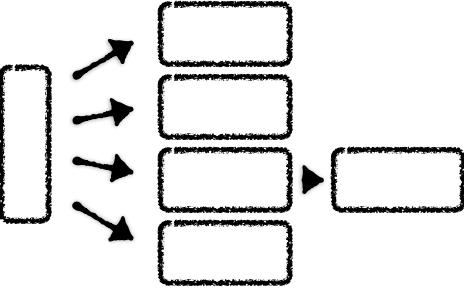


# **Circular Logs & Cuckoo Filters**

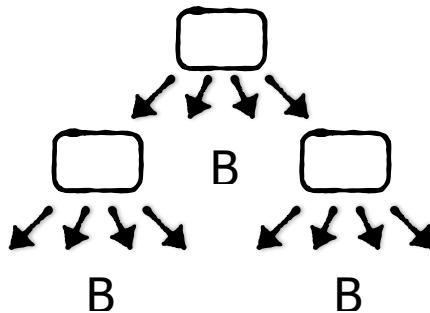
## **CSC443H1 Database System Technology**

**Niv Dayan - March 6, 2023**

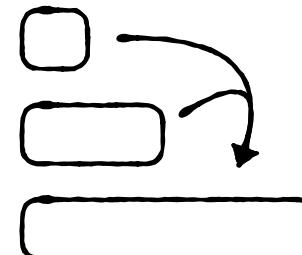
## Extendible Hashing



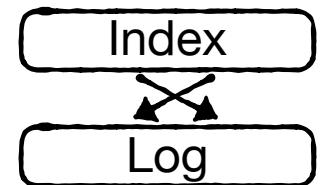
## B-Trees



## LSM-tree



## Circular Log



Cheapest gets  
No scans

Cheap gets  
& scans

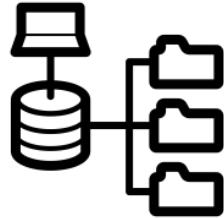
Cheaper writes  
More memory

Cheapest writes  
Most memory  
No scans

## **Circular Log**

Invented in 1992 as a “log structured file system”

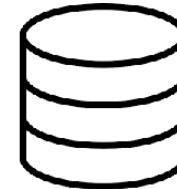
### **File Systems**



### **Flash Translation Layers**



### **KV-stores**



Various names, same data structure:

**Index+Log**

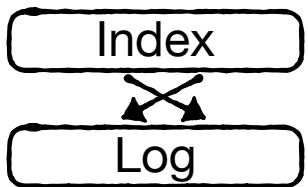
**Circular Log**

**Log-structured Hash Table**

**log structured file system**

# Agenda

Mechanics



Hot/Cold separation



To reduce write-amp

Checkpointing/  
recovery



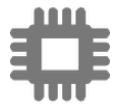
If power fails

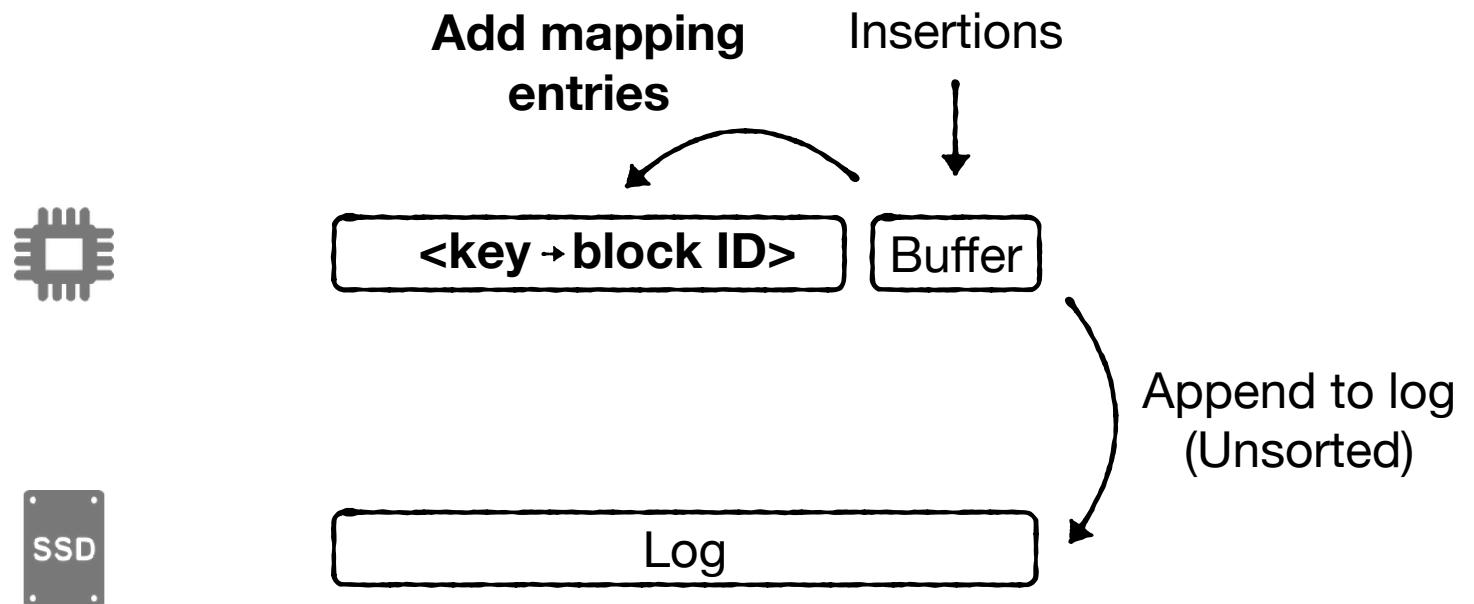
Cuckoo filtering



To reduce memory

# Mechanics





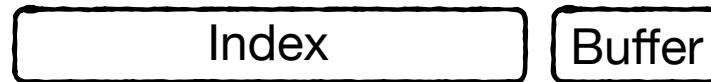
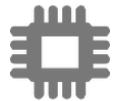
**Write cost**

$O(1/B)$

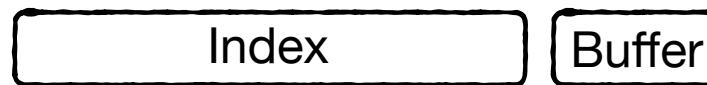
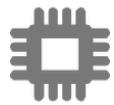
**(Assuming only insertions)**

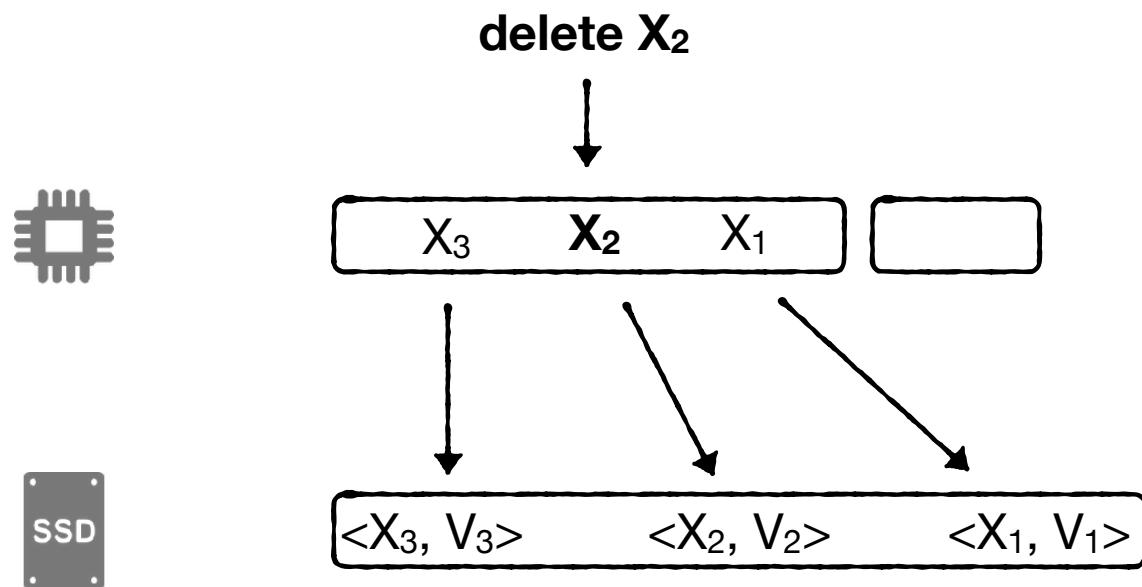
**Get cost**

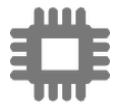
$O(1)$



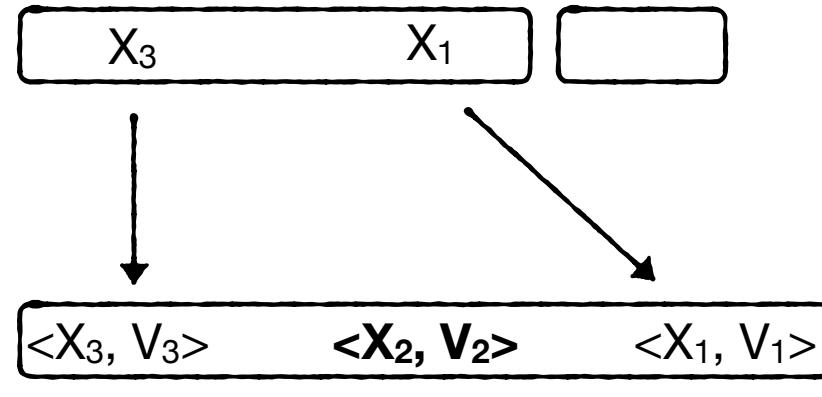
# How to delete?







SSD

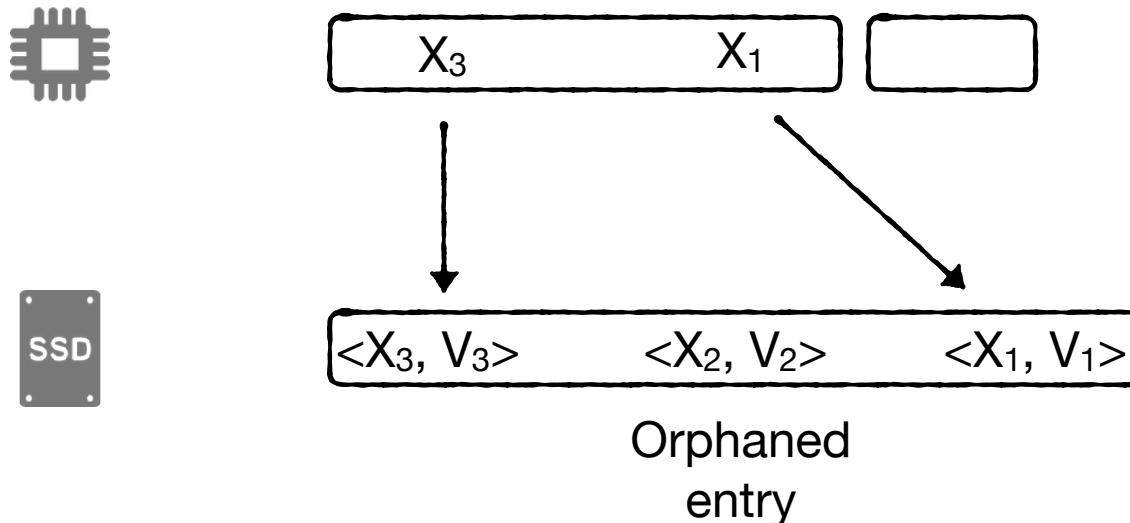


Orphaned  
entry

After many deletes, many orphaned entries accumulate

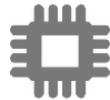
They take up space which we would prefer to use to store valid data

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

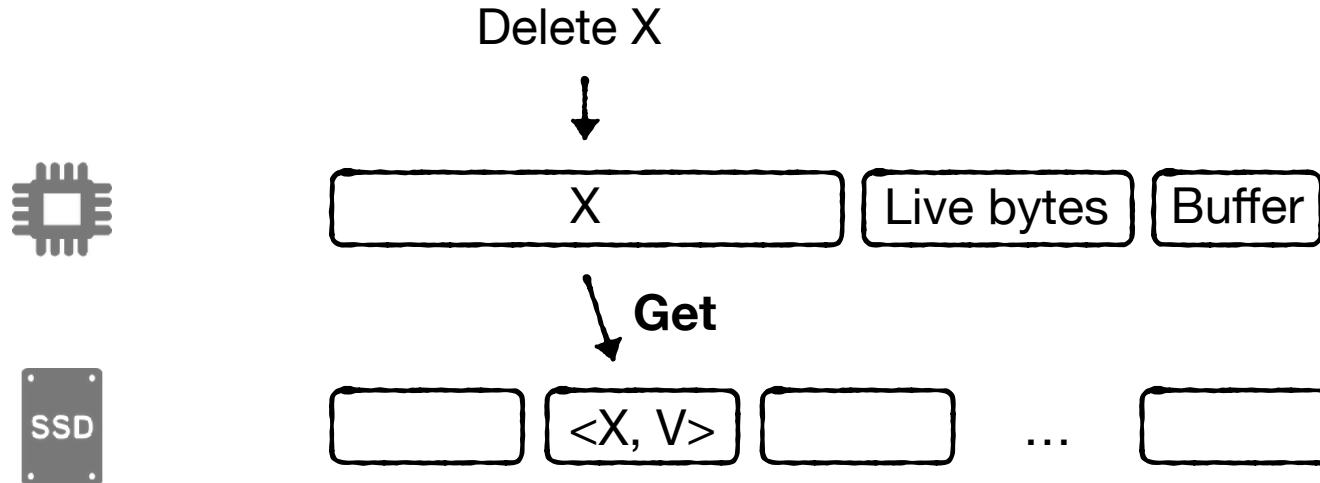


Conceptually divide log into equally-sized areas (can be in order of GBs)

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

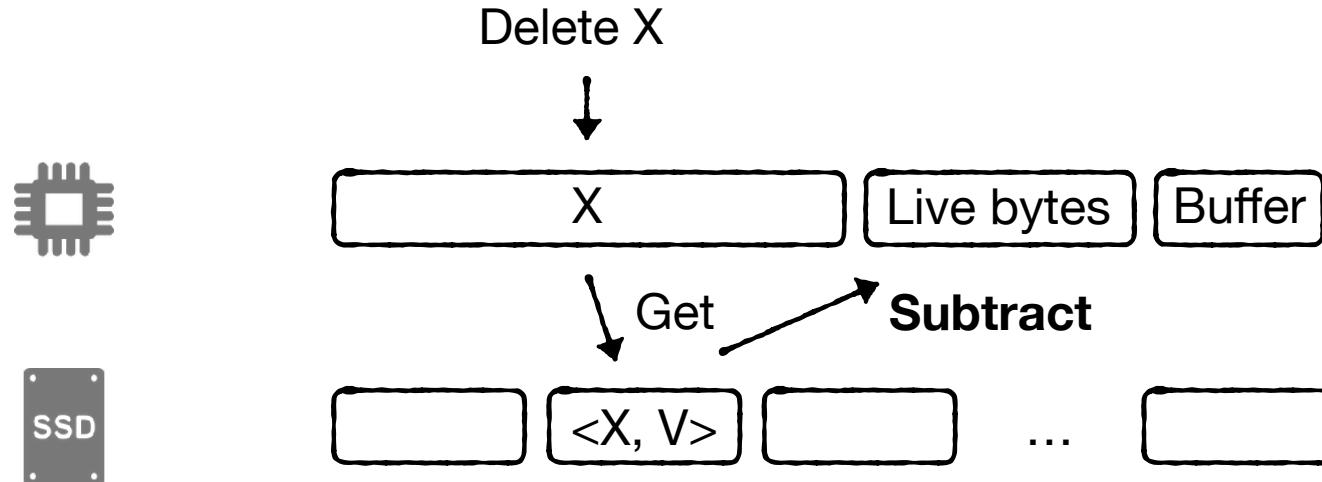


To delete: **(1) get entry from storage to check its size**



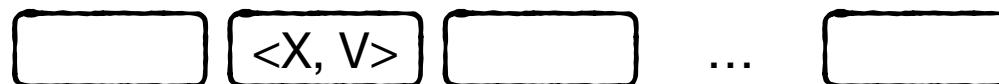
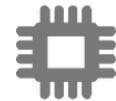
To delete: (1) get entry from storage to check its size

**(2) subtract its size from area's counter**



To delete:

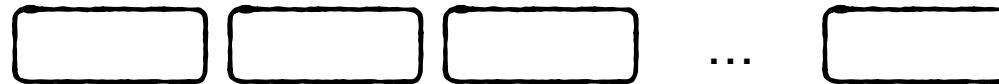
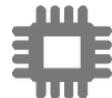
- (1) get entry from storage to check its size
- (2) subtract its size from area's counter
- (3) remove entry from index**



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

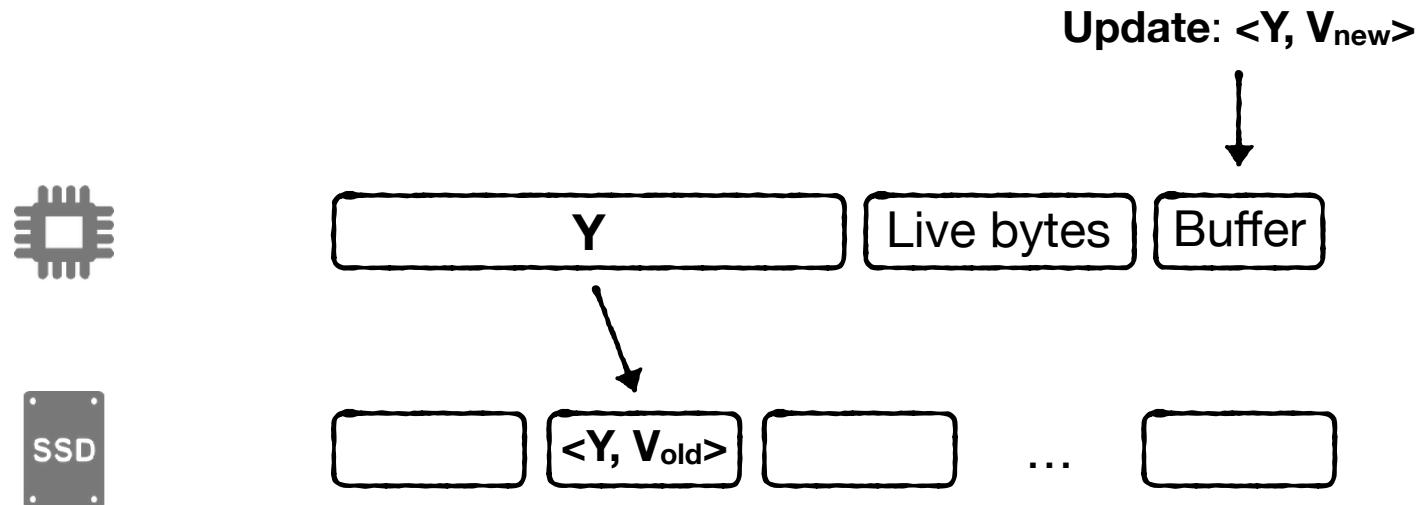
- (2) subtract its size from area's counter
- (3) remove entry from index

