# LocoFS

NHPCC Room 300      Friday, December 8th, 2017      Daniel Shao

# Design and Implementation

## Loosely-Coupled Arch.

Consists of: **Client, DMS, FMS, Object Store**
- **DMS:** Directory Metadata Server. Only 1
  - Enough to hold around 100 million directories in 32GB memory.
  - Simple ACL Management.
- **FMS:** File Metadata Server. Multiple

KV Pattern: **HASHING**
- DMS: full pathname ➡ directory metadata
- FMS: dir_uuid+filename ➡ file metadata

## Rename Discussion

Problem: hashing
- File: only metadata needs relocation
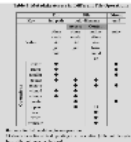- Directory: its metadata as well as all successors' metadata need relocation.

## Flattened Directory Tree

Motivation: DESTROY directory tree
- Backward Directory Entry Organization
- Client Caching: only directories' metadata



## Decoupled File Metadata

Motivation:
- Large-Value access
- (De)Serialization

Tech:
- Fine-grained File Metadata
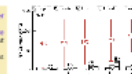- Indexing Metadata Removal
- (De)Serialization



## Motivation
### Problem with FS Directory Tree in DFS



## Motivation
### Gap between FS Metadata and KV Store

Existing file systems have much lower performance than KV stores.
- It has been confirmed that more than half operations are about metadata in file systems.
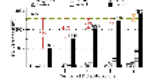- KV Stores have great advantages on small objects.



## Q & A

## Evaluation
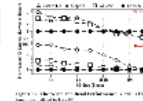### Metadata Performance

- mdtest: 1 million files each time

1. Latency
2. Throughput
3. Bridging gap



## Evaluation
### Full System Performance

Benchmark: not mentioned

# LocoFS
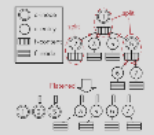
# Design and Implementation

## Loosely-Coupled Arch.

Consists of: **Client, DMS, FMS, Object Store**
- **DMS:** Directory Metadata Server. Only 1
  - Enough to hold around 100 million directories in 32GB memory.
  - Simple ACL Management.
- **FMS:** File Metadata Server. Multiple

KV Pattern: **HASHING**
- DMS: full pathname ➡ directory metadata
- FMS: dir_uuid+filename ➡ file metadata

## Rename Discussion

Problem: hashing
- File: only metadata needs relocation
- Directory: its metadata as well as all successors' metadata need relocation.

## Flattened Directory Tree

Motivation: DESTROY directory tree
- Backward Directory Entry Organization
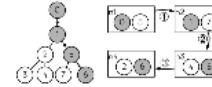- Client Caching: only directories' metadata

## Decoupled File Metadata

Motivation:
- Large-Value access
- (De)Serialization

Tech:
- Fine-grained File Metadata
- Indexing Metadata Removal
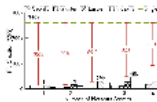- (De)Serialization

## Motivation
### Problem with FS Directory Tree in DFS

## Motivation
### Gap between FS Metadata and KV Store

Existing file systems have much lower performance than KV Stores:
- It has been confirmed that more than half operations are about metadata in file systems.
- KV Stores have great advantages on small objects.

> Q & A

## Evaluation
### Metadata Performance

- mdtest: 1 million files each time

1. Latency
2. Throughput
3. Bridging gap

## Evaluation
### Full System Performance

Benchmark: not mentioned

# LocoFS

NHPCC Room 300          Friday, December 8th, 2017          Daniel Shao

# Design and Implementation

## Loosely-Coupled Arch.

## Flattened Directory Tree

Consists of Client, DMS, FMS, Object Store
- **DMS**: Directory Metadata Server. Only 1
  - Enough to hold around 100 million

Motivation: DESTROY directory tree
- Backward Directory



## Motivation
**Problem with FS Directory Tree in DFS**



## Motivation

# Motivation

## Problem with FS Directory Tree in DFS

# Motivation
## Gap between FS Metadata and KV Store

Existing file systems have much lower performance than KV Stores:

