

Jungle: Towards Dynamically Adjustable Key-Value Store by Combining LSM-Tree and Copy-On-Write B+Tree

Jung-Sang Ahn (junahn@ebay.com),

Mohiuddin Abdul Qader, Woon-Hak Kang, Hieu Nguyen, Guogen Zhang, and Sami Ben-Romdhane

Platform Architecture & Data Infrastructure, eBay Inc.

USENIX HotStorage 2019



Outline

- Background: LSM-tree overview
 - Write amplification
 - Tiering merge: trade-offs
- Jungle
 - (Copy-on-write B+tree overview)
 - Combining LSM-tree and copy-on-write B+tree: why and how?
- Brief evaluation
- What's next?

Log-Structured Merge-Tree (LSM-Tree)

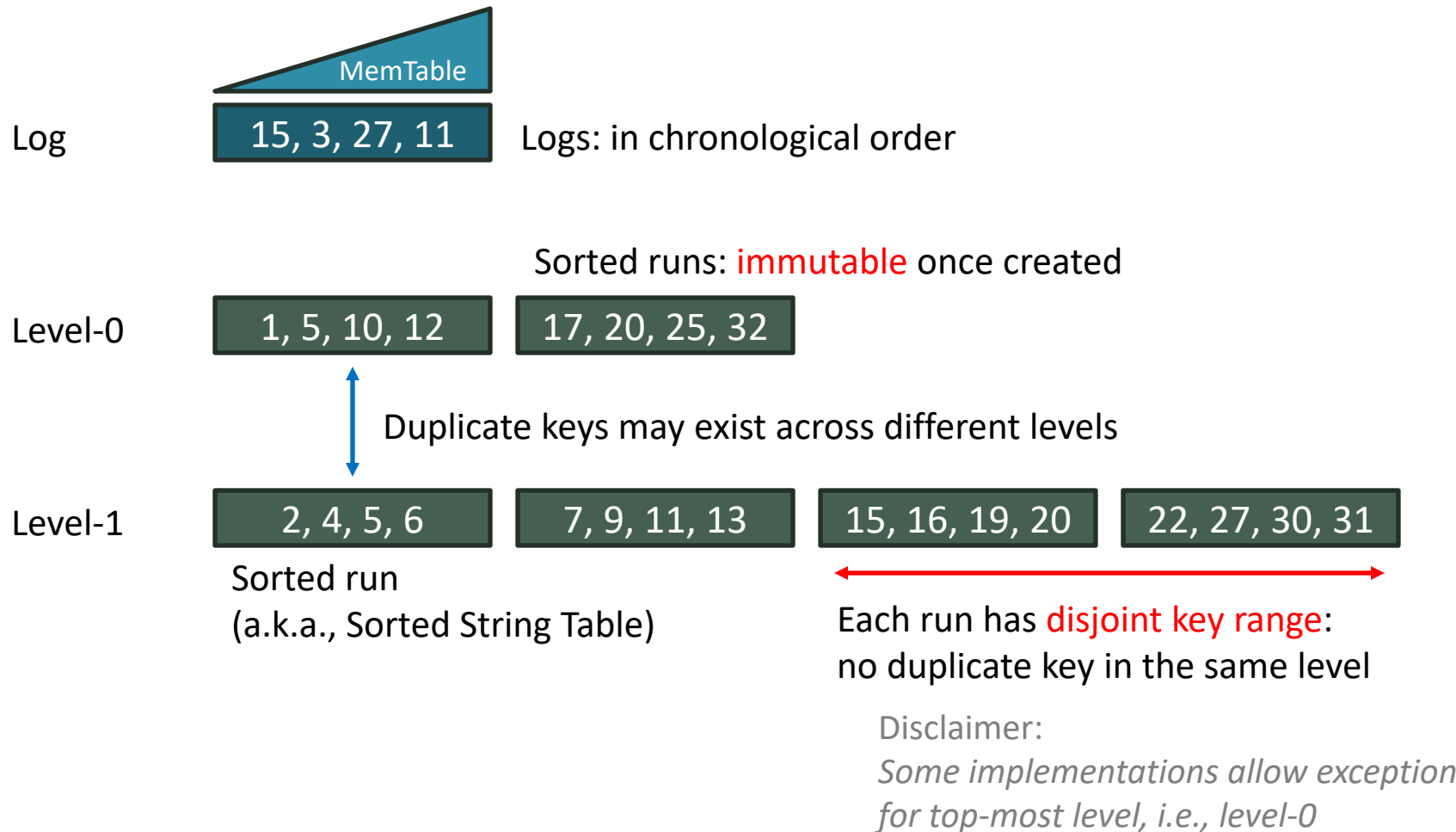
(P. O'Neil et al. 1996)

- Lots of recent key-value stores & databases are using it (or its variants)
- Compared to B+tree
 - Better write performance
 - Relatively degraded (but acceptable) read performance

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

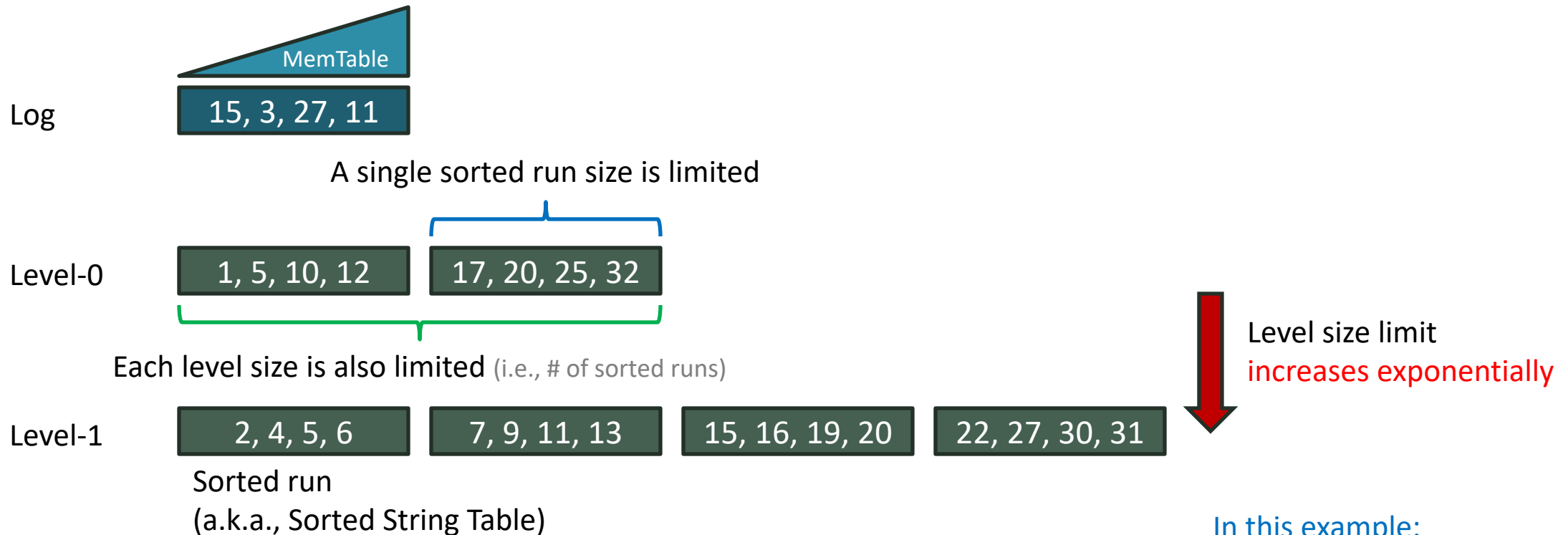
- Basic algorithm – write and merge



Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge



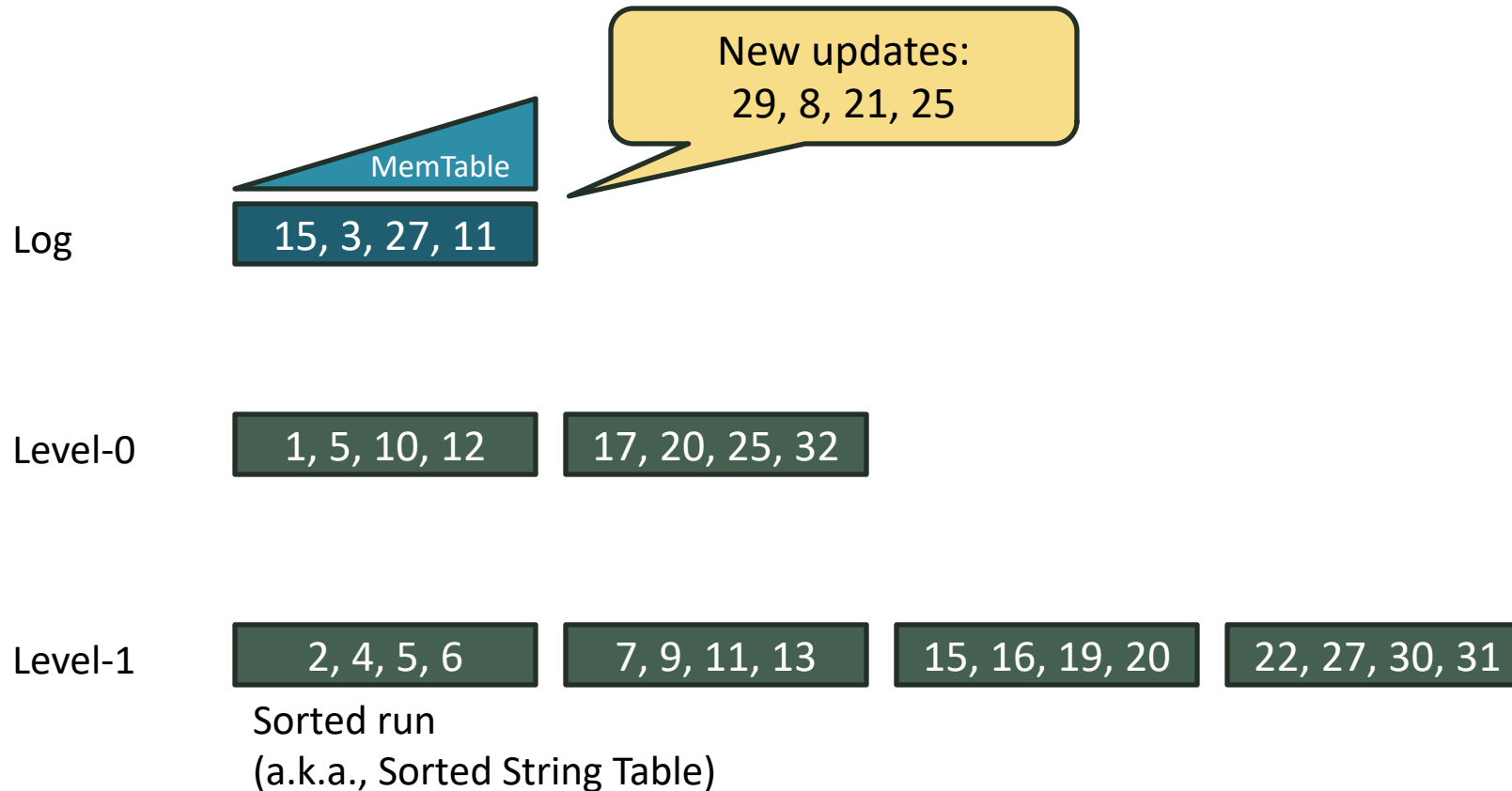
In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge



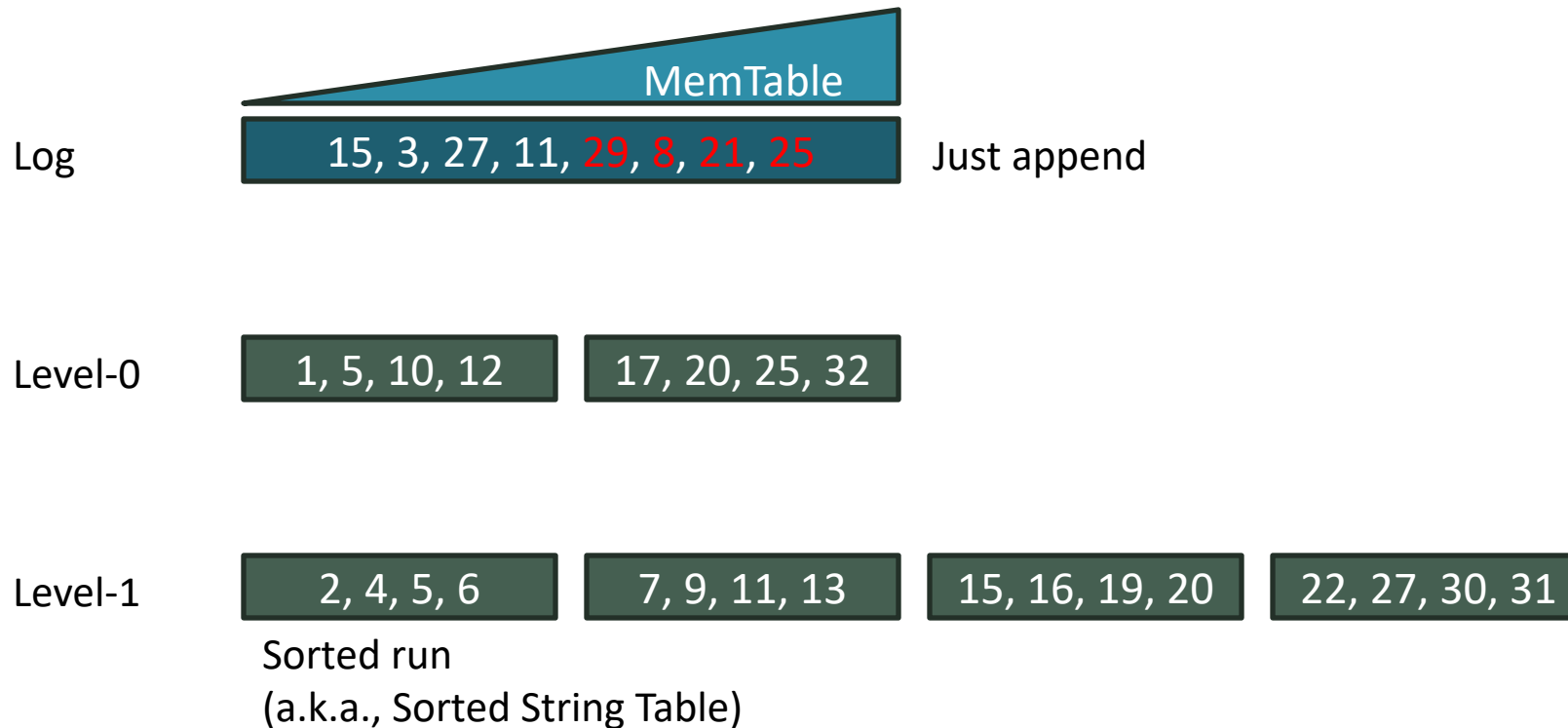
In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge



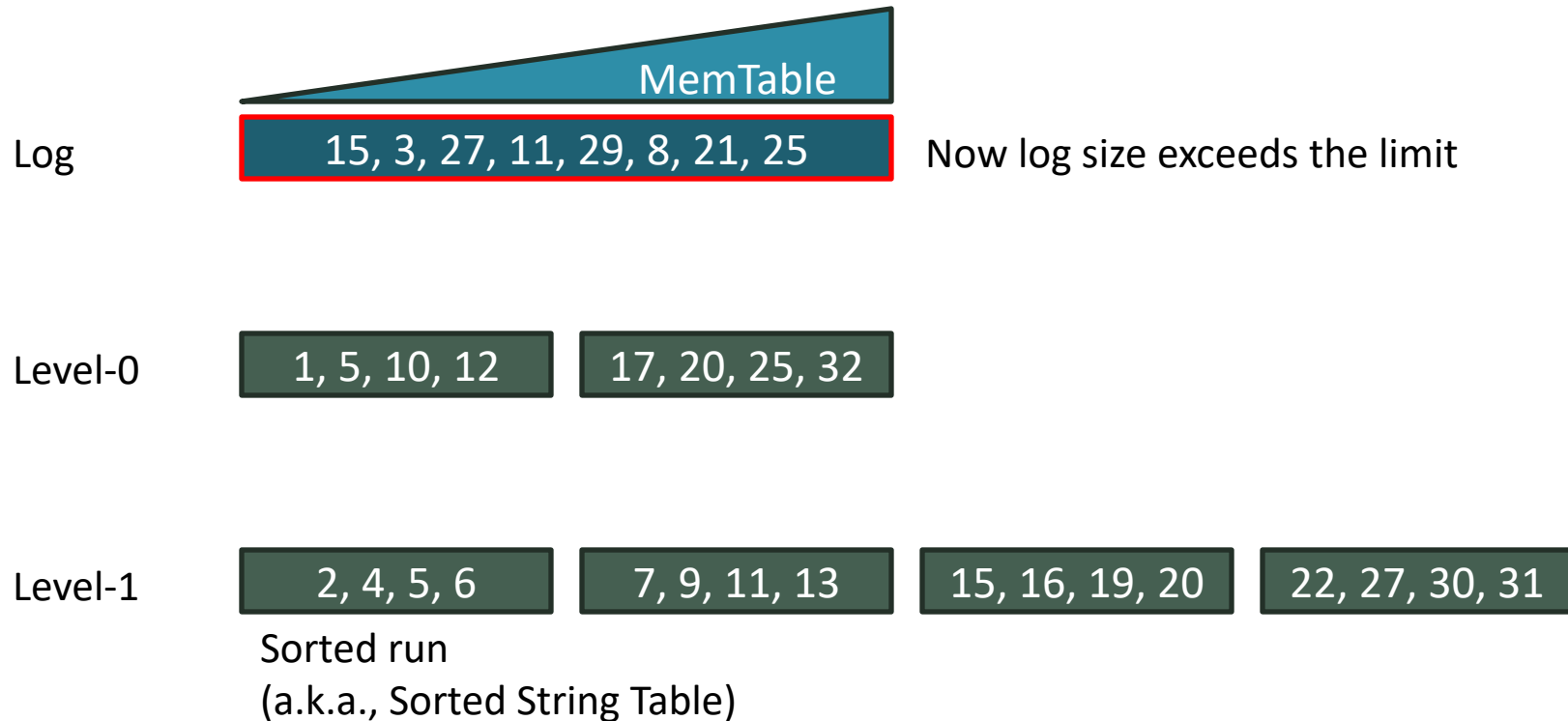
In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge



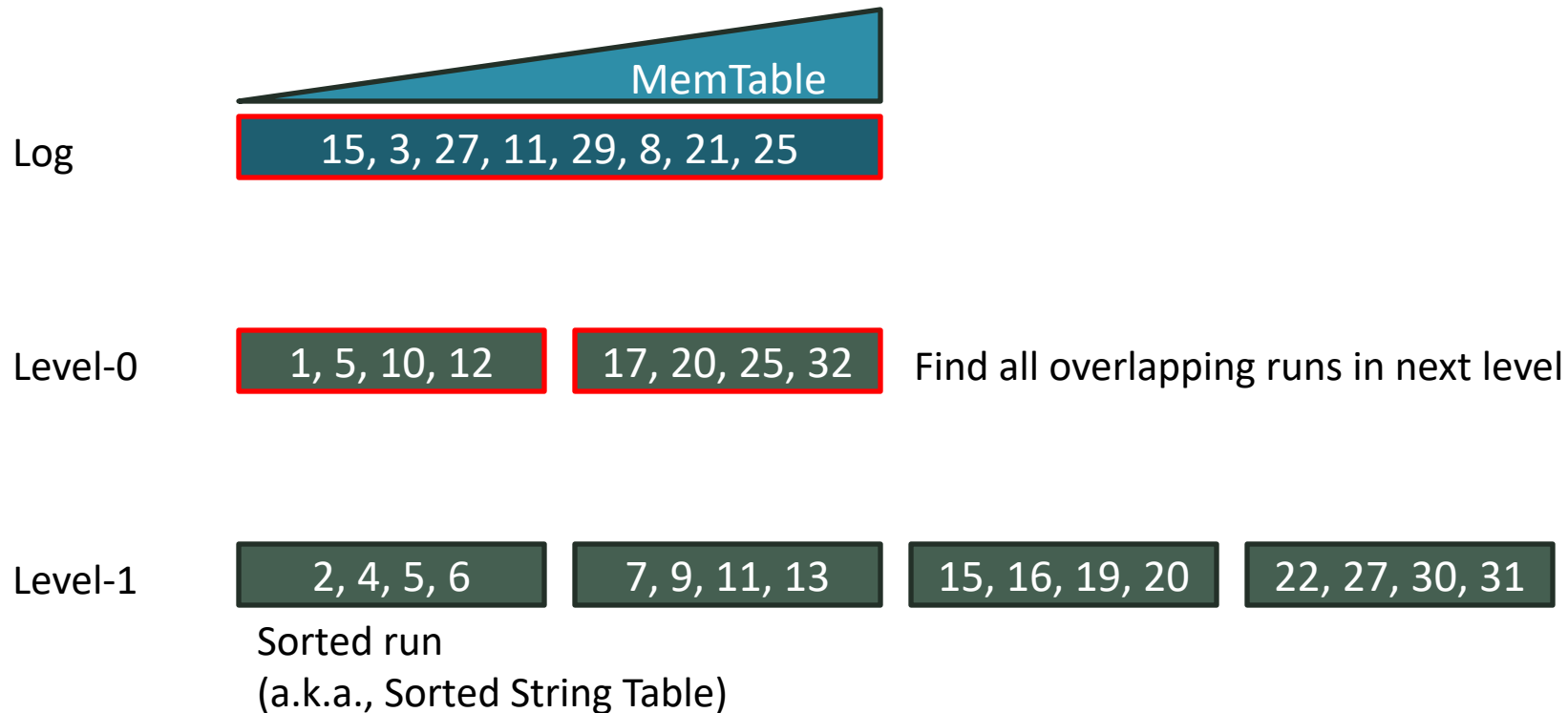
In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge



