



# Key-Value Store: Database for Unstructured Bigdata

<b>Google</b> <ul style="list-style-type: none"><li>- Bigtable, Level DB, Hbase</li><li>- For Web indexing and messaging</li></ul>	<b>Amazon</b> <ul style="list-style-type: none"><li>- Dynamo, SimpleDB</li><li>- For E-commerce</li></ul>	<b>ORACLE</b> <ul style="list-style-type: none"><li>- Oracle</li><li>- NoSQL, Berkeley DB</li><li>- For Configurable KV</li></ul>
<b>Facebook</b> <ul style="list-style-type: none"><li>- Haystack, RocksDB, Cassandra</li><li>- For social network and photo store</li></ul>	<b>Microsoft</b> <ul style="list-style-type: none"><li>- Azure, Cosmos DB</li><li>- For E-commerce</li></ul>	
<b>YAHOO!</b> <ul style="list-style-type: none"><li>- LinkedIn</li><li>- Voldemort</li><li>- For Scalability</li></ul>	<b>Baidu</b> <ul style="list-style-type: none"><li>- Atlas</li><li>- For Cloud data</li></ul>	<b>Basho</b> <ul style="list-style-type: none"><li>- Riak</li><li>- For distributed KV</li></ul>
<b>LinkedIn</b> <ul style="list-style-type: none"><li>- PNUTS</li><li>- For Advertising</li><li>- For Scalability</li></ul>	<b>Open source</b> <ul style="list-style-type: none"><li>- Redis, Memcached</li><li>- For in-memory DB, cache</li></ul>	

September 2, 2021

Jongmoo Choi

Dankook University

<http://embedded.dankook.ac.kr/~choijm>

# Content

---

- What is Key-Value Store?
  - ✓ What is Bigdata?
  - ✓ Structured vs. Unstructured data
  - ✓ Instances and Recent Studies about Key-Value Store
- RocksDB Basic
- Advance Topic 1: Compaction
- Advance Topic 2: Lookup
- Advance Topic 3: WAL
- Advance Topic 4: New approach
- RocksDB Practice
- Conclusion

# What is Key-Value Store? (1/12)

---

## ■ Lecture Objective

- ✓ Understand characteristics of Bigdata.
- ✓ Can explain differences between structured and unstructured data
- ✓ Know about NoSQL, Database for unstructured data
- ✓ Explore real instances of Key-Value Store (a kind of NoSQL)
- ✓ Investigate research issues of Key-Value Store

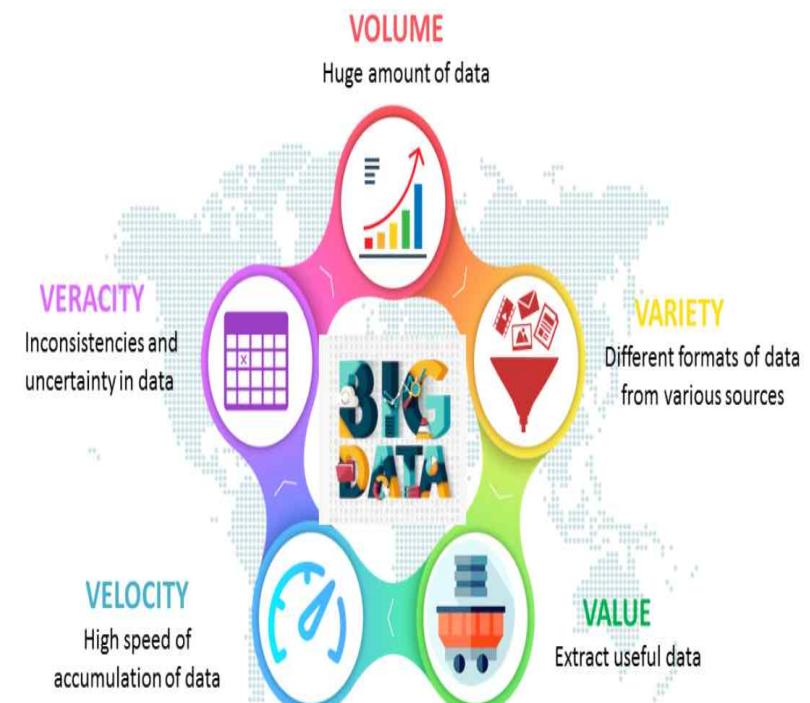
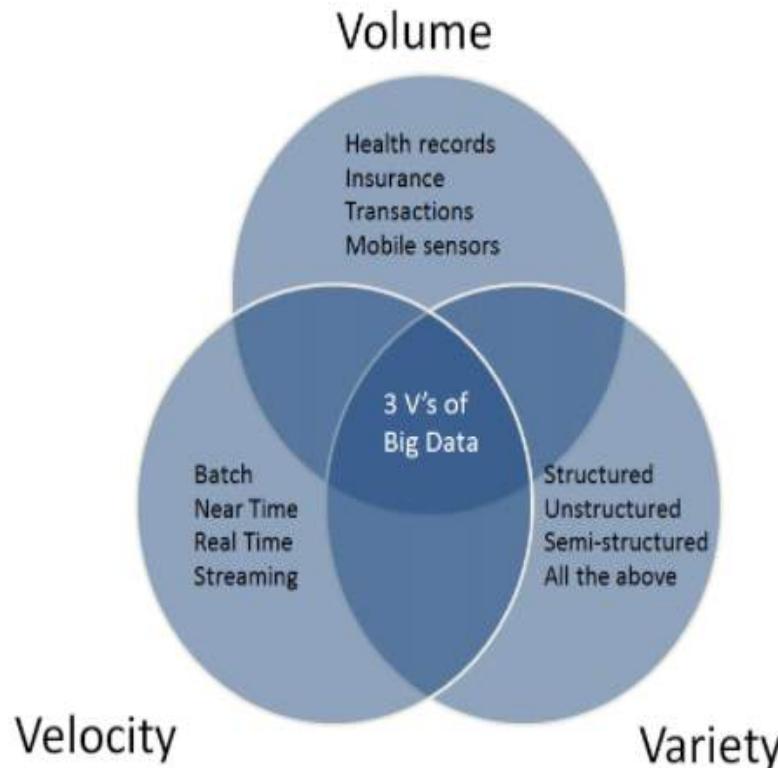


(Source: [www.dreamstime.com](http://www.dreamstime.com))

# What is Key-Value Store? (2/12)

## ■ What is Bigdata?

- ✓ Three features: Volume, Velocity, Variety
  - These days, being extended to Five or Eight features

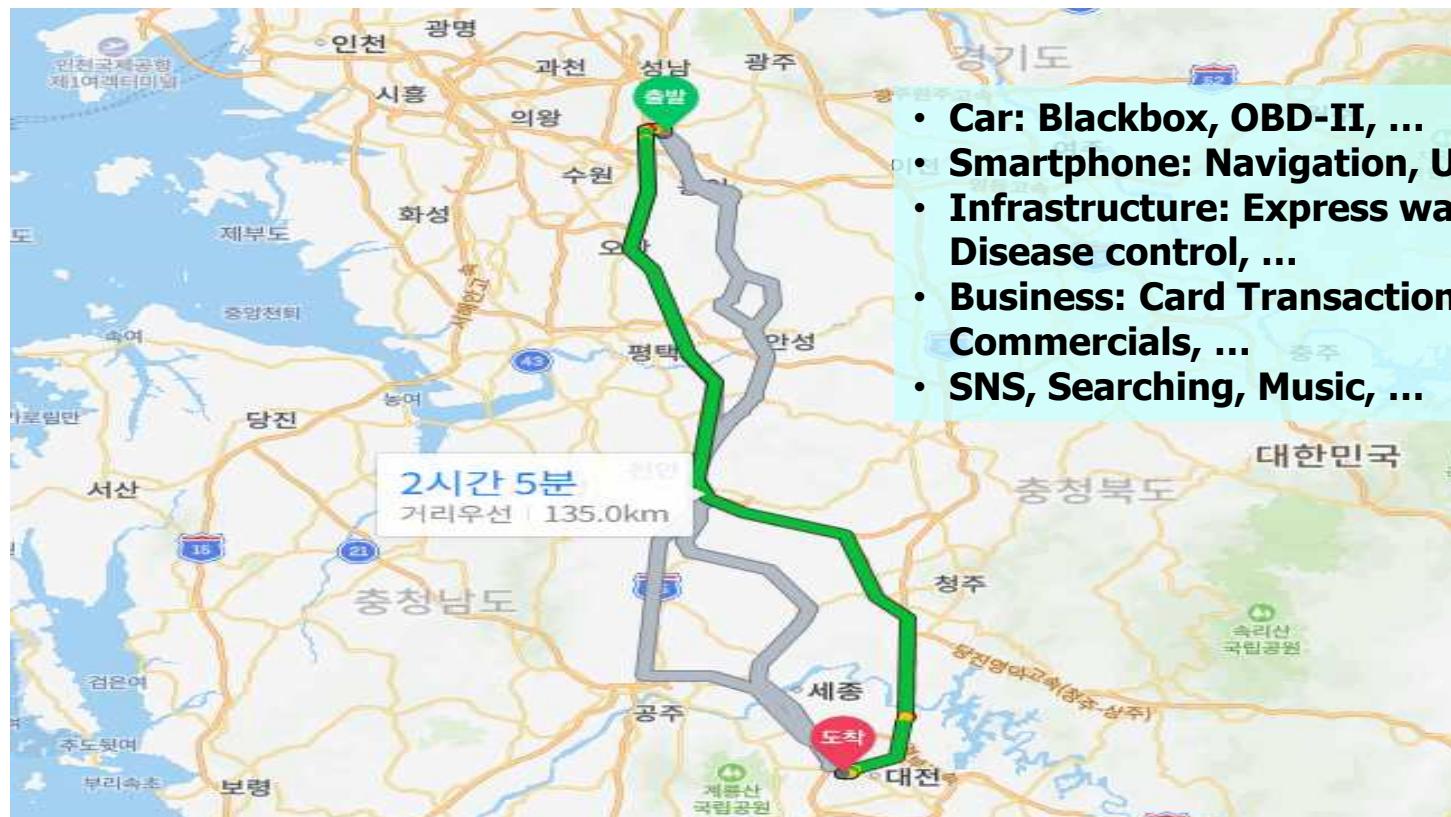


(Source: <https://bigdataldn.com/intelligence/big-data-the-3-vs-explained/>,  
<https://laptrinhx.com/what-is-big-data-a-beginner-s-guide-to-the-world-of-big-data-1936020691/>)

# What is Key-Value Store? (3/12)

## ■ Volume and Velocity Perspective

- ✓ Data gathered while driving



- ✓ 25GB per hour (from IEEE AutoSafety, 2020)

# What is Key-Value Store? (4/12)

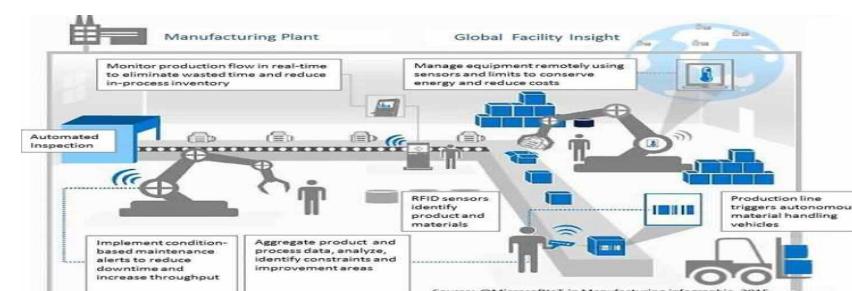
## ■ How about Variety?

- ✓ **Conventional data**: well defined format (schema)
  - Value, Field (Column, Attribute), Record (Row, Tuple), Table, Relation
- ✓ **Bigdata**: various and diverse (schema free)
  - Various fields/frequencies in each record, Diverse format, ...

Field			Table		
id	ISSN-L	ISSNs	PublisherId	Journal_Title	
0	2056-9890	2056-9890		1	Acta Crystallographica Section E Crystallographic Communications
1	2077-0472	2077-0472		2	Agriculture
2	2073-4395	2073-4395		2	Agronomy
3	2076-2615	2076-2615		2	Animals
4	2076-3417	2076-3417		2	Applied Sciences
5	2306-5354	2306-5354		2	Bioengineering
6	2079-7737	2079-7737		2	
7	2079-6374	2079-6374		2	

Relation Name  
↓  
STUDENT  
↓  
Attributes  
↓  
Tuples  
↓  
The attributes and tuples of a relation STUDENT.

	Name	Ssn	Home_phone	Address	Office_phone	Age	Gpa
Benjamin Bayer	305-61-2435	373-1616	2918 Bluebonnet Lane	NULL	19	3.21	
Chung-cha Kim	381-62-1245	375-4409	125 Kirby Road	NULL	18	2.89	
Dick Davidson	422-11-2320	NULL	3452 Elgin Road	749-1253	25	3.53	
Rohan Panchal	489-22-1100	376-9821	265 Lark Lane	749-6492	28	3.93	
Barbara Benson	533-69-1238	839-8461	7384 Fontana Lane	NULL	19	3.25	

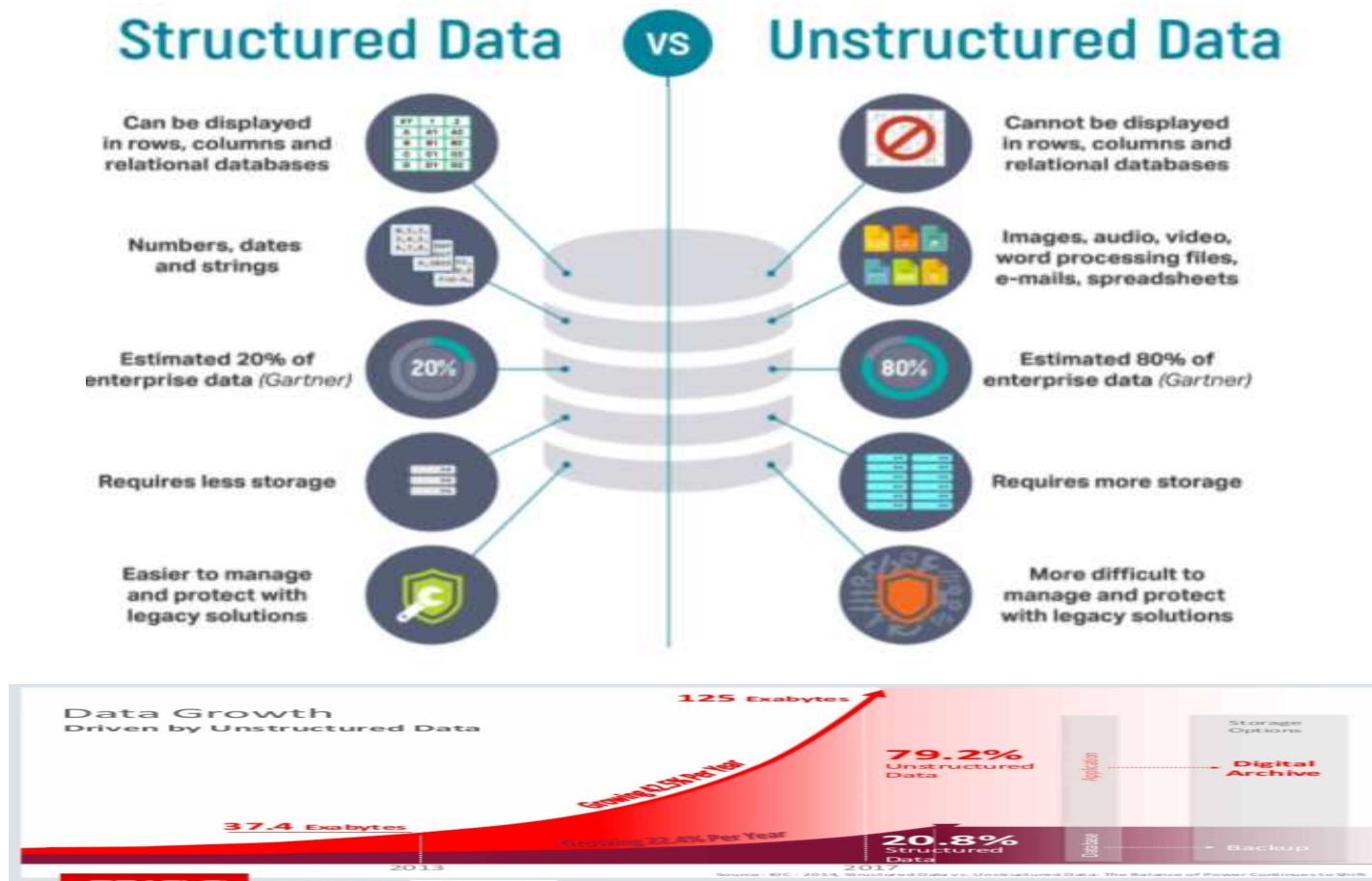


(Source: [www.slideshare.net/b15ku7/chapter-2-relational-data-modelpart1/](http://www.slideshare.net/b15ku7/chapter-2-relational-data-modelpart1/) and [www.researchgate.net/figure/Shows-an-example-of-product-and-data-flow-in-a-Smart-Factory-Product-carries-RFID-tag\\_fig1\\_327884684](http://www.researchgate.net/figure/Shows-an-example-of-product-and-data-flow-in-a-Smart-Factory-Product-carries-RFID-tag_fig1_327884684))

# What is Key-Value Store? (5/12)

## ■ Types of Data

- ✓ Structured vs. Unstructured (vs. Semi-structured)

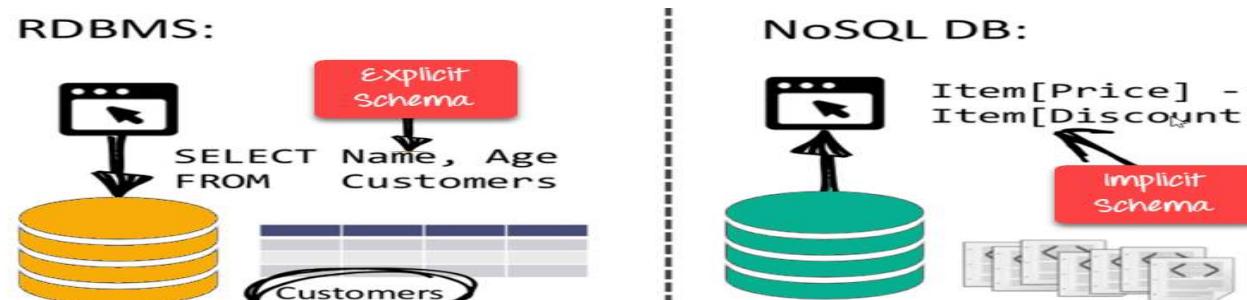


(Source: <https://lawtomated.com/structured-data-vs-unstructured-data-what-are-they-and-why-care/> and <https://www.linkedin.com/pulse/deconstructing-internet-things-oil-gas-keith-moore/>)

# What is Key-Value Store? (6/12)

## ■ Types of Database

- ✓ RDB(Relational DB) for structured vs. Non RDB for unstructured
  - Also known as SQL DB vs. NoSQL DB (Note: **Not-Only SQL!**)



(Source: <https://www.guru99.com/nosql-tutorial.html/>)

- ✓ Example of RDB
  - Infer information via the relation of tables

PubID	Publisher	PubAddress
03-4472822	Random House	123 4th Street, New York
04-7733903	Wiley and Sons	45 Lincoln Blvd, Chicago
03-4859223	O'Reilly Press	77 Boston Ave, Cambridge
03-3920886	City Lights Books	99 Market, San Francisco

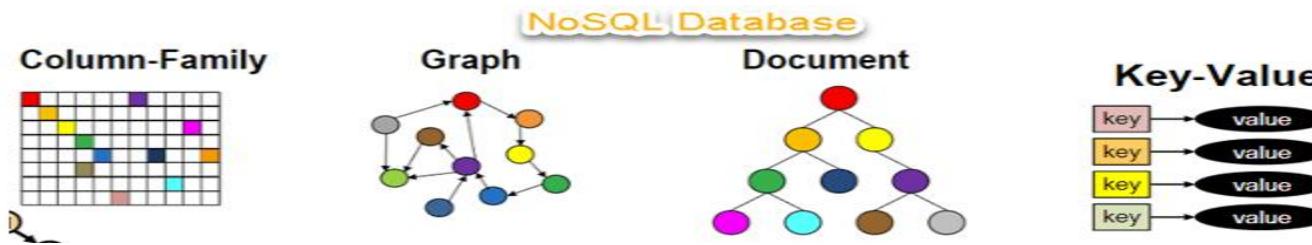
AuthorID	AuthorName	AuthorBDay
345-28-2938	Haile Selassie	14-Aug-92
392-48-9965	Joe Blow	14-Mar-15
454-22-4012	Sally Hemmings	12-Sept-70
663-59-1254	Hannah Arendt	12-Mar-06

ISBN	AuthorID	PubID	Date	Title
1-34532-482-1	345-28-2938	03-4472822	1990	Cold Fusion for Dummies
1-38482-995-1	392-48-9965	04-7733903	1985	Macrame and Straw Tying
2-35921-499-4	454-22-4012	03-4859223	1952	Fluid Dynamics of Aquaducts
1-38278-293-4	663-59-1254	03-3920886	1967	Beads, Baskets & Revolution

# What is Key-Value Store? (7/12)

- Types of Database (cont')
  - ✓ Classification of NoSQL DB



(Source: [www.guru99.com/nosql-tutorial.html/](http://www.guru99.com/nosql-tutorial.html/))

- Key-Value DB (Key Value Store)
  - Manage key and value pair → simple
  - Examples: LevelDB, RocksDB, Redis, ...
- Document DB
  - Extend KV store, KVs are stored in an organized format (like XML)
  - Examples: MongoDB, Couchbase, ...
- Column-family (wide column) DB
  - Column-oriented KV management (sparse)
  - Examples: Bigtable, HBase, Cassandra, ...
- Graph DB
  - KV pairs and their relationships
  - Examples: Neo4j, JanusGraph, ...

```
{
  "id": "webShop635874323254478467",
  "transactionDateUTC": "2016-01-03T15:38:45.44884722",
  "lines": [
    {
      "articleNo": "12345678",
      "articleName": "Bookcase White Rodney 233x100",
      "quantity": 1,
      "itemPrice": 44.45
    }
  ],
  "invoiceAddress": {
    "addressType": "DELIVERY",
    "addressLine1": "The big City Street 2233",
    "addressLine2": "Dept Invoice",
    "city": "Malmö",
    "zipCode": 11111
  },
  "customer": {
    "firstName": "Peter",
    "lastName": "InTheAzureSky",
    "email": "peter@example.com",
    "mobilePhone": "+4642000000"
  }
}
```

# What is Key-Value Store? (8/12)

## ■ Key-Value Store (a.k.a. Key-Value DB)

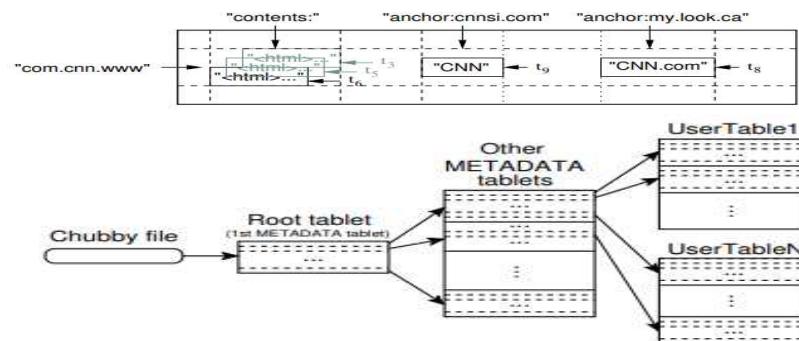
- ✓ A de-facto standard DB for unstructured data
- ✓ Google, Facebook, Amazon, Microsoft, MongoDB, Yahoo, Hbase, LinkedIn, Oracle, Baidu, Basho, In Memory DB (Memcached), ...



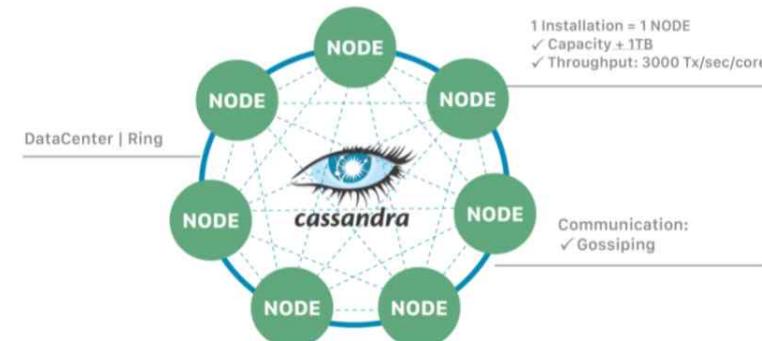
# What is Key-Value Store? (9/12)

## ■ Key-Value Stores: some examples

- ✓ LevelDB
  - By Google, 2011, a subset of Bigtable (Column-oriented DB, OSDI, 2006)
  - Level compaction, Open-source
- ✓ RocksDB
  - By Facebook, 2012, a fork of LevelDB
  - Various algorithms (e.g. Tier compaction), High performance, Diverse applications
- ✓ HBase
  - By Apache, 2008, motivated by Google's Bigtable
  - A distributed data storage system for the Hadoop ecosystem
- ✓ Cassandra
  - By Apache (Facebook), 2008, motivated by Amazon's Dynamo
  - A decentralized architecture, each node is powered by LST-tree based DB



(Source: Bigtable, index by row key, column key, timestamp)



(Source: [cassandra.apache.org/\\_/cassandra-basics.html/](http://cassandra.apache.org/_/cassandra-basics.html/))

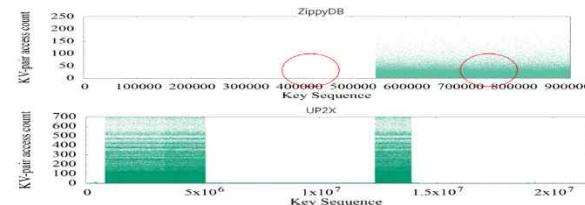
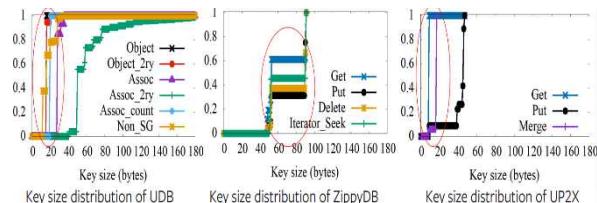
# What is Key-Value Store? (10/12)

---

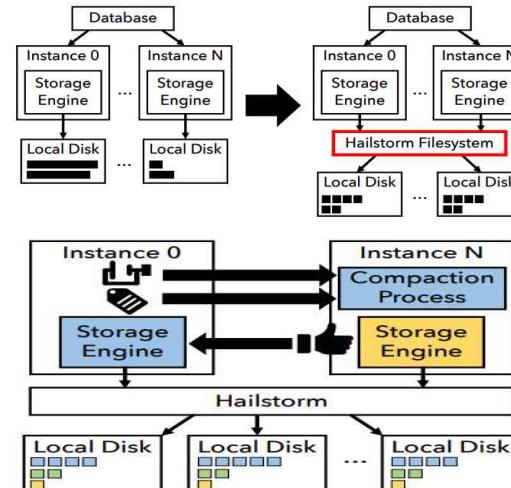
- Key-Value Store: Hot research topics not only academia but also industry
  - ✓ S. Dong et al., "Evolution of Development Priorities in Key-value Stores Serving Large-scale Applications: The [RocksDB](#) Experience", FAST'21.
  - ✓ L. Lu et al., "[WiscKey](#): Separating Keys from Values in SSD-conscious Storage", FAST'16.
  - ✓ Y. Dai et al., "From WiscKey to Bourbon: A [Learned Index](#) for Log-Structured Merge Trees", OSDI'20.
  - ✓ P. Raju et al., "[PebblesDB](#): Building Key-Value Stores using Fragmented Log-Structured Merge Trees", SOSP'17.
  - ✓ O. Balmau et al., "[SILK](#): Preventing Latency Spikes in Log-Structured Merge Key-Value Stores", ATC'19.
  - ✓ W. Zhong et al., "[REMIX](#): Efficient Range Query for LSM-trees", FAST'21.
  - ✓ J. Im et al., "[PinK](#): High-speed In-storage Key-value Store with Bounded Tails", ATC'20.
  - ✓ T. Yao et al., "[MatrixKV](#): Reducing Write Stalls and Write Amplification in LSM-tree Based KV Stores with a Matrix Container in NVM", ATC'20.
  - ✓ H. Chen et al., "[SpanDB](#): A Fast, Cost-Effective LSM-tree Based KV Store on Hybrid Storage", FAST'21
  - ✓ G. DeCandia et al., "[Dynamo](#): Amazon's Highly Available Key-value Store", SOSP'07.
  - ✓ L. Bindschaedler et al., "[Hailstorm](#): Disaggregated Compute and Storage for Distributed LSM-based Databases", ASPLOS'20.
  - ✓ N. Dayan et al., "[Dostoevsky](#): Better Space-Time Trade-Offs for LSM-Tree Based Key-Value Stores via Adaptive Removal of Superfluous Merging", SIGMOD'18.
  - ✓ Y. Matsunobu et al., "[MyRocks](#): LSM-Tree Database Storage Engine Serving Facebook's Social Graph", VLDB'20

# What is Key-Value Store? (11/12)

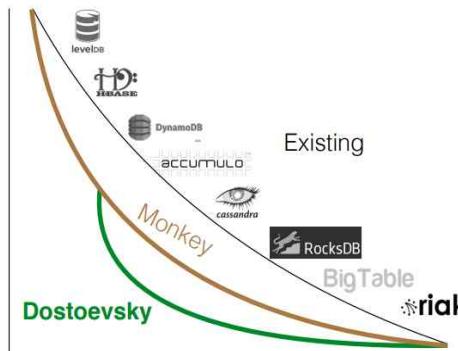
## ■ Interesting recent studies regarding Key-Value Store



(RocksDB analysis, FAST'20)



(Hailstorm, ASPLOS'20)



(Dostoevsky, SIGMOD'18)

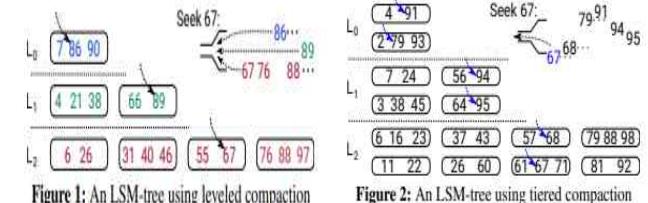
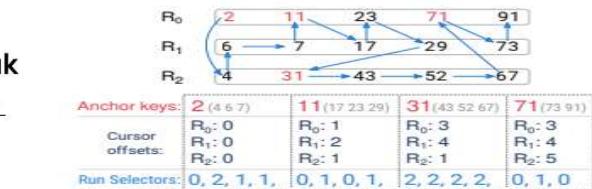
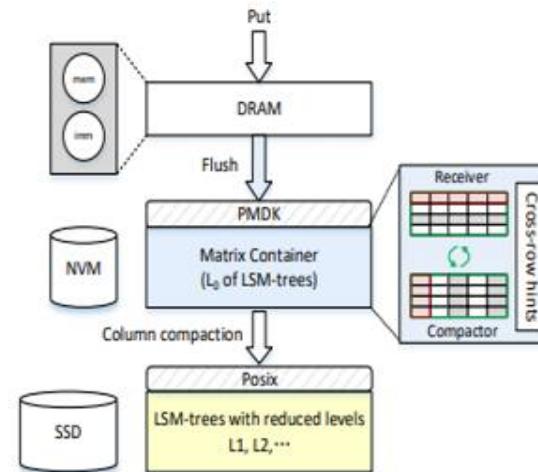


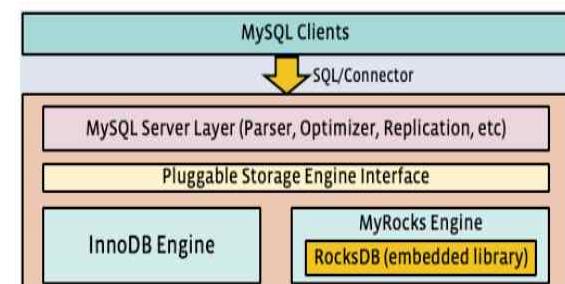
Figure 2: An LSM-tree using tiered compaction



(REMIX, FAST'21)



(MatrixKV, ATC'20)



Engine	Space	CPU seconds/s for writes	CPU seconds/s for reads + writes	Bytes written per second
InnoDB	2187.4GB	0.89	1.83	13.34MB
Myrocks	824.4GB	0.55	1.65	3.42MB

(Myrocks, VLDB'20)

# What is Key-Value Store? (12/12)

---

## ■ Take-away Lessons

- ✓ Understand characteristics of Bigdata
  - Volume, Velocity, Variety → 80% unstructured data
- ✓ NoSQL: Database for unstructured data
  - Key-Value store, Document DB, Column-family DB, Graph DB
- ✓ Key-Value Store examples
  - LevelDB, RocksDB, Bigtable, HBase, Dynamo, ...
  - Wiskey, PebblesDB, SILK, MatrixDB, Hailstorm, Pink, ...
- ✓ Popularly studied from both industry and academy
  - Performance
    - Write: Compaction, Flush, Background threads, ...
    - Read: Bloom filter, Range query, Parallel access, ...
  - Consistency: WAL, Recovery, ...
  - New storage: NVM + SSD, Optane, ...
  - Distributed environment: partition, replication, ...



