




← You are looking at a stream of data.

In networking, we usually talk about **streams of packets**,  
but these questions apply to other domains as well,  
e.g. **search engines and databases**.




← You are looking at a stream of data (packets).  
There are many questions you might ask:

Is a certain element (e.g. ip address) in the stream?

How frequently does an element appear?

How many distinct elements are in the stream?

What are the most frequent elements?



← You are looking at a stream of data (packets).  
There are many questions you might ask:


Is a certain element (e.g. ip address) in the stream?

How frequently does an element appear?

How many distinct elements are in the stream?

What are the most frequent elements?

**In P4, these questions are difficult to answer.**



← You are looking at a stream of data (packets).  
Today, I'll show you how set membership and frequency queries can be realized in P4.

**PART 1**

Is a certain element (e.g. ip address) in the stream?  
→ Bloom filter

**PART 2**

How frequently does an element appear?  
→ CountMin Sketch, Count Sketch, ...

# part 1: set membership queries with Bloom filters

*Is a certain element (e.g. ip address) in the stream?*

(slides by Thomas Holterbach)

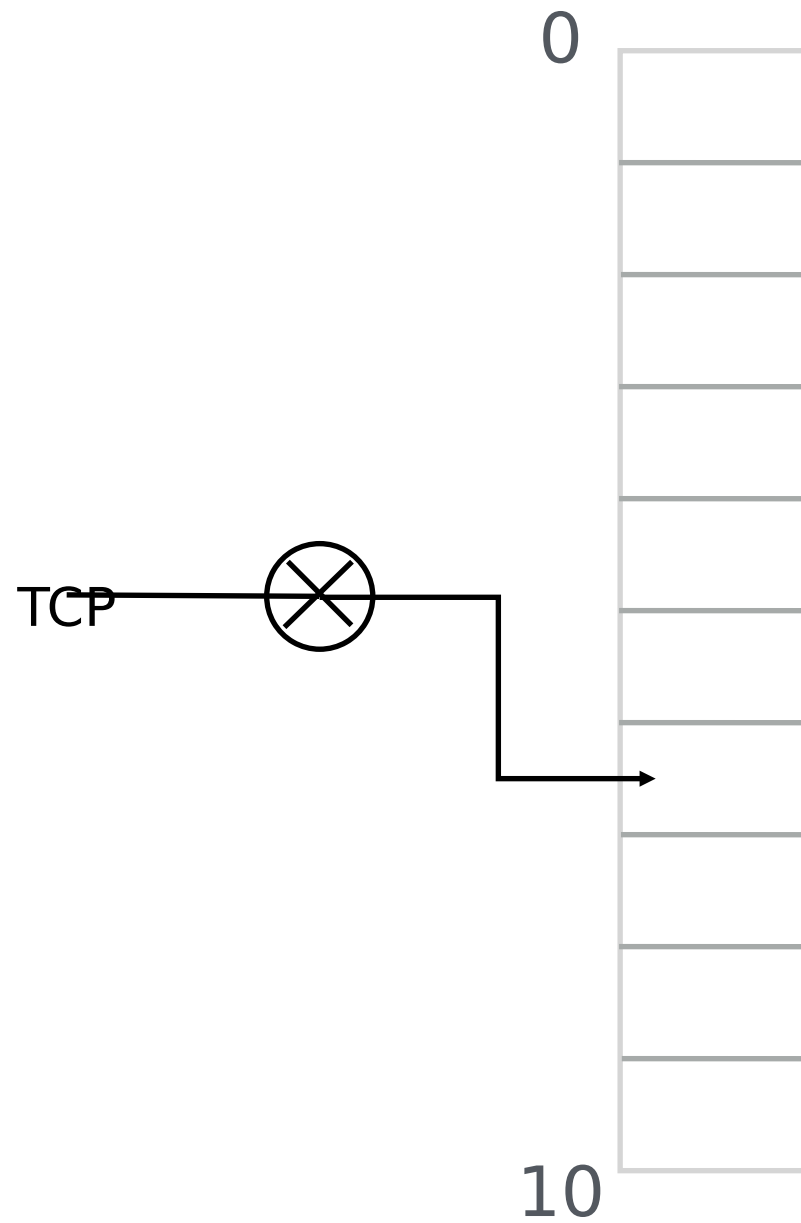
There are two common strategies  
to implement a set

	strategy #1	strategy #2
output	Deterministic	
number of required operations	Probabilistic	

Intuitive implementation of a **set**

# Intuitive implementation of a **set**

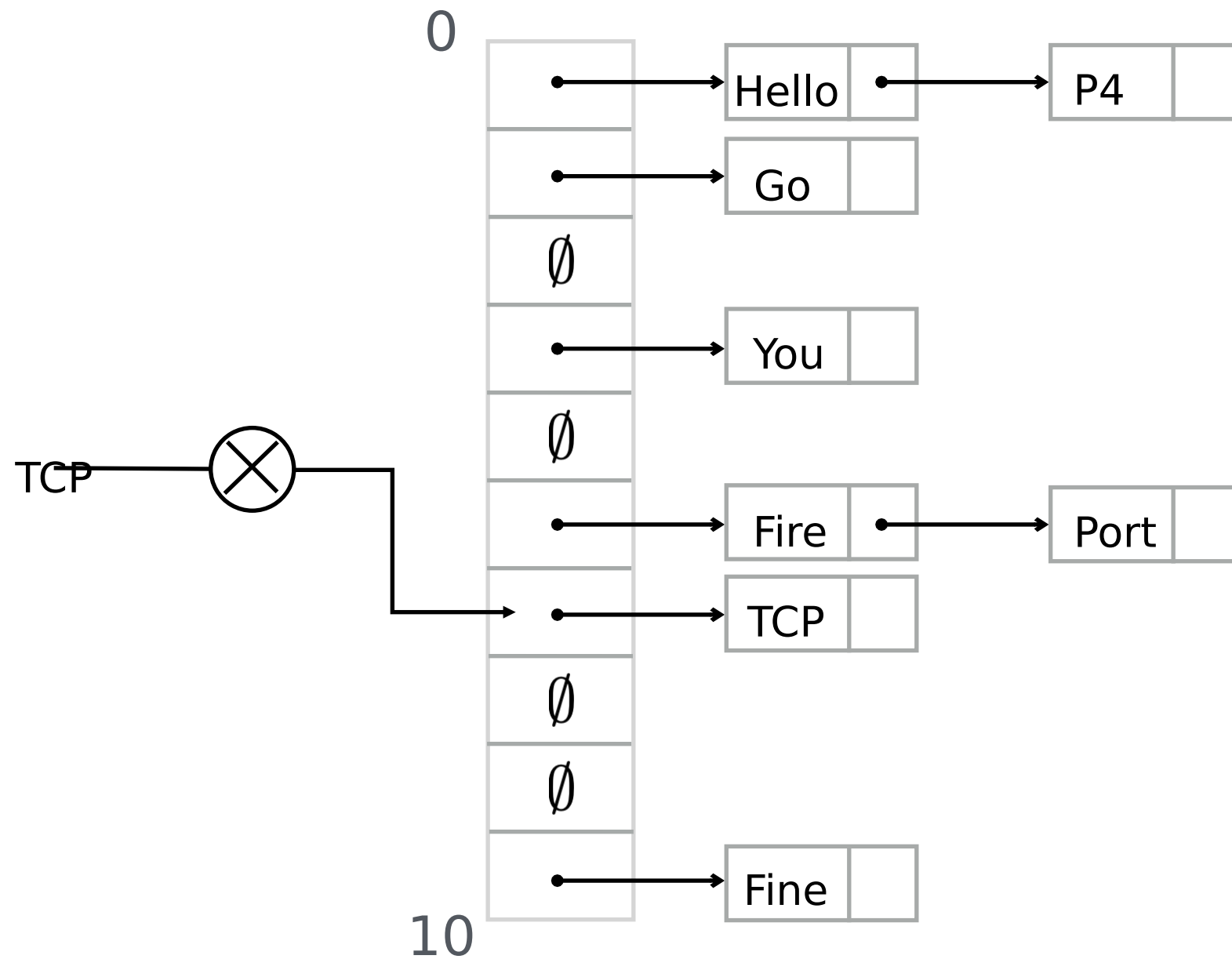
## Separate-chaining





# Intuitive implementation of a **set**

## Separate-chaining



Intuitive implementation of a **set**

Separate-chaining

N elements and M cells

list size

average

$N/M$

worse-case

**N**

Intuitive implementation of a **set**

Separate-chaining

**Pros:** accurate and fast in the average case

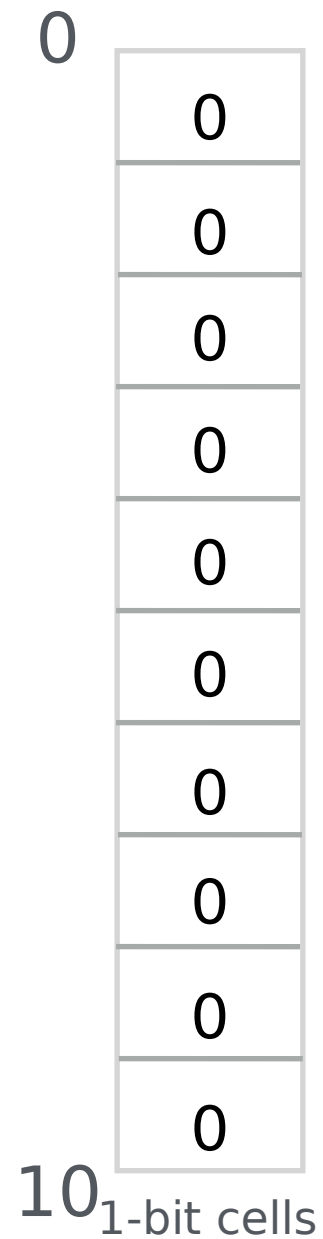
**Con:** only works in hardware if there is a low number of elements (e.g.  $< 100$ )

There are two common strategies  
to implement a set

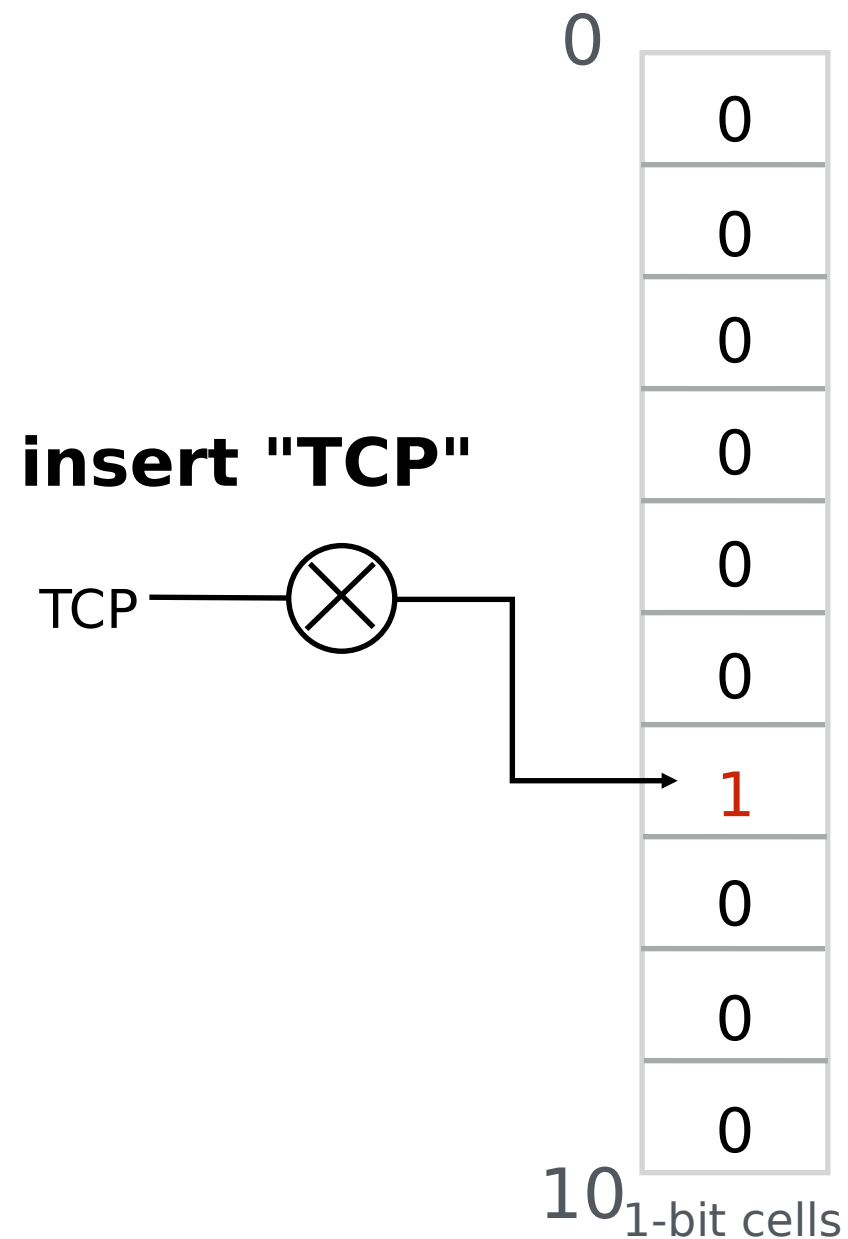
	strategy #1	strategy #2
output	Deterministic	Probabilistic
number of required operations	Probabilistic	Deterministic

*'probabilistic data structures'*

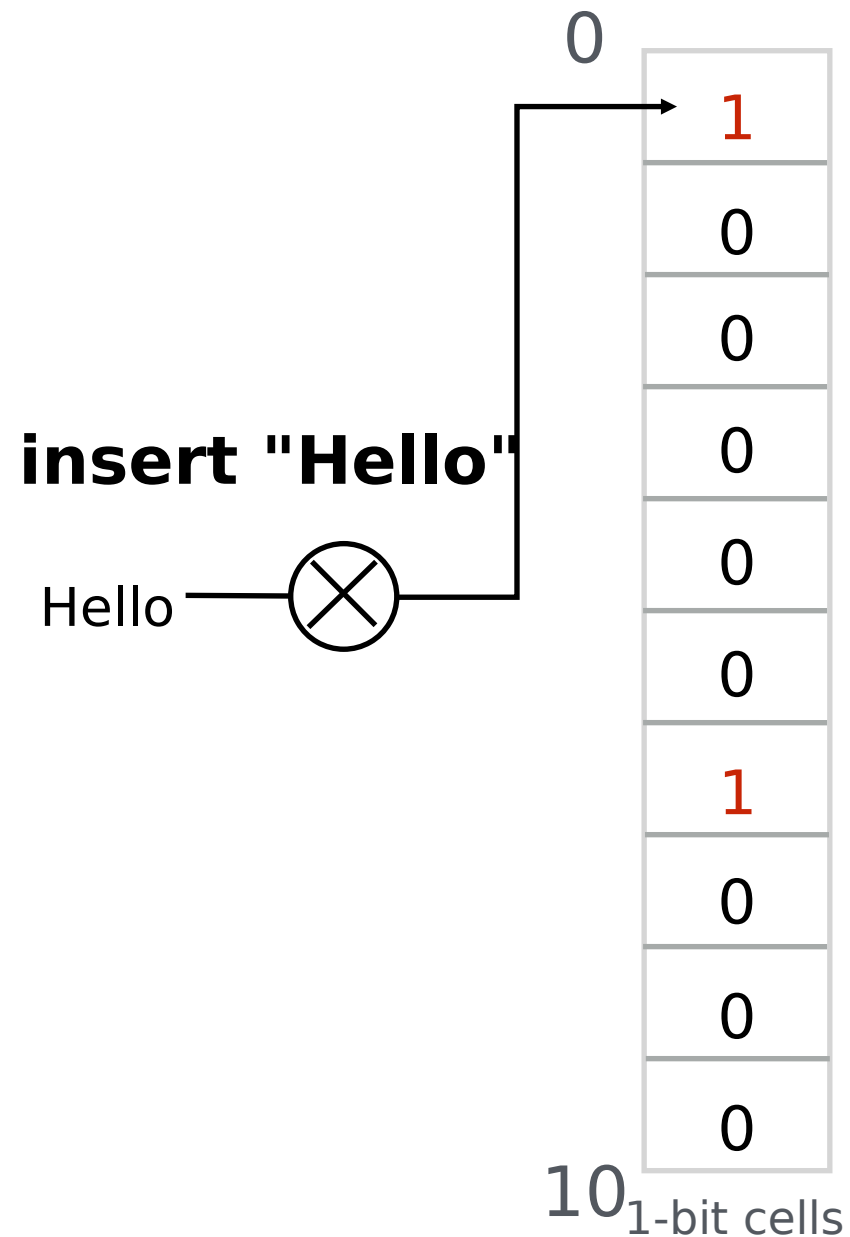
# A simple approach for insertions and membership queries



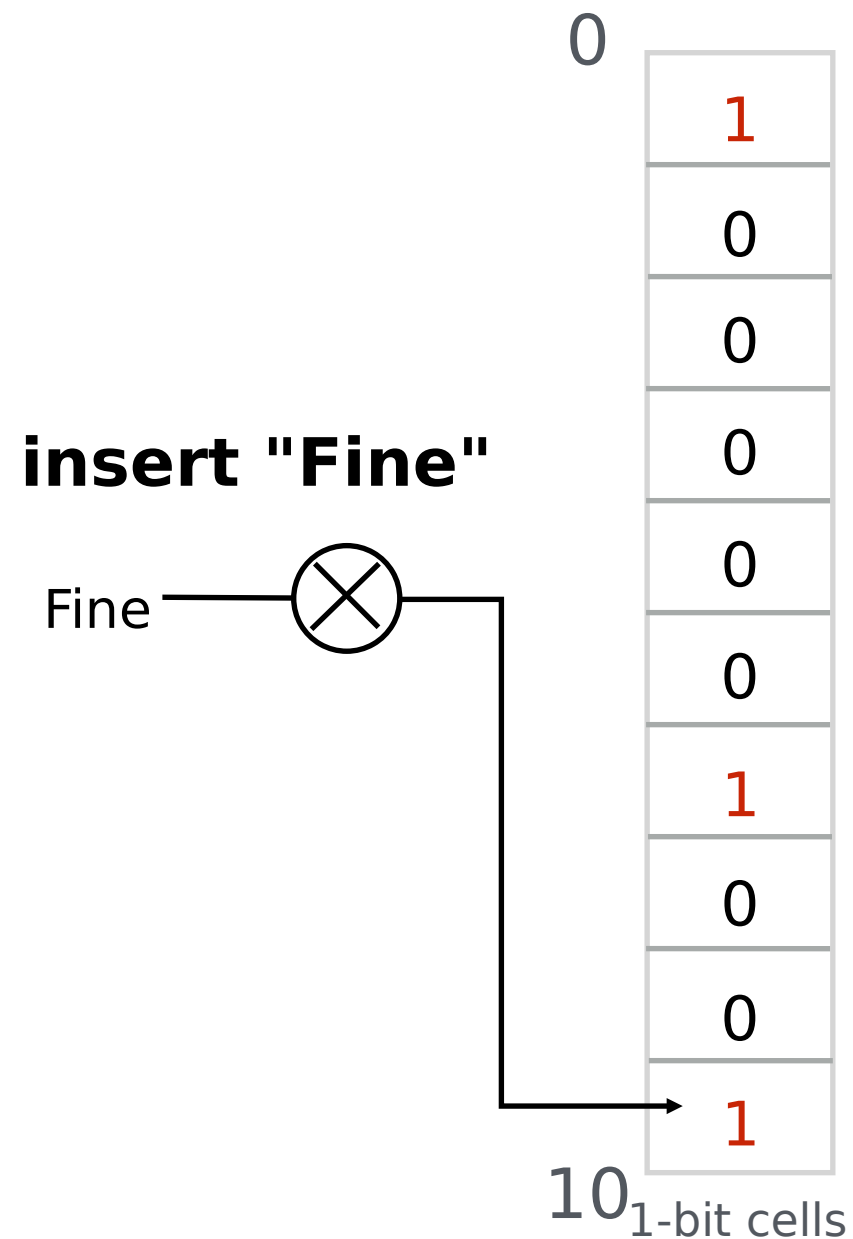
# A simple approach for insertions and membership queries