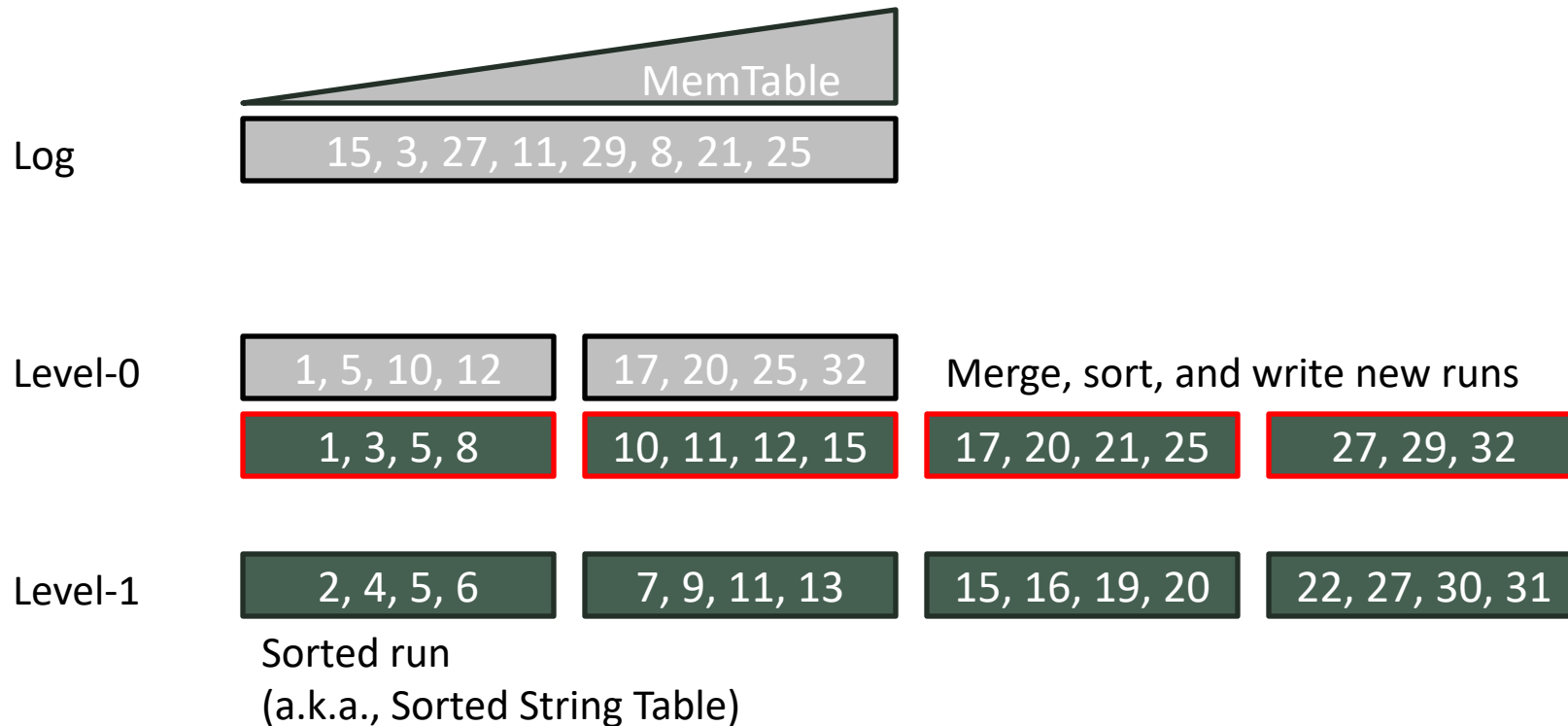
In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge



In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge

Log

Remove old runs & logs

Level-0



Level-1



Sorted run
(a.k.a., Sorted String Table)

In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge

Log

Merge can be **cascaded**:
now level-0 becomes full

Level-0



Level-1



Sorted run
(a.k.a., Sorted String Table)

In this example:

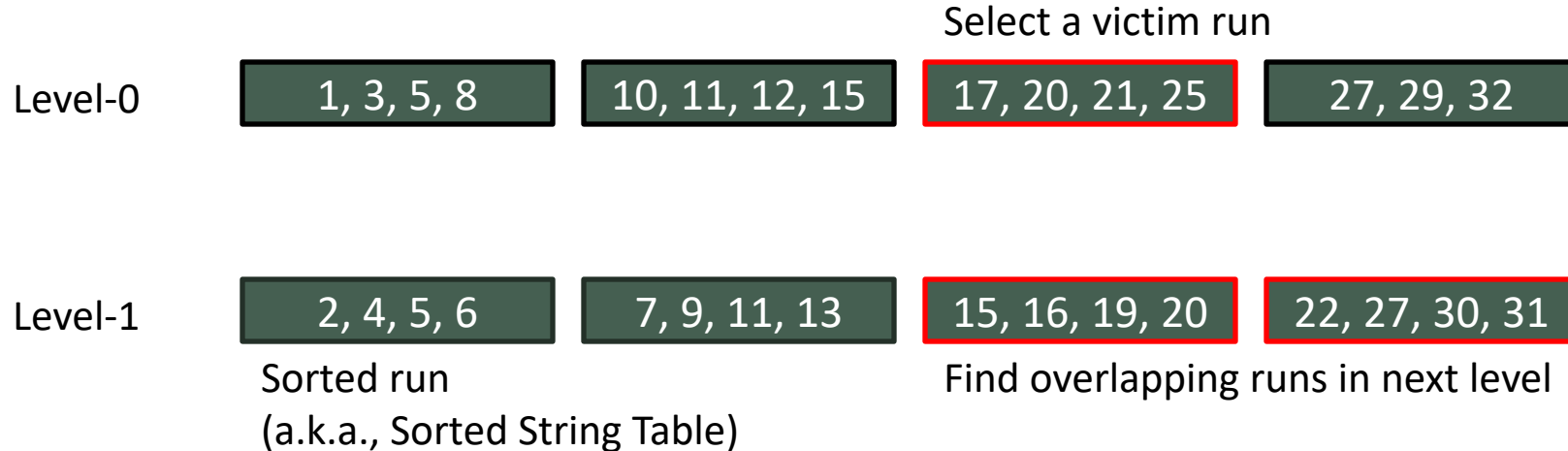
- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge

Log



In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge

Log

Level-0

1, 3, 5, 8	10, 11, 12, 15	17, 20, 21, 25	27, 29, 32
------------	----------------	----------------	------------

Level-1

2, 4, 5, 6	7, 9, 11, 13	15, 16, 19, 20	22, 27, 30, 31
------------	--------------	----------------	----------------

Sorted run
(a.k.a., Sorted String Table)

15, 16, 17, 19

20, 21, 22, 25

Merge, sort, and write new runs

27, 30, 31

In this example:

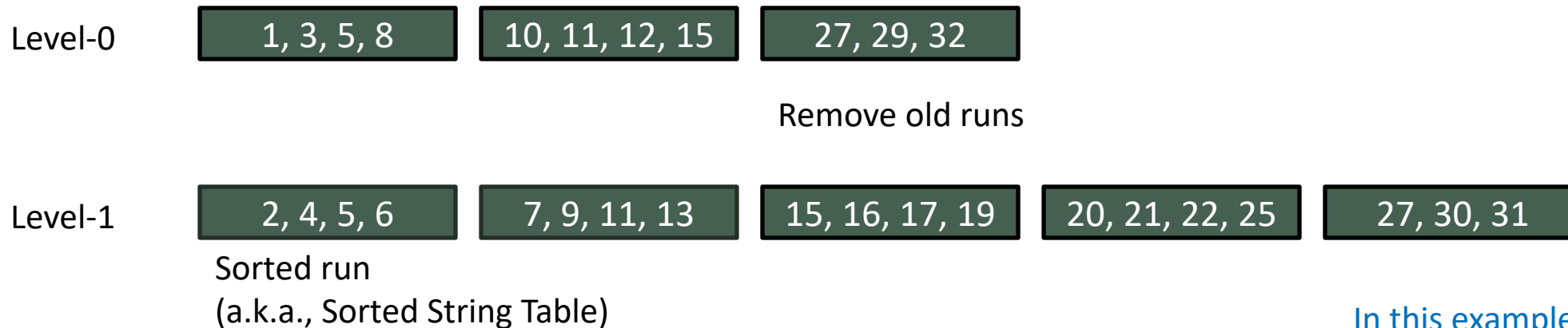
- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge

Log



In this example:

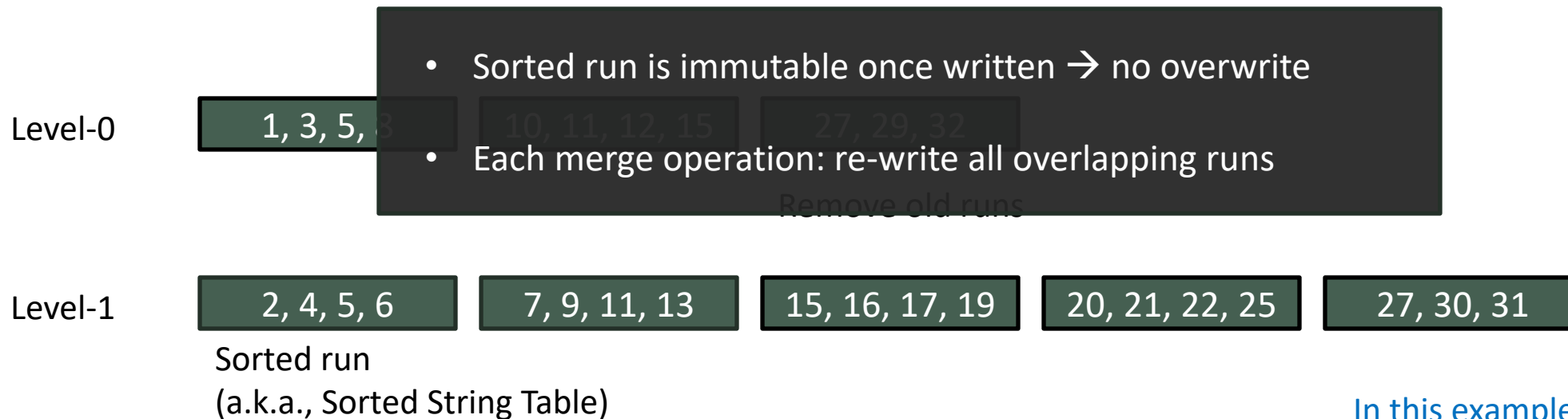
- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – write and merge

Log



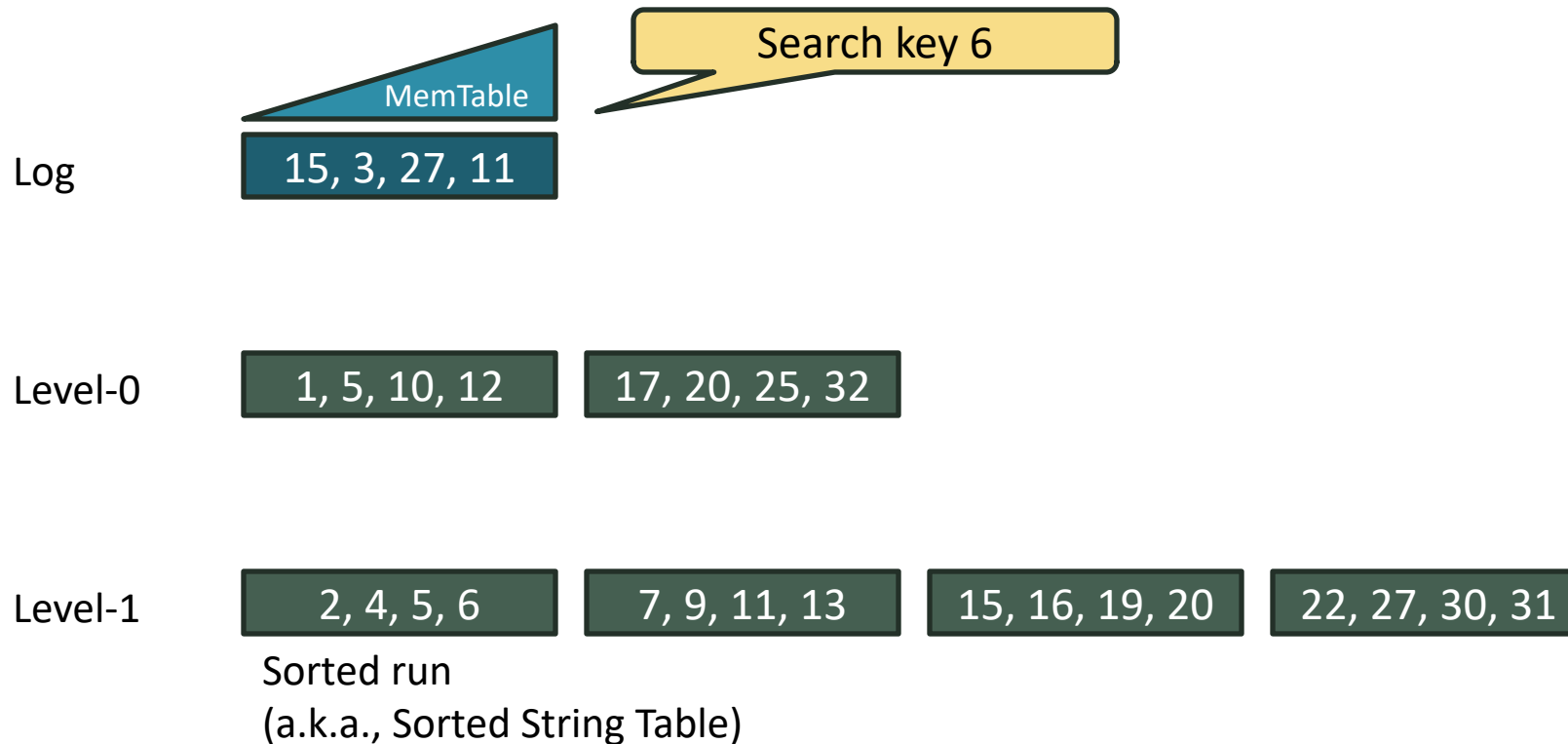
In this example:

- Run size limit: 4 keys
- Level-0: 2 runs
- Level-1: 4 runs (2x)

