sodacon
— Global 2021 — July 13-14
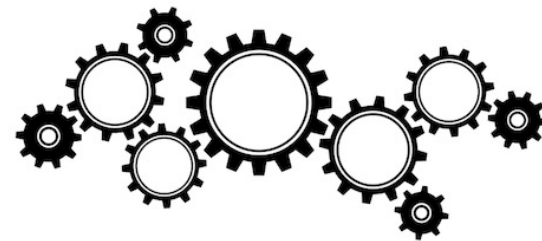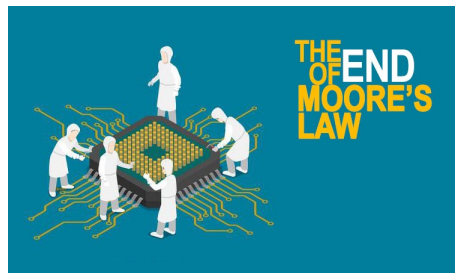
# Innovating IT Infrastructure with Computational Storage – Can you believe (or achieve) the hype?

Tong Zhang, Chief Scientist @ ScaleFlux

soda
foundation

#sodacon2021

# The Rise of Computational Storage



Homogeneous Computing

Heterogenous Computing

| Compute | Networking | Storage | |
|---|---|---|---|
| **FPGA/GPU/TPU** | **SmartNICs** | **Computational Storage** | **Domain Specific Compute** |
| End of Moore's Law | 10 → 100-400Gb/s | **Fast & Big Data Growth** | |

ScaleFlux®

# Beauty of Abstraction

Applications

Video coding

Operating Systems

Compiler

Networking

Instruction Set Architecture (ISA)

Processor

Memory

I/O

SSD

HDD

SmartNICs

Video Codec

Break the principle of abstraction

x **Cost**: Modify one (or more) layers

x **Cost**: Enhance cross-layer interface

x **Cost**: Vendor lock-in

✓ **Benefit**: Improve the system performance/efficiency

**Benefit**

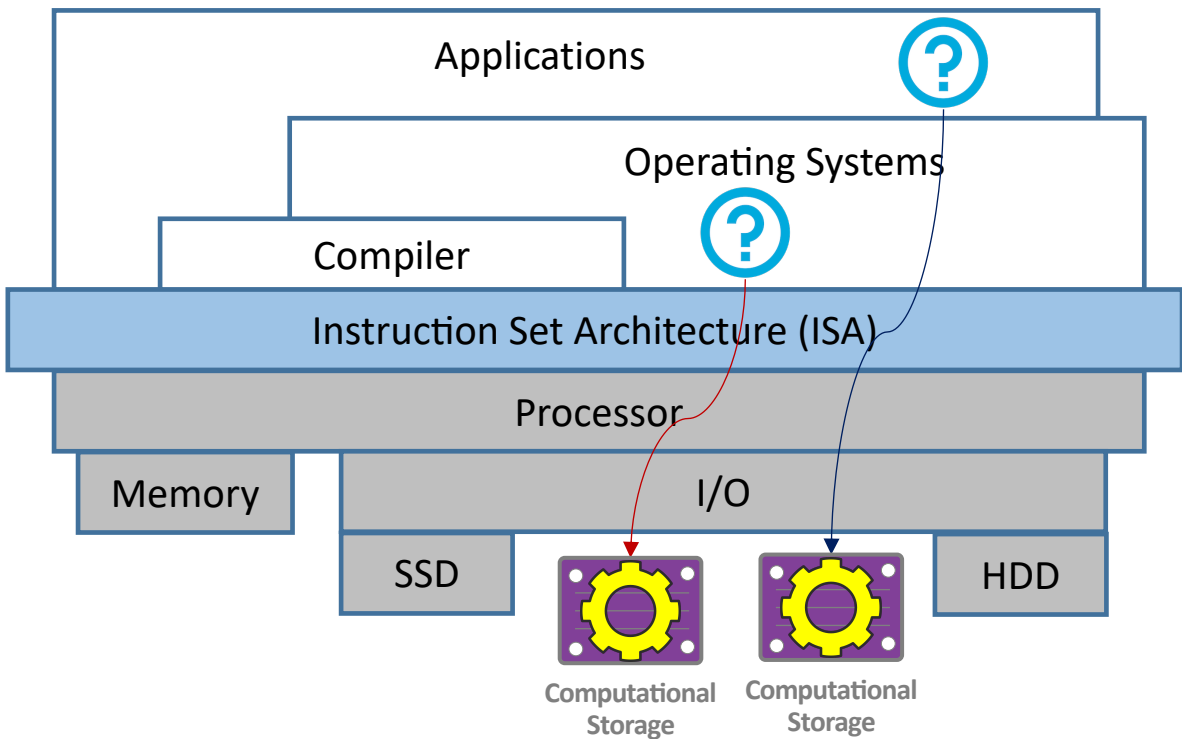**Cost**

ScaleFlux®

# Beauty of Abstraction

Break the principle of abstraction



x **Cost**: Modify one (or more) layers

x **Cost**: Enhance cross-layer interface

x **Cost**: Vendor lock-in

✓ **Benefit**: Improve the system performance/efficiency

# Computational Storage

❑ Academia started to explore this simple concept since 20 years ago

 ➢ Coined many terms: "Intelligent RAM", "Active Disk", "Intelligent SSD", …

 ➢ Not surprisingly, very impressive performance gains were demonstrated **on papers** for applications such as linear algebra, multimedia/graph processing, database, …