(Source: [www.dreamstime.com](http://www.dreamstime.com))

(Special thank you to DKU Embedded Members and IITP's SW StarLab (No.2021-0-01475)

# Content

---

- What is Key-Value Store?
- RocksDB Basic
  - ✓ Definition and Use cases
  - ✓ Architecture: LSM-tree data structure
  - ✓ Interface and Internal Operations
- Advance Topic 1: Compaction
- Advance Topic 2: Lookup
- Advance Topic 3: WAL
- Advance Topic 4: New approach
- RocksDB Practice
- Conclusion

# RocksDB Basic (1/15)

---

## ■ Lecture Objective

- ✓ Understand the features of RocksDB
- ✓ Learn the core algorithm of RocksDB: LSM (Log-Structured Merge)-tree
- ✓ Examine how the key-value interface works: put and get
- ✓ Investigate how the internal operations works: flush, compaction and bloom filter
- ✓ Explore RocksDB development history



(Source: [www.dreamstime.com](http://www.dreamstime.com))

# RocksDB Basic (2/15)

## ■ What is RocksDB? (from RocksDB wiki)

- ✓ 1) Famous KV Store of Facebook, derived from LevelDB
- ✓ 2) A persistent storage engine that supports key/value interface
- ✓ 3) LSM (Log Structured Merge)-Tree based (for SSD)
- ✓ 4) Embedded (C++ library) and Open source
- ✓ 5) Support various algorithms (e.g. compaction, filtering, format), configurations (e.g. size, BG thread), tools and debugging facilities

The image displays two side-by-side screenshots of the GitHub repository for `facebook/rocksdb`.  
The left screenshot shows the `Wiki` page, which includes a **Welcome to RocksDB** section, a **Contents** sidebar with links to the RocksDB Wiki, FAQ, Requirements, and Release Methodology, and a sidebar with links to `rocksdb.org`, `Readme`, `View license`, and `Releases`.  
The right screenshot shows the main repository page, featuring a summary of code statistics (380 issues, 230 pull requests, etc.), a list of recent commits (e.g., "bjlemaire and facebook-github-bot Make mempurge a background process (#837765)", "Add micro-benchmark support (#8493)", etc.), and a sidebar with links to `rocksdb.org`, `Readme`, `View license`, and `Releases`.

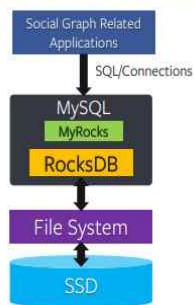
(Source: <https://github.com/facebook/rocksdb/wiki>)

# RocksDB Basic (3/15)

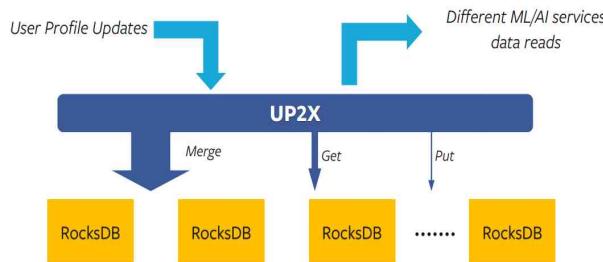
## ■ Use cases

- ✓ 30 applications at Facebook alone (UDB for SQL, ZippyDB for distributed system, UP2X for Logging and ML)
- ✓ Storage engine for other services: LinkedIn, Yahoo, Kafka, Netflix, ...

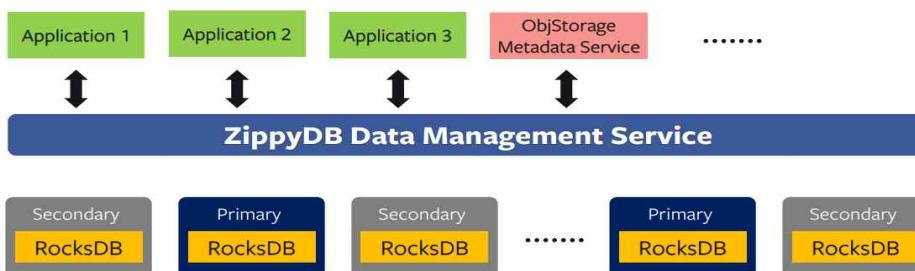
### Use Case 1: UDB



### Use Case 3: UP2X



### Use Case 2: ZippyDB



(Source: Facebook RocksDB paper in FAST'20 and <https://github.com/facebook/rocksdb/blob/master/USERS.md>)

The screenshot shows a GitHub browser window with the URL <https://github.com/facebook/rocksdb/blob/master/USERS.md>. The page title is 'rocksdb/USERS.md at master · facebook/rocksdb'. The content area displays the 'Facebook' section of the file, which discusses the use of RocksDB at Facebook across various services like MyRocks, MongoRocks, ZippyDB, Laser, Dragon, Stylus, and LogDevice. It also links to research publications and code posts. Below the Facebook section is the 'LinkedIn' section, which mentions the use of RocksDB for follow feeds and the Apache Samza framework. At the bottom, it links to a Tech Talk by Ankit Gupta and Naveen Somasundaram.

At Facebook, we use RocksDB as storage engines in multiple data management services and a backend for many different services including:

1. MyRocks -- <https://github.com/MySQLOnRocksDB/mysql-5.6>
2. MongoRocks -- <https://github.com/mongodb-partners/mongo-rocks>
3. ZippyDB -- Facebook's distributed key-value store with Paxos-style replication, built on top of RocksDB.[1] <https://www.youtube.com/watch?v=DfifN7pGODokhtt>
4. Laser -- Laser is a high query throughput, low (millisecond) latency, key-value storage service built on top of RocksDB.
5. Dragon -- a distributed graph query engine. <https://code.facebook.com/posts/1737605303120405/dragon-a-distributed-graph-query-engine>
6. Stylus -- a low-level stream processing framework written in C++.[1]
7. LogDevice -- a distributed data store for logs [2]

[1] <https://research.facebook.com/publications/realtime-data-processing-at-facebook/>  
[2] <https://code.facebook.com/posts/357056558062811/logdevice-a-distributed-data-store-for-logs/>

**LinkedIn**

Two different use cases at LinkedIn are using RocksDB as a storage engine:

1. LinkedIn's follow feed for storing user's activities. Check out the blog post: <https://engineering.linkedin.com/blog/linkedin-s-feed-made-faster-and-smarter>
2. Apache Samza, open source framework for stream processing

Learn more about those use cases in a Tech Talk by Ankit Gupta and Naveen Somasundaram: <http://www.youtube.com/watch?v=120730204806>

**Yahoo**

Yahoo is using RocksDB as a storage engine for their biggest distributed data store Sherpa. Learn more about it here: <http://yahooeng.tumblr.com/post/120730204806/sherpa-scales-new-heights>

# RocksDB Basic (4/15)

## ■ Use cases (cont')

- ✓ Types: DB, Streaming, Logging, Indexing, Caching, ...
- ✓ Core infrastructure for Bigdata service: need engineers

	Read/Write	Read Types	Special Characteristics
Databases	Mixed	Get + Iterator	Transactions and backups
Stream Processing	Write-Heavy	Get or Iterator	Time window and checkpoints
Logging / Queues	Write-Heavy	Iterator	Support HDD too
Index Services	Read-Heavy	Iterator	Bulk loading
Cache	Write-Heavy	Get	Can drop data

Table 1: RocksDB use cases and their workload characteristics

The screenshot shows a web browser displaying a presentation slide from Confluent. The slide is titled "Performance Tuning RocksDB for Kafka Streams' State Stores". It includes a section on "Kafka Summit 2020" and a list of use cases. Below the slide, there is a search result from Google for "RocksDB" in Korean, listing several job titles and descriptions.

Performance Tuning RocksDB for Kafka Streams' State Stores

Kafka Summit 2020

구축 등)에 대한 구축 및 트리플 슈팅, 성능 최적화, DevOps 경험이 있으신 분

- Cloud Front, Akamai등의 CDN을 통한 컨텐츠 전송 경험이 있으신 분

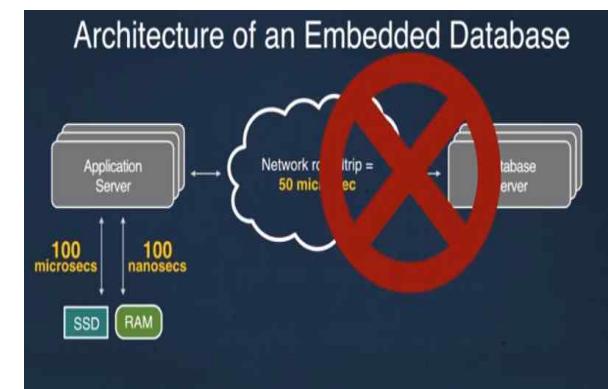
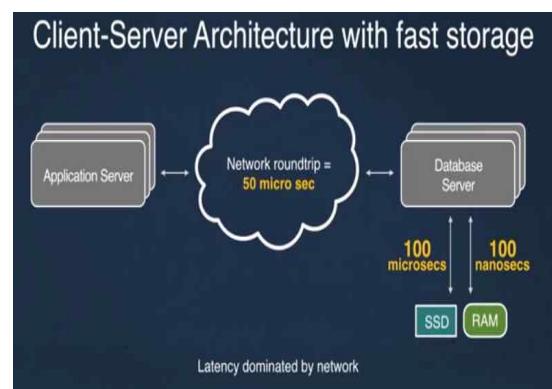
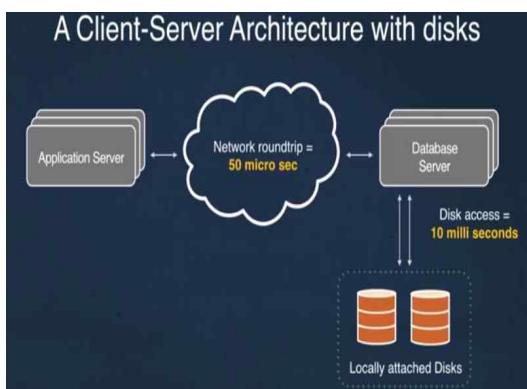
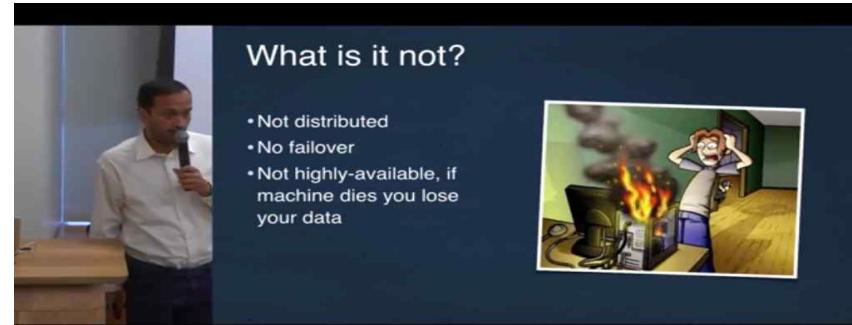
우대사항

- Ceph, Gluster 등의 오픈소스 스토리지에 기여할 수 있는 분
- 대용량 분산 스토리지 구축/운영 경험이 있으신 분
- RocksDB, LevelDB등 스토리지를 위한 Key/Value DB 분석, 최적화 경험이 있으신 분
- OpenStack, Kubernetes 환경 구축, 운영, 최적화, 트리플 슈팅이 가능하신 분
- Ansible 등 DevOps 를 사용이 가능하신 분
- Elasticsearch, Prometheus등의 NoSQL 또는 Time series DB를 활용하여 Grafana등의 대시보드 제작이 가능하신 분
- Hadoop, Spark와 같은 데이터 분석 프레임워크를 이용한 데이터 분석 활용 경험이 가능하신 분
- C/C++ 을 활용한 Linux System Programming 경험 GIT, CD/CD 파이프라인을 통해 소프트웨어 라이프사이클 정책 기반 소프트웨어 배포 경험이 있으신 분
- 외국어로 업무진행이 가능한 분 (영어 또는 일본어)

# RocksDB Basic (5/15)

## ■ Why Embedded?

- ✓ Traditional: separate DB servers from Applications
- ✓ Fast storage: network overhead becomes heavy
- ✓ Decide to place DB as an embedded in an application

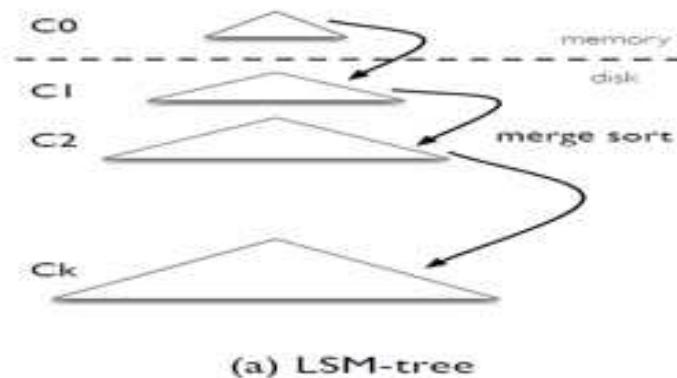


(Source: Dhruba Borthakur, [www.youtube.com/watch?v=hHsxqkcZy7c](https://www.youtube.com/watch?v=hHsxqkcZy7c))

# RocksDB Basic (6/15)

## ■ What is (Log Structured Merge)-tree?

- ✓ By Patrick O'Neil, The Log-Structured Merge Tree, 1996
- ✓ Write optimized data structure
  - **Log-structure:** In-place update → out-of-place update
    - In-place update: good for read, bad for write (due to random writes)
    - Out-of-place update: good for write, possible bad for read (due to multiple locations), need reclaiming mechanism
  - **Merge:** Do merge sort from levels to a next level for deleting old data
    - All data in sorted order
  - **Tree:** larger at lower levels like a tree
    - C0 is in main memory while C1~CK in Storage



(Source: Wisckey paper, FAST'16)

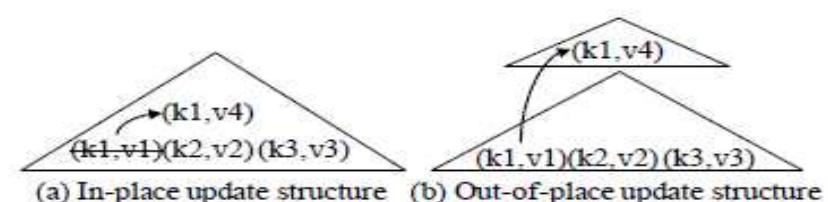


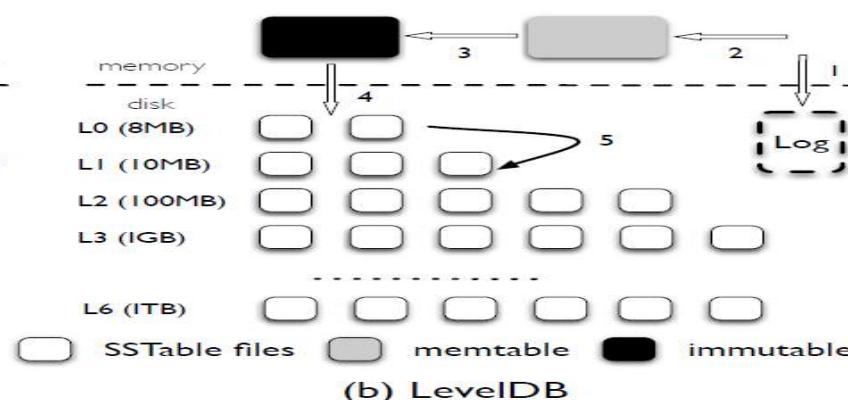
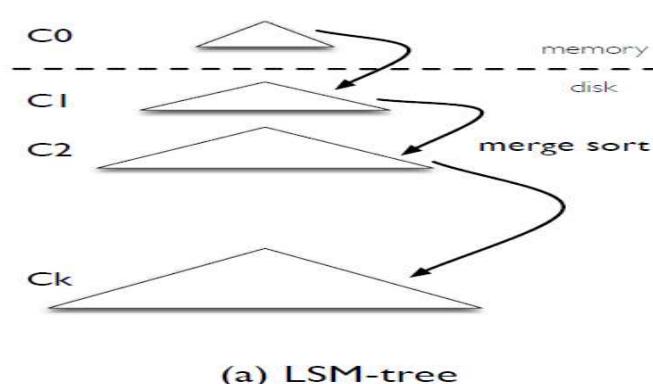
Fig. 1: Examples of in-place and out-of-place update structures: each entry contains a key (denoted as "k") and a value (denoted as "v")

(Source: LSM-based Storage Techniques, VLDB Journal'19)

# RocksDB Basic (7/15)

## ■ What is (Log Structured Merge)-tree?

- ✓ Real implementation in RocksDB (and LevelDB)
  - Memtable for C0
    - Further separate into mutable and immutable
    - Managed by the Skiplist data structure (or hash)
  - A set of SSTables for C1~Ck (multiple Levels, configurable)
    - Default fanout ratio = 10,  $|L_{i+1}| / |L_i|$
    - SSTable internals: data block, index block (logically B+-tree)
  - Log (WAL) for durability
    - A set of records where each record consists of CRC, size, type and payload



(Source: Wisckey paper, FAST'16)

# RocksDB Basic (8/15)

## ■ Key-Value Interface

- ✓ Key/Value: Arbitrary byte streams
- ✓ User visible interfaces
  - `put`, `get`, range scan(iterator), delete, single delete, ...

```
std::string value;
rocksdb::Status s = db->Get(rocksdb::ReadOptions(), key1, &value);
if (s.ok()) s = db->Put(rocksdb::WriteOptions(), key2, value);
if (s.ok()) s = db->Delete(rocksdb::WriteOptions(), key1);
```

(Source: <https://github.com/facebook/rocksdb/wiki>)

## ✓ Internal operations

- **Flush**: when a Memtable is full, it becomes immutable and written into L1, becoming a new SSTable.
- **Compaction**: When L1 is full, compaction is triggered and some SSTables are merged into L1+1



(Source: Dhruba Borthakur and <https://nabillera.tistory.com/>)

# RocksDB Basic (9/15)

## ■ Key-Value Interface

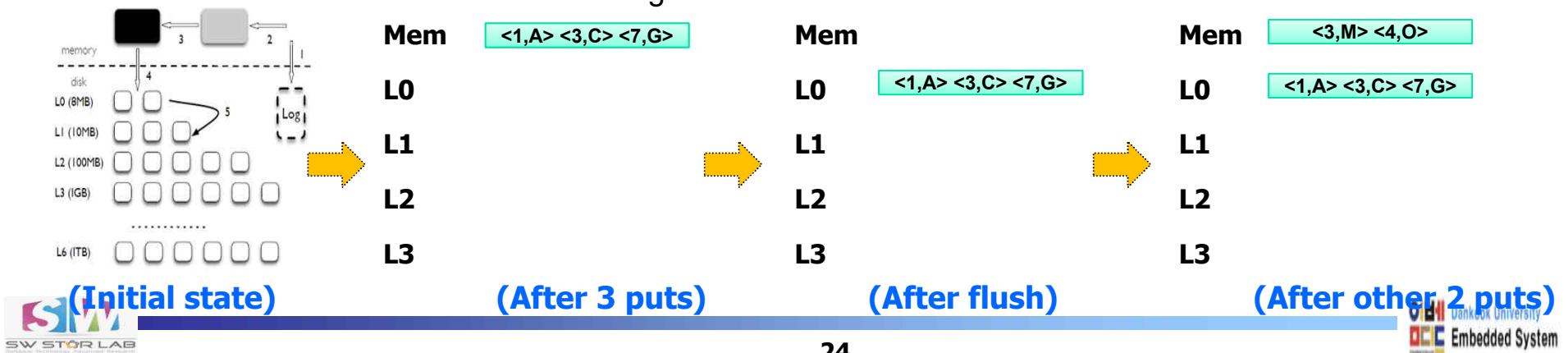
- ✓ Example 1: put/get and flush

- Assumption

- 3 KVs in a Memtable and SSTables
    - $L_0$ : 2 SSTables,  $L_{i+1} = L_i * 2$  (default fanout 10)

- Workload

- put("1", "A"), put("3", "C"), put("7", "G")
      - Memtable → immutable, flushed (Background)
      - A New Memtable is allocated
    - put("3", "M"), put("4", "O")
      - Note: Out-of-place update
    - get(1), get(3)
      - Recent data at higher level



# RocksDB Basic (10/15)

## ■ Key-Value Interface

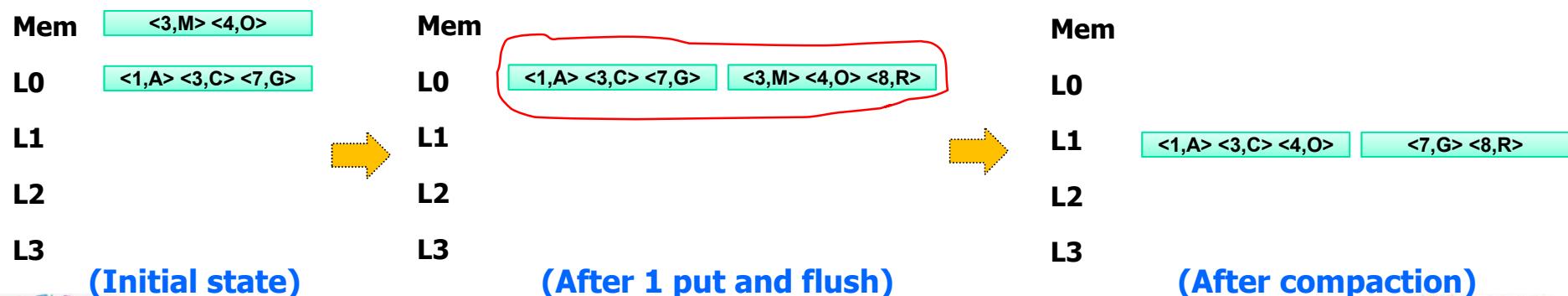
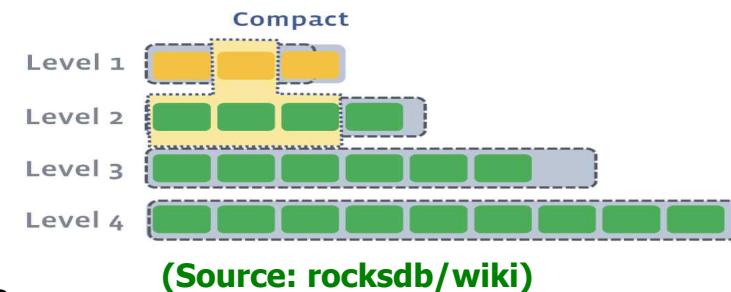
### ✓ Example 2: compaction

#### ▪ Workload

- put("8", "R")
- Memtable → immutable, flushed again
- L0 beyond a threshold → trigger compaction

#### ▪ Compaction Procedure

- 1) select candidate (round robin, least overlapped, ...)
- 2) read overlapped SSTables from current and next level (L0 and L1)
- 3) do merge sort (remove obsolete date)
- 4) write new SSTables to the next level
- Note → L0: can be **overlapped**, L1~Lmax: **no overlapped** among SSTables



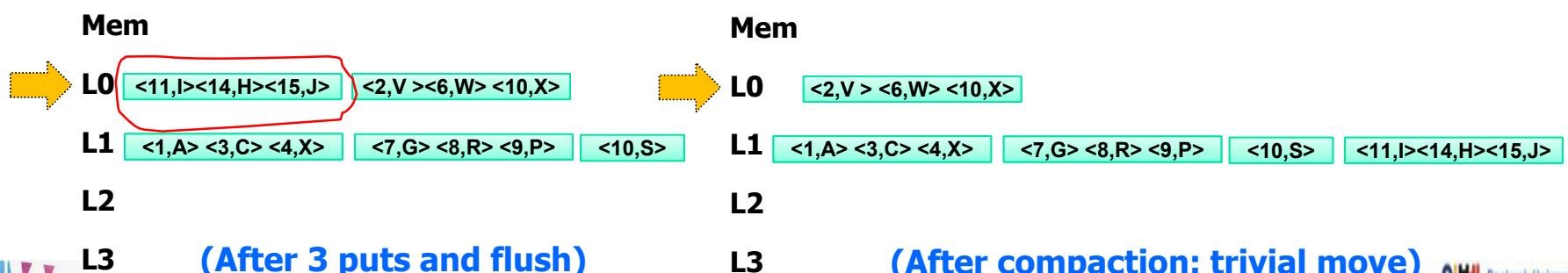
# RocksDB Basic (11/15)

## ■ Key-Value Interface

### ✓ Example 2: compaction (cont')

#### ▪ Workload

- put("10", "S"), put("9", "P"), put("4", "X"),
  - Memtable: managed by Skiplist
- put("15", "J"), put("11", "I"), put("14", "H") → compaction
- put("2", "V"), put("6", "W"), put("10", "X") → compaction
  - Candidate selection affect performance
  - Can trigger compaction recursively (background in actuality)



# RocksDB Basic (12/15)

## ■ Key-Value Interface (cont')

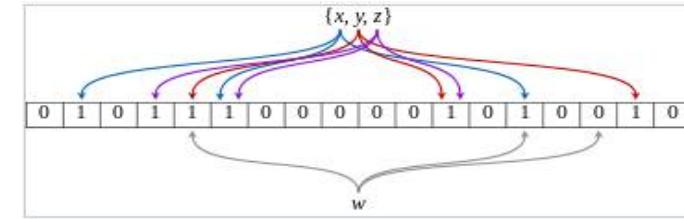
### ✓ Example 3: get with Bloom filter

- Workload

- put("12", "M"), put("14", "Y")
  - get(12) → find at Memtable
  - get(10) → Both L0 and L1 (select higher level)
  - get(8) → at L1 (but unnecessary read at L0)

- Lookup procedure: Memtable first, then SSTables

- Memtable: SkipList
  - SSTable: range check from L0 to the next level until find the request key
  - May do unnecessary read (read amplification) → Employ **Bloom filter**
    - test whether an element is a member of a set, **No false negative**



(Source: Wikipedia)

Mem

L0 <2,V> <6,W> <10,X>

L1 <1,A> <3,C> <4,X> <7,G> <8,R> <9,P> <10,S> <11,I><14,H><15,J>

L2

L3 (Initial state)

Mem <12,M> <14,Y>

L0 <2,V> <6,W> <10,X>

L1 <1,A> <3,C> <4,X> <7,G> <8,R> <9,P> <10,S> <11,I><14,H><15,J>

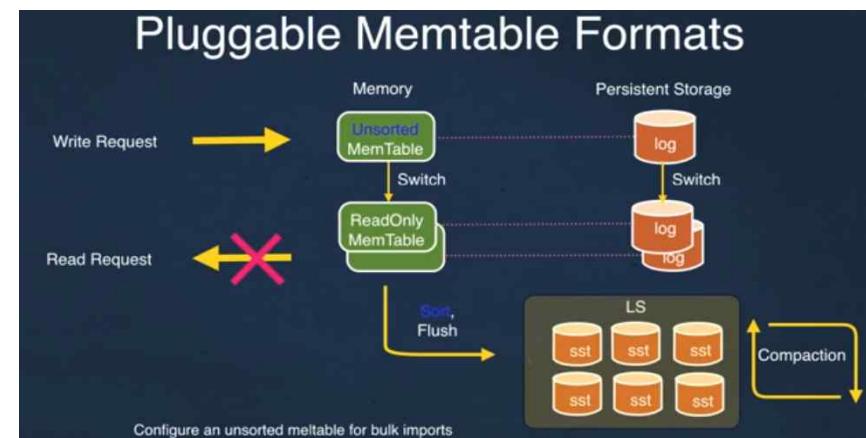
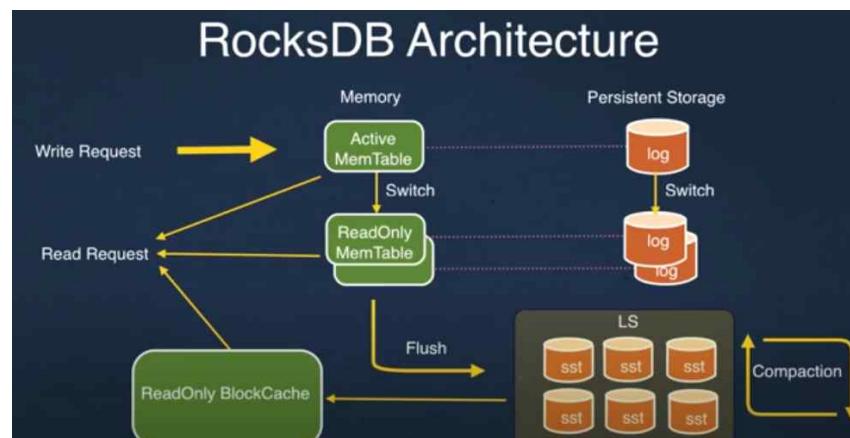
L2

L3 (After 2 puts)

# RocksDB Basic (13/15)

## ■ Features for Performance enhancement

- ✓ Various Compactions: Leveled, Universal (Tiered), FIFO, ...
- ✓ Multi-threaded & Background threads, Compaction filter (e.g. TTL)
- ✓ Bloom filter, Block cache (LRU, Clock, ... → pluggable)
- ✓ Diverse Memtable/SSTable format
  - Memtable format: SkipList, Unsorted (for bulk loading, vectored)
  - SST format: block based vs plain based (less computing, for RAM)
- ✓ New functionalities
  - Column Family (CF), Merge, Transaction, Checksum, Prefix scan, Compression, ...



