# TRIAD: Creating Synergies Between Memory, Disk and Log in Log Structured Key-Value Stores

*USENIX ATC '17, Santa Clara CA*

O. Balmau
EPFL

D. Didona
EPFL

R. Guerraoui
EPFL

W. Zwaenepoel
EPFL

H. Yuan
Nutanix

A. Arora
Nutanix

K. Gupta
Nutanix

P. Konka
Nutanix

# KV Stores

Very **simple** data stores.

**KV pairs.**

Simple operations: **update, read.**
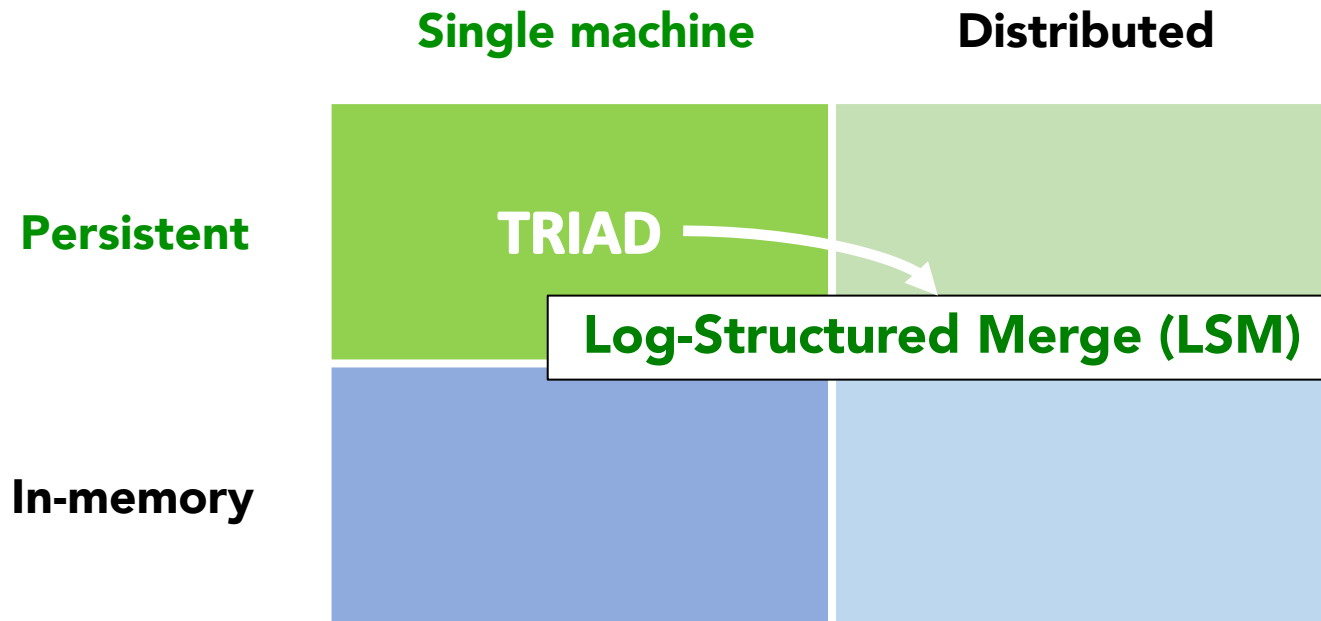
# KV Stores

|  | **Single machine** | **Distributed** |
|---|---|---|
| **Persistent** | | |
| **In-memory** | | |

# KV Stores

|  | **Single machine** | **Distributed** |
|---|---|---|
| **Persistent** | **TRIAD** | |
| **In-memory** | | |

# KV Stores

|  | Single machine | Distributed |
|---|---|---|
| Persistent | TRIAD | |
| In-memory | | |

Log-Structured Merge (LSM)

# TRIAD in a Nutshell

TRIAD LSM KV: achieves 2x throughput on production wklds.

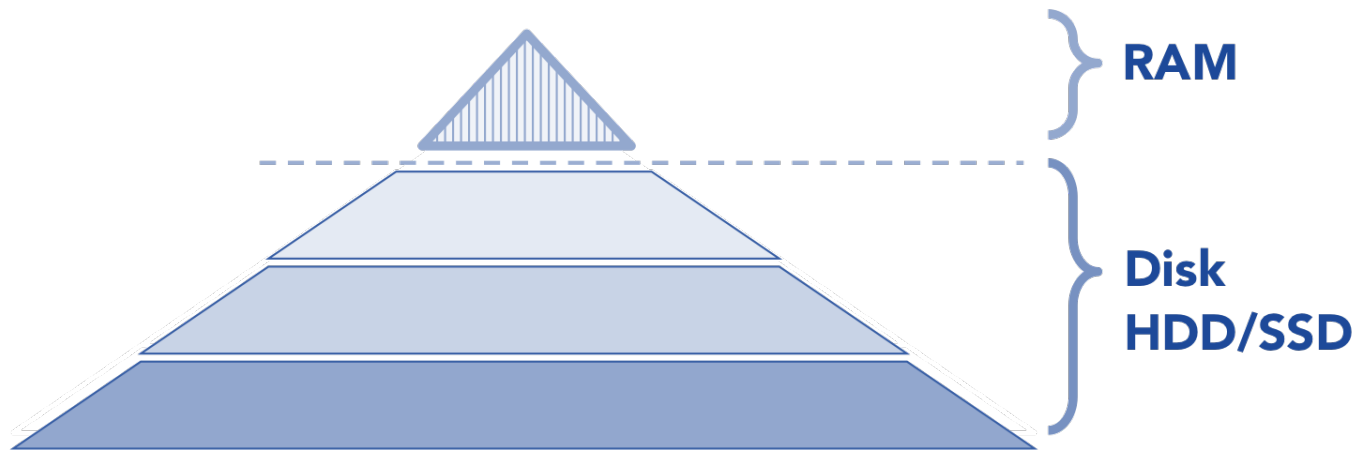Methods:   Reducing background I/O in LSMs.

