# SILK: Preventing Latency Spikes in Log-Structured Merge Key-Value Stores

O. Balmau*, F. Dinu*, W. Zwaenepoel*,

K. Gupta[†], R. Chandhiramoorthi[†], D. Didona[§]

[*] THE UNIVERSITY OF SYDNEY    [†] NUTANIX    [§] IBM Research | Zurich

# Log-Structured Merge (LSM) KVs

✓ **Designed for write-heavy workloads**

✓ **Handle large-scale data**

✓ **Working set does not fit in RAM**

# Log-Structured Merge (LSM) KVs

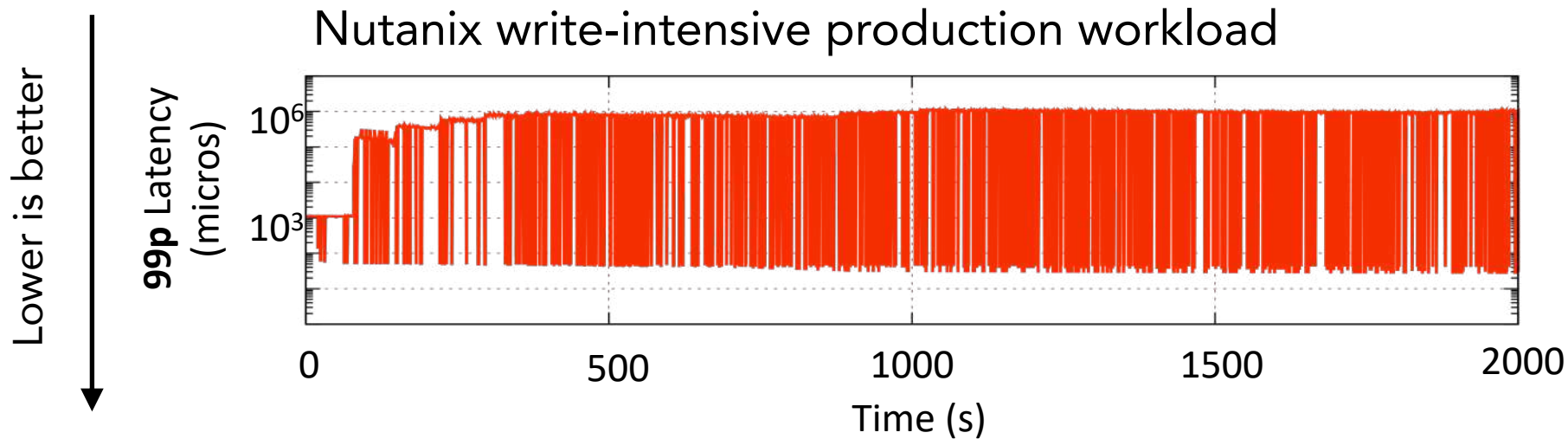✓ **Designed for write-heavy workloads?**

✓ **Handle large-scale data**
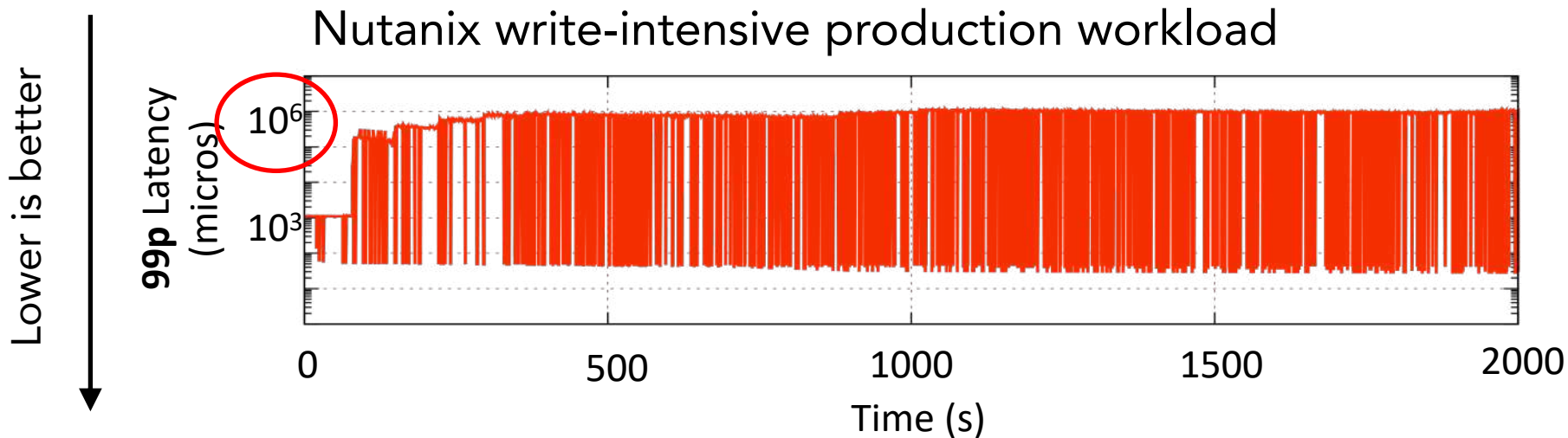
✓ **Working set does not fit in RAM**

# LSM KV Latency Spikes in RocksDB



Nutanix write-intensive production workload

# LSM KV Latency Spikes in RocksDB

Nutanix write-intensive production workload

Lower is better

**99p** Latency (micros)

$10^6$

$10^3$

0          500          1000          1500          2000

Time (s)

**Latency spikes of up to 1s** in write dominated workloads.

# **Latency Spikes** in LSM KVs

## **Why is this important?**

❌  Cannot provide SLA guarantees to clients.

❌  Unpredictable performance when connecting LSM in larger pipelines.

# Our Contribution: The SILK LSM KV

✓ **Solves latency spike problem for write-heavy workloads.**

✓ **No negative side-effects for other workloads.**

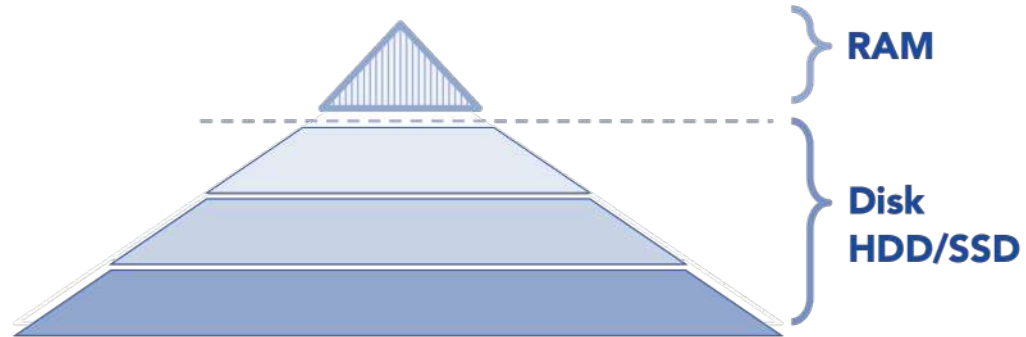✓ **SILK introduces the notion of an I/O scheduler for LSM KVs.**