(Source: Dhruba Borthakur, [www.youtube.com/watch?v=hHsxqkcZy7c](https://www.youtube.com/watch?v=hHsxqkcZy7c))

# RocksDB Basic (14/15)

## ■ History and Recent research

- ✓ Development history: research shift
  - 1) Write amplification: Leveled → Tiered compaction
  - 2) Space amplification: Compression, Compaction frequency (Especially important as shown in Table 2)
  - 3) CPU utilization: Prefix Bloom filter, Hash index
- ✓ Recent focus
  - 1) Distributed system: resource allocation across multiple instances, replication, load balancing
  - 2) Data corruption: need to be detected earlier at every layer
  - 3) Compatibility: for incremental SW upgrade

	CPU	Space Util	Flash Endurance	Read Bandwidth
Stream Processing	11%	48%	16%	1.6%
Logging / Queues	46%	45%	7%	1.0%
Index Services	47%	61%	5%	10.0%
Cache	3%	78%	74%	3.5%

Table 2: System metrics for a typical use case from each application category.

(Source: Facebook RocksDB papers, FAST'21)

A RocksDB Feature Timeline

	Performance	Configurability	Features
2012	• Multi-threaded compactions		• Compaction filters • Locking SSTables from deletion
2013	• Tiered compaction • Prefix Bloom filter • Bloom Filter for MemTables • Separate thread pool for MemTable flush	• Pluggable MemTable • Pluggable file format	• Merge Operator
2014	• FIFO compaction • Compaction rate limiter • Cache-friendly Bloom filters	• String-based config options • Dynamic config changes	• Backup engine • Support for multiple key spaces ("column family") • Physical checkpoints
2015	• Dynamic leveled compaction • File deletion rate limiting • Parallel Level 0 and 1 compaction	• Separate config file • Config compatibility checker	• Bulk loading for SSTable file integration • Optimistic and pessimistic transactions
2016	• Different compression for last level • Parallel MemTable inserts	• MemTable total size caps across instances • Compaction migration tools	• DeleteRange()
2017	• Separate thread pool for bottom-most compactions • Two-level file indices • Level 0 to level 0 compactions	• Single memory limit for both block cache and MemTable	
2018	• Dictionary compression • Hash index into data blocks		• Automatic recovery from out-of-space errors • Query trace and replay tools
2019	• Batched MultiGet() with parallel I/O	• Configure plug-in function using object registry	• Secondary instance
2020	• Multithreaded single file compression		• Entire file checksum • Automatically recover from retrivable errors • Partial support for user-defined timestamps

# RocksDB Basic (15/15)

---

## ■ Take-away Lessons

- ✓ Understand the features of RocksDB
  - 1) Embedded, 2) LSM-tree, 3) Key-Value interface
- ✓ LSM (Log-Structured Merge)-tree
  - 1) Log-structured → Out-of-place update
    - Write-optimized, need compaction
  - 2) Merge
    - Non overlapped SSTables → positive for read
  - 3) Tree
    - Maximum Obsolete data can be bounded
- ✓ Issues
  - How to mitigate the compaction overhead? [Advance topic 1](#)
  - How to enhance read? [Advance topic 2](#)
  - How to support consistency, atomicity, integrity and durability?  
[Advance topic 3](#)
  - How to exploit new storage technology? [Advance topic 4](#)



(Source: [www.dreamstime.com](http://www.dreamstime.com))

# Content

---

- What is Key-Value Store?
- RocksDB Basic
- Advance Topic 1: Compaction
  - ✓ Issue of compaction: Cause latency spike
  - ✓ Compaction styles: Leveled vs Tiered (Universal)
  - ✓ How to reduce compaction overhead?
- Advance Topic 2: Lookup
- Advance Topic 3: WAL
- Advance Topic 4: New approach
- RocksDB Practice
- Conclusion

# Advance Topic 1: Compaction (1/14)

---

## ■ Lecture Objective

- ✓ Understand the role of compaction
- ✓ Learn about the amplification: write, space, and read
- ✓ Examine how to reduce the compaction overhead: scheduling approach
- ✓ Examine how to reduce the compaction overhead: algorithmic approach
- ✓ Explore case studies

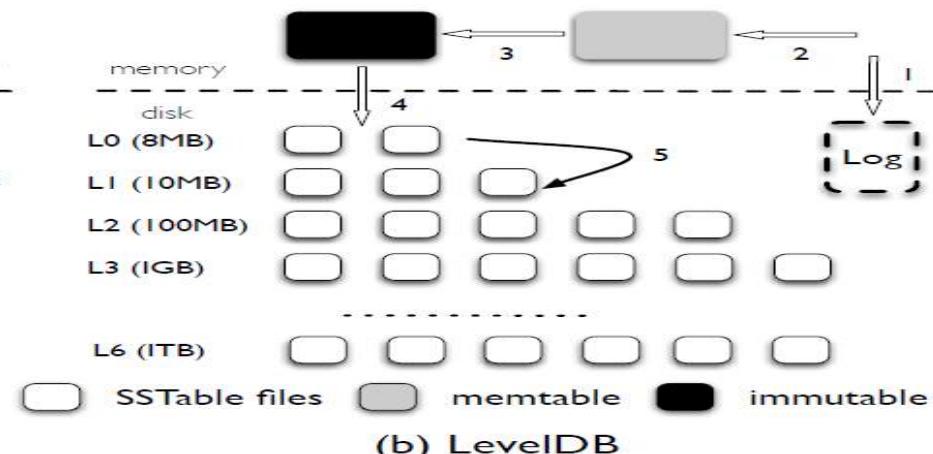
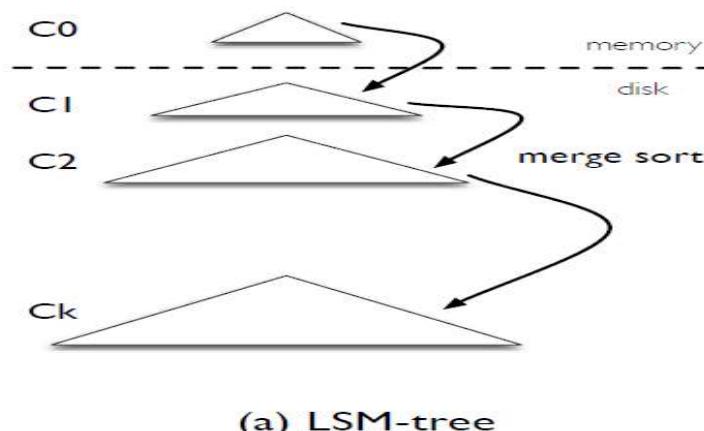


(Source: [www.dreamstime.com](http://www.dreamstime.com))

# Advance Topic 1: Compaction (2/14)

## ■ What is the compaction?

- ✓ RocksDB summary
  - Based on LSM (Log-Structured Merge)-tree
    - Layered: C<sub>0</sub>, C<sub>1</sub>, ..., C<sub>k</sub> (exponentially increasing)
  - Real implementation: Memtable, SSTable
    - Memtable in memory, SSTable in storage (multiple levels: L<sub>0</sub>, L<sub>1</sub>, ..., L<sub>k</sub>)
  - Interface: get, put, range scan, delete, ...
  - put handling
    - 1. WAL, 2. Memtable, 3. immutable Memtable, 4. Flush, 5. Compaction



(Source: Wisckey paper, FAST'16)

# Advance Topic 1: Compaction (3/14)

## ■ What is the compaction?

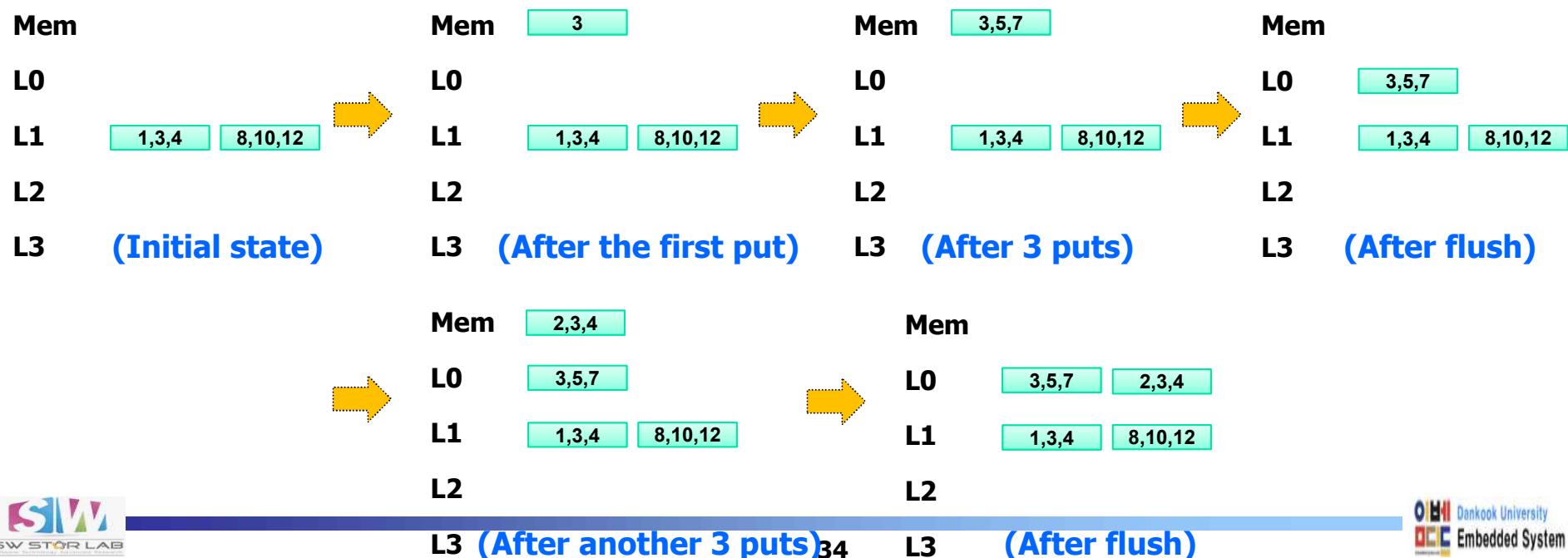
- ✓ Example

- Assumption

- 1) 3 KVs in a Memtable/SSTable, 2)  $L_0$ : 2 SSTables,  $L_{i+1} = L_i * 2$

- Workload

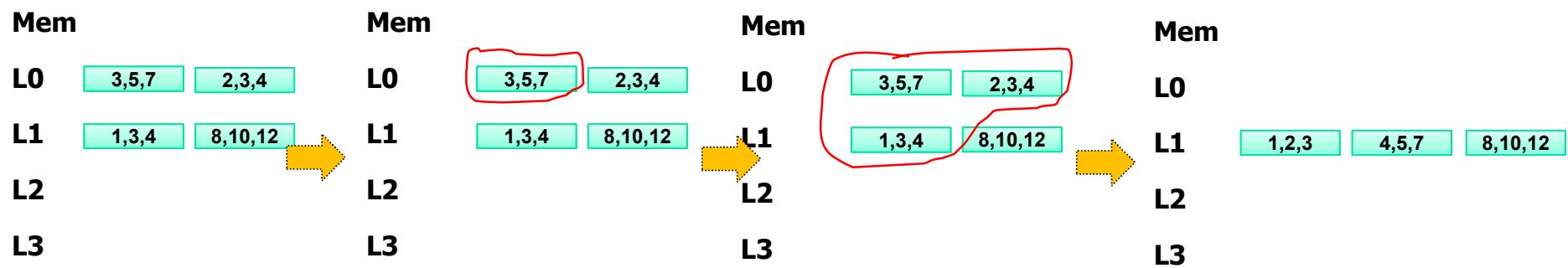
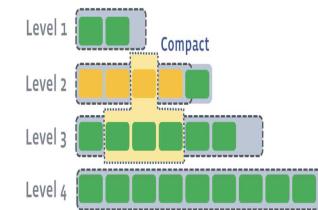
- Initial state: 2 SSTables at  $L_1$  ([show key only here](#))
    - 6 puts: keys = “3”, “7”, “5”, “2”, “3”, “4”
      - Memtable full → trigger flush,  $L_i$  full → trigger compaction



# Advance Topic 1: Compaction (4/14)

## ■ What is the compaction?

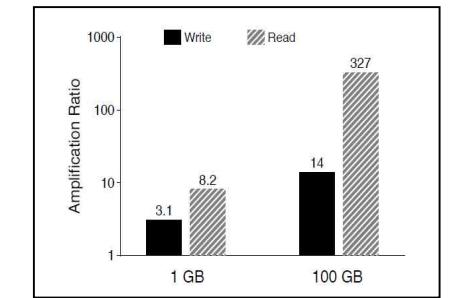
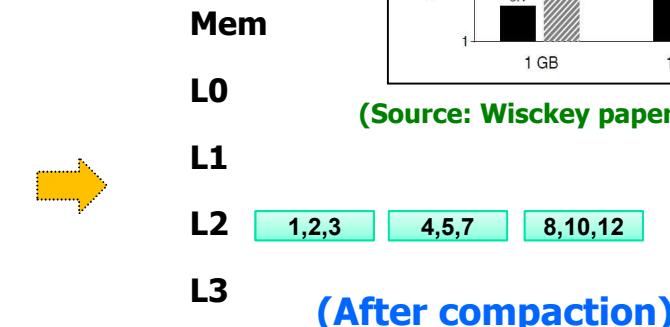
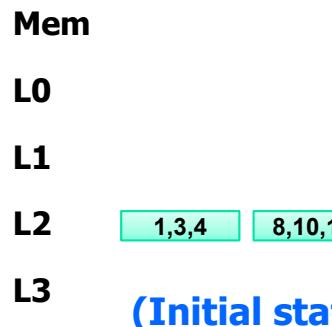
- ✓ Compaction procedure
  - 1. select candidate (FIFO, Least overlapped, ...)
  - 2. read overlapped SSTables from current and next level
  - 3. do merge sort
  - 4. write new SSTables to the next level
- ✓ Compaction effect
  - 1) Remove old data (reclaim), 2) Sort keys at L1, L2, ... (fast lookup)
- ✓ Compaction cost
  - Read/Write SSTables from/into Storage → Heavy operation
  - Cause amplification



# Advance Topic 1: Compaction (5/14)

## ■ What is the amplification?

- ✓ Amplification
  - An undesirable phenomenon where the actual amount of operations/space are more than intended
- ✓ Three types
  - WAF (Write Amplification Factor) = actual writes/intended writes
  - SAF (Space Amplification Factor) = actual used space/required space
  - RAF (Read Amplification Factor) = actual reads/intended reads
- ✓ Example
  - WAF at after compaction = 12/6
  - SAF at after compaction = 9/9, at after flush = 12/9

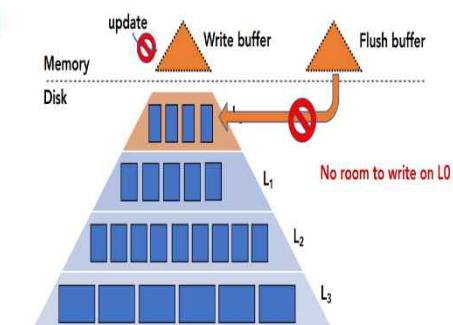
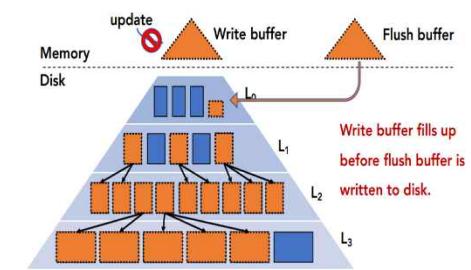
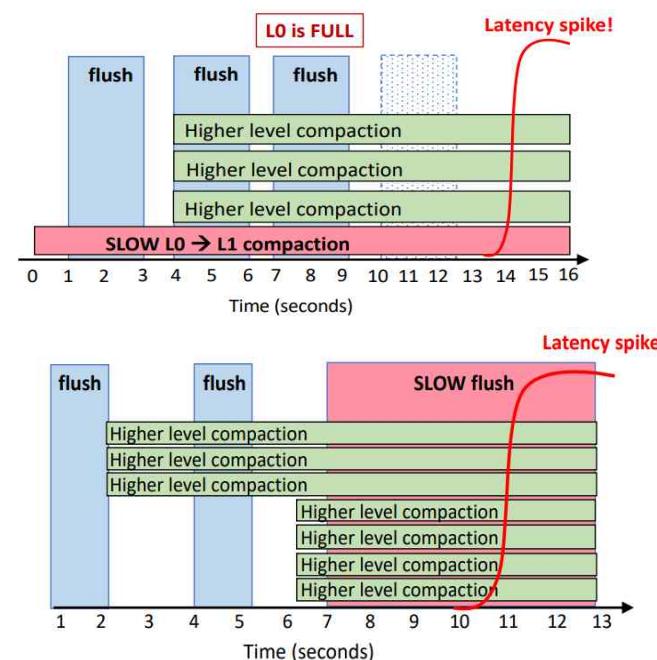
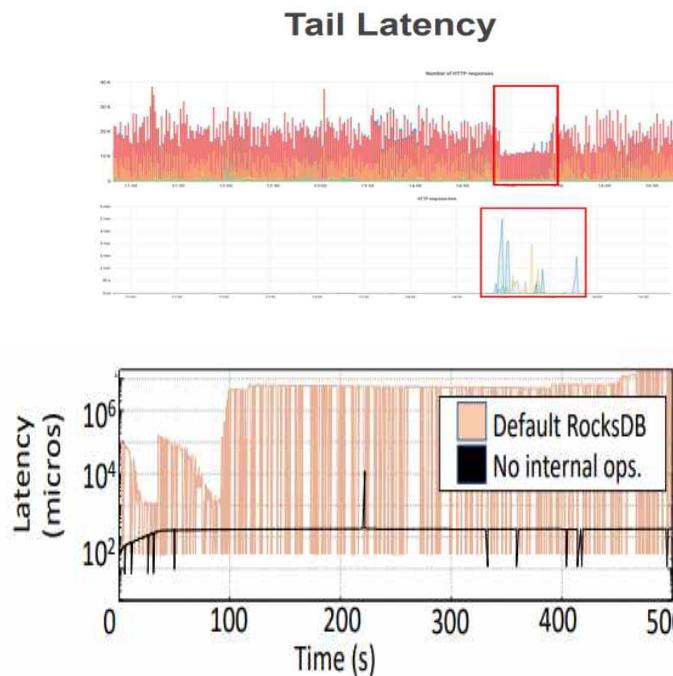


(Source: Wisckey paper, FAST'16)

# Advance Topic 1: Compaction (6/14)

## ■ Concerns about compaction

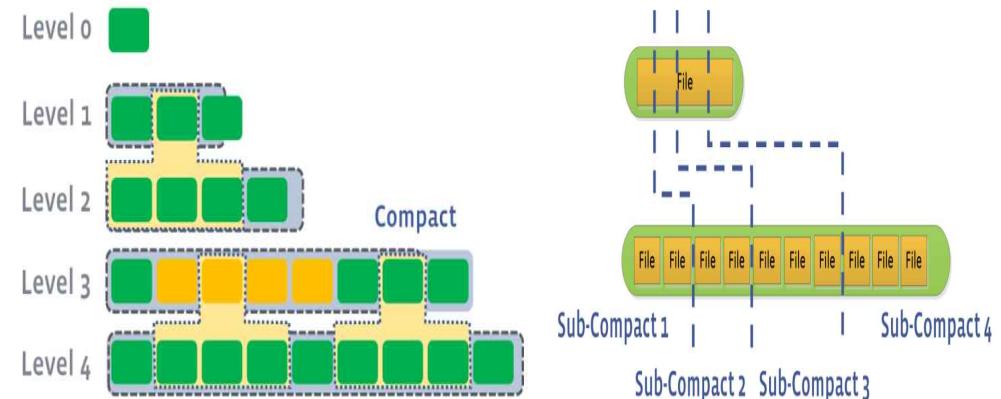
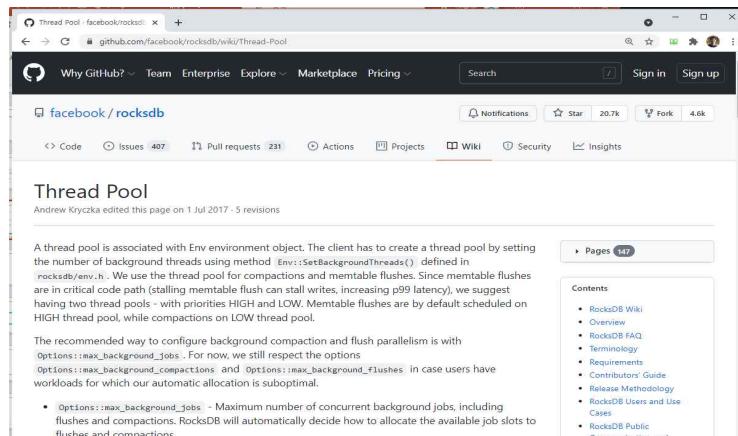
- ✓ 1) Heavy operation (I/O, Sorting), 2) Amplification
- ✓ 3) Delay user requests
  - Question from a developer (domestic company)
  - Unexpected latency Spike → Work Graveyard (밤샘 근무)
  - Key-value DB: Compaction → Interfere Flush → Delay User Request



(Source: SILK papers, ATC'19)

# Advance Topic 1: Compaction (7/14)

- How to reduce compaction overhead: scheduling approach
  - ✓ Thread Pool
    - max\_background\_jobs: # of threads for Flush and Compaction (default:2)
      - max\_background\_compactions and max\_background\_flushes, separately
  - ✓ Background and parallel compaction
    - Multiple compactations: controlled by the max\_background\_compactions
    - Sub compaction: To speed up a compaction job by partitioning it among multiple threads (only for L0 and L1)
  - ✓ Priority control
    - HIGH for flush threads, LOW for compaction threads

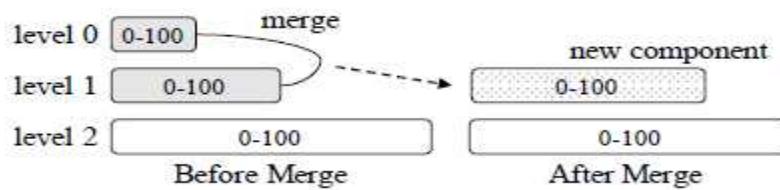


(Source: [rocksdb/wikibook](https://rocksdb.org/wikibook/) and <https://meeeejin.gitbooks.io/rocksdb-wiki-kr/content/leveled-compaction.html>)

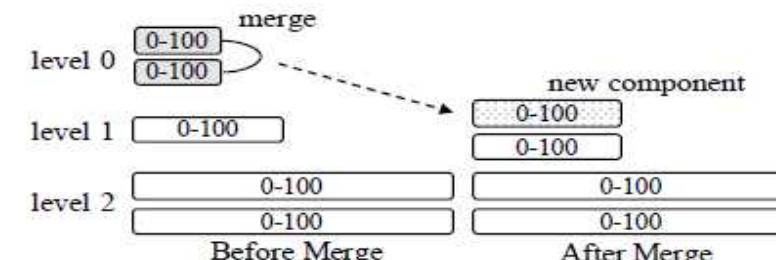
# Advance Topic 1: Compaction (8/14)

## ■ How to reduce compaction overhead: algorithmic approach

- ✓ Leveled compaction
  - Level: consist of one run (component)
  - Merge:  $L_i$  and  $L_{i+1}$
- ✓ Tiered compaction (Universal compaction)
  - Level: consist of T runs (components)
    - . Run: a set of SSTables generated during a time range
  - Merge:  $L_i$  only
- ✓ Difference
  - Leveled: aggressively merge SSTables between levels
  - Tiered: waits for several sorted runs (lazy evaluation)
  - **Tradeoff between WAF and SAF**



(a) Leveling Merge Policy: one component per level



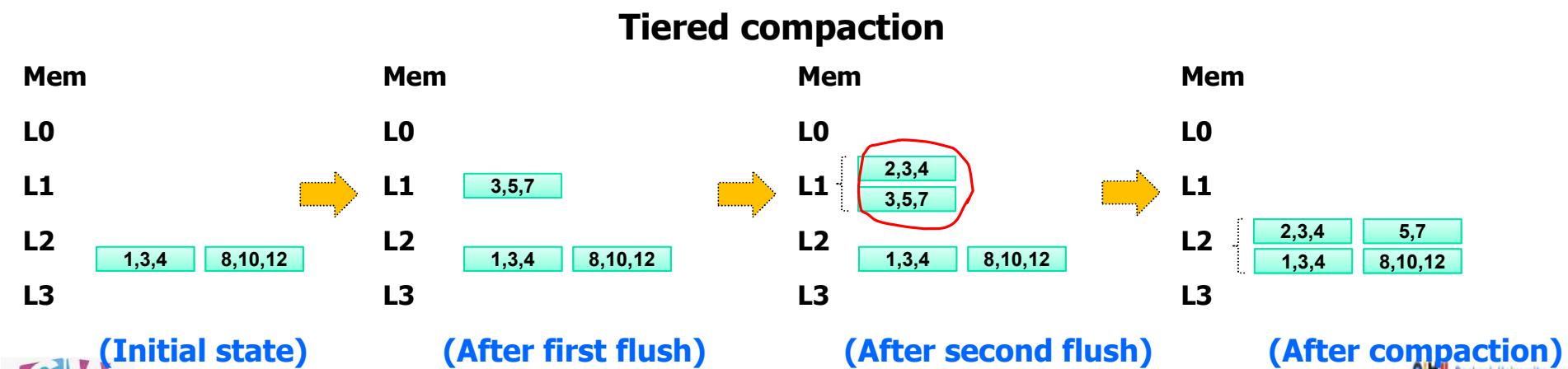
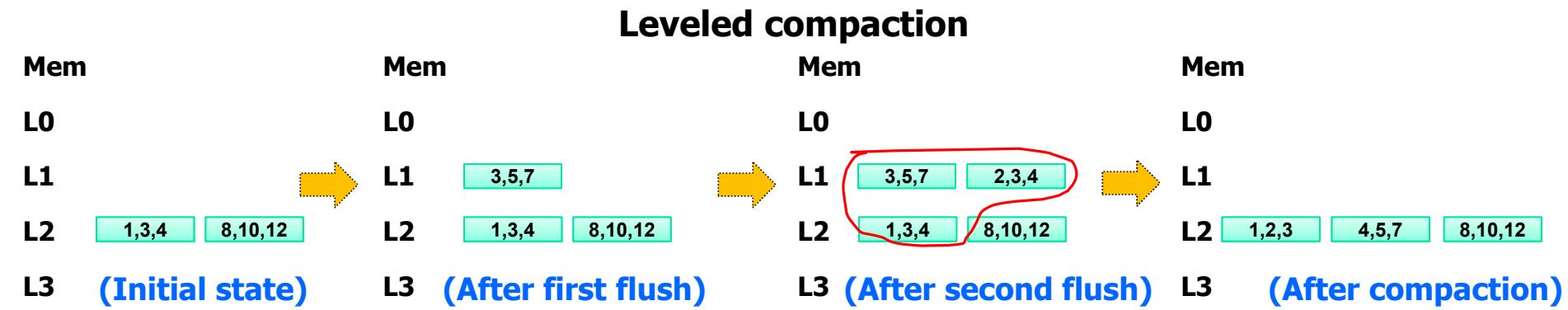
(b) Tiering Merge Policy: up to T components per level

