# Rethinking File Mapping in Persistent Memory

**Ian Neal**[1], Gefei Zuo[1], Eric Shiple[1], Tanvir Ahmed Khan[1],
Youngjin Kwon[3], Simon Peter[2], Baris Kasikci[1]

# Good News! PM has arrived!

# Good News! PM has arrived!

- Great, fast new storage technology called persistent main memory (PM)
  - AKA non-volatile memory (NVM)

# Good News! PM has arrived!

- Great, fast new storage technology
  called persistent main memory (PM)
  - AKA non-volatile memory (NVM)
- 30−40x faster than SSDs

# Good News! PM has arrived!

- Great, fast new storage technology called persistent main memory (PM)
  - AKA non-volatile memory (NVM)
- 30−40x faster than SSDs
- File system IO performance has not been able to keep up with PM performance. Why?

# Good News! PM has arrived!

- Great, fast new storage technology called persistent main memory (PM)
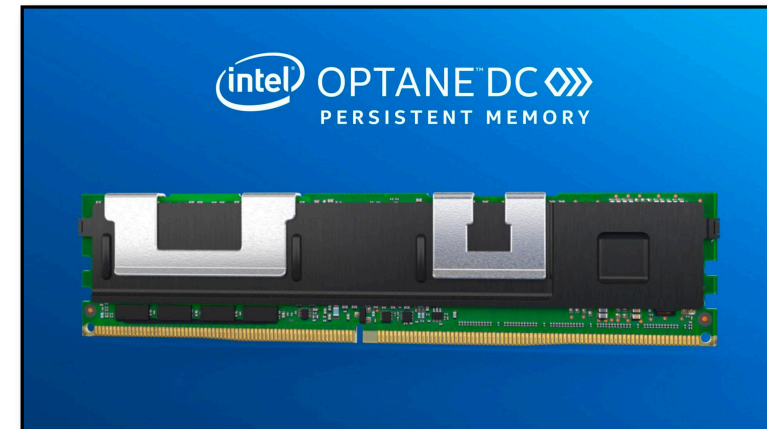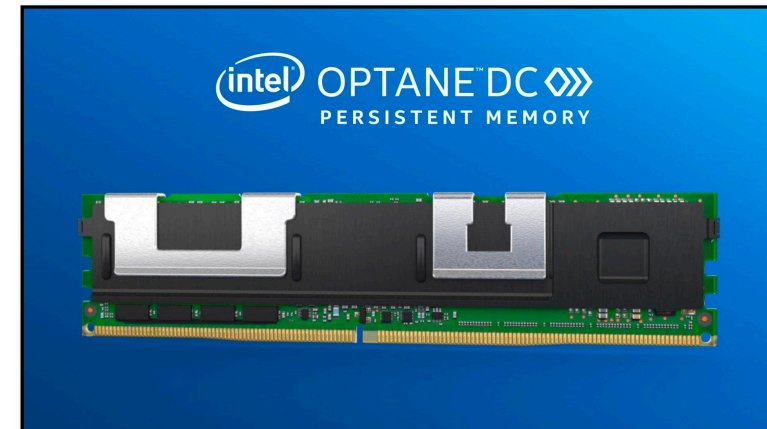  - AKA non-volatile memory (NVM)
- 30−40x faster than SSDs
- File system IO performance has not been able to keep up with PM performance. Why?
- *No rigorous analysis of IO path performance!*

# File IO Before PM

# File IO Before PM



**write()**

| |
|---|
| **file**: /home/ian/notes.txt |
| **offset**: 0x40000000 |
| **data**: "Neque porro quisquam est qui dolorem ..." |

Application

# File IO Before PM



write()

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."

Application

inode num: 0x101
logical block: 0x40000

*File System*

3

# File IO Before PM



**write()**

**file**: /home/ian/notes.txt

**offset**: 0x40000000

**data**: "Neque porro quisquam est qui dolorem ..."

Application

**inode num**: 0x101

**logical block**: 0x40000

*File Mapping*

Block allocator

*File System*

3

# File IO Before PM



```
inode num: 0x101
logical block: 0x40000
```

**File Mapping**

```
write()
file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."
```

Application

Block allocator

Page Cache (DRAM)

DMA

V-NAND SSD
970 EVO Plus
NVMe M.2
SAMSUNG
www.samsung.com/ssd
SAMSUNG ELECTRONICS CO., LTD.    500GB

- *File mapping cache*
- *File data cache*

*File System*

3

# File IO Before PM

# File IO Before PM



write()

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."

Application

inode num: 0x101
logical block: 0x40000

File Mapping

Block allocator

device offset: ...

Page Cache (DRAM)

DMA

V-NAND SSD
970 EVO Plus
NVMe M.2
www.samsung.com/ssd
SAMSUNG ELECTRONICS CO., LTD.
500GB
SAMSUNG

- File mapping cache
- File data cache

*File System*

3

# File IO Before PM

- When devices were slow, FS software overheads were negligible



```
inode num: 0x101
logical block: 0x40000
```

**File Mapping**

```
write()
file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam
est qui dolorem ..."
```
Application

```
device offset: ...
```

Block allocator

Page Cache (DRAM)

DMA

- *File mapping cache*
- *File data cache*

*File System*

# File IO After PM

# File IO After PM



inode num: 0x101

logical block: 0x40000

**File Mapping**

device offset: ...

write()

file: /home/ian/notes.txt

offset: 0x40000000

data: "Neque porro quisquam est qui dolorem ..."

Application

Block allocator

Page Cache (DRAM)

- *File mapping cache*
- *File data cache*

*File System*

4

# File IO After PM



write()

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."

Application

inode num: 0x101
logical block: 0x40000

File Mapping

device offset: ...

Block allocator

Page Cache (DRAM)

- *File mapping cache*
- *File data cache*

*File System*

4

# File IO After PM

- FS components must be optimized for PM



```
inode num: 0x101
logical block: 0x40000
```

**File Mapping**

```
device offset: ...
```

```
write()
file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."
```

Application

Block allocator

Page Cache (DRAM)

- *File mapping cache*
- *File data cache*

*File System*

# File IO After PM

- FS components must be optimized for PM
- Many have in prior work



Application

write()

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."

inode num: 0x101
logical block: 0x40000

File Mapping

device offset: ...

Block allocator

Page Cache (DRAM)

- File mapping cache
- File data cache

Strata (SOSP'17)  ZoFS (SOSP'19)

NOVA (FAST'16)  SplitFS (SOSP'19)

PMFS (EuroSys'14)  ...

File System

4

# File IO After PM

- FS components must be optimized for PM
- Many have in prior work



```
inode num: 0x101
logical block: 0x40000
```

**File Mapping**

```
device offset: ...
```

```
write()

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam
est qui dolorem ..."
```

Application

Block allocator

Page Cache (DRAM)

Strata (SOSP'17)   ZoFS (SOSP'19)

NOVA (FAST'16)   SplitFS (SOSP'19)

PMFS (EuroSys'14)   ...

*File System*

4

# File IO After PM

- FS components must be optimized for PM
- Many have in prior work



write()

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."

Application

inode num: 0x101
logical block: 0x40000

File Mapping

device offset: ...

Block allocator

Page Cache (DRAM)

File System

Strata (SOSP'17)

ZoFS (SOSP'19)

NOVA (FAST'16)

SplitFS (SOSP'19)

PMFS (EuroSys'14)

...

4

# File IO After PM

- FS components must be optimized for PM
- Many have in prior work



**write()**

file: /home/ian/notes.txt

offset: 0x40000000

data: "Neque porro quisquam est qui dolorem ..."

Application

inode num: 0x101

logical block: 0x40000

Block allocator

Page Cache (DRAM)

*File Mapping*

device offset: ...