# A simple approach for **insertions** and **membership queries**

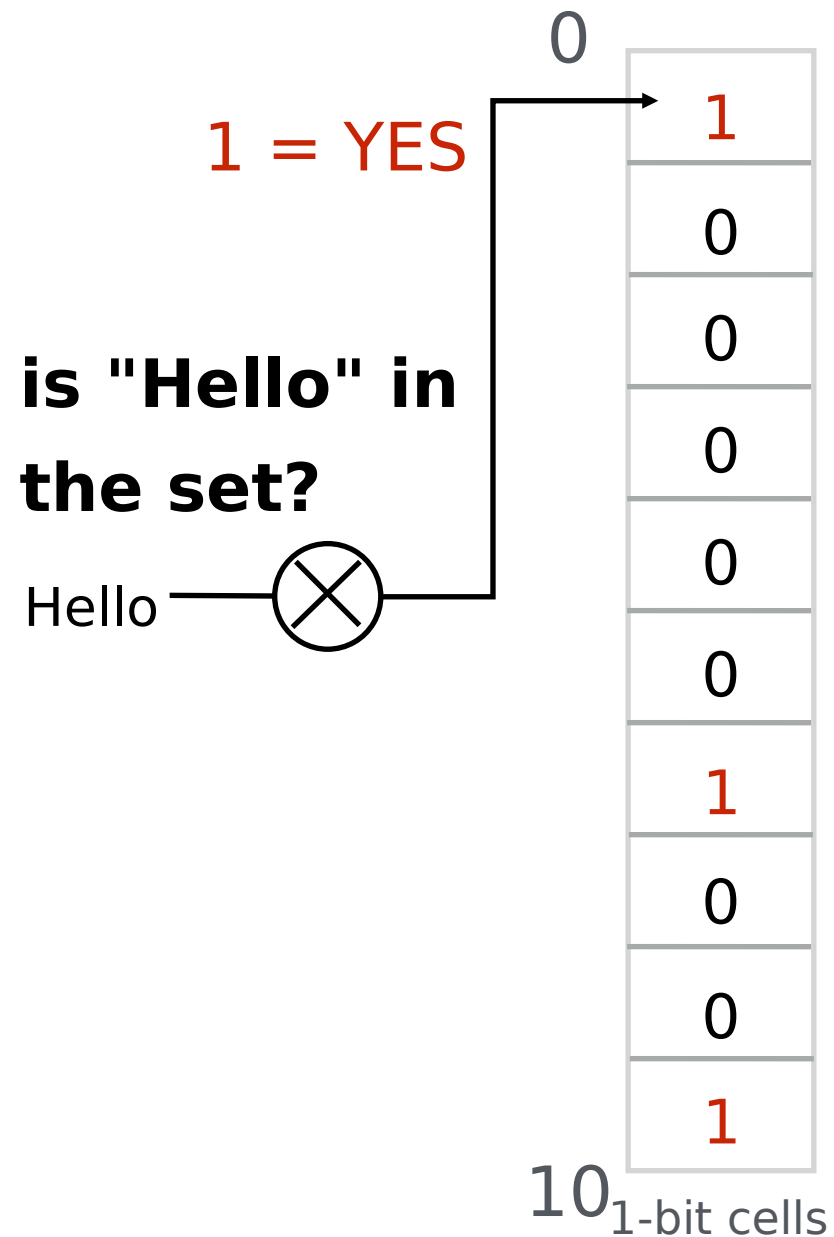


A simple approach for **insertions**  
and **membership queries**

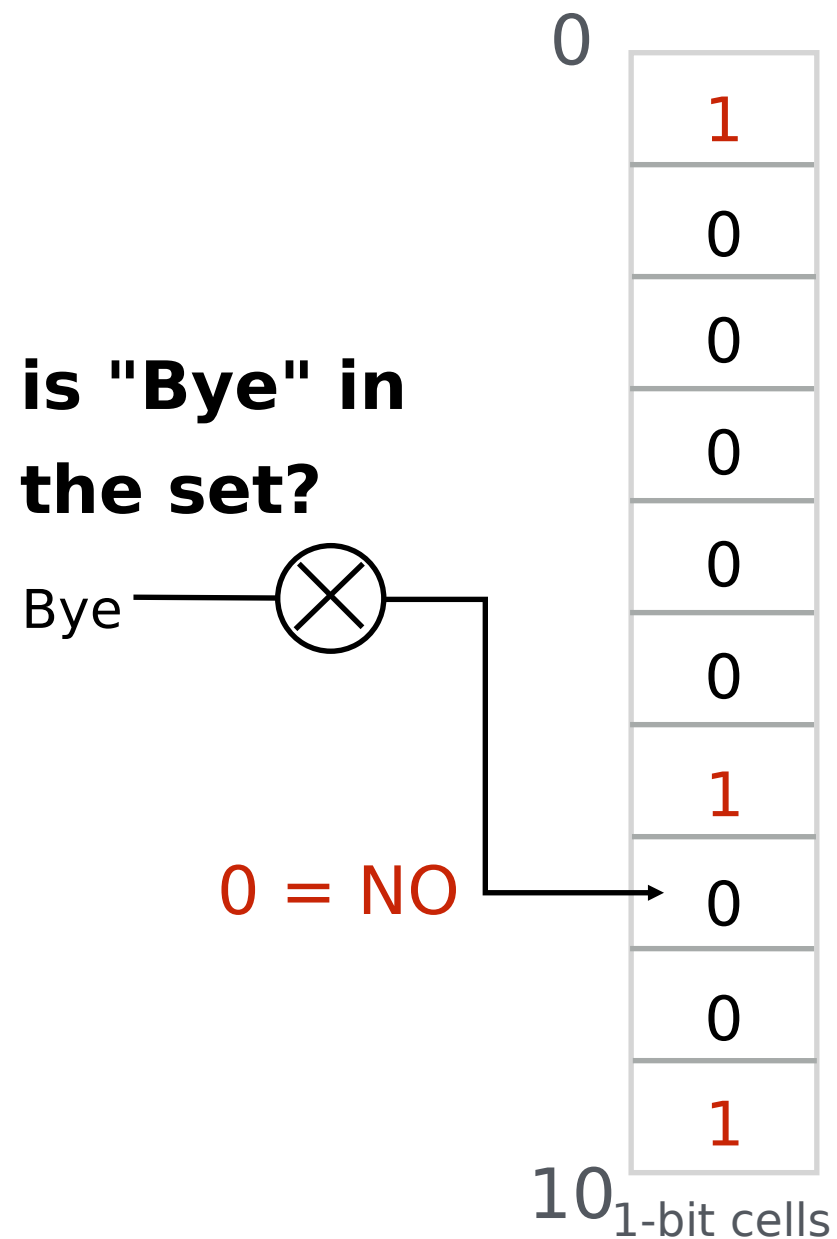




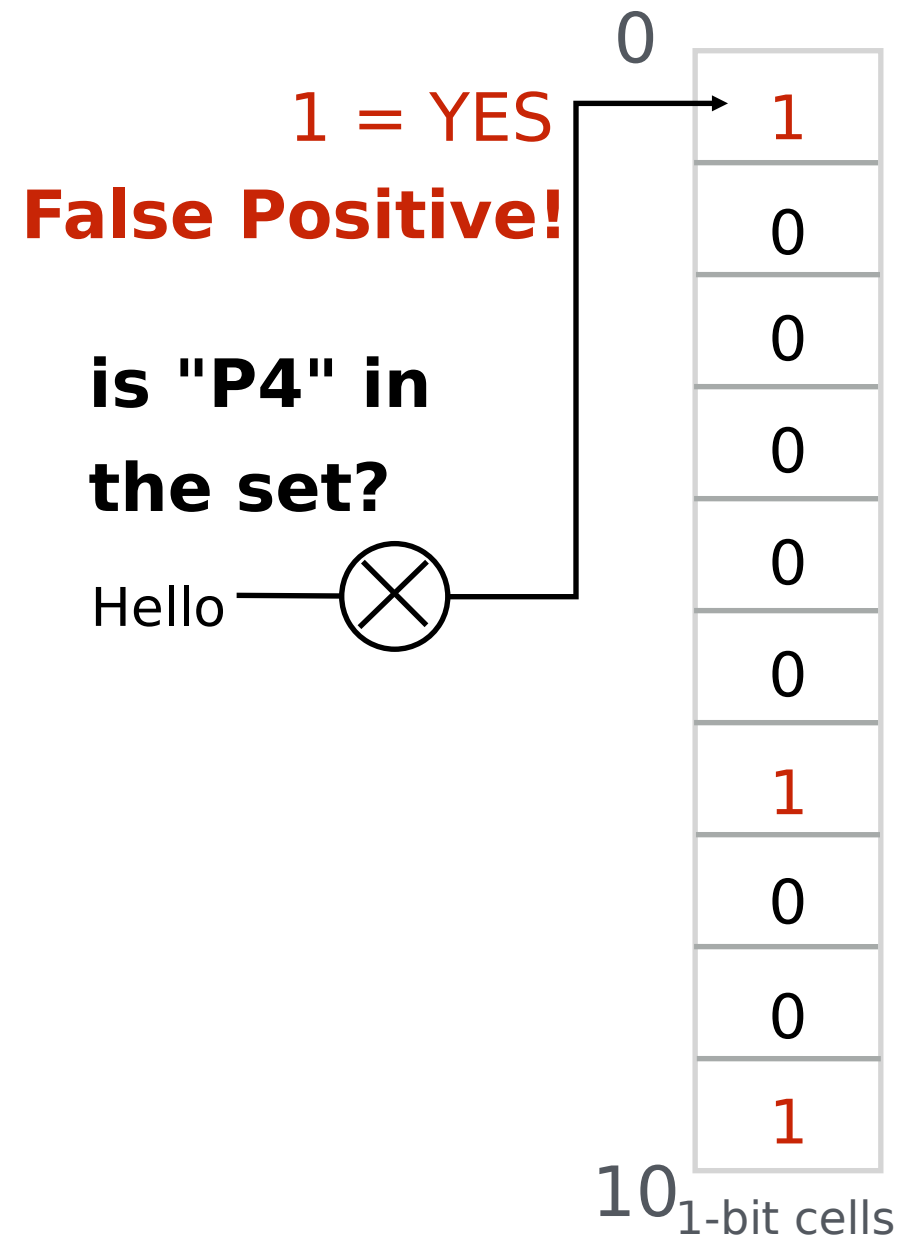
# A simple approach for **insertions** and **membership queries**



# A simple approach for **insertions** and **membership queries**



# A simple approach for **insertions** and **membership queries**



# A simple approach for **insertions** and **membership queries**

N elements and M cells

probability of an element to be  
mapped into a particular cell

$$\frac{1}{M}$$

probability of an element not to  
be mapped into a particular cell

$$1 - \frac{1}{M}$$

probability of a cell to be 0

$$\left(1 - \frac{1}{M}\right)^N$$

false positive rate (FPR)

$$1 - \left(1 - \frac{1}{M}\right)^N$$

false negative rate

$$0$$

# A simple approach for insertions and membership queries

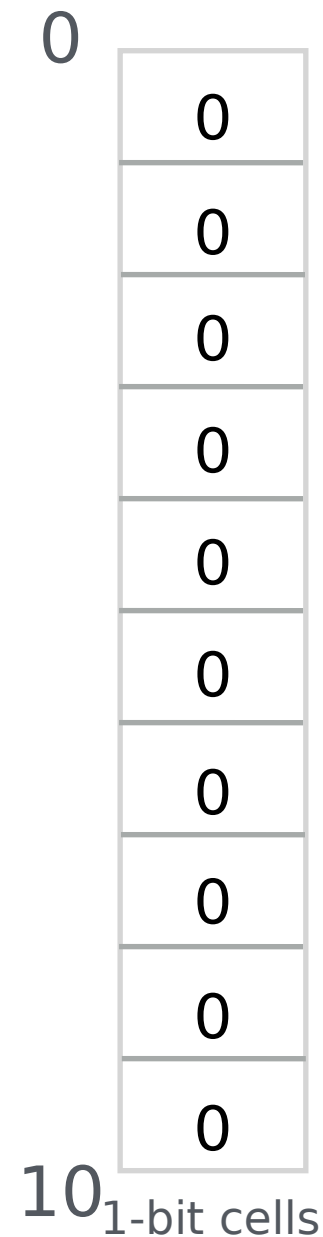
# of elements	# of cells	FPR
1000	10000	9.5%
1000	100000	1%

A simple approach for **insertions**  
and **membership queries**

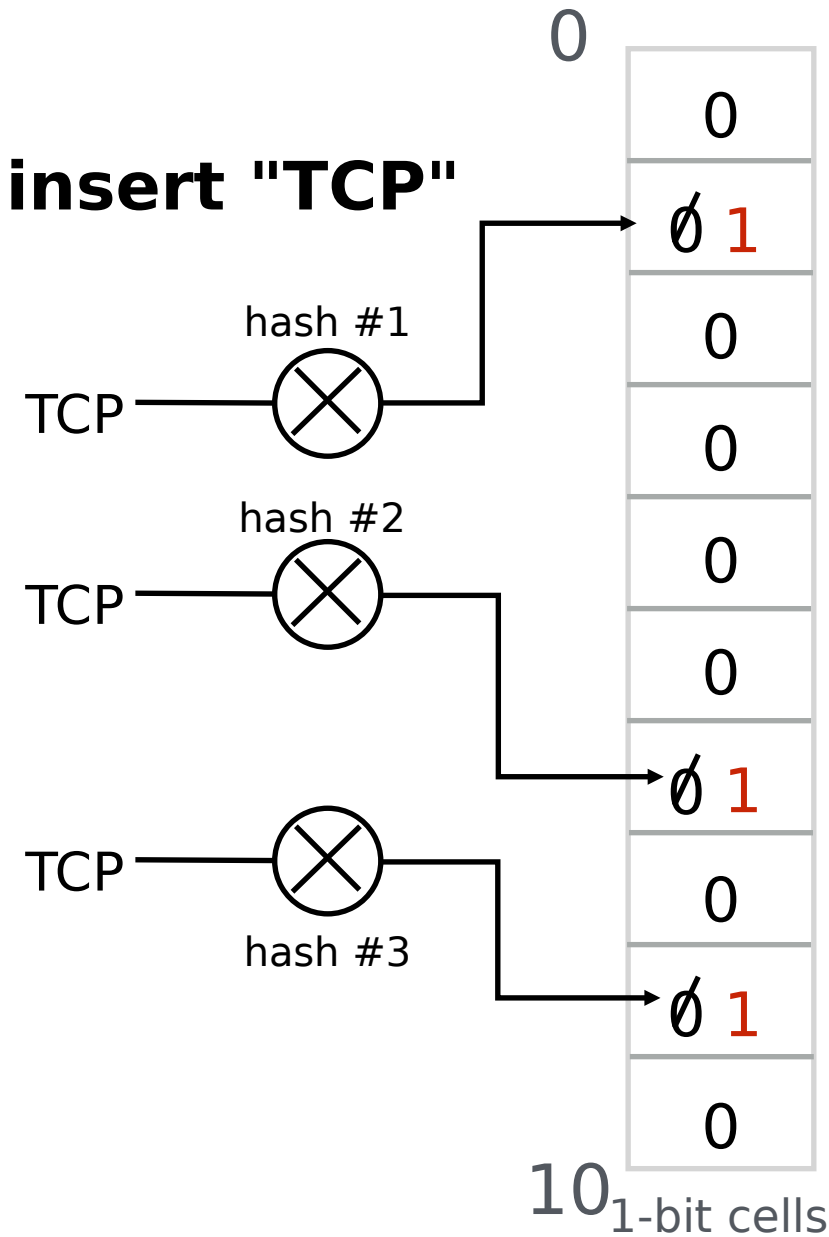
**Pros:** simple and only one operation per  
insertion or query

**Con:** roughly 100x more cells are required than  
the number of element we want to store  
for a 1% false positive rate