**(Source: LSM-based Storage Techniques: A Survey, VLDB Journal'20)**

# Advance Topic 1: Compaction (9/14)

## ■ How to reduce compaction overhead: algorithmic approach

- ✓ Example
  - Leveled compaction : WAF = 12/6, SAF = 9/9 (at after compaction)
  - Tiered compaction : WAF = 11/6, SAF = 11/9 (at after compaction)

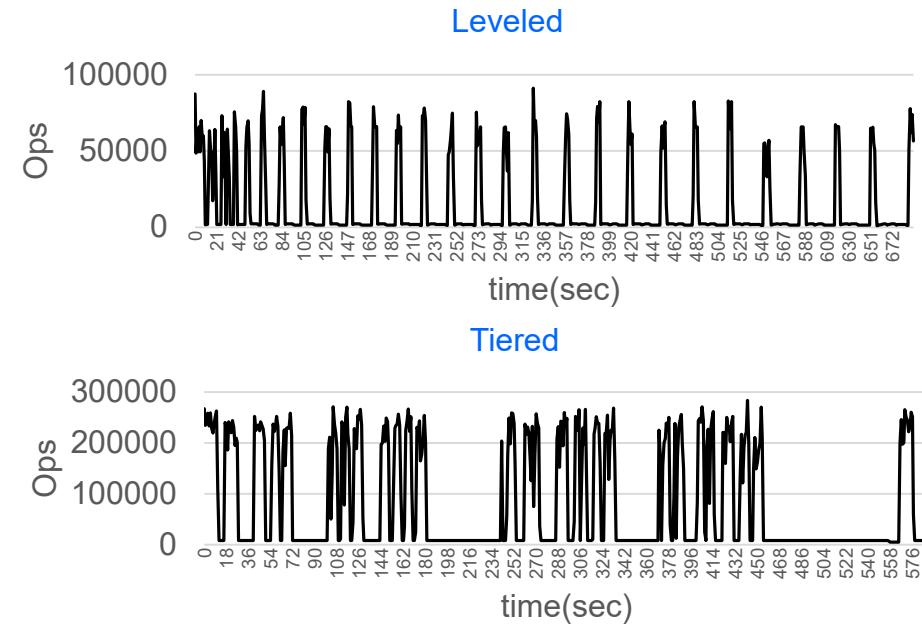
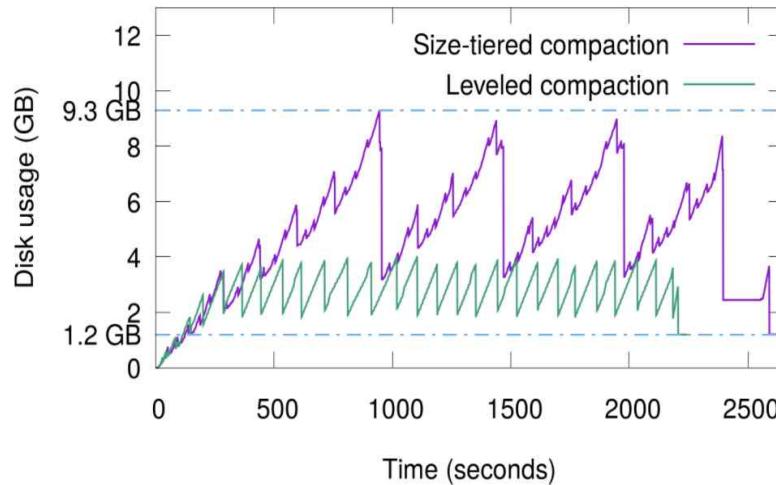


# Advance Topic 1: Compaction (10/14)

## ■ How to reduce compaction overhead: algorithmic approach

### ✓ Comparison:

- Leveled compaction : Good for SAF, Bad for WAF
- Tiered compaction : Good for WAF, Bad for SAF
  - Guarantee write once at a level (c.f. it can be multiples in the Leveled compaction)
  - Might cause intensified compactions due to recursive triggering (lazy evaluation)

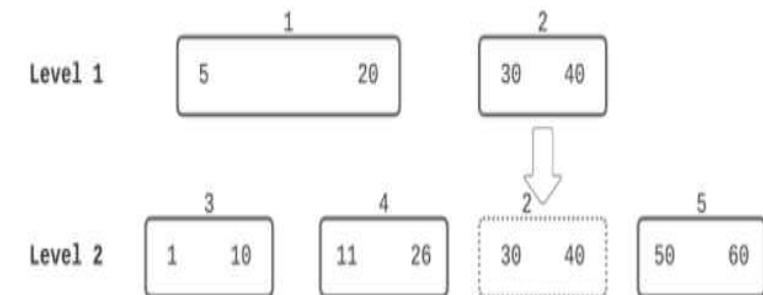
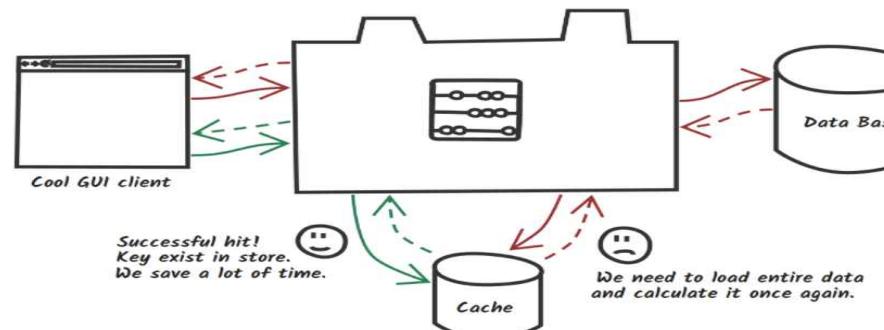


(Source: <https://www.scylladb.com/2018/01/31/compaction-series-leveled-compaction/> and Our own experiment)

# Advance Topic 1: Compaction (11/14)

## ■ Other compaction features in RocksDB

- ✓ FIFO compaction
  - L0 level only, delete SSTables based on TTL (Time-To-Live)
  - For Key-value cache (such as Redis, Memcached)
- ✓ Candidate selection and trivial move
  - Select least overlapped SSTable
- ✓ Hybrid scheme
  - Tiered + Leveled: tiered for the higher levels and leveled for others
    - . Hot data → WAF, Cold data → SAF
  - Leveled-N: Leveled (merge  $L_i$  and  $L_{i+1}$ , but can be N-runs at a level)

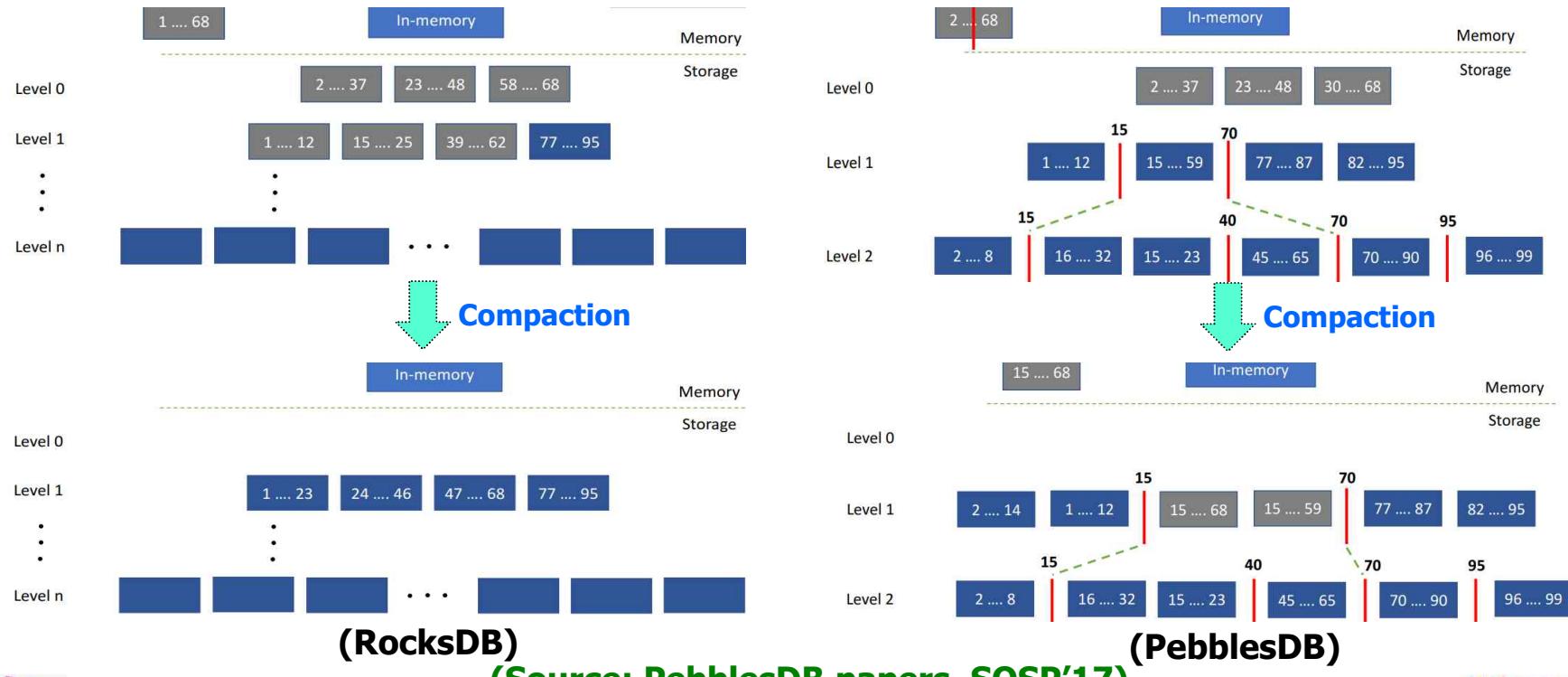


(Source: <https://medium.com/@MatthewFTech/spring-boot-cache-with-redis-56026f7da83a/> and [rocksdb/wiki](#))

# Advance Topic 1: Compaction (12/14)

## ■ Case study: PebblesDB, SOSP'17

- ✓ Goal: Reduce compaction overhead
- ✓ How to?
  - allow overlapping in a level with guard → No need to read  $L_{i+1}$
  - Write once per level vs. May deteriorate read latency

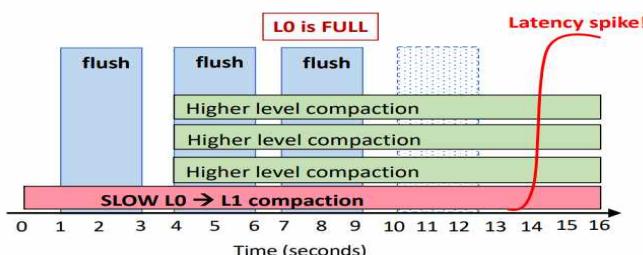


# Advance Topic 1: Compaction (13/14)

## ■ Case study: SILK, ATC'19

- ✓ Goal: Hide compaction overhead
- ✓ How to?
  - I/O scheduling: coordinate I/O bandwidth sharing to minimize interference between internal ops and client ops.
  - Motivate to design the scheduling approach in RocksDB (vice versa)

### 1. Writes Blocked Because L0 is Full.



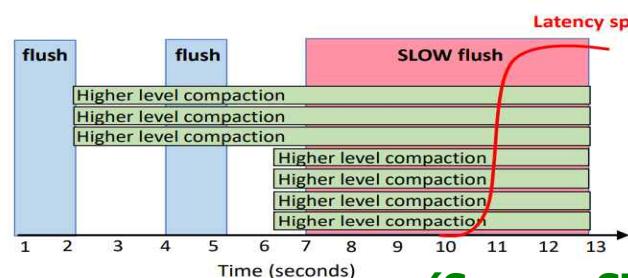
### Lessons Learned

Make sure L0 is never full.

### SILK Design

Prioritize internal operations at lower levels of the tree.

### 2. Writes Blocked Because Flushing is Slow.



Ensure sufficient I/O for flush/compactions on low levels.

Make sure other compactions do not fall behind too much.

Preempt higher level compactions if necessary.

Opportunistically allocate I/O for higher level compactions.

(Source: SILK papers, ATC'19)

# Advance Topic 1: Compaction (14/14)

---

## ■ Take-away Lessons

- ✓ Role of compaction
  - To remove old data generated by updates in a Log-structured style
  - Heavy operation and cause latency spike
- ✓ Three Compaction styles
  - Leveled, Tired, FIFO
  - Various influence in write, space and read amplification
- ✓ Explore case studies
  - PebblesDB: algorithmic approach
  - SILK: scheduling approach
- ✓ Future research
  - Hybrid: Leveled + Tiered
  - Make use of new storage (NVM, Fast SSD, ...)



(Source: [www.dreamstime.com](http://www.dreamstime.com))

**(Special thank you to DKU Embedded Members and IITP's SW StarLab (No.2021-0-01475)**

# Content

---

- What is Key-Value Store?
- RocksDB Basic
- Advance Topic 1: Compaction
- Advance Topic 2: Lookup
  - ✓ Lookup in Memtable: SkipList
  - ✓ Lookup in SSTable: Index block
  - ✓ How to enhance performance? Bloom filter
- Advance Topic 3: WAL
- Advance Topic 4: New approach
- RocksDB Practice
- Conclusion

# Advance Topic 2: Lookup (1/12)

---

## ■ Lecture Objective

- ✓ Understand several data structures for lookup
- ✓ Grasp the RocksDB approach for lookup
- ✓ Learn about details of Memtable and SSTable
- ✓ Investigate how Bloom filter can reduce the read amplification
- ✓ Explore case studies



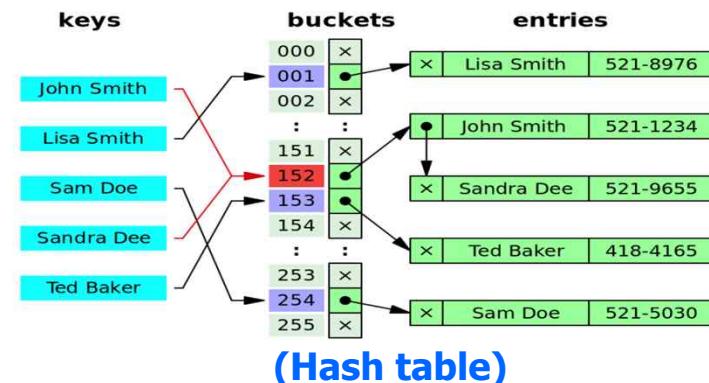
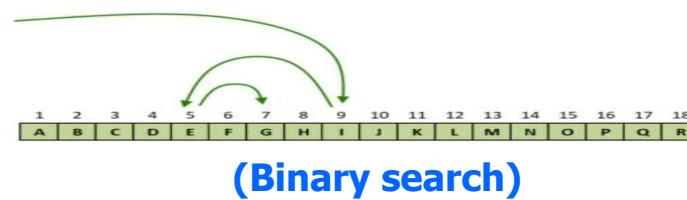
(Source: [www.dreamstime.com](http://www.dreamstime.com))

# Advance Topic 2: Lookup (2/12)

## ■ What is lookup (query)?

- ✓ Two types
  - Point lookup (single query): get a value related to a given key
  - Range lookup (scan): get values related to a range of keys
  - Example: Find a number from an address, Calculate sum or average
- ✓ How to: various data structures (e.g. array, list, sorted, hash, ...)
  - 1) linear search: ( $O(N)$ ), 2) binary search: ( $O(\log N)$ ), 3) hash: ( $O(1)$ )
  - Tradeoffs: search speed, update overhead, scan overhead, ...

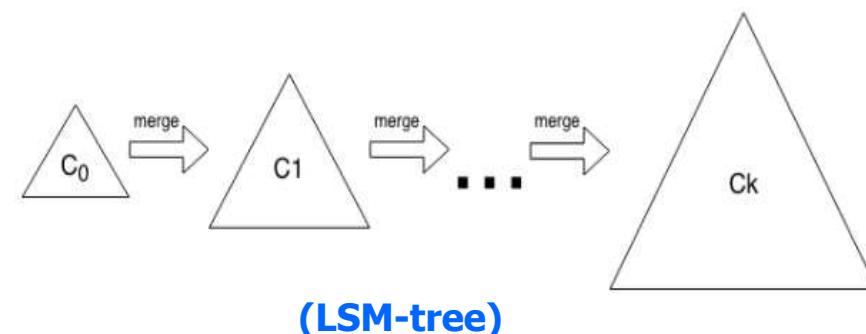
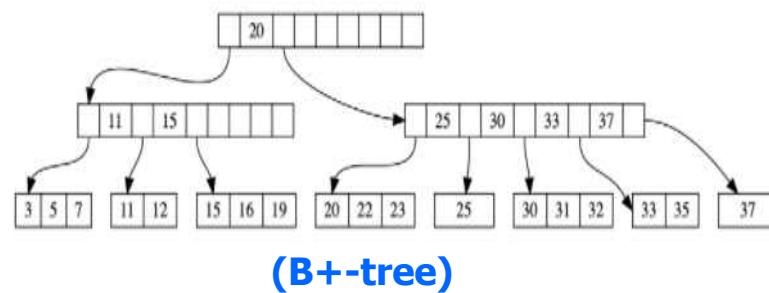
key	value
H. Shin	010-1234-****
S. Kim	010-5678-****
J. Lee	010-1111-****
J. Ahn	010-2222-****
M. Oh	010-3333-****
G. Hong	010-4444-****
J. Choi	010-88**-****
⋮	
C. Kim	010-99**-****



# Advance Topic 2: Lookup (3/12)

## ■ How to materialize lookup in DB?

- ✓ In traditional RDB (e.g. InnoDB)
  - Make use of B-tree (or B+-tree)
    - B-tree: generalized binary tree that has multiple children (n-ary tree)
    - B+-tree: all KV are stored in leaves, all leaves are linked for scan
  - Good for lookup, Bad for update due to tree reconstruction
- ✓ RocksDB
  - Make use of LSM-tree, a write-optimized data structure
    - Employ B+-tree for lookup may deteriorate its original merit (possibly can be used as secondary index)
  - Utilize its own data structures for lookup purpose including SkipList, Index block and Bloom filter



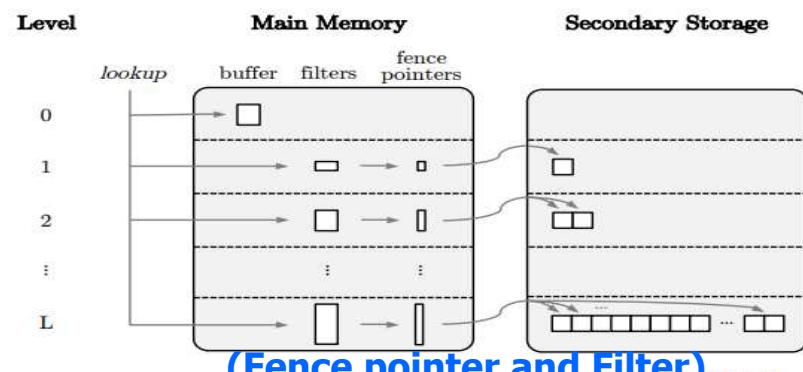
(Source: <https://tikv.org/deep-dive/key-value-engine/b-tree-vs-lsm/>)

# Advance Topic 2: Lookup (4/12)

## ■ RocksDB lookup: Overview

- ✓ Lookup procedure: 1) Memtable, then 2) SSTables (from  $L_0$  to  $L_{max}$ )
- ✓ Revisit the example discussed in the RocksDB basic section
  - get(12): find at Memtable using **Skiplist**
  - get(10): not in Memtable, need to find SSTables
    - Issue 1: many SSTables → using **Fence pointer** in memory (min and max keys in SSTable)
    - Issue 2: Both  $L_0$  and  $L_1$  → select higher level since it contains recent data
    - Issue 3: SSTable is large (default 64MB) → using **Index** for fast retrieval
  - get(8): find in  $L_1$  only
    - Issue 4: Fence pointers suggest both  $L_0$  and  $L_1$  → read  $L_0$  but there is no 8 (read amplification) → using **Bloom filter** to reduce it

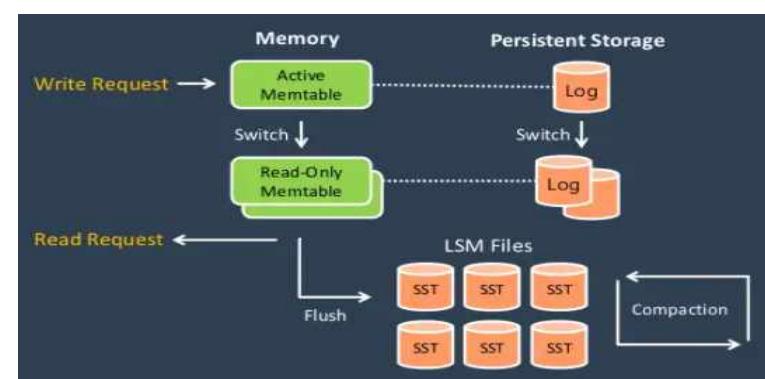
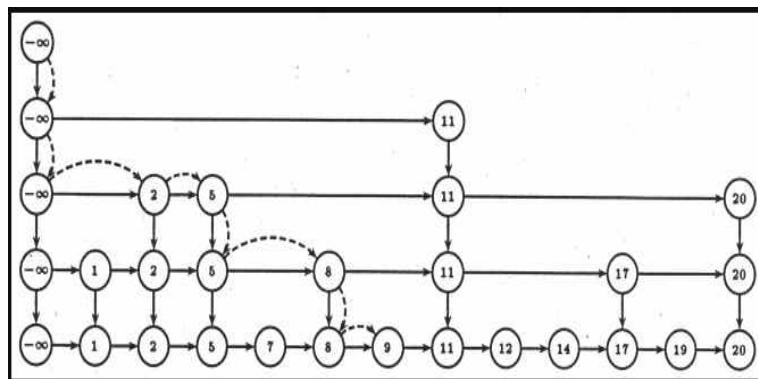
		SSTable: Fence pointers				
Mem		<12,M>	<14,Y>			
$L_0$		<2,V>	<6,W>	<10,X>		
$L_1$		<1,A>	<3,C>	<4,X>	<7,G>	<8,R>
(Initial state)						



# Advance Topic 2: Lookup (5/12)

## ■ RocksDB detail: Memtable

- ✓ KV pairs in memory, managed by SkipList (or Hashlist)
  - SkipList: a data structure with a set of sorted linked lists
    - All keys appear in the last list
    - Some keys also appear in the upper list (for fast search)
  - Good for **both lookup and scan** (get benefit both binary tree and list)
  - Useful in multithreaded system architectures
- ✓ Size options
  - write\_buffer\_size (Memtable size): default 64MB (configurable)
  - max\_write\_buffer\_number: default 2 (also configurable)



(Source: [www.osa.fu-berlin.de/bioinformatics\\_msc/en/exemplary\\_tasks/informatics\\_algorithms/index.html](http://www.osa.fu-berlin.de/bioinformatics_msc/en/exemplary_tasks/informatics_algorithms/index.html)  
and Dhruba Borthakur, [www.youtube.com/watch?v=hHsxqkcZy7c](https://www.youtube.com/watch?v=hHsxqkcZy7c))

# Advance Topic 2: Lookup (6/12)

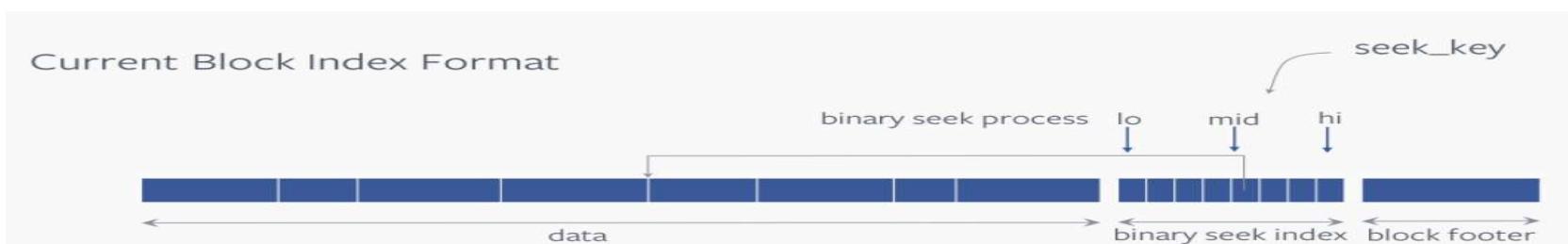
## ■ RocksDB detail: SSTable

- ✓ KV pairs in storage, managed as a file
- ✓ Issue: SSTable is large (default 64 MB)
  - Assume 1KB KV → 64,000 pairs in a file
  - A file is divided into multiple disk blocks
- ✓ Solution: well-defined SSTable format
  - 1) SSTable is divided into data blocks
  - 2) Each KV is searched using **index (B-search)**
  - 3) filter block
  - 4) other meta blocks (e.g. compression)
  - 5) metaindex: one entry for every meta blocks
  - 6) footer: pointers for metaindex & index blocks

### File format

```
<beginning_of_file>
[meta block 1]
[meta block 2]
...
[data block N]
[meta block 1: filter block]
[meta block 2: index block]
[meta block 3: compression dictionary block]
[meta block 4: range deletion block]
[meta block 5: stats block]
...
[meta block K: future extended block] (we m
[metaindex block]
[Footer] (fixed
<end_of_file>
```

(Source: [github.com/facebook/rocksdb/  
wiki/Rocksdb-BlockBasedTable-Format](https://github.com/facebook/rocksdb/wiki/Rocksdb-BlockBasedTable-Format))

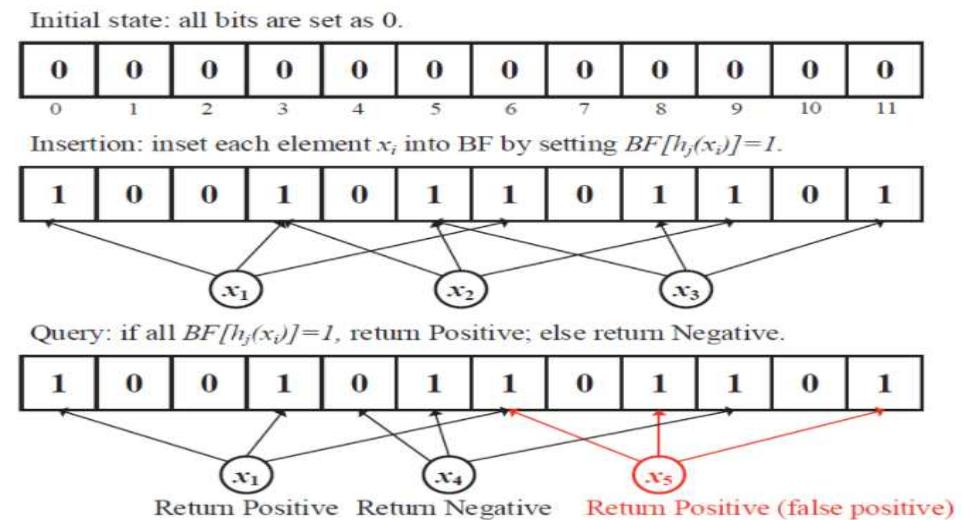
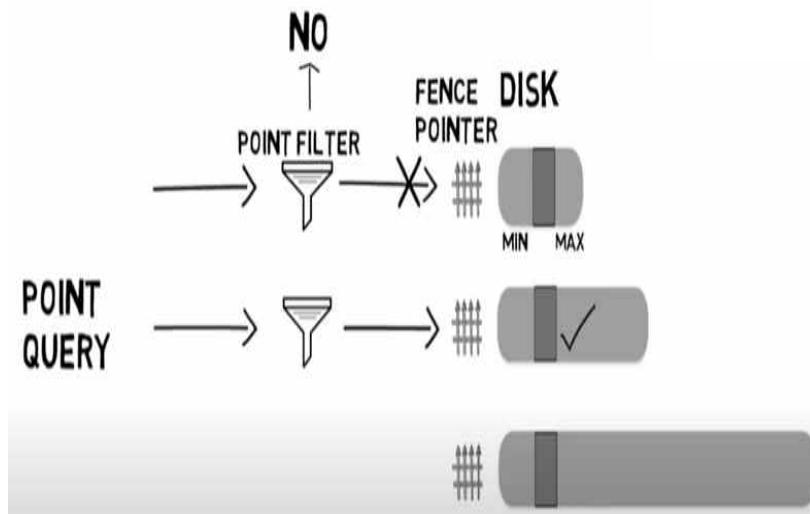


(Source: <https://rocksdb.org/blog/2018/08/23/data-block-hash-index.html>)