- It has been confirmed that more than half operations are about metadata in file systems.
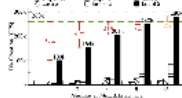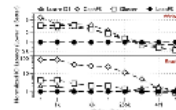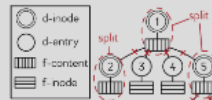- KV Stores have great advantages on small objects

# Design and Implementation

## Loosely-Coupled Arch.
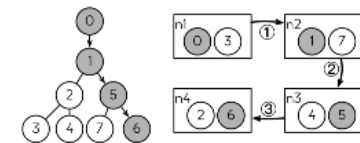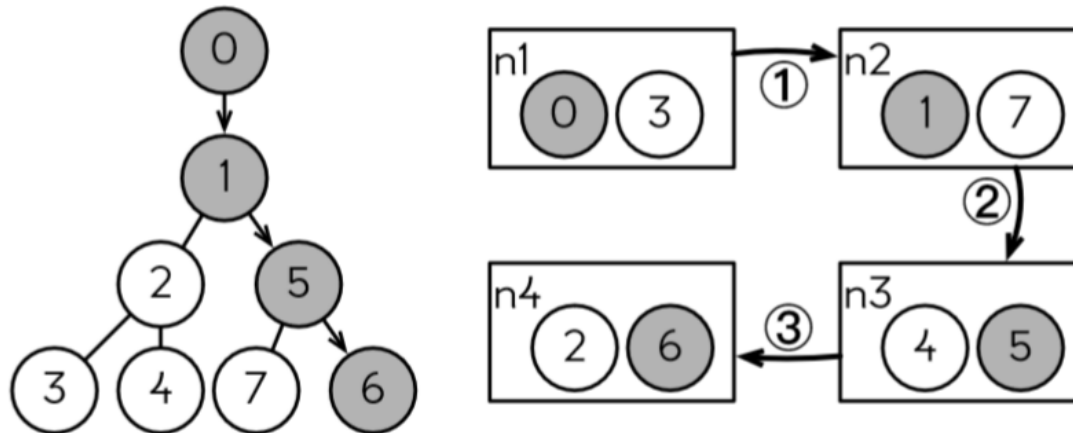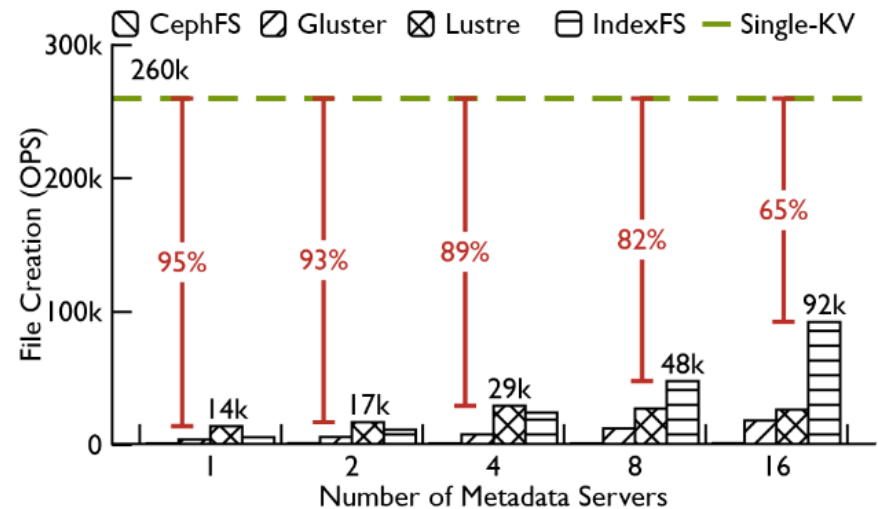
Consists of: **Client, DMS, FMS, Object Store**
- **DMS**: Directory Metadata Server. Only 1
  - Enough to hold around 100 million directories in 32GB memory.
  - Simple ACL Management.
- **FMS**: File Metadata Server. Multiple

KV Pattern: **HASHING**
- DMS: full pathname ➡ directory metadata
- FMS: dir_uuid+filename ➡ file metadata

## Rename Discussion

Problem: hashing
- File: only metadata needs relocation
- Directory: its metadata as well as all successors' metadata need relocation.