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



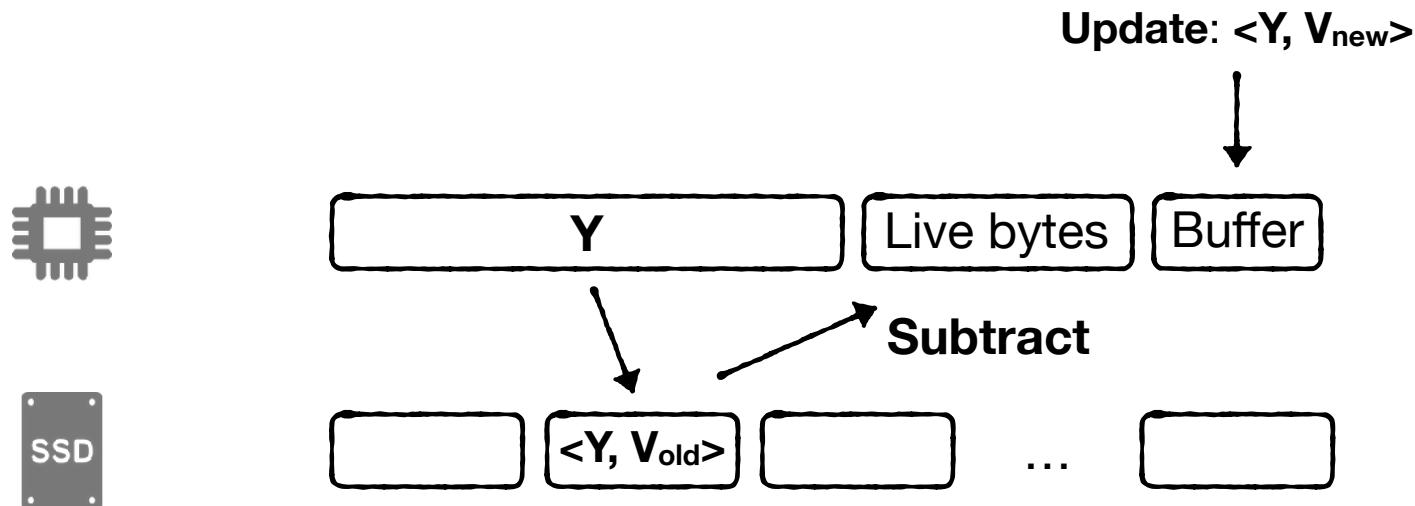
# Updates

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



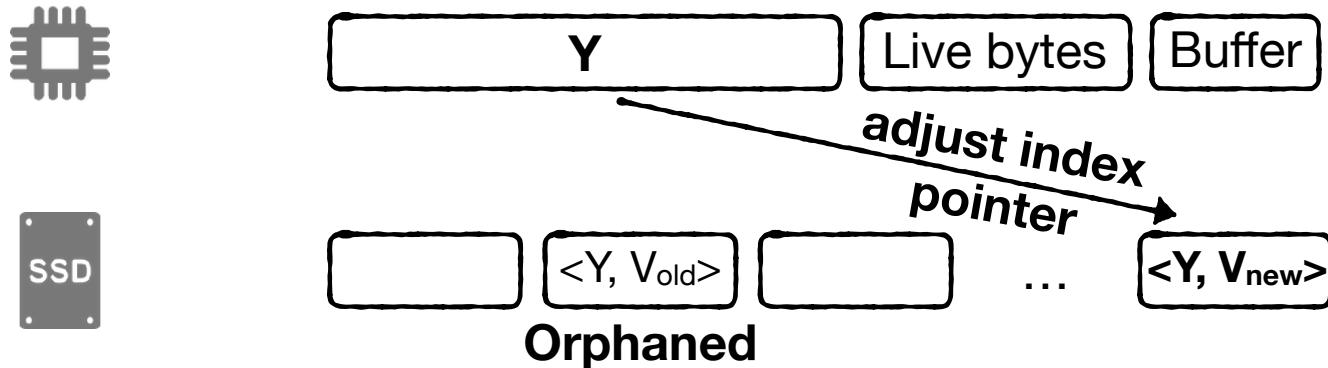
# Updates

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



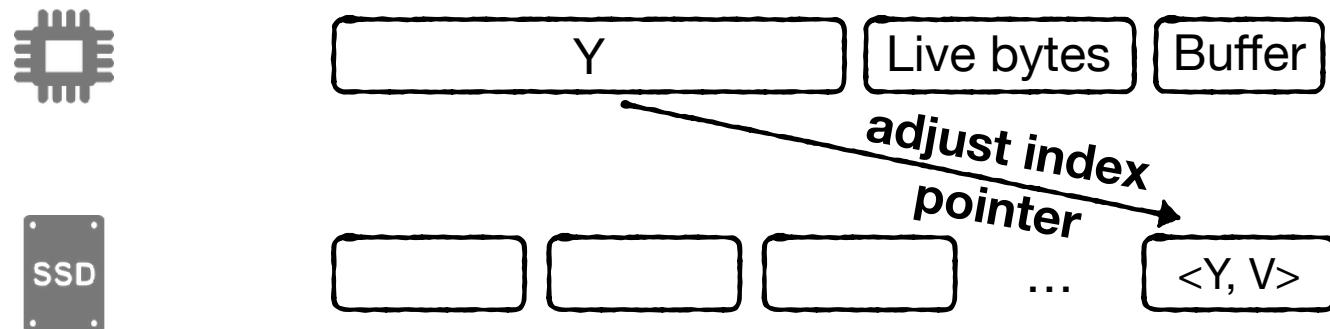
# Updates

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



# Garbage Collection

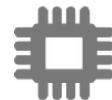
- (1) pick area with least live data left
- (2) Scan area and for each entry
  - (A) If the key is not indexed, the entry had been deleted, so move on
  - (B) If the key is indexed but pointing elsewhere, the entry is outdated, so move on
  - (C) If the key is indexed and pointing to this entry, it is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



## When to trigger garbage-collection?

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

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



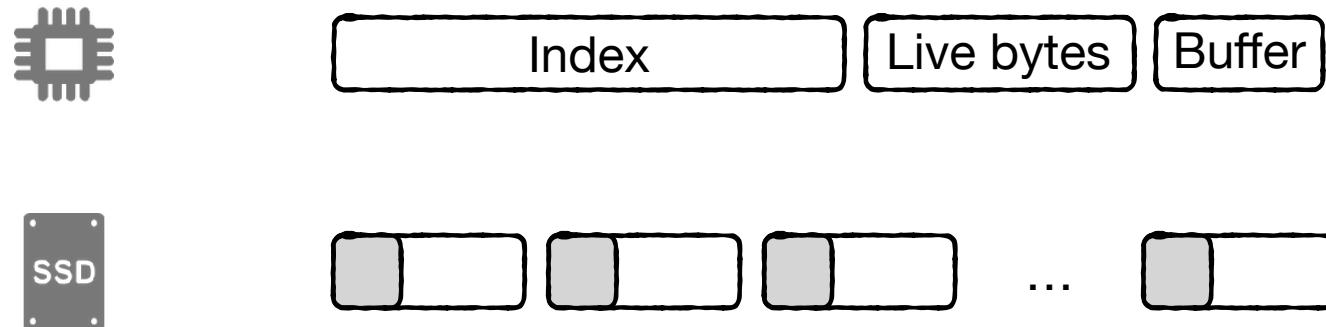
# Garbage-Collection Write-Amplification

How to reason about this?

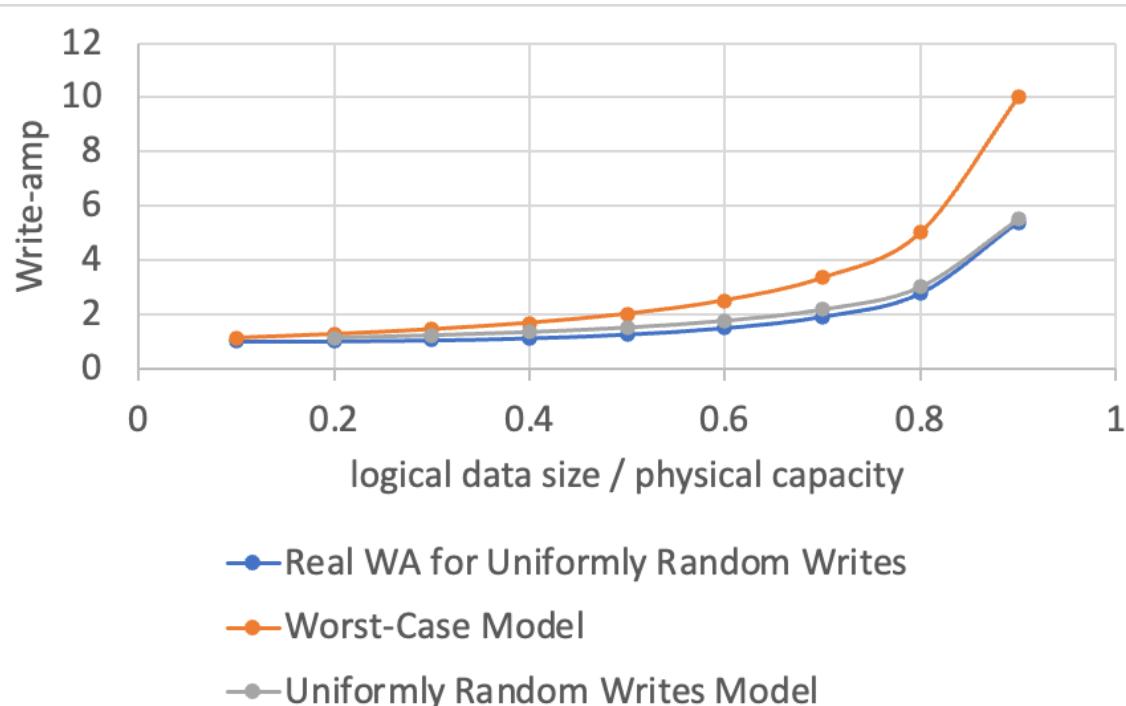
Let  $x$  = avg. % of valid pages in areas we pick to garbage-collect

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

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



# Garbage-Collection Write-Amplification



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

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

**Worst case**

**Uniformly random**

$L$  = logical data size

$P$  = physical data size

## **Hot/Cold Separation**

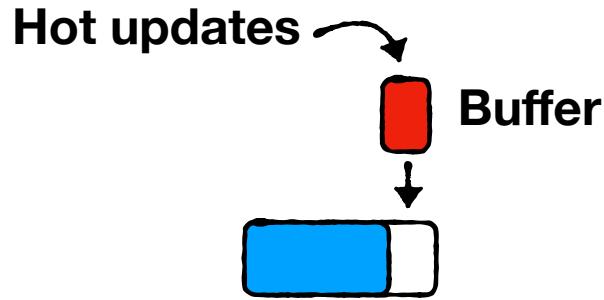
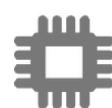
Normal workloads are neither worst-case nor randomly distributed

Normal workloads are neither worst-case nor randomly distributed

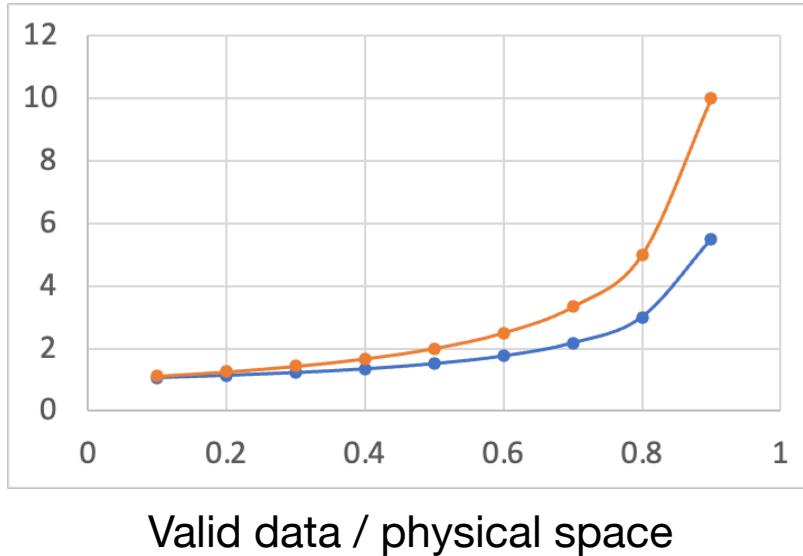
Typically, few entries are frequently updated while most are seldom updated

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.**



Write-amplification



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

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



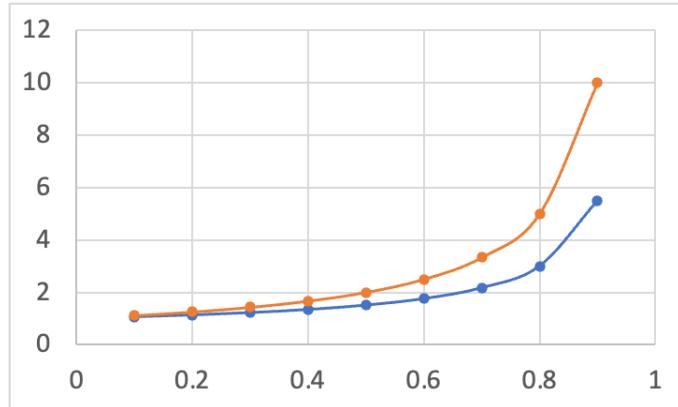
**Mixing hot/cold data brings us towards worst-case**



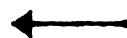
Uniformly random workloads

## Hot vs. Cold Data

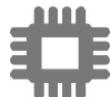
Write-amplification



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



Mixing hot/cold data brings us towards worst-case



How can we avoid garbage-collecting cold data all the time?



# Hot vs. Cold Data Separation

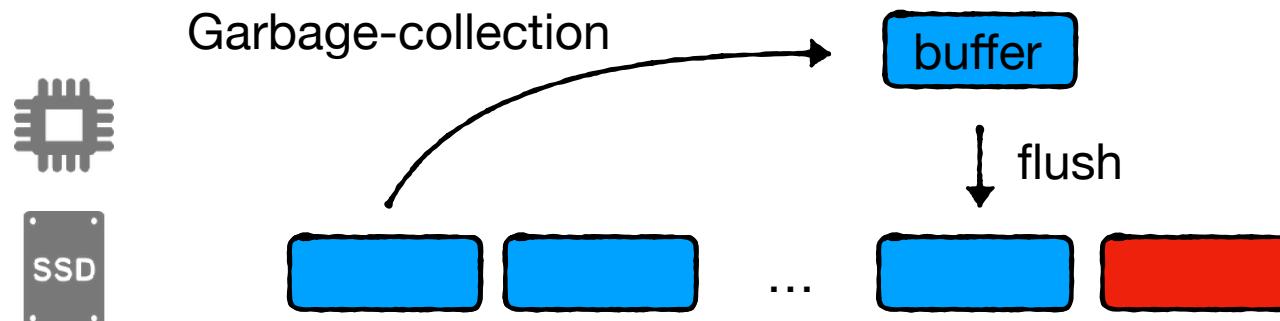
We can separate hot vs. cold data into different areas

**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

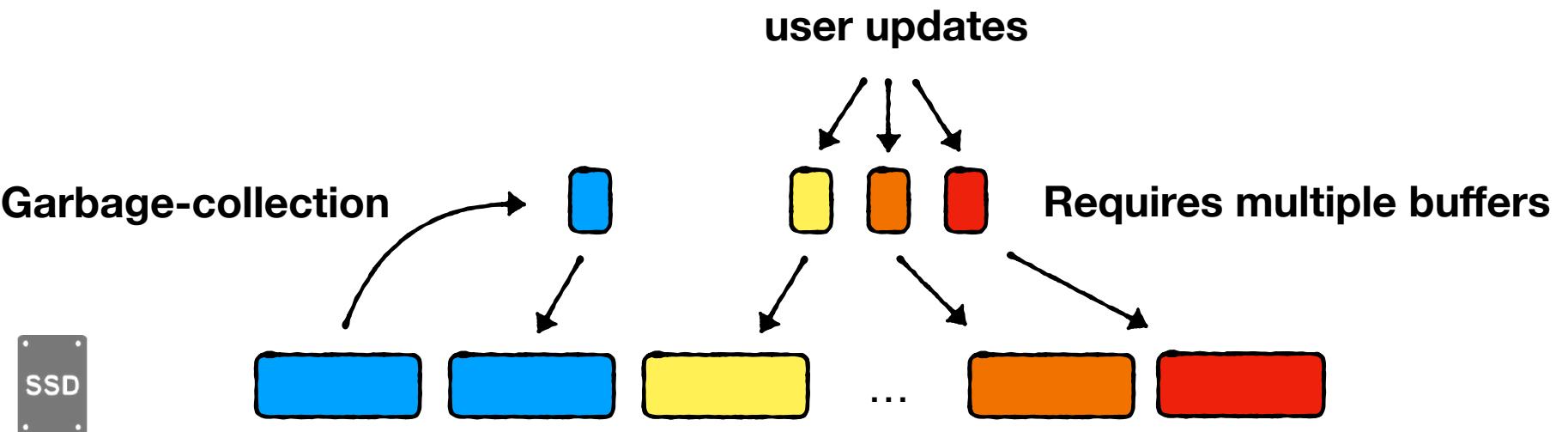
(i.e., it had already existed for a long while without getting updated)



# Hot vs. Cold Data Separation

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

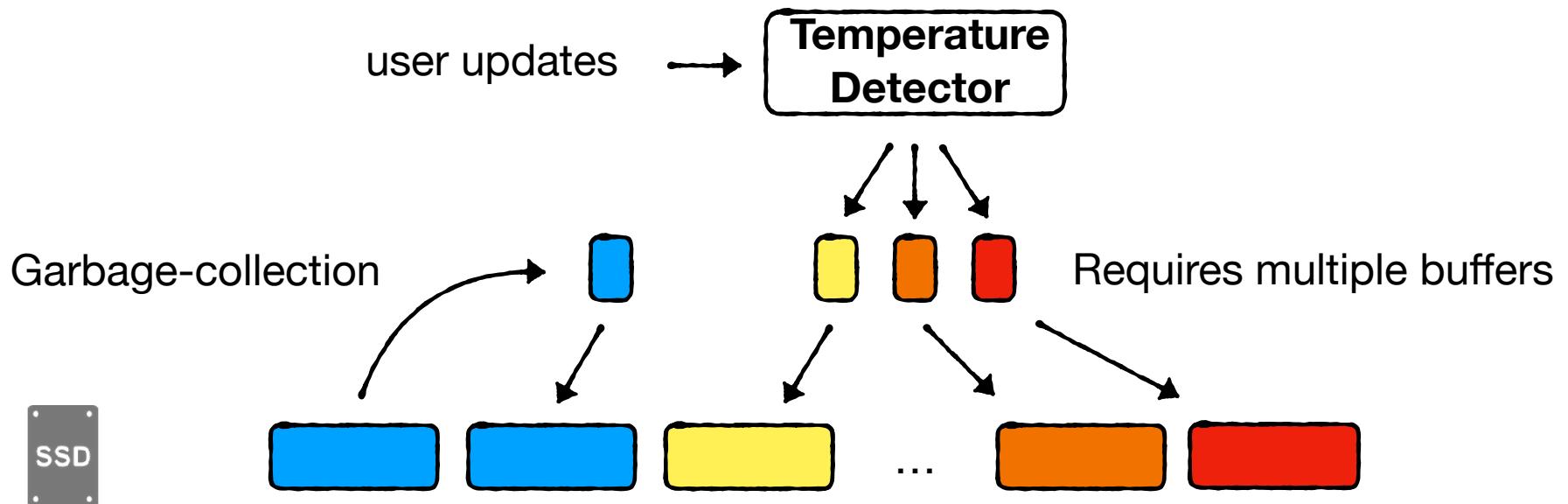
More advanced: separate data with different temperatures into different areas