# TRIAD in a Nutshell
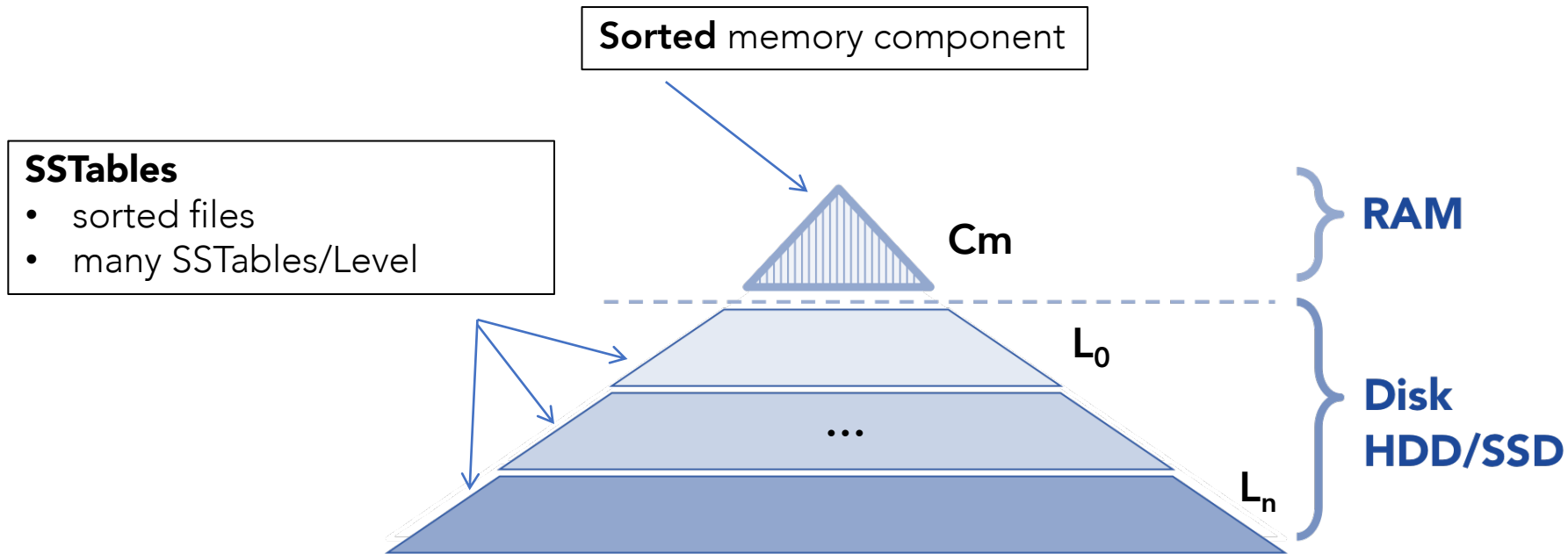
TRIAD LSM KV: achieves 2x throughput on production wklds.

Methods:   Reducing background I/O in LSMs.

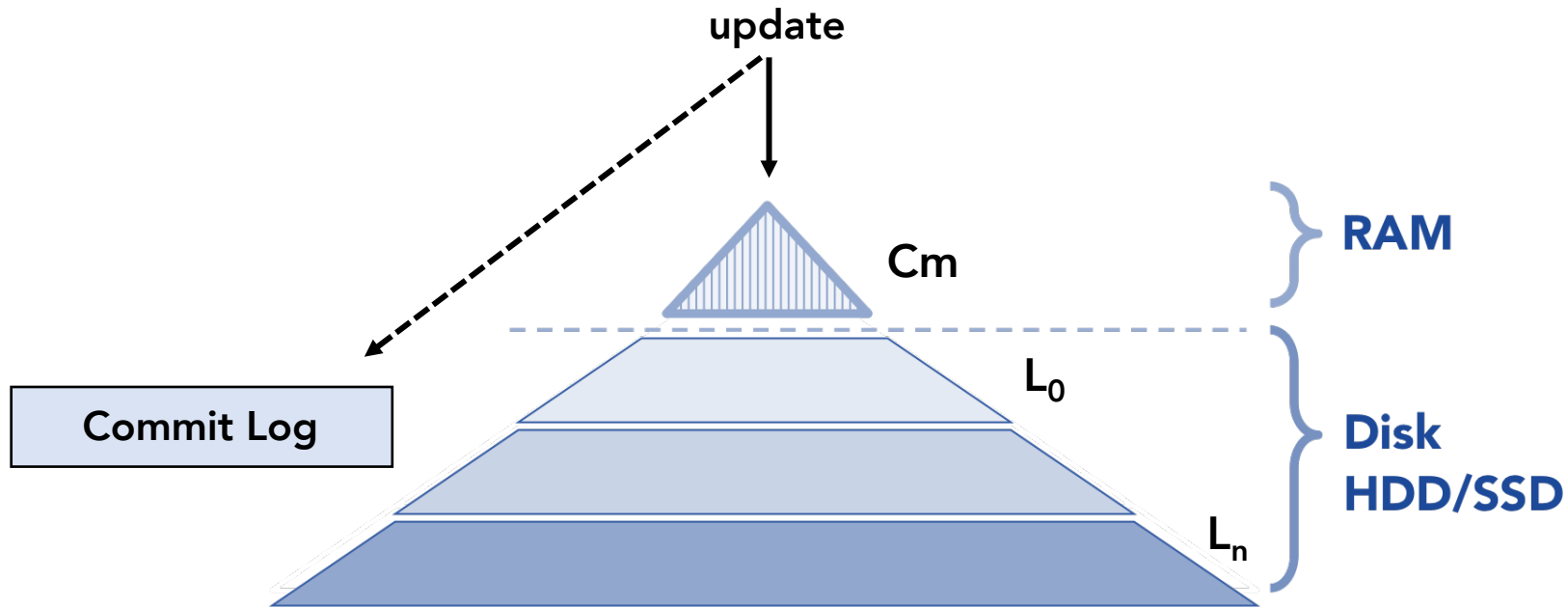☺   No need to know workload a priori.
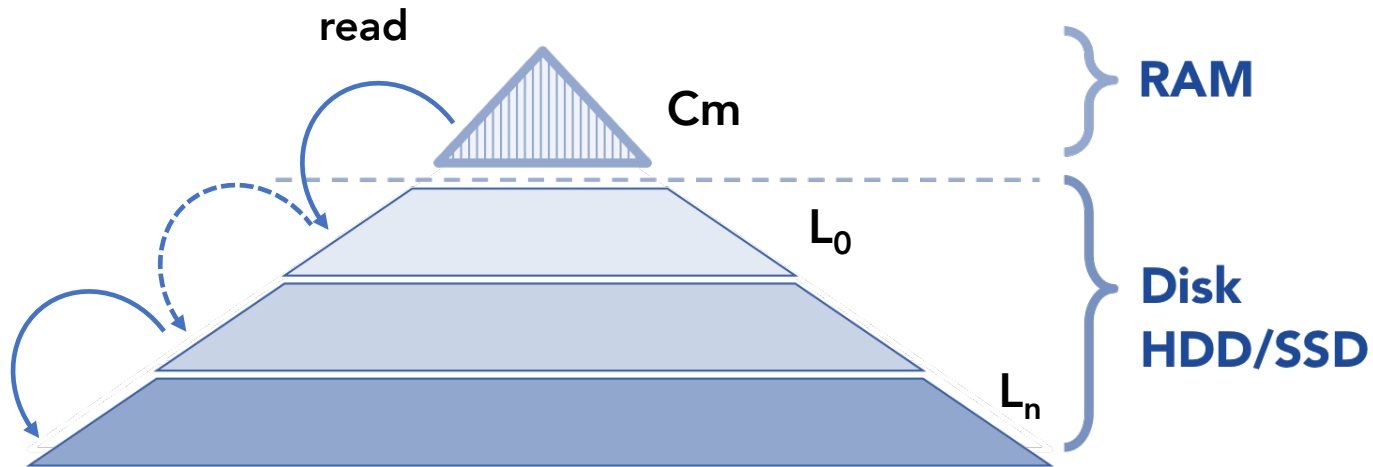
☺   LSM KV semantics preserved.