## Flattened Directory Tree

Motivation: DESTROY directory tree
- Backward Directory Entry Organization
- Client Caching: only directories' metadata



## Decoupled File Metadata

Motivation:
- Large-Value access
- (De)Serialization

Tech:
- Fine-grained File Metadata
- Indexing Metadata Removal
- (De)Serialization

Table 1: Metadata Access in Different File Operations

| | | Dir | File | | Dirent |
|---|---|---|---|---|---|
| Key | | full path | uuid+filename | | uuid |
| | | | Access | Content | |
| Value | | ctime | ctime | mtime | entry |
| | | mode | mode | atime | |
| | | uid | uid | size | |
| | | gid | gid | bsize | |
| | | uuid | | suuid | |
| | | | | sid | |
| Operations | mkdir | ● | | | ● |
| | rmdir | ● | | | ● |
| | readdir | ● | | | |
| | getattr | ● | ● | ● | ● |
| | remove | | ● | | |
| | chmod | ● | ● | | |
| | chown | | ● | | |
| | create | | ● | | ● |
| | open | | | ○ | |
| | read | | | ● | |
| | write | | | ● | |
| | truncate | | | ● | |

● stands for field updating in an operation.
○ stands for optional field updating in an operation (different file system have different implementations).

# Design and Implementat

## Loosely-Coupled Arch.

Consists of: **Client, DMS, FMS, Object Store**
- **DMS**: Directory Metadata Server. Only 1
  - Enough to hold around 100 million directories in 32GB memory.
  - Simple ACL Management.
- **FMS**: File Metadata Server. Multiple

KV Pattern: **HASHING**
- DMS: full pathname ➡ directory metadata
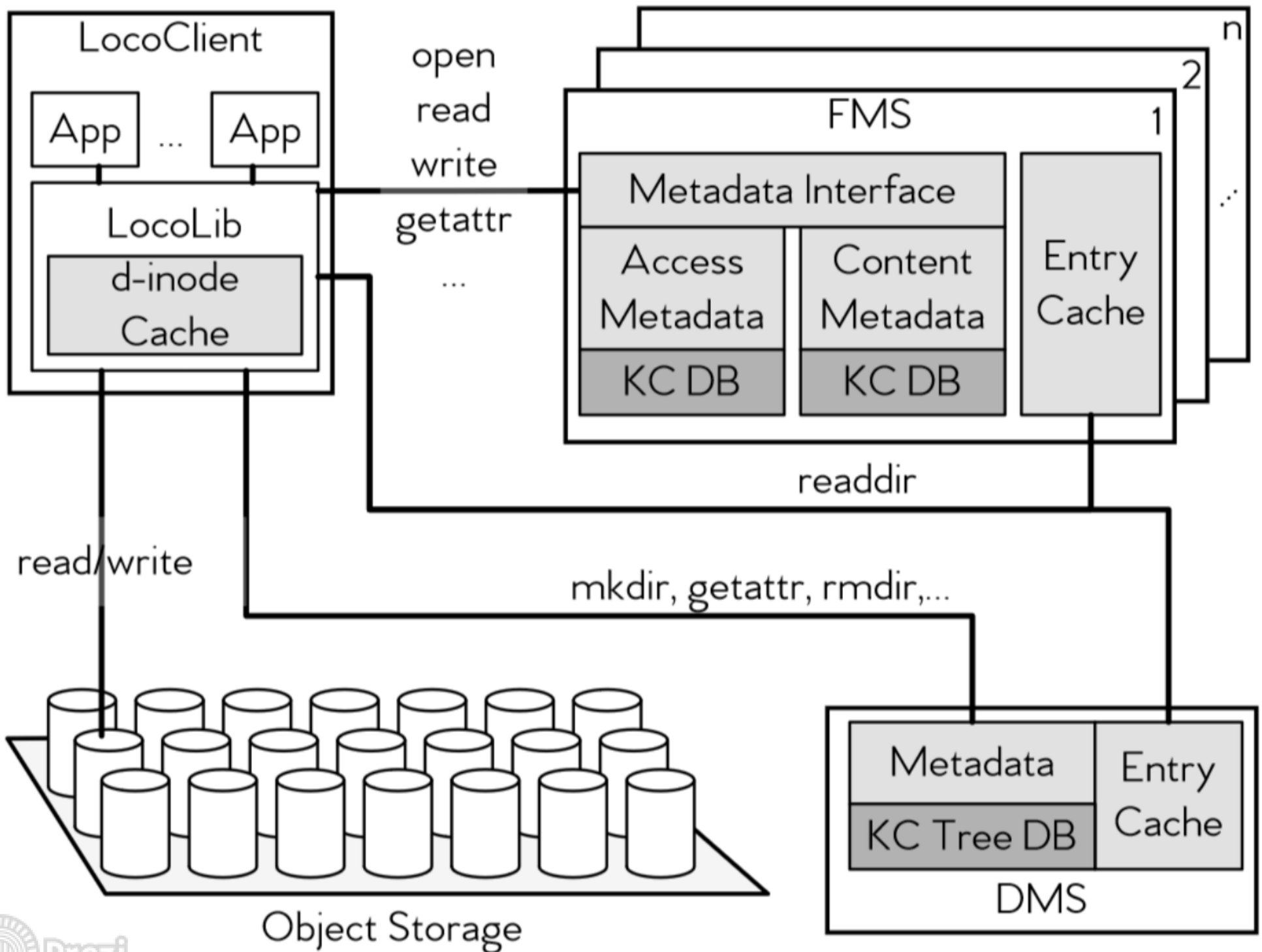- FMS: dir_uuid+filename ➡ file metadata