Strata (SOSP'17)
ZoFS (SOSP'19)
NOVA (FAST'16)
SplitFS (SOSP'19)
PMFS (EuroSys'14)
...

*File System*

4

# File IO After PM

- FS components must be optimized for PM
- Many have in prior work



```
inode num: 0x101
logical block: 0x40000
```

**File Mapping**

```
device offset: ...
```

```
write()

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam
est qui dolorem ..."
```
Application

Block allocator

Page Cache (DRAM)

Strata (SOSP'17)   ZoFS (SOSP'19)

NOVA (FAST'16)   SplitFS (SOSP'19)

PMFS (EuroSys'14)   ...

*File System*

4

# File IO After PM

- FS components must be optimized for PM
- Many have in prior work



**write()**

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."

Application

inode num: 0x101
logical block: 0x40000

*File Mapping*

Block allocator

device offset: ...

Page Cache (DRAM)

Strata (SOSP'17)  ZoFS (SOSP'19)

NOVA (FAST'16)  SplitFS (SOSP'19)

PMFS (EuroSys'14)  ...

*File System*

4

# File IO After PM

- FS components must be optimized for PM
- Many have in prior work

inode num: 0x101

logical block: 0x40000

Can comprise up to *70%* of the IO path overhead!

data: "Neque porro quisquam est qui dolorem ..."

Application

Page Cache (DRAM)

Strata (SOSP'17)
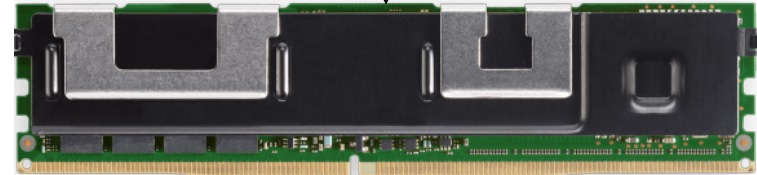
ZoFS (SOSP'19)

NOVA (FAST'16)

SplitFS (SOSP'19)

PMFS (EuroSys'14)

...

*File System*

4

# Our Contributions

# Our Contributions

- Rigorously analyze file mapping in PM

# Our Contributions

- Rigorously analyze file mapping in PM

- Optimize legacy PM file mapping structures

# Our Contributions

- Rigorously analyze file mapping in PM

- Optimize legacy PM file mapping structures

- Design new PM-optimized file mapping approaches

# Our Contributions

- Rigorously analyze file mapping in PM

- Optimize legacy PM file mapping structures

- Design new PM-optimized file mapping approaches

- Evaluate end-to-end performance on real workloads

# Our Contributions

- Rigorously analyze file mapping in PM

- Optimize legacy PM file mapping structures

- Design new PM-optimized file mapping approaches

- Evaluate end-to-end performance on real workloads

# Analysis Overview

# Analysis Overview

| How is file mapping affected by… | Design Question |
| --- | --- |
| Page caching? | Is page caching necessary? |
| File size? | Specialize for different file sizes? |
| IO size? | Optimize for sequential access? |
| Space utilization? | Make file mapping structure elastic? |
| Concurrency? | Is ensuring isolation important? |
| Locality? | Optimize for specific workloads? |
| Fragmentation? | Make robust against file system aging? |
| Storage structures? | Can we reuse PM storage structures? |
| Real workloads? | Are mapping optimizations impactful? |

# Analysis Overview

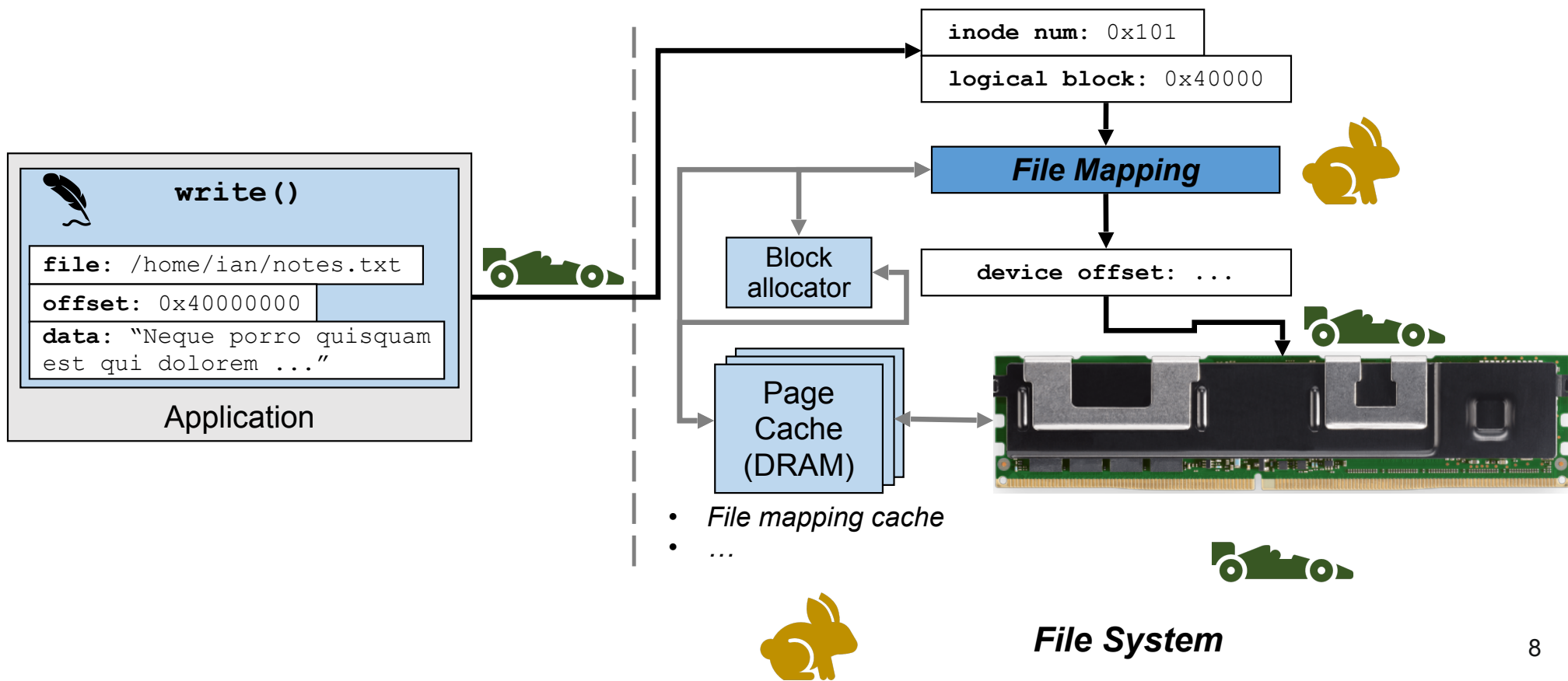| How is file mapping affected by… | Design Question |
|---|---|
| **Page caching?** | **Is page caching necessary?** |
| File size? | Specialize for different file sizes? |
| IO size? | Optimize for sequential access? |
| Space utilization? | Make file mapping structure elastic? |
| Concurrency? | Is ensuring isolation important? |
| Locality? | Optimize for specific workloads? |
| Fragmentation? | Make robust against file system aging? |
| Storage structures? | Can we reuse PM storage structures? |
| **Real workloads?** | **Are mapping optimizations impactful?** |