## Hot vs. Cold Data Separation

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

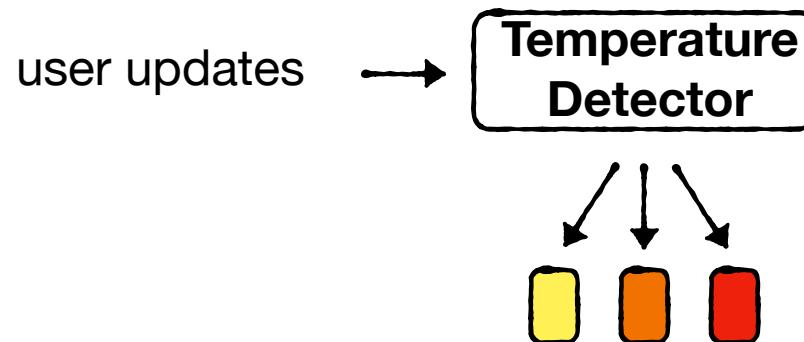
More advanced: separate data with different temperatures into different areas



Estimates how likely a page is to be updated again

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**



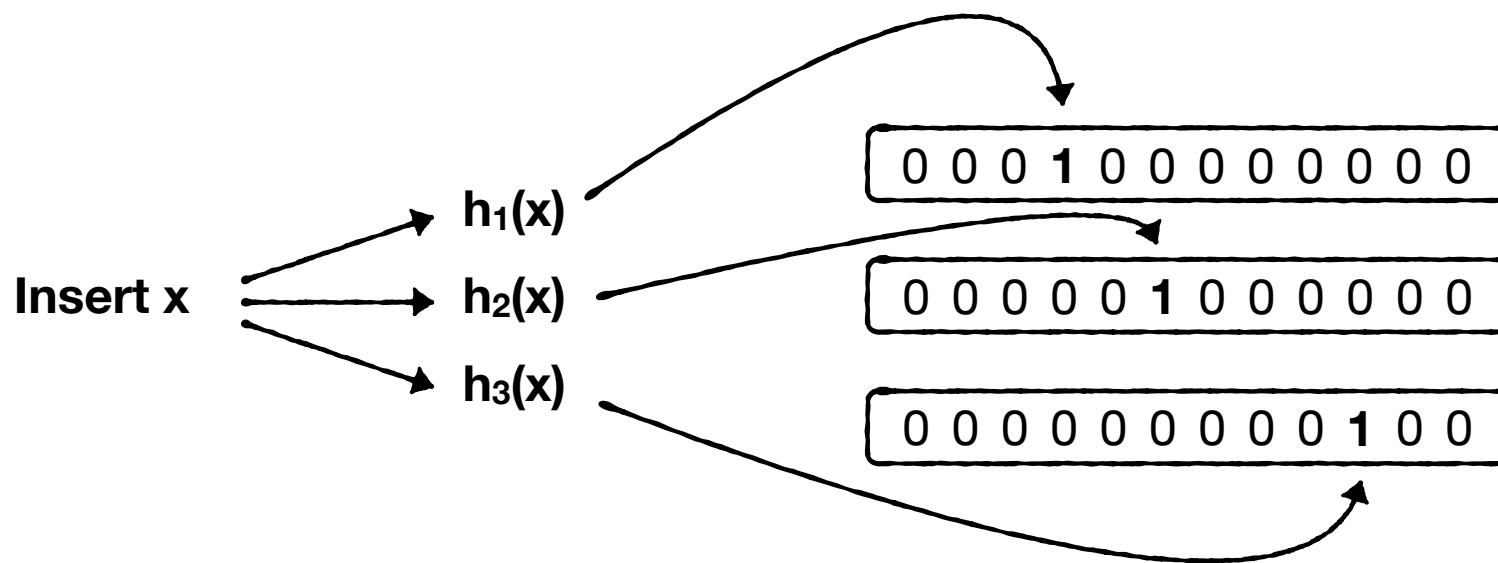
## Count-Min

A data structure that reports the 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**



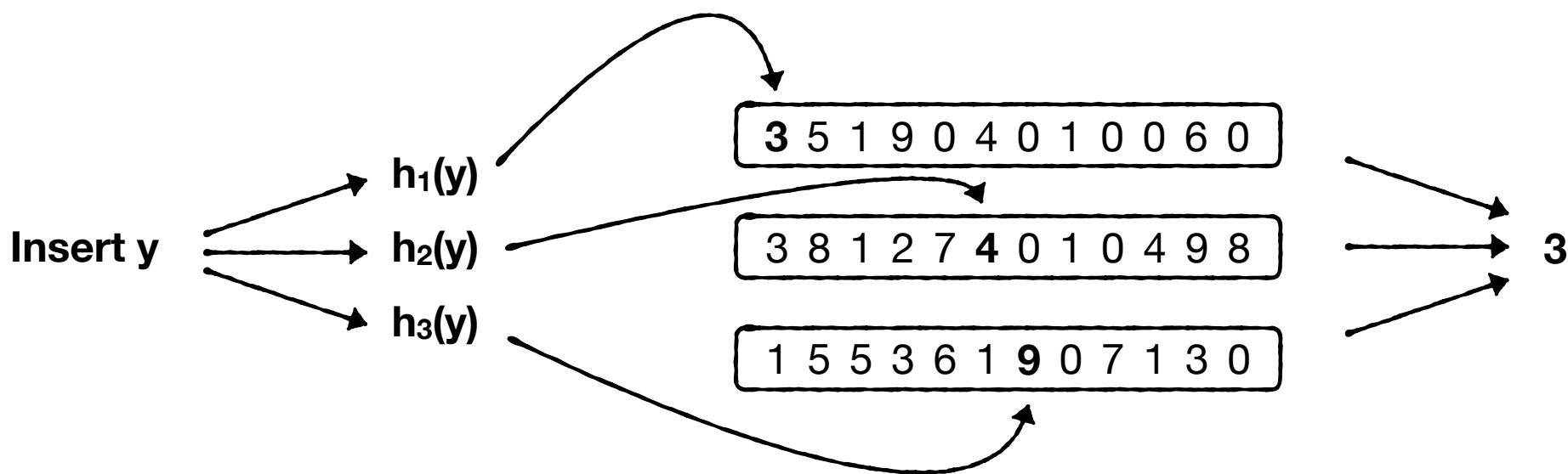
## Count-Min

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**



## Count-Min

Estimated frequency  $\leq$  true frequency +  $\epsilon \cdot$  num insertions

with prob.  $1 - \delta$



3 5 1 9 0 4 0 1 0 0 6 0

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

1 5 5 3 6 1 9 0 7 1 3 0

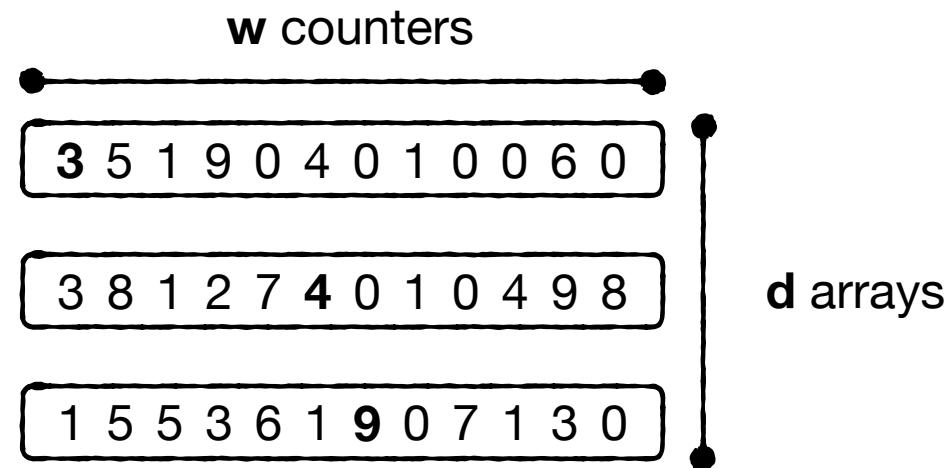
## Count-Min

Estimated frequency  $\leq$  true frequency +  $\epsilon \cdot$  num insertions      with prob.  $1 - \delta$

$$w = \lceil e/\epsilon \rceil$$

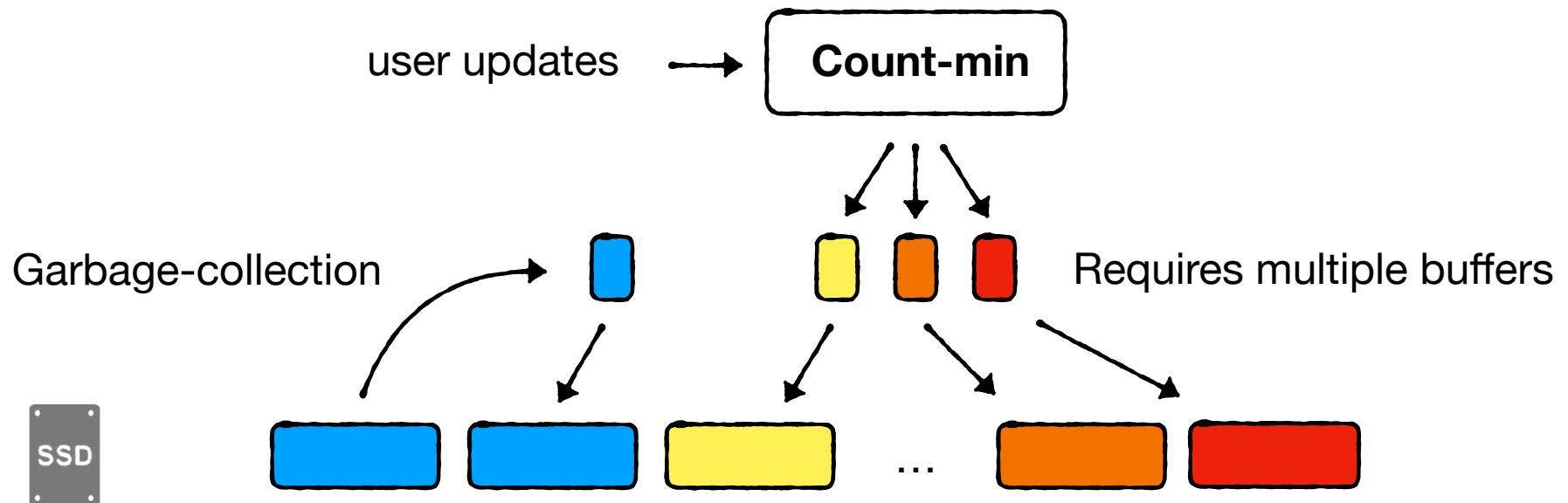
$$d = \lceil \ln 1/\delta \rceil$$

The parameters determine the values of  $w$  and  $d$



We can employ count-min to estimate frequency

**Pitfall:** a value was very hot in the past but became cold.  
Count-min would still tell us its hot.



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

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

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

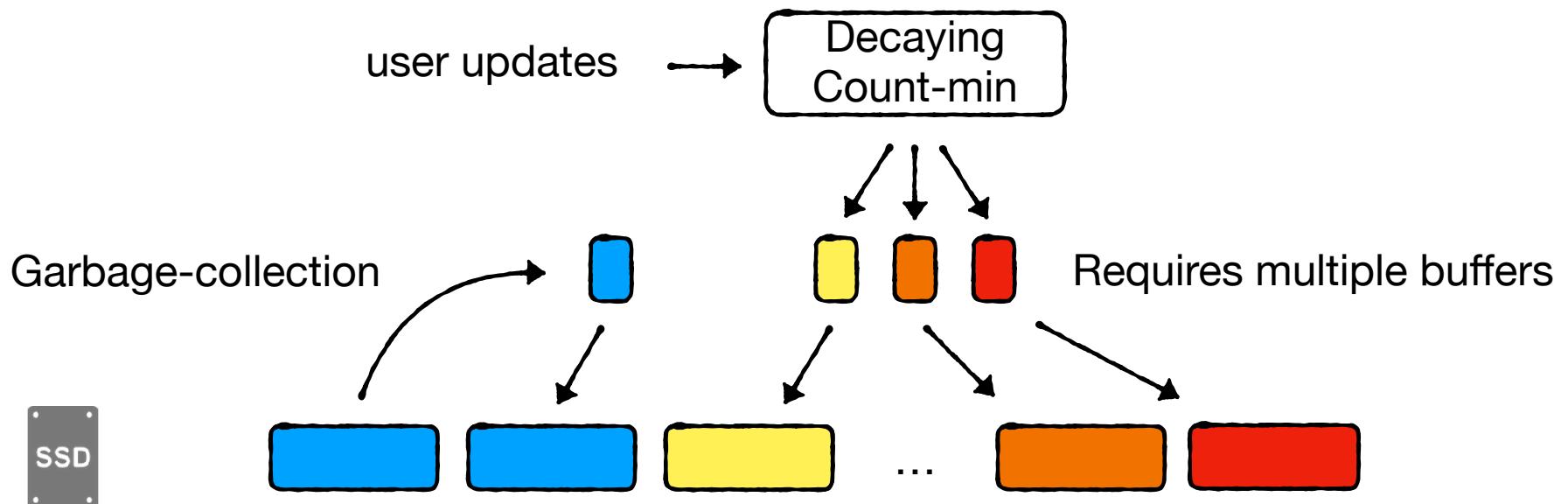
1 2 1 4 0 2 0 0 0 0 3 0

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

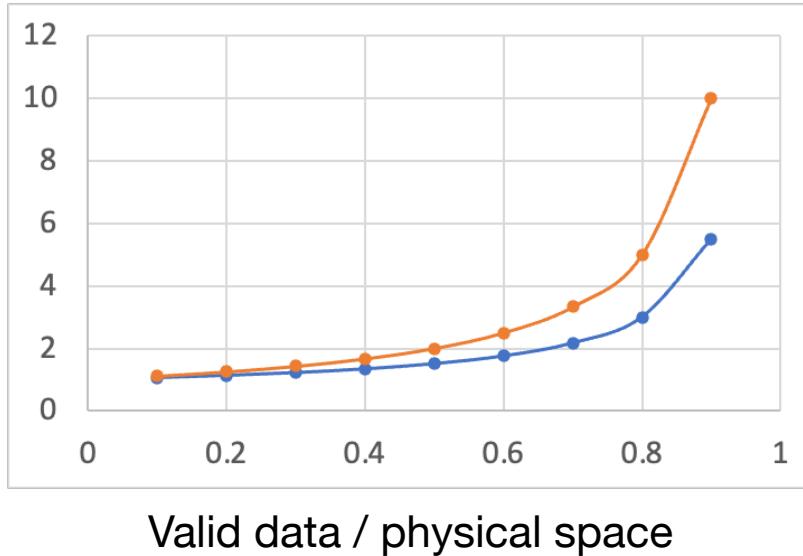
0 2 2 1 3 0 4 0 3 0 1 0

**In reality, hot/cold separation is not widely used in industry.**

**More research & engineering is needed.**



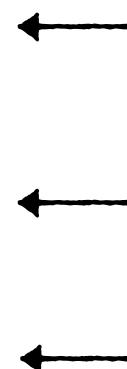
Write-amplification



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

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

...



