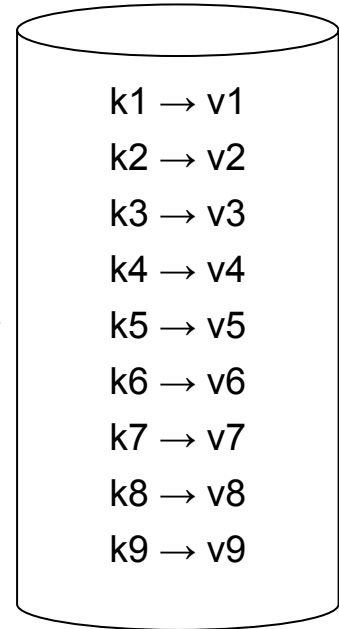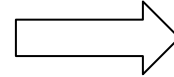# EvenDB: Optimizing Key-Value Storage for Spatial Locality

Eran Gilad, Edward Bortnikov, Anastasia Braginsky, Yonatan Gottesman, Eshcar Hillel (Yahoo Research), Idit Keidar (Technion), Nurit Moscovici (Outbrain), Rana Shahout (Technion)
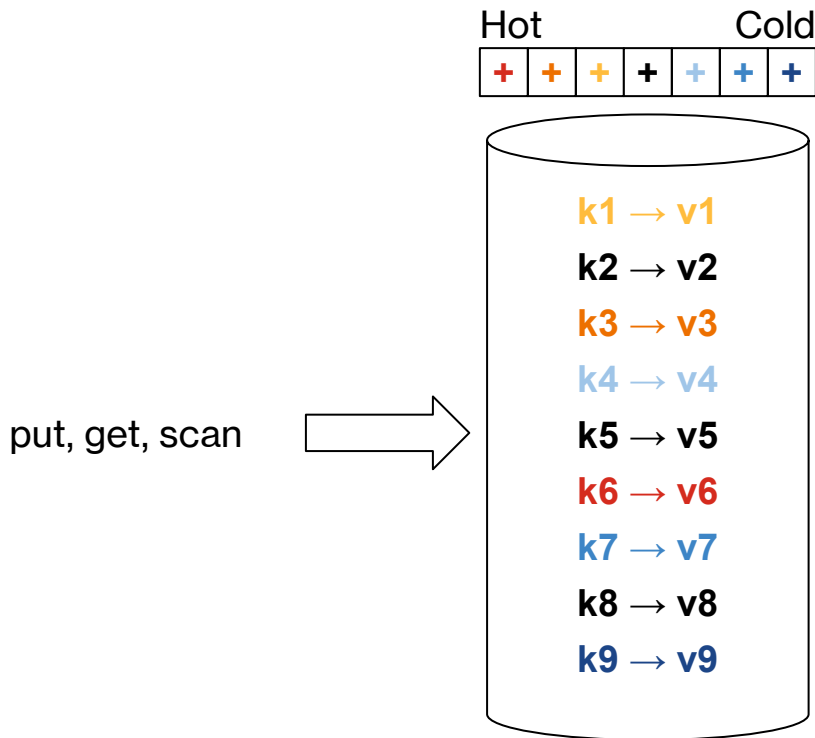
# Key-value stores

- **key -> value mapping**

put, get, scan →

k1 → v1
k2 → v2
k3 → v3
k4 → v4
k5 → v5
k6 → v6
k7 → v7
k8 → v8
k9 → v9

**verizon√ media**

**yahoo! research**

# Key-value stores

Hot                    Cold

- key -> value mapping
- **skewed workload: some *keys* are hotter**

put, get, scan ⟹

k1 → v1
k2 → v2
k3 → v3
k4 → v4
k5 → v5
k6 → v6
k7 → v7
k8 → v8
k9 → v9

verizon√
media

yahoo!
research

# Key-value stores

- key -> value mapping
- skewed workload: some *keys* are hotter
- **spatial locality: some *ranges* are hotter**
  - e.g., complex keys

