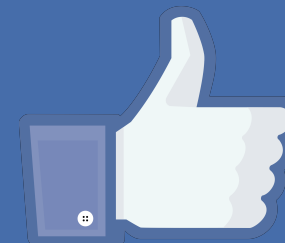facebook

# Analysis of HDFS Under HBase
## A Facebook Messages Case Study

Tyler Harter, Dhruba Borthakur*, Siying Dong*, Amitanand Aiyer*,
Liyin Tang*, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau

University of Wisconsin-Madison          *Facebook Inc.

# Why Study Facebook Messages?

Represents an important type of application.  Universal backend for:

- Cellphone texts
- Chats
- Emails

# Why Study Facebook Messages?

Represents an important type of application.  Universal backend for:

- Cellphone texts

- Chats

- Emails

# Why Study Facebook Messages?

Represents an important type of application.  Universal backend for:

- Cellphone texts
- Chats
- Emails

# Why Study Facebook Messages?

Represents an important type of
application.  Universal backend for:

- Cellphone texts

- Chats

- Emails

# Why Study Facebook Messages?

Represents an important type of application.  Universal backend for:

- Cellphone texts
- Chats
- Emails

## Represents HBase over HDFS

- Common backend at Facebook and other companies
- Similar stack used at Google (BigTable over GFS)

# Why Study Facebook Messages?

Represents an important type of application.  Universal backend for:

- Cellphone texts

- Chats

- Emails

## Represents HBase over HDFS

- Common backend at Facebook and other companies

- Similar stack used at Google (BigTable over GFS)

## Represents layered storage

# Building a Distributed Application (Messages)

We have many machines with many disks.
*How should we use them to store messages?*

Machine 1        Machine 2        Machine 3

# Building a Distributed Application (Messages)

One option: use machines and disks directly.

Messages

Machine 1          Machine 2          Machine 3

# Building a Distributed Application (Messages)

One option: use machines and disks directly.
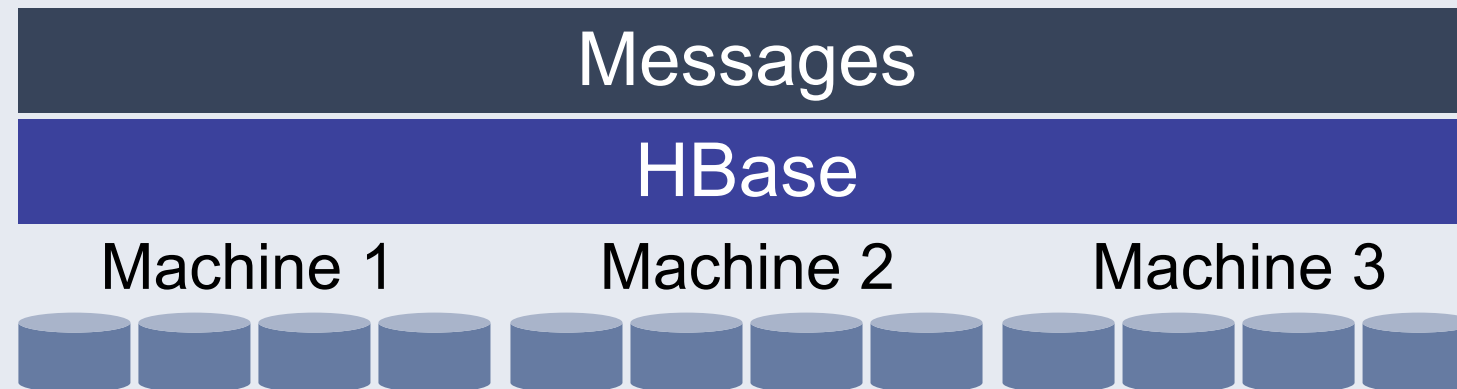Very specialized, but very high development cost.

| Messages |
|---|
| Machine 1 | Machine 2 | Machine 3 |

# Building a Distributed Application (Messages)

**Messages**

Machine 1          Machine 2          Machine 3

# Building a Distributed Application (Messages)

Use HBase for K/V logic

| Messages |
| HBase |

Machine 1          Machine 2          Machine 3

# Building a Distributed Application (Messages)

Use HBase for K/V logic
Use HDFS for replication

| Messages |
| --- |
| HBase |
| Hadoop File System |

| Worker | Worker | Worker |
| --- | --- | --- |
| Machine 1 | Machine 2 | Machine 3 |

# Building a Distributed Application (Messages)

Use HBase for K/V logic
Use HDFS for replication
Use Local FS for allocation

| Messages |
| --- |
| HBase |
| Hadoop File System |

| Worker | Worker | Worker |
| --- | --- | --- |
| Machine 1 | Machine 2 | Machine 3 |

| FS | FS | FS | FS | FS | FS | FS | FS | FS | FS | FS | FS |

# Layered Storage Discussion

## Layering Advantages

- Simplicity (thus fewer software bugs)

- Lower development costs

- Code sharing between systems

## Layering Questions

- Is layering free performance-wise?

- Can layer integration be useful?

- Should there be multiple HW layers?

# Outline

## Intro

- Messages stack overview

- Methodology: trace-driven analysis and simulation

- HBase background

## Results

- Workload analysis

- Hardware simulation: adding a flash layer

- Software simulation: integrating layers

## Conclusions

# Methodology

Actual stack

**Messages**

**HBase**

**HDFS**

**Local FS**

# Methodology

**Actual stack**

| Messages |
|----------|
| HBase |
| HDFS |
| Local FS |

**HDFS Traces** →

## Hadoop Trace FS (HTFS)

- Collects request details
  - Reads/writes, offsets, lengths

- 9 shadow machines

- 8.3 days

# Methodology

Actual stack

Messages

HBase

⟶ **HDFS Traces**

HDFS

Local FS

MapReduce Analysis Pipeline

Workload Analysis

# Methodology

Actual stack

**Messages**

**HBase**

**HDFS**

**Local FS**

**HDFS Traces**

Simulated stack

**HBase+HDFS**

what -ifs

**Local Traces (inferred)**

MapReduce Analysis Pipeline

Workload Analysis

Local Storage

what -ifs

Simulation Results

# Methodology

Actual stack

| |
|---|
| Messages |
| HBase |
| HDFS |
| Local FS |

**HDFS Traces**

# Methodology

Actual stack

**Messages**

**HBase**

**HDFS** → **HDFS Traces**

**Local FS**

**Background: how does HBase use HDFS?**

# Outline

# HBase's HDFS Files

Four activities do HDFS I/O:

HBase memory:

| MemTable |
| --- |

HDFS files:

| LOG |
| --- |

# HBase's HDFS Files

Four activities do HDFS I/O:
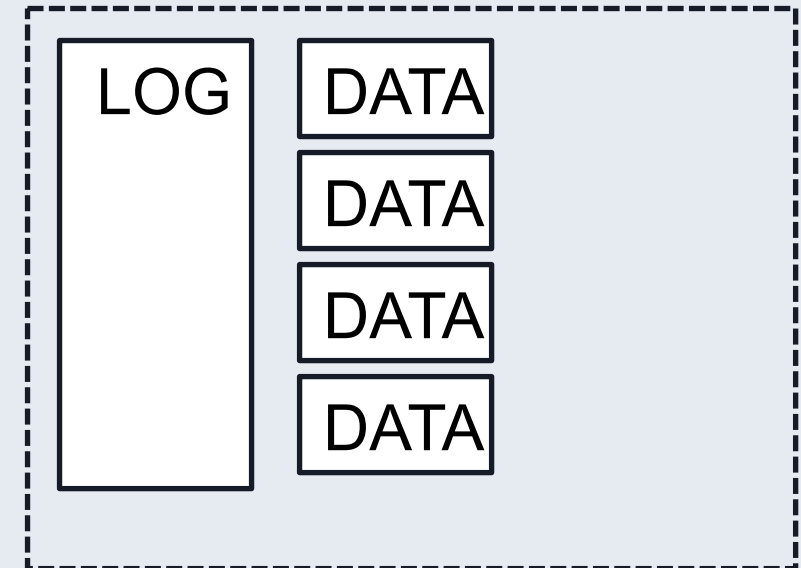
- Logging

HBase memory:

MemTable

HDFS files:

LOG

# HBase's HDFS Files

Four activities do HDFS I/O:

- Logging

HBase memory:

MemTable

HDFS files:

LOG

# HBase's HDFS Files

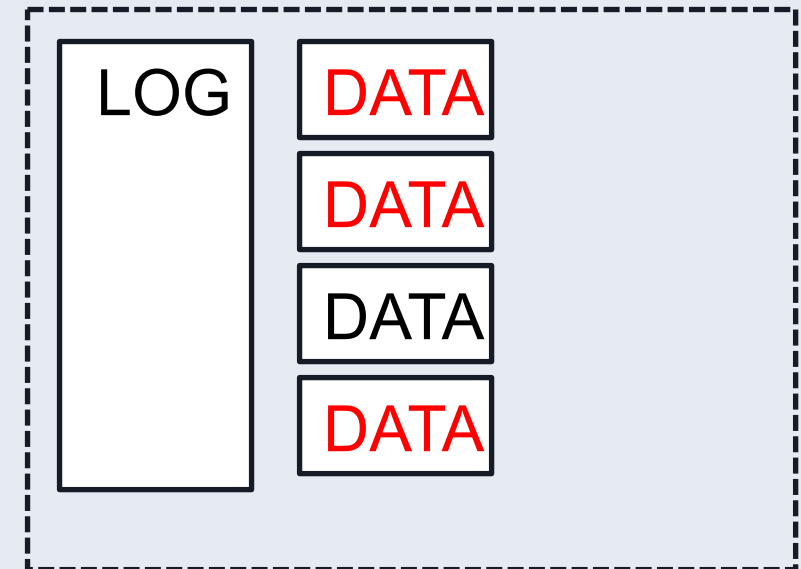Four activities do HDFS I/O :

