

Download this presentation here:

Intro to probabilistic data structures



About Me

- Software Engineer @ AssemblyAI
- Interest in Algorithms & DS and AI

Agenda

- membership
- cardinality
- frequency

In total we'll see five PDS

Probabilistic Data Structures return approximate values:

- Approximate values in case of numbers
- False positives in case of booleans

but they save a lot of memory

TYPE I ERROR: FALSE POSITIVE
TYPE II ERROR: FALSE NEGATIVE
TYPE III ERROR: TRUE POSITIVE FOR
INCORRECT REASONS
TYPE IV ERROR: TRUE NEGATIVE FOR
INCORRECT REASONS
TYPE V ERROR: INCORRECT RESULT WHICH
LEADS YOU TO A CORRECT
CONCLUSION DUE TO
UNRELATED ERRORS
TYPE VI ERROR: CORRECT RESULT WHICH
YOU INTERPRET WRONG
TYPE VII ERROR: INCORRECT RESULT WHICH
PRODUCES A COOL GRAPH
TYPE VIII ERROR: INCORRECT RESULT WHICH
SPARKS FURTHER RESEARCH
AND THE DEVELOPMENT OF
NEW TOOLS WHICH REVEAL
THE FLAW IN THE ORIGINAL
RESULT WHILE PRODUCING
NOVEL CORRECT RESULTS
TYPE IX ERROR: THE RISE OF SKYWALKER

Membership

Checking if an item belongs to a set



Membership use cases

In general when a wrong answer does not involve correctness but more work

Some examples:

- Taken username
- Ad placement: has the user already seen this ad?
- Fraud detection (has the user paid from this location before?)

Membership

- Bloom Filter
- Counting Bloom Filter
- Cuckoo Filter

Bloom Filter

Invented by Burton H. Bloom in 1970

Consists of :

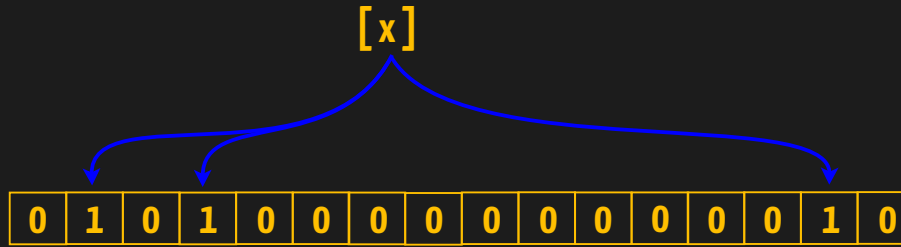
- a bit array of size m
- k different hash functions

Membership condition :

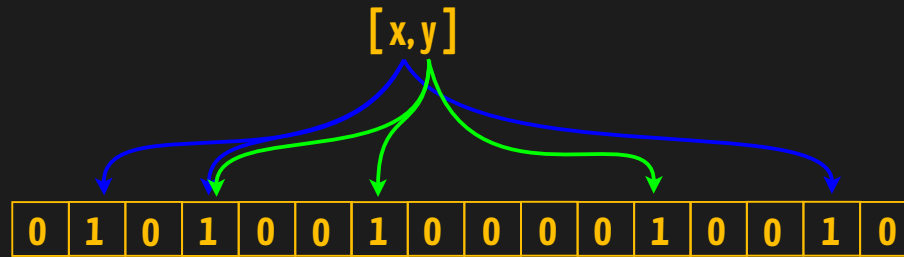
- if all bits are set

Bloom Filter with $m=16, k=3$

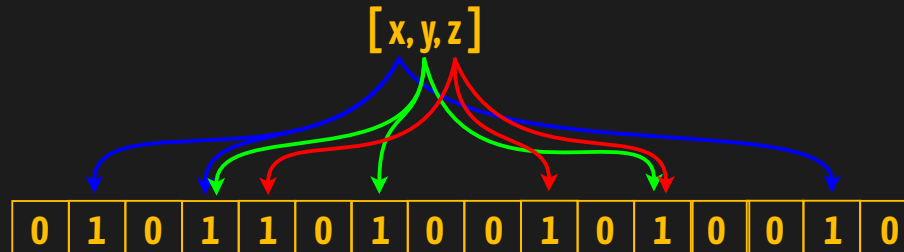
Insert x



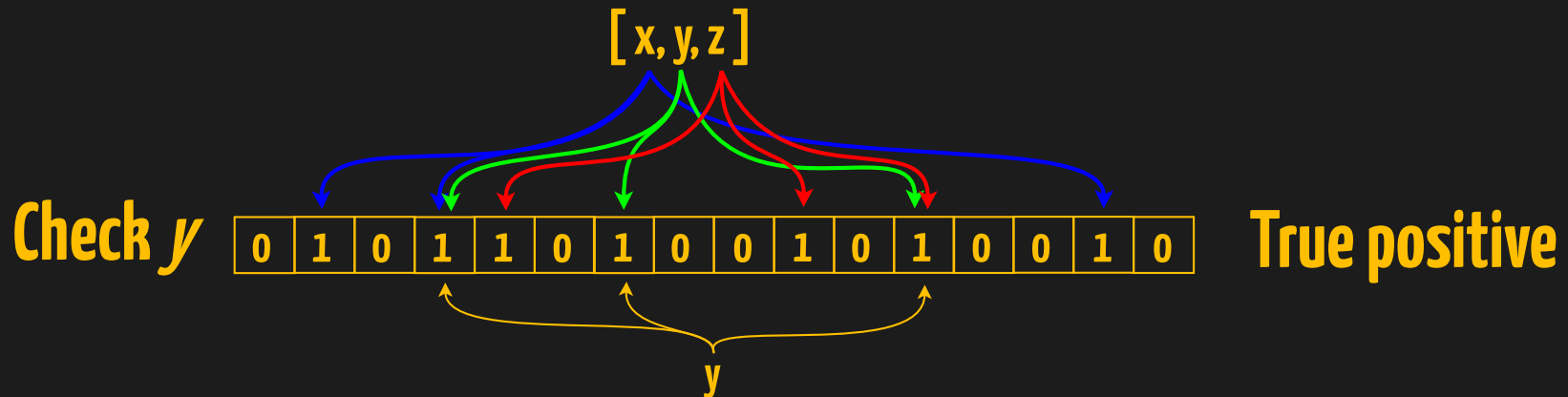
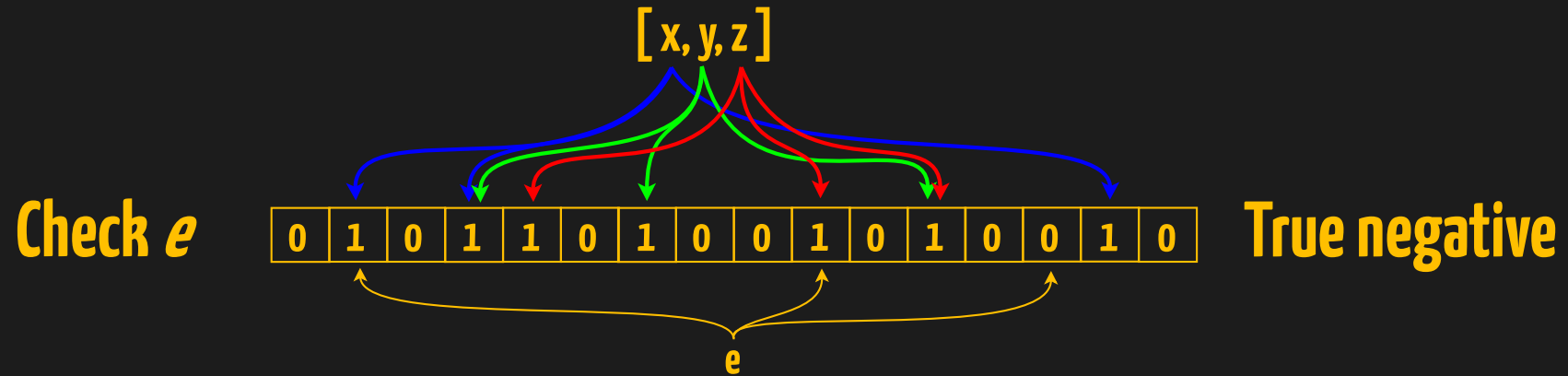
Insert y