**MISSING**
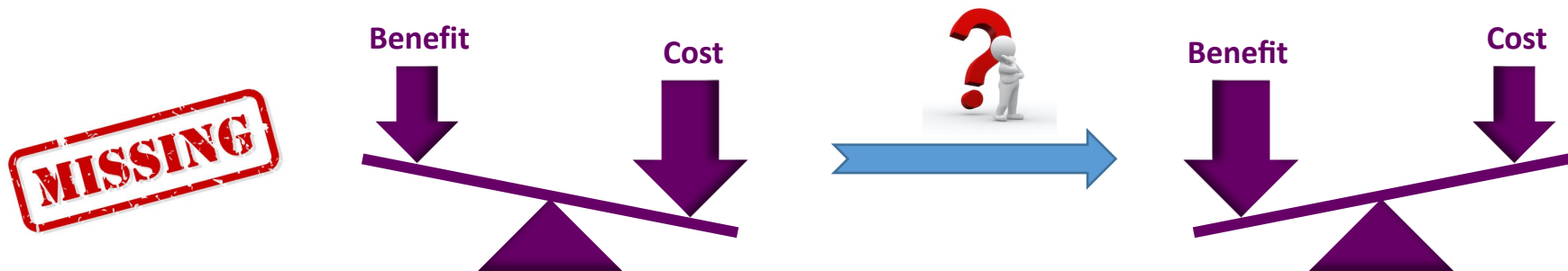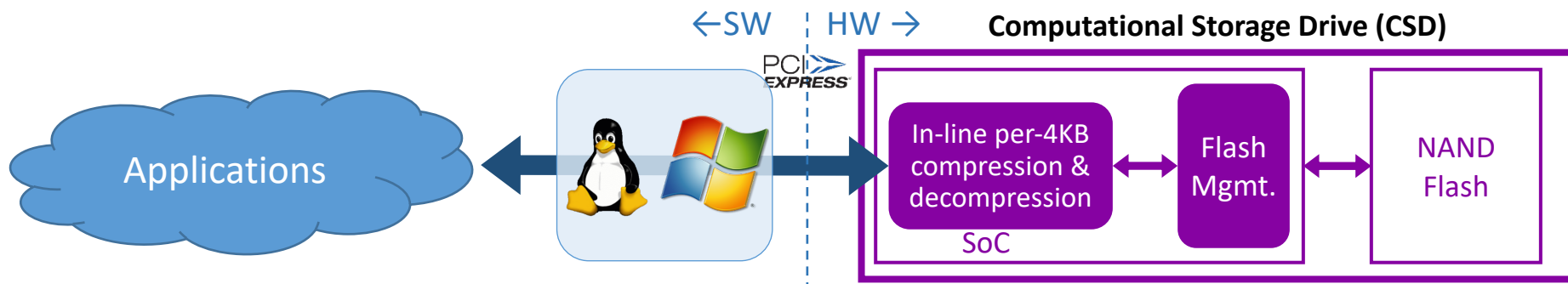
Benefit    Cost

Benefit    Cost

  x  **Cost**: Modify one (or more) layers

  x  **Cost**: Enhance cross-layer interface

  x  **Cost**: Vendor lock-in

**VS**   ✓ **Benefit**: Improve the system performance/efficiency

# Computational Storage

❑ A perfect low-hanging fruit: In-storage **transparent** compression



←SW | HW →

**Computational Storage Drive (CSD)**

Applications

In-line per-4KB compression & decompression

Flash Mgmt.