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – read

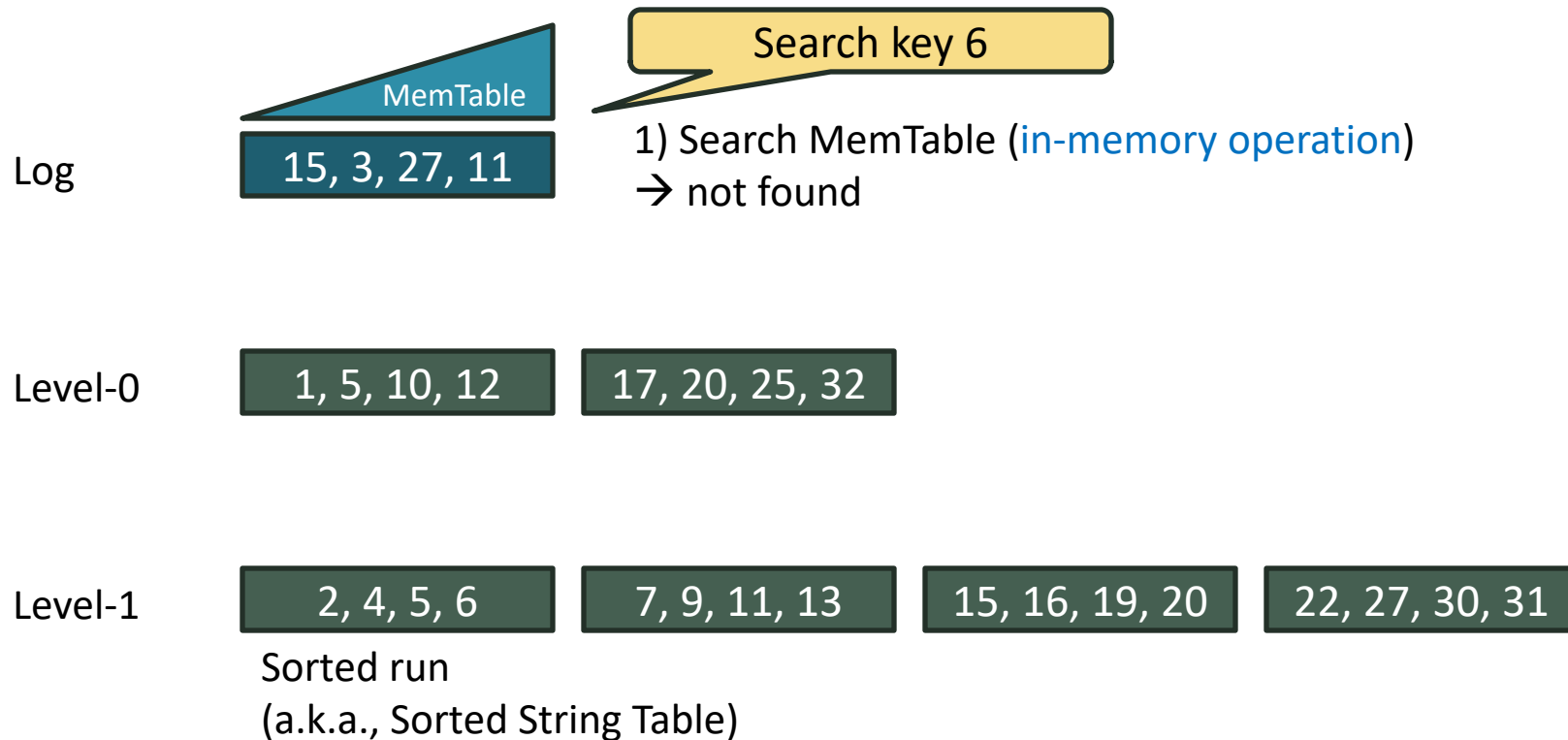


- Inner-level search: **logarithmic**
- Inter-level search: **linear**

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – read

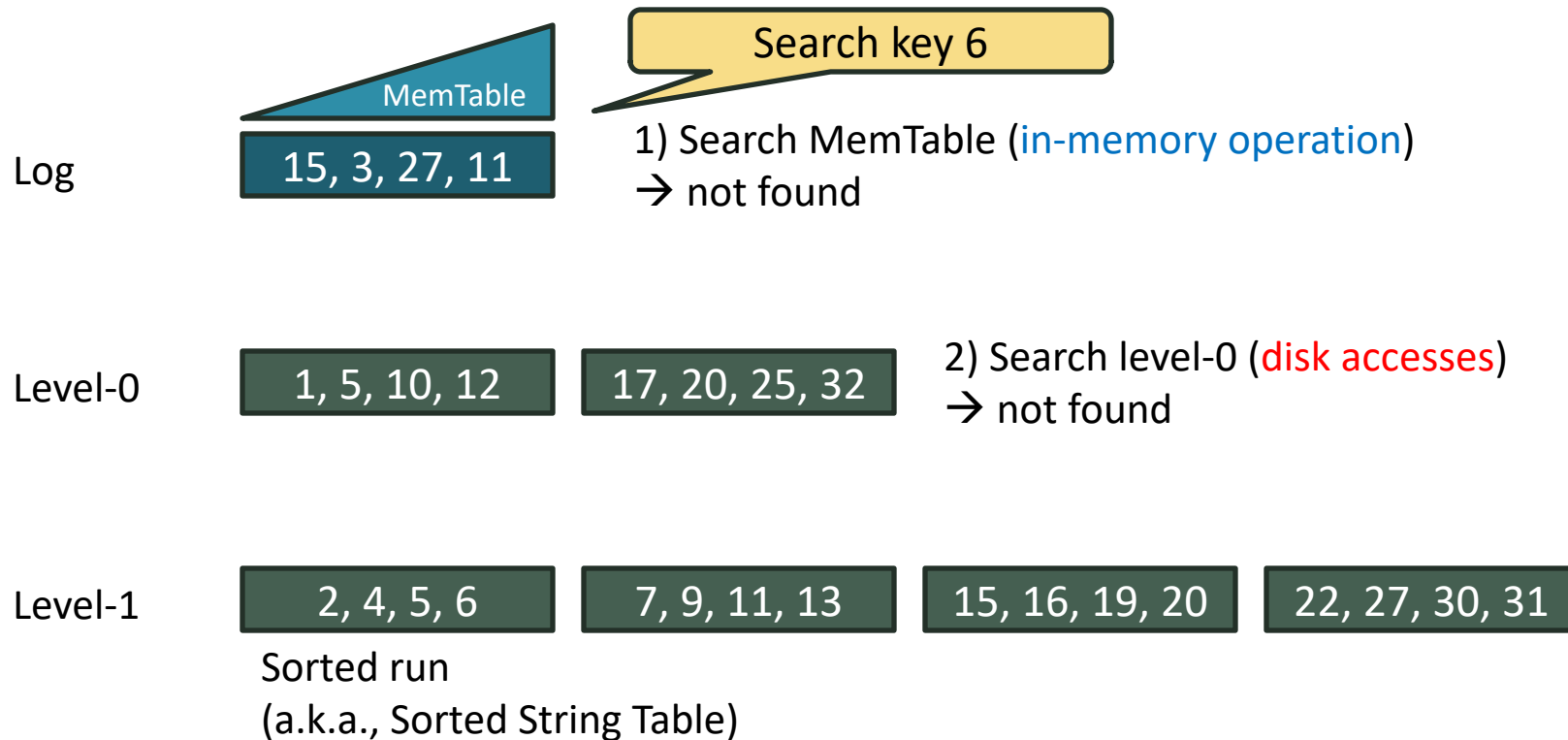


- Inner-level search: **logarithmic**
- Inter-level search: **linear**

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – read

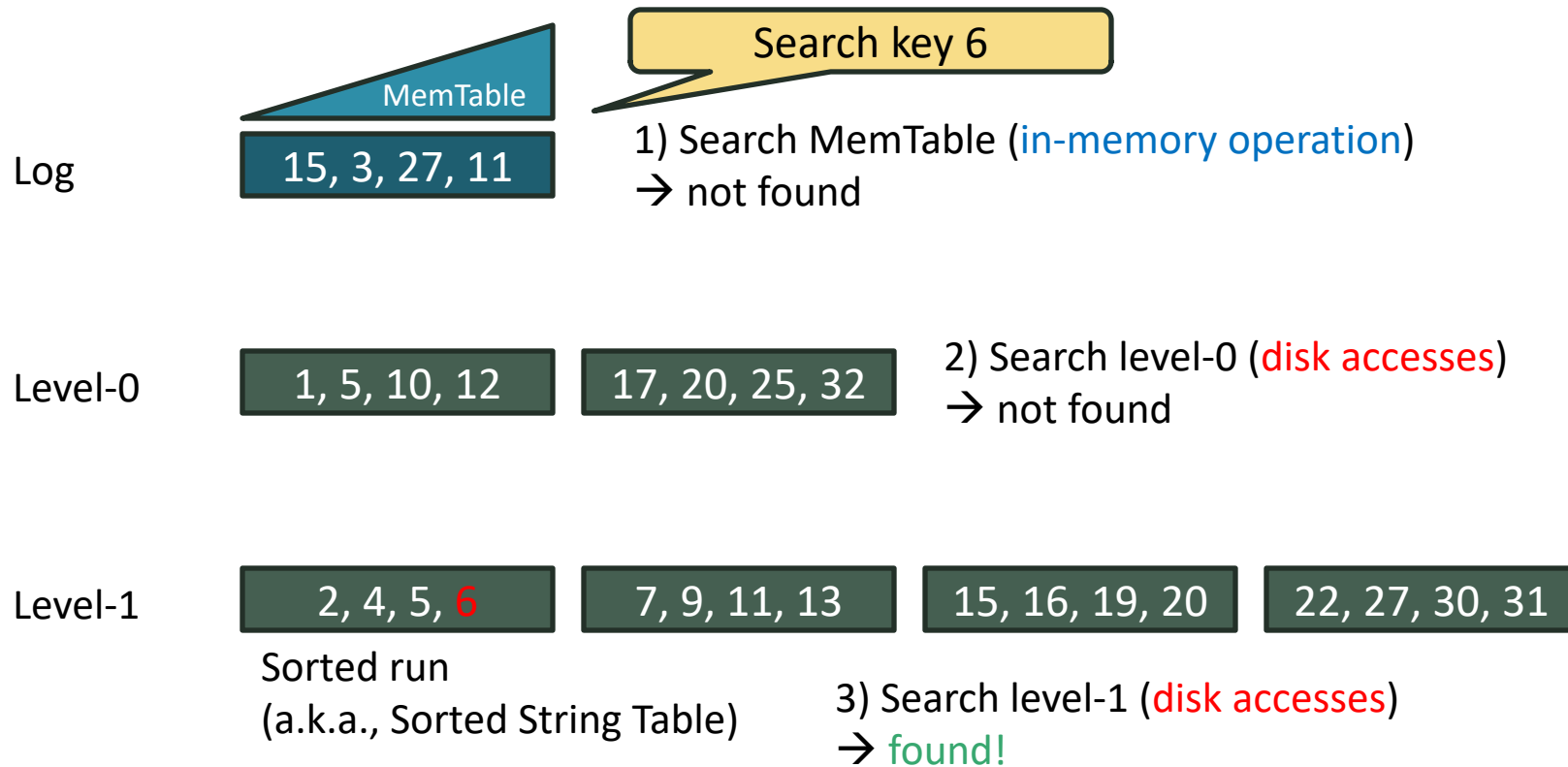


- Inner-level search: logarithmic
- Inter-level search: linear

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – read



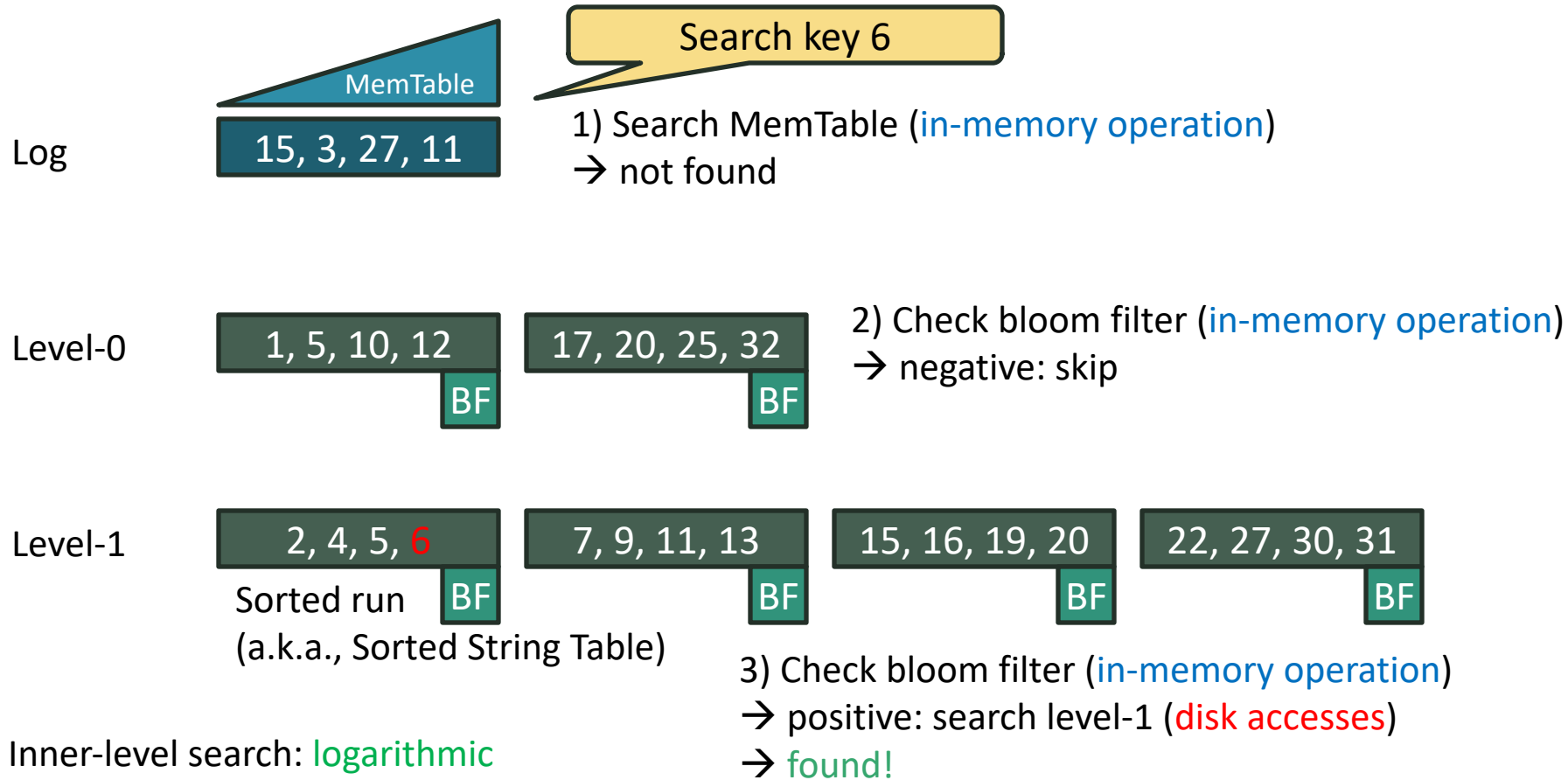
- Inner-level search: logarithmic
- Inter-level search: linear

Search each level one by one

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

- Basic algorithm – read
 - Bloom filter: skip examining unnecessary runs

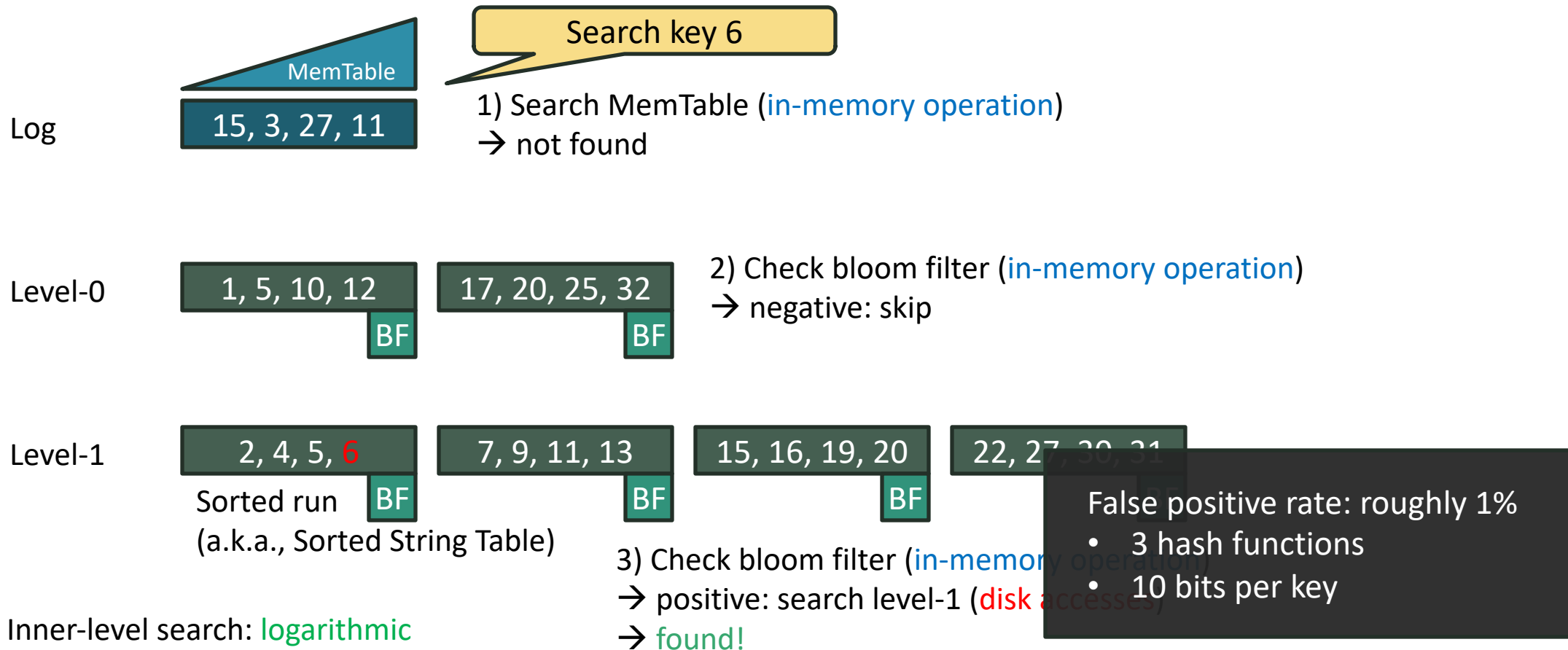


- Inner-level search: logarithmic
- Inter-level search: linear

Log-Structured Merge-Tree (LSM-Tree)

(P. O'Neil et al. 1996)

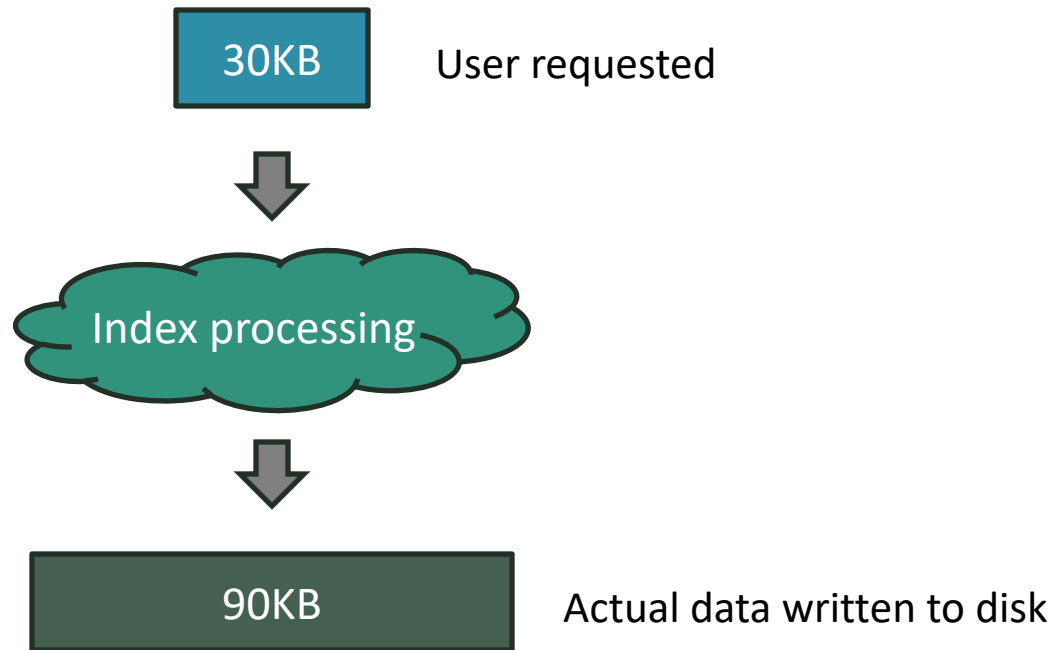
- Basic algorithm – read
 - Bloom filter: skip examining unnecessary runs



- Inner-level search: logarithmic
- Inter-level search: linear

Write Amplification

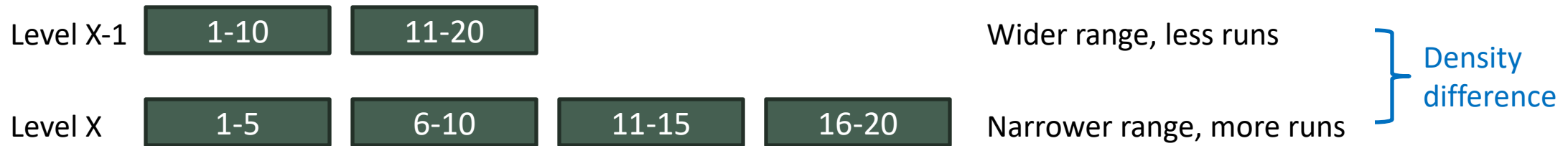
- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk



Write amplification: 3

Write Amplification

- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk
- LSM-tree: merge amplifies amount of writes



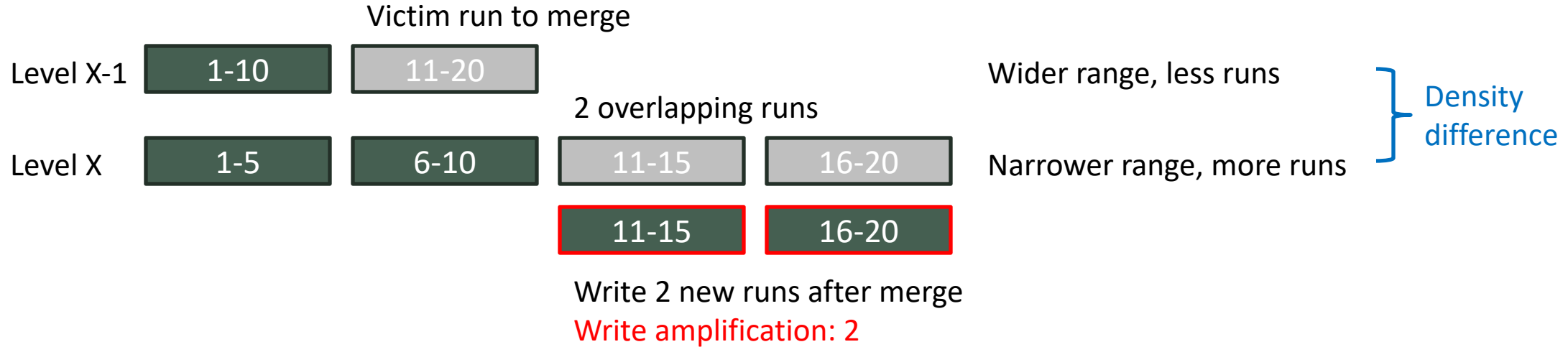
Write Amplification

- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk
- LSM-tree: merge amplifies amount of writes



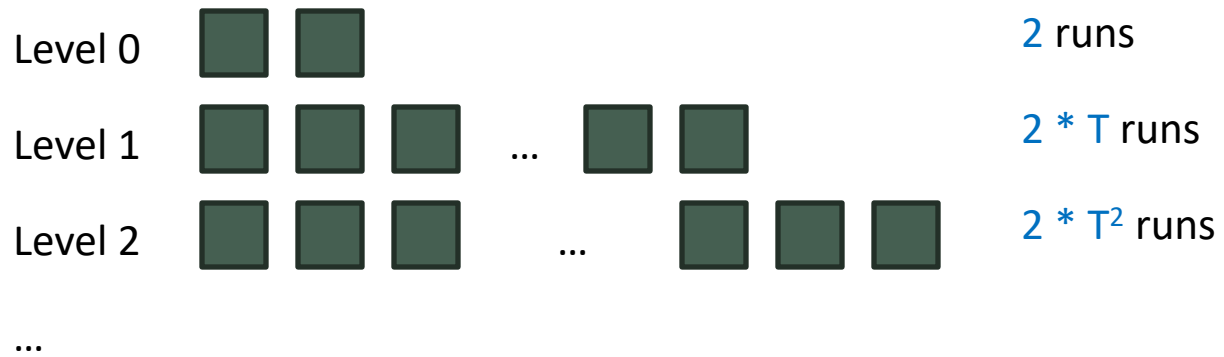
Write Amplification

- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk
- LSM-tree: merge amplifies amount of writes



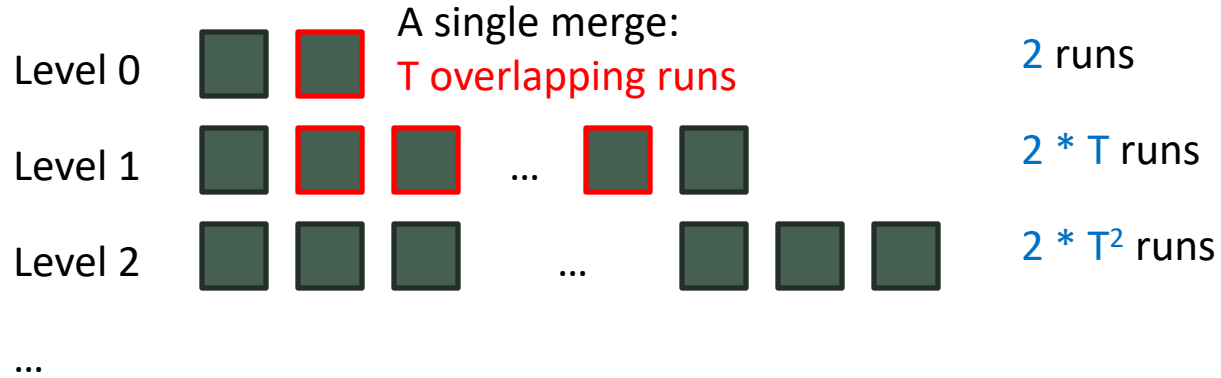
Write Amplification

- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk
- LSM-tree: merge amplifies amount of writes
 - T : size ratio of adjacent levels



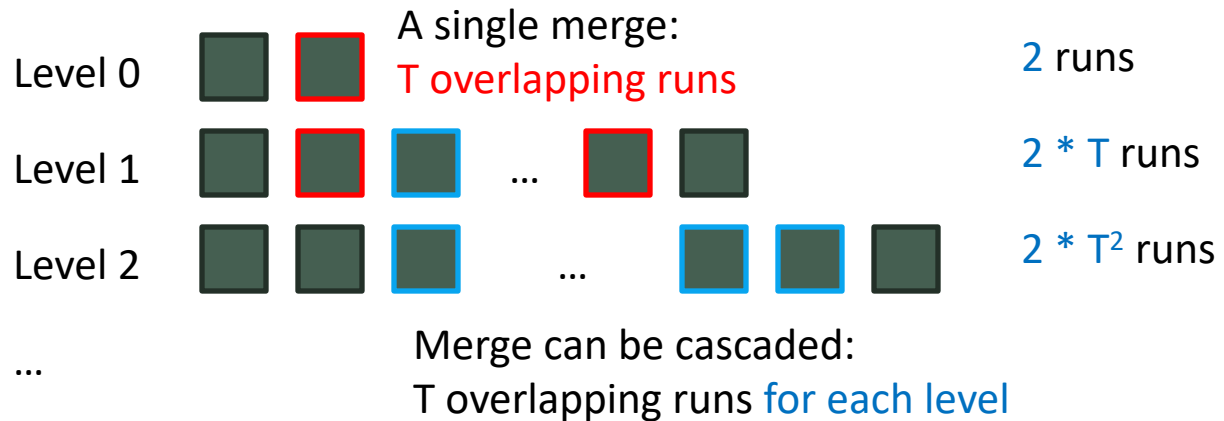
Write Amplification

- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk
- LSM-tree: merge amplifies amount of writes
 - T : size ratio of adjacent levels