# Experimental Study:
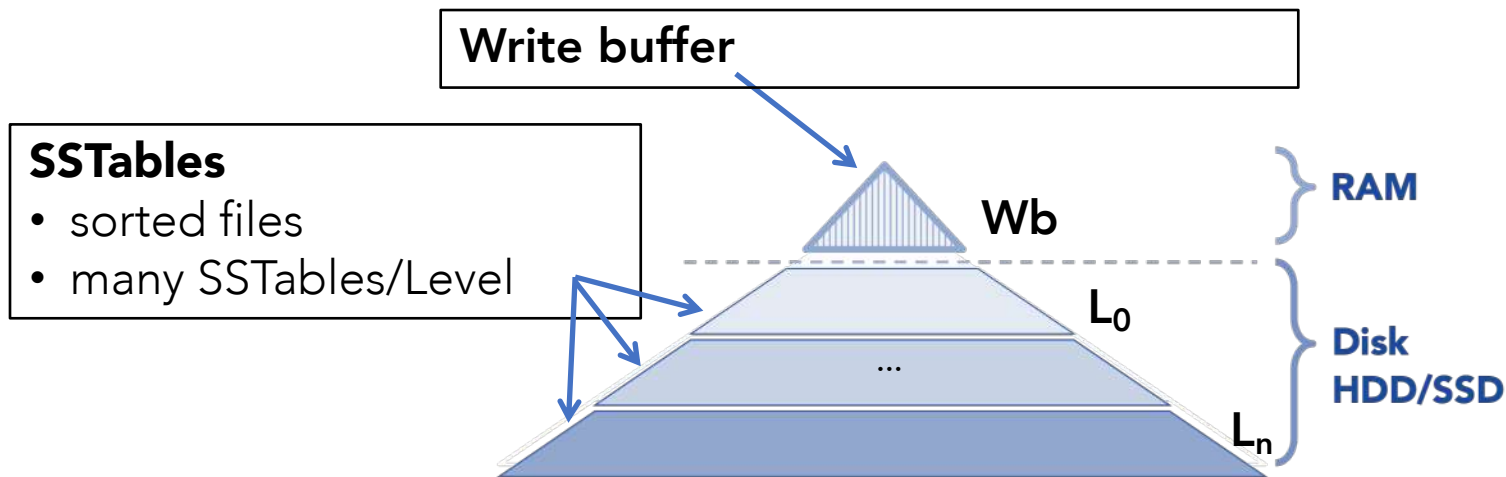# Reason Behind Latency Spikes

# What Causes LSM Latency Spikes?

**Severe competition for I/O bandwidth** between
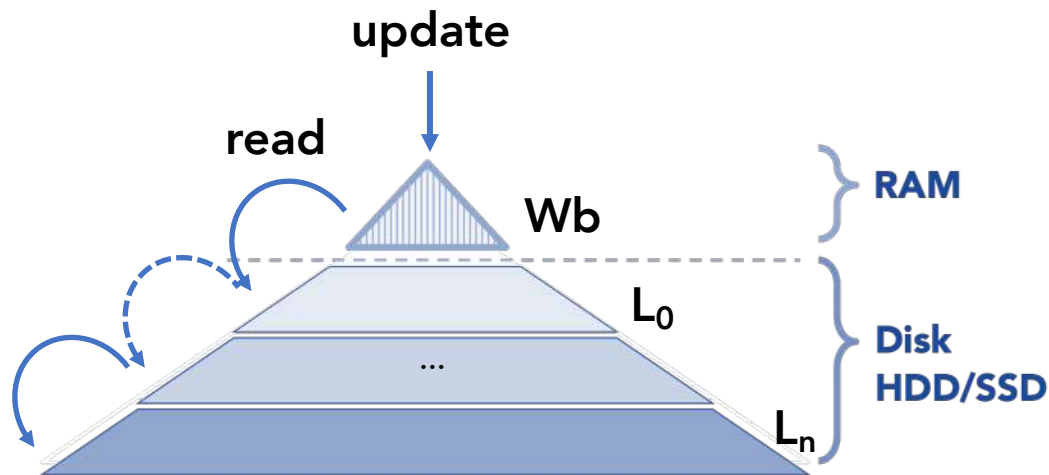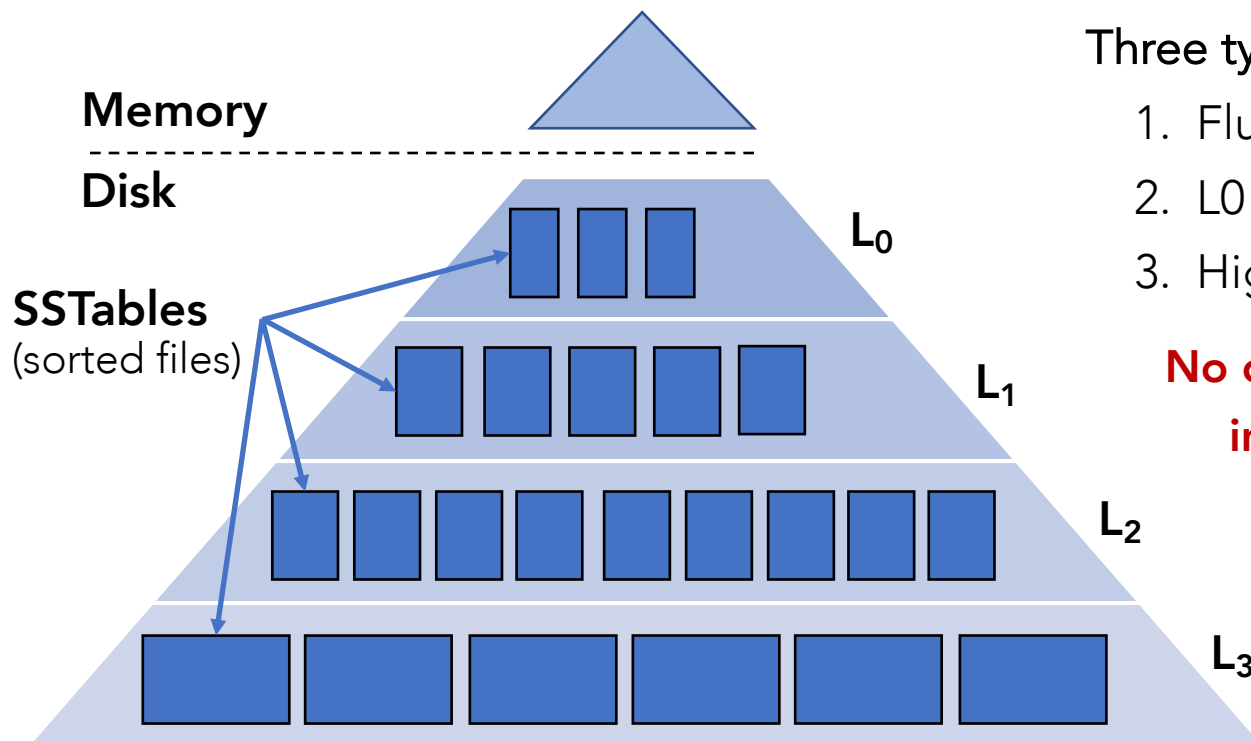client operations and LSM internal operations (~GC).

# LSM KV Overview



RAM

Disk
HDD/SSD

# LSM KV Overview

Write buffer

SSTables
- sorted files
- many SSTables/Level

Wb

RAM

$L_0$

...

$L_n$

Disk
HDD/SSD

# LSM KV Client Operations

# LSM Internal Ops

**Memory**

**Disk**

**SSTables**
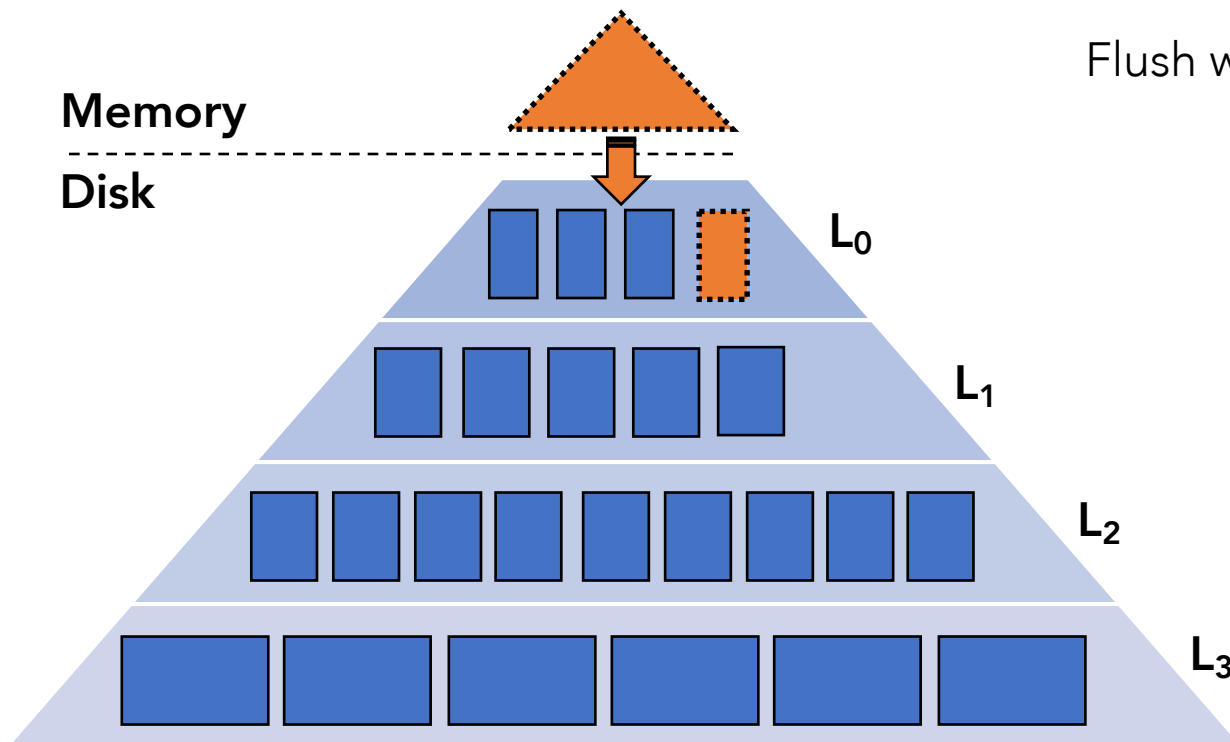(sorted files)

$L_0$

$L_1$

$L_2$

$L_3$

Three types of internal ops:

1. Flushing
2. L0 → L1 compaction
3. Higher level compactions

**No coordination between internal operations.**

# LSM Internal Ops: Flushing

Flush when Write buffer full.
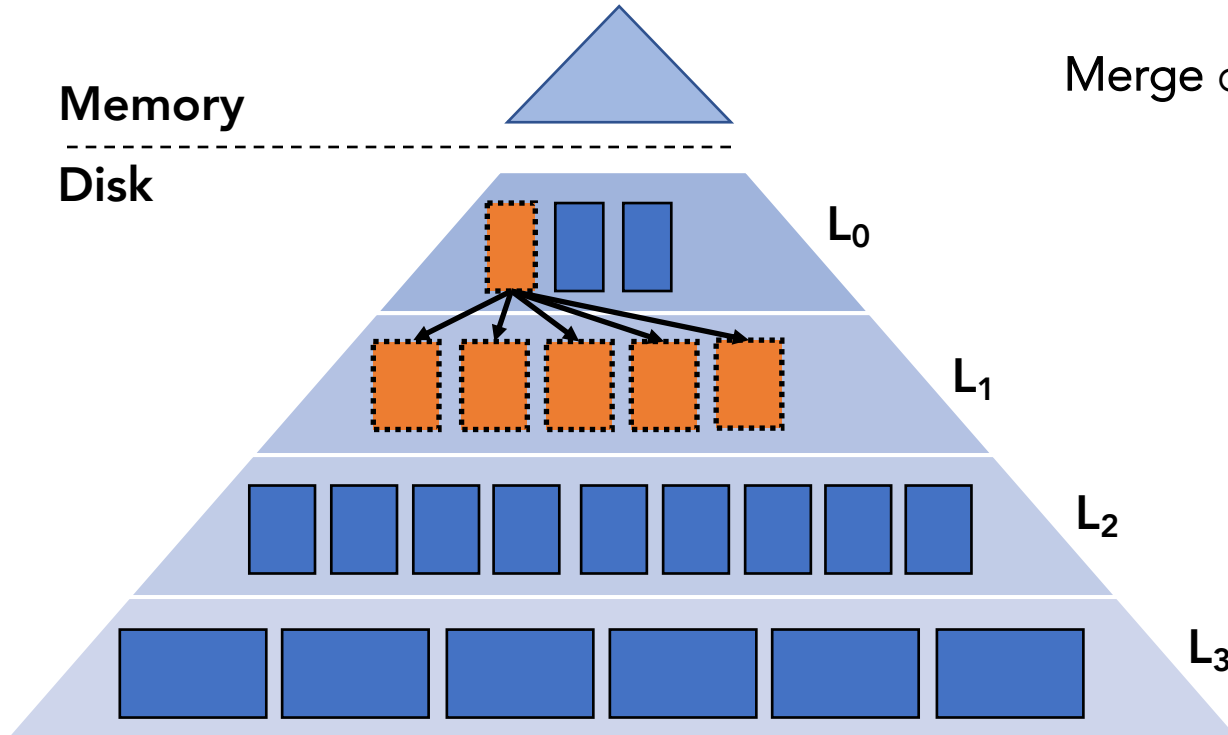


$L_0$

$L_1$

$L_2$

$L_3$

Memory

Disk

# LSM Internal Ops: Flushing



**Memory**

**Disk**

update → New write buffer

Flush buffer

L₁

L₂

L₃

Incoming writes absorbed in new **write buffer**.