# LSM Overview



RAM

Disk
HDD/SSD

# LSM Components

**Sorted** memory component

**SSTables**
- sorted files
- many SSTables/Level

$C_m$

**RAM**

$L_0$

...

**Disk HDD/SSD**

$L_n$

# LSM Updates

update

Cm

RAM

L₀

Commit Log

Disk
HDD/SSD

Lₙ

# LSM Reads

# LSM Background Ops: Flushing



RAM

Cm

Flushing
From memory to L0.

$L_0$

Disk
HDD/SSD

$L_n$

# LSM Background Ops: Flushing



Mem component full

Cm

RAM

Disk          flushing

L0

| Commit Log | |
|---|---|
| K1 | V1 |
| K2 | V2 |
| K1 | V1' |
| K3 | V3 |
|  |  |
|  |  |
|  |  |

The memory component (Cm) contains:

| K1 | V1' |
|---|---|
| K2 | V2 |
| K3 | V3' |
| ... |  |
| Kn | Vn |

# LSM Background Ops: Flushing



Cm

RAM

Disk          flushing

| K1 | V1' |
| K2 | V2 |
| ... | |
| Kn | Vn |

flush$_1$

L0

**Commit Log**

# LSM Background Ops: Flushing



**Cm**

**RAM**

**Disk**          flushing

**L0**

**Commit log full**

**Commit Log**

| K1 | V1 |
|----|-----|
| K2 | V2 |
| K1 | V1'' |
| K3 | V3 |
| ... | |
| K1 | V1' |
| K3 | V3' |

Memtable (Cm):