- Logging

- Flushing

Flush MemTable to sorted file

HBase memory:

MemTable

HDFS files:

LOG    DATA

# HBase's HDFS Files

Four activities do HDFS I/O :

- Logging
- Flushing

HBase memory:

MemTable

HDFS files:

LOG    DATA

# HBase's HDFS Files

Four activities do HDFS I/O :

- Logging

- Flushing

HBase memory:

MemTable

HDFS files:

LOG

DATA

DATA

DATA

DATA

# HBase's HDFS Files

Four activities do HDFS I/O:

- Logging

- Flushing

HBase memory:

MemTable

HDFS files:

LOG

DATA

DATA

DATA

DATA

# HBase's HDFS Files

Four activities do HDFS I/O:

- Logging

- Flushing

- Foreground reads

HBase memory:

MemTable

HDFS files:

LOG

DATA

DATA

DATA

DATA

# HBase's HDFS Files

Four activities do HDFS I/O:

- Logging
- Flushing
- Foreground reads
- Compaction

# HBase's HDFS Files

Four activities do HDFS I/O:

- Logging
- Flushing
- Foreground reads
- Compaction

compaction merge sorts the files

HBase memory:

MemTable

HDFS files:

LOG

DATA

# HBase's HDFS Files

Four activities do HDFS I/O:

- Logging
- Flushing
- Foreground reads
- Compaction

Baseline I/O:

- Flushing and foreground reads are always required

# HBase's HDFS Files

Four activities do HDFS I/O:

- <span style="color:red">Logging</span>
- Flushing
- Foreground reads
- <span style="color:red">Compaction</span>

Baseline I/O:

- <u>Flushing</u> and <u>foreground reads</u> are always required

HBase overheads:

- <u>Logging</u>: useful for crash recovery (but not normal operation)
- <u>Compaction</u>: improves performance (but not required for correctness)

# Outline

Intro

- Messages stack overview

- Methodology: trace-driven analysis and simulation

- HBase background

Results

- Workload analysis

- Hardware simulation: adding a flash layer

- Software simulation: integrating layers

Conclusions

# Workload Analysis Questions

At each layer, what activities read or write?

How large is the dataset?
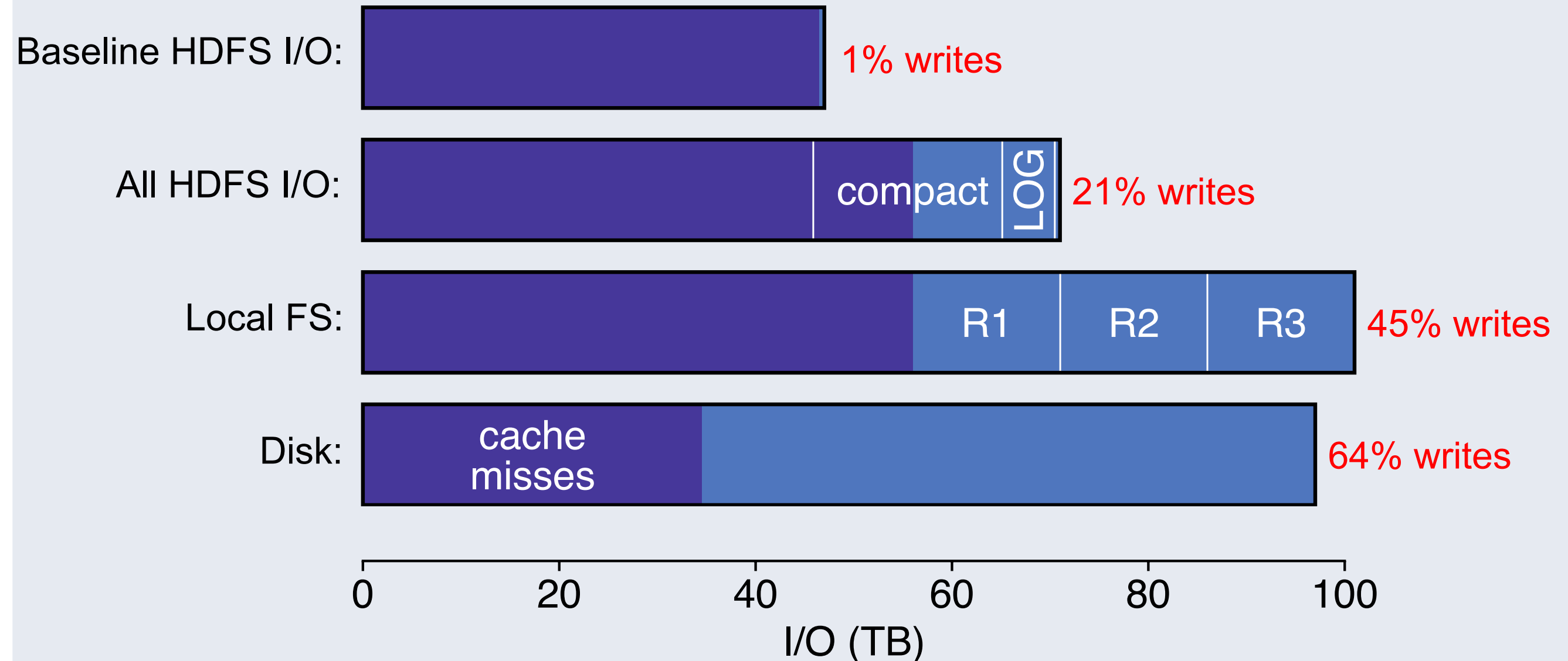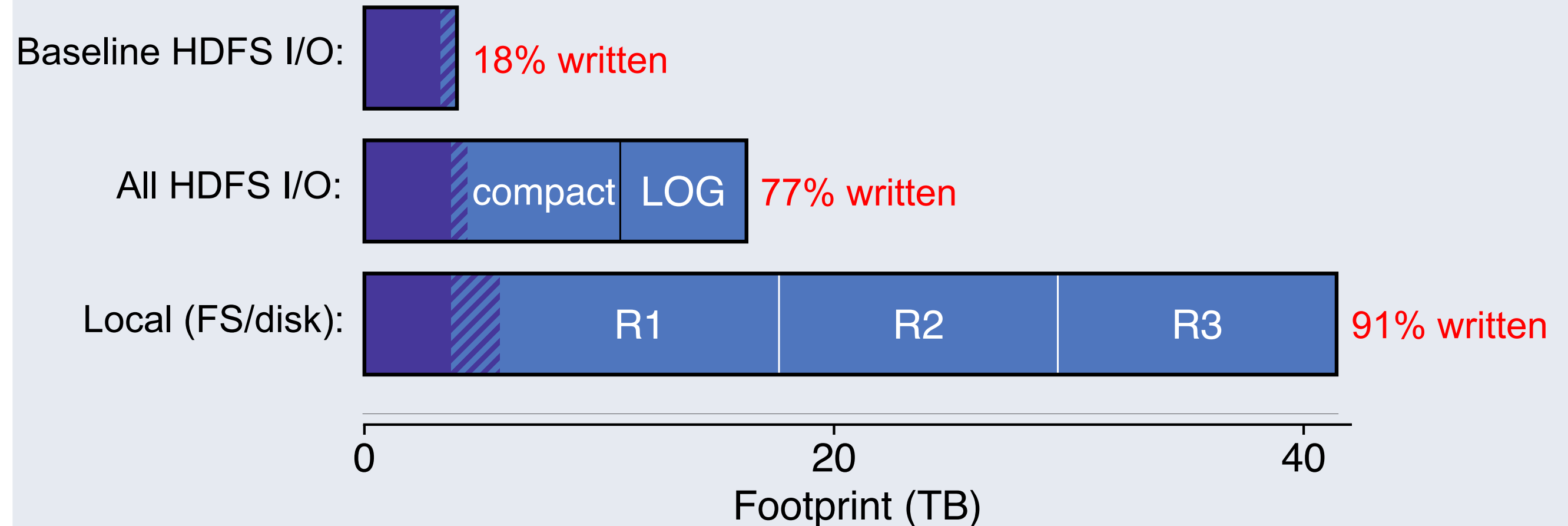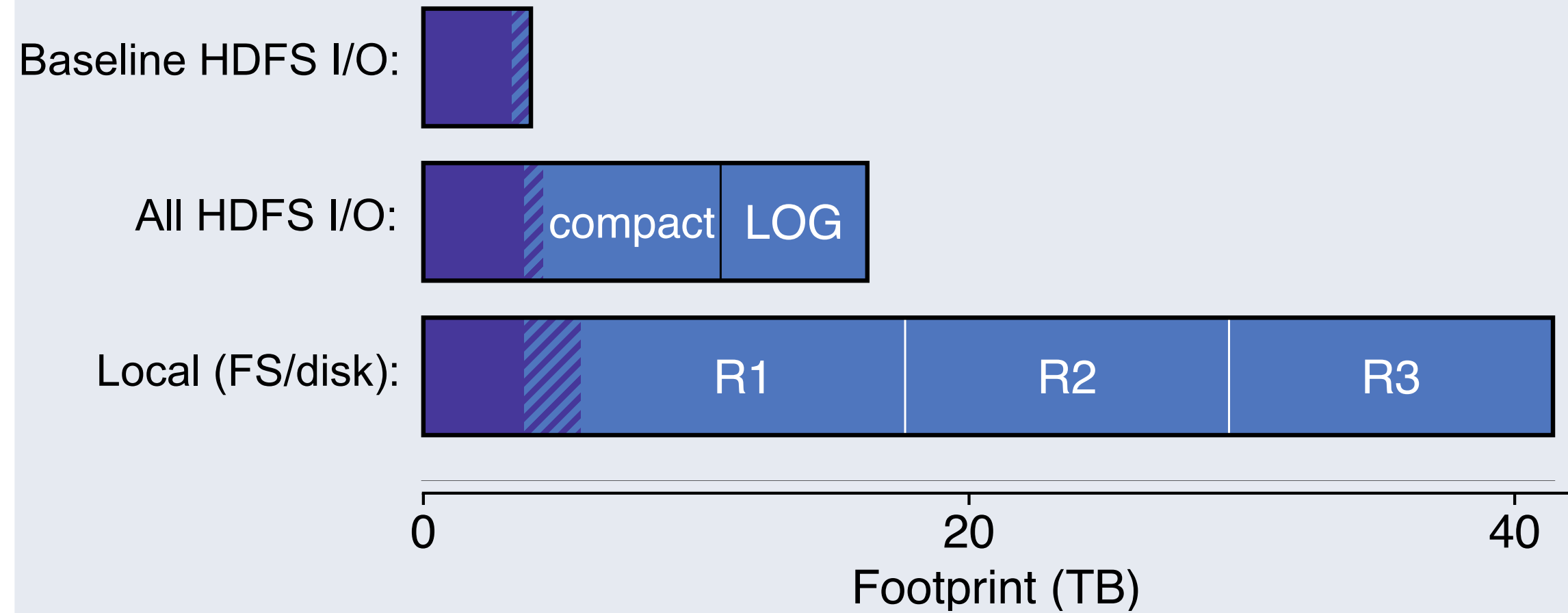
How large are created files?

