# Storage Systems (StoSys) XM_0092

# Lecture 5: Key-Value Stores

Animesh Trivedi
Autumn 2020, Period 2

VU VRIJE UNIVERSITEIT AMSTERDAM

# Reminder: for the coming weeks

We will be gradually transforming to networking and distributed systems

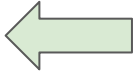It is important you understand networking basics and important concepts such as

- TSO, LRO, Jumbo Frames, Multicore scalability, affinities, and RDMA, etc.

I will only introduce these topics selectively

**Background reading:** Please check out lecture 1, 2 (networking basic), 4 (multicore scalability), and 6 (RDMA networking) from the networking course linked below

- Slides are cross-uploaded in the Canvas for the Storage course
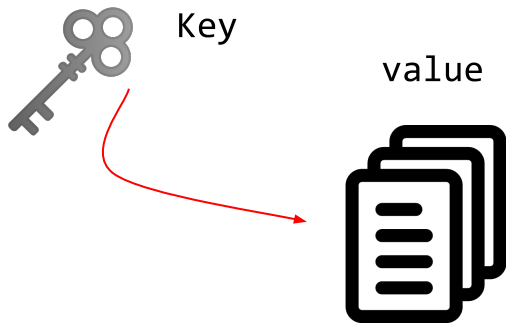- The course page, Advanced Network Programming

# Syllabus outline

1. ~~Welcome and introduction to NVM (today)~~
2. ~~Host interfacing and software implications~~
3. ~~Flash Translation Layer (FTL) and Garbage Collection (GC)~~
4. ~~NVM Block Storage File systems~~
5. NVM Block Storage Key-Value Stores ⬅
6. Emerging Byte-addressable Storage
7. Networked NVM Storage
8. Trends: Specialization and Programmability
9. Distributed Storage / Systems - I
10. Distributed Storage / Systems - II

# So, What is a Key-Value Store

A simplified data structure to store data and identify with a key

Key

value

Examples: associate arrays (array?), dictionaries, hash table

Quite popular with web, scalable services

Isn't a file system suppose to store our data?
- FSes create new files, directories for every object
- Web objects are often small, but basic file system inode overheads per directory/files
  - inodes can be a few kBs, if you want to store 64 bytes of data?
- Files/directories are difficult to iterate over quickly
- Range based queries need further indexing
- Object stores can support flexible consistent models (with FSes typically is a bad idea)
- Performance and feature optimizations, e.g., deduplication, transactions, compression, etc.

# Basic Operations

`put(key, value)`     : saves a value associated with a key

`value = get (key)` : retrieve the value associated with a key

`delete(key)`          : deletes a key (can be equivalent of `put(key, NULL)`)

Batch versions of these commands: `multiget, multiput`

Range based queries: `iterate (start_key, end_key);`

Further helper commands: `replace, add, incr, decr, etc.`

*No single data structure can do all operations efficiently*

# Layout of Coming Slides

**B+ Trees** and what they are good for

- What you need to do for storing them efficiently on NAND flash

**LSM** tree based KV design

- The basic idea
- LSM trees on OCSSD
- Application amplification in LSM

[Optional] A **Hash table**-based KV design (see Backup slides)

- FlashStore (and general topic of {memory $\longleftrightarrow$ I/O} tradeoff)

# B+ Tree

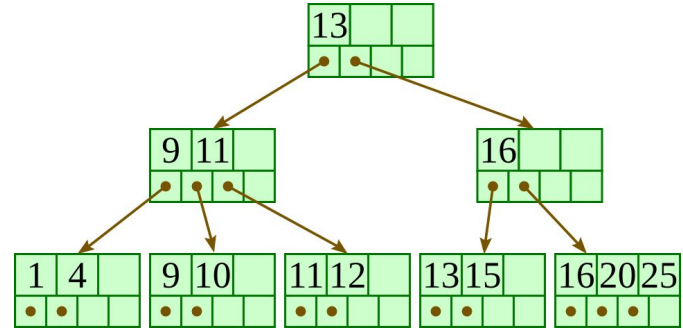M-ary tree with sorted (keys-values) stored in leaves

Useful for block-storage devices as it facilitate on-demand node fetching from the storage in block granularity

d-order tree has "d" keys and (d+1) pointers in non-leaf nodes, non-leaf nodes only contains "keys" for pivoting

Self-balancing (by splitting and merging nodes) and distance to all leaves nodes are equal from the root
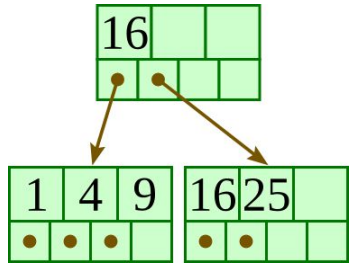
Popular data structure, used in Databases (Oracle, SQL) and file systems (ext4)

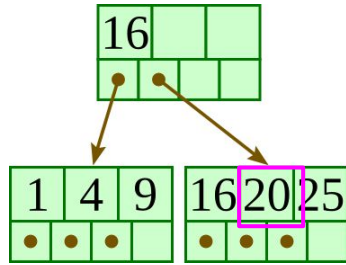Optimized for read-heavy workloads (sorted indexes)
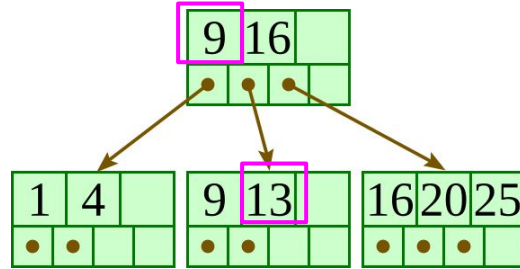


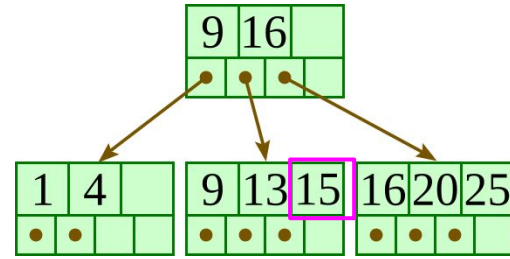http://www.cburch.com/cs/340/reading/btree/index.html
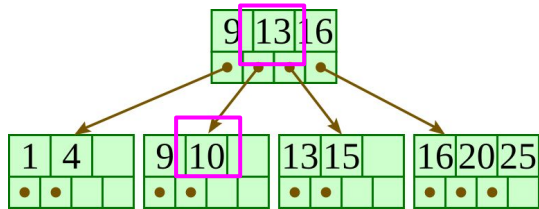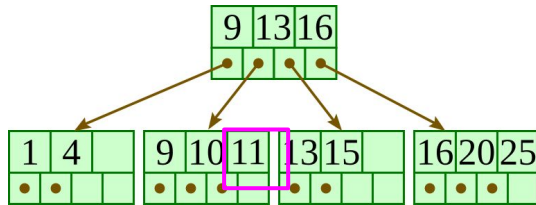
7

# Example: B+ Tree Insertions



Initial

Insert 20

Insert 13

Insert 15

Insert 10

Insert 11

Insert 12

http://www.cburch.com/cs/340/reading/btree/index.html

# Example: B+ Tree Insertions on NAND Flash



NAND flash pages, the same layout used with HDD too
- Whole pages can be read in a single go
- Large sequential transfers, good performance
- All values sorted, so we know which page to load for which node

http://www.cburch.com/cs/340/reading/btree/index.html

# Example: B+ Tree Insertions on NAND Flash

16 | | | | p0

Insert 20

Pick a new page
Copy

1 | 4 | 9 | | 16 | 25 | |

p1 | p2

16 | 20 | 25 | |

p3

NAND pages cannot be in-place updated

# Example: B+ Tree Insertions on NAND Flash

16 | | | **p0**
1 | 4 | 9 | | 16 | 25 | |
**p1** **p2**

**Insert 20**

Pick a new page
Copy

16 | 20 | 25 | |
**p3**

16 | | | **p0**
1 | 4 | 9 | | 16 | 25 | |
**p1** **p2**

16 | 20 | 25 | |
**p3**

*Now we need to update the root page too*

# Example: B+ Tree Insertions on NAND Flash



**Insert 20**

Pick a new page
Copy

*Now we need to update the root page too*

For a simple value insertion we ended up writing 2 new pages (p3 and p4) and generating 2 old (p0 and p2), invalid pages

In general, for a tree "H" height, we will have to update all pages on the path to the root, "H" nodes

It's the same problem what we saw in Log-Structured FS (recursive update problem or also known as Wandering Tree problem)

http://www.cburch.com/cs/340/reading/btree/index.html

# B+ Trees on NAND Flash



## $\mu$-Tree : An Ordered Index Structure for NAND Flash Memory

Dongwon Kang    Dawoon Jung    Jeong-Uk Kang    Jin-Soo Kim
Computer Science Division
Korea Advanced Institute of Science and Technology (KAIST)
Daejeon 305-701, Korea
{dwkang,dwjung,ux}@camars.kaist.ac.kr    jinsoo@cs.kaist.ac.kr

## ABSTRACT

As NAND flash memory becomes increasingly popular as data storage for embedded systems, many file systems and database management systems are being built on it. They require an efficient index structure to locate a particular item quickly from a huge amount of directory entries or database records. This paper proposes $\mu$-Tree, a new ordered index structure tailored to the characteristics of NAND flash memory. $\mu$-Tree is a balanced tree similar to B$^+$-Tree. In $\mu$-Tree, however, all the nodes along the path from the root to the leaf are put together into a single flash memory page in order to minimize the number of flash write operations when a leaf node is updated. Our experimental evaluation shows that $\mu$-Tree outperforms B$^+$-Tree by up to 28% for traces extracted from real workloads. With a small in-memory cache of 8 Kbytes, $\mu$-Tree improves the overall performance by up to 90% compared to B$^+$-Tree with the same cache size.

## Categories and Subject Descriptors

H.3.1 [**Content Analysis and Indexing**]: Indexing methods; D.4.3 [**File Systems Management**]: Directory structures

## General Terms

Algorithms, Design, Performance

## Keywords

B$^+$-Tree, NAND Flash, index structure

## 1. INTRODUCTION

Flash memory is being widely adopted as a storage medium for many portable embedded devices such as PMPs (portable media players), PDAs (personal digital assistants), digital cameras and camcorders, and cellular phones. This is mainly due to the inherent advantageous features of flash memory: non-volatility, small and lightweight form factor, low-power consumption, and solid state reliability.
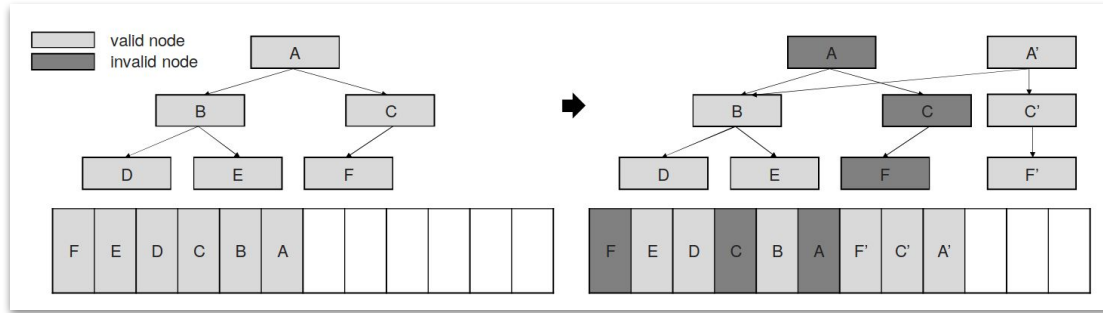
Flash memory comes in two flavors. The NOR type is usually used for storing codes since it can be directly addressable by processors. On the other hand, the NAND type is accessed on a page basis (typically 512 bytes $\sim$ 4 Kbytes) and provides higher cell densities. The NAND type is primarily used for removable flash cards, USB thumb drives, and internal data storage in portable devices.

As the NAND flash technology development continues to double density growth on an average of every 12 months [23], the capacity of a single NAND chip is getting larger at an increasingly lower cost. The declining cost of NAND flash memory has made it a viable and economically attractive alternative to hard disk drives especially in portable embedded systems. As a result, many flash-aware file systems and embedded database management systems (DBMSs) are currently being built on NAND flash memory [2, 7, 9, 13, 24].
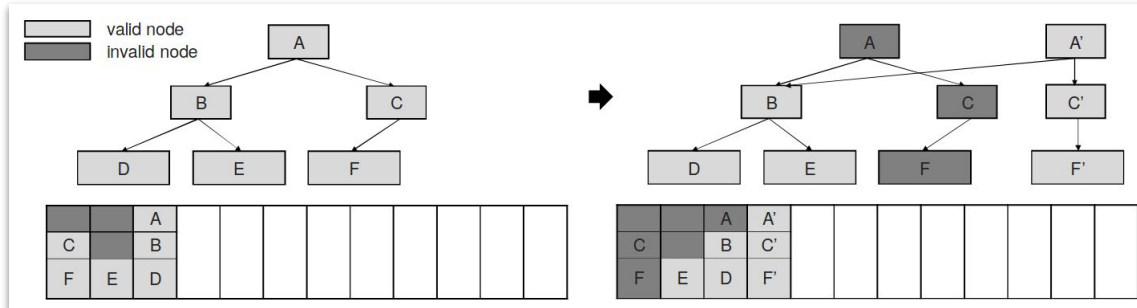
Any file system or DBMS requires an efficient index structure to locate a particular item quickly from a huge amount of directory entries or database records. For small scale systems, the index information can be kept in main memory. For example, JFFS2 keeps the whole index structures in memory that are necessary to find the latest file data on flash memory [24]. Apparently, this approach is not scal-

13

# μ-Tree : The Basic Idea

Rearrange the layout, do not give each nodes its own page. Store multiple nodes on a single page: typically along the path which will be update in case of an insertion



*Basic ("N" writes)*

*Proposed (update in 1 write)*

# How to Pack Nodes in a Page

Should we equally divide space in a page to all levels

Keeps the logic simple, and searchable, we will know exactly which offset in a page a level starts