| K1 | V1' |
|----|-----|
| K2 | V2 |
| K3 | V3' |
| | |
| | |

# LSM Background Ops: Flushing



**Cm**

**RAM**

**Disk**        flushing

**L0**

flush₁

**Commit Log**

# LSM Background Ops: Compaction



RAM

Cm

**Flushing**
From memory to L0.

$L_0$

**Compaction**
Levels on disk.

Disk
HDD/SSD

$L_n$

# LSM Background Ops: Compaction

# LSM Background Ops: Compaction
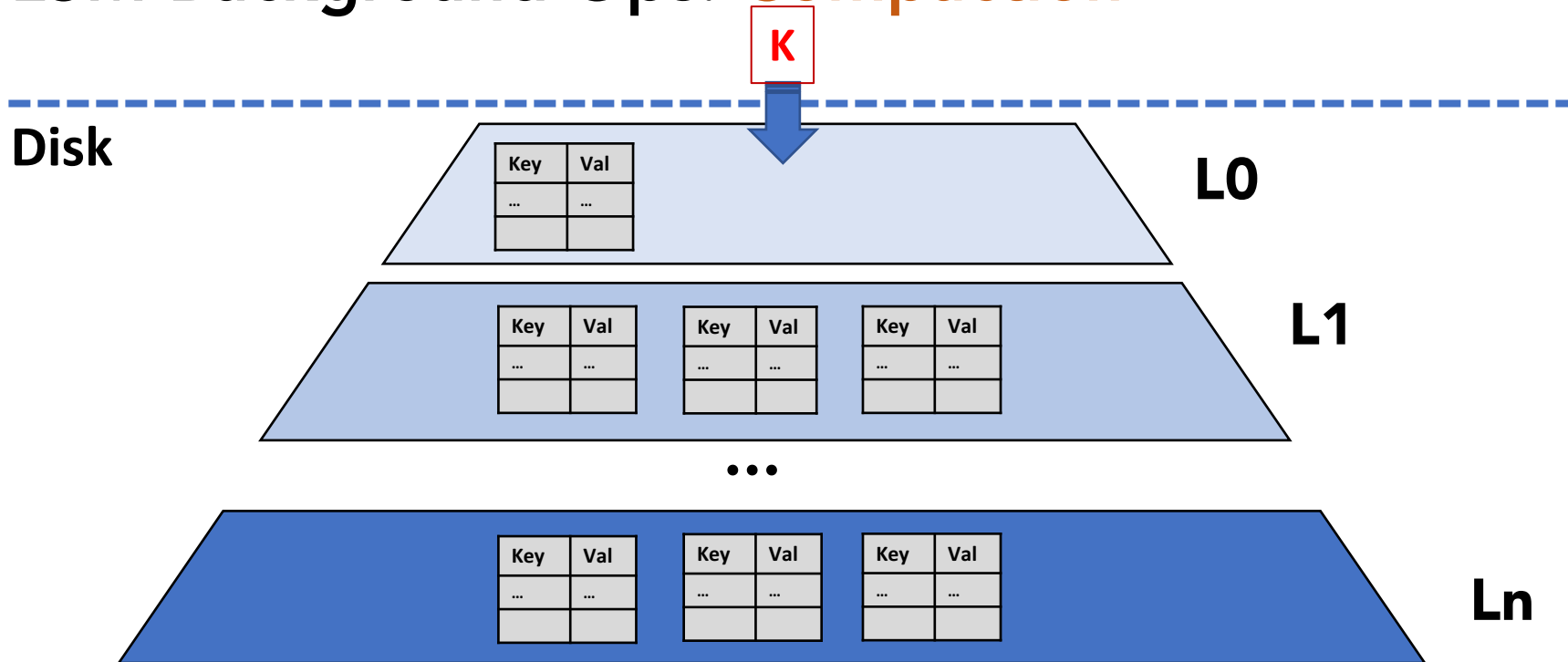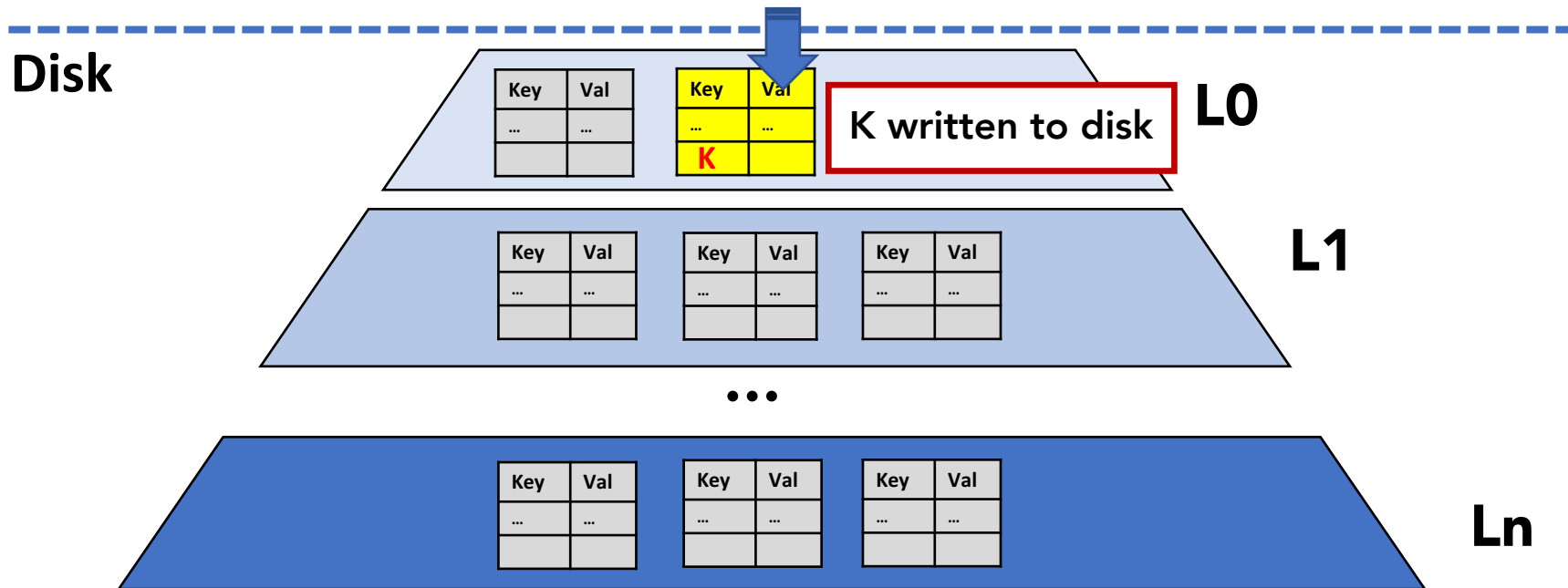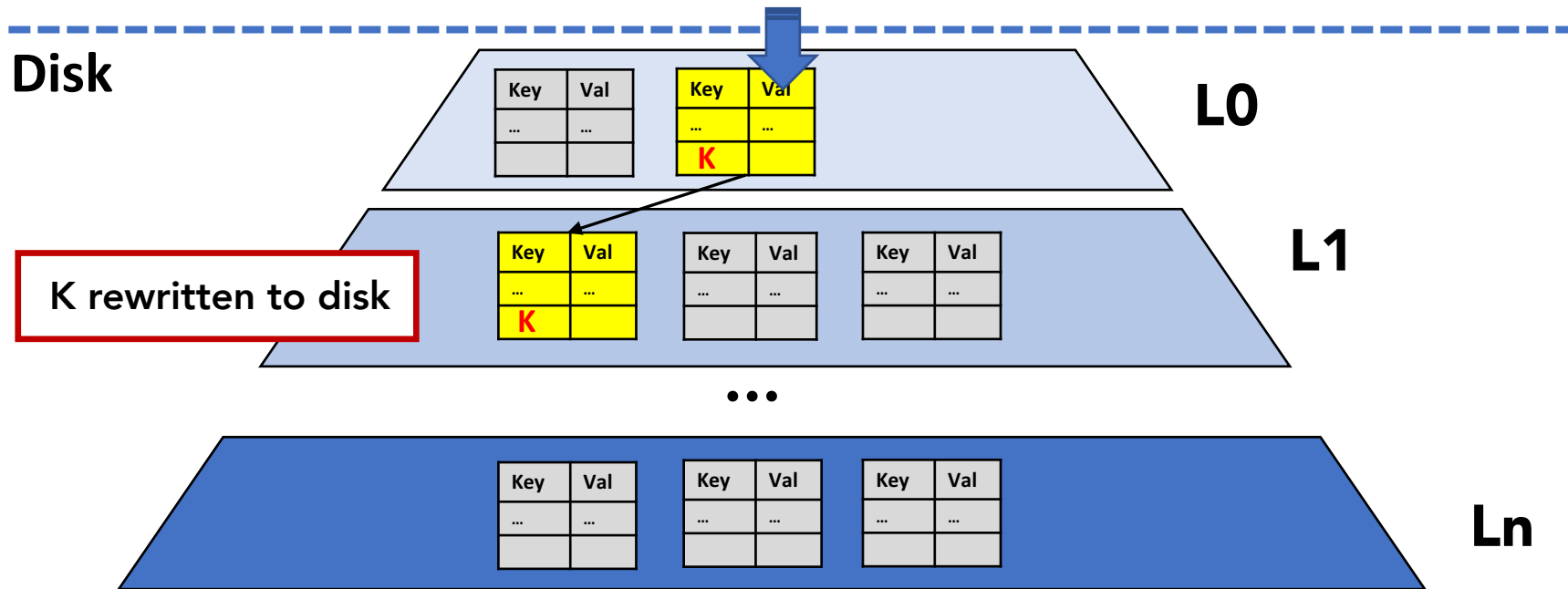
# LSM Background Ops: Compaction

# LSM Background Ops: Compaction

# LSM Background Ops: Compaction

# LSM Background Ops: Compaction



**Disk**

**L0**

Write amplification (WA)

≈ amount of **rewrites** of data to disk

K rewritten $n+1^{th}$ time to disk!

**Ln**

# Insight

Severe competition for compute/storage resources between LSM background ops and user ops.

# Background I/O Overhead

- **Long & slow bg. ops $\longrightarrow$ slowdown of user ops**.