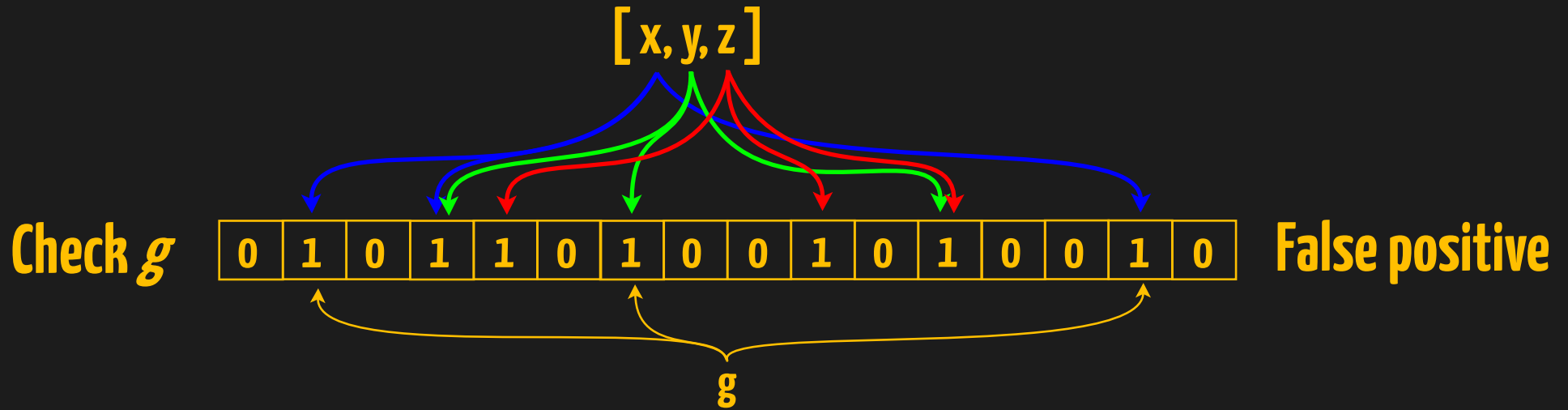
Insert z



Bloom Filter with $m=16, k=3$



Bloom Filter with $m=16, k=3$



Bloom Filter

Approximate optimal size

Given:

- n : the number of expected items
- ε : the wanted false positive rate, where $\varepsilon \in [0, 1]$

Then:

- $m = -\frac{n \cdot \ln(\varepsilon)}{\ln(2)^2}$
- $k = -\log_2(\varepsilon)$

To check 1M items, with 1% error rate:

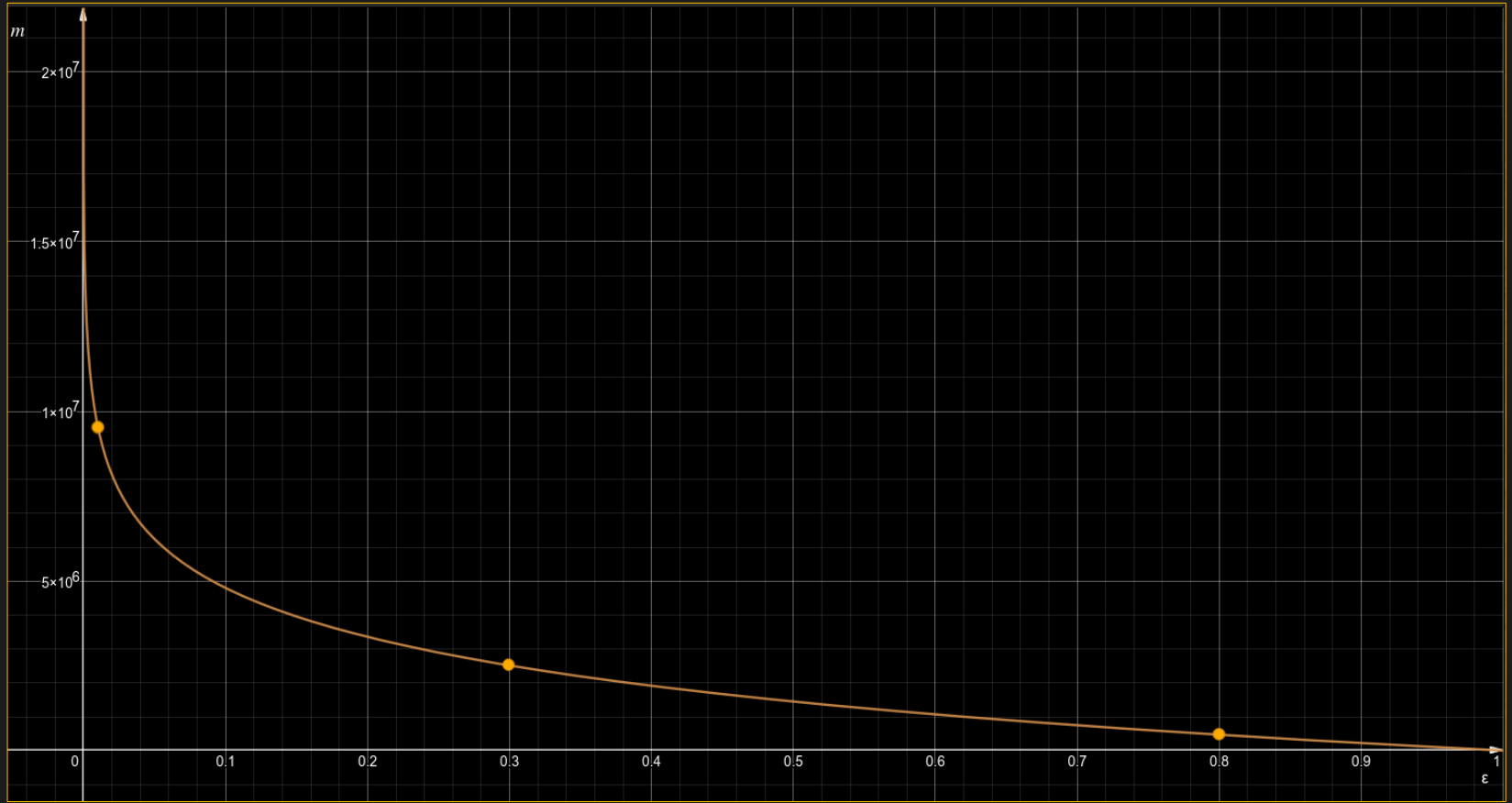
- $m = -\frac{1,000,000 \cdot -4.6}{0.48} \simeq 9.5\text{M} \simeq 1.2\text{ MB}$

- $k \simeq 6.6 = 7$

We have a capped size!

Bloom Filter

Given $n = 1\text{M}$



$$\epsilon = 0.01 \Rightarrow m = 9.5\text{M}$$

$$\epsilon = 0.3 \Rightarrow m = 2.5\text{M}$$

$$\epsilon = 0.8 \Rightarrow m = 460\text{K}$$

Bloom Filter

Fruit flies olfactory neural circuit evolved a variant of a Bloom filter to assess the novelty of odors!

It adds two additional features:

- based on similarity to previously experienced odors
- time elapsed since the odor was last experienced

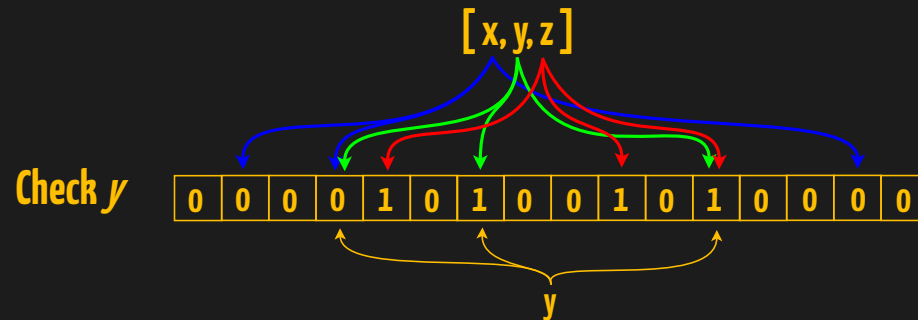
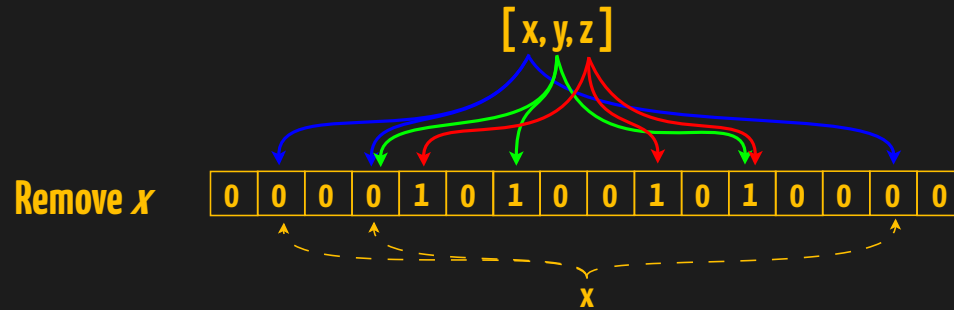
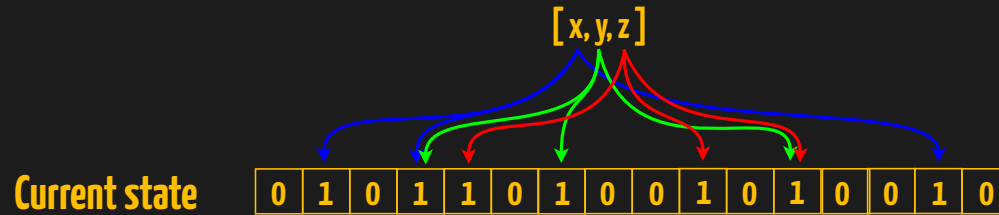