However,

- Then we need to "fix" the maximum height of the tree
- Key space exponentially increases at every level
  - L3 : 2 order tree with 3 pointers
  - L2 : 3 x 3
  - L1 : 3 x 3 x 3 pointers ← this contain data, so we need to proportionally distribute space for different levels with flexibility to increase the level as we increase (or decrease the size of the tree)
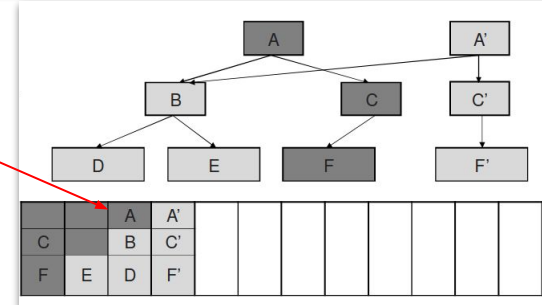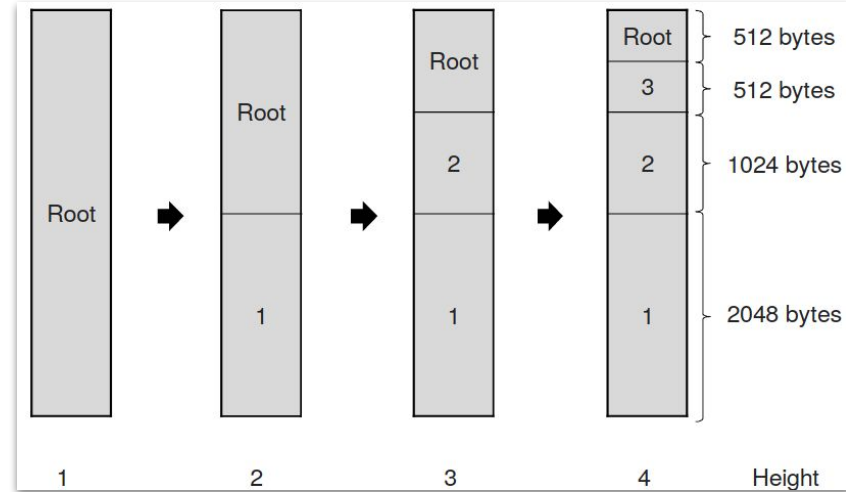
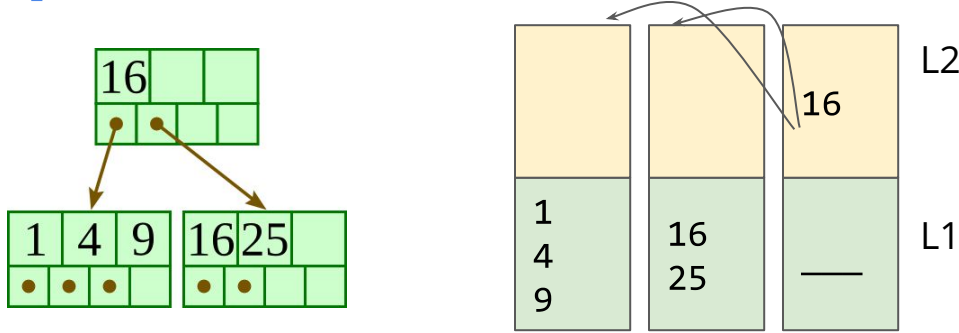| L3 |
| L2 |
| L1 |

# µ-Tree: Proportional Packing

In this setup

- Nodes within a page are still searchable
  - For a given level, and the height of the tree I can calculate which offset the node data starts
- Proportionally distribute space to different levels
- Enables us to do updates in one go, while keeping some date in old pages

The only thing we need to keep track of which page contains the "Root" pointer
- Changed from p2 to p3

# μ-Tree Insertions on NAND Flash
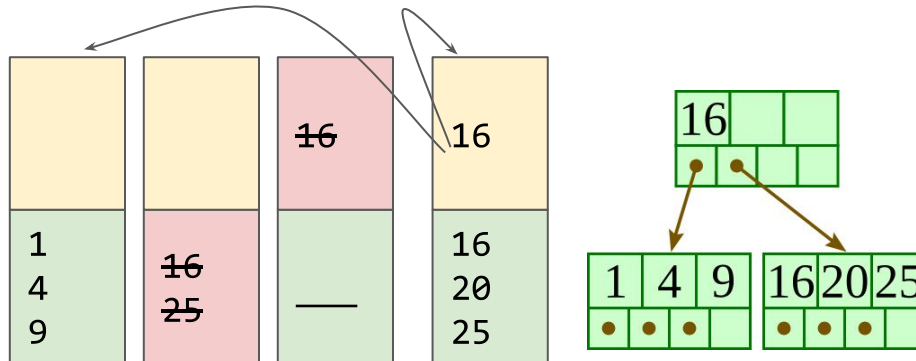
16 | | |
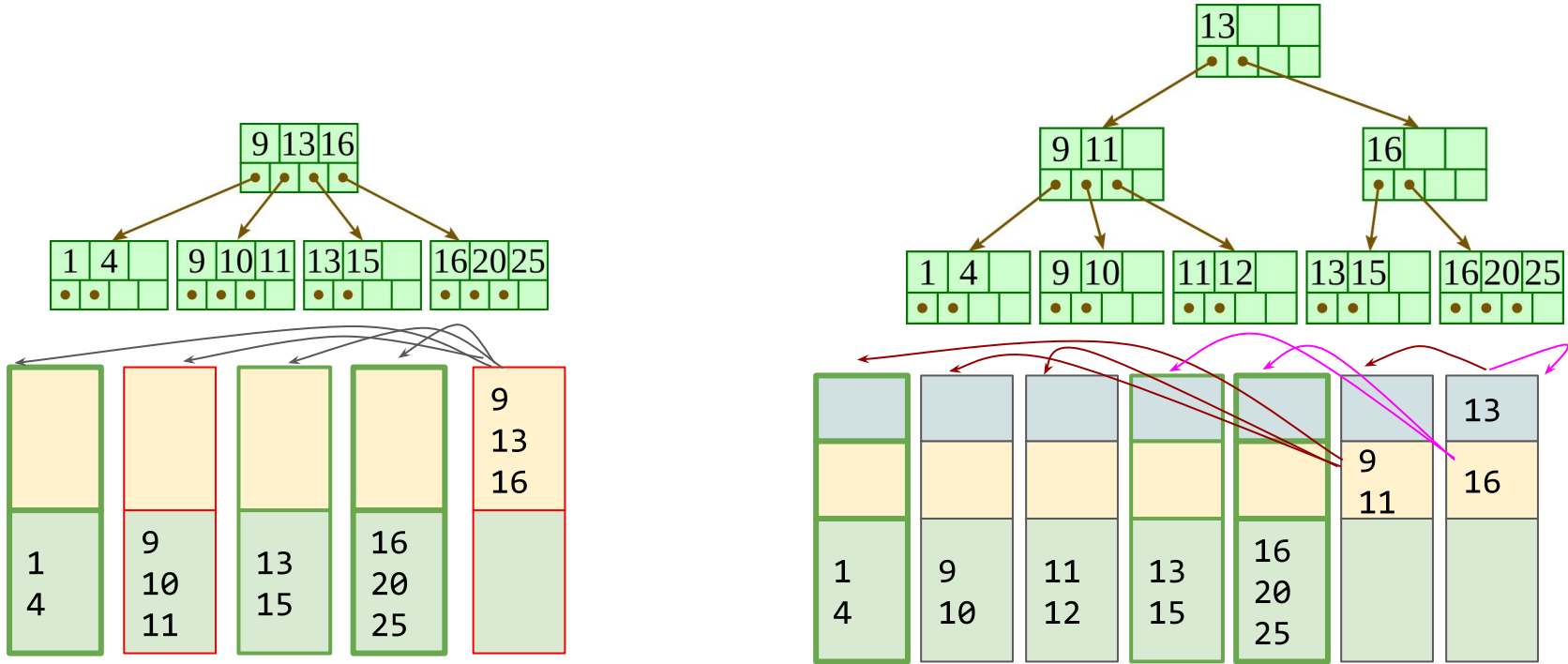1 | 4 | 9 | 16 | 25

L2

16

L1

1
4
9

16
25

___

In this case:
- 2 pages reading
- 1 page writing

In general: H x reading + 1 x writing

**Insert 20**

16

16

1
4
9

~~16~~
~~25~~

___

16
20
25

16 | | |
1 | 4 | 9 | 16 | 20 | 25

# μ-Tree Insertions with Height Increase



*Eventually as you write more, things will be grouped together (the update path) on the same page blocks. A similar logic applies to deletion and tree compaction logic (skipped).*

18

# μ-Tree: Performance (analytical)

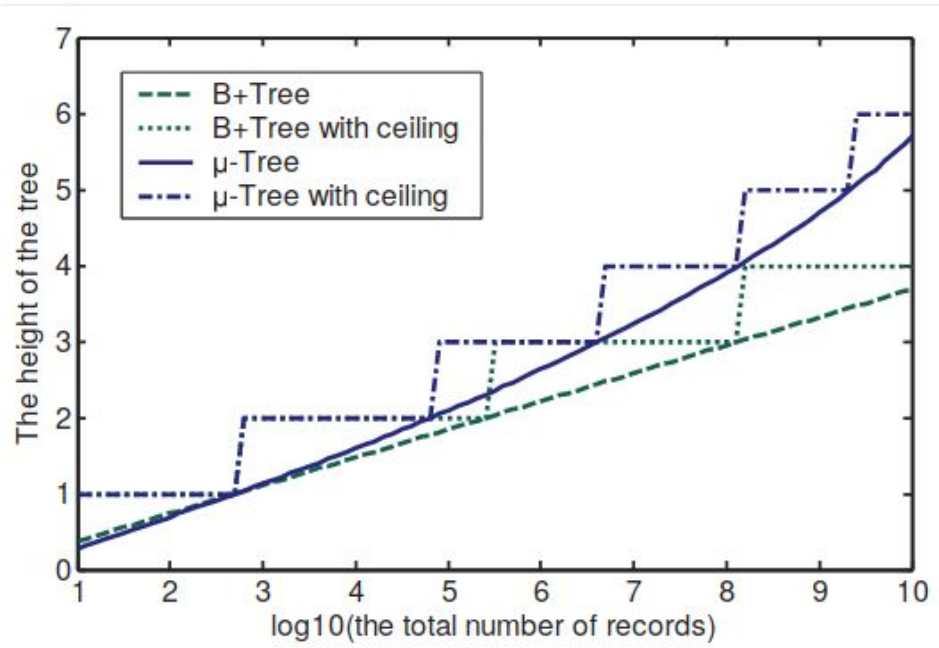Since the number of pointers that can be stored in a single page for a given level is different for μ and B+ Trees
- Height difference, within +1 (upto 1B)
- Takes twice as much flash space

Will results in more reads



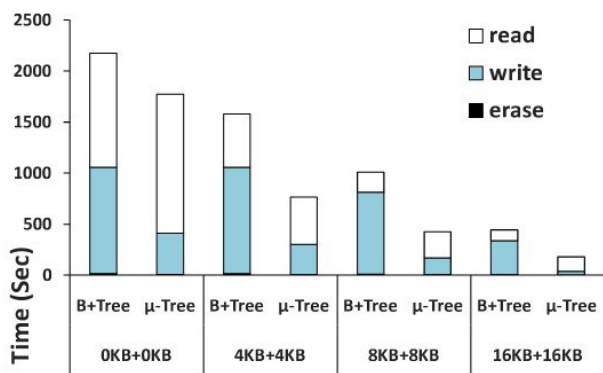| Table 3: The cost of operations | | |
|---|---|---|
| Operations | $B^+$-Tree | $\mu$-Tree |
| Retrieval | $c_r h_B$ | $c_r h_\mu$ |
| Insertion | $(c_r + c_w) h_B$ | $c_r h_\mu + c_w$ |
| Deletion | $(c_r + c_w) h_B$ | $c_r h_\mu + c_w$ |

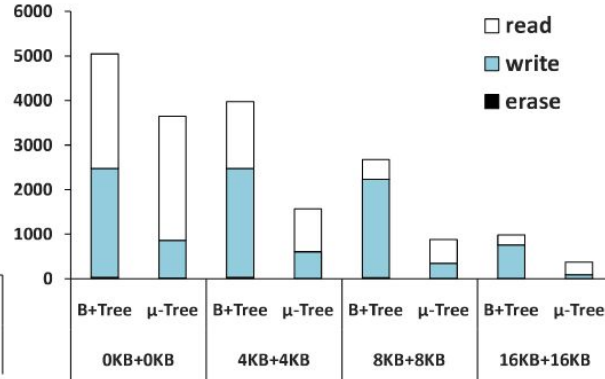*In absence of a split or collapse*

# μ-Tree: Performance

Traces collected from ReiserFS (B+ tree) about node creation, access, deletions
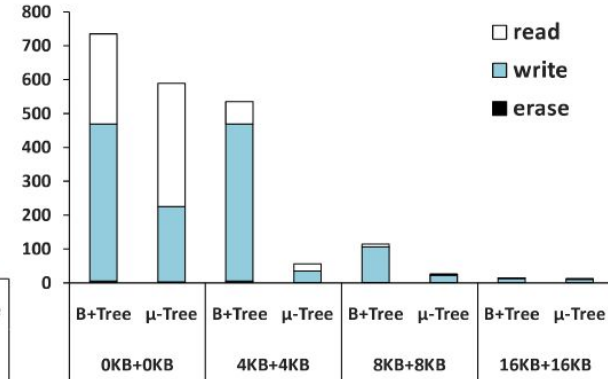
*Could have used some other benchmarks (well!)*



(a) kernel_compile
(b) postmark
(c) mp3

Better performance : decreases the number of writes and with more reads (taller tree)

# There are other works too

## An Efficient B-Tree Layer Implementation for Flash-Memory Storage Systems

CHIN-HSIEN WU and TEI-WEI KUO
National Taiwan University
and
LI PING CHANG
National Chiao-Tung University

With the significant growth of the markets for consumer electronics and various embedded systems, flash memory is now an economic solution for storage systems design. Because index structures require intensively fine-grained updates/modifications, block-oriented access over flash memory could introduce a significant number of redundant writes. This might not only severely degrade the overall performance, but also damage the reliability of flash memory. In this paper, we propose a very different approach, which can efficiently handle fine-grained updates/modifications caused by B-tree index access over flash memory. The implementation is done directly over the flash translation layer (FTL); hence, no modifications to existing application systems are needed. We demonstrate that when index structures are adopted over flash memory, the proposed methodology can significantly improve the system performance and, at the same time, reduce both the overhead of flash-memory management and the energy dissipation. The average response time of record insertions and deletions was also significantly reduced.

