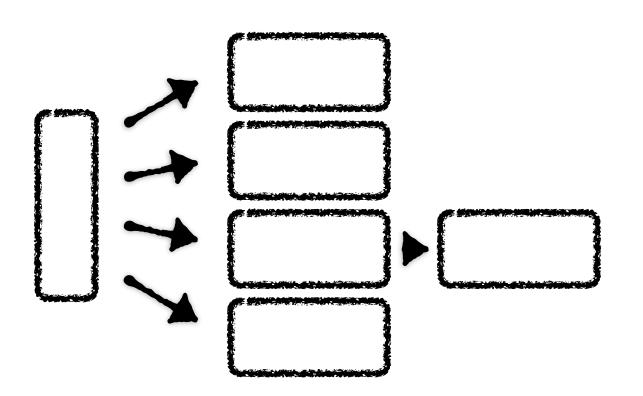
Circular Logs & Cuckoo Filters

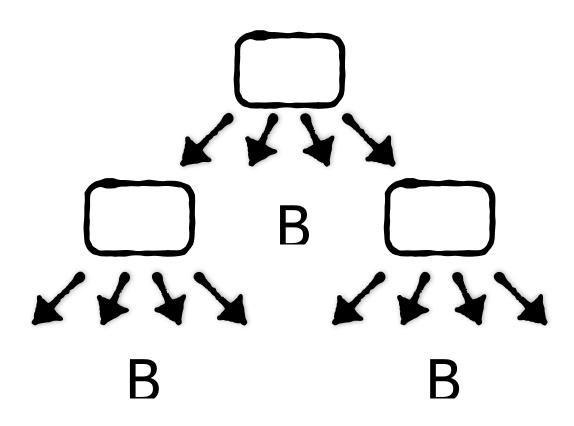
CSC443H1 Database System Technology

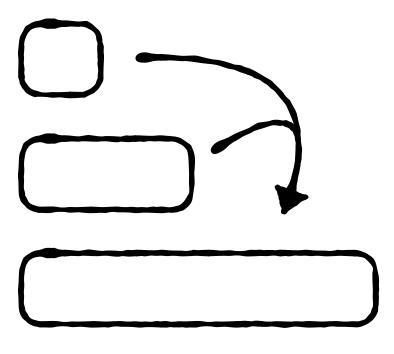
Extendible Hashing



LSM-tree





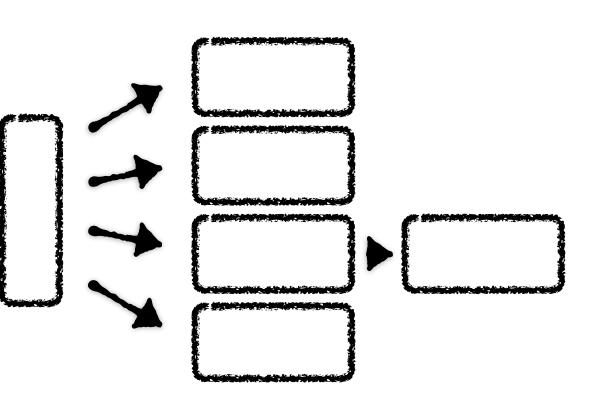


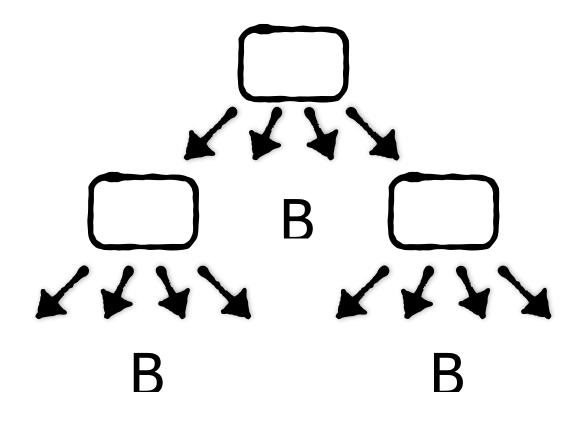
Extendible Hashing

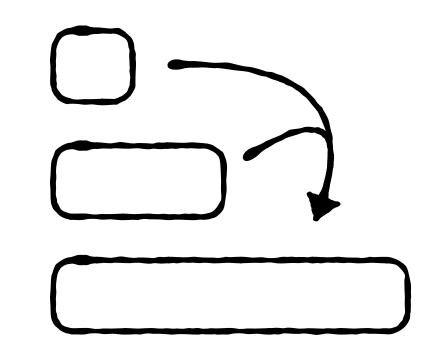
B-Trees

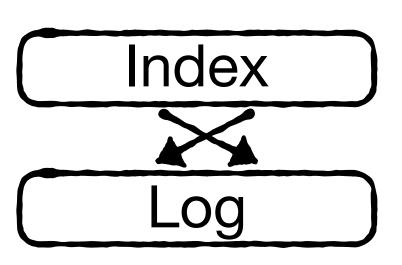
LSM-tree

Circular Log









Cheapest gets
No scans

Cheap gets & scans

Cheaper writes
More memory

Cheapest writes
Most memory
No scans

Circular Log

Invented in 1992 as a "log structured file system"

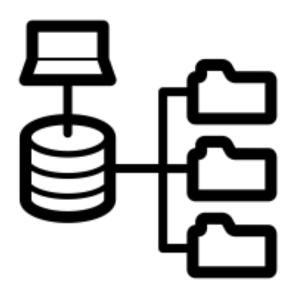
Circular Log

Invented in 1992 as a "log structured file system"

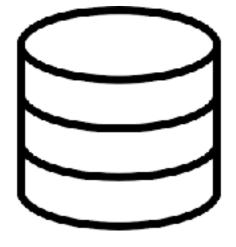
File Systems

Flash Translation Layers

KV-stores







Various names, same data structure:

Index+Log

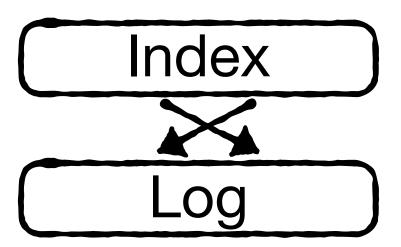
Circular Log

Log-structured Hash Table

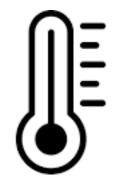
log structured file system

Agenda

Mechanics

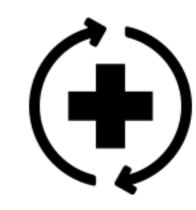


Hot/Cold separation



To reduce write-amp

Checkpointing/recovery



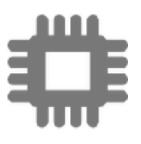
If power fails

Cuckoo filtering



To reduce memory

Mechanics

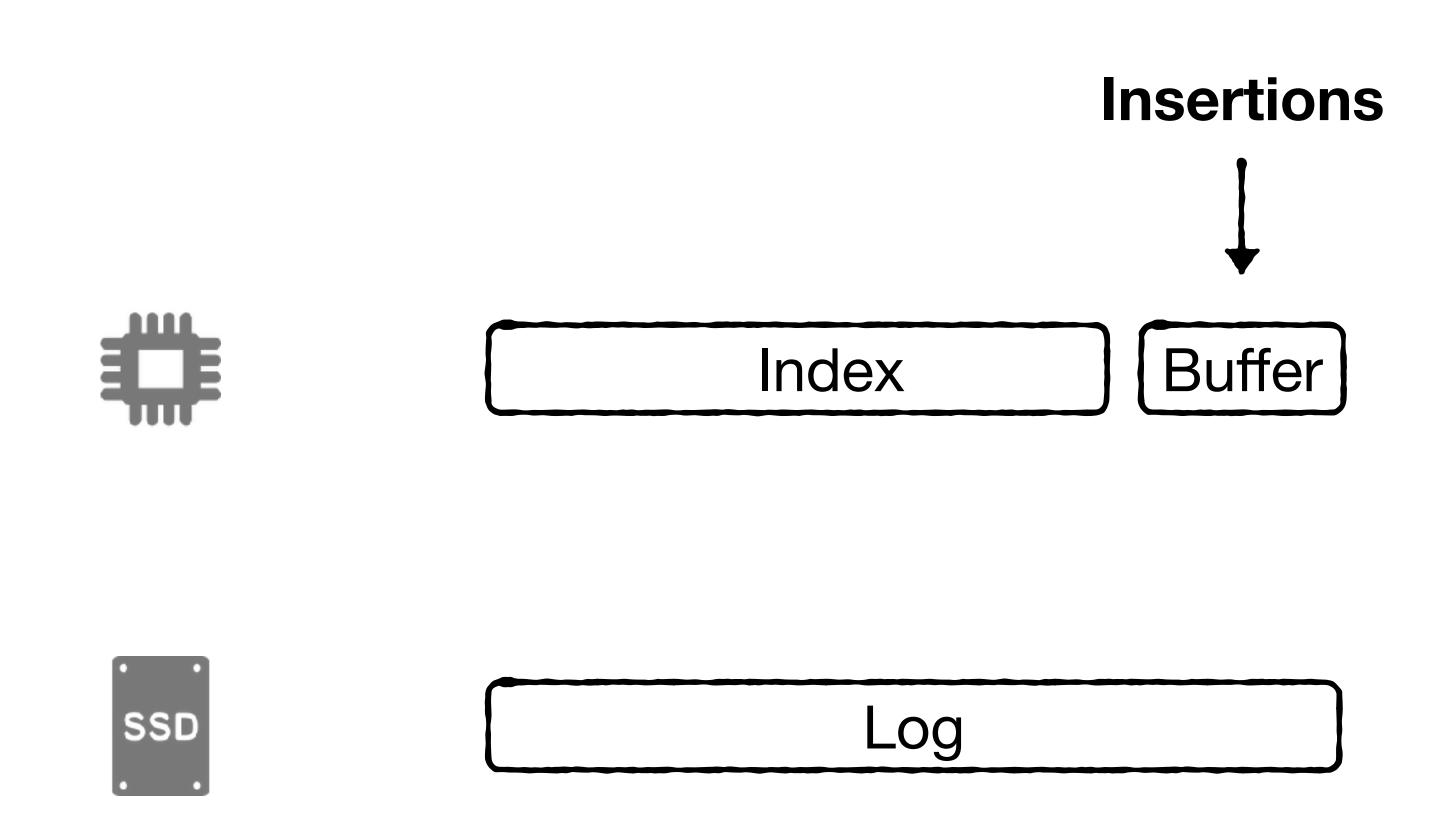


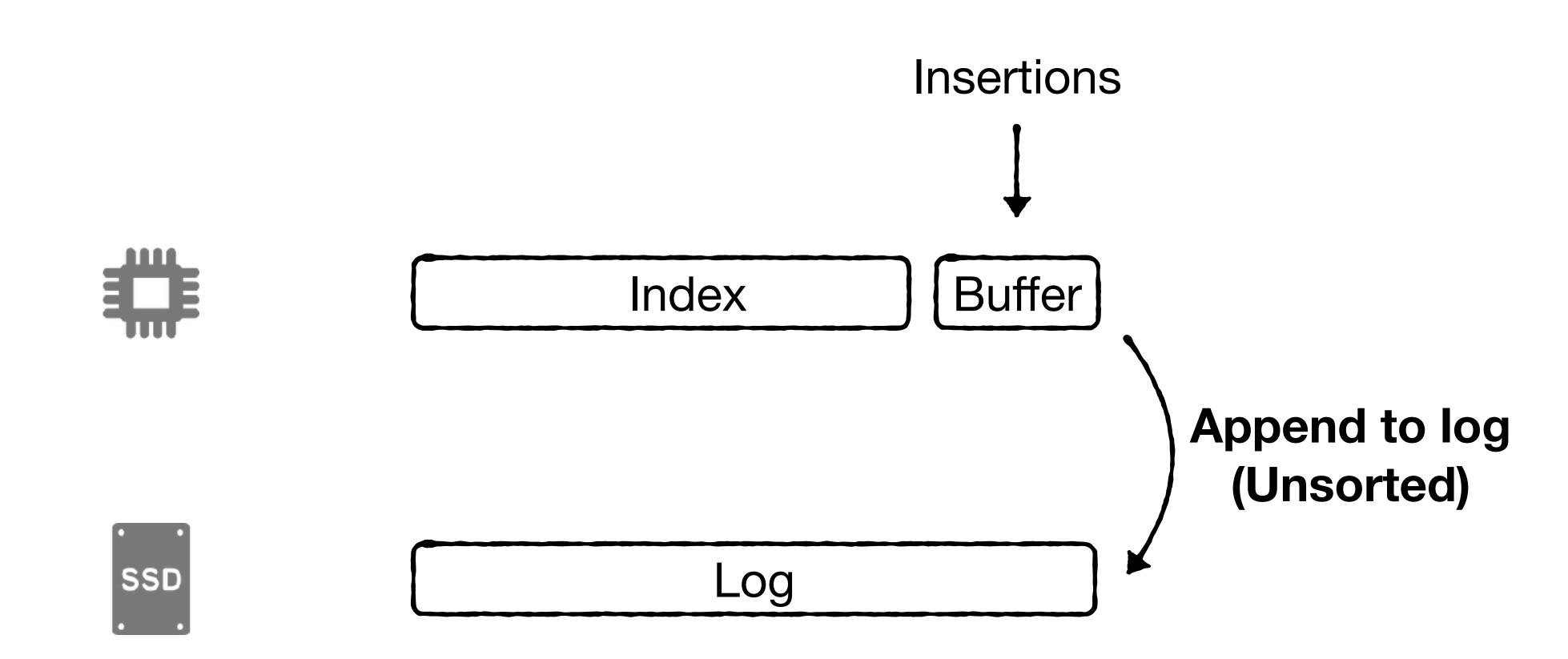
Index

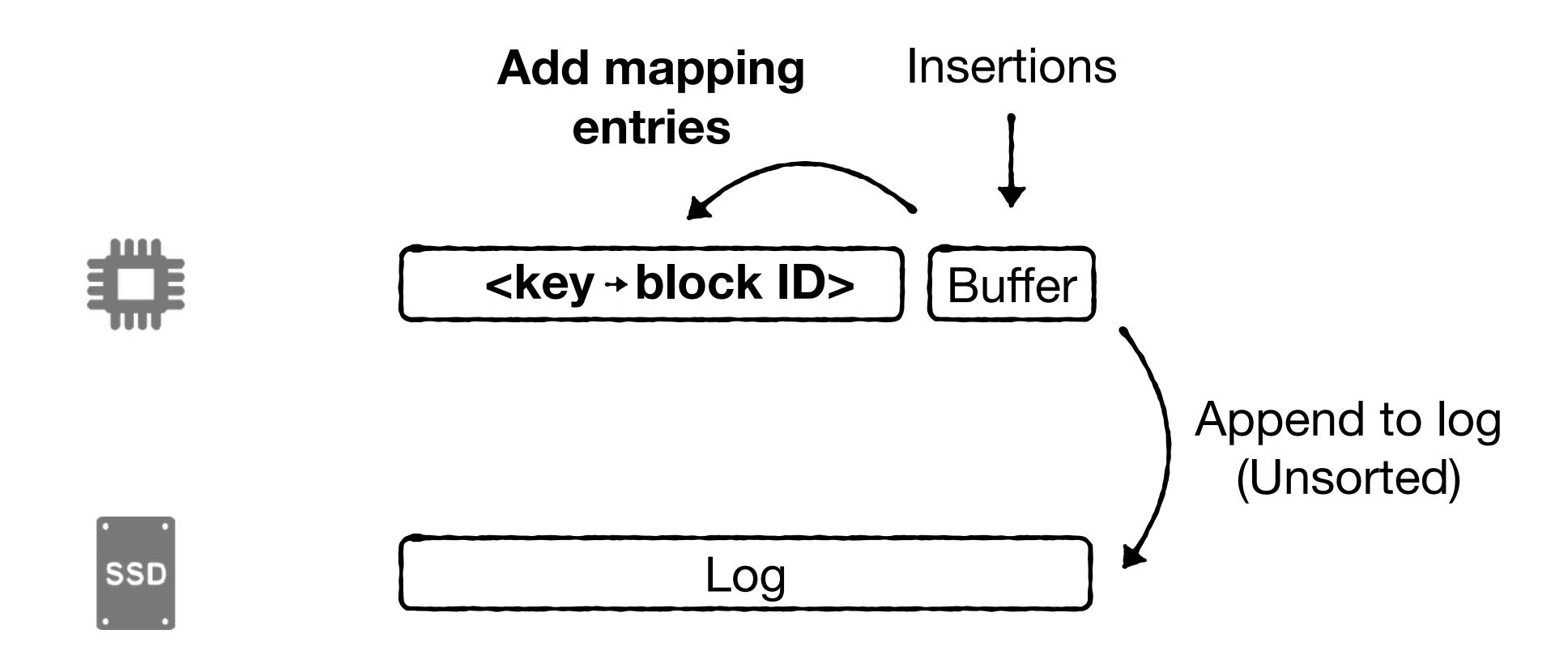
Buffer

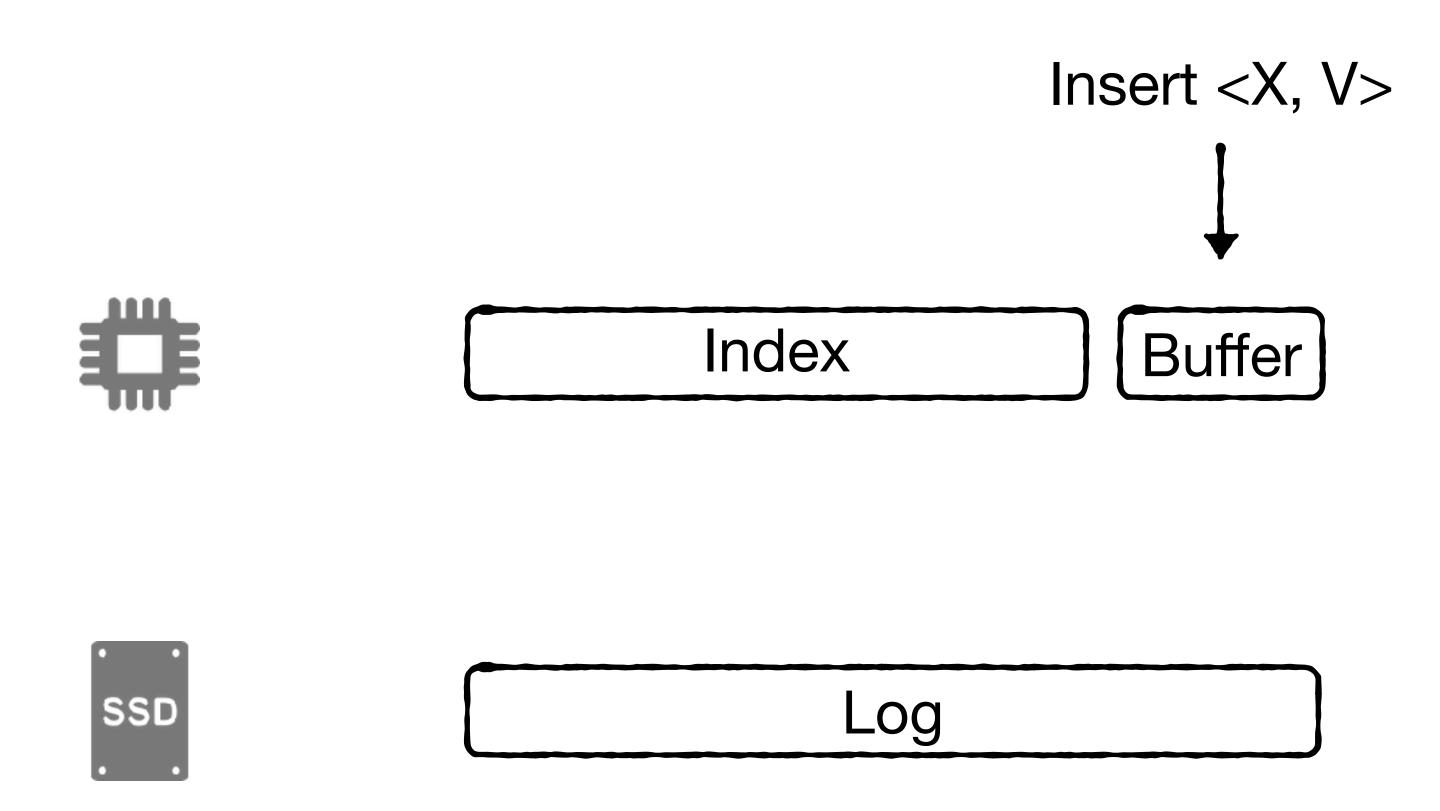


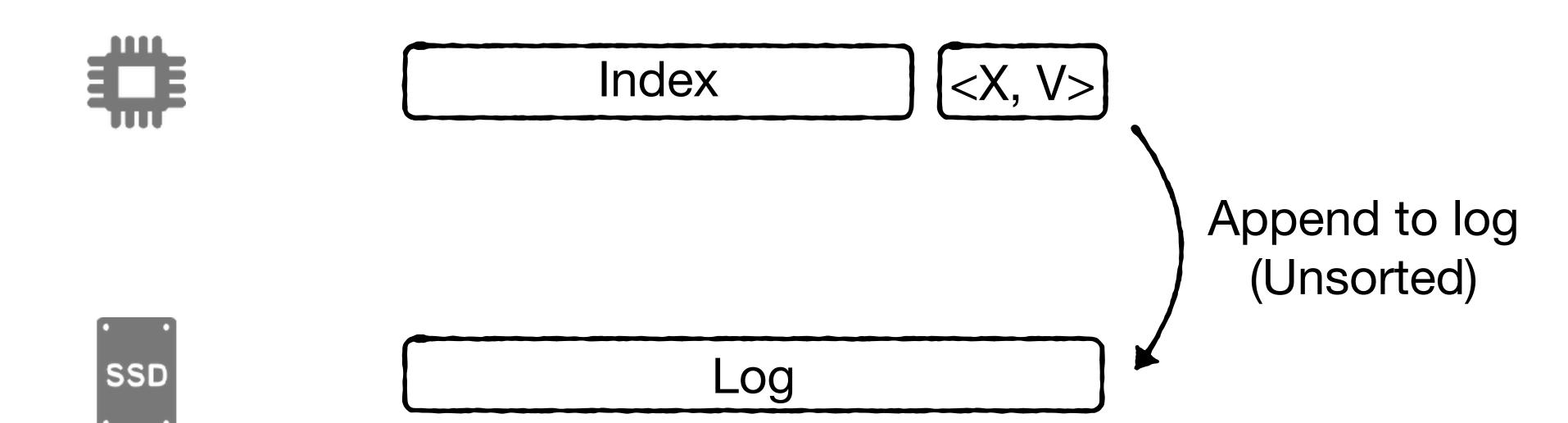
Log



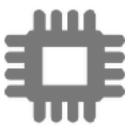


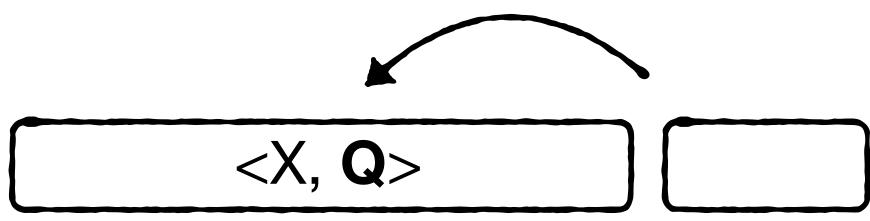






Add mapping entries

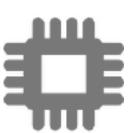


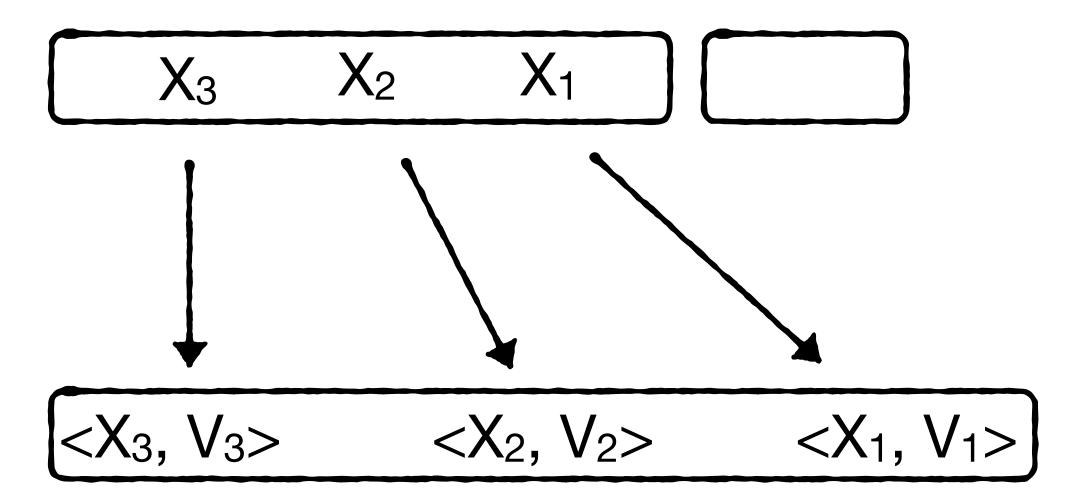






Block Q



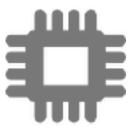




Write cost

(Assuming only insertions)

Get cost



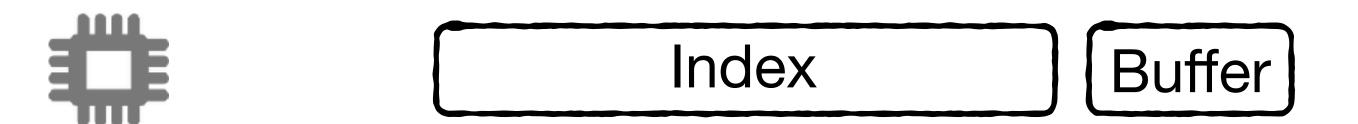
Index

Buffer



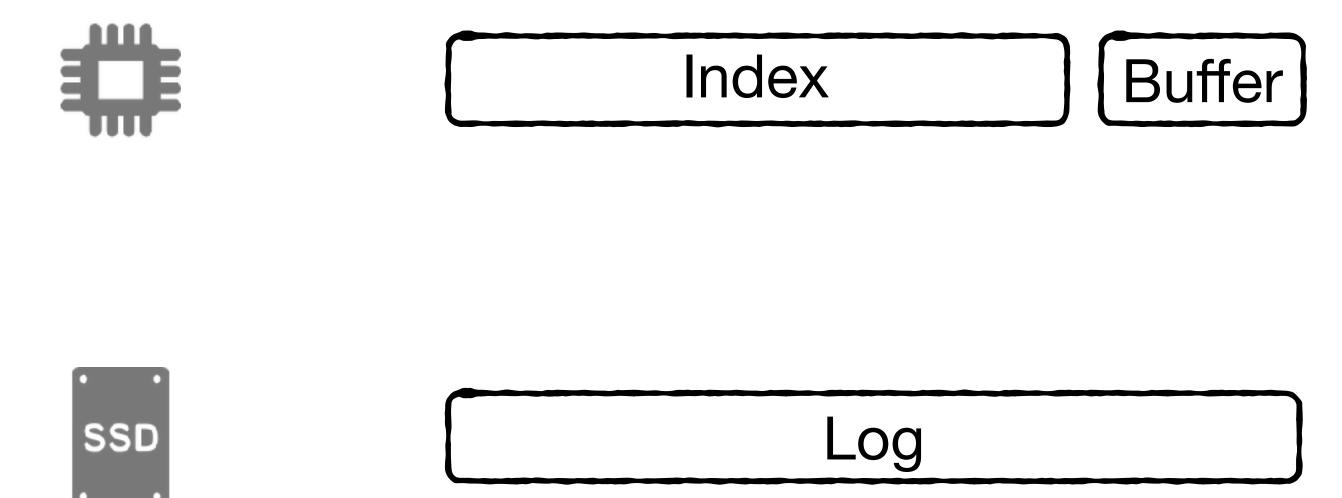
Log

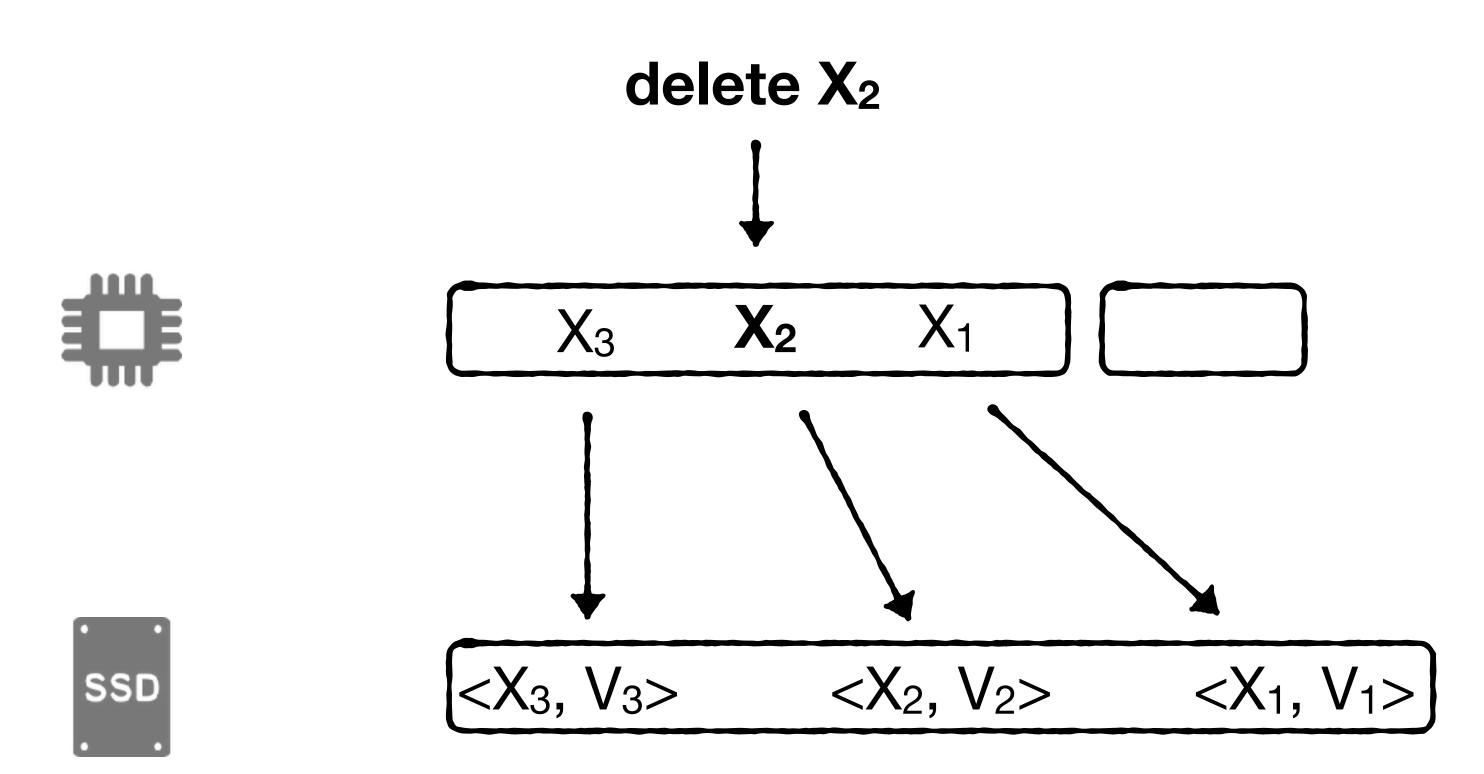
Write cost O(1/B) (Assuming only insertions)
Get cost O(1)

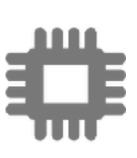


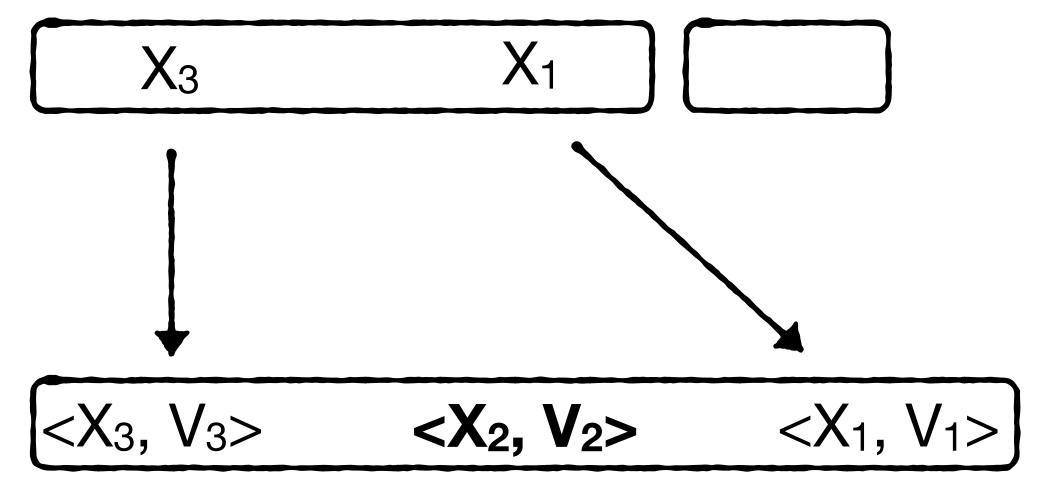
ssD Log

How to delete?





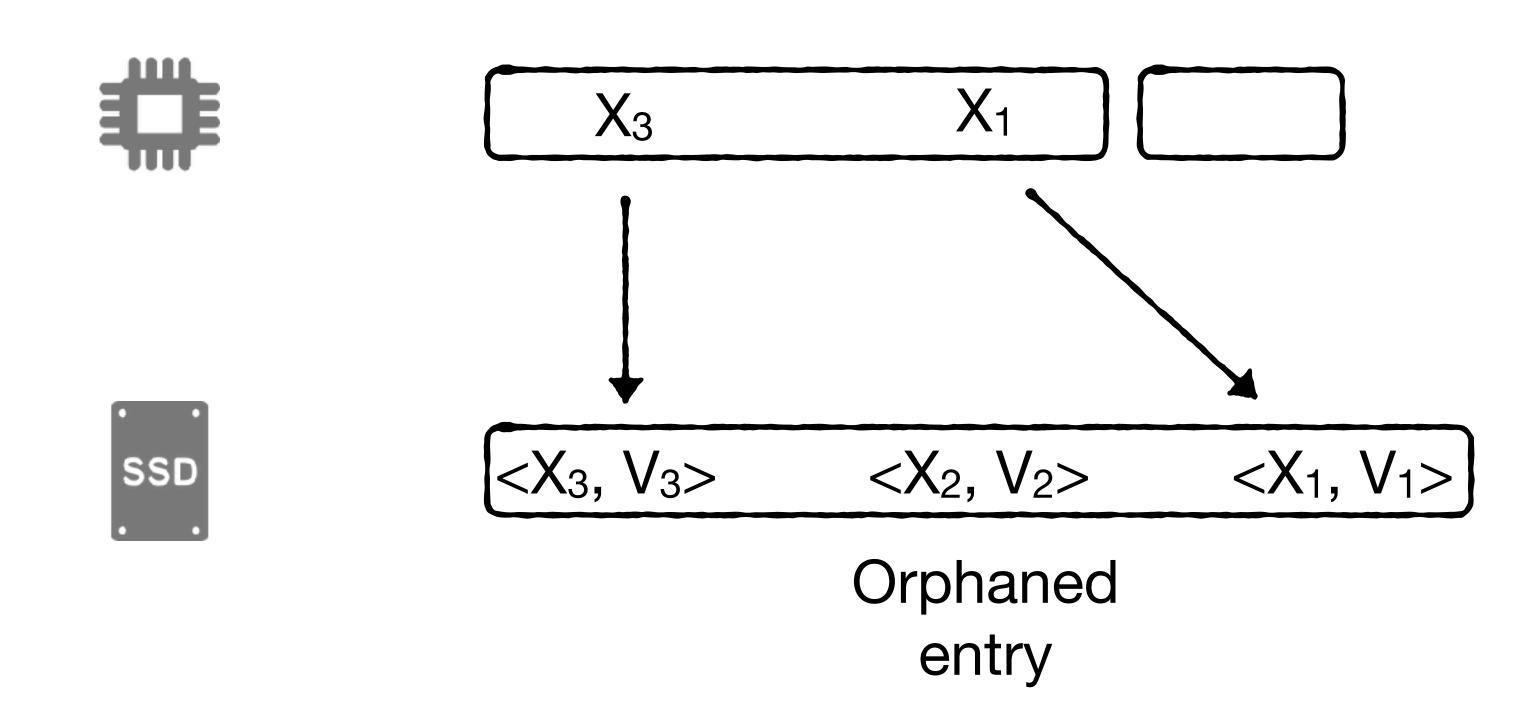






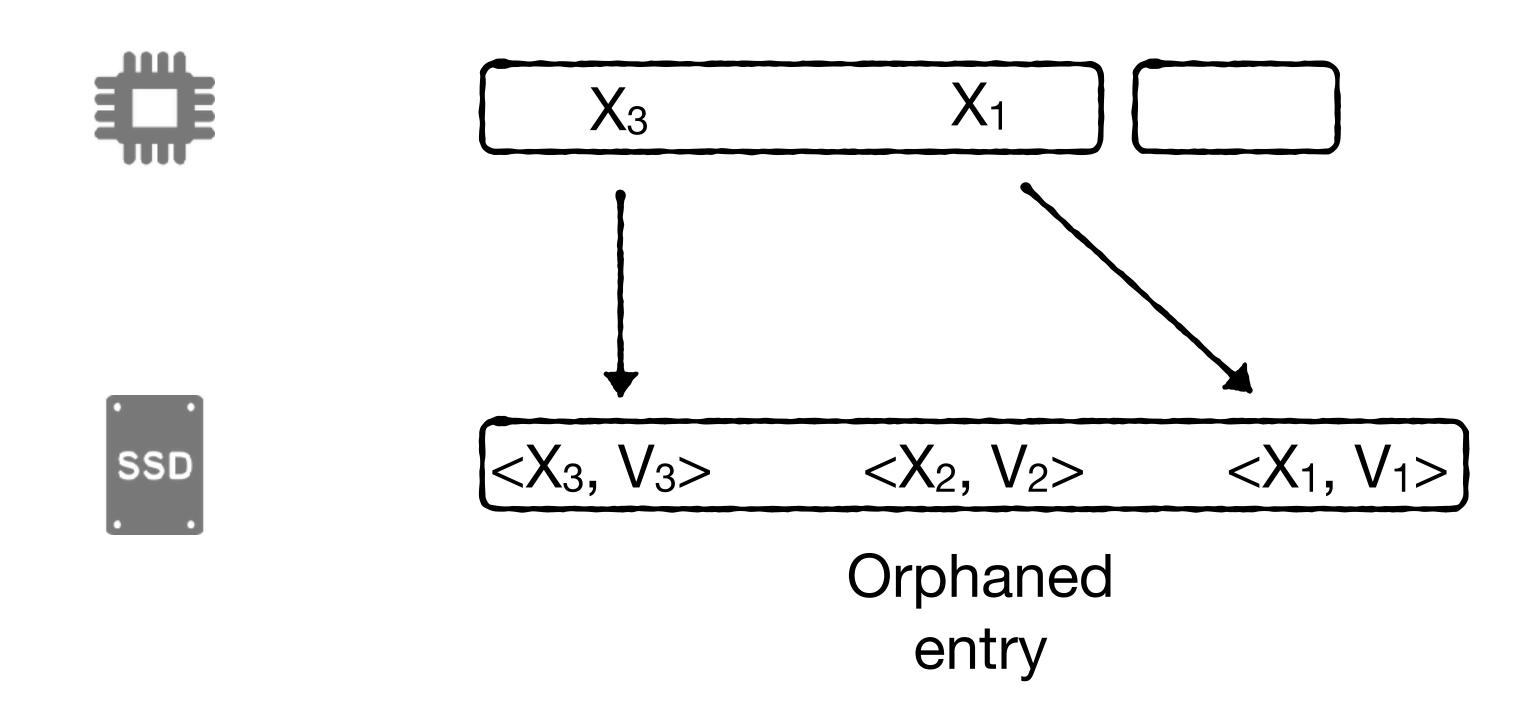
Orphaned entry

After many deletes, many orphaned entries accumulate



After many deletes, many orphaned entries accumulate

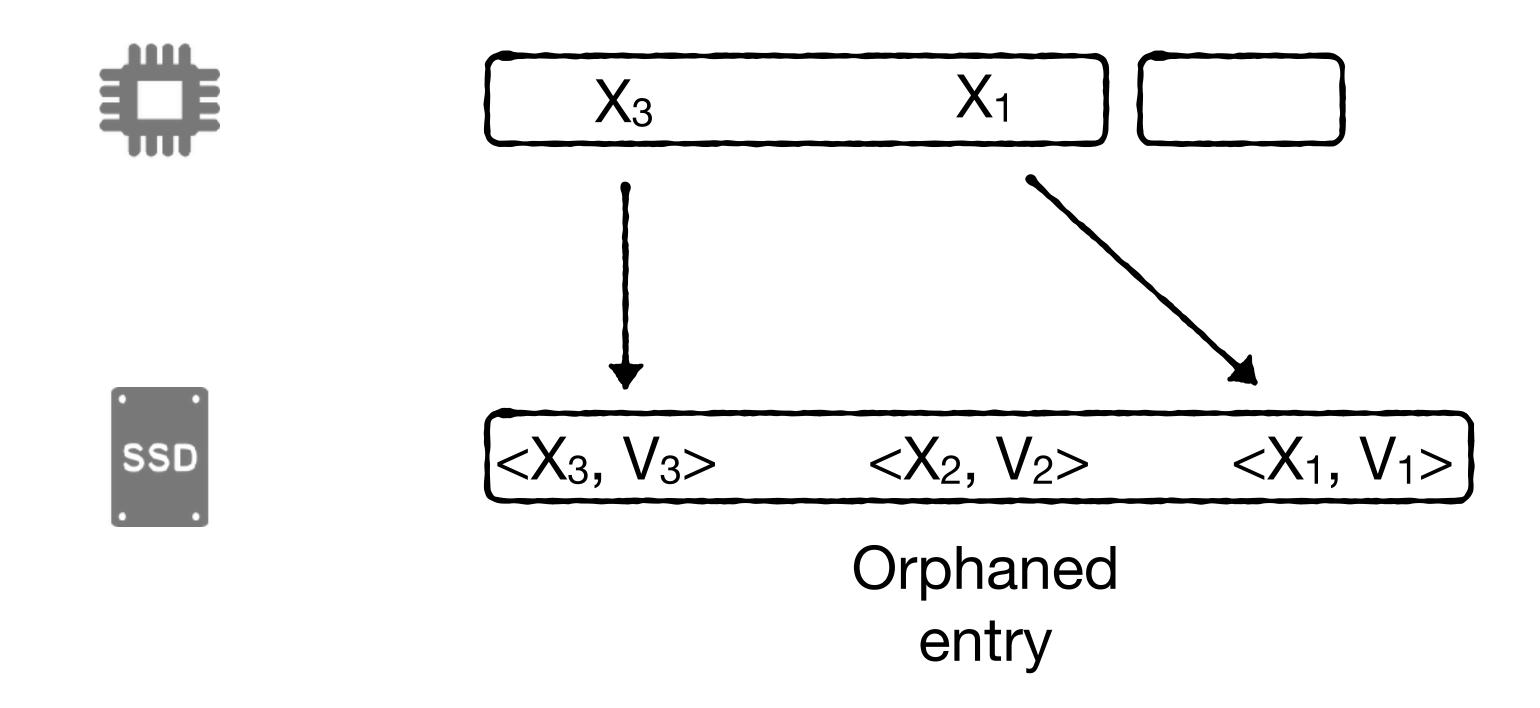
They take up space which we would prefer to use to store valid data



After many deletes, many orphaned entries accumulate

They take up space which we would prefer to use to store valid data

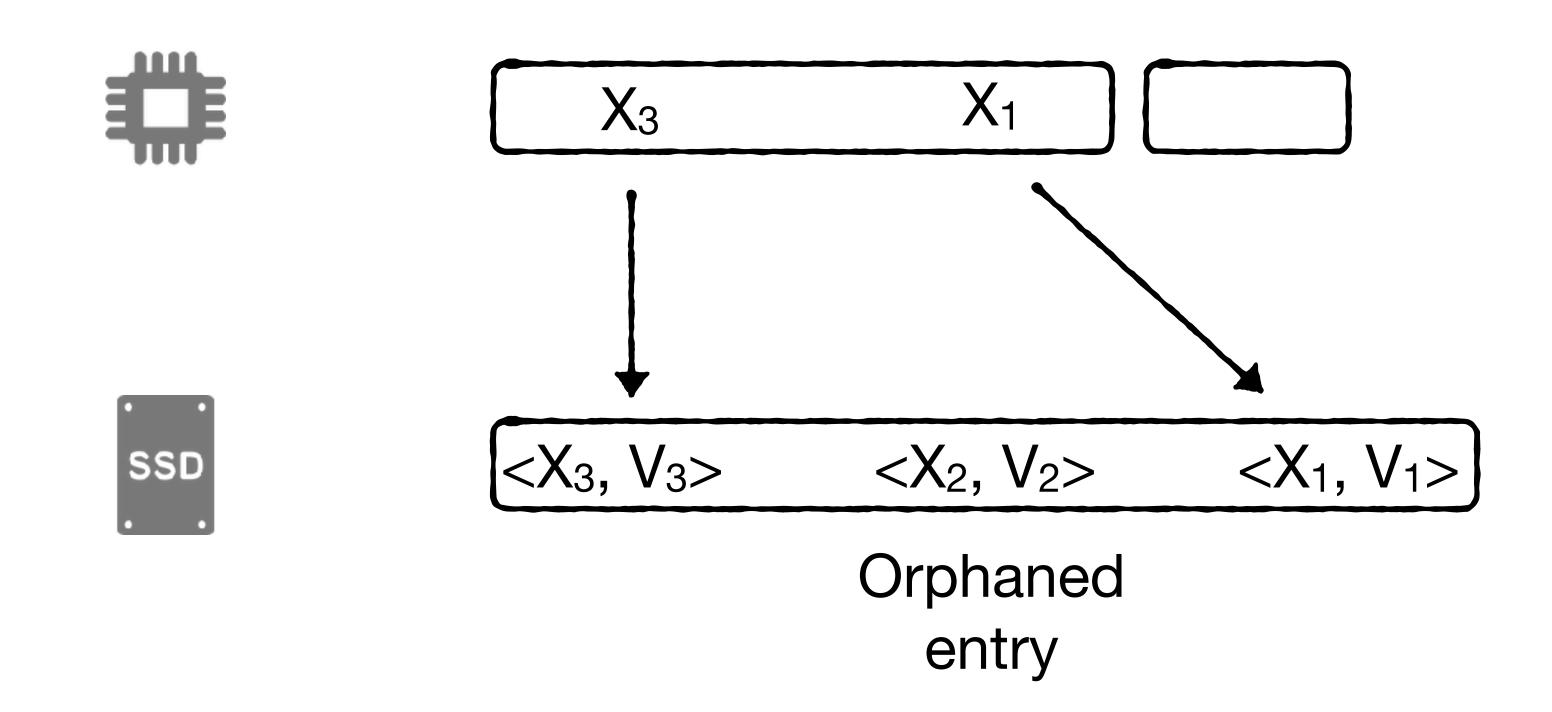
How to fix this?