NAND Flash

SoC
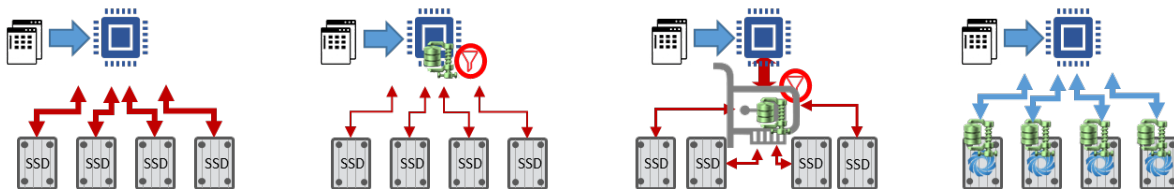
x **Cost**: Modify one (or more) layers

x **Cost**: Enhance cross-layer interface

x **Cost**: Vendor lock in

**VS**

✓ **Benefit**: Improve the system performance/efficiency

ScaleFlux®

# Comparing Compression Options

| | No Compression | Host-Based | Offload Card | CSD |
|---|:---:|:---:|:---:|:---:|
| 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**

**ScaleFlux**®

# Seamless Deployment

➢ Zero changes to application source code

➢ >2x storage cost reduction at zero host CPU usage

➢ Representative use cases: Relational database (e.g., MySQL, Postgre, Oracle), key-value store (e.g., Aerospike), data analytics (e.g., ClickHouse, Apache Kudu)
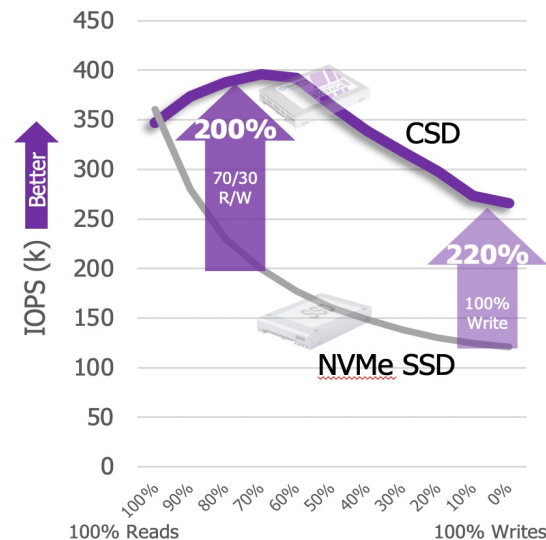
FIO: 8K Random R/W IOPS

❑ >50% storage cost saving

❑ 50% TPS improvement

❑ 38% latency reduction

❑ >50% storage cost saving

❑ 35% TPS improvement

❑ 20% latency reduction

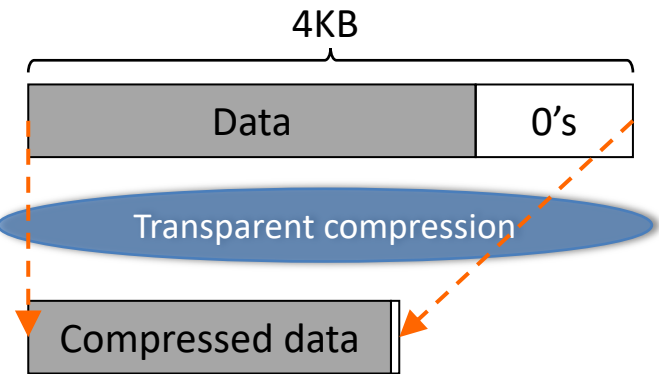❑ >50% storage cost saving

❑ 50% TPS improvement

❑ 98% tail latency reduction

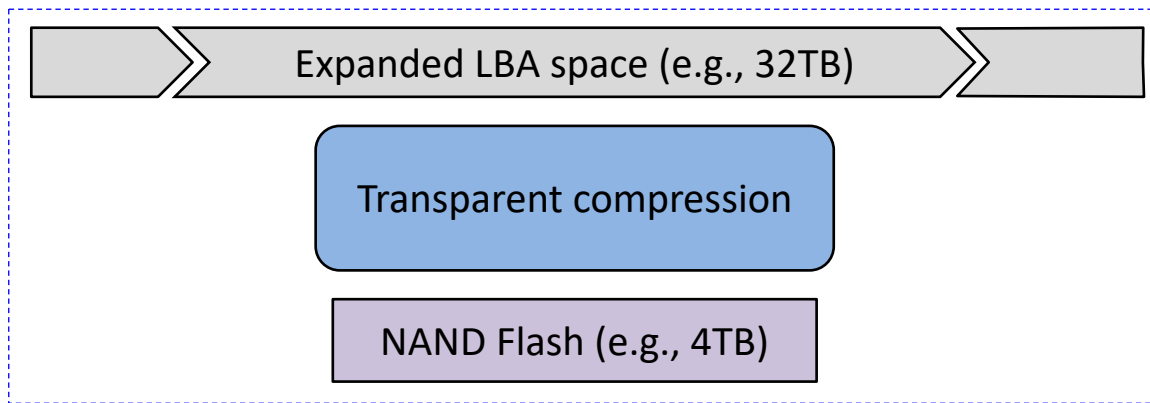# System-level Innovations Enabled by Transparent Compression

4KB

| Data | 0's |
|------|-----|

Transparent compression

| Compressed data |
|-----------------|

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

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

Transparent compression

NAND Flash (e.g., 4TB)