# Advance Topic 2: Lookup (7/12)

## ■ RocksDB detail: Bloom filter

- ✓ Used to reduce the read amplification (unnecessary read)
- ✓ Bloom filter: a data structure for identifying membership
  - Based on bits and multiple hashes
  - Good property: No false negative
  - Issue: can yield false positive → tradeoffs between bits and rate (1% false positive rate with 9.9 bits per key, from RocksDB wiki)
- ✓ Not only SSTable (per SSTable or per block) but also Memtable

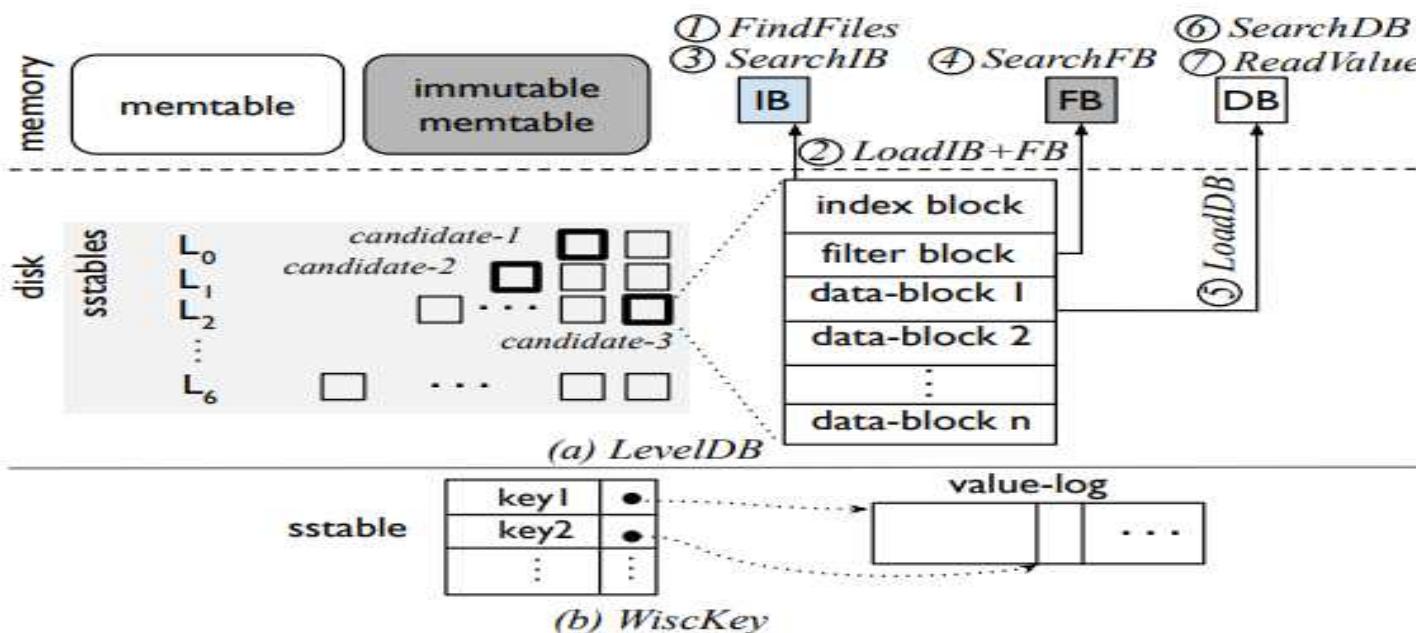


(Source: Rosetta paper, SIGMOD'20 and <https://devopedia.org/bloom-filter>)

# Advance Topic 2: Lookup (8/12)

## ■ RocksDB detail: SSTable and Bloom filter together

- ✓ 1) Find candidates using fence pointers
- ✓ 2) Load IB (Index Block) and FB (Filter Block)
- ✓ 3, 4) Search IB and FB
- ✓ 5) Load DB (Data Block)
- ✓ 6, 7) Search DB and read value



(Source: Bourbon paper, OSDI'20)

# Advance Topic 2: Lookup (9/12)

## ■ RocksDB detail: Range query

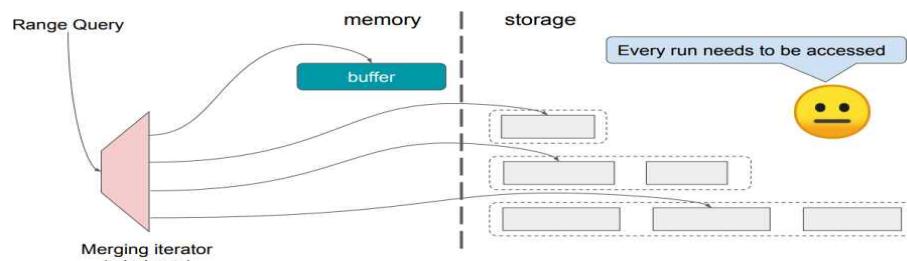
### ✓ Interface

(Source: [rocksdb/wiki](#))

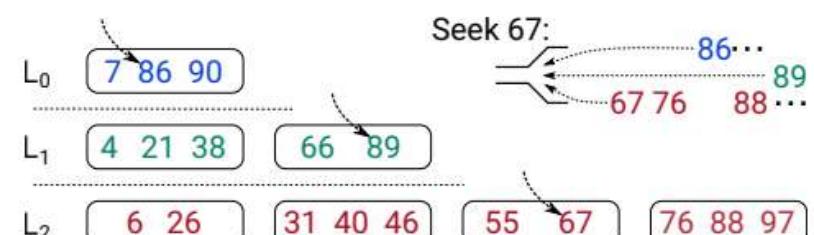
```
rocksdb::Iterator* it = db->NewIterator(rocksdb::ReadOptions());  
for (it->Seek(start);  
     it->Valid() && it->key().ToString() < limit;  
     it->Next()) {  
    ...  
}  
assert(it->status().ok()); // Check for any errors found during the scan
```

### ✓ Internals

- Range query: need to access all levels (c.f. point lookup: stop if found)
- Using min-heap: example of scan(67, 95)
  - 1) find a key larger than seek for every levels (cursor)
  - 2) construct the min-heap using cursors (67, 89, 86)
  - 3) get the smallest key (67), 4) if it is larger than limit, stop, 5) print it, 6) find the next cursor (76), 7) insert it into the min-heap
  - Repeat 3) ~ 7) steps → “67, 76, 86, 88, 89, 90, stop”
- Prefix bloom filter is used to reduce unnecessary read during scan



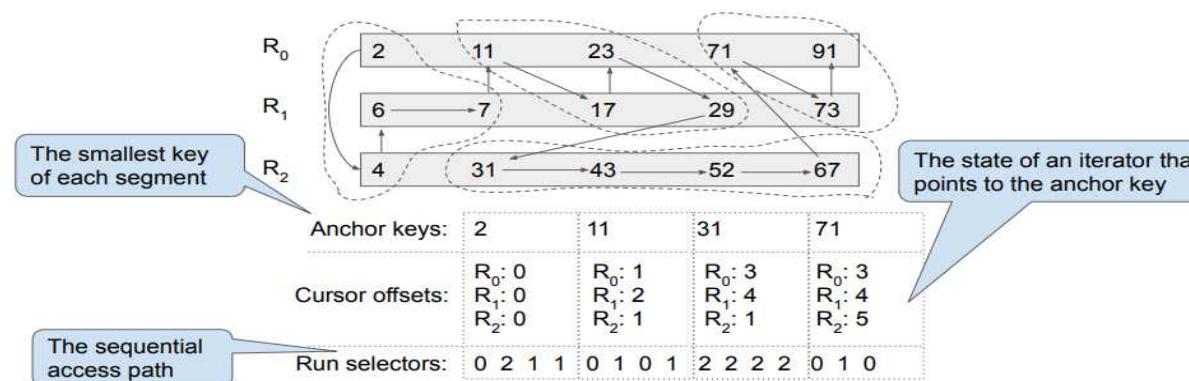
(Source: [REMIN paper, FAST'21](#))



# Advance Topic 2: Lookup (10/12)

## ■ Case study: REMIX, FAST'21

- ✓ Goal: Reduce range query overhead
- ✓ How to?
  - Construct a new data structure that has a global sorted view among runs
    - . 1) Anchor keys: the smallest key in each segment
    - . 2) Cursor offsets: for identifying runs that has keys in each segment
    - . 3) Run selector: for logical sorted view → divide into segments
  - Comparison: example of scan (15, 50)
    - . Traditional: min-heap with “23, 17, 31”, do min-heap operations repeatedly
    - . REMIX: anchor 11, using Run and Cursor for next key (e.g. R0/C1, R1/C2, R0/C2, R1/C3, ...) → no min-heap operations

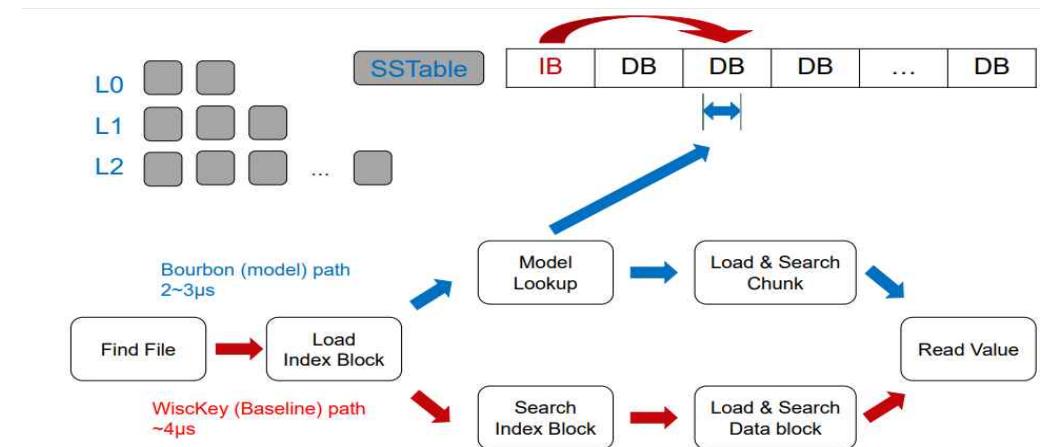
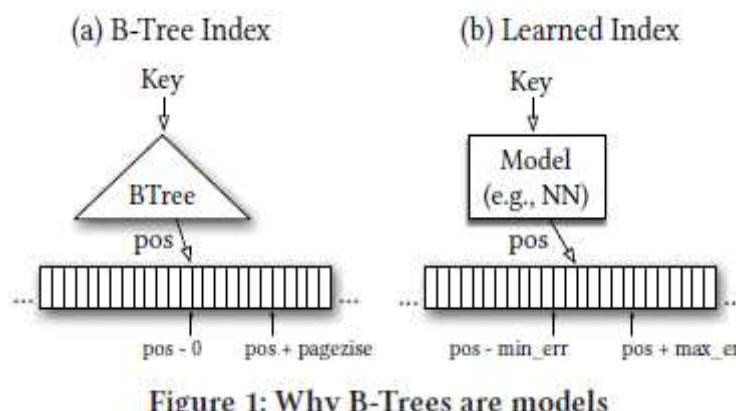


(Source: REMIX papers, FAST'21)

# Advance Topic 2: Lookup (11/12)

## ■ Case study: Bourbon, OSDI'20

- ✓ Goal: Reduce lookup overhead
- ✓ How to? make use of machine learning (Learned index)
  - Idea: Index can be viewed as Learning
    - B-tree: 1) new key → tree reconstruction, 2) lookup: key → position
    - Model: 1) learning from “key -> position” data, 2) lookup: inference from key to expected position (not deterministic, probabilistic)
  - Bourbon: employ Learned index into Key-value store
    - Model lookup instead of Index block access
    - Suitable for Read intensive workloads



(Source: Learned Index, SIGMOD'18 and Bourbon papers, OSDI'21)

# Advance Topic 2: Lookup (12/12)

---

## ■ Take-away Lessons

- ✓ Basic data structures for lookup
  - Linked-List, Binary tree, Hash, B+-tree, ...
- ✓ RocksDB approach for lookup
  - Memtable: Skiplist for both fast lookup and scan
  - SSTable: Index block and Bloom Filter
  - Range scan: min heap (priority queue)
- ✓ Investigate how Bloom filter can reduce the read amplification
  - Tradeoff between Space/CPU overhead and false positive ratio
  - Active used for other purposes (e.g. Prefix Bloom Filter)
- ✓ Explore case studies
  - REMIX: Build extra index for range query
  - Bourbon: Machine learning in Key-Value Store



(Source: [www.dreamstime.com](http://www.dreamstime.com))

(Special thank you to DKU Embedded Members and IITP's SW StarLab (No.2021-0-01475)

# Content

---

- What is Key-Value Store?
- RocksDB basic
- Advance Topic 1: Compaction
- Advance Topic 2: Lookup
- Advance Topic 3: WAL
  - ✓ Durability: WAL and Recovery
  - ✓ Consistency: atomicity and Transaction
  - ✓ Integrity: Checksum
- Advance Topic 4: New approach
- RocksDB Practice
- Conclusion

# Advance Topic 3: WAL (1/11)

---

## ■ Lecture Objective

- ✓ Understand the requirements for DB Reliability
- ✓ Find out the role of WAL (Write Ahead Log)
- ✓ Learn about how to support transaction
- ✓ Explore how to keep integrity
- ✓ Examine the performance effect of WAL



(Source: [www.dreamstime.com](http://www.dreamstime.com))

# Advance Topic 3: WAL (2/11)

## ■ Performance and Reliability

- ✓ Techniques such as Compaction and Bloom filter that we have learnt so far → for performance enhancement
- ✓ DB demands reliability (often degrade performance and vice versa)

## ■ Requirements for high reliability

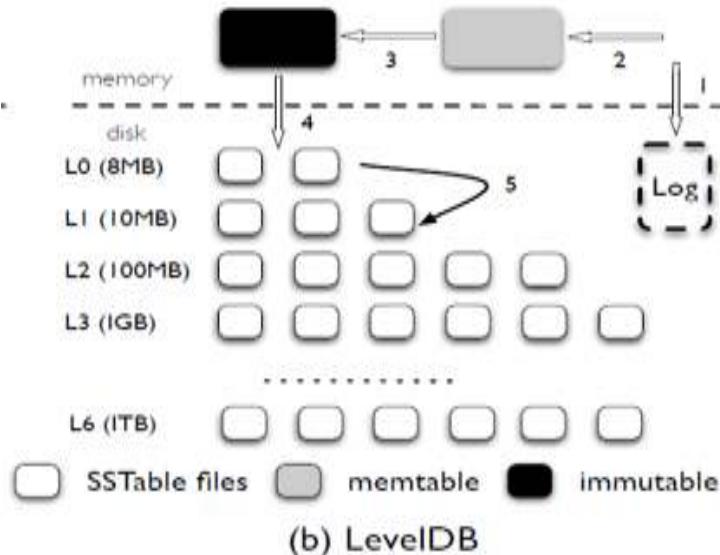
- ✓ Durability: make data persistent
- ✓ Atomicity: all or nothing
- ✓ Consistency: stay in a valid state or agreement among parts
- ✓ Integrity: assurance of data accuracy
- ✓ Recovery: go back to a valid state



# Advance Topic 3: WAL (3/11)

## ■ Reliability matters

- ✓ Put KV: first Memtable, and SSTable via flush
  - By the way, main memory is volatile → What if power failure occurs?
  - Lost user data is a big risk in a DB company → require **durability**
- ✓ RocksDB approach
  - **WAL** (Write Ahead Log, also called as journal or intent log) → in order to support durability and easy recovery



-rw-r--r-- 1 root root 37922501 7월 1 15:44 000214.sst
-rw-r--r-- 1 root root 37920200 7월 1 15:44 000216.sst
-rw-r--r-- 1 root root 37910828 7월 1 15:44 000219.sst
-rw-r--r-- 1 root root 37906740 7월 1 15:44 000221.sst
-rw-r--r-- 1 root root 37905482 7월 1 15:44 000224.sst
-rw-r--r-- 1 root root 37909294 7월 1 15:44 000227.sst
-rw-r--r-- 1 root root 66217640 7월 1 15:44 000228.log
-rw-r--r-- 1 root root 37892964 7월 1 15:44 000229.sst
-rw-r--r-- 1 root root 18621323 7월 1 15:44 000231.log
-rw-r--r-- 1 root root 16 7월 1 15:43 CURRENT
-rw-r--r-- 1 root root 37 7월 1 15:43 IDENTITY
-rw-r--r-- 1 root root 0 7월 1 15:43 LOCK
-rw-r--r-- 1 root root 572956 7월 1 15:44 LOG
-rw-r--r-- 1 root root 14592 7월 1 15:44 MANIFEST-000004
-rw-r--r-- 1 root root 6180 7월 1 15:43 OPTIONS-000007

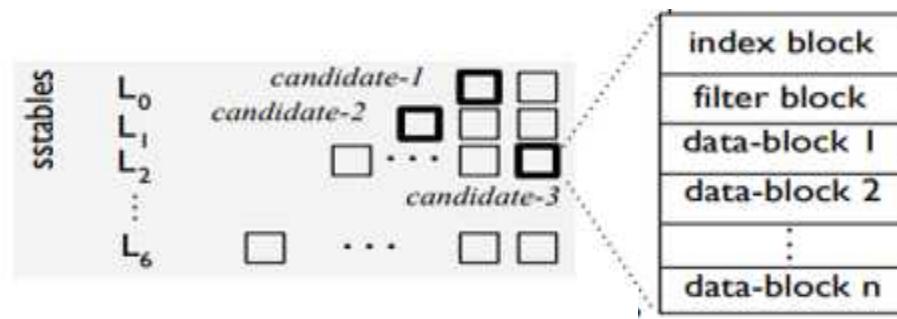
(Source: Wisckey paper, FAST'16 and our experiment)

# Advance Topic 3: WAL (4/11)

## ■ Reliability matters

- ✓ A set of operations are connected
  - Assume that the put() performs while the delete() does not → duplicated
  - Assume that the delete() performs while the put() does not (various causes including scheduling, caching, ...) → Data loss
  - How to overcome? require **atomicity**
- ✓ A set of blocks are connected
  - Assume that data is written while index is not → Space leak
  - Assume that index is written while data is not → Garbage read
  - How to overcome? require **consistency**
- ✓ RocksDB approach
  - WriteBatch and Transaction

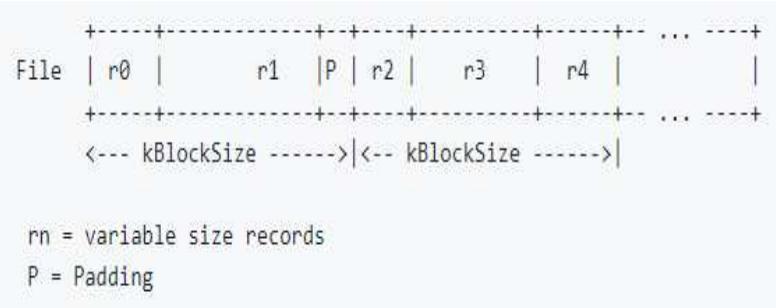
```
std::string value;
rocksdb::Status s = db->Get(rocksdb::ReadOptions(), key1, &value);
if (s.ok()) s = db->Put(rocksdb::WriteOptions(), key2, value);
if (s.ok()) s = db->Delete(rocksdb::WriteOptions(), key1);
```



# Advance Topic 3: WAL (5/11)

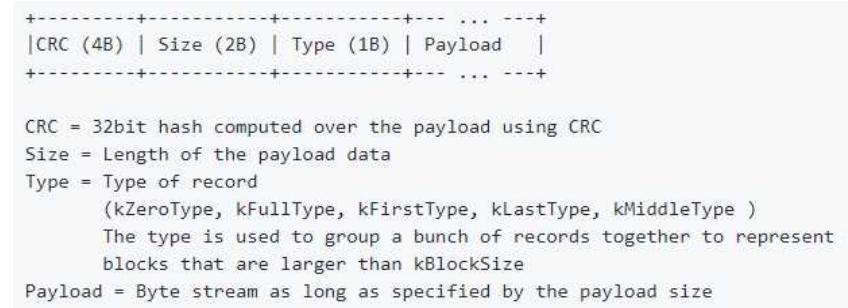
## ■ RocksDB detail: WAL

- ✓ What?
  - Logging every update
- ✓ How?
  - log file: consists of kBlocks (32KB default)
  - kBlock: a set of record (KV) and padding
    - . Type: Full → KV in a record, First/Middle/Last → KV in multiple records
- ✓ Lifecycle
  - A WAL file is created when a new DB is opened
  - A WAL is deleted (or archived if archival is enabled) when all data in the WAL have been persisted to SST files.



(Log file format)

(Source: <https://github.com/facebook/rocksdb/wiki/Write-Ahead-Log-File-Format>)



(Record format)

# Advance Topic 3: WAL (6/11)

## ■ RocksDB detail: WAL

### ✓ Configuration

- DisableWAL, max\_total\_wal\_size, wal\_dir (for different storage), ...
- Sync mode
  - Non-sync mode (sync = false, default)
    - OS-level or manually sync (e.g. call SyncWAL() periodically)
    - Guarantees process crash consistency only
  - Sync mode (sync = true)
    - Do sync() before returning to user
    - Guarantees both process and system crash consistency, but heavy weight

### ✓ Recovery

- Redo WAL records
- Policies: 1) try all (one corrupted record fails recovery), 2) try all except last WAL, 3) until first corrupted record (others from replica), 4) try your best and ignore corrupted records



(Source: Dhruba Borthakur, [www.youtube.com/watch?v=hHsxqkcZy7c](https://www.youtube.com/watch?v=hHsxqkcZy7c))

# Advance Topic 3: WAL (7/11)

## ■ RocksDB detail: Writebatch and Transaction

- ✓ Writebatch
  - For atomicity (in addition for better performance)
- ✓ Transaction
  - For atomicity, consistency (always valid state) and concurrency
  - Add transaction begin and commit into a log

Tr Begin	Put, key2, v2	Delete, key1	Tr Commit
----------	---------------	--------------	-----------

- Recovery
  - Committed → redo, Not committed → undo (rollback)
- Pessimistic (default) vs Optimistic
  - Always locking for thread safe vs validating at commit time

```
#include "rocksdb/write_batch.h"
...
std::string value;
rocksdb::Status s = db->Get(rocksdb::ReadOptions(), key1, &value);
if (s.ok()) {
    rocksdb::WriteBatch batch;
    batch.Delete(key1);
    batch.Put(key2, value);
    s = db->Write(rocksdb::WriteOptions(), &batch);
}
```

### (Writebatch)

(Source: <https://github.com/facebook/rocksdb/wiki/Basic-Operations-and-transactions>)

```
TransactionDB* txn_db;
Status s = TransactionDB::Open(options, path, &txn_db);

Transaction* txn = txn_db->BeginTransaction(write_options, txn_options);
s = txn->Put("key", "value");
s = txn->Delete("key2");
s = txn->Merge("key3", "value");
s = txn->Commit();
delete txn;
```

### (Transaction)

# Advance Topic 3: WAL (8/11)

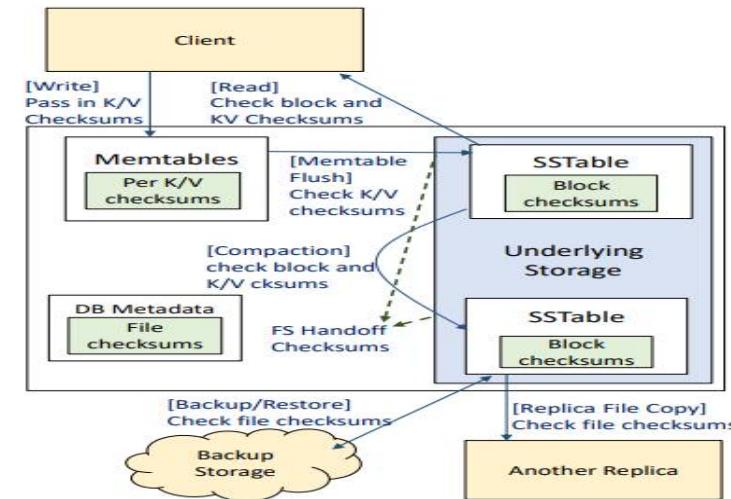
## ■ RocksDB detail: Checksum

- ✓ Why?
  - Atomic and consistent, but what if bit flip occurs (e.g. due to Flash memory error such as bad block, retention error and disturbance)
  - Require **integrity**
- ✓ Checksum: a well-known method for keep Integrity
  - A small-sized block of data derived from main data (e.g. using parity and CRC (Cyclic Redundance Checks))
- ✓ RocksDB approach
  - Four layers
    - Application, Block, File, Handoff

+-----+ <td> CRC (4B) </td> <td>Size (2B)</td> <td>Type (1B)</td> <td>Payload</td> <td>-----+<td>... -----+</td></td>	CRC (4B)	Size (2B)	Type (1B)	Payload	-----+ <td>... -----+</td>	... -----+
---	----------	-----------	-----------	---------	----------------------------	------------

CRC = 32bit hash computed over the payload using CRC  
Size = Length of the payload data  
Type = Type of record  
(kZeroType, kFullType, kFirstType, kLastType, kMiddleType )  
The type is used to group a bunch of records together to represent blocks that are larger than kBlockSize  
Payload = Byte stream as long as specified by the payload size

**(Record format)**



**(Multi-layer checksum)**

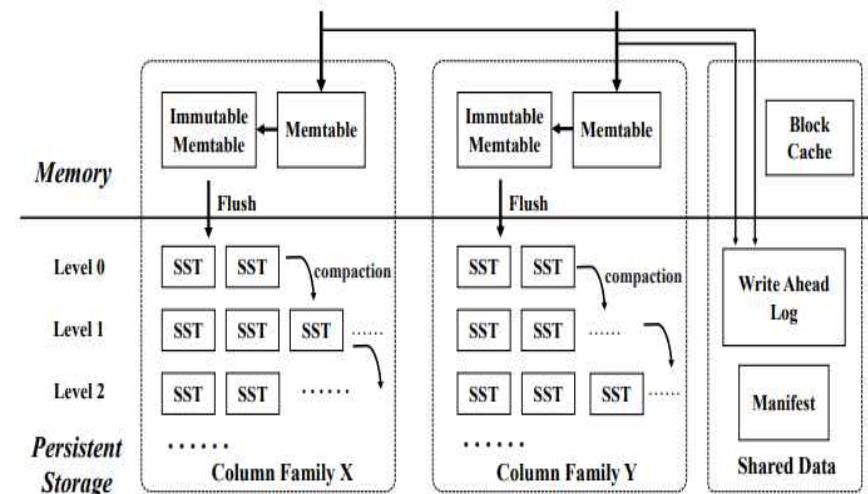
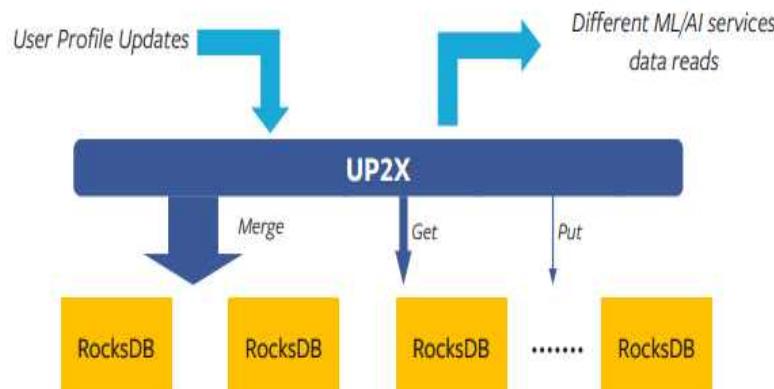
Source: [github.com/facebook/rocksdb/wiki/Write-Ahead-Log-File-Format](https://github.com/facebook/rocksdb/wiki/Write-Ahead-Log-File-Format) & RocksDB Experience papers, FAST'21

# Advance Topic 3: WAL (9/11)

## ■ Other features

- ✓ Merge operator
  - A common pattern: update an existing value → need read-modify-write
  - How about exploit **logging** and **compaction** → write only → combine later (e.g. at the compaction time)
- ✓ Column family
  - A logical partition of DB (e.g. name → phone, name → photo)
  - Support atomic write across column families (share WAL)

## Use Case 3: UP2X

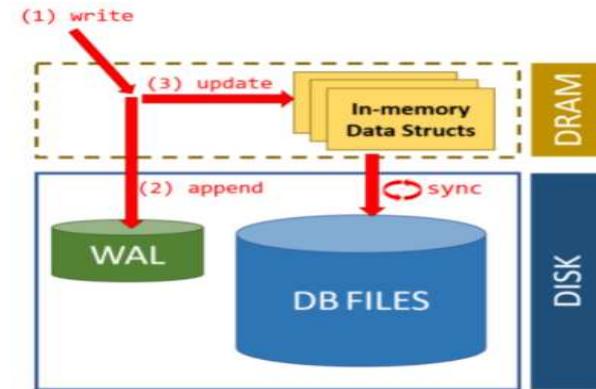


(Source: RocksDB Characterizing papers, FAST'20)