## Rename Discussion

## Flatten

Motivation:
directory tre
- Backwar
  Entry Or
- Client C
  directorie

## Decoup

Motivation:
- Large-V

# Design and Implementati[on]

## Loosely-Coupled Arch.

Consists of: **Client, DMS, FMS, Object Store**
- **DMS**: Directory Metadata Server. Only 1
  - Enough to hold around 100 million directories in 32GB memory.
  - Simple ACL Management.
- **FMS**: File Metadata Server. Multiple

KV Pattern: **HASHING**
- DMS: full pathname ➜ directory metadata
- FMS: dir_uuid+filename ➜ file metadata

## Flattene[d]

Motivation:    D[...]
directory tree [...]
- Backward [...]
  Entry Orga[...]
- Client Cac[...]
  directories'[...]

## Decoupl[ ]

Motivation:
- Large-Valu[...]
  (De)Serial[...]

## Rename Discussion

## Flattened Directory Tree

Motivation: DESTROY directory tree
- Backward Directory Entry Organization
- Client Caching: only directories' metadata



## Decoupled File Metadata

Existing
perform
- It h
  tha

on

① ② ⑤  ③ ④ ⑥ ⑦

# Decoupled File Metadata

Motivation:
- Large-Value access
- (De)Serialization

Tech:
- Fine-grained File Metadata
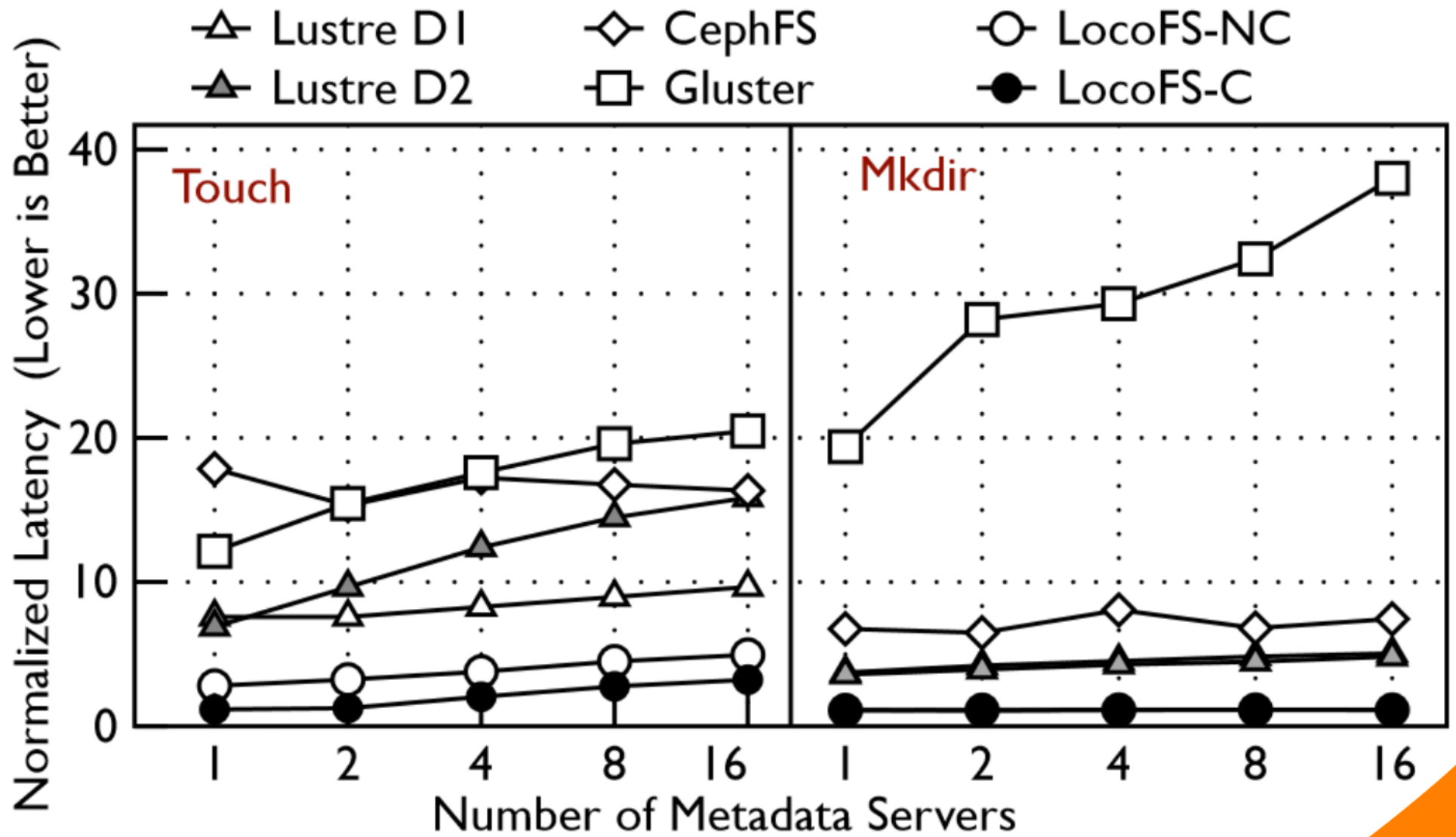- Indexing Metadata Removal
- (De)Serialization

**Table 1: Metadata Access in Different File Operations**

| | | Dir | File | | Dirent |