**Unnecessary** to use all the LBAs

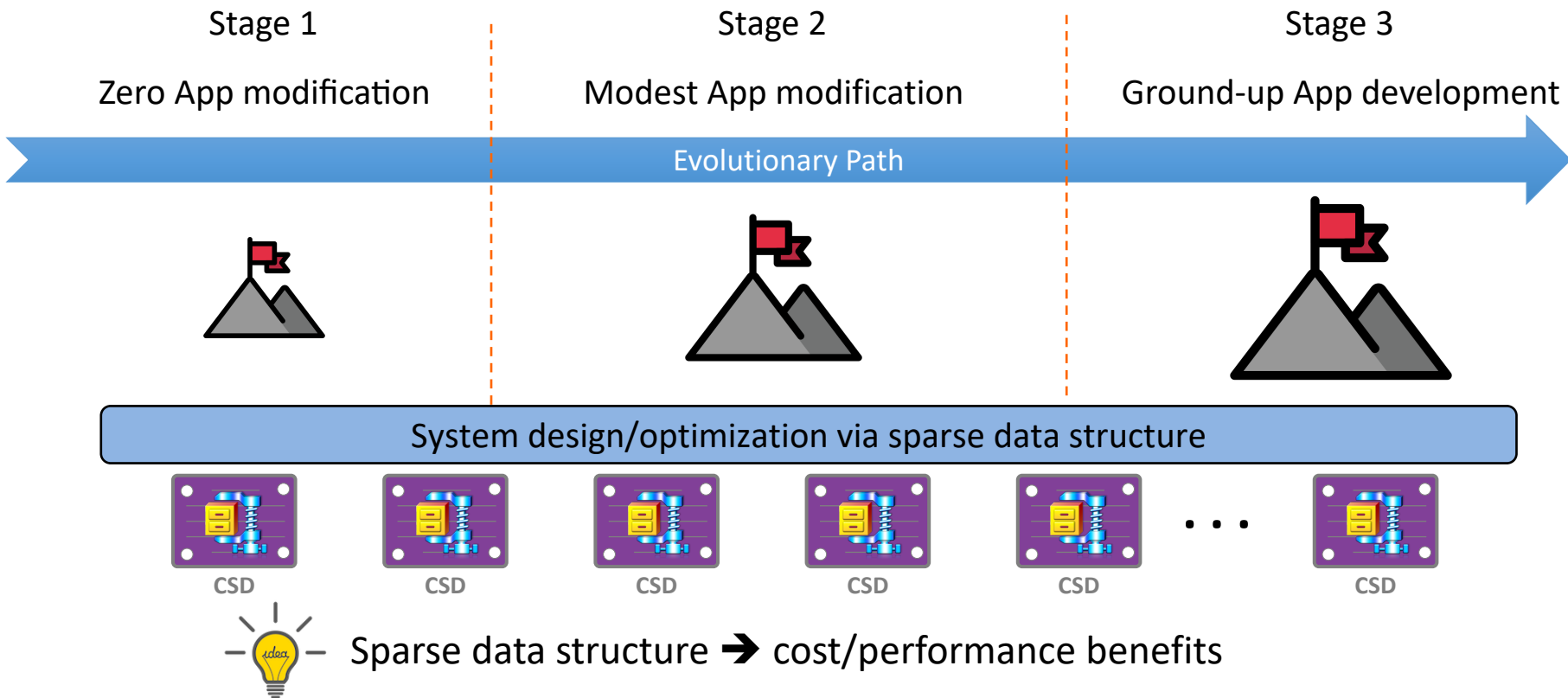OS/Applications can employ *sparse* data structure without sacrificing storage cost