# Analysis Setup

# Analysis Setup

- Evaluation setup:
  - Implemented in Strata (SOSP'17)
  - Baseline mapping structure: Extent trees in the page cache (Strata default)
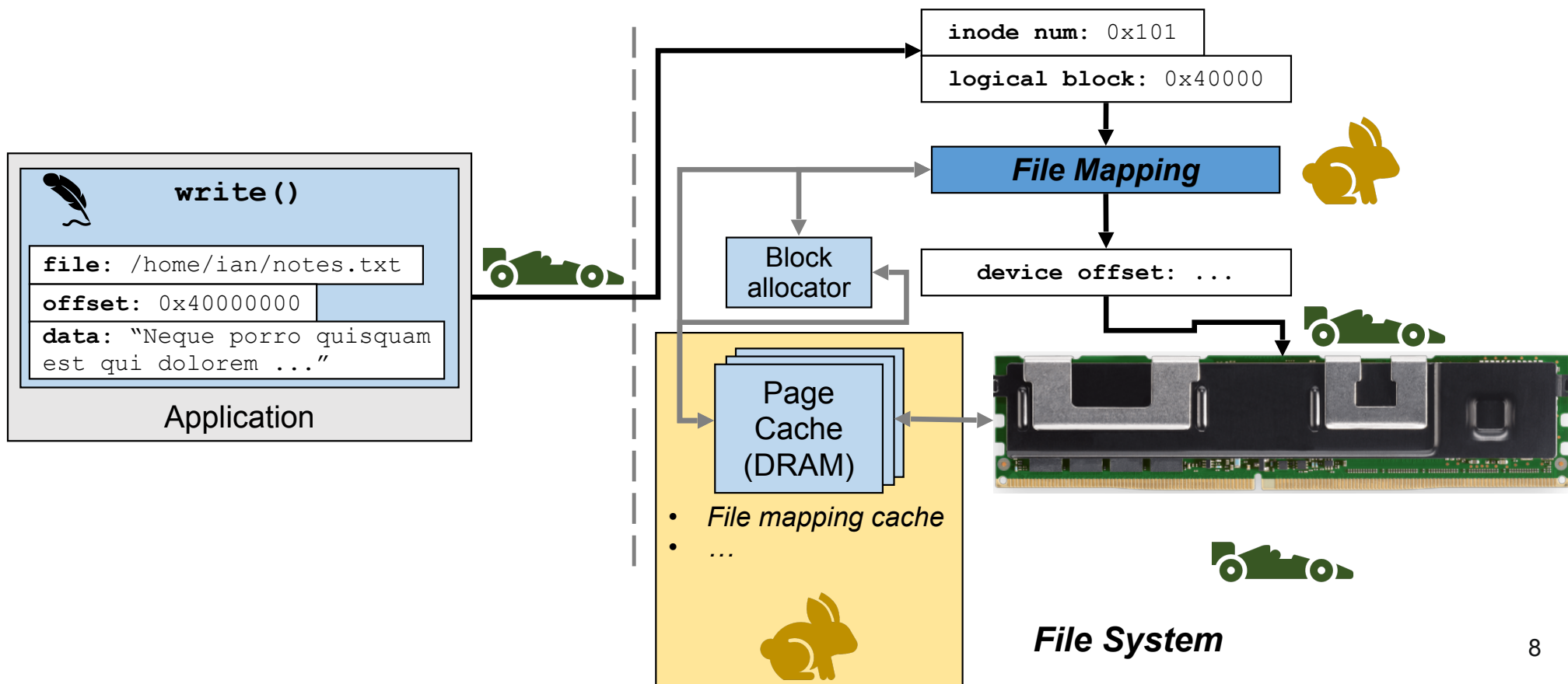  - Evaluated on 256 GB Intel Optane DC NVDIMMs

# Analysis Setup

- Evaluation setup:
  - Implemented in Strata (SOSP'17)
  - Baseline mapping structure: Extent trees in the page cache (Strata default)
  - Evaluated on 256 GB Intel Optane DC NVDIMMs

- Analysis performed using YCSB on LevelDB
  - YCSB: Popular key-value store workload
  - LevelDB: Popular key-value store (used in original Strata evaluation)

# Analysis Setup

- Evaluation setup:
  - Implemented in Strata (SOSP'17)
  - Baseline mapping structure: Extent trees in the page cache (Strata default)
  - Evaluated on 256 GB Intel Optane DC NVDIMMs

- Analysis performed using YCSB on LevelDB
  - YCSB: Popular key-value store workload
  - LevelDB: Popular key-value store (used in original Strata evaluation)

# Analysis Setup

- Evaluation setup:
  - Implemented in Strata (SOSP'17)
  - Baseline mapping structure: Extent trees in the page cache (Strata default)
  - Evaluated on 256 GB Intel Optane DC NVDIMMs

- Analysis performed using YCSB on LevelDB
  - YCSB: Popular key-value store workload
  - LevelDB: Popular key-value store (used in original Strata evaluation)