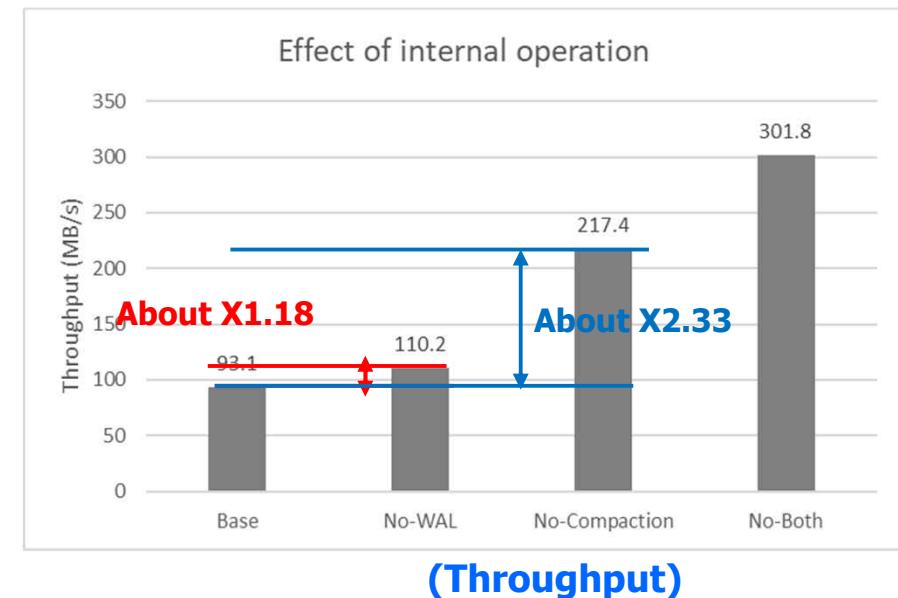
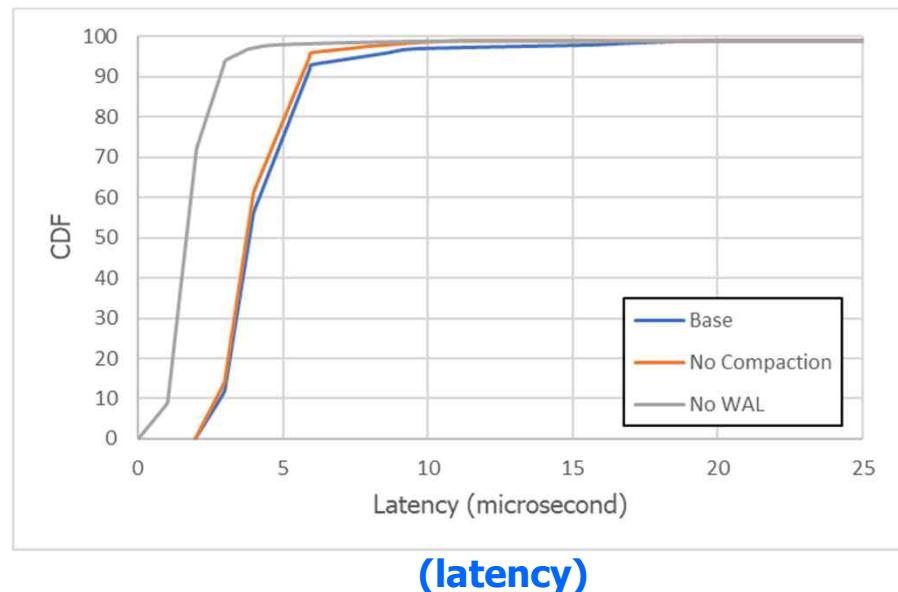
# Advance Topic 3: WAL (10/11)

## ■ Performance consideration

- ✓ WAL incurs additional I/Os
  - Affect latency
- ✓ How to mitigate?
  - Delayed write (No sync)
  - Fine-grained WAL or parallelism
  - Make use of NVM (Non-volatile memory)



(Source: [software.intel.com](http://software.intel.com), "optimizing-WAL-with-intel-optane-persistent-memory")



(Source: our own experiment)

# Advance Topic 3: WAL (11/11)

---

## ■ Take-away Lessons

- ✓ Understand the requirements for DB Reliability
  - Durability, Atomicity, Consistency, Integrity and Recovery
- ✓ RocksDB approach
  - WAL (Write Ahead Log): intent logging before actual actions
  - Writebatch: group a set of updates
  - Transaction: maintain valid states by undo or redo
  - Checksum: at various layers
- ✓ Other features
  - Merge operator: read-modify-write → log and compaction
  - Column family
  - Manifest: Initial point for restart or recovery
  - Checkpoint and snapshot



(Source: [www.dreamstime.com](http://www.dreamstime.com))

(Special thank you to DKU Embedded Members and IITP's SW StarLab (No.2021-0-01475)

# Content

---

- What is Key-Value Store?
- RocksDB basic
- Advance Topic 1: Compaction
- Advance Topic 2: Lookup
- Advance Topic 3: WAL
- Advance Topic 4: New approach
  - ✓ Layout: separate key and value
  - ✓ Make use of New storage and ISP
  - ✓ Disaggregate storage and computing
- RocksDB Practice
- Conclusion

# Advance Topic 4: New approach (1/12)

---

## ■ Lecture Objective

- ✓ Understand the merit of key value separation
- ✓ Explore how to utilize new storage for key-value store
- ✓ Learn about ISP (In-Storge Processing)
- ✓ Find out how to design a distributed key-value store
- ✓ Summarize RocksDB features

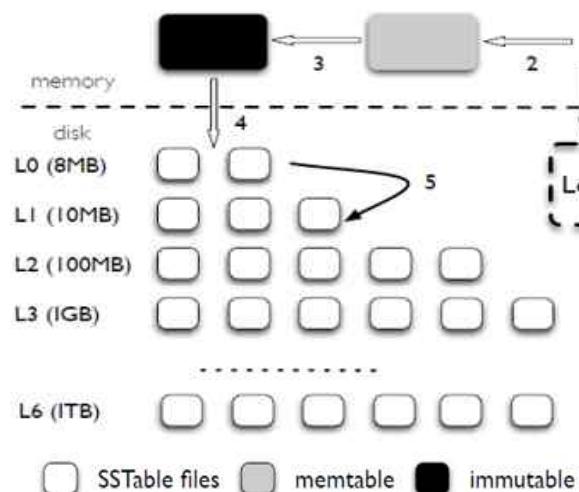


(Source: [www.dreamstime.com](http://www.dreamstime.com))

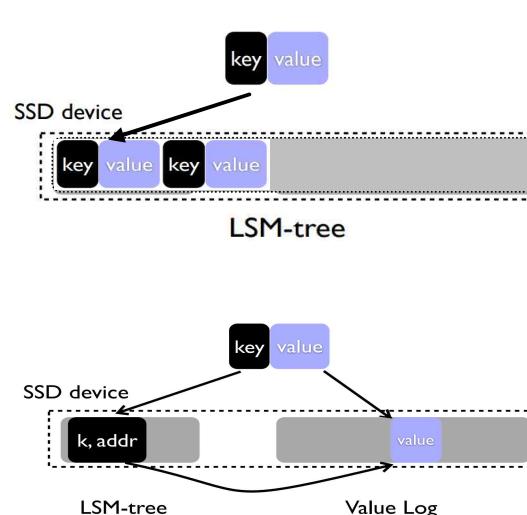
# Advance Topic 4: New approach (2/12)

## ■ Case study: WiscKey, FAST'16

- ✓ Goal: Reduce Compaction Overhead
- ✓ How to?
  - Existing approach: scheduling, algorithmic, ...
  - Proposal: key-value separation
    - Compaction trigger → when the size of  $L_i$  is larger size than threshold
    - LSM-tree: key-value pair → compaction uses key only → How about managing key only in a LSM-tree → smaller LSM-tree → less compaction
    - Also, giving a positive impact on caching and searching



(LSM-tree basic)



(Key-value separation)

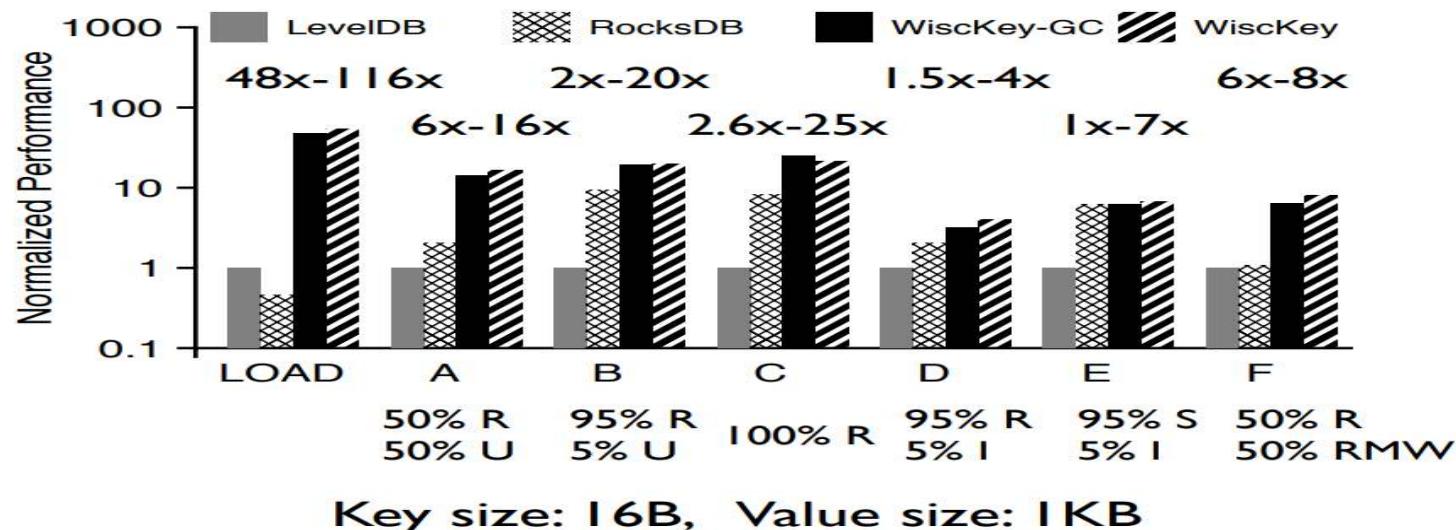
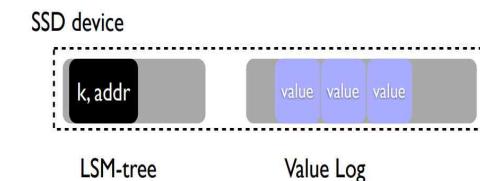
	load 100 GB database	LevelDB	WiscKey
	limits of files	num of files	num of files
L0	9	7	
L1 (5)	30	11	
L2 (50)	365	127	
L3 (500)	2184	460	
L4 (5000)	15752	0	
L5 (50000)	23733	0	
L6 (500000)	0	0	

Small LSM-tree: less compaction, fewer levels to search, and better caching

# Advance Topic 4: New approach (3/12)

## ■ Case study: WiscKey, FAST'16 (cont')

- ✓ How to? (cont')
  - Value log implementation
    - 1) Queue with head (for new insert) and tail (for GC), 2) Lightweight GC
  - Range query
    - thread pool for concurrent processing (make use of SSD parallelism)
- ✓ Lessons
  - Simple idea → Substantial gain (adopted by many companies)



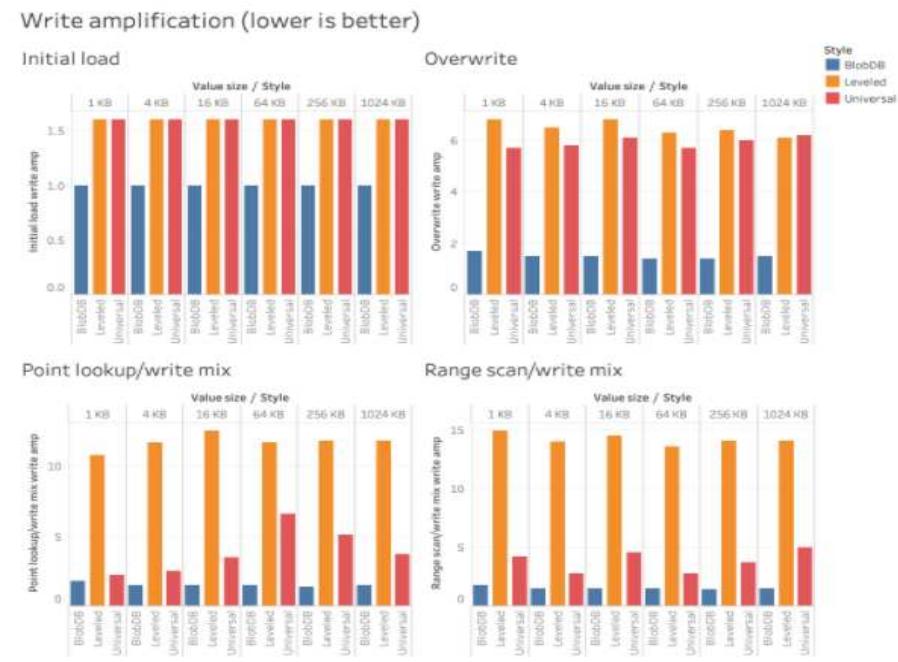
# Advance Topic 4: New approach (4/12)

## ■ Case study: Wisckey, FAST'16 (cont')

### ✓ BlobDB

- Wisckey implementation in RocksDB
- Overall better performance
- Initial stage: being actively added new features (e.g. merge operator, range scan optimization, ...)

The screenshot shows the GitHub repository for BlobDB. The URL is [github.com/facebook/rocksdb/wiki/BlobDB](https://github.com/facebook/rocksdb/wiki/BlobDB). The page title is "BlobDB". It contains sections for "Legacy implementation" and "Integrated implementation". The "Legacy implementation" section discusses the original API based on the `StackableDB` interface, noting it is geared towards FIFO/TLI use cases and lacks recovery mechanisms. The "Integrated implementation" section notes it uses the well-known `rocksdb::DB` API. A sidebar on the right lists various contents such as Overview, Requirements, and Basic Operations.



(Source: [github.com/facebook/rocksdb/wiki/BlobDB](https://github.com/facebook/rocksdb/wiki/BlobDB), [rocksdb.org/blog/2021/05/26/integrated-blob-db.html](https://rocksdb.org/blog/2021/05/26/integrated-blob-db.html))

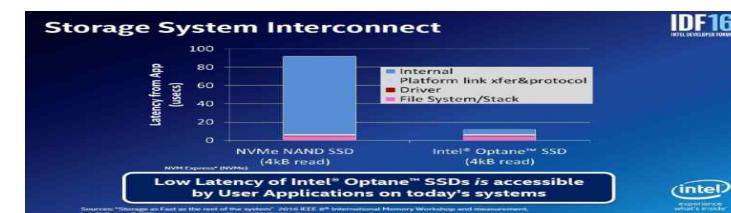
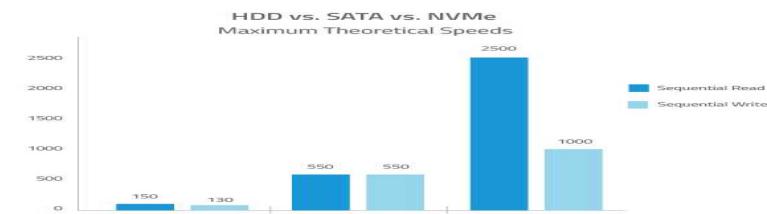
# Advance Topic 4: New approach (5/12)

## ■ Case study: MatrixKV, ATC'20

- ✓ Goal: Make use of NVM (Non-Volatile Memory) or faster storage (e.g. NVMe SSD, Optane SSD)
  - NVM
    - Characteristics: both non-volatile and byte-addressable, fast
    - Types: PRAM (PCM), MRAM, RRAM, NVDIMM, 3D-Xpoint, Optane, ...
  - SSD
    - SATA, NVMe, Optane SSD (SSD based on NVM) → hybrid storage
  - Accelerate compaction via fast NVM/SSD



(NVM characteristics)



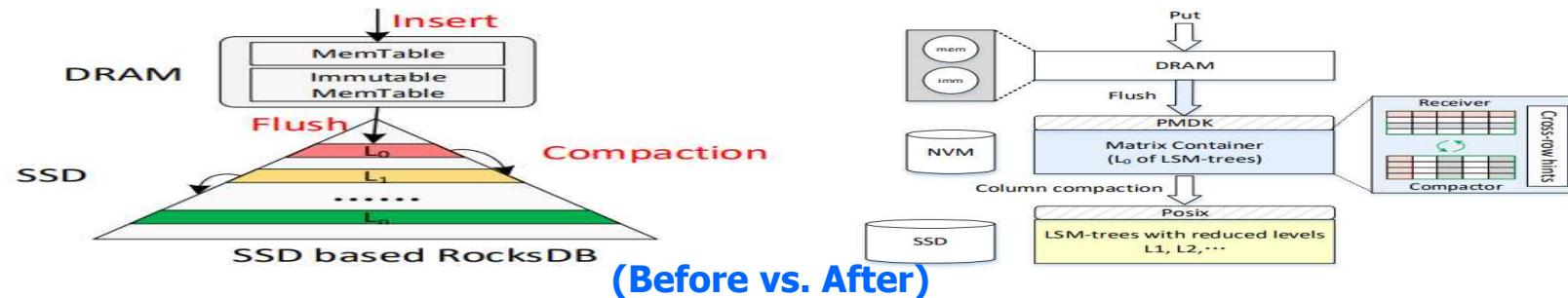
(SATA, NVMe, Optane SSD)

# Advance Topic 4: New approach (6/12)

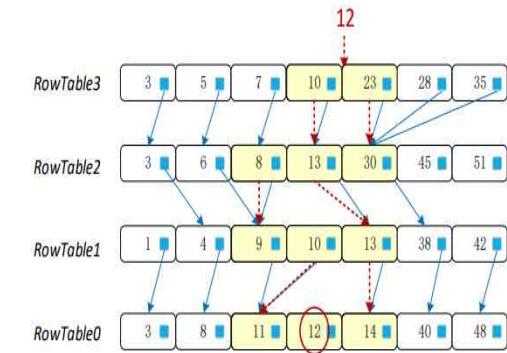
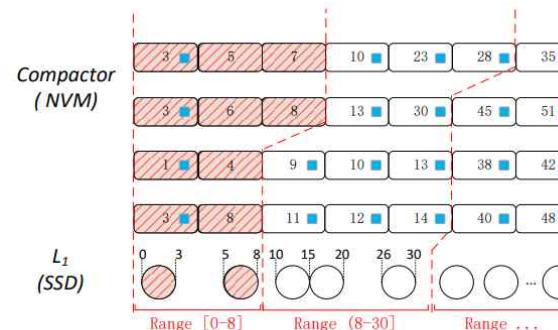
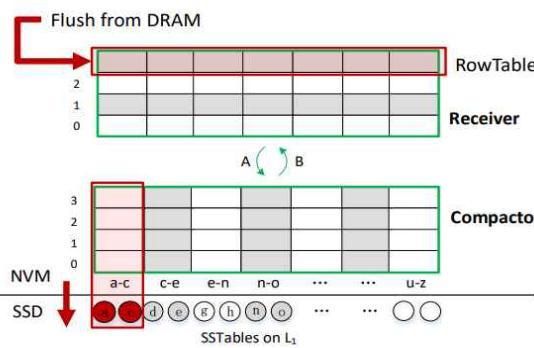
## ■ Case study: MatrixKV, ATC'20

- ✓ How to?

- Place L0 in NVM (can extend to L1, L2, ...), others in SSD



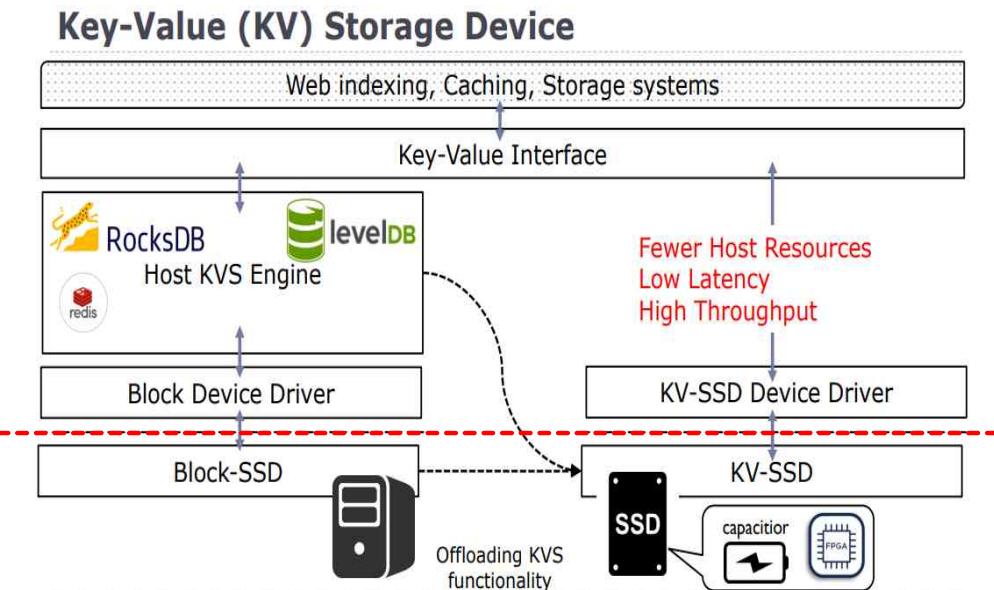
- Design several novel techniques
  - Matrix form in NVM: 1) receiver (like tiered using rows), 2) compactor
  - Fine-grained compaction: based on non-overlapped compaction
  - Cross-row hint search: using forward pointers for fast search



# Advance Topic 4: New approach (7/12)

## ■ Case study: Pink, ATC'20

- ✓ Goal: Make use of ISP (In-Storage Processing)
  - SSD: equip with CPUs (multicores), large DRAM, HW accelerator, ...
    - Why CPUs in SSD? For running SW called FTL (Flash Translation Layer) to handle unique Flash features such as erase-before-write and endurance
  - KVSSD
    - Conventional architecture: APP+KV Store in host, FTL in SSD, Block interface
    - KVSSD architecture: APP in host, KV Store + FTL in SSD, KV interface



(Source: [www.quietpc.com/multicore-pcie-ssd](http://www.quietpc.com/multicore-pcie-ssd), [m.blog.naver.com/ki630808/221968329122](http://m.blog.naver.com/ki630808/221968329122), and Samsung KVSSD paper)

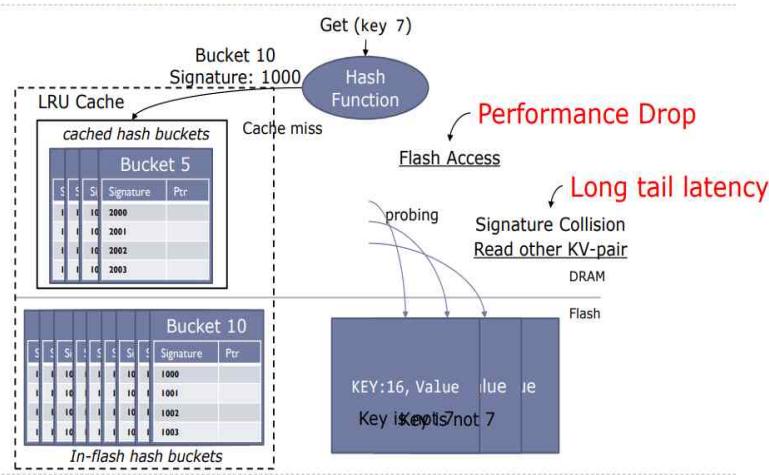
# Advance Topic 4: New approach (8/12)

## ■ Case study: PinK, ATC'20

### ✓ How to?

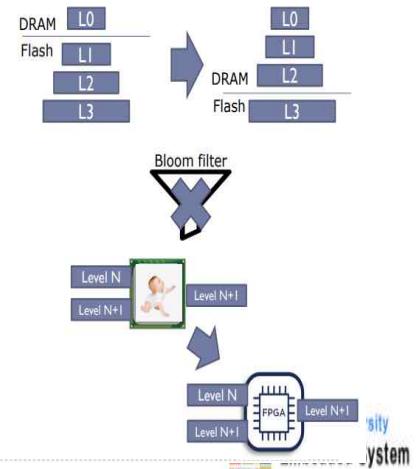
- KV store in SSD: Hash based vs. LSM-tree based → Hash causes performance drop and long-tail latency due to hash collision and overflow → LSM-tree based approach
- LSM-tree issues
  - Long tail due to SSTable accesses → Memory pinning (L0~L2)
  - Bloom filter overhead (due to embedded CPUs) → No Bloom Filter (pinning makes it feasible)
  - Compaction overhead (lead to I/O overhead) → HW accelerator

### Problem of Hash-based KV-SSD



### PinK : New LSM-tree-based KV-SSD

- ▶ Long tail latency?
  - ▶ Using "Level-pinning"
- ▶ CPU overhead?
  - ▶ "No Bloom filter"
  - ▶ "HW accelerator" for compaction
- ▶ I/O overhead?
  - ▶ Reducing compaction I/O by level-pinning
  - ▶ Optimizing GC by reinserting valid data to LSM-tree

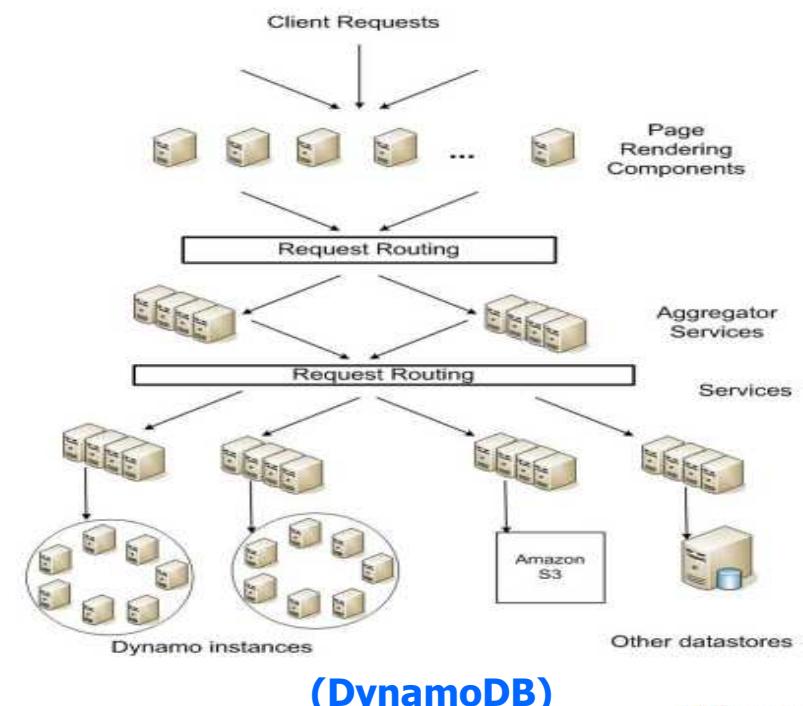
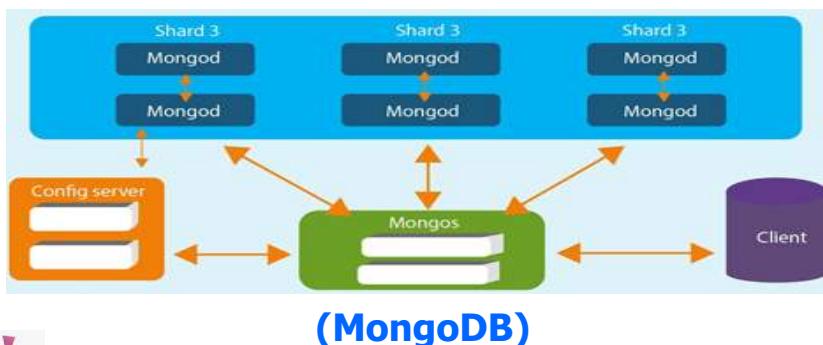
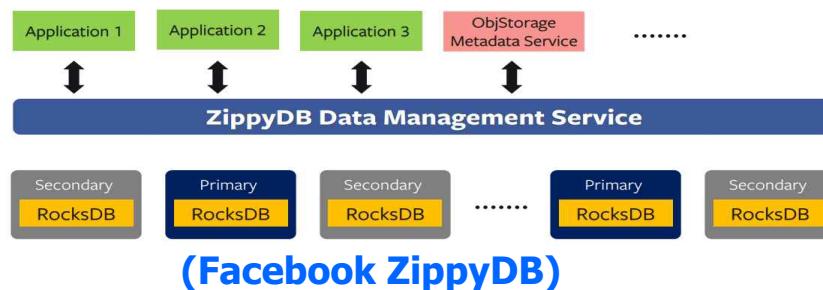


# Advance Topic 4: New approach (9/12)

## ■ Case study: Hailstorm, ASPLOS'20

- ✓ Goal: Load balancing for Distributed Key-Value Store
  - Distributed key-value store: core infrastructure for data service companies to obtain availability and aggregated performance
  - Issue 1. Key distribution: mapping table, hash, circular, ...
  - Issue 2. Load imbalance: main topic in this study

