# The True Value of Storage Drives with Built-in Transparent Compression: Far Beyond Lower Storage Cost

**Tong Zhang**

ScaleFlux Inc.

San Jose, CA

# The Rise of Computational Storage

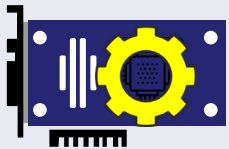THE END OF MOORE'S LAW

Homogeneous Computing

Heterogenous Computing

**Compute** ⚙️

FPGA/GPU/TPU
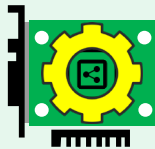
End of Moore's Law
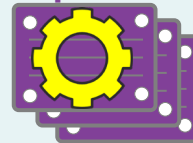
**Networking**

SmartNICs

10 → 100-400Gb/s

**Storage**

**Computational**

**Fast & Big Data Growth**

**Domain Specific Compute**

# Computational Storage: A Very Simple Idea

☐ End of Moore's Law ➔ heterogeneous computing

**Computational Storage**

Low-hanging fruits

FPGA/GPU/TPU

intel Xeon processor

SmartNICs

←SW | HW →

PCI EXPRESS

MySQL.

APACHE Spark™

ORACLE®

TensorFlow

MariaDB®

In-line per-4KB zlib compression & decompression

Flash Control

NAND Flash

FPGA

Computational Storage Drive (CSD) with Data Path Transparent Compression

# ScaleFlux Computational Storage Drive: CSD 2000



**SSD**

Flash Controller

Flash   Flash

**Multiple, discrete components for Compute and SSD Functions**

**CSD**

FPGA   FC

Flash   Flash

**Single FPGA combines Compute and SSD Functions**

- ✓ Complete, validated solution
  - ✓ Pre-Programmed FPGA
  - ✓ Hardware
  - ✓ Software
  - ✓ Firmware
- ✓ No FPGA knowledge or coding
- ✓ Field upgradeable
- ✓ Standard U.2 & AIC form factors

# CSD 2000: Data Path Transparent Compression



**FIO:  4K Random R/W IOPS**

IOPS (k) / Better

700
600
500
400
300
200
100
0

170%
70/30 R/W

230%
100% Write

CSD 2000

NVMe SSD

100%  90%  80%  70%  60%  50%  40%  30%  20%  10%  0%

100% Reads                                          100% Writes

2.5:1 Compressible Data, 8 jobs, 32 QD, steady state after preconditioning

**FIO:  16K Random R/W IOPS**

IOPS (k) / Better

250
200
150
100
50
0

220%
70/30 R/W

220%
100% Write

CSD 2000

NVMe SSD

100%  90%  80%  70%  60%  50%  40%  30%  20%  10%  0%

100% Reads                                          100% Writes

2.5:1 Compressible Data, 8 jobs, 32 QD, steady state after preconditioning

# Comparing Compression Options



| | No Compression | Host-Based | Offload Card | CSD 2000 |
|---|:---:|:---:|:---:|:---:|
| No CPU Overhead | ✔ | ✘ | ✔ | ✔ |
| Reduced $/User GB | ✘ | ✔ | ✔ | ✔ |
| Performance scales with capacity | ✔ | ✘ | ✘ | ✔ |
| Transparent App Integration | - | ✘ | ✘ | ✔ |
| Zero App Latency | ✔ | ✘ | ✘ | ✔ |
| No incremental power usage | ✔ | ✘ | ✘ | ✔ |
| No incremental physical footprint | ✔ | ✔ | ✘ | ✔ |

**Scalable CSD-based compression reduces Cost/GB without choking the CPU**

# In-Storage Transparent Compression: Why is It Hard to Build?



Logical Block Address (LBA)

Flash Memory

Proc. 1

Proc. 2

Proc. n

**CORE**

Flash Translation Layer

(w/o compression)

4KB LBA ➜ 4KB flash block mapping

**Regularity & uniformity**

✓ Relatively simple FTL implementation

✓ Relatively easy to achieve high speed

✓ Relatively easy to ensure storage stability

# In-Storage Transparent Compression: Why is It Hard to Build?

# CSD 2000: Highest OLTP TPS, Lowest $/User GB

- Sysbench (MySQL 5.7.25, InnoDB)
- 50M records, 64 Threads
- 1hr Test run
- Intel(R) Xeon(R) CPU E5-2667 v4 @ 3.20GHz, 256GB DRAM