How sequential is I/O?

# Workload Analysis Questions

At each layer, what activities read or write?

How large is the dataset?

How large are created files?
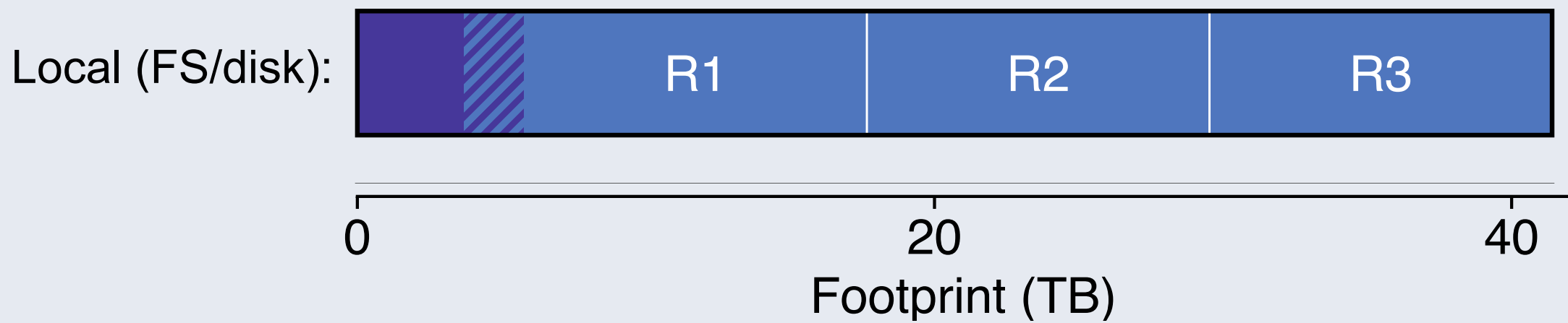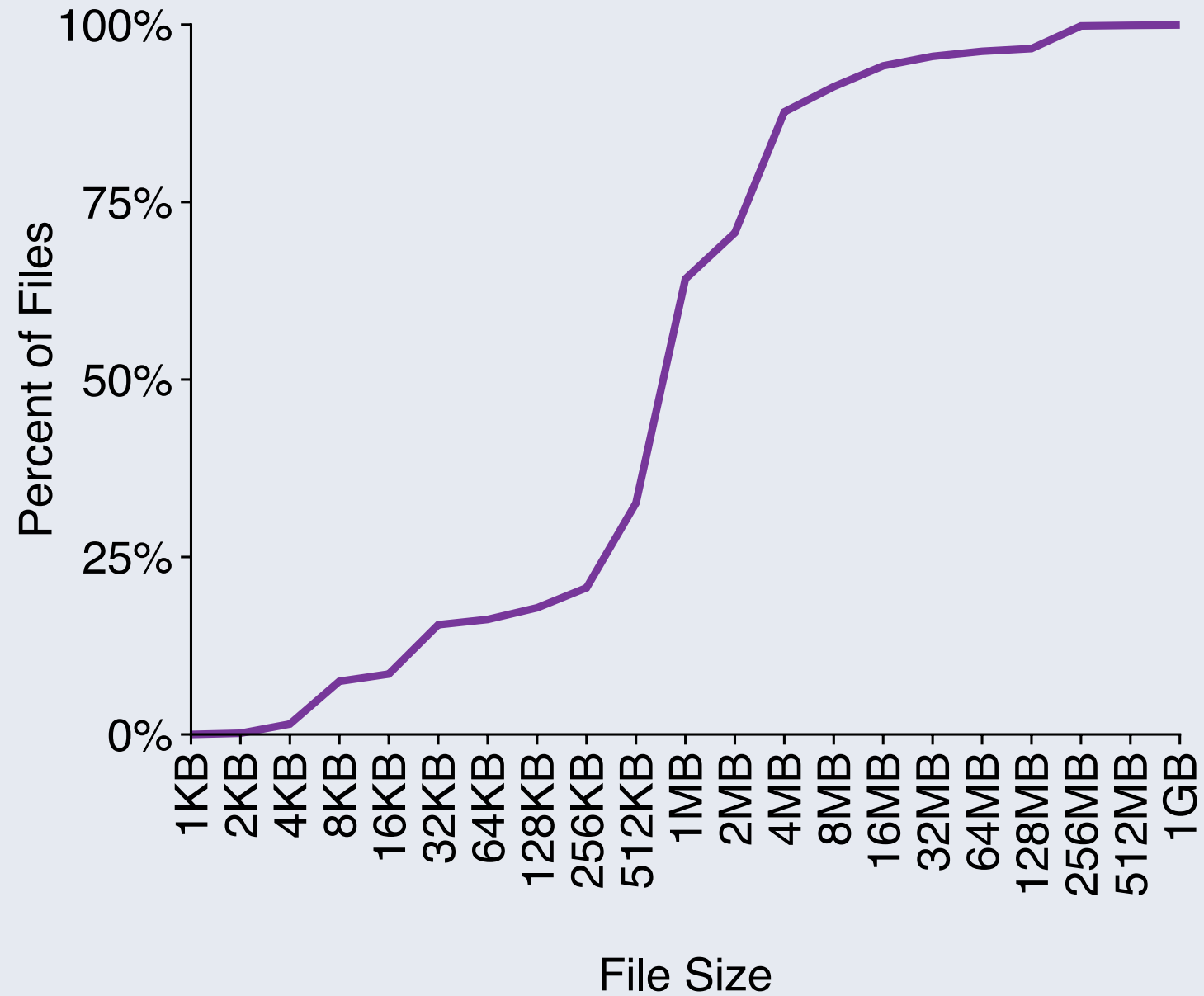
How sequential is I/O?

# Cross-layer R/W Ratios

Baseline HDFS I/O:

reads    writes

1% writes

I/O (TB)

# Cross-layer R/W Ratios



Baseline HDFS I/O: — 1% writes

All HDFS I/O: compact LOG — 21% writes

Local FS: R1 R2 R3 — 45% writes

Disk: cache misses — 64% writes

I/O (TB)

0   20   40   60   80   100

# Workload Analysis Conclusions

① Layers amplify writes: <span style="color:red">1% => 64%</span>

  ◆ Logging, compaction, and replication <u>increase writes</u>

  ◆ Caching <u>decreases reads</u>

# Workload Analysis Questions

At each layer, what activities read or write?

How large is the dataset?

How large are created files?