# Background I/O Overhead

- **Long & slow bg. ops $\longrightarrow$ slowdown of user ops.**



Chart: K Operations/s

- RocksDB
- RocksDB No BG I/O

**up to 3x throughput gap** ☹

Uniform 50r-50w

Skewed 50r-50w

# Goal

**Decrease background ops overhead
to increase user throughput.**

# TRIAD

# TRIAD

| | Workload | Improve WA in |
|---|---|---|
| **TRIAD-MEM** | Skewed workloads | Flushing and Compaction |
| **TRIAD-DISK** | In-between | Compaction |
| **TRIAD-LOG** | Uniform workloads | Flushing |

**Three techniques work together and are complementary.**

# TRIAD

|  | Workload | Improve WA in |
|---|---|---|
| **TRIAD-MEM** | Skewed workloads | Flushing and Compaction |
| **TRIAD-LOG** | Uniform workloads | Flushing |

# TRIAD-MEM

|  | Workload | Improve WA in |
|---|---|---|
| **TRIAD-MEM** | Skewed workloads | Flushing and Compaction |

# Problem: Flushing with Skewed Workloads



Cm

RAM

Disk    flushing

Commit Log

L0

# Problem: Flushing with Skewed Workloads



RAM

Disk

flushing

**Cm**

**L0**

**Commit Log**

| K1 | V1$_1$ |
|----|--------|
|    |        |
|    |        |
|    |        |
|    |        |
|    |        |
|    |        |

33

# Problem: Flushing with Skewed Workloads



**Cm**

**RAM**

**Disk**         flushing

**Commit Log**

| K1 | V1$_1$ |
|----|--------|
| K1 | V1$_2$ |
|    |        |
|    |        |
|    |        |
|    |        |
|    |        |

**L0**

# Problem: Flushing with Skewed Workloads



**RAM**

**Disk**   flushing

**Cm**

**L0**

| K1 | V1$_2$ |
|----|--------|
| K2 | V2 |
| | |
| | |
| | |

**Commit Log**

| K1 | V1$_1$ |
|----|--------|
| K1 | V1$_2$ |
| K2 | V2 |
| | |
| | |
| | |
| | |

# Problem: Flushing with Skewed Workloads

**Cm**

**RAM**

**Disk**

flushing

**L0**

| K1 | V1$_3$ |
|----|--------|
| K2 | V2 |
|  |  |
|  |  |
|  |  |

**Commit Log**

| K1 | V1$_1$ |
|----|--------|
| K1 | V1$_2$ |
| K2 | V2 |
| K1 | V1$_3$ |
|  |  |
|  |  |
|  |  |

# Problem: Flushing with Skewed Workloads



RAM

Disk

flushing

Cm

L0

**Commit Log**

| K1 | $V1_1$ |
|----|--------|
| K1 | $V1_2$ |
| K2 | V2 |
| K1 | $V1_3$ |
| K1 | $V1_4$ |
| | |

# Problem: Flushing with Skewed Workloads

**RAM**

**Disk**

**Cm**

| K1 | $V1_n$ |
|----|--------|
| K2 | V2 |
| | |
| | |
| | |

flushing

**L0**

**Commit Log**

| K1 | $V1_1$ |
|----|--------|
| K1 | $V1_2$ |
| K2 | V2 |
| K1 | $V1_3$ |
| K1 | $V1_4$ |
| ... | |
| K1 | $V1_n$ |

# Problem: Flushing with Skewed Workloads



**Cm**

**RAM**

**Disk**

**flushing**

**Commit Log**

**L0**

K1  V1$_n$

K2  V2

flush$_1$

# Problem: Flushing with Skewed Workloads



**Cm**

**RAM**

**Disk**

flushing

**L0**

**High data skew**

- Flush because commit log is full
- Flush mostly empty mem comp

**Commit Log**

| | |
|---|---|
| K1 | $V1_1$ |
| K2 | V2 |
| K1 | $V1_2$ |
| K1 | $V1_3$ |
| K1 | $V1_4$ |
| ... | |
| K1 | $V1_n$ |

# Problem: Compaction with Skewed Workloads

# Problem: Compaction with Skewed Workloads

# Problem: Compaction with Skewed Workloads



L0

L1

# Problem: Compaction with Skewed Workloads

# Problem: Compaction with Skewed Workloads



L0

L1

Rewritten to disk

# Problem: Compaction with Skewed Workloads

File on L1 **rewritten to disk twice** because of one key ☹

Rewritten to disk

**L1**

# TRIAD-MEM: Hot-cold key separation



**Cm**

**RAM**

**Disk**    flushing

**L0**

**Commit Log**

| K1 | $V1_1$ |
|----|--------|
| K2 | V2 |
| K1 | $V1_2$ |
| K1 | $V1_3$ |
| K1 | $V1_4$ |
| ... | |
| K1 | $V1_n$ |

# TRIAD-MEM: Hot-cold key separation

**Idea:**

Keep **hot keys in memory**

**Flush** only **cold keys**

Keep **hot keys in CL**

| K1 | $V1_n$ |
|----|--------|
| K2 | V2 |
| K3 | V3 |
| ... | |
| Kn | Vn |

**RAM**

**Disk**

**flushing**

**L0**

**Commit Log**

| K1 | $V1_1$ |
|----|--------|
| K2 | V2 |
| K1 | $V1_2$ |
| K1 | $V1_3$ |
| K1 | $V1_4$ |
| ... | |
| K1 | $V1_n$ |

48

# TRIAD-MEM: Hot-cold key separation

**Idea:**

Keep **hot keys in memory**

**Flush** only **cold keys**

Keep **hot keys in CL**

CM

| K1 | V1$_n$ |
|----|--------|
|    |        |
|    |        |
|    |        |
|    |        |

**RAM**

**Disk**

flushing

L0

| K2 | V2 |
|----|----|
| K3 | V3 |
| ... |   |
| Kn | Vn |

**Commit Log**

| K1 | V1$_n$ |
|----|--------|
|    |        |
|    |        |
|    |        |
|    |        |
|    |        |
|    |        |

# TRIAD-MEM Summary

✓ **Good for skewed workloads.**

✓ **Reduce flushing WA: less data written from memory to disk.**

✓ **Reduce compaction WA: avoid repeatedly compacting hot keys.**

# TRIAD-LOG

| | Workload | Improve WA in |
|---|---|---|
| **TRIAD-LOG** | Uniform workloads | Flushing |

