | A | | 6 | B | 1 | | | |
| 7 | C | 2 | 8 | D | 2 | 3 | bin | | | | |
| 10 | F | 3 | 4 | lib | 0 | 11 | | | | | 4 |
| | ... | | | ... | | | ... | | | | |

DB Node 0          DB Node 1          DB Node 3

**Distributed Index Scan Operations**
**Scale Poorly**
**210 K ops/sec (590 times better than DFTS)**

```
Select * from Inodes where
PID = 1
WITH(INDEX(…))
```

ls – l /home

CPU          Network          Latency

# HopsFS

- Uses NewSQL Relational Database that allows

1.  **User Defined Partitioning (UDP):** Take control of how the data is distributed across difference database nodes

2.  **Distribution Aware Transactions (DAT):** Take control over which Transaction Coordinator handles which file system operations.

# Solution (User Defined Partitioning)

DB Node 0

| ID | Name | PID |
|----|------|-----|
| 1  | /    | 0   |
| 9  | E    | 4   |
| 10 | F    | 4   |
| 15 | K    | 8   |
| 15 | L    | 8   |

DB Node 1

| ID | Name | PID |
|----|------|-----|
| 2  | home | 1   |
| 3  | usr  | 1   |
| 4  | bin  | 1   |
| 11 | G    | 5   |
| 12 | H    | 5   |

DB Node 2

| ID | Name | PID |
|----|------|-----|
| 5  | A    | 2   |
| 6  | B    | 2   |
| 16 | M    | 10  |
| 17 | N    | 10  |
|    |      |     |

DB Node 3

| ID | Name | PID |
|----|------|-----|
| 7  | C    | 3   |
| 8  | D    | 3   |
| 13 | I    | 7   |
| 14 | J    | 7   |
|    |      |     |

**Hash Fn**
PID % 4 = Partition No



29

# Solution (Partition Pruned Index Scan Operations)

drwxrwx--- ... ome/A
drwxrw... ...me/B

| TC | Name | PID | | TC | Name | PID | | TC | | Name | |
|---|------|-----|---|---|------|-----|---|---|---|------|---|
| 1 | / | 0 | | 2 | home | 1 | | 5 | | | |
| 9 | E | 4 | | 3 | usr | 1 | | | | | |
| 10 | F | 4 | | 4 | bin | 1 | | | | | 7 |
| 15 | K | 8 | | 11 | | | | | J | | 7 |
| 15 | L | 8 | | | | | | | | | |

DB Node 0     DB Node 2     DB Node 3

**Partition Pruned Index Scan**

**Scales Well**

```
Select * from Inodes where
PID = 2
WITH(INDEX(…))
```

ls – l /home


CPU    Network    Latency

30

# Solution (Distribution Aware Transactions)

```
drwxrwx---  /home/A
drwxrwx---  /home/B
```

TC  TC  TC  TC

| ID | Name | PID |
|----|------|-----|
| 1 | / | |
| 9 | E | 4 |
| 10 | F | 4 |
| 15 | K | 8 |
| 15 | L | 8 |

DB Node 0

| ID | Name | PID |
|----|------|-----|
| 2 | home | 1 |
| 3 | usr | 1 |
| 4 | bin | 1 |
| 11 | | |

| ID | Name | PID |
|----|------|-----|
| 5 | A | |

| | | |
|----|------|-----|
| | | 7 |

DB Node 3

Distribution Aware Transactions &
Partition Pruned Index Scan
**Scales Very Well**
~ 1Million ops/ sec (5X DIS)

```
Start Transaction on Node 2
Select * from Inodes where
PID = 2
WITH(INDEX(name))
```