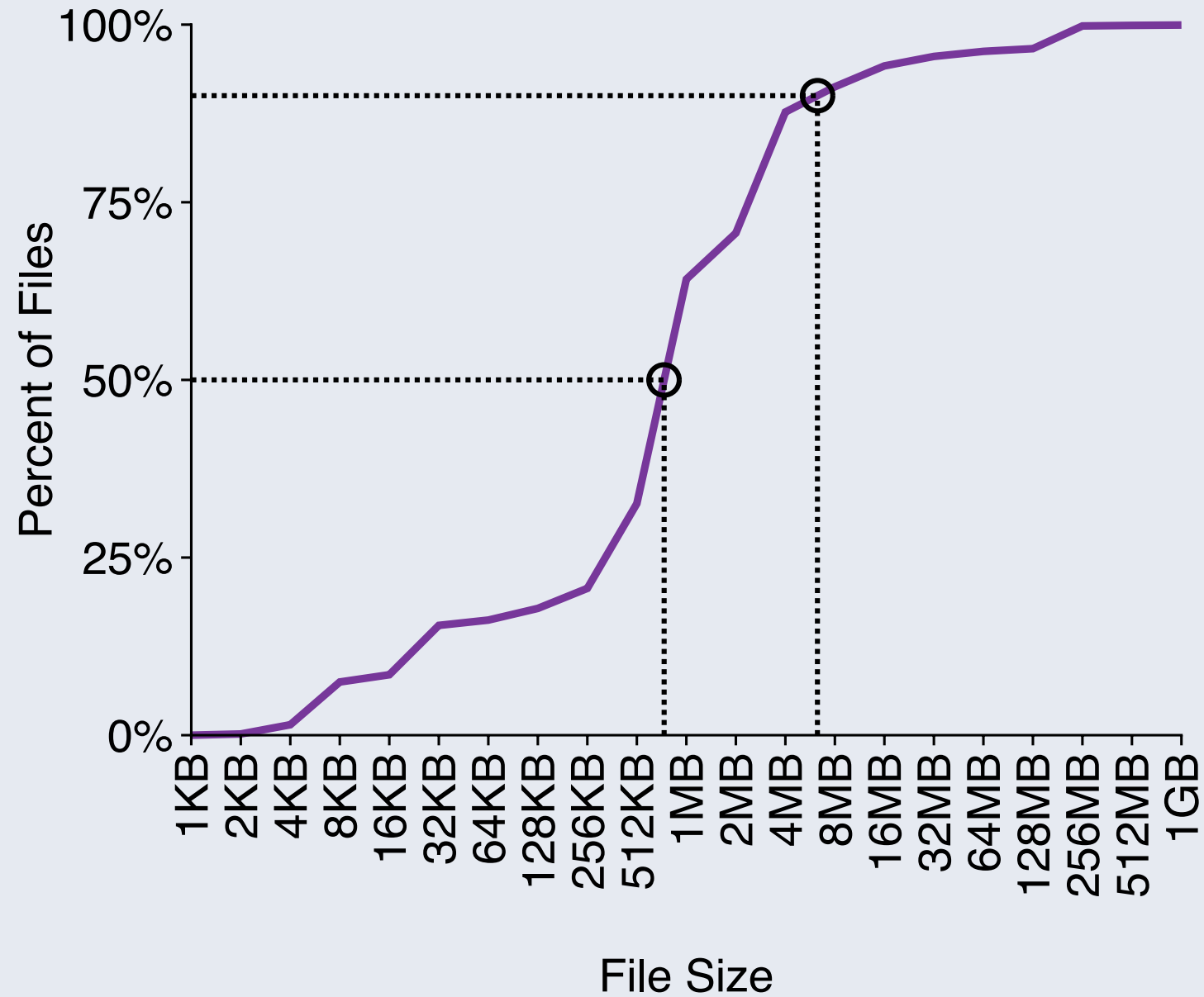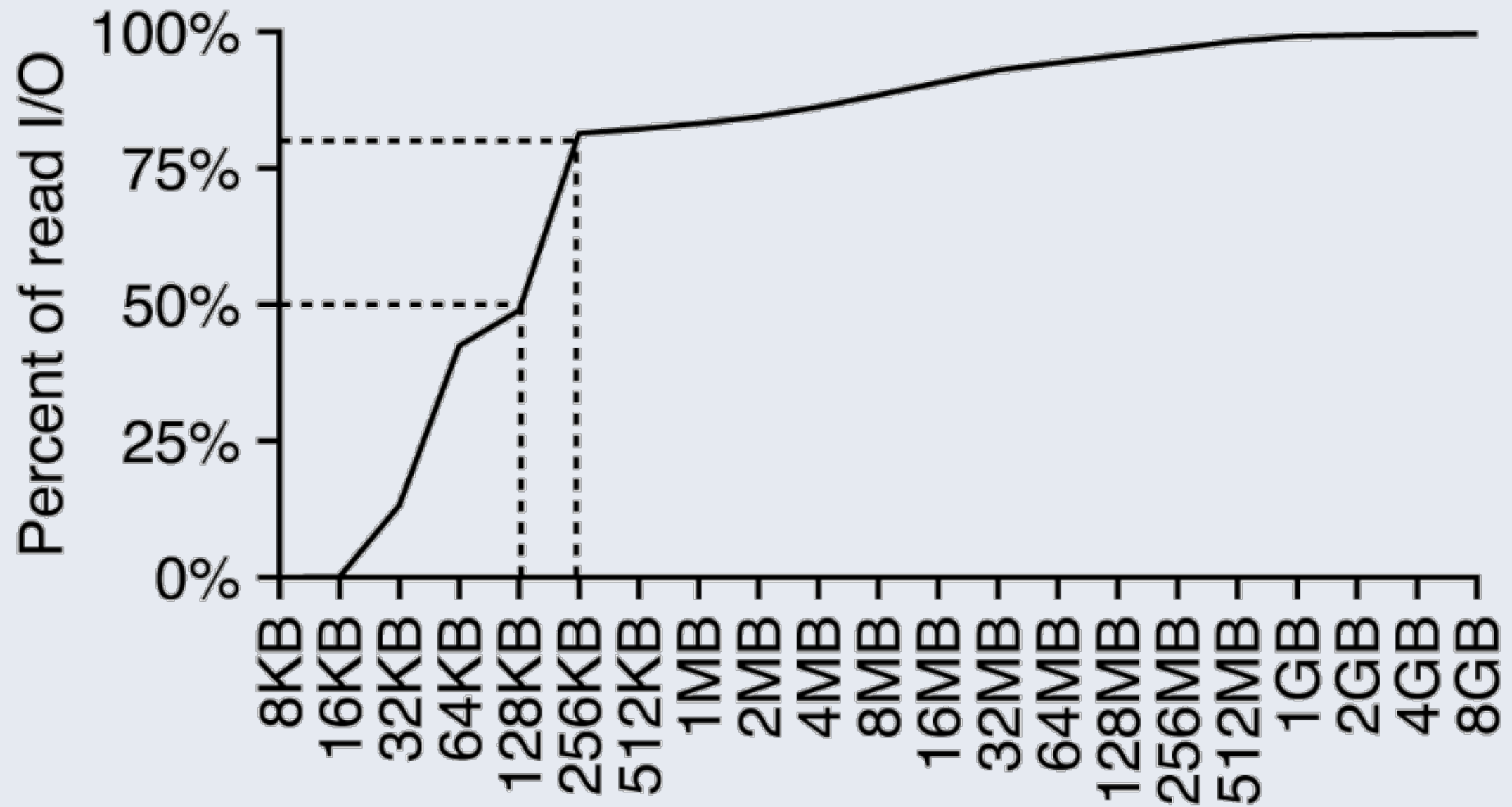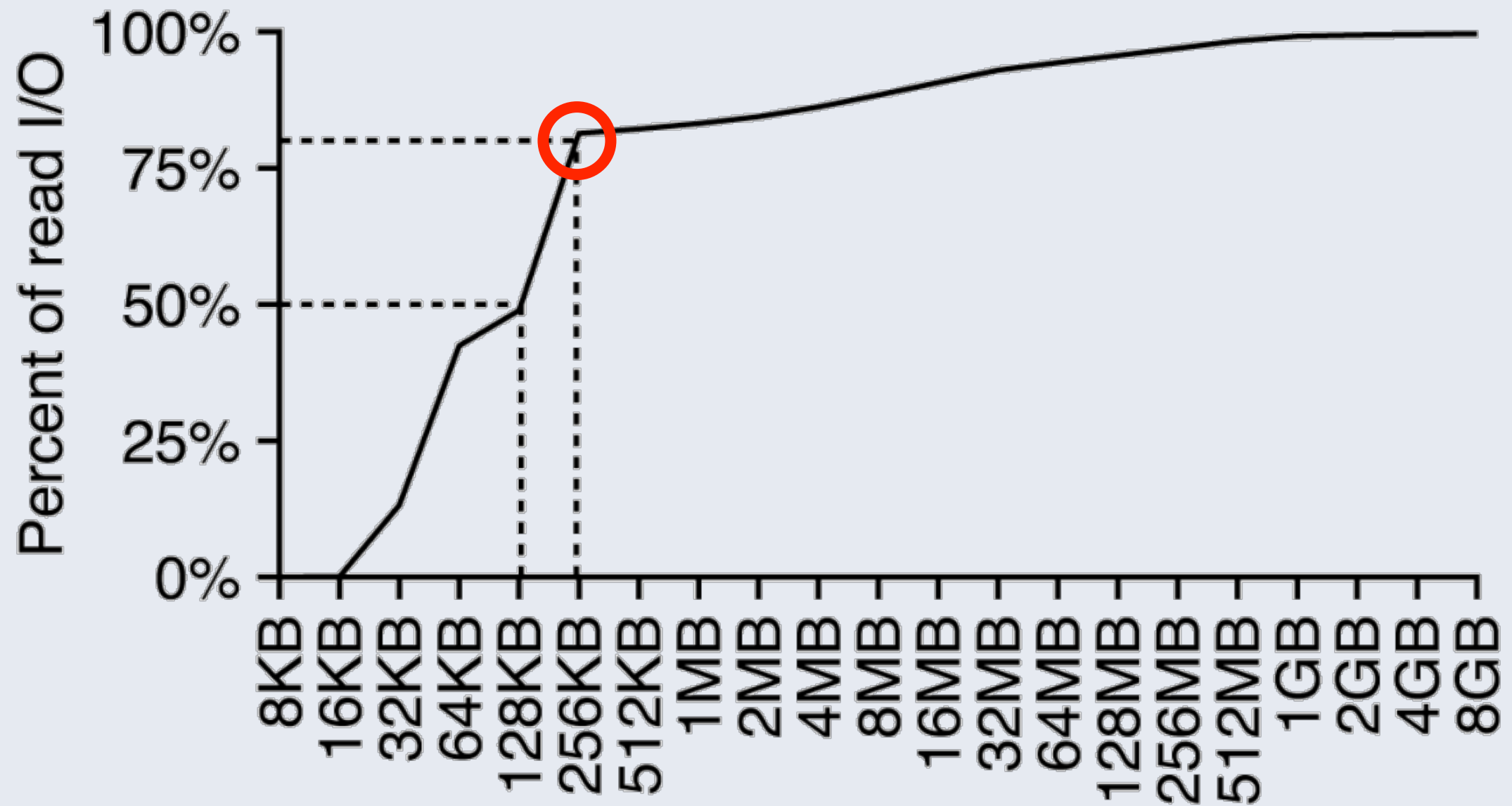How sequential is I/O?

# Cross-layer Dataset (Accessed Data)

Baseline HDFS I/O: — 18% written

All HDFS I/O: compact LOG — 77% written

Local (FS/disk): R1 R2 R3 — 91% written

Footprint (TB)

0    20    40

# Workload Analysis Conclusions

① Layers amplify writes: <span style="color:red">1% => 64%</span>

② Most touched data is <span style="color:red">only written</span>

Cold Data

Baseline HDFS I/O:

All HDFS I/O: compact LOG

Local (FS/disk): R1 R2 R3

Footprint (TB)
0          20          40

# Cold Data

# Cold Data

# Workload Analysis Conclusions

① Layers amplify writes: 1% => 64%

② Most touched data is only written

③ The dataset is large and cold: 2/3 of 120TB never touched

# Workload Analysis Questions

At each layer, what activities read or write?