put, get, scan ⟹

Hot                                    Cold
| + | + | + | + | + | + | + |

k1_l1 → v1
k1_l2 → v2
k1_l3 → v3
k2_l1 → v4
k2_l2 → v5
k3_l1 → v6
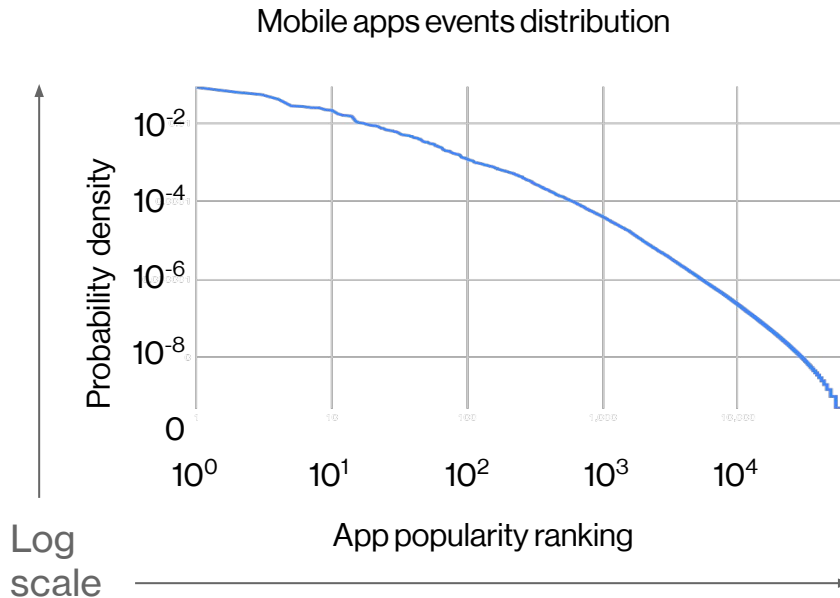k3_l2 → v7
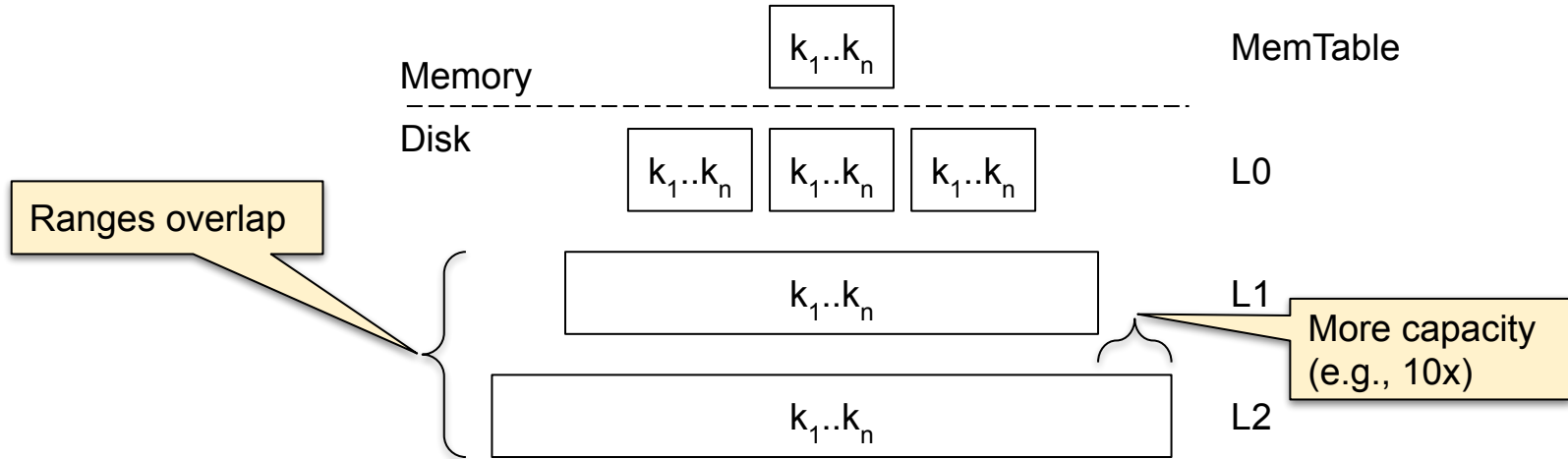k3_l3 → v8
k3_l4 → v9