# Problem: Flushing with Uniform Workloads



**Cm**

| Key | Val |
|-----|-----|
| K1 | V1' |
| K2 | V2 |
| ... | |
| Kn | Vn |

**RAM**

**Disk**

flushing

**L0**

**Commit Log**

| K1 | V1 |
|-----|-----|
| K2 | V2 |
| K1 | V1' |
| K3 | V3 |
| K3 | V3' |
| ... | |
| Kn | Vn |

# Problem: Flushing with Uniform Workloads

**Cm**

**RAM**

**Disk**     flushing

| Key | Val |
|-----|-----|

| Key | Val |
|-----|-----|
| K1 | V1' |
| K2 | V2 |
| ... | |
| Kn | Vn |

**flush**

**L0**

**Commit Log**

| K1 | V1 |
|----|-----|
| K2 | V2 |
| K1 | V1' |
| K3 | V3 |
| K3 | V3' |
| ... | |
| Kn | Vn |

# Problem: Flushing with Uniform Workloads



**Cm**

**RAM**

**Disk**        flushing        **Commit Log**

**L0**

**flush**

# Problem: **Flushing with Uniform Workloads**



**RAM**

**Disk**

**Insight:**
Flushed data already written **to commit log.**

**Idea:**
Use commit logs as SSTables. **Avoid bg I/O due to flushing.**

| Key | Val |
|-----|-----|
| K1  | V1' |
| K2  | V2  |
| ... |     |
| Kn  | Vn  |

**flush**

**L0**

| K1  | V1  |
|-----|-----|
| K2  | V2  |
| K1  | V1' |
| K3  | V3  |
| K3  | V3' |
| ... |     |
| Kn  | Vn  |



55

# TRIAD-LOG

| Key | Val | CL Index |
|-----|-----|----------|
| K1 | V1' | 3 |
| K2 | V2 | 2 |
| ... | | |
| Kn | Vn | n |

**Cm**

**RAM**

**Disk**   flushing

Point to most recent entry in CL.

**L0**

**Commit Log**

| K1 | V1 |
|----|----|
| K2 | V2 |
| K1 | V1' |
| K3 | V3 |
| K3 | V3' |
| ... | |
| Kn | Vn |

# TRIAD-LOG

| Key | Val | CL Index |
|-----|-----|----------|
| K1 | V1' | 3 |
| K2 | - | 2 |
| ... | | |
| Kn | Vn | n |

**Cm**

**RAM**

**Disk**

**flushing**

**L0**

**Commit Log**

| K1 | V1 |
|----|----|
| K2 | V2 |
| K1 | V1' |
| K3 | V3 |
| K3 | V3' |
| ... | |
| Kn | Vn |

# TRIAD-LOG

| Key | Val | CL Index |
|-----|-----|----------|
|     |     | K1: 3    |
|     |     | K2: 2    |
|     |     |          |
|     |     | Kn: n    |

**Cm**

| CL Index |
|----------|
| K1: 3    |
| K2: 2    |
|          |
| Kn: n    |

Keep index in memory for further reads.

**CL-SSTable**
Only flush CL Index from memory and couple it with the current Commit Log.

**flushing**

**Commit Log**

| CL Index | K1  | V1  |
|----------|-----|-----|
| K1: 3    | K2  | V2  |
| K2: 2    | K1  | V1' |
| ...      | ... |     |
| Kn: n    | Kn  | Vn  |

**L0**

# TRIAD-LOG Summary

- ✓ **Good for uniform workloads.**

- ✓ **Reuse Commit Log as L0 SST.**

- ✓ **No more flushing of mem component to disk.**

# TRIAD Summary

**TRIAD-<span style="color:red">MEM</span>, TRIAD-<span style="color:blue">DISK</span>, TRIAD-<span style="color:green">LOG</span>:**

- **Complementary**, targeting different wklds.

- Working **simultaneously**.

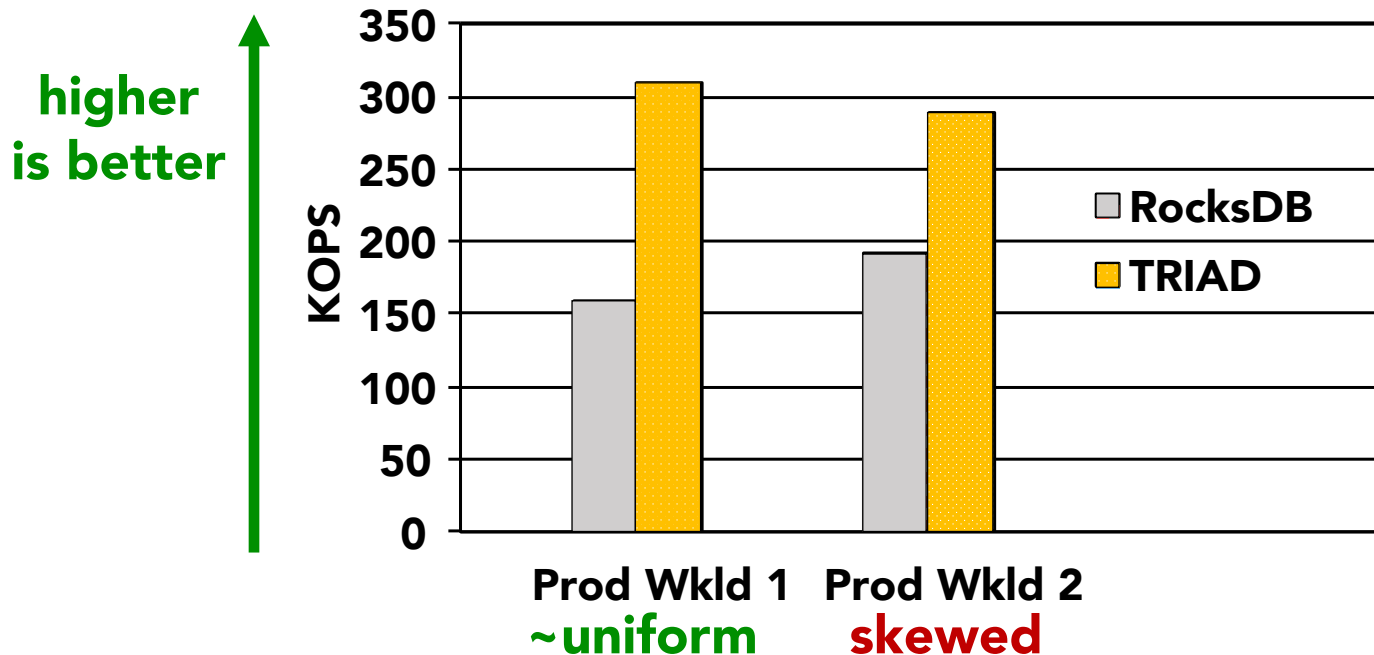- **Transparent** to the workloads; no a priori knowledge needed.