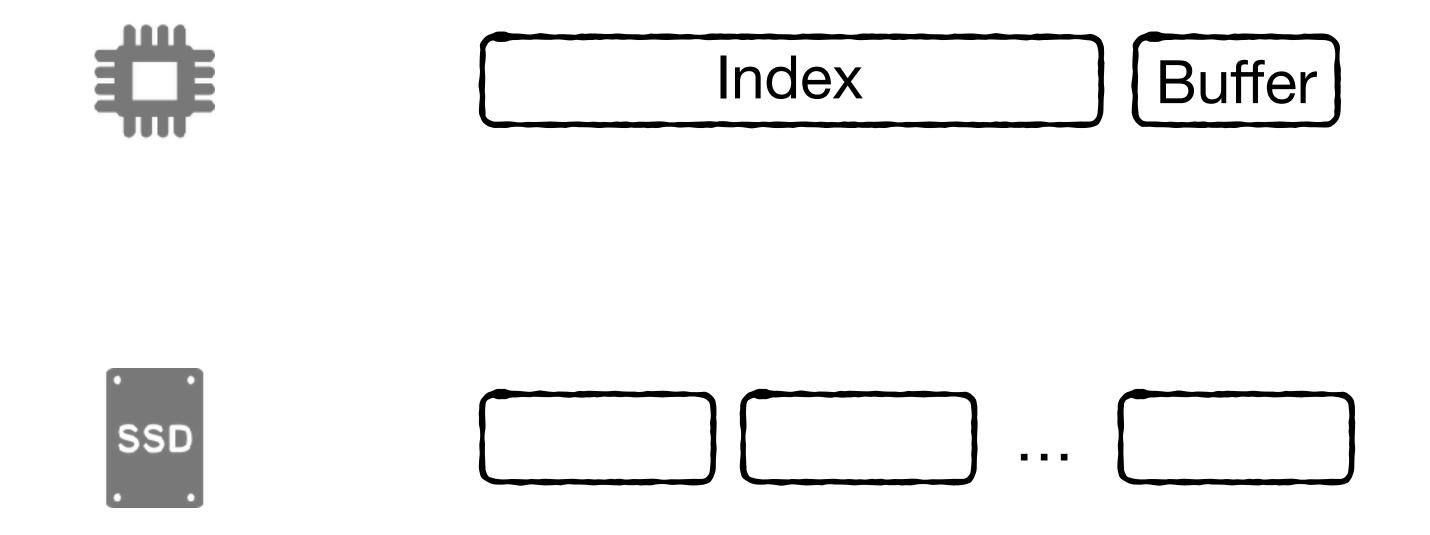
After many deletes, many orphaned entries accumulate

They take up space which we would prefer to use to store valid data

How to fix this? Garbage Collection (as we saw in SSDs)

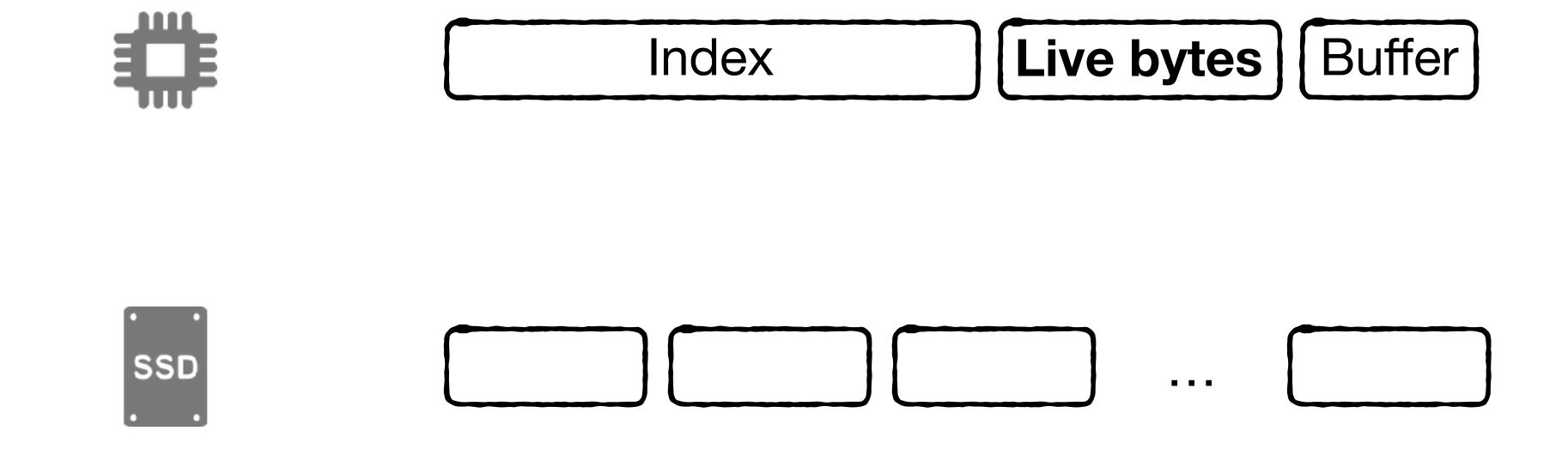


Conceptually divide log into equally-sized areas (can be in order of GBs)

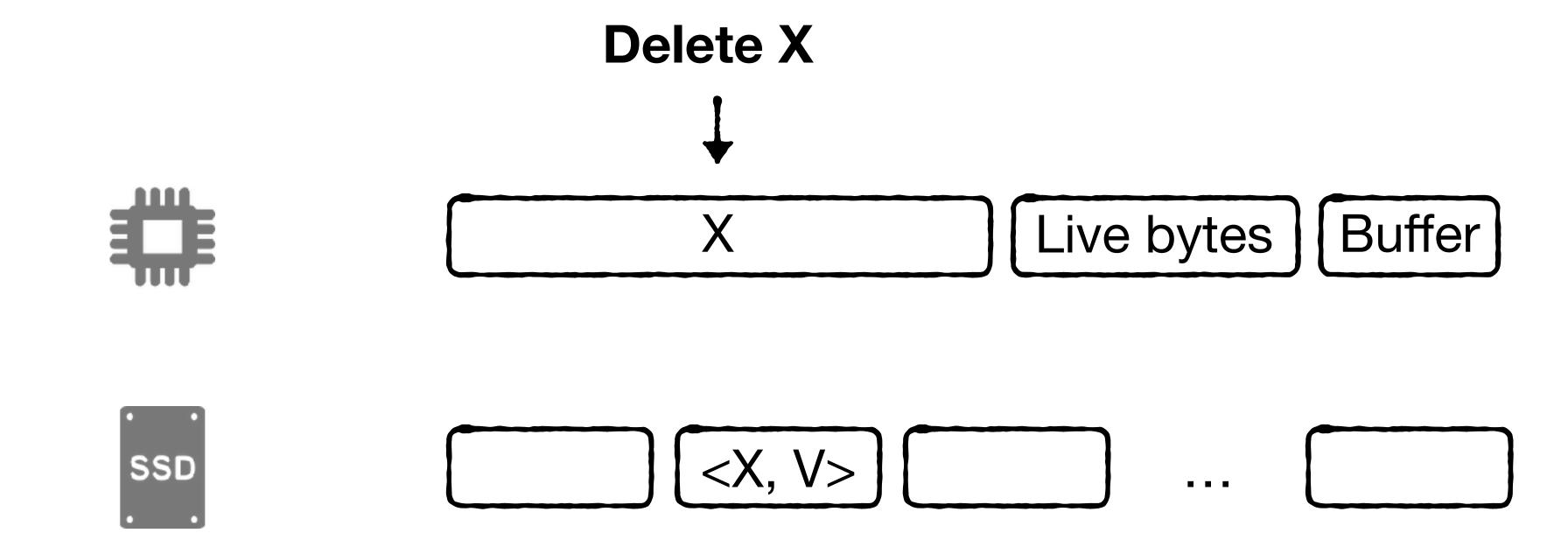


Conceptually divide log into equally-sized areas (can be in order of GBs)

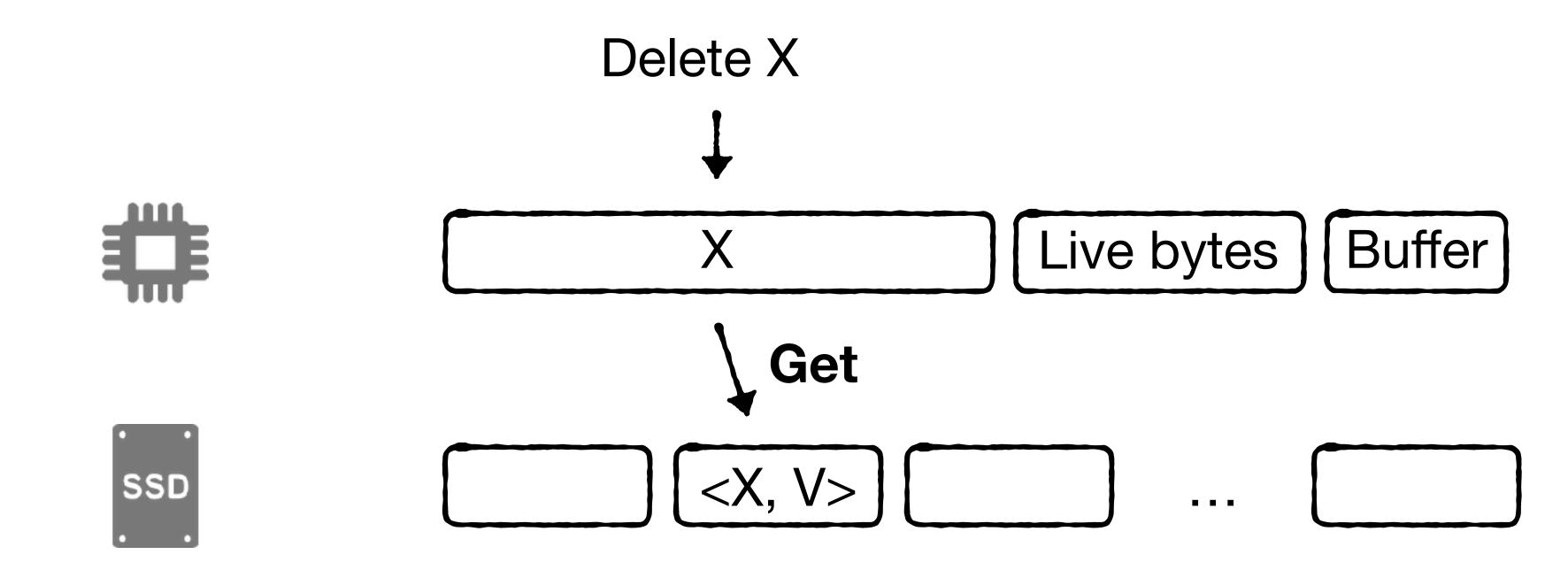
For each area, maintain a counter of the number of bytes representing live data



To delete:

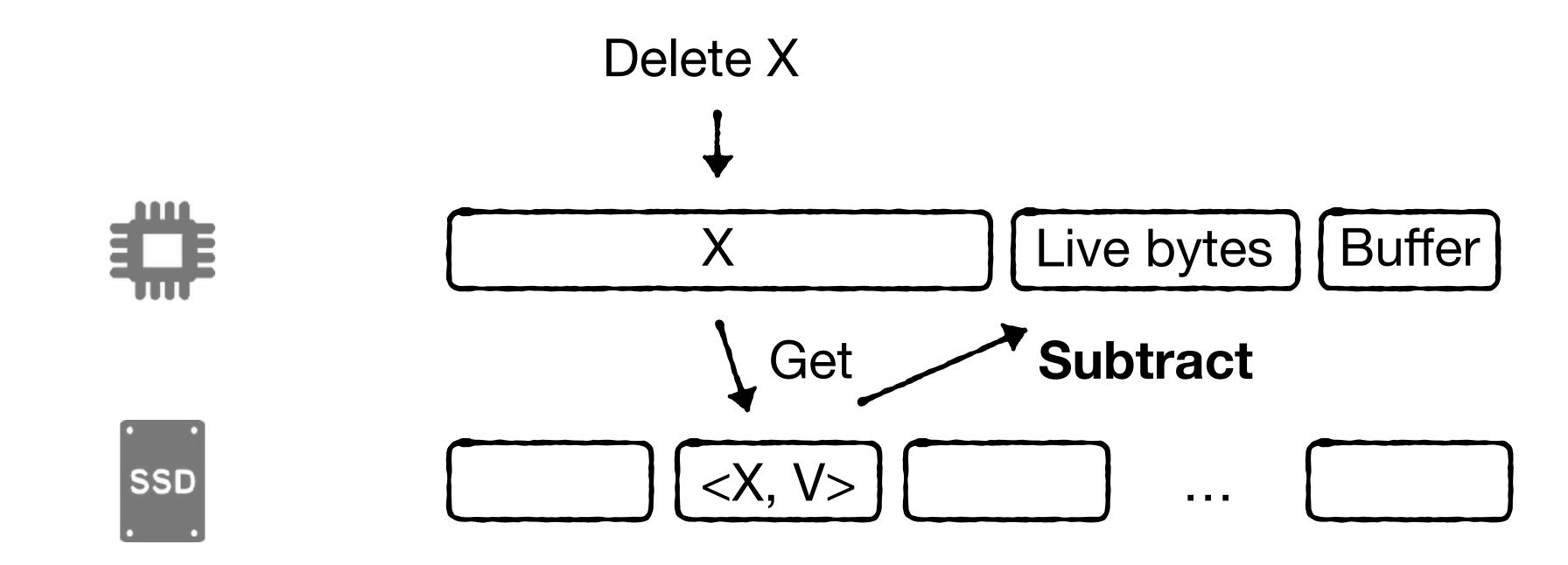


To delete: (1) get entry from storage to check its size



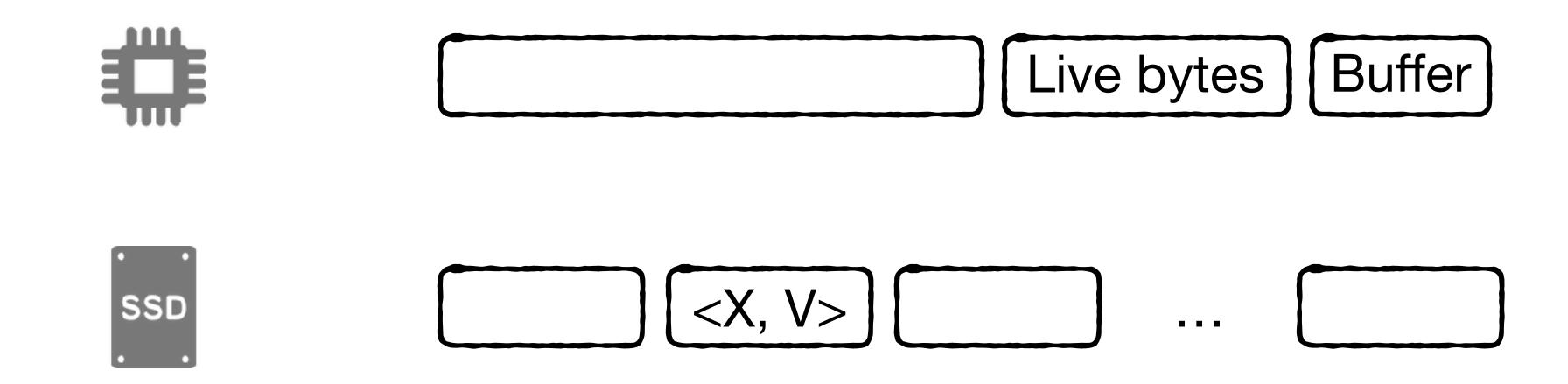
To delete: (1) get entry from storage to check its size

(2) subtract its size from area's counter



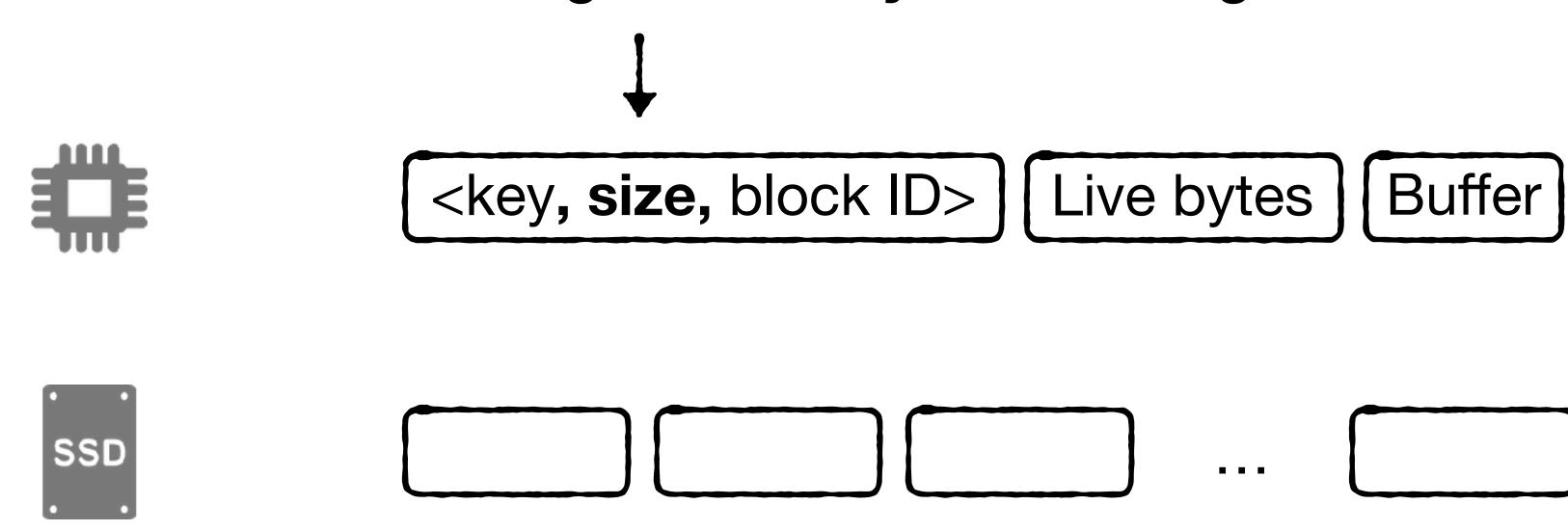
To delete: (1) get entry from storage to check its size

- (2) subtract its size from area's counter
- (3) remove entry from index



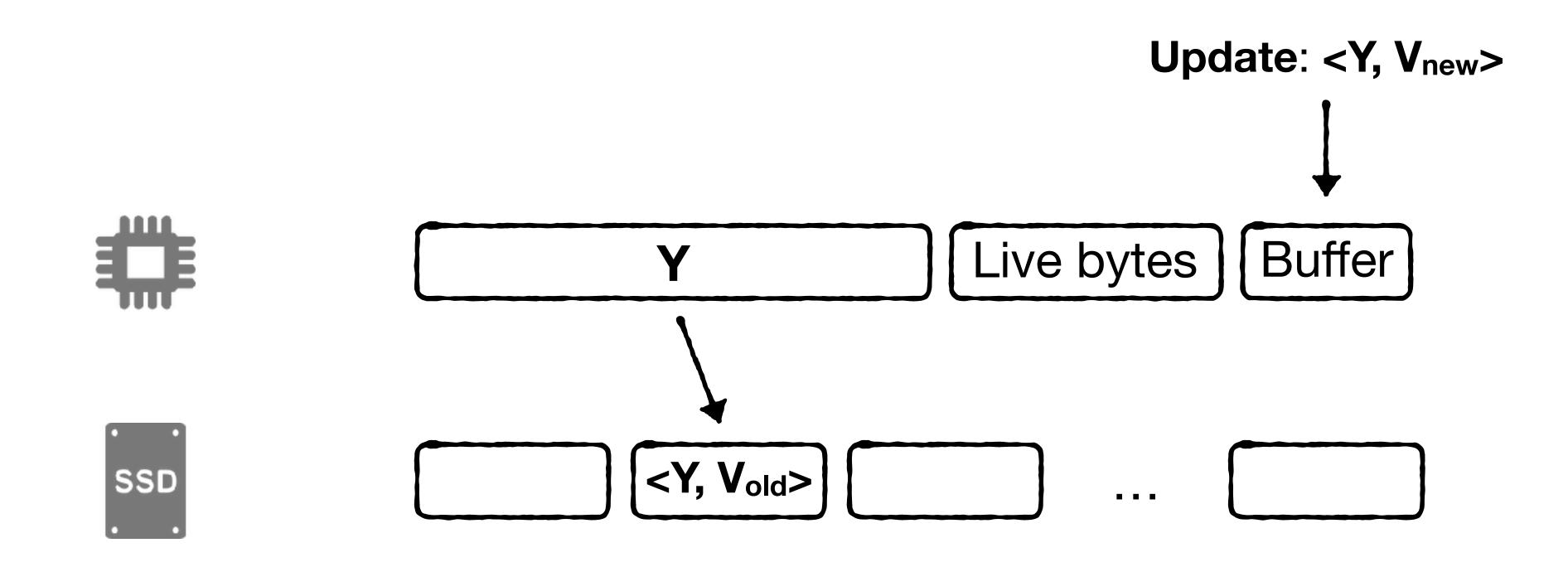
- To delete: (1) get entry from storage to check its size -
 - (2) subtract its size from area's counter
 - (3) remove entry from index

Can also store each entry's size in the index, so we not have to get the entry from storage



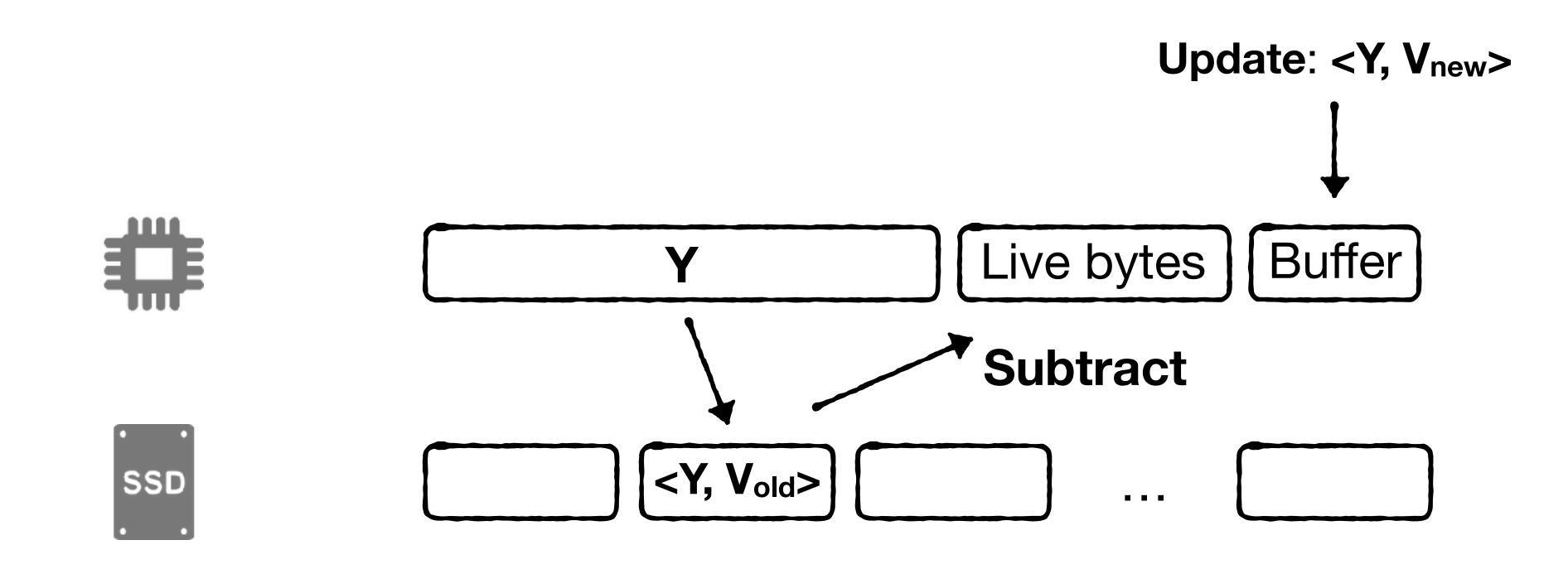
Updates

An update is a delete followed by an insertion of an entry with the same key



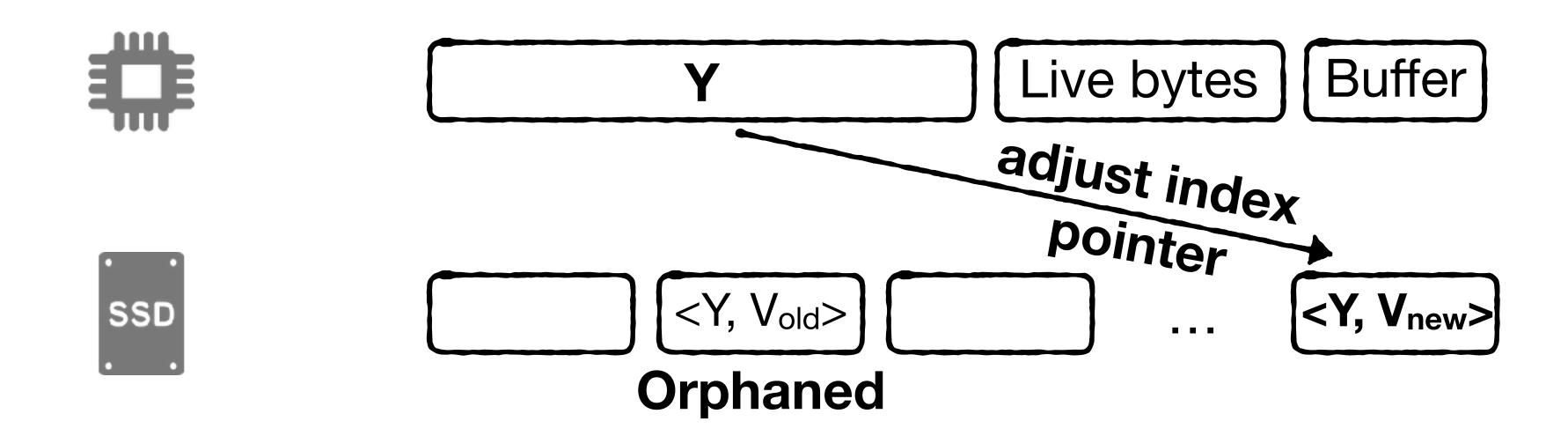
Updates

An update is a delete followed by an insertion of an entry with the same key

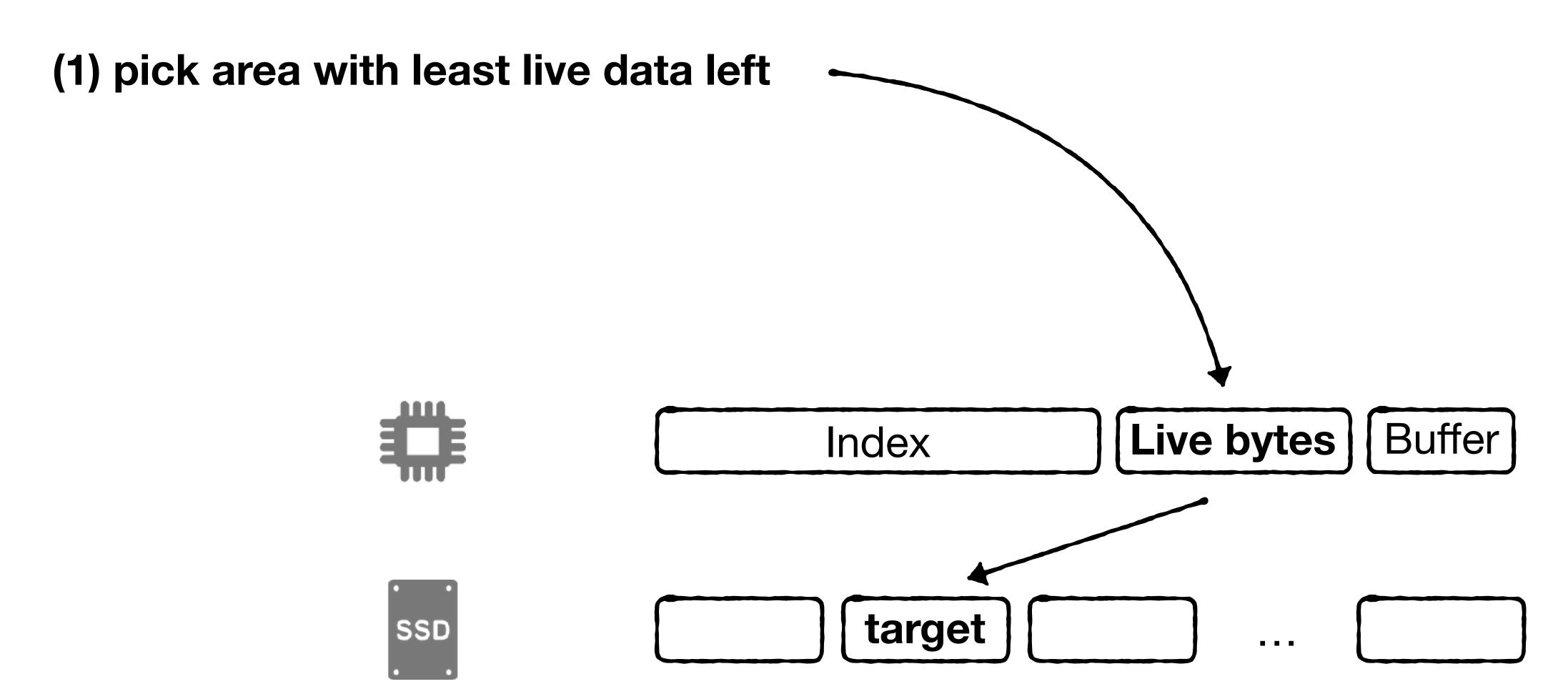


Updates

An update is a delete followed by an insertion of an entry with the same key

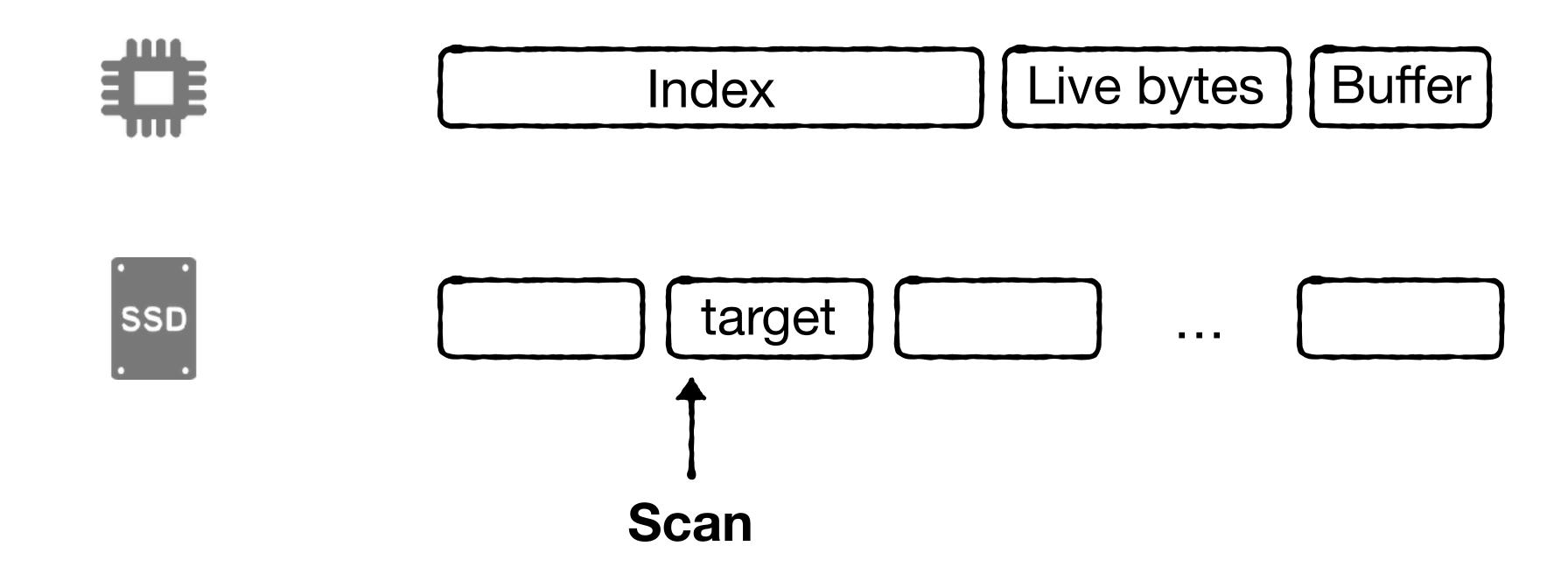


Garbage Collection



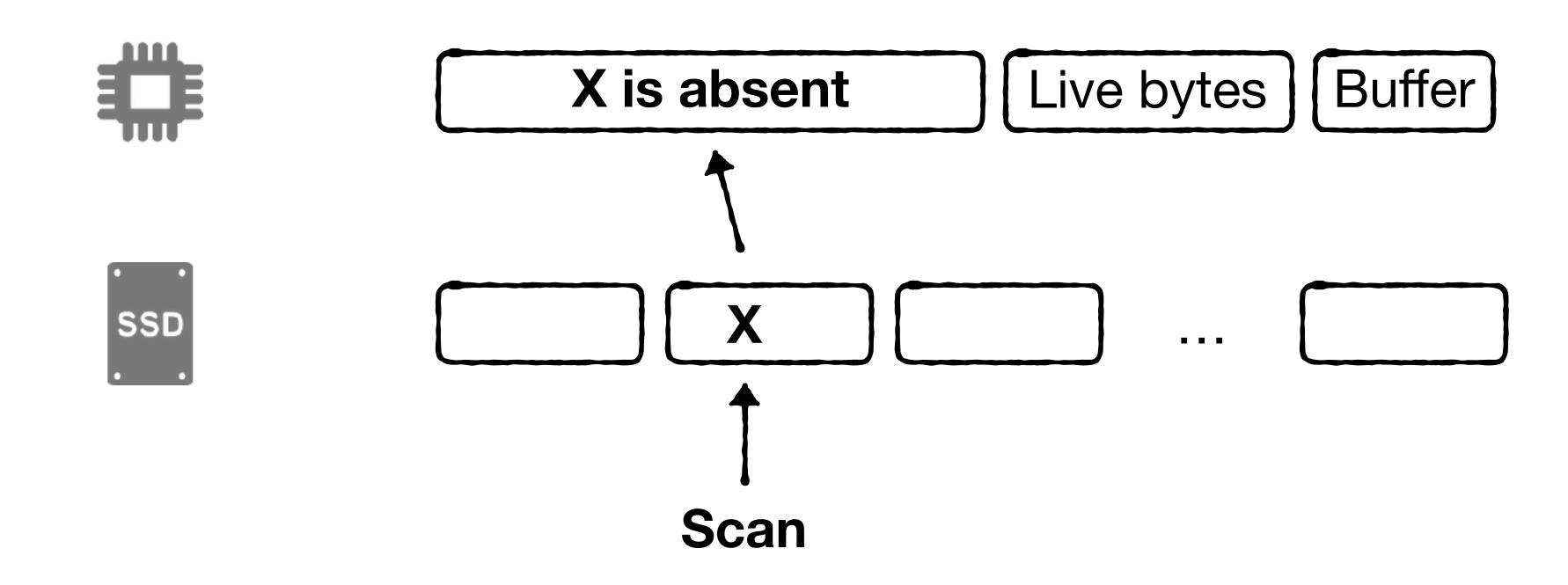
Garbage Collection

- (1) pick area with least live data left
- (2) Scan area and for each entry

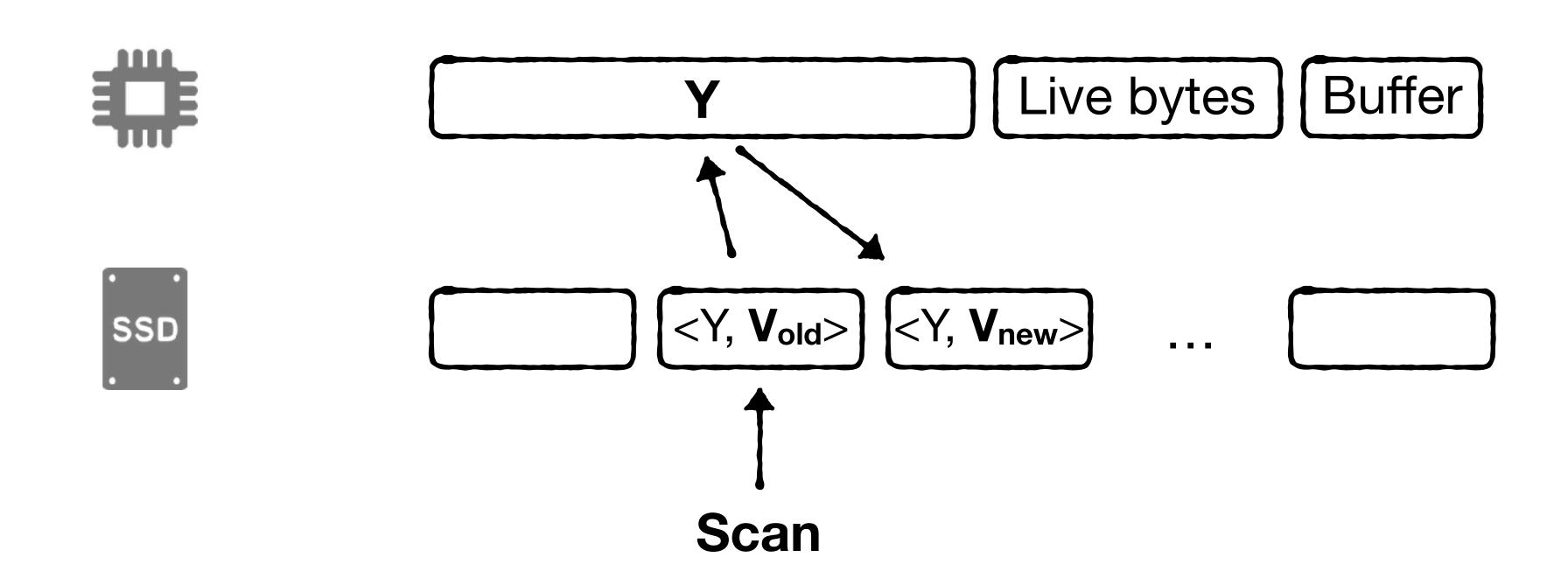


- (1) pick area with least live data left
- (2) Scan area and for each entry

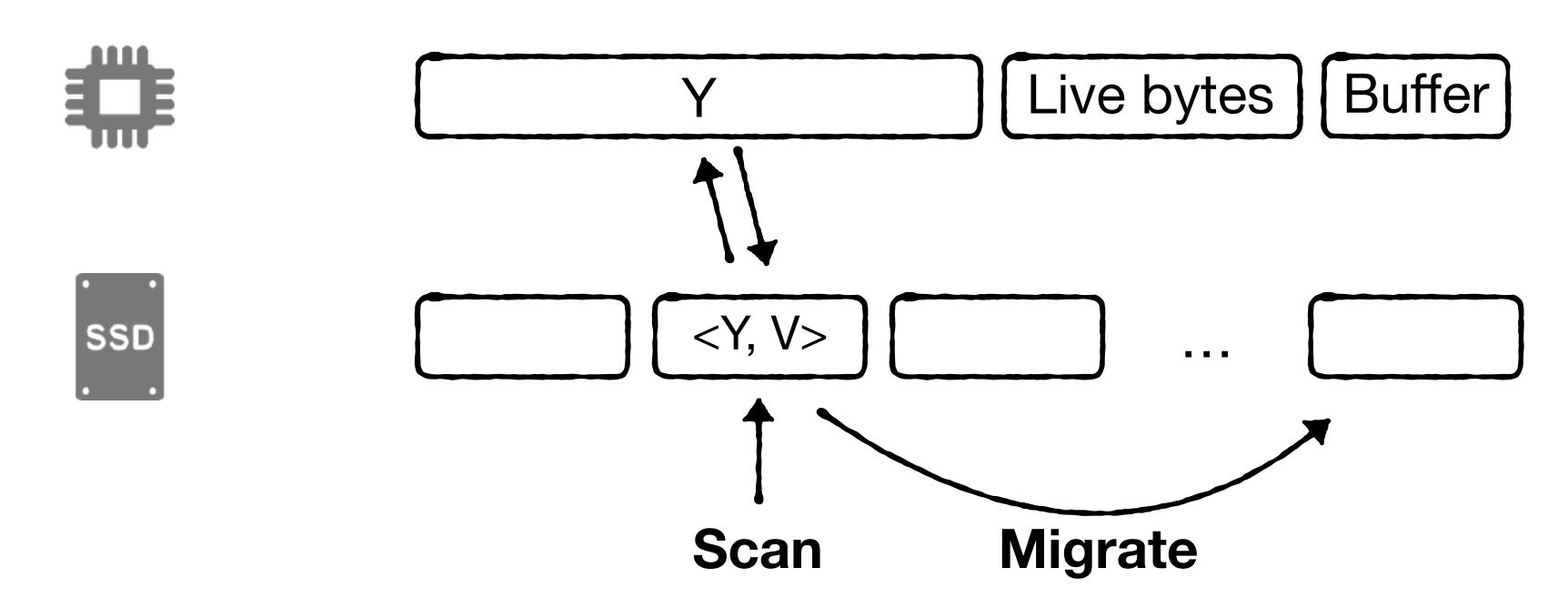
(A) If the key is not indexed, the entry had been deleted, so move on



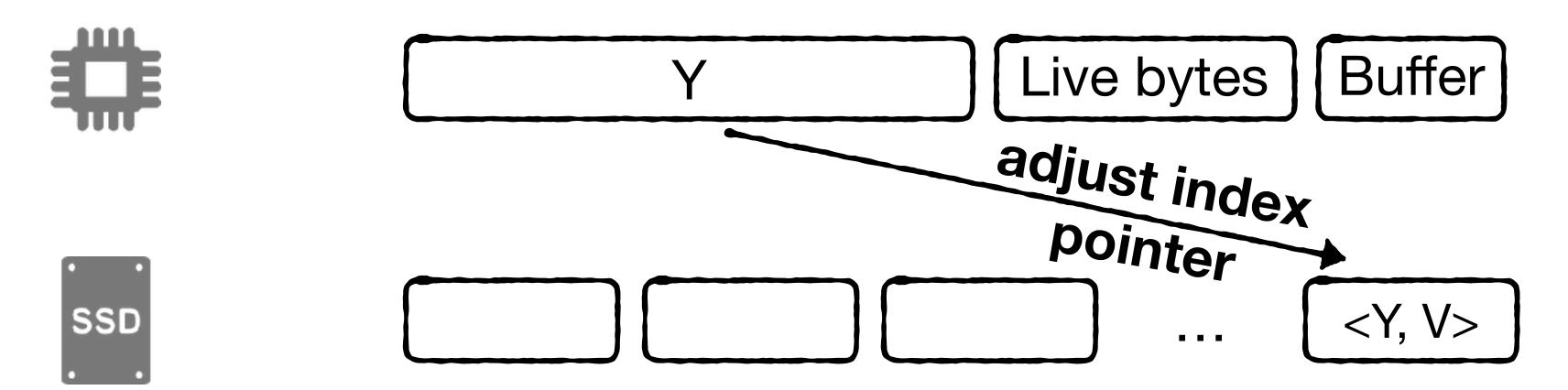
- (1) pick area with least live data left
- (2) Scan area and for each entry
 - (A) If the key is not indexed, the entry had been deleted, so move on
 - (B) If the key is indexed but pointing elsewhere, the entry is outdated, so move on



- (1) pick area with least live data left
- (2) Scan area and for each entry
 - (A) If the key is not indexed, the entry had been deleted, so move on
 - (B) If the key is indexed but pointing elsewhere, the entry is outdated, so move on
 - (C) If the key is indexed and pointing to this entry, it is valid, so migrate it

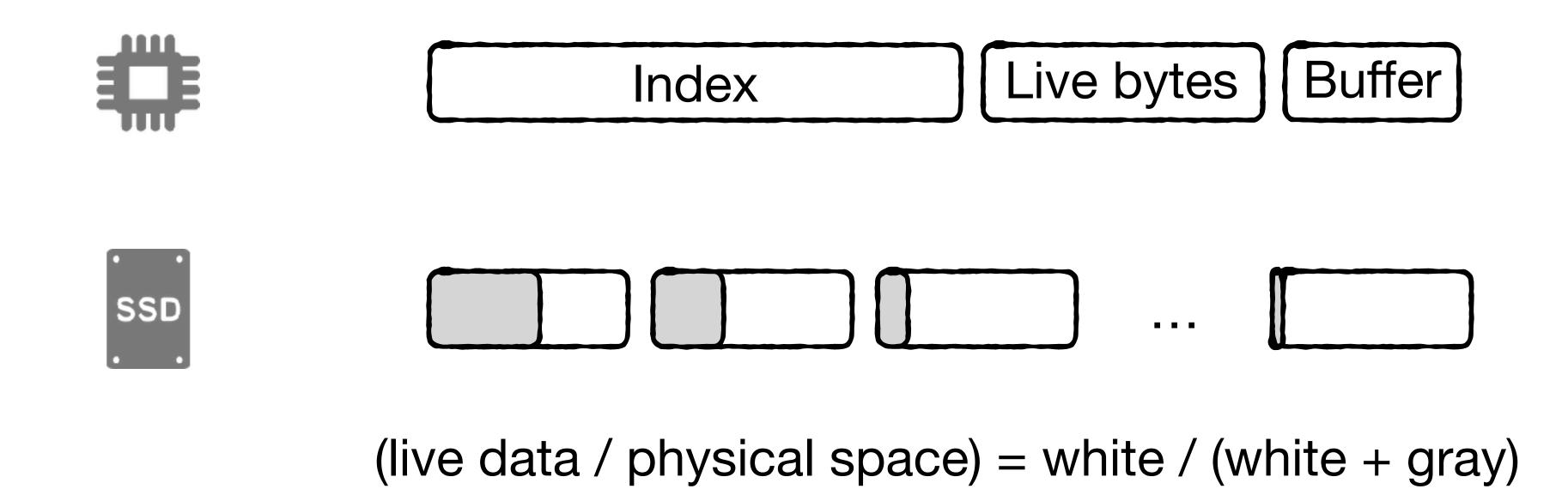


- (1) pick area with least live data left
- (2) Scan area and for each entry
 - (A) If the key is not indexed, the entry had been deleted, so move on
 - (B) If the key is indexed but pointing elsewhere, the entry is outdated, so move on
 - (C) If the key is indexed and pointing to this entry, it is valid, so migrate it



When to trigger garbage-collection?

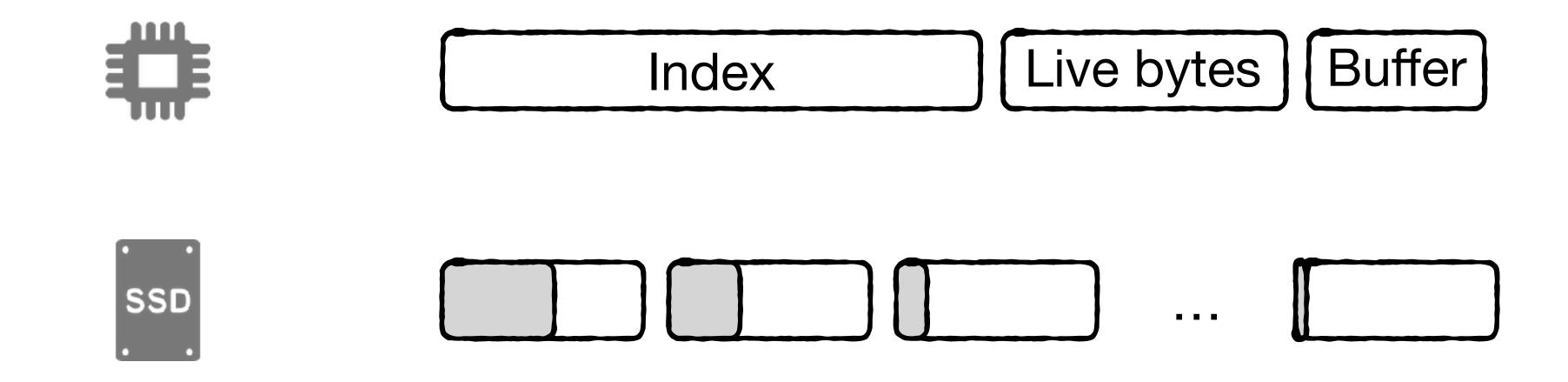
Define a global threshold of (live data L / physical space P)



When to trigger garbage-collection?

Define a global threshold of (live data L / physical space P)

When this threshold is reached, trigger garbage-collection to free space

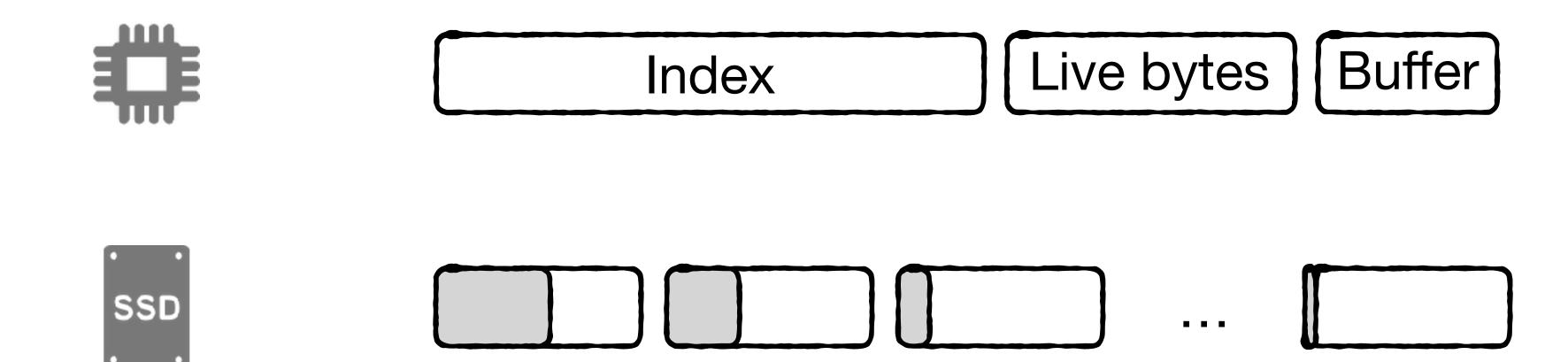


Garbage-Collection Write-Amplification

How to reason about this?

Let x = avg. % of valid pages in areas we pick to garbage-collect

$$\mathbf{WA} = 1 + \frac{x}{1 - x}$$



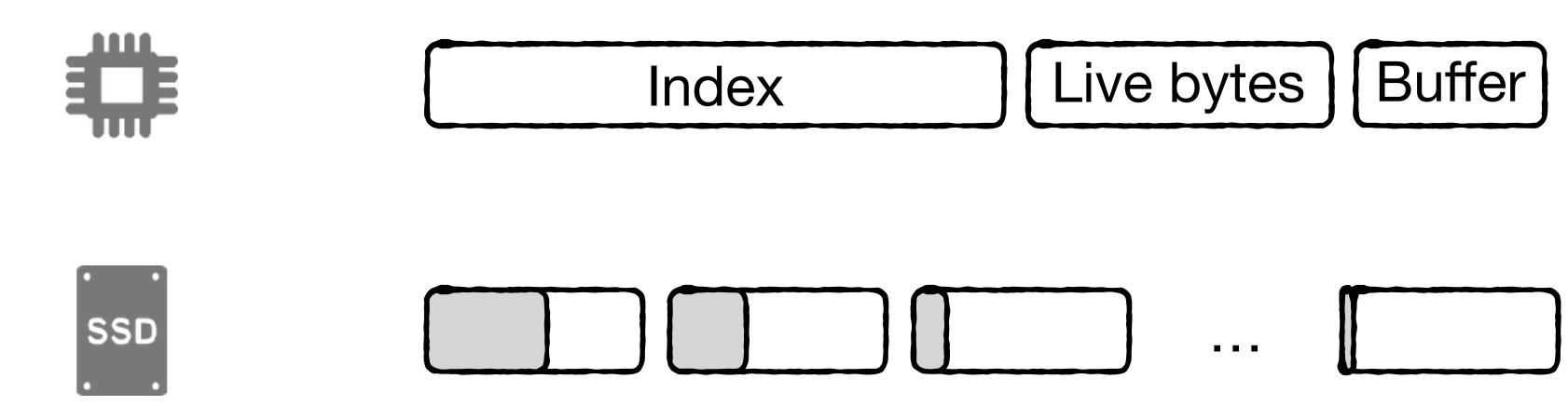
Garbage-Collection Write-Amplification

How to reason about this?

Let x = avg. % of valid pages in areas we pick to garbage-collect

$$WA = 1 + \frac{x}{1 - x}$$

Same analysis as for garbage-collection in SSDs, so refer to that:)



Garbage-Collection Write-Amplification



$$+\frac{L/P}{1 I/D}$$

 $1 + \frac{1}{2} \cdot \frac{L/P}{1 - I/P}$

Worst case

Uniformly random

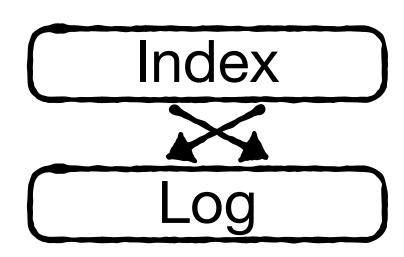
L = logical data size

P = physical data size

Mechanics

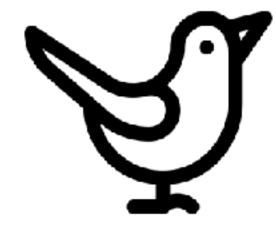
Hot/Cold Separation Checkpointing/recovery

Cuckoo filtering









To reduce write-amp

If power fails

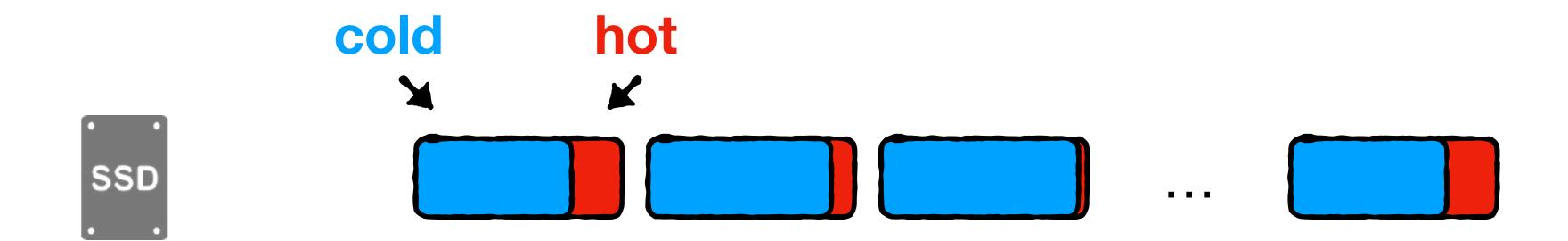
To reduce memory

Hot/Cold Separation

Normal workloads are neither worst-case nor randomly distributed

Normal workloads are neither worst-case nor randomly distributed

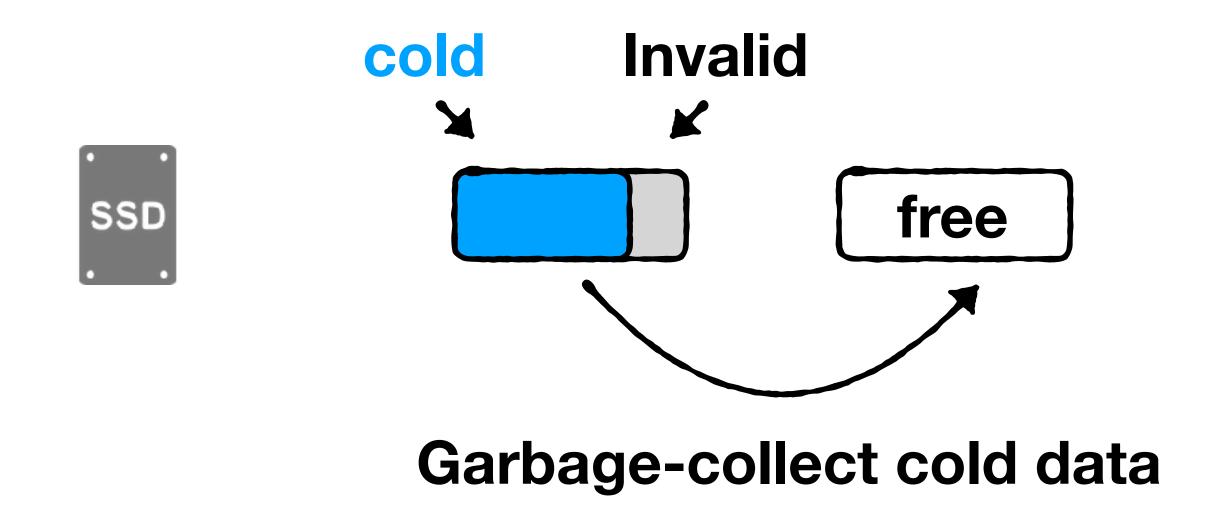
Typically, few entries are frequently updated while most are seldom updated



Normal workloads are neither worst-case nor randomly distributed

Typically, few entries are frequently updated while most are seldom updated

Hot entries are invalidated quickly, so by the time we garbage-collect, there is usually only cold data left

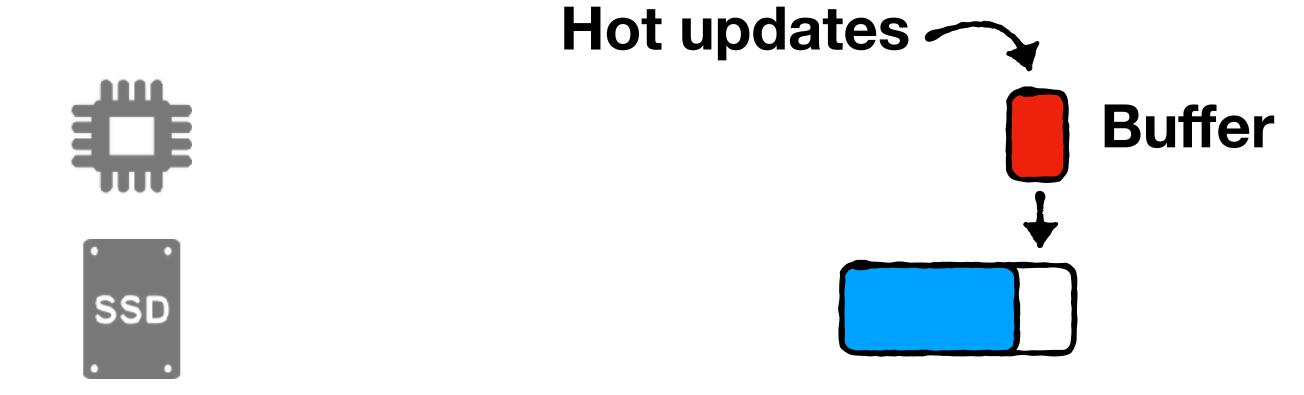


Normal workloads are neither worst-case nor randomly distributed

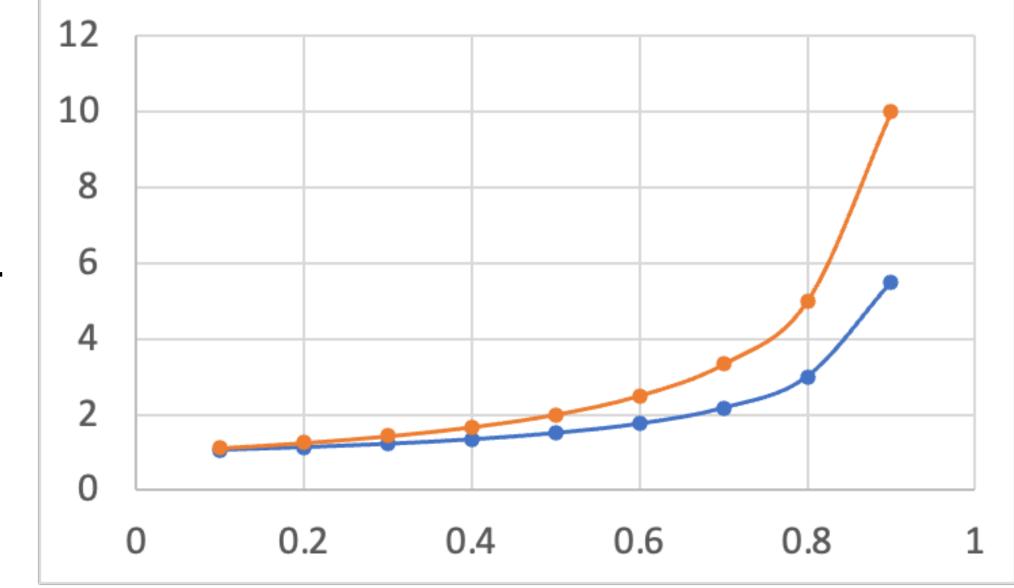
Typically, few entries are frequently updated while most are seldom updated

Hot entries are invalidated quickly, so by the time we garbage-collect, there is usually only cold data left

This migrated cold data gets mixed with more hot data, and the cycle repeats.





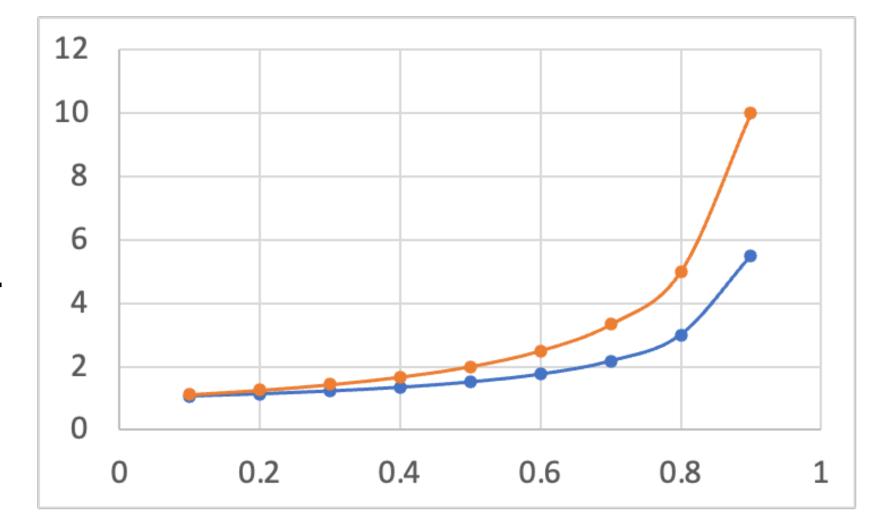


$$+\frac{L/P}{1-L/P}$$
 Mixing hot/cold data brings us towards worst-case

$$+\frac{L/P}{2 \cdot (1 - L/P)}$$
 Uniformly random workloads

Hot vs. Cold Data

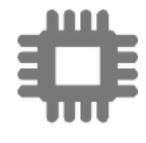




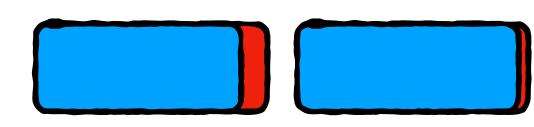
$$+\frac{L/P}{1-L/P}$$

Mixing hot/cold data brings us towards worst-case

How can we avoid garbagecollecting cold data all the time?