|---|---|---|---|---|---|
| | **Key** | full path | uuid+filename | | uuid |
| | | | **Access** | **Content** | |
| | **Value** | ctime<br>mode<br>uid<br>gid<br>uuid | ctime<br>mode<br>uid<br>gid | mtime<br>atime<br>size<br>bsize<br>suuid<br>sid | entry |
| **Operations** | mkdir | ● | | | ● |
| | rmdir | ● | | | ● |
| | readdir | ● | | | ● |
| | getattr | ● | ● | ● | |
| | remove | | ● | ● | ● |
| | chmod | ● | ● | | |
| | chown | ● | ● | | |
| | create | | ● | | ● |
| | open | | ● | ○ | |
| | read | | | ● | |
| | write | | | ● | |
| | truncate | | | ● | |

● stands for field updating in an operation.
○ stands for optional field updating in an operation (different file system have different implementations).

- Enough to hold around 100 million directories in 32GB memory.
- Simple ACL Management.
- **FMS**: File Metadata Server. Multiple

KV Pattern: **HASHING**
- DMS: full pathname ➡ directory metadata
- FMS: dir_uuid+filename ➡ file metadata

- Backward Di
  Entry Organiza
- Client Caching
  directories' me

# Decouplec

# Rename Discussion

Problem: hashing
- File: only metadata needs relocation
- Directory: its metadata as well as all successors' metadata need relocation.

