# KVell: the Design and Implementation of a Fast Persistent Key-Value Store

Baptiste Lepers

Oana Balmau

Karan Gupta

Willy Zwaenepoel

THE UNIVERSITY OF SYDNEY

NUTANIX™

Made in Australia from 0% Australian ingredients
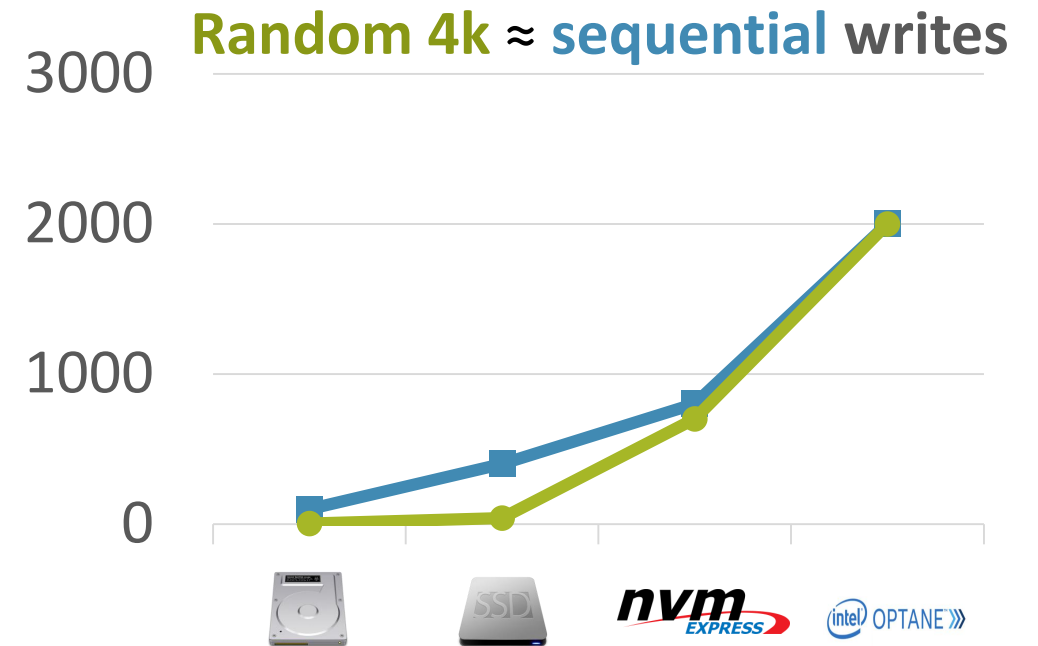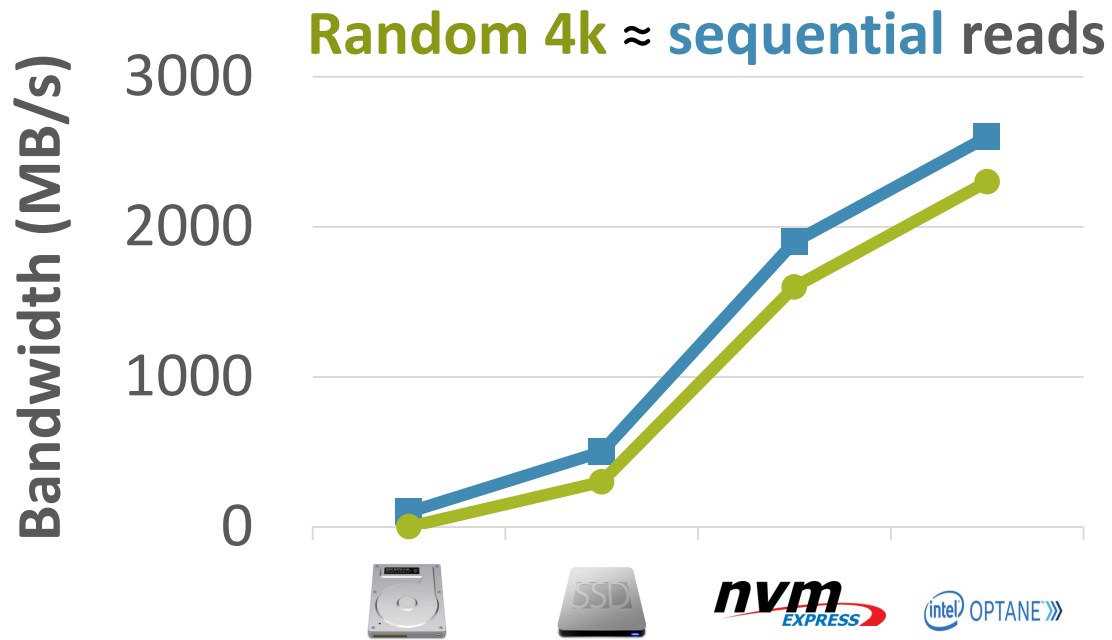
# Single Machine Persistent KVs