(source: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6304992/>)

Bloom Filter Code

```
01 class StringBloomFilter(expectedSize: Int, errorRate: Double) {
02
03     private val m: Int = -(expectedSize * ln(errorRate) / ln2squared).toInt() // bitset size
04     private val k = ceil(-log2(errorRate)).toInt() // number of hash functions
05     private val bitSet = BitSet(m)
06     private val hashers = (1..k).map { n → Hasher(primes[n + 4], primes[3 * n + 3]) }
07
08     fun contains(item: String)=hashers.all { hasher→bitSet.get(abs(hasher.hashCode(item)) % m) }
09
10     fun add(item: String)=hashers.forEach { hasher→bitSet.set(abs(hasher.hashCode(item)) % m, true)}
11
12     private class Hasher(private val base: Int, private val multiplier: Int) {
13         fun hashCode(value: String): Int = value
14             .map { it.code }
15             .fold(base) { acc, curr → acc * multiplier + curr }
16     }
17 }
```

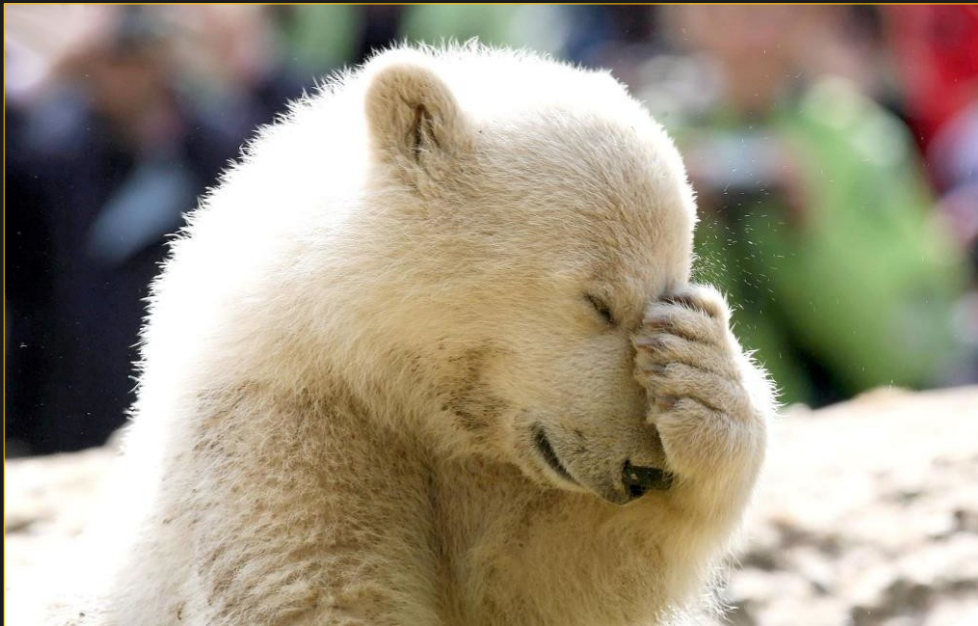
What if we remove an element?



False negative!

Bloom Filter

What's the use of a DS that may return both false positives and false negatives?



Counting Bloom Filter

Consists of:

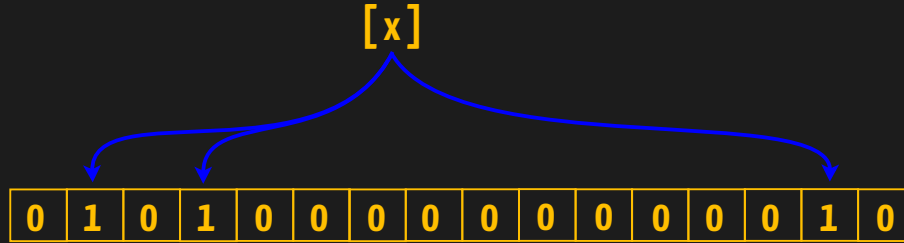
- a count array of size m
- k different hash functions

Membership condition:

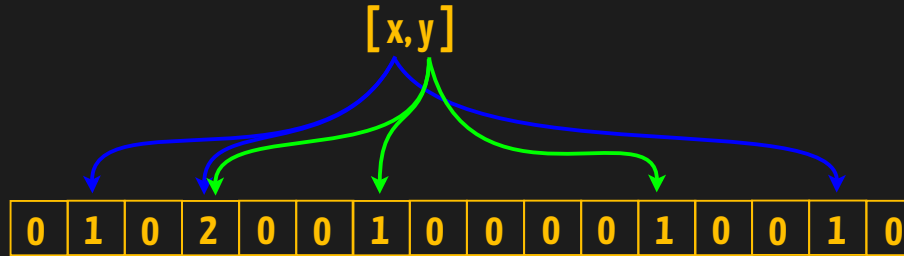
- if all counts are > 0

Counting Bloom Filter

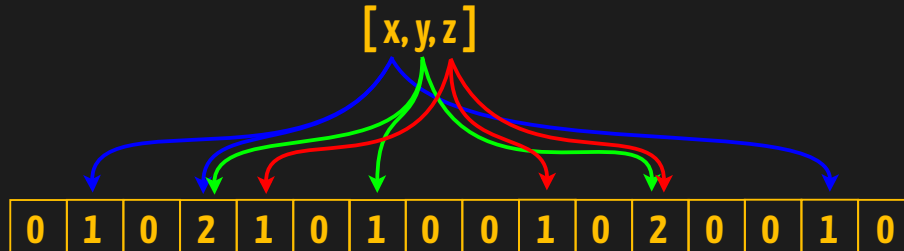
Insert x



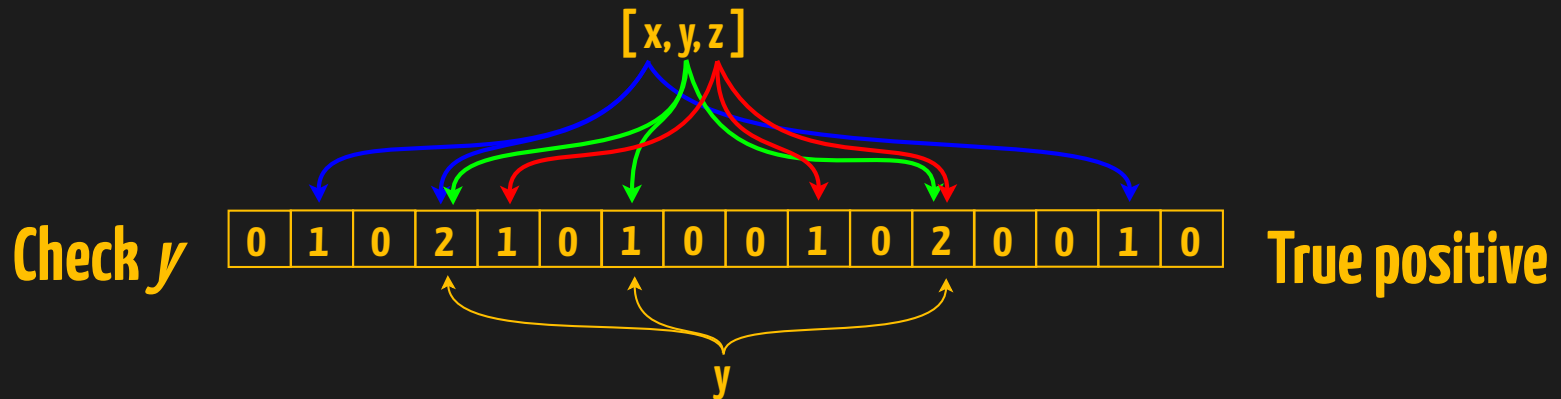
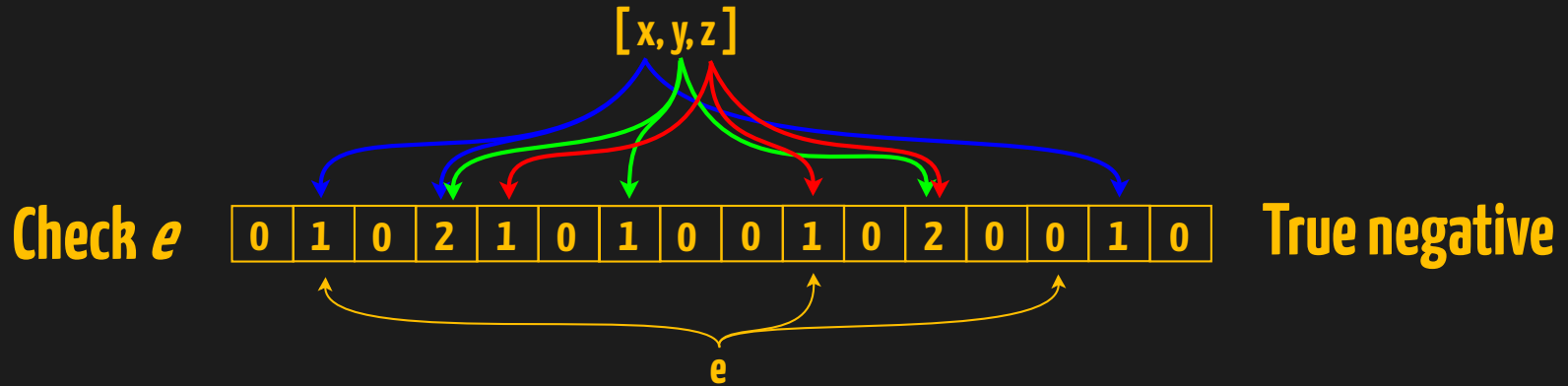
Insert y