## FlashDB: Dynamic Self-tuning Database for NAND Flash

Suman Nath
Microsoft Research
sumann@microsoft.com

Aman Kansal
Microsoft Research
kansal@microsoft.com

**ABSTRACT**

FlashDB is a self-tuning database optimized for sensor networks using NAND flash storage. In practical systems flash is used in different packages such as on-board flash chips, compact flash cards, secure digital cards and related formats. Our experiments reveal non-trivial differences in their access costs. Furthermore, databases may be subject to different types of workloads. We show that existing databases for flash are not optimized for *all* types of flash devices or for *all* workloads and their performance is thus suboptimal in many practical systems. FlashDB uses a novel self-tuning index that dynamically adapts its storage structure to workload and underlying storage device. We formalize the self-tuning nature of an index as a two-state task system and propose a 3-competitive online algorithm that achieves the theoretical optimum. We also provide a framework to determine the optimal size of an index node that minimizes energy and latency for a given device. Finally, we propose optimizations to further improve the performance of our index. We prototype and compare different indexing schemes on multiple flash devices and workloads, and show that our indexing scheme outperforms existing schemes under *all* workloads and flash devices we consider.

example includes sensor networks of mobile devices which have significant local processing power [4, 12]. In these cases rather than uploading the entire raw data stream, one may save energy and bandwidth by processing queries locally at a cluster-head or a more capable node and uploading only the query response or the compressed or summary data. Storage centric networks have also been discussed in [6, 7].

In most cases where the storage is part of the sensor network, the storage device used is flash based rather than a hard disk due to shock resistance, node size, and energy considerations. Additionally, flash is also common in many mobile devices such as PDA's, cell-phones, music players, and personal exercise monitors. These devices can benefit from a having light weight database.

Our objective is to design storage and retrieval functionality for flash storage. A simple method is to archive data without an index, and that is in fact efficient in many scenarios. However, as we show in section 6, for scenarios where the number of queries is more than a small fraction ($\approx 1\%$) of the number of data items, having an index is useful. Hence, we focus on indexed storage. Prior work on flash storage provides file systems (e.g., ELF [5]) and other useful data structures such as stacks, queues and limited indexes (e.g., Capsule [14], MicroHash [22]). Our goal is to extend the functionality provided by those methods to B+-tree based indexing to support useful queries such as lookups, range-queries, multi-dimensional range-queries, and joins.

Existing database products are not well suited for sensor networks due to several reasons. Firstly, existing products, including
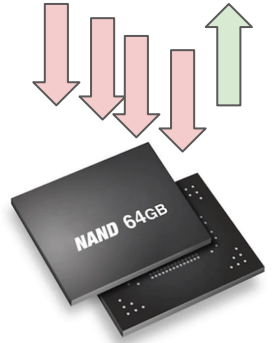
# Now, what about write-heavy workloads?

Write heavy workloads on flash can be really bad
- Key-Values can be really small (32-64-128 bytes). Hash checksums are also small (512/1024 bits)

The best solution so far we have seen is a log (FTL, file system)
- Append small writes to a log and read from there (search)

How can we improve searching the log?
- We can build a hash table (key) → {flash offset}
  - But will need a lot of memory for the hash table
    - Simply 8 bytes per key (similar to FTL)
- Does not allow doing fast range-based queries and lookups   ideas?

# Back to the Future: LSM Trees

**Log-Structured Merge (LSM)** Tree data structure
Invented and optimized for HDD, why?

- Same logic as LogFS
  - Disks have fast sequential performance
  - Disks have poor random, small I/O performance
- Read/Write large chunks to disk
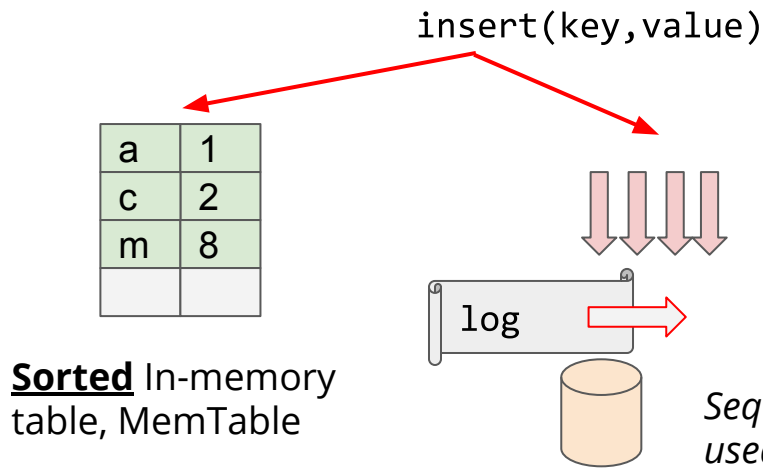- Eliminates random insertions, updates and deletions

*Patrick E. O'Neil, Edward Cheng, Dieter Gawlick, Elizabeth J. O'Neil:
The Log-Structured Merge-Tree (LSM-Tree). Acta Informatica 33(4):
351-385 (**1996**)*

Very popular data structure:  Bigtable, HBase, LevelDB,
SQLite4, Tarantool, RocksDB, WiredTiger, Apache Cassandra,
InfluxDB, and ScyllaDB



DOI:10.1145/3209210

Article development led by [acm]queue
queue.acm.org

**Different uses for read-optimized
B-trees and write-optimized LSM-trees.**

BY ALEX PETROV

# Algorithms Behind Modern Storage Systems

THE AMOUNTS OF data processed by applications are constantly growing. With this growth, scaling storage becomes more challenging. Every database system has its own trade-offs. Understanding them is crucial, as it helps in selecting the right one from so many available choices.

Every application is different in terms of read/write workload balance, consistency requirements, latencies, and access patterns. Familiarizing yourself with database and storage internals facilitates

when they arise, and fine-tunes the database for your workload.

It is impossible to optimize a system in all directions. In an ideal world there would be data structures guaranteeing the best read and write performance with no storage overhead but, of course, in practice that is not possible.

This article takes a closer look at two storage system design approaches used in a majority of modern databases —read-optimized B-trees[1] and write-optimized LSM (log-structured merge)-trees[2]—and describes their use cases and trade-offs.

https://queue.acm.org/detail.cfm?id=3220266

23

# LSM Tree Basics

`insert(key,value)`

| | |
|---|---|
| a | 1 |
| c | 2 |
| m | 8 |
| | |

**Sorted** In-memory table, MemTable

`log`

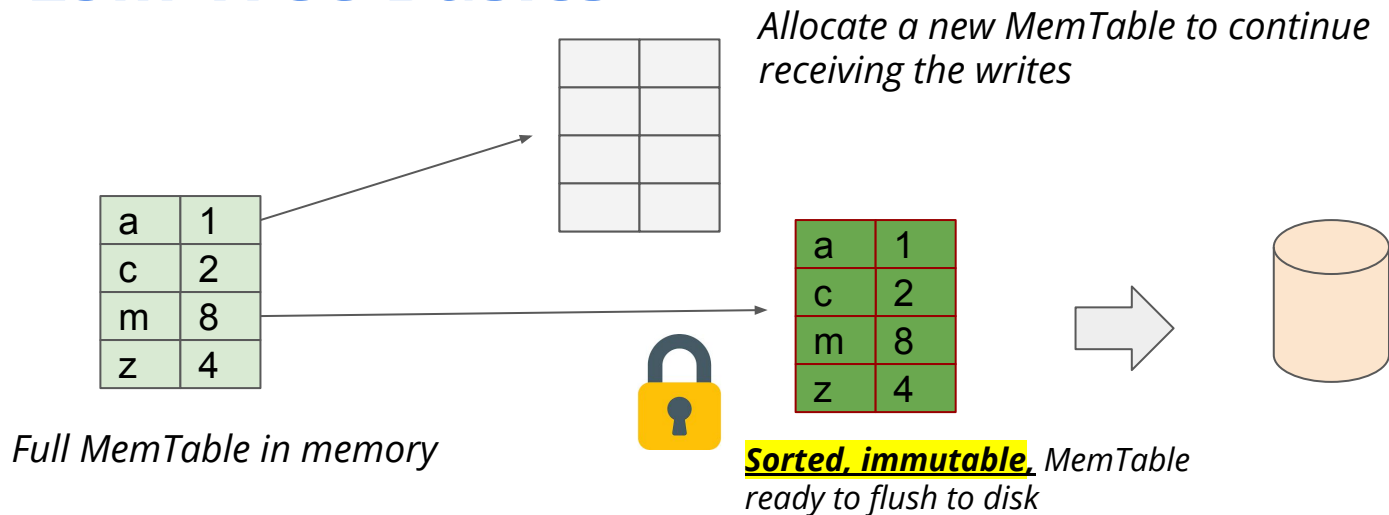*Sequential log on disk, only used in for failure recovery*

At insertion, (key,value) is
- written to the disk-resident log (large sequential performance)
- Inserted in the sorted MemTable to enable fast lookup with a range based query

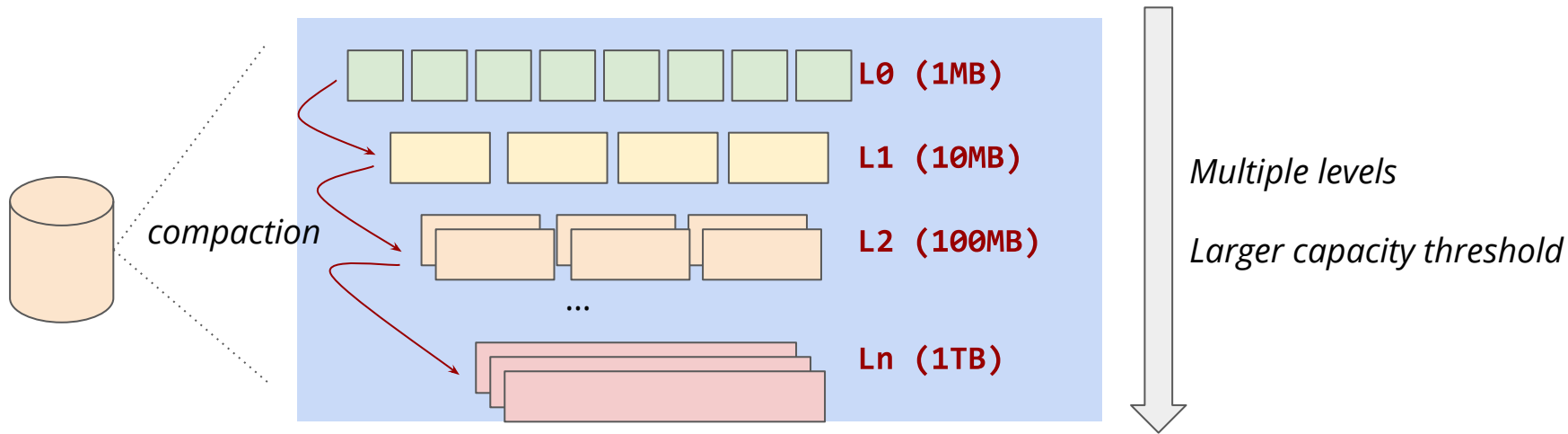What happens when the in-memory data structure is full?

# LSM Tree Basics

*Allocate a new MemTable to continue receiving the writes*



| a | 1 |
|---|---|
| c | 2 |
| m | 8 |
| z | 4 |

*Full MemTable in memory*

| a | 1 |
|---|---|
| c | 2 |
| m | 8 |
| z | 4 |

***Sorted, immutable,*** *MemTable ready to flush to disk*

- Once the in-memory table is full : the MemTable is marked immutable and *flushed* to disk
- Key lookup requires searching in the MemTable + looking up on the disk
  - *(we will see how this can be made efficient)*
- If data is present in both location, use timestamps to reconcile which is the newest write

Challenge now is how to (a) manage and (b) search TBs of data on disk to look for a key

# LSM Tree Basics



Data is stored in a multi-level, large, immutable files on the disk (no holes/gaps). Each level has a fixed size that increases as you go to the higher levels

A new table flush is written always written to **L0**

Just like in-memory table, once, a preconfigured size of file is reached, a files are level i can be merged with (i+1). This process is known as **compaction.** Since files written are sorted, the compaction is essentially an N-way merge sort from level (i) to (i+1)

# On-Disk File Format (SSTables)

**Sorted String Tables (SSTables)**

| Index | Bloom Filters | Data Block b0 | Data Block b1 | Data Block b2 |
|---|---|---|---|---|
| {K101:201, b0} | K101 → K201 | | | |
| {K321:350, b1} | K321 → K350 | | | |
| {K500:624, b2} | K500 → K624 | | | |
| {K876:900, b3} | K876 → K900 | | | |

When searching : find a value in the index range, then check in the bloom filter

Then go fetch the "block" for reading and scan the value inside

All files are immutables, hence, a delete is a new insertion with a "NULL" value

# Recap: Bloom Filters

Set key1



Bitmap or an array (any size)

Set key2

A bunch of hash function, `h1, h2, h3`

*collision*



0  1  2  3  4  5  6  7

Now if we were to check for key3 and key4
- lookup (key3) = 1, 4, 7 // all set, but they key3 was never set, <mark>false positive</mark>
- lookup (key4) = 0, 2, 5 // nope, this key was never set, always accurate! <mark>Cannot have false negative</mark>

The rate of false positive depends upon the size of the filter and the quality of the hash functions

For more fun read see https://blog.cloudflare.com/when-bloom-filters-dont-bloom/

# Example Compaction Process

**L0**  | 11, 21 |

**L1**

**L2**

# Example Compaction Process

*L0 can have duplicates keys in different files*

**L0**  `11, 21`    `11, 21`  `11, 13`                              `1, 15`    `12, 99`

**L1**                              `11, 13, 21`              `11, 13, 21`

**L2**                                                    *Pick all files which have overlapping ranges*

# Example Compaction Process

*L0 can have duplicates keys in different files*

**L0**    `11, 21`          `11, 21`    `11, 13`                          `1, 15`       `12, 99`

**L1**                                              `11, 13, 21`                `11, 13, 21`

**L2**                                                          *Pick all files which have overlapping ranges*

---

**L0**                                      `12, 99`        *You can check how many segments this compaction will touch*

**L1**        `1, 11, 13, 15`    `21`              `1, 11, 13, 15`    `21, 53, 65, 90`

**L2**

# Example Compaction Process

*L0 can have duplicates keys in different files*

**L0**   `11, 21`      `11, 21`   `11, 13`              `1, 15`      `12, 99`

**L1**                              `11, 13, 21`              `11, 13, 21`

**L2**

*Pick all files which have overlapping ranges*

**L0**                        `12, 99`     *You can check how many segments this compaction will touch*

**L1**   `1, 11, 13, 15`   `21`        →    `1, 11, 13, 15`   `21, 53, 65, 90`

**L2**

**L0**   `2, 16, 95`

**L1**   `1, 11, 12, 13`   `15, 21, 53, 65`   `90, 99`    →    `1, 2, 11, 12`   `13, 15, 16, 21`   `53, 65, 90, 95`

**L2**                                                                              `99`

32

# How to Optimize for Searching Files?

**Look in**: (i) mutable MemTable (ii) look at all the files at L0

- L0 files can contain overlapping key ranges, hence, **<u>all files</u>** need to be searched at **L0**