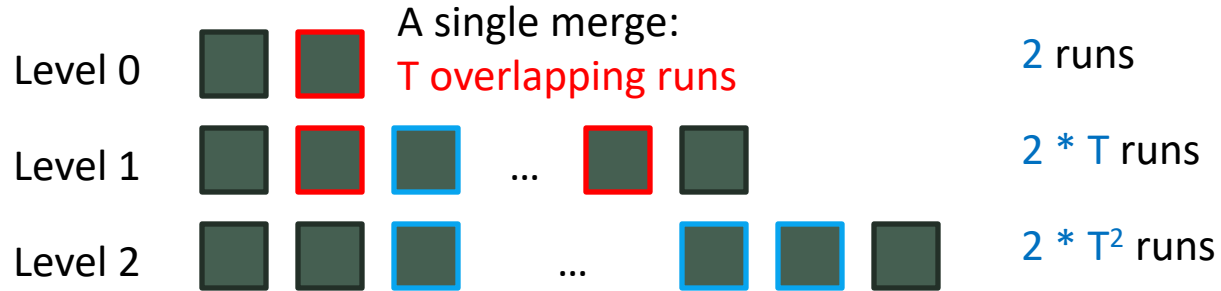
Write Amplification

- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk
- LSM-tree: merge amplifies amount of writes
 - T: size ratio of adjacent levels



Write Amplification

- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk
- LSM-tree: merge amplifies amount of writes
 - T: size ratio of adjacent levels



... Merge can be cascaded:
T overlapping runs for each level

$$\text{Write amplification: } T + T + \dots + T = O(L * T)$$

The number of levels (L)

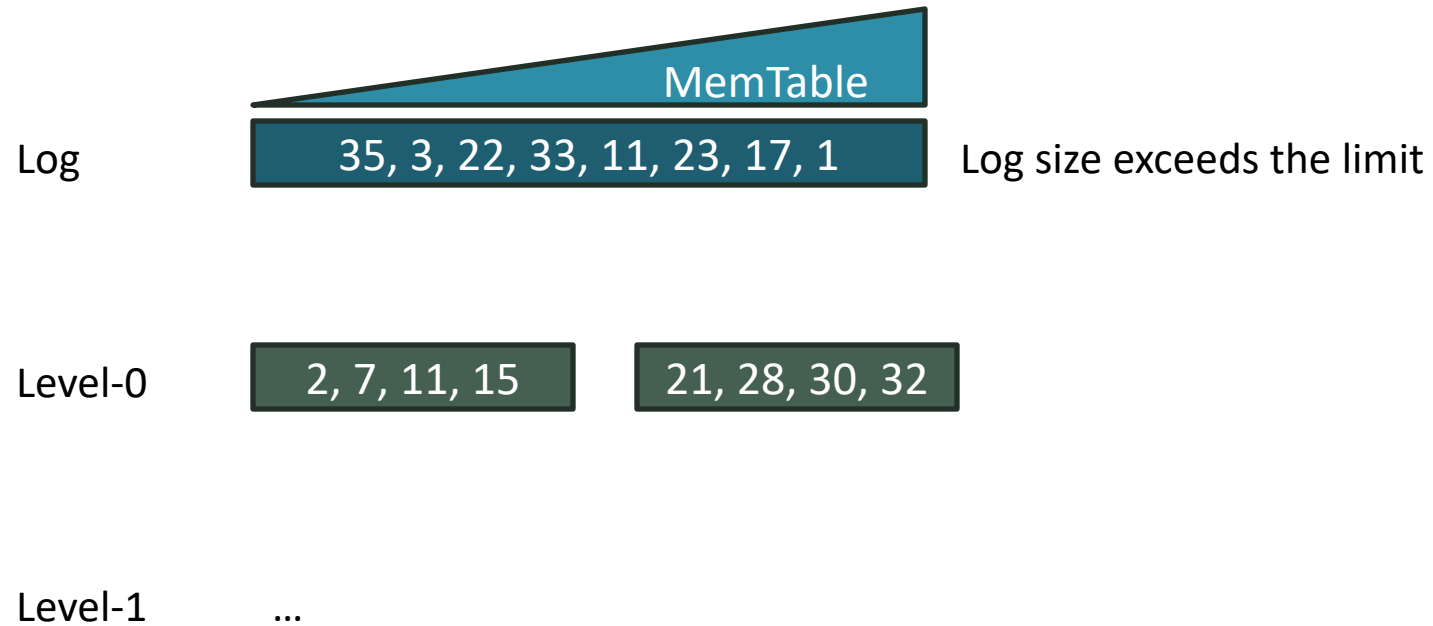
Write Amplification

- Ratio between
 - Amount of data requested by user
 - Amount of data actually written to disk
- LSM-tree: merge amplifies amount of writes
 - T: size ratio of adjacent levels
 - Write amplification: $O(L * T)$
- Why write amplification matters?
 - User: put 2KB records at 10,000 records/sec rate (20MB/s)
 - Write amplification: 20
 - Actual traffic to disk: 400MB/s
 - Easily hit the upper bound of disk bandwidth
 - Affect read latency or throughput as well

Tiering Merge

(H.V. Jagadish et al. VLDB 1997, PebblesDB SOSP 2017, Dostoevsky SIGMOD 2018)

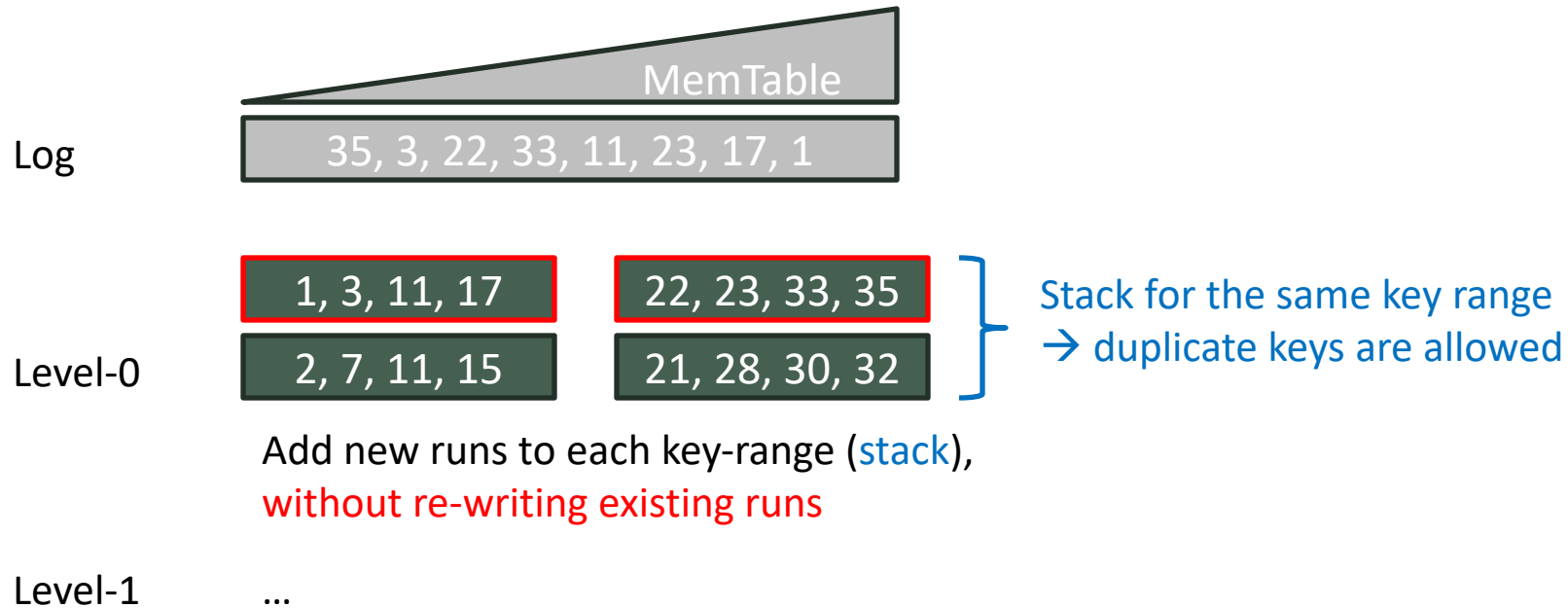
- Delay merge, and keep stacks of runs (for the same key range)



Tiering Merge

(H.V. Jagadish et al. VLDB 1997, PebblesDB SOSP 2017, Dostoevsky SIGMOD 2018)

- Delay merge, and keep stacks of runs (for the same key range)

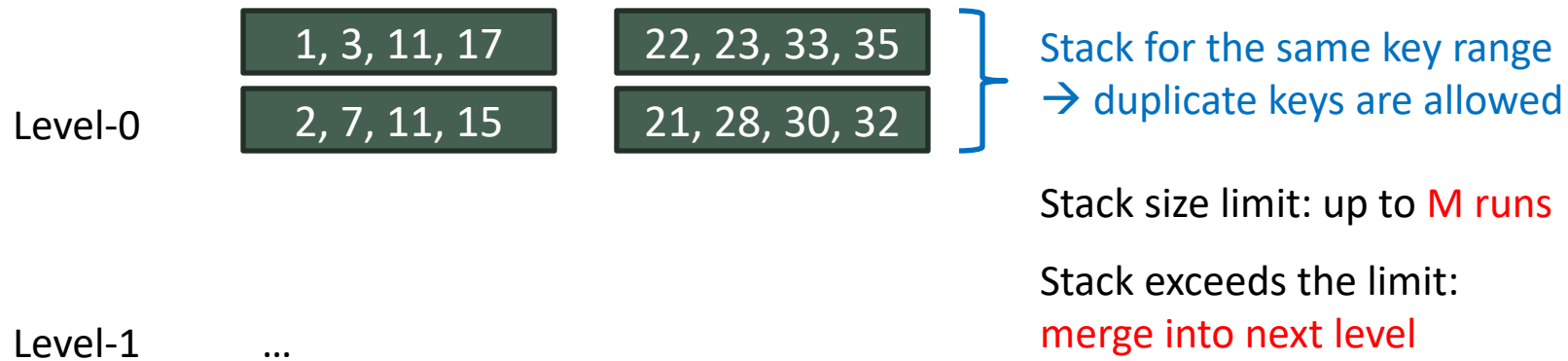


Tiering Merge

(H.V. Jagadish et al. VLDB 1997, PebblesDB SOSP 2017, Dostoevsky SIGMOD 2018)

- Delay merge, and keep stacks of runs (for the same key range)

Log

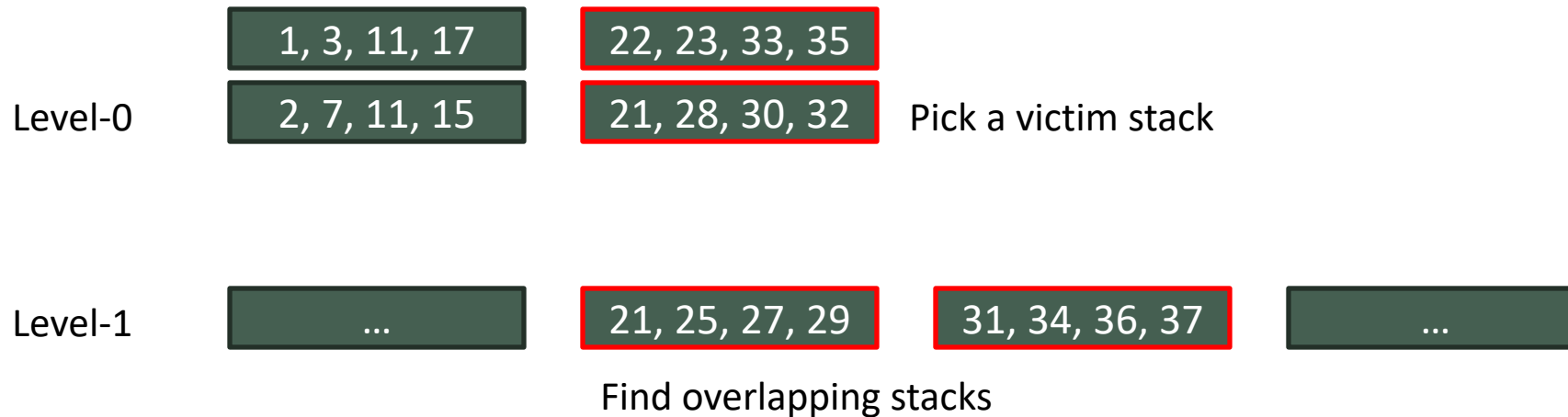


Tiering Merge

(H.V. Jagadish et al. VLDB 1997, PebblesDB SOSP 2017, Dostoevsky SIGMOD 2018)

- Delay merge, and keep stacks of runs (for the same key range)

Log

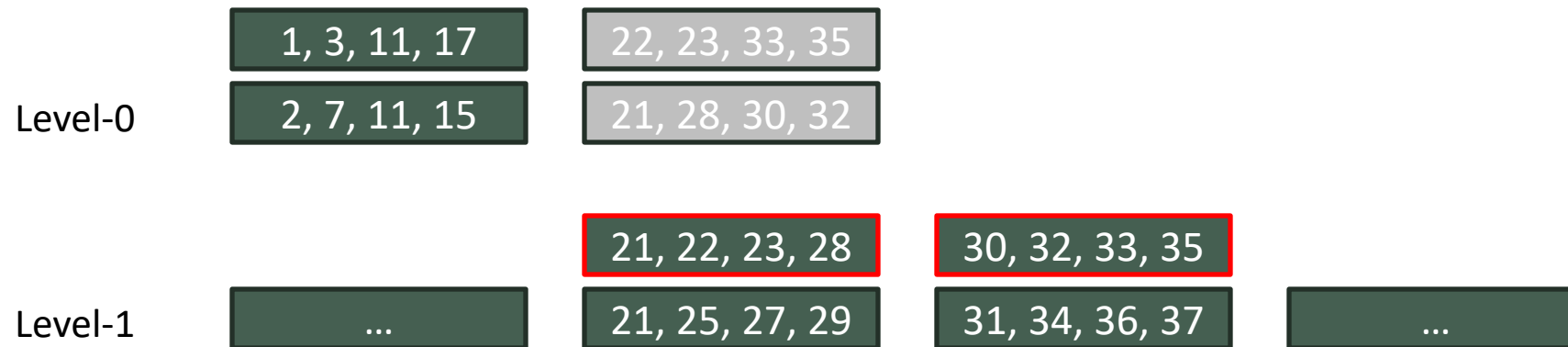


Tiering Merge

(H.V. Jagadish et al. VLDB 1997, PebblesDB SOSP 2017, Dostoevsky SIGMOD 2018)

- Delay merge, and keep stacks of runs (for the same key range)

Log

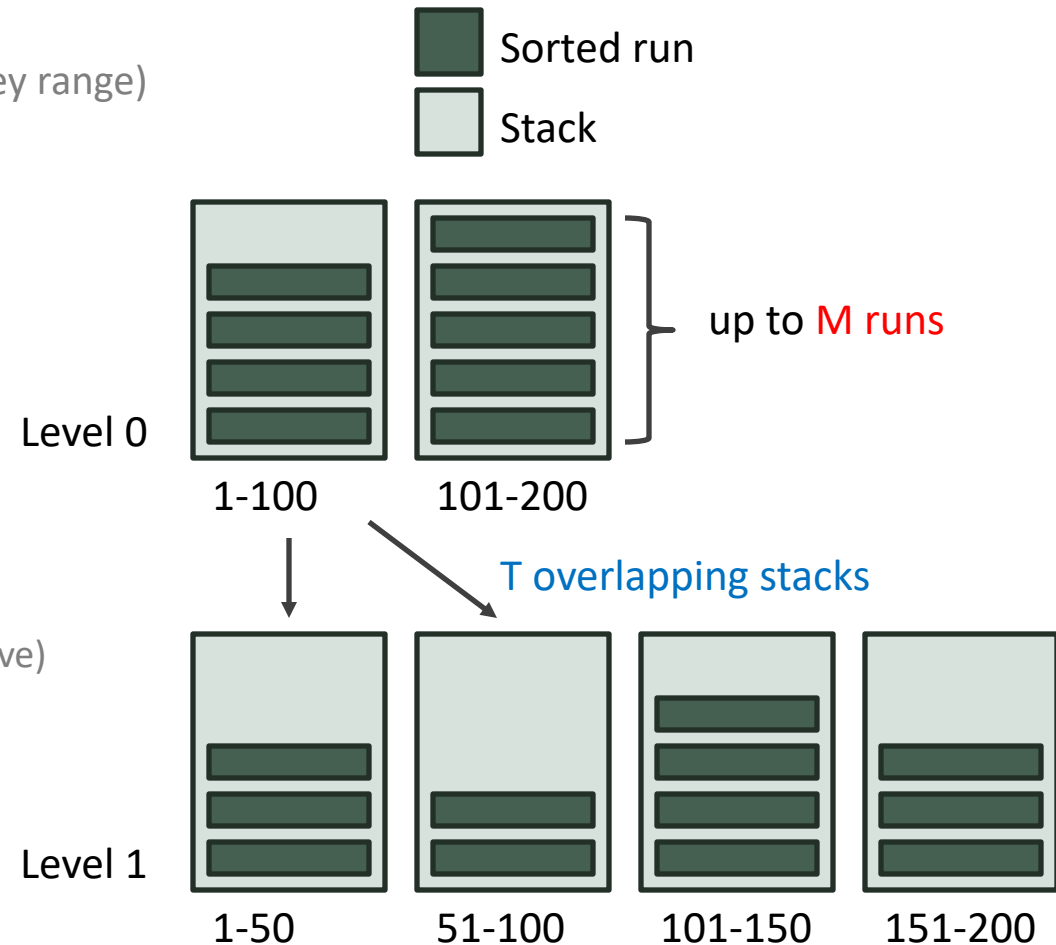


Add new runs to each stack,
without re-writing existing runs

Tiering Merge

(H.V. Jagadish et al. VLDB 1997, PebblesDB SOSP 2017, Dostoevsky SIGMOD 2018)

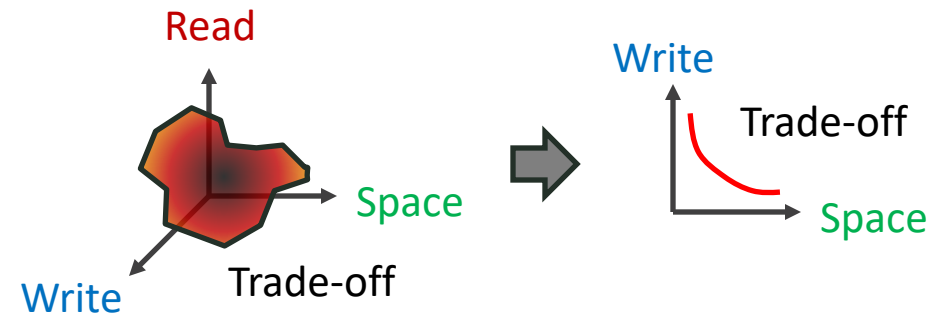
- Delay merge, and keep stacks of runs (for the same key range)
- Trade-offs
 - Stack size limit: up to M runs
 - Write cost
 - Can delay merge M times
 - M times smaller write amplification
 - Read cost
 - M times more searches
 - (even with Bloom filter, more searches upon false positive)
 - Space cost
 - M times more space



↓ Write cost vs. read cost and space cost ↑

Jungle: Motivation

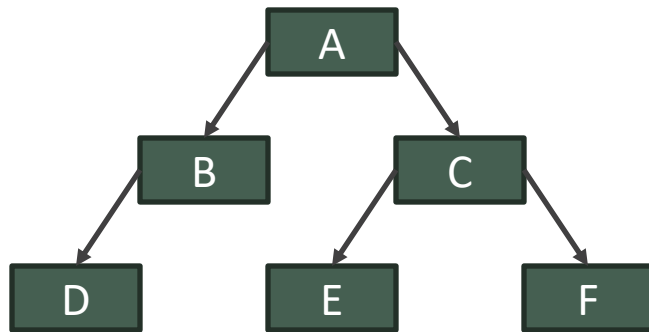
- What if we can
 - Reduce write cost
 - Without sacrificing read cost
- Trade-off change: 3-dimensional \rightarrow 2-dimensional
- How?
 - Replace each stack with copy-on-write B+tree



Copy-On-Write (CoW) B+Tree

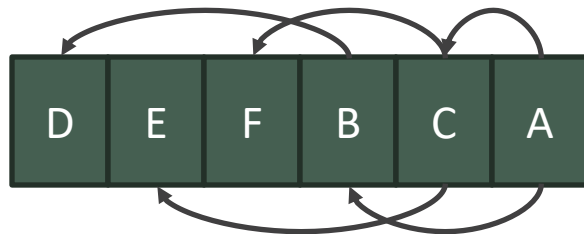
- Logically the same, but about how to write B+tree nodes to disk
 - Out-of-update manner: no overwrite, but append
- All nodes are immutable, written in chronological order

Logical view



 B+tree node

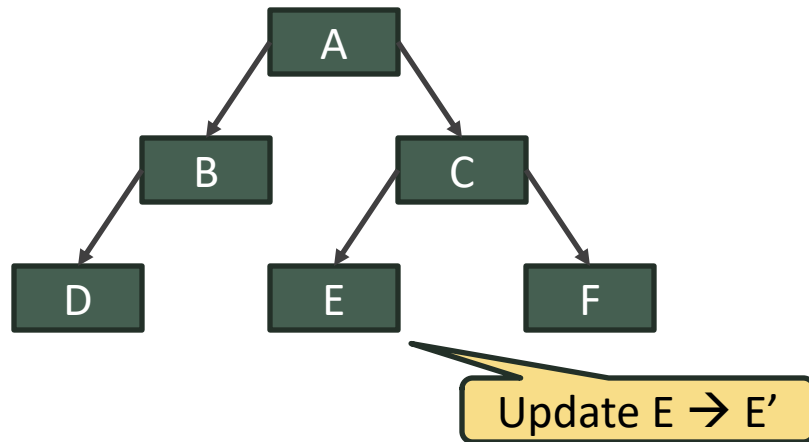
Flattened view
(on disk)