- Evaluate on Filebench
  - fileserver (1:2 read/write ratio)
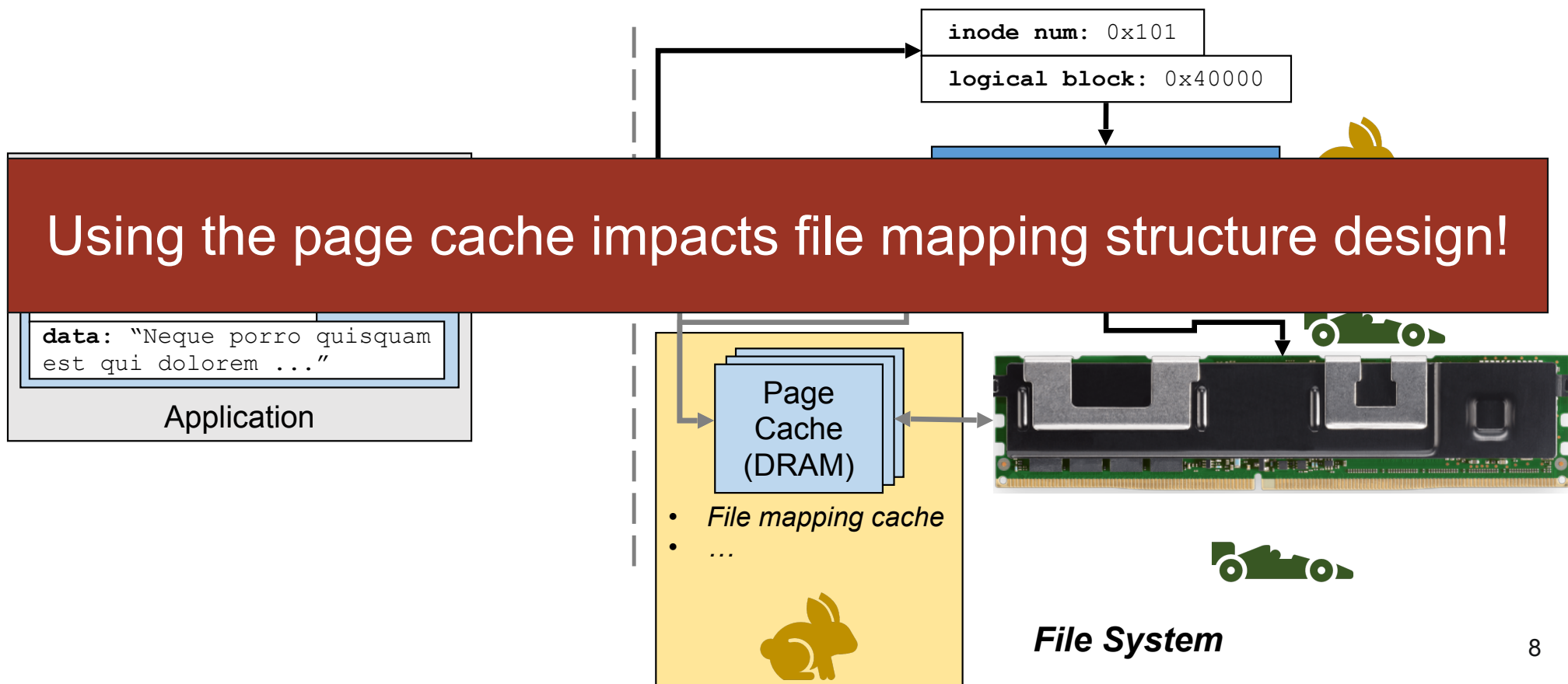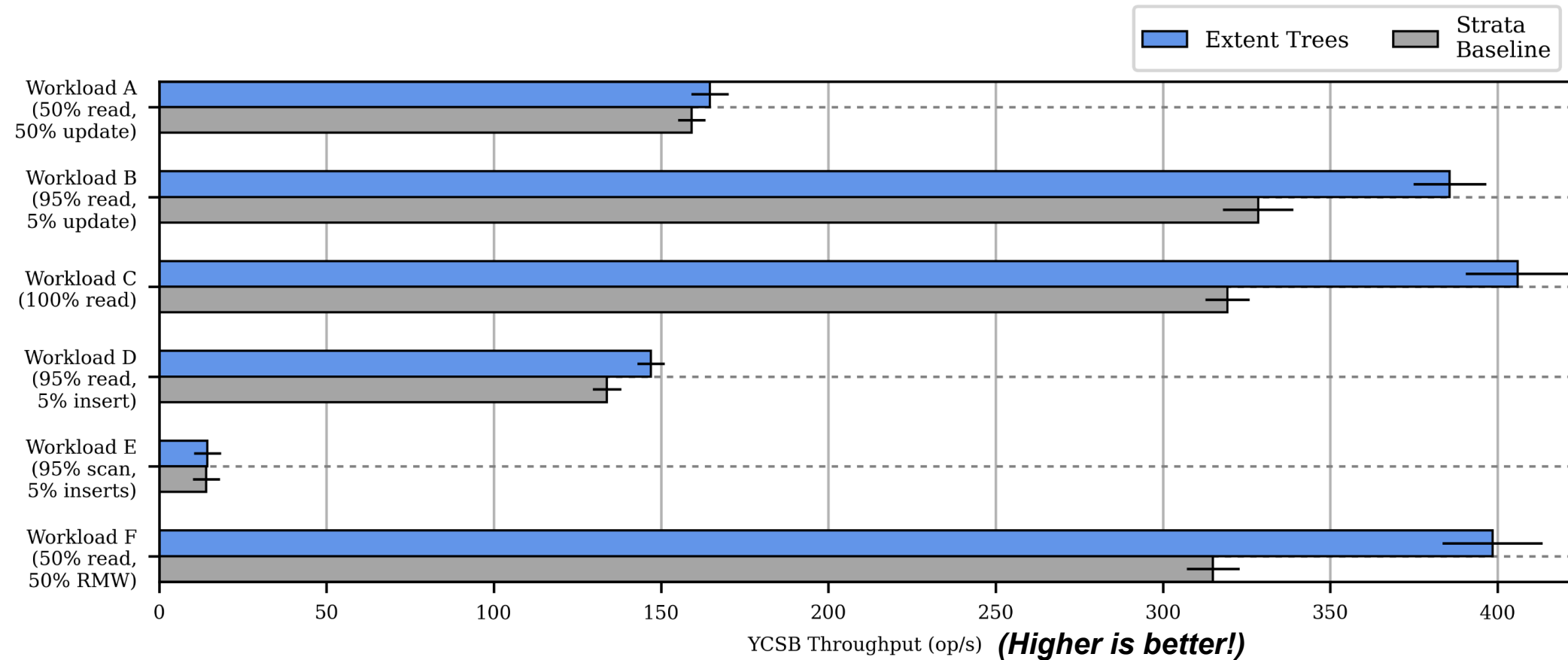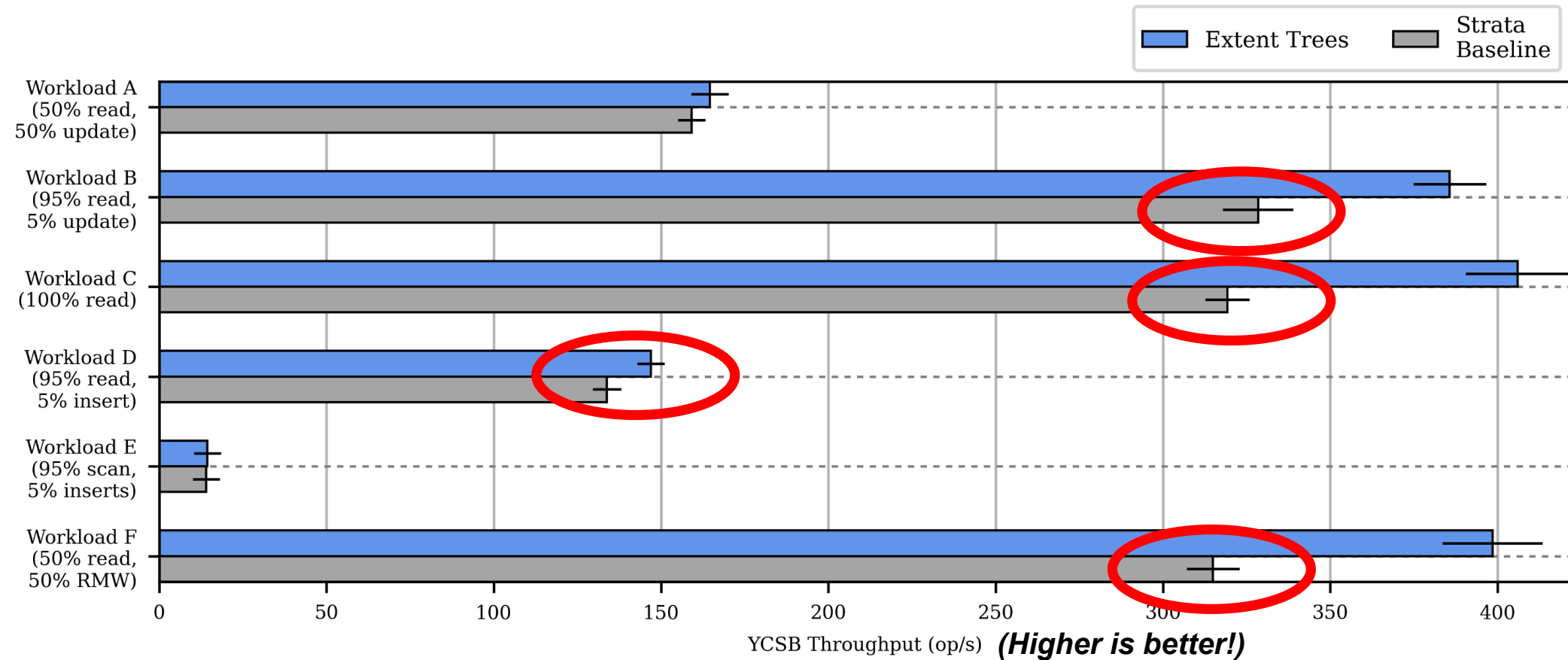  - webproxy (5:1 read/write ratio)

# Is Page Caching Necessary?



write()

file: /home/ian/notes.txt
offset: 0x40000000
data: "Neque porro quisquam est qui dolorem ..."

Application

inode num: 0x101
logical block: 0x40000

File Mapping

device offset: ...