Further down, it can be a bit simpler as

- Files at L1 onwards **<u>do not have overlapping ranges</u>** (they are built that way)
- Hence, for each level, only need to check the range block and the bloom filter, not need to have read the file
- Lower levels contain fresher data (e.g., data at L3 would be newer than at L5)

Also, since indexes are sorted and immutable, support range-based queries

# General LSM Considerations

What are the size threshold for each level

What are the block sizes

When to do compaction
- Will result in decreasing the number of files
- Which level should be compacted to which next level
- Also: as L0 fills up the speed of writes will be stalled (in the end it will stop completely)

When to do garbage collection
- Deletion of old values which have been deleted
- Typically read the keys from the tree, and insert them back in the system

# An Efficient Design and Implementation of LSM-Tree based Key-Value Store on Open-Channel SSD (2014)



*Yes, this is your milestone 4 - just copy the design and implement it :)*
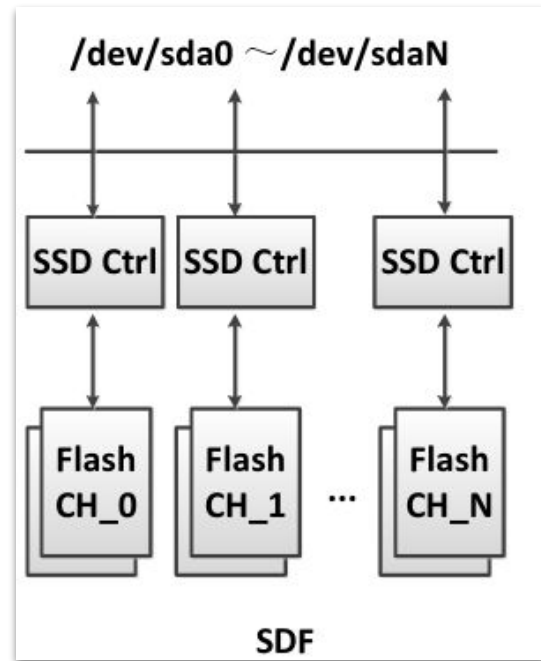
# Challenges with the Basic LSM Design

Recall: Open-Channel SSD is similar to SDF where all device internals and placement information is exposed -- high parallelism

1. Single head writing of immutable SSTable
2. Operation unaware scheduling (read, write, erase)
3. Placement and parallelism unaware scheduling

This work: **LOCS**
"LSM-tree-based KV store on Open-Channel SSD"

They retain the basic LSM design, but optimize it for OCSSD



/dev/sda0 ~/dev/sdaN

SSD Ctrl  SSD Ctrl  SSD Ctrl

Flash CH_0  Flash CH_1  ... Flash CH_N

SDF

# Idea 1: Enable Concurrent Accesses
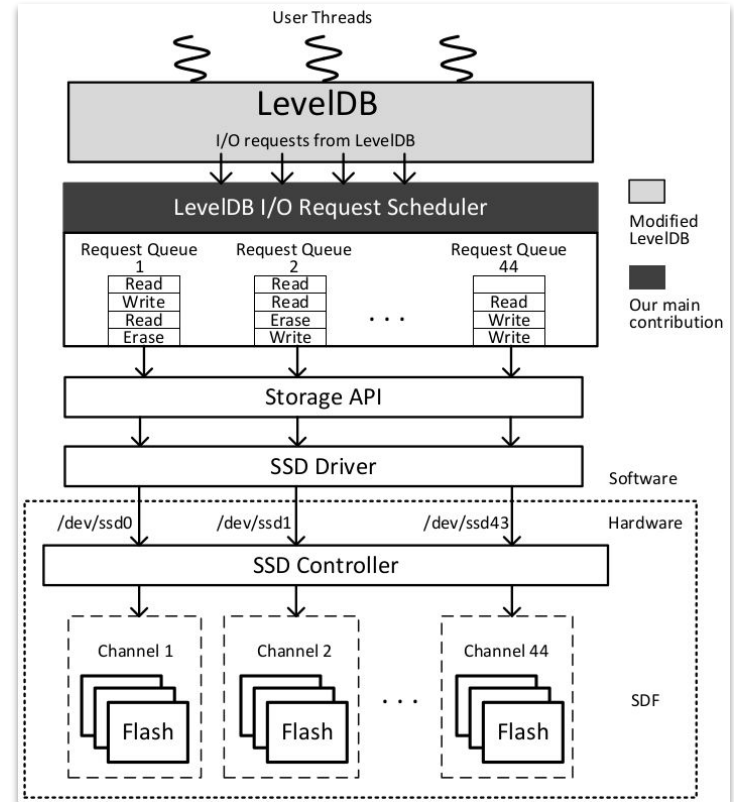
There is still a single <u>mutable</u> MemTable

Number of <u>immutable</u> in-memory MemTables are increased to 44
- Can absorb write bursts

Run multiple parallel compaction at the same time
- Was not possible with HDD because there is only single read/write head

# Idea 2: Scheduling Optimization

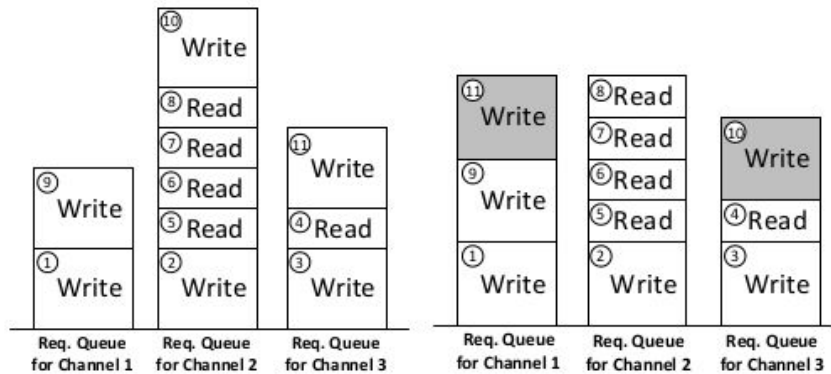Question: How should you pick which channel an SSTable should be flushed?
- Writes decides read workload too

**Strategy 1:** Round-Robin

**Strategy 2:** Least Weighted Queue Length Write dispatching
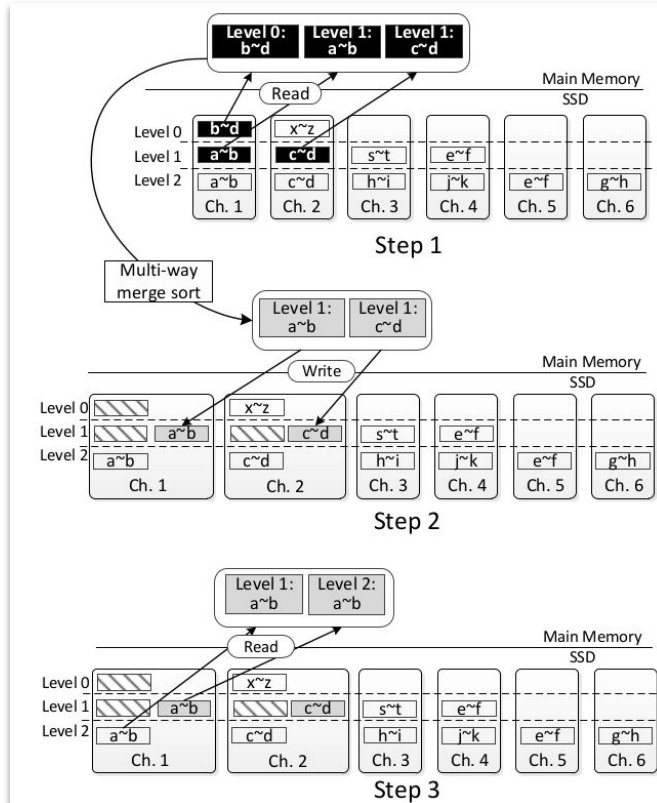
- Weight is read/write/erase cost

$$Length_{weight} = \sum_1^N W_i \times Size_i$$



(a) Round-Robin    (b) Least Weighted-Queue-Length

| Trace # | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---------|---|---|---|---|---|---|---|---|---|----|----|
| Op. | W | W | W | R | R | R | R | R | W | W | W |
| Channel | – | – | – | 3 | 2 | 2 | 2 | 2 | – | – | – |

(c) Trace of the example

# Idea 3: Placement Aware Compaction



Recall that LSM trees need compaction

**Here**: L0 file (b-d) is being pushed to L1

At L1 it overlaps with two files (a-b),(c-d)
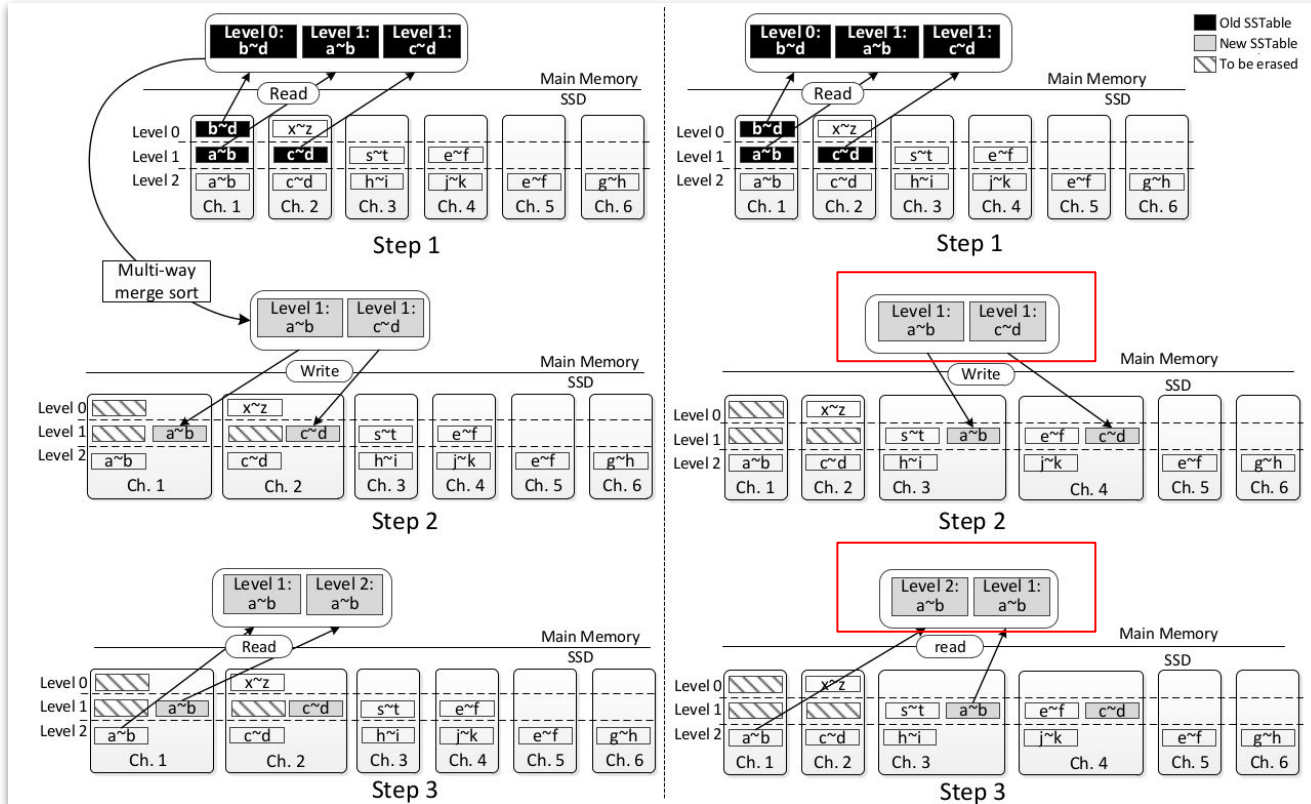
We first read those two files in DRAM

Do a multi-way merge sort with the three files

Then write out the L1 files (a-b) and (c-d)

Next-level of compaction at level L1 and L2 for key ranges of (a-b)

**Problem?**

# Idea 3: Placement Aware Compaction
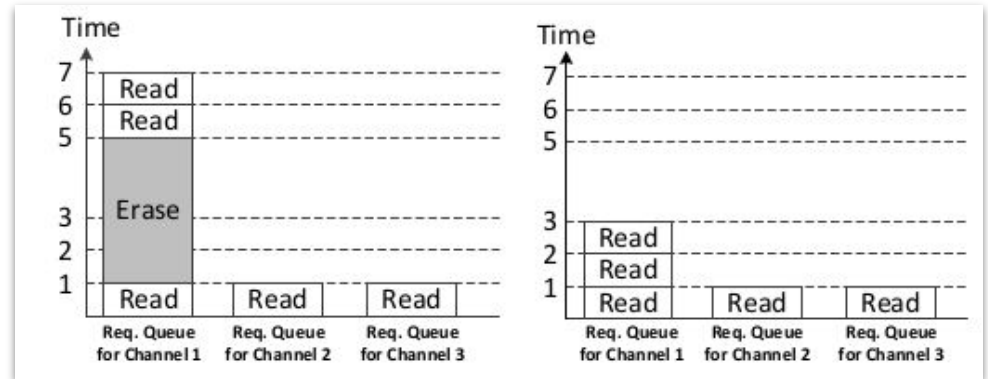
# Idea 4: Erase Aware Scheduling

Once the compaction is done, then one must erase blocks

Unlike read/write, erase can be scheduled by the KV when it is most opportune, when is that?