Copy-On-Write (CoW) B+Tree

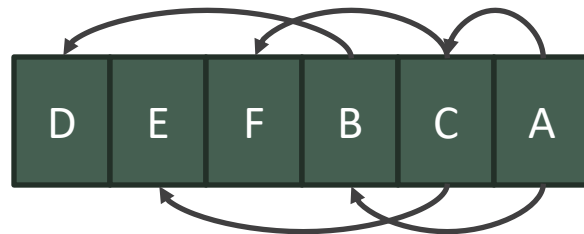
- Logically the same, but about how to write B+tree nodes to disk
 - Out-of-update manner: no overwrite, but append
- All nodes are immutable, written in chronological order

Logical view



 B+tree node

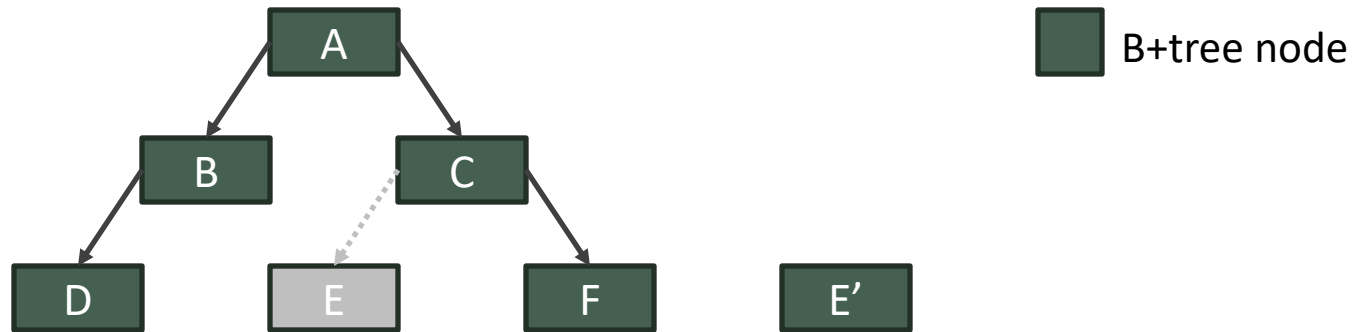
Flattened view
(on disk)



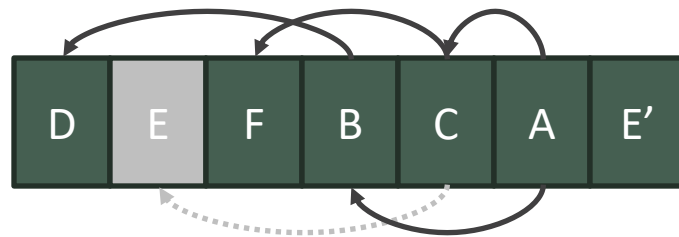
Copy-On-Write (CoW) B+Tree

- Logically the same, but about how to write B+tree nodes to disk
 - Out-of-update manner: no overwrite, but append
- All nodes are immutable, written in chronological order

Logical view



Flattened view
(on disk)

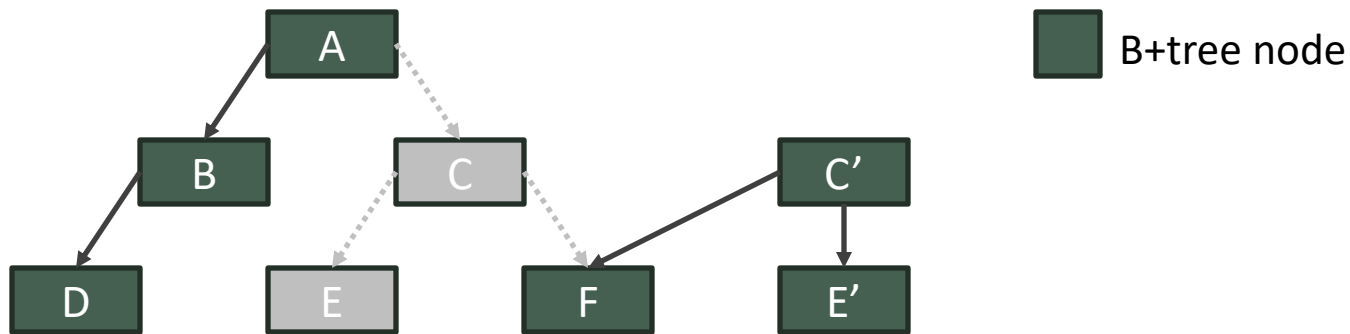


Append E'
→ parent C also needs to be updated

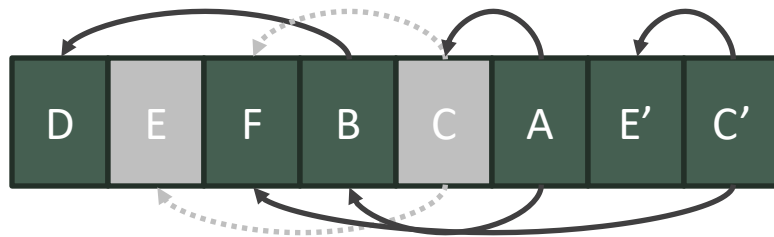
Copy-On-Write (CoW) B+Tree

- Logically the same, but about how to write B+tree nodes to disk
 - Out-of-update manner: no overwrite, but append
- All nodes are immutable, written in chronological order

Logical view



Flattened view
(on disk)

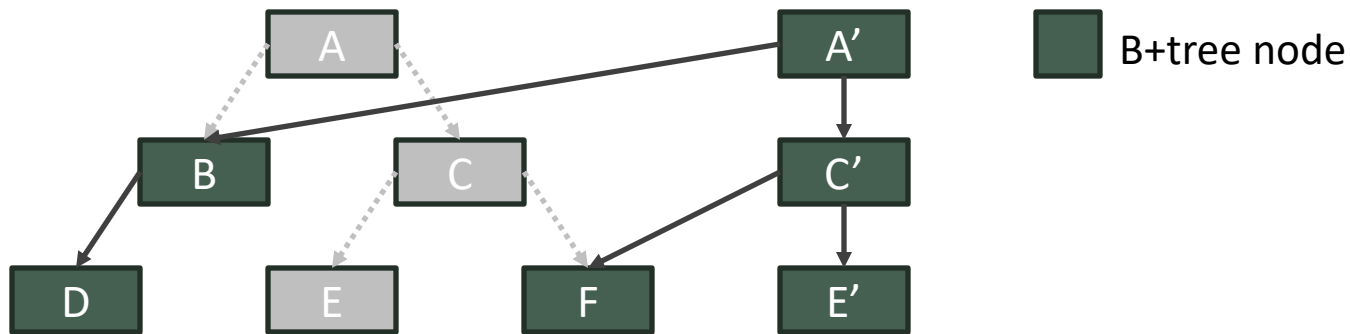


Append C'
→ parent A also needs to be updated

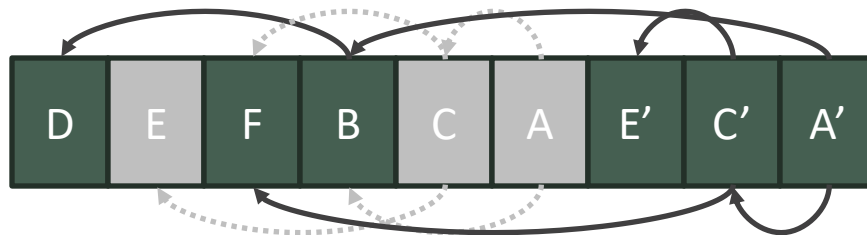
Copy-On-Write (CoW) B+Tree

- Logically the same, but about how to write B+tree nodes to disk
 - Out-of-update manner: no overwrite, but append
- All nodes are immutable, written in chronological order

Logical view



Flattened view
(on disk)

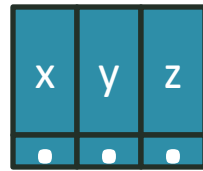


Copy-On-Write (CoW) B+Tree

- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs



K: value of key x



B+tree node with key x, y, and z

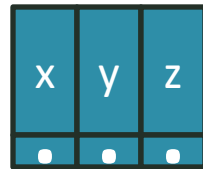
Insert {10, A}, {15, B}, {27, C}, {50, D}

Copy-On-Write (CoW) B+Tree

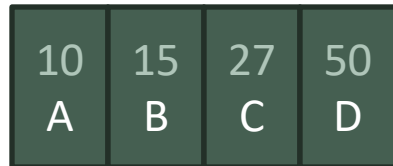
- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs



K: value of key x



B+tree node with key x, y, and z

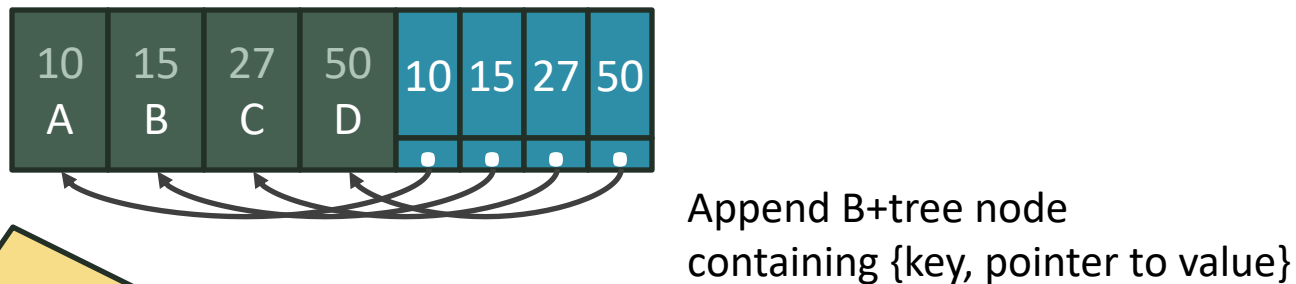
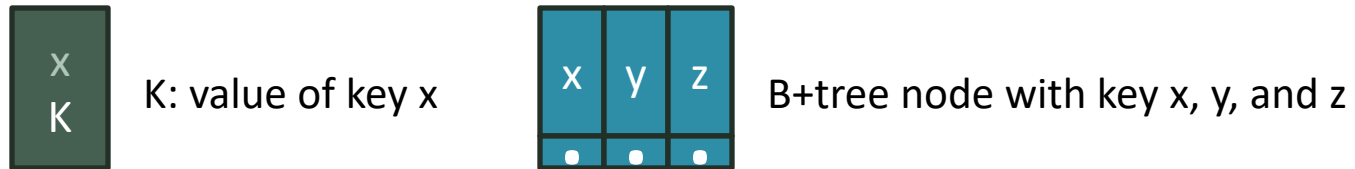


Append values (in key order)

Insert {10, A}, {15, B}, {27, C}, {50, D}

Copy-On-Write (CoW) B+Tree

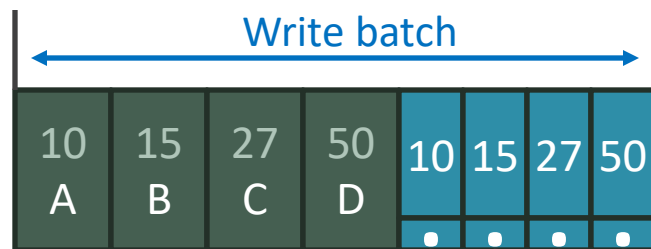
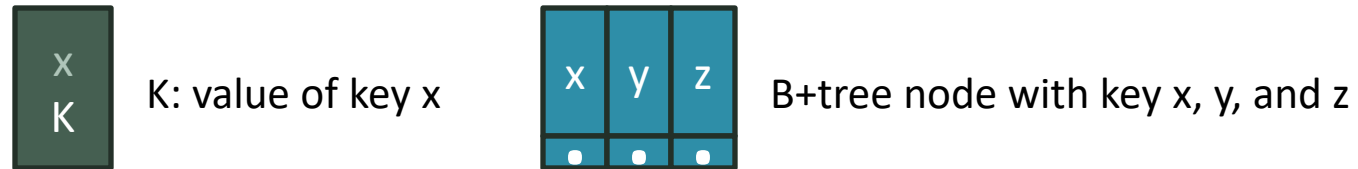
- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs



Insert {10, A}, {15, B}, {27, C}, {50, D}

Copy-On-Write (CoW) B+Tree

- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs

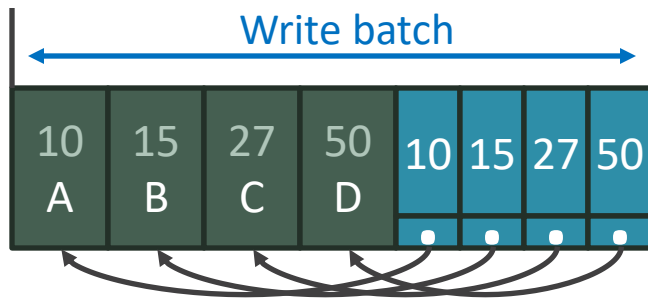
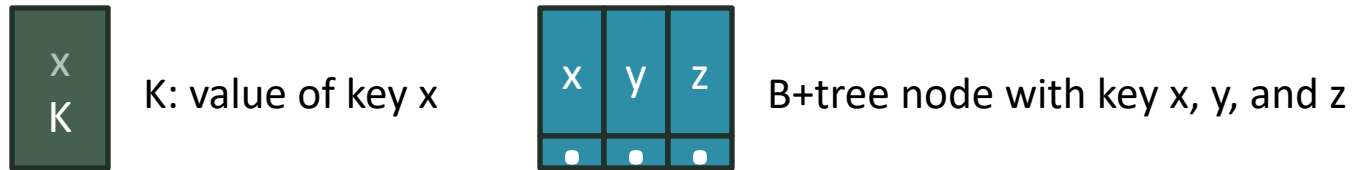


Append B+tree node containing {key, pointer to value}

Insert {10, A}, {15, B}, {27, C}, {50, D}

Copy-On-Write (CoW) B+Tree

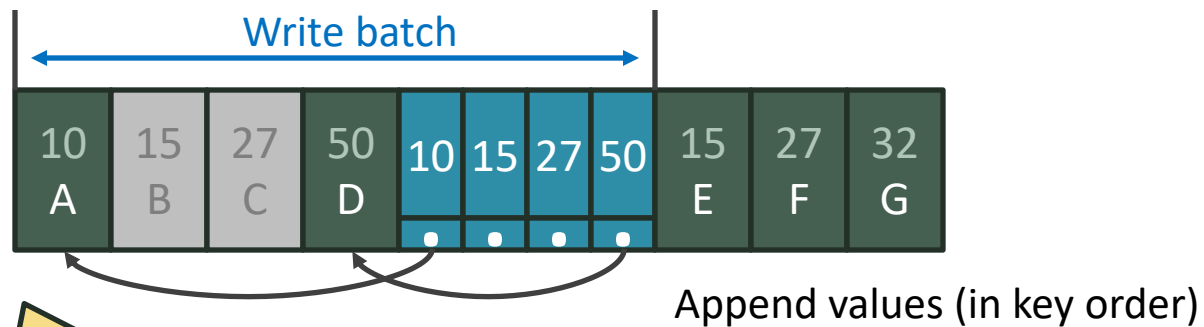
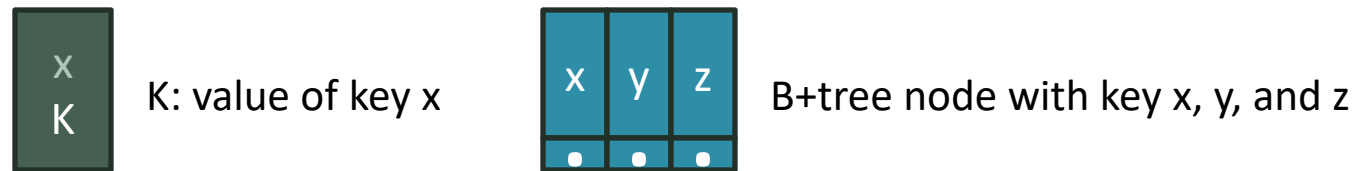
- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs



Set {15, E}, {27, F}, {32, G}

Copy-On-Write (CoW) B+Tree

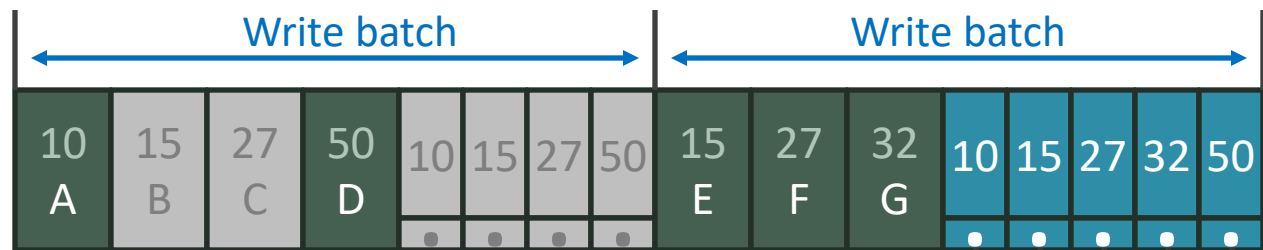
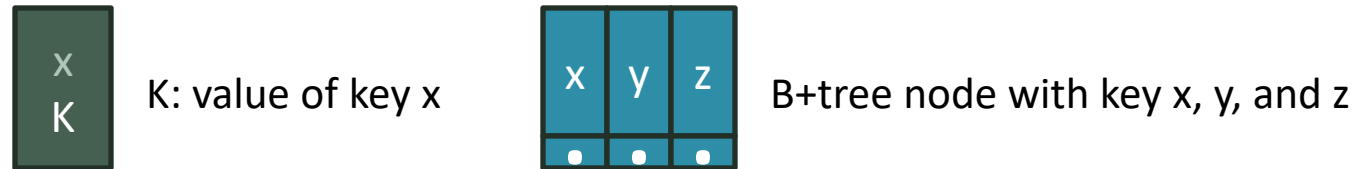
- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs



Set {15, E}, {27, F}, {32, G}

Copy-On-Write (CoW) B+Tree

- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs



Append updated B+tree node

Set {15, E}, {27, F}, {32, G}

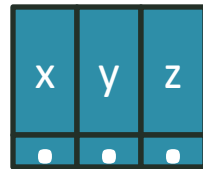
Keys in the same node should be copied together although they are not updated

Copy-On-Write (CoW) B+Tree

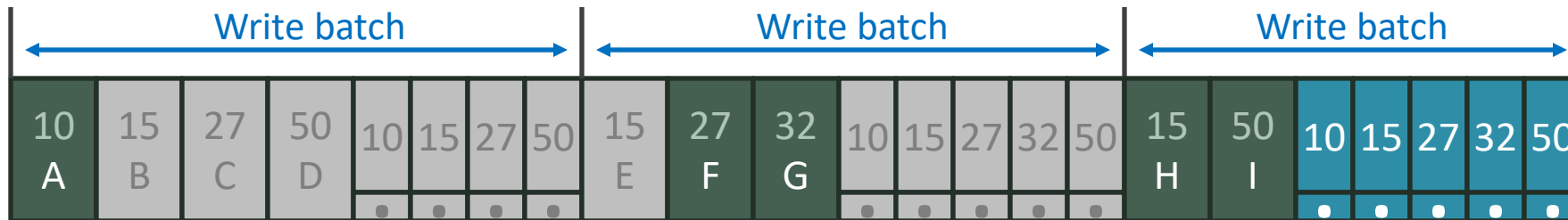
- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs



K: value of key x



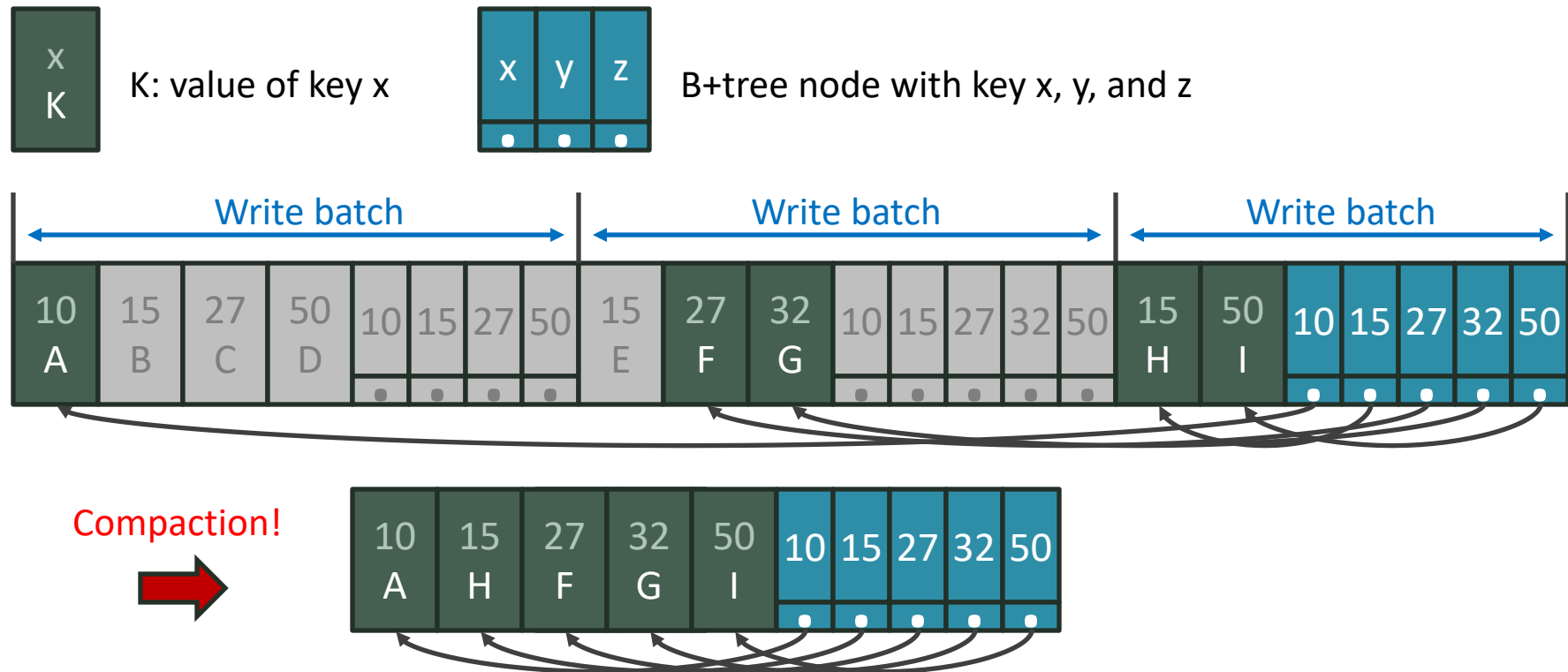
B+tree node with key x, y, and z



Set {15, H}, {50, I}

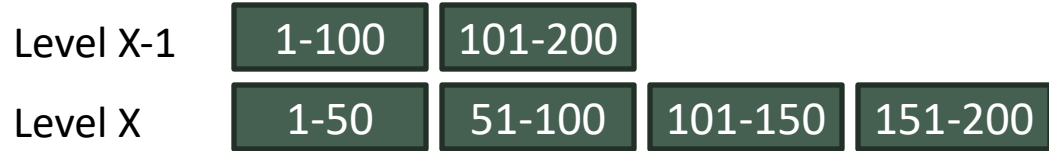
Copy-On-Write (CoW) B+Tree

- Batching + decoupling value from B+tree node
 - To reduce write amplification + to make B+tree compact
 - Append values first, and then append updated B+tree nodes
 - B+tree leaf nodes: contain {key, pointer to value} pairs

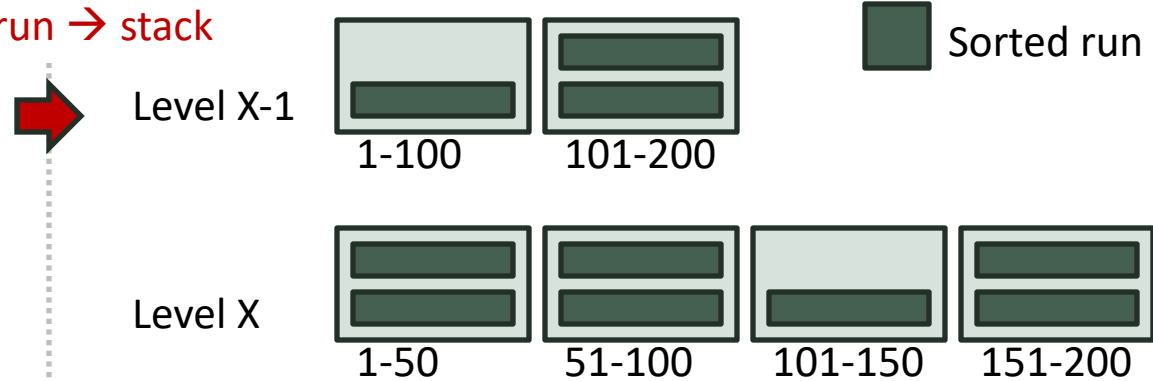


Jungle: Replacing Stack with CoW B+Tree

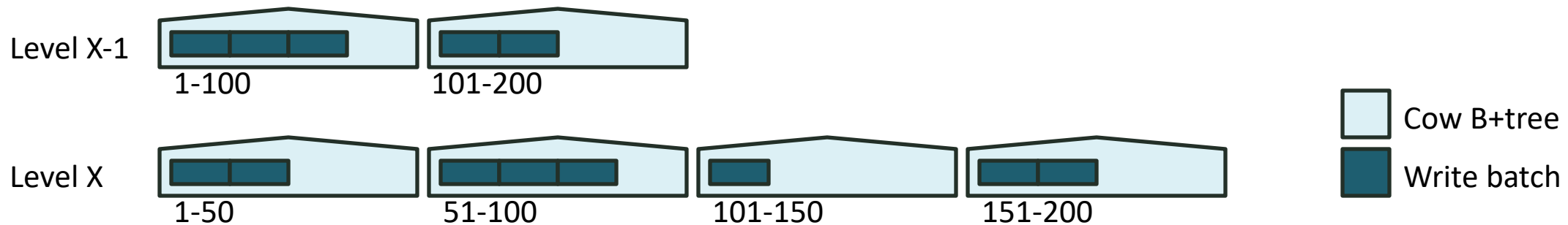
- Original LSM-tree (leveling)



- Tiering merge LSM-tree



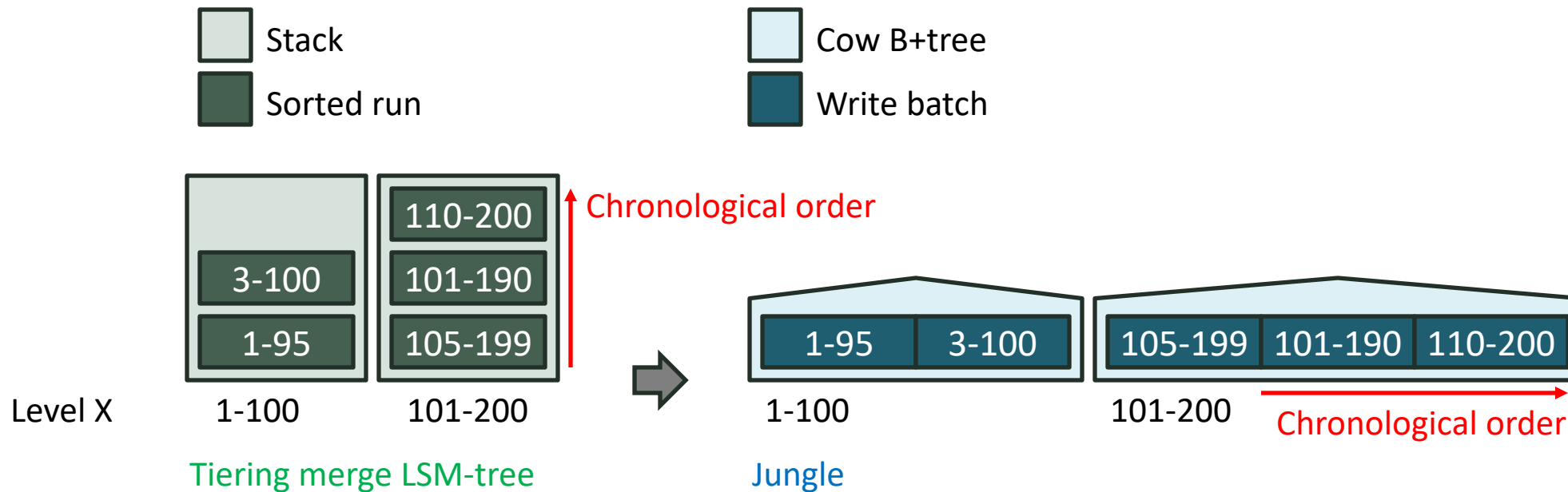
- Our approach: Jungle



Stack → cow B+tree

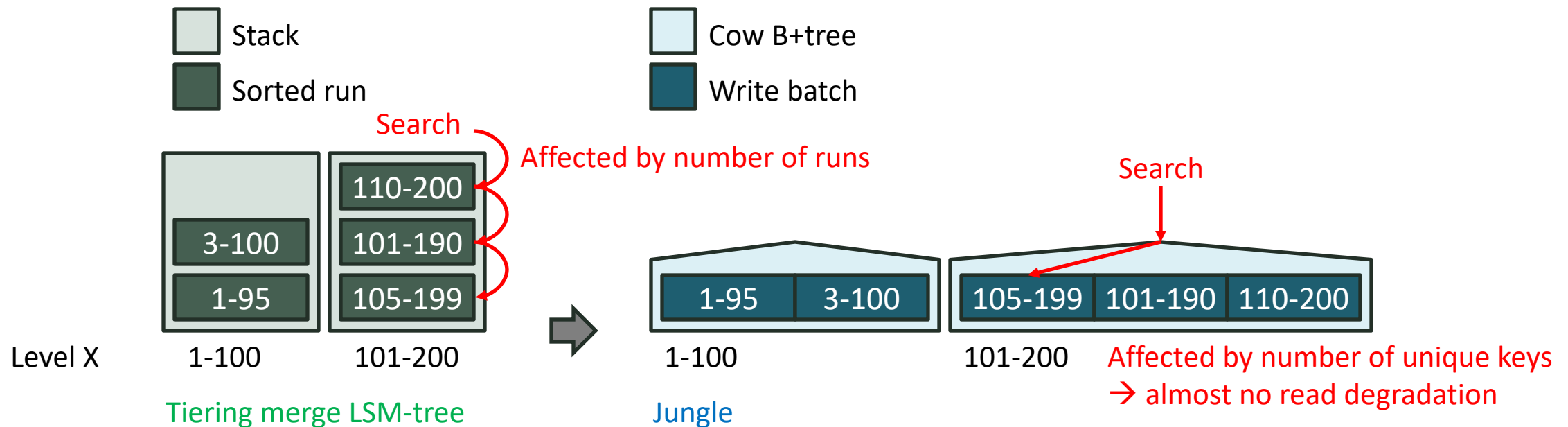
Jungle: Replacing Stack with CoW B+Tree

- Similarities between **tiering stack** vs. **cow B+tree**
 - **Stack** → **cow B+tree**
 - **Sorted run** → **write batch** (if locally sorted in key order)
 - Immutable: no overwrite
 - Appended in chronological order



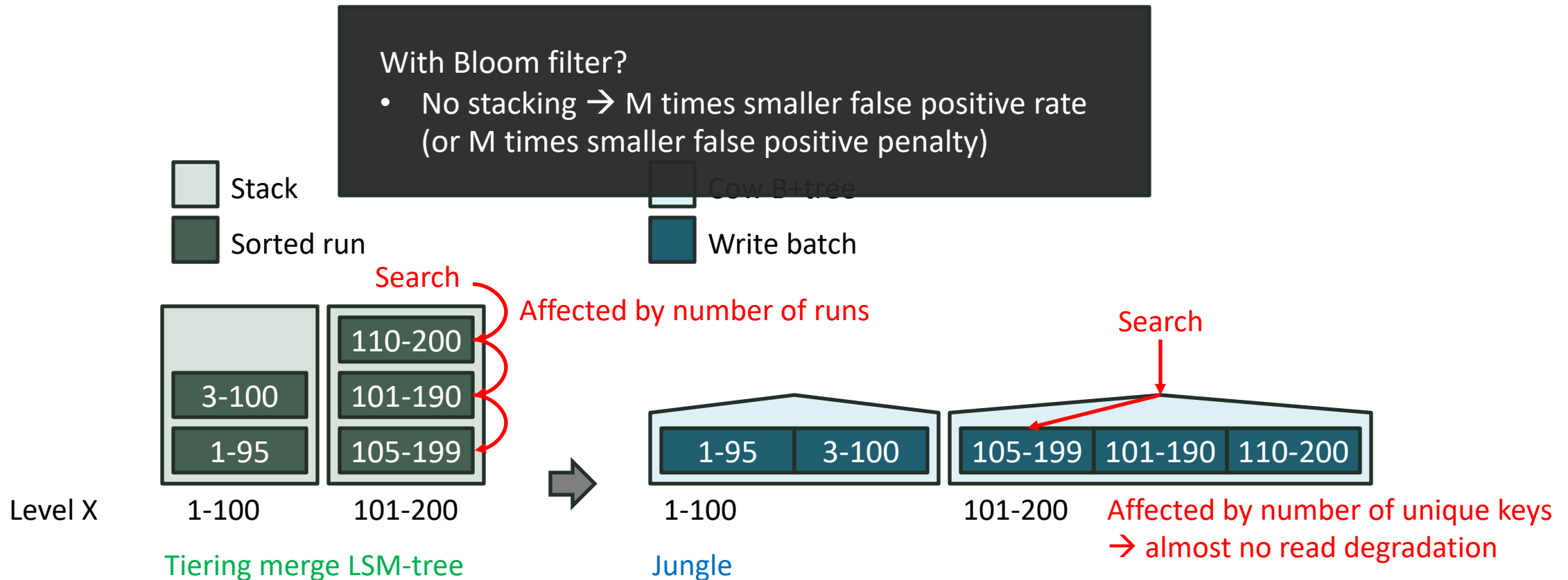
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)



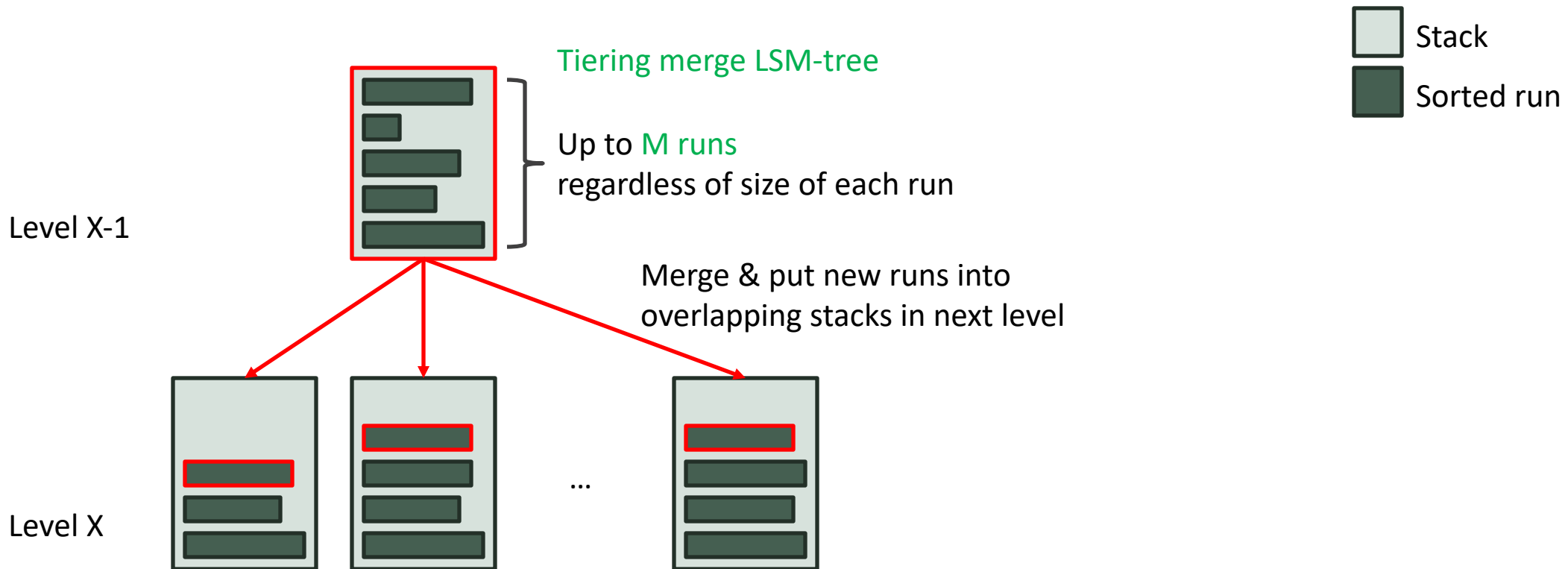
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)



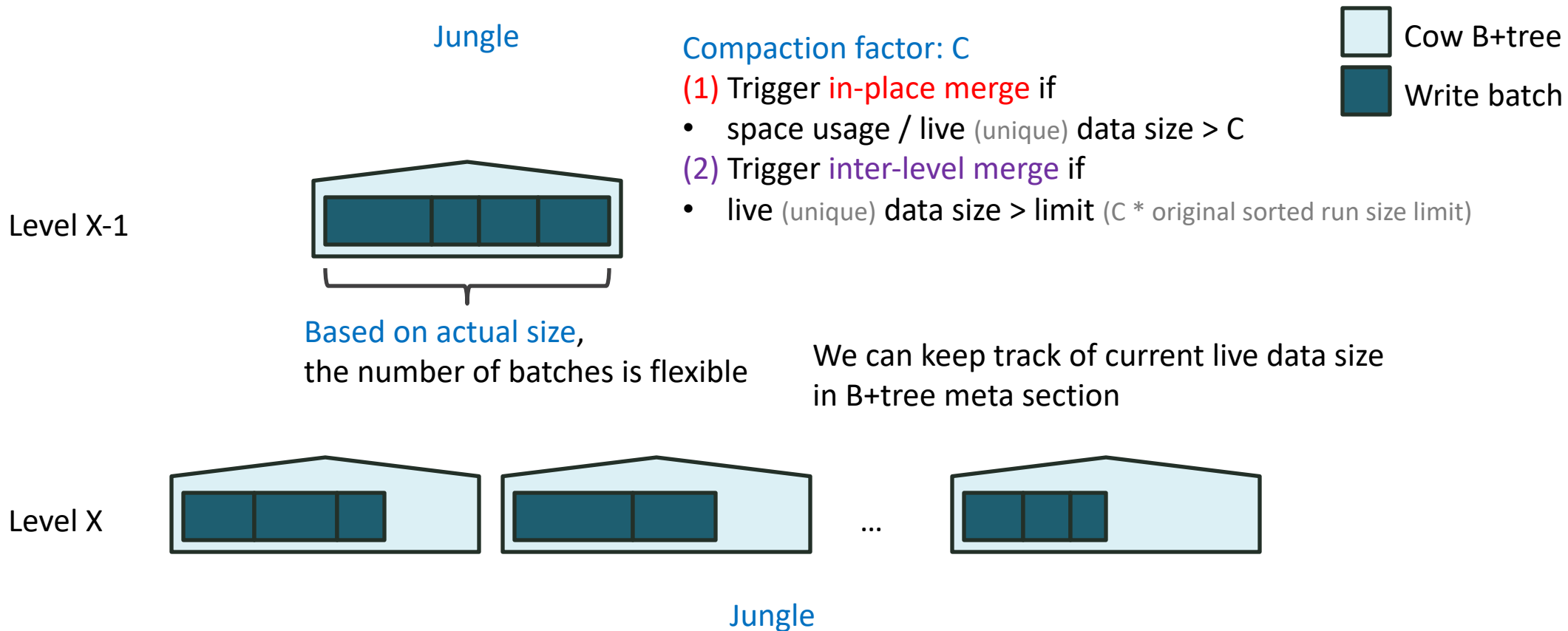
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



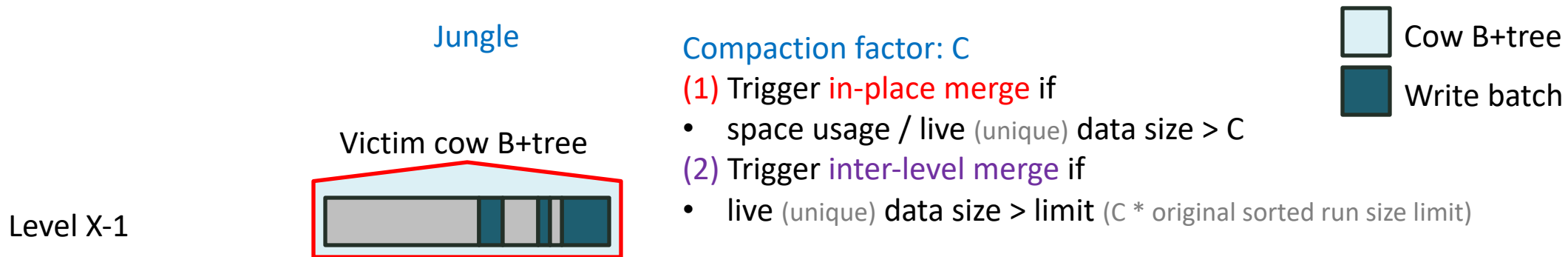
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