Block allocator

Page Cache (DRAM)

- *File mapping cache*
- *...*

*File System*

# Is Page Caching Necessary?



**write()**

file: /home/ian/notes.txt

offset: 0x40000000

data: "Neque porro quisquam est qui dolorem ..."

Application

inode num: 0x101

logical block: 0x40000

File Mapping

device offset: ...

Block allocator

Page Cache (DRAM)

- *File mapping cache*
- ...

*File System*

8

# Is Page Caching Necessary?



inode num: 0x101

logical block: 0x40000

Using the page cache impacts file mapping structure design!

data: "Neque porro quisquam est qui dolorem ..."

Application

Page Cache (DRAM)

- File mapping cache
- …

File System

8

# Is Page Caching Necessary?

# Is Page Caching Necessary?



YCSB Throughput (op/s) *(Higher is better!)*

Legend: Extent Trees, Strata Baseline

Workload A (50% read, 50% update)
Workload B (95% read, 5% update)
Workload C (100% read)
Workload D (95% read, 5% insert)
Workload E (95% scan, 5% inserts)
Workload F (50% read, 50% RMW)

9

# Is Page Caching Necessary?



**No!** Copying mappings to DRAM is costly

Legend: Extent Trees | Strata Baseline

Y-axis categories:
- Workload A (50% read, 50% update)
- Workload B (95% read, 5% ...)
- Workload ... (10...)
- Workload D (95% read, 5% insert)
- Workload E (95% scan, 5% inserts)
- Workload F (50% read, 50% RMW)

X-axis: YCSB Throughput (op/s) *(Higher is better!)*
X-axis values: 0, 50, 100, 150, 200, 250, 300, 350, 400

9

# Contributions

- Rigorously analyze file mapping in PM

- Optimize legacy PM file mapping structures

- Design new PM-optimized file mapping approaches

- Evaluate end-to-end performance on real workloads

# Analysis Overview

| How is file mapping affected by… | Design Question |
| --- | --- |
| Page caching? | Is page caching necessary? |
| **File size?** | **Specialize for different file sizes?** |
| **IO size?** | **Optimize for sequential access?** |
| **Space utilization?** | **Make file mapping structure elastic?** |
| **Concurrency?** | **Is ensuring isolation important?** |
| **Locality?** | **Optimize for specific workloads?** |
| **Fragmentation?** | **Make robust against file system aging?** |
| **Storage structures?** | **Can we reuse PM storage structures?** |
| Real workloads? | Are mapping optimizations impactful? |

# Analysis Highlights

# Analysis Highlights

- Analyze 4 different file mapping approaches optimized for PM

- **Optimize legacy PM file mapping structures**
  - Extent trees (Strata, ext4-DAX)
  - Radix trees (page cache mapping, NOVA)

# Analysis Highlights

- Analyze 4 different file mapping approaches optimized for PM

- **Optimize legacy PM file mapping structures**
  - Extent trees (Strata, ext4-DAX)
  - Radix trees (page cache mapping, NOVA)

- Legacy structures suffered performance degradation on large files, update operations expensive

# Analysis Highlights

- Analyze 4 different file mapping approaches optimized for PM

- **Optimize legacy PM file mapping structures**
  - Extent trees (Strata, ext4-DAX)
  - Radix trees (page cache mapping, NOVA)

- Legacy structures suffered performance degradation on large files, update operations expensive

# Analysis Highlights

- Analyze 4 different file mapping approaches optimized for PM

- **Optimize legacy PM file mapping structures**
  - Extent trees (Strata, ext4-DAX)
  - Radix trees (page cache mapping, NOVA)

- Legacy structures suffered performance degradation on large files, update operations expensive

- **Design new PM-optimized file mapping approaches**
  - Cuckoo hashing
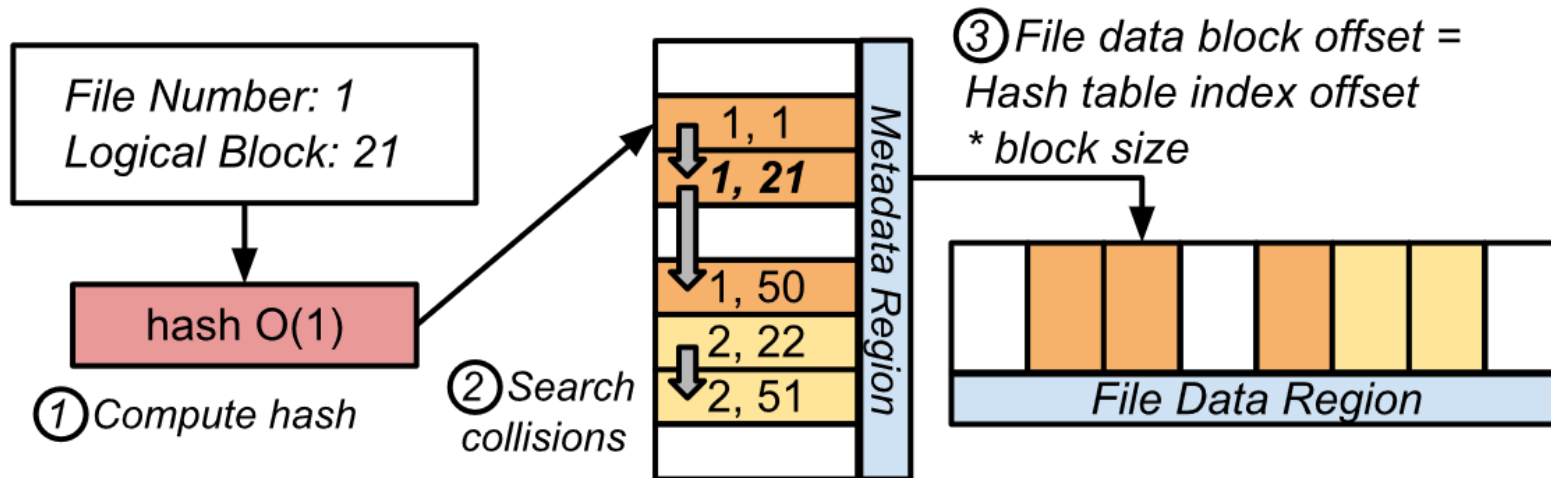  - HashFS

# HashFS

# HashFS

- Hash table structure (linear probing)
  - Possible due to no page cache + PM byte addressability
    - Makes sparse, random updates efficient

# HashFS

- Hash table structure (linear probing)
  - Possible due to no page cache + PM byte addressability
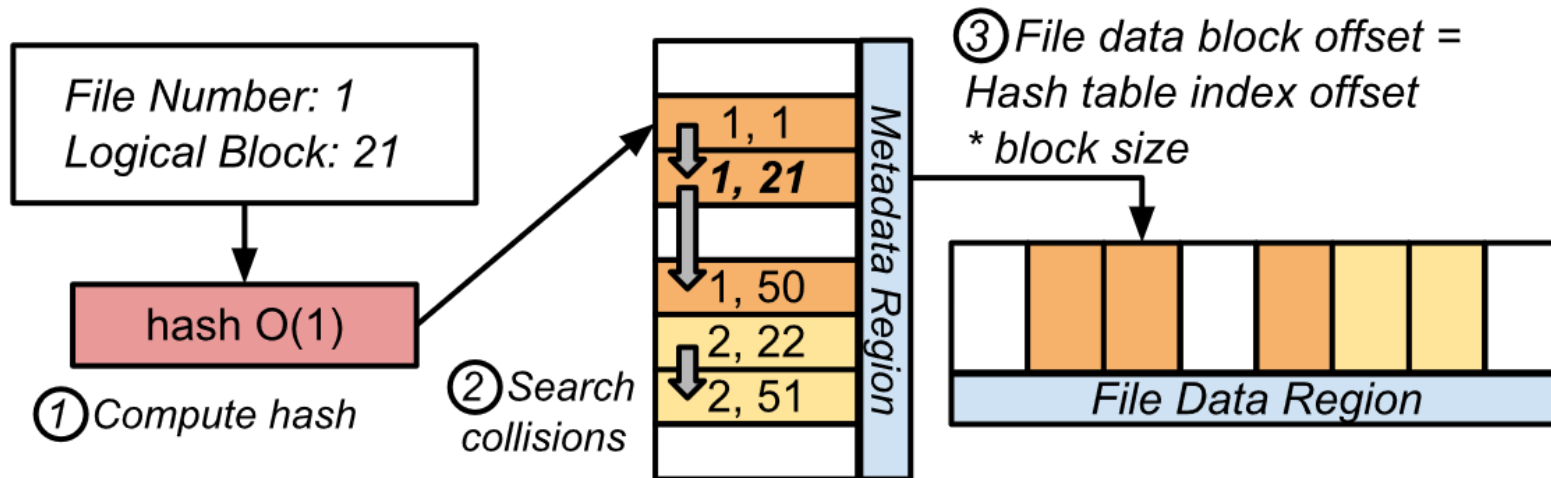    - Makes sparse, random updates efficient