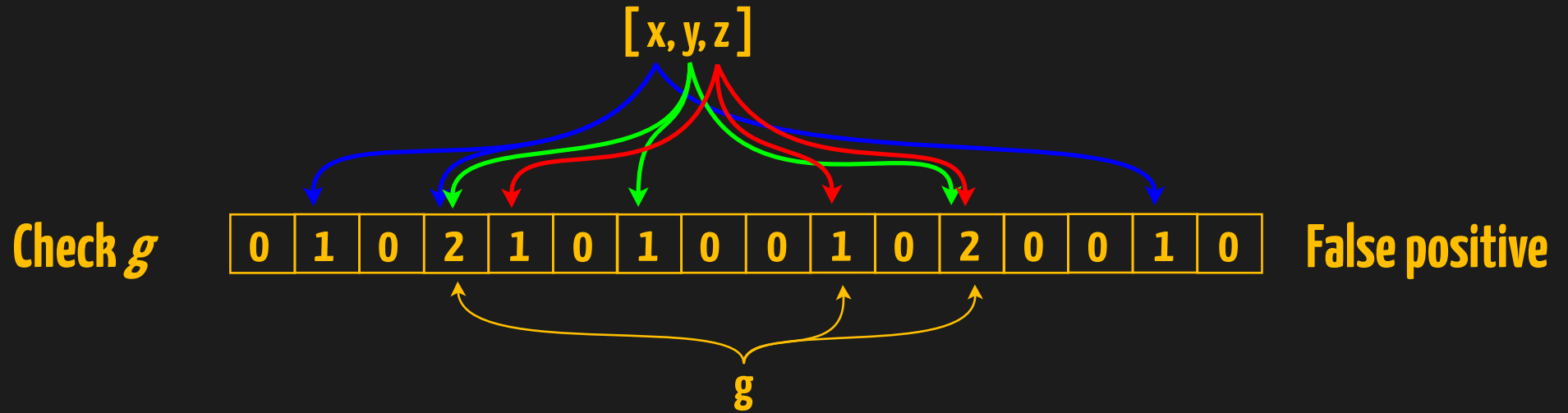
Insert z



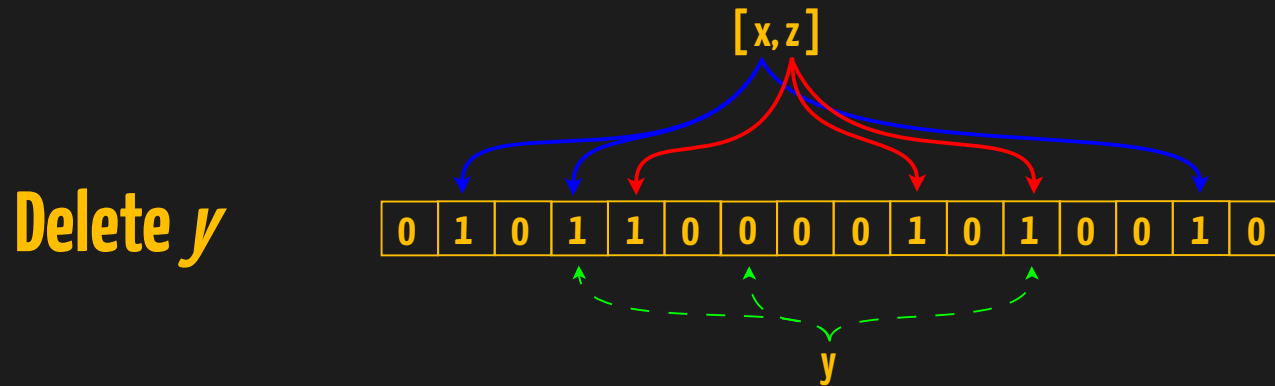
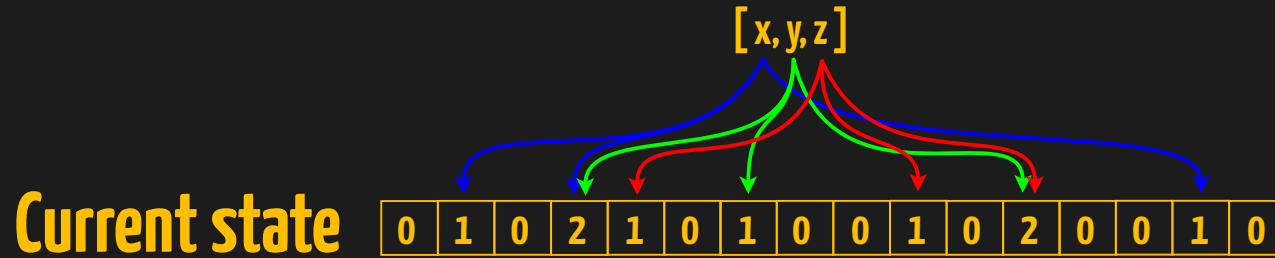
Counting Bloom Filter



Counting Bloom Filter



Counting Bloom Filter



Counting Bloom Filter Code

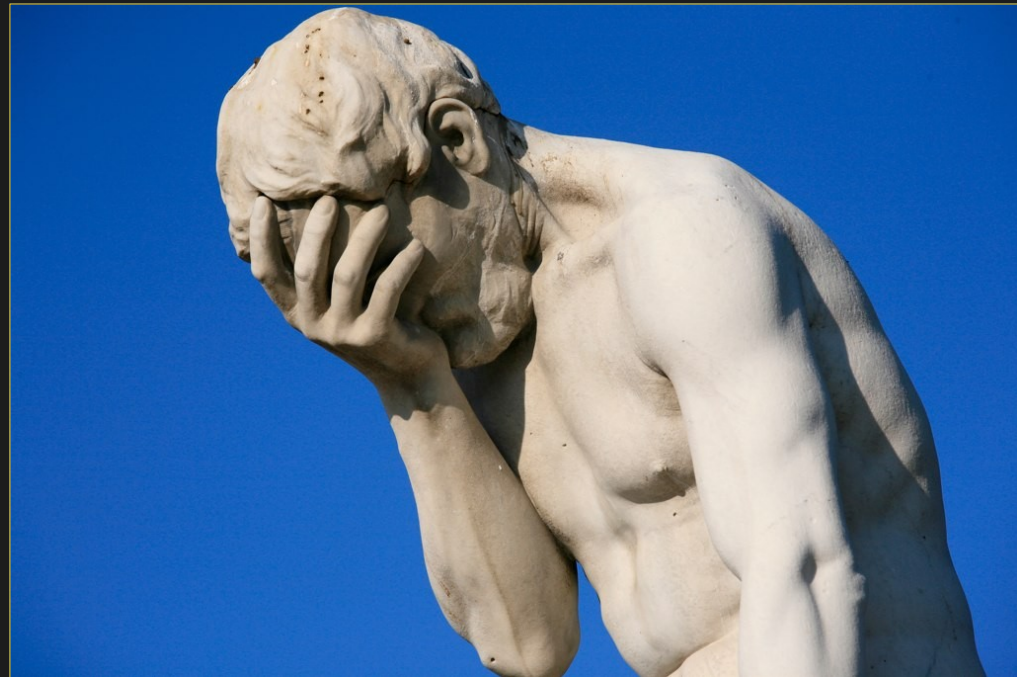
```
01 class StringCountingBloomFilter(expectedSize: Int, errorRate: Double) {
02   private val m: Int = -(expectedSize * ln(errorRate) / ln2squared).toInt()
03   private val k = ceil(-log2(errorRate)).toInt()
04   private val counters = ByteArray(m) { Byte.MIN_VALUE }
05   private val hashers = (1..k).map { n → Hasher(primes[n + 4], primes[3 * n + 3]) }
06
07   fun add(item: String)=hashers.forEach { counters[abs(it.hashCode(item)) % m]++ }
08   fun delete(item: String)=hashers.forEach { counters[abs(it.hashCode(item)) % m]-- }
09   fun contains(item: String)=hashers.all { counters[abs(it.hashCode(item)) % m] > MIN_VALUE }
10
11   private class Hasher(private val base: Int, private val mult: Int) {
12     fun hashCode(value: String) = value.map { it.code }
13       .fold(base) { acc, curr → acc * mult + curr }
14   }
15 }
```