### Use Case 2: ZippyDB

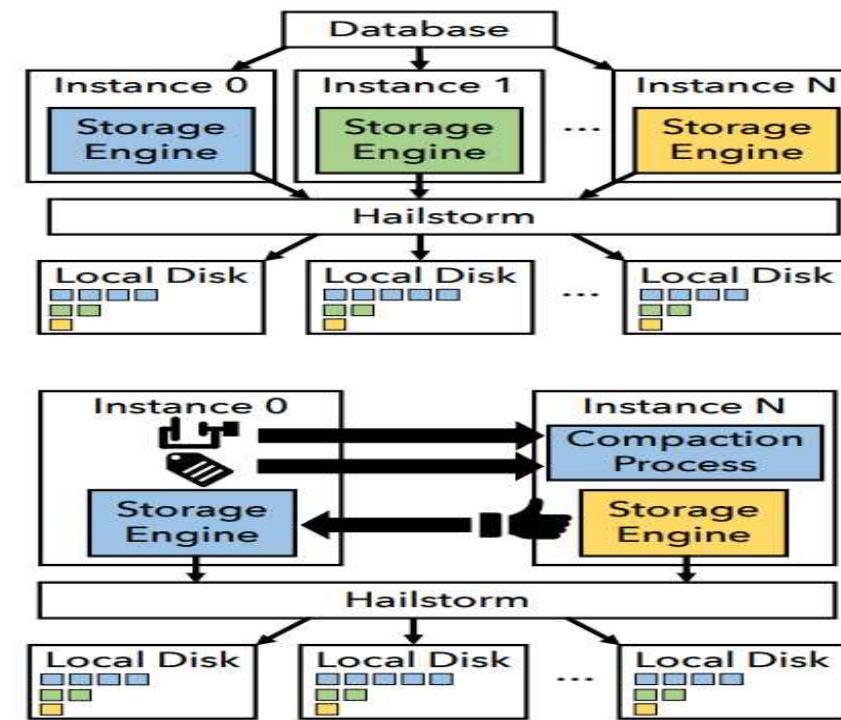
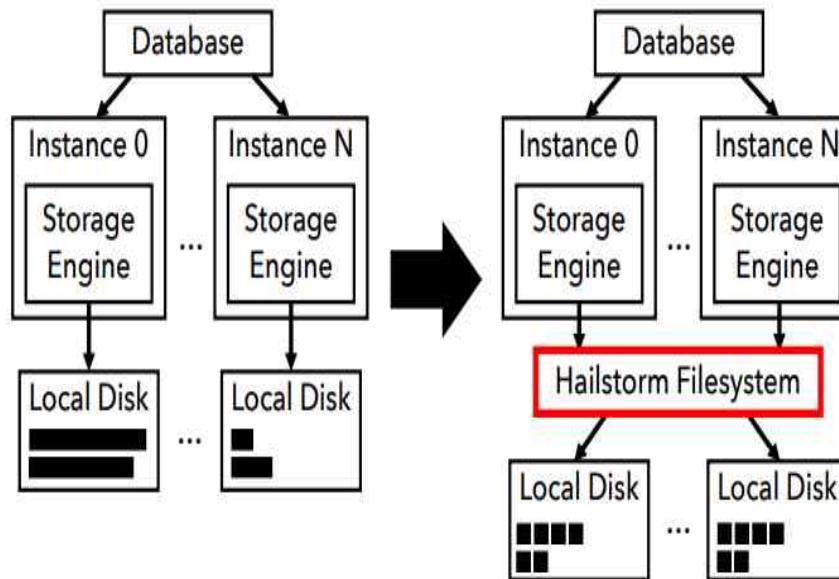


# Advance Topic 4: New approach (10/12)

## ■ Case study: Hailstorm, ASPLOS'20

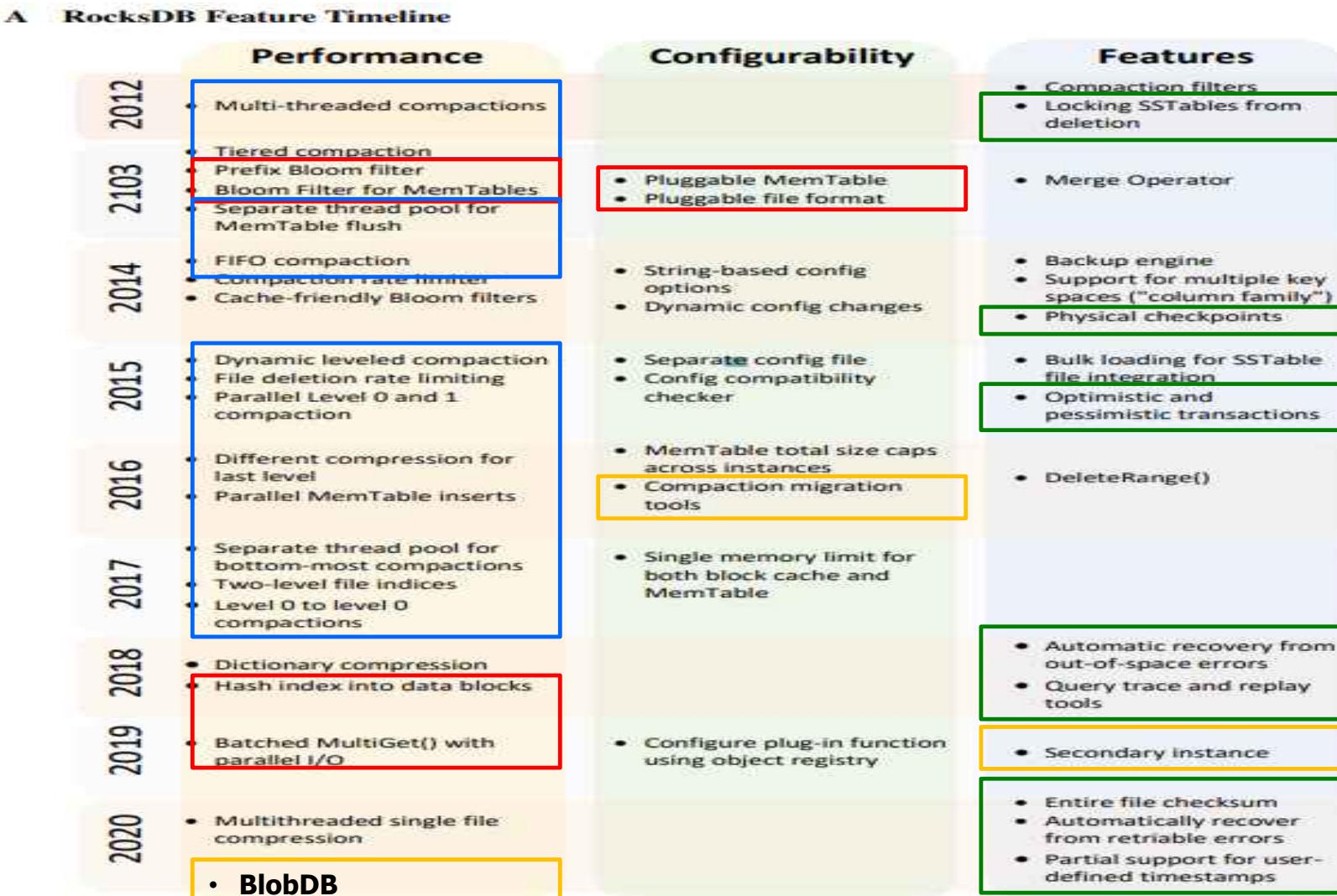
### ✓ How to?

- Causes of load imbalance
  - 1) Data skewness, 2) Compaction
- Data skewness
  - Employ a middleware (file system in this study) for load balancing
- Interference due to Compaction
  - Compaction offloading



# Advance Topic 4: New approach (11/12)

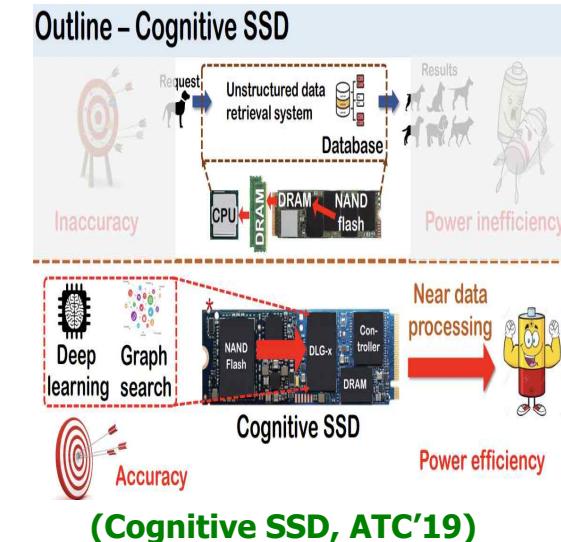
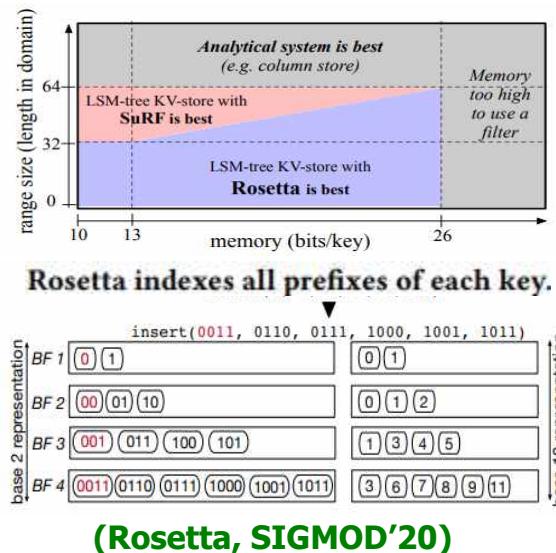
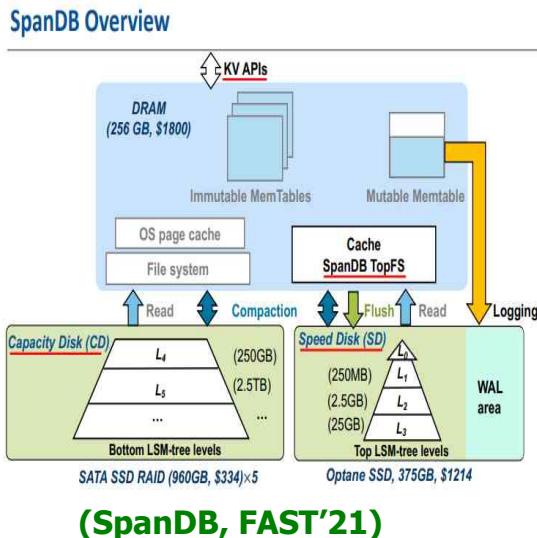
## ■ RocksDB Features (what we have learned at the second lecture)



# Advance Topic 4: New approach (12/12)

## ■ Take-away Lessons

- ✓ Outstanding researches are conducted
  - Key value separation for smaller LSM-tree
  - Utilize new storage (e.g. NVM, Optane SSD) and hybrid storage
  - Exploit computing capability in storage
  - Computing offload in a distributed key-value store
- ✓ Still a hot research topic



(Special thank you to DKU Embedded Members and IITP's SW StarLab (No.2021-0-01475)

# Content

---

- What is Key-Value Store?
- RocksDB basic
- Advance Topic 1: Compaction
- Advance Topic 2: Lookup
- Advance Topic 3: WAL
- Advance Topic 4: New approach
- RocksDB Practice
  - ✓ Install RocksDB in your computer
  - ✓ Run db\_bench
  - ✓ How to analyze outputs?
- Conclusion

# RocksDB Practice (1/12)

---

## ■ Lecture Objective

- ✓ Learn how to set up RocksDB in your computer
- ✓ Understand how to run a benchmark on RocksDB
- ✓ Investigate representative options for RocksDB
- ✓ Examine how to analyze evaluation results
- ✓ Explore the source code of RocksDB

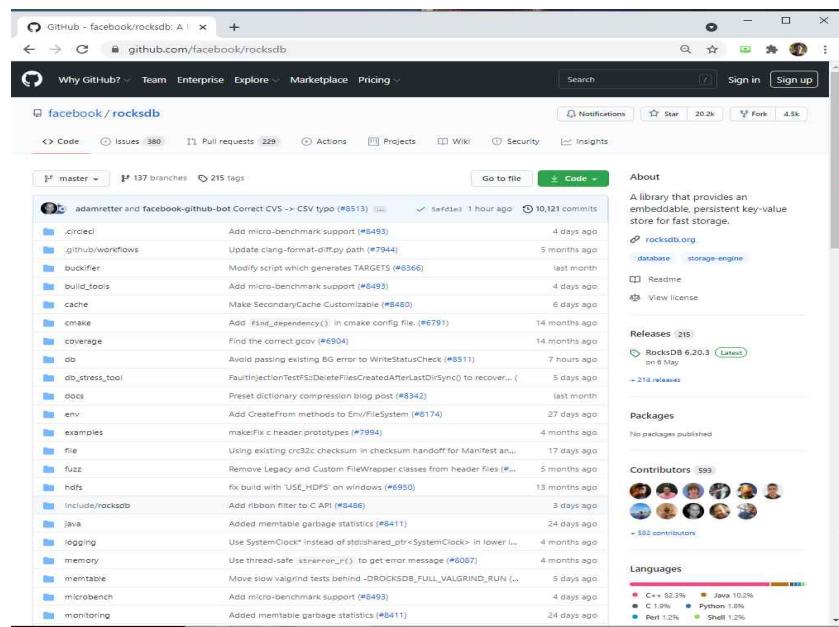


(Source: [www.dreamstime.com](http://www.dreamstime.com))

# RocksDB Practice (2/12)

## ■ Set up RocksDB in your computer

- ✓ Preparation: Linux machine, ubuntu on Virtual box, WSL, Cloud
- ✓ Download RocksDB from github
  - What is Github?
    - . Git: distributed version control system
    - . Github: open source hosting service based on git
  - How to? \$git clone → \$ls rocksdb



```
root@shin96:/home/shin96/RocksDB_Festival# git clone https://github.com/facebook/rocksdb.git
Cloning into 'rocksdb'...
remote: Enumerating objects: 99699, done.
remote: Counting objects: 100% (589/589), done.
remote: Compressing objects: 100% (296/296), done.
remote: Total 99699 (delta 308), reused 465 (delta 272), pack-reused 99110
Receiving objects: 100% (99699/99699), 168.23 MiB | 22.59 MiB/s, done.
Resolving deltas: 100% (75591/75591), done.
Checking connectivity... done.

root@shin96:/home/shin96/RocksDB_Festival# ls
rocksdb
root@shin96:/home/shin96/RocksDB_Festival# cd rocksdb/
root@shin96:/home/shin96/RocksDB_Festival/rocksdb# ls
appveyor.yml          db          HISTORY.md
AUTHORS               dist        include
buckifier              stress_tool
build_tools            DEFAULT_OPTIONS_HISTORY.md
cache                  docs        INSTALL.md
cmake                 DUMP_FORMAT.md
CMakeLists.txt         env         LICENSE.Apache
CODE_OF_CONDUCT.md    examples    LICENSE.leveldb
CONTRIBUTING.md       file        logging
COPYING                fuzz        Makefile
coverage               hdfs        memory
memtable               monitoring
memtable               options
memtable               plugin
memtable               tools
memtable               trace_replay
memtable               USERS.md
memtable               util
memtable               utilities
memtable               Vagrantfile
memtable               WINDOWS_PORT.md
memtable               TARGETS

The files listed are:
- appveyor.yml
- AUTHORS
- buckifier
- build_tools
- cmake
- CMakeLists.txt
- CODE_OF_CONDUCT.md
- CONTRIBUTING.md
- COPYING
- coverage
- db
- dist
- DEFAULT_OPTIONS_HISTORY.md
- docs
- env
- examples
- fuzz
- HISTORY.md
- include
- INSTALL.md
- issue_template.md
- java
- LICENSE.Apache
- LICENSE.leveldb
- logging
- Makefile
- memory
- memtable
- monitoring
- options
- plugin
- tools
- trace_replay
- USERS.md
- util
- utilities
- Vagrantfile
- WINDOWS_PORT.md
- TARGETS
```

# RocksDB Practice (3/12)

## ■ Build and Run benchmark

### ✓ Compile

- Make: A build tool that controls the generation of executables via Makefile
- How to? `$make db_bench`
  - It takes quite a long time (depend on your system)

### ✓ Run benchmark

- `db_bench`: a program for testing RocksDB performance
- How to? `./db_bench`

```
root@shin96:/home/shin96/RocksDB_Festival/rocksdb# make db_bench
$DEBUG_LEVEL is 1
Makefile:176: Warning: Compiling in debug mode. Don't use the resulting binary in production
$DEBUG_LEVEL is 1
Makefile:176: Warning: Compiling in debug mode. Don't use the resulting binary in production
CC      tools/db_bench.o
CC      tools/db_bench_tool.o
CC      tools/simulated_hybrid_file_system.o
CC      test_util/testutil.o
CC      cache/cache.o
CC      cache/cache_entry_roles.o
CC      cache/clock_cache.o
CC      cache/lru_cache.o
CC      cache/sharded_cache.o
CC      db/arena_wrapped_db_iter.o
CC      db/blob/blob_fetcher.o
```

```
root@shin96:/home/shin96/RocksDB_Festival/rocksdb# ls
appveyor.yml          file
AUTHORS                fuzz
buildifier              hdfs
build_tools             HISTORY.md
cache                  DEFAULT_OPTIONS_HISTORY.md
cmake                  include
CMakeLists.txt          INSTALL.md
CODE_OF_CONDUCT.md     issue_template.md
CONTRIBUTING.md        java
COPYING                LANGUAGE_BINDINGS.md
examples               librocksdb.debug.a
COPYRIGHT              monitoring
LICENSE.Apache          options
LICENSE.level0           port
LICENSE.level1          README.md
logging                 ROCKSDB_LITE.md
make_config.mk          src.mk
Makefile                utilities
memory                 table
mmt                    TARGETS
monitoring             test_util
options                thirdparty
port                   Vagrantfile
README.md               WINDOWS_PORT.md
thirdparty.inc          tools
trace_replay            USERS.md
util                   utilities
Vagrantfile             windows_port.md

root@shin96:/home/shin96/RocksDB_Festival/rocksdb# ./db_bench --db=/mnt --benchmarks=fillsq --num=1000000
Initializing RocksDB Options from the specified file
Initializing RocksDB Options from command-line flags
RocksDB: version 6.22
Date:   Fri Jul 2 22:02:50 2021
CPU:    4 * Intel(R) Core(TM) i5-6600 CPU @ 3.30GHz
CPU Cache: 6144 KB
Keys:   16 bytes each (+ 0 bytes user-defined timestamp)
Values: 100 bytes each (50 bytes after compression)
Entries: 1000000
Prefix: 0 bytes
Keys per prefix: 0
RawSize: 110.6 MB (estimated)
Filesize: 62.9 MB (estimated)
Write rate: 0 bytes/second
Read rate: 0 ops/second
Compression: Snappy
Compression sampling rate: 0
Memtablerep: skip_list
Perf Level: 1
WARNING: Assertions are enabled; benchmarks unnecessarily slow
-----
Initializing RocksDB Options from the specified file
Initializing RocksDB Options from command-line flags
DB path: [/mnt]
fillsq : 2.606 micros/op 383782 ops/sec 42.5 MB/s
root@shin96:/home/shin96/RocksDB_Festival/rocksdb#
```

# RocksDB Practice (4/12)

## ■ Configurability

- ✓ There are a lot of options (parameters for configuration)
  - Balance between **customizability** and **self-adaptability**
- ✓ 1) Benchmark related options
  - Benchmark: a series of well-defined test to measure performance
  - Workload type: sequential, random, # of threads, ...
  - Workload size: num (# of KV pairs), key\_size, value\_size, ...

The image shows two browser windows side-by-side. The left window displays the help output for the db\_bench command, listing various flags and their descriptions. The right window shows the full list of configuration parameters for RocksDB, which are categorized into several groups such as IO, Compression, and Logging and Monitoring.

**Left Tab (db\_bench help output):**

```
$ ./db_bench -help
db_bench:
USAGE:
./db_bench [OPTIONS]...
```

Flags from tools/db\_bench\_tool.cc:

- advise\_random\_on\_open (Advise random access on table file open) type: bool default: true
- allow\_concurrent\_memtable\_write (Allow multi-writers to update mem tables in parallel.) type: bool default: true
- base\_background\_compactions (The base number of concurrent background compactions to occur in parallel.) type: int32 default: 1
- batch\_size (Batch size) type: int64 default: 1
- benchmark\_read\_rate\_limit (If non-zero, db\_bench will rate-limit the reads from RocksDB. This is the global rate in ops/second.) type: uint64 default: 0
- benchmark\_write\_rate\_limit (If non-zero, db\_bench will rate-limit the writes going into RocksDB. This is the global rate in bytes/second.) type: uint64 default: 0
- benchmarks (Comma-separated list of operations to run in the specified order. Available benchmarks:
  - fillseq -- write N values in sequential key order in async mode
  - fillseqdeterministic -- write N values in the specified key order and keep the shape of the LSM tree
  - fillrandom -- write N values in random key order in async mode
  - filluniquerandomdeterministic -- write N values in a random key order and keep the shape of the LSM tree
  - overwrite -- overwrite N values in random key order in async mode
  - fillsync -- write N/100 values in random key order in sync mode
  - fill10K -- write N/1000 1000 values in random order in async mode
  - deleteseq -- delete N keys in sequential order
  - deletarandom -- delete N keys in random order
  - readseq -- read N times sequentially)

**Right Tab (Full list of configuration parameters):**

- key\_id\_range (Range or possible value of key id (used in timeseries only.)) type: int32 default: 100000
- key\_size (Size of each key) type: int32 default: 16
- keys\_per\_prefix (Control average number of keys generated per prefix, 0 means no special handling of the prefix, i.e. use the prefix comes with the generated random number.) type: int64 default: 0
- level0\_file\_num\_compaction\_trigger (Number of files in level-0 when compactions start) type: int32 default: 4
- level0\_slowdown\_writes\_trigger (Number of files in level-0 that will slow down writes.) type: int32 default: 20
- level0\_stop\_writes\_trigger (Number of files in level-0 that will trigger put stop.) type: int32 default: 36
- level\_compaction\_dynamic\_level\_bytes (Whether level size base is dynamic) type: bool default: false
- max\_background\_compactions (The maximum number of concurrent background compactions that can occur in parallel.) type: int32 default: 1
- max\_background\_flushes (The maximum number of concurrent background flushes that can occur in parallel.) type: int32 default: 1
- max\_bytes\_for\_level\_base (Max bytes for level-1) type: uint64 default: 268435456
- max\_bytes\_for\_level\_multiplier (A multiplier to compute max bytes for level-N ( $N > 1$ )) type: double default: 10
- max\_bytes\_for\_level\_multiplier\_additional (A vector that specifies additional fanout per level) type: string default: ""
- max\_compaction\_bytes (Max bytes allowed in one compaction) type: uint64 default: 0
- max\_num\_range\_tombstones (Maximum number of range tombstones to insert.) type: int64 default: 0
- max\_successive\_merges (Maximum number of successive merge operations on a key in the memtable) type: int32 default: 0
- max\_total\_wal\_size (Set total max WAL size) type: uint64 default: 0
- max\_write\_buffer\_number (The number of in-memory memtables. Each memtable is of size write\_buffer\_size.) type: int32 default: 2
- max\_write\_buffer\_number\_to\_maintain (The total maximum number of write buffers to maintain in memory including copies of buffers that have already been flushed. Unlike max\_write\_buffer\_number, this parameter does not affect flushing. This controls the minimum amount of write history)

# RocksDB Practice (5/12)

## ■ Configurability (cont')

### ✓ 2) Memtable related options

- write\_buffer\_size: Size of a single Memtable (default: 64MB).
- Max\_write\_buffer\_number: # of Memtables (default: 2)
  - 1 for mutable and others for immutable that wait for being flushed

### ✓ 3) SSTable related options

- Level0\_file\_num\_compaction\_trigger (default 4): # of files to trigger
- All non-0 levels: target\_size, max\_bytes\_for\_level\_multiplier (default 10)
  - if (level\_compaction\_dynamic\_level\_bytes == false)
    - target\_size( $L_1$ ) = **max\_bytes\_for\_level\_base** (default=256MB, 64\*4 from L0)
    - target\_size( $L_{n+1}$ ) = target\_size( $L_n$ ) \* max\_bytes\_for\_level\_multiplier
  - else
    - determine the last level size as the actual size
    - then, target\_size( $L_{n-1}$ ) = target\_size( $L_n$ ) / max\_bytes\_for\_level\_multiplier



# RocksDB Practice (6/12)

---

## ■ Configurability (cont')

- ✓ 4) Compaction related options
  - Compaction\_style: kCompactionStyleLevel (default), Universal, FIFO, ...
  - disable\_auto\_compactions: useful in the case of bursty case
  - max\_background\_jobs, max\_background\_compactions, level\_compaction\_dynamic\_level\_bytes, ...
- ✓ 5) Block cache related options
  - Block cache size, policy (LRU or Clock), allow\_OS\_buffer,, ...
- ✓ 6) Storage related options
  - Block\_size (default 4KB): I/O unit
    - . When reading a KV pair from a file, an entire block is loaded into memory
  - I/O rate limiter (e.g. rate\_bytes\_per\_sec)
    - . Throttle maximum write speed to give a chance to read
- ✓ 7) Compression related options
- ✓ 8) ...

```
cf_options.level_compaction_dynamic_level_bytes = true;
options.max_background_compactions = 4;
options.max_background_flushes = 2;
options.bytes_per_sync = 1048576;
options.compaction_pri = kMinOverlappingRatio;
table_options.block_size = 16 * 1024;
table_options.cache_index_and_filter_blocks = true;
table_options.pin_l0_filter_and_index_blocks_in_cache = true;
table_options.format_version = <the latest version>;
```

# RocksDB Practice (7/12)

## ■ Evaluation examples: from RocksDB wiki

- ✓ Set up
  - AWS instance: 8 CPU, 32 GB Memory, 1 x 300 NVMe SSD
  - Linux kernel: 4.14, FS: XFS with discard enabled, RocksDB: 6.10
- ✓ Tests (being added)
  - Test1 (bulkload): Bulk load of Keys in Random Order (write only)
  - Test2 (overwrite): Random write, started from Test 1
  - Test3 (readwhilewriting): randomly read keys and ongoing updates to existing keys, Based on Test 2
  - Test4 (randomread): random read performance
- ✓ Result 1
  - 100M Keys, 2KB value, Default block size (4KB), bulkload time = 2018s

Test	ops/sec	mb/sec	Size-GB	L0_GB	Sum_GB	W-Amp	W-MB/s	usec/op	p50	p75	p99
bulkload	272448	537.3	0.1	85.3	85.3	1.0	242.6	3.7	0.7	1.1	2.1
overwrite	22940	45.2	0.1	229.3	879.4	3.8	169.0	1394.9	212.9	350.7	760.1
readwhilewriting	87093	154.2	0.1	5.4	162.6	30.1	31.0	367.4	369.2	491.9	2200.1
readrandom	95666	169.9	0.1	0.0	0.0	0	0	334.5	411.1	498.7	740.1

# RocksDB Practice (8/12)

## ■ Evaluation examples: from Rocksdb wiki (cont')

### ✓ Result 2

- NUM\_KEYS=9x10<sup>8</sup>, NUM\_THREADS=32, CACHE\_SIZE=6442450944
- Analyze while change an option (block size in this case)
  - With different block size: 8KB, 4KB, 16KB, 8KB with DIO, 8KB with RL
  - RL (Rate Limiting): a timed rate-limited operation was place before the reported operation → give a chance to do background jobs → creating more predictable results
- Results are given per each test

### Test Case 1 : benchmark.sh bulkload

8K: Complete bulkload in 4560 seconds 4K: Complete bulkload in 5215 seconds 16K: Complete bulkload in 3996 seconds DIO:  
Complete bulkload in 4547 seconds 8K RL: Complete bulkload in 4388 seconds

Block	ops/sec	mb/sec	Size-GB	L0_GB	Sum_GB	W-Amp	W-MB/s	usec/op	p50	p75	p99	p99.9
8K	924468	370.3	0.2	157.1	157.1	1.0	167.5	1.1	0.5	0.8	2	4
4K	853217	341.8	0.2	165.3	165.3	1.0	165.9	1.2	0.5	0.8	2	4
16K	1027567	411.6	0.1	149.0	149.0	1.0	181.6	1.0	0.5	0.8	2	3
DIO	921342	369.0	0.2	156.6	156.6	1.0	167.0	1.1	0.5	0.8	2	4
8K RL	989786	396.5	0.2	159.4	159.4	1.0	179.5	1.0	0.5	0.8	2	4

### Test Case 2 : benchmark.sh overwrite

Block	ops/sec	mb/sec	Size-GB	L0_GB	Sum_GB	W-Amp	W-MB/s	usec/op
8K	85756	34.3	0.1	161.4	739.9	4.5	142.2	373.1
4K	79856	32.0	0.2	166.0	716.9	4.3	136.3	400.7
16K	93678	37.5	0.1	174.4	825.0	4.7	156.8	341.6
DIO	85655	34.3	0.1	163.9	734.9	4.4	140.7	373.6
8K RL	85542	34.3	0.1	161.2	757.8	4.7	143.6	748.1