# Bloom Filters: a more memory-efficient approach for insertions and membership queries

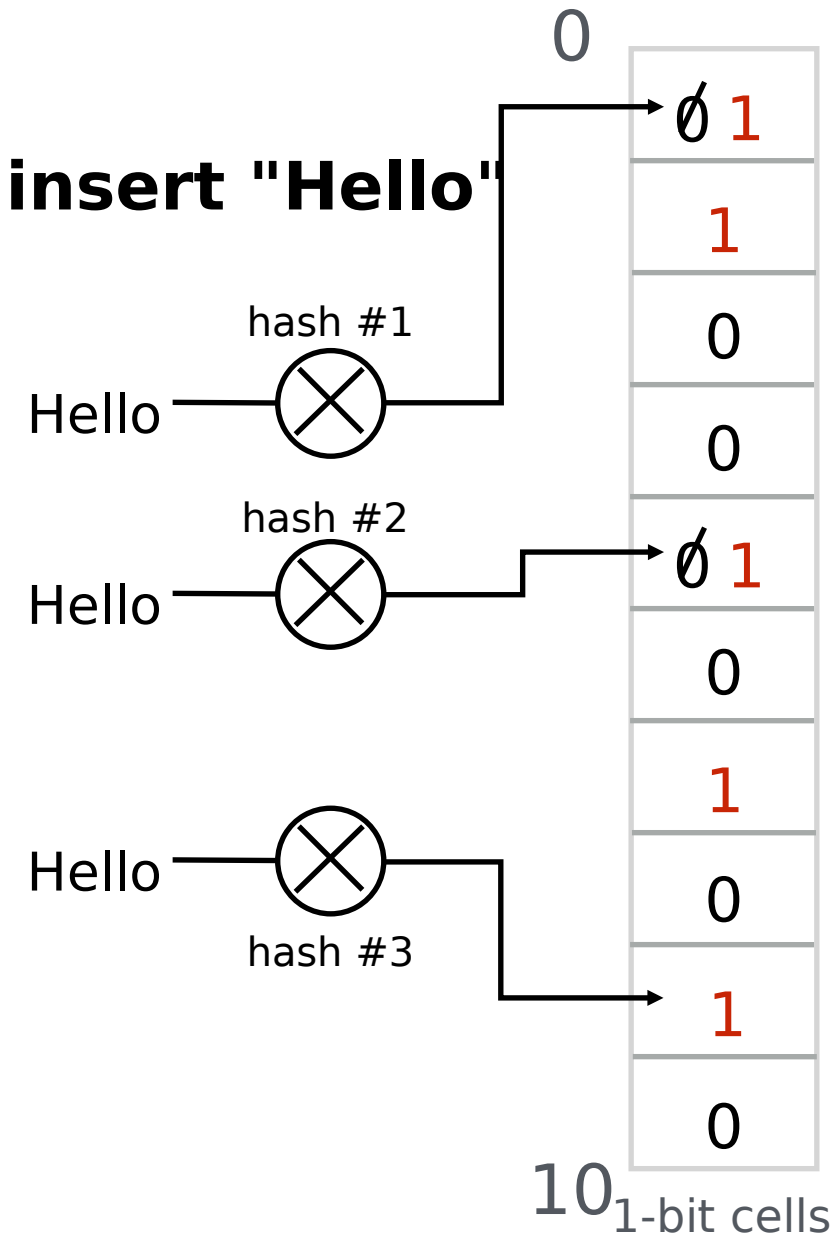


# Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

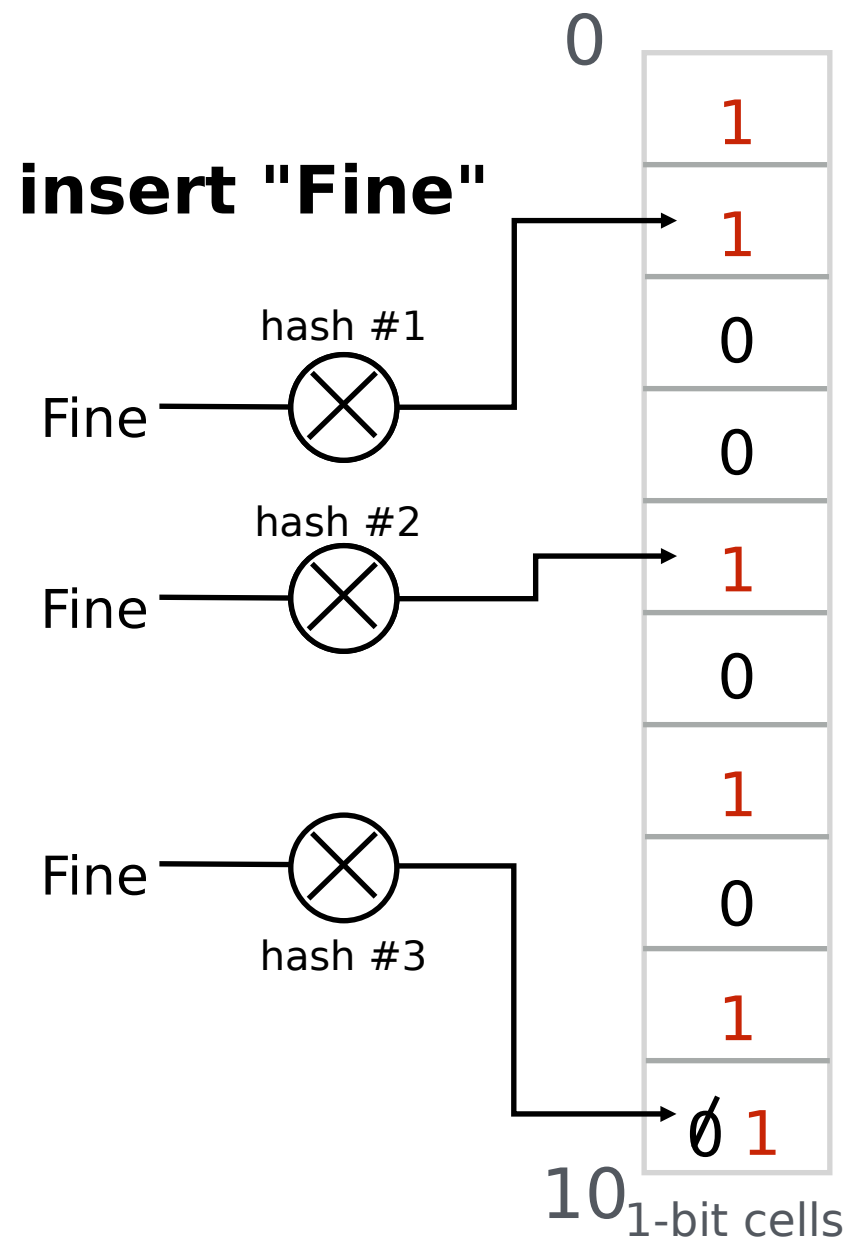




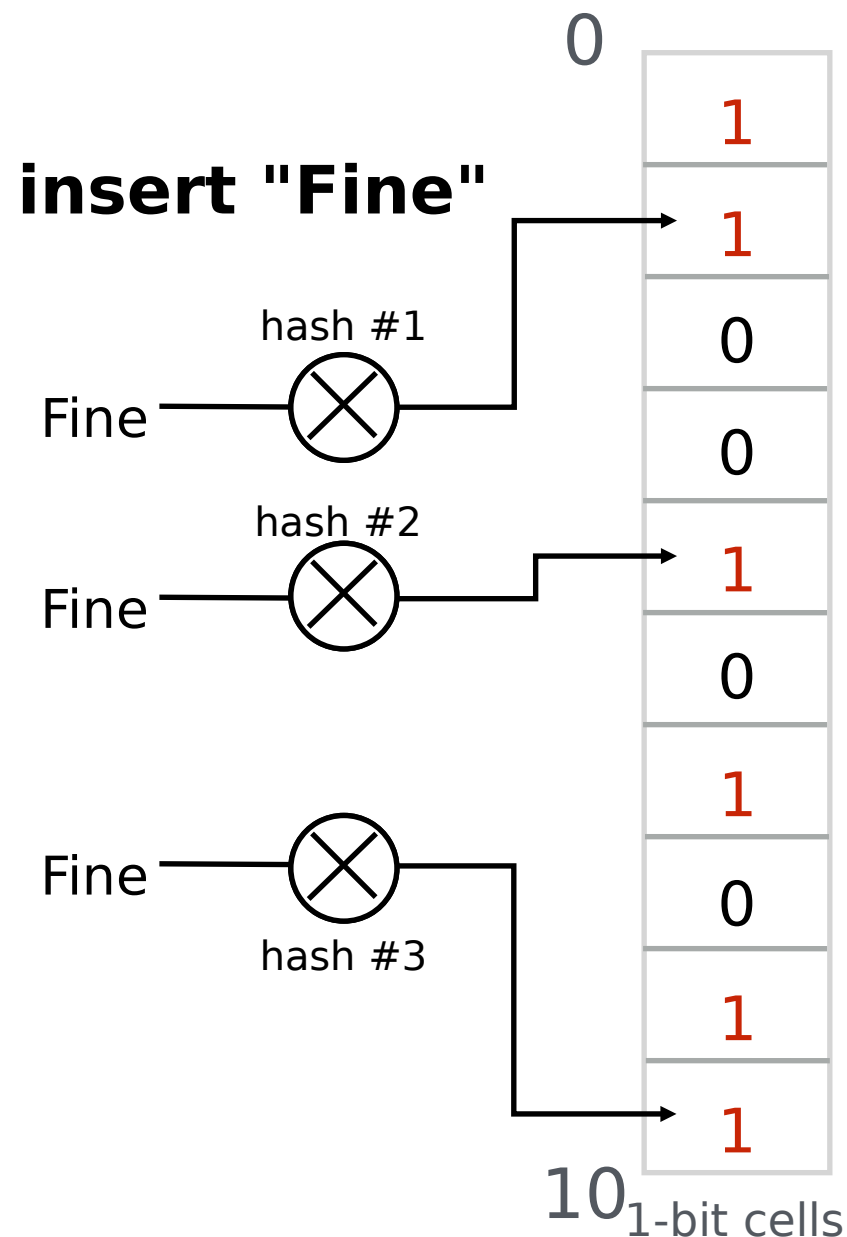
# Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**



# Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**



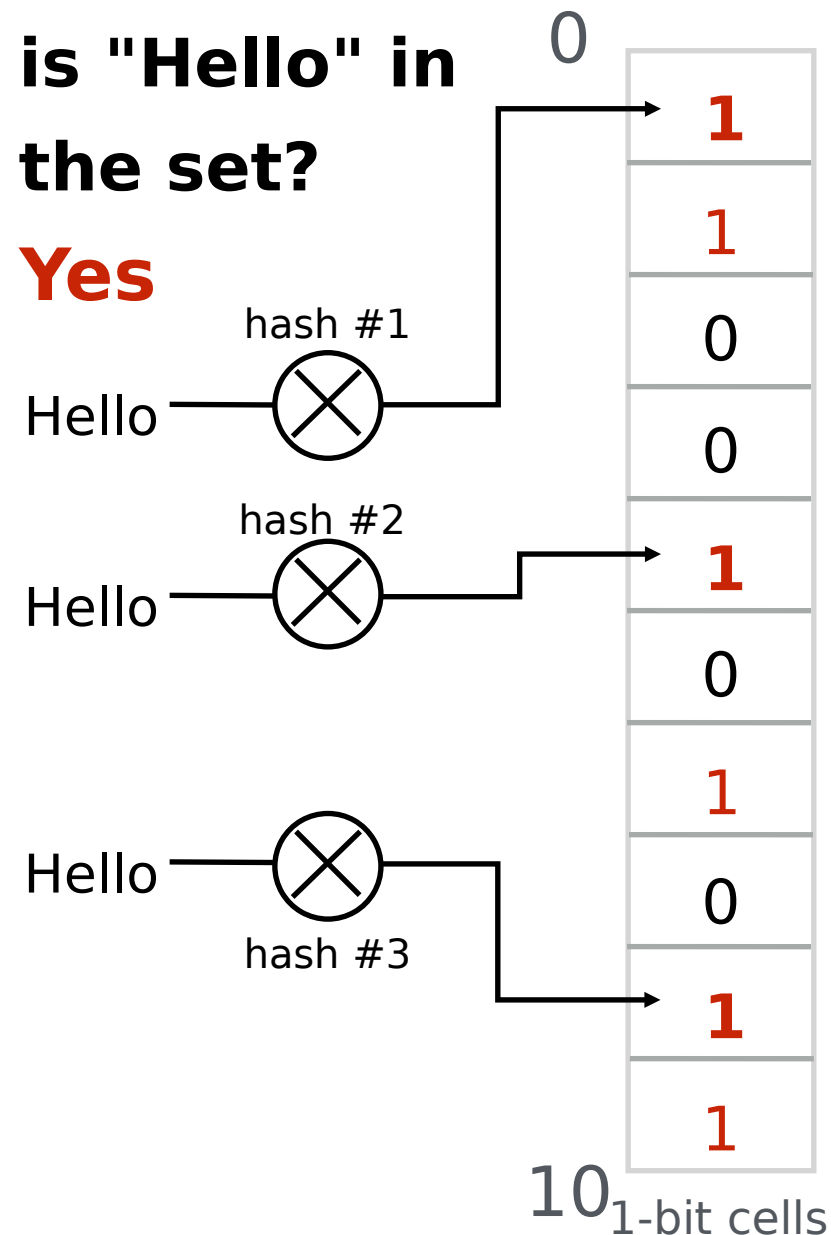
# Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**



An element is considered in the set if **all** the hash values map to a cell with 1

An element is not in the set if **at least** one hash value maps to a cell with 0

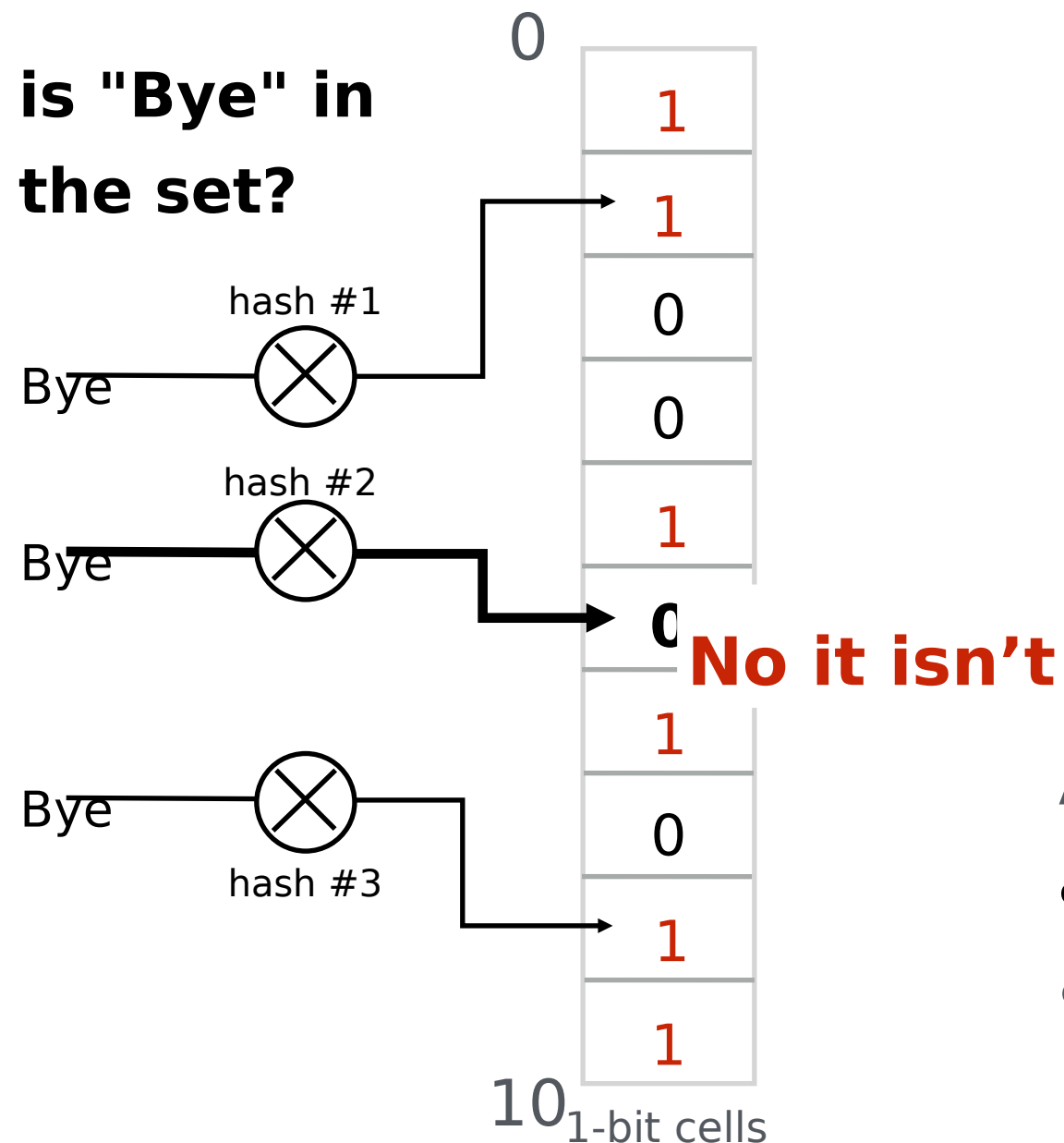
# Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**



An element is considered in the set if **all** the hash values map to a cell with 1

An element is not in the set if **at least** one hash value maps to a cell with 0

# Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

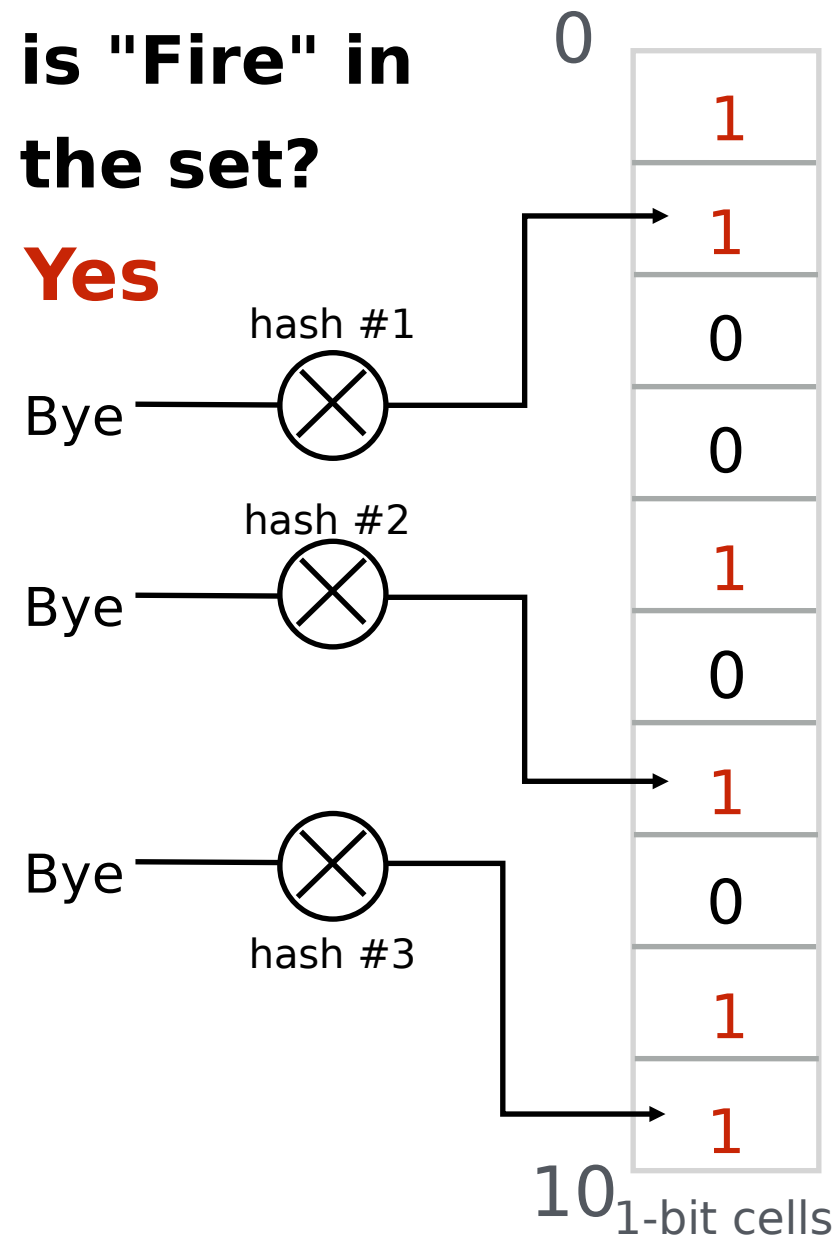


An element is considered in the set if **all** the hash values map to a cell with 1

An element is not in the set if **at least** one hash value maps to a cell with 0

# Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

## False Positive!



An element is considered in the set if **all** the hash values map to a cell with 1

An element is not in the set if **at least** one hash value maps to a cell with 0

# Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

N elements, M cells and K *independent* hash functions

probability that one hash function  
returns the index of a particular cell  $\frac{1}{M}$

probability that one hash function does  
not return the index of a particular cell  $1 - \frac{1}{M}$

probability of a cell to be 0  $(1 - \frac{1}{M})^{KN}$

false positive rate  $(1 - (1 - \frac{1}{M})^{KN})^K$

false negative rate 0

Bloom Filters: a more memory-efficient approach  
for **insertions** and **membership queries**

# of elements	# of cells	# hash functions	FPR
1000	10000	7	0.82%
1000	100000	7	$\approx 0\%$

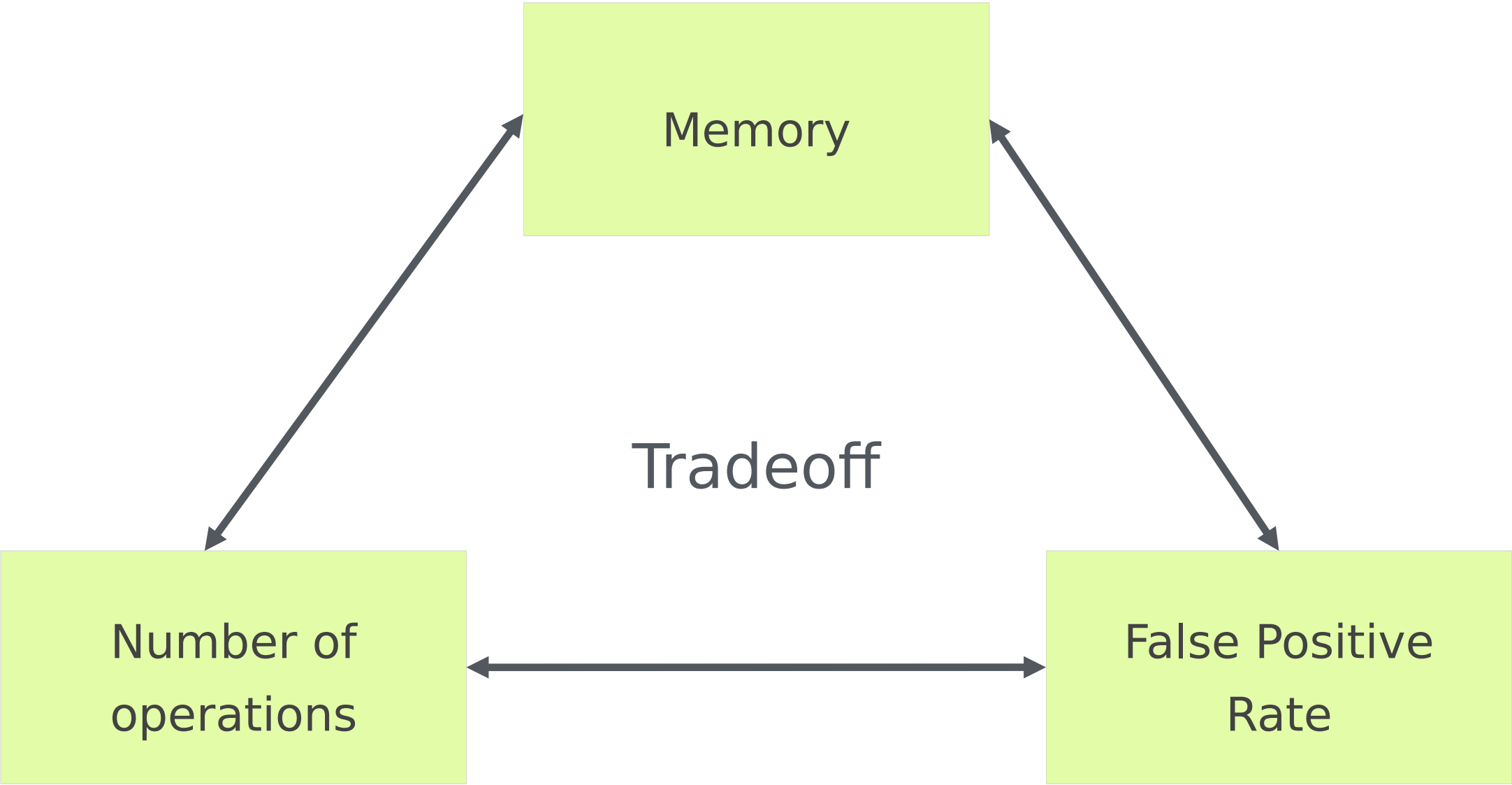


Bloom Filters: a more memory-efficient approach for **insertions** and **membership queries**

**Pro:** consumes roughly 10x less memory than the simple approach

**Con:** Requires slightly more operations than the simple approach (7 hashes instead of just 1)

# Dimension your Bloom Filter



Dimension your Bloom Filter

$N$  elements

$M$  cells

$K$  hash functions

FP false positive rate

Dimension your Bloom Filter

N elements

M cells

K hash functions

FP false positive rate

asymptotic approx.

$$FP = \left(1 - \left(1 - \frac{1}{M}\right)^{KN}\right)^K \approx \left(1 - e^{-KN/M}\right)^K$$

with calculus you can  
dimension your bloom filter

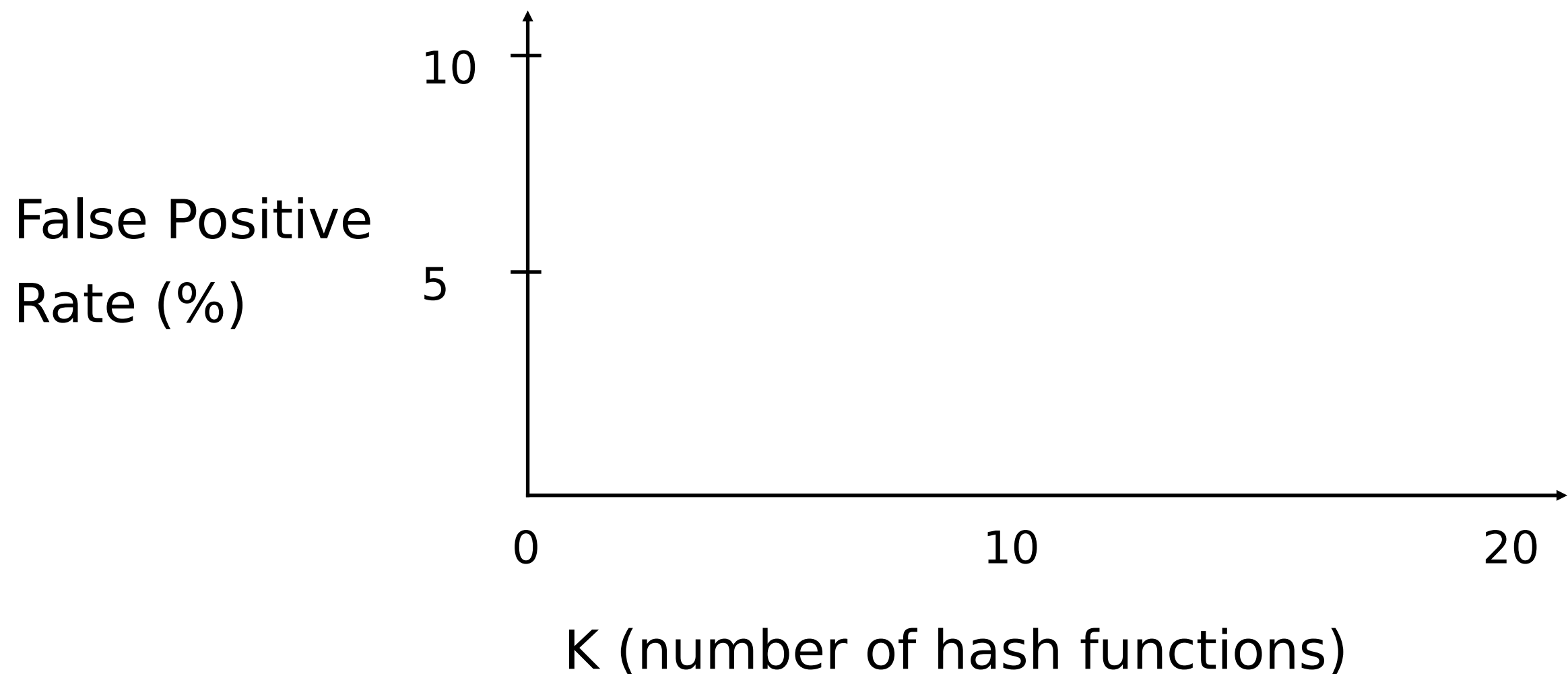
# Dimension your Bloom Filter

N elements

M cells

K hash functions

FP false positive rate



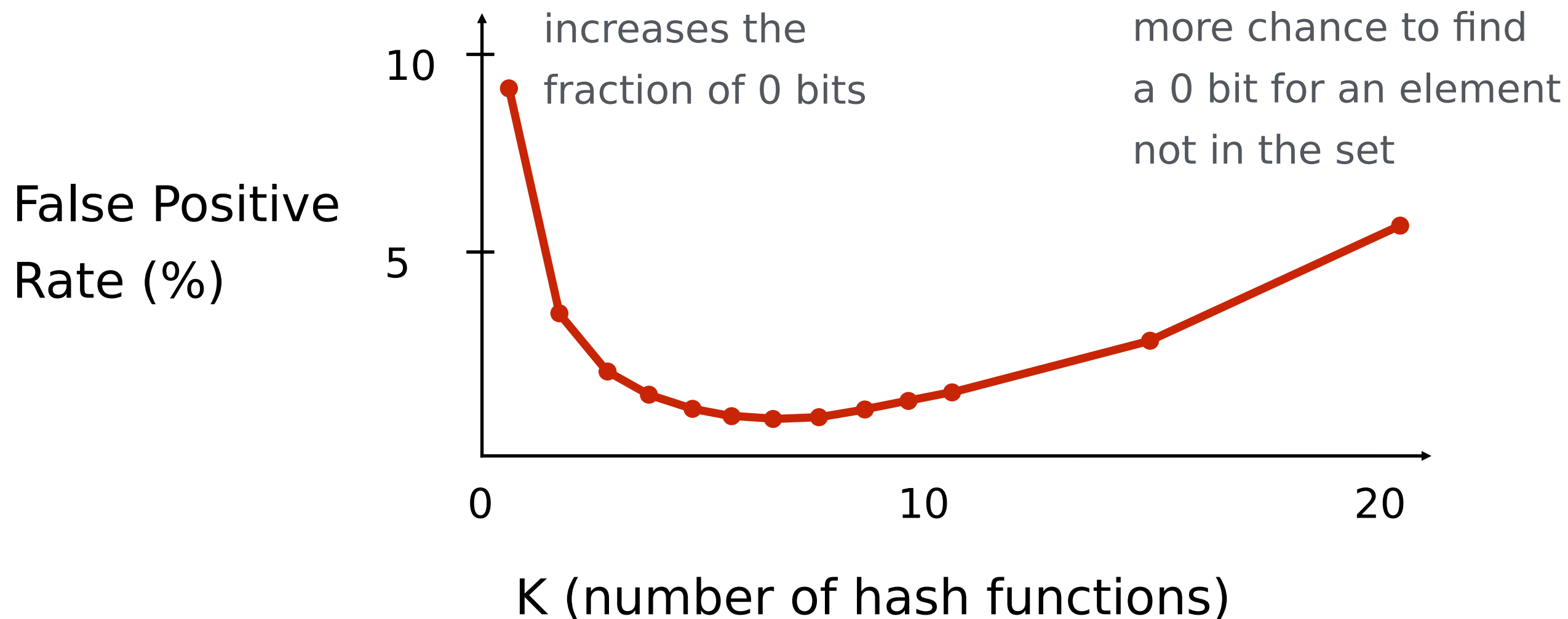
# Dimension your Bloom Filter

N elements

M cells

K hash functions

FP false positive rate



# Dimension your Bloom Filter

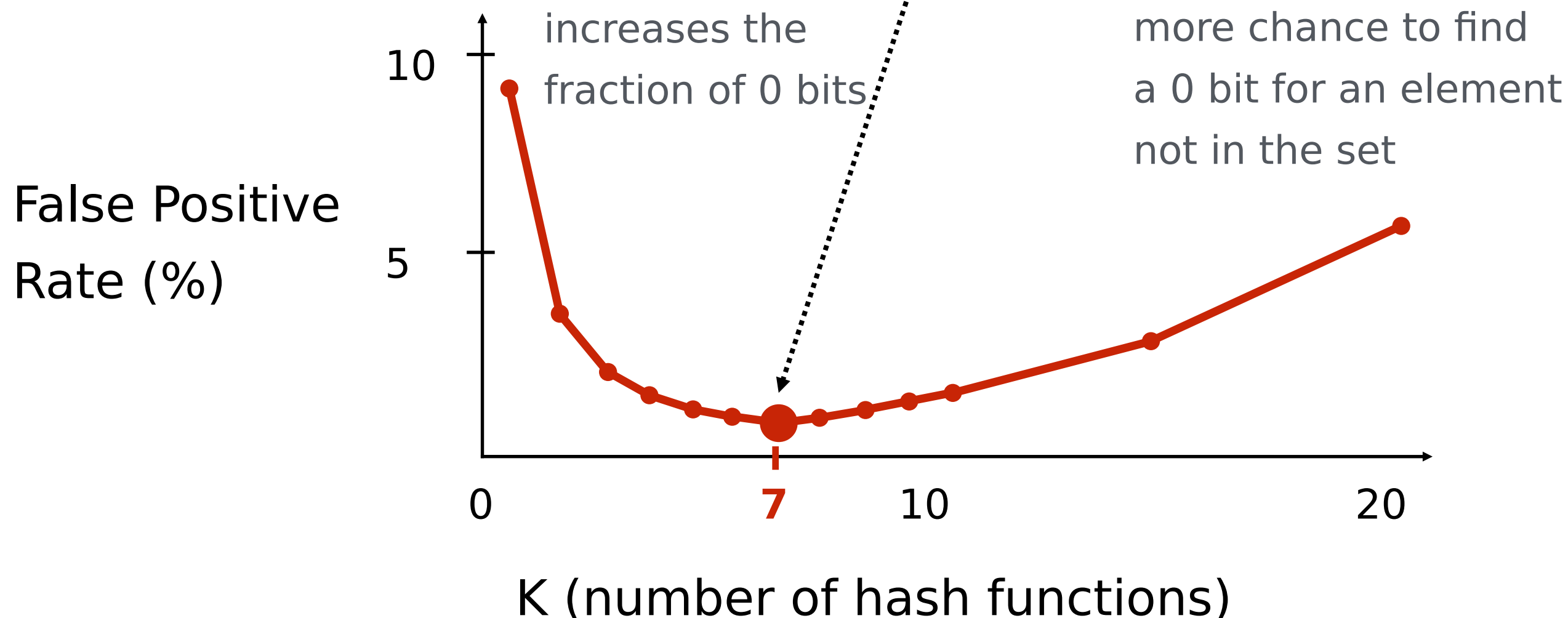
N elements

M cells

K hash functions

FP false positive rate

there is always a  
global minimum when  
 $K = \ln 2 * (M/N)$  found  
by taking the derivative  
of  $\approx (1 - e^{-KN/M})^K$



# Implementation of a Bloom Filter in P4<sub>16</sub>

You will have to use hash functions

v1model

```
enum HashAlgorithm {  
    crc32,  
    crc32_custom,  
    crc16,  
    s,  
    random,  
    identity,  
    csum16,  
    xor16  
}
```

```
extern void hash<O, T, D, M>(out O result,  
    in HashAlgorithm algo, in T base, in D data, in M  
    max);
```

more info

<https://github.com/p4lang/p4c/blob/master/p4include/v1model.p4>



# Implementation of a Bloom Filter in P4<sub>16</sub>

You will have to use hash functions, as well as registers

v1model

```
extern register<T> {  
  
    register(bit<32> size);  
  
    void read(out T result, in bit<32> index);  
    void write(in bit<32> index, in T value);  
}
```

# Implementation of a Bloom Filter in P4<sub>16</sub> with 2 hash functions

```
control MyIngress(...) {  
    register<bit<1>>(NB_CELLS) bloom_filter;
```

# Implementation of a Bloom Filter in P4<sub>16</sub> with 2 hash functions

```
control MyIngress(...) {  
  register<bit<1>>(NB_CELLS) bloom_filter;  
  apply {  
    hash(meta.index1, HashAlgorithm.my_hash1, 0,  
          {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);  
    hash(meta.index2, HashAlgorithm.my_hash2, 0,  
          {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);  
  }  
}
```

# Implementation of a Bloom Filter in P4<sub>16</sub> with 2 hash functions

```
control MyIngress(...) {
  register register<bit<1>>(NB_CELLS) bloom_filter;
  apply {
    hash(meta.index1, HashAlgorithm.my_hash1, 0,
         {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);
    hash(meta.index2, HashAlgorithm.my_hash2, 0,
         {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);

    if (meta.to_insert == 1) {
      bloom_filter.write(meta.index1, 1);
      bloom_filter.write(meta.index2, 1);
    }
    if (meta.to_query == 1) {
      bloom_filter.read(meta.query1, meta.index1);
      bloom_filter.read(meta.query2, meta.index2);

      if (meta.query1 == 0 || meta.query2 == 0) {
        meta.is_stored = 0;
      }
      else {
        meta.is_stored = 1;
      }
    }
  }
}
```