Counting Bloom Filter

It uses a lot more memory than Bloom Filter:

- 8x with byte
- 16x with short
- 32x with int

It also might overflow



Cuckoo Filter

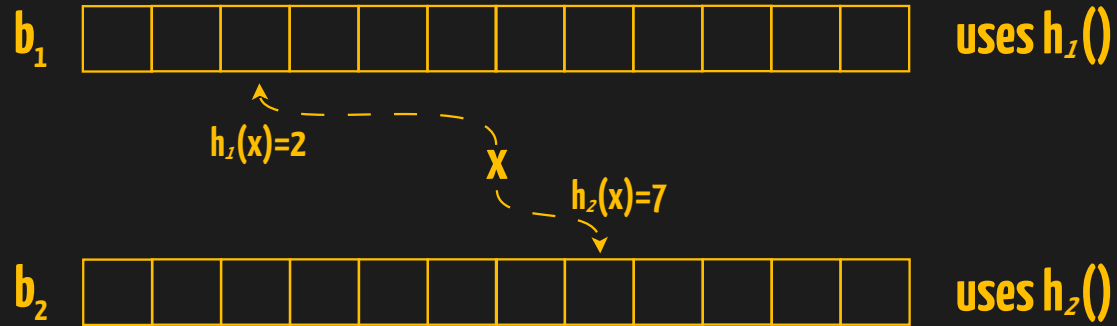
Described by Fan, Andersen, Kaminsky, and Mitzenmacher in 2014

Based on Cuckoo Hashing (a *classic* Data Structure), which consists of:

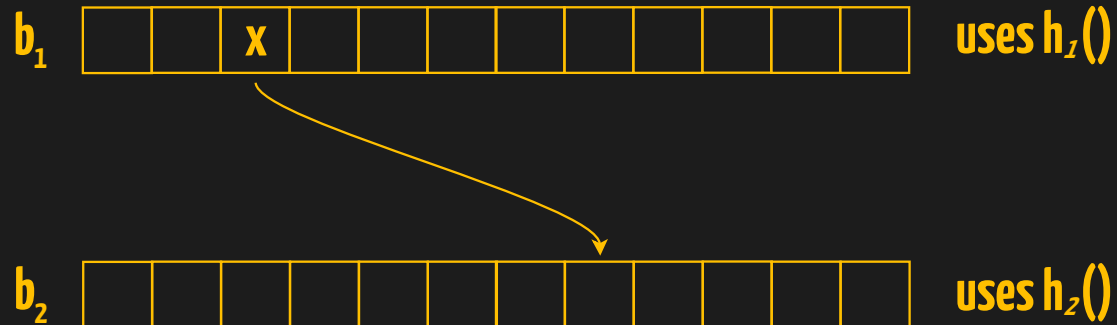
- 2 arrays: b_1 and b_2
- 2 hash functions: $h_1()$ and $h_2()$