**Flush buffer** written to L0.

# LSM Internal Ops: L0 → L1 compactions

Merge one L0 SSTable with L1.

**Memory**

**Disk**

$L_0$

$L_1$

$L_2$

$L_3$

# LSM Internal Ops: L0 → L1 compactions



**Memory**

**Disk**

$L_0$

$L_1$

$L_2$

$L_3$

**Merge** one L0 SSTable with L1.

Makes room on L0 for flushing.

# LSM Internal Ops: Higher Level Compactions



**Memory**

**Disk**

$L_0$

$L_1$

$L_2$

$L_3$

~GC in the LSM tree.

Discard duplicates & delete values.

Less urgent than L0➔L1 compactions.

… *but* need to complete.

I/O bandwidth intensive.

# LSM Internal Ops: Higher Level Compactions

**Memory**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Disk**

$L_0$

$L_1$

$L_3$

~GC in the LSM tree.

Discard duplicates & delete values.

Less urgent than L0➔L1 compactions.

… *but* need to complete.

Can have **many** higher level
compactions running **in parallel**.

# LSM Review

Internal operations:

1. **Flushing**. From memory to disk.

2. **L0 ➔ L1 compaction**. Make room to flush new files.

3. **Higher level compactions**. ~GC, I/O intensive.

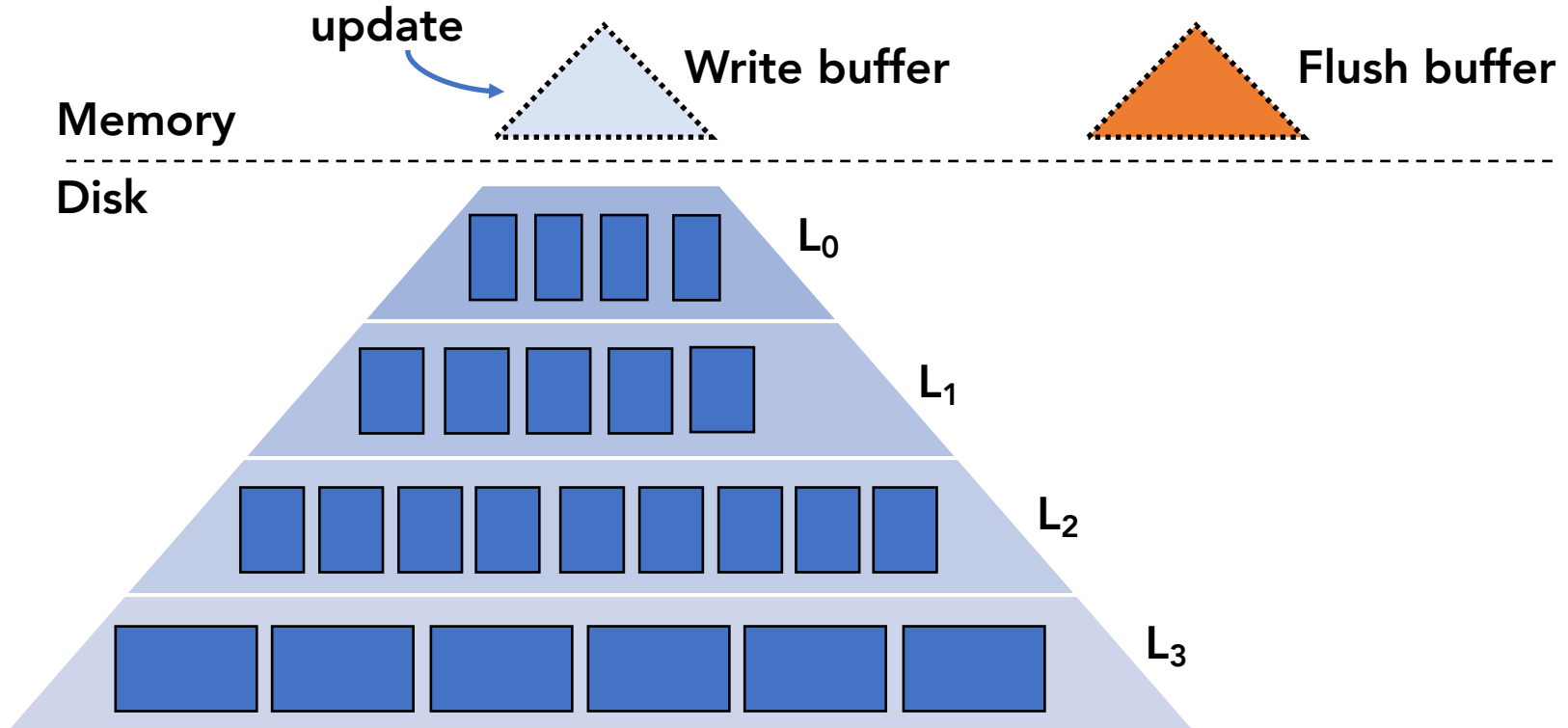⚠️ **No coordination between internal ops and client ops.**

# What Causes LSM Latency Spikes?

Both reads and writes experience latency spikes.
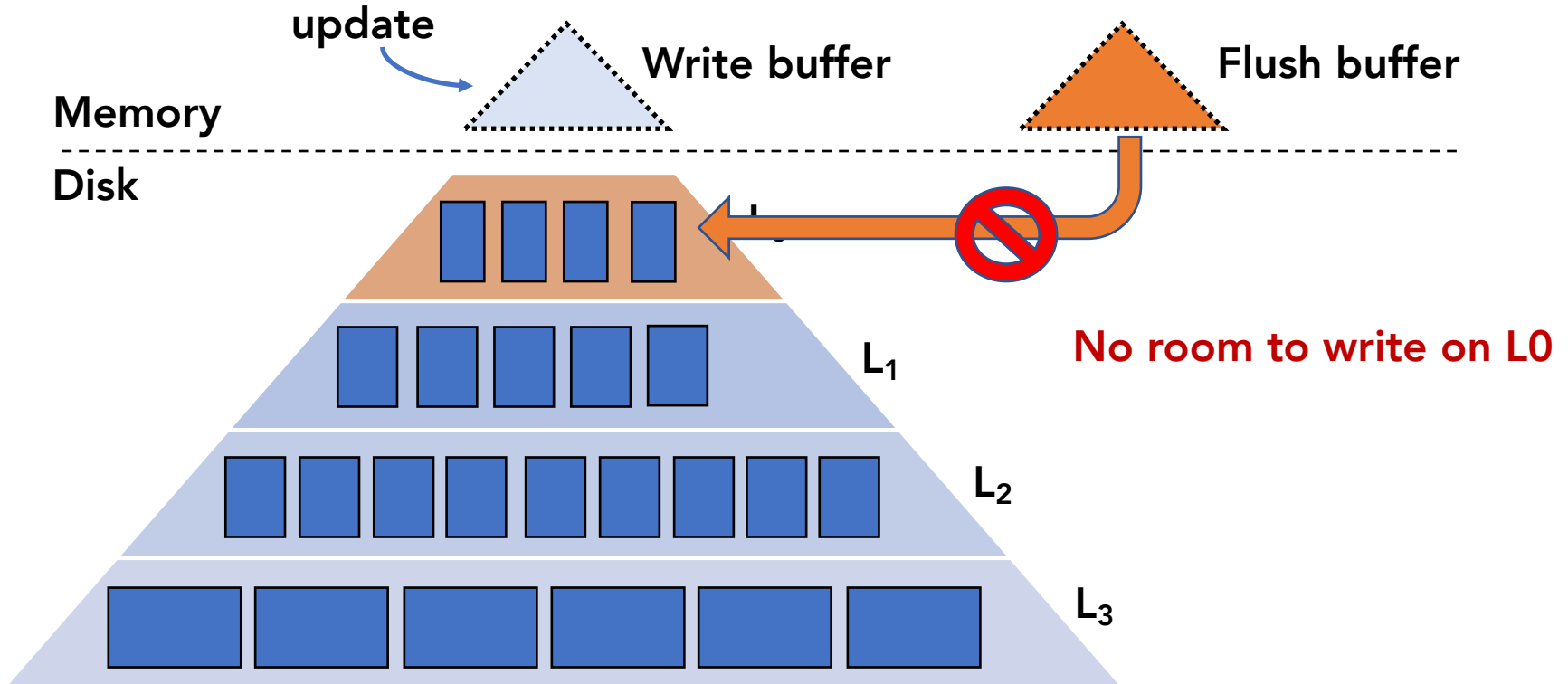
Focus on writes. Less intuitive.

Writes finish in memory. **Why do we have 1s latencies?**

# L0 Full, Cannot Flush

update

Write buffer

Flush buffer

Memory

Disk

$L_0$

$L_1$

$L_2$

$L_3$

# L0 Full, Cannot Flush

update

Write buffer

Flush buffer

Memory

Disk

L₀

L₁

L₂

L₃

No room to write on L0

# L0 Full, Cannot Flush

update

Write buffer

Flush buffer

Memory

Disk

$L_0$

$L_1$

No room to write on L0

$L_2$

$L_3$

# L0 Full, Cannot Flush



update

Write buffer

Flush buffer

Memory

Disk

L1

No room to write on L0

L2

L3

# L0 Full, Cannot Flush

update

Write buffer

Flush buffer

Memory

Disk

L$_1$

L$_2$

L$_3$

No room to write on L0

# L0 Full, Cannot Flush

update

Write buffer

Flush buffer

Memory

Disk

L₀

L₁

L₂

L₃

**No room to write on L0**

# 1. Writes Blocked Because L0 is Full.

**No coordination between internal ops.**

↓

Higher level compactions take over I/O.