Mixing hot/cold data brings us towards worst-case

Uniformly random workloads

**Hot/cold separation improves WA beyond uniform case**

**GC overheads =**

$$\left\{ \begin{array}{l} 1 + \frac{L/P}{1 - L/P} \\ 1 + \frac{L/P}{2 \cdot (1 - L/P)} \\ \dots \end{array} \right.$$

Mixing hot/cold data brings us towards worst-case

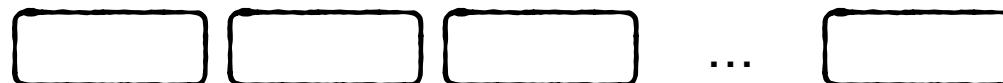
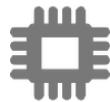
Uniformly random workloads

**Hot/cold separation improves WA beyond uniform case**

# Overall Cost Analysis

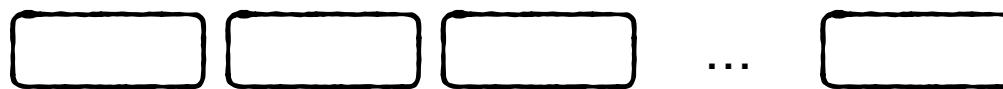
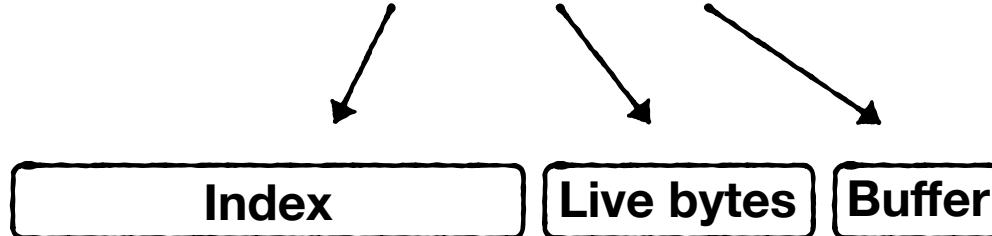
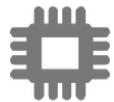
**Write cost: O( GC / B )**

Reads: 1 read I/O

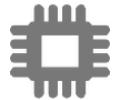


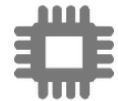
# Checkpointing & Recovery

If power fails, we lose these



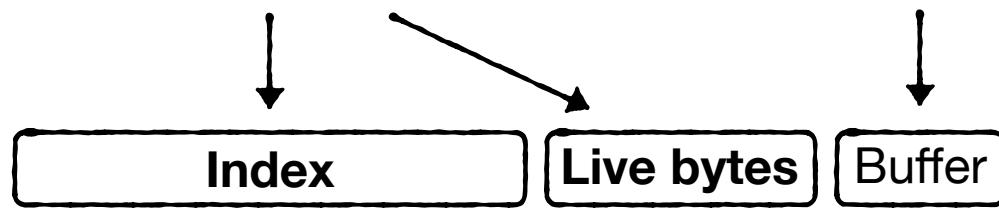
**Let's say it's ok to  
lose the buffer**



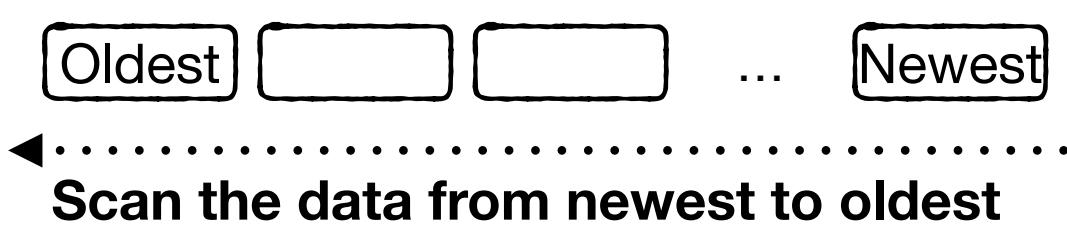
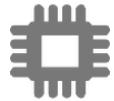


We need to recover  
these

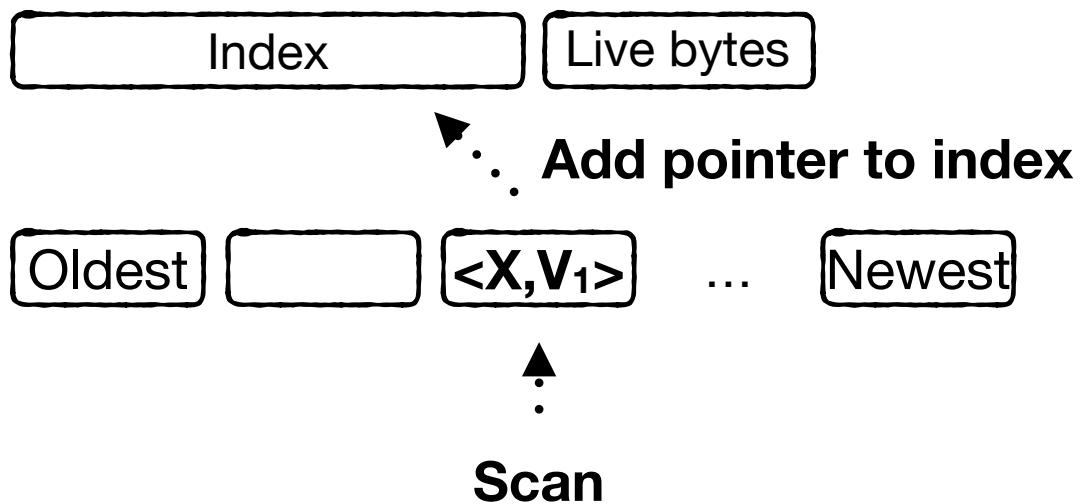
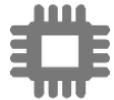
Let's say it's ok to  
lose the buffer



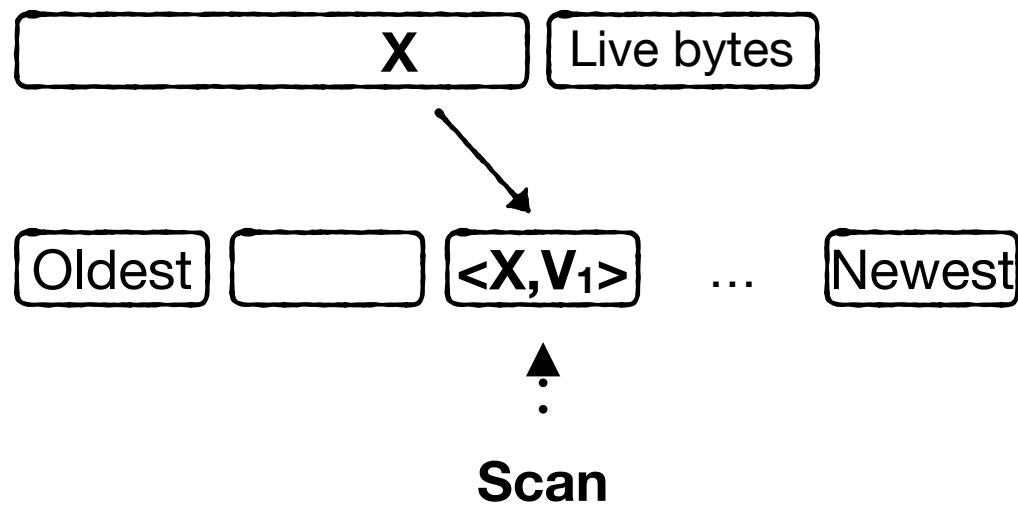
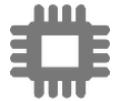
## Simple recovery algorithm?



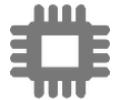
## Simple recovery algorithm?



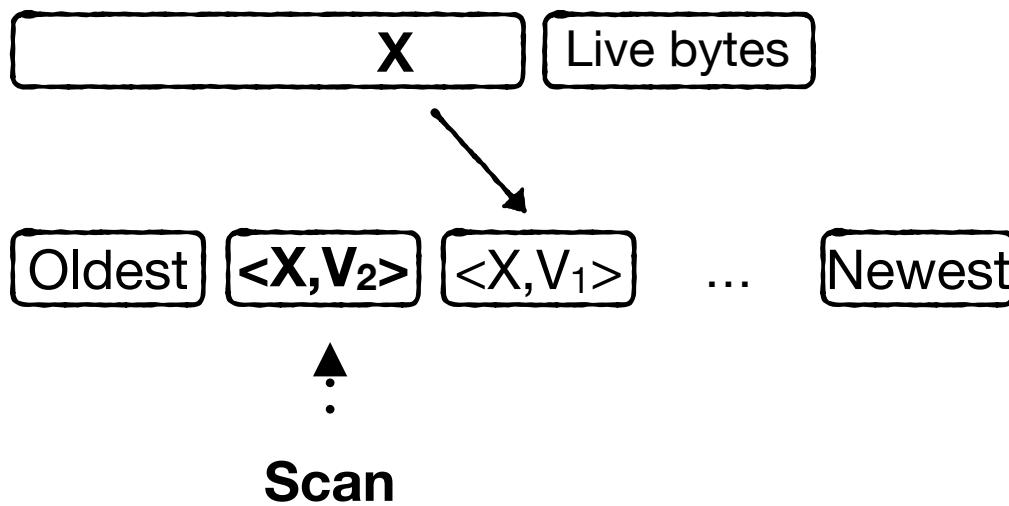
## Simple recovery algorithm?



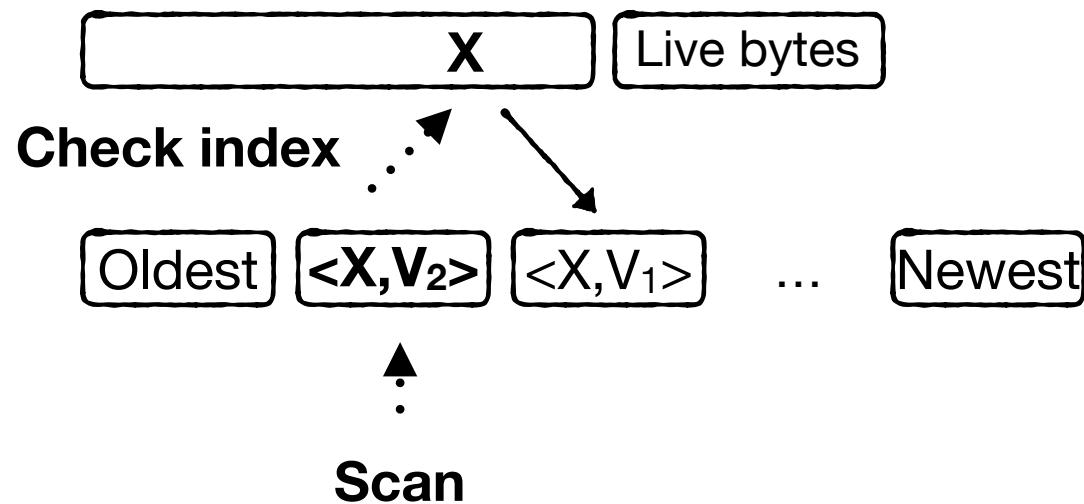
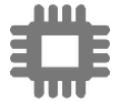
## Simple recovery algorithm?



SSD

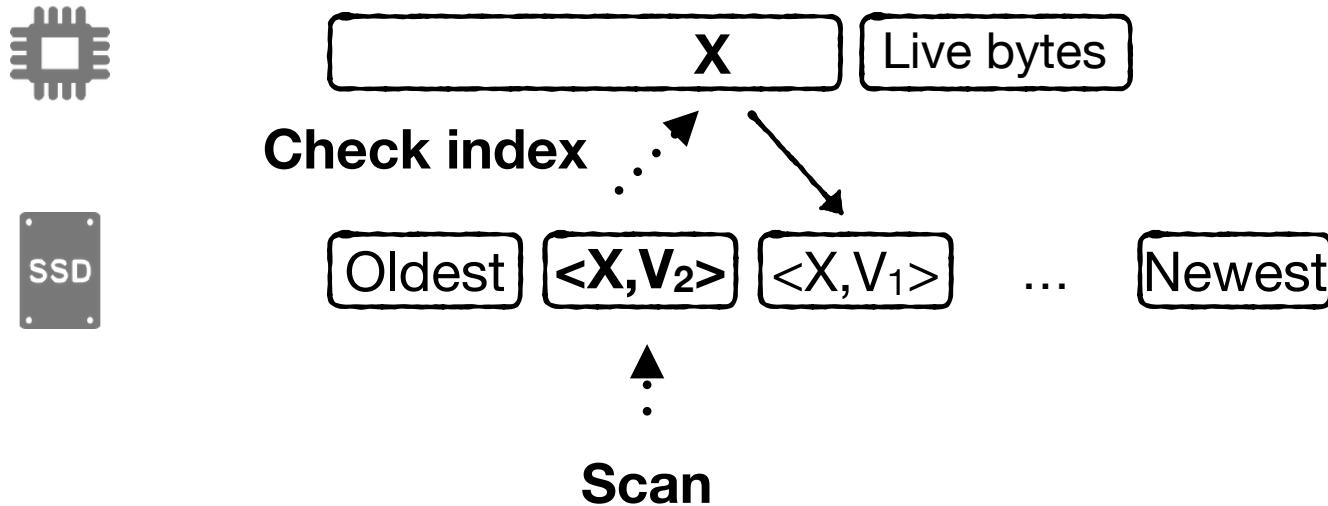


## Simple recovery algorithm?



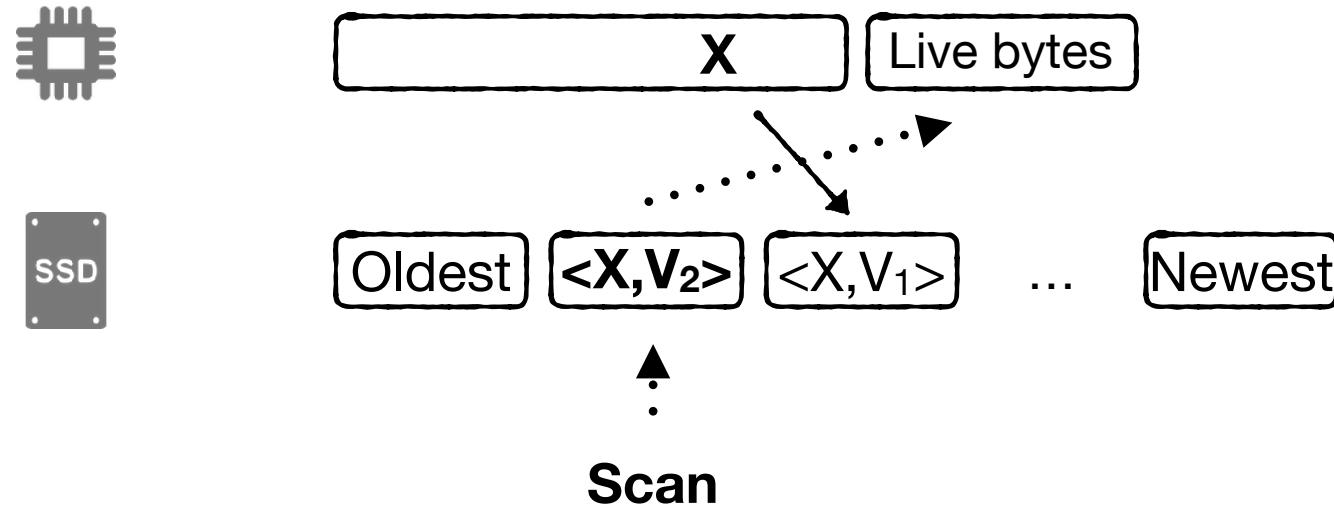
Simple recovery algorithm?

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



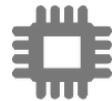
Simple recovery algorithm?

**So we instead subtract size of entry from live bytes.**



We can recover with one pass over the data

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



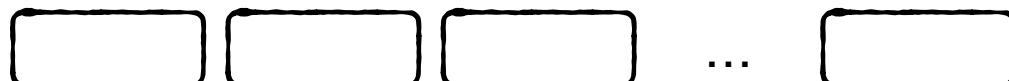
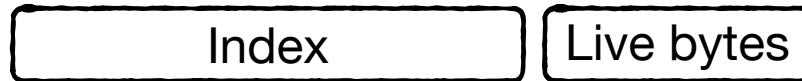
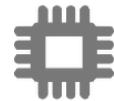
e.g., 200 MB/s sequential  
throughput

1 TB of data takes 1 hour to recover

We can recover with one pass over the data

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

**Any good ideas?**



e.g., 200 MB/s sequential  
throughput

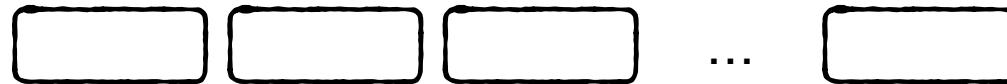
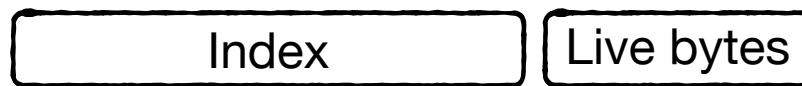
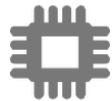