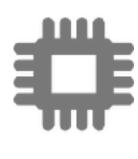




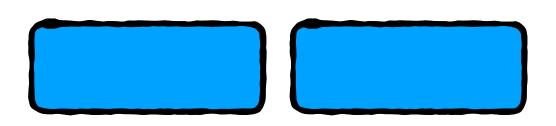








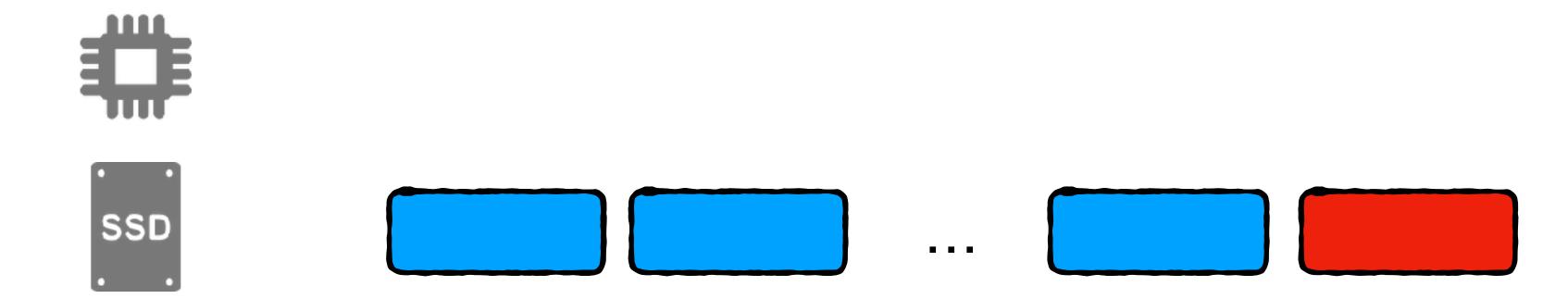






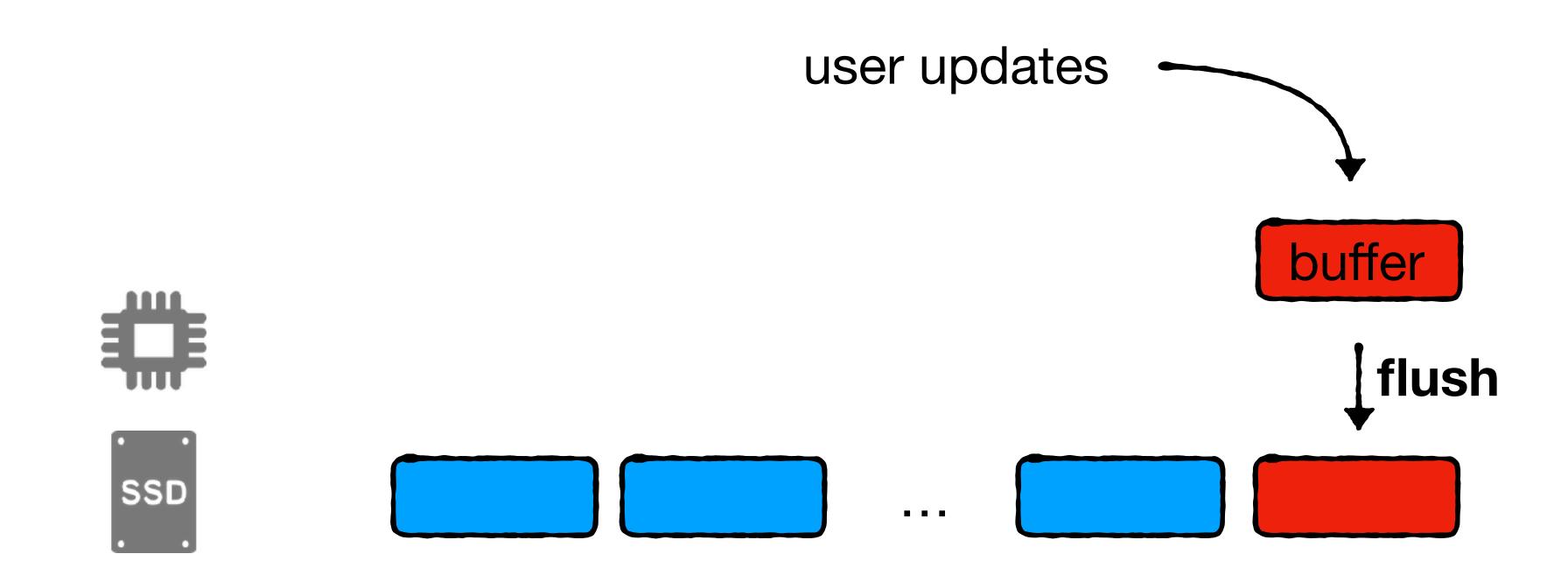


We can separate hot vs. cold data into different areas



We can separate hot vs. cold data into different areas

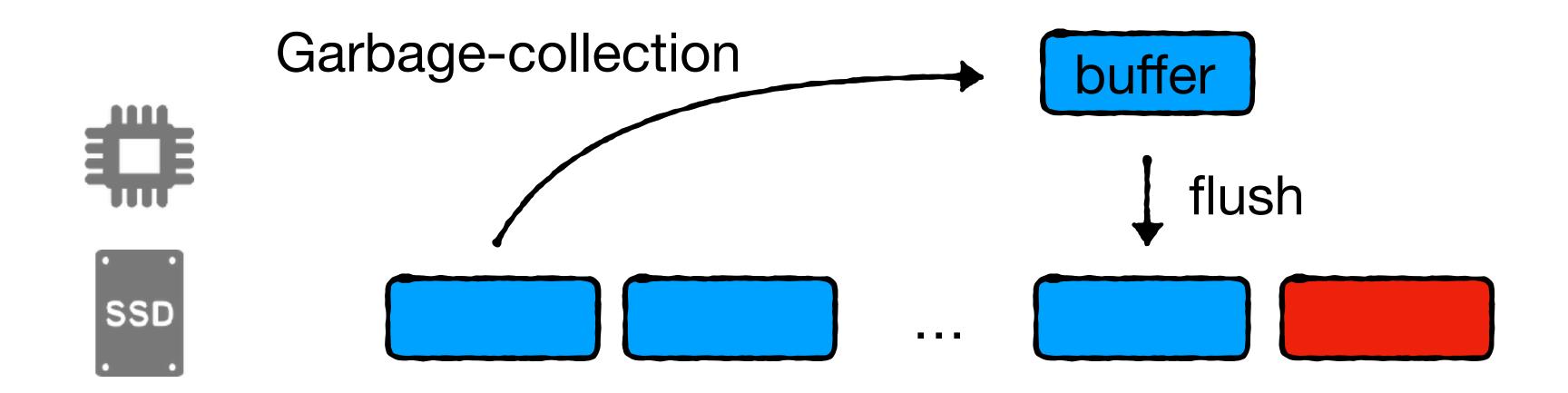
Insight 1: user updates are generally hot (i.e., data recently written is likely to be written again)



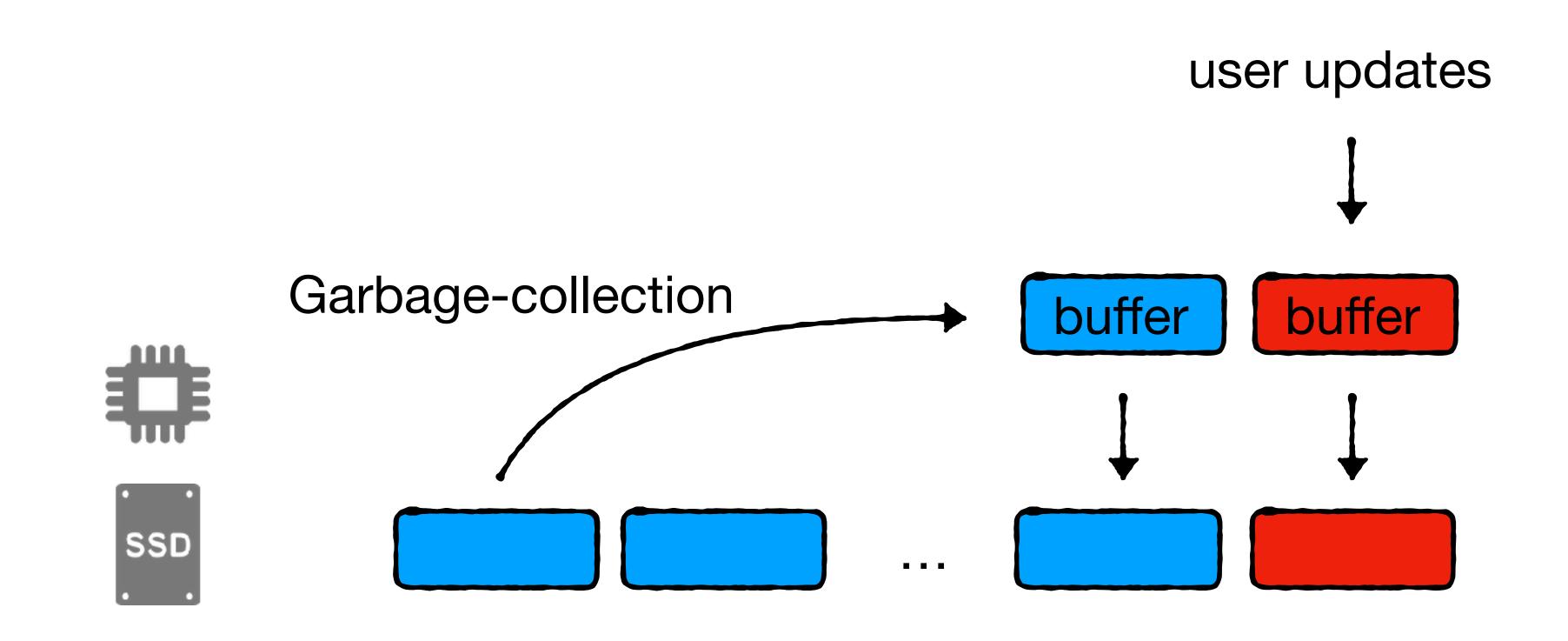
We can separate hot vs. cold data into different areas

Insight 1: user updates are generally hot (i.e., data recently written is likely to be written again)

Insight 2: garbage-collected data is generally cold (i.e., it had already existed for a long while without getting updated)



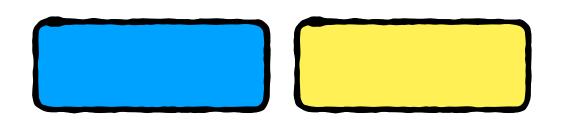
Simplest solution: separate user updates and cold data using different buffers into different areas



Simplest solution: separate user updates and cold data using different buffers into different areas

More advanced: separate data with different temperatures into different areas



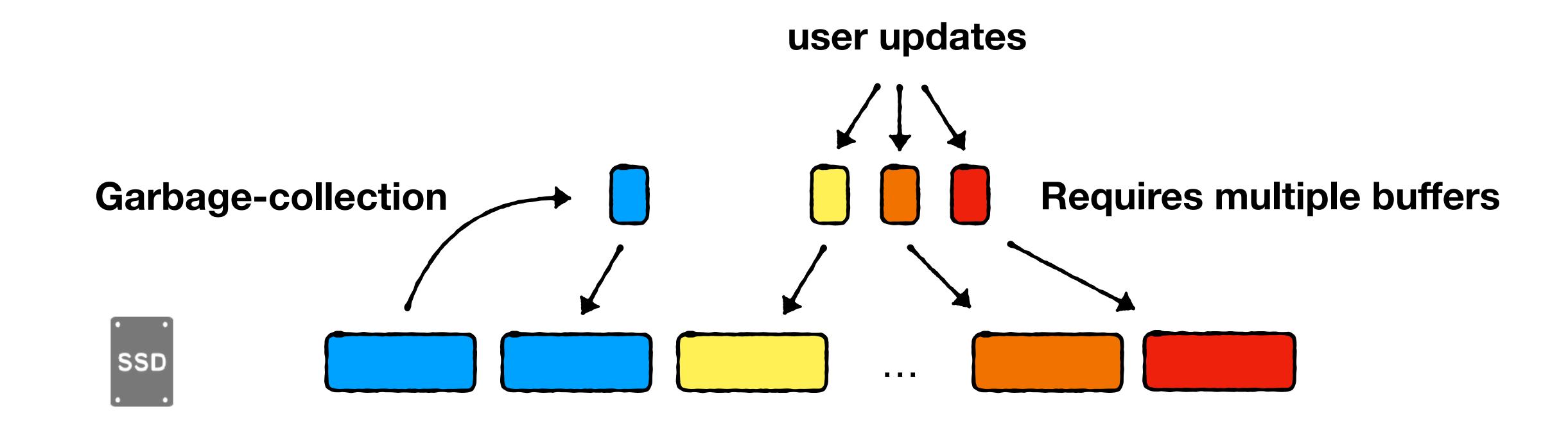






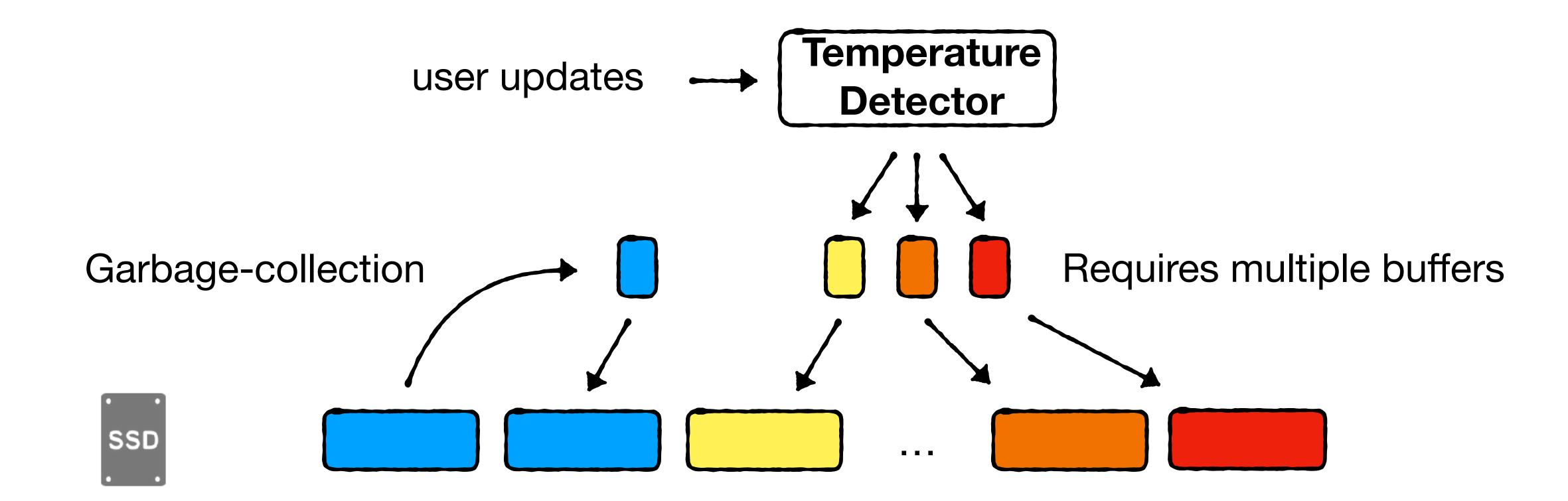
Simplest solution: separate user updates and cold data using different buffers into different areas

More advanced: separate data with different temperatures into different areas

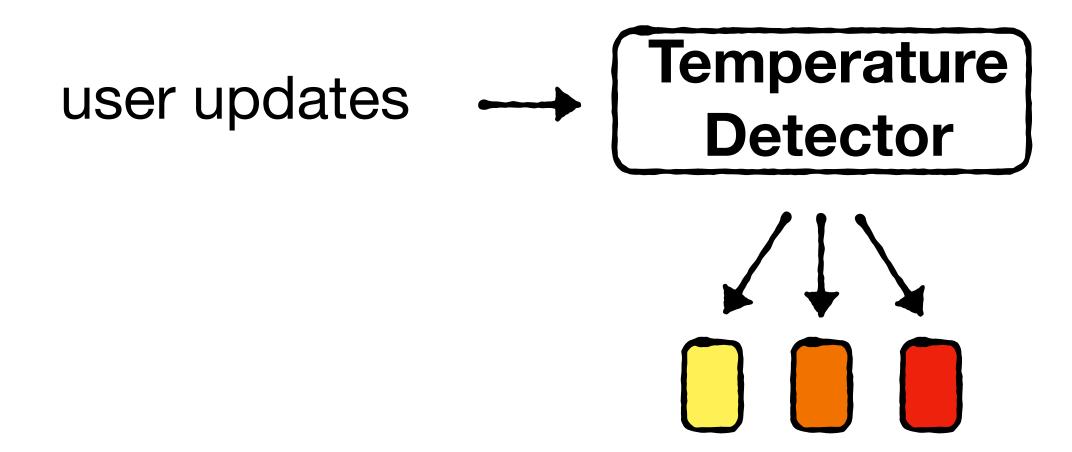


Simplest solution: separate user updates and cold data using different buffers into different areas

More advanced: separate data with different temperatures into different areas

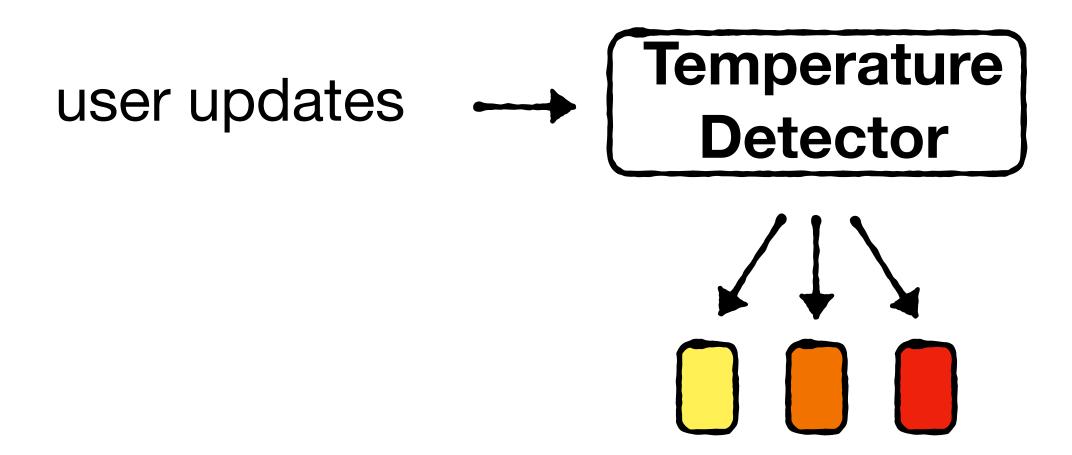


Estimates how likely a page is to be updated again



Estimates how likely a page is to be updated again

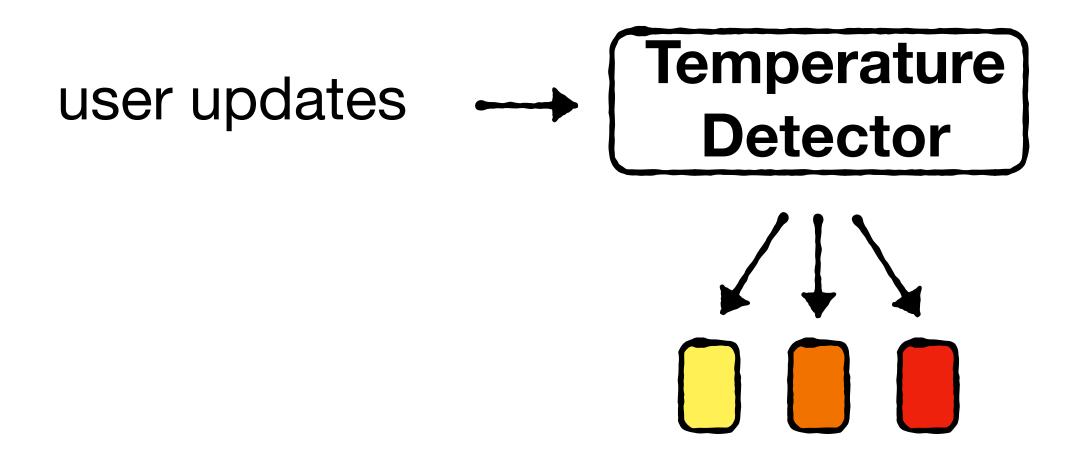
Should be a light-weight data structure that can fit in memory



Estimates how likely a page is to be updated again

Should be a light-weight data structure that can fit in memory

There is a lot of research on this, but we'll explore just one solution relying on a cool data structure called count-min



A data structure that reports that frequency of elements in a data stream

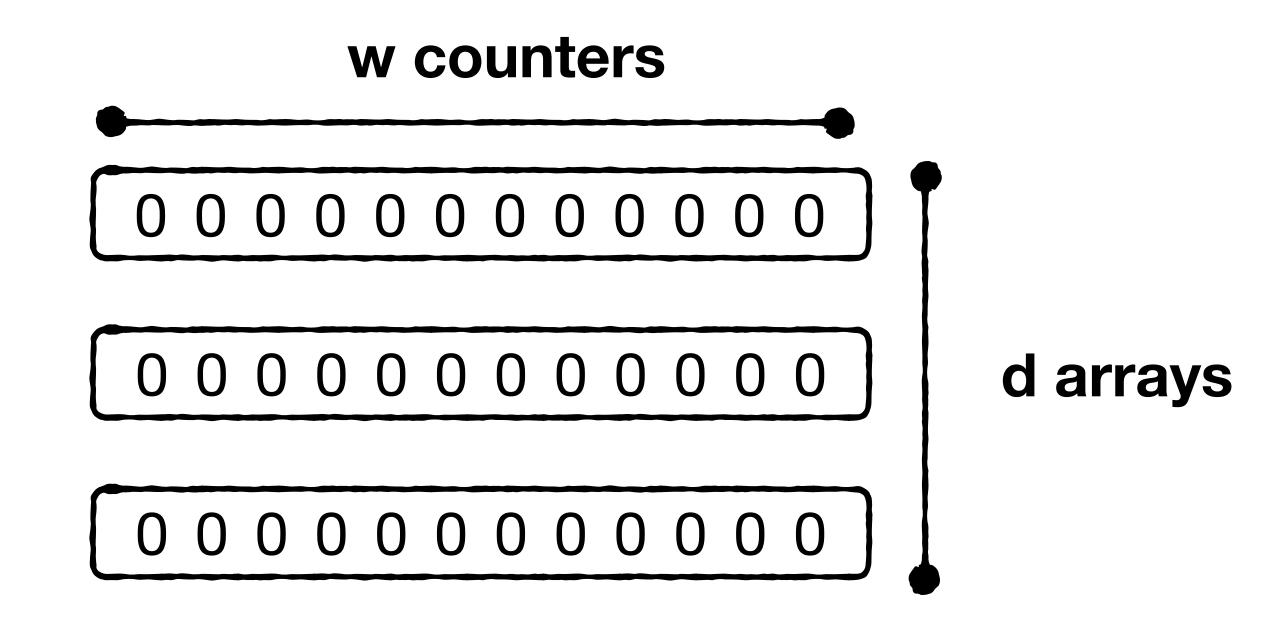
0000000000

0 0 0 0 0 0 0 0 0 0

0000000000

A data structure that reports that frequency of elements in a data stream

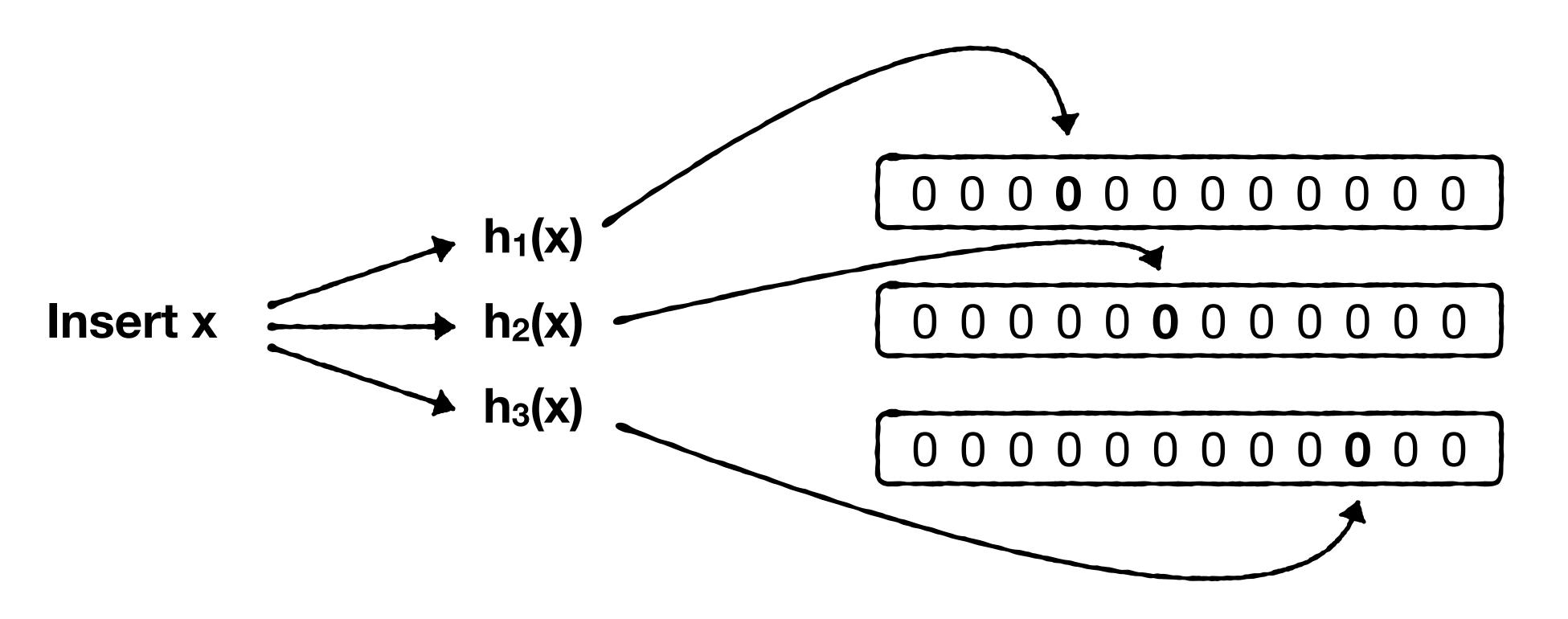
Consists of d arrays of w counters each



A data structure that reports that frequency of elements in a data stream

Consists of d arrays of w counters each

Insert an entry by hashing it to one counter in each array using different hash function

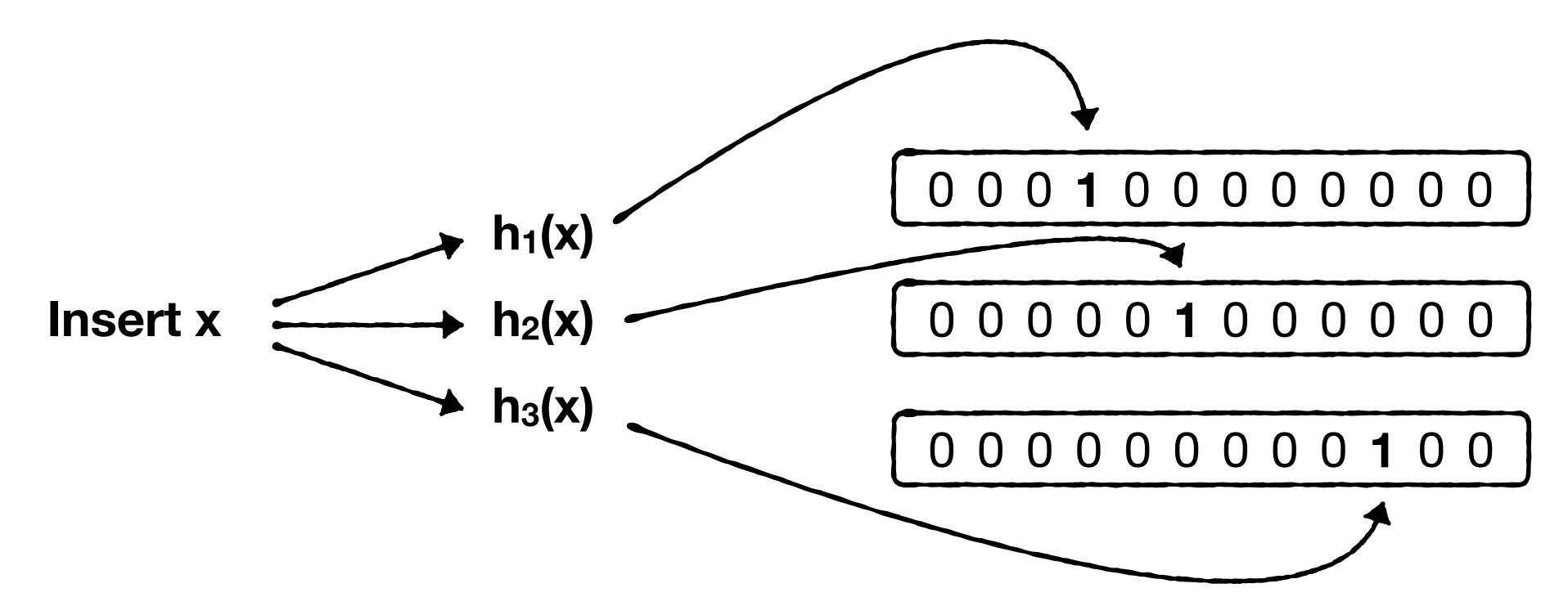


A data structure that reports that frequency of elements in a data stream

Consists of d arrays of w counters each

Insert an entry by hashing it to one counter in each array using different hash function

And increment that counter



After a while counters have a wide distribution of values

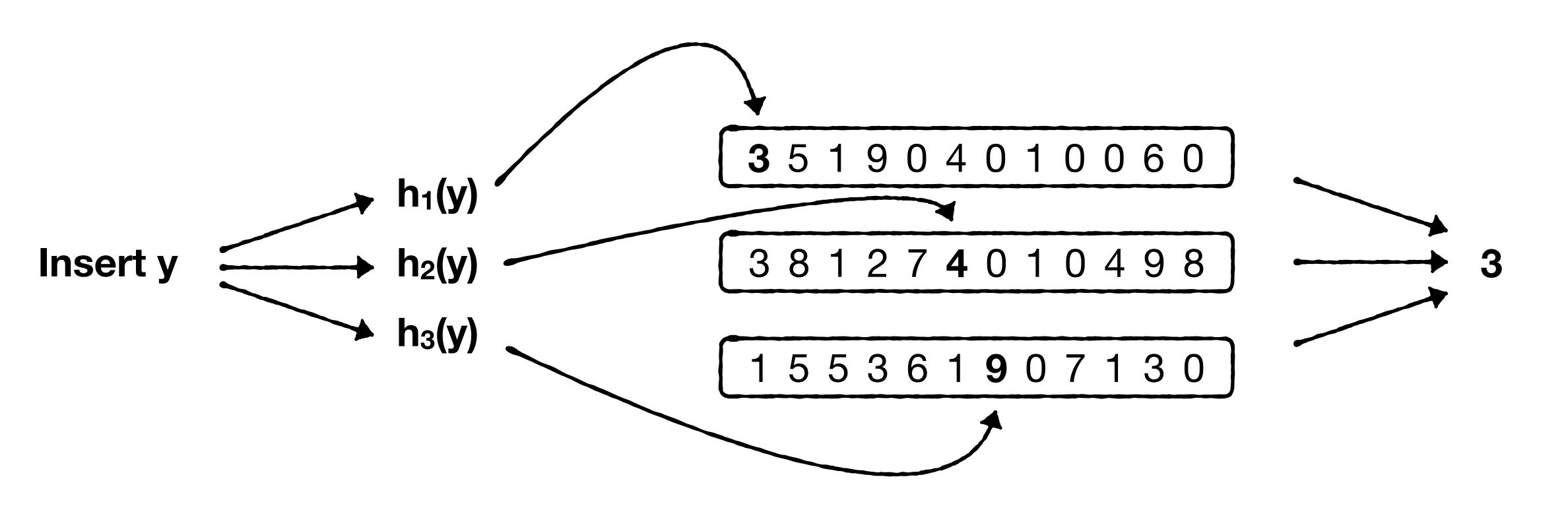
351904010060

3 8 1 2 7 4 0 1 0 4 9 8

155361907130

After a while counters have a wide distribution of values

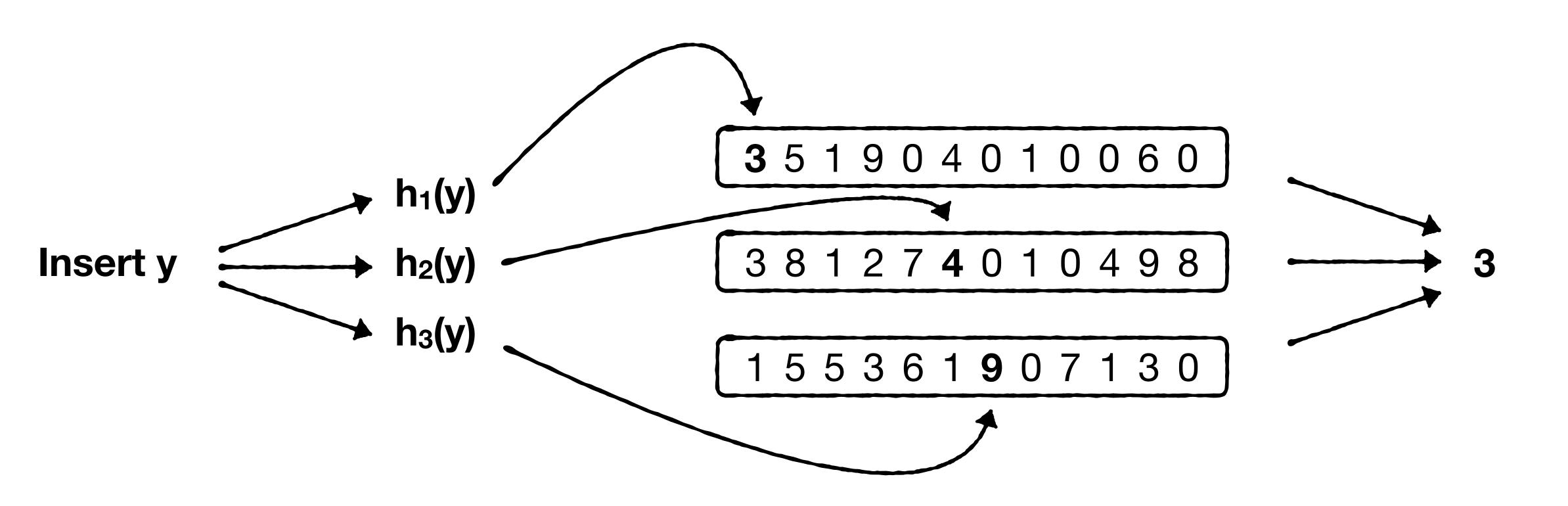
To query for the frequency of a key, hash it to each array and return minimum



After a while counters have a wide distribution of values

To query for the frequency of a key, hash it to each array and return minimum

This result is a guaranteed upper bound of the real count

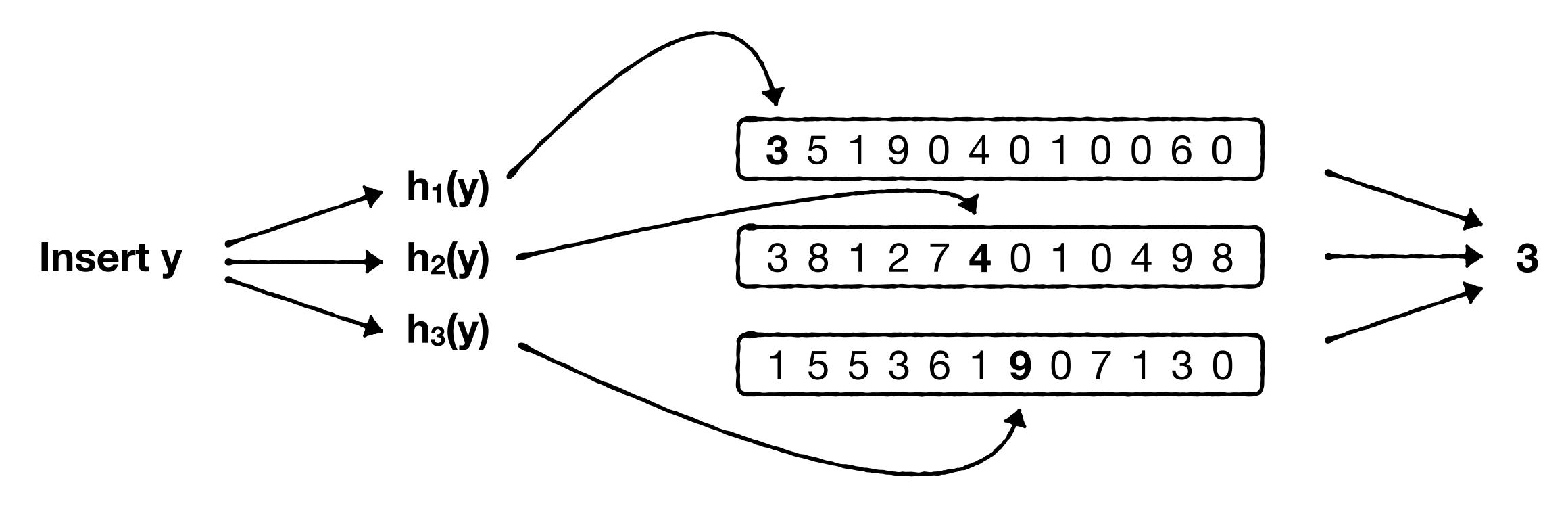


After a while counters have a wide distribution of values

To query for the frequency of a key, hash it to each array and return minimum

This result is a guaranteed upper bound of the real count

Result might overestimate the true answer due to hash collisions



Estimated frequency ≤ true frequency + ε · num insertions

with prob. 1- δ

351904010060

3 8 1 2 7 4 0 1 0 4 9 8

155361907130

Count-Min

Estimated frequency ≤ true frequency + ε · num insertions



Parameters

351904010060

3 8 1 2 7 4 0 1 0 4 9 8

155361907130

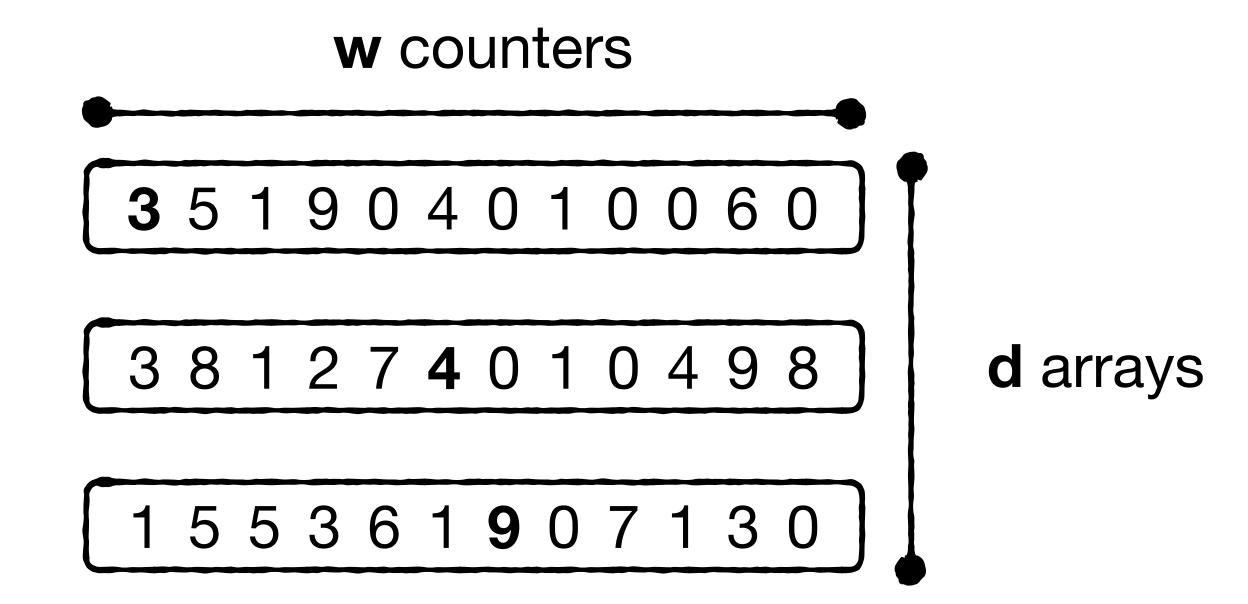
Count-Min

Estimated frequency ≤ true frequency + ε · num insertions

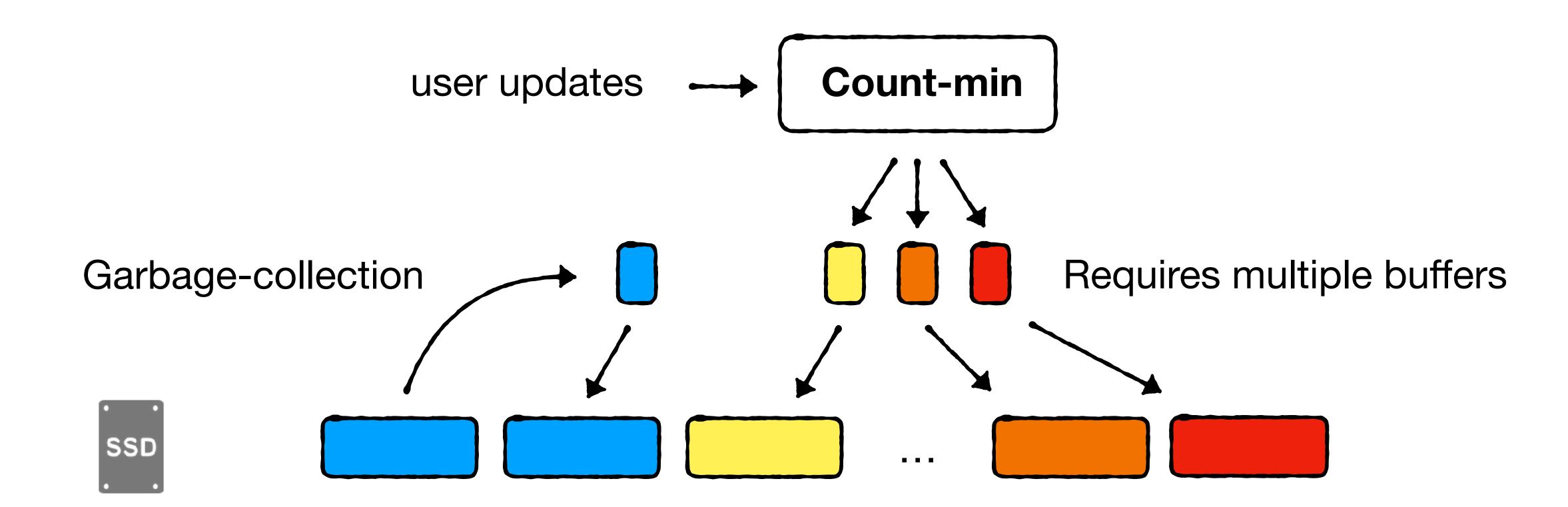
with prob. 1- δ

$$w = \lceil e/\epsilon \rceil$$
$$d = \lceil \ln 1/\delta \rceil$$

The parameters determine the values of w and d



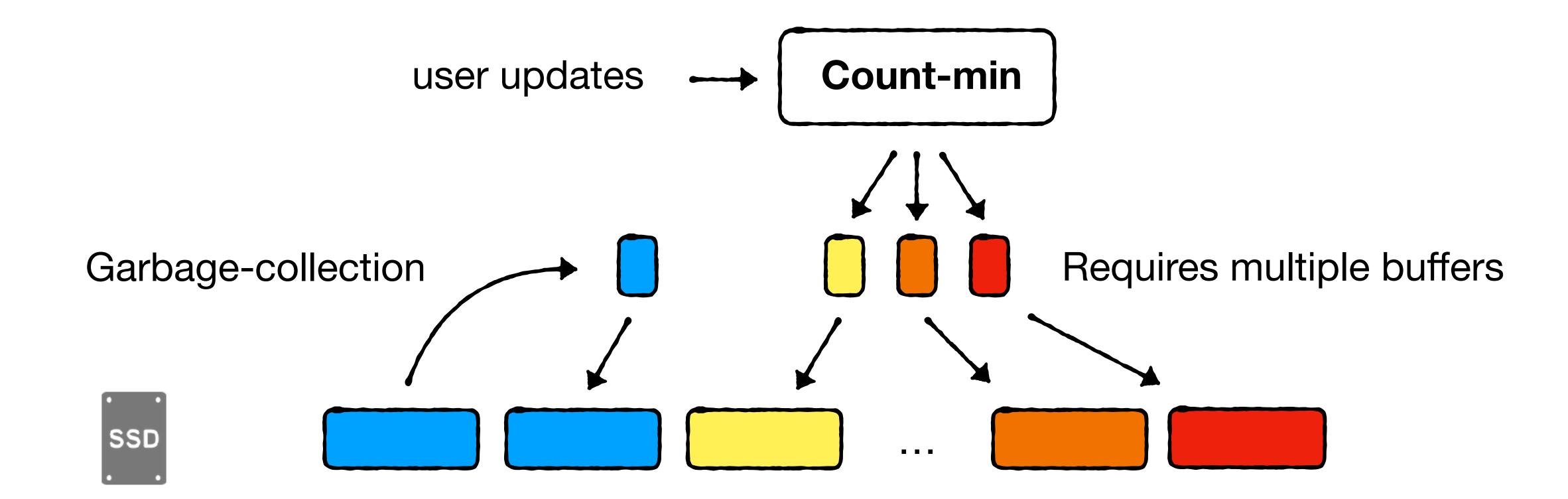
We can employ count-min to estimate frequency



We can employ count-min to estimate frequency

Pitfall: a value was very hot in the past but became cold.

Count-min would still tell us its hot. Solutions?



Decay: every x insertions, we can divide all counters by 2.

351904010060

3 8 1 2 7 4 0 1 0 4 9 8

155361907130

Decay: every x insertions, we can divide all counters by 2.

12140200030

1 4 0 1 3 2 0 0 0 2 4 4

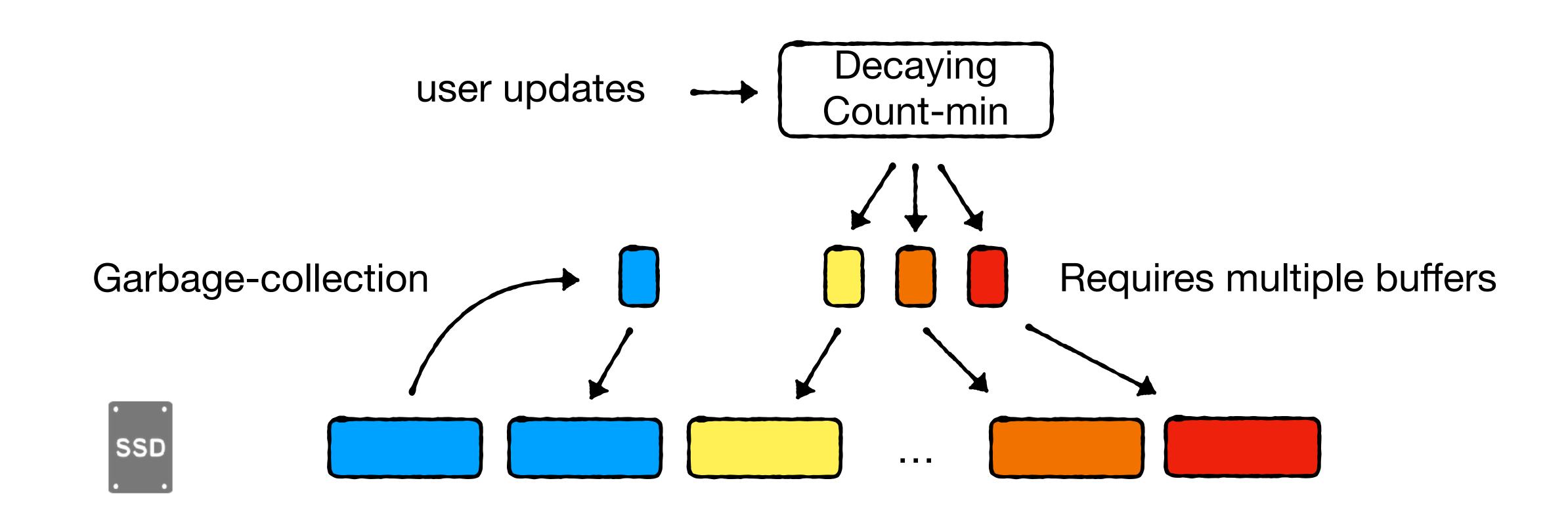
0 2 2 1 3 0 4 0 3 0 1 0

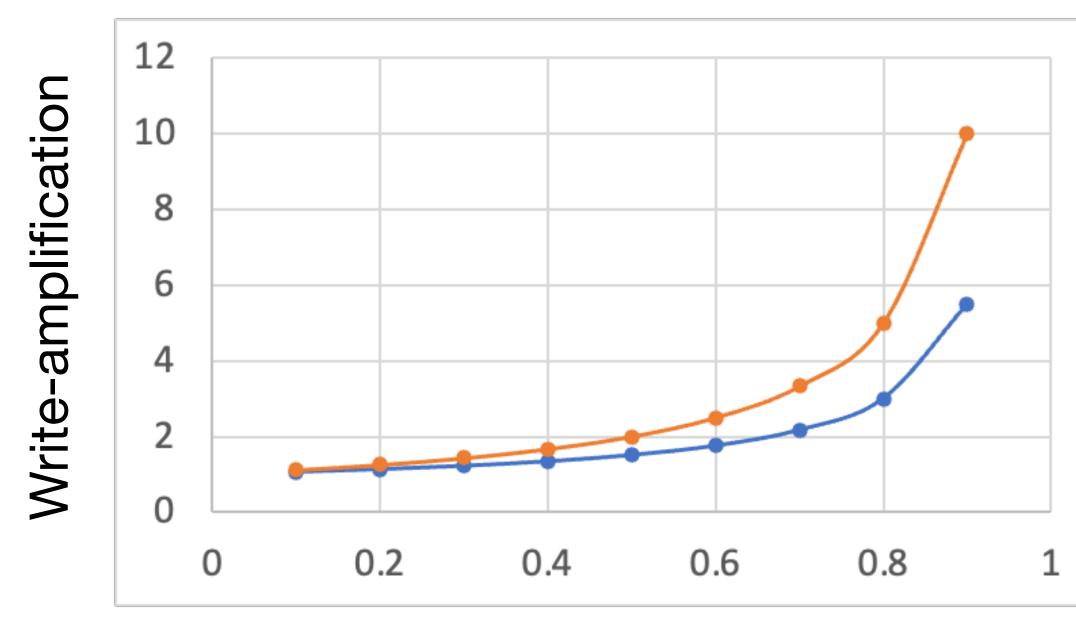
In addition, we can safeguard against counter overflows by not incrementing counters that have reached their maximum value.