- Eager, as soon as possible

Erase is a long operation

Can lead to interferences with read operation (poor perf)



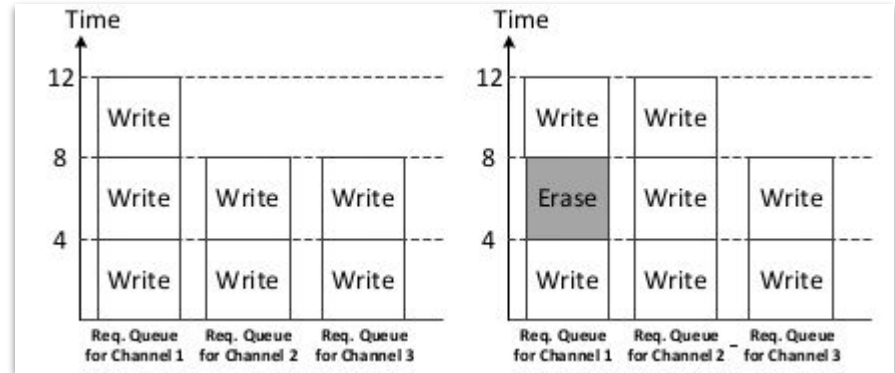*Eager scheduling of erase might be bad for read performance*

# Idea 4: Erase Aware Scheduling

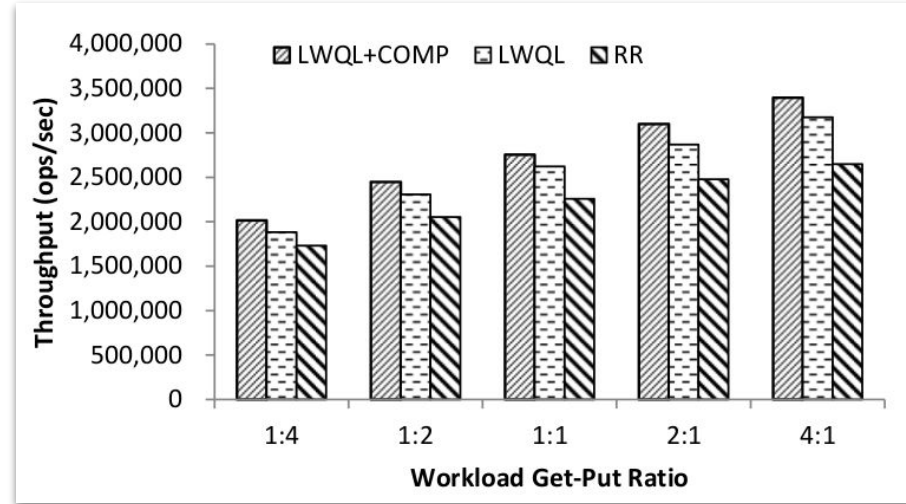The trick here is to schedule Erase with Writes, not with Read, why?
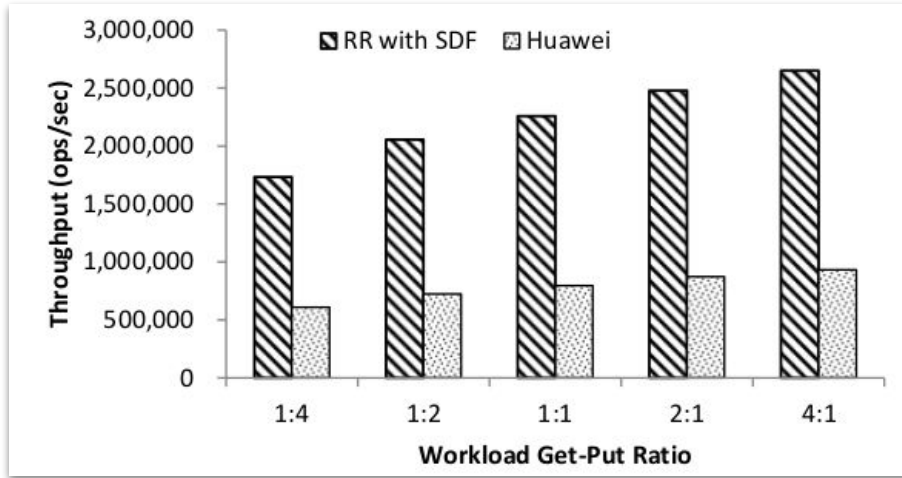
- Because writes can be put to any channel (flexible)
  - Reads cannot be moved around because they need to a read a given address from that channel
- Erase + Write can be used to balance out work among channels

In this example, we can insert Erase with write operations to maintain A balanced LWQL queue

E.g., with Erase in write it will take 19 units, where as Erase in write takes 15 units

# Performance: LOCS



Basic idea of software-managed parallelism over channels make sense

RR delivers good performance, LWQL even better, LWQL with Compaction aware optimizations the best of the three

43

# WiscKey: Separating Keys from Values in SSD-Conscious Storage (2016)

Lanyue Lu, Thanumalayan Sankaranarayana Pillai,
Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau

*University of Wisconsin, Madison*

## Abstract

We present WiscKey, a persistent LSM-tree-based key-value store with a performance-oriented data layout that separates keys from values to minimize I/O amplification. The design of WiscKey is highly SSD optimized, leveraging both the sequential and random performance characteristics of the device. We demonstrate the advantages of WiscKey with both microbenchmarks and YCSB workloads. Microbenchmark results show that WiscKey is 2.5×–111× faster than LevelDB for loading a database and 1.6×–14× faster for random lookups. WiscKey is faster than both LevelDB and RocksDB in all six YCSB workloads.

## 1 Introduction

Persistent key-value stores play a critical role in a variety of modern data-intensive applications, including web indexing [16, 48], e-commerce [24], data deduplication [7, 22], photo stores [12], cloud data [32], social networking [9, 25, 51], online gaming [23], messaging [1, 29], software repository [2] and advertising [20]. By enabling efficient insertions, point lookups, and range queries, key-value stores serve as the foundation for this growing group of important applications.

For write-intensive workloads, key-value stores based on Log-Structured Merge-Trees (LSM-trees) [43] have become the state of the art. Various distributed and local stores built on LSM-trees are widely deployed in large-scale production environments, such as BigTable [16] and LevelDB [48] at Google, Cassandra [33], HBase [29] and RocksDB [25] at Facebook, PNUTS [20] at Yahoo!, and Riak [4] at Basho. The main advantage of LSM-

throughout its lifetime; as we show later (§2), this I/O amplification in typical LSM-trees can reach a factor of 50x or higher [39, 54].

The success of LSM-based technology is tied closely to its usage upon classic hard-disk drives (HDDs). In HDDs, random I/Os are over 100× slower than sequential ones [43]; thus, performing additional sequential reads and writes to continually sort keys and enable efficient lookups represents an excellent trade-off.

However, the storage landscape is quickly changing, and modern solid-state storage devices (SSDs) are supplanting HDDs in many important use cases. As compared to HDDs, SSDs are fundamentally different in their performance and reliability characteristics; when considering key-value storage system design, we believe the following three differences are of paramount importance. First, the difference between random and sequential performance is not nearly as large as with HDDs; thus, an LSM-tree that performs a large number of sequential I/Os to reduce later random I/Os may be wasting bandwidth needlessly. Second, SSDs have a large degree of internal parallelism; an LSM built atop an SSD must be carefully designed to harness said parallelism [53]. Third, SSDs can wear out through repeated writes [34, 40]; the high write amplification in LSM-trees can significantly reduce device lifetime. As we will show in the paper (§4), the combination of these factors greatly impacts LSM-tree performance on SSDs, reducing throughput by 90% and increasing write load by a factor over 10. While replacing an HDD with an SSD underneath an LSM-tree does improve performance, with current LSM-tree technology, the SSD's true potential goes largely unrealized.

# So, What is the Problem?

We briefly referenced that reading
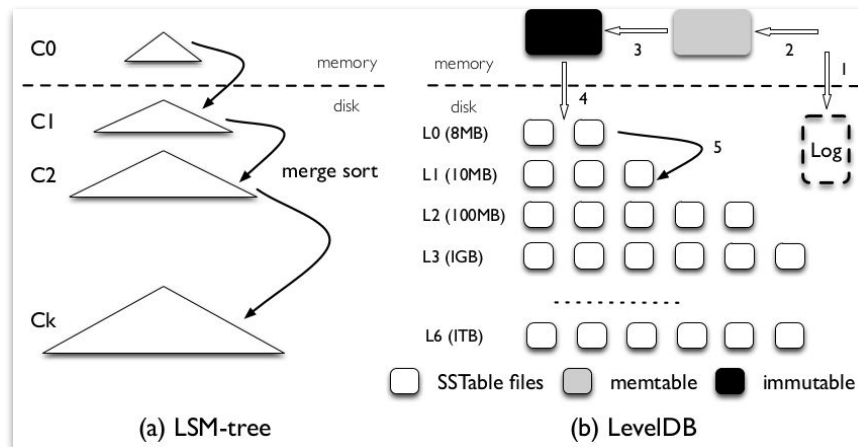performance on LSM can be problematic

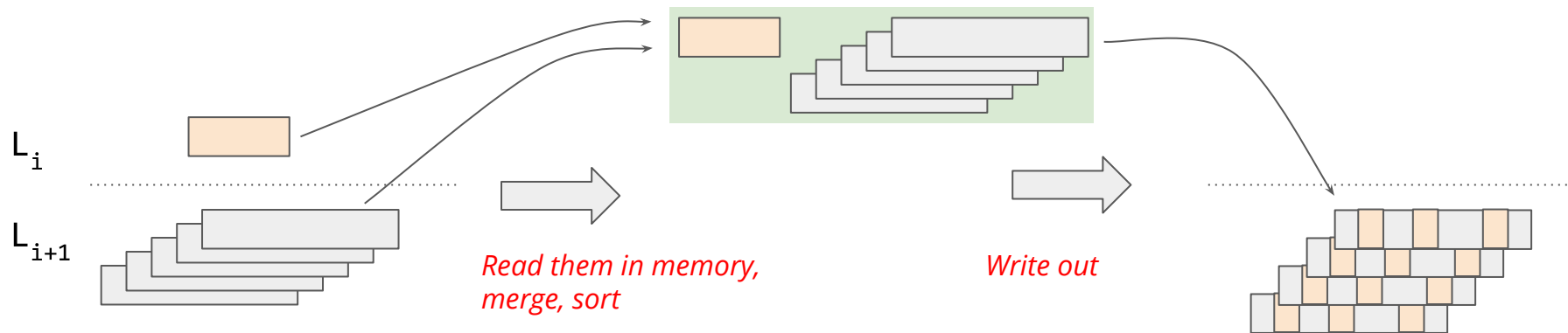<mark>Any guesses why?</mark>

What was the read path order?

- MemTable → L0 → L1 ... L6 (here)

So, if you were to read simple 1 byte key-value,
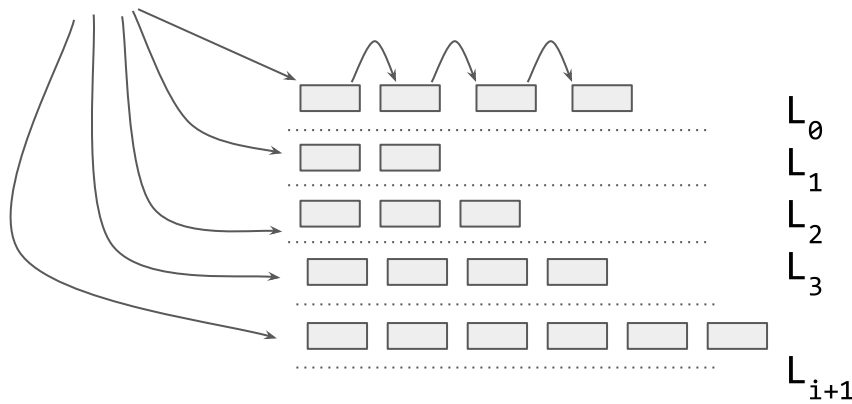how much data you have to read before you can find a 1 byte result?

We have looked this type of problem before in FTL for writes
(**recall**: write-amplification)



(a) LSM-tree    (b) LevelDB

# LSM has Read and Write Amplifications

$L_i$

$L_{i+1}$

*Read them in memory, merge, sort*

*Write out*

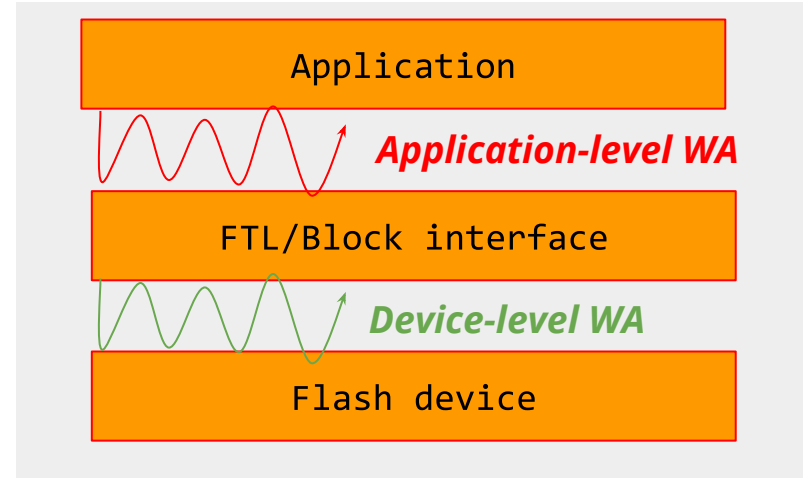read,lookup

$L_0$

$L_1$

$L_2$

$L_3$

$L_{i+1}$

# Analysis : Write/Read Amplification (RA/WA)

Compaction can result in
- Reading "n" times data from the next level to merge from the current level
  - For LevelDB this is 10x between levels
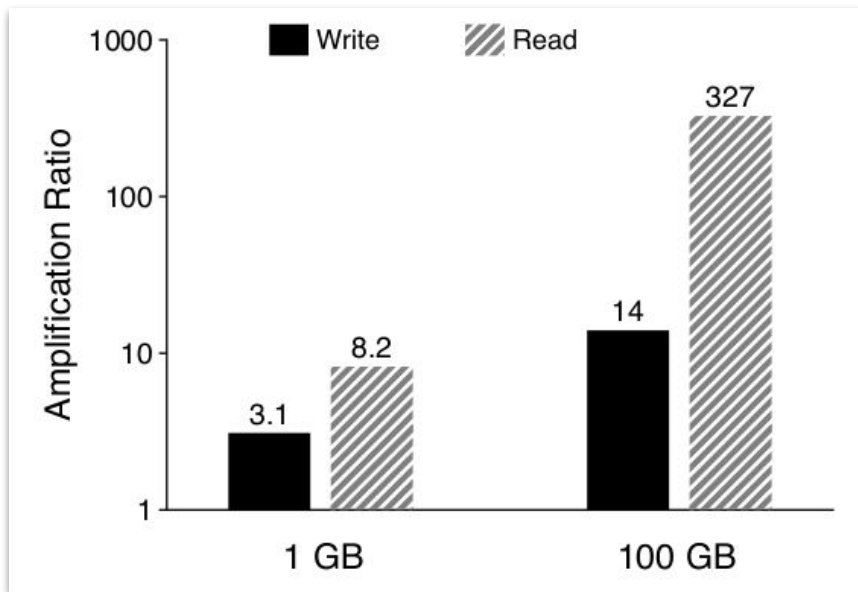  - For 6 levels, it could be 50x