# Evaluation

# Evaluation

- **Compare TRIAD with RocksDB**

- **Workloads:** Production, Synthetic

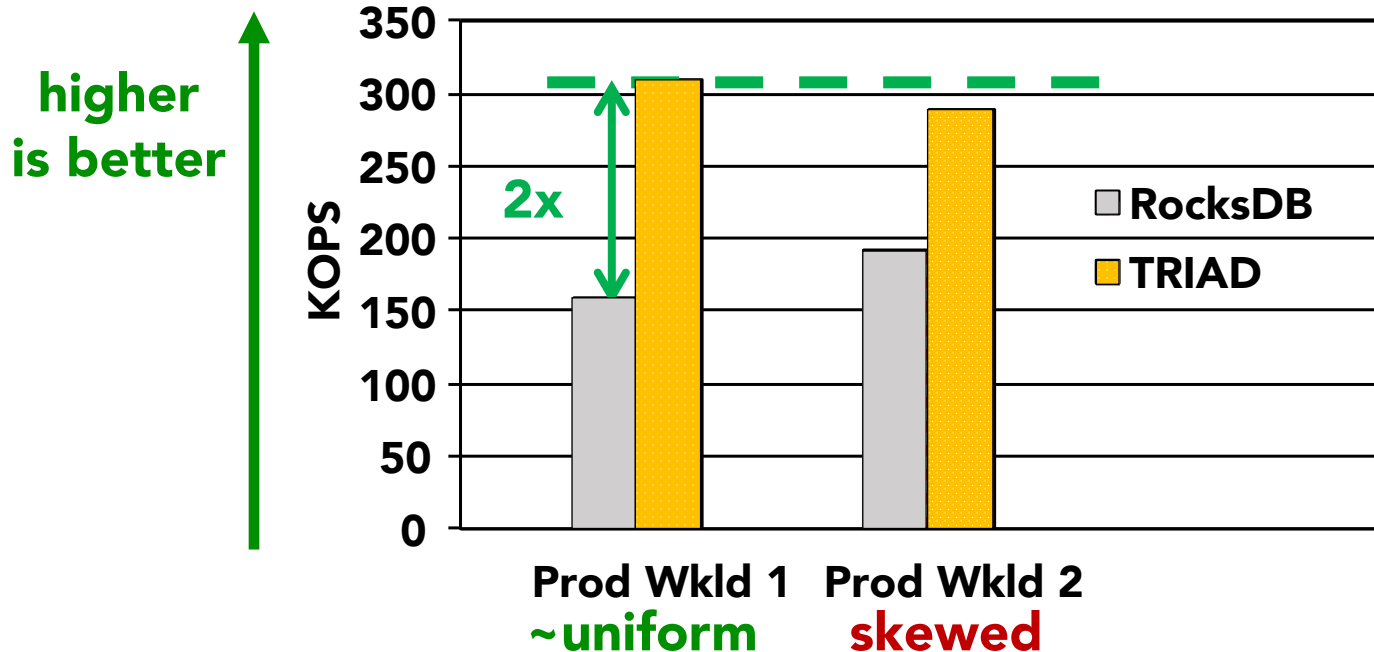- **Metrics:** Throughput, Write Amplification (WA)

- **Code:** https://github.com/epfl-labos/TRIAD

# Write Amplification (WA)

$$WA = \frac{\text{total data written to storage}}{\text{data written by app}}$$
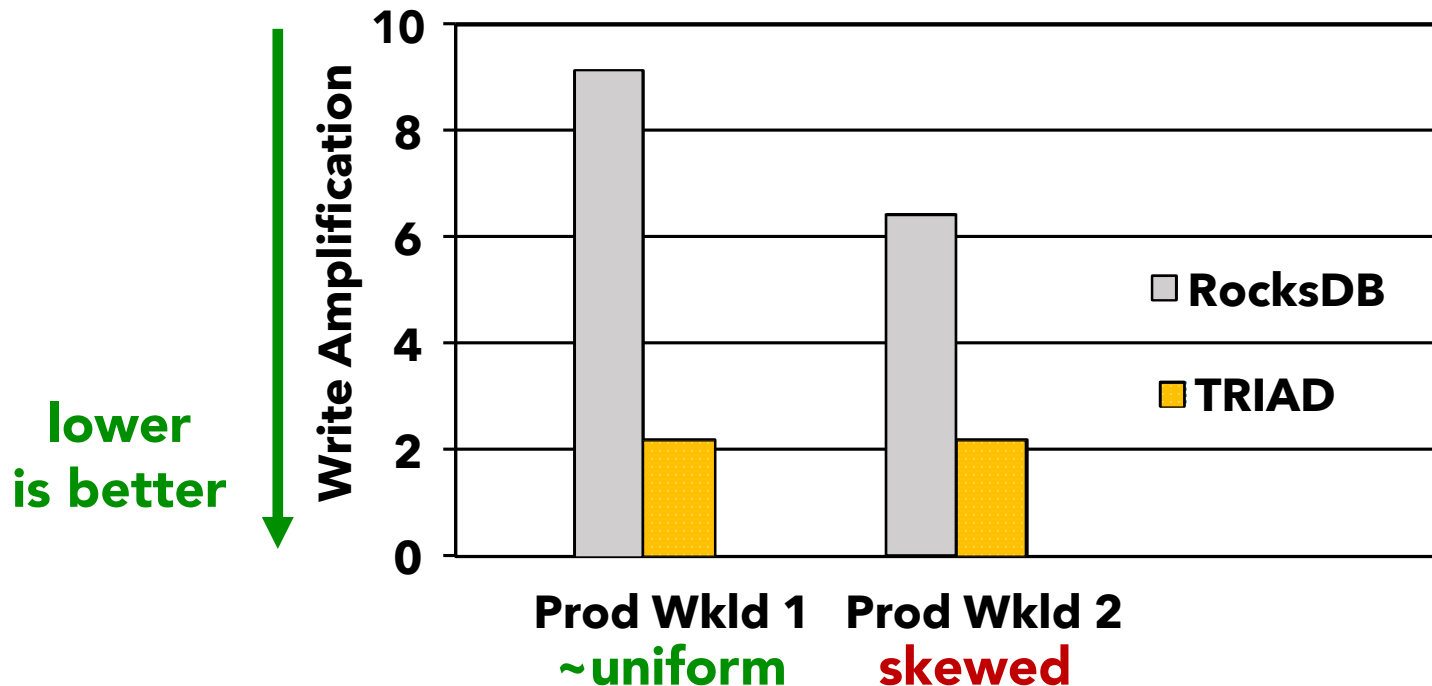
# Production Workloads: Throughput