ls – l /home

CPU        Network        Latency

31

Transactional FS Operations

- File System Operation
  - **Distributed Transaction**
    - START DISTRIBUTION AWARE TRANSACTION
      - Primary Key Ops
      - Partition Pruned Index Scan Ops
      - Batching and Caching
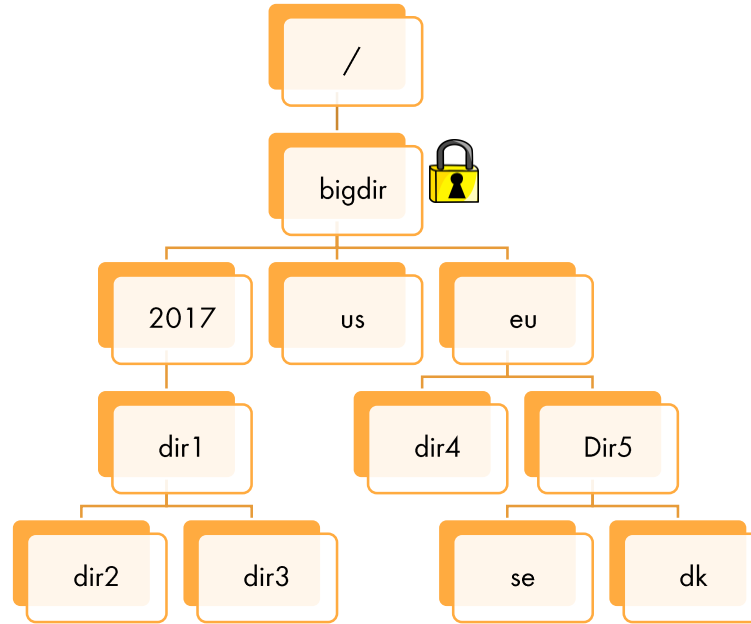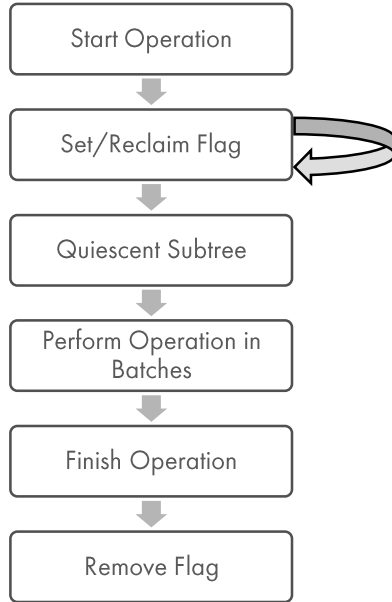    - COMMIT TRANSACTION

# Solution (Contd.)

- ~40 Tables & hundreds of FS operations.

- Most NewSQL databases do not provide serializable transaction isolation level

  - HopsFS uses **read committed** transaction isolation level and **row level locks**

- Avoiding deadlocks in file system operation.

  - Total Order Locking

- Implementing large file system operations that do not fit in a transaction

# Large File System Operations

# Subtree Operations

**Subtree operation Stages**

- Start Operation
- Set/Reclaim Flag
- Quiescent Subtree
- Perform Operation in Batches
- Finish Operation
- Remove Flag

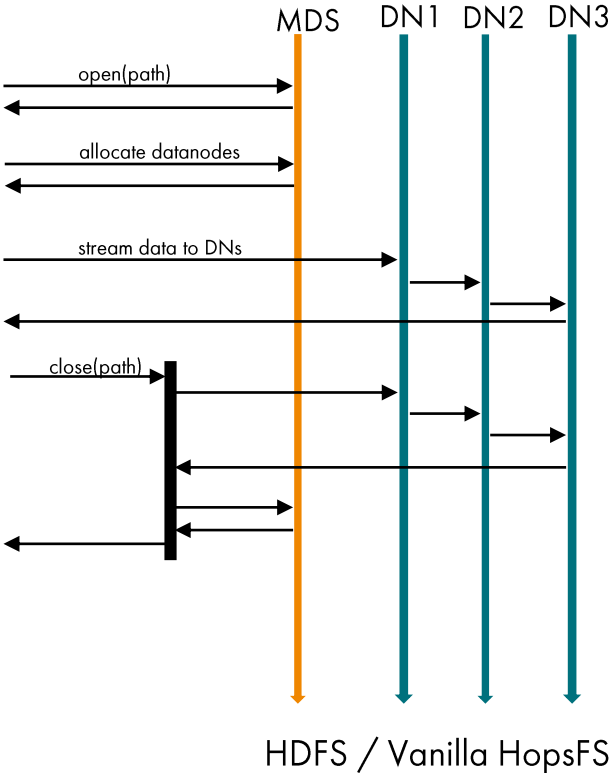# Failures during Subtree Operations

All subtree operations are implemented in such a way that if the operations fail halfway, then the namespace is not left in an inconsistent state.
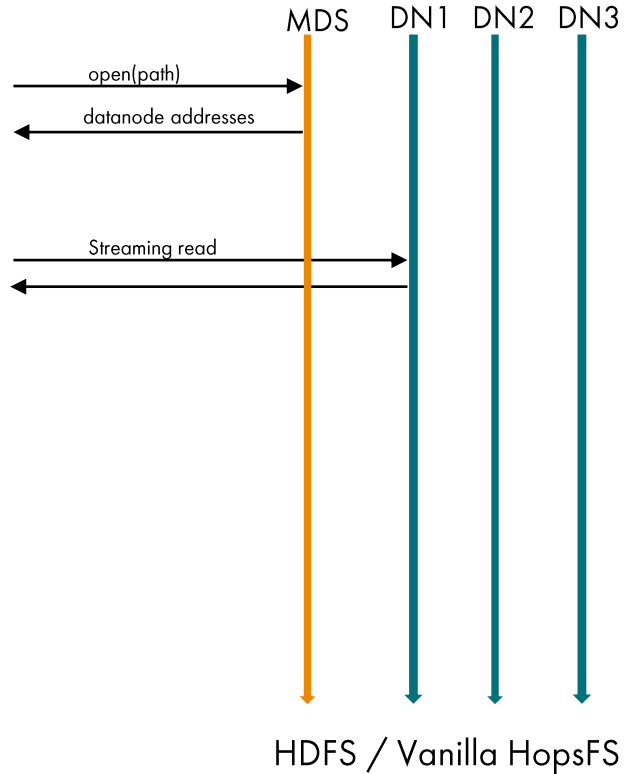
# HopsFS Performance