How large is the dataset?

How large are created files?

How sequential is I/O?

# Created Files: Size Distribution

# Created Files: Size Distribution



50% of files are <750KB

# Created Files: Size Distribution



Percent of Files

File Size

90% of files are <6.3MB

# Workload Analysis Conclusions

① Layers amplify writes: 1% => 64%

② Most touched data is only written

③ The dataset is large and cold: 2/3 of 120TB never touched

④ Files are very small: 90% smaller than 6.3MB

# Workload Analysis Questions

At each layer, what activities read or write?

How large is the dataset?

How large are created files?

How sequential is I/O?

Reads: Run Size

# Reads: Run Size



50% of runs (weighted by I/O) <130KB

Reads: Run Size

80% of runs (weighted by I/O) <250KB

# Workload Analysis Conclusions

①  Layers amplify writes: 1% => 64%

②  Data is read or written, but rarely both

③  The dataset is large and cold: 2/3 of 120TB never touched

④  Files are very small: 90% smaller than 6.3MB

⑤  Fairly random I/O: 130KB median read run

# Outline

# Hardware Architecture: Workload Implications

Option 1: pure disk      Option 2: pure flash      Option 3: hybrid
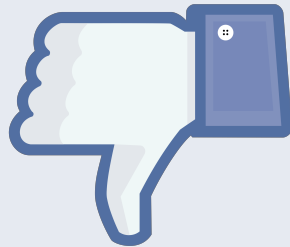
# Hardware Architecture: Workload Implications

Option 1: <span style="color:red">pure disk</span>

- Very random reads
- Small files

Option 2: <span style="color:red">pure flash</span>

Option 3: <span style="color:red">hybrid</span>

# Hardware Architecture: Workload Implications

Option 1: <span style="color:red">pure disk</span>

- Very random reads
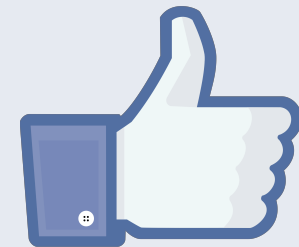
- Small files

Option 2: <span style="color:red">pure flash</span>

Option 3: <span style="color:red">hybrid</span>

# Hardware Architecture: Workload Implications

Option 1: pure disk
- Very random reads
- Small files

Option 2: pure flash
- Large dataset
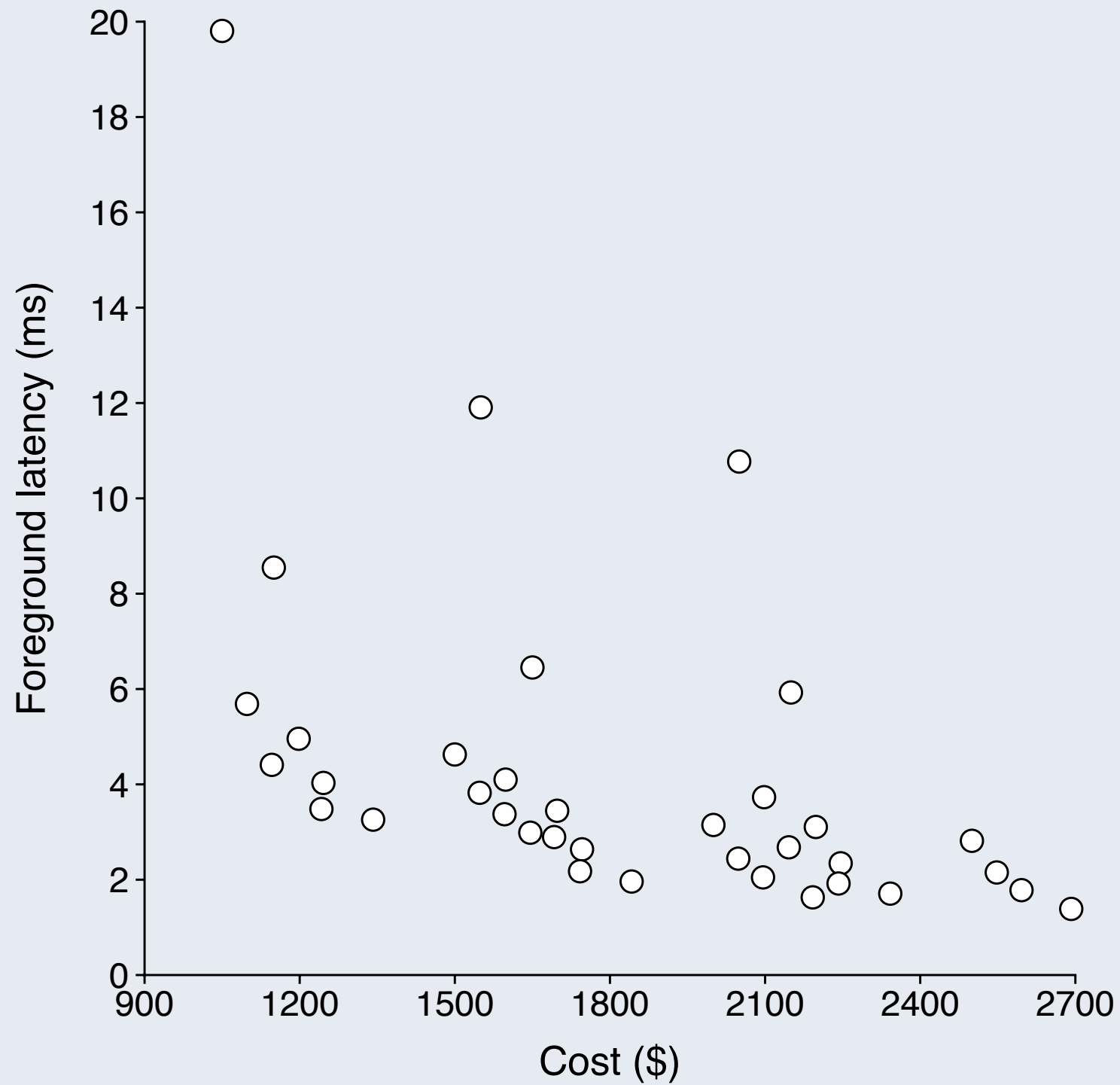- Mostly very cold
- >$10K / machine

Option 3: hybrid

# Hardware Architecture: Workload Implications

Option 1: pure disk

- Very random reads
- Small files

Option 2: pure flash

- Large dataset
- Mostly very cold
- >$10K / machine

Option 3: hybrid

# Hardware Architecture: Workload Implications

## Option 1: <span style="color:red">pure disk</span>

- Very random reads
- Small files

## Option 2: <span style="color:red">pure flash</span>

- Large dataset
- Mostly very cold
- <span style="color:red">>$10K / machine</span>

## Option 3: <span style="color:red">hybrid</span>

- Process of elimination

# Hardware Architecture: Simulation Results

Evaluate <u>cost</u> and <u>performance</u> of 36 hardware combinations (3x3x4)

- Disks: 10, 15, or 20
- RAM (cache): 10, 30, or 100GB
- Flash (cache): 0, 60, 120, or 240GB

Assumptions:

| Hardware | Cost | Performance |
|---|---|---|
| HDD | $100/disk | 10ms seek, 100MB/s |
| RAM | $5/GB | zero latency |
| Flash | $0.8/GB | 0.5ms |

Cost/performance tradeoff for 36 hardware combinations
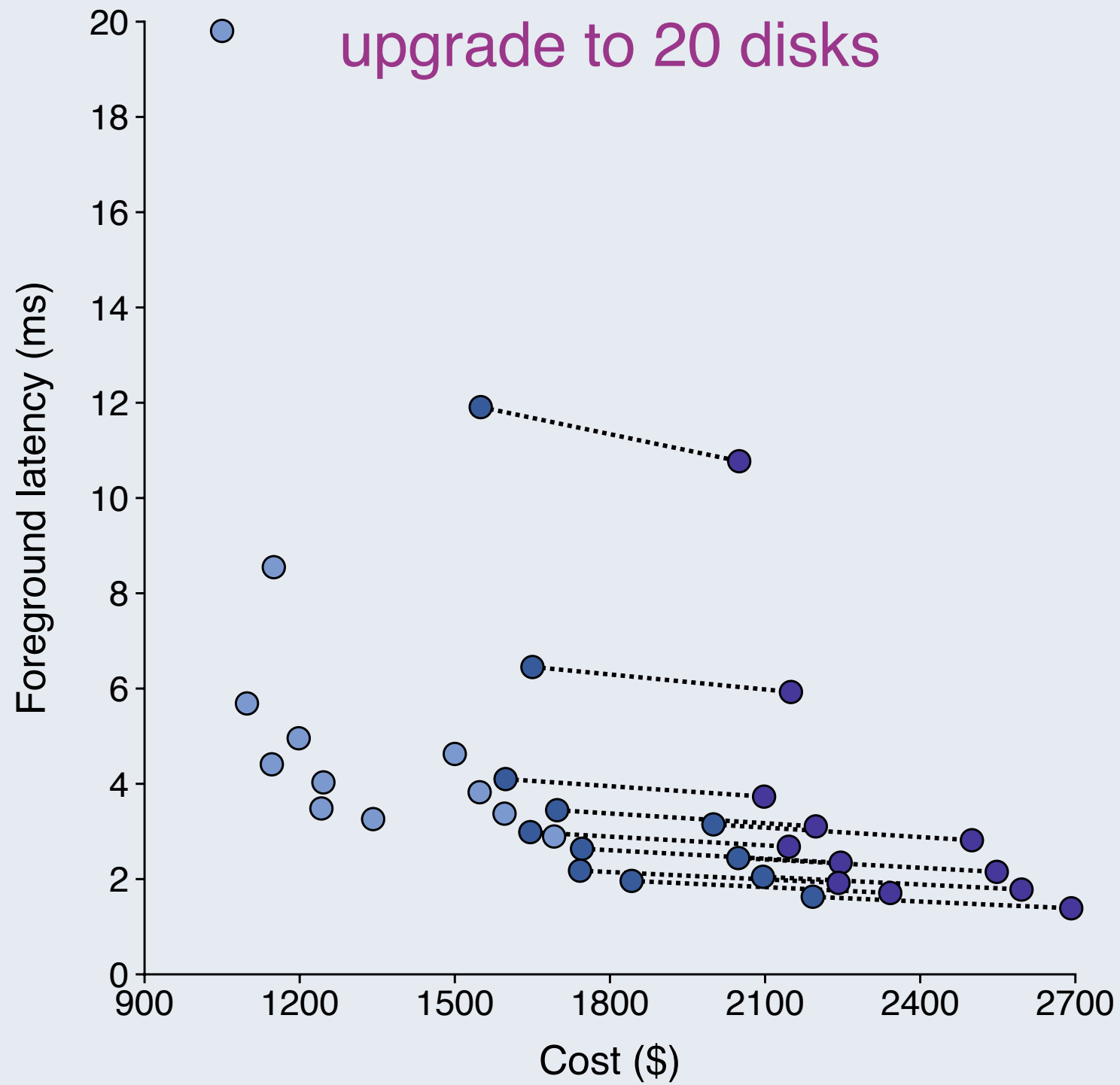
Upgrades decrease latency but increase cost