higher
is better

~uniform   skewed

RocksDb
TRIAD

10

Write Amplification

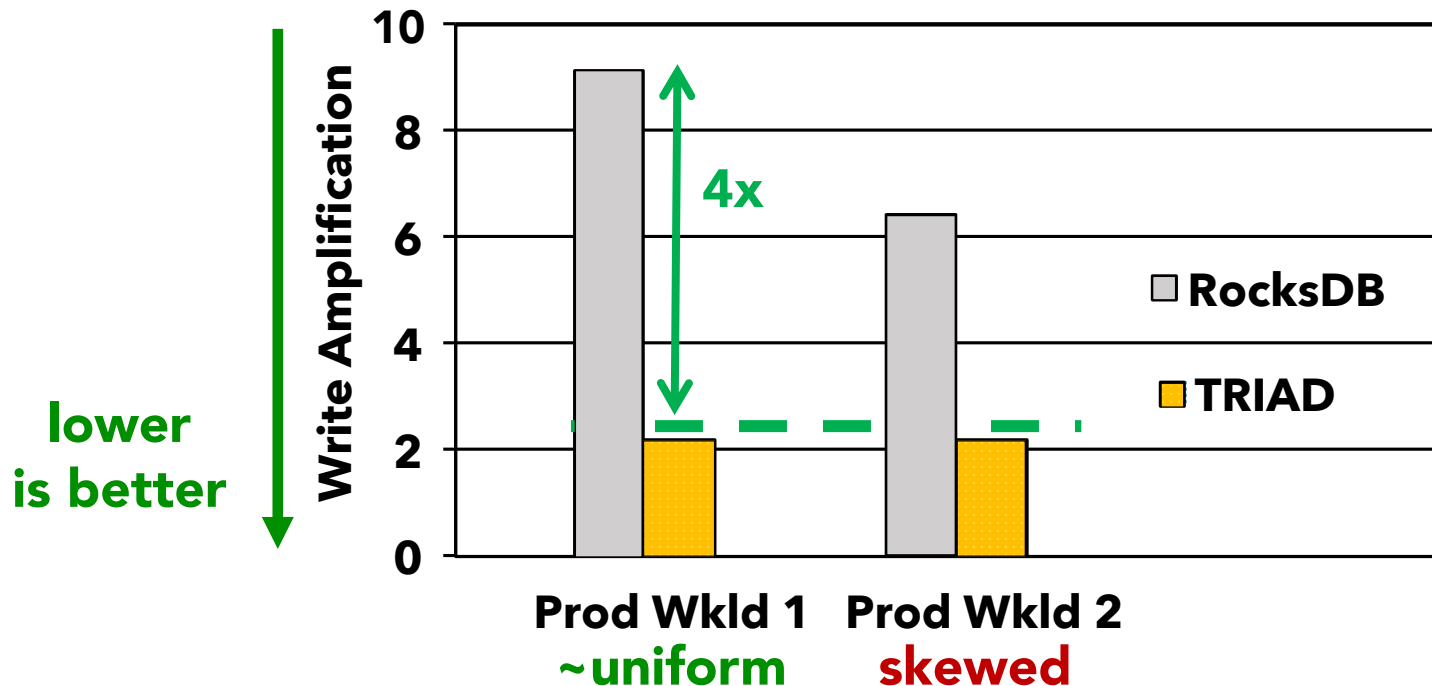Prod Wl

# Production Workloads: Throughput

# Production Workloads: Write Amplification

# Production Workloads: Write Amplification



TRIAD: low and <u>uniform</u> WA.

lower is better

RocksDB

TRIAD

4x

Prod Wkld 1
~uniform

Prod Wkld 2
skewed

# TRIA... etic Workloads



higher is better

**Skew 1%-99%**

KOPS

| TRIAD | TRIAD-MEM | TRIAD-LOG | RocksDB |

**No Skew**

KOPS

| TRIAD | TRIAD-MEM | TRIAD-LOG | RocksDB |

# TRIA... etic Workloads

# More in Our Paper

o **More production workloads**

o **More synthetic workloads**

o **Detailed breakdown of TRIAD techniques**

o **TRIAD-DISK**

# Related work

o **LevelDB**: first LSM-based KV store. No attempts to reduce WA.

# Related work

o **LevelDB**: first LSM-based KV store. No attempts to reduce WA.

o A number of systems attempt to **reduce WA.**
  - **LSMs:** RocksDB, *bLSM* (SIGMOD/PODS '12), *VT-tree* (FAST '13), *HyperLevelDB, LSM-trie* (USENIX ATC '15), Cassandra, WiscKey (FAST '16).

# Related work

o **LevelDB**: first LSM-based KV store. No attempts to reduce WA.

o A number of systems attempt to **reduce WA.**

- **LSMs:** RocksDB, *bLSM* (SIGMOD/PODS '12), *VT-tree* (FAST '13), *HyperLevelDB, LSM-trie* (USENIX ATC '15), Cassandra, WiscKey (FAST '16).

- **B-epsilon trees:** Tucana (USENIX ATC '16), BetrFS (FAST '15, '16)

# Related work

o **LevelDB**: first LSM-based KV store. No attempts to reduce WA.

o A number of systems attempt to **reduce WA.**

- **LSMs:** RocksDB, *bLSM* (SIGMOD/PODS '12), *VT-tree* (FAST '13), *HyperLevelDB, LSM-trie* (USENIX ATC '15), Cassandra, WiscKey (FAST '16).

- **B-epsilon trees:** Tucana (USENIX ATC '16), BetrFS (FAST '15, '16)

- But no hot/cold key separation, not use commit log as a pseudo-SST.

# Related work

o **LevelDB**: first LSM-based KV store. No attempts to reduce WA.

o A number of systems attempt to **reduce WA.**
  - **LSMs:** RocksDB, *bLSM* (SIGMOD/PODS '12), *VT-tree* (FAST '13), *HyperLevelDB, LSM-trie* (USENIX ATC '15), Cassandra, WiscKey (FAST '16).

  - **B-epsilon trees:** Tucana (USENIX ATC '16), BetrFS (FAST '15, '16)

  - But no hot/cold key separation, not use commit log as a pseudo-SST.

o Hot/cold key separation idea used in different context.
  - **FLASH storage systems:** dual-pool algorithm (SAC '07), Application-Managed Flash (FAST '16)

# Take-home Messages

✓ **TRIAD: LSM key-value store with high throughput and low WA.**

# Take-home Messages

✓ **TRIAD: LSM key-value store with high throughput and low WA.**

✓ **Complementary techniques transparent to workload types.**

# Take-home Messages

✓ **TRIAD: LSM key-value store with high throughput and low WA.**

✓ **Complementary techniques transparent to workload types.**

✓ **Impact of LSM I/O on user throughput reduced.**