Sparse data structure ➔ cost/performance benefits

**ScaleFlux**®

# Beauty of Abstraction
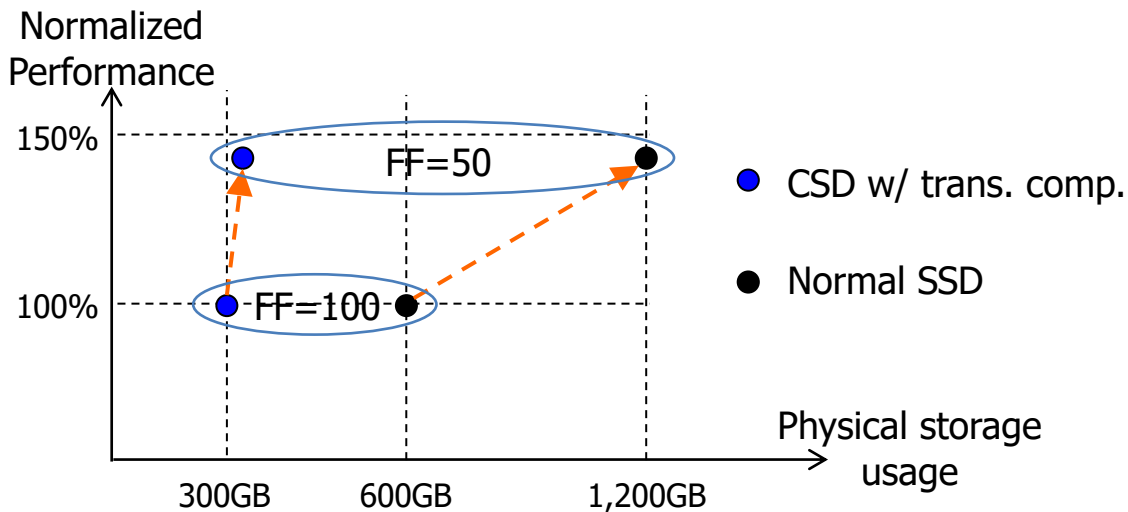
Data management with sparse data structure

| |
|---|
| Applications |
| Operating Systems |
| Compiler |
| Instruction Set Architecture (ISA) |
| Processor |
| Memory | I/O |
| SSD |

Computational Storage   . . .   Computational Storage

x **Cost**: Modify one (or more) layers

x ~~**Cost**: Enhance cross-layer interface~~

x ~~**Cost**: Vendor lock-in~~

✓ **Benefit**: Improve the system performance/efficiency

**Benefit**          **Cost**

ScaleFlux®

# System-level Innovations Enabled by Transparent Compression

Stage 1

Zero App modification

Stage 2

Modest App modification

Stage 3

Ground-up App development

Evolutionary Path

System design/optimization via sparse data structure

CSD    CSD    CSD    CSD    CSD    • • •    CSD

Sparse data structure ➜ cost/performance benefits

# Stage-1 Example: PostgreSQL



8KB/page

Data | Reserved for future update

Fillfactor (FF)

FF ↓ → Performance ↑

→ Storage space ↑

8KB/page

Data | 0's

Data path compression

Compressed data

Normalized Performance

150% — FF=50

100% — FF=100

● CSD w/ trans. comp.

● Normal SSD

300GB    600GB    1,200GB

Physical storage usage

# Stage-1 Example: Experiments (Sysbench-TPCC)



**740GB** — Normalized TPS chart

- Normal SSD: FF100 (740GB) = 100.0%, FF75 (905GB) ≈ 133%
- CSD: FF100 (178GB) ≈ 99%, FF75 (189GB) ≈ 135%

**1.4TB** — Normalized TPS chart

- Normal SSD: FF100 (1,433GB) = 100.0%, FF75 (1,762GB) ≈ 134%
- CSD: FF100 (342GB) ≈ 100%, FF75 (365GB) ≈ 137%

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

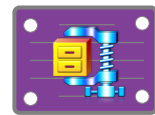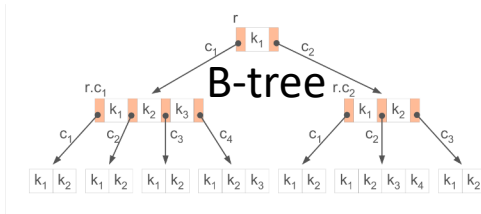| Fillfactor | Logical size (GB) | Drive | Physical size (GB) | Comp Ratio |
|---|---|---|---|---|
| 100 | **1,433** | Normal SSD | 1,433 | 1.00 |
| | | CSD | **342** | 4.19 |
| 75 | **1,762** | Normal SSD | 1,762 | 1.00 |
| | | CSD | **365** | 4.82 |

ScaleFlux®

# Stage-2 Example: Reduce B-tree Write Amplification

❑ Log-structured merge tree (LSM-tree)

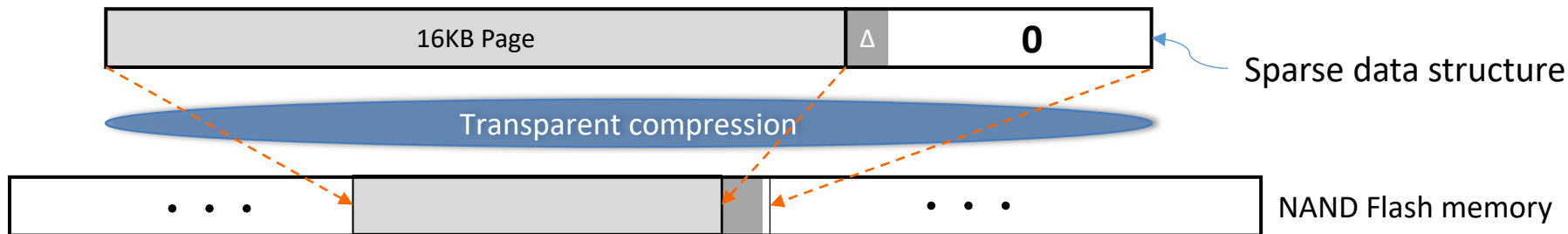➢ Significant recent interest: RocksDB, Cassandra, HBase, ScyllaDB, …