# Implementation of a Bloom Filter in P4<sub>16</sub>

## with 2 hash functions

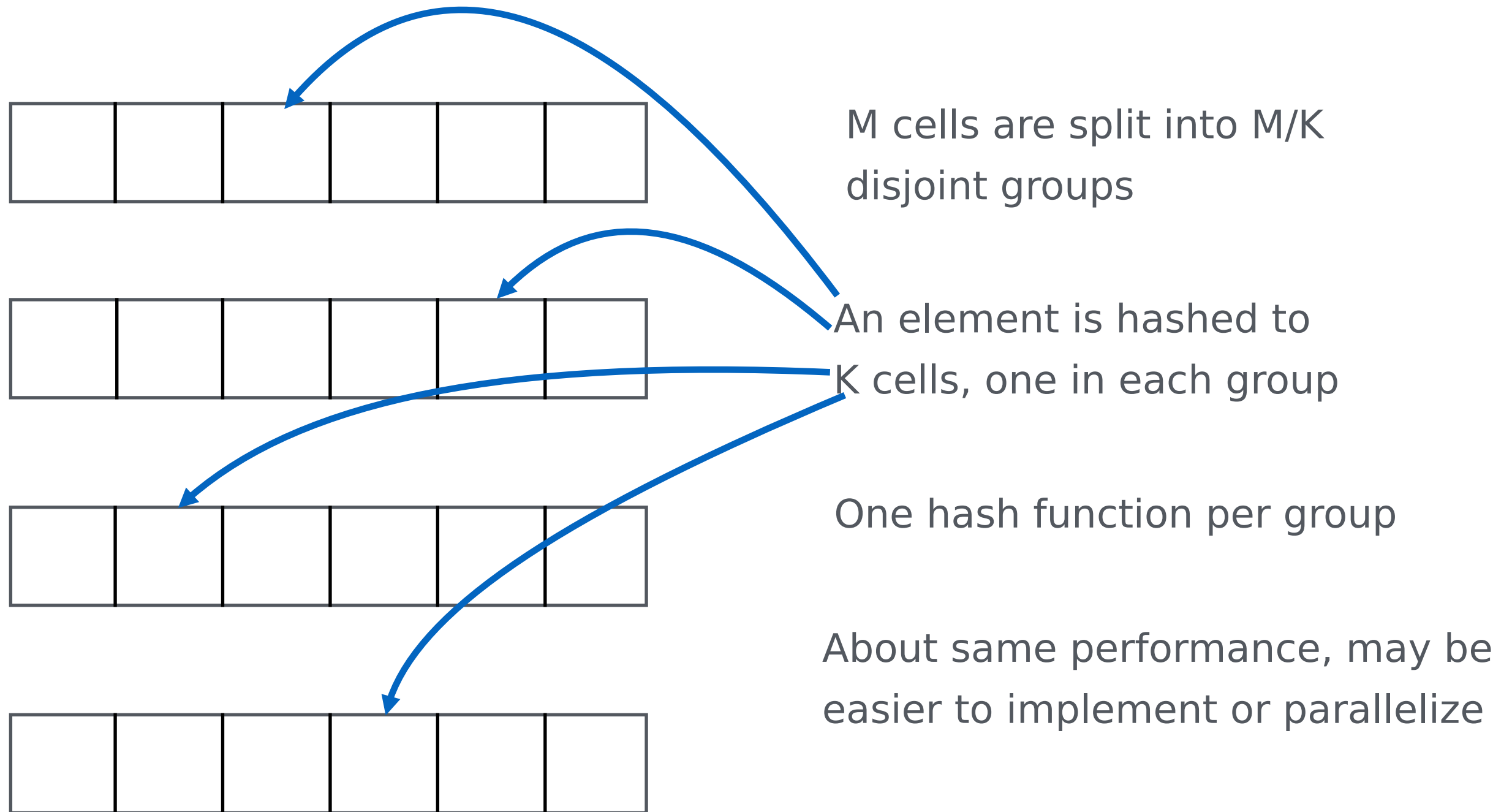
```
control MyIngress(...) {
  register register<bit<1>>(NB_CELLS) bloom_filter;
  apply {
    hash(meta.index1, HashAlgorithm.my_hash1, 0,
         {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);
    hash(meta.index2, HashAlgorithm.my_hash2, 0,
         {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);

    if (meta.to_insert == 1) {
      bloom_filter.write(meta.index1, 1);
      bloom_filter.write(meta.index2, 1);
    }
    if (meta.to_query == 1) {
      bloom_filter.read(meta.query1, meta.index1);
      bloom_filter.read(meta.query2, meta.index2);

      if (meta.query1 == 0 || meta.query2 == 0) {
        meta.is_stored = 0;
      }
      else {
        meta.is_stored = 1;
      }
    }
  }
}
```

Everything in bold red must be adapted for your program

Depending on the hardware limitations,  
splitting the bloom filter might be required



Because deletions are not possible, the controller may need to regularly **reset** the bloom filters

Resetting a bloom filter takes some time during which it is not usable

Common trick: use two bloom filters and use one when the controller resets the other one

Bloom filters may be extended to allow deletions and to list the filter content.

If you are curious, check out the extended slides for:

- counting Bloom filters (allow deletions)
- invertible Bloom filters (allow to list content)

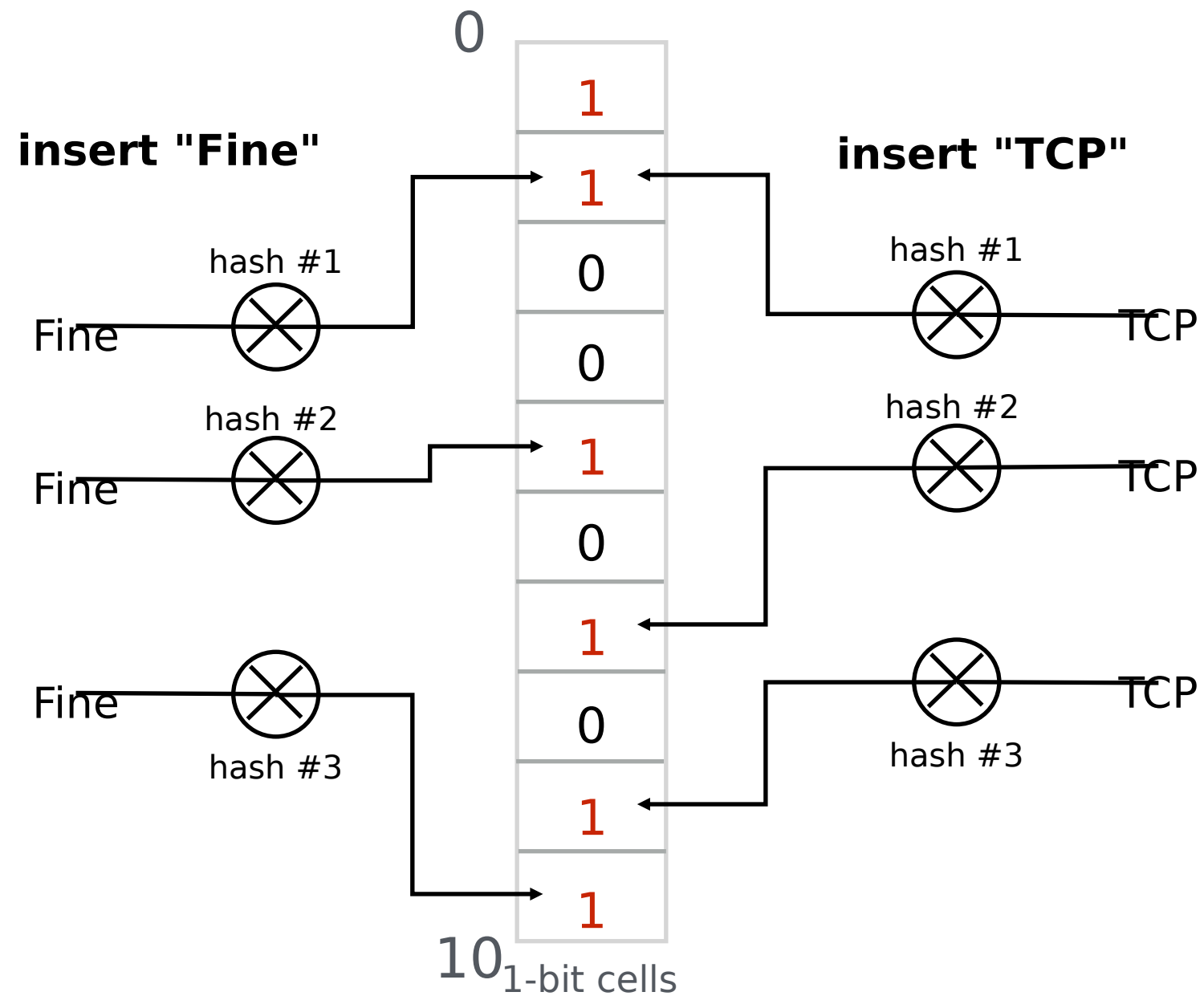


So far we have seen how to do insertions and membership queries

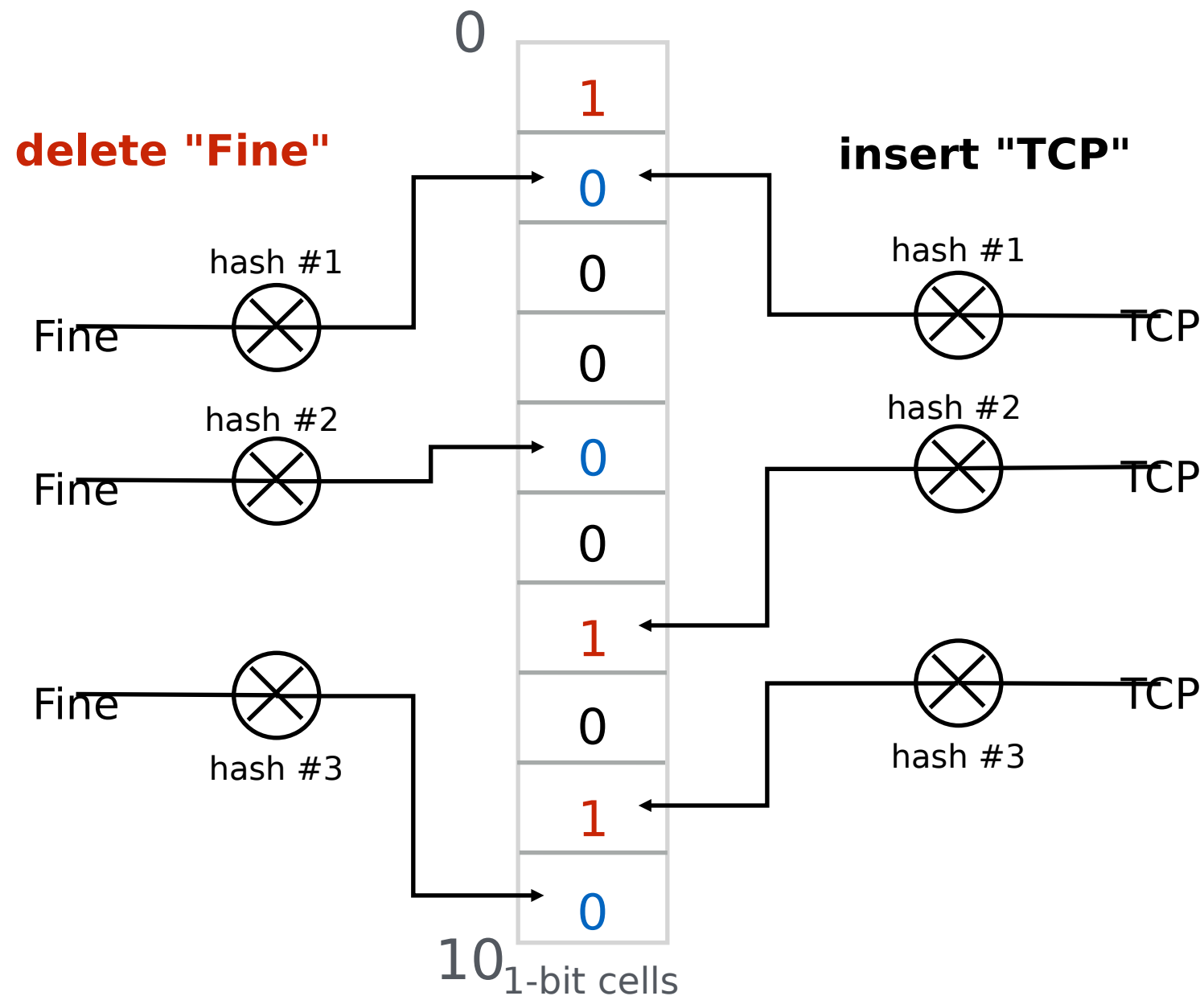
	strategy #1	strategy #2
output	Deterministic	Probabilistic
number of required operations	Probabilistic	Deterministic

Bloom Filters

However Bloom Filters do not handle deletions



However Bloom Filters do not handle **deletions**

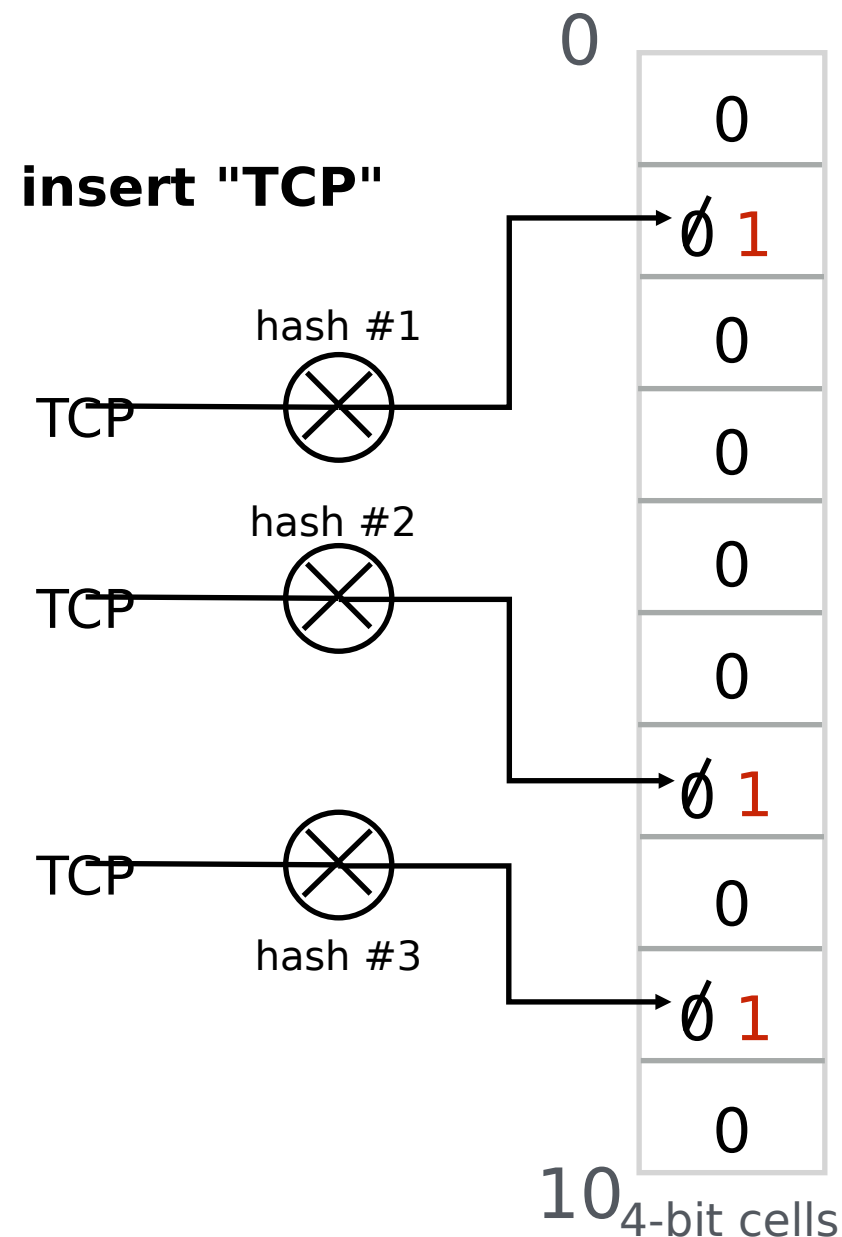


If deleting an element means resetting 1s to 0s, then deleting "Fine" also deletes "TCP"

But we can easily extend them to handle deletions

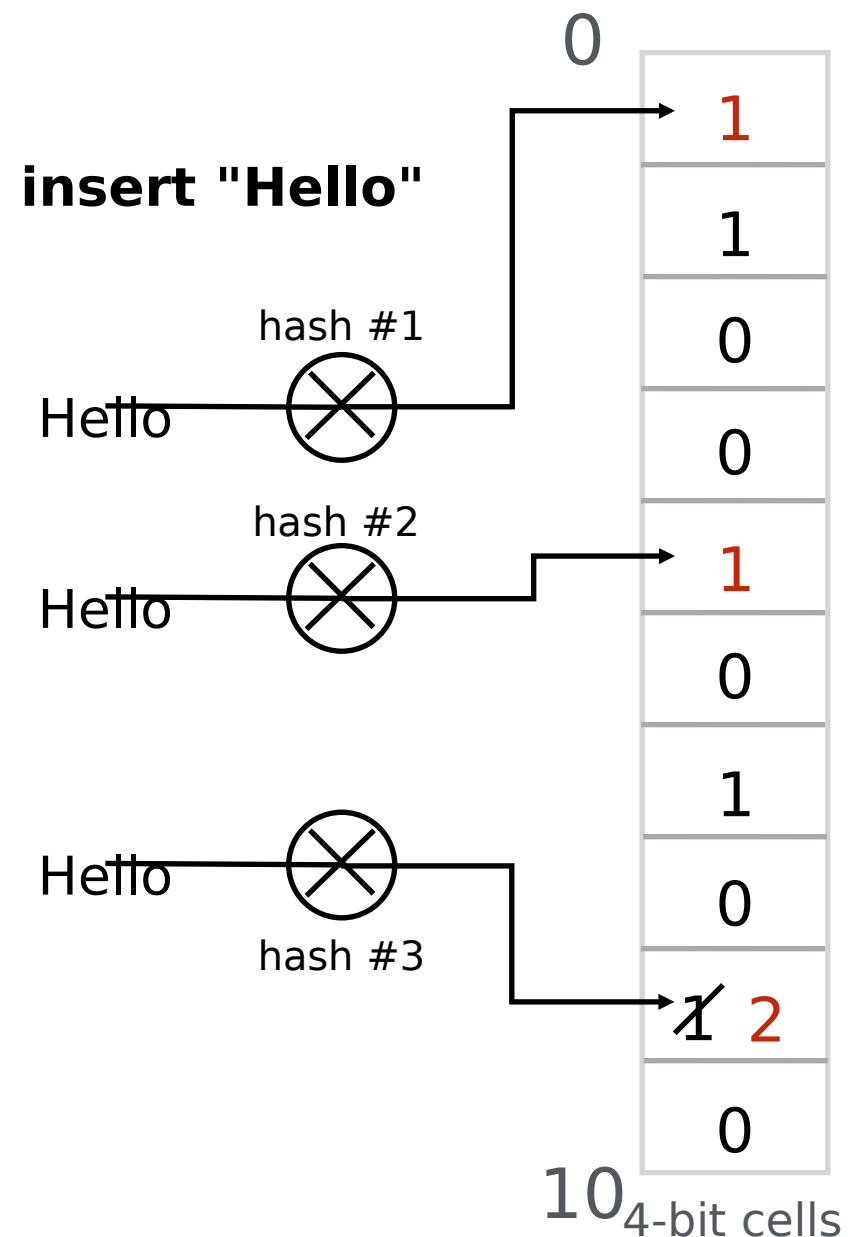
This extended version is called a **Counting Bloom Filter**

But we can easily extend them to handle deletions  
This extended version is called a **Counting Bloom Filter**