File Number: 1
Logical Block: 21

# HashFS

- Hash table structure (linear probing)
  - Possible due to no page cache + PM byte addressability
    - Makes sparse, random updates efficient
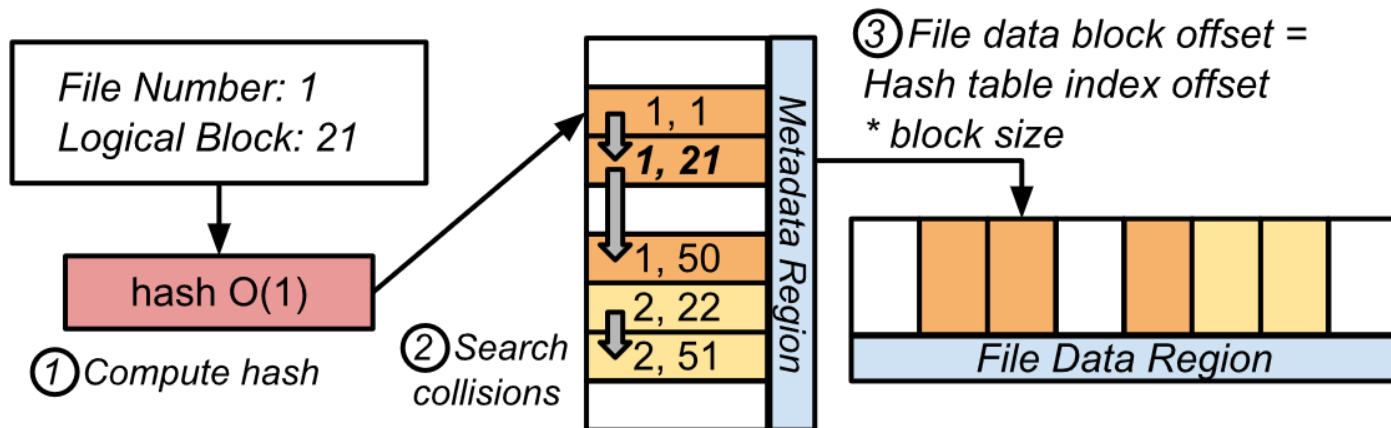
File Number: 1
Logical Block: 21

① Compute hash

# HashFS

- Hash table structure (linear probing)
  - Possible due to no page cache + PM byte addressability
    - Makes sparse, random updates efficient

File Number: 1
Logical Block: 21

hash O(1)

① Compute hash

# HashFS

- Hash table structure (linear probing)
  - Possible due to no page cache + PM byte addressability
    - Makes sparse, random updates efficient

# HashFS

- Hash table structure (linear probing)
  - Possible due to no page cache + PM byte addressability
    - Makes sparse, random updates efficient

# HashFS

- Hash table structure (linear probing)
  - Possible due to no page cache + PM byte addressability
    - Makes sparse, random updates efficient

# HashFS

- Hash table structure (linear probing)
  - Possible due to no page cache + PM byte addressability
    - Makes sparse, random updates efficient
- Combined block-allocation and file-mapping scheme
  - Insert into hash table implicitly allocates block at corresponding offset
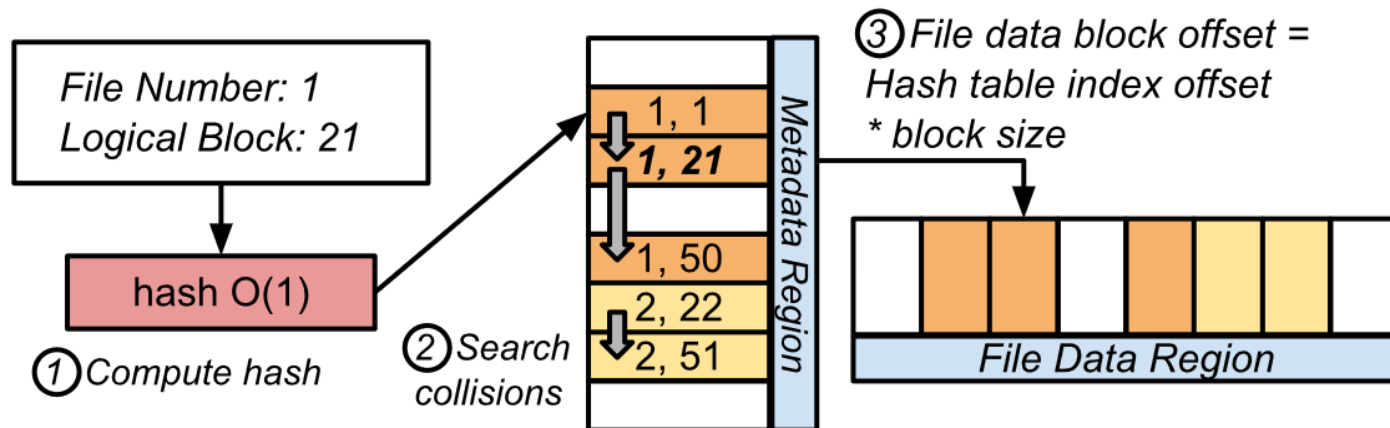  - **Bypasses expensive block allocator management** (cf. our paper)

# HashFS (cont.)



File Number: 1
Logical Block: 21

hash O(1)