## Performance: 150% TPS

Sysbench (MySQL): Read-Write Transactions per Second (TPS)
2.4TB Dataset



2.4TB Dataset Physical Flash consumed on NVMe A; 0.9TB on CSD 2000

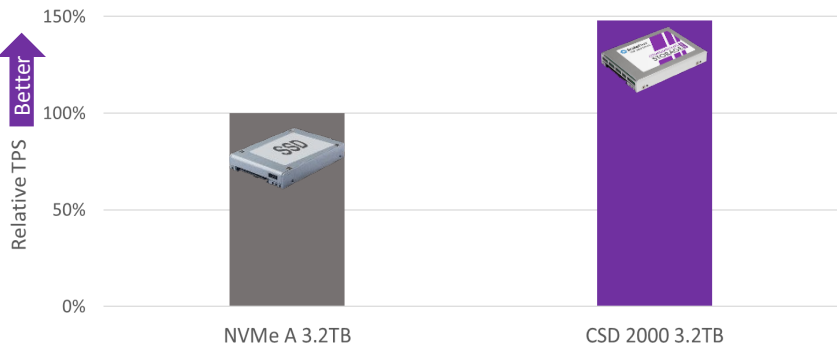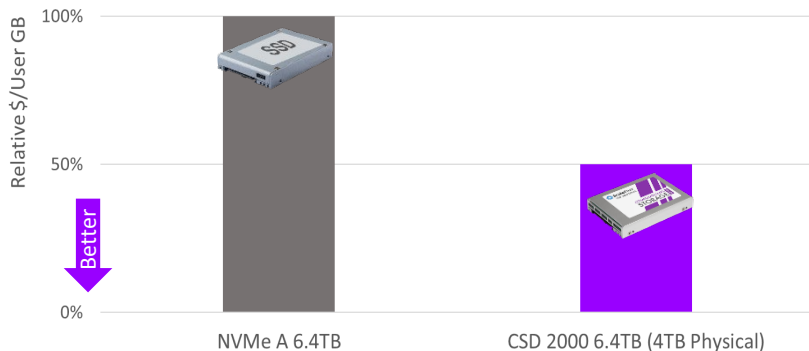## Cost: 50% Less $/User GB

Flash Storage $/User GB Comparison
4.8TB Dataset



4.8TB Dataset Physical Flash consumed on NVMe A; 1.6TB on CSD 2000
CSD 2000 delivers 30% higher Read-Write TPS in this cost comparison

## Flexible Drive Capacity Enables the Best Performance ↔ Cost

# Open a Door for System Innovation

Transparent compression

**Logical** storage space utilization efficiency

**Physical** storage space utilization efficiency

Exposed LBA space (e.g., 32TB)

FTL with transparent compression

SSD

NAND Flash (e.g., 4TB)

**Unnecessary** to use all the LBAs

4KB

| Valid user data | 0's |

Transparent compression

Compressed data

**Unnecessary** to fill each 4KB sector with user data

OS/Applications can **_purposely waste_** logical storage space to gain benefits

# Case Study 1: PostgreSQL



8KB/page

Data

Fillfactor (FF)

Reserved for future update

8KB/page

Data | 0's

Transparent compression

Compressed data

FF ↓

Performance ↑

Storage space ↑

Normalized Performance

$R_L \approx 4KB$

200%

100%

$R_L \approx 0$

● SFX CSD 2000

● Commodity SSD

Physical storage usage

300GB  600GB  1.2TB

# Case Study 1: PostgreSQL



**740GB**

Normalized TPS

Vendor-A | CSD 2000

FF100 (740GB) → FF100 (178GB)
FF75 (905GB) → FF75 (189GB)

| Fillfactor | Drive | Logical size (GB) | Physical size (GB) | Comp Ratio |
|---|---|---|---|---|
| 100 | Vendor-A | **740** | 740 | 1.00 |
| | CSD 2000 | | **178** | 4.12 |
| 75 | Vendor-A | **905** | 905 | 1.00 |
| | CSD 2000 | | **189** | 4.75 |