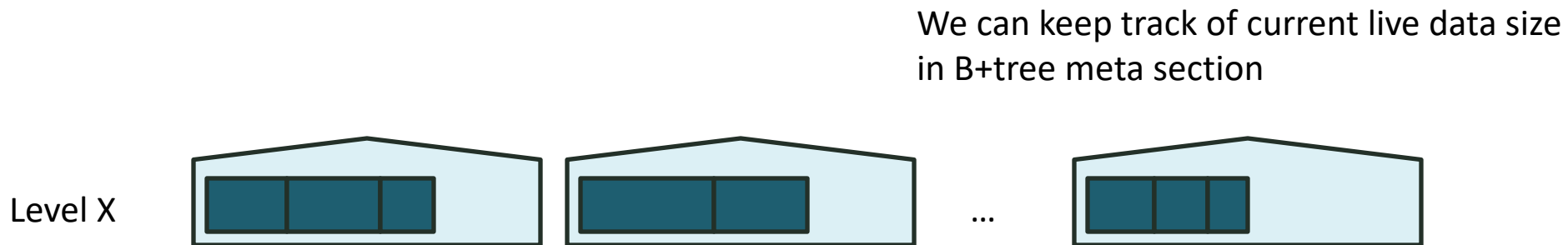
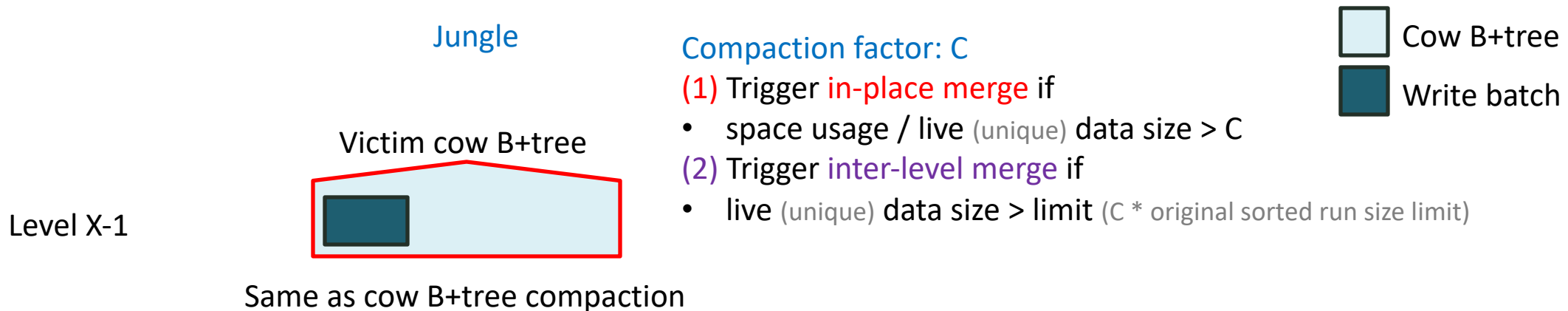
(1) Not much unique data \rightarrow **in-place merge**

We can keep track of current live data size in B+tree meta section



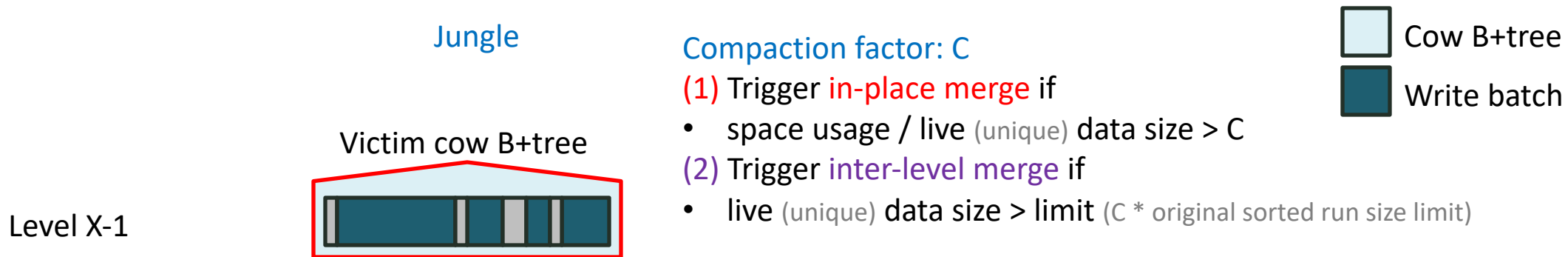
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



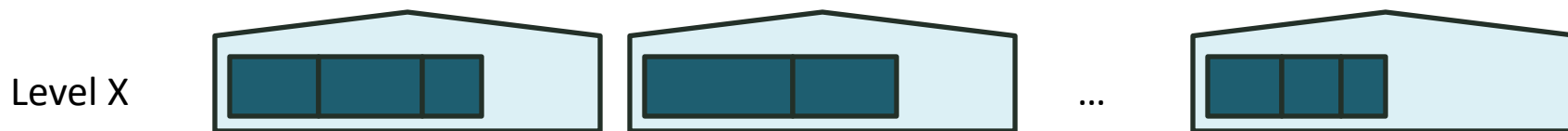
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



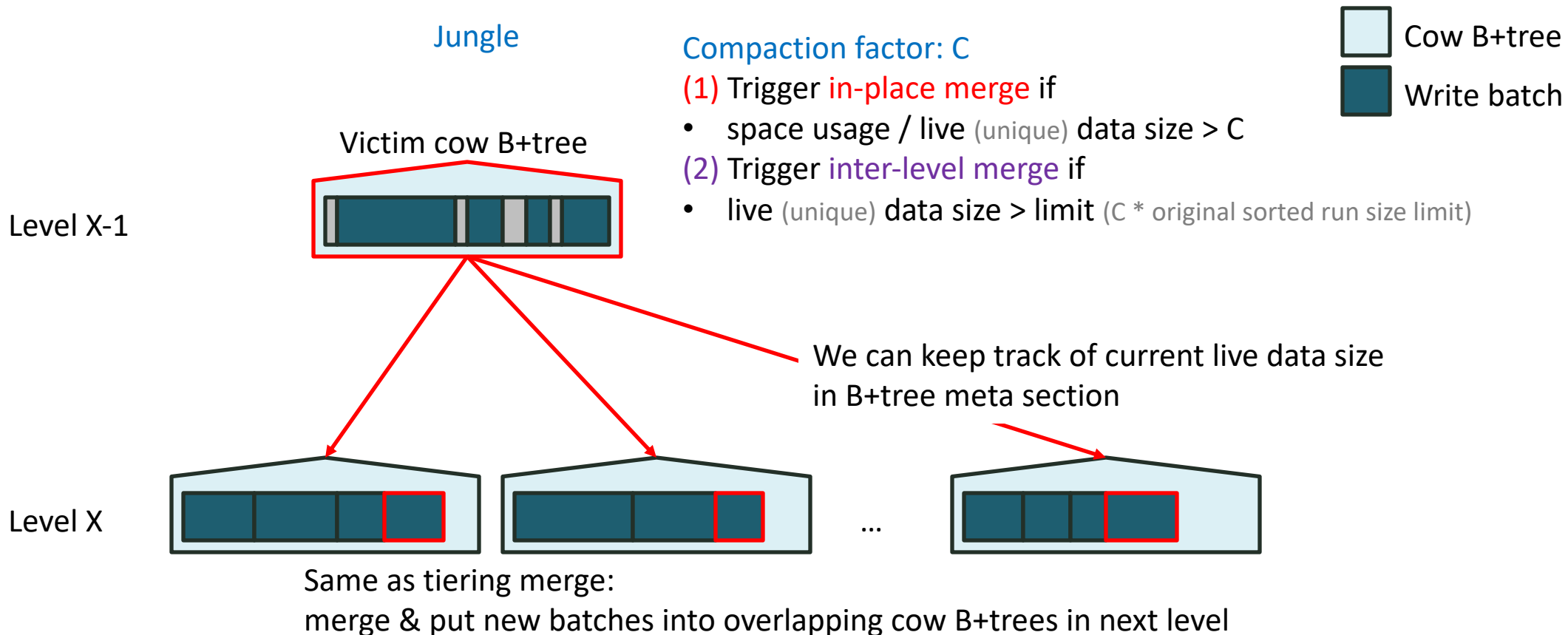
(2) Almost all data are unique → **inter-level merge**

We can keep track of current live data size in B+tree meta section



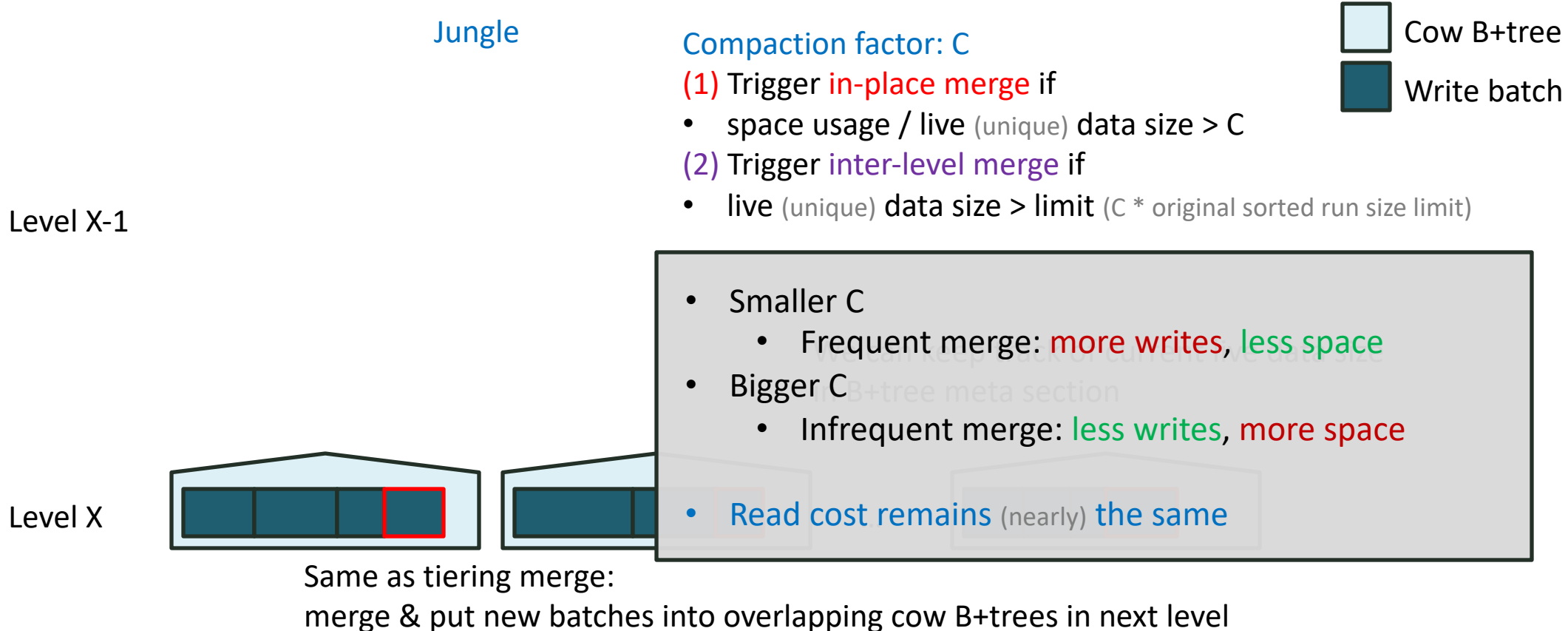
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



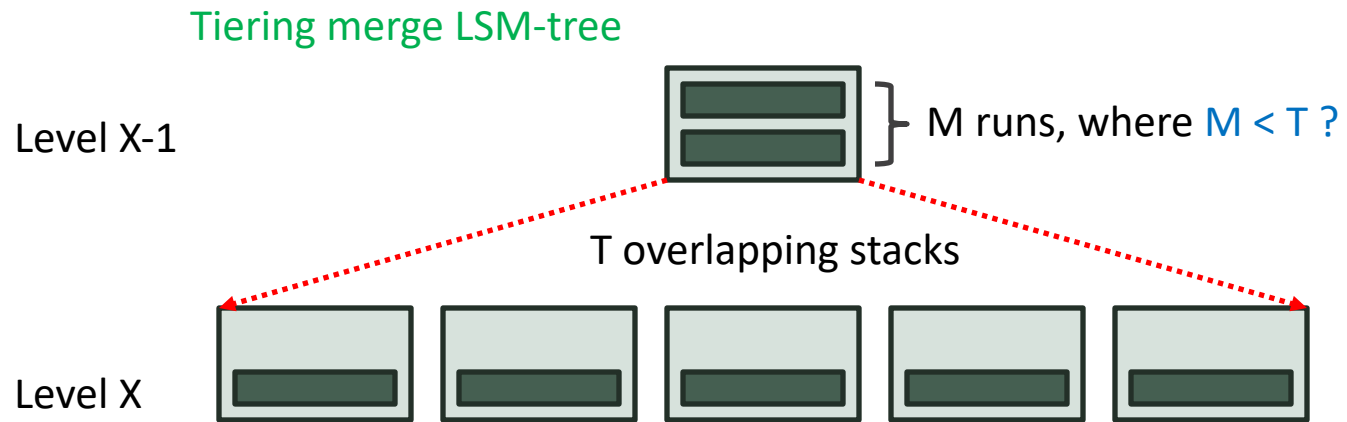
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



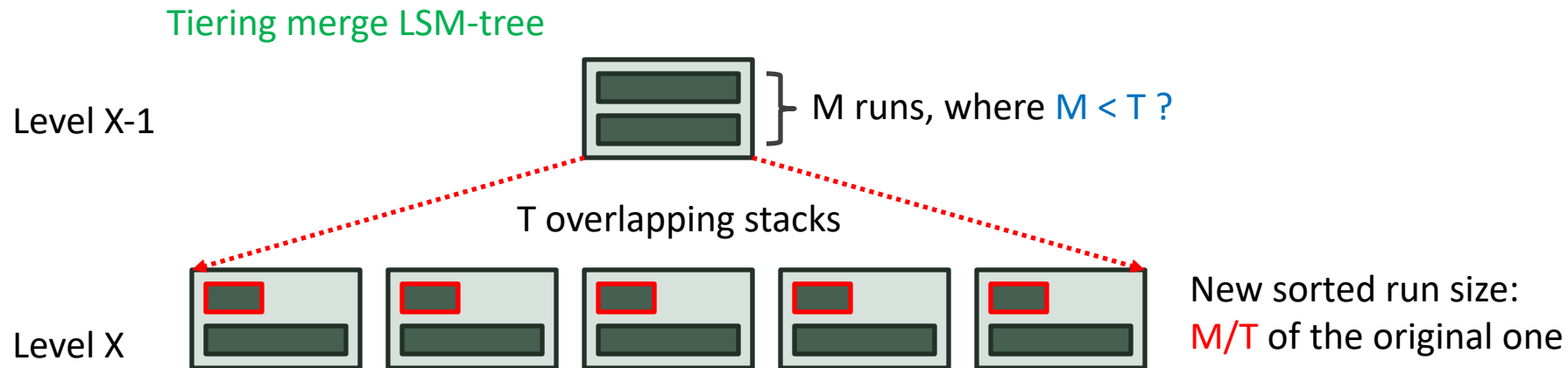
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



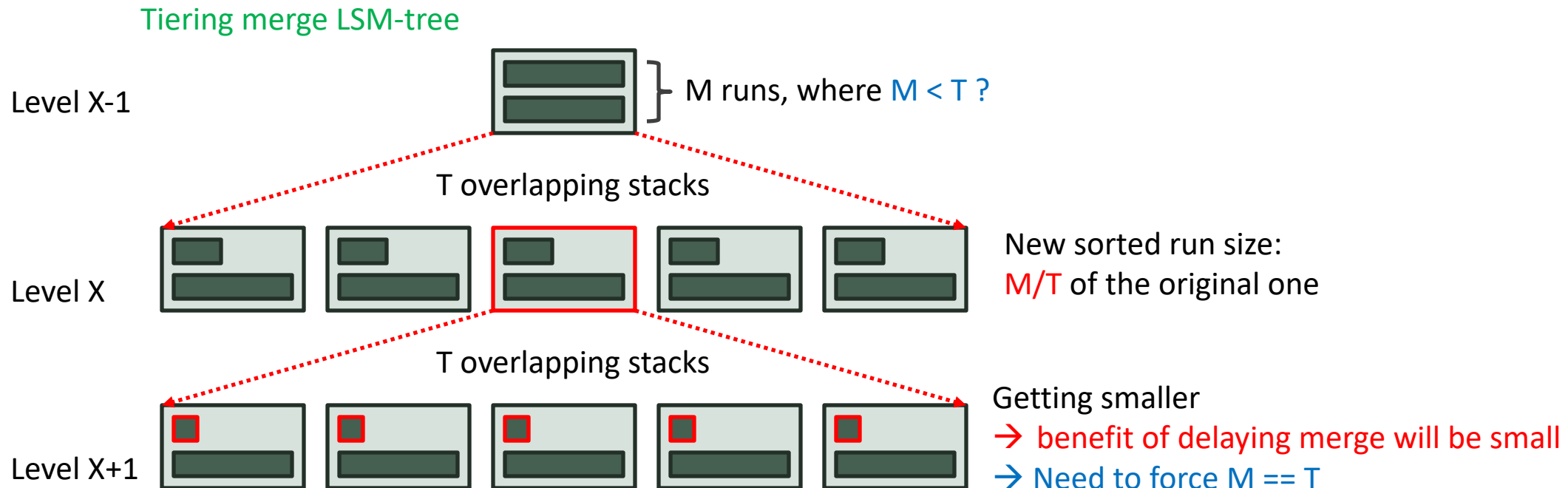
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



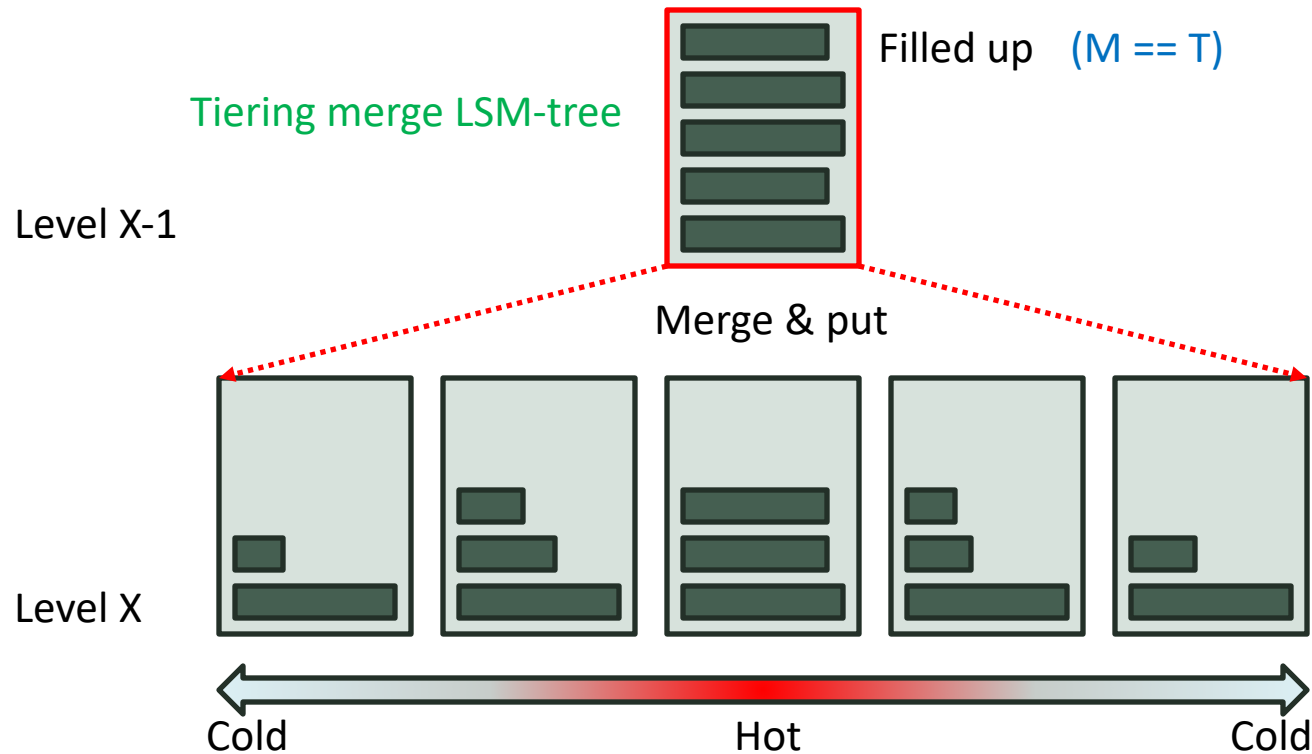
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)



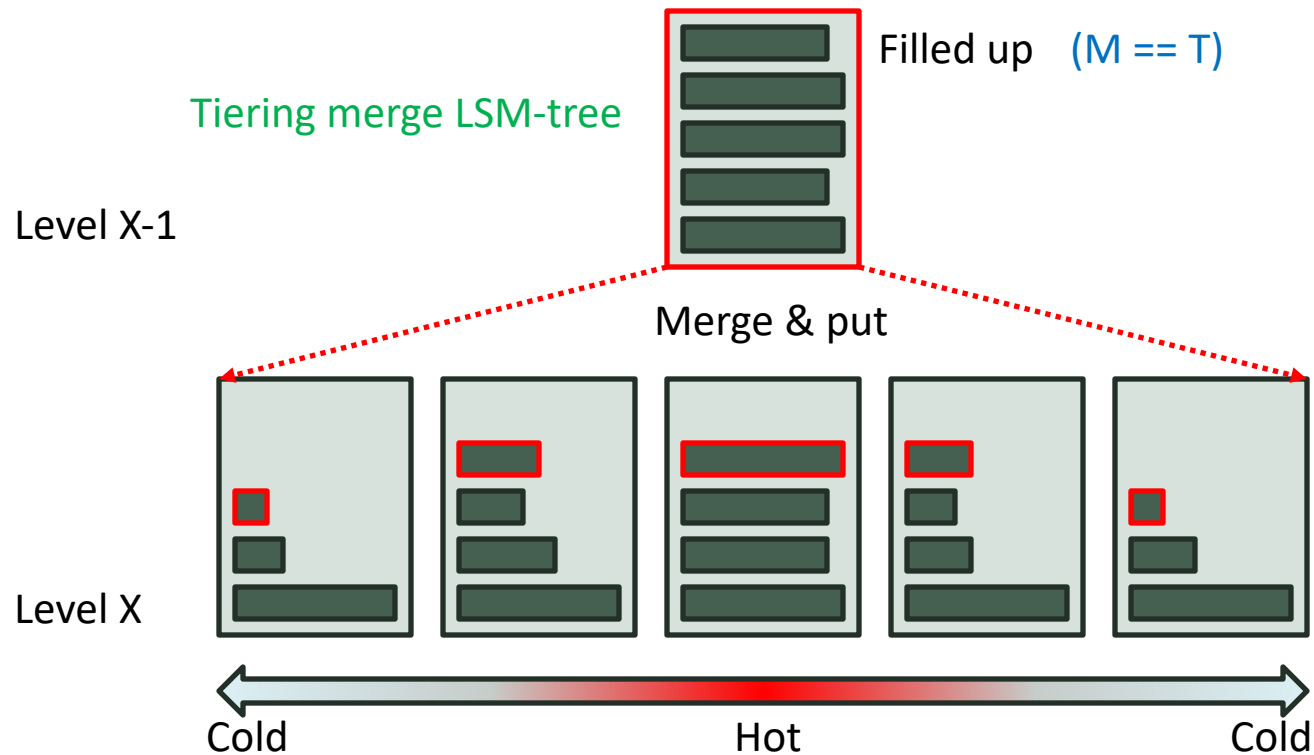
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)
 - Locality?