Cuckoo Hashing

Choose location for x

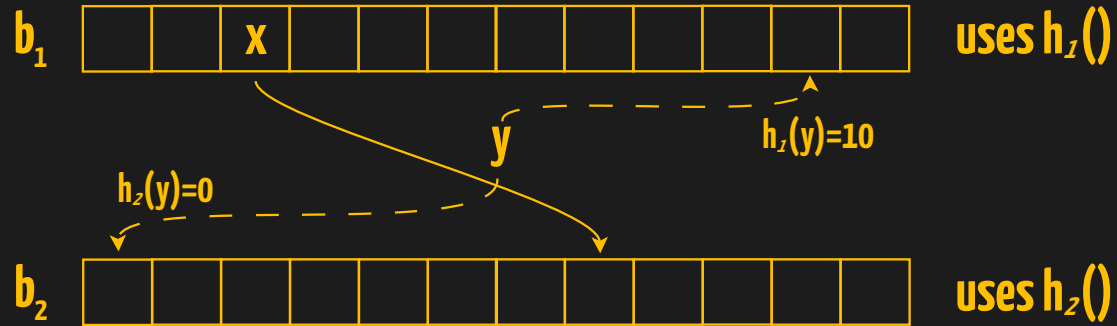


Insert

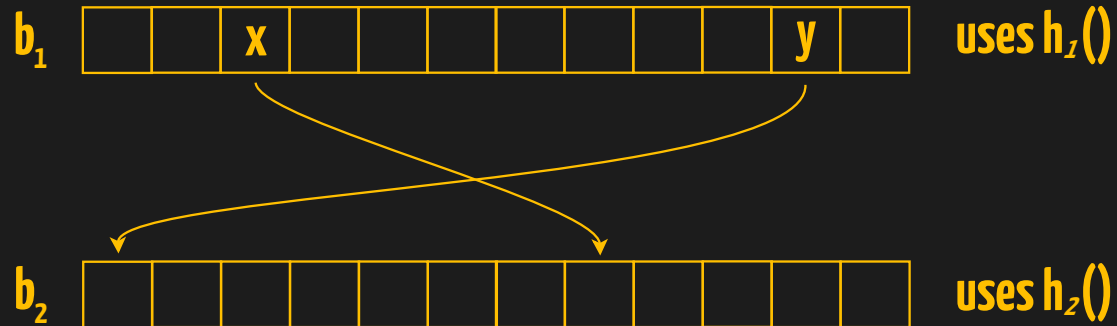


Cuckoo Hashing

Choose location for y

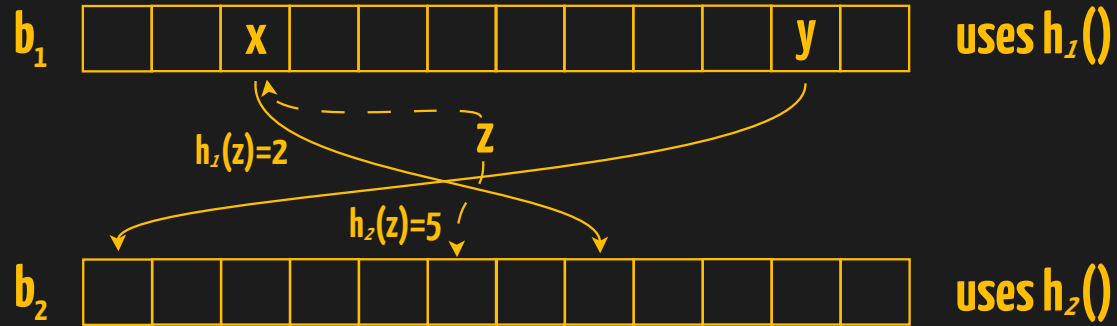


Insert

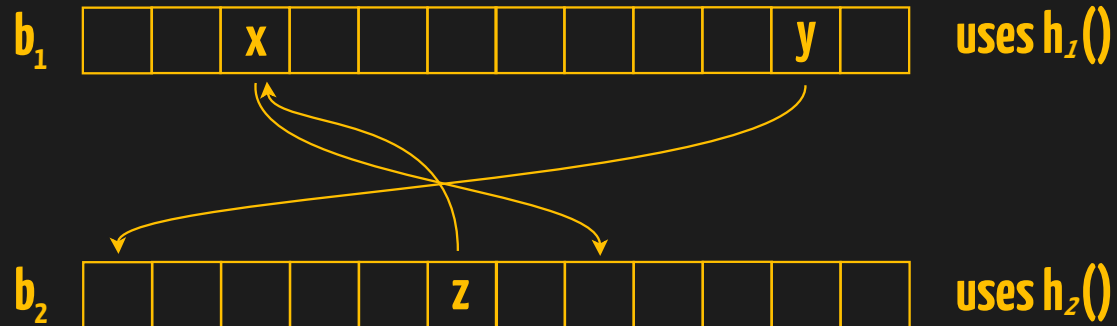


Cuckoo Hashing

Choose location for z

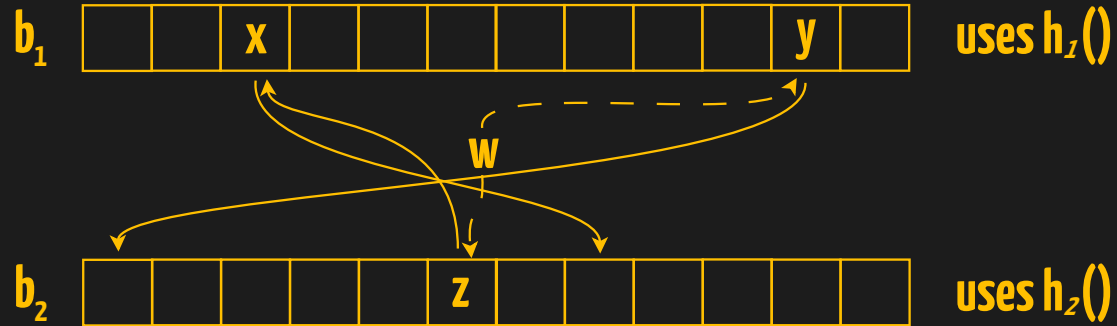


Insert

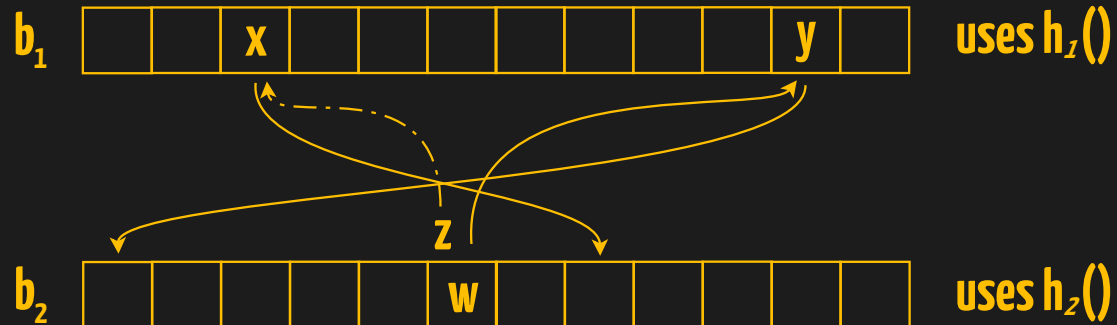


Cuckoo Hashing

Choose location for w

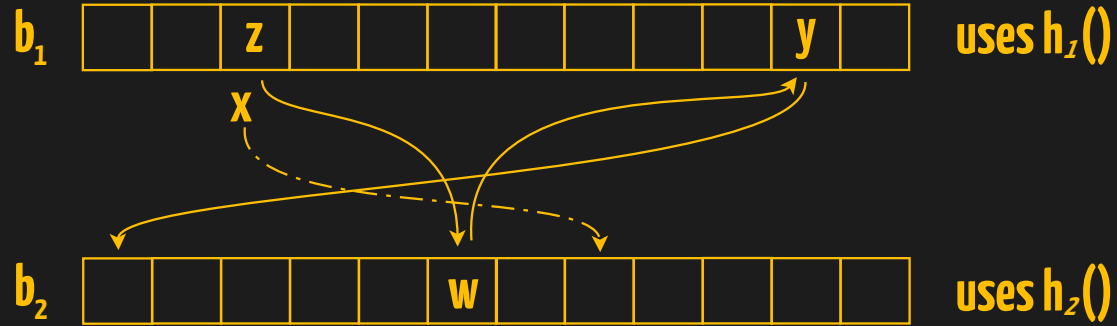


Randomly choose b_2 and kick out z !

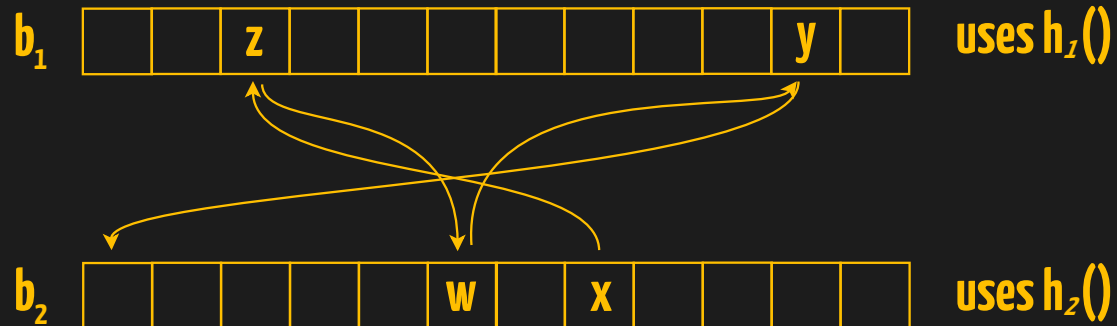


Cuckoo Hashing

z kicks out x



Insert x into its alternate location



Cuckoo Hashing

Performance

- $O(1)$ lookup runtime complexity
- $O(1)$ delete runtime complexity
- Insert: if it fails (cycle or threshold) the two hash tables need to be enlarged and rehashed
- $O(1)$ insert amortized runtime complexity

Cuckoo Filter

- to save memory instead of storing the whole item we store a fingerprint (f bit)
- how can we get the alternate location if we don't have the kicked out item?

Solution: create a hash that keeps track of fingerprint

Cuckoo Filter

- a fingerprint function returning a f -bit value

- two hash functions
$$\begin{cases} h_1(x) = \text{hash}(x) \\ h_2(x) = h_1(x) \oplus \text{hash}(\text{fingerprint}(x)) \end{cases}$$

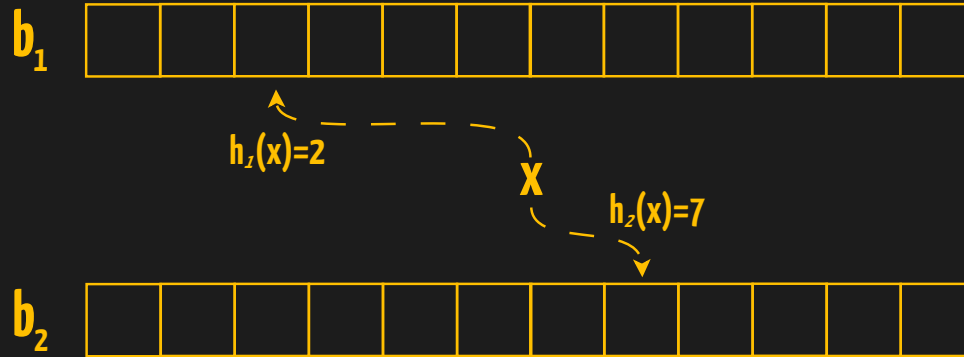


$$h_1(x) = h_2(x) \oplus \text{hash}(\text{fingerprint}(x))$$

(thanks to XOR symmetry)

Cuckoo Filter

Choose location for x



Insert fingerprint(x)



$$\begin{cases} h_1(x) = \text{hash}(x) \\ h_2(x) = h_1(x) \oplus \text{hash}(\text{fingerprint}(x)) \end{cases}$$

Cuckoo Filter

x is kicked out by z



Move fingerprint(x) and insert fingerprint(z)



Cuckoo Filter

Lookup of x

$$f(x) == b_1[h_1(x)] \quad || \quad f(x) == b_2[h_2(x)]$$

When a matching fingerprint is found in any of the b arrays, the entry might be in the filter. We can have false positives when a different entry has the same fingerprint as the searched one.

Cuckoo Filter

Cuckoo filter requires $\frac{\log_2 \left(\frac{1}{\epsilon} \right) + 1 + \log_2 (b)}{\alpha}$ bits per key

where α is the load factor of the cuckoo hash table.