Motivation:
- Large-Value a
- (De)Serializat
Tech:
- Fine-grained
  Metadata
- Indexing Meta
  Removal
- (De)Serializat

Evaluation

Evaluation

# Evaluation
## Metadata Performance

- mdtest: 1 million files each time

1. Latency
2. Throughput
3. Bridging gap

# Evaluation
## Metadata Performance

- mdtest: 1 million files each time

1. Latency
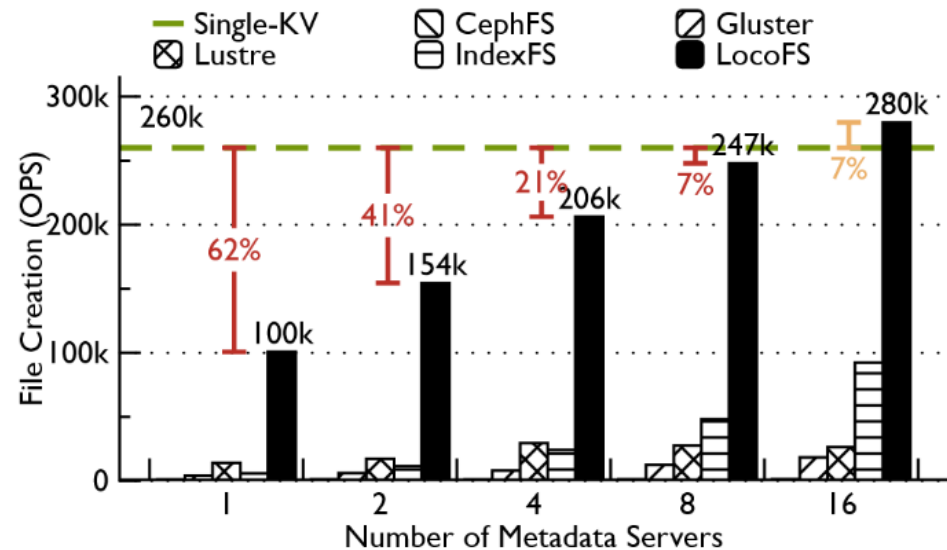2. Throughput
3. Bridging gap