↓
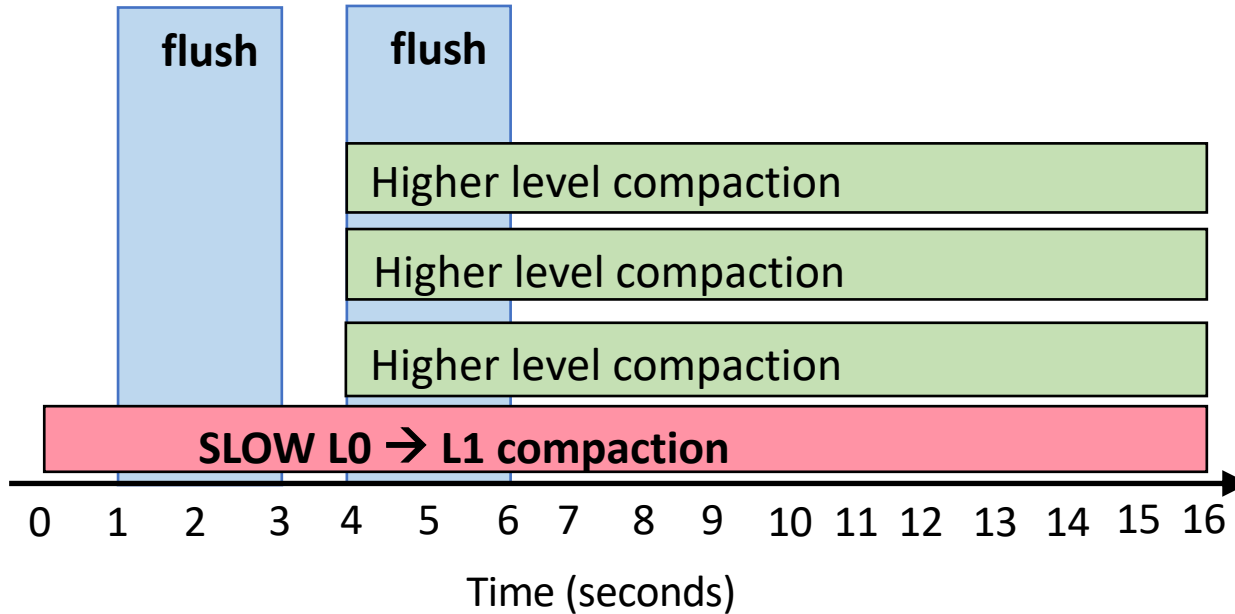
L0 → L1 compaction is too slow.
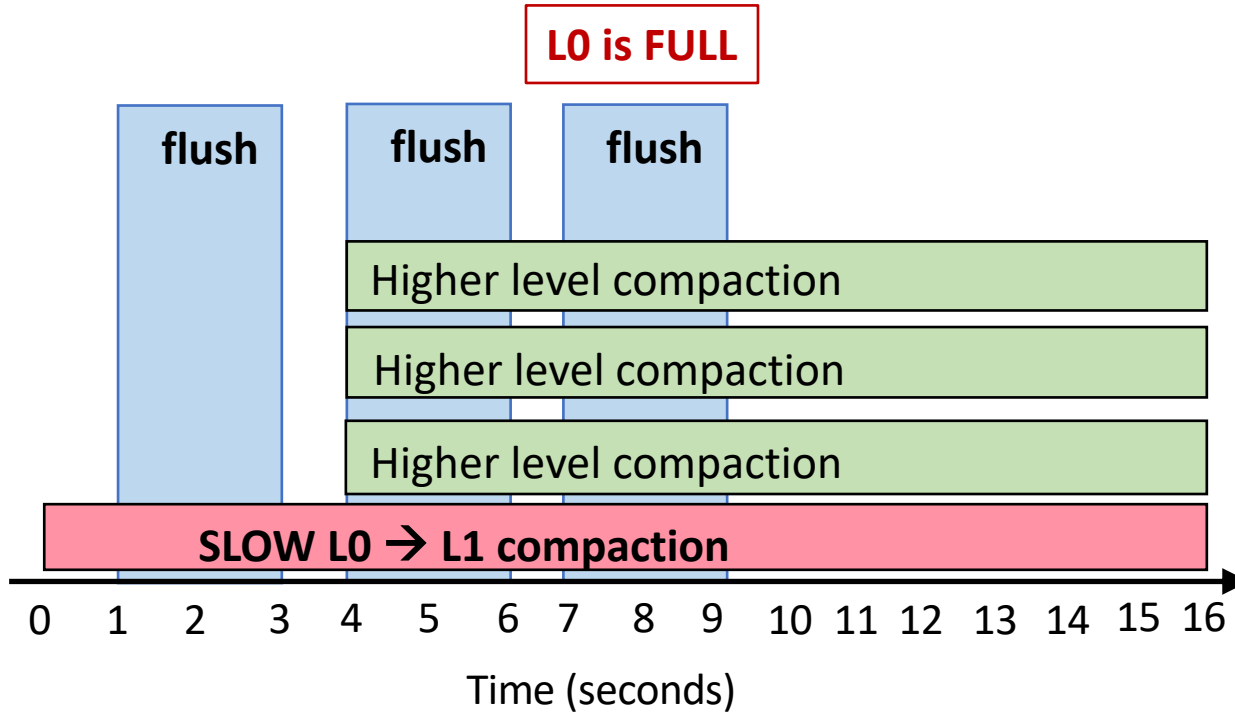
↓

Not enough space on L0.

↓

**Cannot flush memory component.**

# 1. Writes Blocked Because L0 is Full.

# 1. Writes Blocked Because L0 is Full.



L0 is FULL

flush     flush     flush

Higher level compaction

Higher level compaction

Higher level compaction

SLOW L0 → L1 compaction

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15  16

Time (seconds)

# 1. Writes Blocked Because L0 is Full.



L0 is FULL

Cannot flush.
No space on L0

flush

flush

flush

Higher level compaction

Higher level compaction

Higher level compaction

SLOW L0 → L1 compaction

0   1   2   3   4   5   6   7   8   9   10   11   12   13   14   15   16

Time (seconds)

# 1. Writes Blocked Because L0 is Full.



**L0 is FULL**

**Latency spike!**

flush   flush   flush

Higher level compaction

Higher level compaction

Higher level compaction

**SLOW L0 → L1 compaction**

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15  16

Time (seconds)

# Flushing is Slow

update

Write buffer

Flush buffer

**Memory**

**Disk**

$L_0$

$L_1$

$L_2$

$L_3$

# Flushing is Slow

update

Write buffer

Flush buffer

Memory

Disk

L₁

L₂

L₃

# Flushing is Slow

update

Write buffer

Flush buffer

Memory

Disk

L₁

L₂

L₃

# Flushing is Slow

update

Write buffer

Flush buffer

**Memory**

**Disk**

$L_0$

$L_1$

$L_2$

$L_3$

**Write buffer fills up before flush buffer is written to disk.**

# Flushing is Slow



update

Write buffer

Flush buffer

Memory

Disk

$L_0$

$L_1$

$L_2$

$L_3$

Write buffer fills up before flush buffer is written to disk.

# Flushing is Slow



update 🚫

Memory

Write buffer

Flush buffer

Disk

$L_0$

$L_1$

$L_2$

$L_3$

Write buffer fills up before flush buffer is written to disk.

# 2. Writes Blocked Because Flushing is Slow.

**No coordination between internal ops.**

↓

Higher level compactions take over I/O.

↓

Flushing does not have enough I/O.

↓

Flushing is very slow.