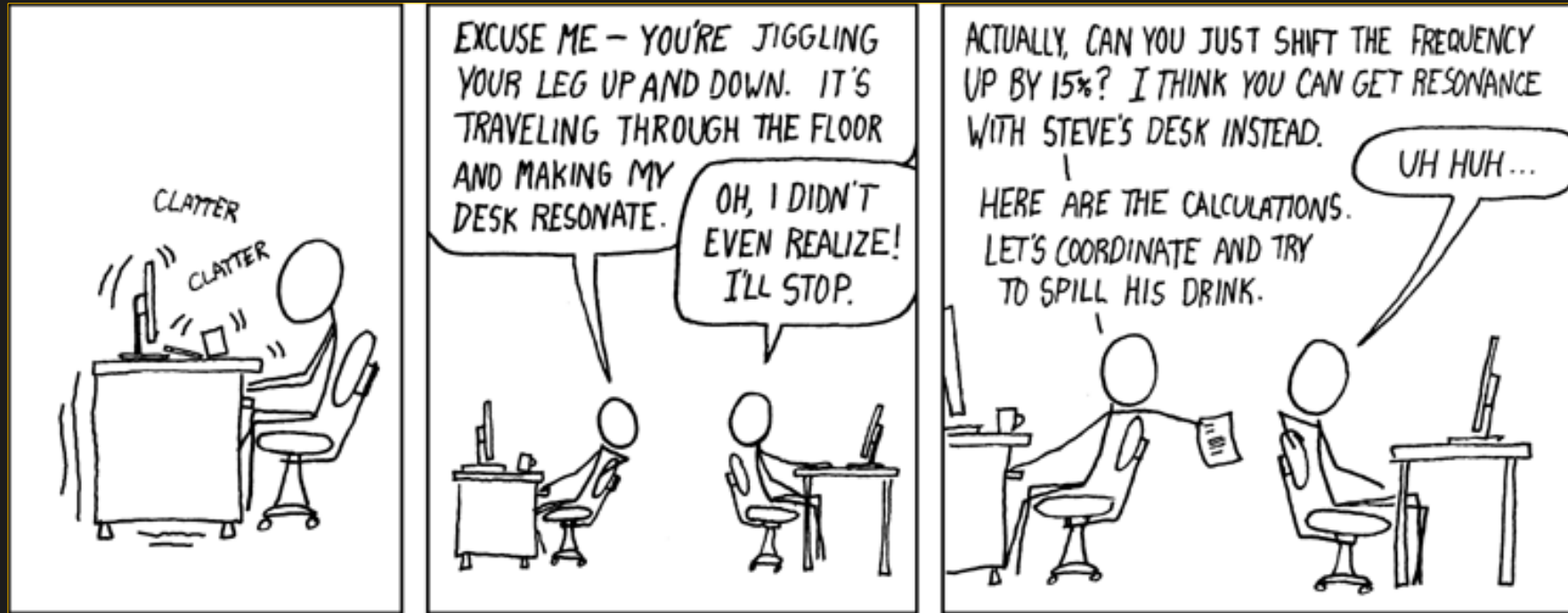
Cuckoo Filter Code

```
01 fun insert(item: String): Boolean {
02     var f = fingerprint(item)
03     val index1 = hash(item)
04     if (bucket[index1] == EMPTY) {
05         bucket[index1] = f
06         return true
07     }
08     val index2 = index1 xor hash(f)
09     if (bucket[index2] == EMPTY) {
10         bucket[index2] = f
11         return true
12     }
13
14     // relocate existing items
15     var i = if (Random.nextBoolean()) index1 else index2
16     for (n in 0..MAX_KICKS) {
17         f = bucket[i].also { bucket[i] = f }
18         i = i xor hash(f)
19         if (bucket[i] == EMPTY) {
20             bucket[i] = f
21         }
22     }
23     return false // Max kicks reached
24 }
```

Membership Libraries

- **Guava:** <https://github.com/google/guava> (Bloom Filter)
- **Hadoop** <https://hadoop.apache.org/> (Bloom and Counting Bloom Filters)
- **Fastfilter** https://github.com/FastFilter/fastfilter_java (Bloom, Counting Bloom, Cuckoo and other filters)
- **Boom Filters:** <https://github.com/tylertreat/BoomFilters> (Bloom, Counting Bloom, Cuckoo and other filters)

Frequency



(source: <https://xkcd.com/228/>)

Frequency

Counting the number of times an item appears in a stream

- Count-Min Sketch

Frequency Use Cases

Some examples:

- Unique users of a service
- Top k items (e.g. top players in online gaming)
- Network traffic analysis

Count-Min Sketch

Invented by G. Cormode et al in 2005

- d hash functions (pairwise independent)
- a bidimensional array *count* of size $d \times w$ for counting the items, where:

$$d = \ln\left(\frac{1}{\delta}\right)$$

ε is the error rate and δ its probability

$$w = \left(\frac{e}{\varepsilon}\right)$$

If we want an error of at most 0.1% (of the sum of all frequencies) with 99.9% certainty, then:

$$w = \left(\frac{e}{0.001}\right) \simeq 2,718 \quad d = \ln\left(\frac{1}{0.001}\right) \simeq 6.9 = 7$$

Using 32 bit counters, `sizeof(count)` = $w \cdot d \cdot 4 \simeq 76$ KB

Count-Min Sketch

$d = 4, w = 16$

Initial state

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Hash Function #1

Hash Function #2

Hash Function #3

Hash Function #4

Count-Min Sketch

$d = 4, w = 16$

Insert x

$$h_1(x) = 2$$

$$h_2(x) = 13$$

$$h_3(x) = 12$$

$$h_4(x) = 7$$

0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0

Hash Function #1

Hash Function #2

Hash Function #3

Hash Function #4

Count-Min Sketch

$d = 4, w = 16$

Insert y

$$h_1(y) = 6$$

$$h_2(y) = 13$$

$$h_3(y) = 8$$

$$h_4(y) = 2$$

0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0
0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0
0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0

Hash Function #1

Hash Function #2

Hash Function #3

Hash Function #4

Count-Min Sketch

$d = 4, w = 16$

Insert x

$$h_1(x) = 2$$

$$h_2(x) = 13$$

$$h_3(x) = 12$$

$$h_4(x) = 7$$

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

Hash Function #1

Hash Function #2

Hash Function #3

Hash Function #4

Count-Min Sketch

$d = 4, w = 16$

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

Hash Function #1

Hash Function #2

Hash Function #3

Hash Function #4

$$\begin{aligned}\text{Count}(x) &= \min(\text{count}[0][h_1(x)], \\ &\quad \text{count}[1][h_2(x)], \\ &\quad \text{count}[2][h_3(x)], \\ &\quad \text{count}[3][h_4(x)]) \\ &= \min(2, 3, 2, 2) = 2\end{aligned}$$

$$\begin{aligned}\text{Count}(y) &= \min(\text{count}[0][h_1(y)], \\ &\quad \text{count}[1][h_2(y)], \\ &\quad \text{count}[2][h_3(y)], \\ &\quad \text{count}[3][h_4(y)]) \\ &= \min(1, 3, 1, 1) = 1\end{aligned}$$

Count-Min Sketch

If there are collisions on hashes we get higher counts

Count-Min Sketch Code

```
01 class CountMinSketch(d: Int, private val w: Int) {
02     private val count = Array(d) { LongArray(w) }
03     private val hashers = (1..k).map { n → Hasher(primes[n + 4], primes[3 * n + 3]) }
07
08     fun add(item: String) = hashers.forEachIndexed { idx, hasher →
09         count[idx][hasher.hashCode(item)]++
10     }
11
12     fun count(item: String) = hashers
13         .mapIndexed { idx, hasher → count[idx][hasher.hashCode(item)] }
14         .min()
15
16     private class Hasher(val size: Int, val base: Long, val multiplier: Long) {
17         fun hashCode(value: String) = value.map { it.code }
18             .fold(base){ acc, curr → acc * multiplier + curr }
19             .mod(size)
20     }
21 }
```

Frequency Libraries

- **Apache Spark:** <https://spark.apache.org/>
- **Boom Filters:** <https://github.com/tylertreat/BoomFilters>

Cardinality



Cardinality

Counting the number of distinct items in a collection

- HyperLogLog

Cardinality Use Cases

Some examples:

- **Google BigQuery:** APPROX_COUNT_DISTINCT, APPROX_QUANTILES, APPROX_TOP_COUNT, APPROX_TOP_SUM
- **Snowflake:** HLL, HLL_ACCUMULATE, HLL_COMBINE, HLL_ESTIMATE
- **Number of users in search engines (where ads are paid per user)**
- **Materialized views in Data Warehouses**

HyperLogLog

Flajolet et al. 2007

Binary representation of random numbers:

- $\frac{1}{2}$ start with “0”
 - $\frac{1}{2}$ start with “1”
 - $\frac{1}{4}$ start with “00”
 - $\frac{1}{4}$ start with “01”
 - $\frac{1}{4}$ start with “10”
 - $\frac{1}{4}$ start with “11”
 - $\frac{1}{8}$ start with “000”
 - $\frac{1}{8}$ start with “001”
 - $\frac{1}{8}$ start with “010”
 - $\frac{1}{8}$ start with “011”
 - $\frac{1}{8}$ start with “111”
 - $\frac{1}{8}$ start with “101”
 - $\frac{1}{8}$ start with “110”
 - $\frac{1}{8}$ start with “111”
- General rule: $p(\text{first } k \text{ bits}) = 2^{-k}$

HyperLogLog

Underlying idea: let's imagine we have a big set of items hashed to numbers; if we add to the set the hash of a new item starting with "0000" we can say that it's likely that the set has size $2^4 = 16$.