Good upgrade
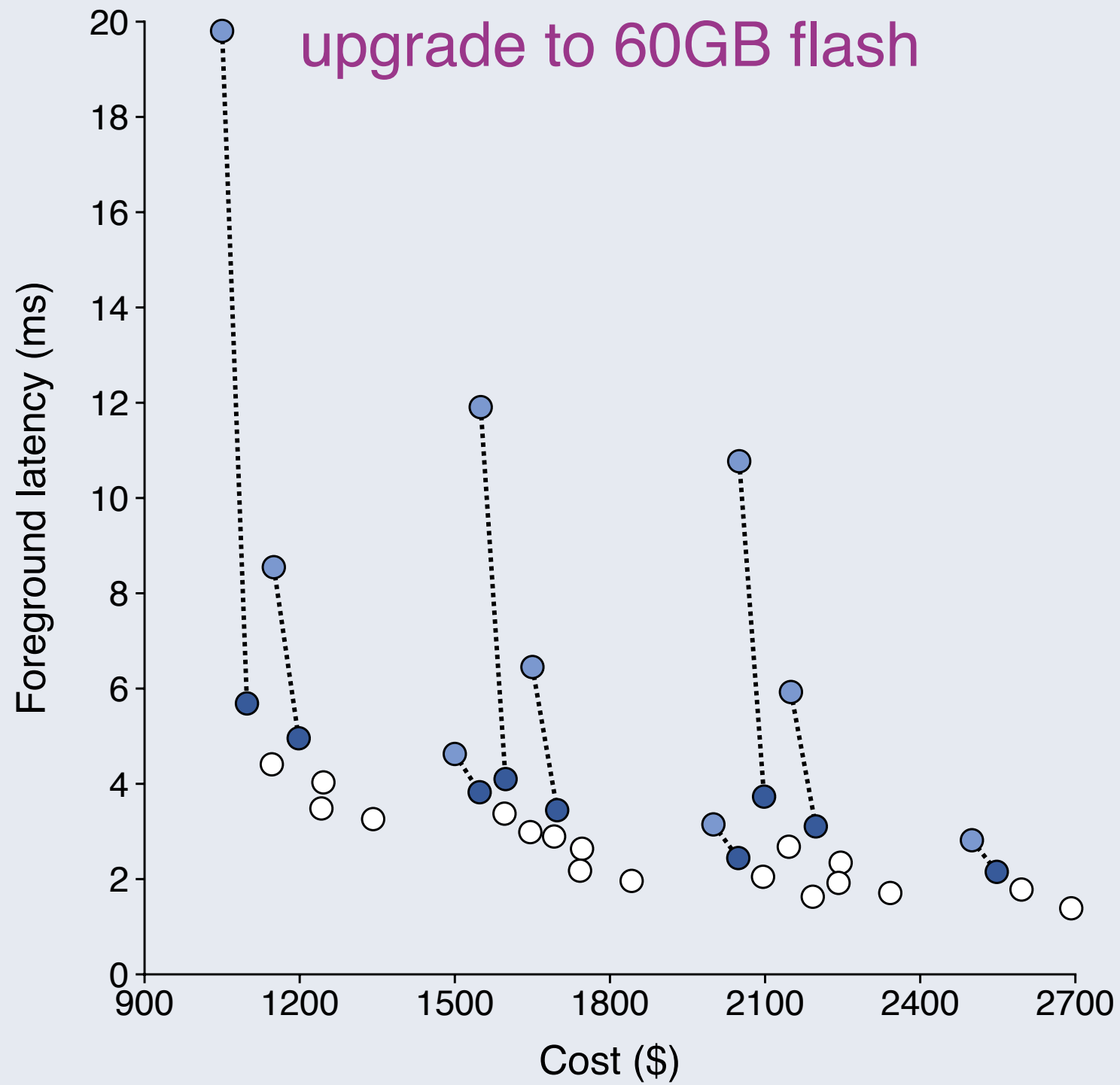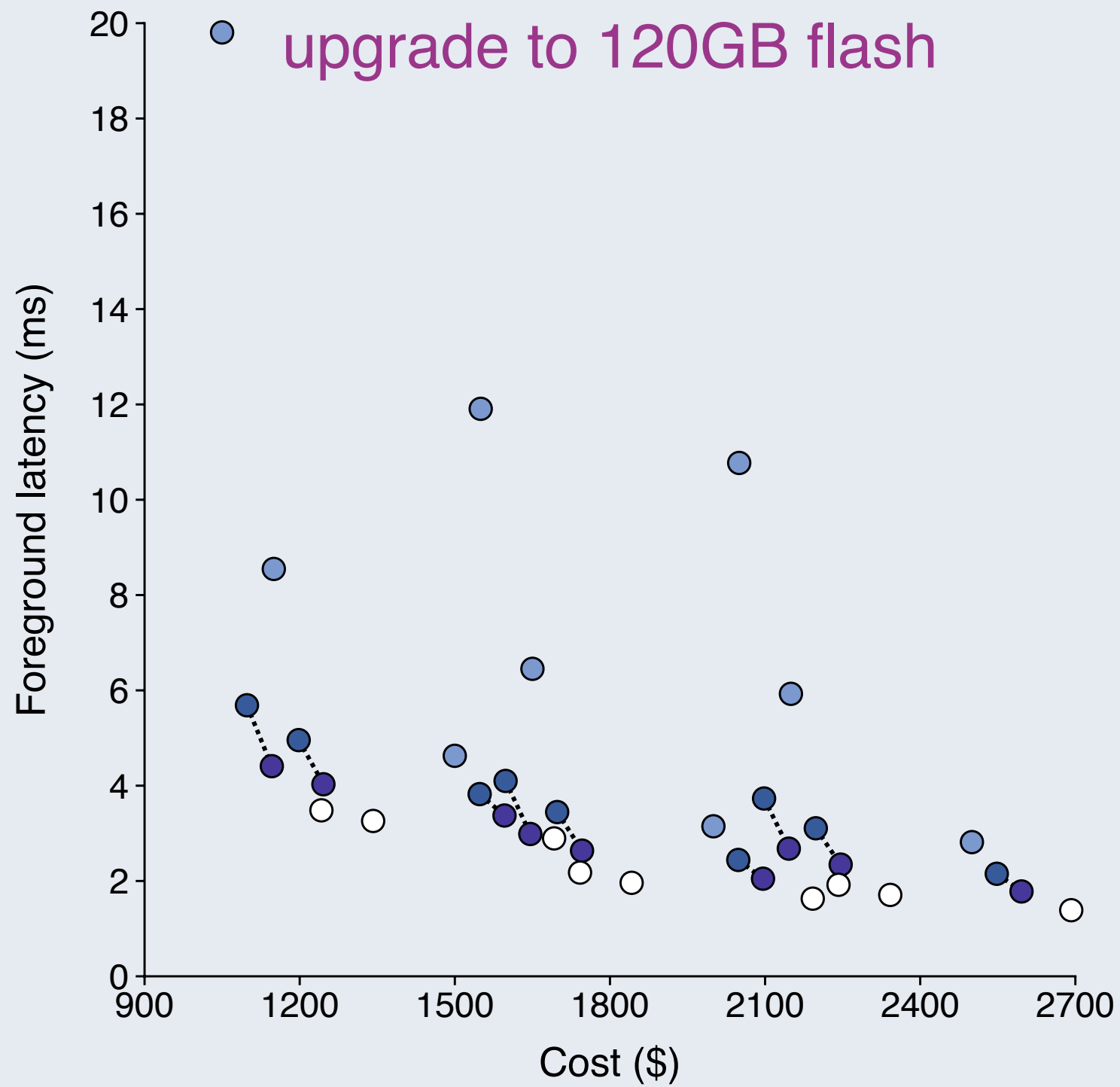
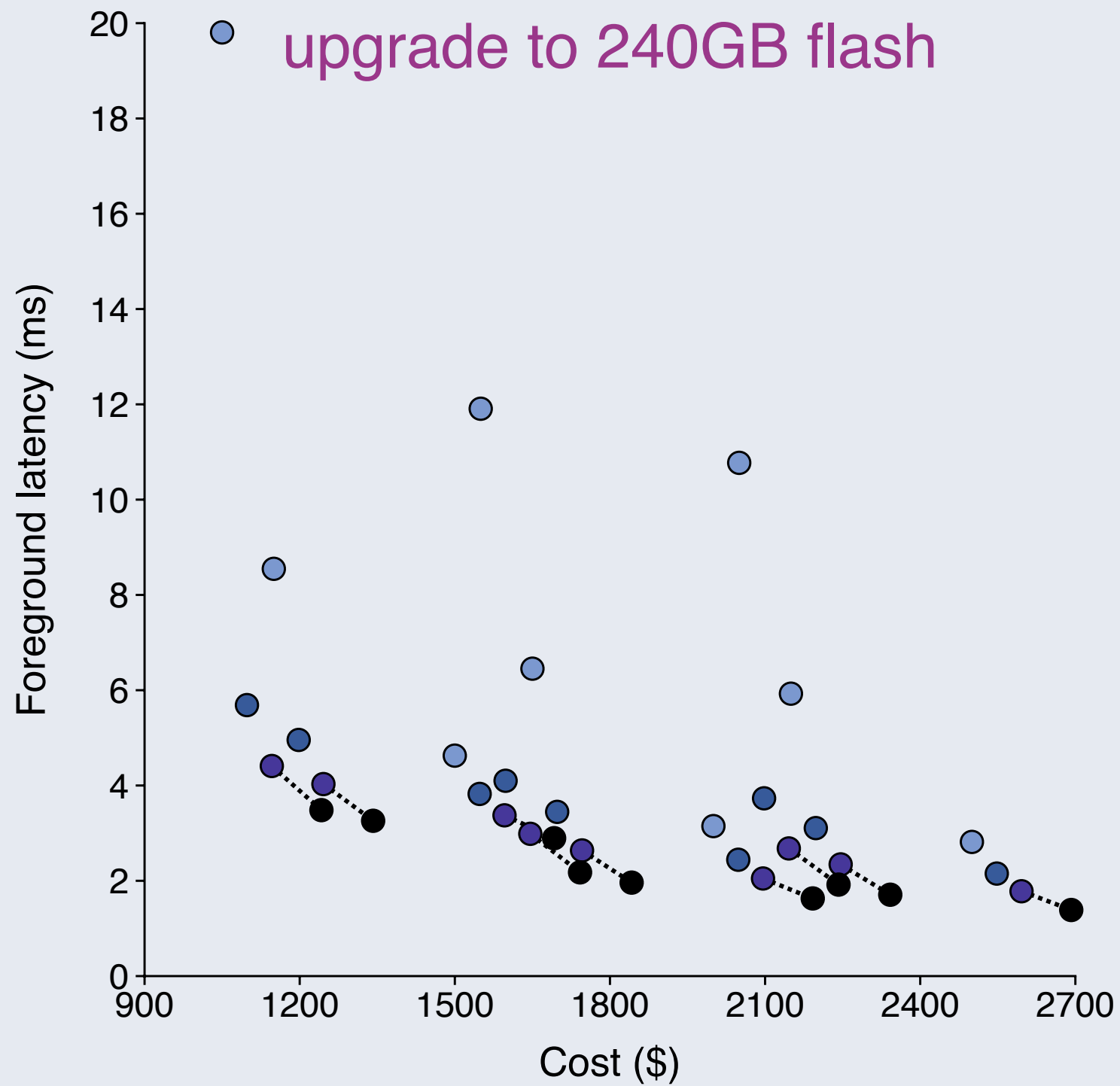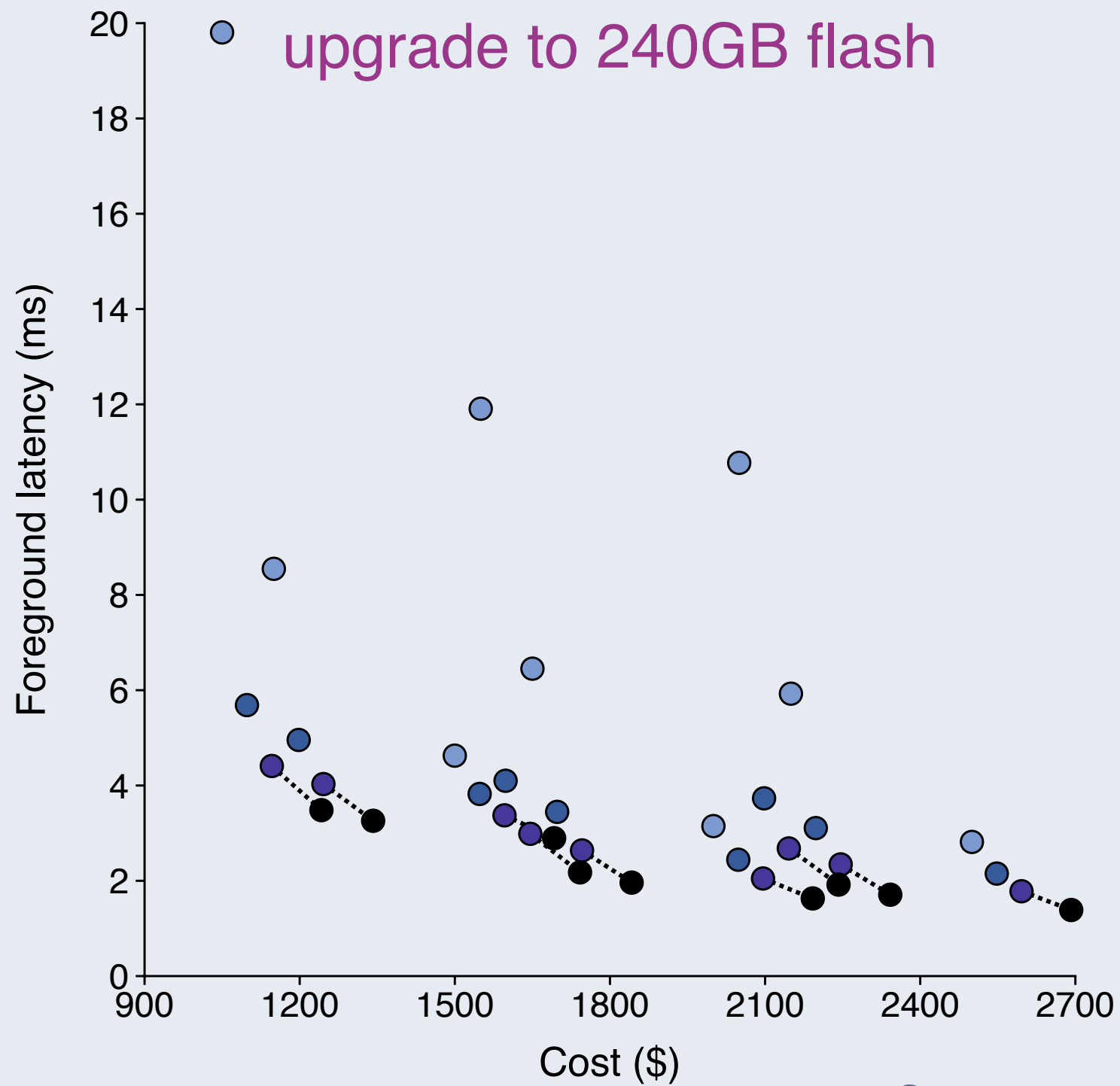Bad upgrade