↓

**Memory component becomes full.**

# 2. Writes Blocked Because Flushing is Slow.
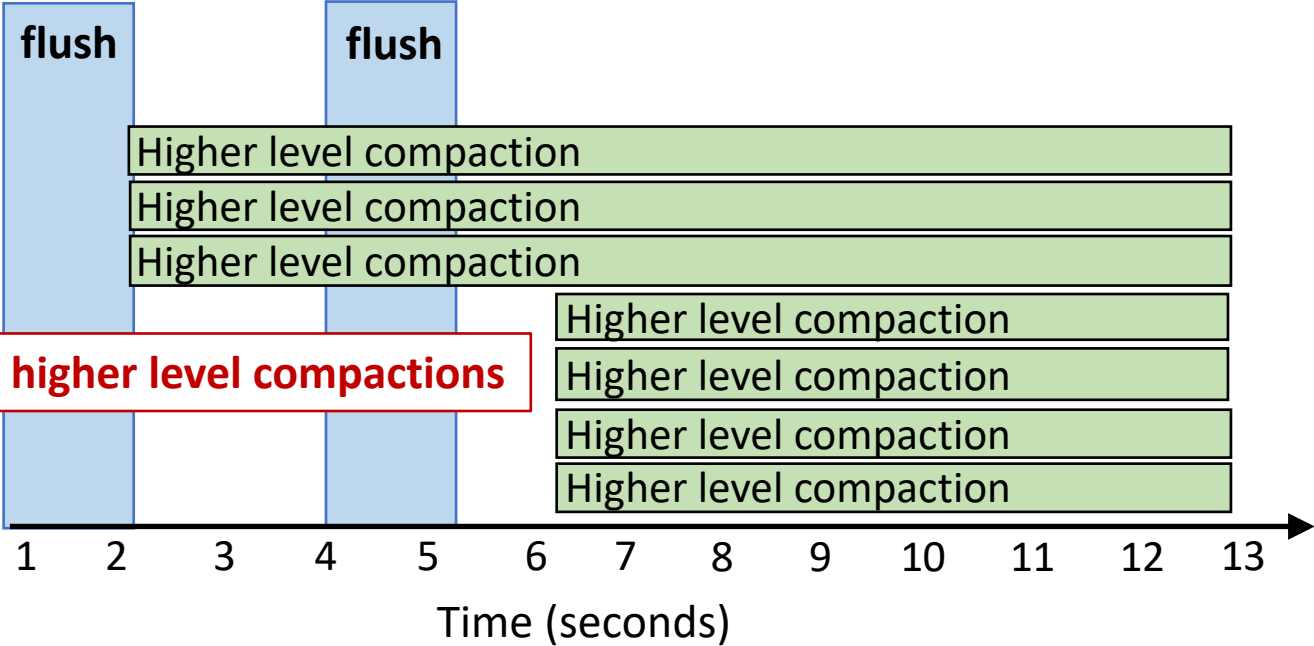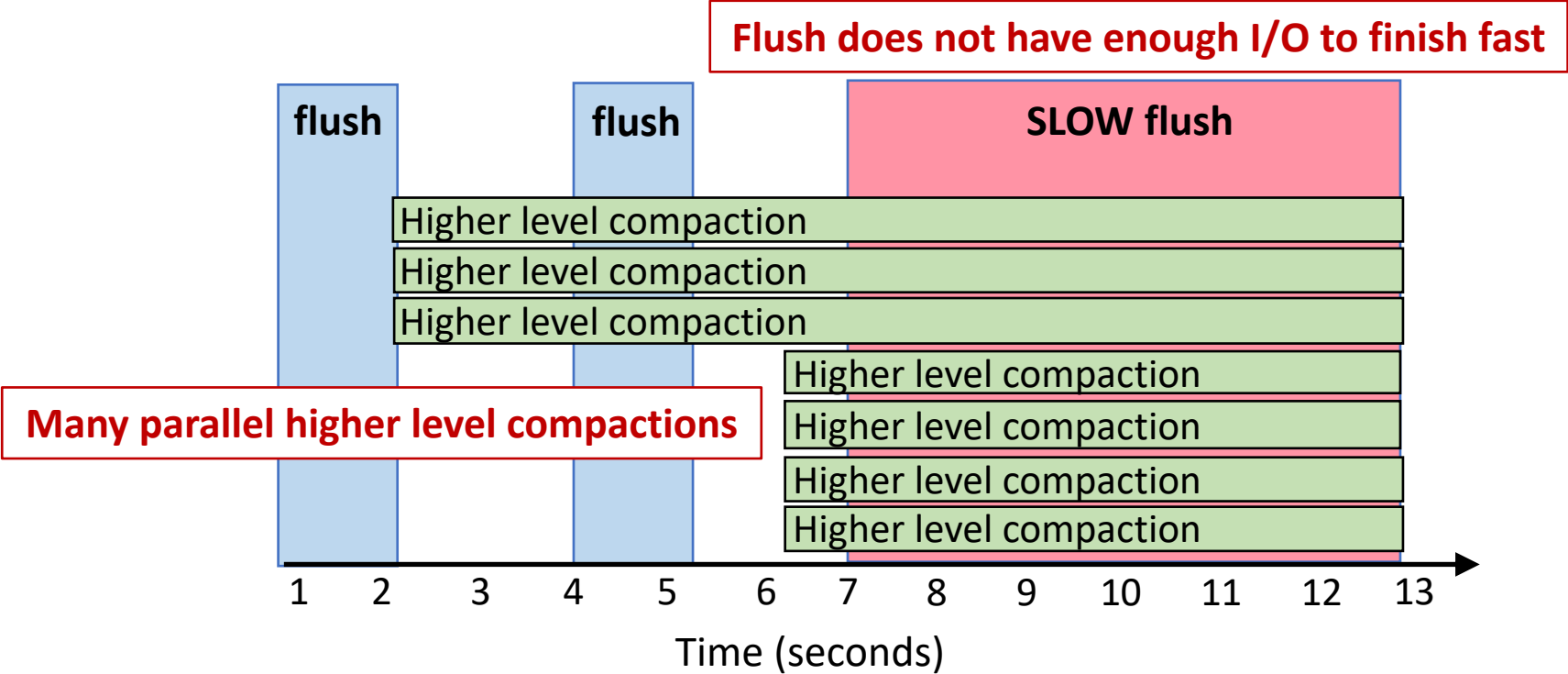
# 2. Writes Blocked Because Flushing is Slow.



Flush does not have enough I/O to finish fast

| flush | | flush | | SLOW flush |

Higher level compaction

Higher level compaction

Higher level compaction

Many parallel higher level compactions

Higher level compaction

Higher level compaction

Higher level compaction

Higher level compaction

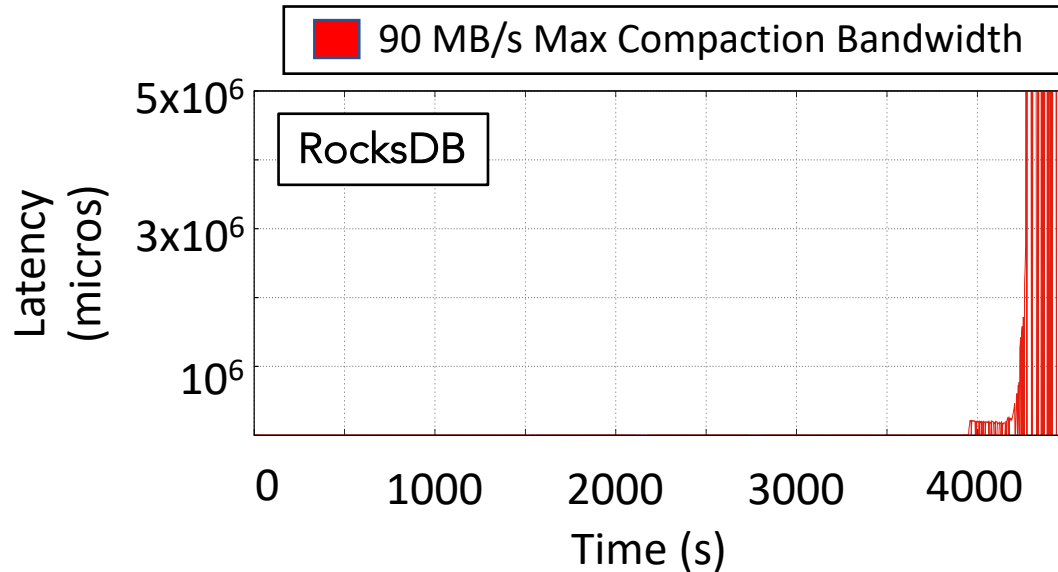1  2  3  4  5  6  7  8  9  10  11  12  13

Time (seconds)

# 2. Writes Blocked Because Flushing is Slow.

# Naïve Solution 1: Compaction Rate Limiting
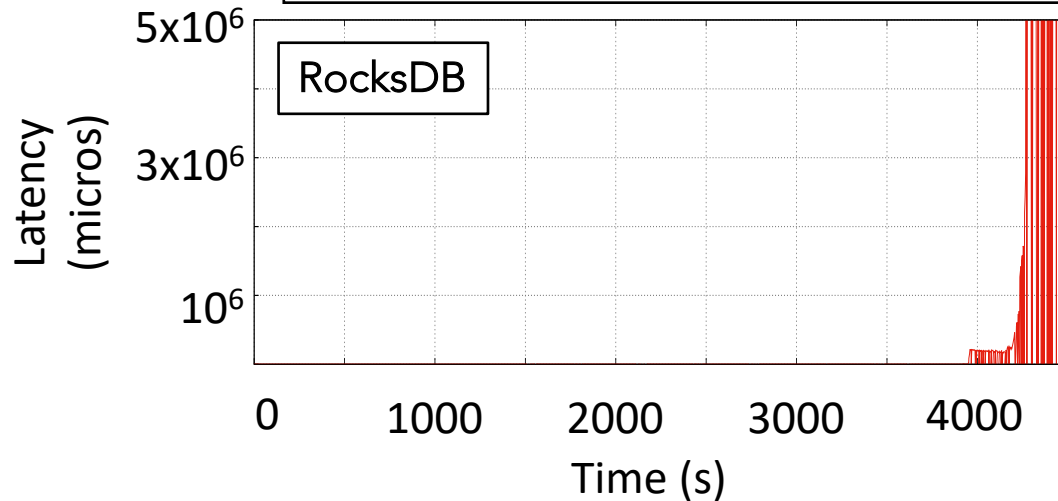
**Rate Limiting**: simple attempt to coordinate between internal and external ops.

# Naïve Solution 1: Compaction Rate Limiting

Rate l **Static compaction rate limiting does not work in the long term.** s.
Chance to run many parallel high level compactions increases.

# Naïve Solution 2: Delay Compaction Work

Selective/Delayed Compaction (TRIAD [USENIX ATC '17], PebblesDB [SOSP '17]).

# Naïve Solution 2: Delay Compaction Work

Select

**Being selective about compactions does not avoid interference.**
Eventually need to do the delayed compaction work.



Latency (micros) vs Time (s), labeled TRIAD. Y-axis: $10^3$ to $10^6$. X-axis: 0 to 2000.

# Lessons Learned

1. **Make sure L0 is never full.**

# Lessons Learned

1. **Make sure L0 is never full.**

2. **Ensure sufficient I/O for flush/compactions on low levels.**

# Lessons Learned

1. Make sure L0 is never full.

2. Ensure sufficient I/O for flush/compactions on low levels.

3. Higher level compactions should not fall behind too much.

# The SILK I/O Scheduler

# SILK Key Idea

**I/O scheduler** for LSM KVs: **coordinate I/O bandwidth sharing** to **minimize interference** between internal ops and client ops.

# Lessons Learned

Make sure L0 is never full.

Ensure sufficient I/O for flush/ compactions on low levels.

Make sure other compactions do not fall behind too much.

# SILK Design

# Lessons Learned

# SILK Design

**Make sure L0 is never full.**

**Prioritize internal operations at lower levels of the tree.**

Ensure sufficient I/O for flush/ compactions on low levels.

Make sure other compactions do not fall behind too much.

53

# Lessons Learned

Make sure L0 is never full.

**Ensure sufficient I/O for flush/ compactions on low levels.**

Make sure other compactions do not fall behind too much.

# SILK Design

Prioritize internal operations at lower levels of the tree.