Caveats:

- only works on big cardinalities
- hash function must return uniformly distributed binary representations

HyperLogLog

Consists of:

- a hash function
- an array of $\text{int } M$

With only one hash function, we risk that one single item can skew the result, so we could use multiple hash functions.

HyperLogLog

Insert x

$$h_1(x) = 0000100100101001$$

$$h_2(x) = 0011100110011010$$

$$h_3(x) = 0101011110100011$$

$$h_4(x) = 1111100111101001$$

$$h_5(x) = 0010110110000110$$

...

$$h_{16}(x) = 0001100001001100$$

Array M

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

HyperLogLog

Insert y

$$h_1(y) = 1100010101110101$$

$$h_2(y) = 0101110010010101$$

$$h_3(y) = 0001100010101010$$

$$h_4(y) = 0101010100111101$$

$$h_5(y) = 0111001100110100$$

...

$$h_{16}(y) = 0110011010100101$$

Array M

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

HyperLogLog

Computing multiple hash functions is computationally expensive!

Can we do it with only one function?

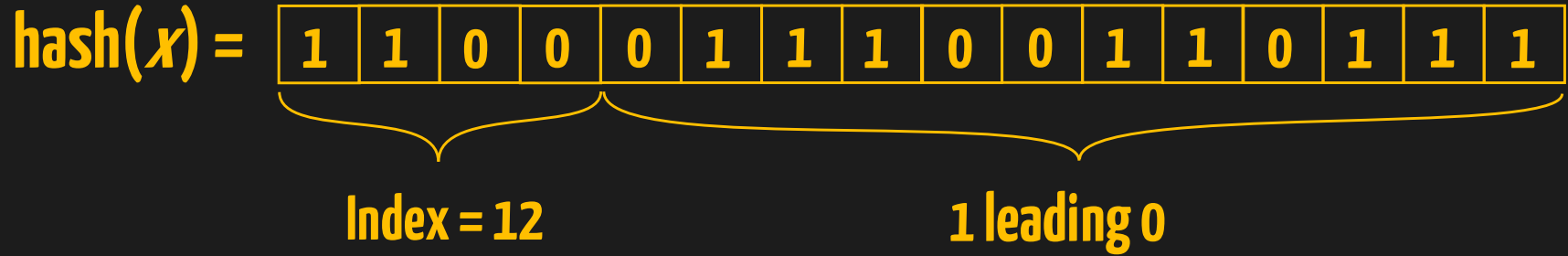
We can split the outcome of one hash function in two and use one part as the index of the array and the other as the value.



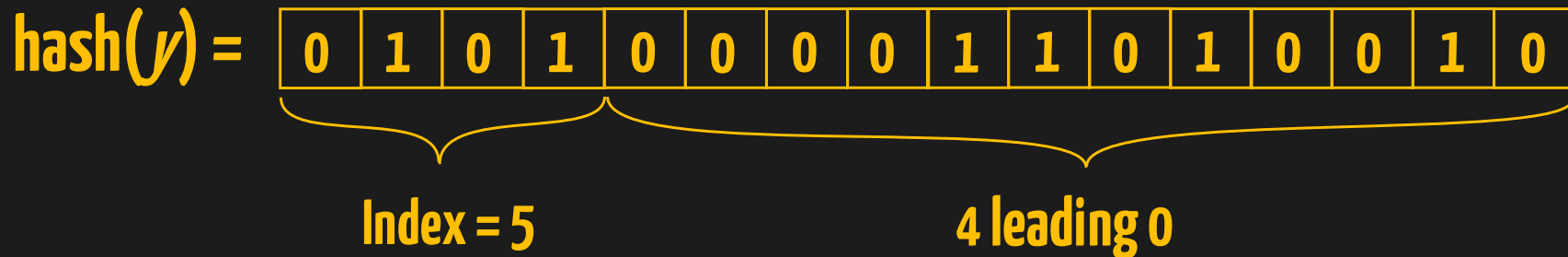
array size $m = 2^4 = 16$

Remaining 12 bits as
the hash value

HyperLogLog



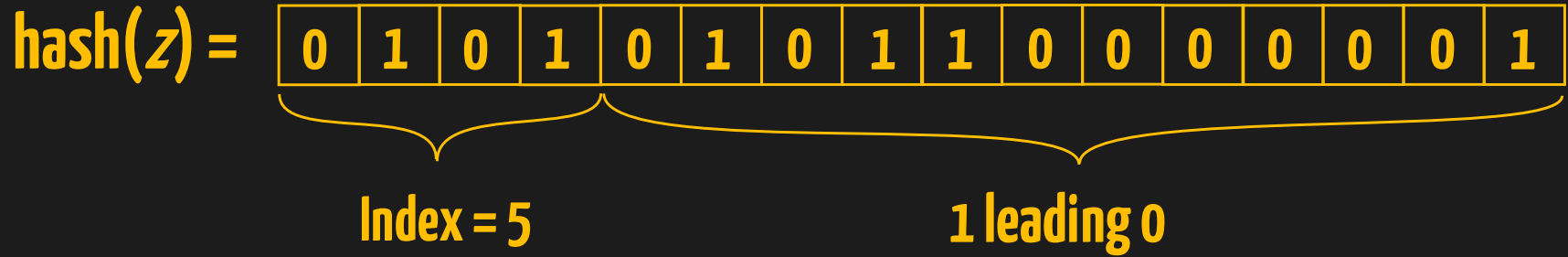
HyperLogLog



Array M

0	0	0	0	0	4	0	0	0	0	0	0	1	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

HyperLogLog



Array M

0	0	0	0	0	4	0	0	0	0	0	0	1	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

HyperLogLog

Get the cardinality

- compute a normalized harmonic mean $Z = \left(\sum_{i=1}^m 2^{-M[i]} \right)^{-1}$ of the values of M
- return Z with a correction factor α : $\text{cardinality} = Z \cdot |M|^2 \cdot \alpha_m$

where $\alpha_m = \begin{cases} \alpha_{16} = 0.673 \\ \alpha_{32} = 0.697 \\ \alpha_{64} = 0.709 \\ \alpha_m = \frac{0.7213}{1 + \frac{1.079}{m}} \text{ for } m \geq 128 \end{cases}$

$$\text{Error rate} = \frac{1.04}{\sqrt{(m)}}$$

HyperLogLog

Array M

4	5	2	3	5	4	7	2	6	5	4	5	3	6	2	5
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

$$\text{cardinality} = Z \cdot M^2 \cdot \alpha \simeq 0.72 \cdot 256 \cdot 0.673 \simeq 124.59$$

Array M

4	5	12	3	5	4	7	2	6	5	4	5	3	6	2	5
---	---	----	---	---	---	---	---	---	---	---	---	---	---	---	---

$$\text{cardinality} = Z \cdot M^2 \cdot \alpha \simeq 0.88 \cdot 256 \cdot 0.673 \simeq 152.06$$

HyperLogLog Code

```
01 class HyperLogLog {
02     private val p: Int = 4
03     private val m: Int = 2.0.pow(p.toDouble()).toInt()
04     private val alpha = 0.673
05     private val buckets = IntArray(m)
06
07     fun add(item: String) {
08         val hash = item.hashCode()
09         val index = hash.ushr(32 - p)
10         val leadingZeroes = Integer.numberOfLeadingZeros(hash shl p) + 1
11         buckets[index] = max(buckets[index], leadingZeroes)
12     }
13
14     fun count(): Long {
15         val harmonicMean = buckets.sumOf { 1.0 / (1 shl it) }
16         val estimate = alpha * (m * m).toDouble() / harmonicMean
17         return if (estimate ≤ 5.0 / 2.0 * m) {
18             (m * Ln(m.toDouble() / estimate)).toLong()
19         } else {
20             estimate.toLong()
21         }
22     }
23 }
```

Cardinality Libraries

- **Apache DataSketches:** <https://datasketches.apache.org>
- **Boom Filters:** <https://github.com/tylertreat/BoomFilters>

Probabilistic Data Structures

Pros: - Save huge amount of memory

Cons: - Uncertainty
- Cold start problem: need snapshots
- Distributed?

Original Papers

- **Bloom filter:** <https://dl.acm.org/doi/pdf/10.1145/362686.362692>
- **Cuckoo Filter:** <https://www.cs.cmu.edu/~dga/papers/cuckoo-conext2014.pdf>
- **Count-min Sketch:** <http://dimacs.rutgers.edu/~graham/pubs/papers/cm-full.pdf>
- **HyperLogLog:** <https://algo.inria.fr/flajolet/Publications/FlFuGaMe07.pdf>

Code shown in this presentation: <https://github.com/andreaiacono/TalkProbabilisticDataStructures>

Slides of this presentation:



Questions?