**Put**(k, v)

**Get**(k)$\rightarrow$ v

**Scan**($k_X$, $k_Y$) $\rightarrow$ [$k_X$ $v_X$ , ... , $k_Y$ $v_Y$]

# Disks are much faster

# Random as fast as sequential
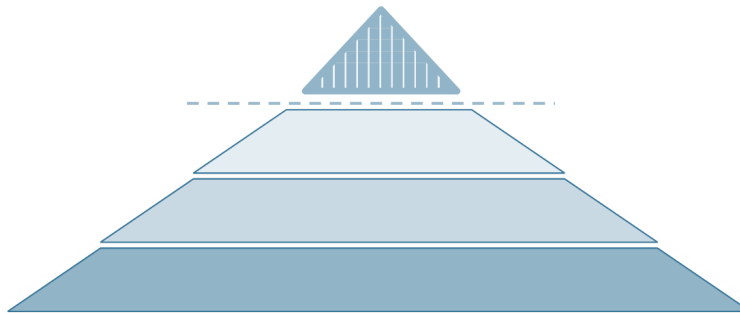
# This Talk

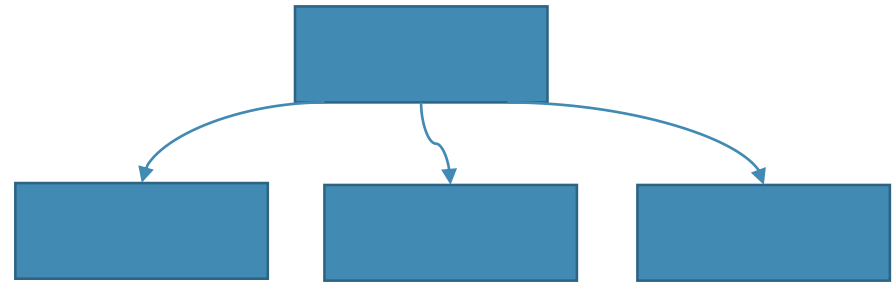**Existing KVs <span style="color:red">not designed for fast drives</span>**