**1.4TB**

Normalized TPS

Vendor-A | CSD 2000

FF100 (1,433GB) → FF100 (342GB)
FF75 (1,762GB) → FF75 (365GB)

| Fillfactor | Drive | Logical size (GB) | Physical size (GB) | Comp Ratio |
|---|---|---|---|---|
| 100 | Vendor-A | **1,433** | 1,433 | 1.00 |
| | CSD 2000 | | **342** | 4.19 |
| 75 | Vendor-A | **1,762** | 1,762 | 1.00 |
| | CSD 2000 | | **365** | 4.82 |

# Case Study 2: Sparse Write-Ahead Logging

❑ Write-ahead logging (WAL)

➢ Universally used by data management systems to achieve _atomicity_ and _durability_

# Case Study 2: Sparse Write-Ahead Logging

# Case Study 2: Sparse Write-Ahead Logging

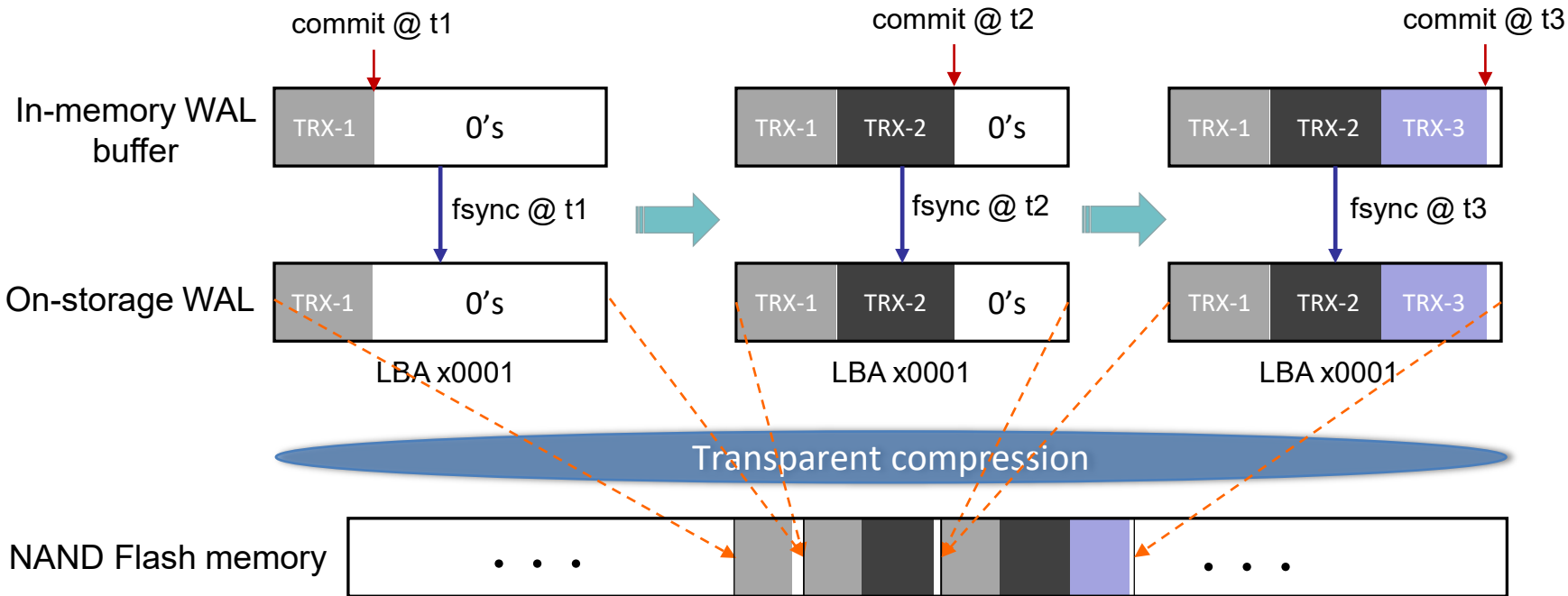❑ Sparse WAL: Allocate a new 4KB sector per transaction commit

✓ Waste logical storage space ➜ reduce WAL-induced write amplification

commit @ t1          commit @ t2          commit @ t3

In-memory WAL
buffer

| TRX-1 | 0's |     | TRX-2 | 0's |     | TRX-3 | 0's |

fsync @ t1          fsync @ t2          fsync @ t3

On-storage WAL

| TRX-1 | 0's |     | TRX-2 | 0's |     | TRX-3 | 0's |

LBA x0001          LBA x0002          LBA x0003

Transparent compression

NAND Flash memory    . . .    . . .