yahoo!
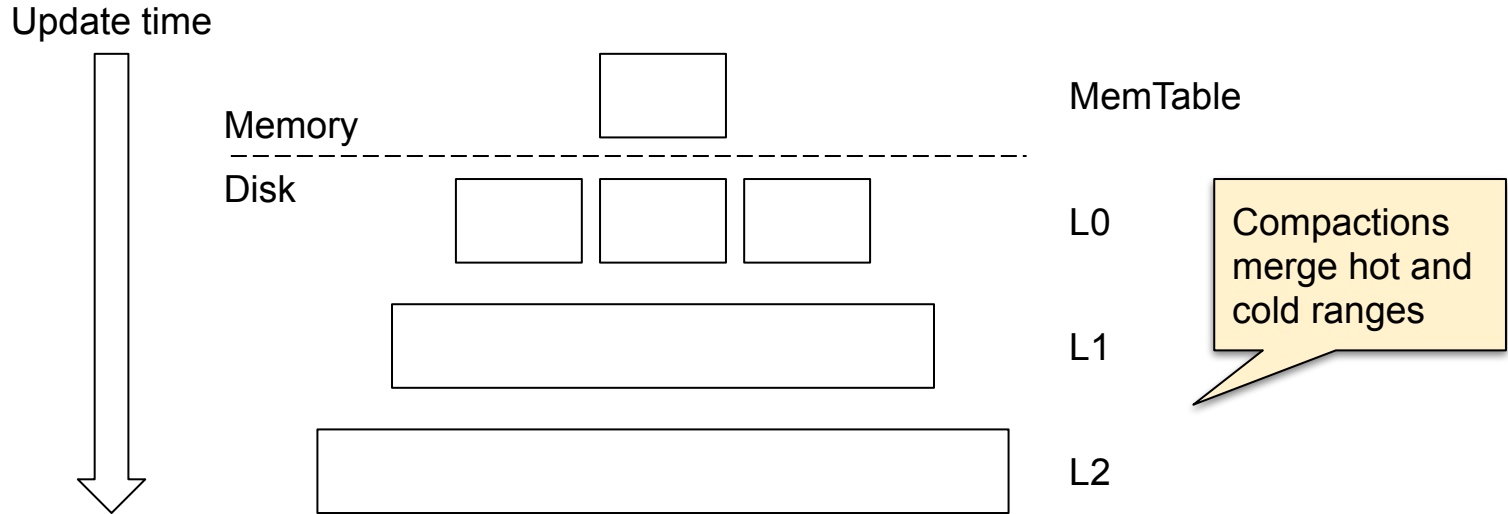research

# Key-value stores

- key -> value mapping
- skewed workload: some *keys* are hotter
- spatial locality: some *ranges* are hotter
  - e.g., complex keys
- **Sample production trace:**
  - appname_timestamp
  - 1% of apps ⇒ 1% key prefixes ⇒ 94% of events

Mobile apps events distribution



Probability density: $10^{-2}$, $10^{-4}$, $10^{-6}$, $10^{-8}$, 0

App popularity ranking: $10^0$, $10^1$, $10^2$, $10^3$, $10^4$

Log scale

# LSM-trees

Memory

$k_1..k_n$   MemTable

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

Disk

$k_1..k_n$   $k_1..k_n$   $k_1..k_n$   L0

Ranges overlap

$k_1..k_n$   L1

More capacity (e.g., 10x)