# Outline

Intro

- Messages stack overview

- Methodology: trace-driven analysis and simulation

- HBase background

Results

- Workload analysis

- Hardware simulation: adding a flash layer

- Software simulation: integrating layers

Conclusions

# Software Architecture: Workload Implications

Writes are amplified

- 1% at HDFS (excluding overheads) to 64% at disk (given 30GB RAM)

- <u>We should optimize writes</u>

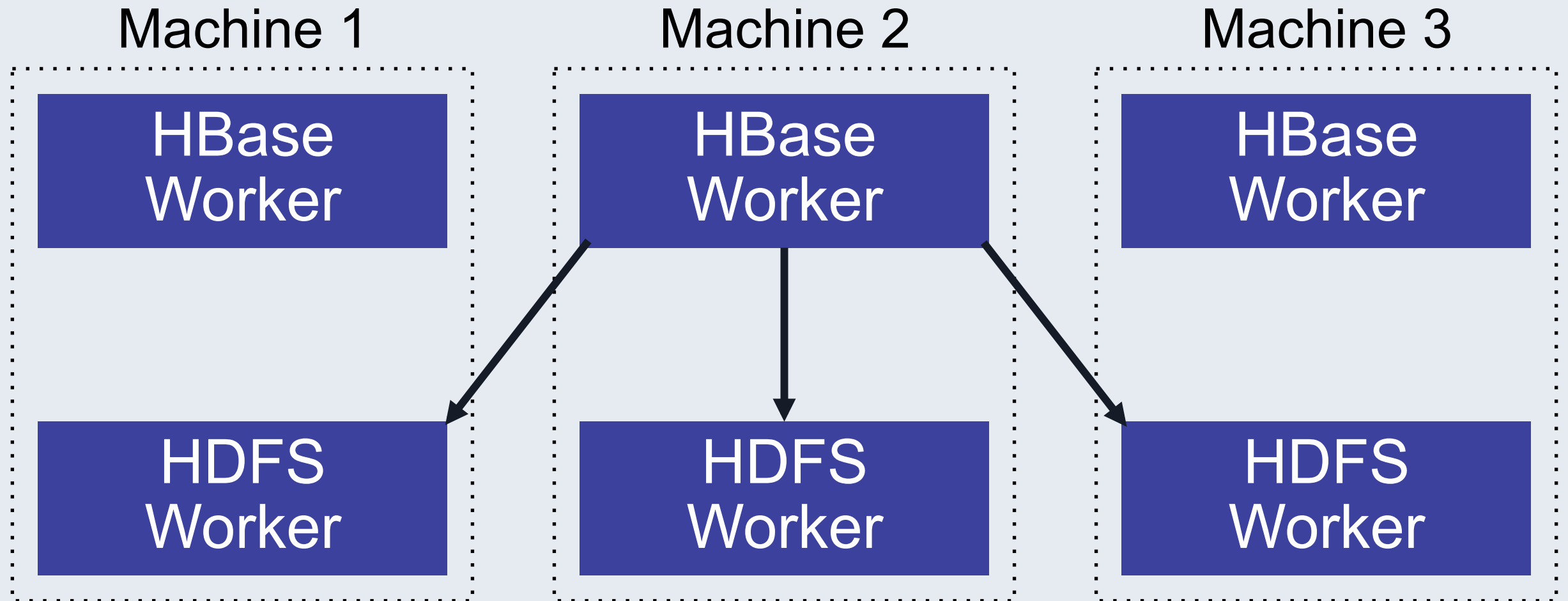# Software Architecture: Workload Implications