**Preempt higher level compactions if necessary.**

# Lessons Learned

Make sure L0 is never full.

Ensure sufficient I/O for flush/compactions on low levels.

**Make sure other compactions do not fall behind too much.**

# SILK Design

Prioritize internal operations at lower levels of the tree.

Preempt higher level compactions if necessary.

**Opportunistically allocate I/O for higher level compactions.**

# Prioritize & Preempt

Prioritize internal ops at **lower tree levels:**

**1**   **Flushing**

**2**   **L0 → L1 compactions**

**3**   **Higher level compactions**

# Prioritize & Preempt

Prioritize internal ops at **lower tree levels:**

**1** **Flushing** *– dedicated flush operation queue.*

**2** **L0 → L1 compactions**

**3** **Higher level compactions**

# Prioritize & Preempt

Prioritize internal ops at **lower tree levels:**

**1** **Flushing** *– dedicated flush operation queue.*

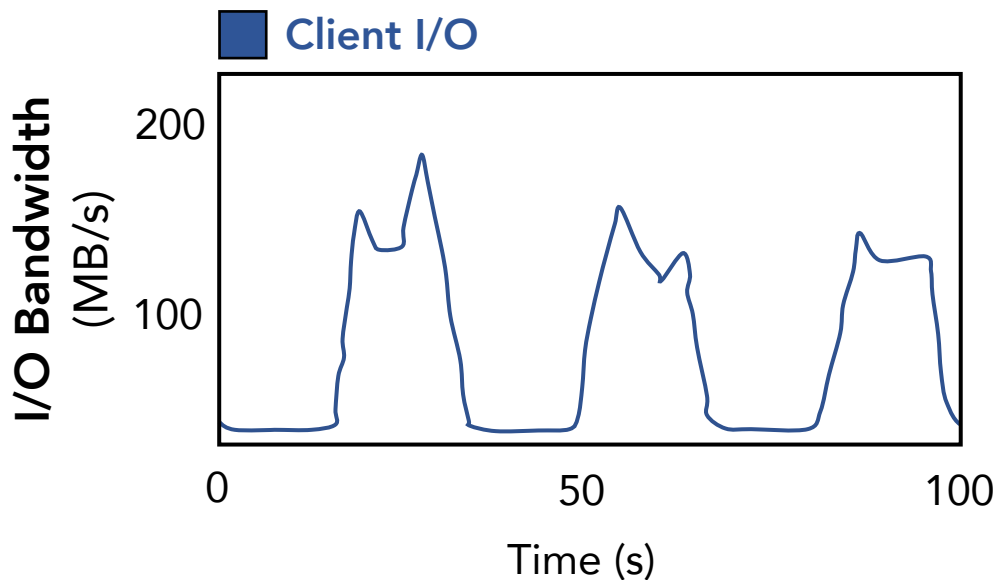**2** **L0 → L1 compactions**

**3** **Higher level compactions**

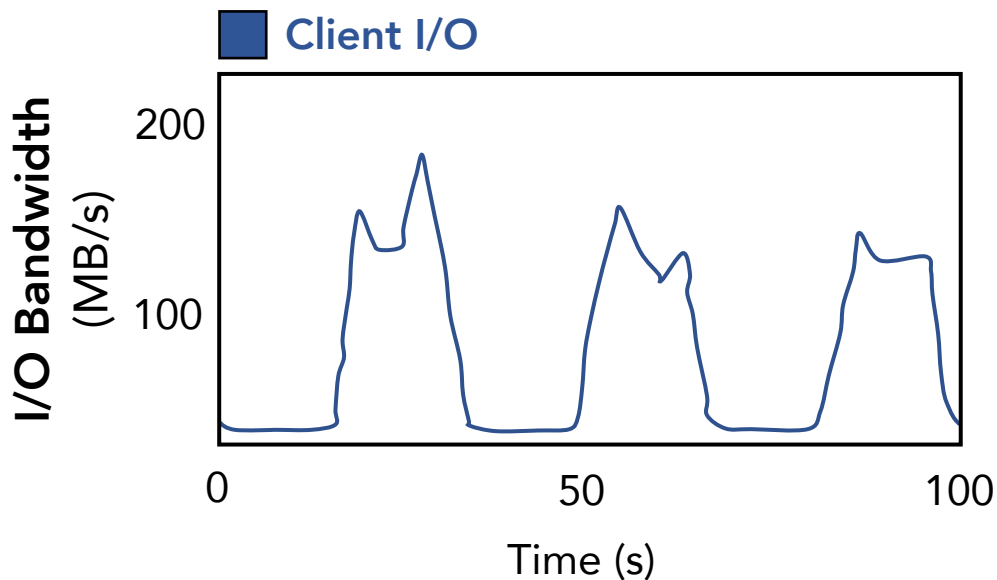*L0 → L1 compaction preempts higher level compactions.*

# Opportunistically allocate I/O for compactions

**Real Nutanix client load example**
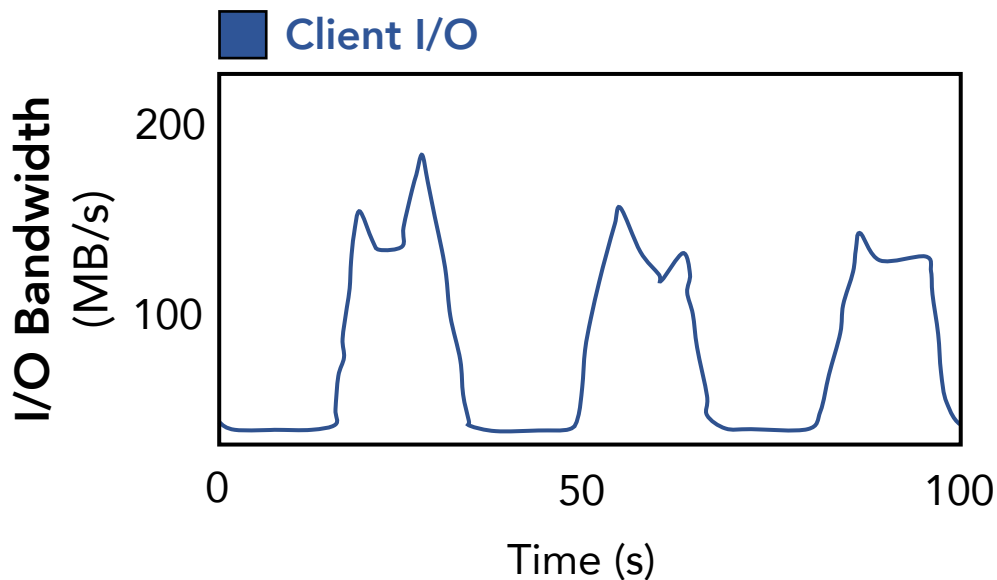
# Opportunistically allocate I/O for compactions

**Real Nutanix client load example**

■ **Client I/O**

Client workload is **not constant**.



I/O Bandwidth (MB/s)

200

100

0        50       100

Time (s)

# Opportunistically allocate I/O for compactions

**Real Nutanix client load example**



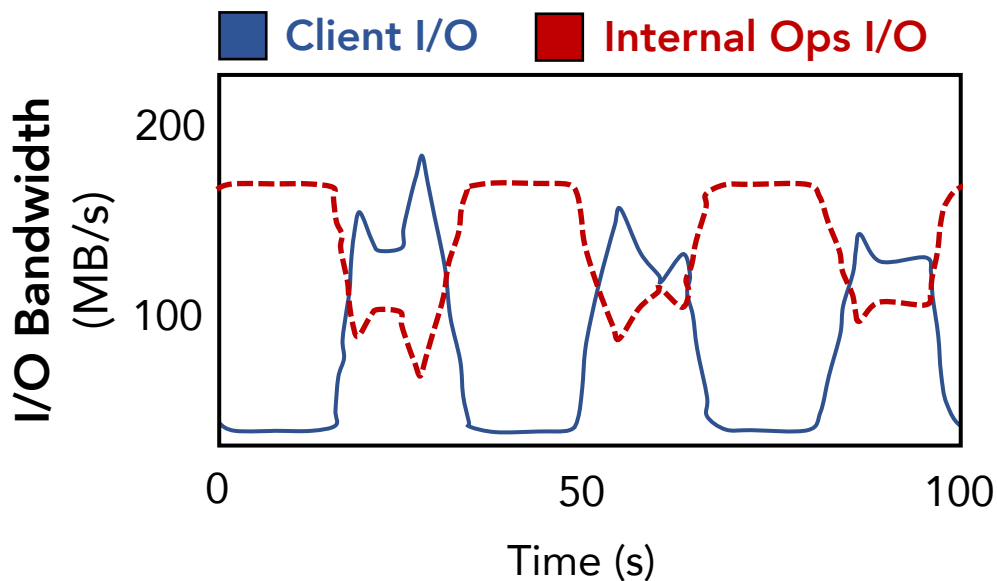Client workload is **not constant**.

SILK **continuously monitors** client **I/O bandwidth** use.

# Opportunistically allocate I/O for compactions

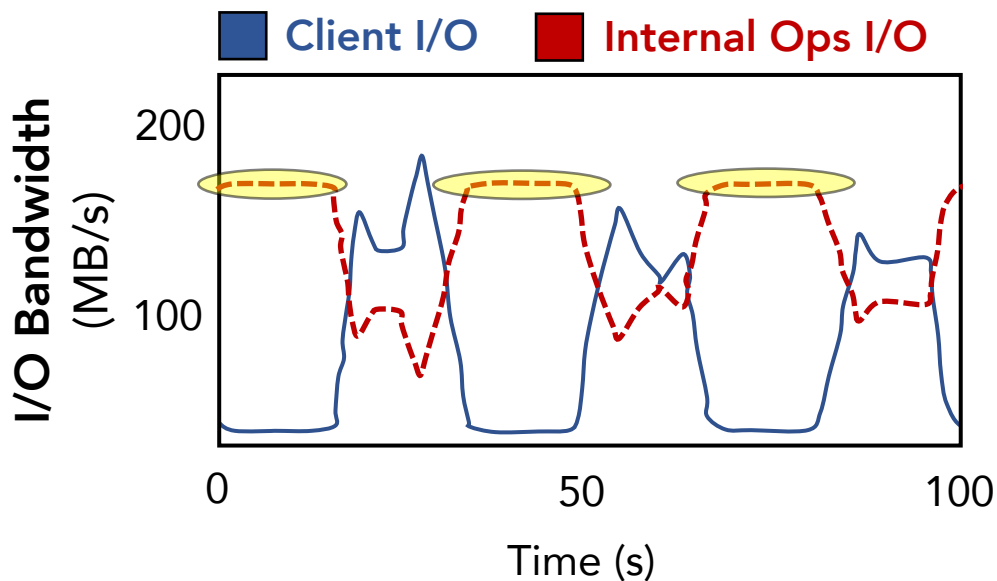**Real Nutanix client load example**



Allocate **less I/O to compactions** during **client load peaks.**

Allocate **more I/O to compactions** during **low client load.**

# Opportunistically allocate I/O for compactions

**Real Nutanix client load example**



More I/O to high level compactions during low load → **don't fall behind**.