$k_1..k_n$   L2

verizon✓
media

**6**

yahoo!
research

# LSM-trees are designed for temporal locality

Update time

Memory
Disk

MemTable

L0

L1

L2

Compactions merge hot and cold ranges

verizon✓
media

yahoo!
research

# LSM-trees are less suited for spatial locality

*scan(...):*

MemTable

Memory

Disk

L0

Ranges are
fragmented

L1

L2

**verizon✓
media**
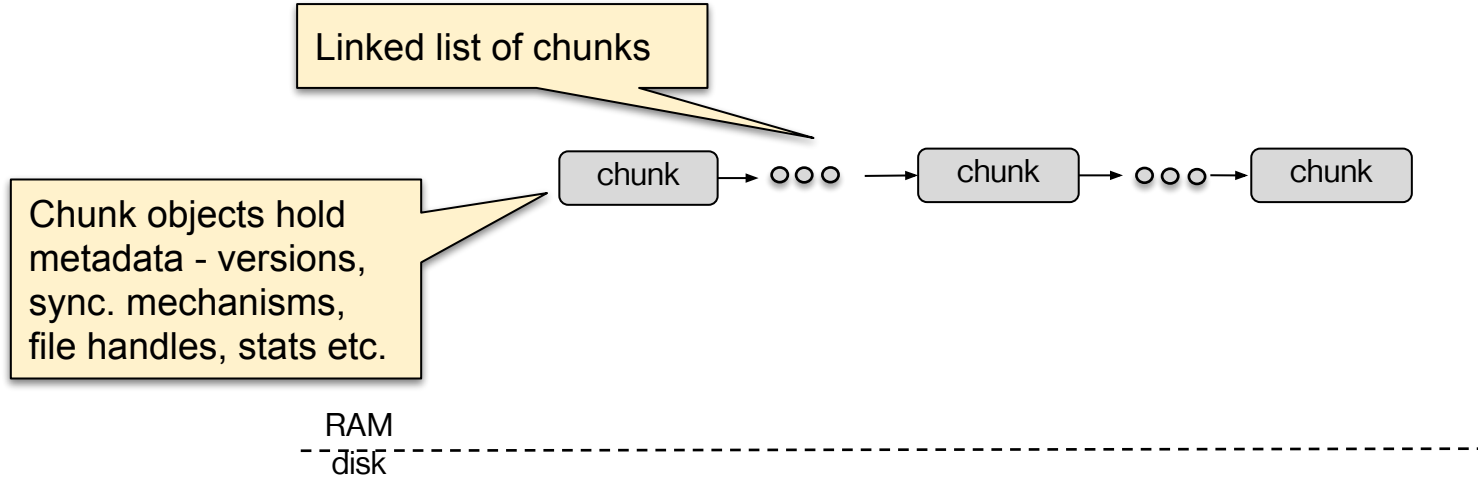
8

**yahoo!
research**

# EvenDB

- **Ordered key-value store**

- **Optimized for spatial locality**

- **Low write amplification**

- **Persistent, fast recovery**

- **Atomic operations, including scan**

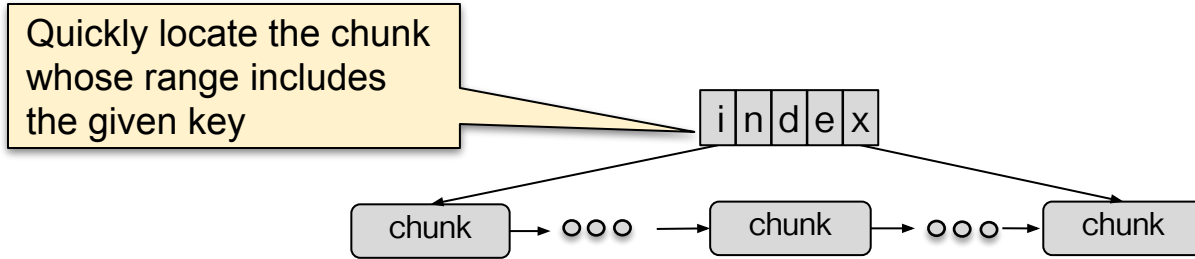**verizon**✓
**media**

yahoo!
research