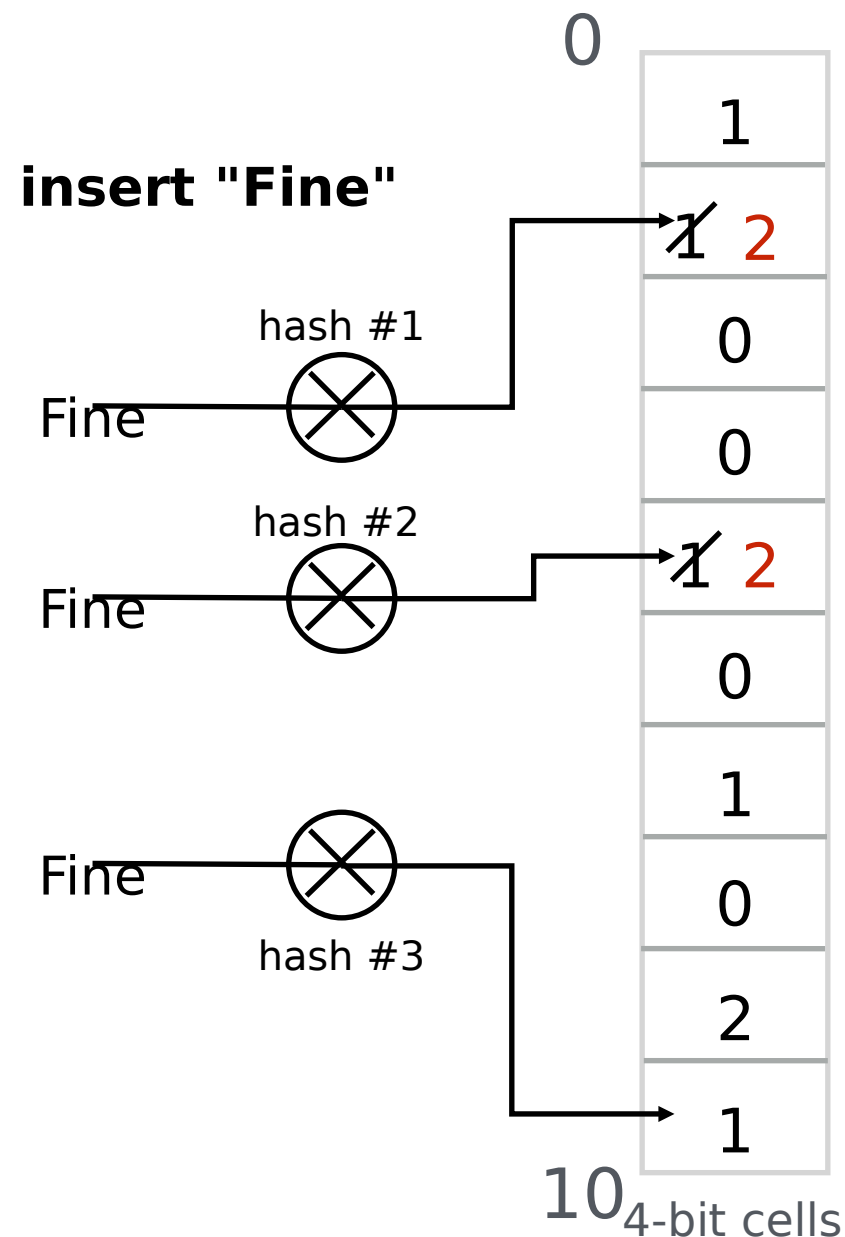
To add an element, increment the corresponding counters

But we can easily extend them to handle deletions  
This extended version is called a **Counting Bloom Filter**



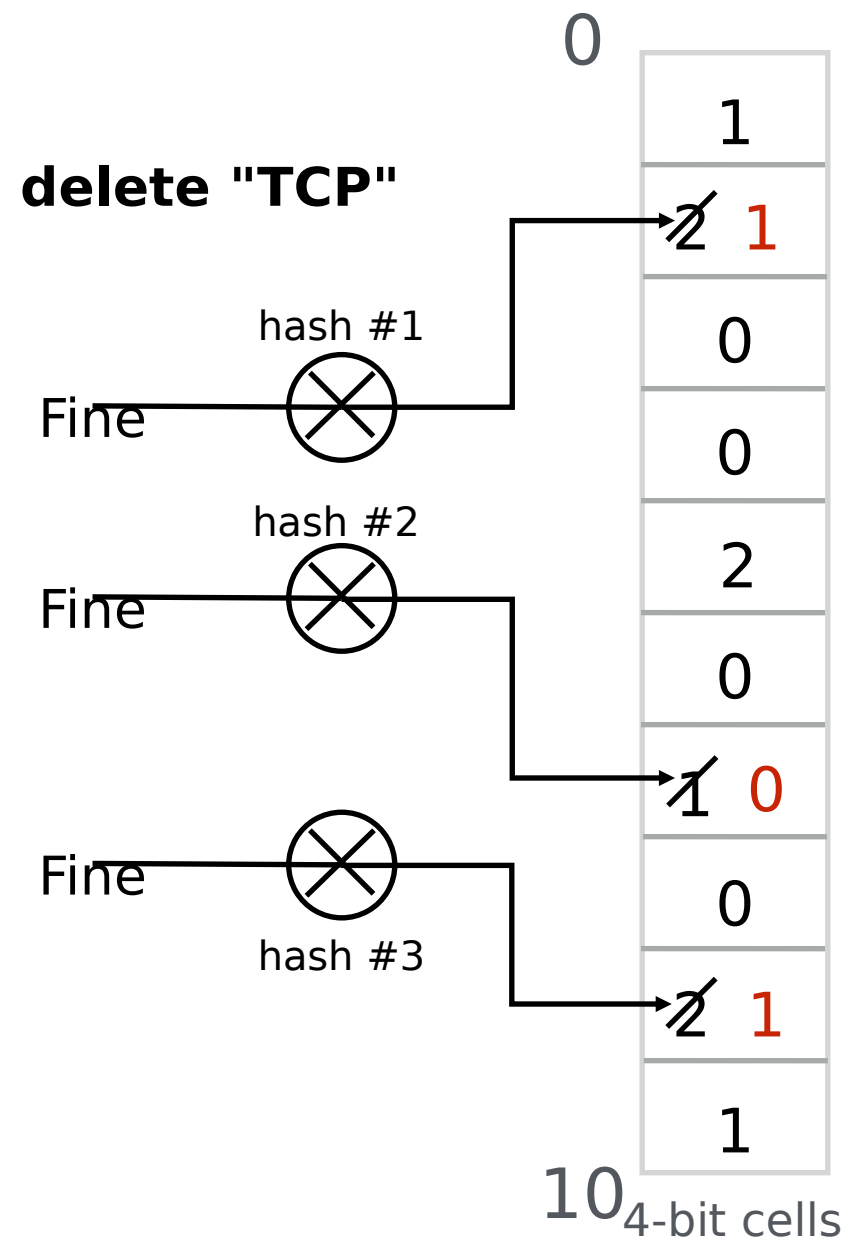
To add an element, increment  
the corresponding counters

But we can easily extend them to handle deletions  
This extended version is called a **Counting Bloom Filter**



To add an element, increment  
the corresponding counters

But we can easily extend them to handle deletions  
This extended version is called a **Counting Bloom Filter**

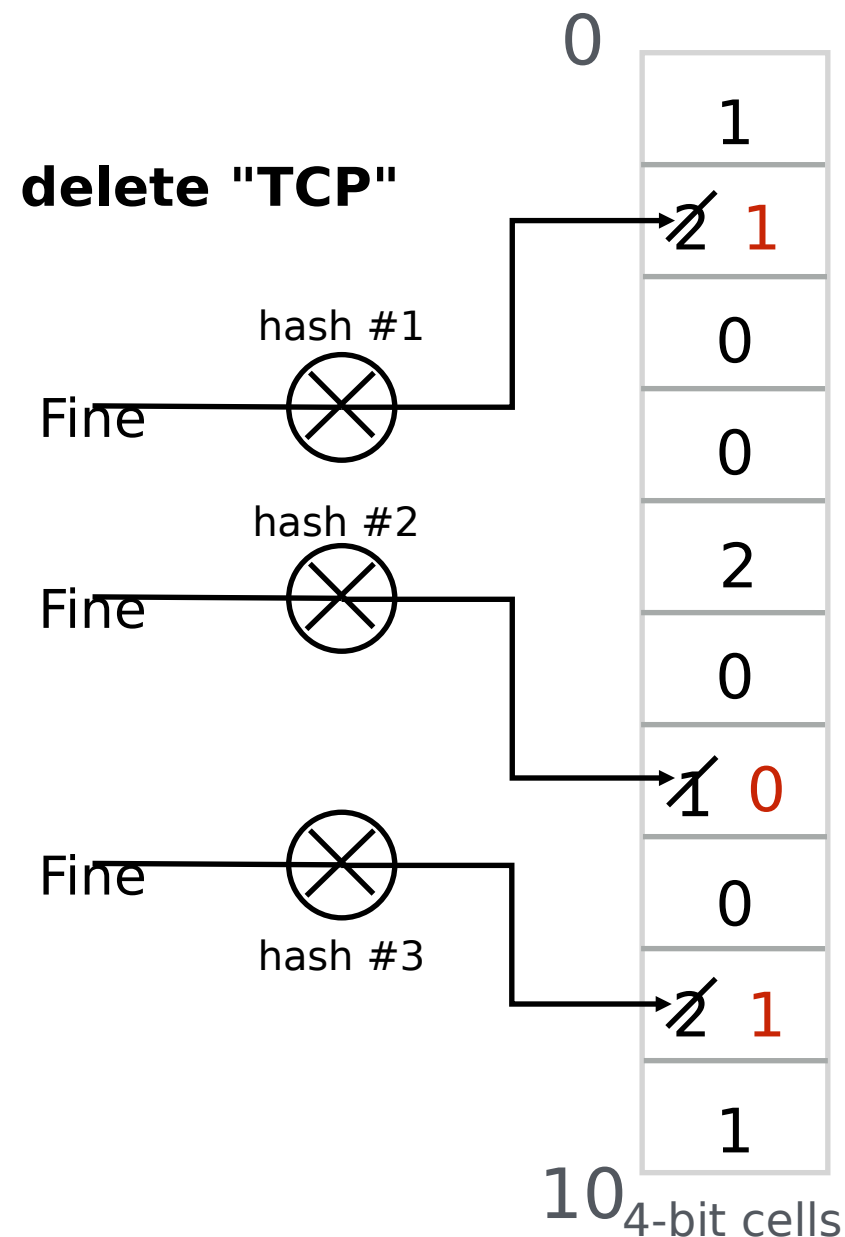


To add an element, increment  
the corresponding counters

To delete an element, decrement  
the corresponding counters



But we can easily extend them to handle deletions  
This extended version is called a **Counting Bloom Filter**



To add an element, increment the corresponding counters

To delete an element, decrement the corresponding counters

All of our prior analysis for standard bloom filters applies to counting bloom filters

Counting Bloom Filters do handle deletions  
at the price of using more memory

Counting Bloom Filters do handle deletions  
at the price of using more memory

Counters must be large enough to avoid overflow  
If a counter eventually overflows, the filter may yield  
false negatives

Counting Bloom Filters do handle **deletions**  
at the price of using **more memory**

Counters must be large enough to avoid overflow  
If a counter eventually overflows, the filter may yield  
false negatives

Poisson approximation suggests 4 bits/counter

The average load (i.e.,  $\frac{NK}{M}$ ) is  $\ln 2$  assuming  $K = \ln 2 * (M/N)$

With  $N=10000$  and  $M=80000$  the probability that some  
counter overflows if we use  $b$ -bit counters is at most

$$M * Pr(Poisson(\ln 2) \geq 2^b) = 1.78e-11$$

# Implementation of a Counting Bloom Filter in P4<sub>16</sub>

## with 2 hash functions

### Add a new element

```
control MyIngress(...) {
    register<bit<1>>(NB_CELLS) bloom_filter;
    apply {
        hash(meta.index1, HashAlgorithm.my_hash1, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);
        hash(meta.index2, HashAlgorithm.my_hash2, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);

        // Add a new element if not yet in the set
        bloom_filter.read(meta.query1, meta.index1);
        bloom_filter.read(meta.query2, meta.index2);

        if (meta.query1 == 0 || meta.query2 == 0) {
            bloom_filter.write(meta.index1, meta.query1 + 1);
            bloom_filter.write(meta.index2, meta.query2 + 1);
        }
    }
}
```

# Implementation of a Counting Bloom Filter in P4<sub>16</sub>

## with 2 hash functions

### Delete an element

```
control MyIngress(...) {
    register<bit<1>>(NB_CELLS) bloom_filter;
    apply {
        hash(meta.index1, HashAlgorithm.my_hash1, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);
        hash(meta.index2, HashAlgorithm.my_hash2, 0,
            {meta.dstPrefix, packet.ip.srcIP}, NB_CELLS);

        // Delete a element only if it is in the set
        bloom_filter.read(meta.query1, meta.index1);
        bloom_filter.read(meta.query2, meta.index2);

        if (meta.query1 > 0 && meta.query2 > 0) {
            bloom_filter.write(meta.index1, meta.query1 - 1);
            bloom_filter.write(meta.index2, meta.query2 - 1);
        }
    }
}
```

So far we have seen how to do insertions, deletions and membership queries

	strategy #1	strategy #2
output	Deterministic	Probabilistic
number of required operations	Probabilistic	Deterministic

Bloom Filters  
Counting Bloom Filters

Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**



Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

Each cell contains three fields

**count** which counts the number of entries mapped to this cell

**keySum** which is the sum of all the keys mapped to this cell

**valueSum** which is the sum of all the values mapped to this cell

Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

**Add** a new key-value pair  
(assuming it is not in the set yet )

**For each** hash function  
**hash** the key to find the  
index

Then at this index  
**increment** the count by one  
**add** key to keySum  
**add** value to valueSum

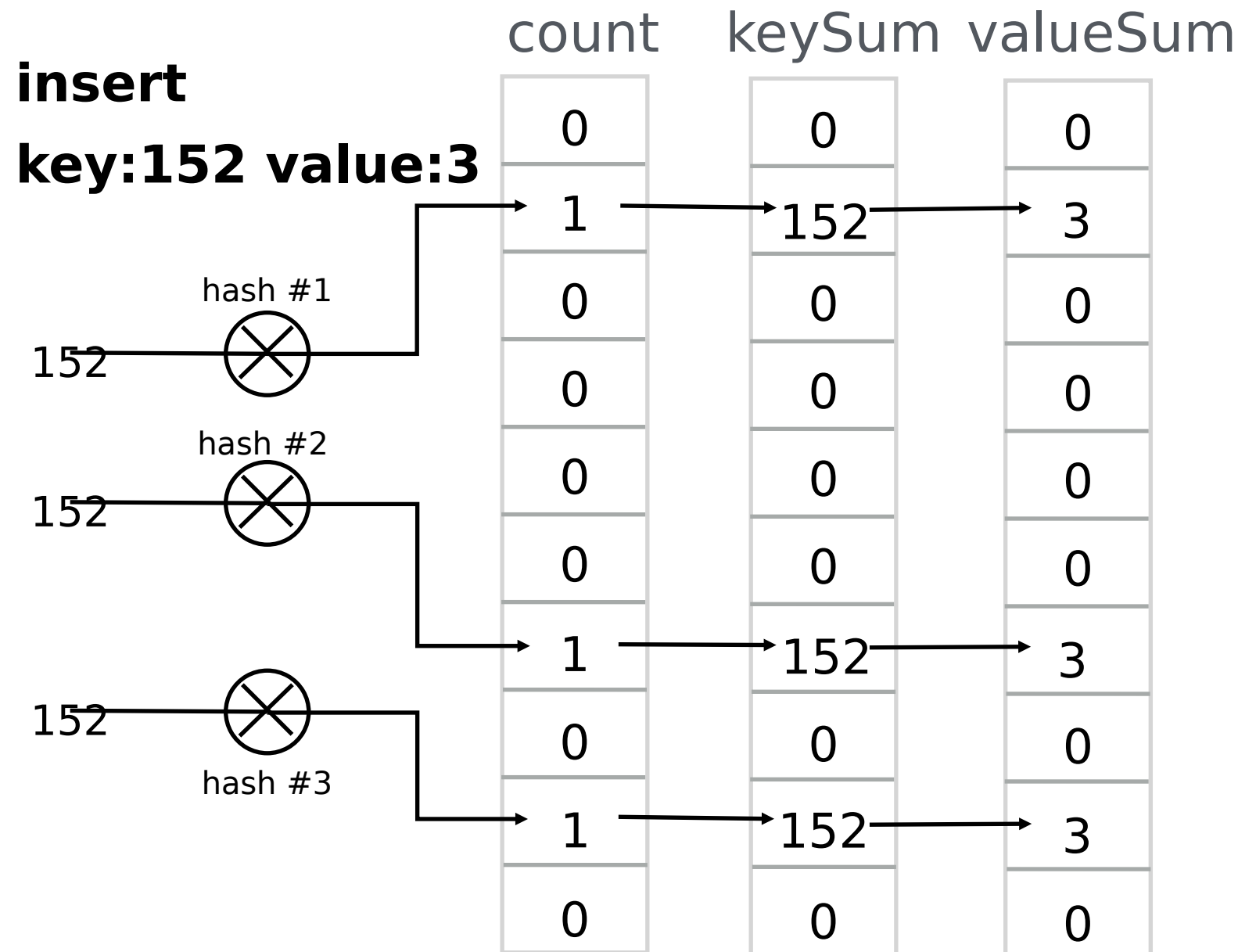
Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

**Delete** a key-value pair  
(assuming it is in the set)

**For each** hash function  
**hash** the key to find the  
index

Then at this index  
**subtract** one to the count  
**subtract** key to keySum  
**subtract** value to valueSum

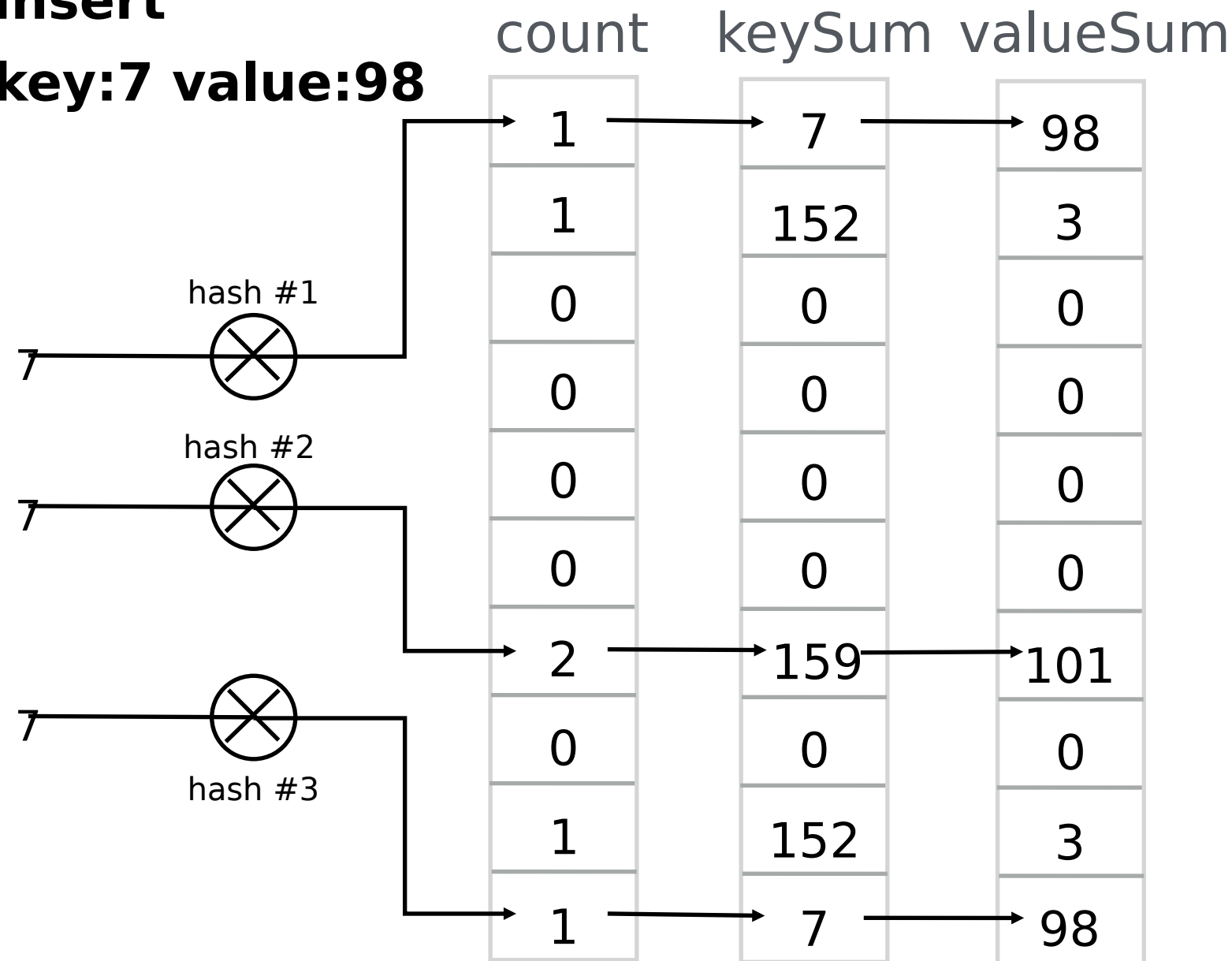
Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**