- ## Throughput

  - **16X** the throughput of HDFS (Spotify Workload).

  - **38X** the throughput of HDFS for 20% write intensive workload

- ## Low Latency

  - Identical avg op latency (~3ms) for small number (50) of clients

  - **10X** lower latency for large number (6500) of clients

- ## Metadata Scaling

  - **37X** more metadata.
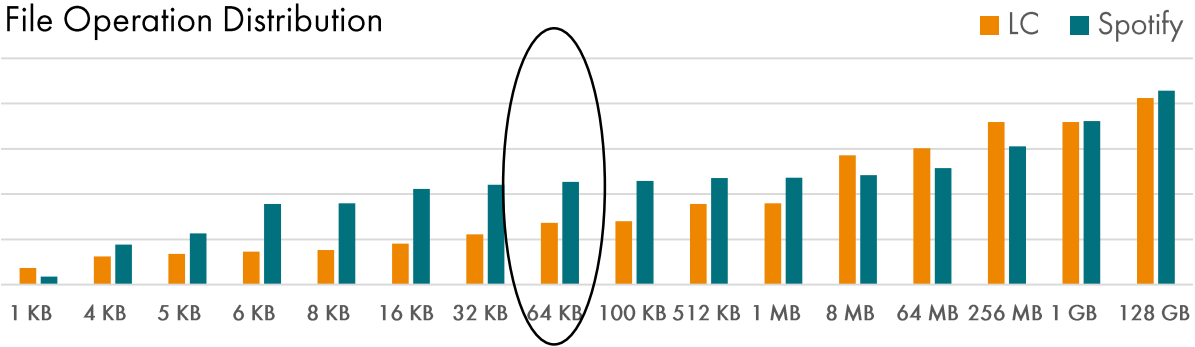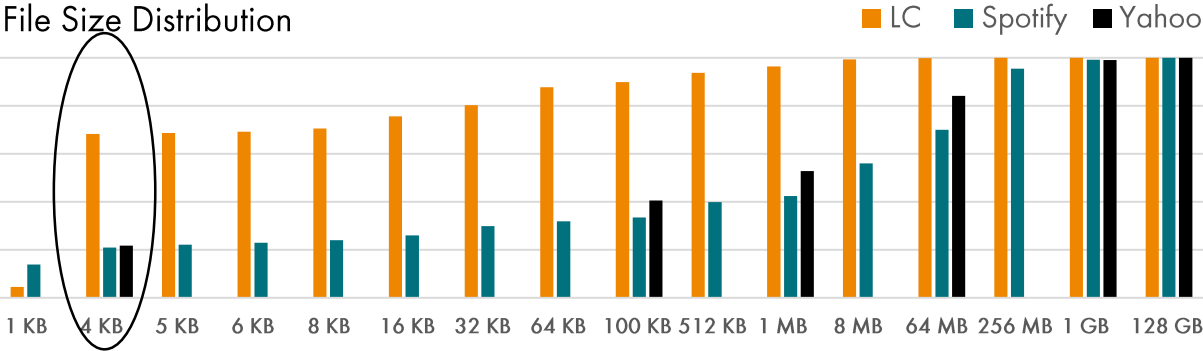
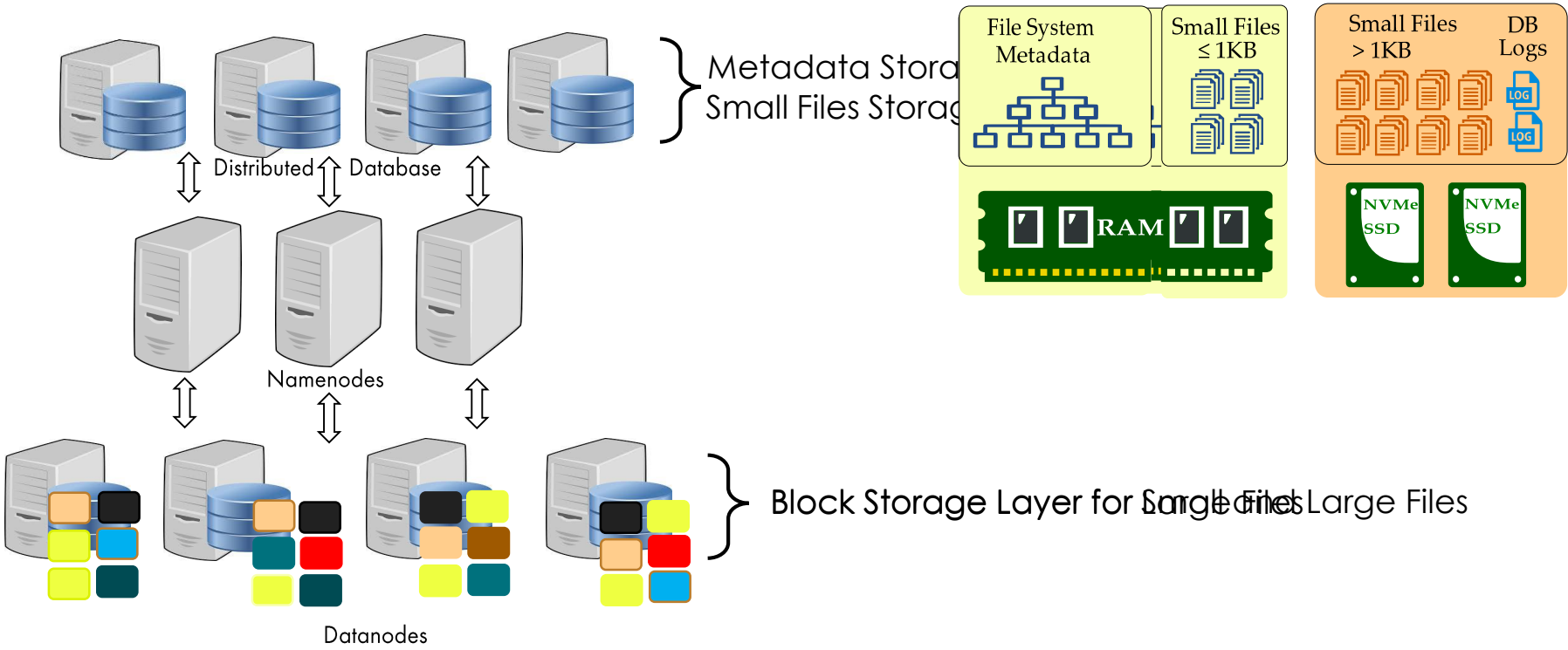# Support for Small Files