1 TB of data takes 1 hour to recover

# Checkpointing

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

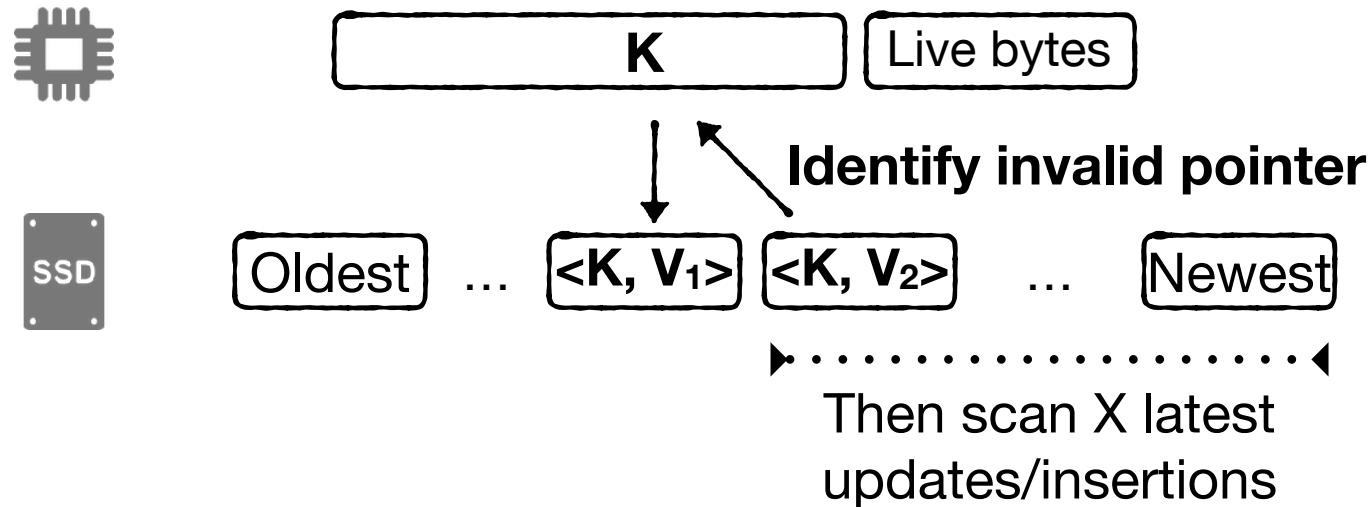
**After recovery, load copy**



# Checkpointing

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

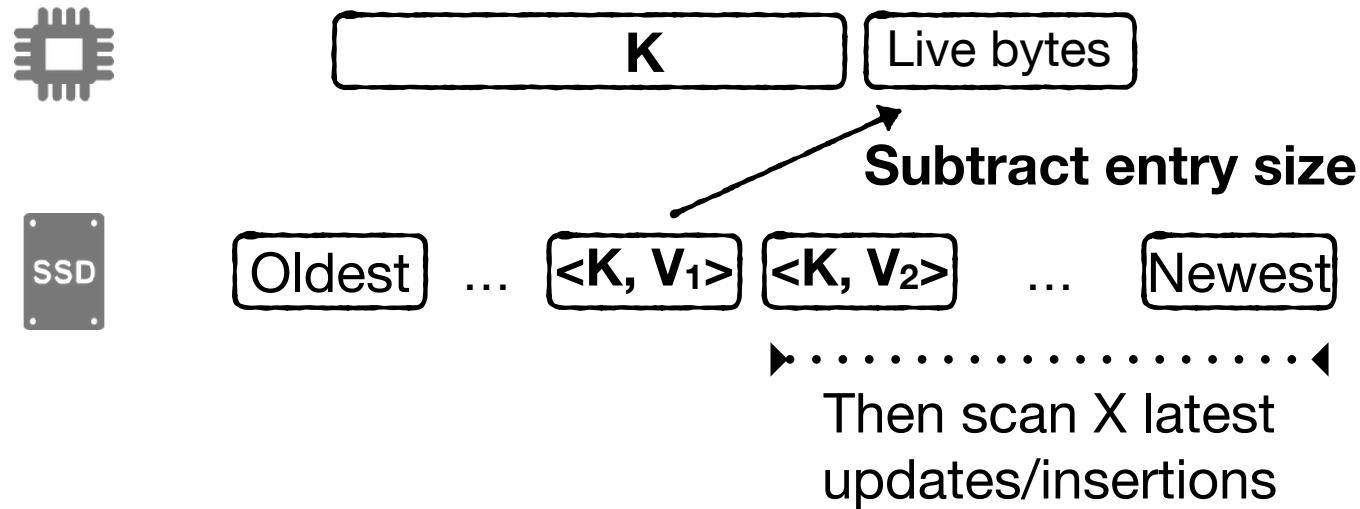
**After recovery, load copy**



# Checkpointing

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

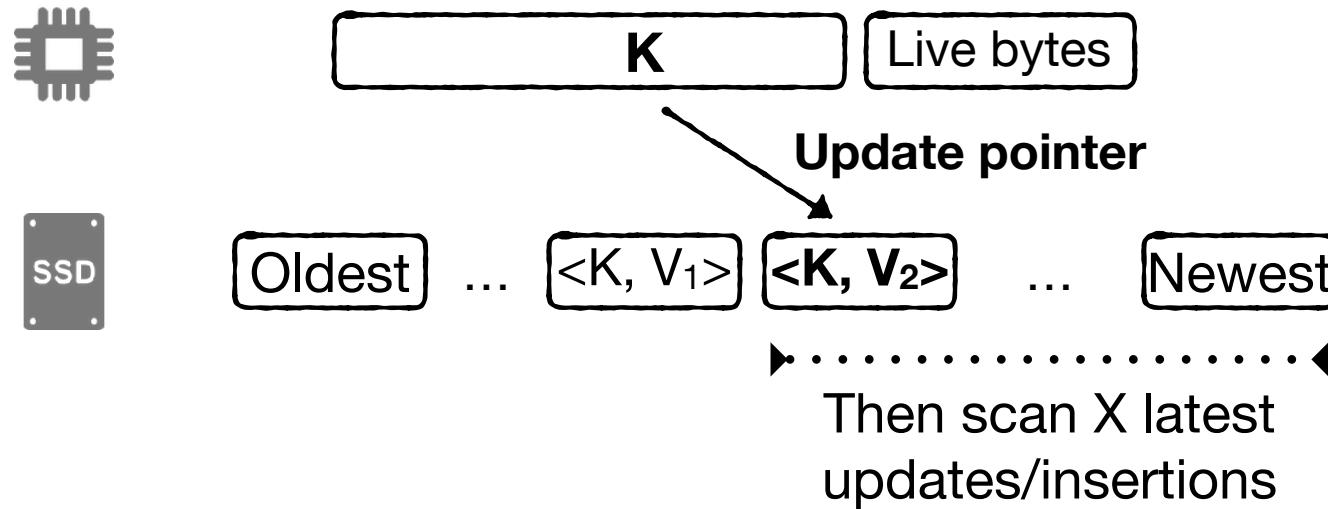
**After recovery, load copy**



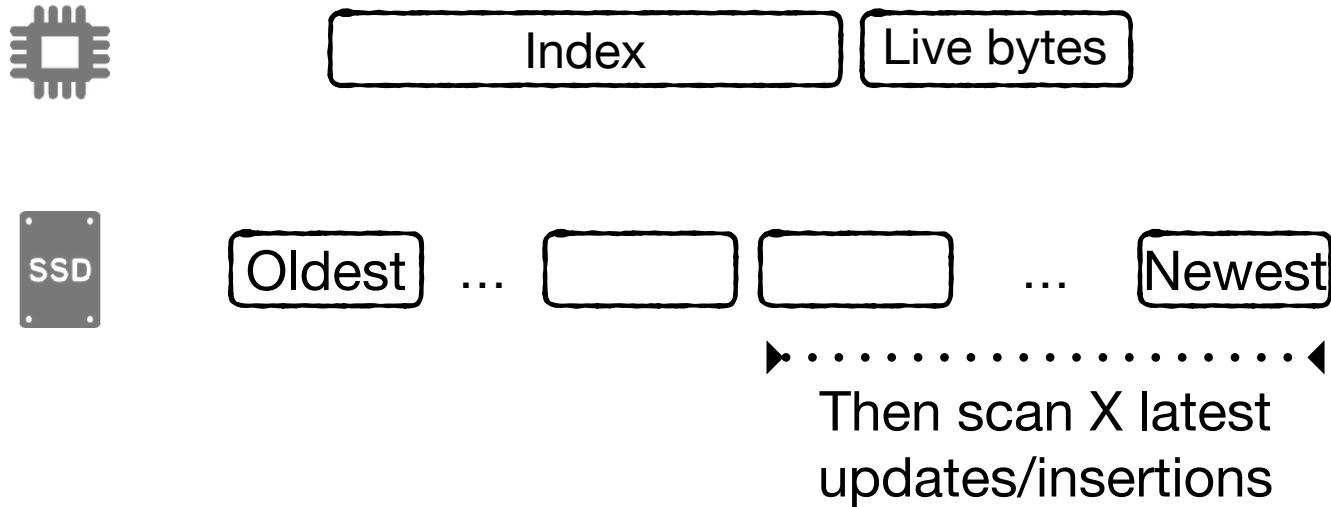
# Checkpointing

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

**After recovery, load copy**

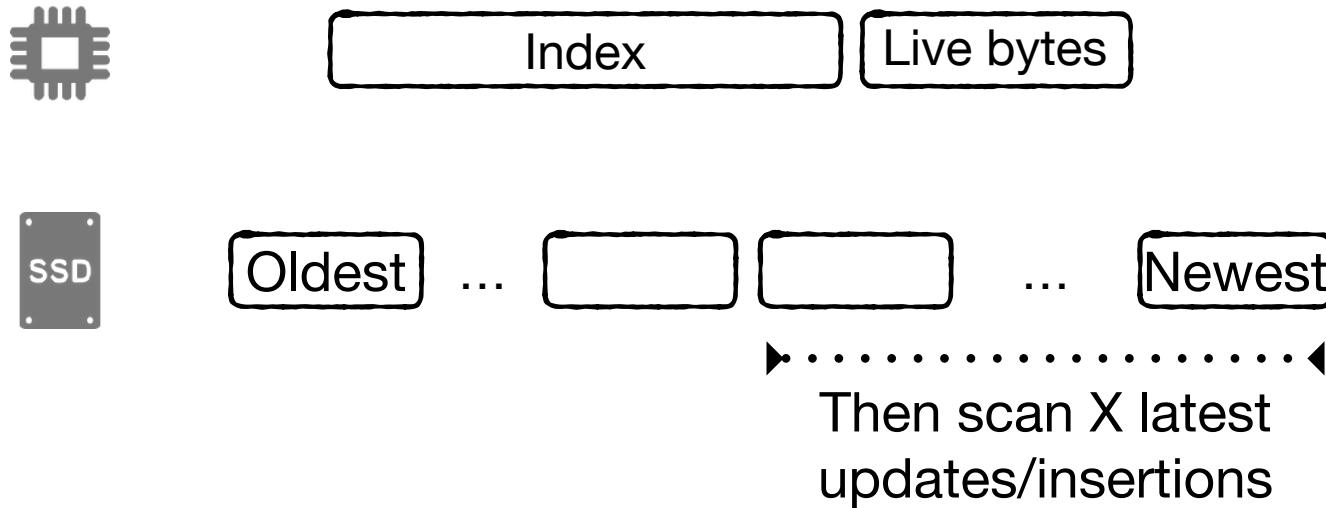


**The frequency of checkpointing controls a trade-off between write cost and recovery time.**



**Additional write cost:**  $O(\text{index size} / X)$

**Recovery time:**  $O(X)$  reads for backwards scan

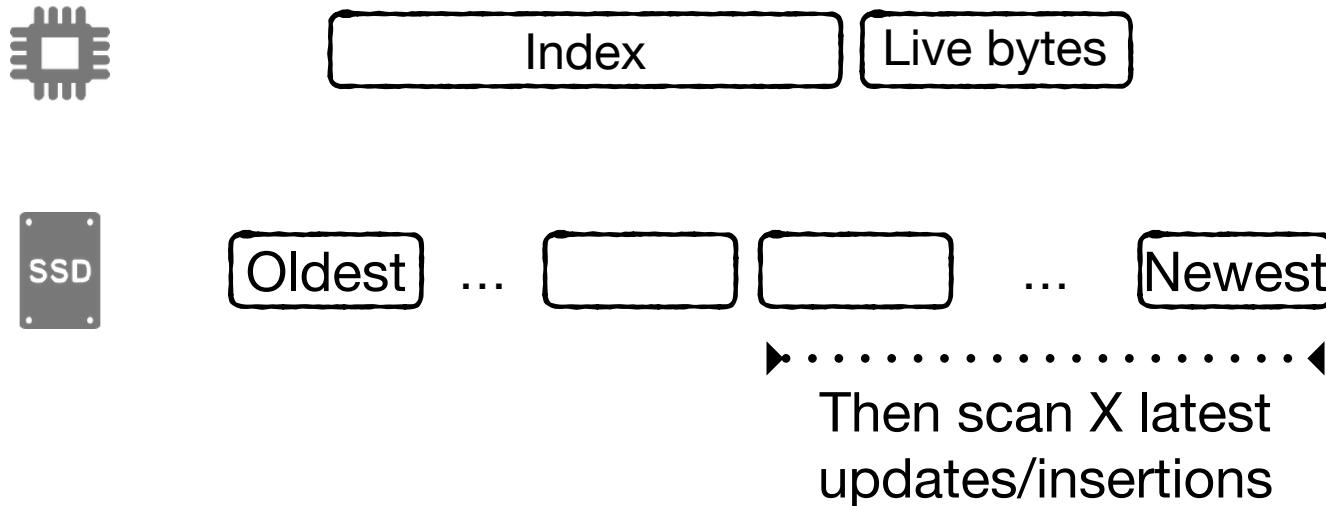


# Costs Summary

**write cost:**  $O(GC/B + \text{index size} / X)$

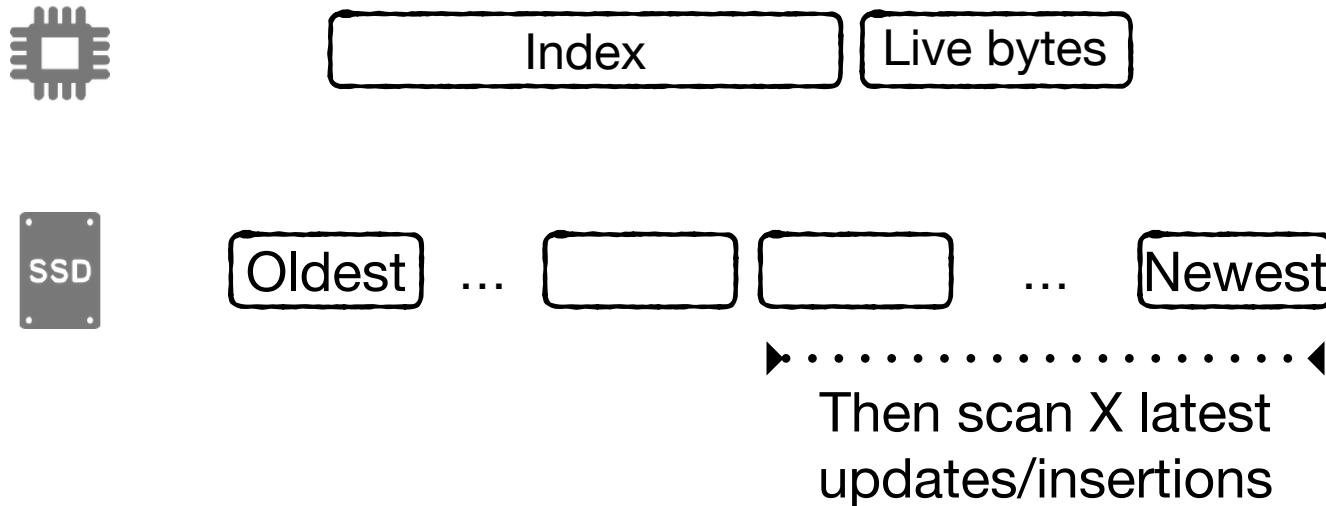
**Recovery time:**  $O(X)$  reads for backwards scan

**Get cost:**  $O(1)$

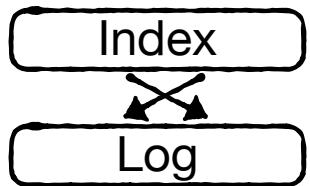


## Note on Deletes & Tombstones for Recovery

After recovery, any key-value pair that was physically present in the log but deleted in the index should remain deleted in the index. We can achieve this by adding a tombstone for each entry. During recovery, if the first instance of a key we see is a tombstone, we ignore all subsequent entries with this key.



## Mechanics



## Hot/Cold separation



To reduce write-amp

## Checkpointing & Recovery

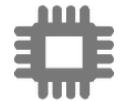


If power fails

## Cuckoo filtering



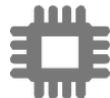
To reduce memory



Can be large

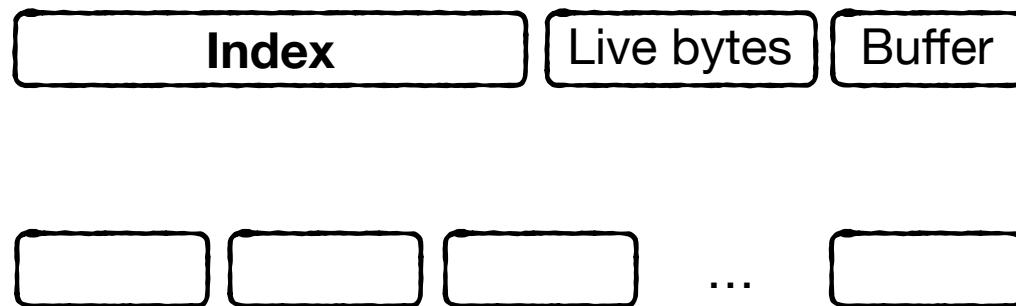


**Ultimately attack this  
with cuckoo filters**

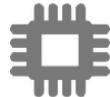


$$\text{Index size} = N \cdot (P + K) / \alpha$$

→       $N = \text{data size}$   
 $P = \text{pointer size} = O(\log_2 N/B)$   
**K = key size =  $\Omega(\log_2 N)$**   
 $\alpha = \text{collision resolution overheads} \approx 0.8$



**But first attack this  
with cuckoo hashing**

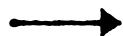


$$\text{Index size} = N \cdot (P + K) / \alpha$$

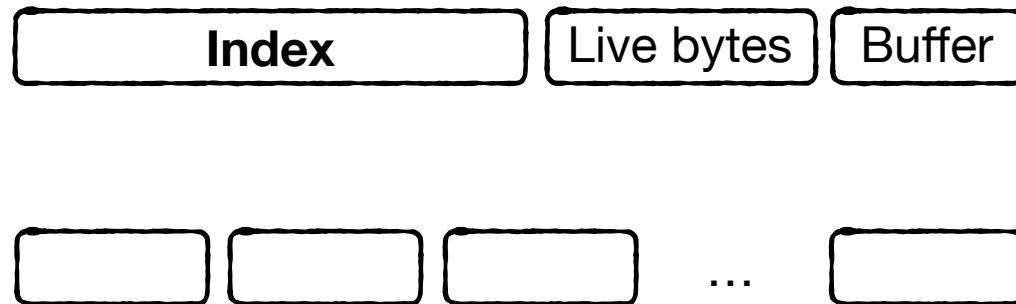
$N$  = data size

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

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



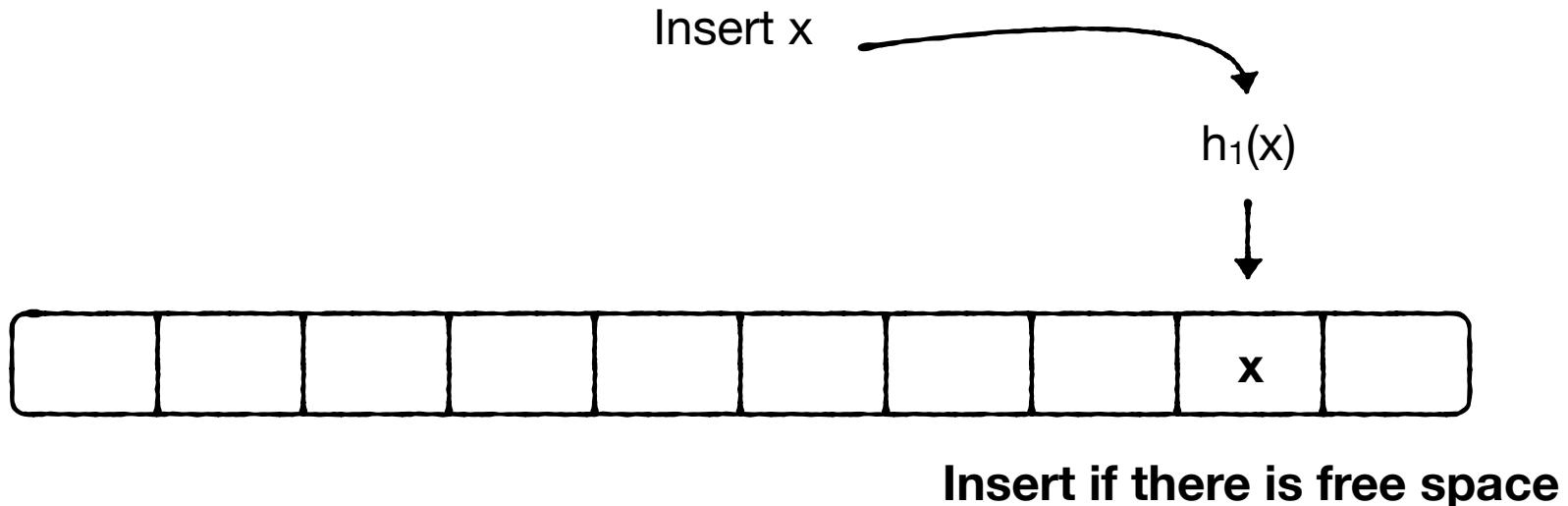
$\alpha$  = collision resolution overheads  $\approx 0.8$