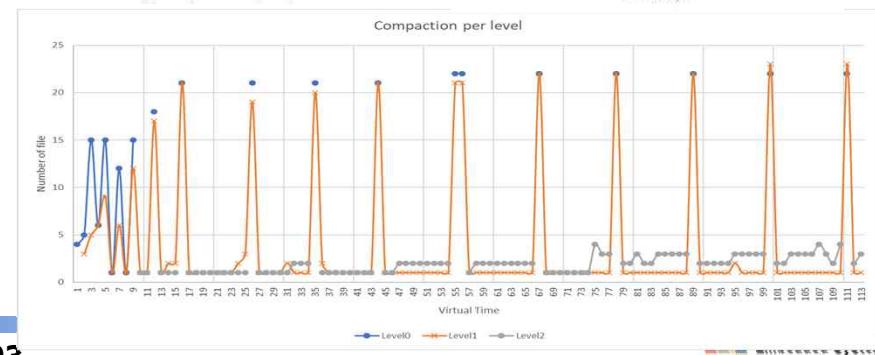
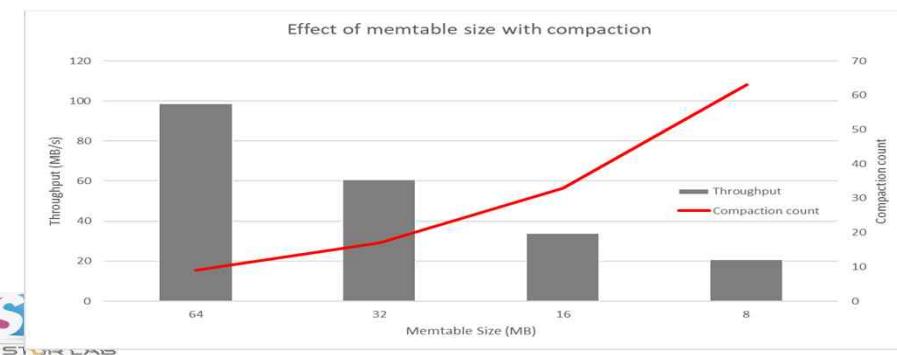
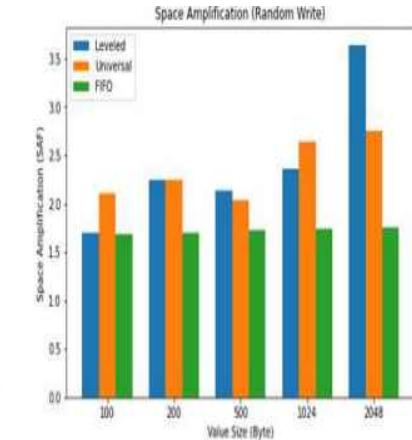
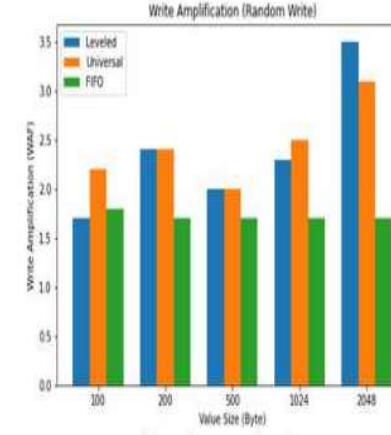
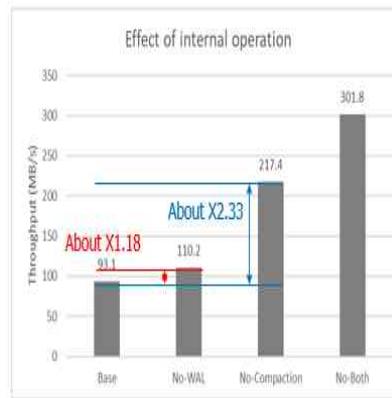
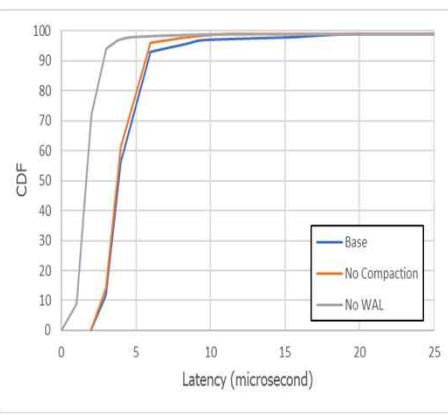
(Source: <https://github.com/facebook/rocksdb/wiki/Performance-Benchmarks>)

# RocksDB Practice (9/12)

## ■ Evaluation examples: from our experiments

### ✓ Result 3

- Effects of internal operations
- Effect of Memtable size: smaller size incurs higher compaction
- Effects of Compaction Policies (styles)
- # of SSTables involved in each compaction



# RocksDB Practice (10/12)

## ■ How to modify RocksDB?

- ✓ Browse source code
- ✓ Modify a file (e.g. \$vi tools/db\_bench\_tool.cc for Run())
- ✓ Build and run the db\_bench

```
root@linux-server-93:/home/choigunhee/hojin/RocksDB_Explorer# ls
appveyor.yml      db_stress_tool      issue_template.md    parsing_csv      third-party
AUTHORS          DEFAULT_OPTIONS_HISTORY.md   java           plugin        thirdparty.inc
buckifier        defs.bz1             LANGUAGE-BINDINGS.md  PLUGINS.md    tools
build_tools       docs               librocksdb_debug.a  port          trace_replay
cache            DUMP_FORMAT.md      LICENSE.Apache     python_parser  USERS.md
cmake            env                LICENSE.leveldb    README.md    util
CMakeLists.txt    examples          logging          result_txt    utilities
CODE_OF_CONDUCT.md file              make_config.mk  RocksDB_explorer_sh Vagrantfile
CONTRIBUTING.md  fuzz              Makefile         ROCKSDB_LITE.md WINDOWS_PORT.md
COPYING          hdfs              memory          src.mk
coverage         HISTORY.md       memtable        table
db               include          monitoring      TARGETS
db_bench        INSTALL.md      options         test_util
root@linux-server-93:/home/choigunhee/hojin/RocksDB_Explorer#
```

```
void Run() {
    if (!SanityCheck()) {
        ErrorExit();
    }
    Open(&open_options_);
    PrintHeader();
    fprintf(stdout, "Source Code Modifying!!!\n");
    std::stringstream benchmark_stream(FLAGS_benchmarks);
```



```
root@shin96:/home/shin96/RocksDB_Festival/rocksdb# make db_bench
$DEBUG_LEVEL is 1
Makefile:176: Warning: Compiling in debug mode. Don't use the resulting binary in production
$DEBUG LEVEL is 1
Makefile:176: Warning: Compiling in debug mode. Don't use the resulting binary in production
CC      tools/db_bench_tool.o
CCLD    db_bench
root@shin96:/home/shin96/RocksDB_Festival/rocksdb#
```

```
root@shin96:/home/shin96/RocksDB_Festival/rocksdb# ./db_bench --db=
Initializing RocksDB Options from the specified file
Initializing RocksDB Options from command-line flags
RocksDB: version 6.22
Date:   Fri Jul  2 22:06:32 2021
CPU:   4 * Intel(R) Core(TM) i5-6600 CPU @ 3.30GHz
CPU Cache: 6144 KB
Keys: 16 bytes each (+ 0 bytes user-defined timestamp)
Values: 100 bytes each (50 bytes after compression)
Entries: 1000000
Prefix: 0 bytes
Keys per prefix: 0
RawSize: 110.6 MB (estimated)
FileSize: 62.9 MB (estimated)
Write rate: 0 bytes/second
Read rate: 0 ops/second
Compression: Snappy
Compression sampling rate: 0
Memtablerep: skip_list
Perf Level: 1
WARNING: Assertions are enabled; benchmarks unnecessarily slow
-----
Source Code Modifying!!!
Initializing RocksDB Options from the specified file
Initializing RocksDB Options from command-line flags
DB path: [/mnt]
fillseq : 2.649 micros/op 377562 ops/sec, 41.8 MB/s
root@shin96:/home/shin96/RocksDB_Festival/rocksdb#
```

# RocksDB Practice (11/12)

## ■ RocksDB explorer

- ✓ By Hojin Shin (at DKU Lab.)

- Logging for source code analysis, Scripts for various experiments, ...

**RocksDB Explorer**

**RocksDB Festival : Quantitative Approaches by DKU StarLab**

This is a place for stdyuing RocksDB (Facebook) by Dankook University

- Writer : Hojin Shin
- Laboratory : Embedded System Lab. in Dankook University
- 2021 DKU RocksDB Festival Lecture Information [link](#)

**How to use This Github**

**1. Setting Environment for experiments**

Follow the instructions below

```
git clone https://github.com/DKU-StarLab/RocksDB_Explorer.git
cd RocksDB_Explorer
make db_bench (It will take some time)
```

You need to edit your DEV\_PATH before starting experiments.

**2. Policy-related experiments**

2.1 What if we turn off the WAL (Write-Ahead-Log)

- Throughput

```
root@linux-server-93:/home/choigunhee/hojin/RocksDB_Explorer/RocksDB_explorer sh# ls
compaction_latency.sh    configuration_threads.sh policy_wal.sh workload_key_value_size.sh
configuration_memtable.sh policy_compaction.sh wal_latency.sh workload_value_size_distribution.sh
(scripts)
```

```
Keys:      16 bytes each (+ 0 bytes user-defined timestamp)
Values:    1024 bytes each (1024 bytes after compression)
Entries:   5000000
Prefix:    0 bytes
Keys per prefix: 0
RawSize:  4959.1 MB (estimated)
FileSize: 4959.1 MB (estimated)
Write rate: 0 bytes/second
Read rate: 0 ops/second
Compression: Snappy
Compression sampling rate: 0
Memtablerep: skip_list
Perf Level: 1
WARNING: Assertions are enabled; benchmarks unnecessarily slow
```

```
Initializing RocksDB Options from the specified file
Initializing RocksDB Options from command-line flags
db_bench Write Flow - MaybeScheduleFlushOrCompaction() in db_impl_compaction_flush.cc
db_bench Write Flow - MaybeScheduleFlushOrCompaction() in db_impl_compaction_flush.cc
DB path: [../../mnt/rocksdb_test/]
db_bench Write Flow - DoWrite() in db_bench_tool.cc
db_bench Write Flow - Write() in db_impl_write.cc
db_bench Write Flow - WriteImpl() in db_impl_write.cc
db_bench Write Flow - PipelinedWriteImpl() in db_impl_write.cc
db_bench Write Flow - PreprocessWrite() in db_impl_write.cc
db_bench Write Flow - Write() in db_impl_write.cc
db_bench Write Flow - WriteImpl() in db_impl_write.cc
```

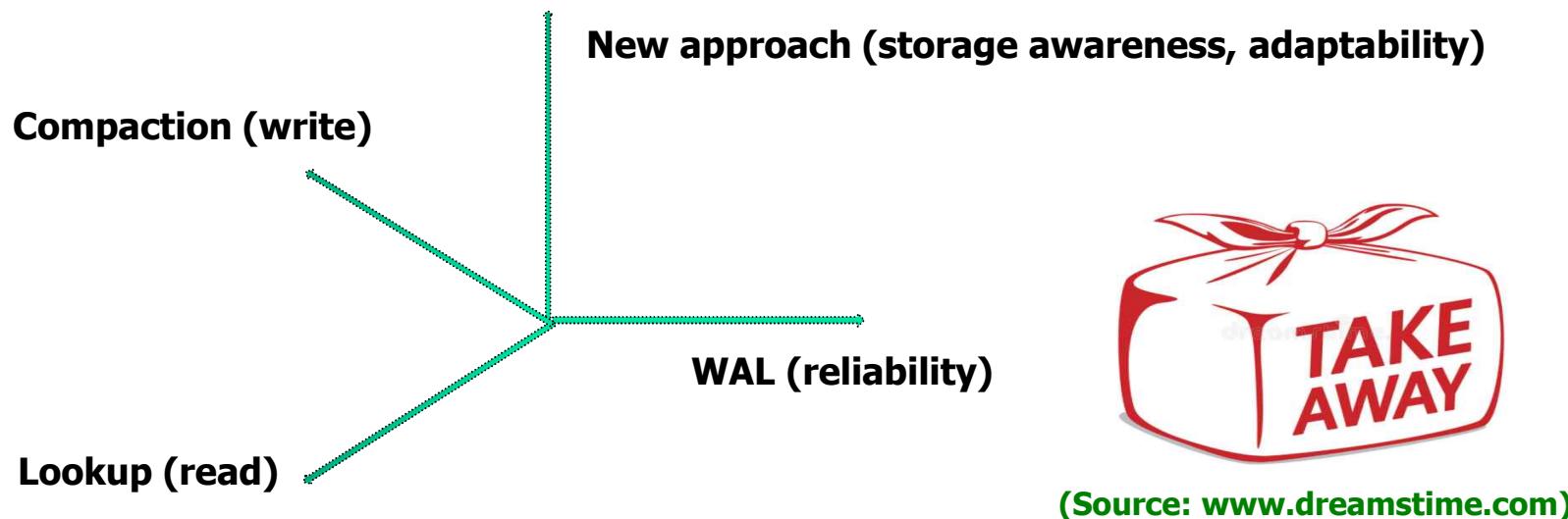
**(control flow analysis)**

**(Source: [https://github.com/DKU-StarLab/RocksDB\\_Explorer](https://github.com/DKU-StarLab/RocksDB_Explorer))**

# RocksDB Practice (12/12)

## ■ Take-away Lessons

- ✓ Set up, run, analyze RocksDB is straightforward :-)
  - Open source (**waiting for your contribution**)
  - Provide a variety of options
  - Being actively developed by amazing experts
- ✓ Let's do our own research!!



(Special thank you to DKU Embedded Members and IITP's SW StarLab (No.2021-0-01475)

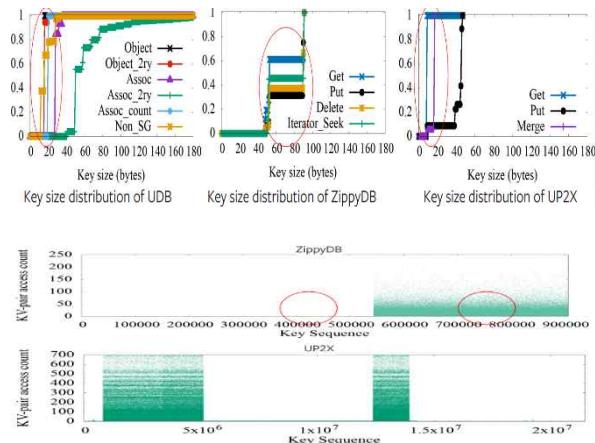
# Discussion

---

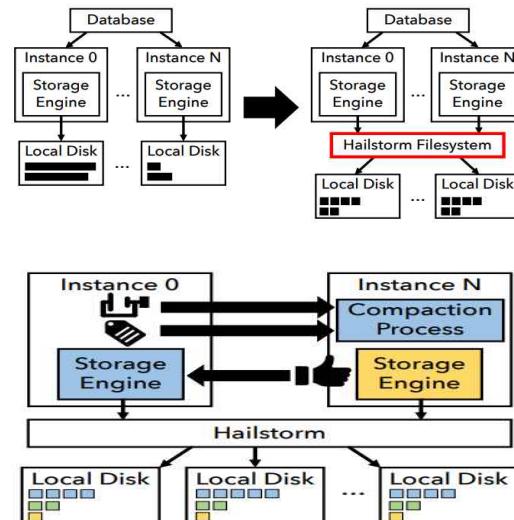


# Appendix

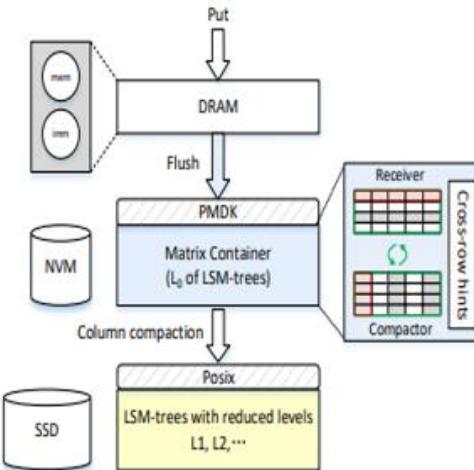
## ■ Interesting recent studies regarding Key-Value Store



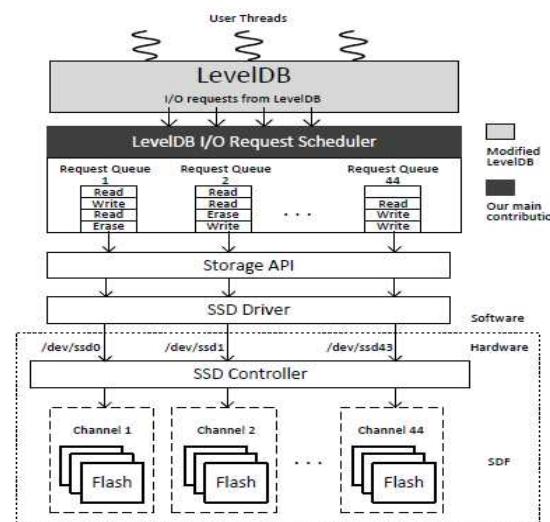
(RocksDB analysis, FAST'20)



 (Hailstorm, ASPLOS'20)

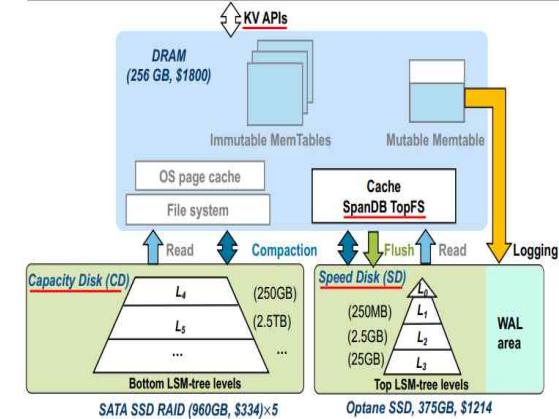


(MatrixKV, ATC'20)



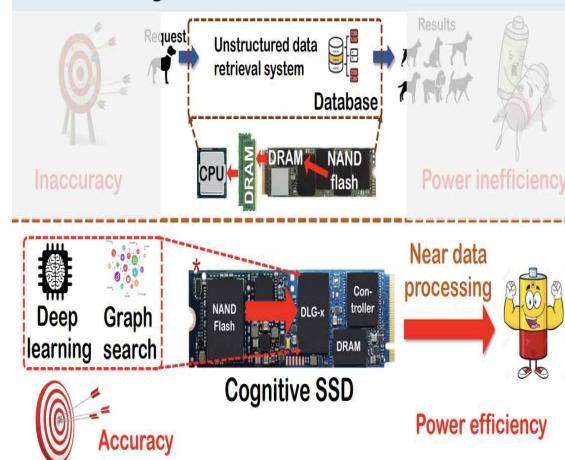
(SDF, ASPLOS'14)

SpanDB Overview



(SpanDB, FAST'21)

Outline - Cognitive SSD



(Cognitive SSD, ATC'19)

# Appendix

## ■ Interesting recent studies regarding Key-Value Store

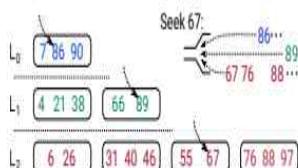
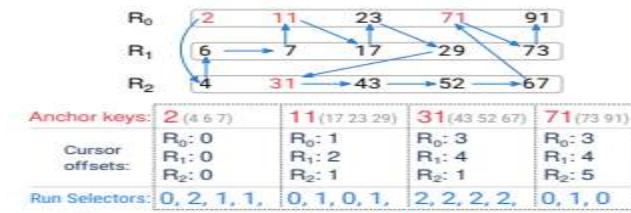
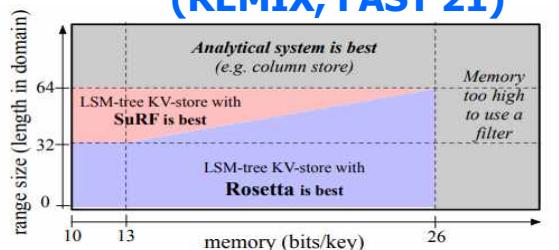


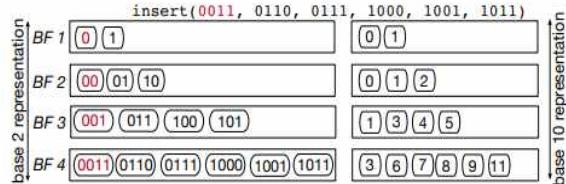
Figure 1: An LSM-tree using leveled compaction



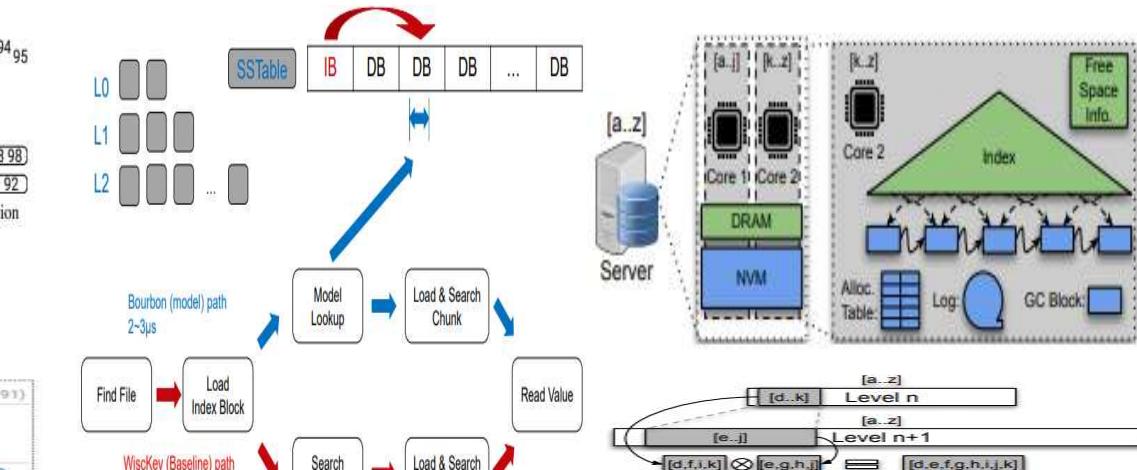
(REMIX, FAST'21)



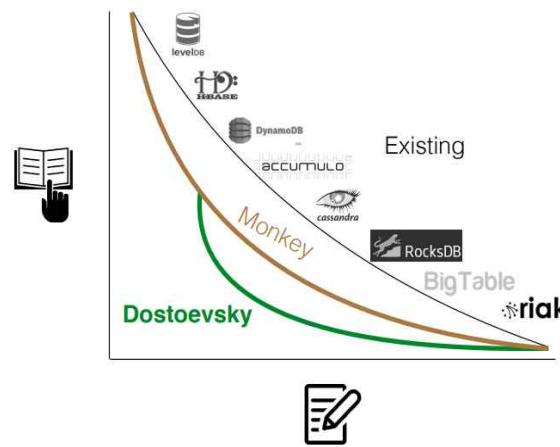
Rosetta indexes all prefixes of each key.



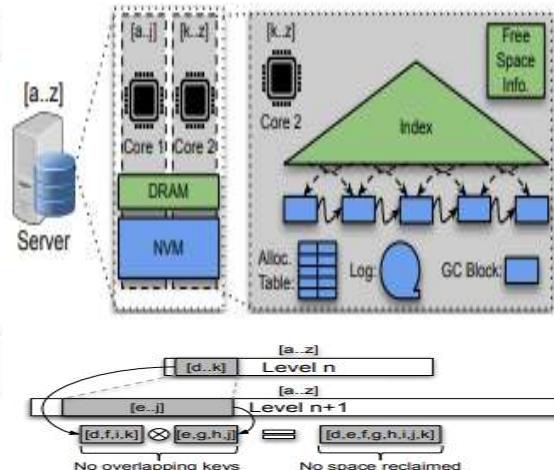
(Rosetta, SIGMOD'20)



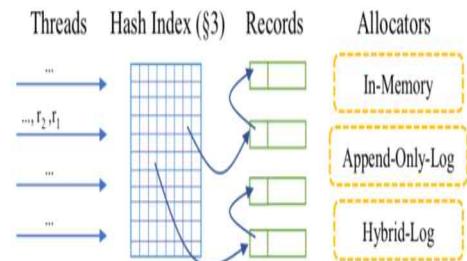
(Bourbon, OSDI'20)



(Dostoevsky, SIGMOD'18)



(RStore, VLDB'20)



	Record Allocator	Concurrent & Latch-free	Larger-Than-Memory	In-Place Updates
\$4	In-Memory	✓		✓
\$5	Append-Only-Log	✓	✓	
\$6	Hybrid-Log	✓	✓	✓

(Faster, SIGMOD'18)

# Appendix

## ■ YCSB (Yahoo! Cloud Serving Benchmark)

- ✓ A de facto standard benchmark for key-value stores
- ✓ 1) download, 2) build, 3) run workload
  - git clone <https://github.com/brianfrankcooper/YCSB.git>
  - ./bin/ycsb run rocksdb -s -P workloads/workloada -p rocksdb.dir=[PATH] -p operationcount=[num]

```
root@shin96:/home/shin96/RocksDB_Festival# git clone https://github.com/brianfrankcooper/YCSB.git
Cloning into 'YCSB'...
remote: Enumerating objects: 20648, done.
remote: Total 20648 (delta 0), reused 0 (delta 0), pack-reused 20648
Receiving objects: 100% (20648/20648), 31.68 MiB | 17.69 MiB/s, done.
Resolving deltas: 100% (8016/8016), done.
Checking connectivity... done.
root@shin96:/home/shin96/RocksDB_Festival# ls
rocksdb YCSB
root@shin96:/home/shin96/RocksDB_Festival# cd YCSB/
root@shin96:/home/shin96/RocksDB_Festival/YCSB# ls
accumulo_0.9 cassandra distribution googledatastore LICENSE.txt pom.xml s3 zookeeper
aerospike checkstyle.xml doc griddb maprdb postgresql scylla
arangodb cloudspanner dynamodb hbase1 maprjsondb rados seaweedfs
asyncbase CONTRIBUTING.md elasticsearch hbase2 memcached README.ind solr7
azurecosmos core elasticsearch5 ignite mongodb redis tablestore
azurereatablestorage couchbase foundationdb infinispan nosqlb rest tarantool
bin couchbase2 geode jdb orientdb redis voltdb
binding-parent crail googlebigtable kudu orientdb rocksdb workloads
root@shin96:/home/shin96/RocksDB_Festival/YCSB#
```

(git clone)

```
root@shin96:/home/shin96/RocksDB_Festival/YCSB# mvn -pl site.ycsb:rocksdb-binding -am clean package
[INFO] Scanning for projects...
[INFO] Reactor Build Order:
[INFO] 1. YCSB Root
[INFO] 2. Core YCSB
[INFO] 3. Per Datastore Binding descriptor
[INFO] 4. RocksDB Binding Parent
[INFO] 5. RocksDB Java Binding
[INFO] 6. Building YCSB Root 0.18.0-SNAPSHOT
[INFO] 7. Building Core YCSB 0.18.0-SNAPSHOT
[INFO] 8. Building RocksDB Binding 0.18.0-SNAPSHOT
[INFO] 
[INFO] -- maven-clean-plugin:2.5:clean (default-clean) @ root ---
[INFO] Deleting /home/shin96/RocksDB_Festival/YCSB/target
[INFO] -- maven-enforcer-plugin:3.0.0-M1:enforce (enforce-maven) @ root ---
[INFO] -- maven-checkstyle-plugin:2.16:check (validate) @ root ---
[INFO] -- maven-resources-plugin:2.6:resources (default-resources) @ core ---
[INFO] Using 'UTF-8' encoding to copy filtered resources.
[INFO] Copying 1 resource
[INFO] -- maven-compiler-plugin:3.7.0:compile (default-compile) @ core ---
[INFO] Changes detected - recompiling the module!
[INFO] Compiling 63 source files to /home/shin96/RocksDB_Festival/YCSB/core/target/classes
[INFO] -- maven-resources-plugin:2.6:testResources (default-testResources) @ core ---
[INFO] Using 'UTF-8' encoding to copy filtered resources.
[INFO] skip non-existing resourceDirectory /home/shin96/RocksDB_Festival/YCSB/core/src/test/resources
```

(build with maven)

```
Command line: -db site.ycsb.db.rocksdb.RocksDBClient -s -P workloads/workloada -p rocksdb.dir=/tmp/rocksdb-data-ycsb -t
YCSB Client 0.18.0-SNAPSHOT

Loading workload...
Starting test.
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
2021-07-02 00:56:53.0 sec: 0 operations, est completion in 0 second
Configuration: report latency for each error is false, specific error codes to track for latency are: []
2021-07-02 00:56:53.0 sec: 0 operations, 1000 operations/sec. [READ: Count=521, Max=184, Min=3, Avg=13.95, 90=28, 99=74, 99.9=140, 99.99=184] {CLEANUP: Count=1, Max=1123, Min=1123, Avg=1123, 99=1123, 99.9=1123, 99.99=1123} [UPDATE: Count=479, Max=4065, Min=19, Avg=55.09, 90=86, 99=239, 99.9=4065, 99.99=4065]
[OVERALL: RunTime(ms), 229
[OVERALL: Throughput(ops/sec), 4366.812227074236
[TOTAL_GCS_PS_Scavengel], Count, 0
[TOTAL_GC_TIME_PS_Scavengel], Time(ms), 0
[TOTAL_GC_TIME_%_PS_Scavengel], Time(%), 0.0
[TOTAL_GCS_PS_MarkSweep], Count, 0
[TOTAL_GC_TIME_PS_MarkSweep], Time(ms), 0
[TOTAL_GC_TIME_%_PS_MarkSweep], Time(%), 0.0
[TOTAL_GCS], Count, 0
[TOTAL_GC_TIME], Time(ms), 0
[TOTAL_GC_TIME_%], Time(%), 0.0
[READ], Operations, 521
[READ], AverageLatency(us), 13.953934740882918
[READ], MinLatency(us), 3
[READ], MaxLatency(us), 18
[READ], 95thPercentileLatency(us), 41
[READ], 99thPercentileLatency(us), 74
[READ], Return-OK, 521
[CLEANUP], Operations, 1
[CLEANUP], AverageLatency(us), 1123.0
[CLEANUP], MinLatency(us), 1123
[CLEANUP], MaxLatency(us), 1123
[CLEANUP], 95thPercentileLatency(us), 1123
[CLEANUP], 99thPercentileLatency(us), 1123
[UPDATE], Operation, 479
[UPDATE], AverageLatency(us), 55.02713987473904
[UPDATE], MinLatency(us), 19
[UPDATE], MaxLatency(us), 4065
[UPDATE], 95thPercentileLatency(us), 153
[UPDATE], 99thPercentileLatency(us), 233
[UPDATE], Return-OK, 479
root@shin96:/home/shin96/RocksDB_Festival/YCSB#
```

(run workloada)

```
root@shin96:/home/shin96/RocksDB_Festival/YCSB/workloads# ls
```

```
tsworkload tsworkload_template workload workloadb workloadc workloadd workloade workloadf workloadf_template
root@shin96:/home/shin96/RocksDB_Festival/YCSB/workloads#
```

(type of workloads)