① Compute hash

② Search collisions

1, 1
**1, 21**
1, 50
2, 22
2, 51

Metadata Region

③ File data block offset =
Hash table index offset
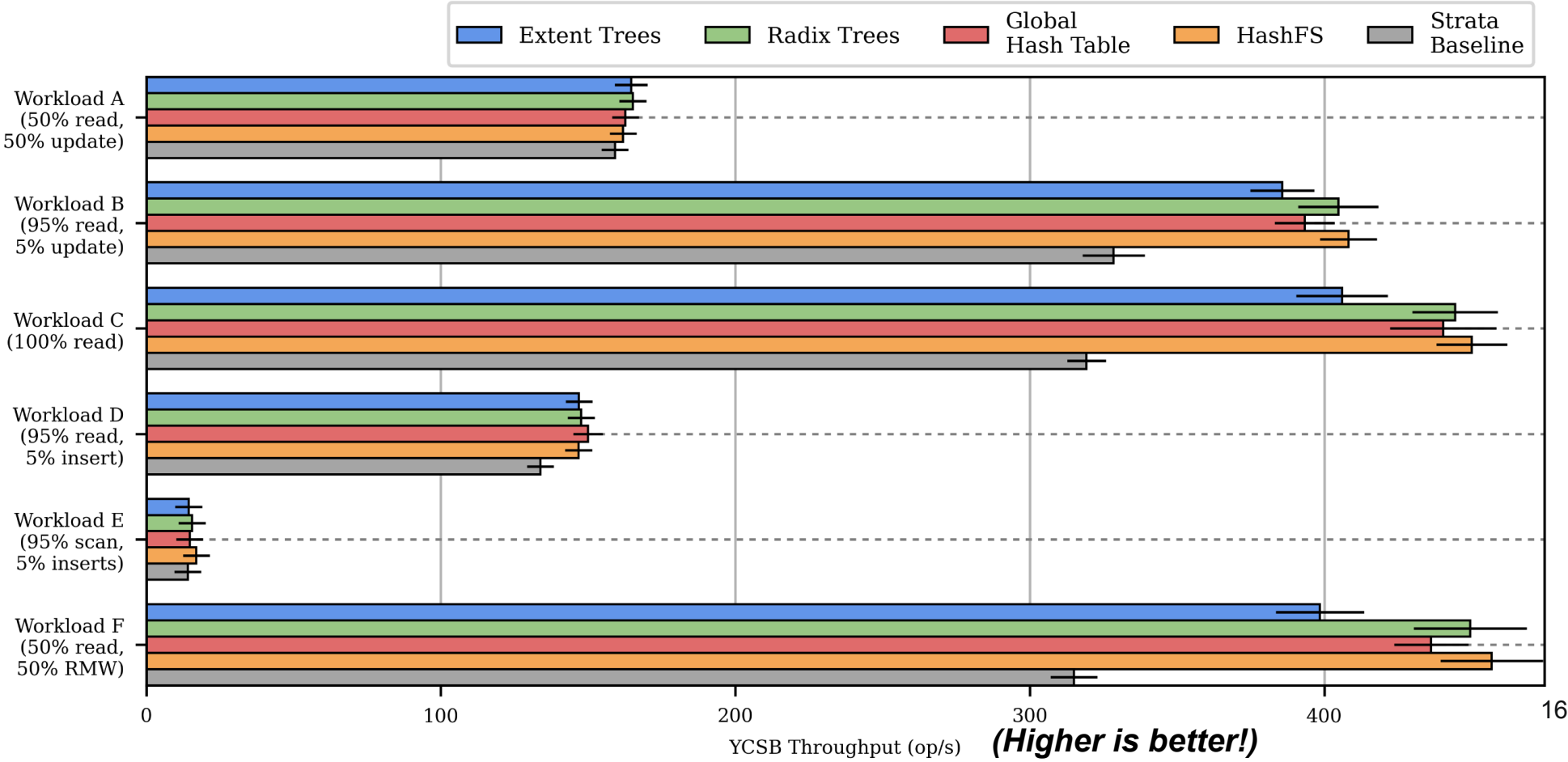* block size

File Data Region

14

# HashFS (cont.)

- Need to avoid all resizing (incurs high update latency)
  - Must also be a global structure (one structure for all files)
  - Statically allocate max size at creation



① Compute hash

② Search collisions

③ File data block offset = Hash table index offset * block size

File Number: 1
Logical Block: 21

hash O(1)

1, 1
1, 21
1, 50
2, 22
2, 51

Metadata Region

File Data Region

# HashFS (cont.)

- Need to avoid all resizing (incurs high update latency)
  - Must also be a global structure (one structure for all files)
  - Statically allocate max size at creation
- FS optimization: use SIMD for large IO operations
  - Many file system workloads perform large IO operations
  - For efficiency, mapping structures must return ranges of mappings
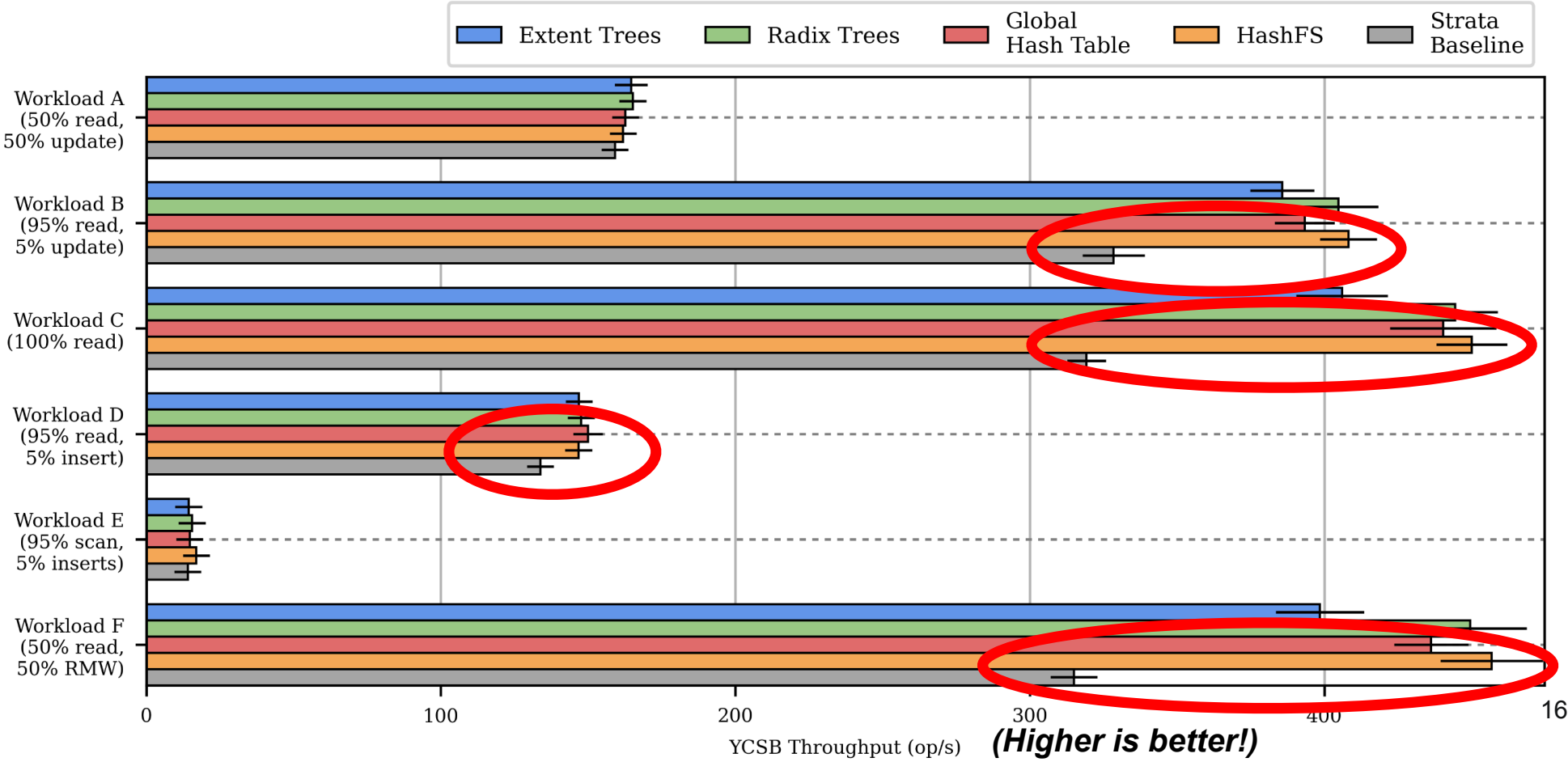  - **Perform hash table operations in parallel**

# Contributions

- Rigorously analyze file mapping in PM

- Optimize legacy PM file mapping structures

- Design new PM-optimized file mapping approaches

- Evaluate end-to-end performance on real workloads
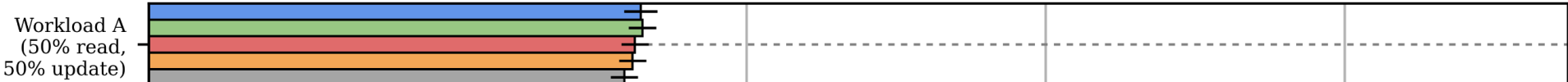
# Are File Mapping Optimizations Impactful?