# SILK Evaluation

# SILK Implementation

**Extends RocksDB.**



**Open Source** `https://github.com/theoanab/SILK-USENIXATC2019`
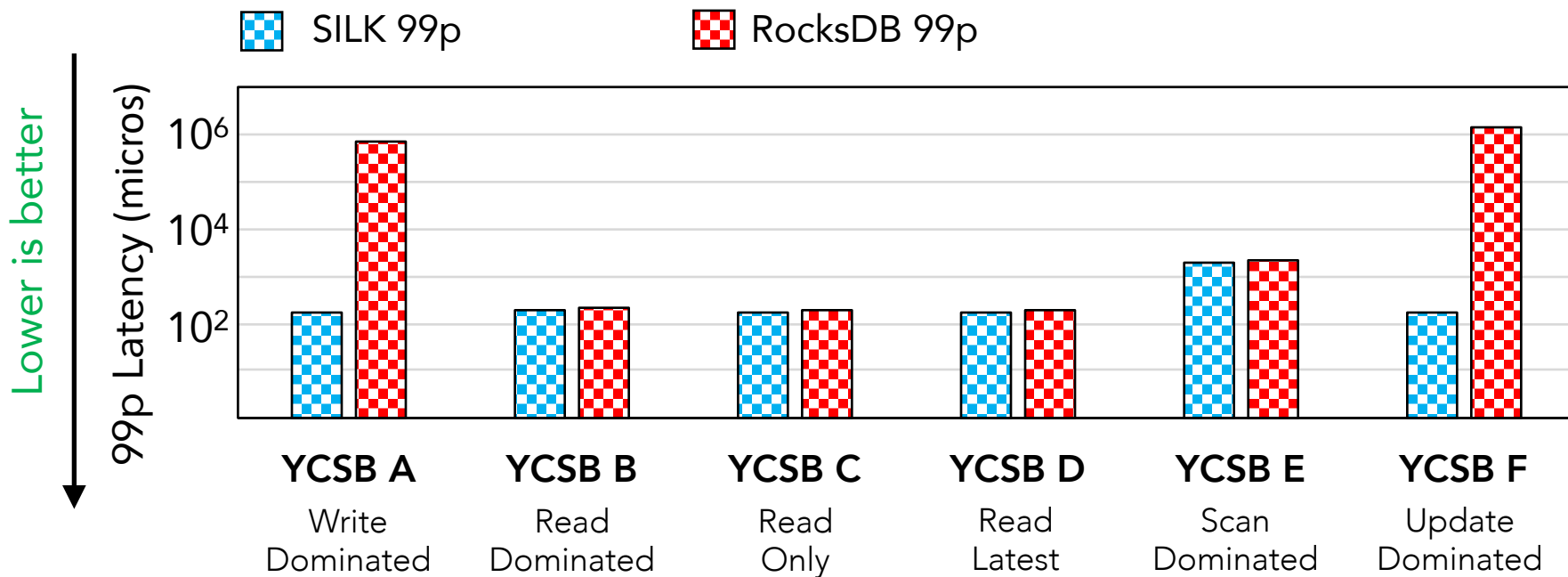
# YCSB

Benchmark with different workloads:

write-intensive, read-intensive, scan-intensive.
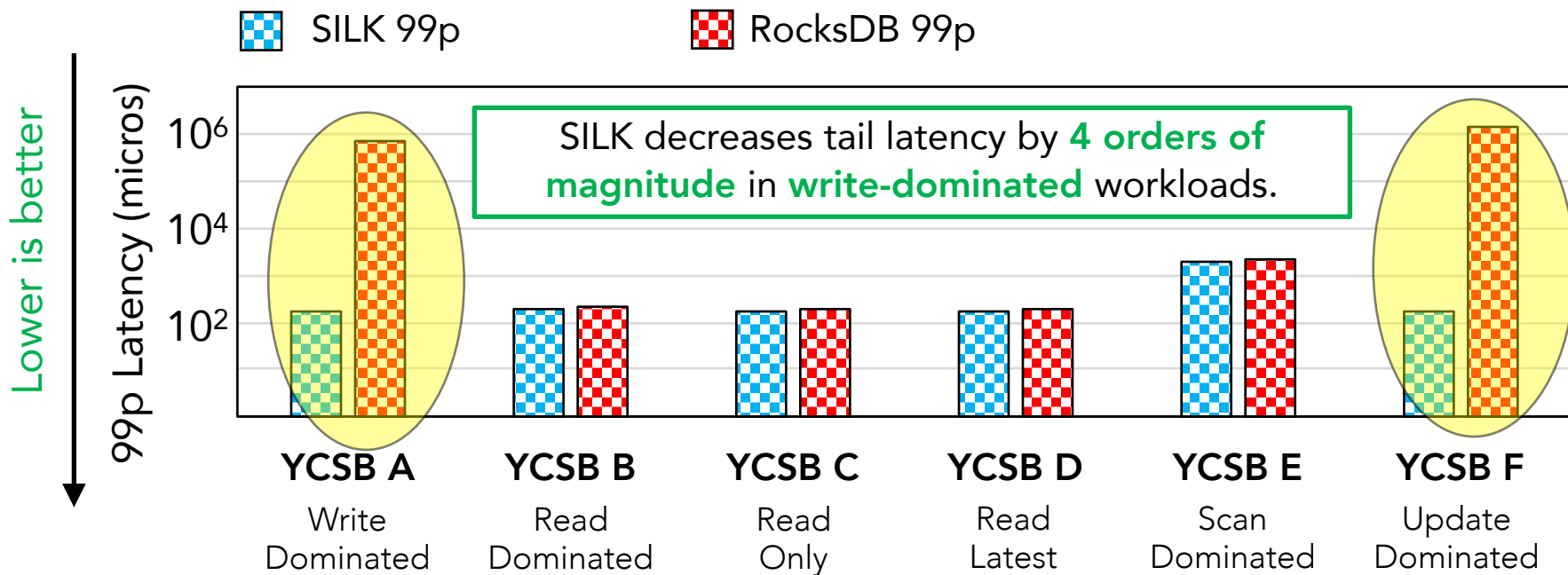
Show:

1. Write-heavy workloads: SILK is much better for tail latency.

2. Other workloads: SILK is not detrimental.

# YCSB

# YCSB



**SILK 99p**  **RocksDB 99p**

Lower is better

99p Latency (micros)

$10^6$

$10^4$

$10^2$

SILK decreases tail latency by **4 orders of magnitude** in **write-dominated** workloads.

| YCSB A | YCSB B | YCSB C | YCSB D | YCSB E | YCSB F |
|---|---|---|---|---|---|
| Write Dominated | Read Dominated | Read Only | Read Latest | Scan Dominated | Update Dominated |

# YCSB



SILK 99p      RocksDB 99p

SILK **does not affect read/scan dominated** workloads

Lower is better

99p Latency (micros)

$10^6$

$10^4$

$10^2$

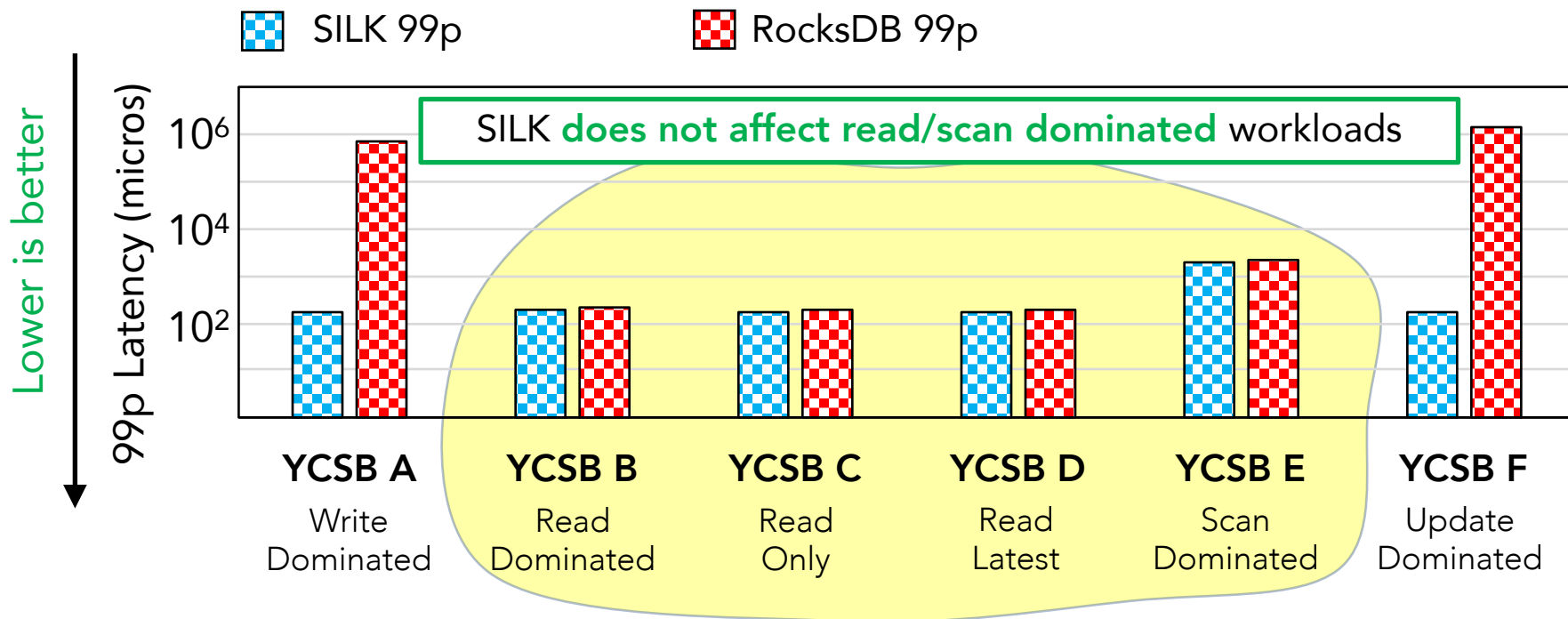| YCSB A | YCSB B | YCSB C | YCSB D | YCSB E | YCSB F |
|--------|--------|--------|--------|--------|--------|
| Write Dominated | Read Dominated | Read Only | Read Latest | Scan Dominated | Update Dominated |