# Case Study 2: Sparse Write-Ahead Logging

❑ Sparse WAL: Allocate a new 4KB sector per transaction commit
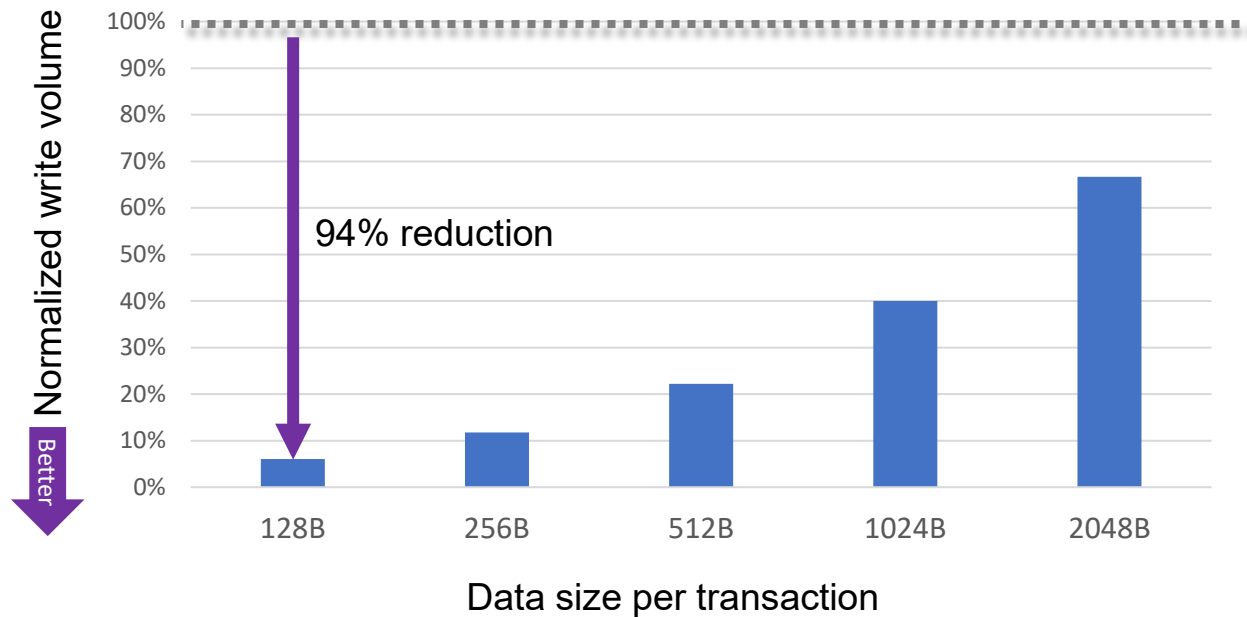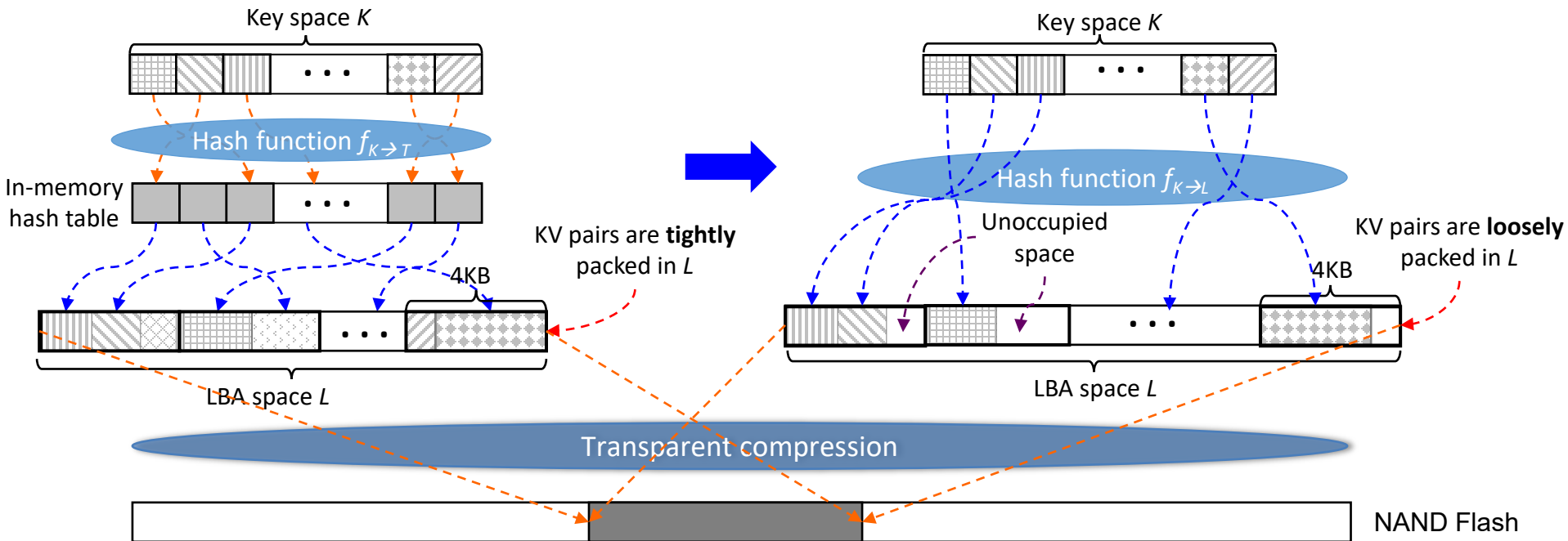
  ✓ Waste logical storage space ➔ reduce WAL-induced write amplification