1 2 1 4 0 2 0 0 0 0 3 0

1 4 0 1 3 2 0 0 0 2 4 4

0 2 2 1 3 0 4 0 3 0 1 0





Valid data / physical space

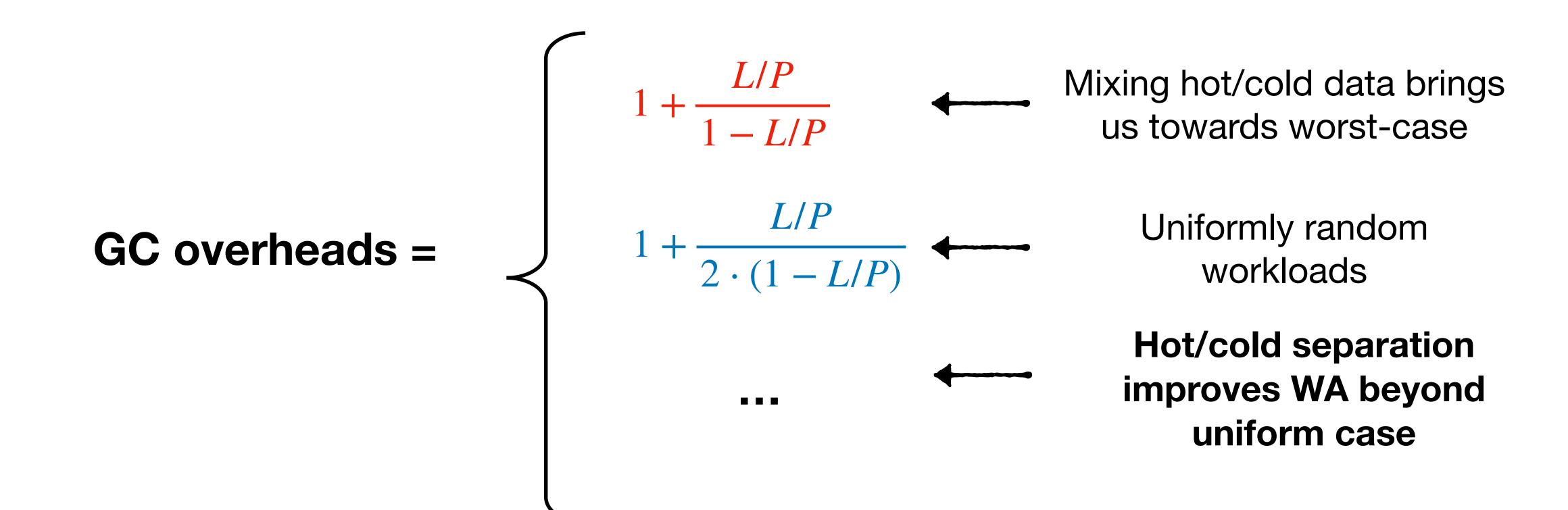


$$1 + \frac{L/P}{2 \cdot (1 - L/P)} \longleftarrow$$

Uniformly random

workloads

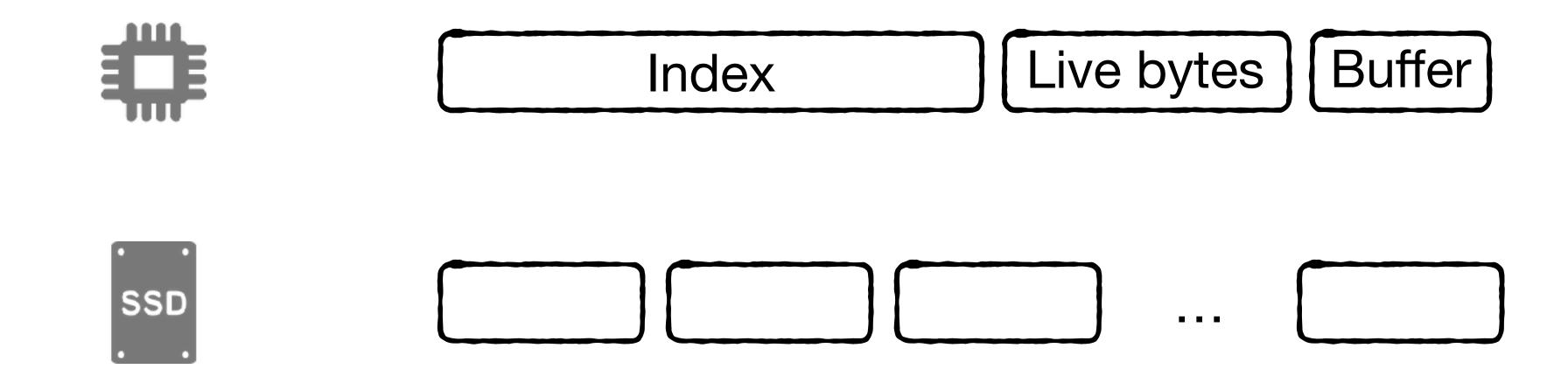
Hot/cold separation improves WA beyond uniform case



Overall Cost Analysis

Write cost: O(GC/B)I/O

Reads: 1 read I/O

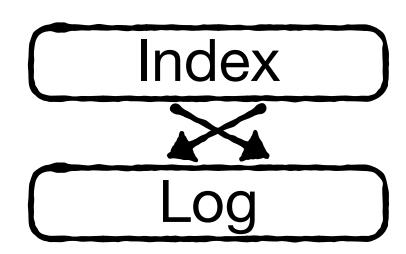


Mechanics

Hot/Cold separation

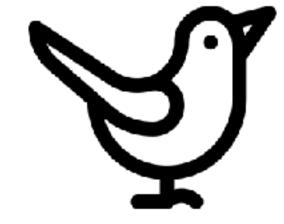
Checkpointing & Recovery

Cuckoo filtering







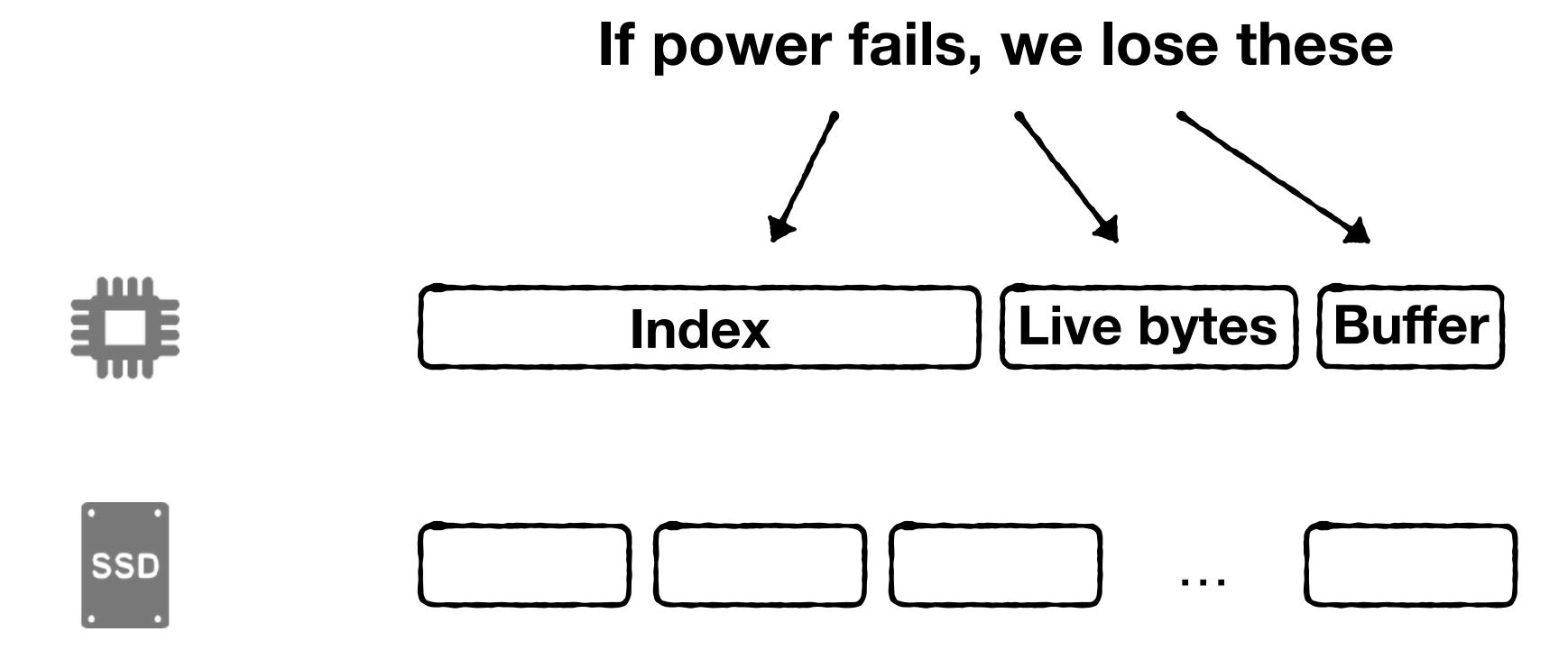


To reduce write-amp

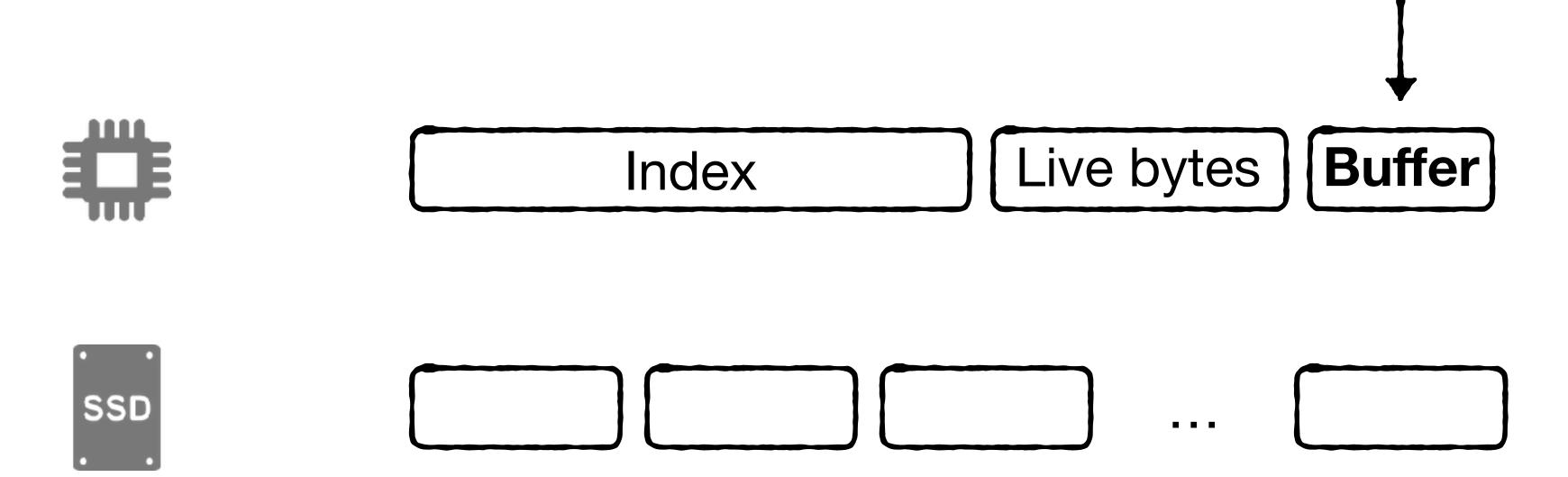
If power fails

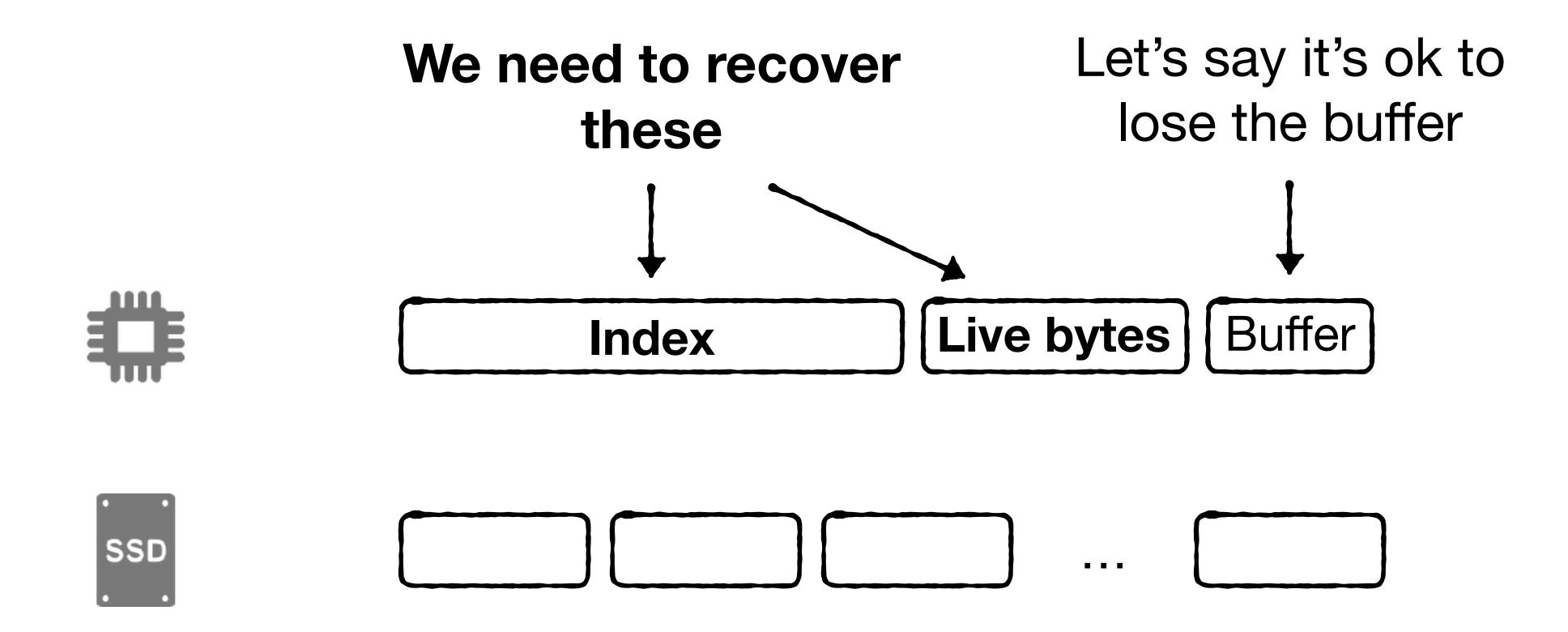
To reduce memory

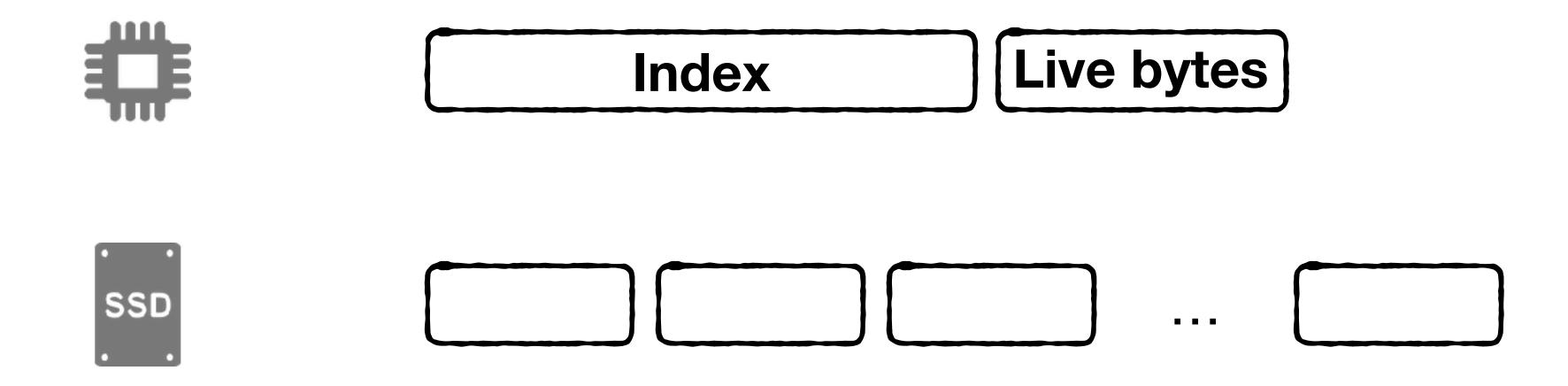
Checkpointing & Recovery

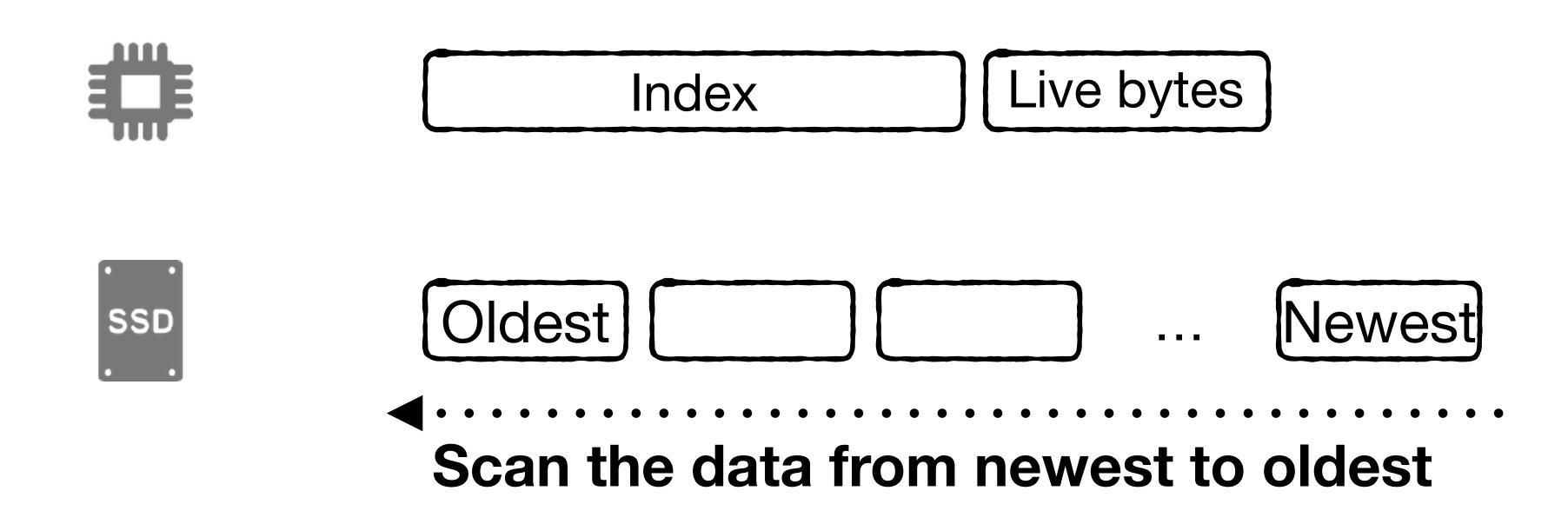


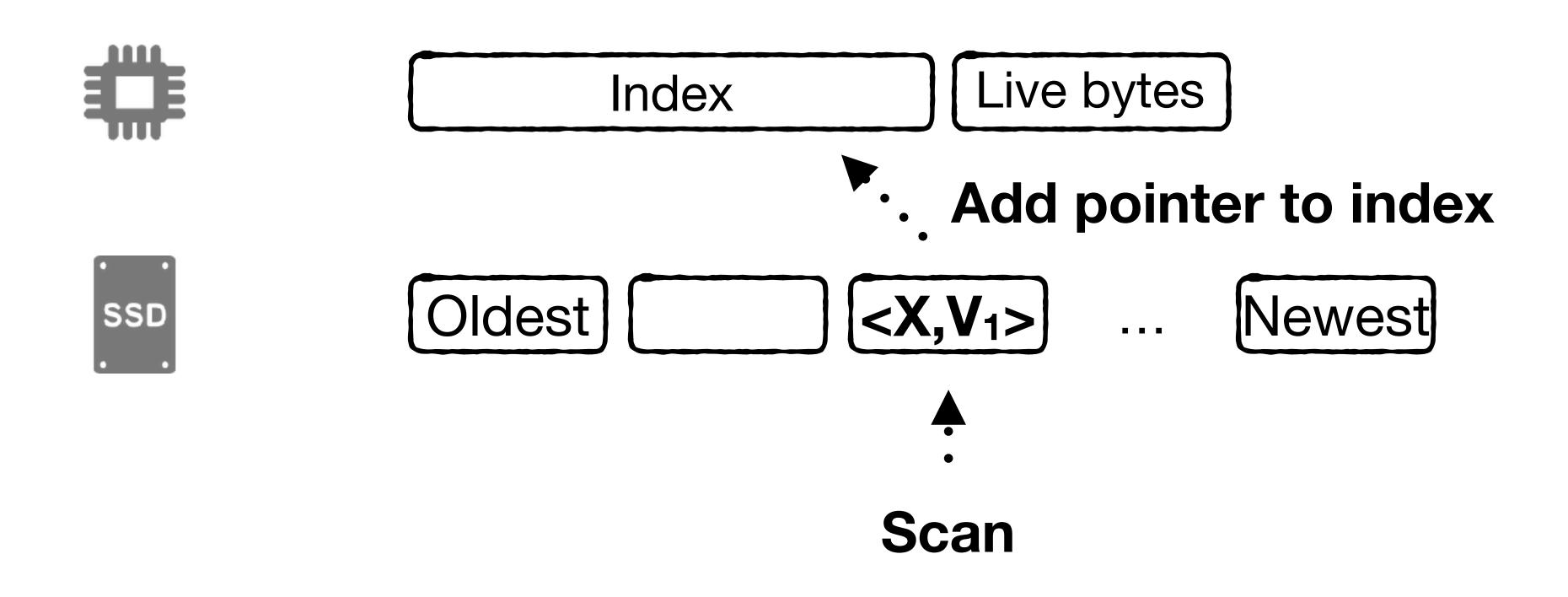
Let's say it's ok to lose the buffer

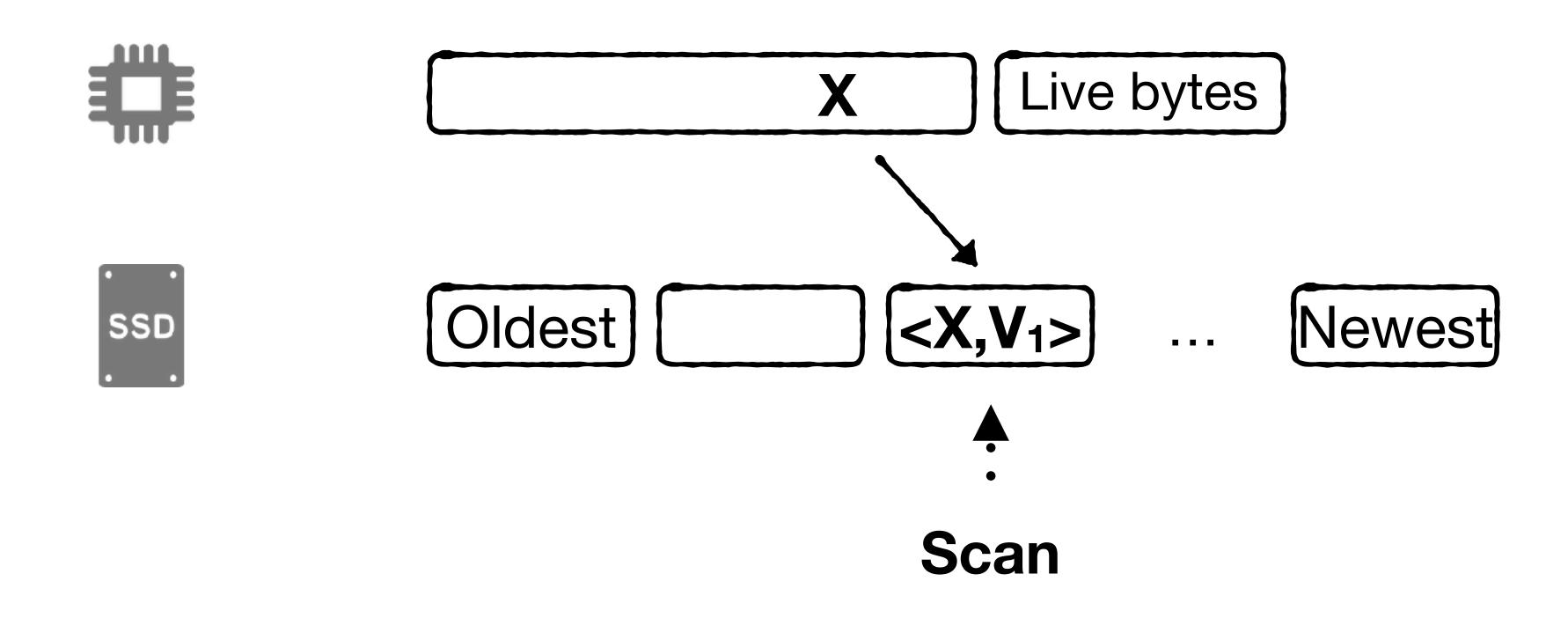


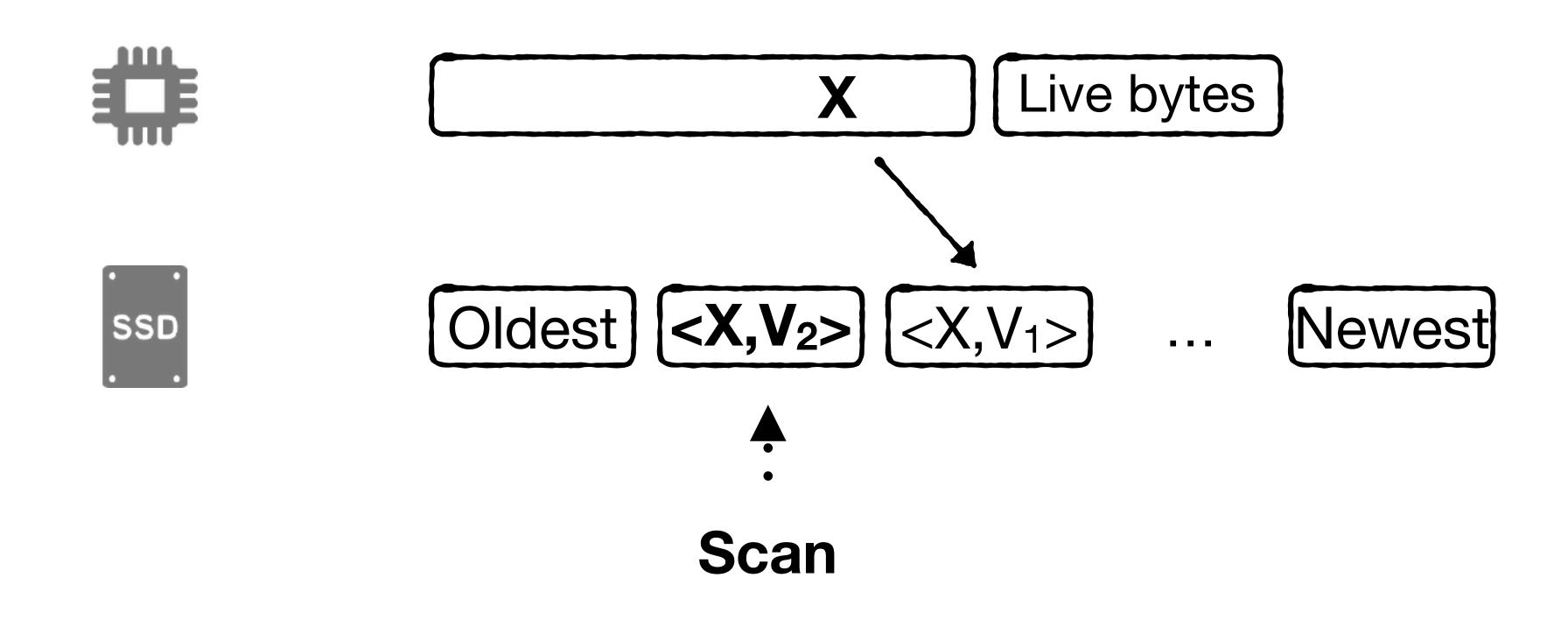


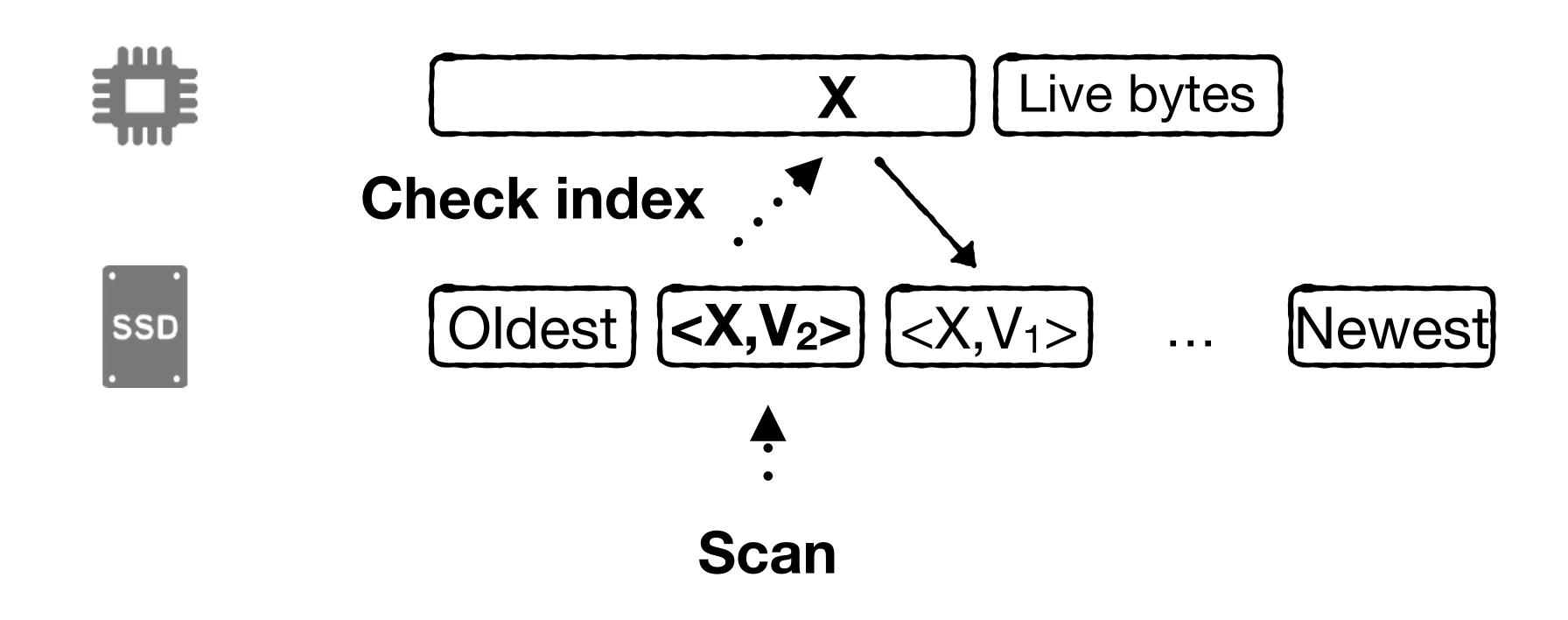




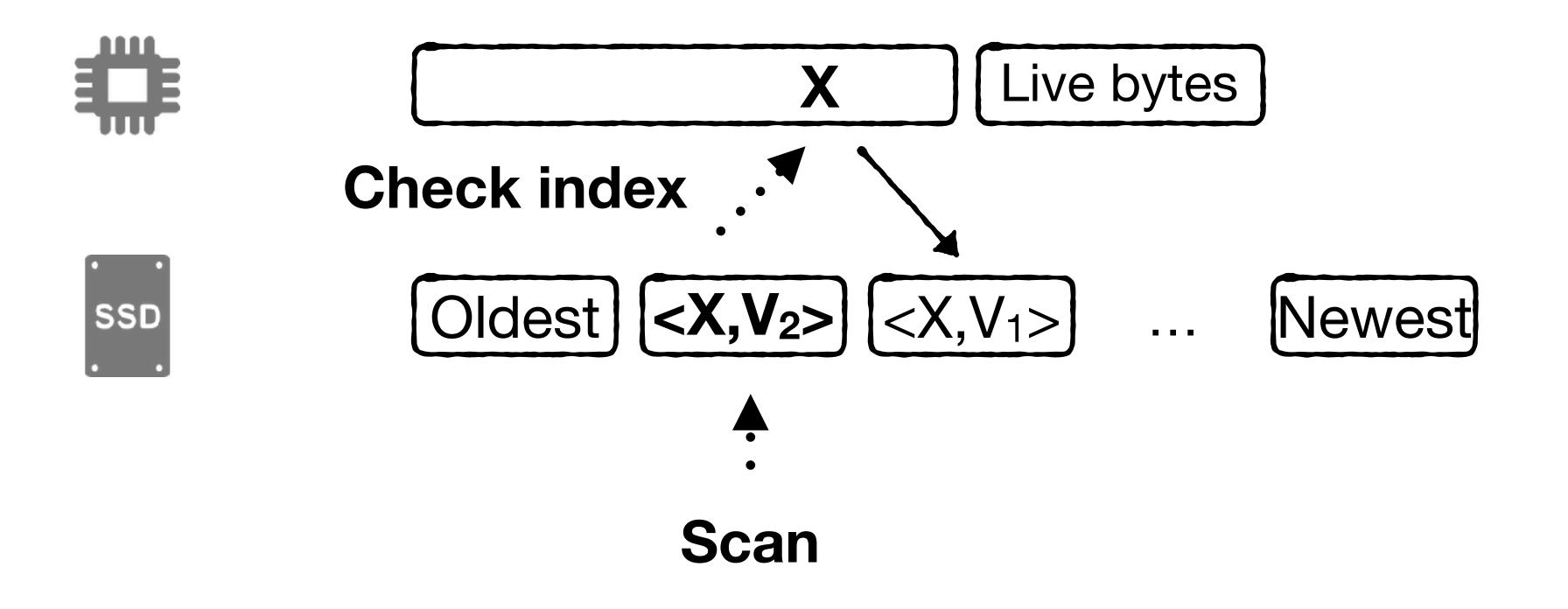




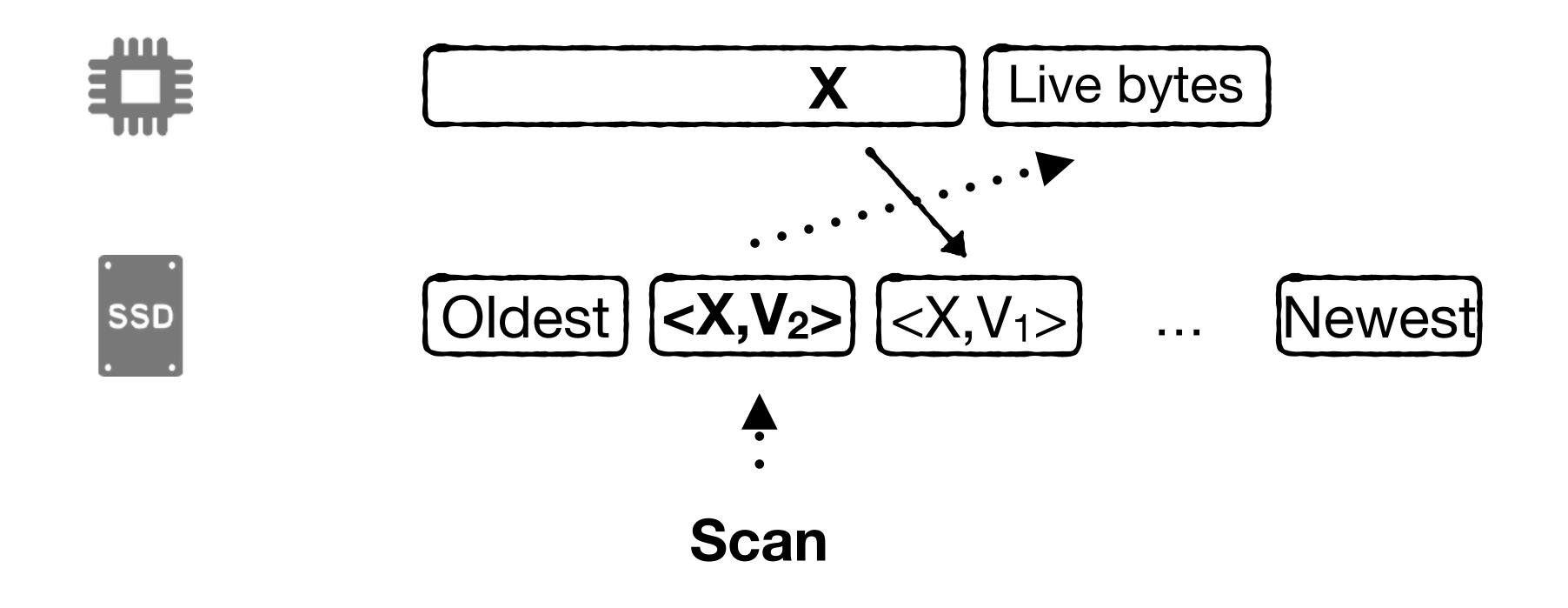




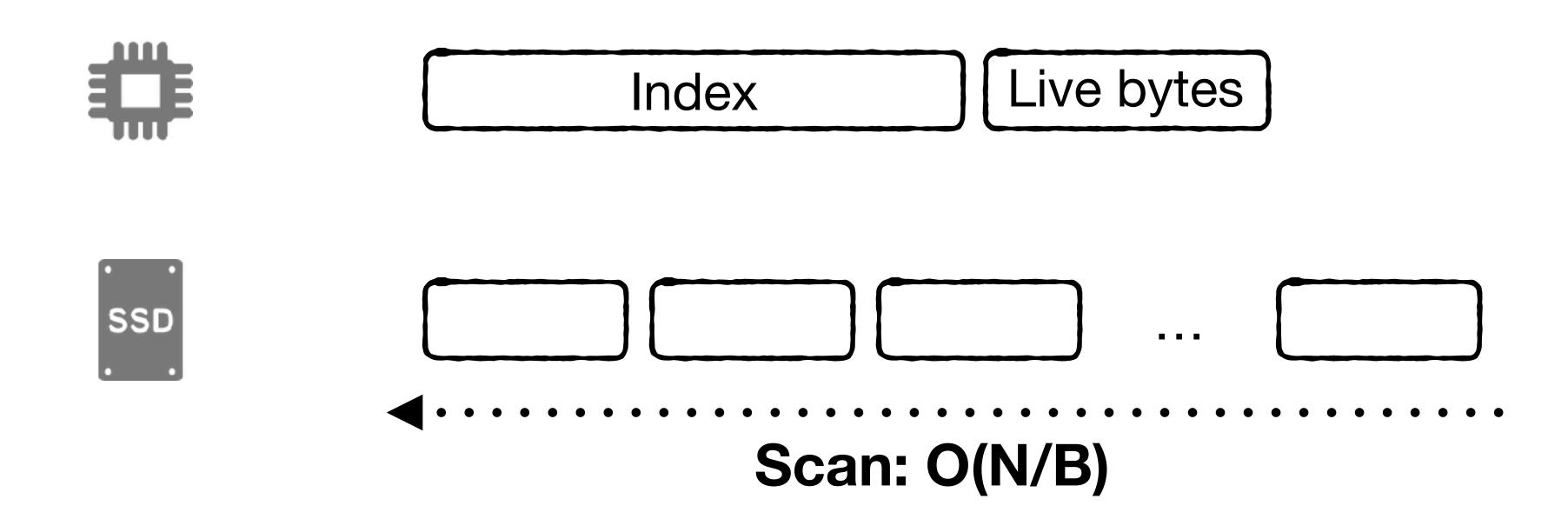
Since a mapping entry with the same key already exists, the current entry is obsolete.



So we instead subtract size of entry from live bytes.

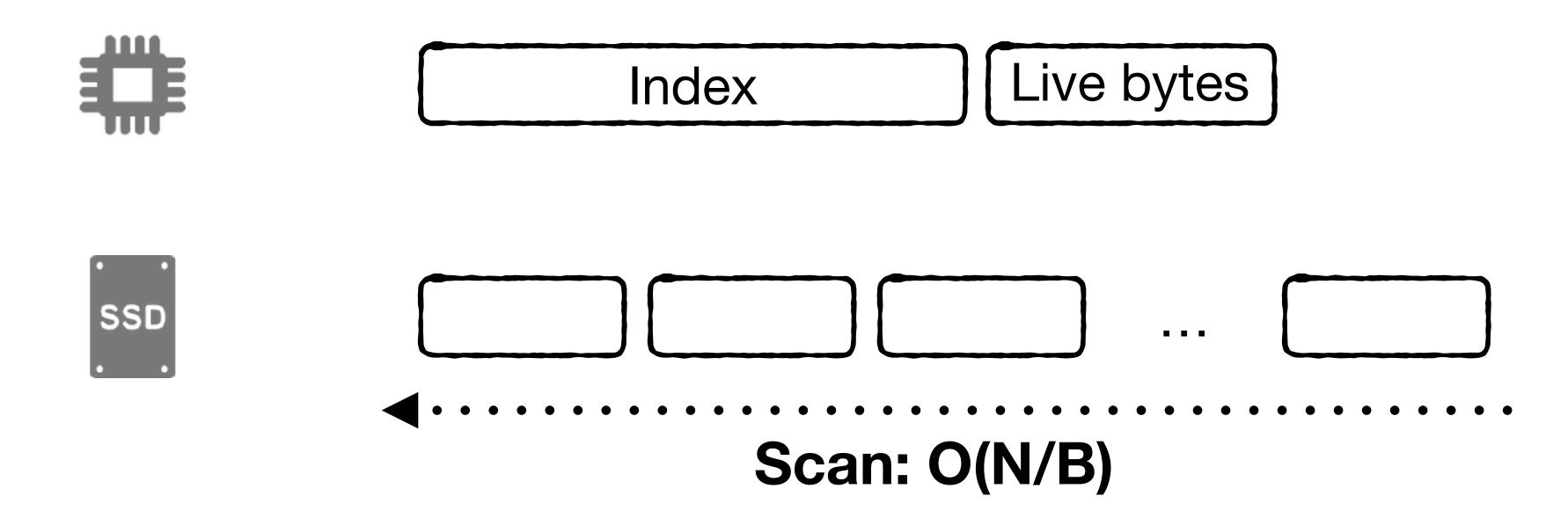


We can recover with one pass over the data



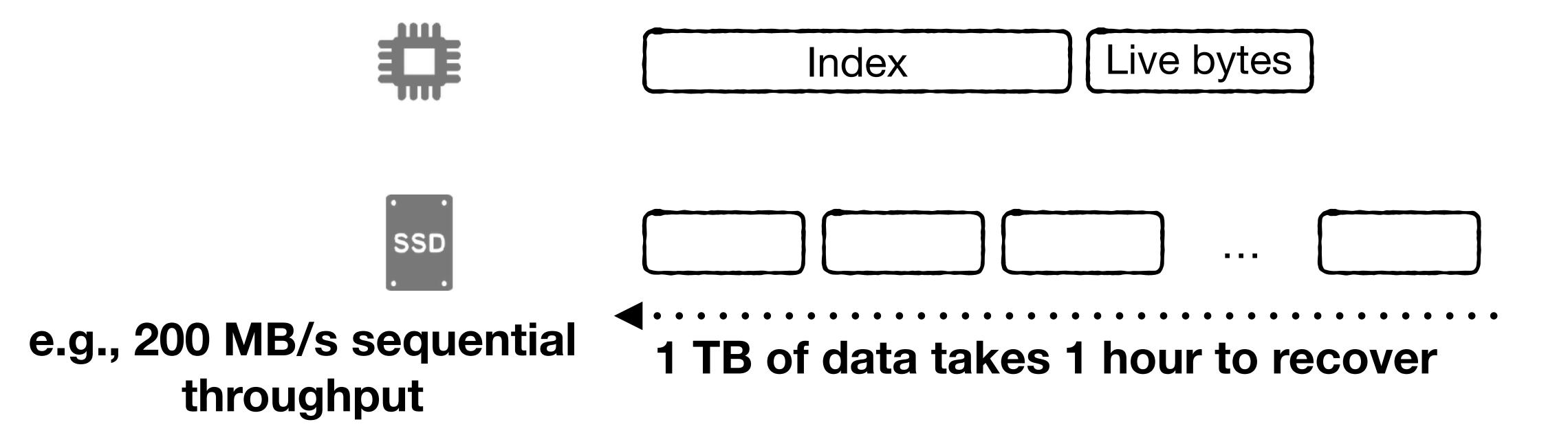
We can recover with one pass over the data

But if the data is huge, this can take a while.



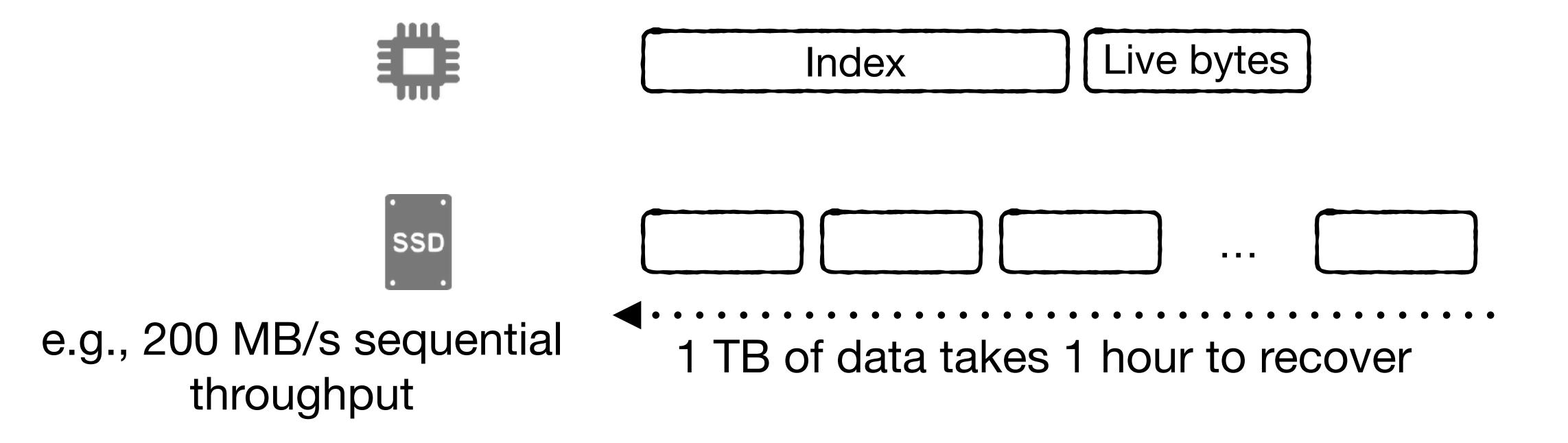
We can recover with one pass over the data

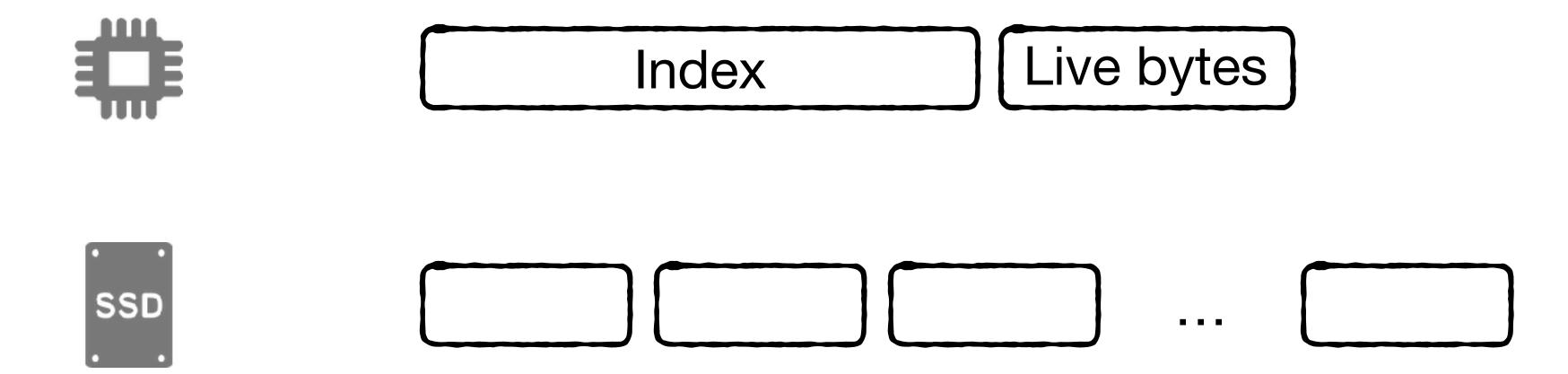
But if the data is huge, this can take a while.



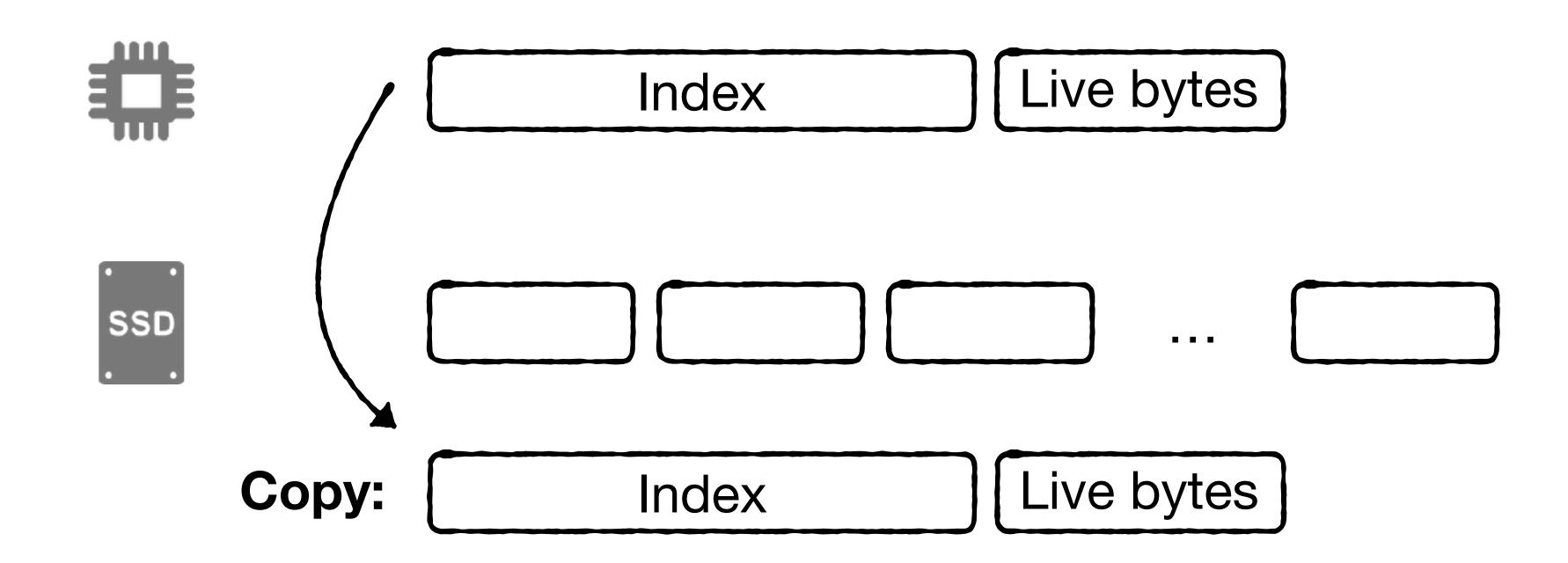
We can recover with one pass over the data But if the data is huge, this can take a while.

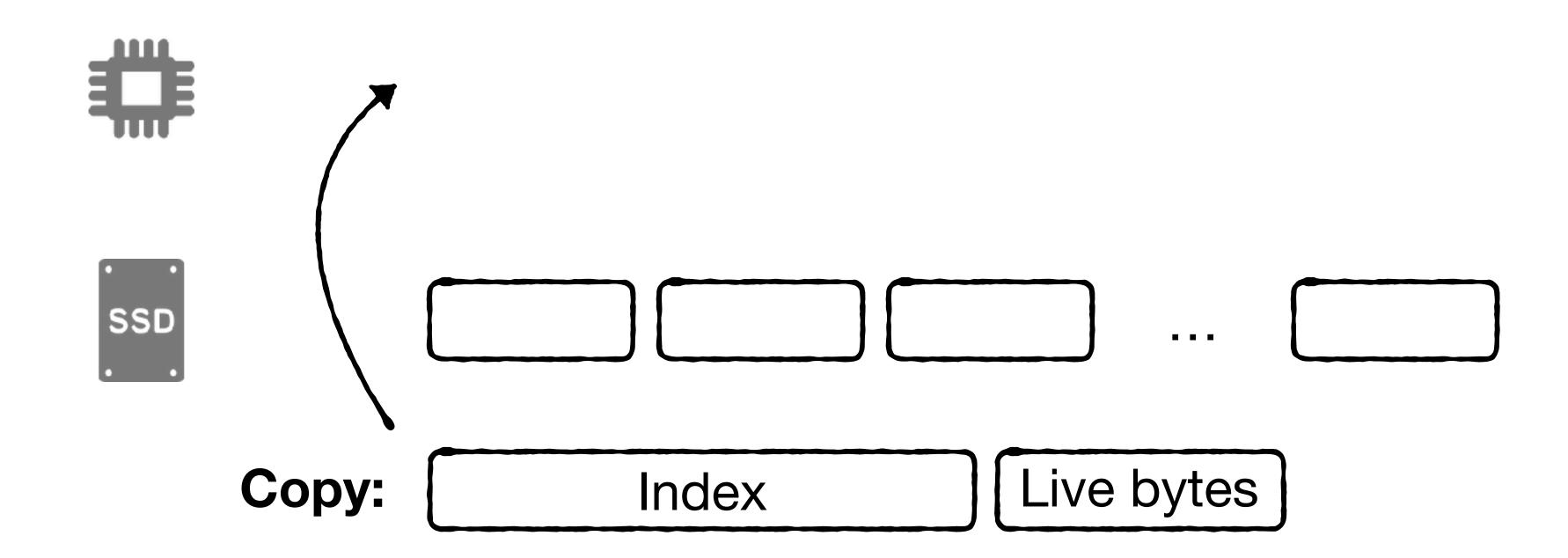
Any good ideas?