# Evaluation

## Metadata Performance

- mdtest: 1 million files each time

1. Latency
2. Throughput
3. Bridging gap

# Evaluation
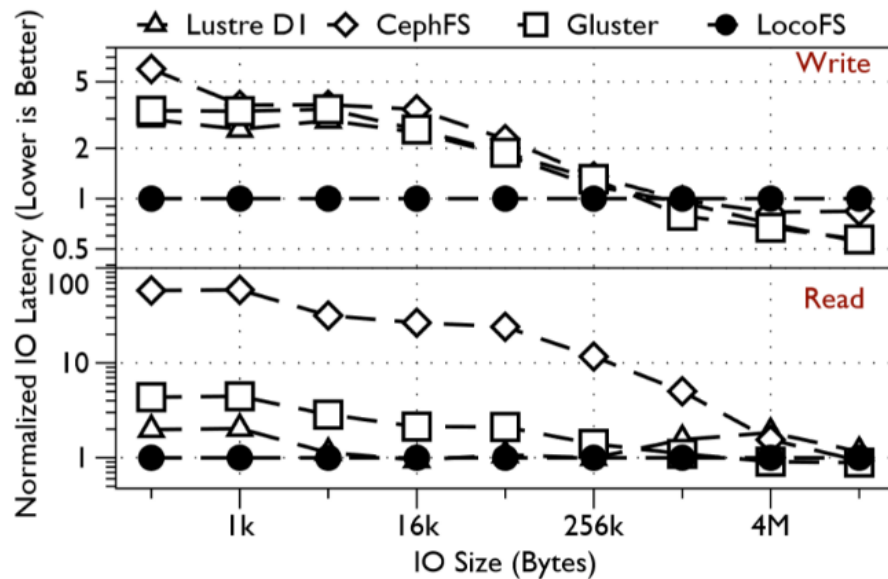
## Full System Performance

Benchmark:
not mentioned



**Figure 12: The Write and Read Performance.** Y-axis is the latency normalized to LocoFS.
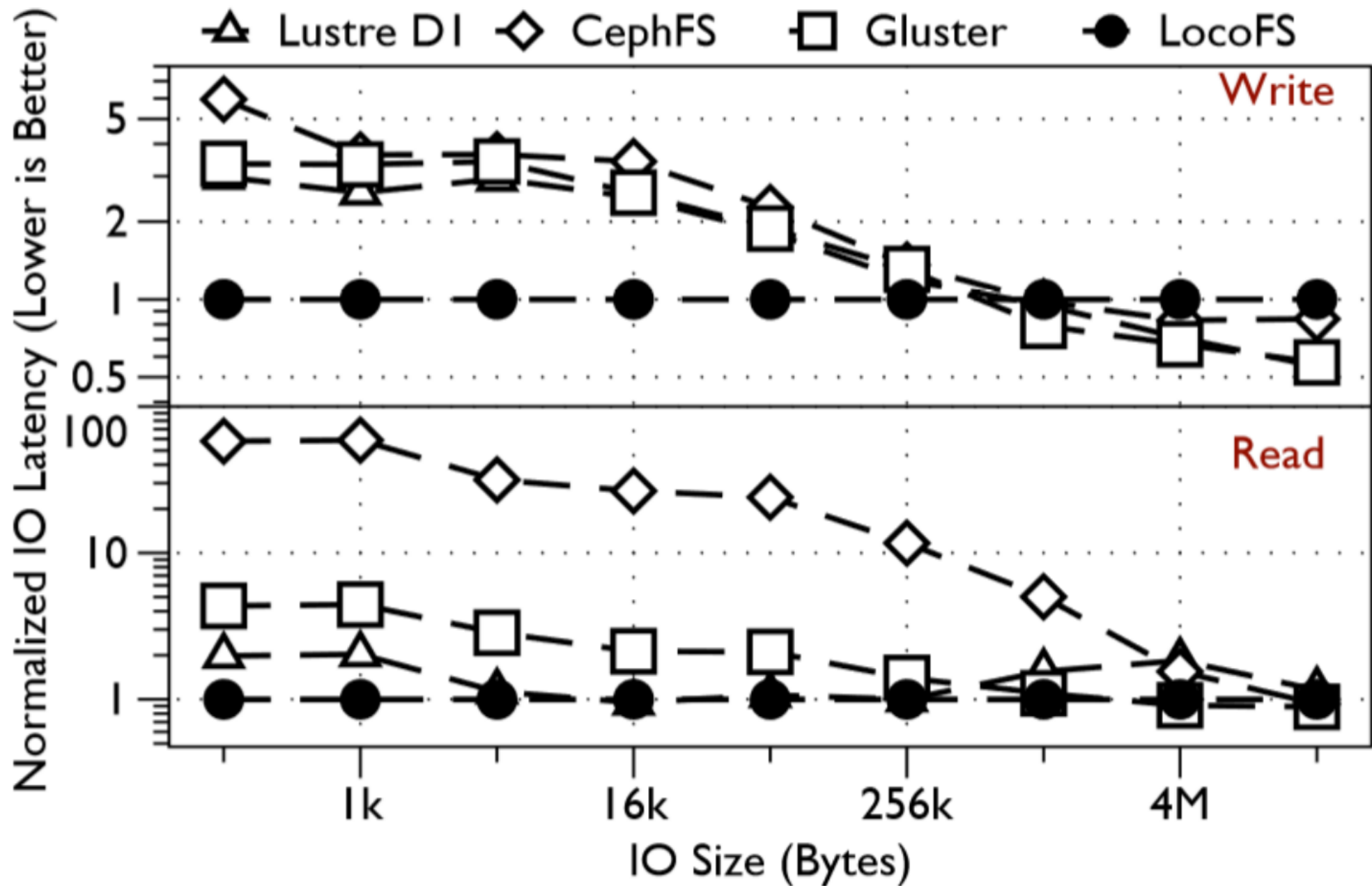
**Figure 12: The Write and Read Performance.** Y-axis is the latency normalized to LocoFS.