# HDFS/HopsFS Write Operation



HDFS / Vanilla HopsFS

# HDFS/HopsFS Read Operation



HDFS / Vanilla HopsFS

# Prevalence of Small Files In Hadoop



File Size Distribution

File Operation Distribution

# Small Files' Support in HopsFS

Metadata Storage &
Small Files Storage

Distributed  Database

Namenodes

File System
Metadata

Small Files
≤ 1KB

Small Files
> 1KB

DB
Logs

LOG

LOG

RAM

NVMe
SSD

NVMe
SSD

Block Storage Layer for Small and Large Files

Datanodes

# Optimizing Write Operation For Small Files



HDFS / Vanilla HopsFS

MDS   DN1   DN2   DN3

open(path)

allocate datanodes

stream data to DNs

close(path)

Optimized HopsFS

MDS   DN1   DN2   DN3

open

close (path, data)

**61X** Throughput

**Using only 6 NVMe SSDs**

43

# Optimizing Read Operations For Small Files



open(path)

datanode addresses

Streaming read

MDS   DN1   DN2   DN3

HDFS / Vanilla HopsFS

open(path)

return data

MDS   DN1   DN2   DN3

**4X** Throughput

**Using only 6 NVMe SSDs**

Optimized HopsFS

# Questions

http://www.hops.io
http://github.com/hopshadoop
@hopshadoop

Read More

Scaling Distributed Hierarchical File Systems Using NewSQL Databases
Salman Niazi. Ph.D. Thesis. KTH Royal Institute of Technology

Scaling hierarchical file system metadata using newsql databases
S Niazi, M Ismail, S Haridi, J Dowling, S Grohsschmiedt, M Ronström
15th USENIX Conference on File and Storage Technologies (FAST 17), 89-104

Scaling HDFS to more than 1 million operations per second with HopsFS
M Ismail, S Niazi, M Ronström, S Haridi, J Dowling
2017 17th IEEE/ACM International Symposium on Cluster, Cloud and Grid …

Size Matters: Improving the Performance of Small Files in Hadoop
S Niazi, M Ronström, S Haridi, J Dowling
Proceedings of the 19th International Middleware Conference, 26-39