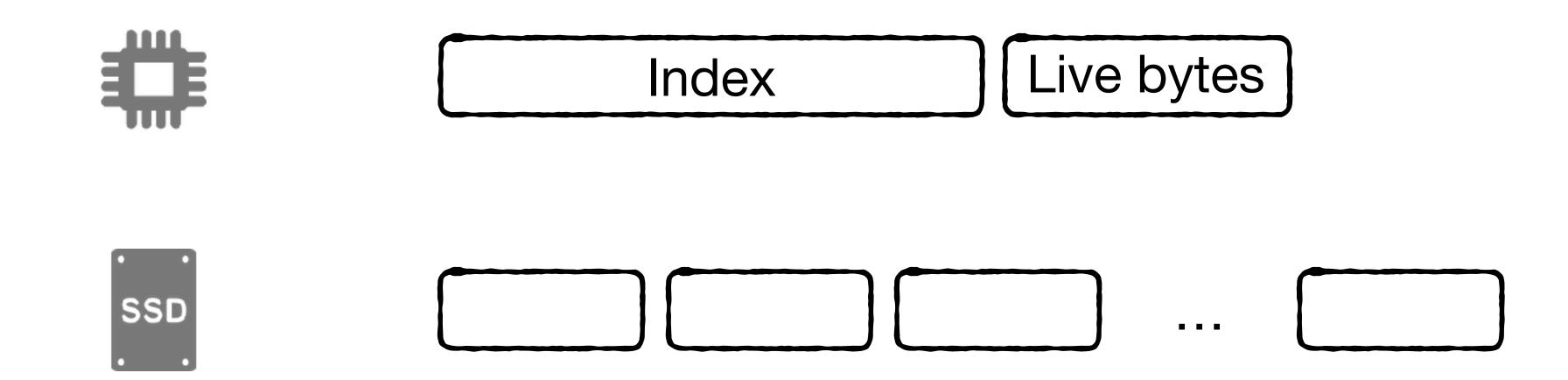


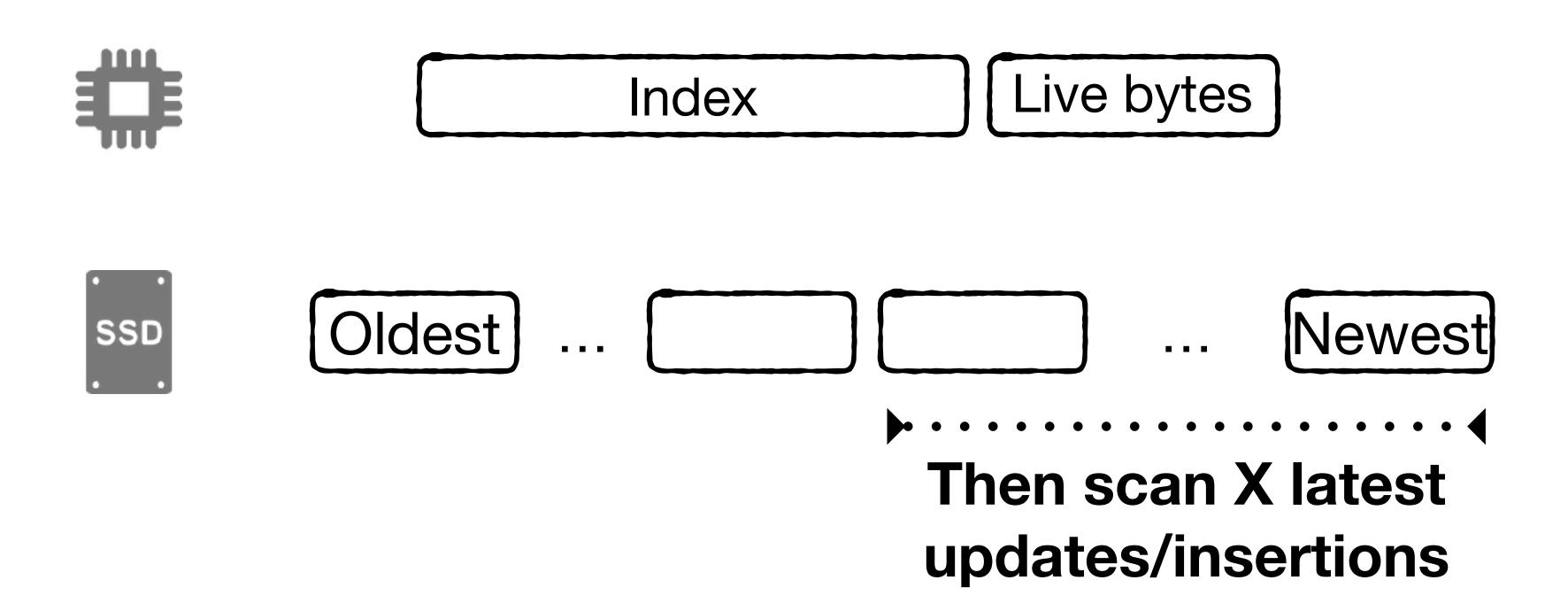


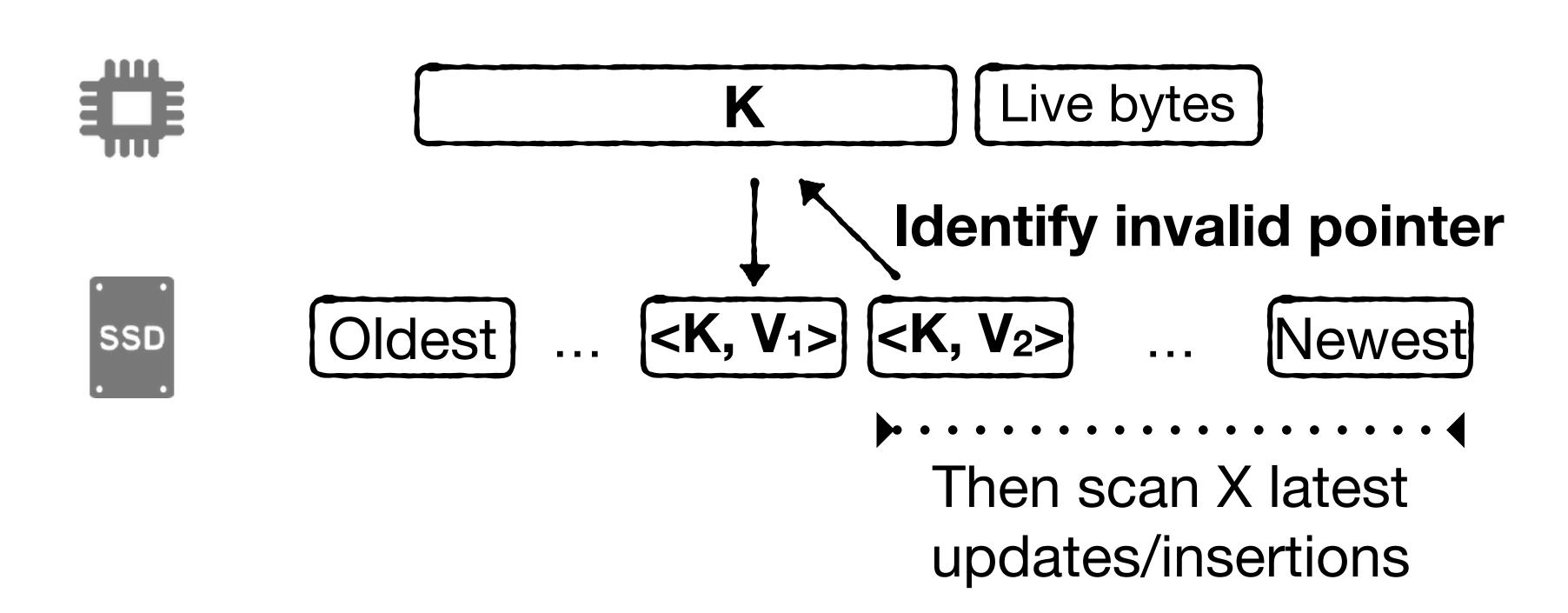
Every X updates/insertions, store copy of index & live bytes

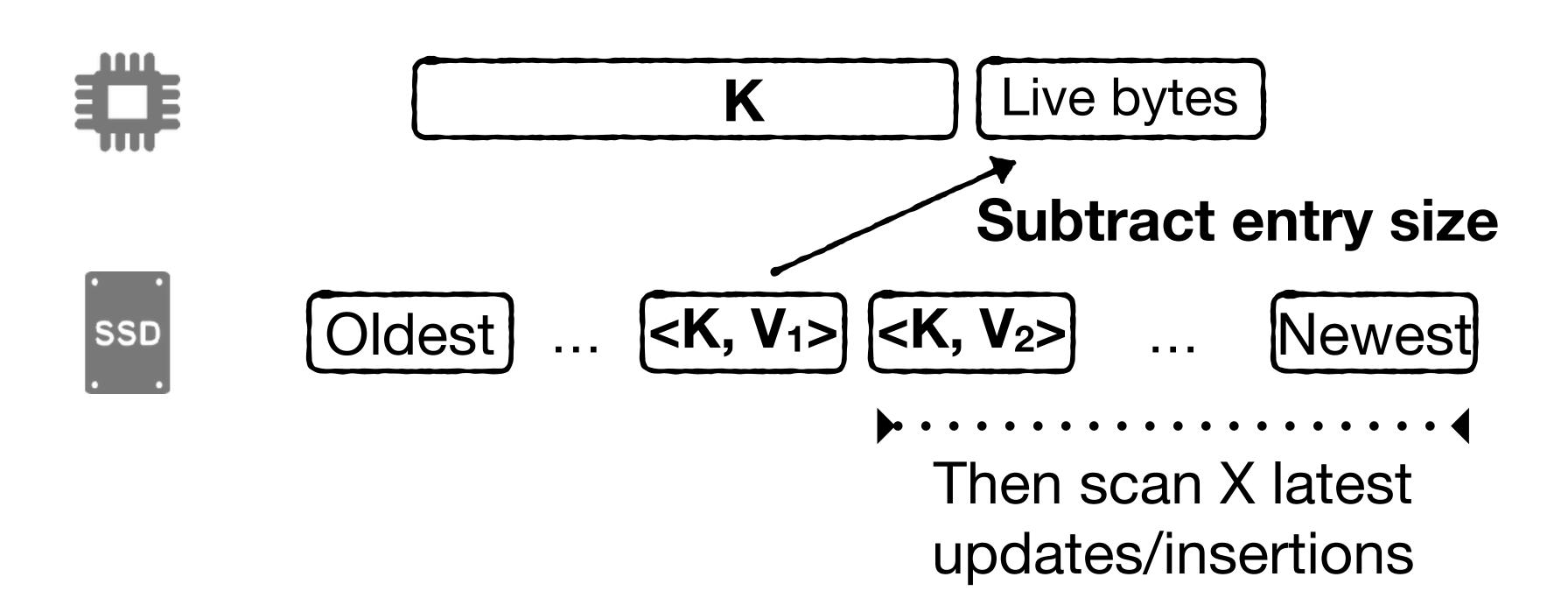


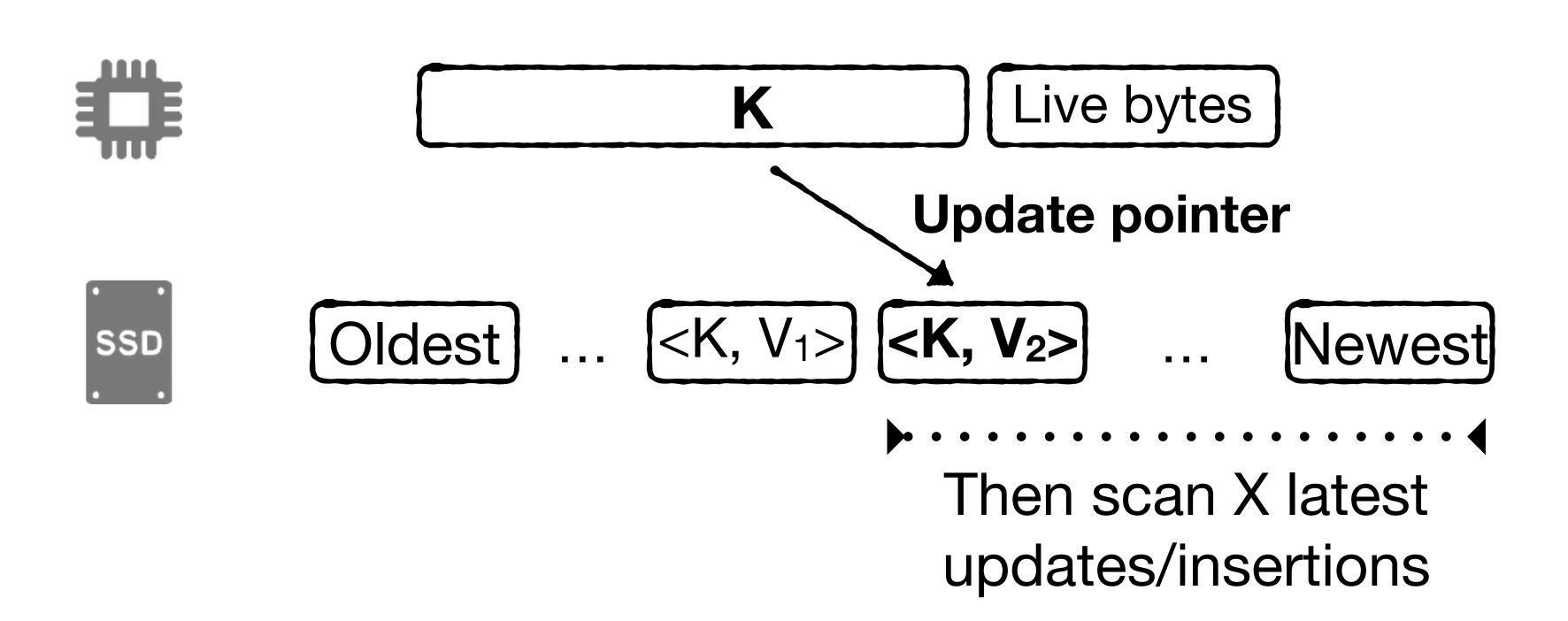




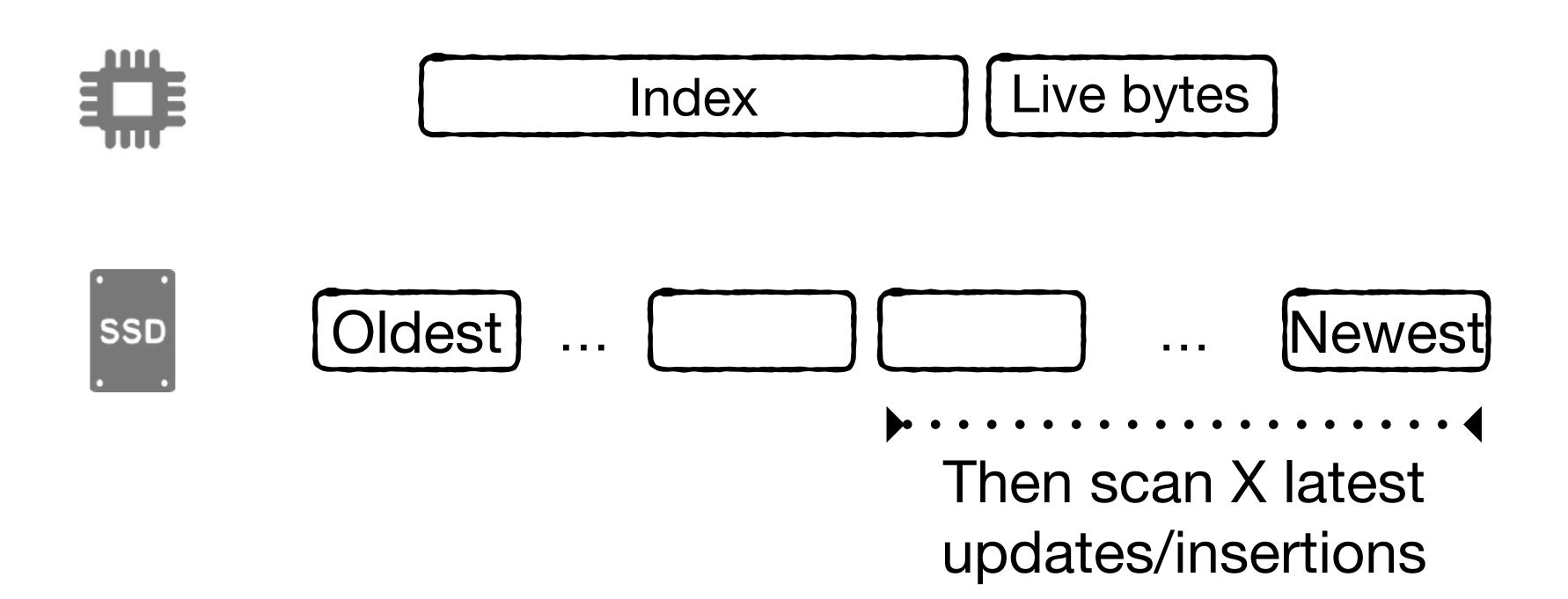






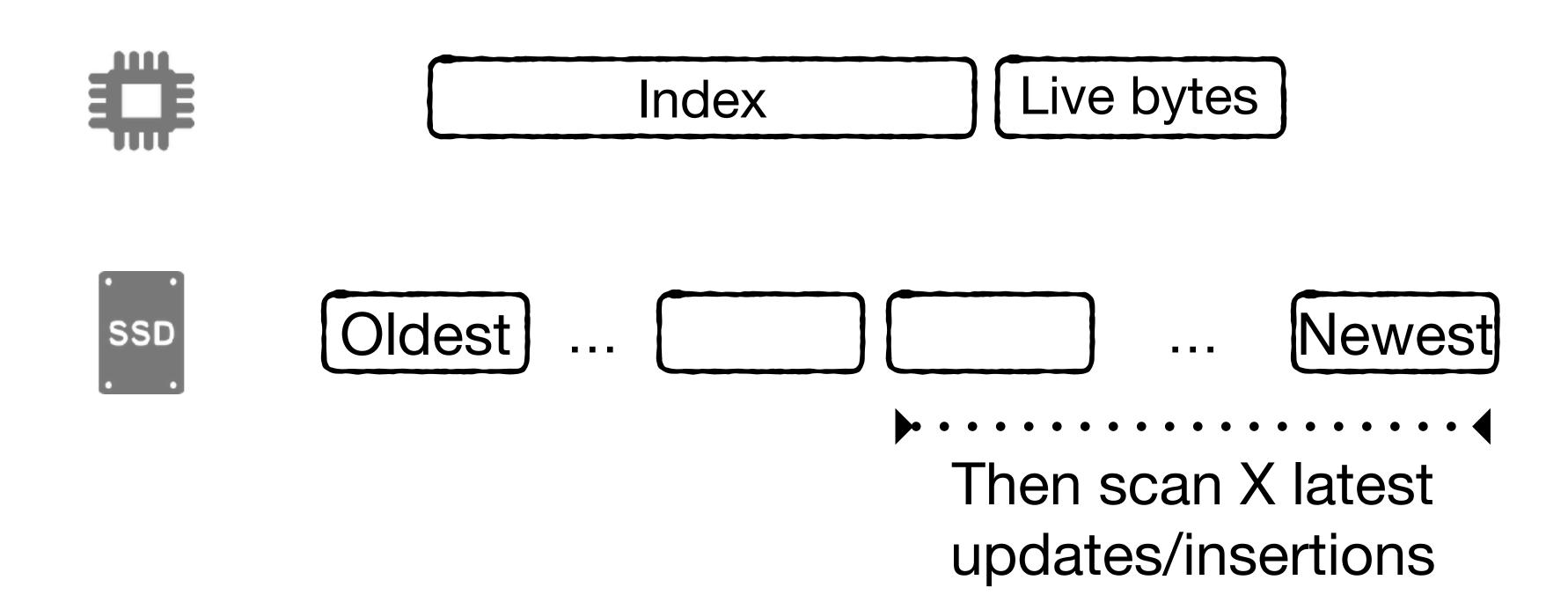


The frequency of checkpointing controls a trade-off between write cost and recovery time.



Additional write cost: O(index size / X)

Recovery time: O(X/B) reads for backwards scan

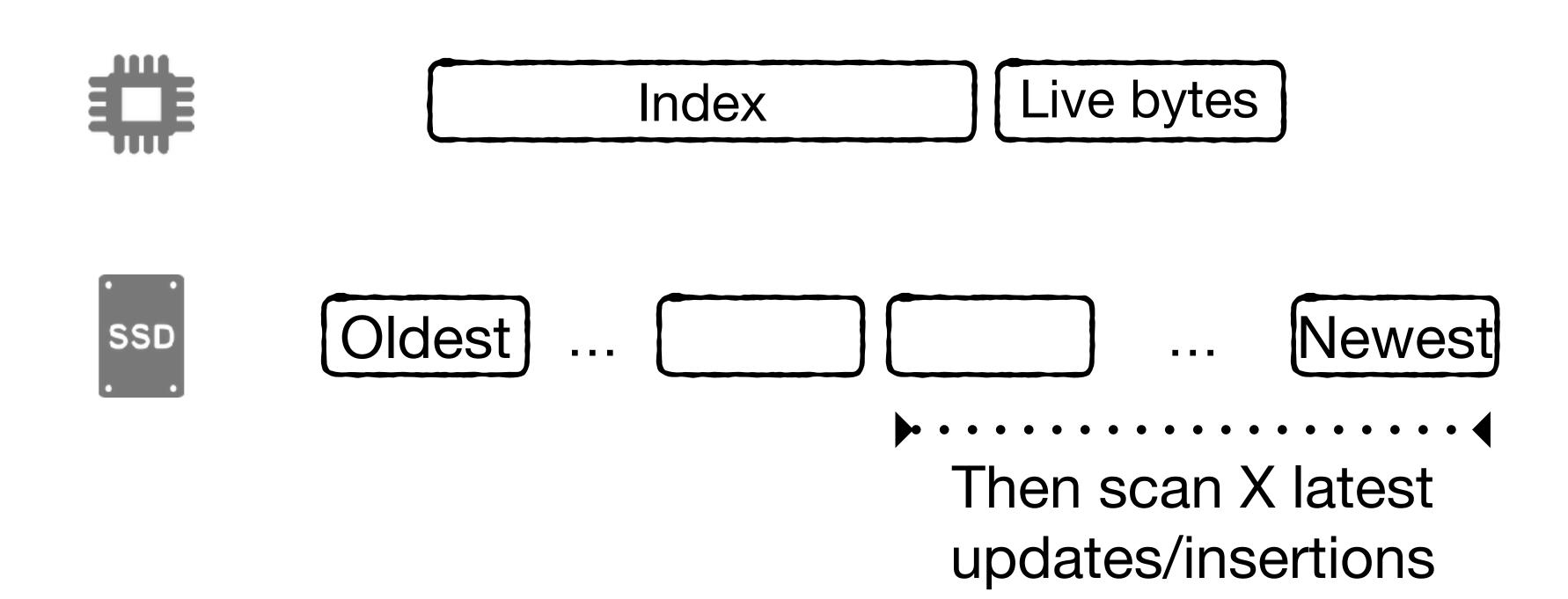


Costs Summary

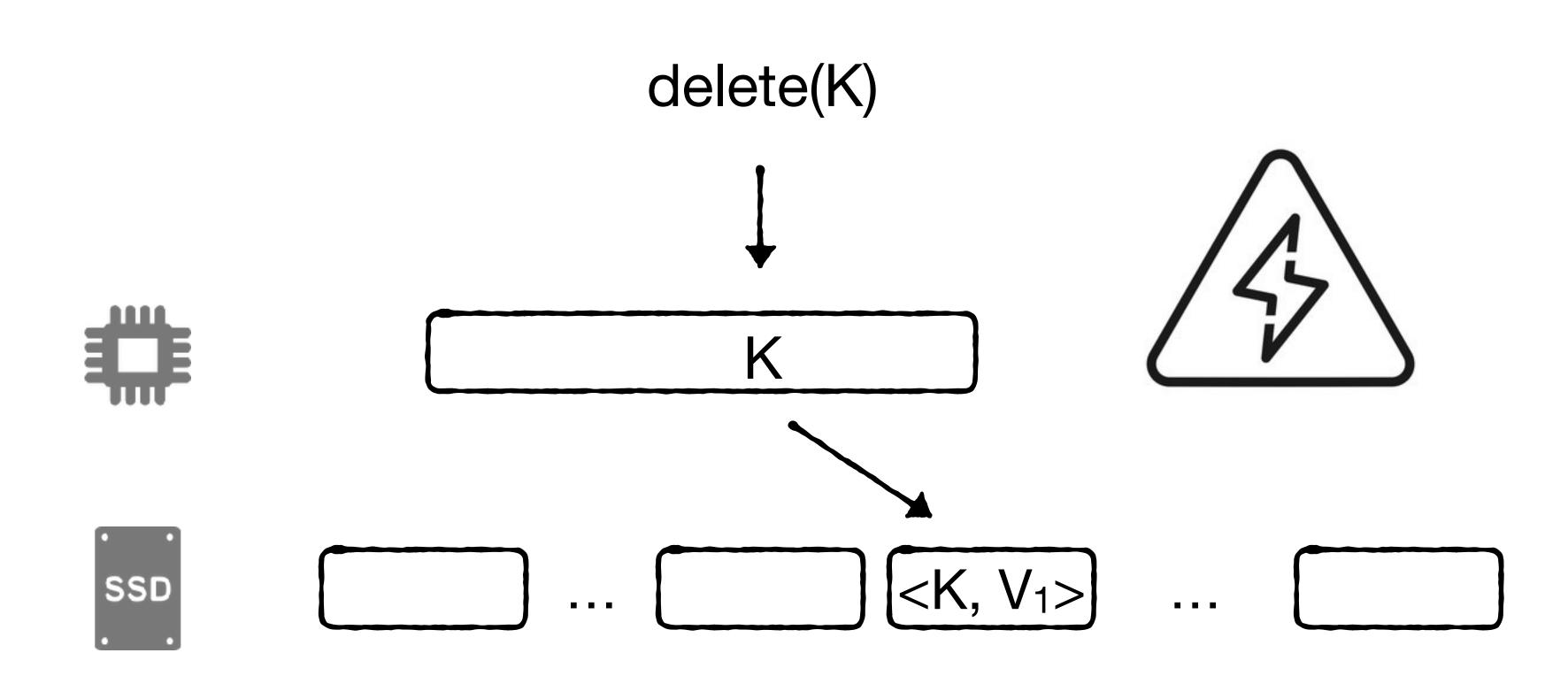
write cost: O(GC/B + index size / X)

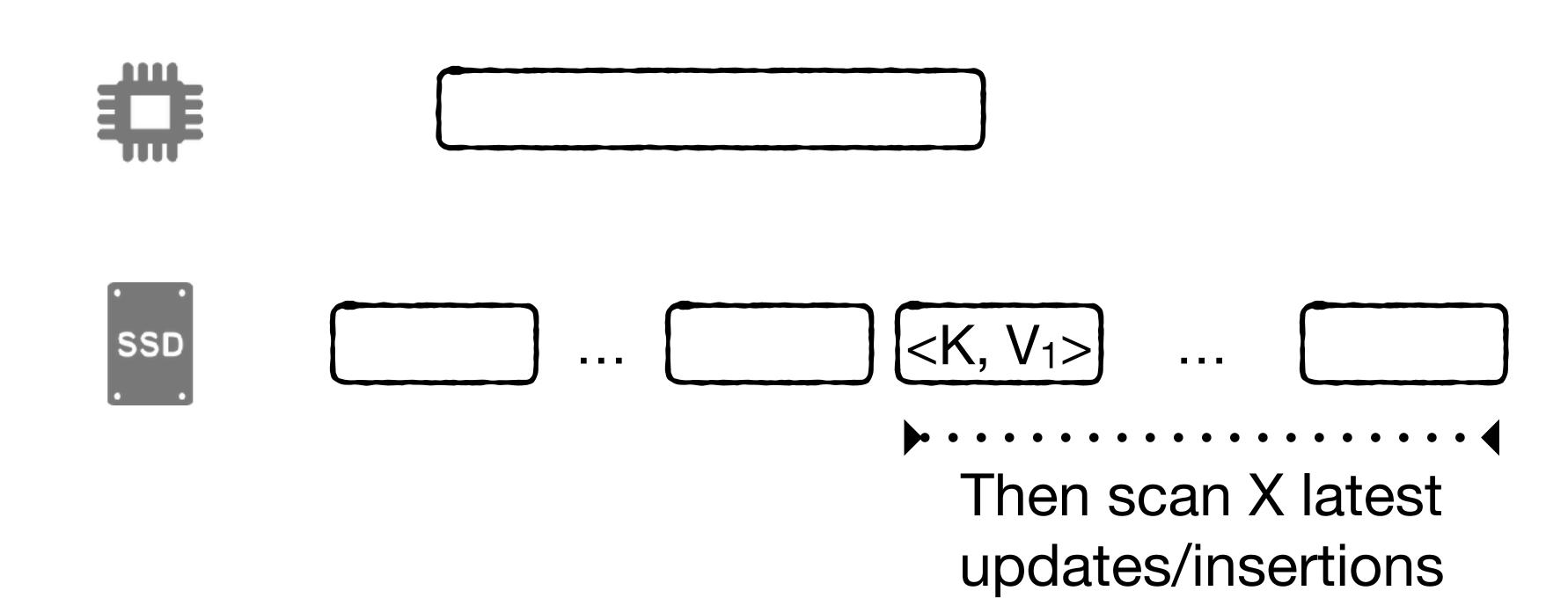
Recovery time: O(X/B) reads for backwards scan

Get cost: O(1)

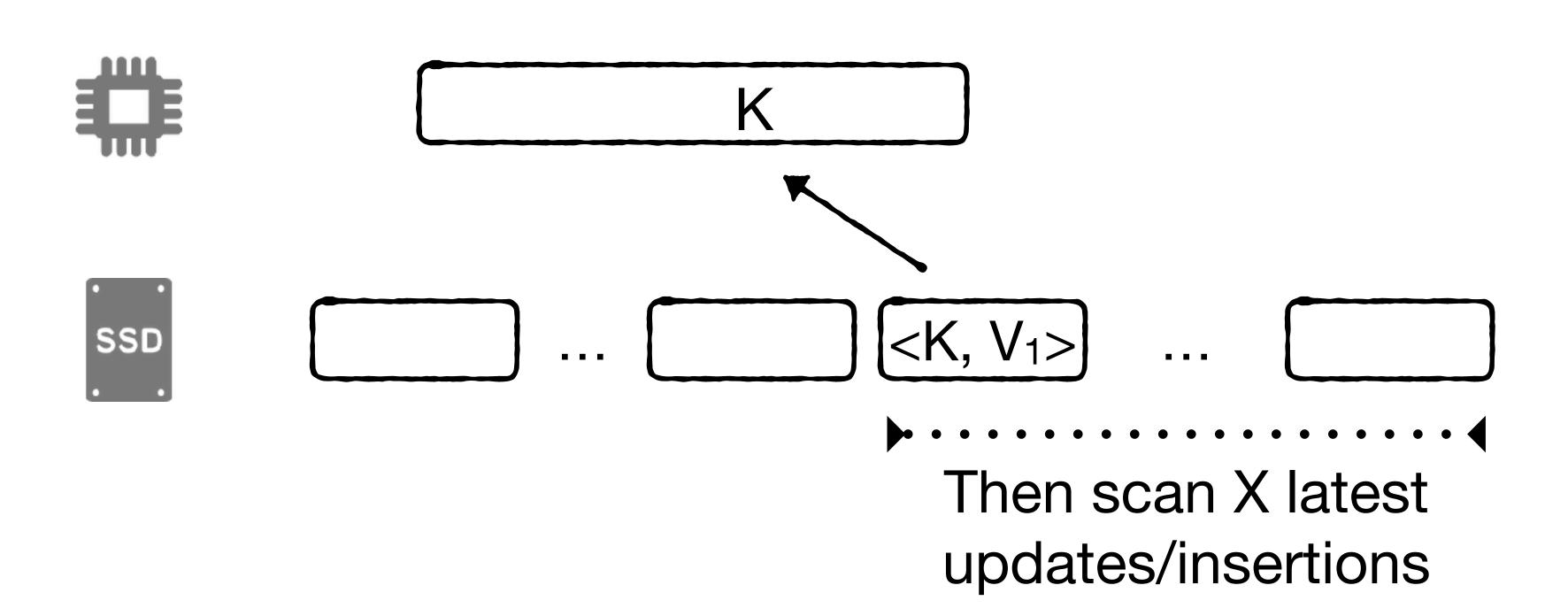


Suppose we delete an entry and then power fails. Can you spot a problem?

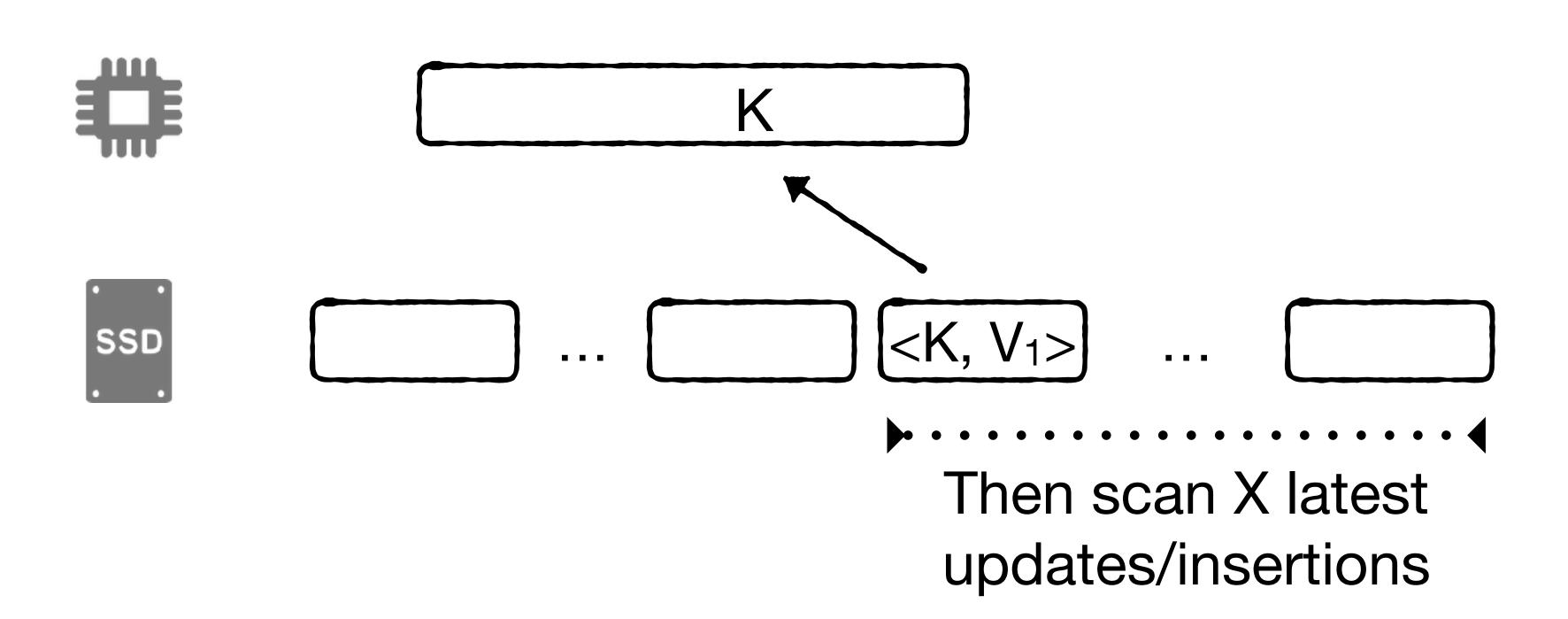




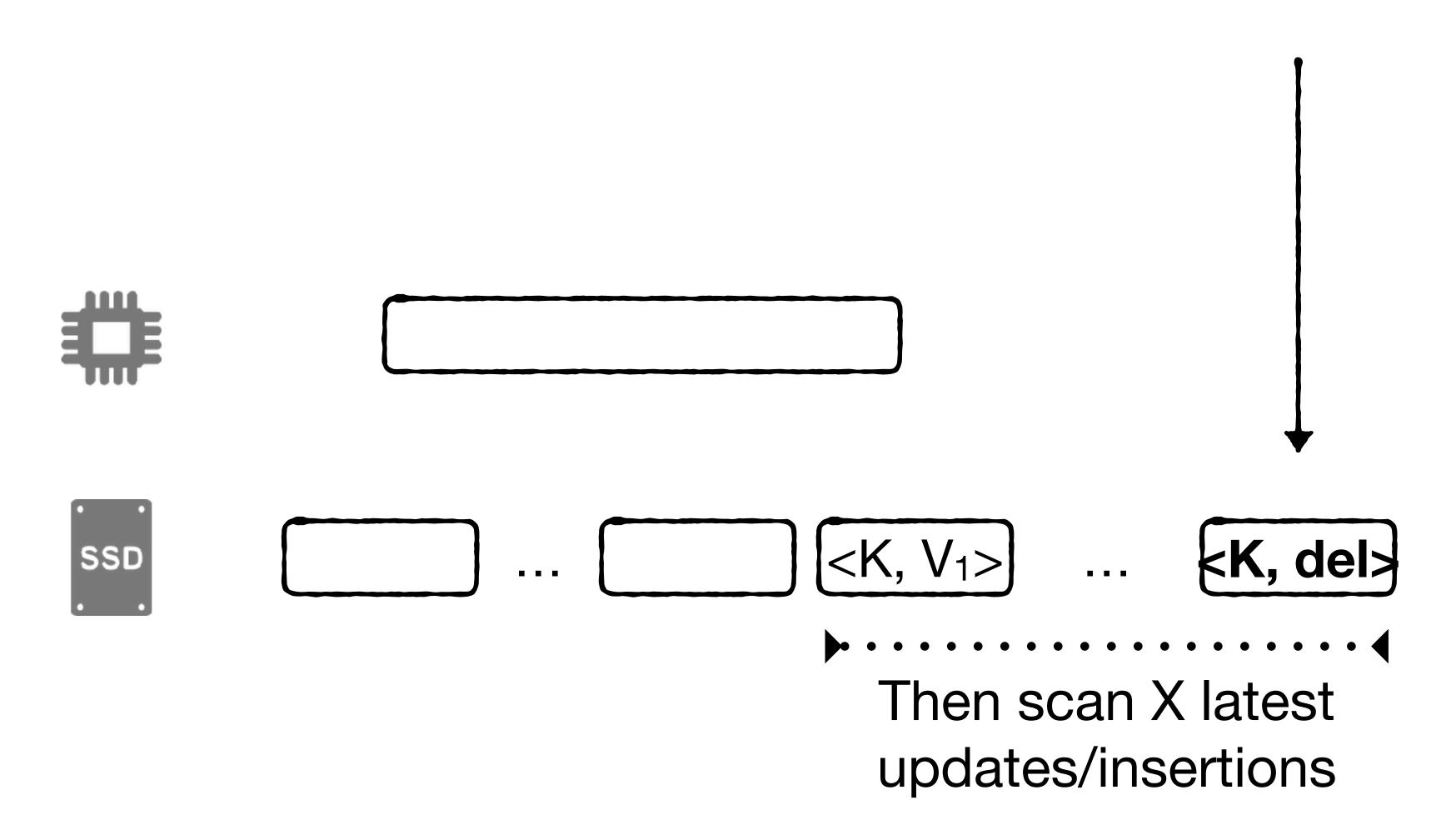
May load deleted entry back to index during recovery



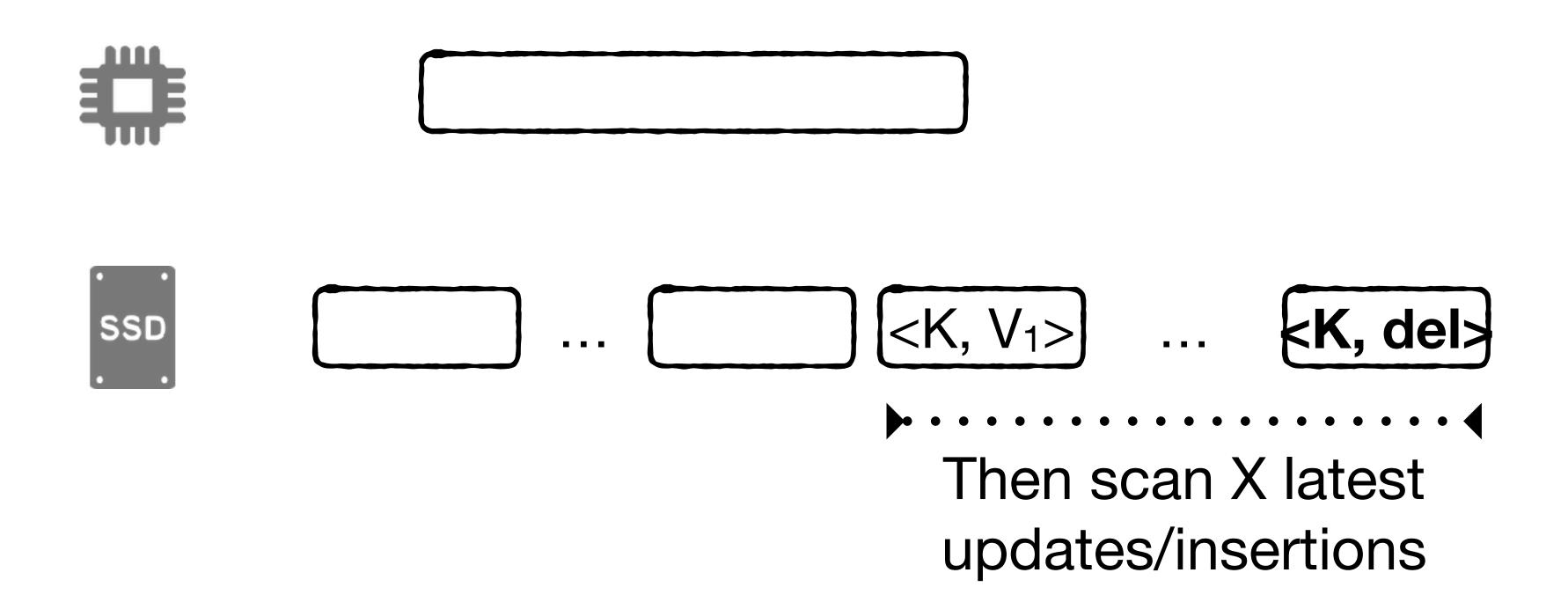
May load deleted entry back to index during recovery **Solution?**



Insert tombstone when deleting



While recovering, if the first instance of a key we see is a tombstone, we ignore all subsequent entries with this key.

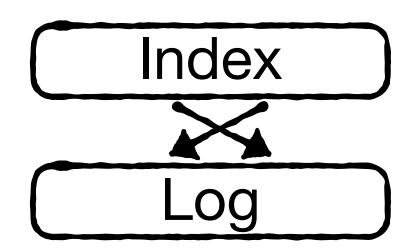


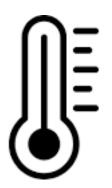
Mechanics

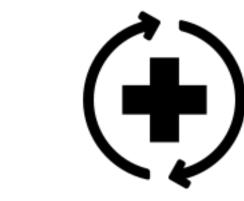
Hot/Cold separation

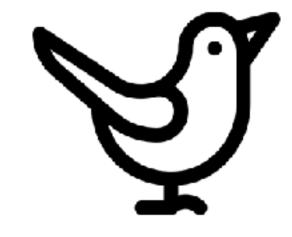
Checkpointing & Recovery

Cuckoo filtering





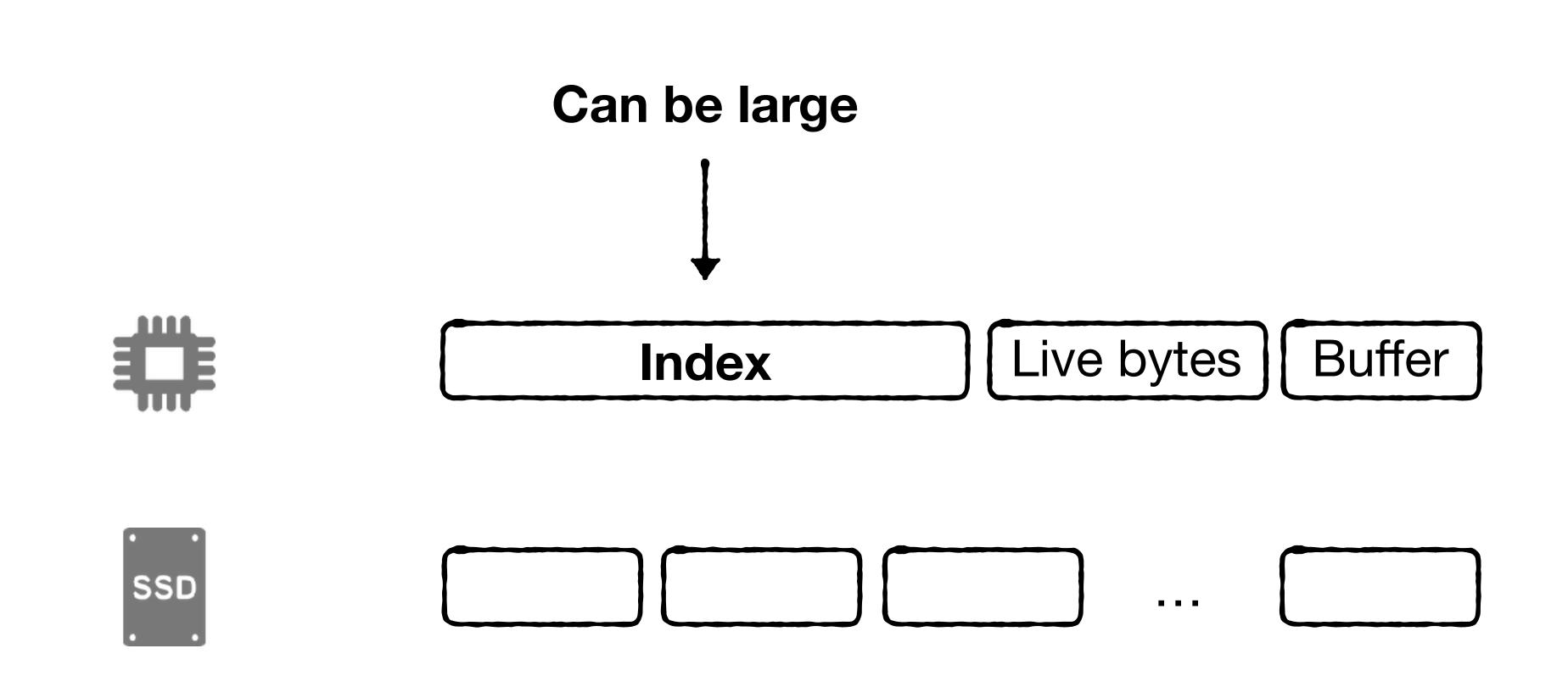


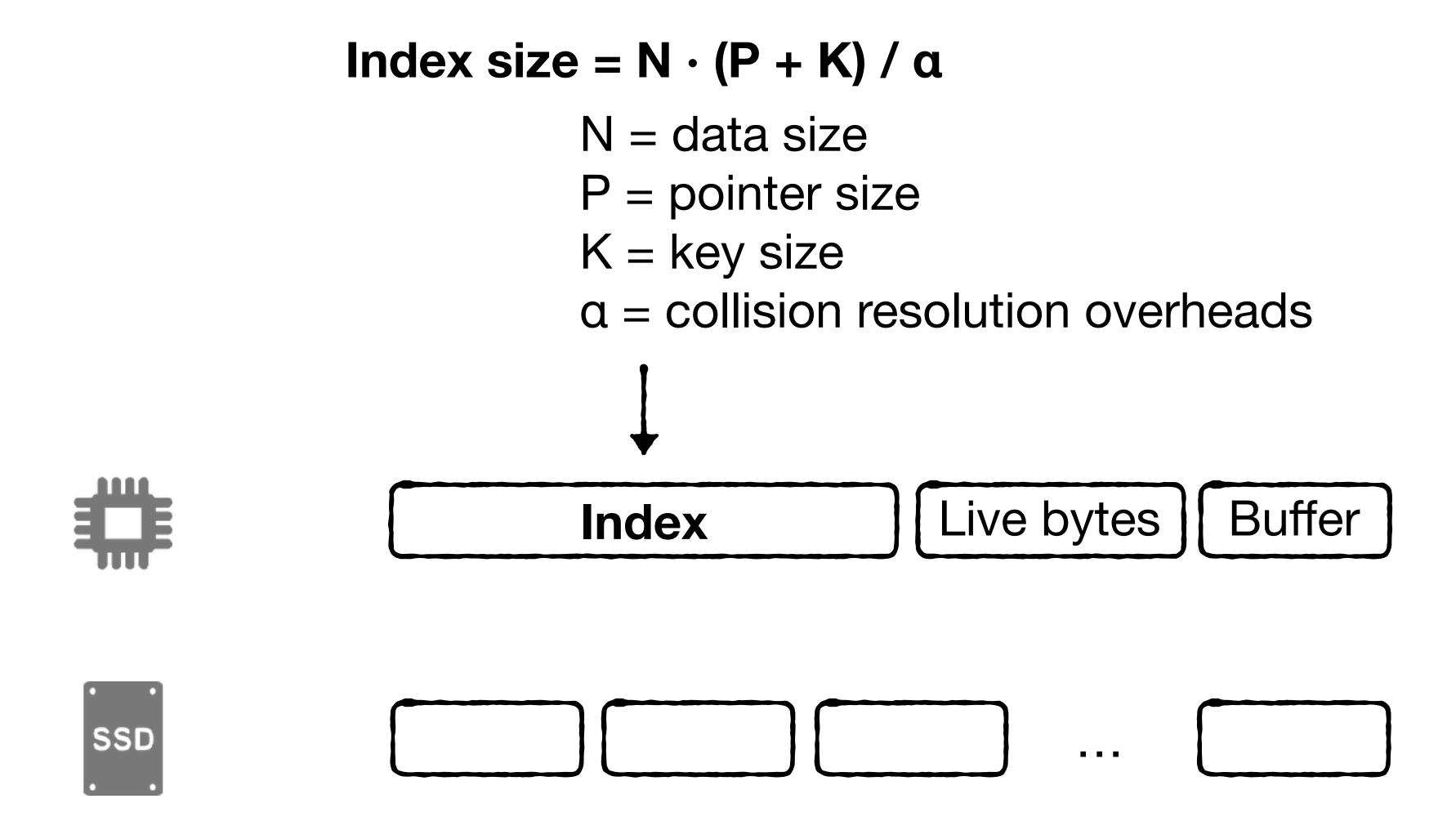


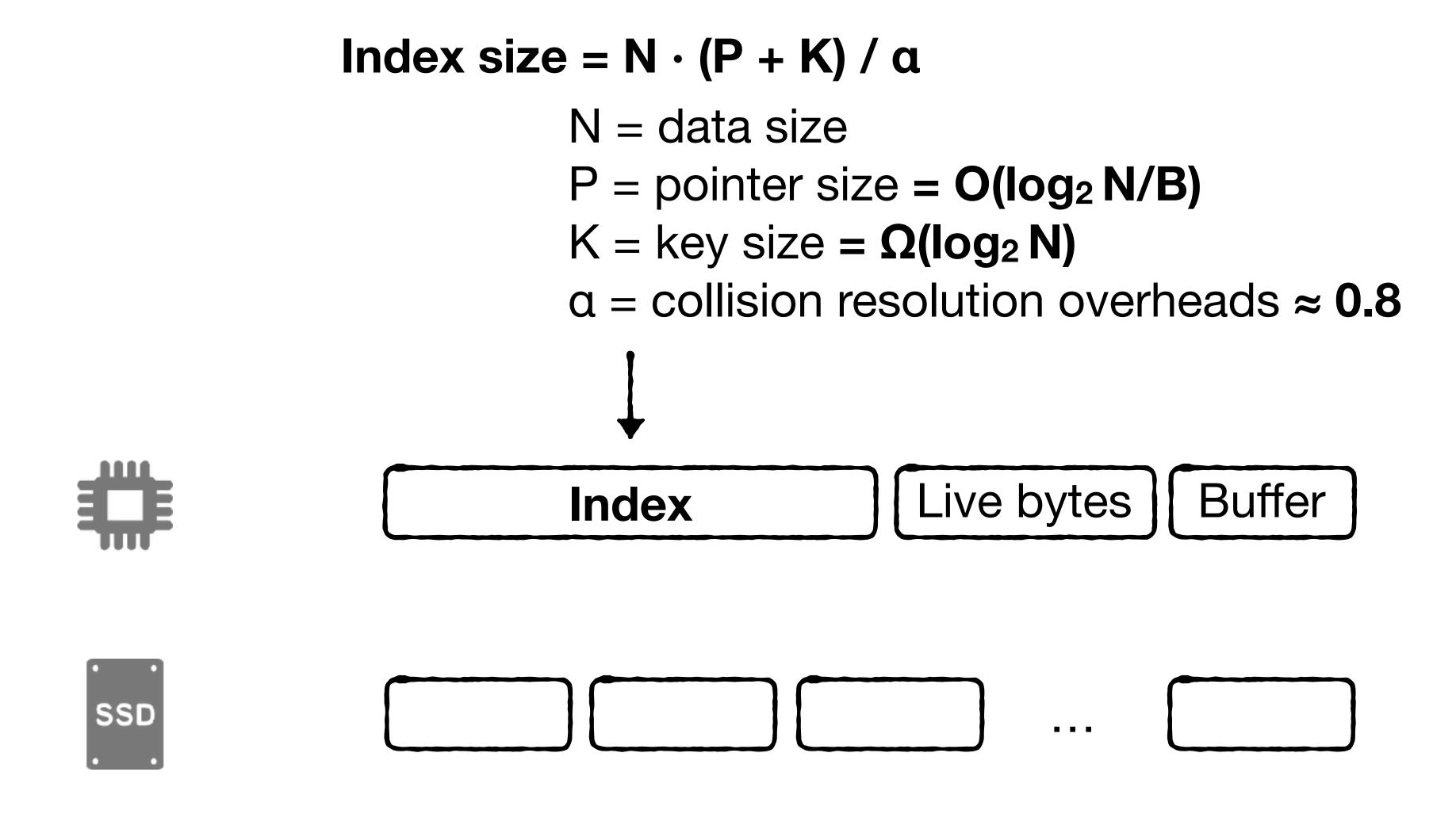
To reduce write-amp

If power fails

To reduce memory







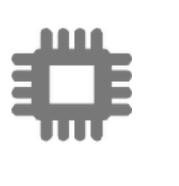
Index size = $N \cdot (P + K) / \alpha$

N = data size

 $P = pointer size = O(log_2 N/B)$

 \leftarrow K = key size = $\Omega(\log_2 N)$

 α = collision resolution overheads \approx 0.8



Ultimately attack this

with cuckoo filters

Index

Live bytes

Buffer





. . .

Index size = $N \cdot (P + K) / \alpha$

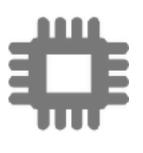
N = data size

 $P = pointer size = O(log_2 N/B)$

 $K = \text{key size} = \Omega(\log_2 N)$

But first attack this with cuckoo hashing

 α = collision resolution overheads \approx 0.8



Index

Live bytes

Buffer





. . .