Reading can result in
- Reading "n" files on L0 and then 1 file on following level
  - LevelDB, 8 files (at L0) + 6 files (L1-L6) = 14 files
    - Within the file we need to read the "index" + "bloom filter" + data block
    - For level DB index (16kB), bloom (4kB) + data (4kB)
    - So, if we are looking for a 1kB file: 14 files x (24 kb) = 336 kb ⇒ 336x **RA**
  - Determined by how many files do you have to touch and read to find a value

**LSM Trees trade high "amplification" for having "sequential performance" → Why this was ok with HDDs?**



Application

*Application-level WA*

FTL/Block interface

*Device-level WA*

Flash device

# Quantify and Justify
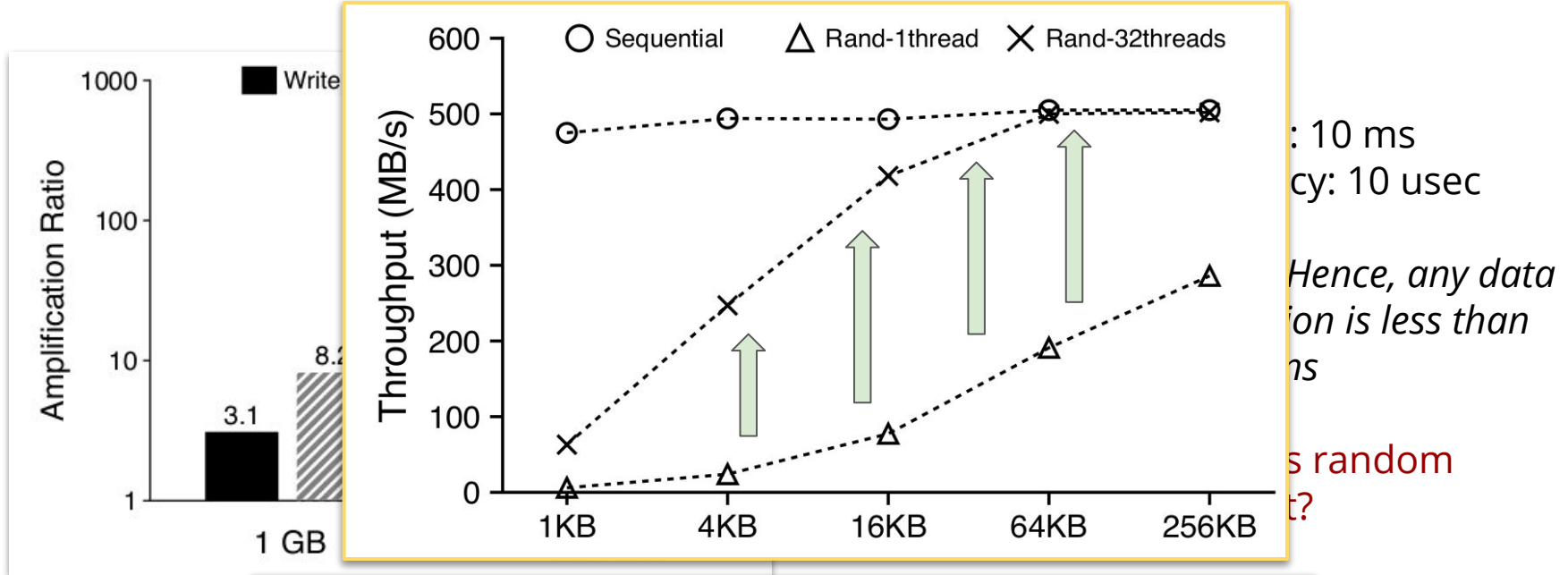


Key size: 16 bytes, value size : 1024 bytes

**Justification for HDD**
- Random 1kB latency: 10 ms
- Sequential 1kB latency: 10 usec

Ratio is seq:rand **1:1000**. *Hence, any data structure where amplification is less than 1000, sequential access wins*

On SSD? Are sequential vs random accesses are 1:1000 apart?

# Quantify and Justify



: 10 ms
cy: 10 usec
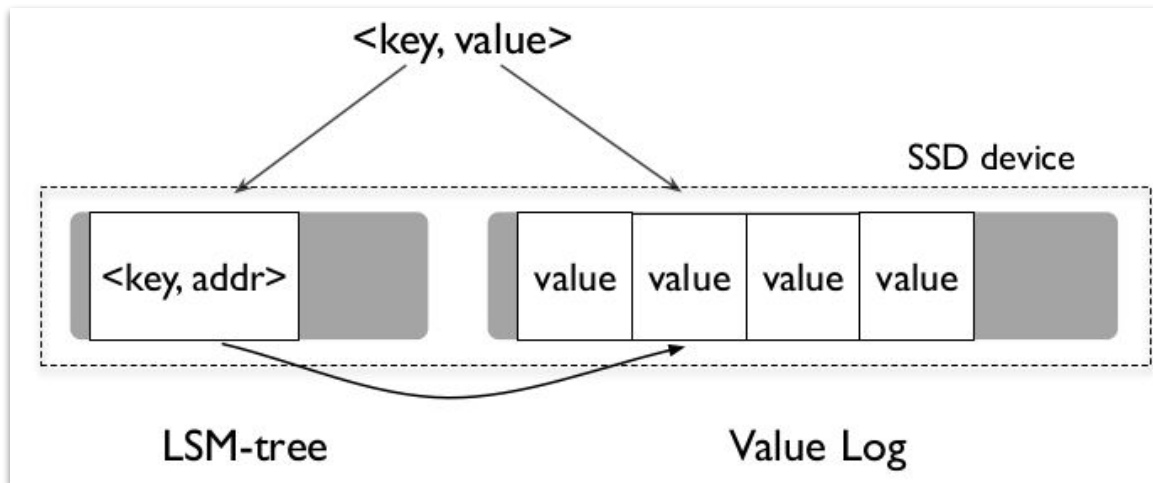
*Hence, any data
ion is less than*
*ns*

s random
t?

There exists a gap between random and sequential performance, but
- Not for large values
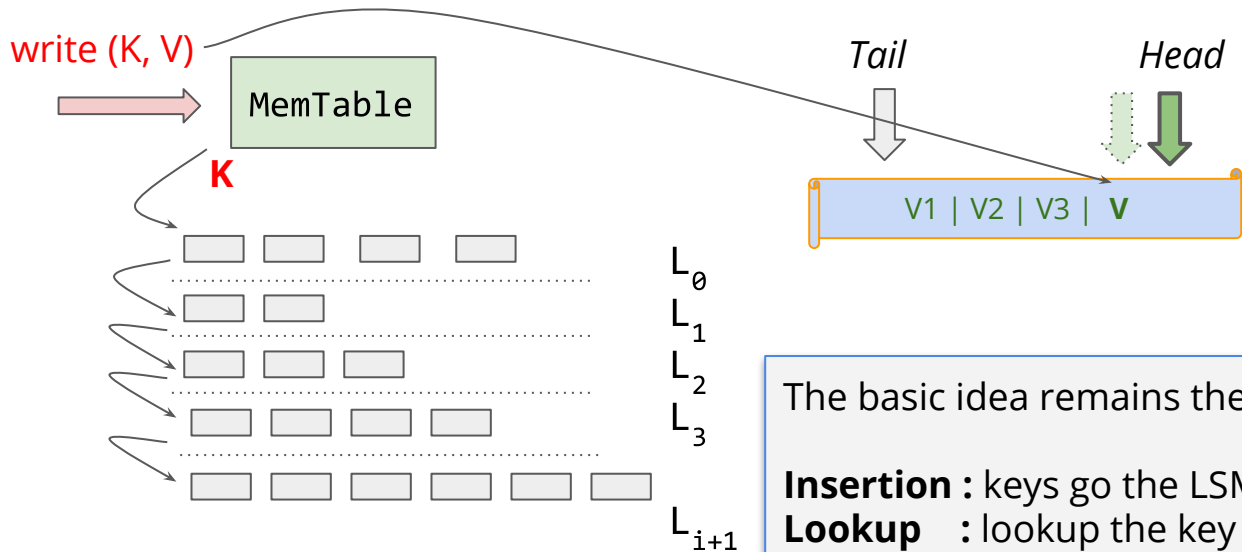- The gap can be closed by issuing multiple parallel requests

# What does WiscKey Proposes

**Key Idea:** separate keys from the values

- Maintain keys in the LSM tree
- Maintain value in a sequential append value log

# Key-Value Insertion and Lookup

write (K, V)

MemTable

**K**

$L_0$
$L_1$
$L_2$
$L_3$

$L_{i+1}$

*Tail*          *Head*

V1 | V2 | V3 | **V**

The basic idea remains the same

**Insertion :** keys go the LSM tree, values to the log
**Lookup    :** lookup the key in the LSM tree, then read the offset from the log

For **range-based queries**, the log can be read in parallel
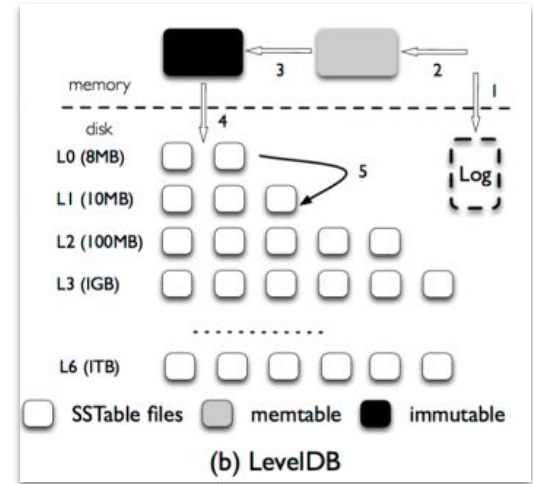
# WiscKey: LSM Tree made out of Keys

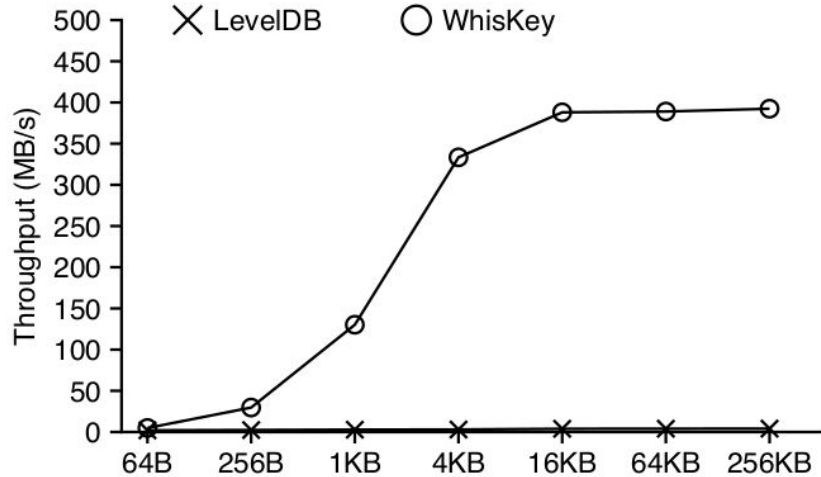What advantages a <u>key-only</u> LSM tree brings
- [with assumptions] keys are small and values are big
- Much improved write-amplification
  - Before WA was: ~10-50x
  - Now (10 x key_size) + value_size / (key + value size)
  - E.g., (10 x 16 + 1024) / (1024 + 16) = **1.14** (not 10x)
  - Worse case : (50 x 16 + 1024) / (1024+16) = **1.76** (not 50x)
- Lower write amplification means longer device life time
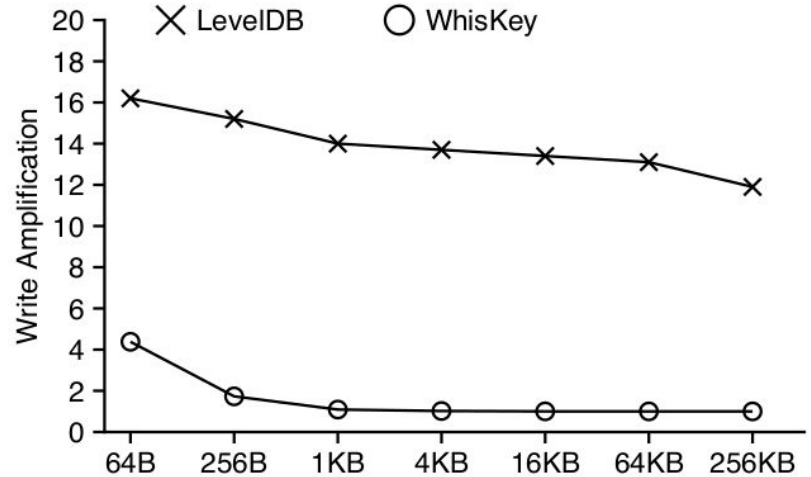
Also, the size of the tree can be small (small keys)
- Less levels than a comparable key-value LSM tree
- Small tree can be cached in the memory for fast lookups