# Evaluation
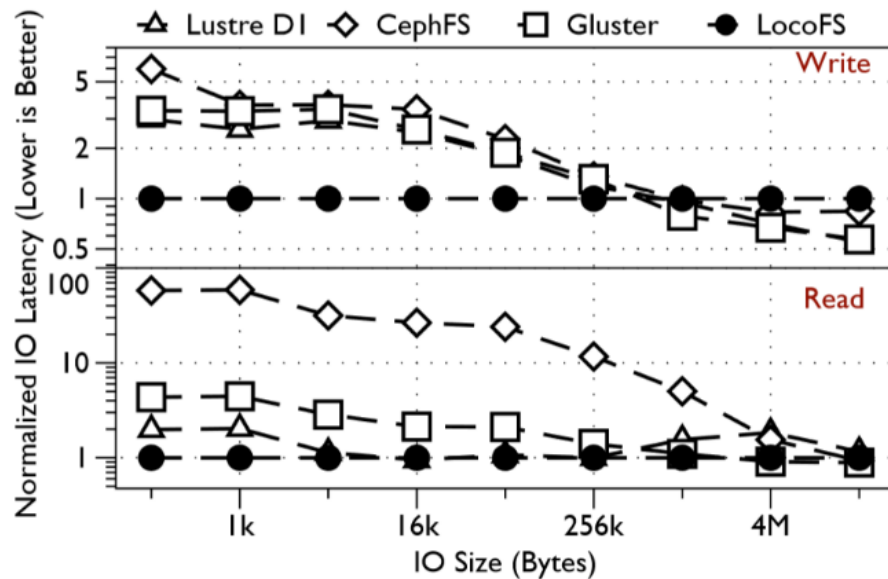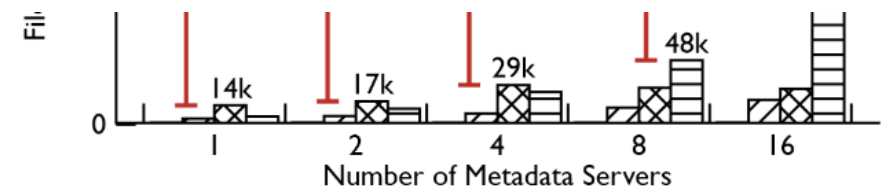
## Full System Performance

Benchmark:
not mentioned



**Figure 12: The Write and Read Performance.** Y-axis is the latency normalized to LocoFS.

- KV Stores have great advantages on small objects



14k   17k   29k   48k

Number of Metadata Servers

Q & A

# LocoFS

NHPCC Room 300 — Friday, December 8th, 2017 — Daniel Shao

# Design and Implementation

## Loosely-Coupled Arch.

Consists of: **Client, DMS, FMS, Object Store**
- **DMS**: Directory Metadata Server. Only 1
  - Enough to hold around 100 million directories in 32GB memory.
  - Simple ACL Management.
  - **FMS**: File Metadata Server. Multiple
- KV Pattern: **HASHING**
  - DMS: full pathname ➞ directory metadata
  - FMS: dir_uuid+filename ➞ file metadata

## Rename Discussion

Problem: hashing
- File: only metadata needs relocation
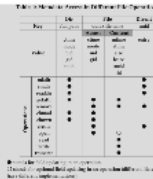- Directory: its metadata as well as all successors' metadata need relocation.

## Flattened Directory Tree

Motivation:  DESTROY directory tree
- Backward Directory Entry Organization
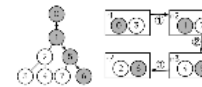- Client Caching: only directories' metadata

## Decoupled File Metadata

Motivation:
- Large-Value access
- (De)Serialization
Tech:
- Fine-grained File Metadata
- Indexing Metadata Removal
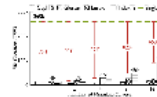- (De)Serialization

## Motivation
### Problem with FS Directory Tree in DFS

## Motivation
### Gap between FS Metadata and KV Store

Existing file systems have much lower performance than KV Stores:
- It has been confirmed that more than half operations are about metadata in file systems.
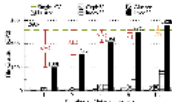- KV Stores have great advantages on small objects.

> Q & A

## Evaluation
### Metadata Performance

- mdtest: 1 million files each time

1. Latency
2. Throughput
3. Bridging gap

## Evaluation
### Full System Performance

Benchmark:
not mentioned