## Cuckoo Hashing

Table of buckets, each containing 1 entry

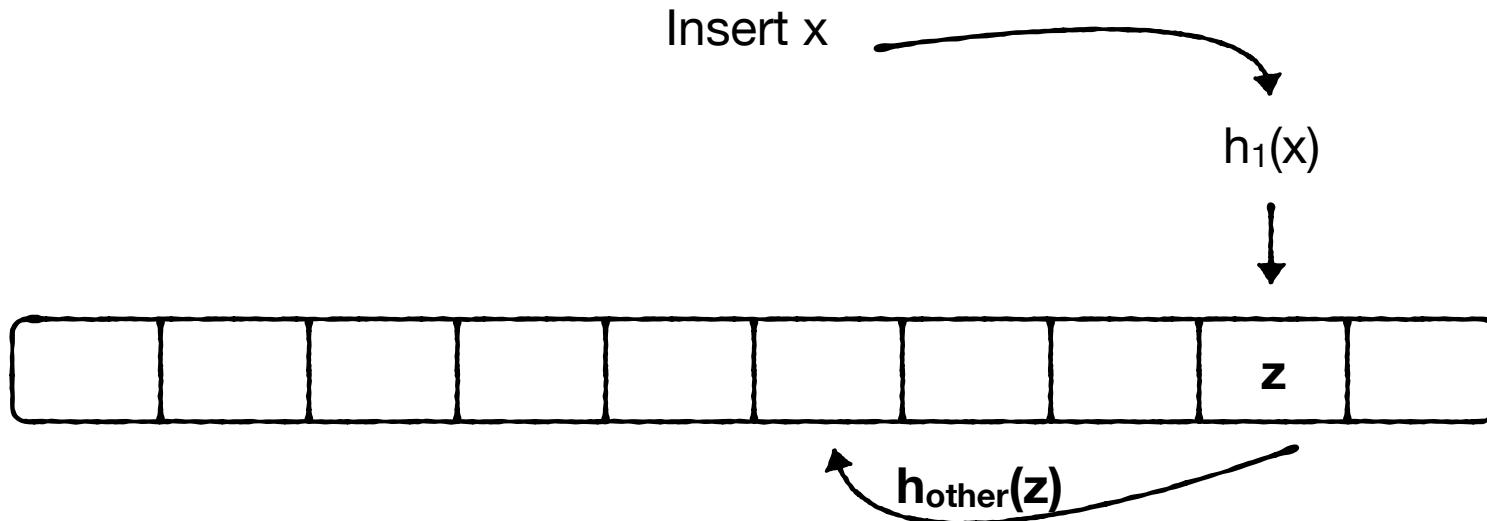
Two hash functions



## Cuckoo Hashing

Table of buckets, each containing 1 entry

Two hash functions

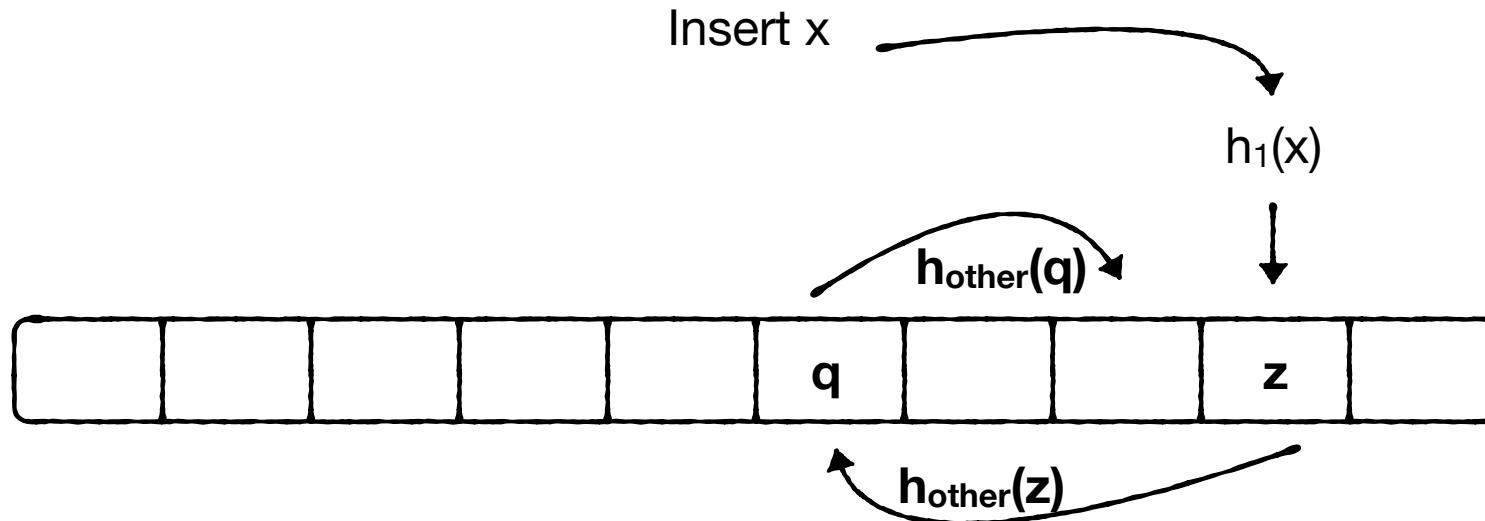


**Otherwise, evict existing entry using its alternative hash function to an alternative slot**

## Cuckoo Hashing

Table of buckets, each containing 1 entry

Two hash functions

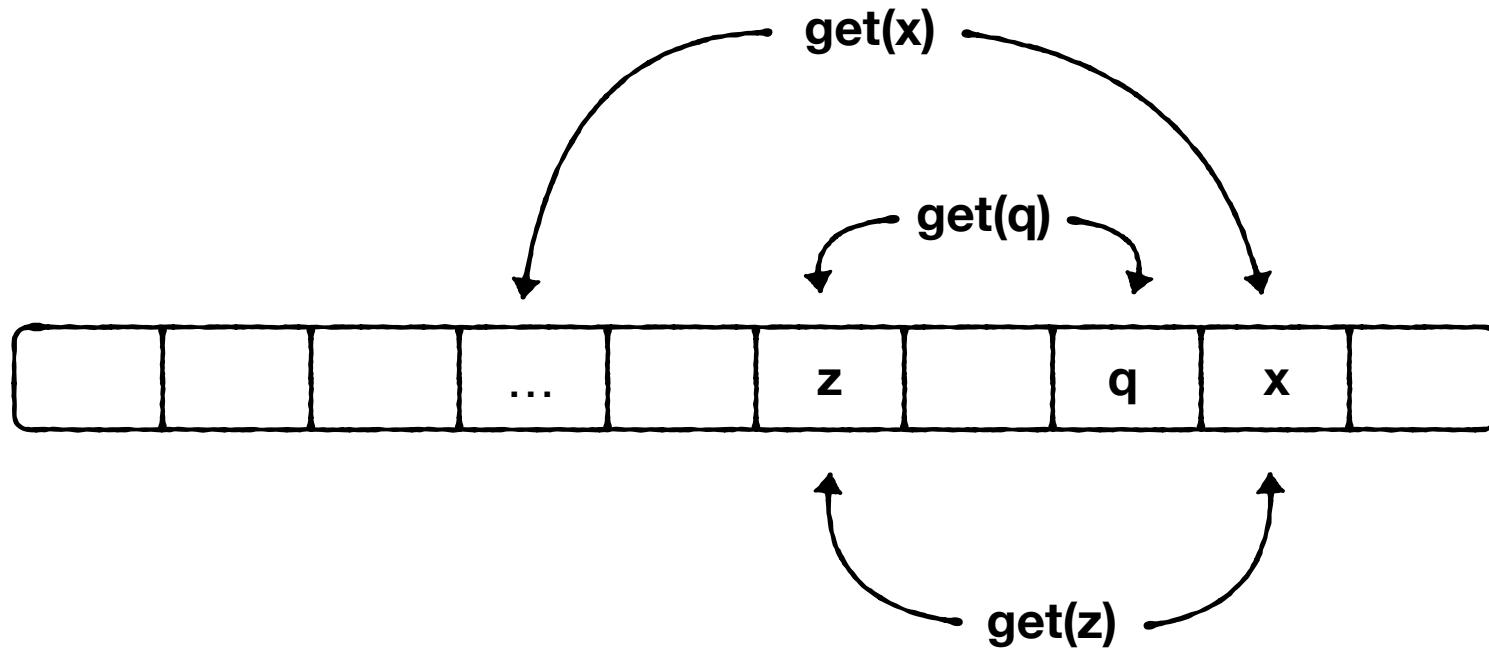


Continue doing this recursively until all entries are mapped to an empty bucket

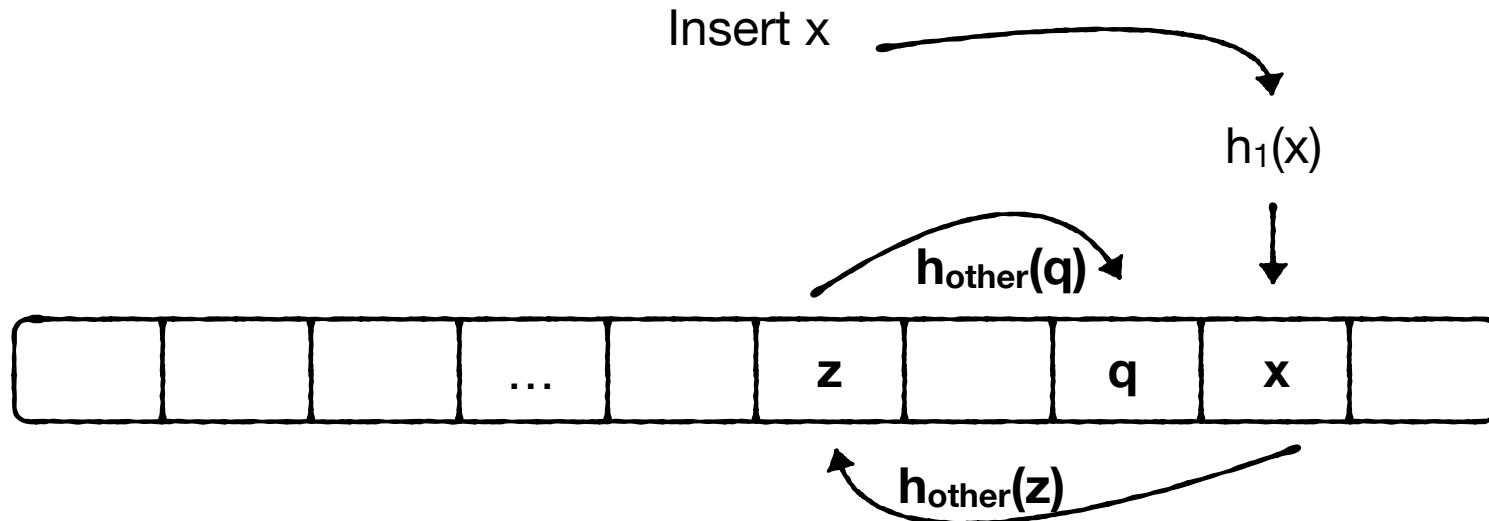


**Continue doing this recursively until all entries are mapped to an empty bucket**

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



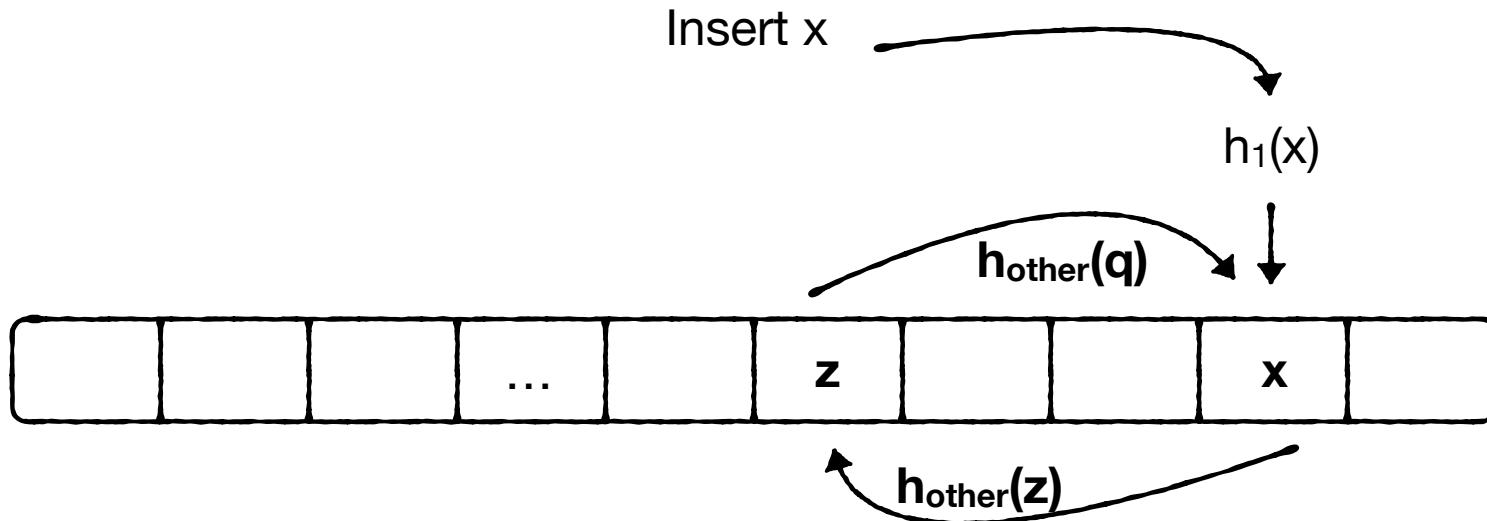
**Insertions may endure multiple evictions and swapping of keys across buckets**



Insertions may endure multiple evictions and swapping of keys across buckets

Worse: an infinite loop is technically possible

**We can alleviate this by storing multiple keys per bucket**



Insertions may endure multiple evictions and swapping of keys across buckets

Worse: an infinite loop is technically possible

**We can alleviate this by storing multiple keys per bucket**



**Theory: insertions succeed in  $O(1)$  expected time with high probability**

Assuming **load factor < 50% for bucket size 1**

**load factor < 84% for bucket size 2**

**load factor < 95% for bucket size 4**

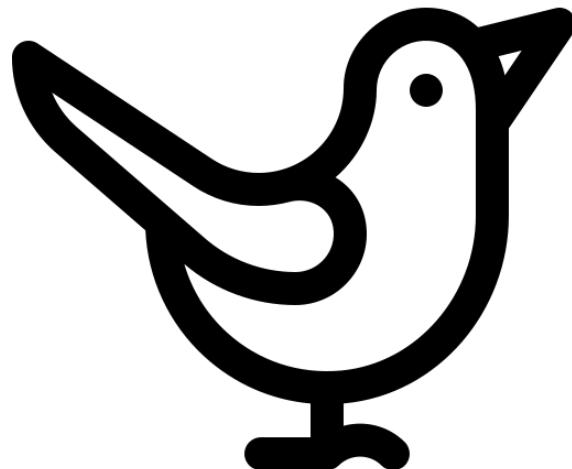


**Theory: insertions succeed in  $O(1)$  expected time with high probability**

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.

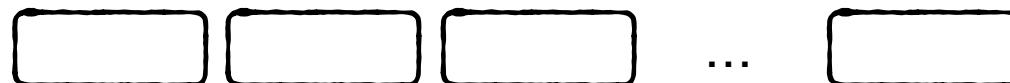
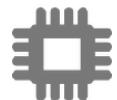
$$\text{Index size} = N * (P + K) / a$$

N = data size

P = pointer size

K = key size

**Addressed with  
cuckoo hashing** → **a = collision resolution overheads ≈0.8 → ≈0.95**



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

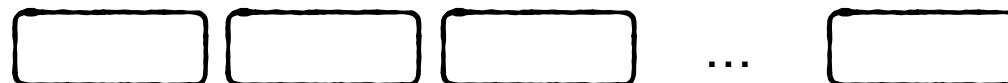
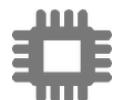
$N$  = data size

$P$  = pointer size

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

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

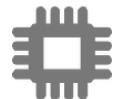
Now lets attack this  
With a Cuckoo Filter



## 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 )



## Cuckoo Filter

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

**A fingerprint is a hash digest derived by hashing a key**

$$\text{FP}(\text{KEY}) = \text{FINGERPRINT}$$

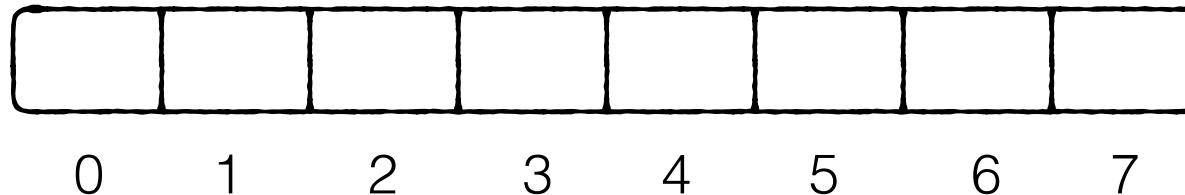


## Cuckoo Filter

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

A fingerprint is a hash digest derived by hashing a key

Example:       $\text{FP}(X) = \overbrace{0100}^{\text{M bits}}$        $h_1(X)$  = Bucket address

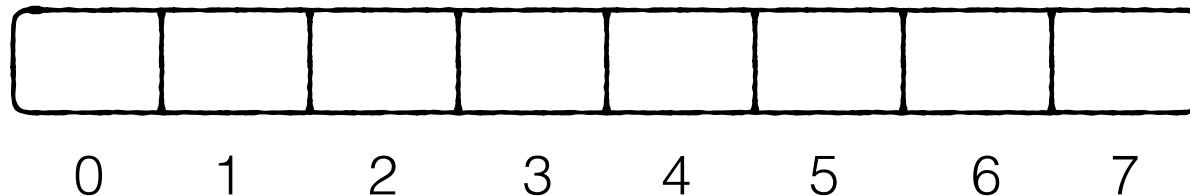


## Cuckoo Filter

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

A fingerprint is a hash digest derived by hashing a key

Example:       $\text{FP}(X) = \overbrace{0100}^{\text{M bits}}$        $h_1(X) = 3$

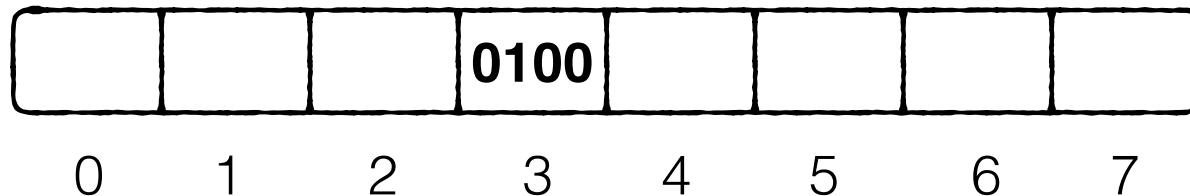


## Cuckoo Filter

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

A fingerprint is a hash digest derived by hashing a key

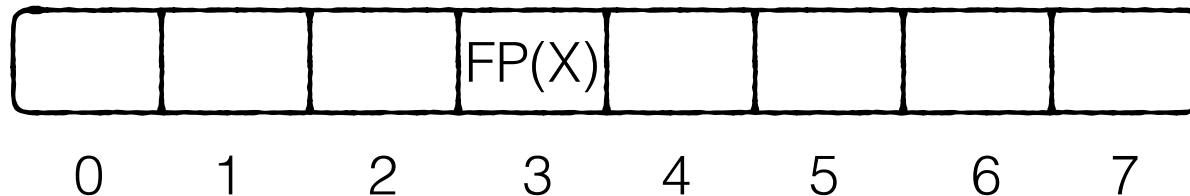
Example:       $\text{FP}(X) = \overbrace{\quad\quad\quad}^{\text{M bits}}$        $h_1(X) = \mathbf{3}$