**KVell: a new design for fast drives**
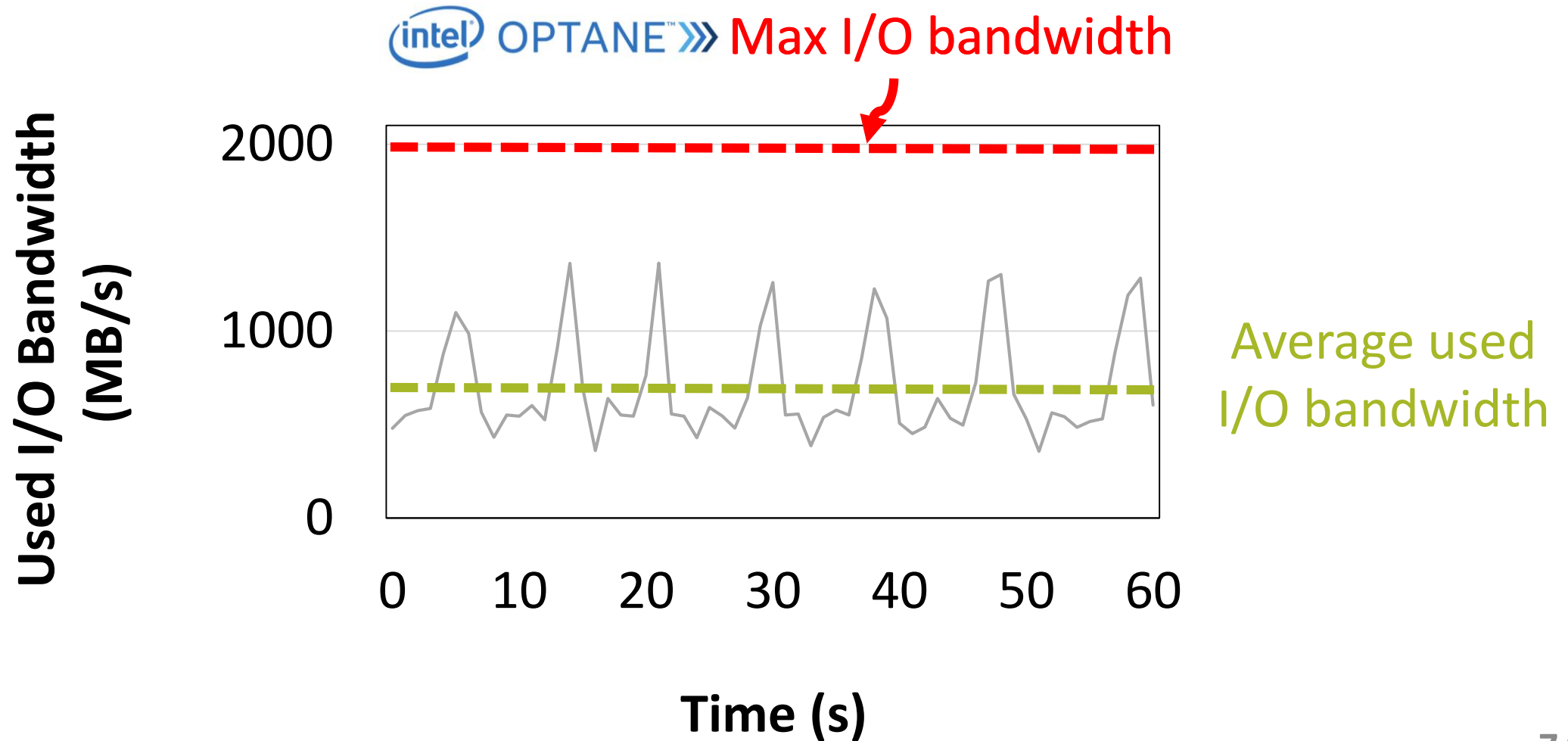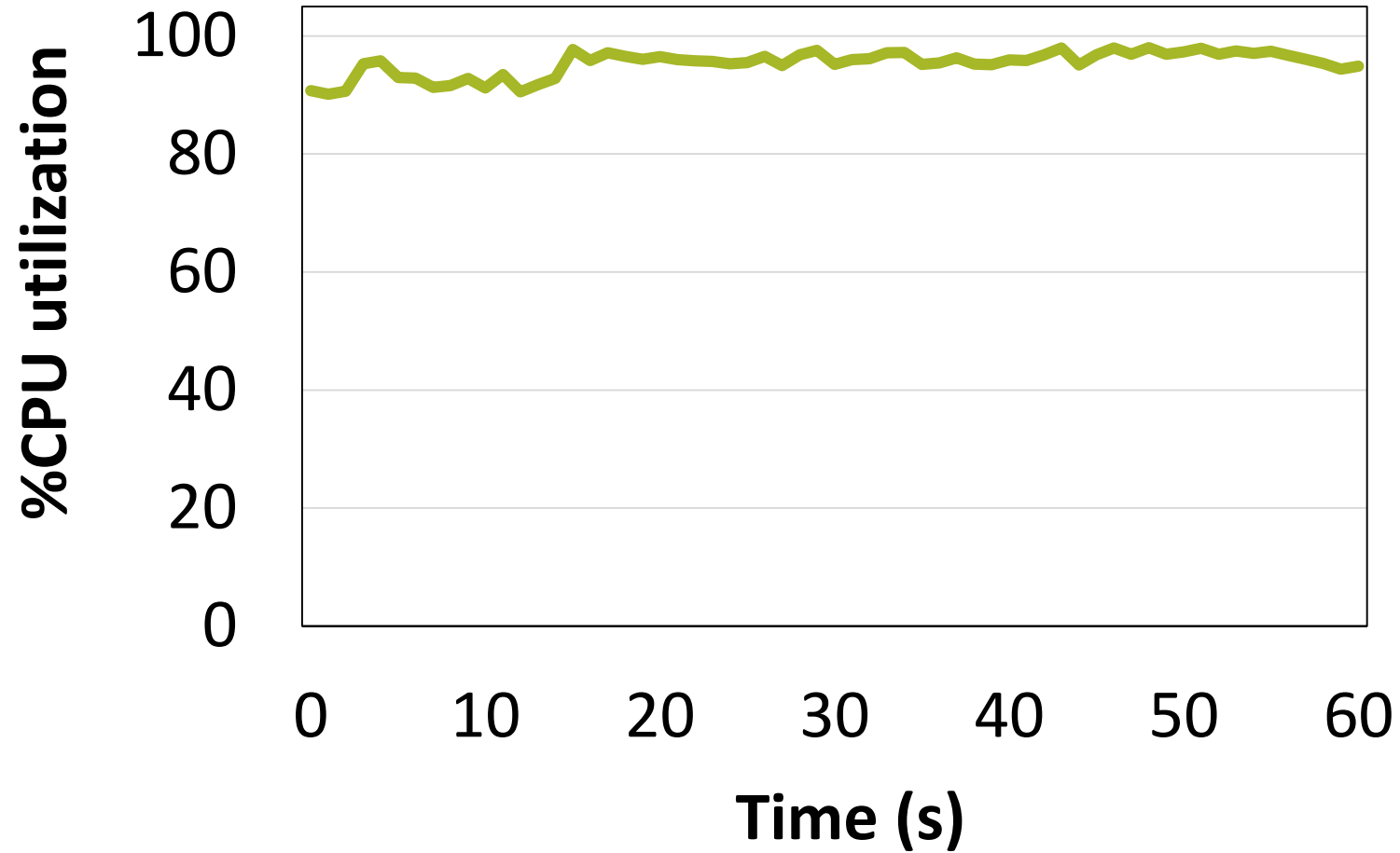
# Popular designs