Table of buckets, each containing 1 entry Two hash functions

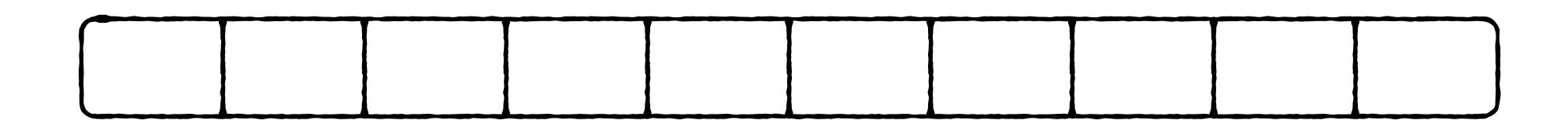


Table of buckets, each containing 1 entry Two hash functions

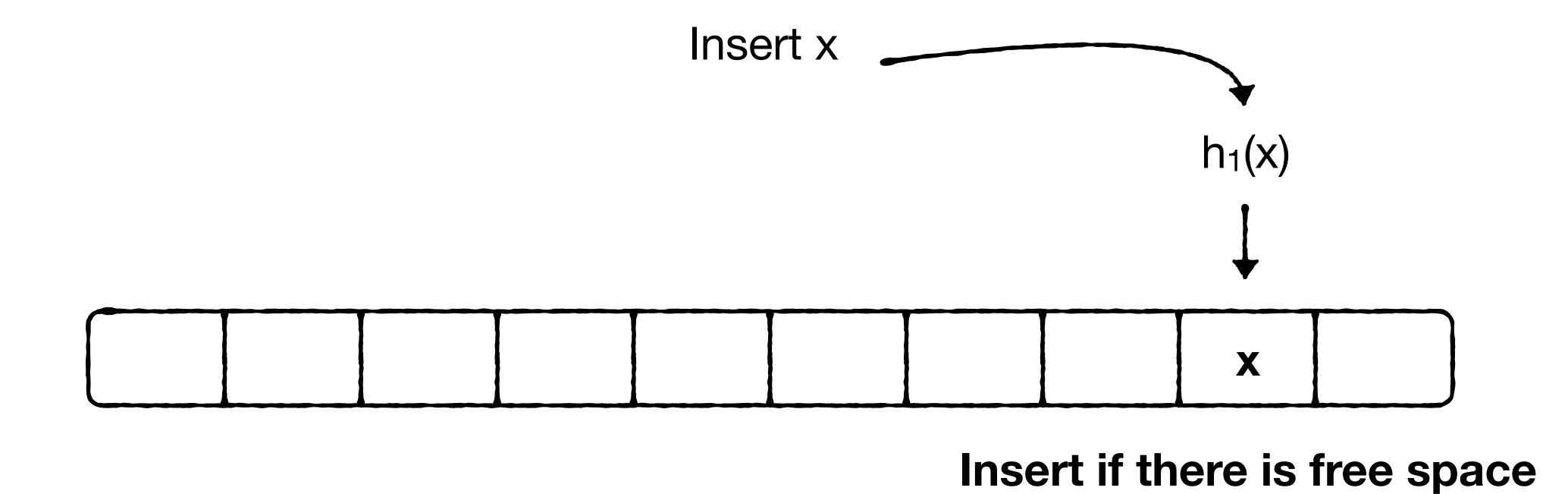
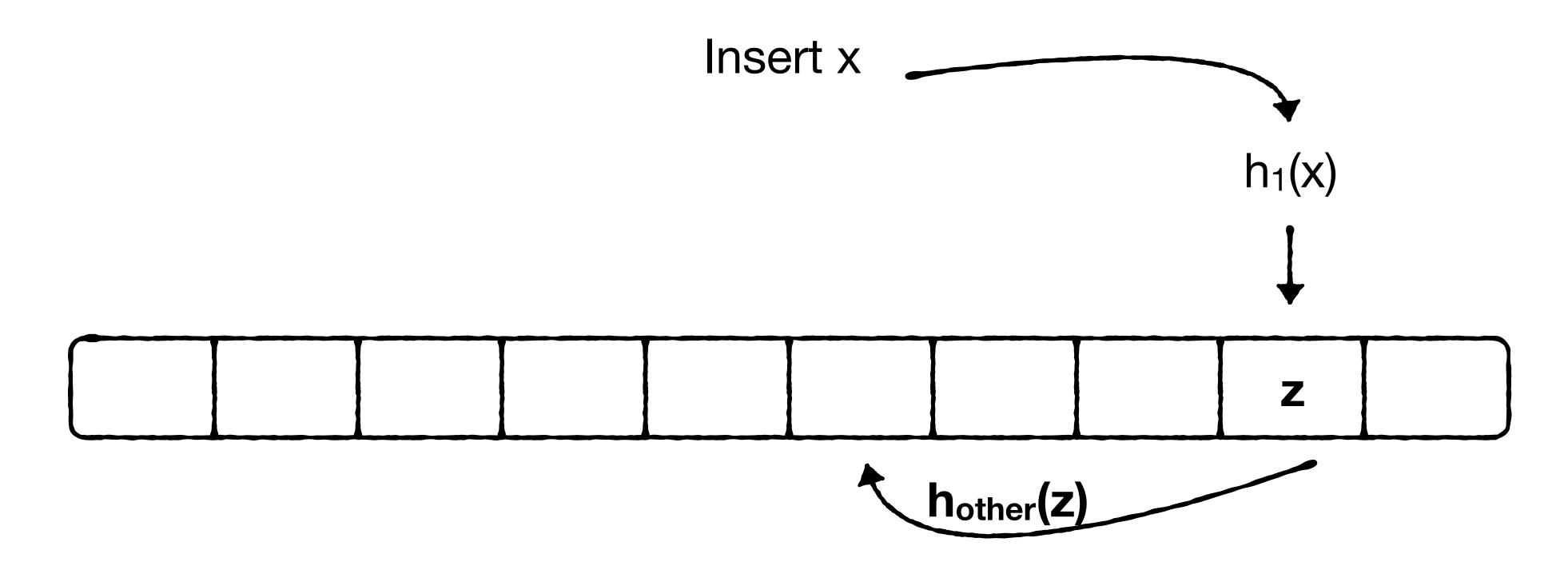
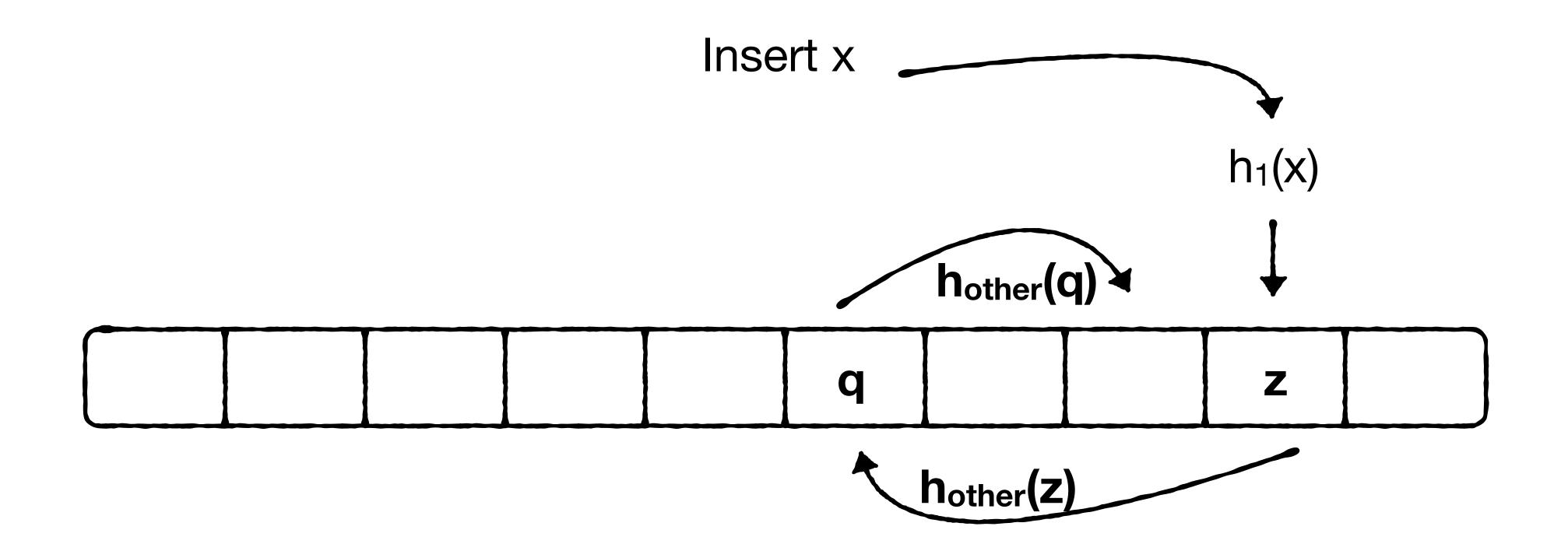


Table of buckets, each containing 1 entry Two hash functions

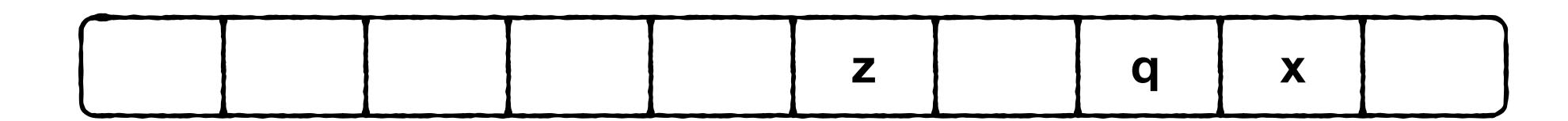


Otherwise, evict existing entry using its alternative hash function to an alternative slot

Table of buckets, each containing 1 entry Two hash functions



Continue doing this recursively until all entries are mapped to an empty bucket

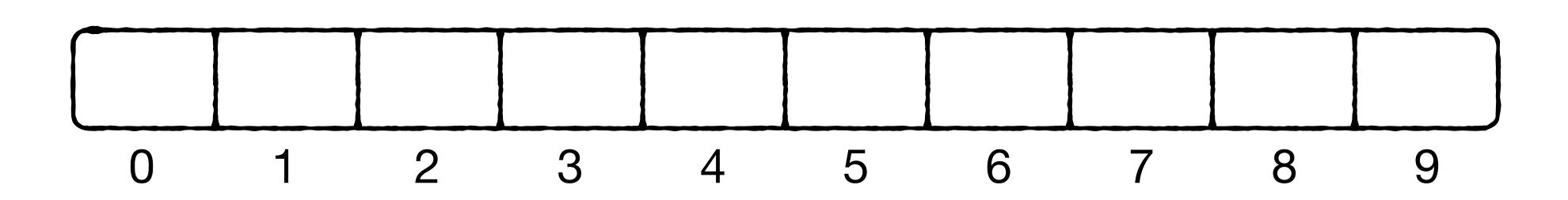


Continue doing this recursively until all entries are mapped to an empty bucket

$$h_1(z) = 7$$
 $h_2(z) = 2$

$$h_1(q) = 7$$
 $h_2(q) = 5$

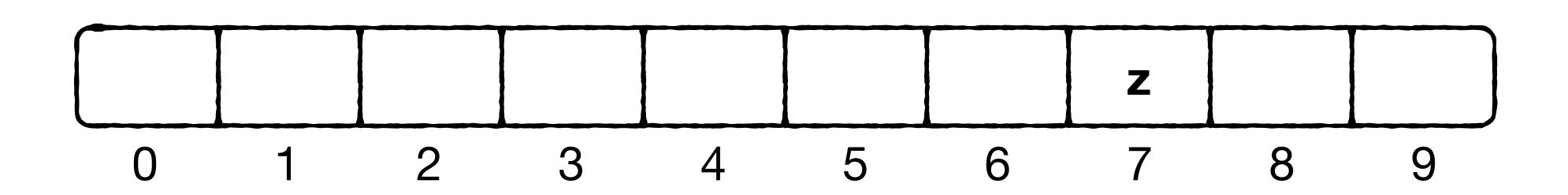
$$h_1(x) = 2$$
 $h_2(x) = 9$



$$h_1(z) = 7$$
 $h_2(z) = 2$

$$h_1(q) = 7$$
 $h_2(q) = 5$

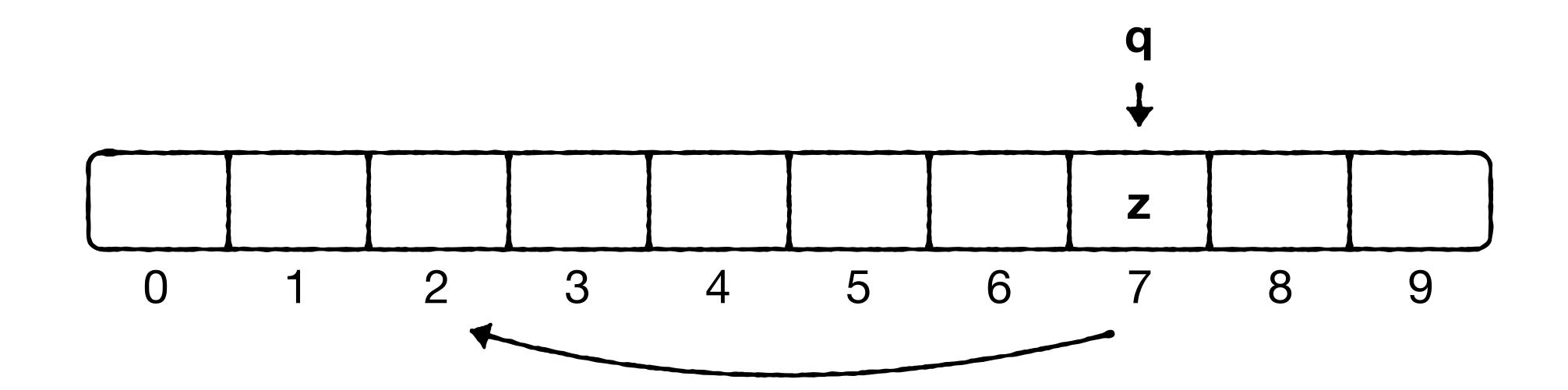
$$h_1(x) = 2$$
 $h_2(x) = 9$



$$h_1(z) = 7$$
 $h_2(z) = 2$

$$h_1(q) = 7$$
 $h_2(q) = 5$

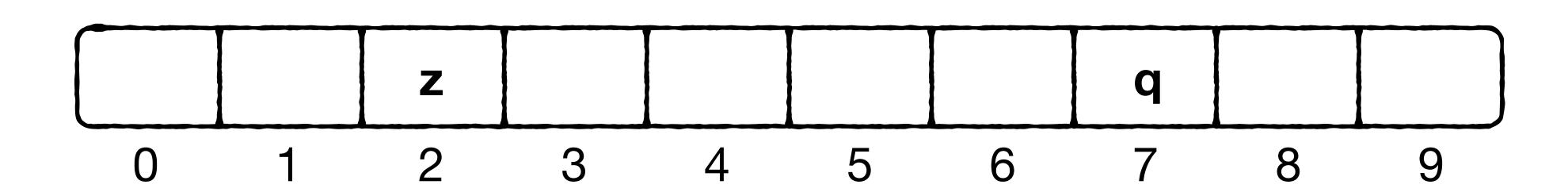
$$h_1(x) = 2$$
 $h_2(x) = 9$



$$h_1(z) = 7$$
 $h_2(z) = 2$

$$h_1(q) = 7$$
 $h_2(q) = 5$

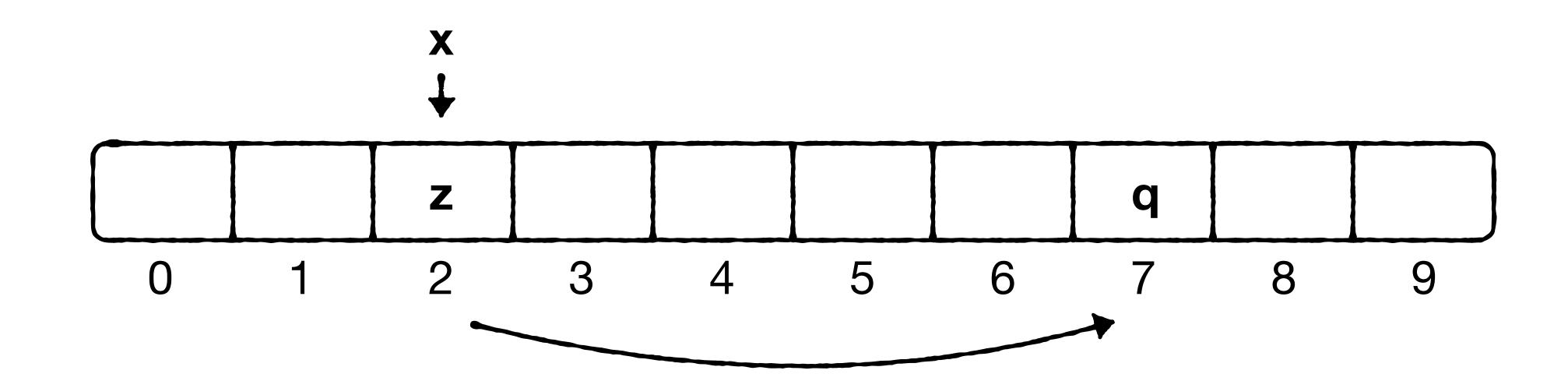
$$h_1(x) = 2$$
 $h_2(x) = 9$



$$h_1(z) = 7$$
 $h_2(z) = 2$

$$h_1(q) = 7$$
 $h_2(q) = 5$

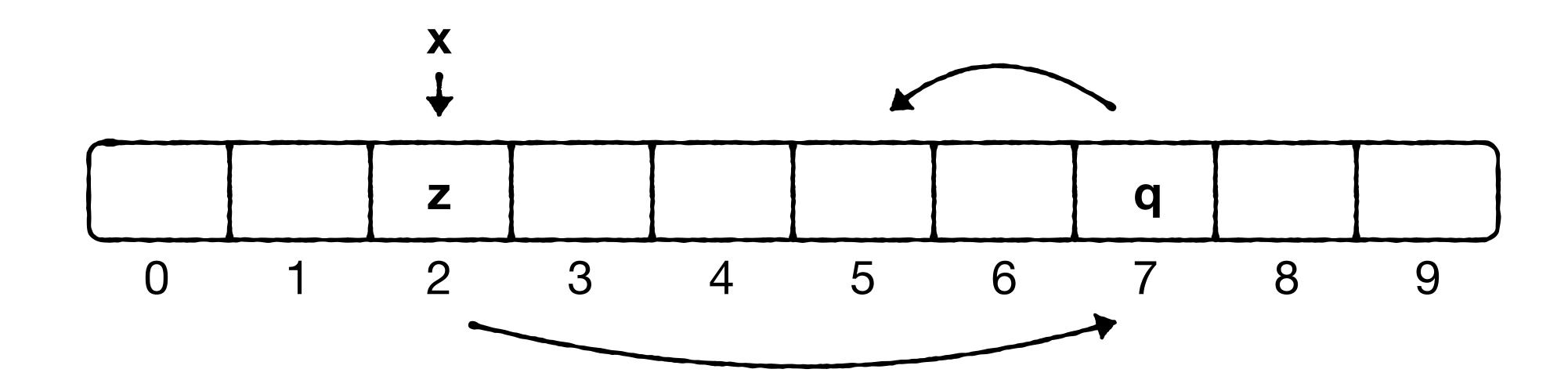
$$h_1(x) = 2$$
 $h_2(x) = 9$



$$h_1(z) = 7$$
 $h_2(z) = 2$

$$h_1(q) = 7$$
 $h_2(q) = 5$

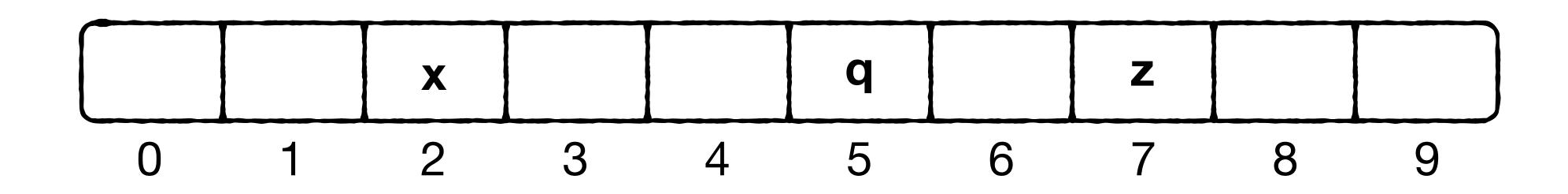
$$h_1(x) = 2$$
 $h_2(x) = 9$



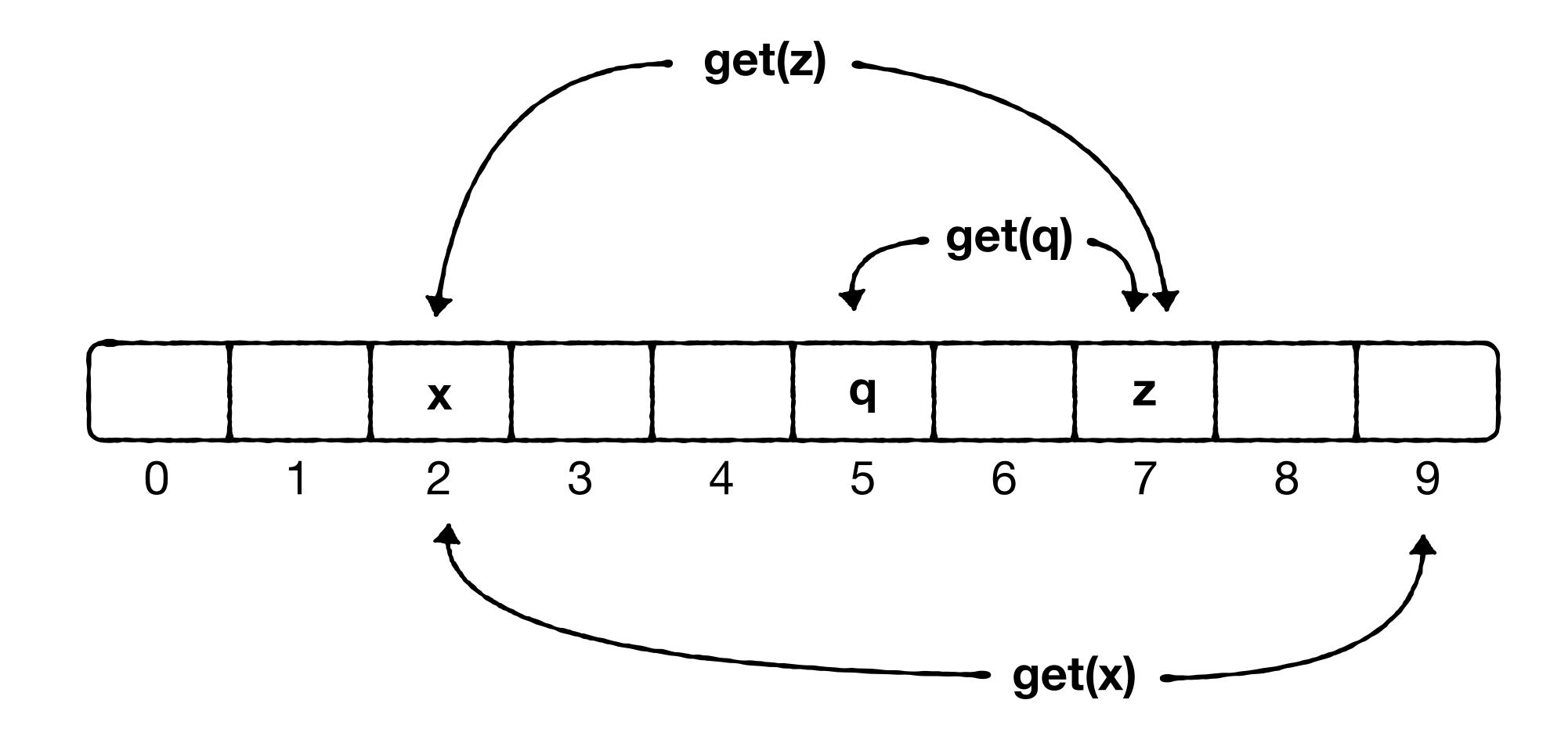
$$h_1(z) = 7$$
 $h_2(z) = 2$

$$h_1(q) = 7$$
 $h_2(q) = 5$

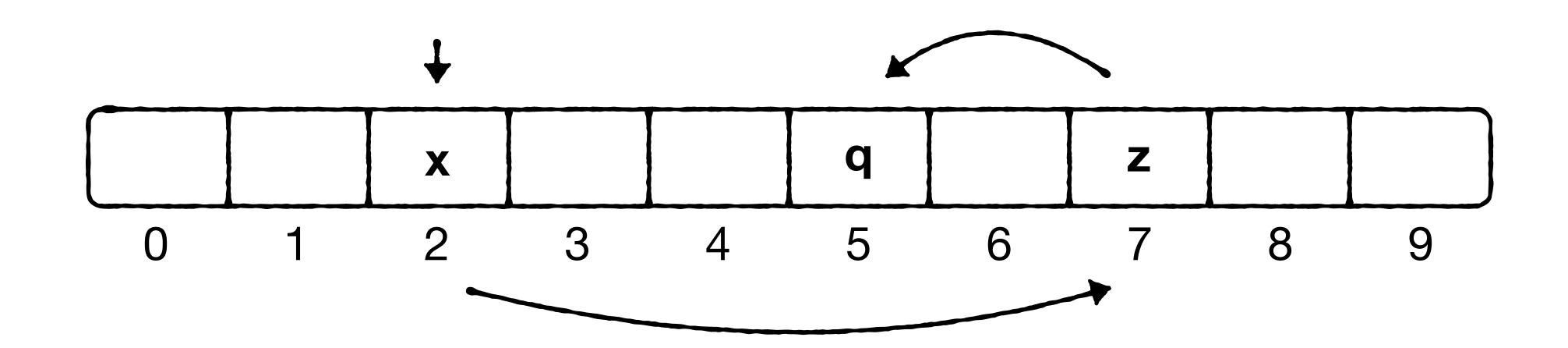
$$h_1(x) = 2$$
 $h_2(x) = 9$



A query has to check at most two locations to find a key



Insertions may endure multiple evictions and swapping of keys across buckets



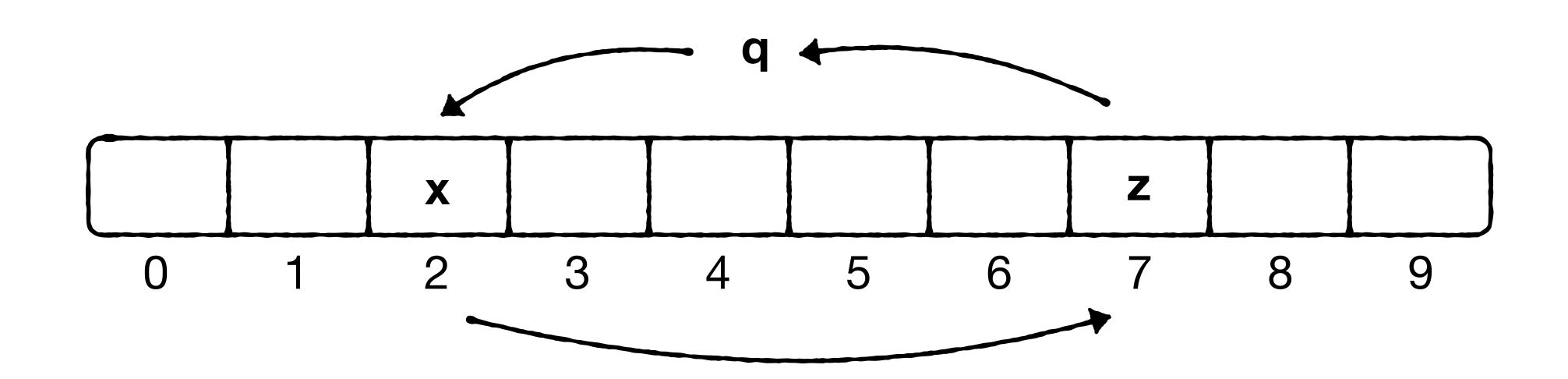
Insertions may endure multiple evictions and swapping of keys across buckets

Worse: an infinite loop is technically possible

$$h_1(z) = 7$$
 $h_2(z) = 2$

$$h_1(q) = 2$$
 $h_2(q) = 7$

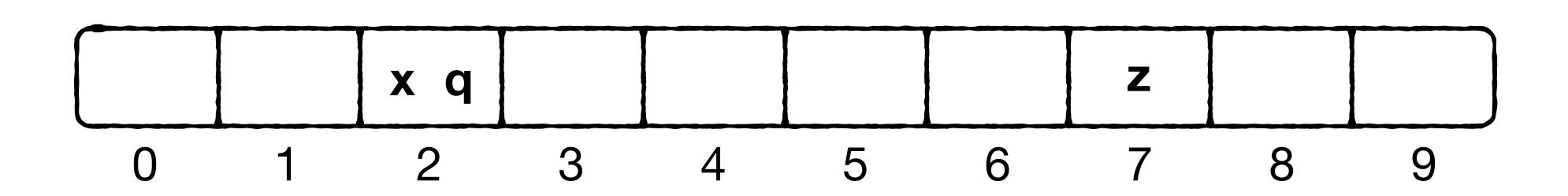
$$h_1(x) = 7$$
 $h_2(x) = 2$



Insertions may endure multiple evictions and swapping of keys across buckets

Worse: an infinite loop is technically possible

alleviate by allowing multiple keys per bucket



Assuming...



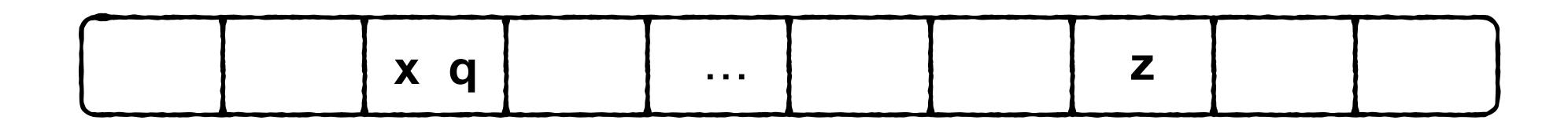
Assuming load factor < 50% for bucket size 1



Assuming load factor < 50% for bucket size 1 load factor < 84% for bucket size 2



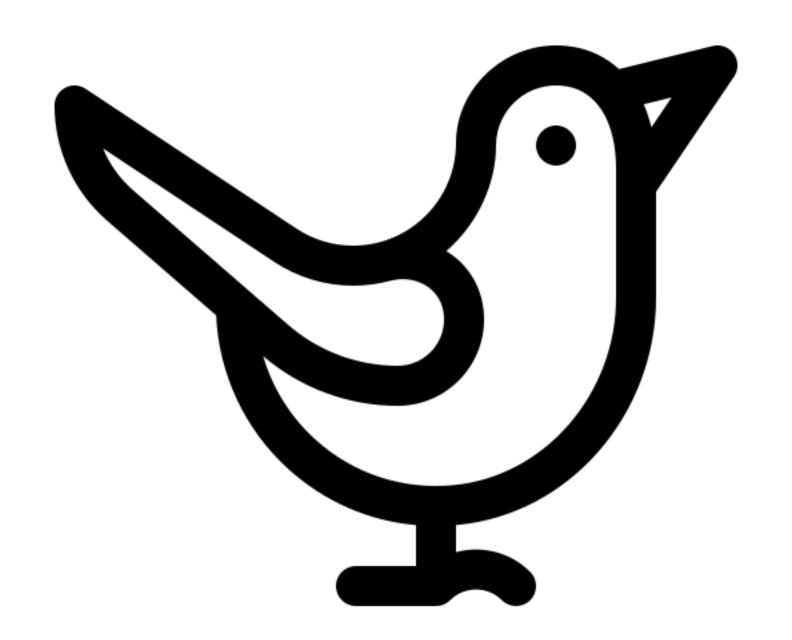
Assuming load factor < 50% for bucket size 1 load factor < 84% for bucket size 2 load factor < 95% for bucket size 4



We'll use this → load factor < 95% for bucket size 4



Why is it called Cuckoo hashing?



Lay eggs in other birds' nests, and the hatchlings "evict" the other birds' eggs.

```
Index size = N * (P + K) / \alpha
                       N = data size
                       P = pointer size
                       K = key size
                       \alpha = collision resolution overheads
                                          Live bytes
                                                        Buffer
                        index
SSD
```

Index size = $N * (P + K) / \alpha$

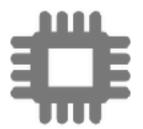
N = data size

P = pointer size

K = key size

Addressed with cuckoo hashing

 \rightarrow α = collision resolution overheads



Cuckoo index

Live bytes

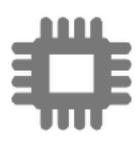
Buffer





Index size = N * (P + K) / α N = data size
P = pointer size
K = key size

a = collision resolution overheads $\approx 0.8 \rightarrow \approx 0.95$



Cuckoo index

Live bytes

Buffer





Index size =
$$N * (P + K) / \alpha$$

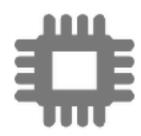
N = data size

P = pointer size

Now lets attack this With a Cuckoo Filter

 \rightarrow K = key size = $\Omega(\log_2 N)$

 $\alpha = \text{collision resolution overheads } \approx 0.8 \longrightarrow \approx 0.95$



Cuckoo index

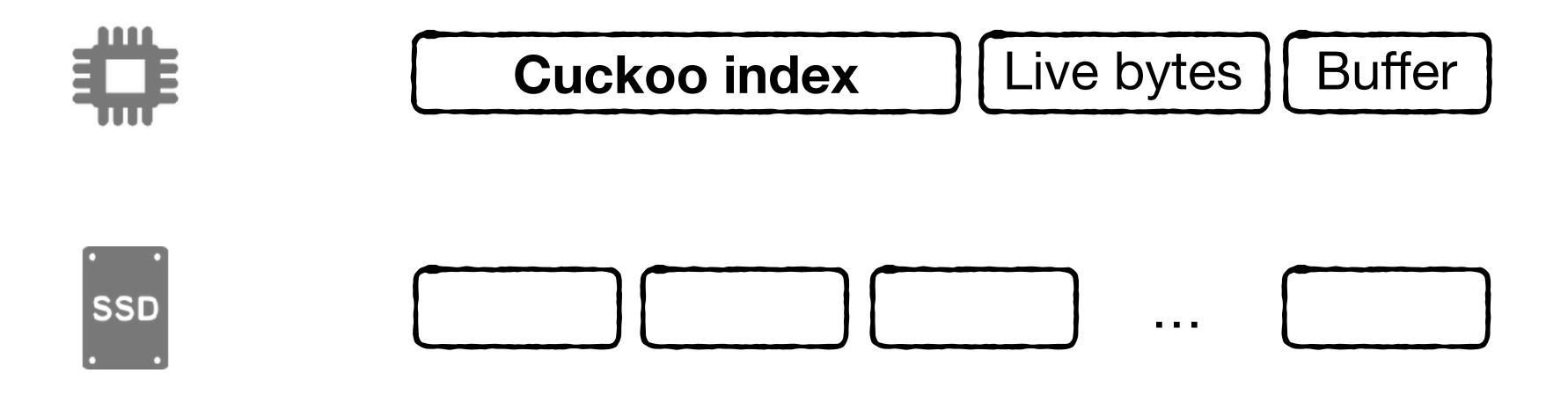
Live bytes

Buffer



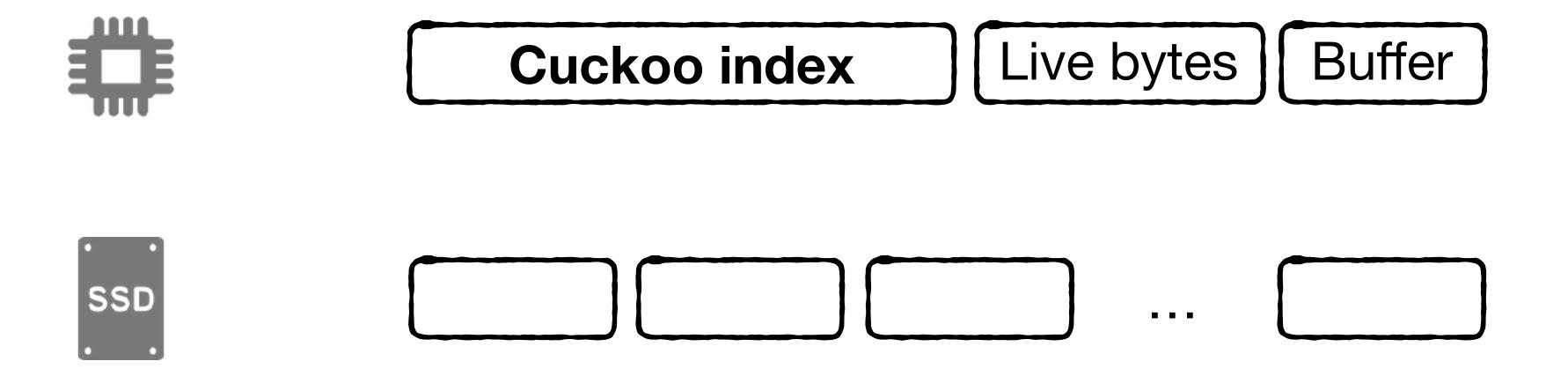
Problems with storing full keys in a hash table

(1) Keys may be arbitrarily large

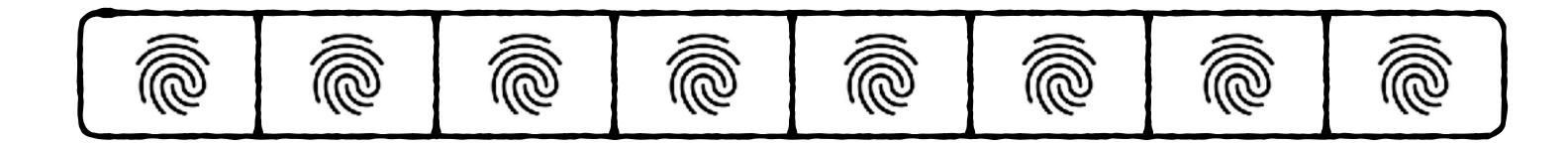


Problems with storing full keys in a hash table

- (1) Keys may be arbitrarily large
- (2) Keys can be variable-length (requires additional metadata and CPU cycles to encode)

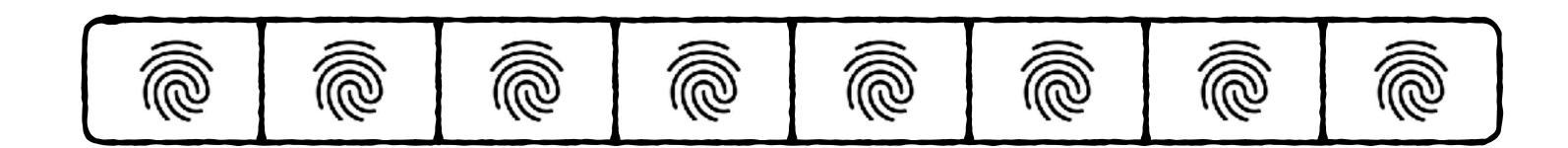


Same as Cuckoo hash tables, but store fingerprints instead of keys



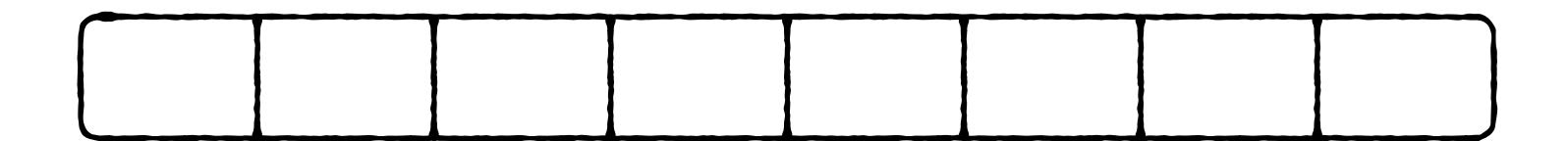
Same as Cuckoo hash tables, but store fingerprints instead of keys

A fingerprint is a hash digest derived by hashing a key



Same as Cuckoo hash tables, but store fingerprints instead of keys A fingerprint is a hash digest derived by hashing a key

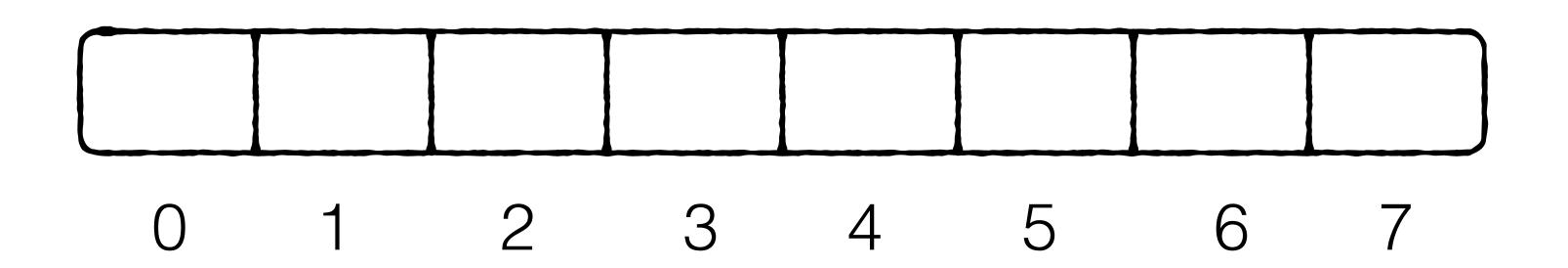
Example:
$$FP(X) = 0100$$



Same as Cuckoo hash tables, but store fingerprints instead of keys A fingerprint is a hash digest derived by hashing a key

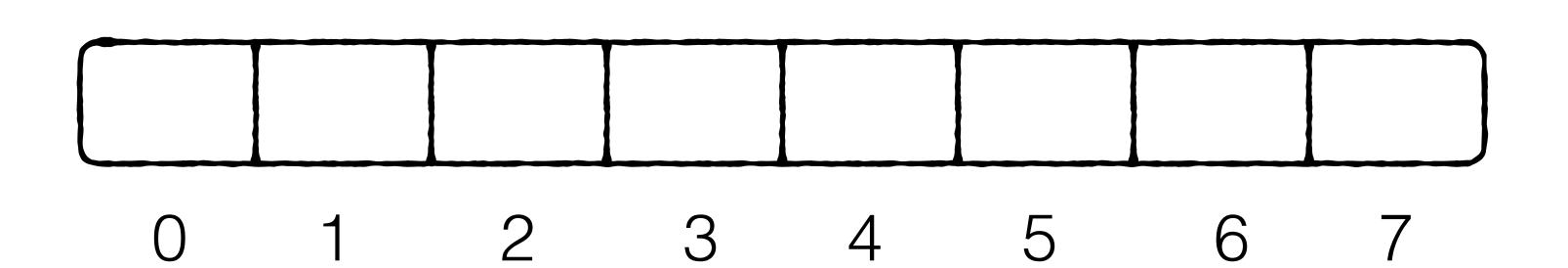
Example:
$$\mathbf{FP}(\mathbf{X}) = 0100$$

 $h_1(X)$ = Bucket address



Same as Cuckoo hash tables, but store fingerprints instead of keys A fingerprint is a hash digest derived by hashing a key

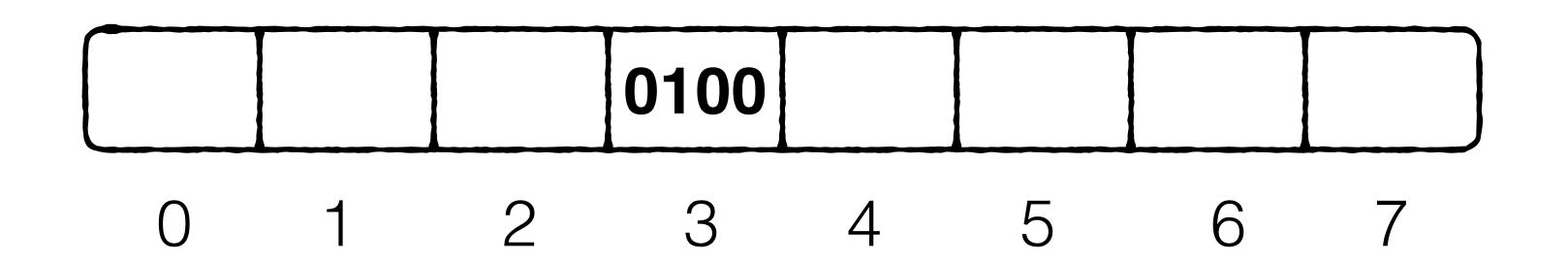
Example:
$$\mathbf{FP}(\mathbf{X}) = 0100$$



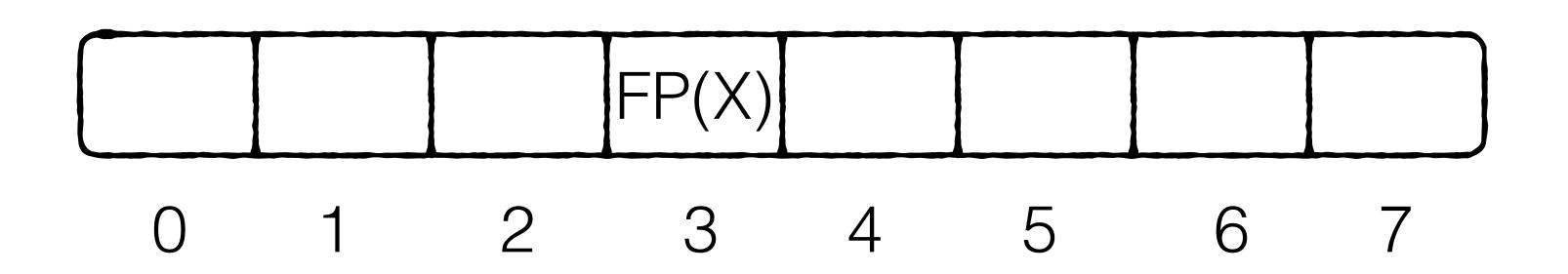
 $h_1(X) = 3$

Same as Cuckoo hash tables, but store fingerprints instead of keys A fingerprint is a hash digest derived by hashing a key

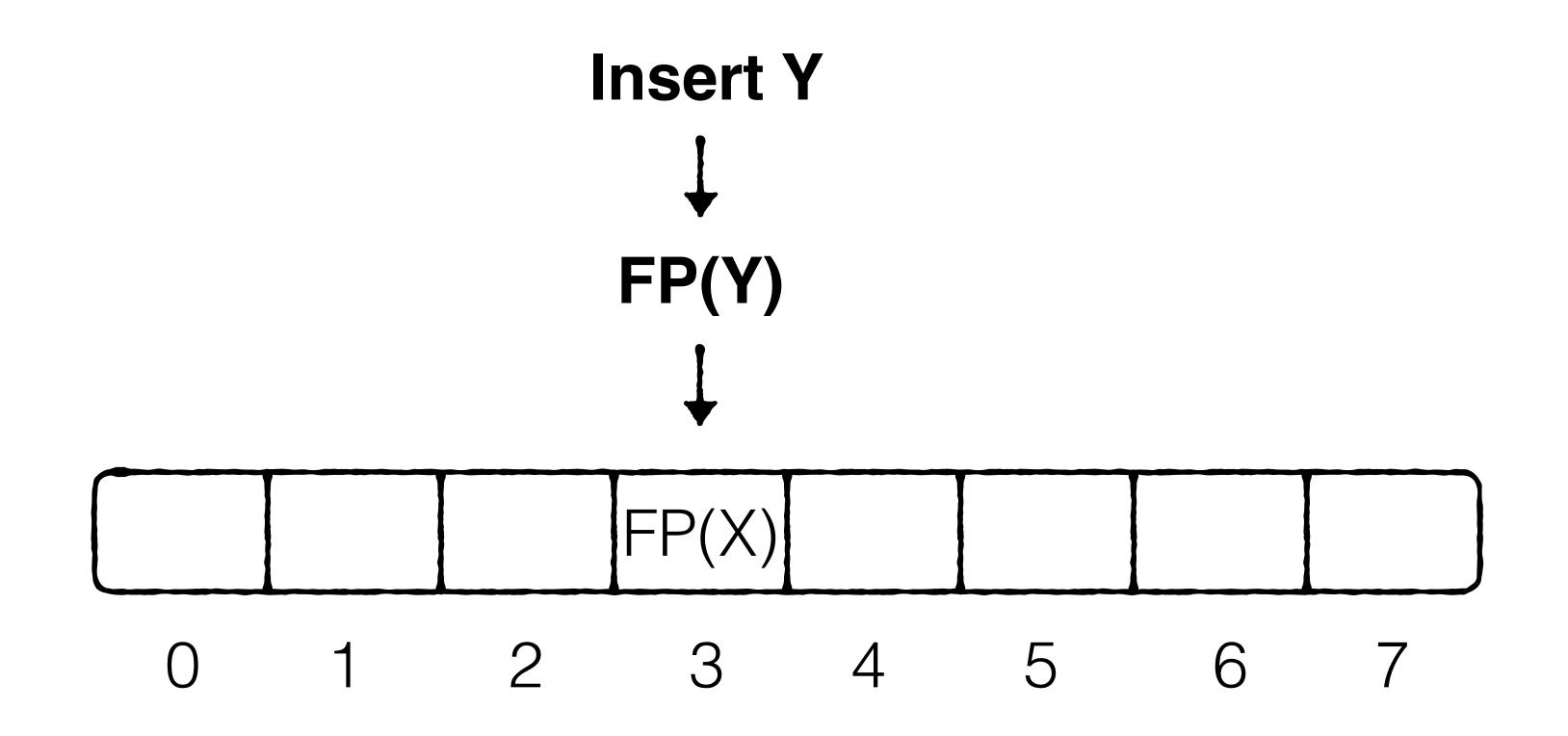
Example:
$$FP(X) = Mbits$$
 $h_1(X) = 3$



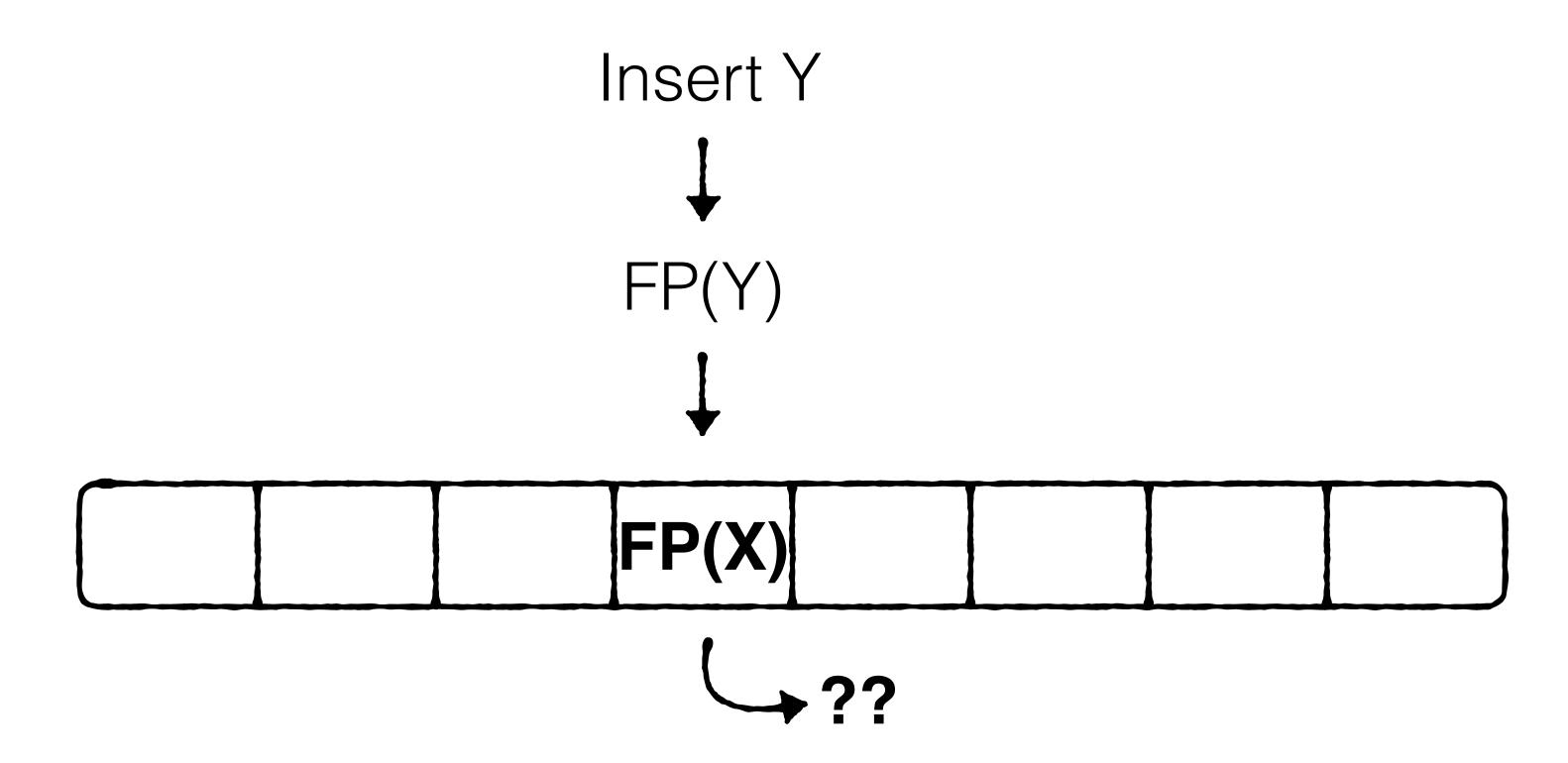
Same as Cuckoo hash tables, but store fingerprints instead of keys A fingerprint is a hash digest derived by hashing a key



Suppose we then insert another fingerprint to the same bucket



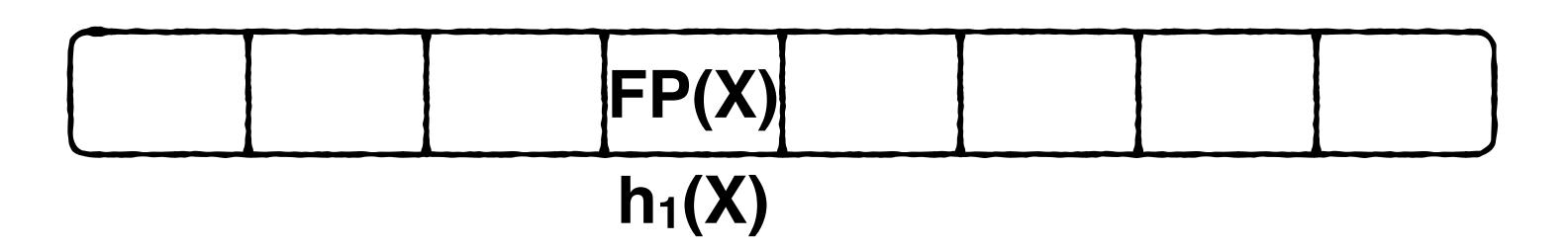
Suppose we then insert another fingerprint to the same bucket



Where to evict X? We no longer have its key! How to derive an alternative bucket?