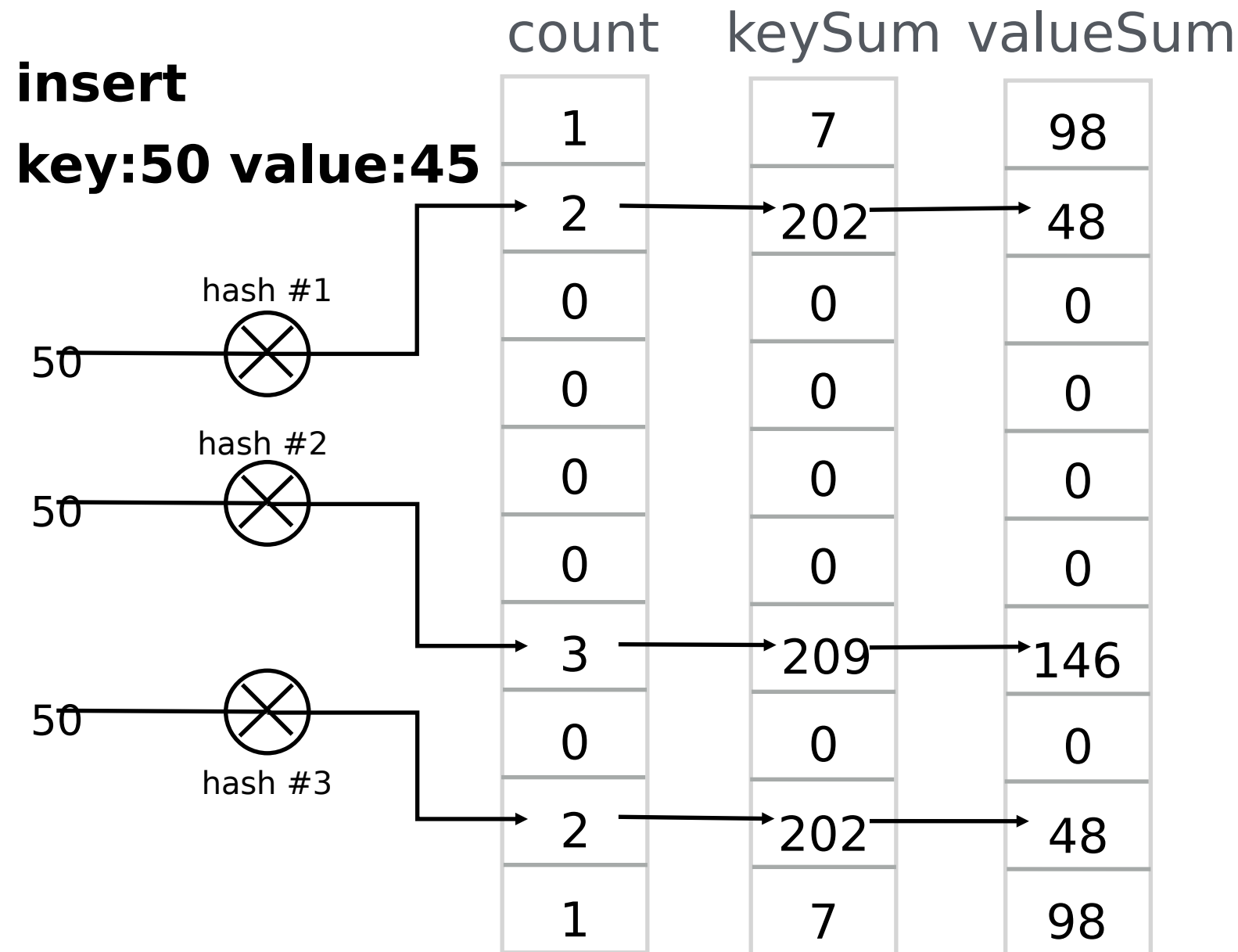
Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

**insert**

**key:7 value:98**



Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**



Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

Key-value pair **lookup**

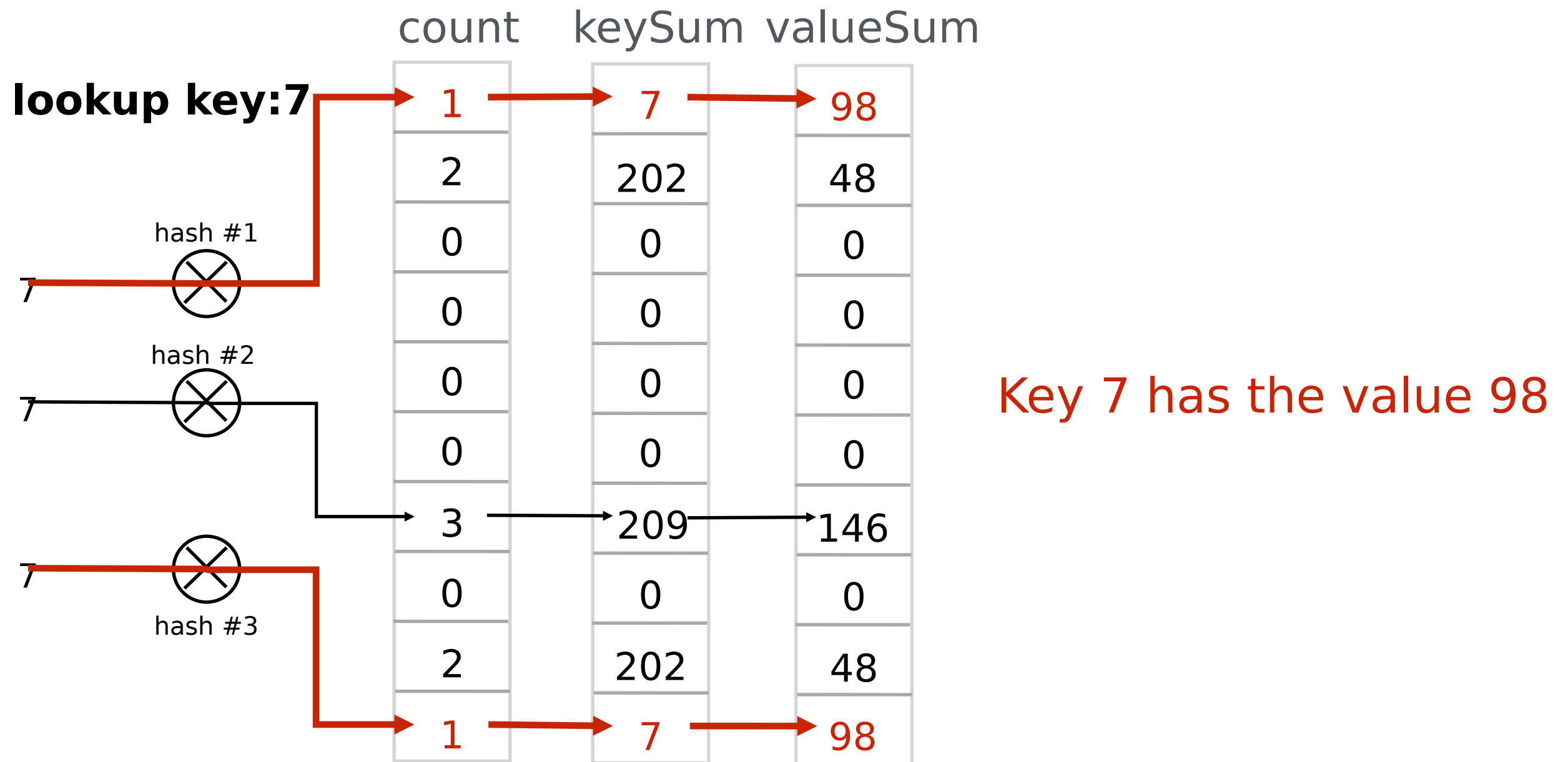
The value of a key can be found if the key is associated to **at least** one cell with a count = 1

Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

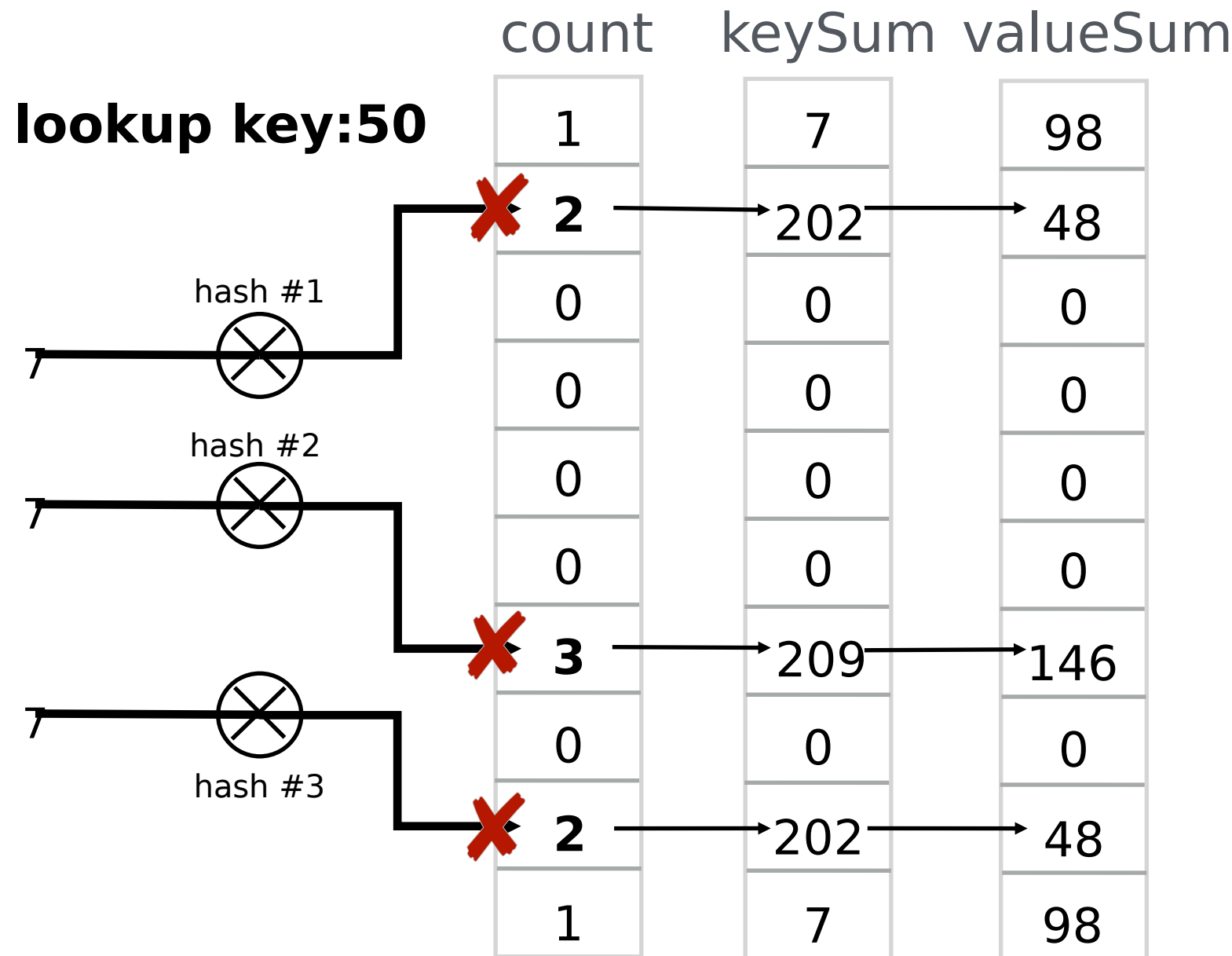
count	keySum	valueSum
1	7	98
2 →	202 →	48
0	0	0
0	0	0
0	0	0
0	0	0
3 →	209 →	146
0	0	0
2 →	202 →	48
1	7	98



Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**



Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**



The value for the key 50  
can't be found

Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

**Listing** the IBLT

**While** there is an index for which count = 1

**Find** the corresponding key-value pair and return it

**Delete** the corresponding key-value pair

Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

## **Listing** the IBLT

**While** there is an index for which count = 1  
**Find** the corresponding key-value pair and return it  
**Delete** the corresponding key-value pair

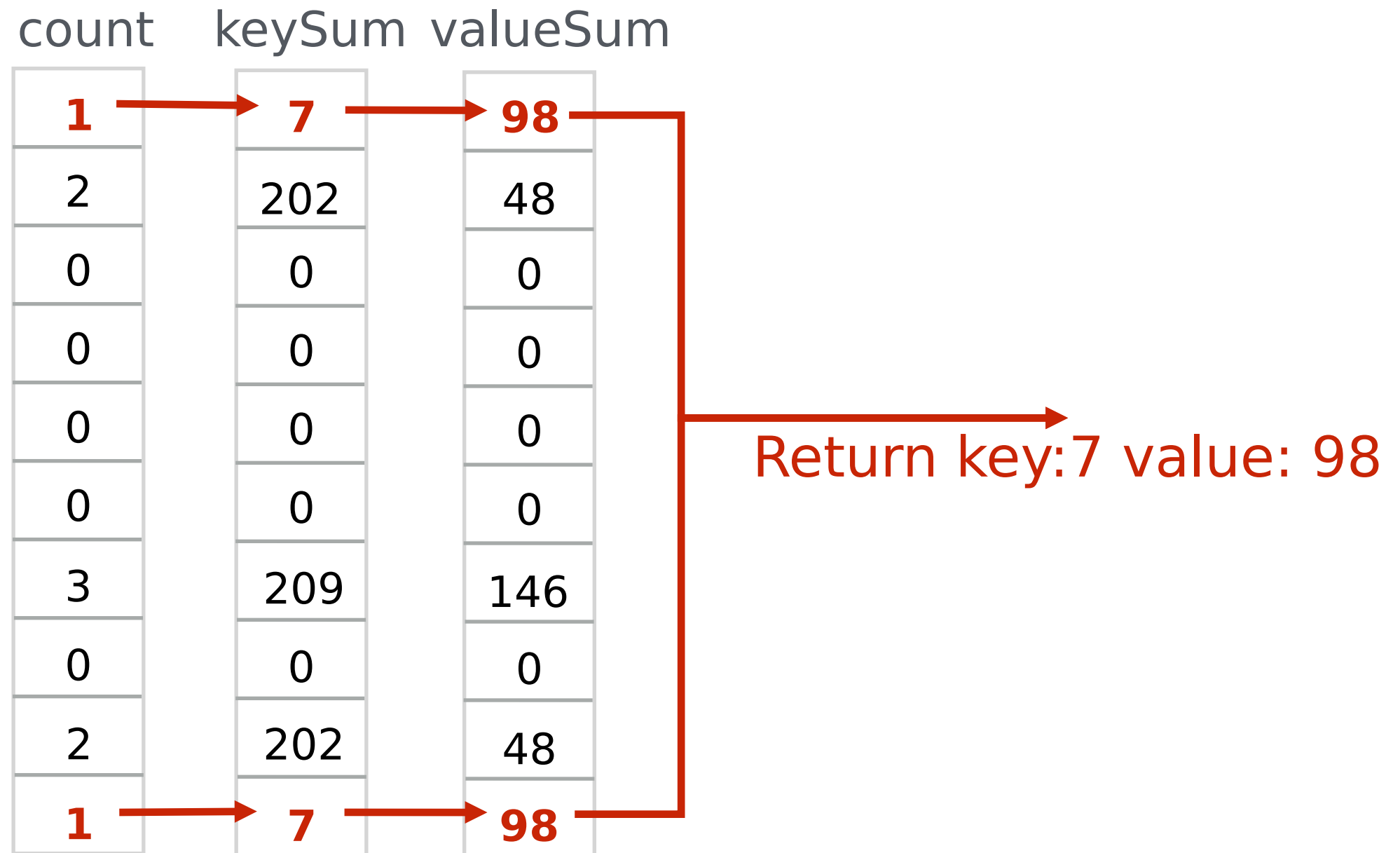
Unless the number of iterations is very low, loops can't be implemented in hardware