## Log Structured Merge Tree (LSM)

## B+ Tree

# RocksDB 50% GET, 50% PUT



Used I/O Bandwidth (MB/s) vs Time (s). Intel Optane Max I/O bandwidth shown as red dashed line at ~2000 MB/s. Average used I/O bandwidth shown as green dashed line at ~700 MB/s.
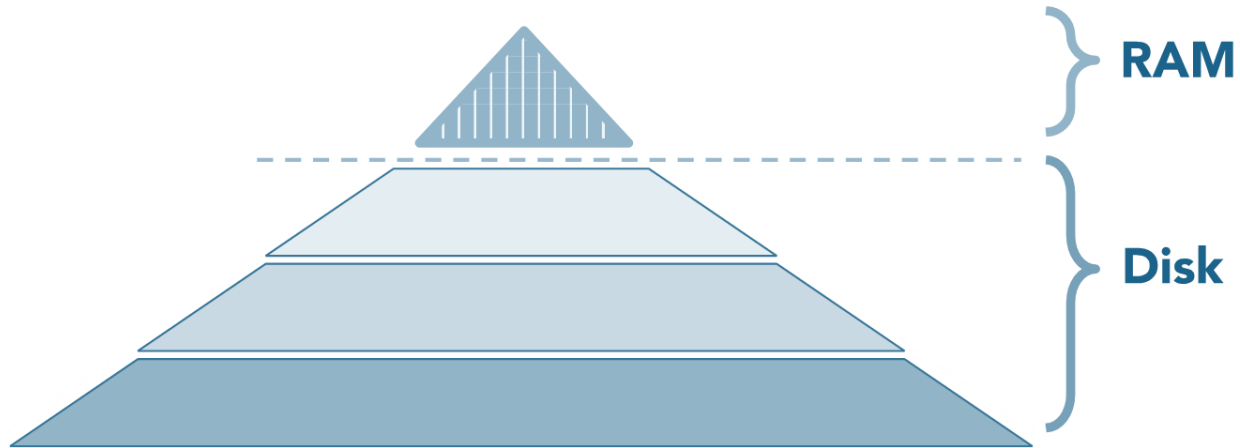
# RocksDB is CPU-bound