94% reduction

Normalized write volume — Better

| Data size per transaction | 128B | 256B | 512B | 1024B | 2048B |

# Case Study 3: Table-less Hash-based KV Store

❑ Very simple idea

➢ Hash *key space* directly onto *logical storage space* ➜ eliminate the in-memory hash table

➢ Transparent compression eliminates the "unoccupied space" from physical storage space



Key space $K$

Hash function $f_{K \to T}$

In-memory hash table

KV pairs are **tightly** packed in $L$

4KB

LBA space $L$

Key space $K$

Hash function $f_{K \to L}$

Unoccupied space

KV pairs are **loosely** packed in $L$

4KB

LBA space $L$

Transparent compression

NAND Flash

# Case Study 3: Table-less Hash-based KV Store

❑ Eliminate in-memory hash table
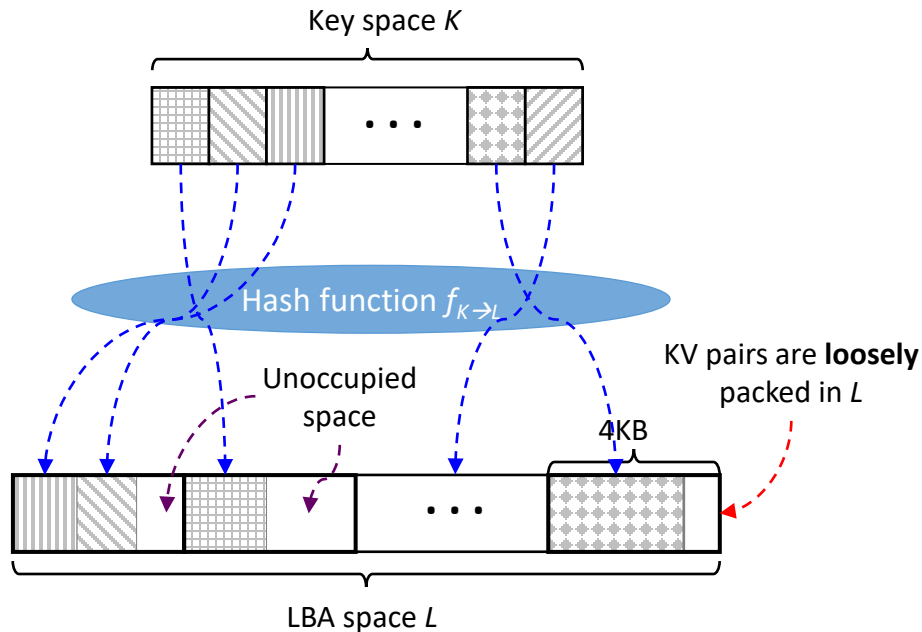
    ✓  Very small memory footprint

    ✓  High operational parallelism

    ✓  Short data access data path

    ✓  Very simple code base

Key space $K$

Hash function $f_{K \to L}$

Unoccupied space

KV pairs are **loosely** packed in $L$

4KB

LBA space $L$

❑ Under-utilize logical storage space

    ✓  Obviate frequent background operations (e.g., GC and compaction)