## Cuckoo Filter

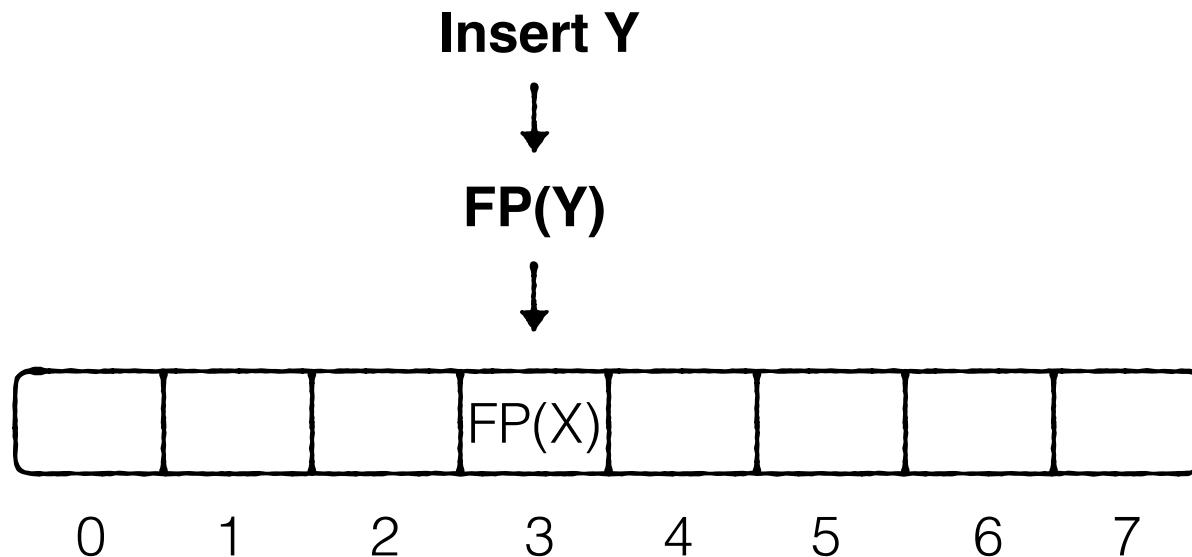
Same as Cuckoo hash tables, but store fingerprints instead of keys

A fingerprint is a hash digest derived by hashing a key



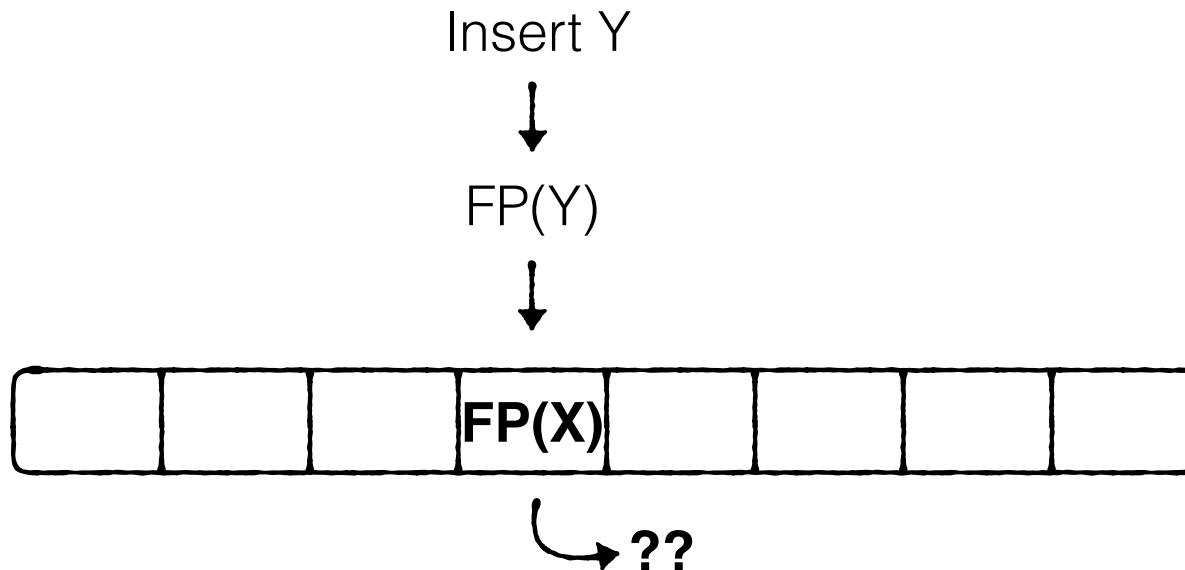
## Cuckoo Filter

Suppose we then insert another fingerprint to the same bucket



## Cuckoo Filter

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?**

## How to derive an alternative bucket?

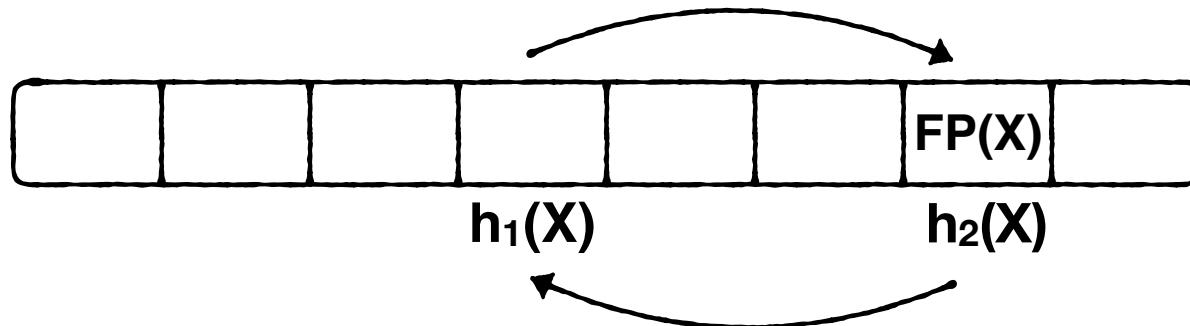
We have two pieces of information about X

(1) Bucket address:  $h_1(X)$

(2) fingerprint  $FP(X)$

We want to combine them to give alternative bucket address  $h_2(X)$

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



**Any ideas?**

# XOR Operator

<b>Input 1</b>	0	0	1	1
$\oplus$				
<b>Input 2</b>	0	1	0	1
=				
<b>Parity</b>	0	1	1	0

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

## Parity can help recover any input

Suppose we  
lost input 2

Input 1	0	0	1	1
$\oplus$				
Parity	0	1	1	0
=				
Input 2	0	1	0	1

**Recovered**

## Parity can help recover any input

Or suppose we  
lost input 1

Input 2	0	1	0	1
		⊕		
Parity	0	1	1	0
		=		
Input 1	0	0	1	1

**Recovered**

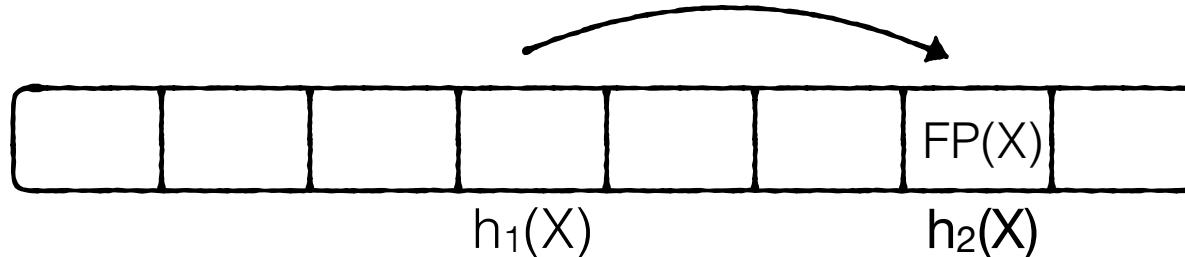
## Using XOR in Cuckoo filters

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**

$$h_2(X) = h_1(X) \text{ xor } \text{hash}(FP(X))$$



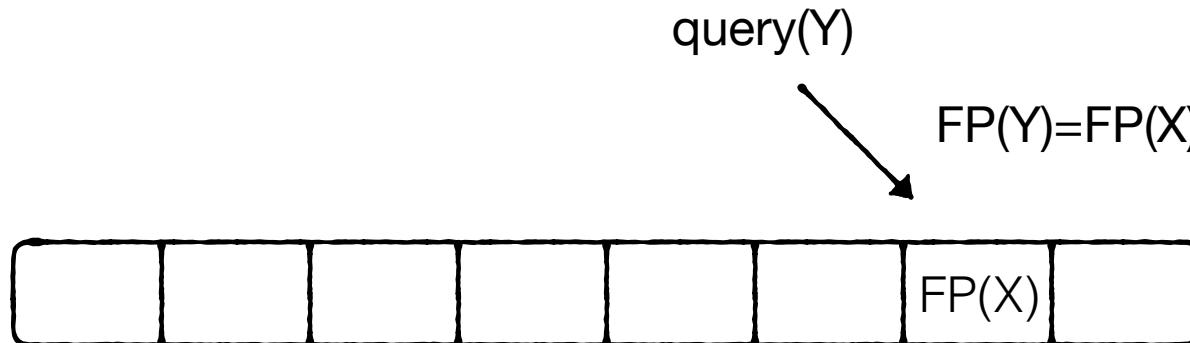
$$h_1(X) = h_2(X) \text{ xor } \text{hash}(FP(X))$$

## Using XOR in Cuckoo filters

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.



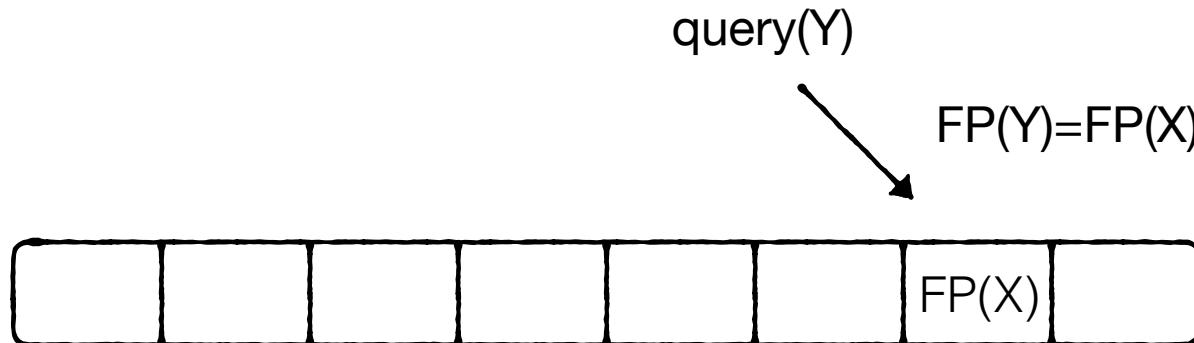
**Can we have false negatives?**

## Using XOR in Cuckoo filters

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?

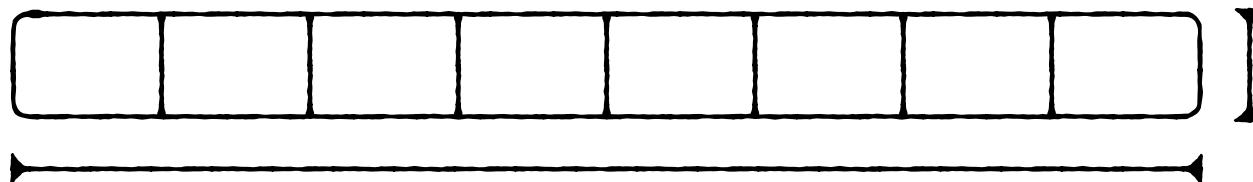
**No, because a fingerprint for a key that had been inserted is in one of only two possible buckets, both of which we search.**

**Since a query searches  
two buckets**



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

M bits per  
fingerprint



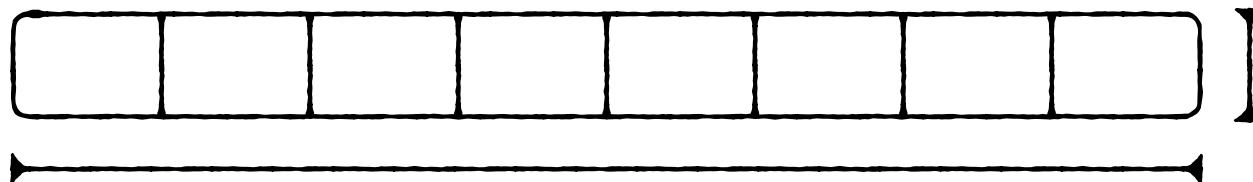
β slots per  
bucket

**When at capacity**



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

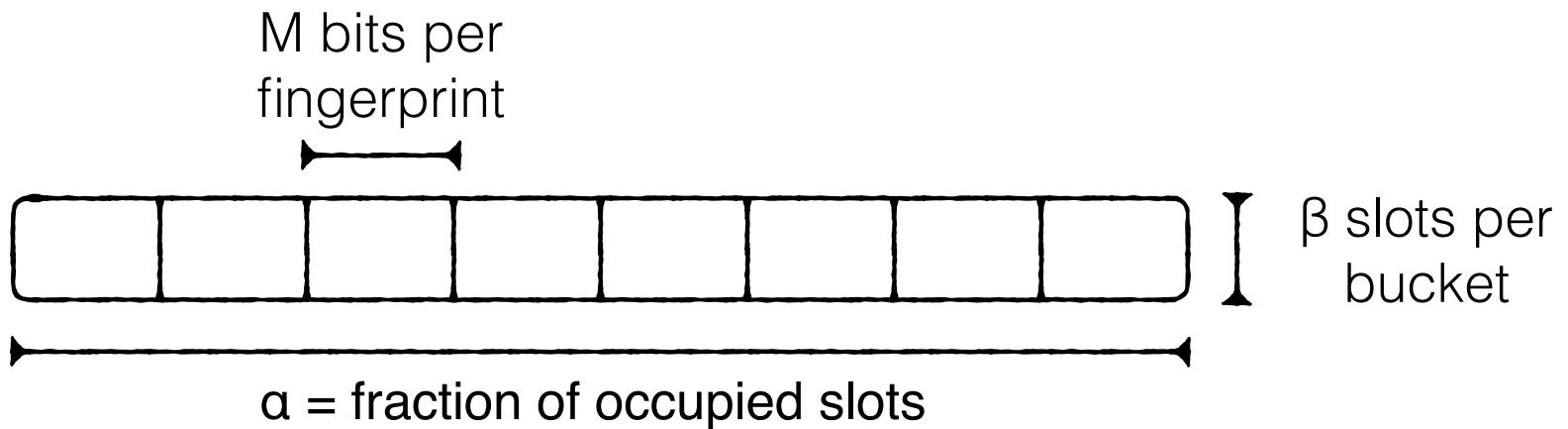
M bits per  
fingerprint



a = fraction of occupied slots

β slots per  
bucket

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



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

**For Bloom filter  
when  $M \geq 10$**



## Cuckoo Filter      Bloom filter

false positive rate:       $\approx 2^{-M+3}$        $2^{-M} \cdot \ln(2)$

**Memory accesses for  
positive query**      1.5       $M \cdot \ln(2)$



## Cuckoo Filter      Bloom filter

false positive rate:       $\approx 2^{-M+3}$        $2^{-M} \cdot \ln(2)$

Memory accesses for  
positive query      1.5       $M \cdot \ln(2)$

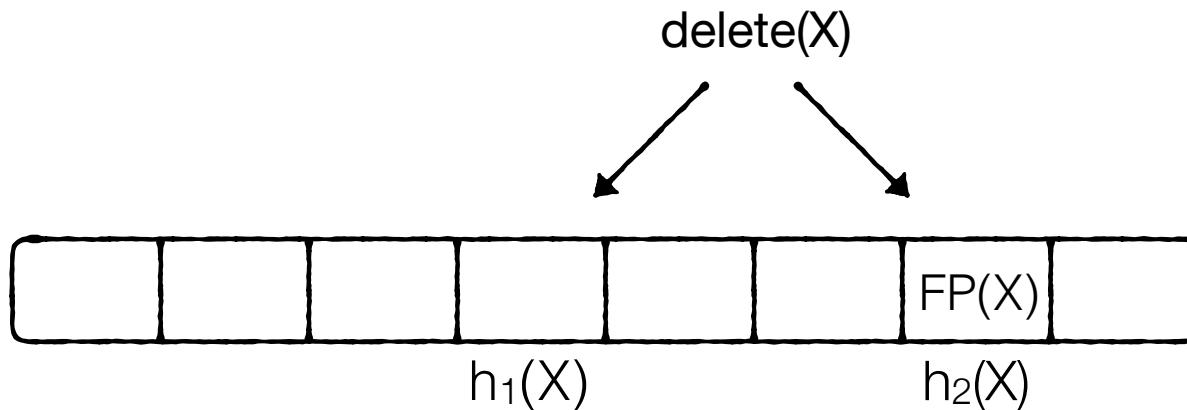
Memory accesses for  
negative query      2       $\approx 2$

Insertion cost      Exp. O(1)       $M \cdot \ln(2)$

**Deletes?**      **N/A**

**Storing payloads?**      **N/A**

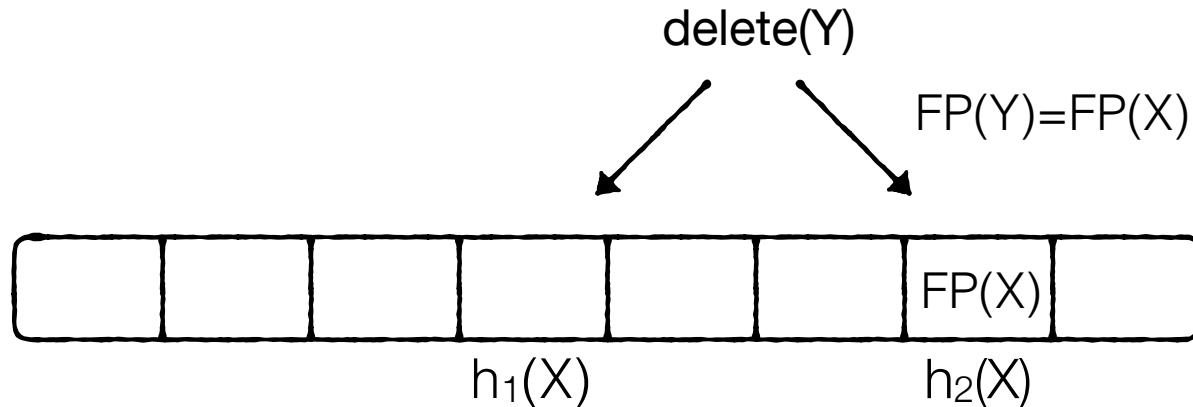
## Deletes in a Cuckoo Filter



## Deletes in a Cuckoo Filter

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?**



## Deletes in a Cuckoo Filter

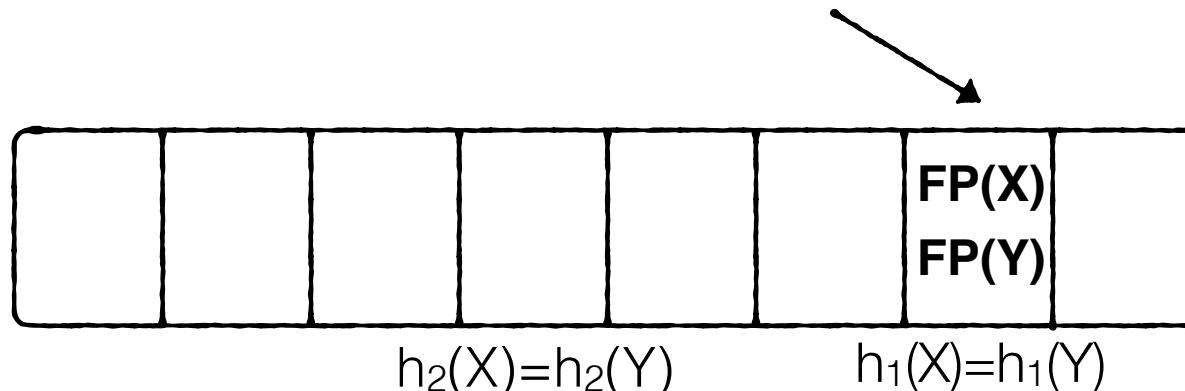
A delete searches both buckets for a key and removes a matching fingerprints.