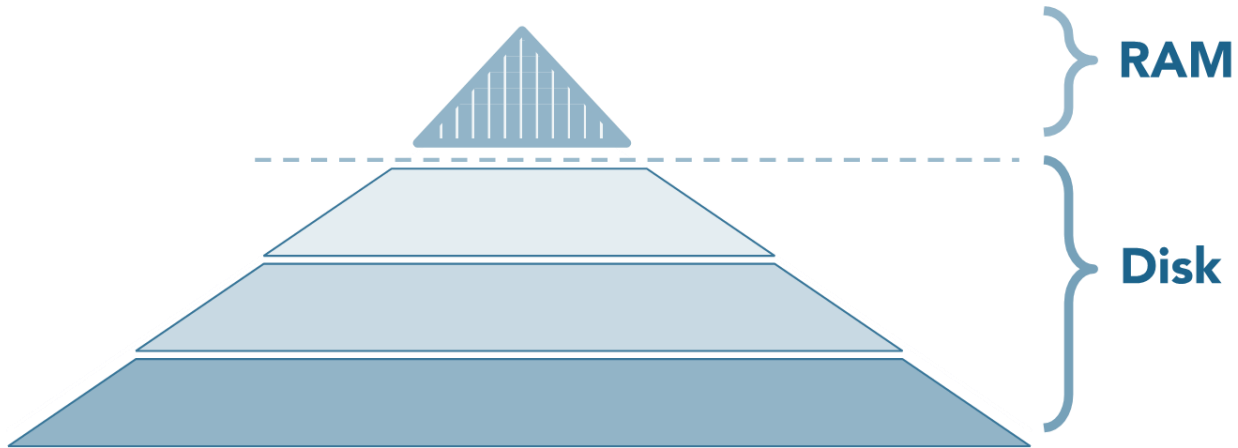# Popular design #1: LSM



RAM

Disk

# Popular design #1: LSM

RAM

Disk

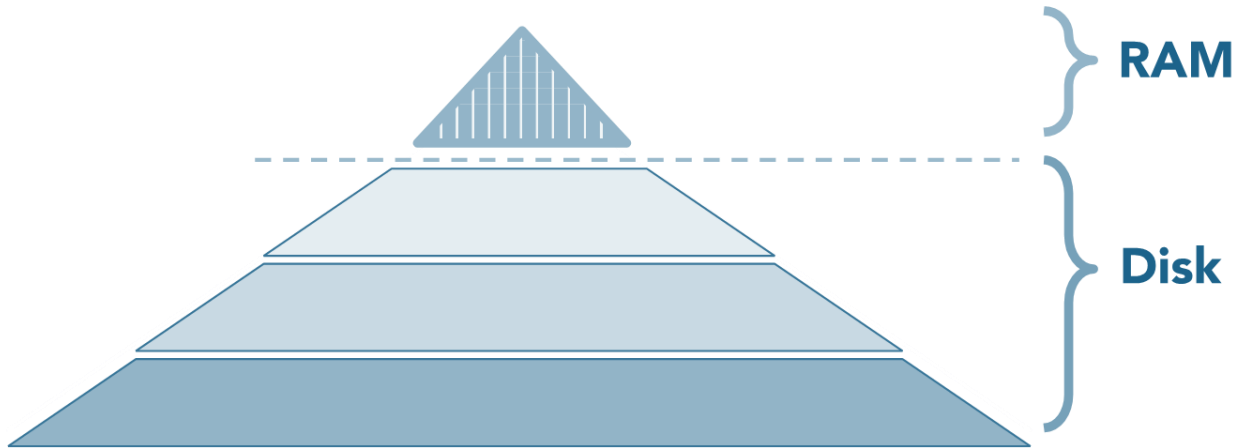Data **ordered by key**
in RAM and on disk

# Popular design #1: LSM



Updates **buffered** in RAM.