Advantages of LSM-tree

~~➢ Less storage space usage~~
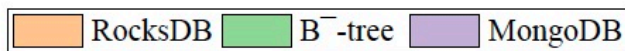
~~➢ Much smaller write amplification~~



B-tree

**Computational Storage**
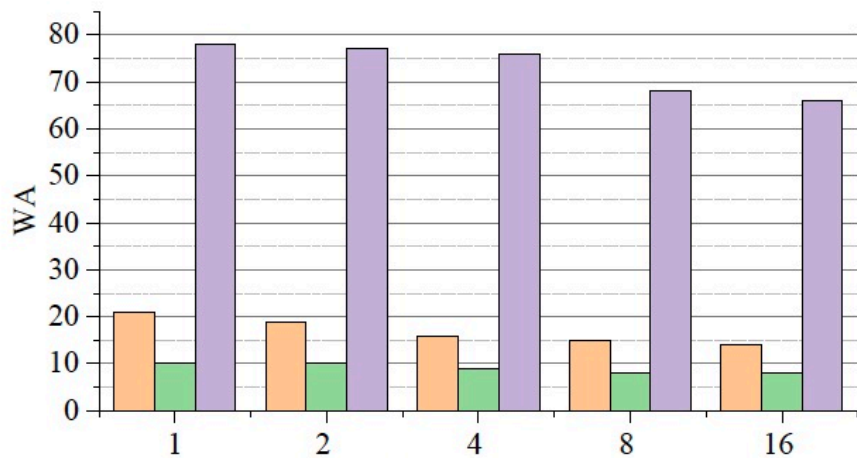
❑ Localized page modification logging

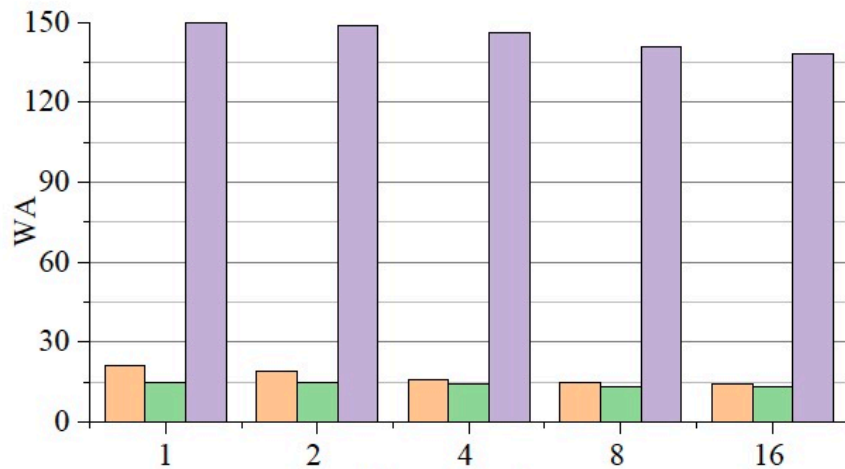➢ Very simple idea: Log page modification instead of re-writing the entire page every time



| 16KB Page | Δ | 0 |
|---|---|---|

Sparse data structure

Transparent compression

NAND Flash memory

ScaleFlux®

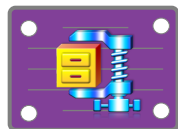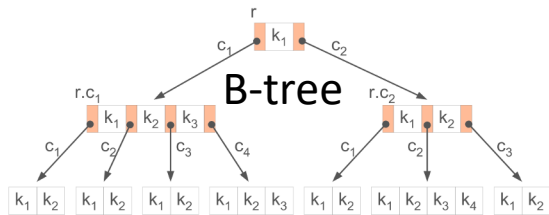# Stage-2 Example: Reduce B-tree Write Amplification