# Chunk-based organization

- **Dynamically partitioned key space into *chunks***
  - Much smaller than shards
  - Much larger than blocks
- **Chunks are the basic unit for**
  - Disk I/O
  - Compaction
  - Memory caching
  - Concurrency control

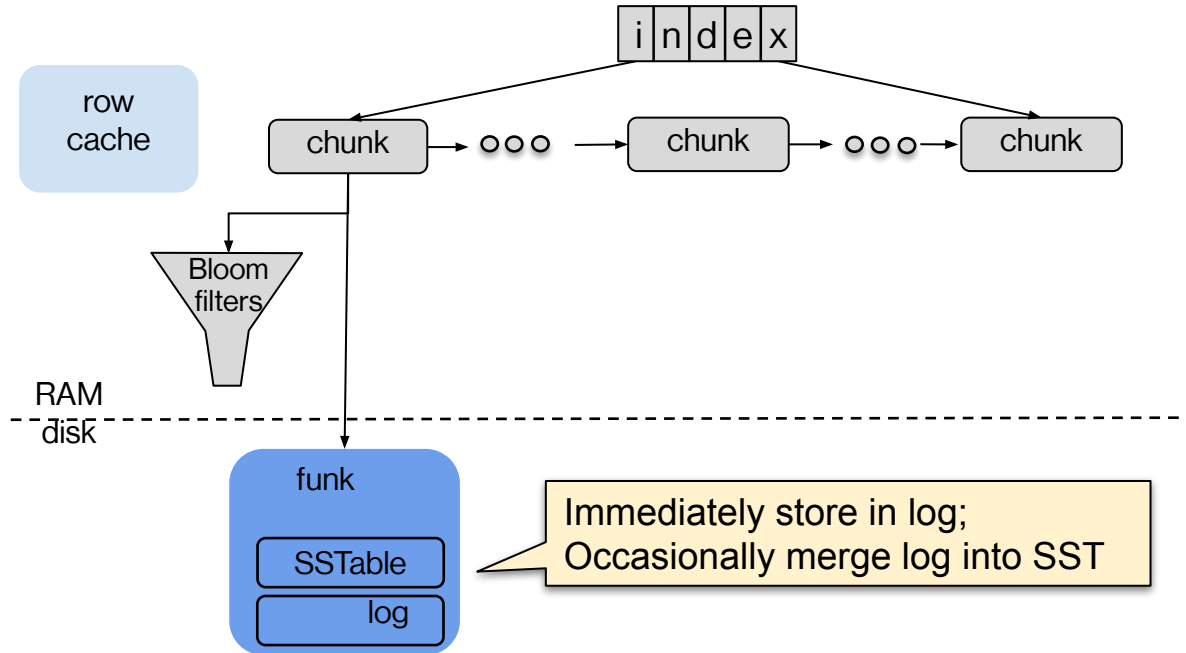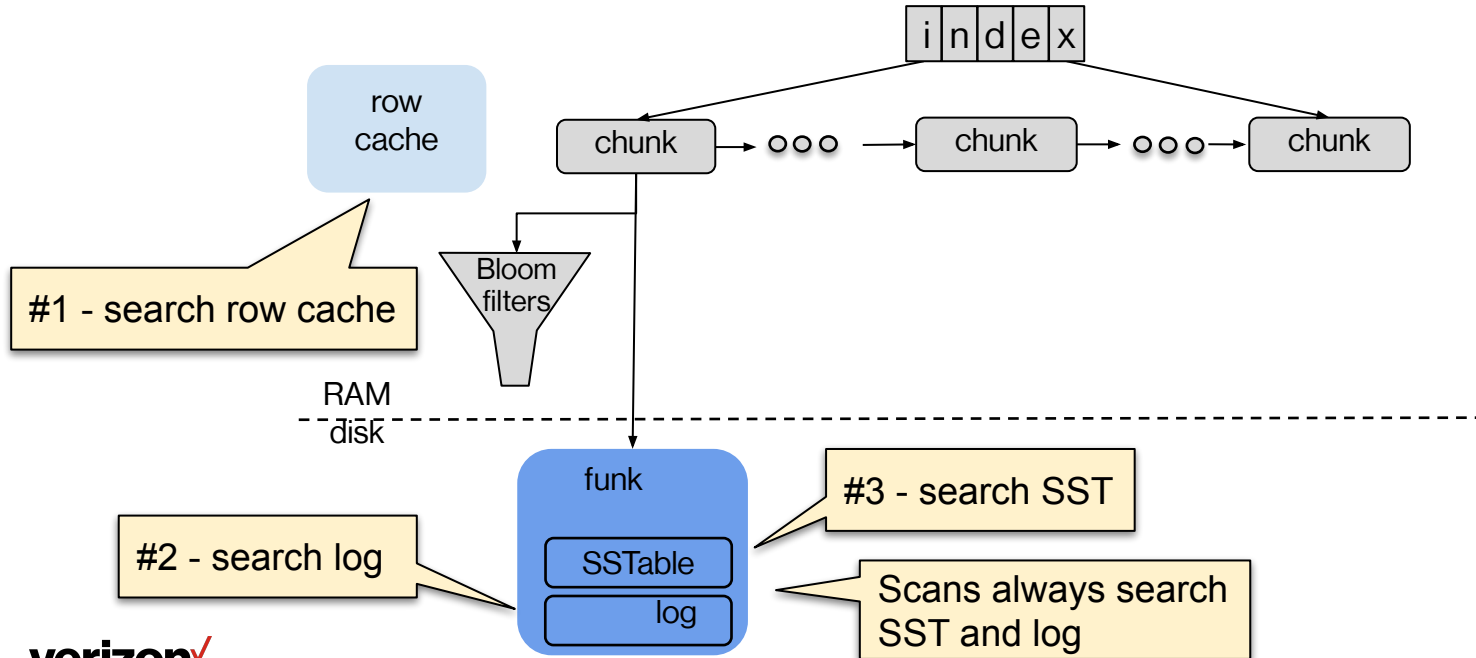**verizon√**
**media**