RAM flushed to disk
➔ **Large sequential IO**

# Popular design #1: LSM
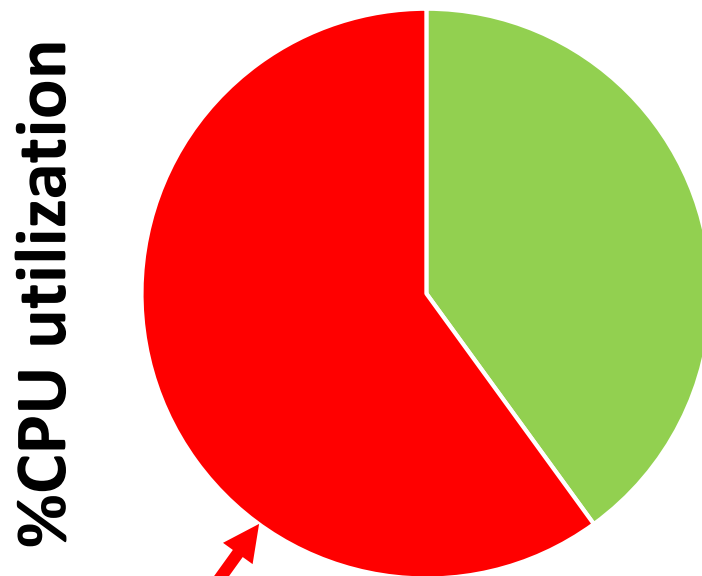


Updates **buffered** in RAM.

RAM flushed to disk
**merged** in the **ordered** main
structure (**compaction**)