# Nutanix Production Workload
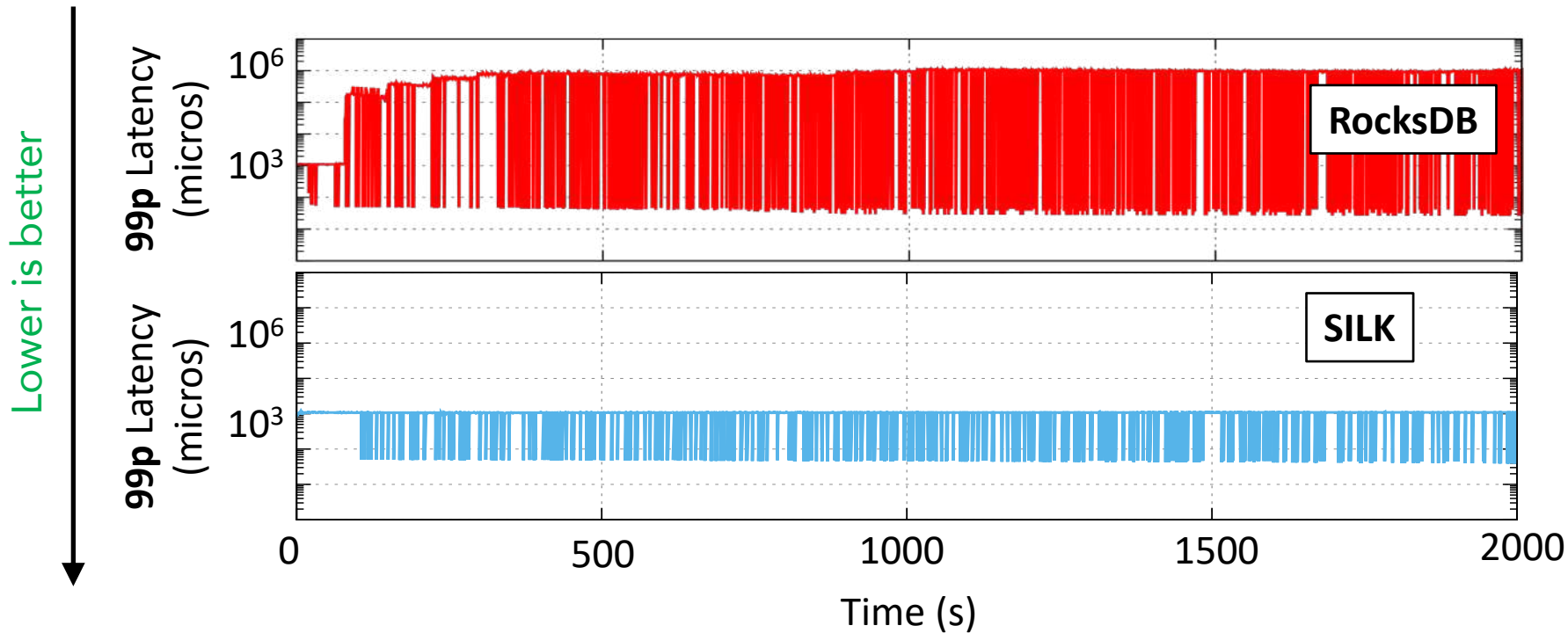
**Write dominated:**

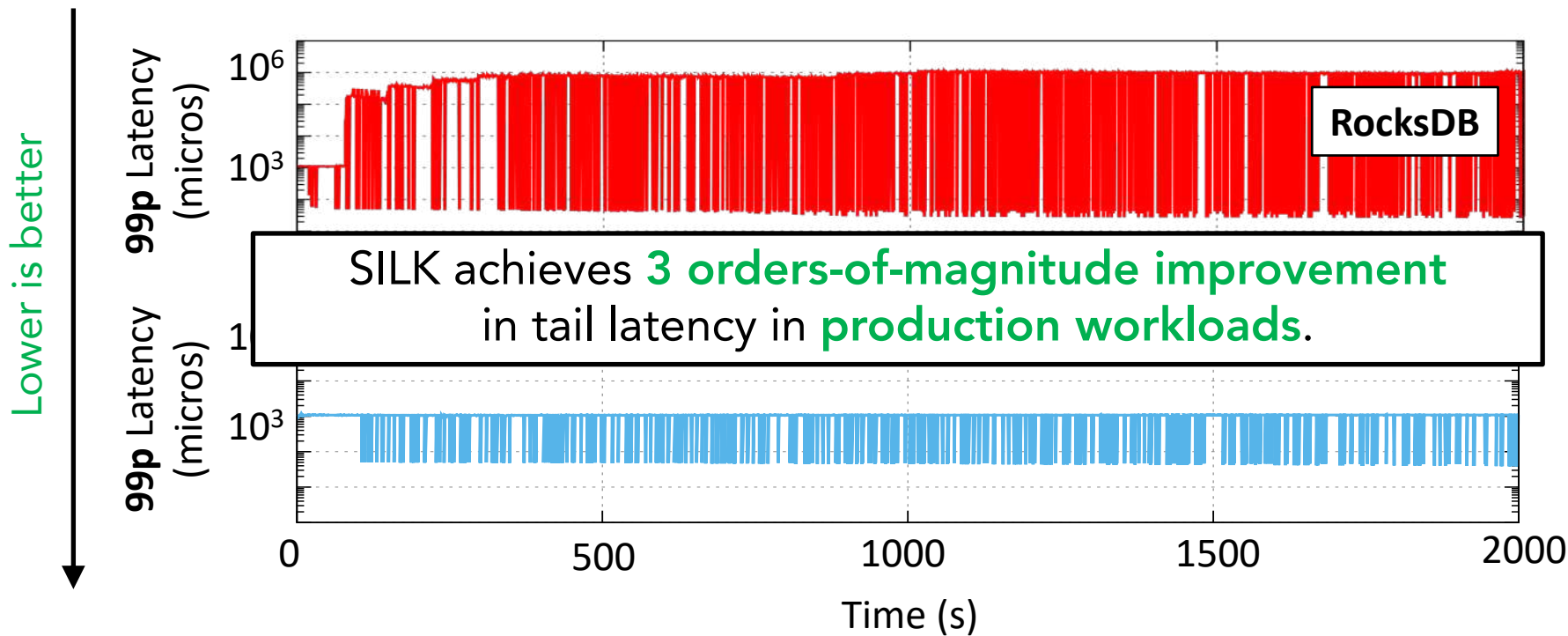   57% writes, 41% reads, 2% scans.

**Bursty (open loop):**

   Peaks and valleys in client load.
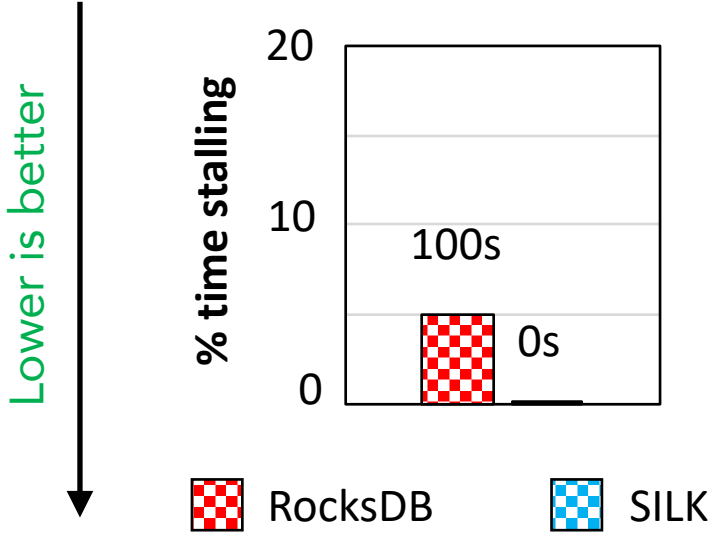
**Dataset size:**  500GB, KV tuple size 400B on average.

# SILK vs RocksDB Tail Latency 99P



Lower is better

# SILK vs RocksDB Tail Latency 99P



Lower is better

99p Latency (micros)

$10^6$

$10^3$

RocksDB

SILK achieves **3 orders-of-magnitude improvement** in tail latency in **production workloads**.

99p Latency (micros)

1

$10^3$

0   500   1000   1500   2000

Time (s)

# SILK vs RocksDB Stalling

# SILK vs RocksDB Stalling

Lower is better

% time stalling

20

10

0

100s

0s

SILK never stalls because it can always do timely flushing.

RocksDB    SILK

# More in the paper…

✓ **More experiments and workloads.**

✓ **With SILK, throughput is steady and close to the client load.**

✓ **Comparison with more state-of-the-art LSMs (TRIAD, PebblesDB).**

# SILK Take-Home Message

- We introduce the **new concept** of an **I/O scheduler for LSM**.

- **Coordinate I/O sharing** to avoid latency spikes.

- **Three orders-of-magnitude improvements** on tail latency.

**Thank you! Questions?**