We have two pieces of information about X (1) Bucket address: h₁(X)

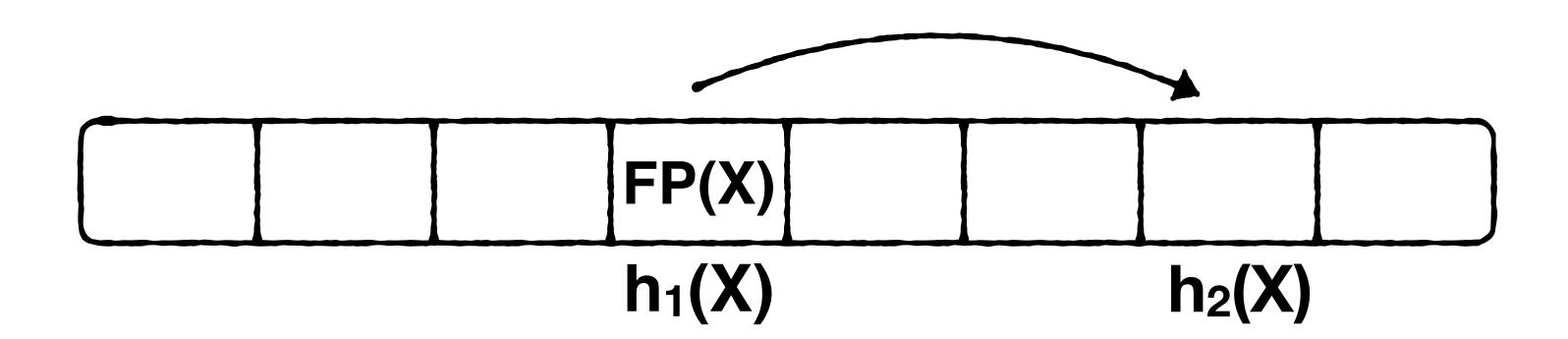
(2) fingerprint FP(X)



We have two pieces of information about X

- (1) Bucket address: h₁(X)
- (2) fingerprint FP(X)

We want to combine them to give alternative bucket address h₂(X)

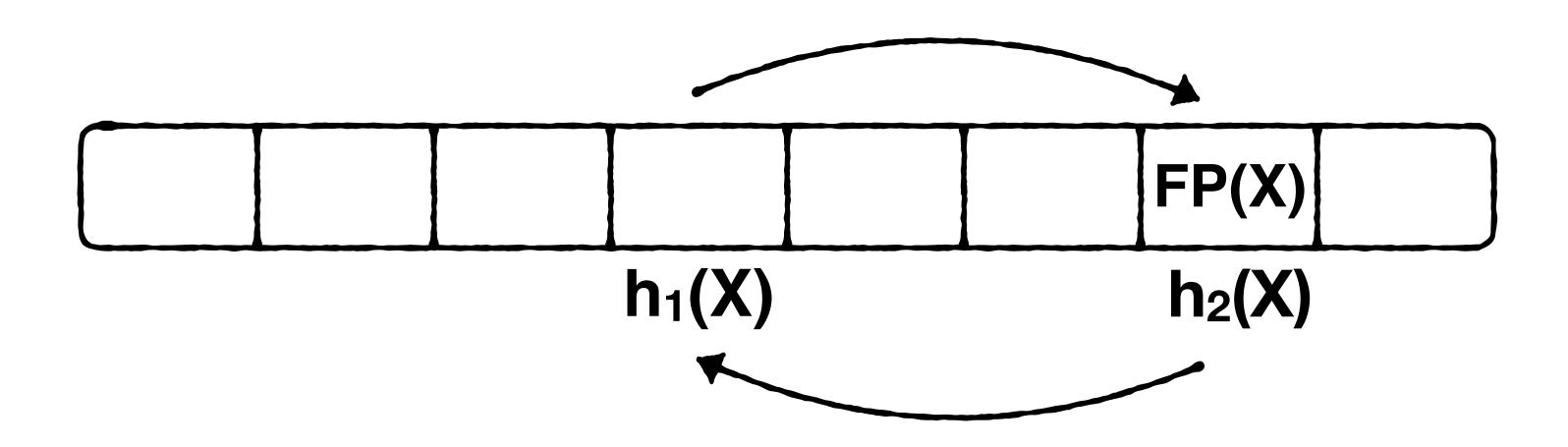


We have two pieces of information about X

- (1) Bucket address: h₁(X)
- (2) fingerprint FP(X)

We want to combine them to give alternative bucket address h₂(X)

This mapping must be reversible, so we can derive h₁(X) from h₂(X) and FP(X)

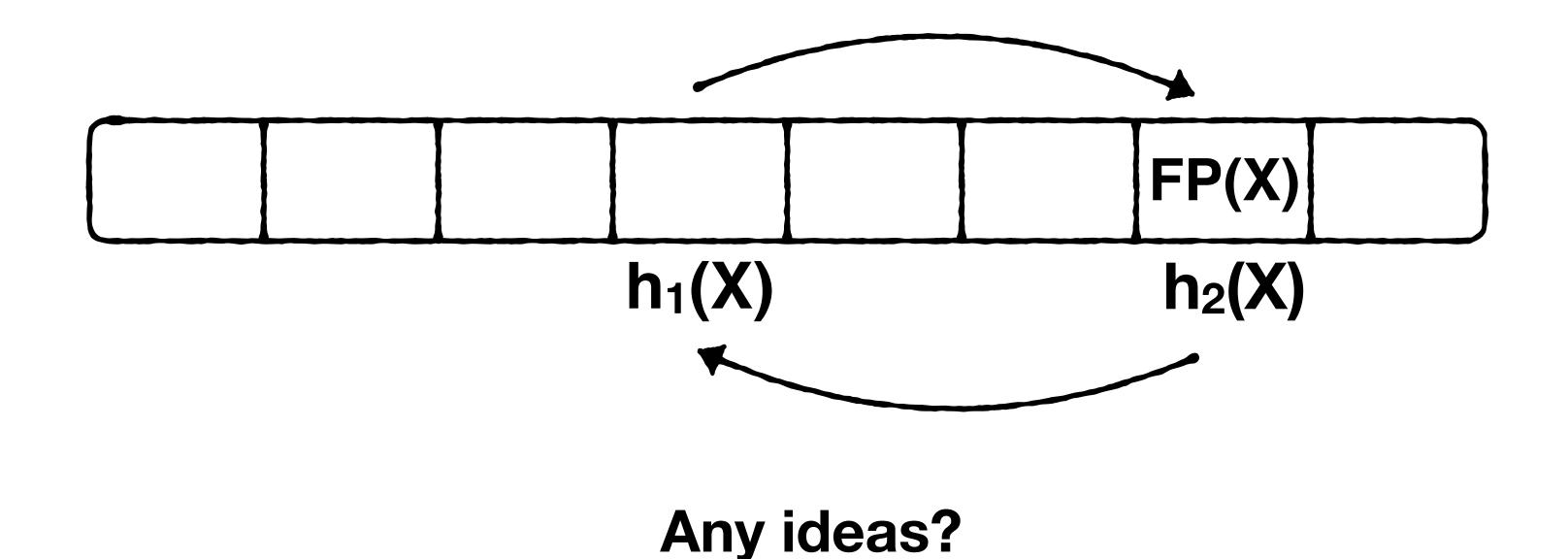


We have two pieces of information about X

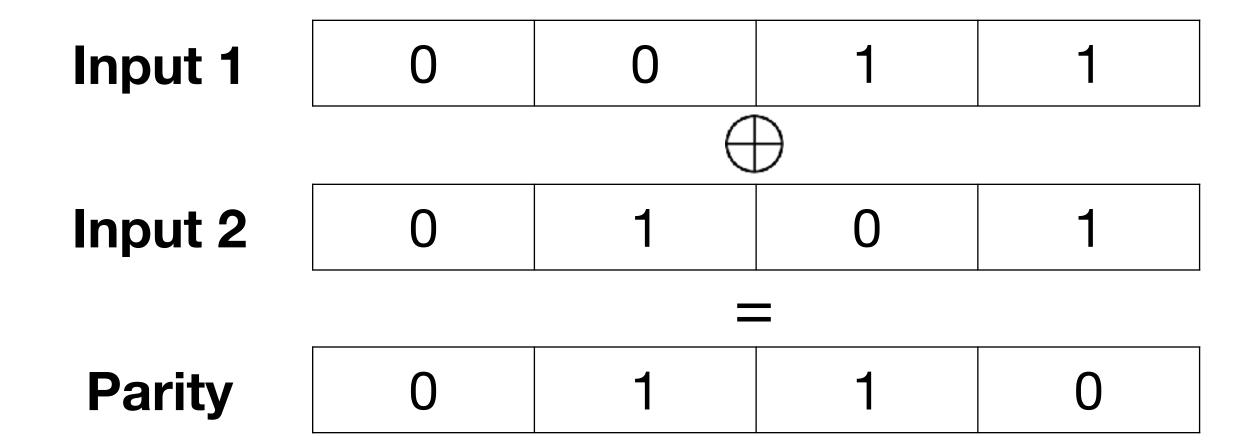
- (1) Bucket address: h₁(X)
- (2) fingerprint FP(X)

We want to combine them to give alternative bucket address h₂(X)

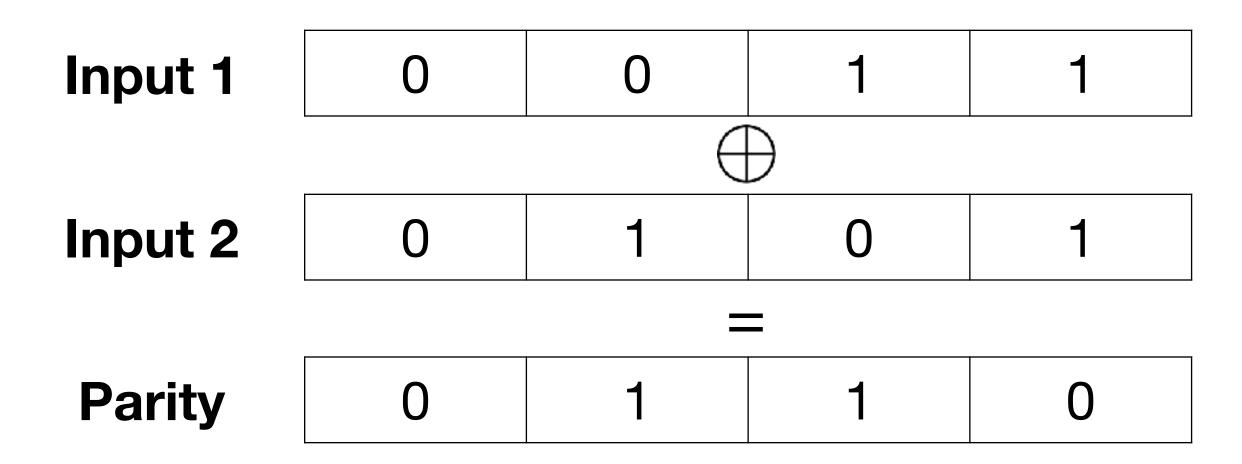
This mapping must be reversible, so we can derive $h_1(X)$ from $h_2(X)$ and FP(X)



XOR Operator



XOR Operator



If number of 1s in an input column is even, the result is 0. If odd, it is 1.

Suppose we lost input 2

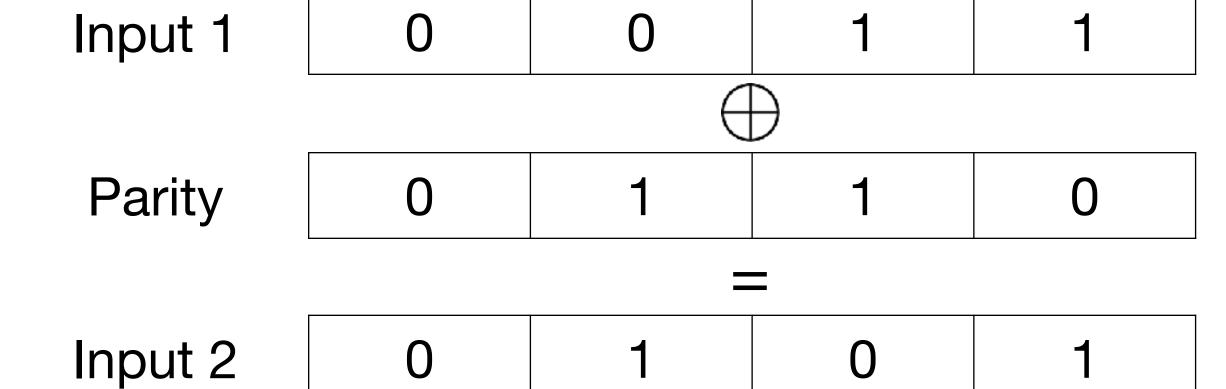
Input 1

0 0 1 1

Parity

0 1 1 0





Recovered

Or suppose we lost input 1

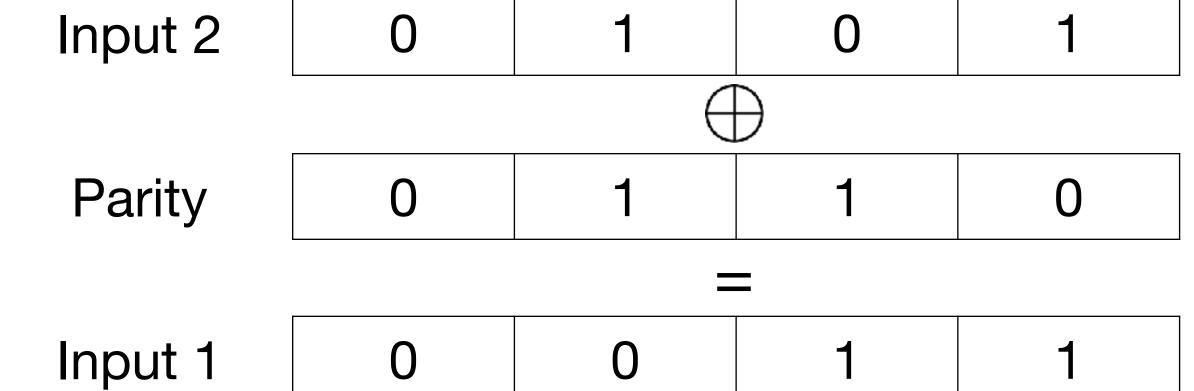
Input 2

0 1 0 1

Parity

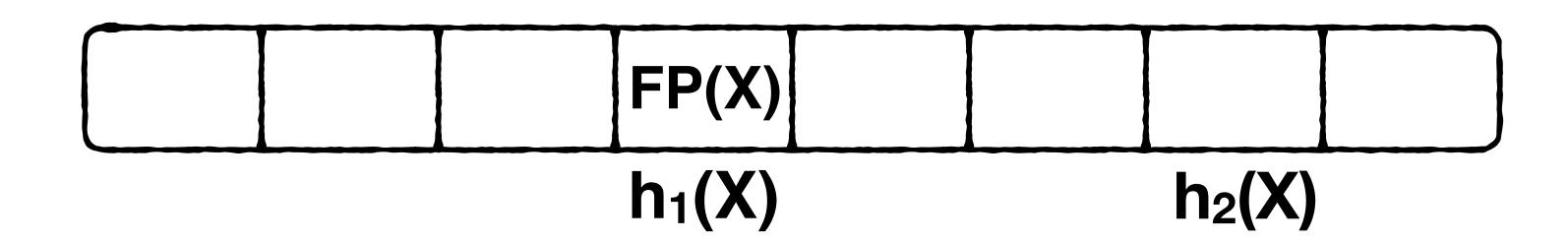
0 1 1 0



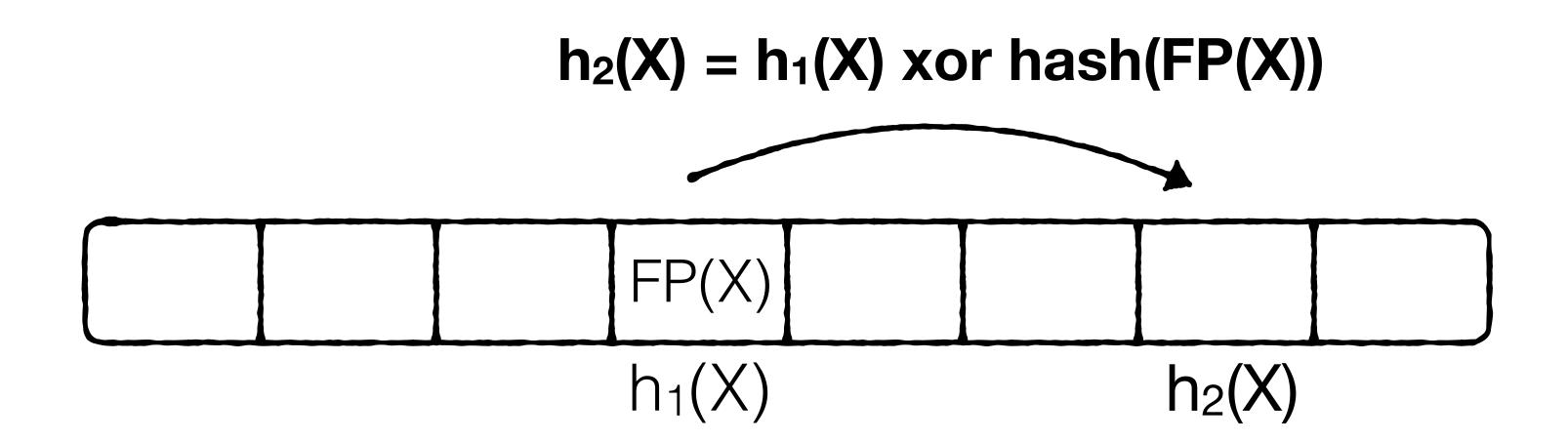


Recovered

Let's XOR the bucket address and a hash of the fingerprint:

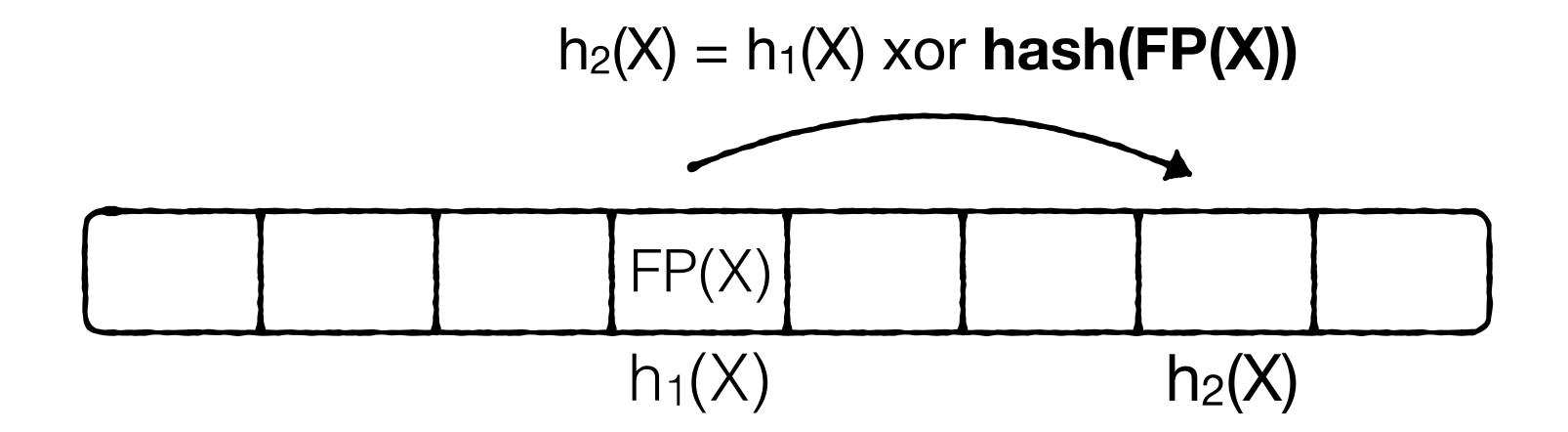


Let's XOR the bucket address and a hash of the fingerprint:



Let's XOR the bucket address and a hash of the fingerprint:

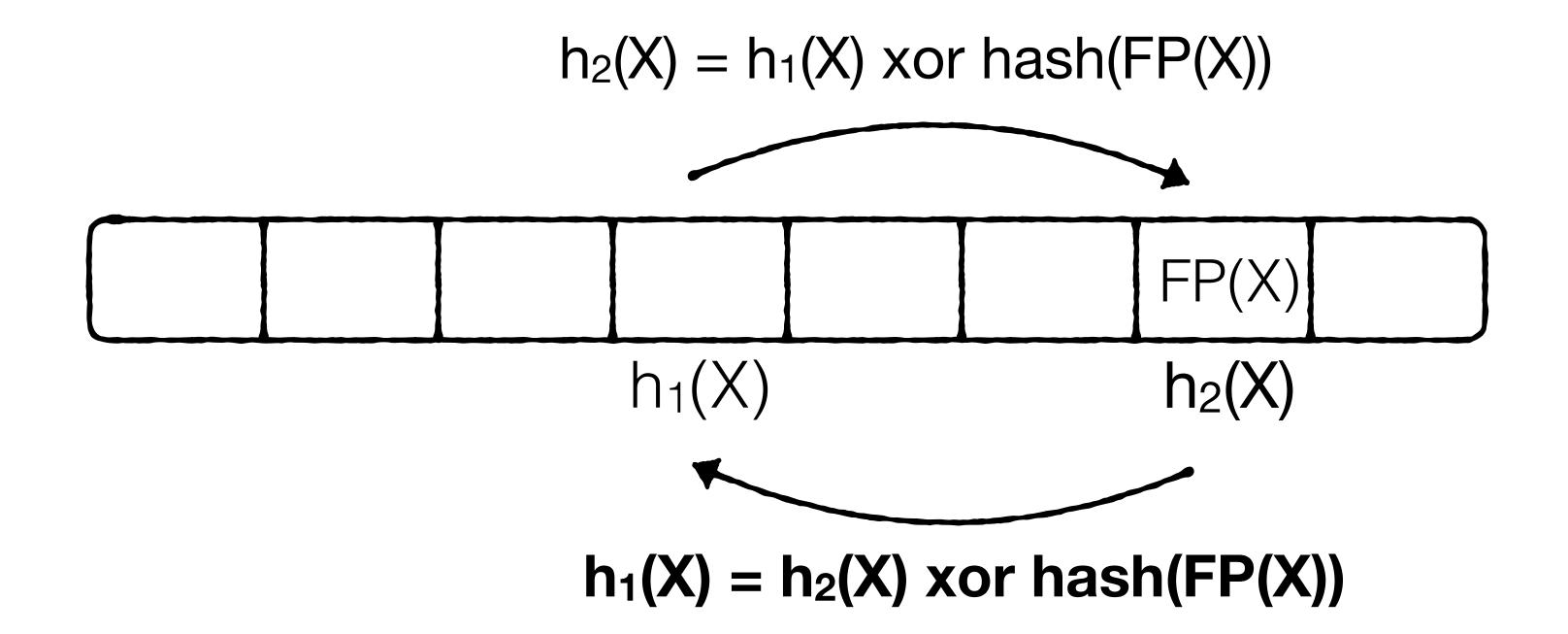
(We hash the fingerprint to map it to the same address space size as the buckets)



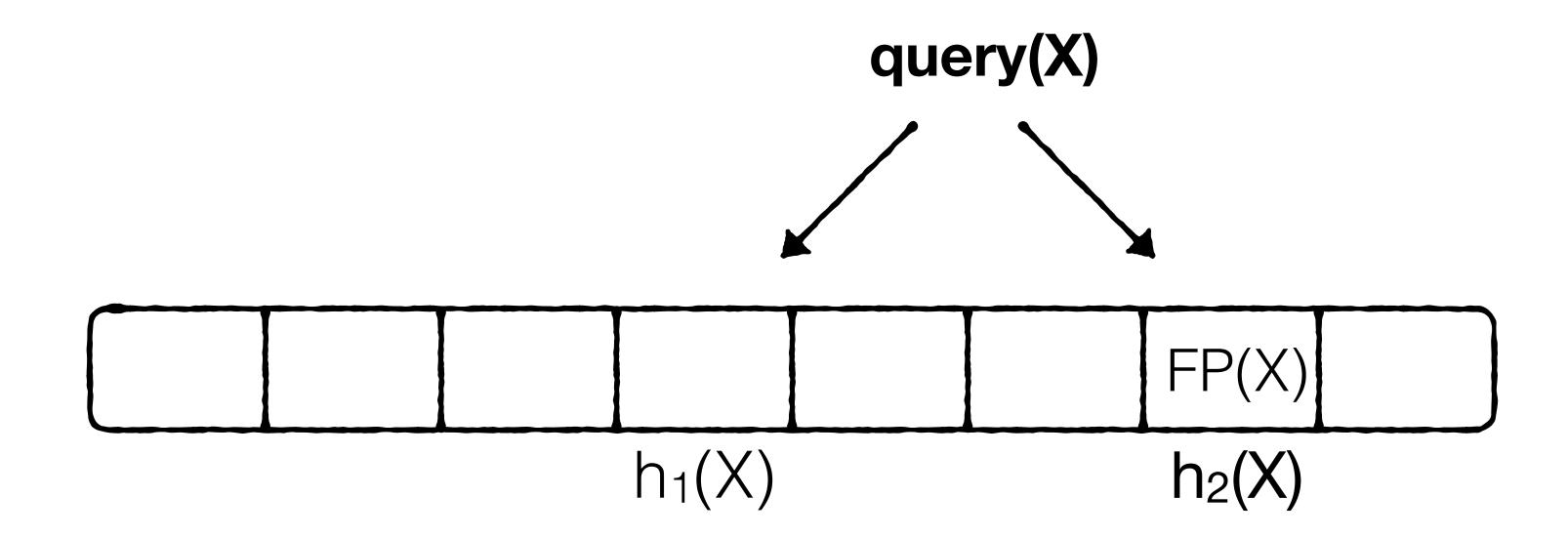
Let's XOR the bucket address and a hash of the fingerprint:

(We hash the fingerprint to map it to the same address space size as the buckets)

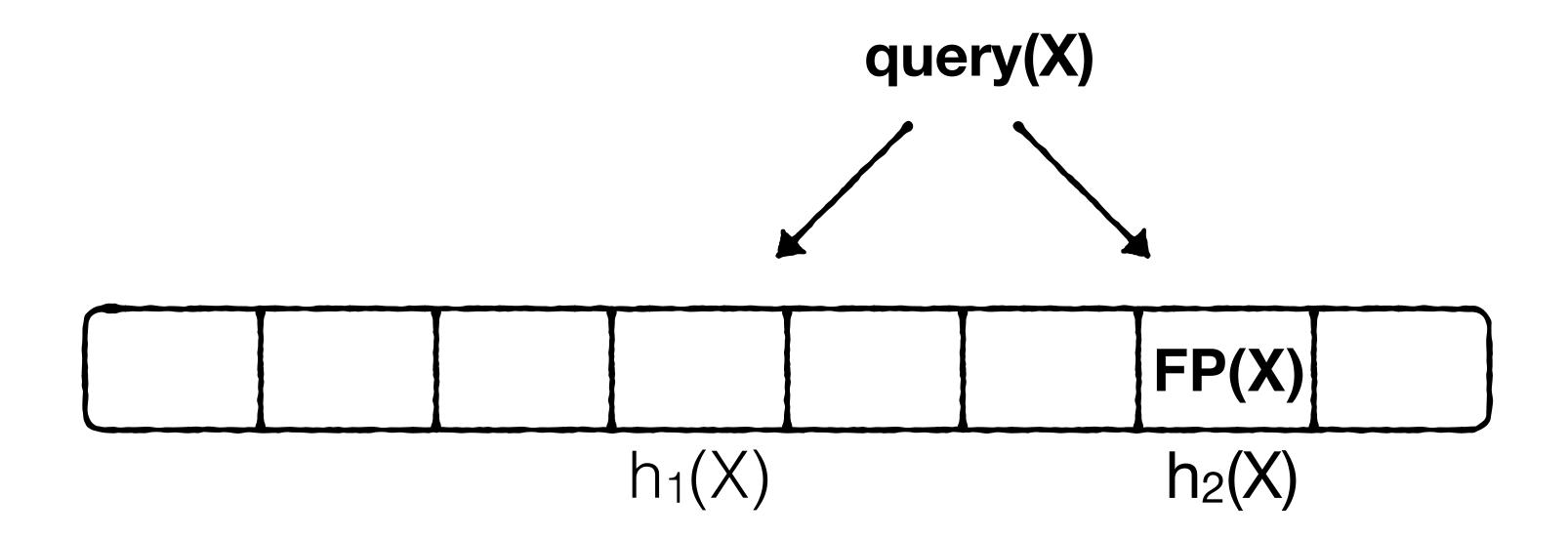
The resulting mapping is reversible



Thus, we must search only two buckets to find an entry's fingerprint

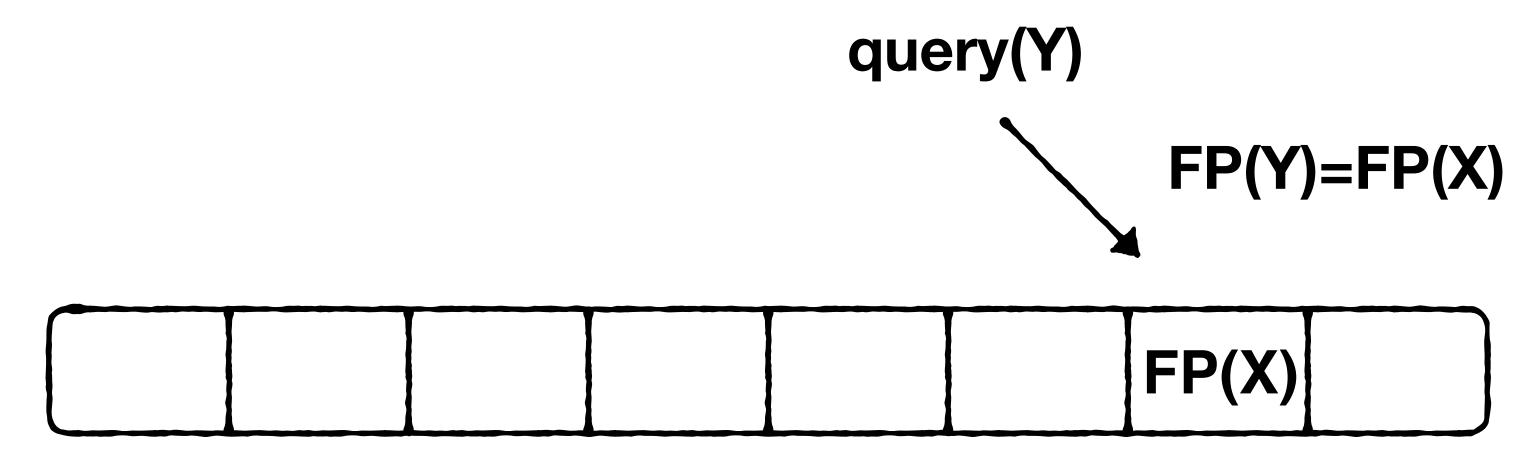


Thus, we must search only two buckets to find an entry's fingerprint If we find a matching fingerprint, we report a positive.



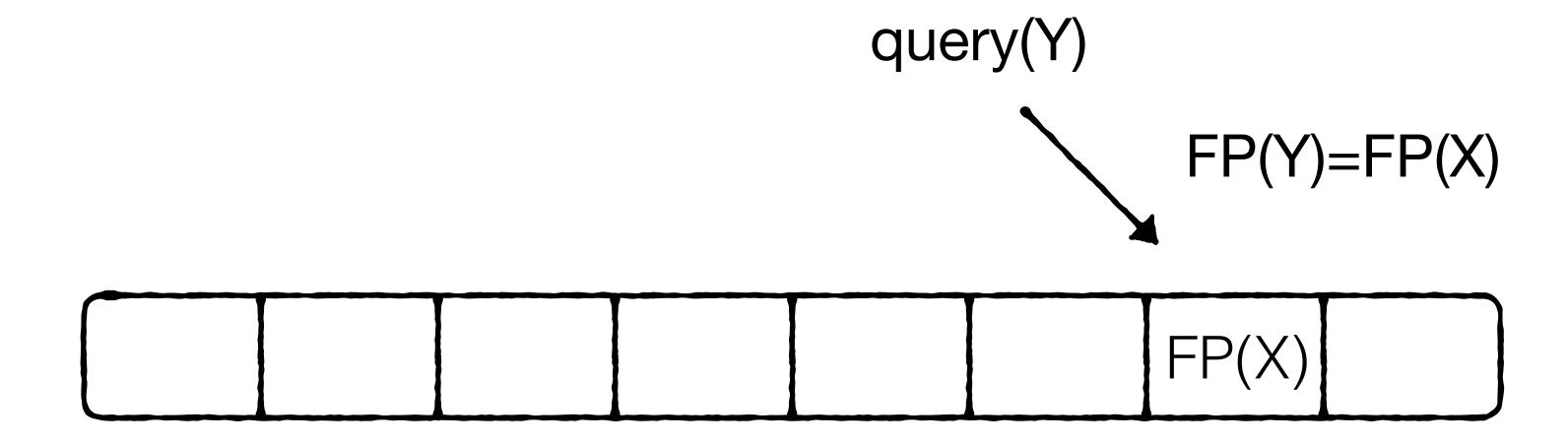
Thus, we must search only two buckets to find an entry's fingerprint If we find a matching fingerprint, we report a positive.

However, we can have false positives.



Thus, we must search only two buckets to find an entry's fingerprint If we find a matching fingerprint, we report a positive.

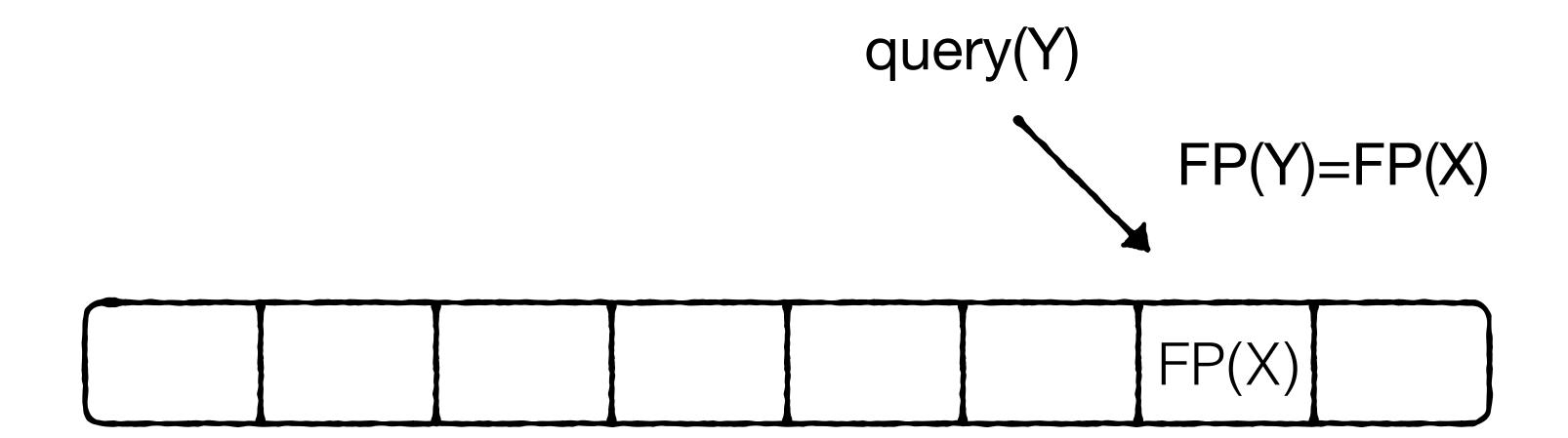
However, we can have false positives.



Can we have false negatives?

Thus, we must search only two buckets to find an entry's fingerprint If we find a matching fingerprint, we report a positive.

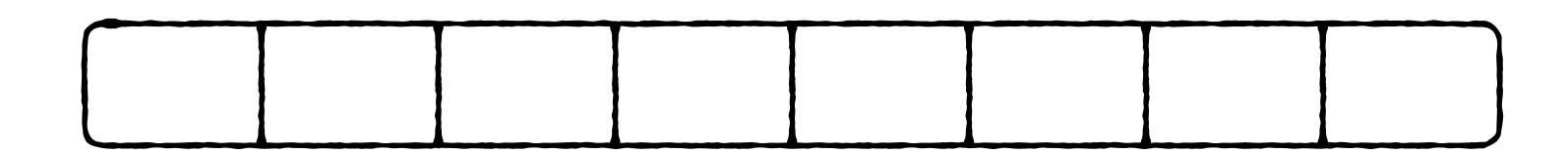
However, we can have false positives.



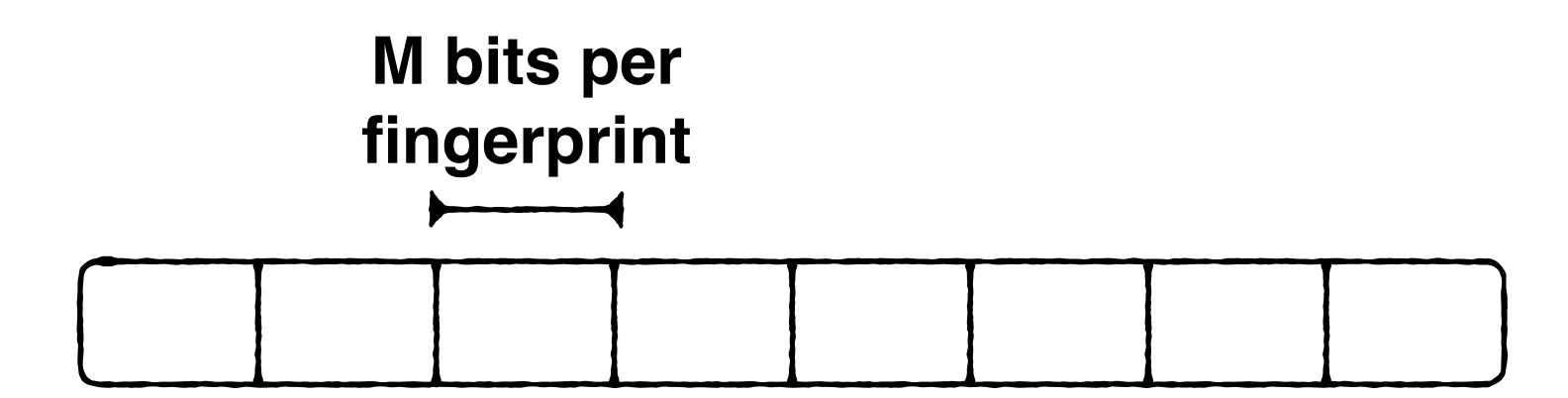
Can we have false negatives?

No, because a fingerprint for a key that had been inserted is in one of only two possible buckets, both of which we search.

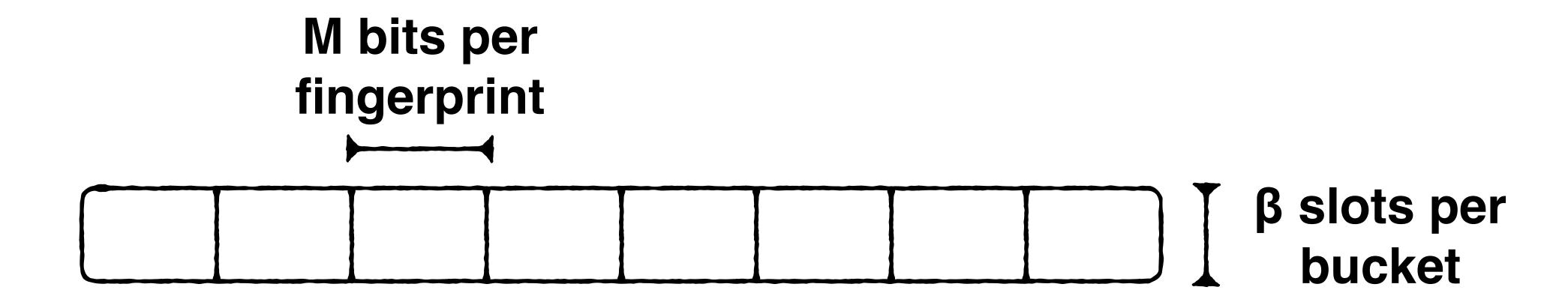
What's the false positive rate?



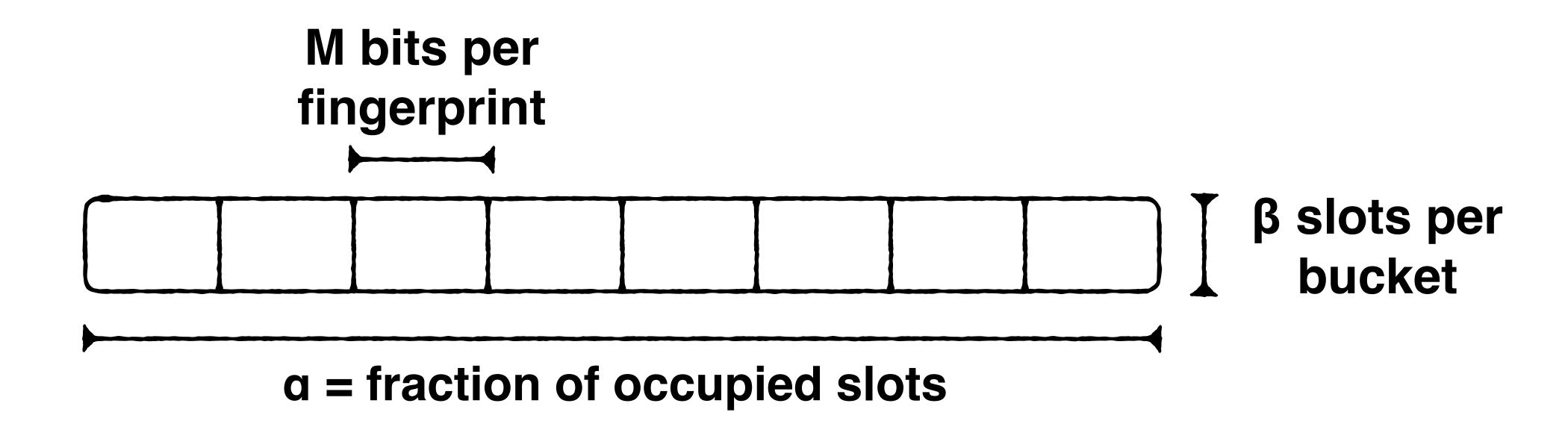
What's the false positive rate?



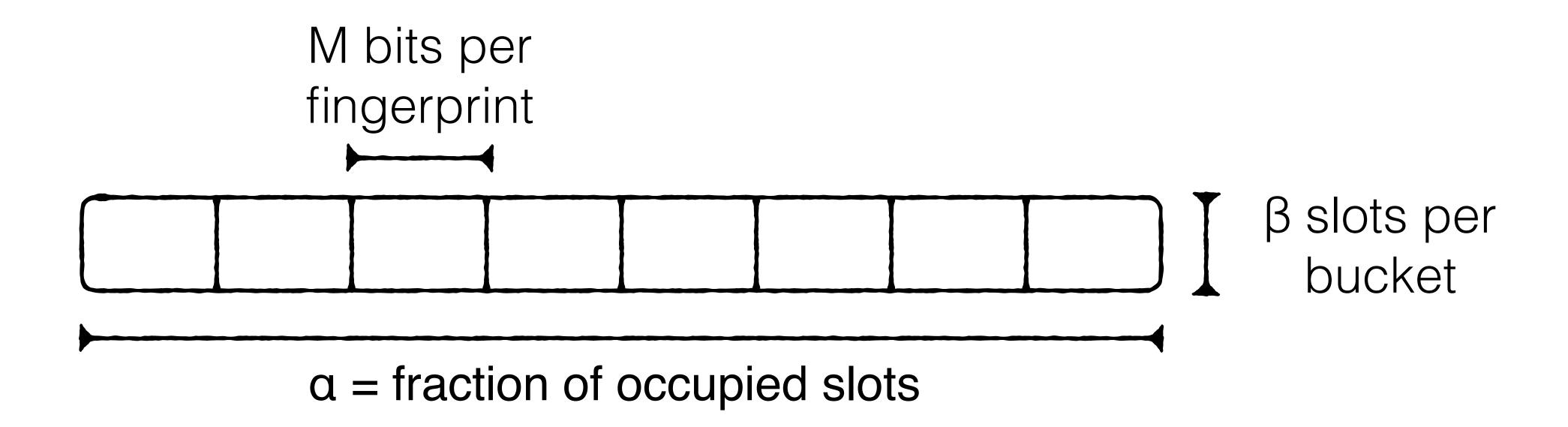
What's the false positive rate?



What's the false positive rate?

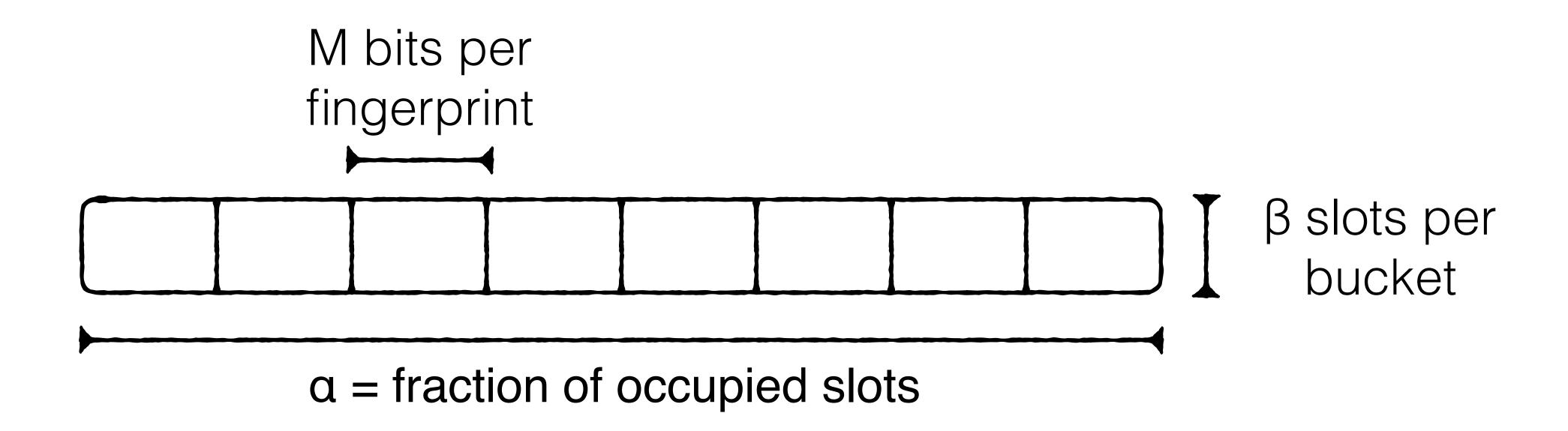


false positive rate $\approx 2 \cdot 2^{-M} \cdot \beta \cdot \alpha$



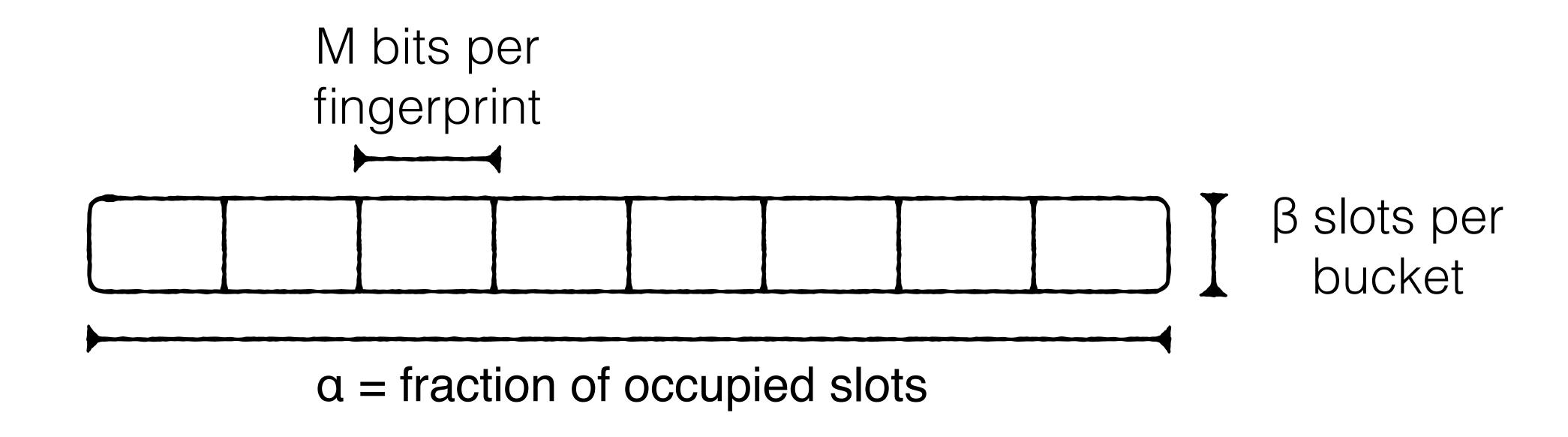
Since a query searches two buckets

false positive rate $\approx 2 \cdot 2^{-M} \cdot \beta \cdot \alpha$

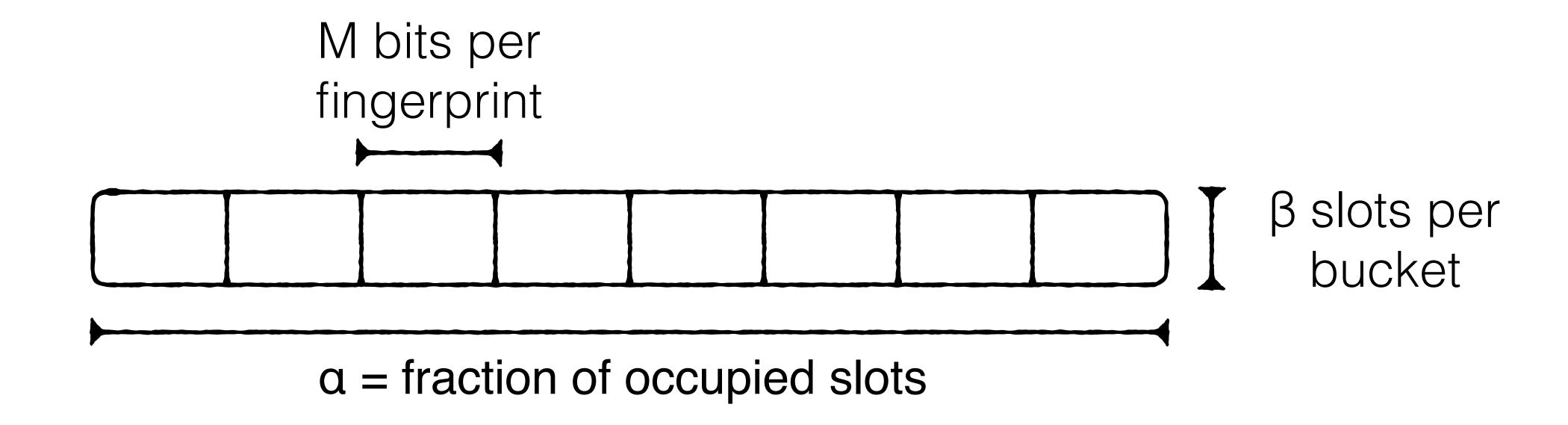


When at capacity

false positive rate $\approx 2 \cdot 2^{-M} \cdot 4 \cdot \textbf{0.95}$

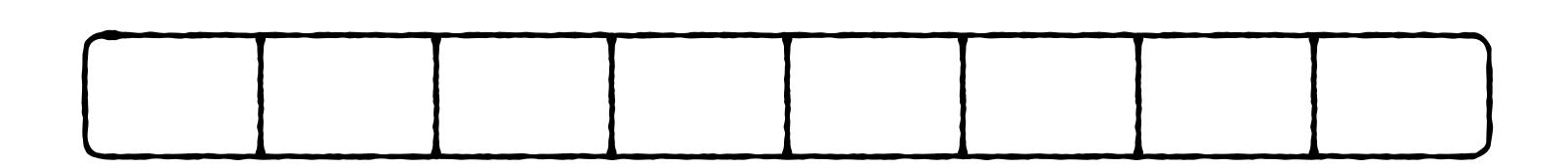


false positive rate $\approx 2^{-M+3}$



false positive rate $\approx 2^{-M+3}$ < 2-M·In(2)

For Bloom filter when M ≥ 10

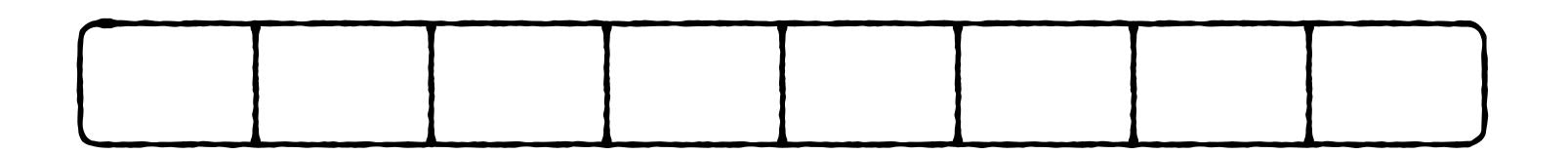


false positive rate:

 $\approx 2^{-M+3}$

2-M · In(2)

Memory accesses for positive query?



false positive rate:

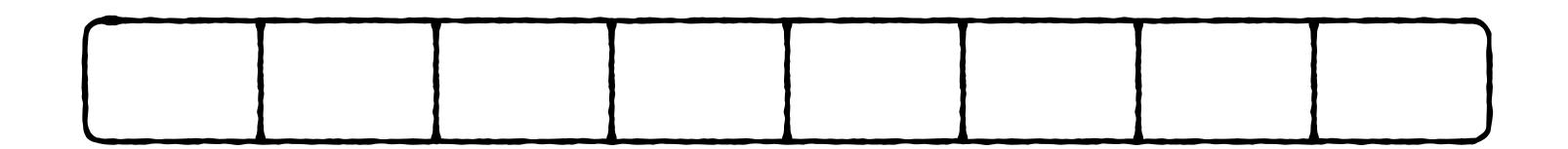
 $\approx 2^{-M+3}$

2-M · In(2)

Memory accesses for positive query

1.5

M · In(2)



false positive rate:

 $\approx 2^{-M+3}$

 $2-M \cdot ln(2)$

Memory accesses for positive query

1.5

 $M \cdot ln(2)$

Memory accesses for negative query?

false positive rate:

 $\approx 2^{-M+3}$

2-M · In(2)

Memory accesses for positive query

1.5

 $M \cdot ln(2)$

Memory accesses for negative query

2

≈2

Cuckoo Filter

Bloom filter

false positive rate:

 $\approx 2\text{-M}+3$

2-M · In(2)

Memory accesses for positive query

1.5

 $M \cdot ln(2)$

Memory accesses for negative query

2

 ≈ 2

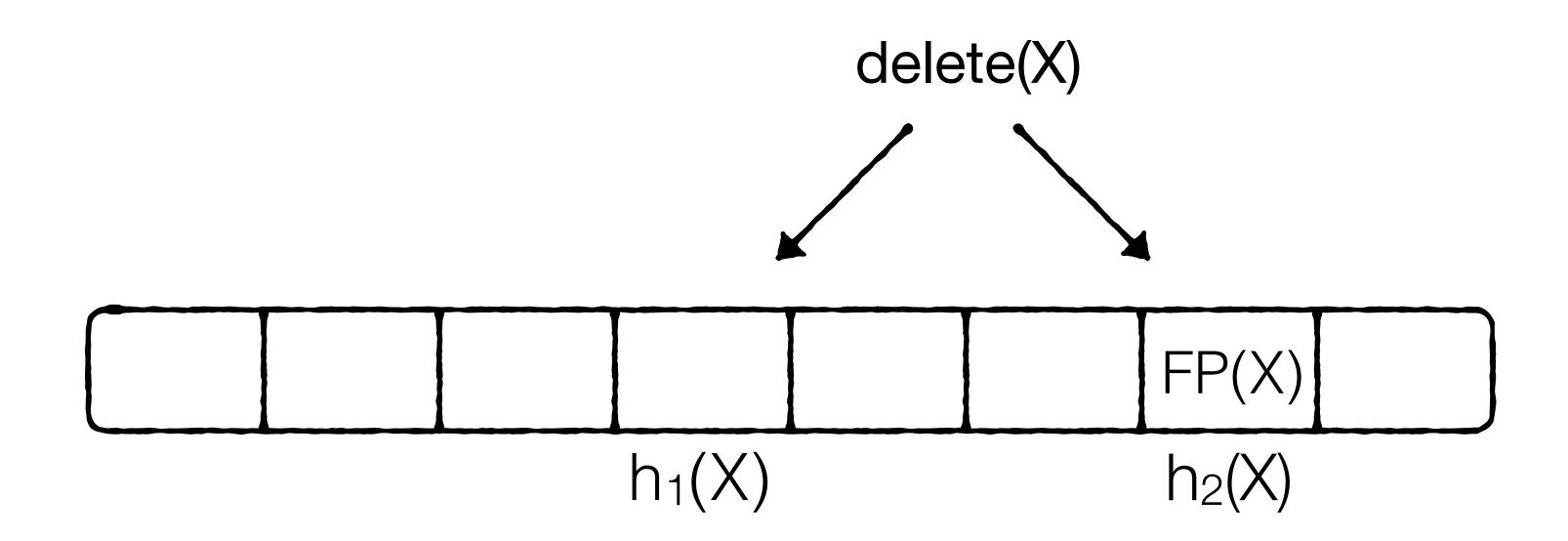
Insertion cost

Exp. O(1)

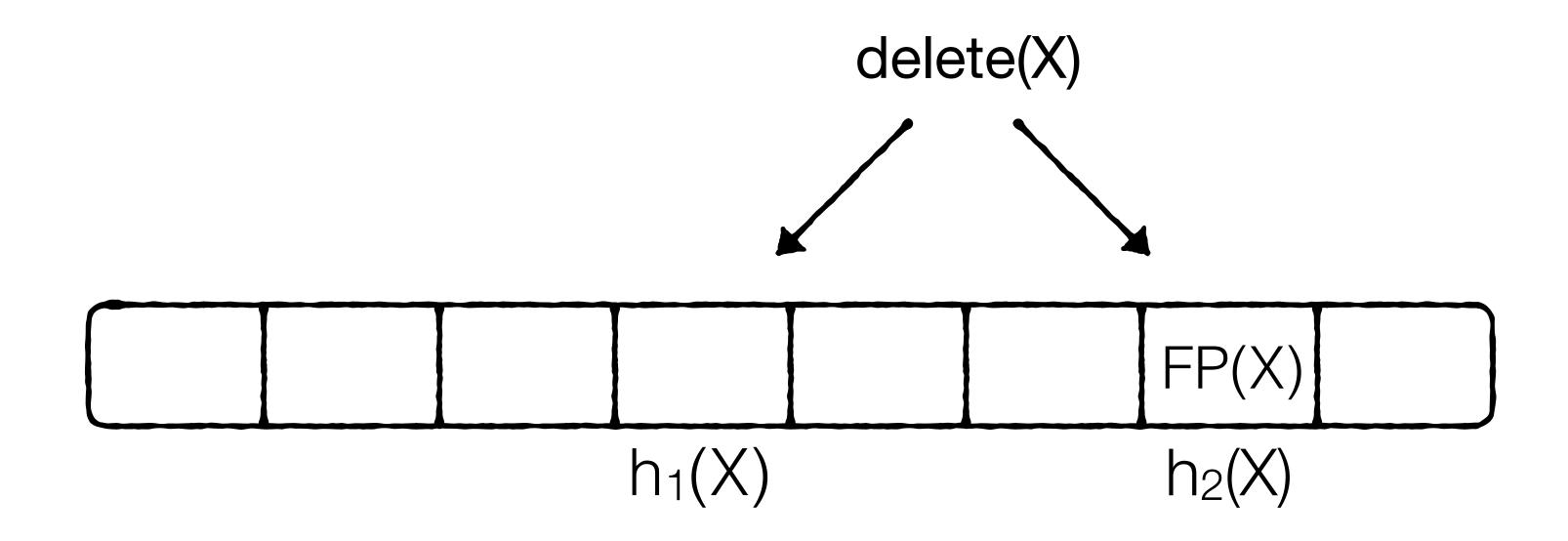
 $M \cdot ln(2)$

| | Cuckoo Filter | Bloom filter |
|------------------------------------|-----------------|-----------------|
| false positive rate: | ≈ 2 -M+3 | 2-M · In(2) |
| Memory accesses for positive query | 1.5 | M · In(2) |
| Memory accesses for negative query | 2 | ≈2 |
| Insertion cost | Exp. O(1) | $M \cdot ln(2)$ |
| Deletes? | | N/A |

Storing payloads? N/A

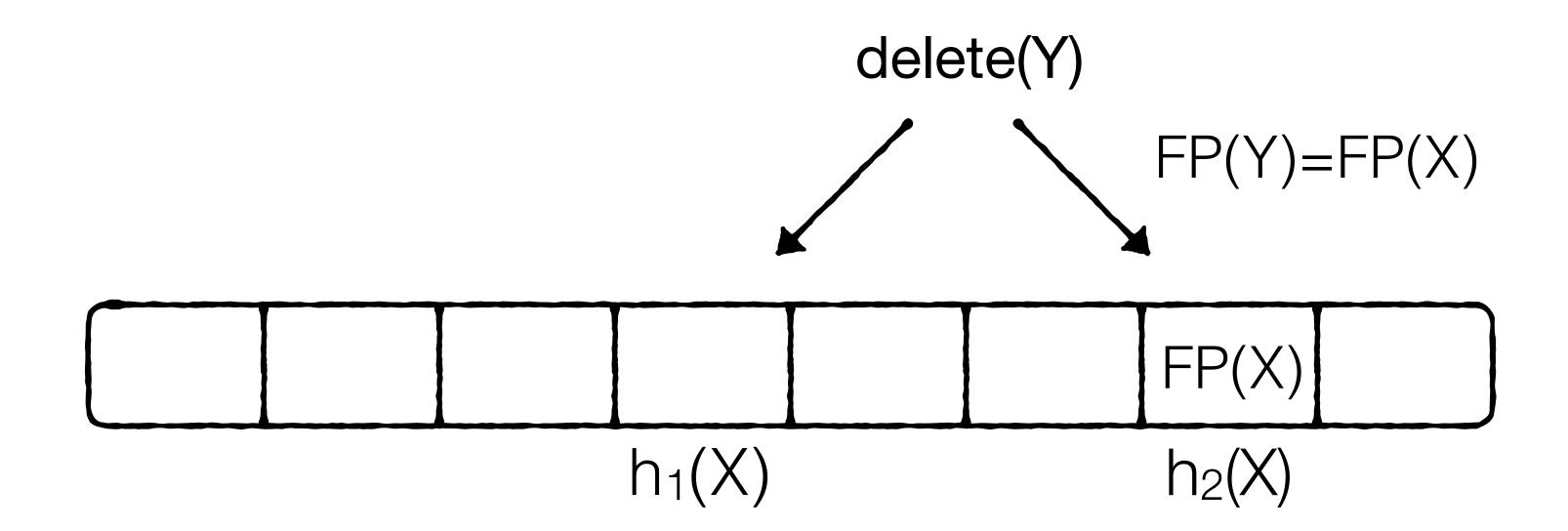


A delete searches both buckets for a key and removes a matching fingerprints.