Writes are greatly amplified

- 1% at HDFS (excluding overheads) to 64% at disk
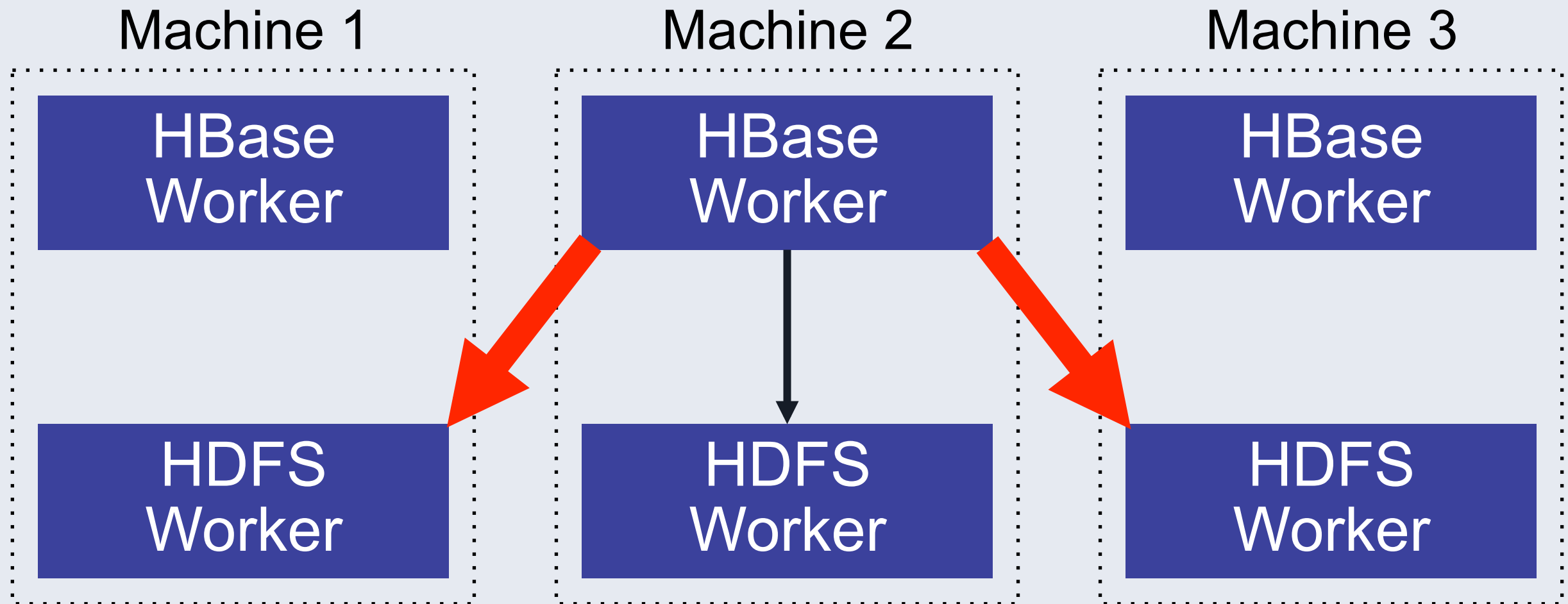
- <u>We should optimize writes</u>

61% of writes are for compaction

- <u>We should optimize compaction</u>

- Compaction interacts with replication inefficiently

# Solution: Ship Computation to Data

# Solution: do Local Compaction

Machine 1

Machine 2

Machine 3

HBase Worker

HBase Worker

HBase Worker

do compact

do compact

HDFS Worker

HDFS Worker

HDFS Worker

# Solution: do Local Compaction

# Local Compaction



Normally 3.5TB of network I/O

net (normal)

# Local Compaction

Normally 3.5TB of network I/O

Local comp: 62% reduction

# Local Compaction

Read I/O (TB) vs Cache size (GB)
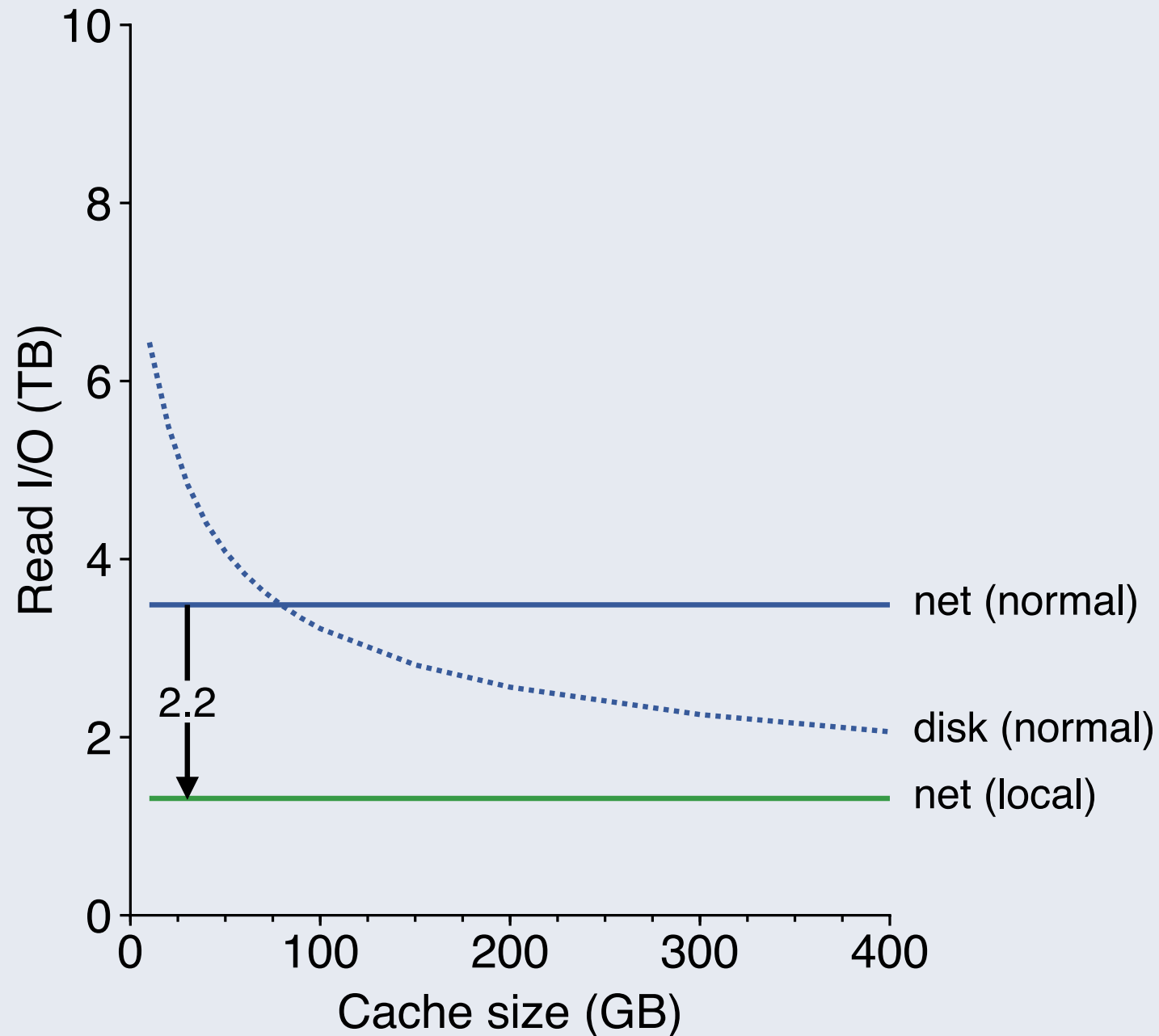
- net (normal) — 3.5
- disk (normal)
- net (local)
- 2.2

Normally 3.5TB of network I/O

Local comp: 62% reduction

# Local Compaction



Normally 3.5TB of network I/O

Local comp: 62% reduction

Network I/O becomes disk I/O

- 9% overhead (30GB cache)

- Compaction reads: (a) usually misses, (b) pollute cache

# Local Compaction



Normally 3.5TB of network I/O

Local comp: 62% reduction

Network I/O becomes disk I/O

- 9% overhead (30GB cache)

- Compaction reads: (a) usually misses, (b) pollute cache

Still good!

- Disk I/O is cheaper than network

# Outline

Intro

- Messages stack overview

- Methodology: trace-driven analysis and simulation

- HBase background

Results

- Workload analysis

- Hardware simulation: adding a flash layer

- Software simulation: integrating layers

Conclusions

# Conclusion 1: Messages is a New HDFS Workload

Original GFS paper:

- "*high sustained bandwidth is more important than low latency*"

- "*multi-GB files are the common case*"

We find <u>files are small and reads are random</u>

- 50% of files <750KB

- >75% of reads are random

# Conclusion 2: Layering is Not Free

Layering "*proved to be vital for the verification and logical soundness*" of the THE operating system ~ Dijkstra

We find <u>layering is not free</u>

- Over half of network I/O for replication is unnecessary

Layers can amplify writes, multiplicatively

- E.g., logging overhead (10x) with replication (3x) => 30x write increase

Layer integration can help

- Local compaction reduces network I/O caused by layers

# Conclusion 3: Flash Should not Replace Disk

Jim Gray predicted (for ~2012) that "*tape is dead, disk is tape, flash is disk*"

We find <u>flash is a poor disk replacement</u> for Messages

- Data is very large and mostly cold

- Pure flash would cost >$10K/machine

However, small flash tier is useful

- A 60GB SSD cache can double performance for a 5% cost increase

# Thank you!  Any questions?

University of Wisconsin-Madison          Facebook Inc.