**yahoo!**
**research**

# Chunks metadata

Linked list of chunks

chunk → ○○○ → chunk → ○○○ → chunk

Chunk objects hold metadata - versions, sync. mechanisms, file handles, stats etc.

RAM
disk

# Chunks index

Quickly locate the chunk whose range includes the given key

i n d e x

chunk → ○○○ → chunk → ○○○ → chunk

RAM
disk

**verizon**√
**media**

yahoo!
research

# Disk storage - updates



Immediately store in log;
Occasionally merge log into SST

# Disk storage - lookups



index

row cache

chunk → ○○○ → chunk → ○○○ → chunk

#1 - search row cache

Bloom filters

RAM
disk

funk

#2 - search log

SSTable

log

#3 - search SST

Scans always search SST and log

verizon√
media

yahoo!
research

# Memory cache - updates



index

row cache

chunk ○○○ chunk ○○○ chunk

Bloom filters

munk

munk

munk cache

#2 - Store in munk

RAM
disk

#3 - Occasionally rebalance munk

funk

funk

funk

#4 - Rarely create SST from munk

SSTable

SSTable

SSTable

log

log

log

#1 - Store in log

verizon√
media

yahoo!
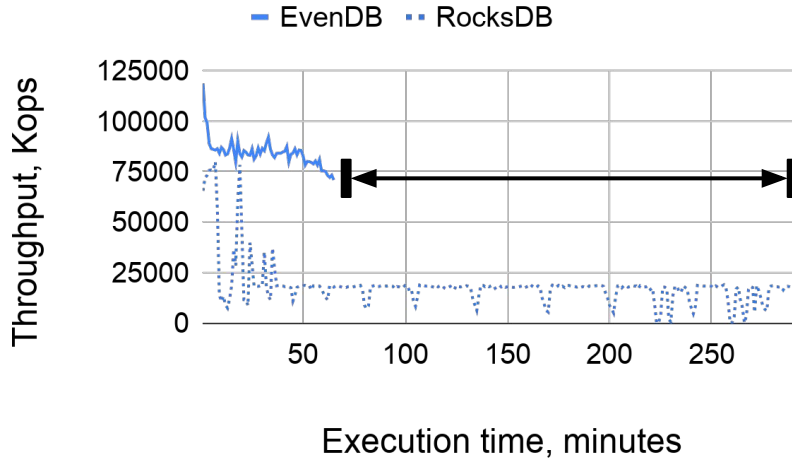research

15

# Memory cache - lookups

# Evaluation

- ## 3 benchmark suites
  - Traces from internal production system, 256GB DB - some presented next
  - Standard and extended YCSB benchmarks - results in paper

- ## State-of-the-art LSM: RocksDB

**verizon**√
**media**

yahoo!
research

# Real dataset ingestion

Throughput dynamics - 256GB DB creation

**EvenDB 4.4x faster,**

**write amp. 4x lower (better)**

verizon✓
media

yahoo!
research

# Compactions impact

Throughput dynamics - 256GB DB creation



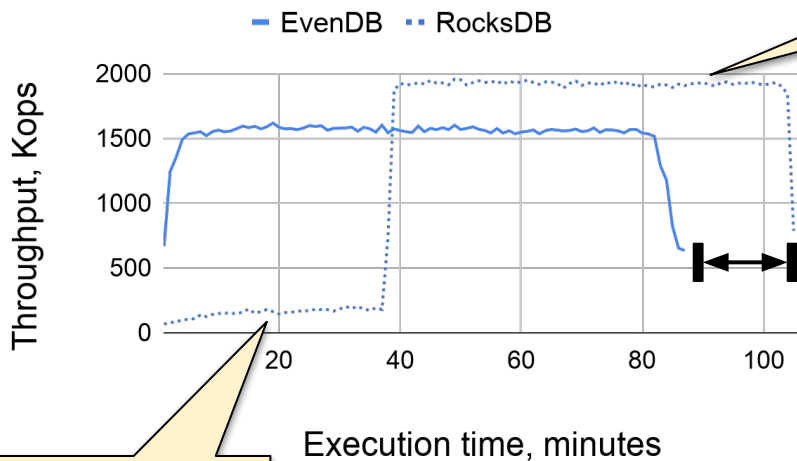Space amp.: DB size during ingestion



RocksDB throughput drops during compaction

EvenDB runs much smoother

# Real dataset scans

Scan throughput dynamics, 256GB



RocksDB faster after storage optimized

**EvenDB 1.2x faster than RocksDB**

~38 minutes stall after DB creation

verizon√
media

**20**

yahoo!
research

# Summary

Thank you! Qs?

- **EvenDB introduces a novel key-value store architecture**
- **Chunk arrangement better suited for spatially-local workloads than LSM:**
    - Lower write amplification
    - Single level of storage, no overlapping
    - Memory serves reads and writes
- **EvenDB outperforms RocksDB when:**
    - Workload is spatially-local or most working set fits in RAM
    - In par otherwise
    - Demonstrated in real workload and synthetic YCSB benchmarks

**verizon**√
**media**

yahoo!
research