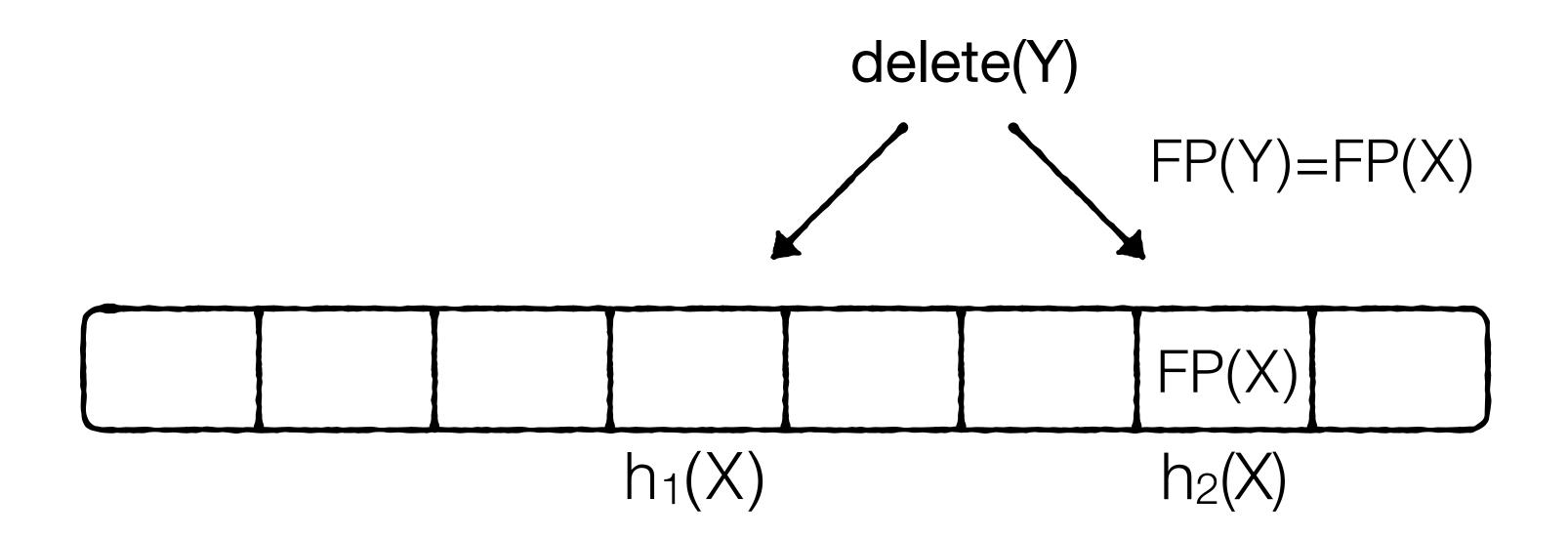
A delete searches both buckets for a key and removes a matching fingerprints.

What if we delete a fingerprint for a key that was never inserted?



A delete searches both buckets for a key and removes a matching fingerprints.

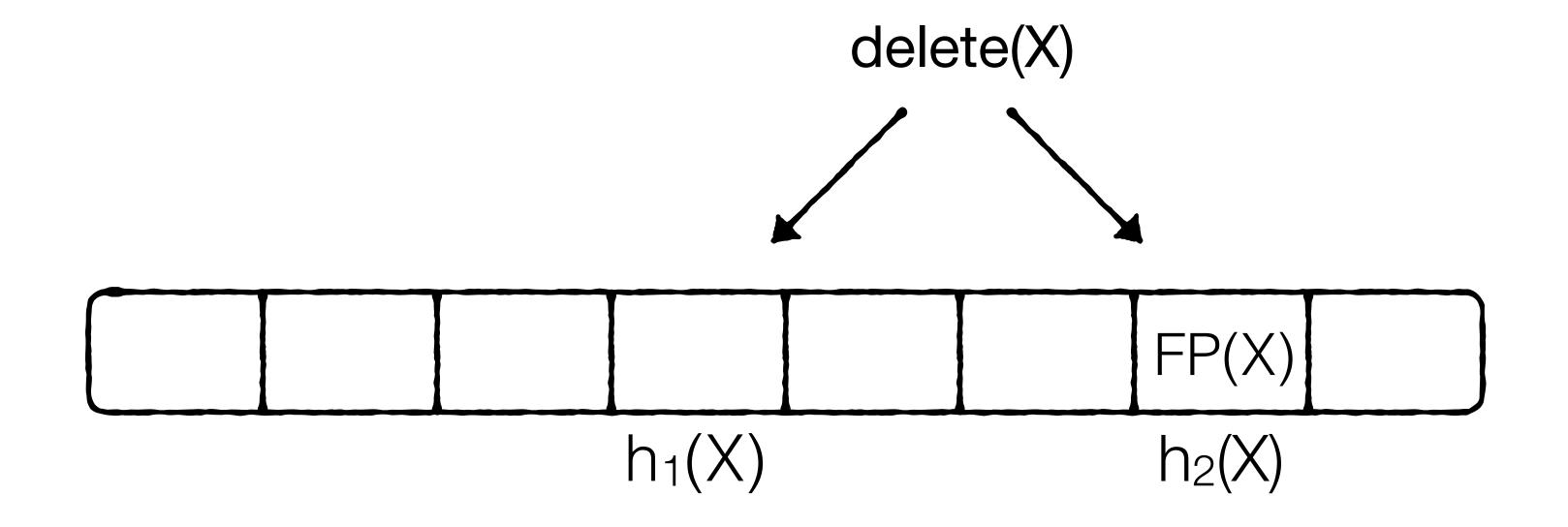
What if we delete a fingerprint for a key that was never inserted? False negatives...



A delete searches both buckets for a key and removes a matching fingerprints.

What if we delete a fingerprint for a key that was never inserted? False negatives...

As long as we delete keys we know for sure have been inserted, no false negatives can occur

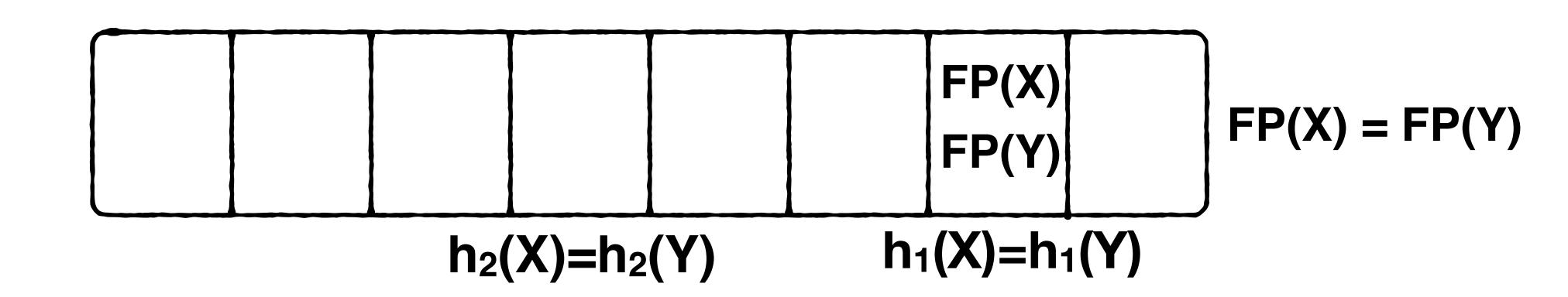


A delete searches both buckets for a key and removes a matching fingerprints.

What if we delete a fingerprint for a key that was never inserted? False negatives...

As long as we delete keys we know for sure have been inserted, no false negatives can occur

Works even if we have matching fingerprints for different keys in the same bucket



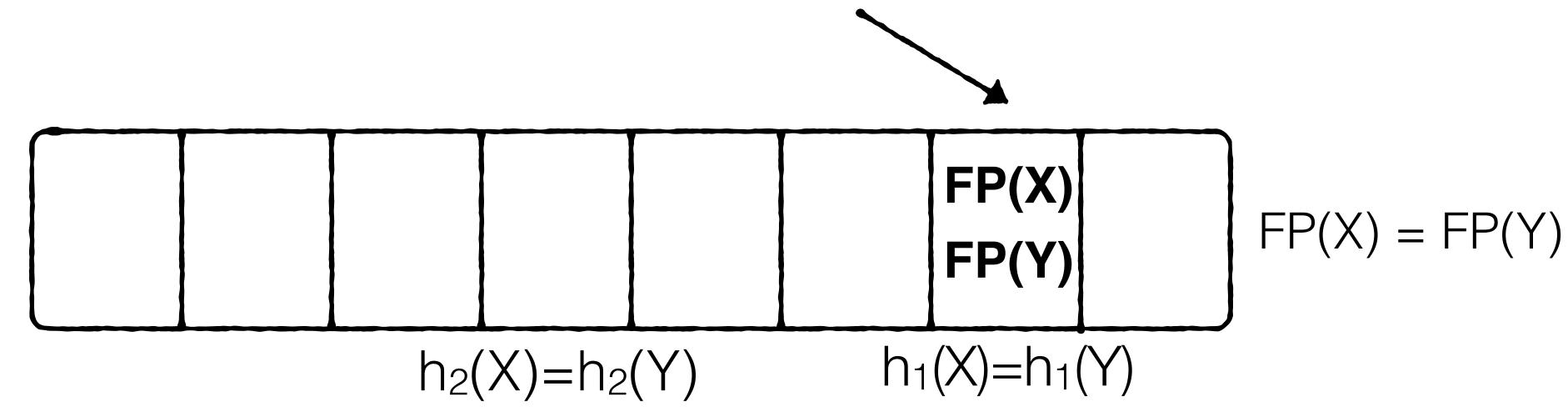
A delete searches both buckets for a key and removes a matching fingerprints.

What if we delete a fingerprint for a key that was never inserted? False negatives...

As long as we delete keys we know for sure have been inserted, no false negatives can occur

Works even if we have matching fingerprints for different keys in the same bucket

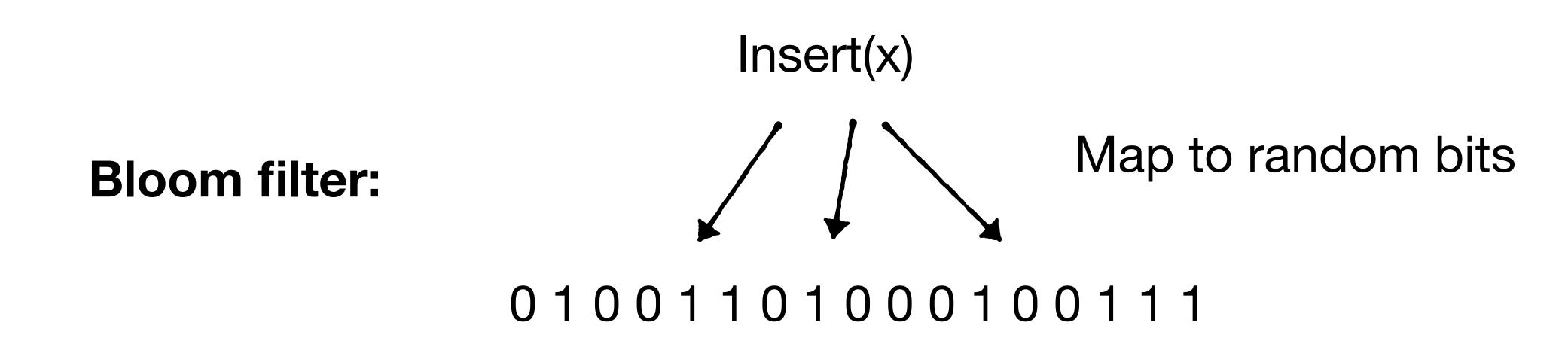
After delete(X), get(Y) will still succeed, whichever fingerprint we delete.



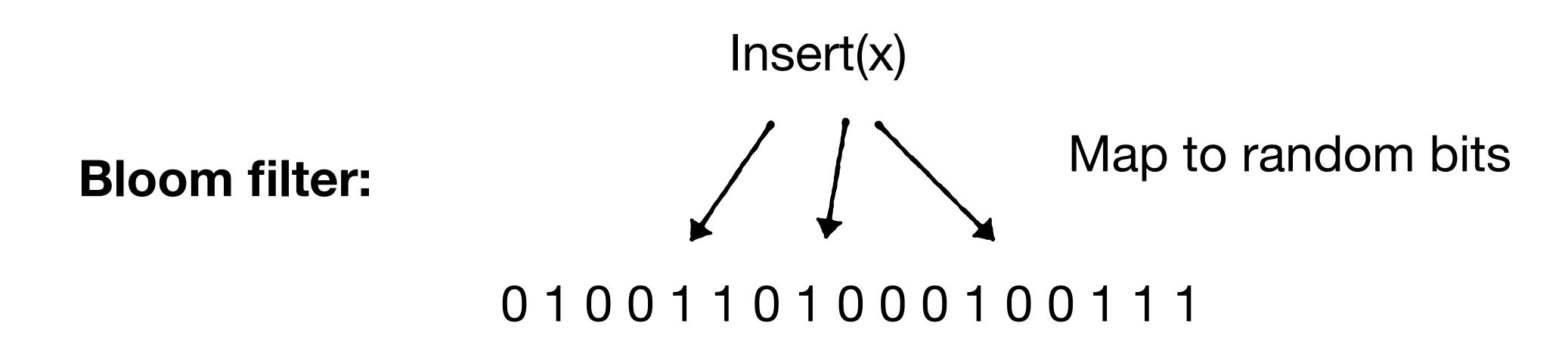
| | Cuckoo Filter | Bloom filter |
|------------------------------------|-----------------|---------------------|
| false positive rate: | ≈ 2 -M+3 | 2-M · In(2) |
| Memory accesses for positive query | 1.5 | M · In(2) |
| Memory accesses for negative query | 2 | 2 |
| Insertion cost | O(1) | M · In(2) |
| Deletes? | 1.5 | N/A |
| Storing Payloads? | | N/A |

Can we store a payload associated with each key and retrieve it on positive query?

Can we store a payload associated with each key and retrieve it on positive query?



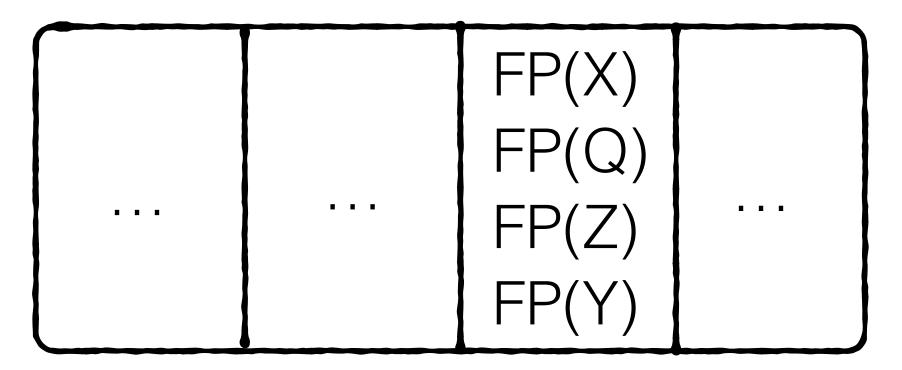
Can we store a payload associated with each key and retrieve it on positive query?



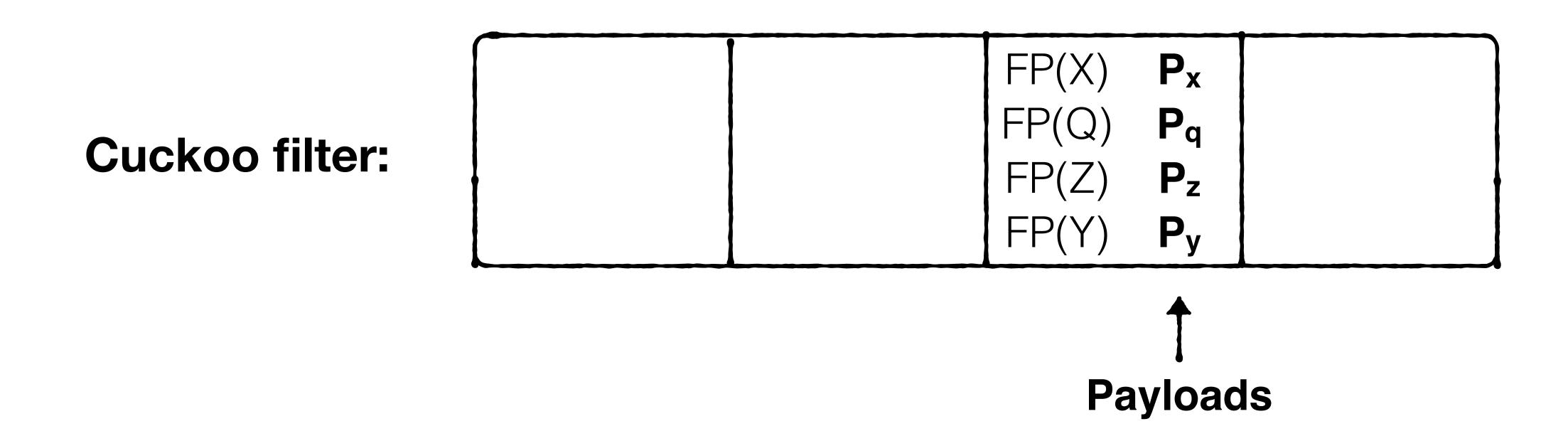
There is nowhere obvious to associate a payload with each key

Can we store a payload associated with each key and retrieve it on positive query?

Cuckoo filter:

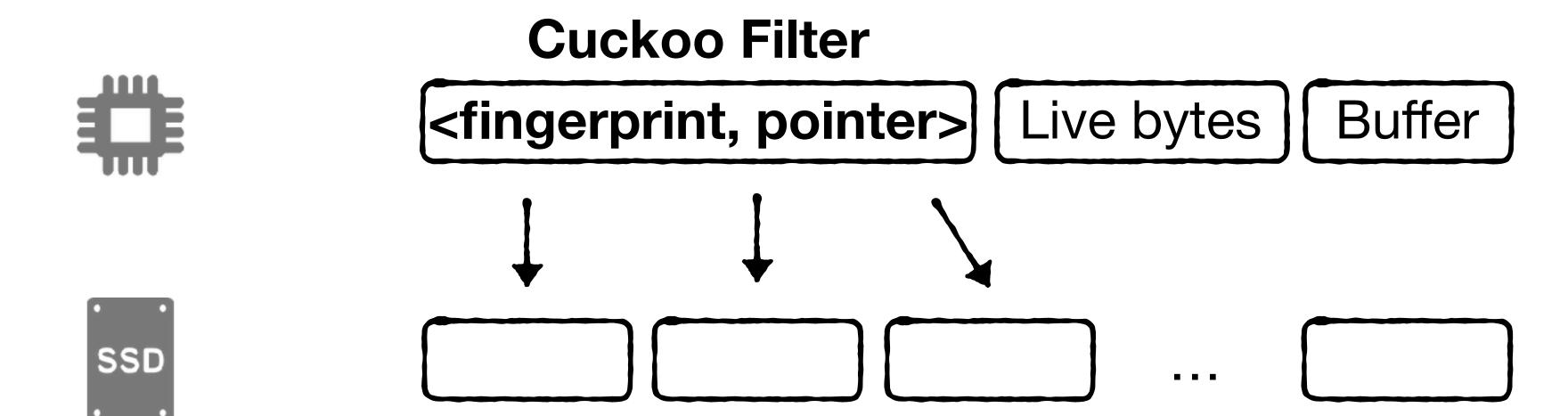


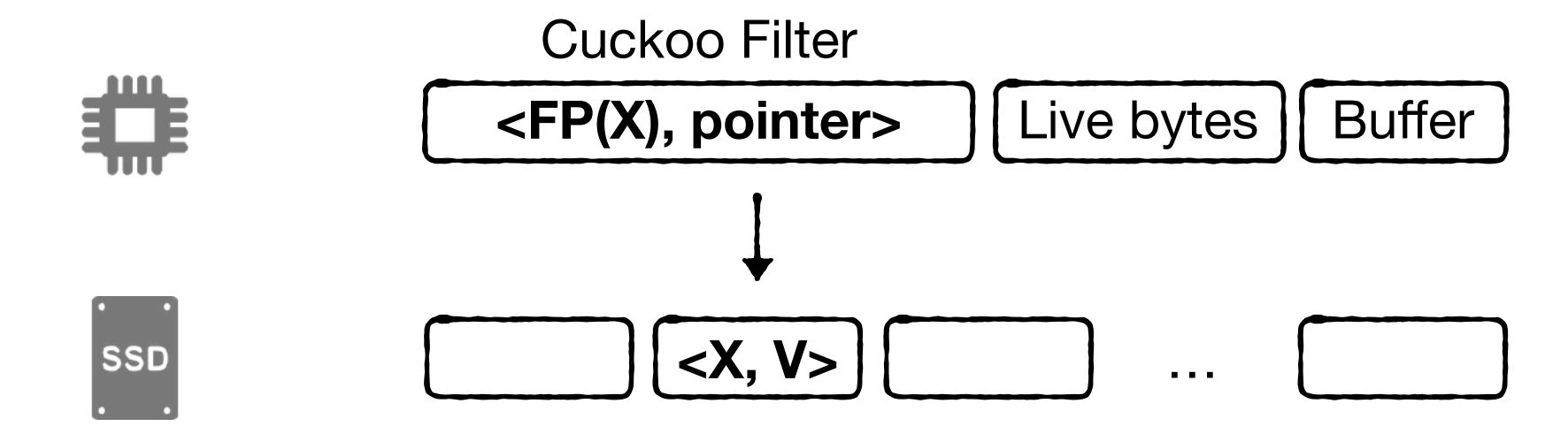
Can we store a payload associated with each key and retrieve it on positive query?



Just store a payload alongside each fingerprint

| | Cuckoo Filter | Bloom filter |
|------------------------------------|-----------------|--------------|
| false positive rate: | ≈ 2 -M+3 | 2-M · In(2) |
| Memory accesses for positive query | 1.5 | M·In(2) |
| Memory accesses for negative query | 2 | 2 |
| Insertion cost | O(1) | M · In(2) |
| Deletes? | 1.5 | N/A |
| Storing Payloads? | Yes | N/A |





Index size = $N * (P + K) * \alpha$

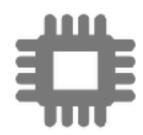
N = data size

 $P = pointer size = O(log_2 N/B)$

Our goal was reducing this

 \longrightarrow K = key size = $\Omega(\log_2 N)$

 $\alpha = \text{collision resolution overheads } \approx 0.8 \longrightarrow \approx 0.95$



Cuckoo Filter

Live bytes

Buffer



Index size = $N * (P + M) / \alpha$

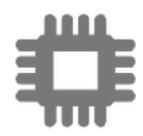
N = data size

 $P = pointer size = O(log_2 N/B)$

Our goal was reducing this

 \longrightarrow K = key size = $\Omega(\log_2 N)$ \longrightarrow M bits / entry

 $\alpha = \text{collision resolution overheads } \approx 0.8 \longrightarrow \approx 0.95$

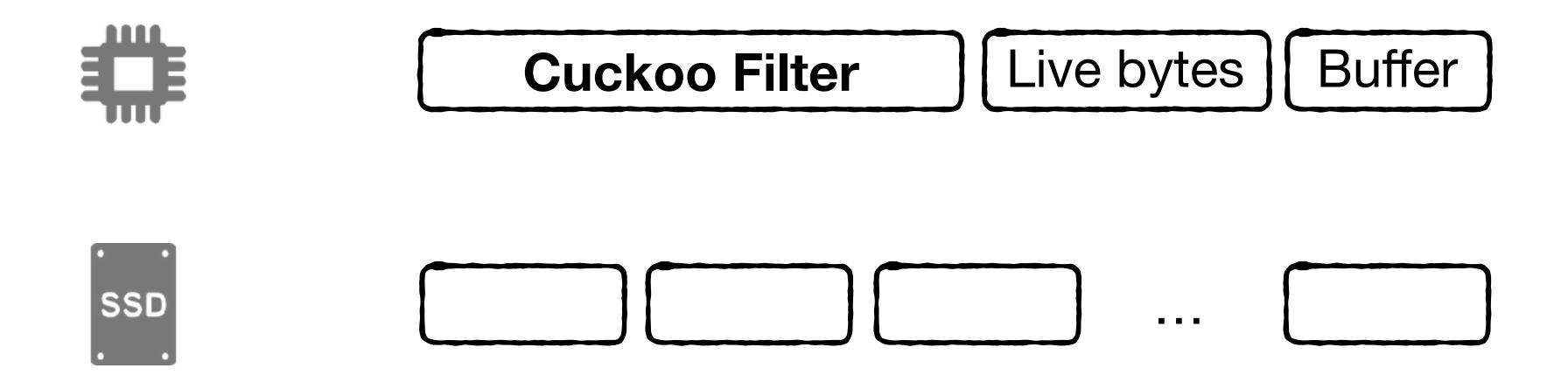


Cuckoo Filter

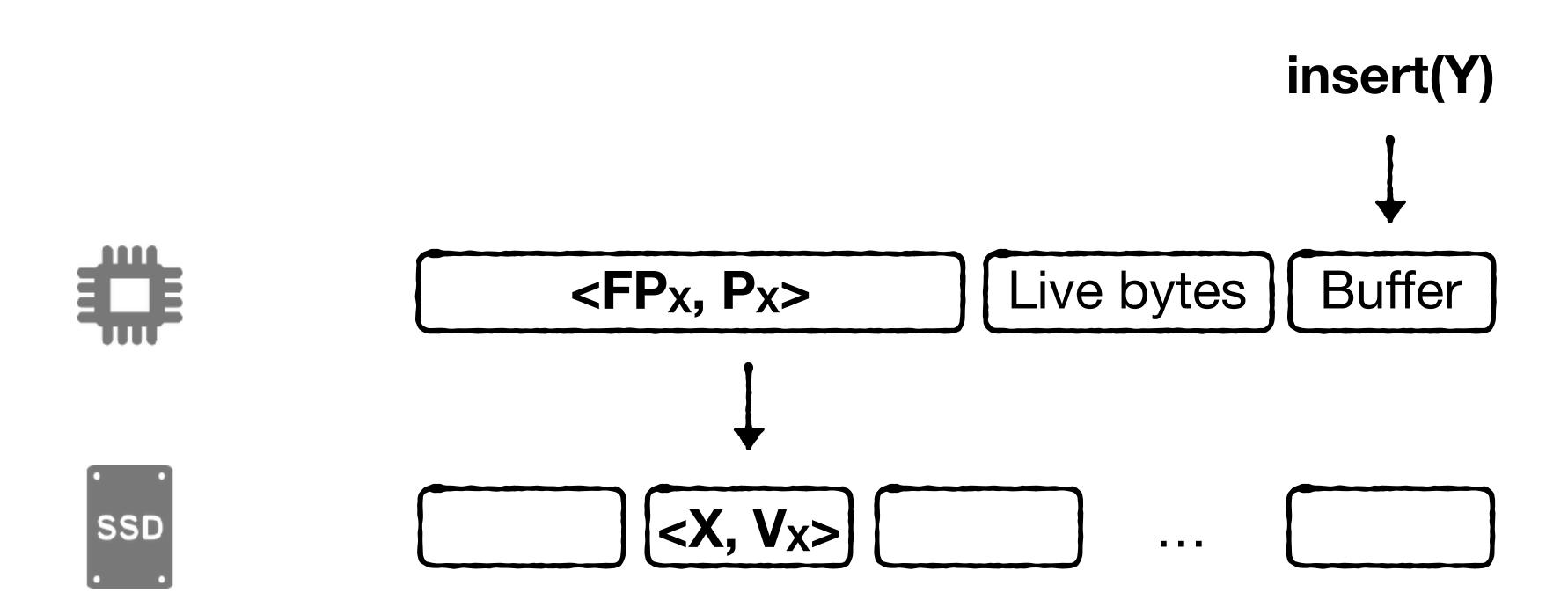
Live bytes

Buffer

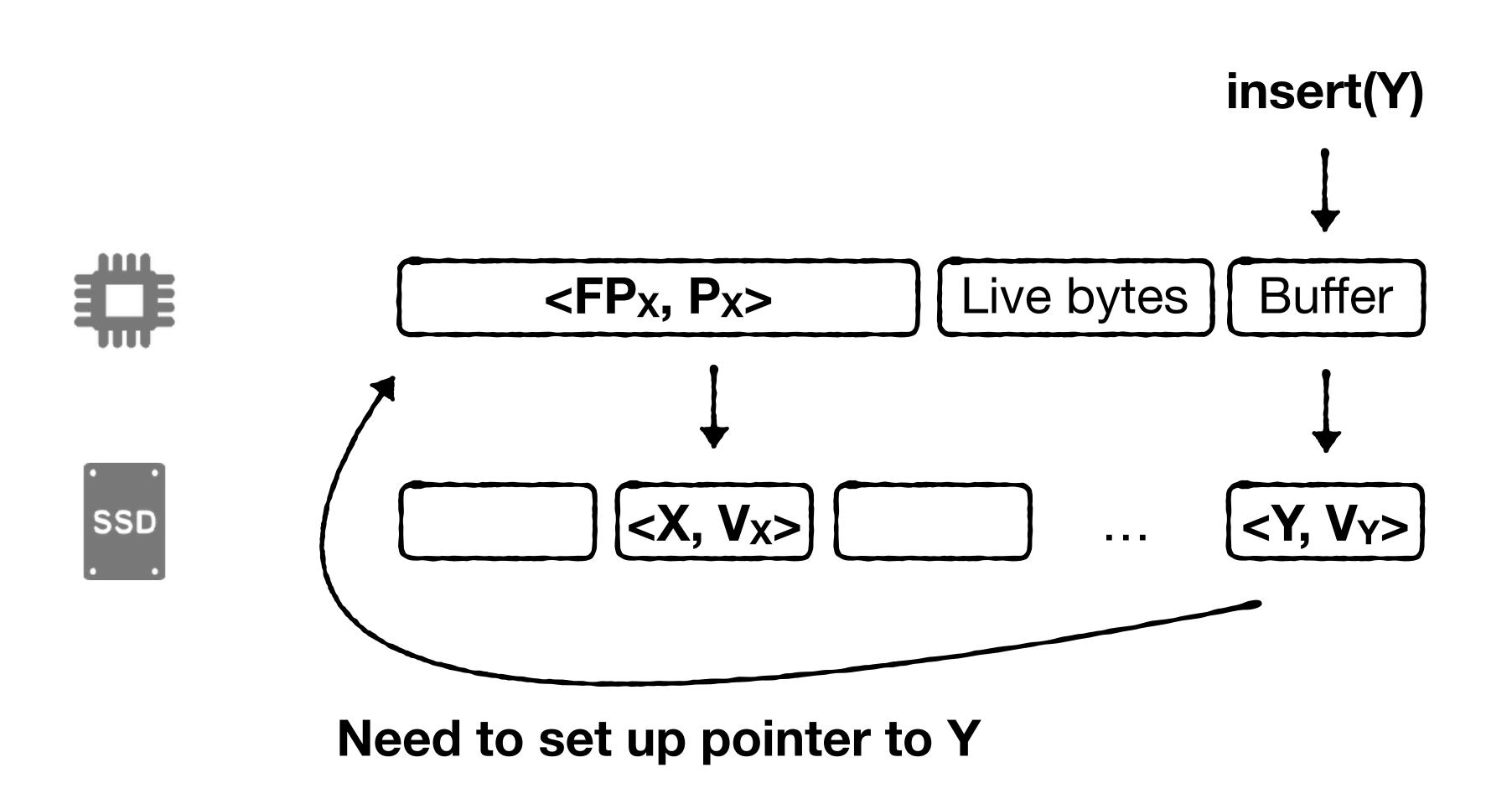




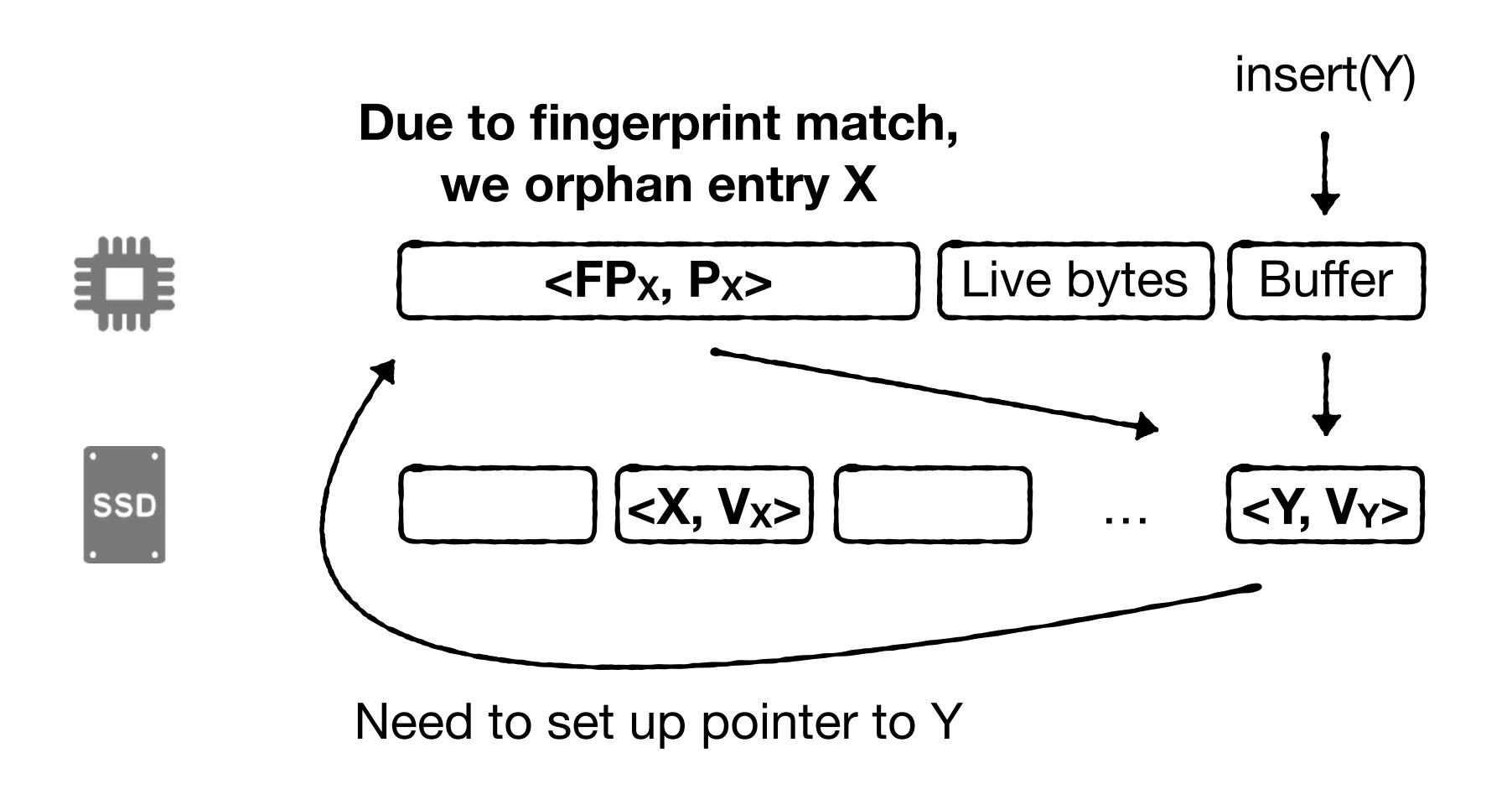
Suppose we insert key Y that has a matching fingerprint to existing key X (FP_{X=}FP_y)



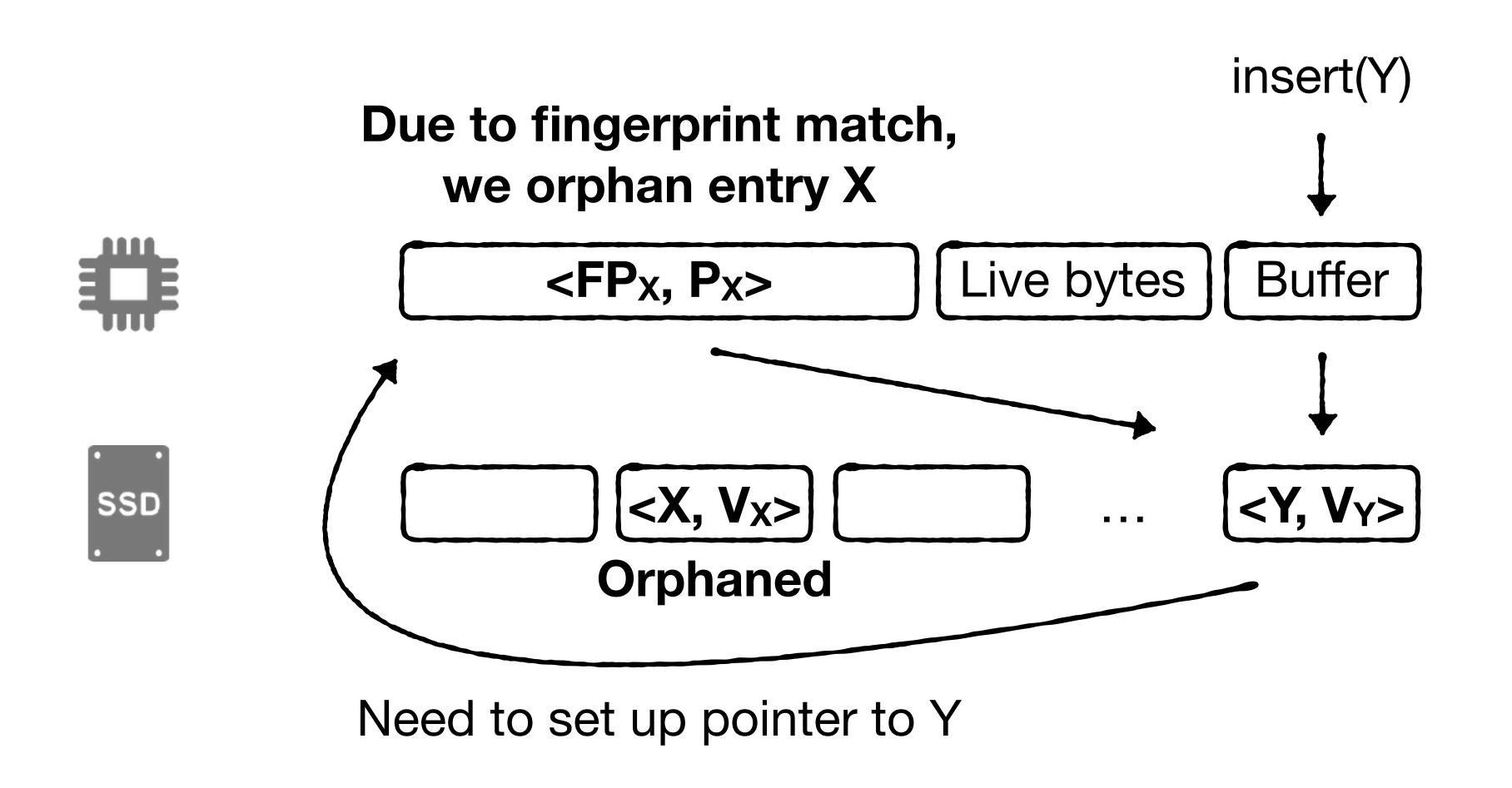
Suppose we insert key Y that has a matching fingerprint to existing key X



Suppose we insert key Y that has a matching fingerprint to existing key X

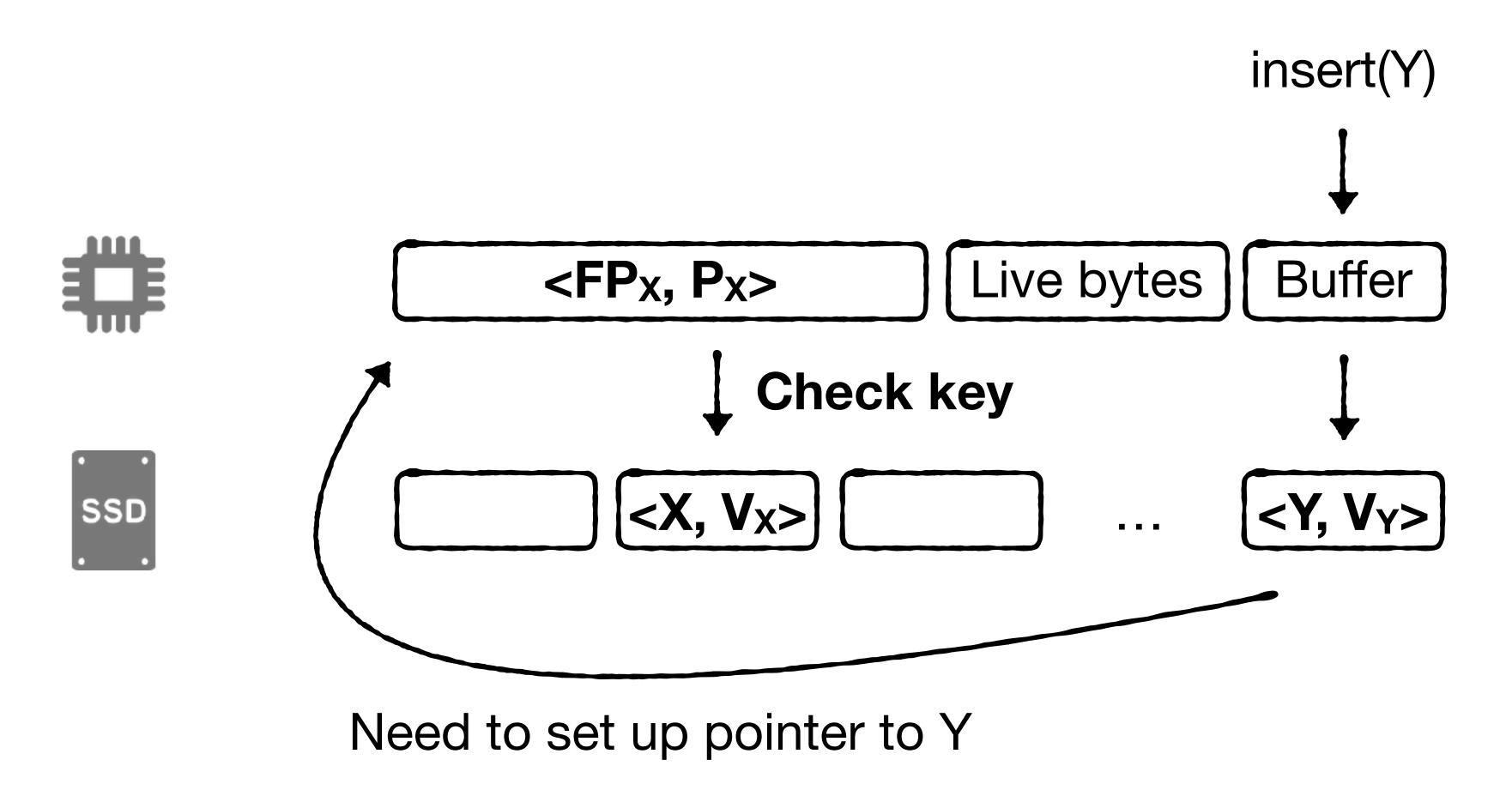


Suppose we insert key Y that has a matching fingerprint to existing key X

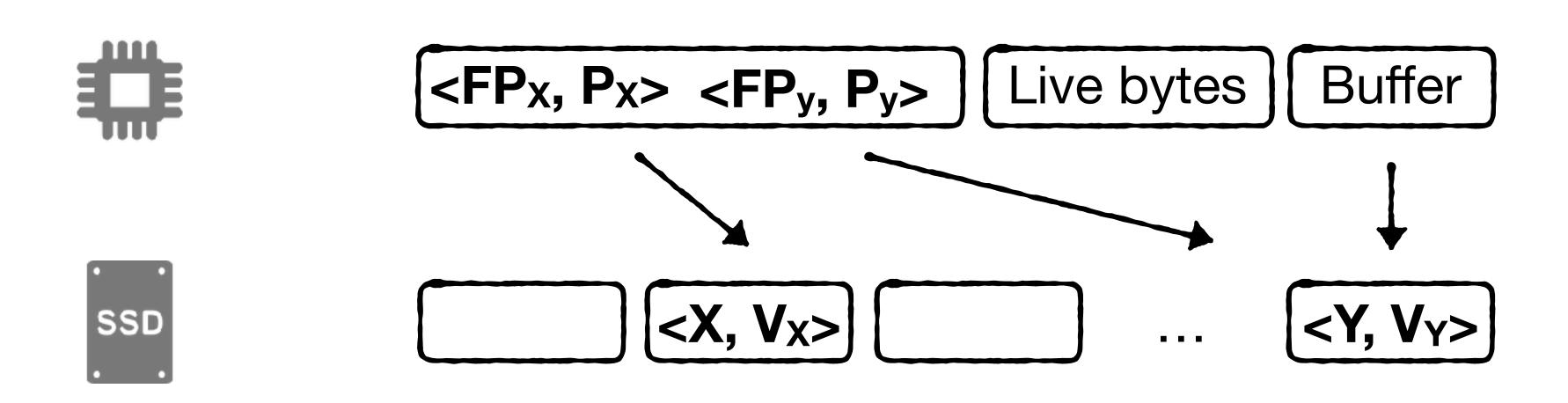


Suppose we insert key Y that has a matching fingerprint to existing key X

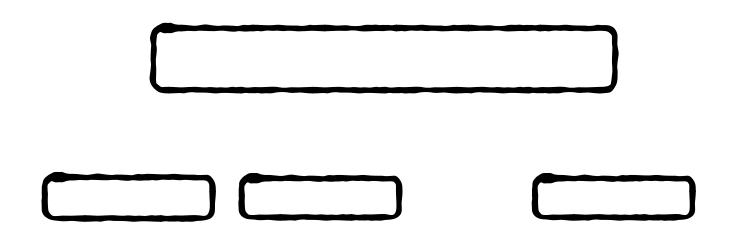
To safeguard against orphaning, we must issue read-before-write



Suppose we insert key Y that has a matching fingerprint to existing key X To safeguard against orphaning, we must issue read-before-write



Circular Log w. Cuckoo Filter



Updates/

Deletes:

O(1+2-M+3) reads & O(GC/B) writes

O(2-M+3) read & O(GC/B) writes

Gets:

Insert:

O(1+2-M+3)

Scan:

O(N/B)

Memory

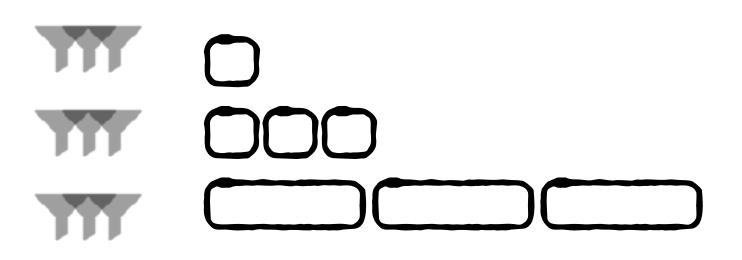
(bits / entry)

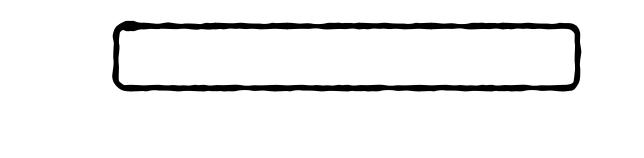
 $O(log_2(N/B) + M))$

(excluding checkpointing)

Basic LSM-tree w. Monkey

Circular Log w. Cuckoo Filter





Updates/ Deletes:

O(log₂(N/P) / B) reads & writes

O(1+2-M+3) reads & O(GC/B) writes

Insert:

O(log₂(N/P) / B) reads & writes

O(2-M+3) read & O(GC/B) writes

Gets:

O(1+2-M*In(2))

O(1+2-M+3)

Scan:

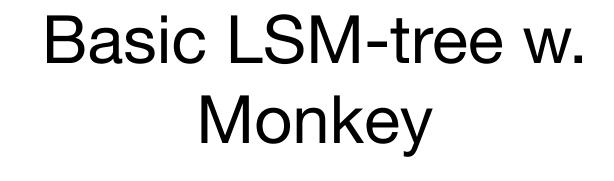
 $O(log_2 N/P + S/B)$

O(N/B)

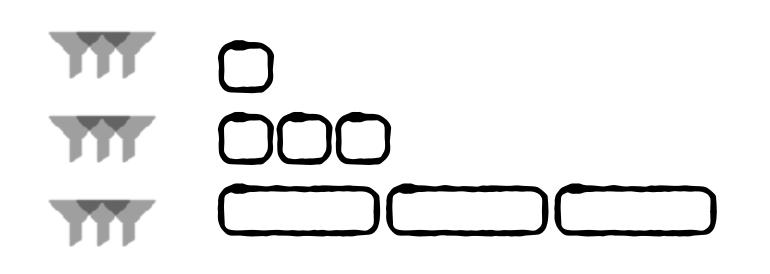
Memory (bits / entry)

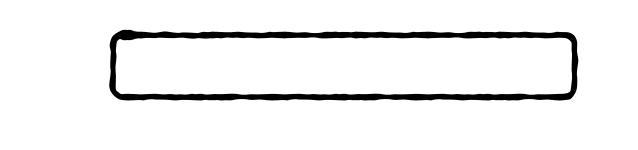
O(K/B + M)

 $O(log_2(N/B) + M))$



Circular Log w. Cuckoo Filter





Updates/ Deletes:

O(log₂(N/P) / B) reads & writes

O(1+2-M+3) reads & O(GC/B) writes

Insert:

O(log₂(N/P) / B) reads & writes

 $O(2^{-M+3})$ read & O(GC/B) writes $O(1+2^{-M+3})$

Gets:

O(1+2-M*In(2))

 $O(log_2 N/P + S/B)$

O(N/B)

Scan:

O(K/B + M)

 $O(log_2(N/B) + M))$

(bits / entry)
Recovery

Memory

Swift

Long

And now: office hours