➡ High performance, low memory cost, and low CPU usage
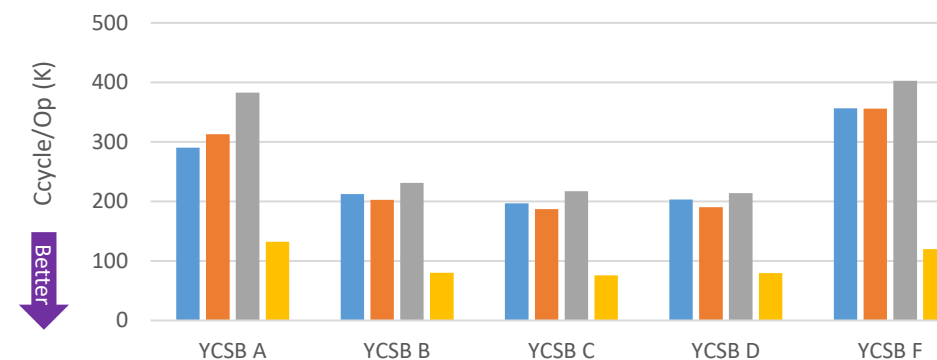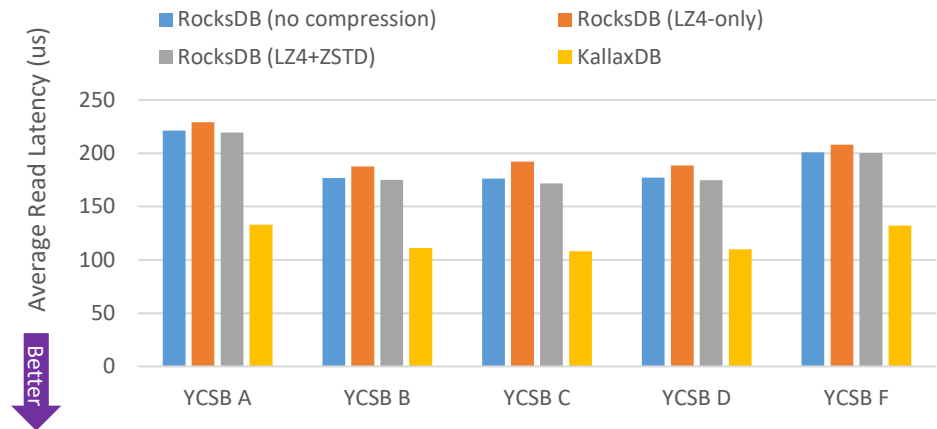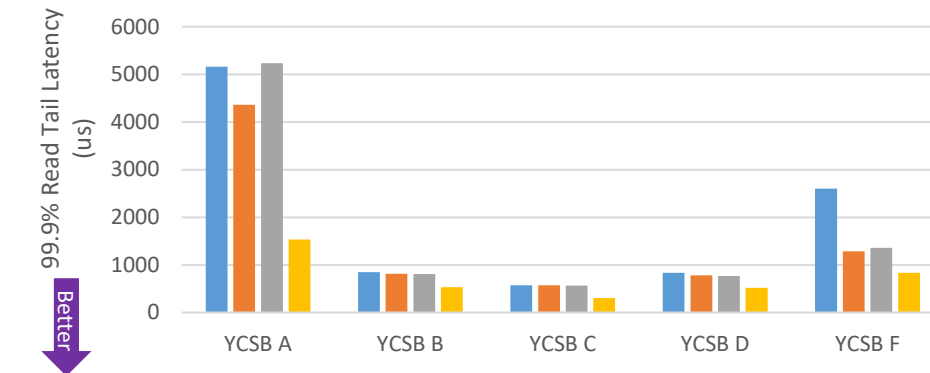
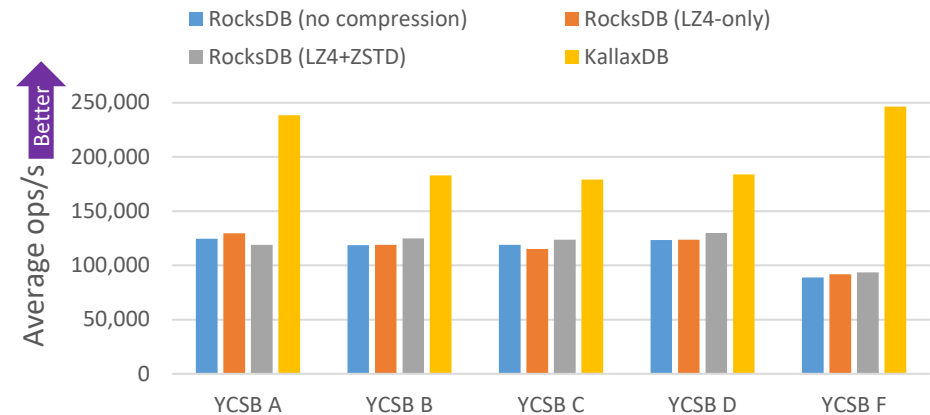# Case Study 3: Table-less Hash-based KV Store