150GB dataset & 1GB cache

(a) 128B record, 8KB page

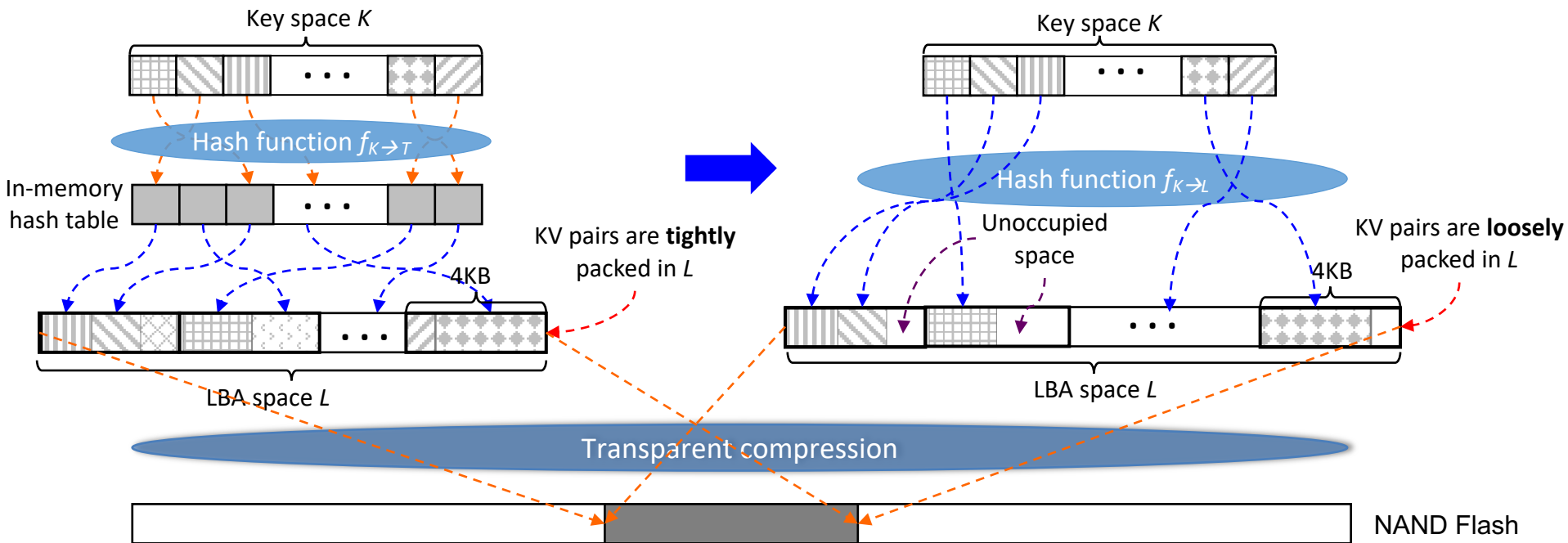(b) 128B record, 16KB page

B-tree **+** Computational Storage → Make LSM-tree (e.g., RocksDB) **much less** appealing

# Stage-3 Example: 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

# Stage-3 Example: Table-less Hash-based KV Store
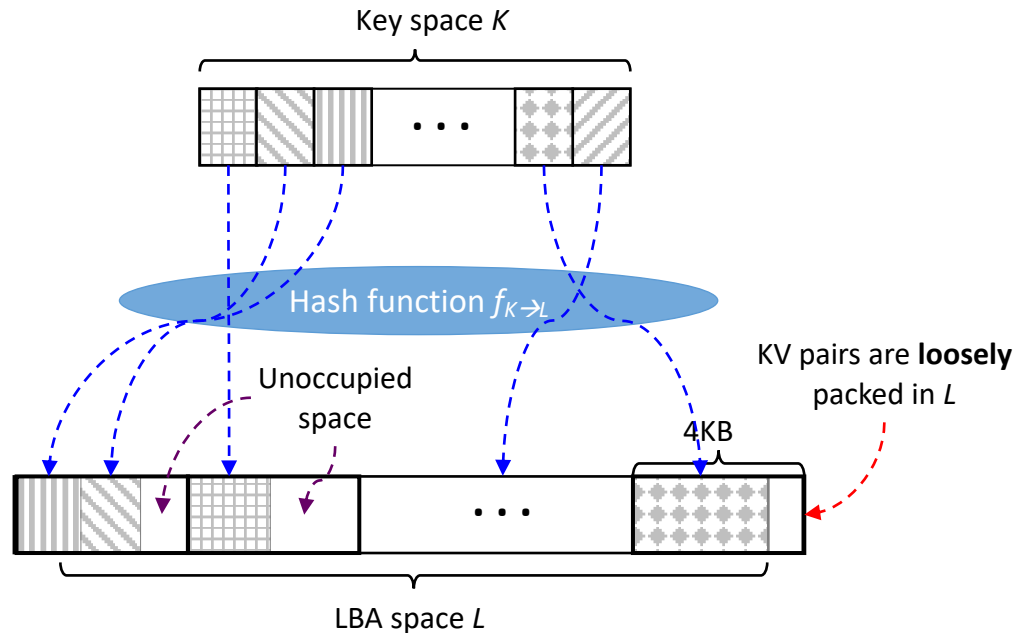
- ❑ Eliminate in-memory hash table