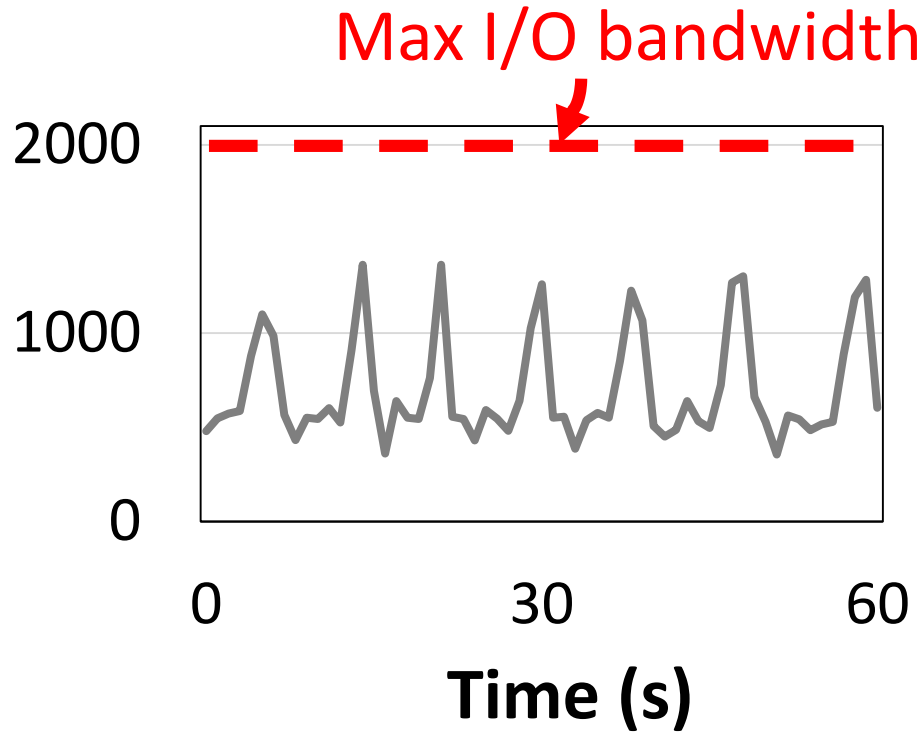# RocksDB is CPU-bound



%CPU utilization

**60% - merging + creating indexes of the disk structure**

# RocksDB's performance **fluctuates**

Max I/O bandwidth

# RocksDB's performance **fluctuates**

Max I/O bandwidth

1 flush = large backlog of work

# Popular design #2: B+ Trees



**%CPU utilization**

**60% - Contention on shared data structures**
**➜ low average throughput**

**Throughput (KOps/s)**

120

60

0

0    20    40    60

**Time (s)**

**Large buffers**
**➜ fluctuations**

# Lessons learned

❌ **Ordering**

❌ **Contention** ➜ **low average throughput**

❌ **Large buffers** ➜ **fluctuations**

# How to design an efficient KV for very fast drives?

# Key ideas

❌ ~~Ordering~~

✔️ **Data unsorted on disk (but sorted in memory)**

# Key ideas

❌ ~~Ordering~~

❌ ~~Contention~~

✅ **Data unsorted on disk (but sorted in memory)**

✅ **Shared-nothing**

# Key ideas

❌ ~~Ordering~~

✅ **Data unsorted on disk (but sorted in memory)**

❌ ~~Contention~~

✅ **Shared-nothing**

❌ ~~Large buffers~~

✅ **No buffering**

# Key idea #1 – data unsorted on disk

# Key idea #1 – data unsorted on disk

Unsorted data on disk

Put( **k**, **v** )

**k**, **v**

File 1

File 2

# Key idea #1 – data ordered in memory

**In-memory *ordered* index *(for scans)***

**Unsorted data on disk**

Prefix(K0)

Prefix(K3)

Prefix(k) ➜ **[file, idx]**

Prefix(K4) | Prefix(K5)

Prefix(K2)

**RAM**

file

**k, v**   idx

SSD

# Key idea #2 – no sharing

Sharding (static partitioning) - N independent workers

**Worker 1**

Key % 3 == 0

**Worker 2**

Key % 3 == 1

**Worker 3**

Key % 3 == 2

# Key idea #2 – no sharing

Workers have their own index and files

# Key idea #3 – no buffering

**Traditionally**   Put(k, v) —**write**→ Page Cache **RAM** ····**delayed write**····> SSD

# Key idea #3 – no buffering

**Traditionally**

Put(k, v) —write→ Page Cache (RAM) --delayed write--> SSD

**KVell**

Put(k, v) ——→ Page Cache (RAM) —write→ SSD

# Implementation challenges

**Syscall** cost

**Data structures**

**Manage disk queue length**

# Evaluation

**Machines:**

4 cores, 32GB RAM, Optane 905P drive (**500K IOPS**, 2GB/s)

**Benchmark:**

**YCSB – 1KB items, 100M elements (100GB)**

**Competition:**

# Evaluation – YCSB

# Evaluation – YCSB



Legend: RocksDB, PebblesDB, TokuMX, WiredTiger, Kvell

Chart axis: Throughput (KOps/s) — 0, 200, 400, 600, 800; right axis 0, 5, 20

Uniform

**KVell runs at disk BW
(75% of CPU time idle)**

YCSB A — 50/50 read/write
YCSB B — 95/5 read/write
YCSB C — 100% read
YCSB D — 50/50 read/write
YCSB F — 95/5 read/write
YCSB E — 95/5 scans/write

# Evaluation – YCSB – Scans



YCSB E

95/5
scans/write

Throughput (KOps/s)

Time (s)

■ RocksDB   ▤ PebblesDB   ☐ TokuMX   ☐ WiredTiger   ■ Kvell

# Evaluation – YCSB – Scans



**RocksDB drops to 1.8K scans/s even on a read mostly workload**

20

15

Th                0

YCSB E

95/5
scans/write

Throughput (KOps/s)

60

40

20

0

0        60       120      180      240

Time (s)

■ RocksDB    ▤ PebblesDB    ☐ TokuMX    ☐ WiredTiger    ■ Kvell

# In the paper

- Limitations:
  - Indexes have to fit in memory
  - Suboptimal scans for small items

- AWS machine, 15GB/s, 5TB dataset

- Production workload

- Recovery time

...

# Conclusions & take away messages

- **Ordering data is expensive**
- **Buffering** creates big fluctuation

- Optimizing for **CPU utilization is key**

**To kvell: to feel happy and proud**

https://github.com/BLepers/KVell
Code and scripts to reproduce results on AWS