- ❑ Experimental Setup
  - ➢ 24-core 2.6GHz Intel CPU, 32GB DDR4 DRAM, and a 3.2TB SFX CSD2000
  - ➢ RocksDB 6.10 (12 compaction threads and 4 flush threads)
  - ➢ 400-byte KV pair size, 1 billion KVs ➔ 400GB raw data
  - ➢ Memory usage: RocksDB (5GB), KallaxDB (600MB)

| YCSB A | 50% reads, 50% updates |
|--------|-------------------------|
| YCSB B | 95% reads, 5% updates |
| YCSB C | 100% reads |
| YCSB D | 95% reads, 5% inserts |
| YCSB F | 50% reads, 50% read-modify-writes |

| | Storage Usage |
|--------|---------------|
| RocksDB (no compression) | 428GB |
| RocksDB (LZ4-only) | 235GB |
| RocksDB (LZ4+ZSTD) | 201GB |
| KallaxDB | 216GB |

# Case Study 3: Experimental Results (24 clients)

# Open a Door for System Innovation

Transparent compression

**Logical** storage space utilization efficiency

**Physical** storage space utilization efficiency

Exposed LBA space (e.g., 32TB)

FTL with transparent compression

NAND Flash (e.g., 4TB)

SSD

**Unnecessary** to use all the LBAs

4KB

| Valid user data | 0's |

Transparent compression

Compressed data

**Unnecessary** to fill each 4KB sector with user data

OS/Applications can **_purposely waste_** logical storage space to gain benefits

# Open a Door for System Innovation

8KB/page

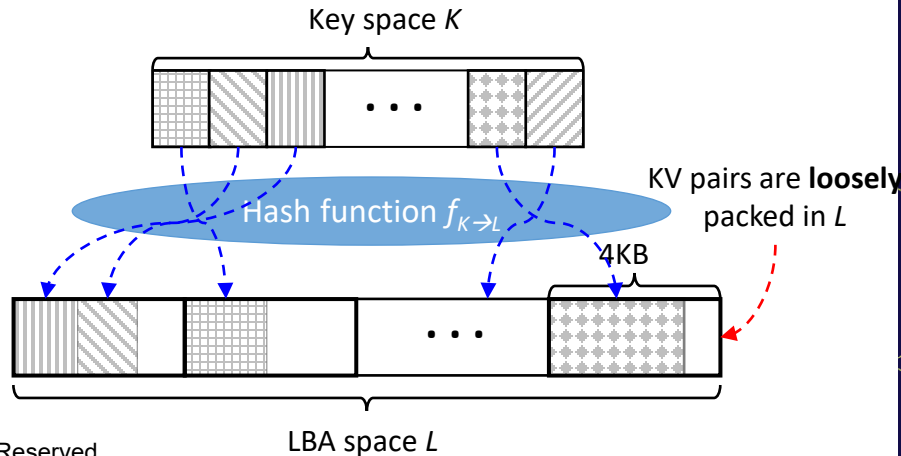PostgreSQL

| Data | Reserved for future update |

Reserve more space for future update to improve performance @ zero storage overhead

Sparse WAL ⟹ Reduce WAL-induced write amplification @ zero storage overhead

Table-less hash-based KV store

⬇

High performance, low memory/CPU usage
@ zero storage overhead

Key space $K$

· · ·

Hash function $f_{K \rightarrow L}$

KV pairs are **loosely** packed in $L$

4KB

· · ·

LBA space $L$

# Thank You

www.scaleflux.com

info@scaleflux.com

tong.zhang@scaleflux.com