(b) LevelDB

# WiscKey: Performance
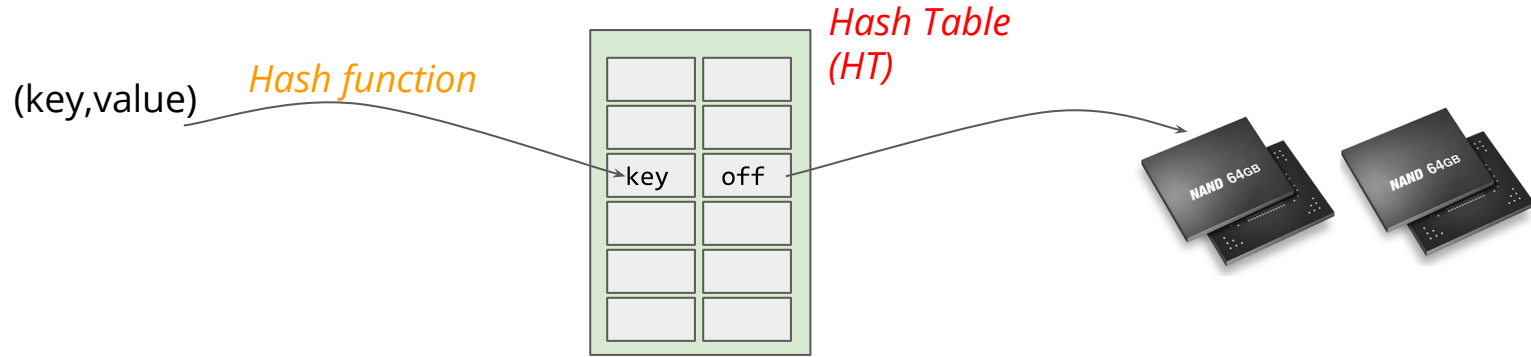


Key: 16B, Value: 64B to 256KB

Key: 16B, Value: 64B to 256KB

LevelDB is at 2-4MB/sec whereas WiscKey is at 350 MB/sec (46-111x)

Significantly better write amplification performance

# Hash Tables on Flash



(key,value) — *Hash function* → | key | off | — *Hash Table (HT)*

This simple hash table based schema works, but it needs to deal with
- Small writes (multiple writes must be packed together)
- Can do fast get and put, but no range-based queries (without additional indexes)
- Trade off {DRAM size of the HT } ← → {number of I/O operations}
  - The same tradeoff as FTL design, how much memory do we need to store a hash table with 1 TB of values
  - Can store the table in flash itself, to decrease the memory size, then multiple I/O

https://en.wikipedia.org/wiki/Hash_table#Choosing_a_hash_function

# Alternate Hash Table Designs

## SkimpyStash: RAM Space Skimpy Key-Value Store on Flash-based Storage

Biplob
Mi
EM

**ABSTRACT**

We present SkimpyStash, a RAM space s
on flas -based storage, designed for high t
server applications. The distinguishing fea
the design goal of extremely low RAM fo
0.5) byte per key-value pair, which is mor
lier designs. SkimpyStash uses a hash tabl
index key-value pairs stored in a log-struc
To break the barrier of a flas pointer (say, 4
overhead per key, it "moves" most of the p
key-value pair from RAM to flas itself. I
resolving hash table collisions using linear
tiple keys that resolve (collide) to the sam
chained in a linked list, and (ii) storing the
self with a pointer in each hash table buck
the beginning record of the chain on flash .
ple flas reads per lookup. Two further tec
prove performance: (iii) two-choice based l
wide variation in bucket sizes (hence, chain
lookup times), and a bloom filte in each h
in RAM to disambiguate the choice during
paction procedure to pack bucket chain rec
flas pages so as to reduce flas reads durin
bucket size is the critical design parameter t
ful knob for making a continuum of tradeo
usage and low lookup latencies. Our eval
server platforms with real-world data center
SkimpyStash provides throughputs from fe
of 100,000 get-set operations/sec.

## FlashStore: High Throughput Persistent Key-Value Store

Biplob Debnath*
University of Minnesota
Twin Cities, USA
biplob@umn.edu

Sudipta Sengupta
Microsoft Research
Redmond, USA
sudipta@microsoft.com

Jin Li
Microsoft Research
Redmond, USA
jinl@microsoft.com

**ABSTRACT**

We present FlashStore, a high throughput persistent key-value store, that uses flash memory as a non-volatile *cache* between RAM and hard disk. FlashStore is designed to store the working set of key-value pairs on flash and use one flash read per key lookup. As the working set changes over time, space is made for the current working set by destaging recently unused key-value pairs to hard disk and recycling pages in the flash store. FlashStore organizes key-value pairs in a log-structure on flash to exploit faster sequential write performance. It uses an in-memory hash table to index them, with hash collisions resolved by a variant of cuckoo hashing. The in-memory hash table stores compact key signatures instead of full keys so as to strike tradeoffs between RAM usage and false flash read operations.

FlashStore can be used as a high throughput persistent key-value storage layer for a broad range of server class applications. We compare FlashStore with BerkeleyDB, an embedded key-value store application, running on hard disk and flash separately, so as to bring out the performance gain of FlashStore in not only using flash as a cache above hard disk but also in its use of flash aware algorithms. We use real-world data traces from two data center applications, namely, Xbox LIVE Primetime online multi-player game and inline storage deduplication, to drive and evaluate the design of FlashStore on traditional and low power server platforms. FlashStore outperforms BerkeleyDB by up to 60x on throughput (ops/sec), up to 50x on energy efficiency (ops/Joule), and up to 85x on cost efficiency (ops/sec/dollar) on the evaluated datasets.

A high throughput persistent key-value store can help t
improve the performance of such applications. Flash mem
ory is a natural choice for such a store, providing persis
tency and 100-1000 times lower access times than hard disk
Compared to DRAM, flash access times are about 100 time
higher. Flash stands in the middle between DRAM and dis
also in terms of cost – it is 10x cheaper than DRAM, whil
20x more expensive than disk – thus, making it an ideal ga
filler between DRAM and disk.

There are two types of popular flash devices, NOR an
NAND flash. NAND flash architecture allows a denser lay
out and greater storage capacity per chip. As a resul
NAND flash memory has been significantly cheaper tha
DRAM, with cost decreasing at faster speeds. NAND flas
characteristics have lead to an explosion in its usage i
consumer electronic devices, such as MP3 players, phone
caches and Solid State Disks (SSDs). In the rest of the pa
per, we use NAND flash based SSDs as the architectura
choice and simply refer to it as flash memory. We describ
SSDs in detail in Section 2. To get the maximum perfo
mance per dollar out of SSDs, it is necessary to use flas
aware data structures and algorithms to avoid small randor
writes that not only have a higher latency but also reduc
flash device lifetimes through increased page wearing.

In this paper, we present the design and evaluation o
FlashStore, a high performance key-value storage system us
ing flash as a cache between RAM and hard disk. When a
key-value blob is written, it is sequentially logged in flash.
A specialized RAM-space efficient hash table index using a
variant of cuckoo hashing [32] and compact key signatures
is used to index the key-value blobs stored in flash mem-

## SILT: A Memory-Efficient, High-Performance Key-Value Store

Hyeontaek Lim[1], Bin Fan[1], David G. Andersen[1], Michael Kaminsky[2]

[1]Carnegie Mellon University, [2]Intel Labs

**ABSTRACT**

SILT (Small Index Large Table) is a memory-efficient, high-performance key-value store system based on flash storage that scales to serve billions of key-value items on a single node. It re-quires only 0.7 bytes of DRAM per entry and retrieves key/value pairs using on average 1.01 flash reads each. SILT combines new algorithmic and systems techniques to balance the use of memory, storage, and computation. Our contributions include: (1) the design of three basic key-value stores each with a different emphasis on memory-efficiency and write-friendliness; (2) synthesis of the basic key-value stores to build a SILT key-value store system; and (3) an analytical model for tuning system parameters carefully to meet the needs of different workloads. SILT requires one to two orders of magnitude less memory to provide comparable throughput to cur-rent high-performance key-value systems on a commodity desktop system with flash storage.

**Categories and Subject Descriptors**

D.4.2 [**Operating Systems**]: Storage Management; D.4.7 [**Operating Systems**]: Organization and Design; D.4.8 [**Operating Systems**]: Performance; E.1 [**Data**]: Data Structures; E.2 [**Data**]: Data Storage Representations; E.4 [**Data**]: Coding and Information Theory

**General Terms**

Algorithms, Design, Measurement, Performance

**Keywords**

Algorithms, design, flash, measurement, memory efficiency, performance

## 1. INTRODUCTION

Key-value storage systems have become a critical building block for today's large-scale, high-performance data-intensive applications.

| Metric | 2008 → 2011 | Increase |
|---|---|---|
| CPU transistors | 731 → 1,170 M | 60 % |
| DRAM capacity | 0.062 → 0.153 GB/$ | 147 % |
| Flash capacity | 0.134 → 0.428 GB/$ | 219 % |
| Disk capacity | 4.92 → 15.1 GB/$ | 207 % |

Table 1: From 2008 to 2011, flash and hard disk capacity increased much faster than either CPU transistor count or DRAM capacity.
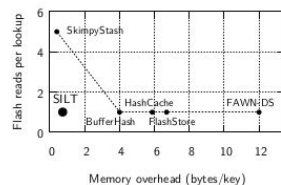
Figure 1: The memory overhead and lookup performance of SILT and the recent key-value stores. For both axes, smaller is better.

e-commerce platforms [21], data deduplication [1, 19, 20], picture stores [7], web object caching [4, 30], and more.
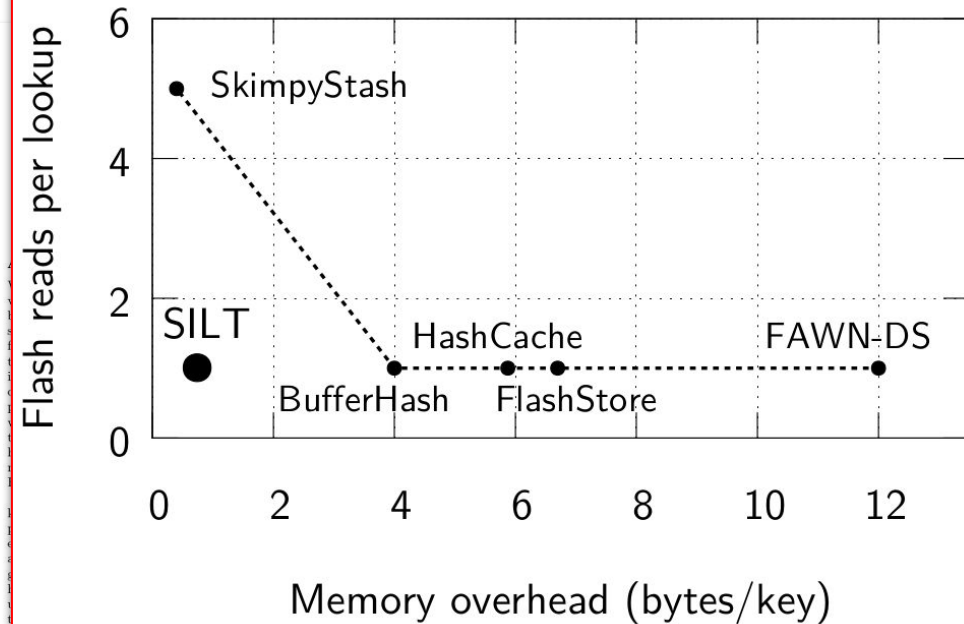
To achieve low latency and high performance, and make best use of limited I/O resources, key-value storage systems require efficient indexes to locate data. As one example, Facebook engineers recently created a new key-value storage system that makes aggressive use of DRAM-based indexes to avoid the bottleneck caused by multiple disk operations when reading data [7]. Unfortunately, DRAM is up to 8X more expensive and uses 25X more power per bit than flash, and as Table 1 shows, is growing more slowly than the capacity of

# Alternate Hash Table Designs

# Summary of Data Structures

- B+ Tree (read-optimized)
  - Fast, bounded lookup for read/get (log(n))
  - Efficient range based queries
  - But poor performance for write-heavy workloads, update bubbling (also small updates)

- Log-structured Merge (LSM) Tree (write-optimized)
  - Good performance for write-heavy workloads, large sequential log based updates
  - Ranged based queries possible
  - Read/Write amplification is a problem

- Simple hash table (hash like md5 on the key → map to a location)
  - [Typically uses] Log-based writing
  - Easy and fast lookup and retrieval (O(1))
  - Limited range based query support (need additional indexing)
  - Tradeoff between (memory usage, and flash I/O)

# What you should know from this lecture

1. The idea of B+ Tree, LSM Tree, and Hash Tables
2. Choices these data structures (B+ Tree, LSM, and Hash Table)
3. What advantages and disadvantages they offer when implementing them over NAND flash
4. Key problem and solution: uTree
5. Key problem and solution: LOCS
6. Key problem and solution: WiscKey
7. What is read/write amplification in LSM tree

# Further References

- Hyeontaek Lim, Bin Fan, David G. Andersen, and Michael Kaminsky. 2011. SILT: a memory-efficient, high-performance key-value store. In Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles (SOSP '11).
- Assaf Eisenman, Asaf Cidon, Evgenya Pergament, Or Haimovich, Ryan Stutsman, Mohammad Alizadeh, and Sachin Katti. 2019. Flashield: a hybrid key-value cache that controls flash write amplification. In Proceedings of the 16th USENIX Conference on Networked Systems Design and Implementation (NSDI'19). USENIX Association, USA, 65–78.
- Oana Balmau, Diego Didona, Rachid Guerraoui, Willy Zwaenepoel, Huapeng Yuan, Aashray Arora, Karan Gupta, and Pavan Konka. 2017. TRIAD: creating synergies between memory, disk and log in log structured key-value stores. In Proceedings of the 2017 USENIX Conference on Usenix Annual Technical Conference (USENIX ATC '17). USENIX Association, USA, 363–375.
- PinK: High-speed In-storage Key-value Store with Bounded Tails, USENIX ATC 2020.
- Jiacheng Zhang, Youyou Lu, Jiwu Shu, and Xiongjun Qin. 2017. FlashKV: Accelerating KV Performance with Open-Channel SSDs. ACM Trans. Embed. Comput. Syst. 16, 5s, Article 139 (October 2017), 19 pages.
- Lanyue Lu, Thanumalayan Sankaranarayana Pillai, Andrea C. Arpaci-Dusseau,  and Remzi H. Arpaci-Dusseau, WiscKey: Separating Keys from Values  in SSD-conscious Storage, USENIX FAST 2016.
- Leonardo Marmol, Swaminathan Sundararaman, Nisha Talagala, and Raju Rangaswami. 2015. NVMKV: a scalable, lightweight, FTL-aware key-value store. In Proceedings of the 2015 USENIX Conference on Usenix Annual Technical Conference (USENIX ATC '15). USENIX Association, USA, 207–219.
- Zhaoyan Shen, Feng Chen, Yichen Jia, and Zili Shao. 2018. DIDACache: An Integration of Device and Application for Flash-based Key-value Caching. ACM Trans. Storage 14, 3, Article 26 (November 2018), 32 pages.