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 positives 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.**



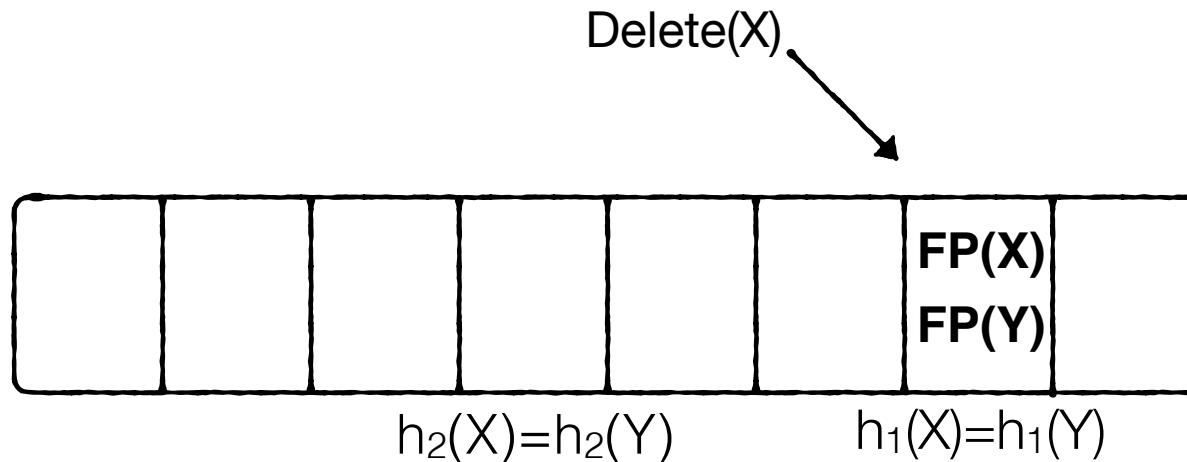
## Deletes in a Cuckoo Filter

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 positives can occur

**Example: insert keys X and Y that happen to have matching buckets & fingerprints**



## Cuckoo Filter      Bloom filter

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      2      2

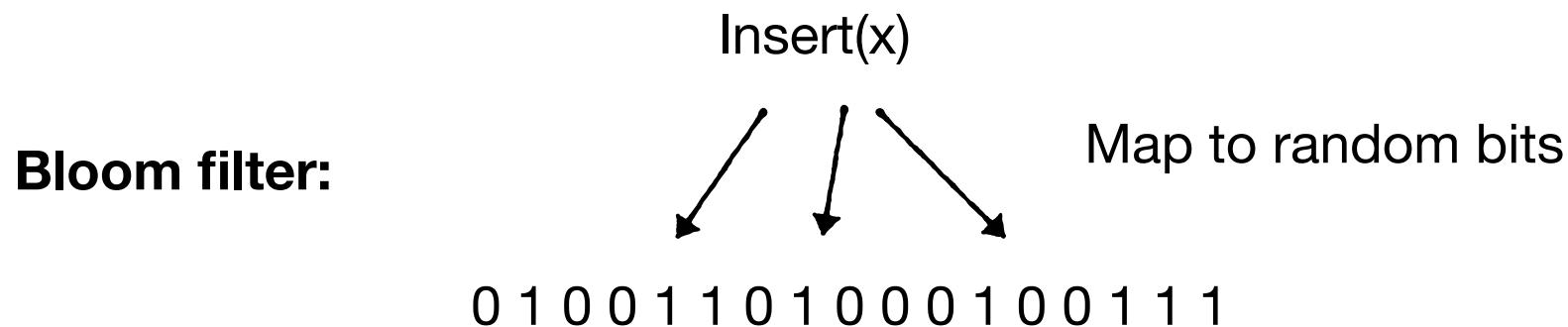
Insertion cost       $O(1)$        $M \cdot \ln(2)$

**Deletes?**      **1.5**      **N/A**

Storing Payloads?      N/A

## Storing Payloads in a Filter

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

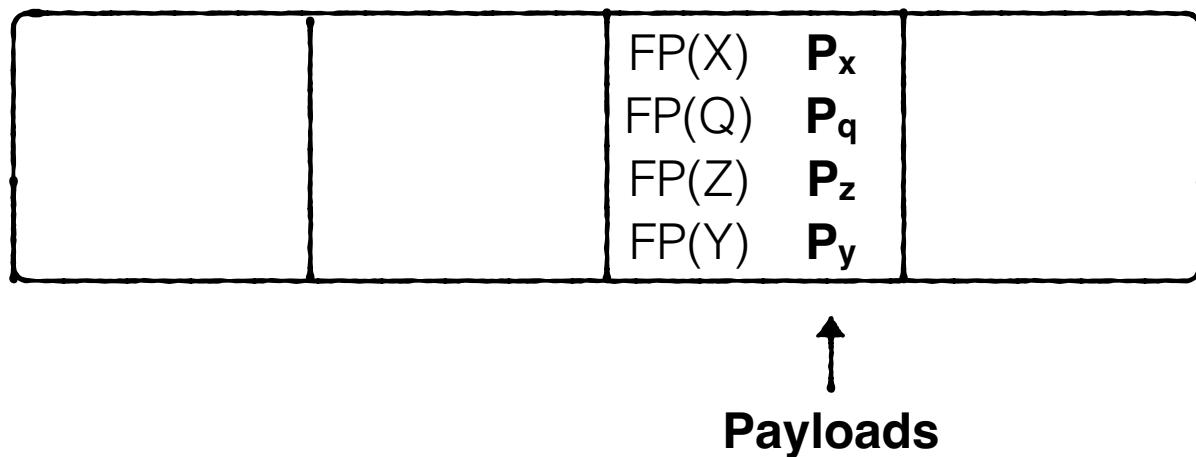


**There is nowhere obvious to associate a payload with each key**

## Storing Payloads in a Filter

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

**Cuckoo filter:**



**Just store a payload alongside each fingerprint**

## Cuckoo Filter      Bloom filter

false positive rate:       $\approx 2^{-M+3}$        $2^{-M} \cdot \ln(2)$

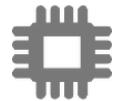
Memory accesses for  
positive query      1.5       $M \cdot \ln(2)$

Memory accesses for  
negative query      2      2

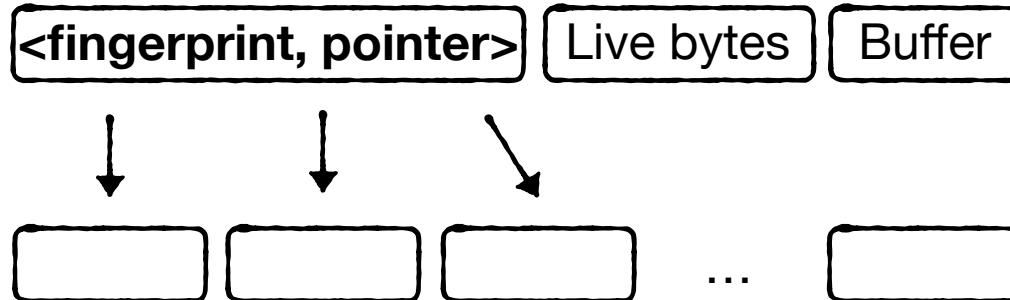
Insertion cost       $O(1)$        $M \cdot \ln(2)$

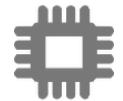
Deletes?      1.5      N/A

**Storing Payloads?**      Yes      N/A



## Cuckoo Filter



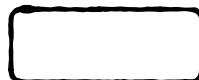


Cuckoo Filter

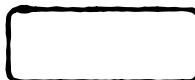
**<FP(X), pointer>**

Live bytes

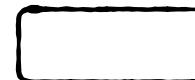
Buffer



**<X, V>**



...



**Our goal was  
reducing this**

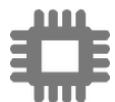
$$\text{Index size} = N * (P + M) / a$$

$N$  = data size

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

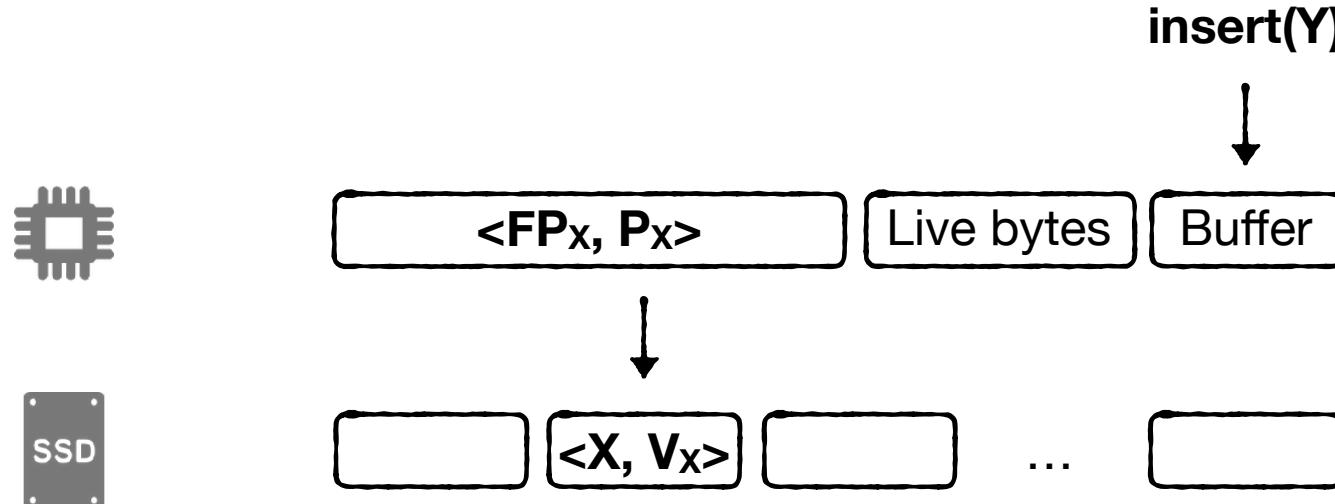
→  $K = \text{key size} = \Omega(\log_2 N)$  →  $M$  bits / entry

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



## Implication of Using a Filter on Insertions/Deletes/Updates

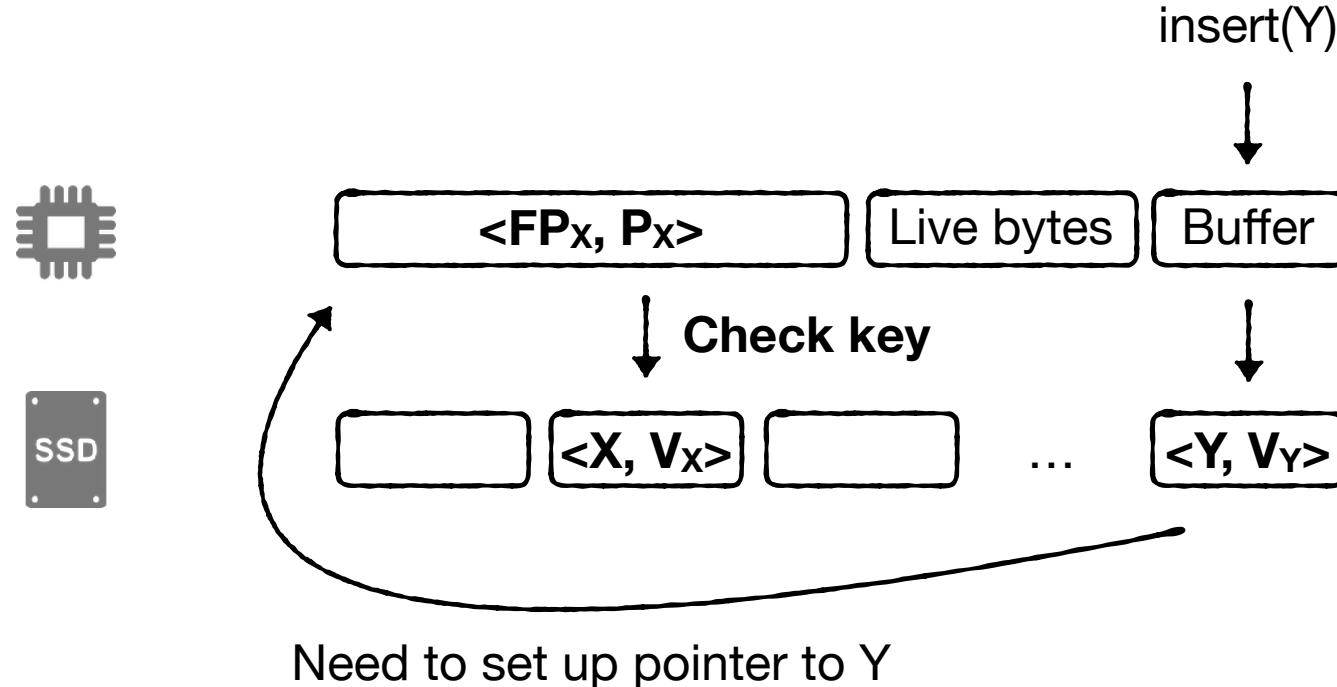
Suppose we insert key Y that has a matching fingerprint to existing key X ( $FP_X = FP_Y$ )



## Implication of Using a Filter on Insertions/Deletes/Updates

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

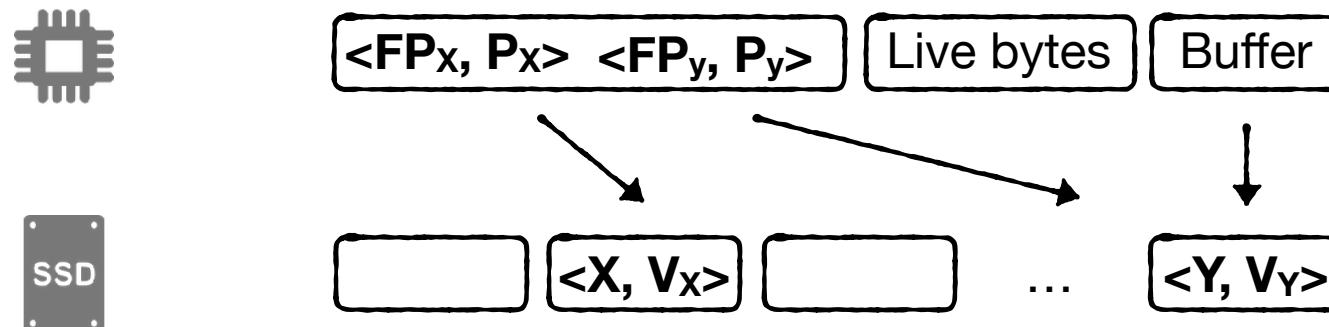
**To safeguard against orphaning, we must issue read-before-write**



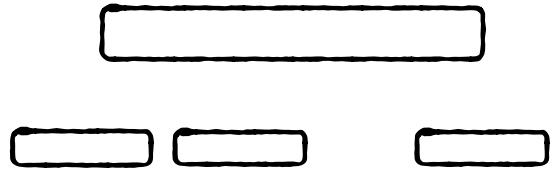
## Implication of Using a Filter on Insertions/Deletes/Updates

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

(excluding checkpointing)

Insert:

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

Gets:

$O(1+2^{-M+3})$

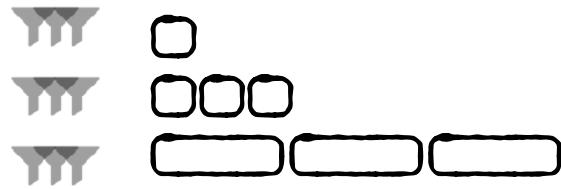
Scan:

$O(N/B)$

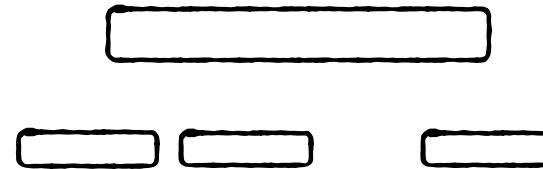
Memory  
(bits / entry)

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

## Basic LSM-tree w. Monkey



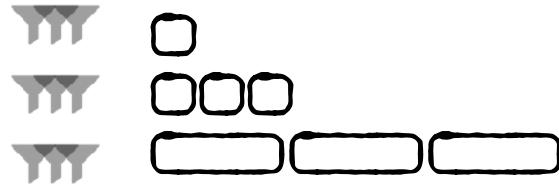
## Circular Log w. Cuckoo Filter



Updates/ Deletes:	$O(\log_2(N/P) / B)$ reads & writes
Insert:	$O(\log_2(N/P) / B)$ reads & writes
Gets:	$O(1+e^{-M} * \ln(2))$
Scan:	$O(\log_2 N/P)$
Memory (bits / entry)	$O(K/B + M)$

$O(1+2^{-M+3})$ reads & $O(GC/B)$ writes
$O(2^{-M+3})$ read & $O(GC/B)$ writes
$O(1+2^{-M+3})$
$O(N/B)$
$O(\log_2(N/B) + M))$

## Basic LSM-tree w. Monkey



	Basic LSM-tree w. Monkey	Circular Log w. Cuckoo Filter
Updates/ Deletes:	$O(\log_2(N/P) / B)$ reads & writes	$O(1+2^{-M+3})$ reads & $O(GC/B)$ writes
Insert:	$O(\log_2(N/P) / B)$ reads & writes	$O(2^{-M+3})$ read & $O(GC/B)$ writes
Gets:	$O(1+e^{-M} * \ln(2))$	$O(1+2^{-M+3})$
Scan:	$O(\log_2 N/P)$	$O(N/B)$
Memory (bits / entry)	$O(K/B + M)$	$O(\log_2(N/B) + M))$
Recovery	<b>Swift</b>	<b>Long</b>