YCSB Throughput (op/s) *(Higher is better!)*

# Are File Mapping Optimizations Impactful?



YCSB Throughput (op/s)  *(Higher is better!)*

# Are File Mapping Optimizations Impactful?


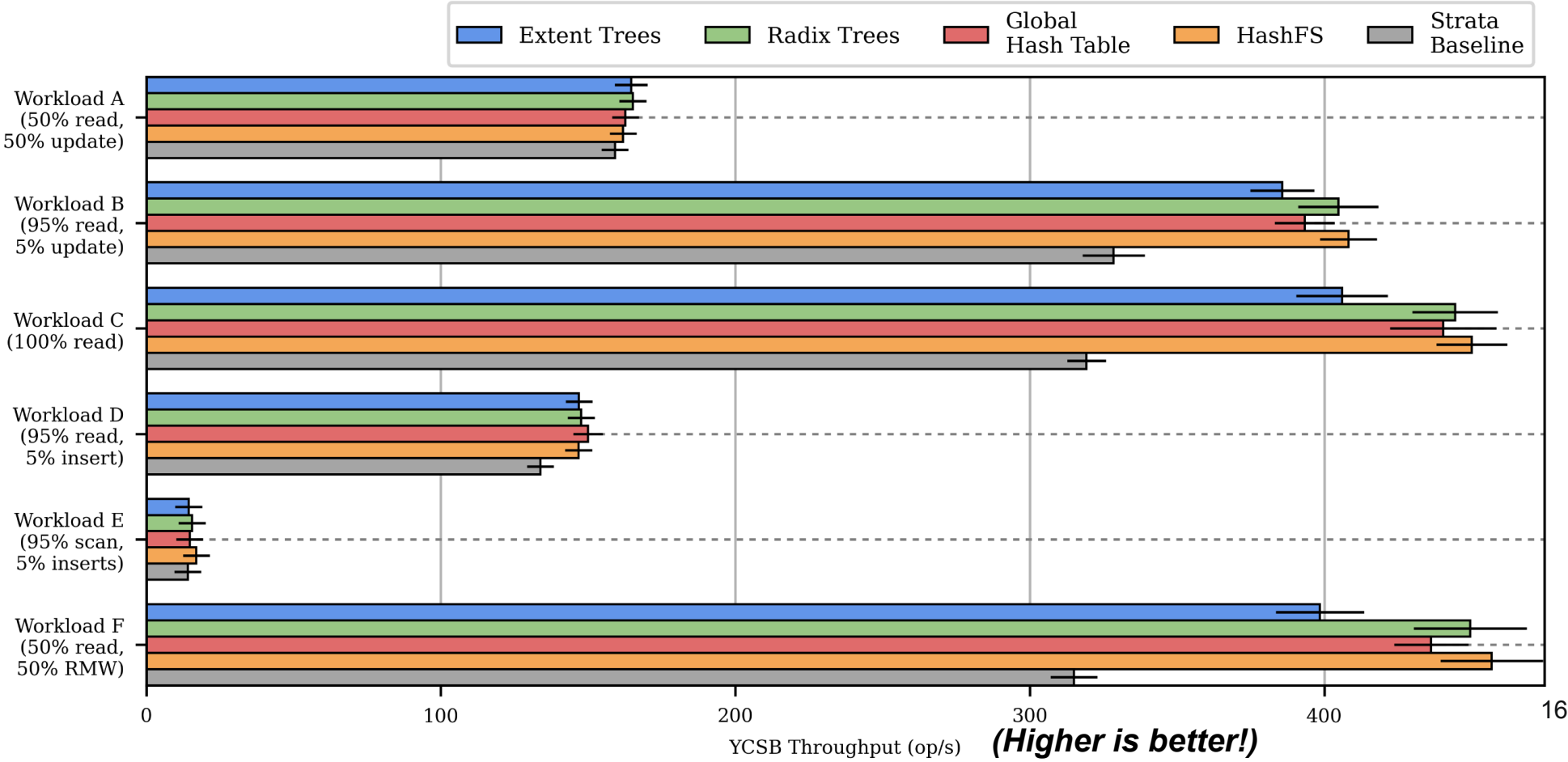
HashFS provides +10–45% throughput!

YCSB Throughput (op/s)  *(Higher is better!)*

# Are File Mapping Optimizations Impactful?



Legend: Extent Trees, Radix Trees, Global Hash Table, HashFS, Strata Baseline

- Workload A (50% read, 50% update)
- Workload B (95% read, 5% update)
- Workload C (100% read)
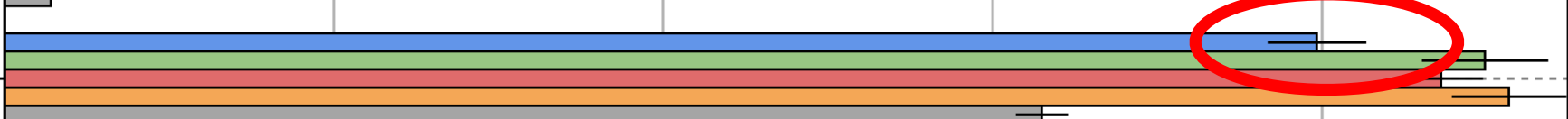- Workload D (95% read, 5% insert)
- Workload E (95% scan, 5% inserts)
- Workload F (50% read, 50% RMW)

X-axis: YCSB Throughput (op/s) *(Higher is better!)* — 0, 100, 200, 300, 400

16

# Are File Mapping Optimizations Impactful?



**YCSB Throughput (op/s)** *(Higher is better!)*

Legend: Extent Trees, Radix Trees, Global Hash Table, HashFS, Strata Baseline

Workload A (50% read, 50% update)
Workload B (95% read, 5% update)
Workload C (100% read)
Workload D (95% read, 5% insert)
Workload E (95% scan, 5% inserts)
Workload F (50% read, 50% RMW)

16

# Are File Mapping Optimizations Impactful?



PM-optimized extent trees have 13% lower throughput!

# Are File Mapping Optimizations Impactful?

# Are File Mapping Optimizations Impactful?



Filebench Throughput (mb/s) *(Higher is better!)*

HashFS provides +26% throughput!

# Conclusions

# Conclusions

- We must re-examine file mapping specifically for PM file systems!

# Conclusions

- We must re-examine file mapping specifically for PM file systems!

- A rigorous analysis yields insights into performance-shortfalls of existing mapping approaches

# Conclusions

- We must re-examine file mapping specifically for PM file systems!

- A rigorous analysis yields insights into performance-shortfalls of existing mapping approaches

- We design two new, global file mapping approaches (cuckoo hashing, HashFS)

# Conclusions

- We must re-examine file mapping specifically for PM file systems!

- A rigorous analysis yields insights into performance-shortfalls of existing mapping approaches

- We design two new, global file mapping approaches (cuckoo hashing, HashFS)

- HashFS (our new PM-optimized file mapping approach) outperforms the state-of-the-art by up to 45% in real workloads

# Thank you!

Corresponding Author: **Ian Neal**

iangneal@umich.edu

https://about.iangneal.io