  - ✓ Very small memory footprint

  - ✓ High operational parallelism

  - ✓ Short data access data path

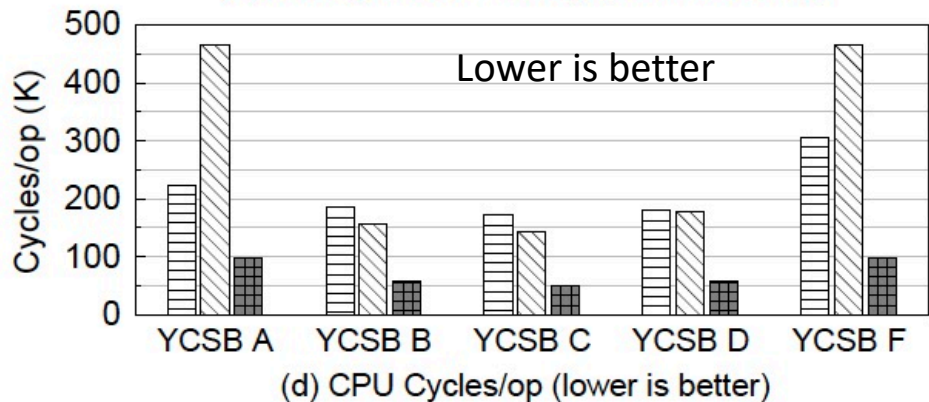  - ✓ Very simple code base
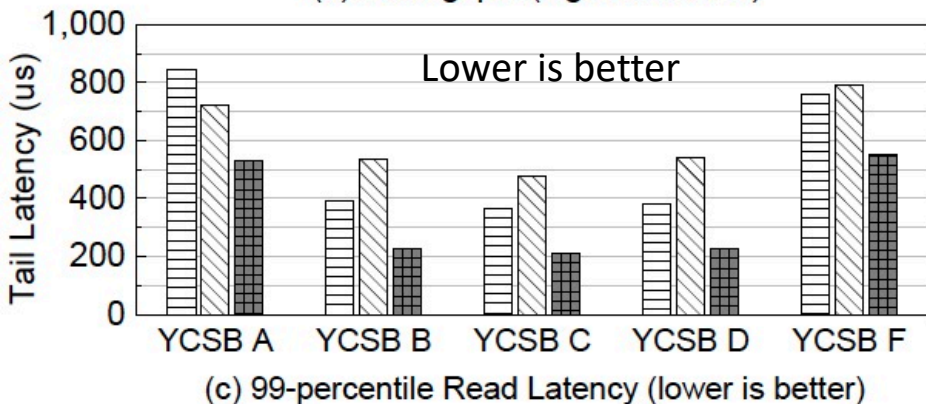
- ❑ Under-utilize logical storage space

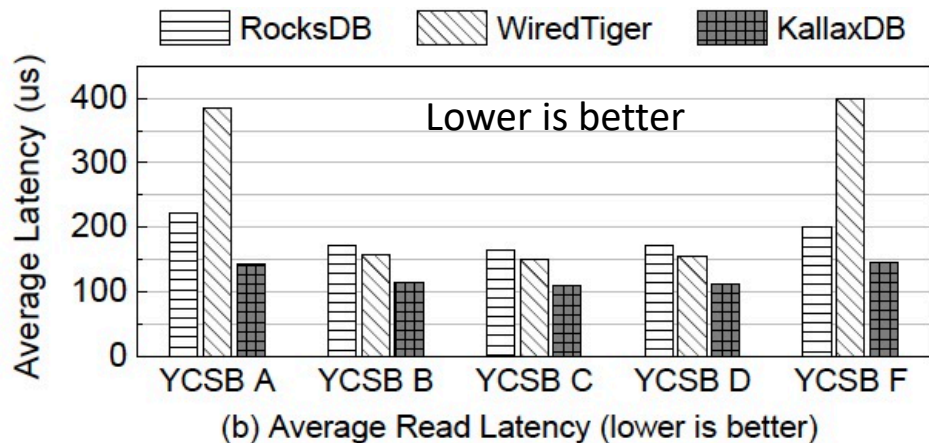  - ✓ Obviate frequent background operations (e.g., GC and compaction)

  ➡ High performance, low memory cost, and low CPU usage

Key space $K$

Hash function $f_{K \to L}$

Unoccupied space

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

4KB

LBA space $L$

ScaleFlux®

# Stage-3 Example: Results (400-byte KV, 400GB dataset)



(a) Throughput (higher is better)

(b) Average Read Latency (lower is better)

(c) 99-percentile Read Latency (lower is better)

(d) CPU Cycles/op (lower is better)

# Conclusion

☐ Yes, computational storage indeed has a real potential to innovate IT infrastructure

☐ The best starting point: Computational storage drive with transparent compression

| Stage 1 | Stage 2 | Stage 3 |
|---------|---------|---------|
| Zero App modification | Modest App modification | Ground-up App development |

Exciting IT Infrastructure Innovation Potential Enabled by Transparent Compression

CASE STUDY — PostgreSQL

CASE STUDY — Modify B-tree to reduce write amplification

CASE STUDY — A table-less hash-based key-value store

☐ Beyond performance gain, two additional HUGE benefits to end customers:

✓ Avoid modifying interface across applications↔OS/filesystem↔driver ↔HW

✓ Avoid computational storage drive vendor lock-in

ScaleFlux®