# Example 2: HashTable on Flash

## FlashStore: High Throughput Persistent Key-Value Store

Biplob Debnath[*]
University of Minnesota
Twin Cities, USA
biplob@umn.edu

Sudipta Sengupta
Microsoft Research
Redmond, USA
sudipta@microsoft.com

Jin Li
Microsoft Research
Redmond, USA
jinl@microsoft.com

### ABSTRACT

We present FlashStore, a high throughput persistent key-value store, that uses flash memory as a non-volatile *cache* between RAM and hard disk. FlashStore is designed to store the working set of key-value pairs on flash and use one flash read per key lookup. As the working set changes over time, space is made for the current working set by destaging recently unused key-value pairs to hard disk and recycling pages in the flash store. FlashStore organizes key-value pairs in a log-structure on flash to exploit faster sequential write performance. It uses an in-memory hash table to index them, with hash collisions resolved by a variant of cuckoo hashing. The in-memory hash table stores compact key signatures instead of full keys so as to strike tradeoffs between RAM usage and false flash read operations.

FlashStore can be used as a high throughput persistent key-value storage layer for a broad range of server class applications. We compare FlashStore with BerkeleyDB, an embedded key-value store application, running on hard disk and flash separately, so as to bring out the performance gain of FlashStore in not only using flash as a cache above hard disk but also in its use of flash aware algorithms. We use real-world data traces from two data center applications, namely, Xbox LIVE Primetime online multi-player game and inline storage deduplication, to drive and evaluate the design of FlashStore on traditional and low power server platforms. FlashStore outperforms BerkeleyDB by up to 60x on throughput (ops/sec), up to 50x on energy efficiency (ops/Joule), and up to 85x on cost efficiency (ops/sec/dollar) on the evaluated datasets.

A high throughput persistent key-value store can help to improve the performance of such applications. Flash memory is a natural choice for such a store, providing persistency and 100-1000 times lower access times than hard disk. Compared to DRAM, flash access times are about 100 times higher. Flash stands in the middle between DRAM and disk also in terms of cost – it is 10x cheaper than DRAM, while 20x more expensive than disk – thus, making it an ideal gap filler between DRAM and disk.

There are two types of popular flash devices, NOR and NAND flash. NAND flash architecture allows a denser layout and greater storage capacity per chip. As a result, NAND flash memory has been significantly cheaper than DRAM, with cost decreasing at faster speeds. NAND flash characteristics have lead to an explosion in its usage in consumer electronic devices, such as MP3 players, phones, caches and Solid State Disks (SSDs). In the rest of the paper, we use NAND flash based SSDs as the architectural choice and simply refer to it as flash memory. We describe SSDs in detail in Section 2. To get the maximum performance per dollar out of SSDs, it is necessary to use flash aware data structures and algorithms to avoid small random writes that not only have a higher latency but also reduce flash device lifetimes through increased page wearing.

In this paper, we present the design and evaluation of FlashStore, a high performance key-value storage system using flash as a cache between RAM and hard disk. When a key-value blob is written, it is sequentially logged in flash. A specialized RAM-space efficient hash table index using a variant of cuckoo hashing [32] and compact key signatures is used to index the key-value blobs stored in flash mem-

60

# FlashStore: Data Structures

Many workloads are _read-heavy_ and do not need indexing (B+ tree a bit of an overkill) - restrictive layout how the keys can be stores

- Microsoft wanted to have flash SSDs as a KV cache in front of their HDDs

If we just do a simple hash(key) → location, that would be good enough

- Hash has O(1) lookup time, not O(Log(n)) like B+ tree

But the "small write" problem. We cannot store each key in its own page (in efficient) and cannot do small writes to just to update the key

**Goal**: fast KV cache with a single flash I/O read to locate data

# Design Goals and Issues

1.  Deliver low-latency, high-throughput operations
    a.  For small key looks up
    b.  Values can be in DRAM cache or on Flash

2.  Use flash-aware data structures
    a.  Do not do small page updates

3.  Low RAM footprint for indexing to lookup on flash
    a.  Technically you can use 8 bytes per key and 64 bytes of value
    b.  So for a 1 TB of flash drive, you will need 1 TB / (64 + 8) x 8 bytes = 122 GB of DRAM (!)
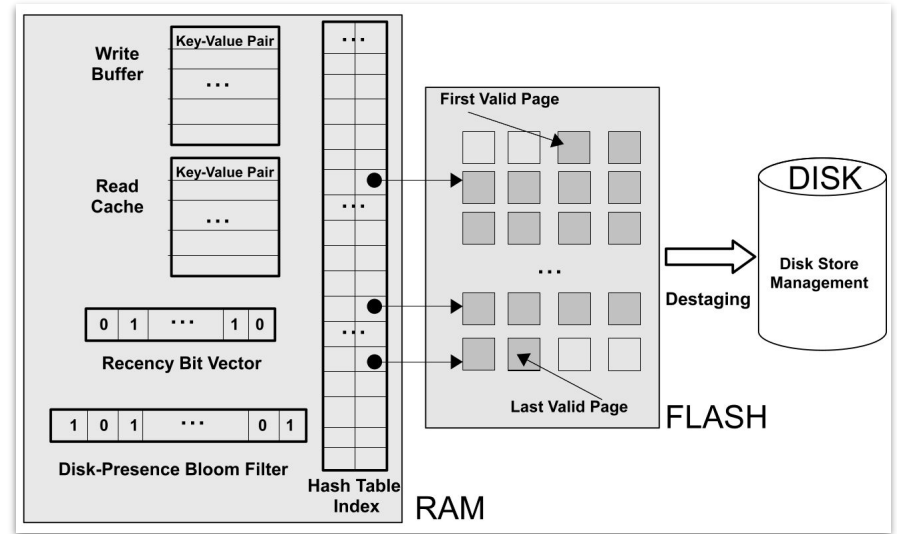    c.  Same problem as with the FTL

# Architecture

**RAM Write buffer :** buffer until the flash page size

**Read cache:** fixed-size read cache for recently used items (LRU)

**Recency Bit Vector:** maintains access information for staging data between flash and disk

**Bloom filter:** probabilistic "false positive", but never "false negative" *(it's not there when it is there)*

**HashTable:** The primary data structure to look for key → flash location in one flash read
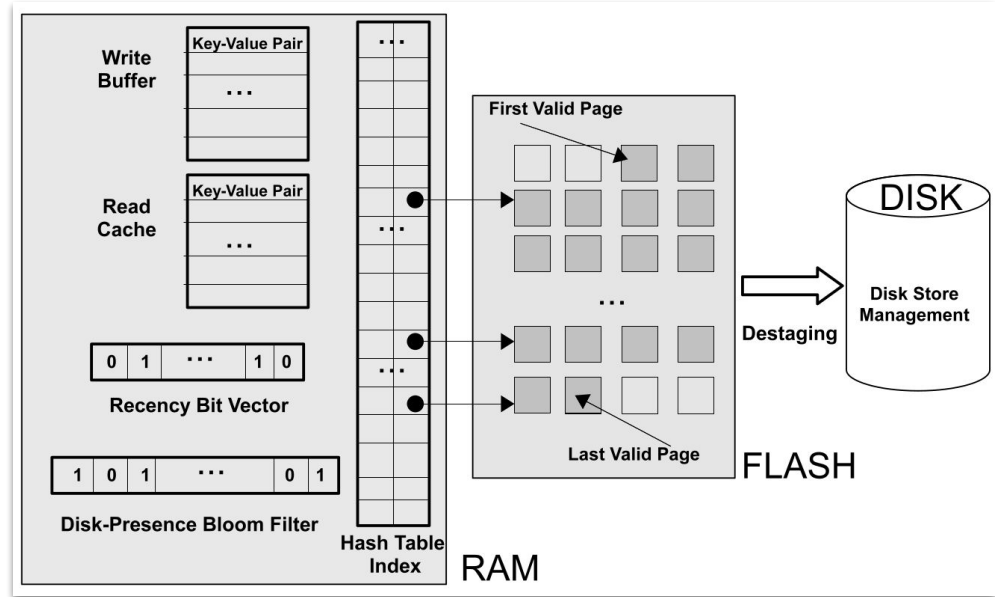
# Key Lookup and Insertion Operations

**Insert** (with timestamps):
1. Into the write buffer
2. Wait until full
3. Write out to flash
4. Update the HT index

**Lookup**

1. In RAM read cache
2. In RAM write cache
3. Lookup in HT index to find on flash
4. Lookup bloom filter
   a. No: return NULL
   b. Yes: disk search (B+ tree)
5. Update recency bit
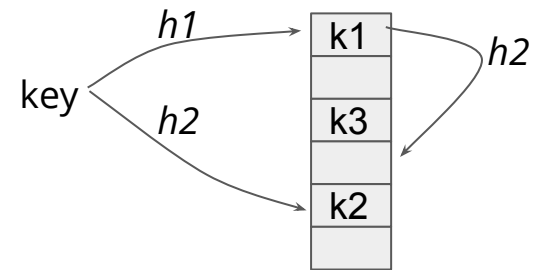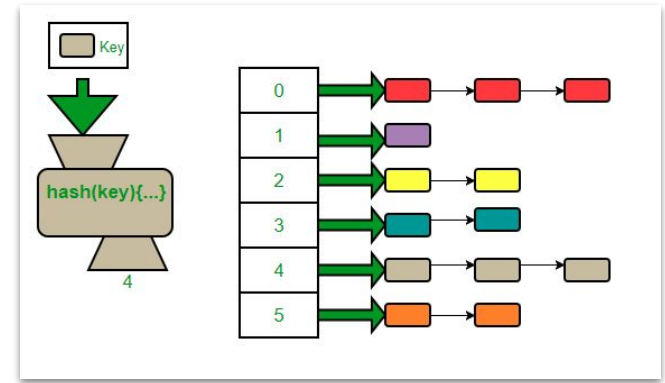6. (Optional) put in RAM read cache

# Hash Table Design

In a simple hash table, we can do something like



- Hash(key) → HT slot → check if the key stored there matches
  - OK, then follow the flash page pointer (8bytes)
  - Collision: then follow the link list of collision pointers

Uses **Cuckoo hashing** : use "n" hash functions and find the first free location to put the key. No need to scan any linear list in case of high collision
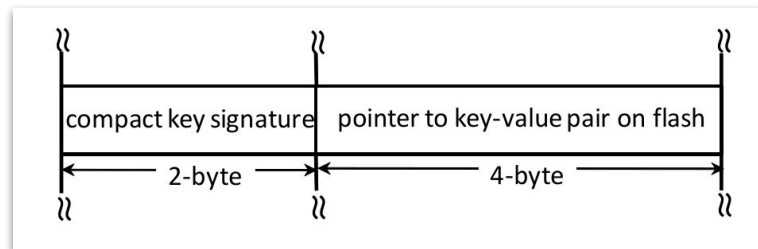
*What to store in these hash table slots?* Full key and flash page address? (lots of data)

# Hash Table Memory Usage: What to Store?

Compact key signature (instead of full key and hash):



- A full key can be of any size, hashes are large too (160-512 bits)
- If the key used $i^{th}$ hash function then used the top-order **16 bits** as a compact signature

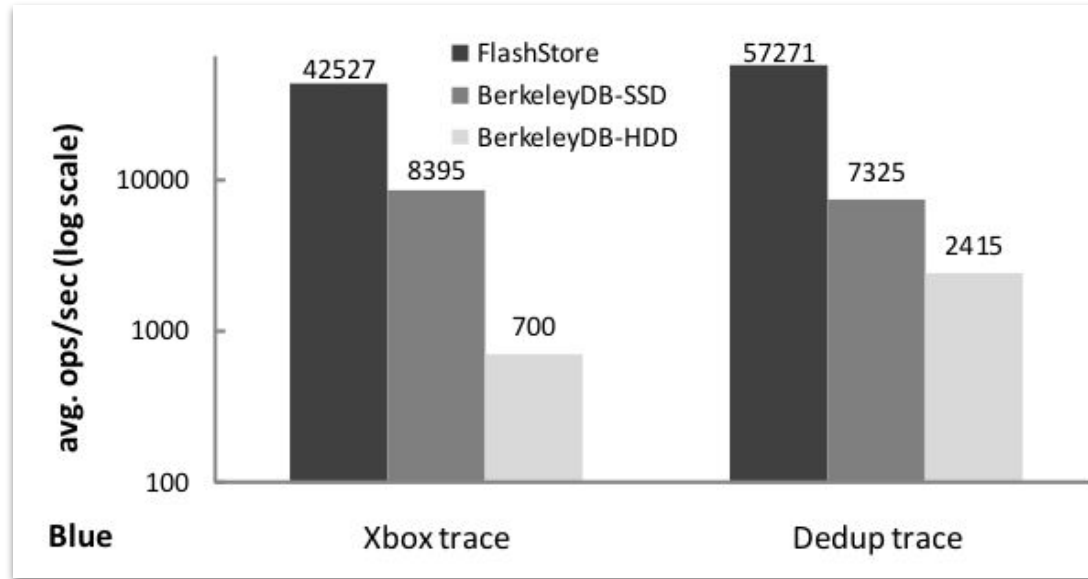Flash page offset as 4 byte pointers (not 8 bytes) : maximum size = $2^{32}$ x 4KB = 8TB

- How many bits to use, can be optimized for the given size of the device
- For example, 160GB device (what they used), 160GB/4KB = 26 bits only
  - Rest of the (32 - 26) = 6 bits, can be used for in-page offsets of 128 bytes
  - Hence, 128 bytes becomes the minimum packing granularity

Broadly speaking: a memory-efficient HT table design is an active research problem (many papers are out there in the field, we are only covering one trick)

# Flash Specific Concerns

- Filled flash pages are written in a log-append order (lookup is done using the in-memory HT table)
  - Log garbage collection for entries that have been overwritten or deleted (similar logic)

- After certain HT table occupancy and Flash usage - trigger destaging from flash to HDD
  - Pick pages and check the recency bitmap in memory to find if they have been accessed recently
    - Yes, put them in write buffer (back in the circulation)
    - No, push them to HDD and make space
- At crash
  - Default option: build HT by scanning flash logs
  - Options 2: checkpointing

# Performance



Delivers performance for two important workloads for Microsoft (xbox, and dedup)

Compared with running BerkeleyDB (B+Tree) on SSD and HDD

# WiscKey: Doing garbage collection in vLog

A **native way** would be : to scan the LSM key tree to identify all valid values and then remove them.



Value Log

**Better way**: to keep a back reference to the keys in the value log as well

Once GC kicks in, values from the tail are read, validated by querying the LSM tree, and then move to the head

The new tail, and addresses are then inserted in the LSM tree before cleaning values