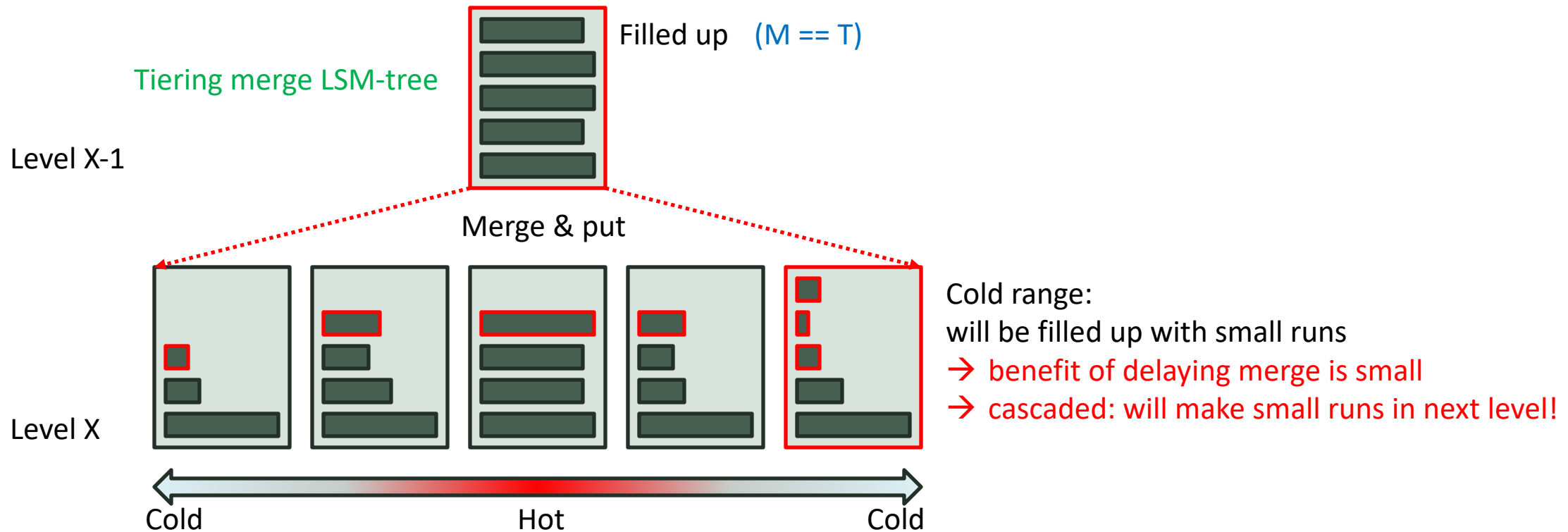
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)
 - Locality?



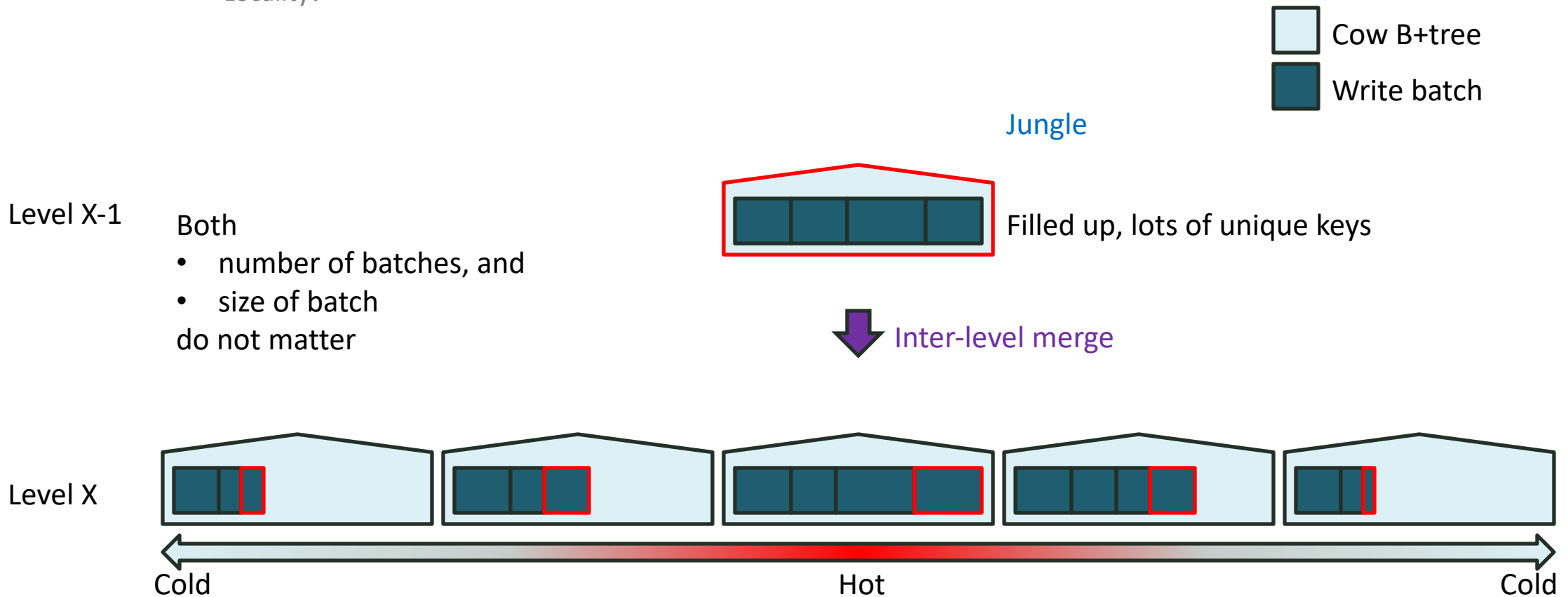
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)
 - Locality?



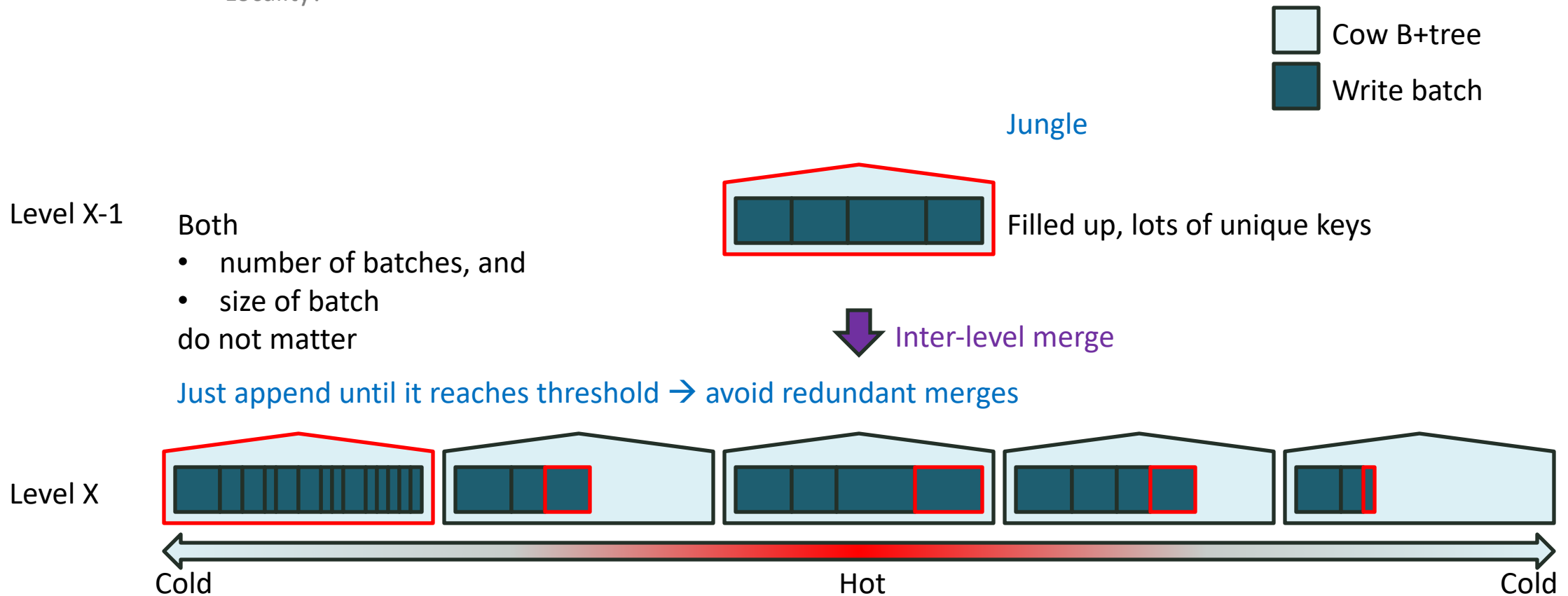
Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)
 - Locality?



Jungle: Replacing Stack with CoW B+Tree

- Differences between **tiering stack** vs. **cow B+tree**
 - Search: **linear** (stack) vs. **logarithmic** (cow B+tree)
 - Unit of limit: **the number of runs** (stack) vs. **actual data size** (cow B+tree)
 - Locality?



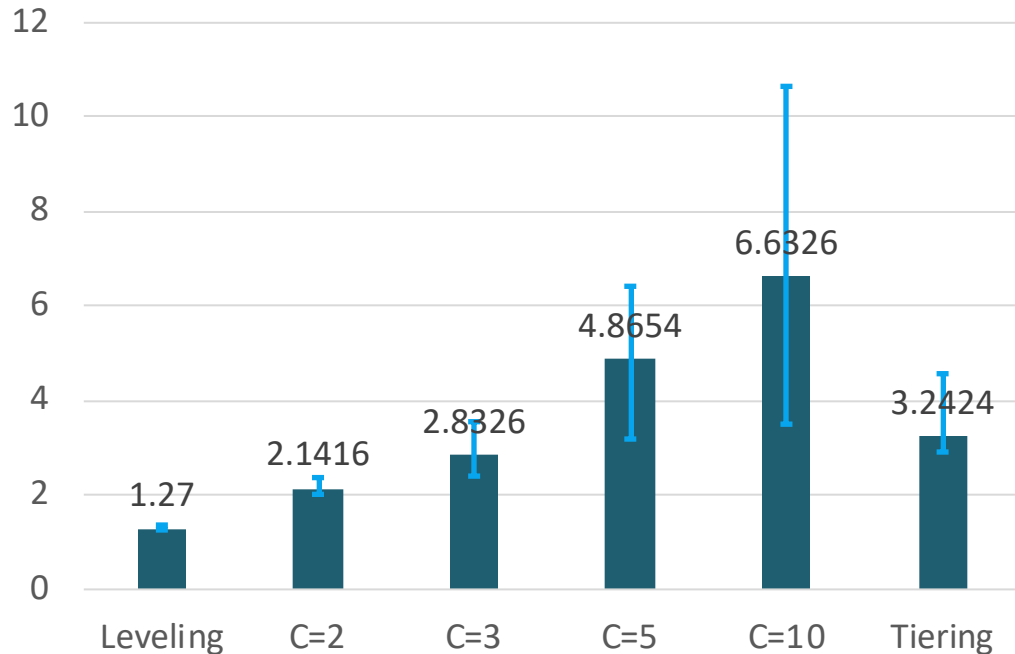
Brief Evaluation

- To prove the concept of Jungle
 - Comparison against leveling (original LSM-tree) and tiering
 - Widely used LSM-based approaches: leveling
- Environment
 - Samsung 860 QVO 1TB, Ext4
 - 20M random key-value pairs
 - Key: 8 bytes, value 1024 KB
 - RAM size: limited to 2.5 GB
 - LSM-tree settings
 - Max sorted run size: 64MB, L0 size limit: 256MB
 - Size ratio between levels (T): 10
 - Max stack size of tiering (M): 10 runs (==T)
 - Bloom filter: 10 bits per key (~1% false positive rate)

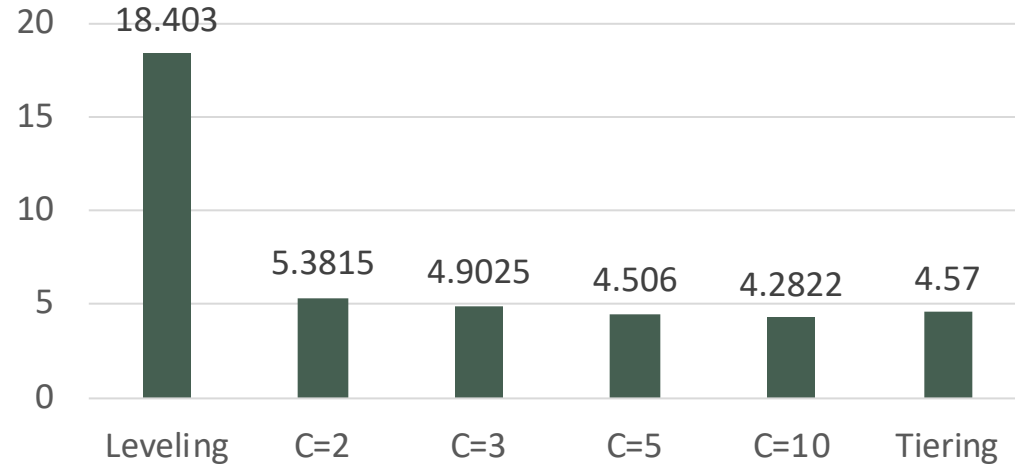
Brief Evaluation

- Write and space amplification
 - Issued 400M uniform random updates (no locality → best case for tiering)

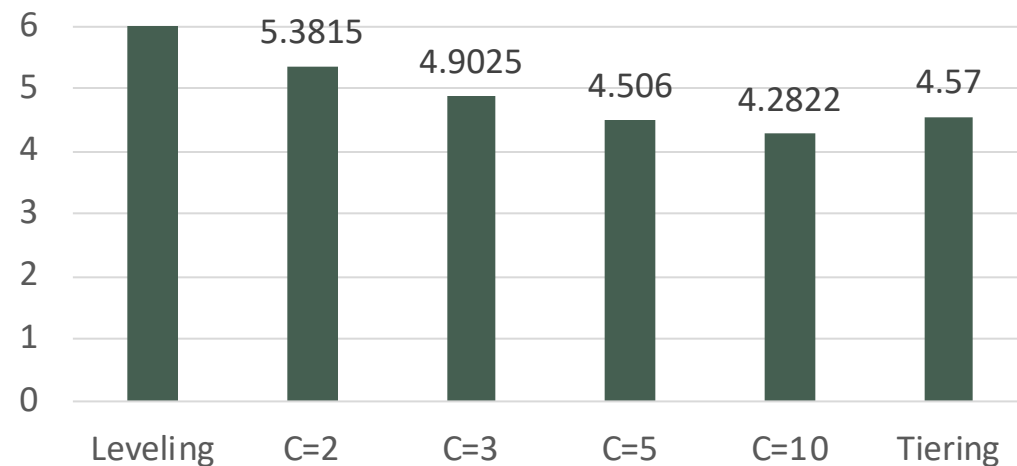
Space amplification



Write amplification



Write amplification (zoomed in)

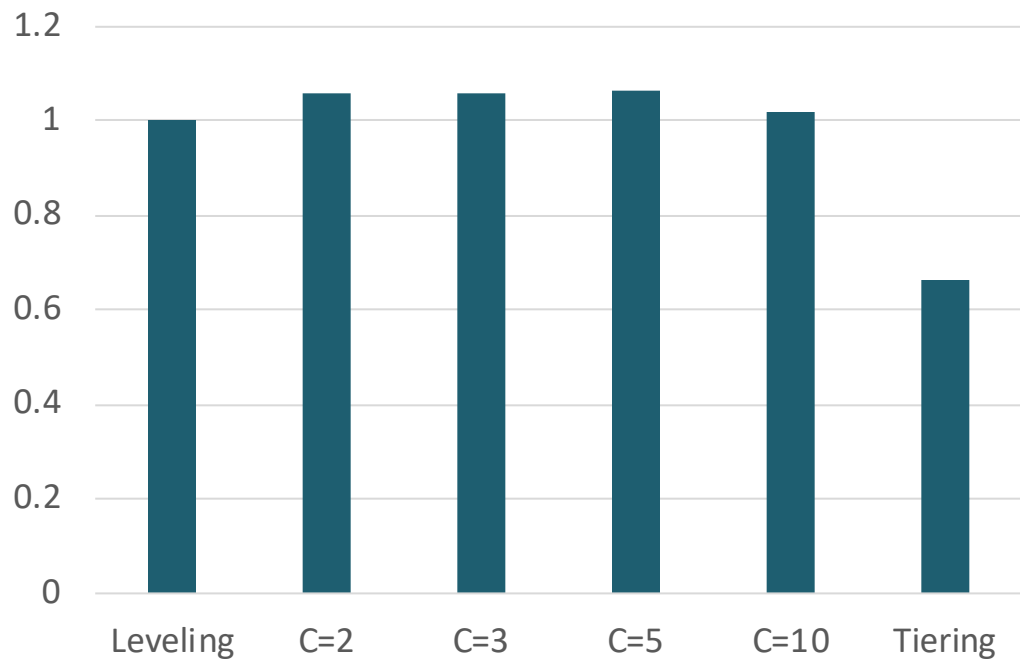


Jungle with different parameters

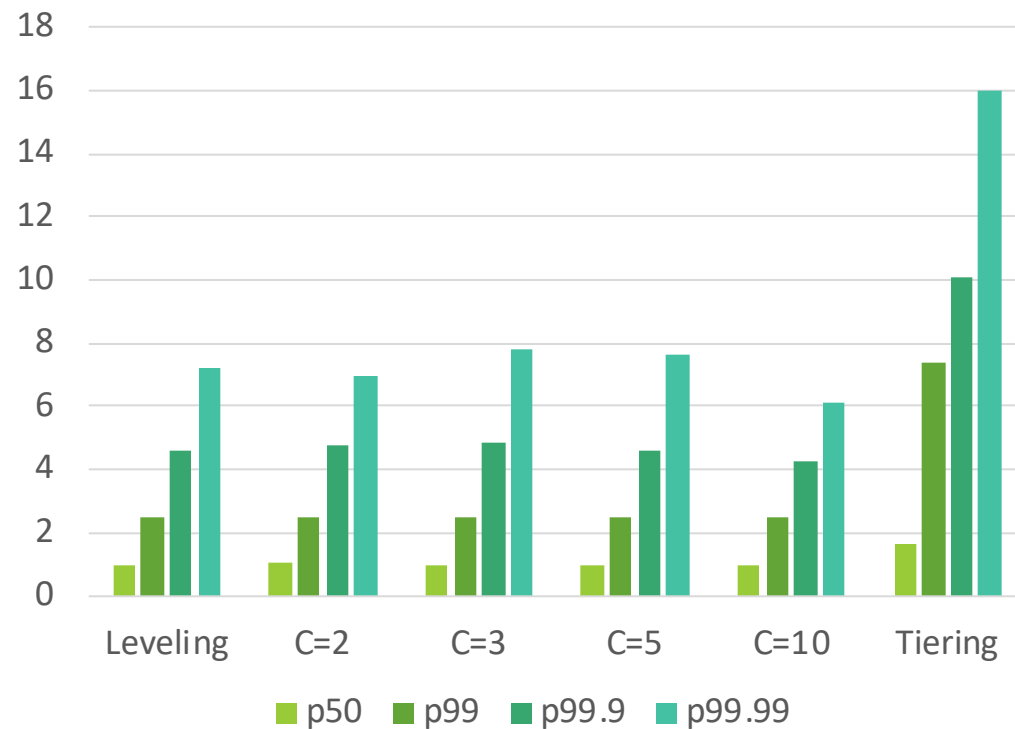
Brief Evaluation

- Point read throughput and latencies
 - Issued random read operation on aged indexes

Normalized read throughput

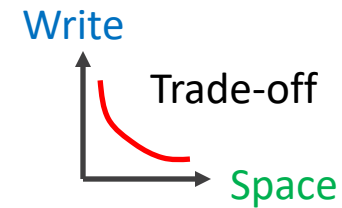
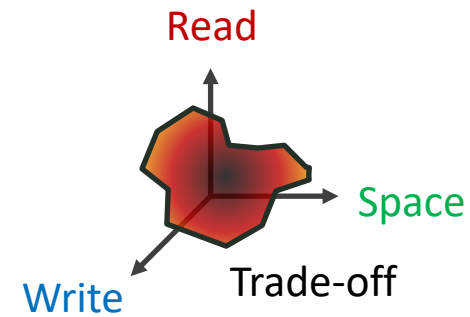


Normalized read latencies



Summary

- Traditional LSM-tree trade-offs
 - Read, write, and space
- Jungle
 - Transplant copy-on-write B+tree into LSM-tree
 - Get rid of read cost from the trade-off
 - Introduce a practical parameter to adjust write and space
- What's next?
 - Fundamental ways to re-think LSM-tree structure
 - More chances to deform or optimize



Thank You