**The listing is done by the controller**

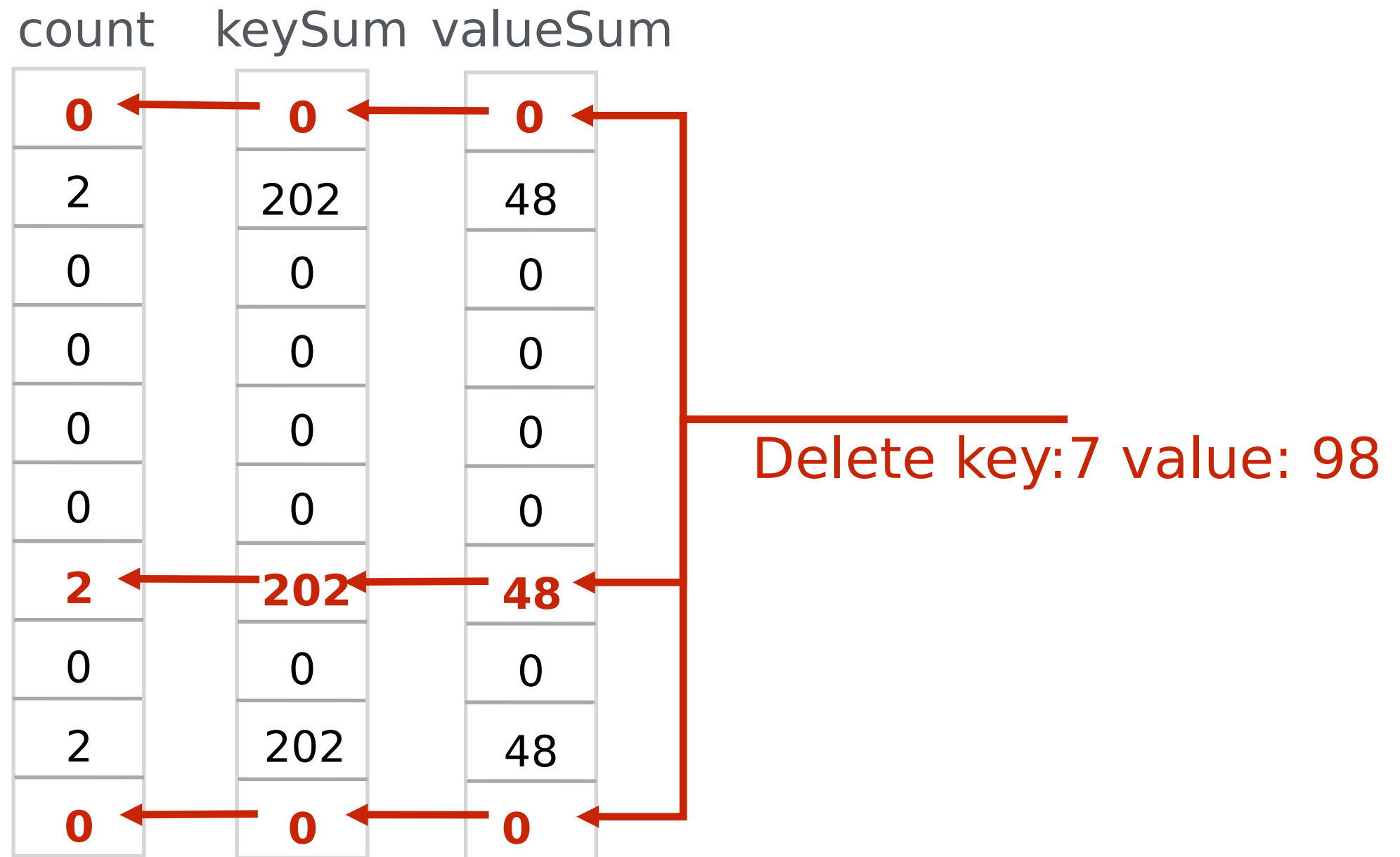
Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

count	keySum	valueSum
1	7	98
2	202	48
0	0	0
0	0	0
0	0	0
0	0	0
3	209	146
0	0	0
2	202	48
1	7	98

Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**



Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**



Invertible Bloom Lookup Tables (IBLT) stores key-value pairs and allows for **lookups** and a complete **listing**

count	keySum	valueSum
0	0	0
<b>2</b>	202	48
0	0	0
0	0	0
0	0	0
0	0	0
<b>2</b>	202	48
0	0	0
<b>2</b>	202	48
0	0	0

In this example, a complete listing is not possible



In many settings, we can use XORs in place of sums  
For example to avoid overflow issues

count	keySum	XOR
1	7	
2	202	
0	0	
0	0	
0	0	
0	0	
3	209	
0	0	
2	202	
1	7	


Bloom filters are **probabilistic data structures** for **set membership queries**. For more info, see:

Space/Time Trade-offs in Hash Coding with Allowable Errors. Burton H. Bloom. 1970.

Network Applications of Bloom Filters: A Survey. Andrei Broder and Michael Mitzenmacher. 2004.

Invertible Bloom Lookup Tables. Michael T. Goodrich and Michael Mitzenmacher. 2015.

FlowRadar: A Better NetFlow for Data Centers Yuliang Li et al. NSDI 2016.



← You are looking at a stream of data (packets).  
Today, I'll show you how set membership and frequency queries can be realized in P4.

**PART 1**

Is a certain element (e.g. ip address) in the stream?  
→ Bloom filter

**PART 2**

How frequently does an element appear?  
→ CountMin Sketch, Count Sketch, ...

## part 2: **counting with sketches**

How frequently does an element appear?

(slides by yours truly)

We are going to look at **frequencies**,  
i.e. **how often** an element occurs in a data stream.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \quad \begin{array}{l} \text{vector of frequencies (counts)} \\ \text{of all *distinct* elements } x_i \end{array}$$

We are going to look at **frequencies**,  
i.e. **how often** an element occurs in a data stream.

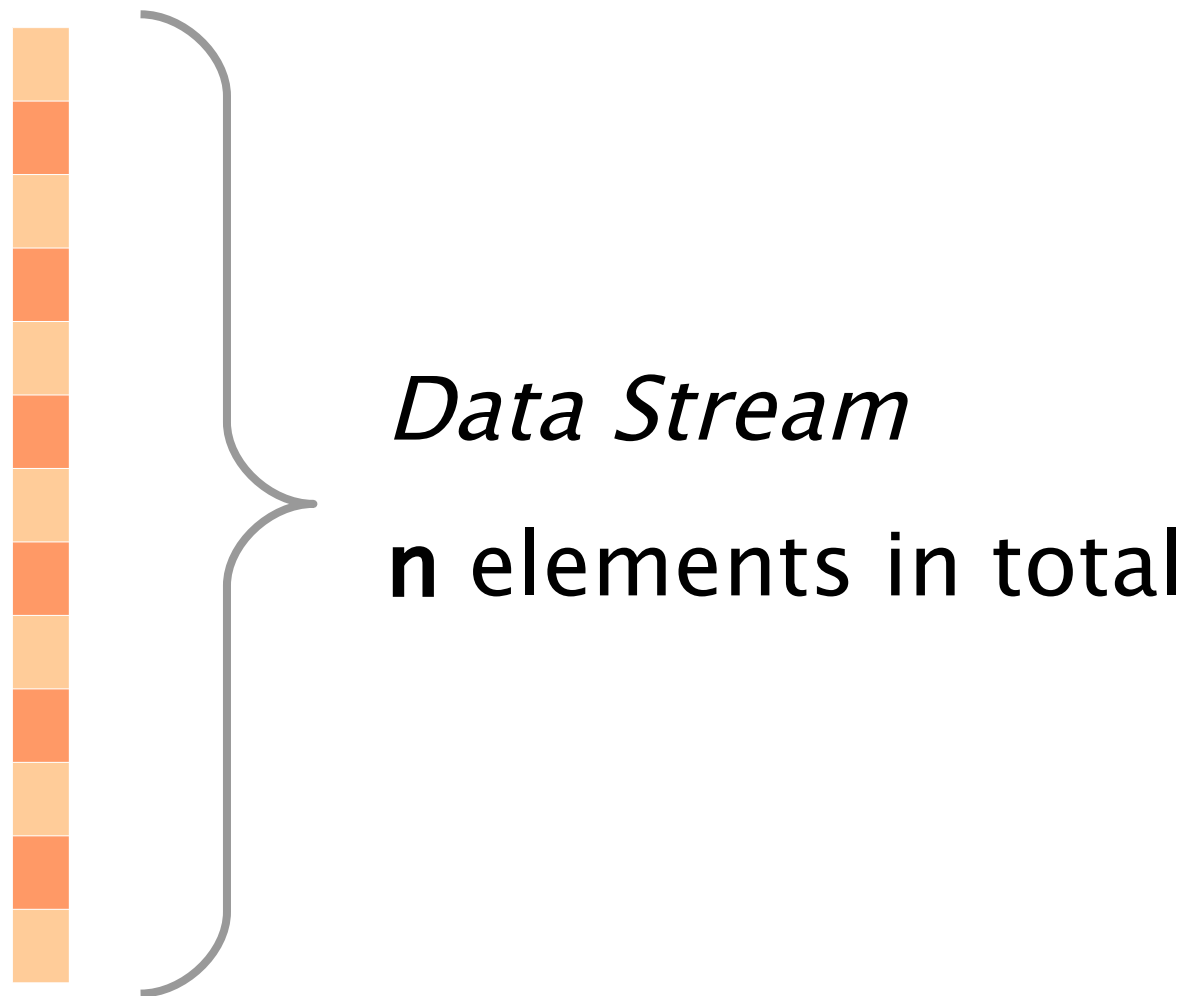
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

*vector of frequencies (counts)  
of all **distinct elements**  $x_i$*

*e.g. flows, ip addresses, ...*

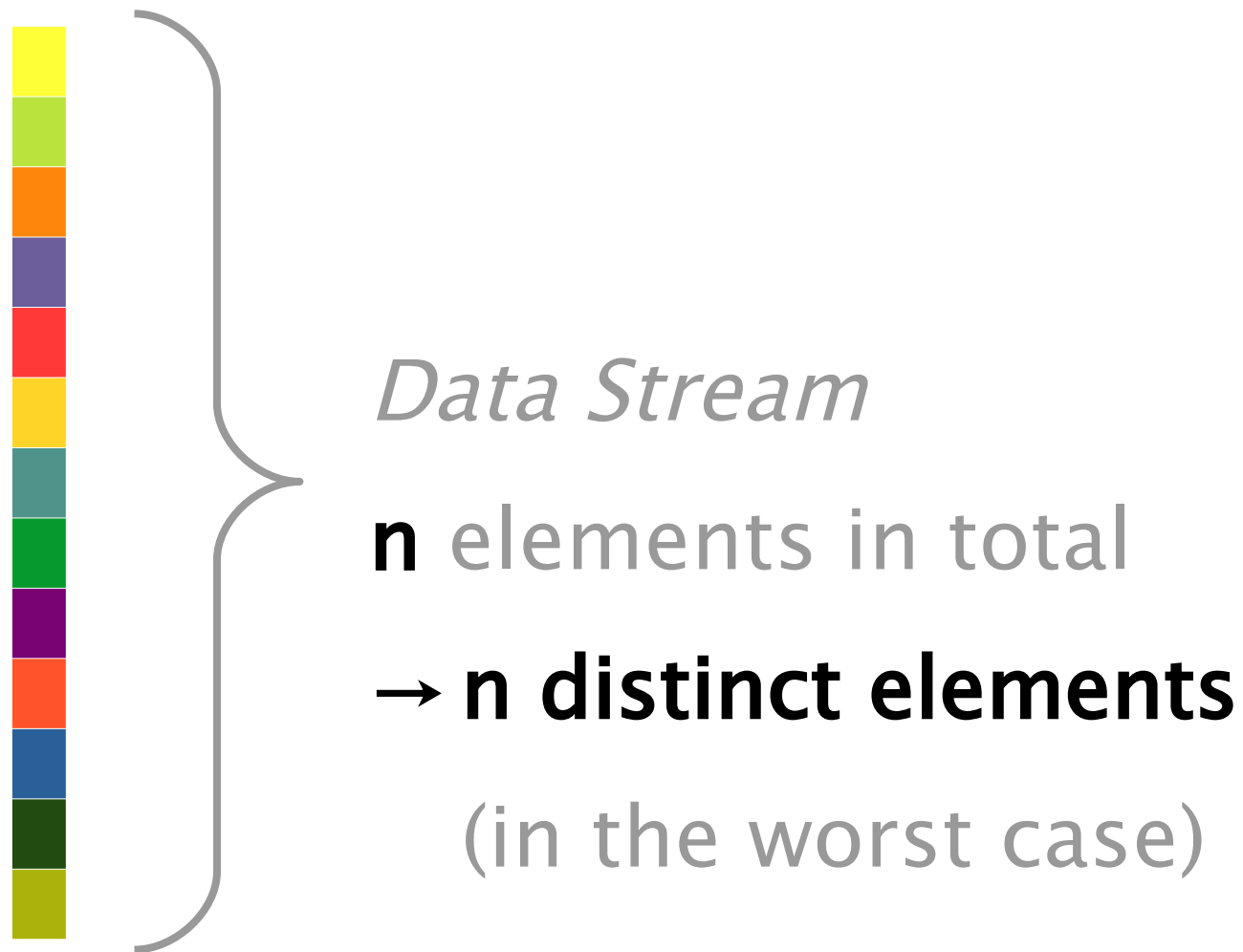
In the worst case, an algorithm providing **exact frequencies** requires **linear space**.

In the worst case, an algorithm providing  
**exact frequencies** requires **linear space**.

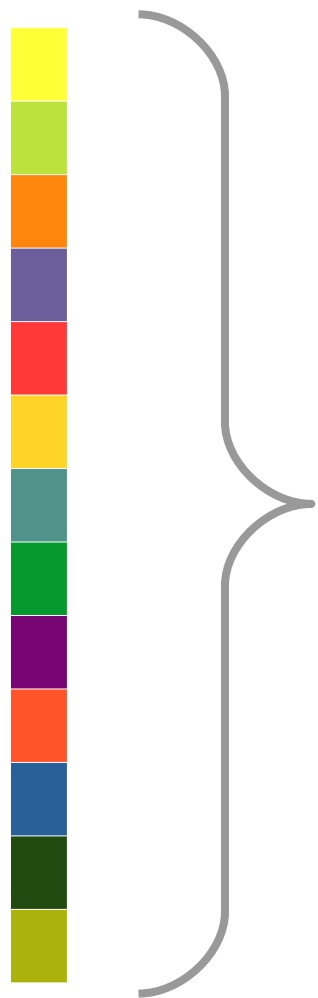




In the worst case, an algorithm providing exact frequencies requires linear space.



In the worst case, an algorithm providing exact frequencies requires linear space.



*Data Stream*

**n** elements in total

→ **n** distinct elements

(in the worst case)

→ **n counters** required? :(

# Probabilistic datastructures can help again!

## Bloom Filters

quickly “filter” only those elements that might be in the set

More efficient by allowing false positives.

# Probabilistic datastructures can help again!

## Bloom Filters

quickly “filter” only those elements that might be in the set

More efficient by allowing false positives.

## Sketches

provide a approximate frequencies of elements in a data stream.

More efficient by allowing mis-counting.

A **CountMin sketch** uses the same principles as a counting bloom filter, but is **designed** to have **provable L1 error bounds** for frequency queries.

A **CountMin sketch** uses the same principles as a counting bloom filter, but is **designed** to have **provable L1 error bounds** for frequency queries.

Notation reminder:

vector of frequencies (counts)  
of all **distinct elements**  $x_i$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

A CountMin sketch uses the same principles as a counting bloom filter, but is **designed** to have **provable L1 error bounds** for frequency queries.

$$Pr \left[ \underbrace{\hat{x}_i}_{\text{estimated frequency}} - \underbrace{x_i}_{\text{true frequency}} \geq \underbrace{\varepsilon \|\mathbf{x}\|_1}_{\text{sum of frequencies}} \right] \leq \delta$$

*The estimation error **exceeds**  $\varepsilon \|\mathbf{x}\|_1$   
with a **probability smaller than**  $\delta$*



$$Pr \left[ \underbrace{\hat{x}_i}_{\text{estimated frequency}} - \underbrace{x_i}_{\text{true frequency}} \geq \underbrace{\varepsilon \|\mathbf{x}\|_1}_{\substack{\text{relative to L1 norm} \\ \text{sum of frequencies}}} \right] \leq \delta$$

*The estimation error **exceeds**  $\varepsilon \|\mathbf{x}\|_1$   
with a **probability smaller than**  $\delta$*

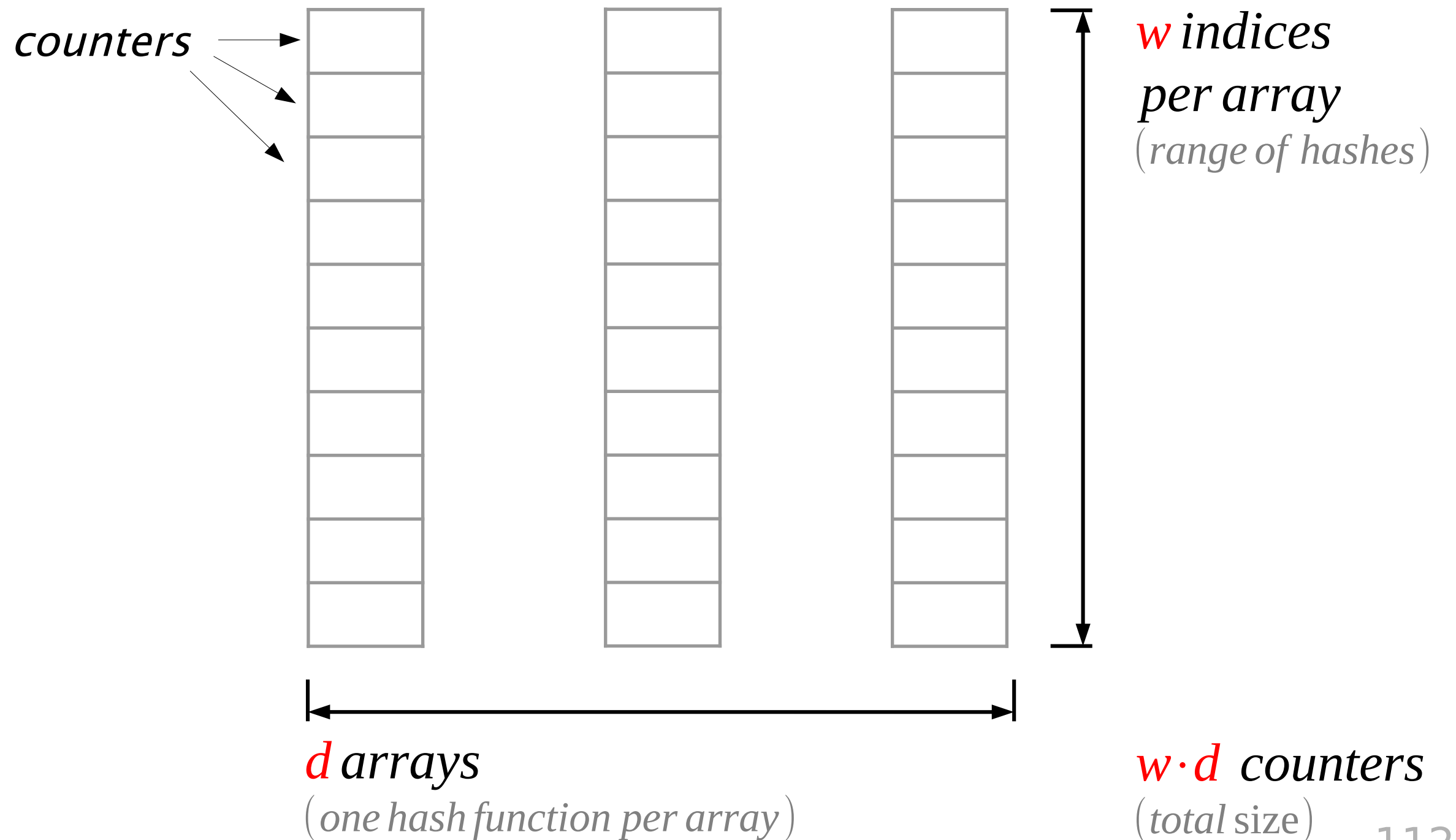
$$Pr \left[ \underset{\substack{\text{estimated} \\ \text{frequency}}}{\hat{x}_i} - \underset{\substack{\text{true} \\ \text{frequency}}}{x_i} \geq \underset{\substack{\text{sum of} \\ \text{frequencies}}}{\varepsilon \|\mathbf{x}\|_1} \right] \leq \delta$$

*Let  $\|\mathbf{x}\|_1 = 10000$  ,  $\varepsilon = 0.01$  ,  $\delta = 0.05$*

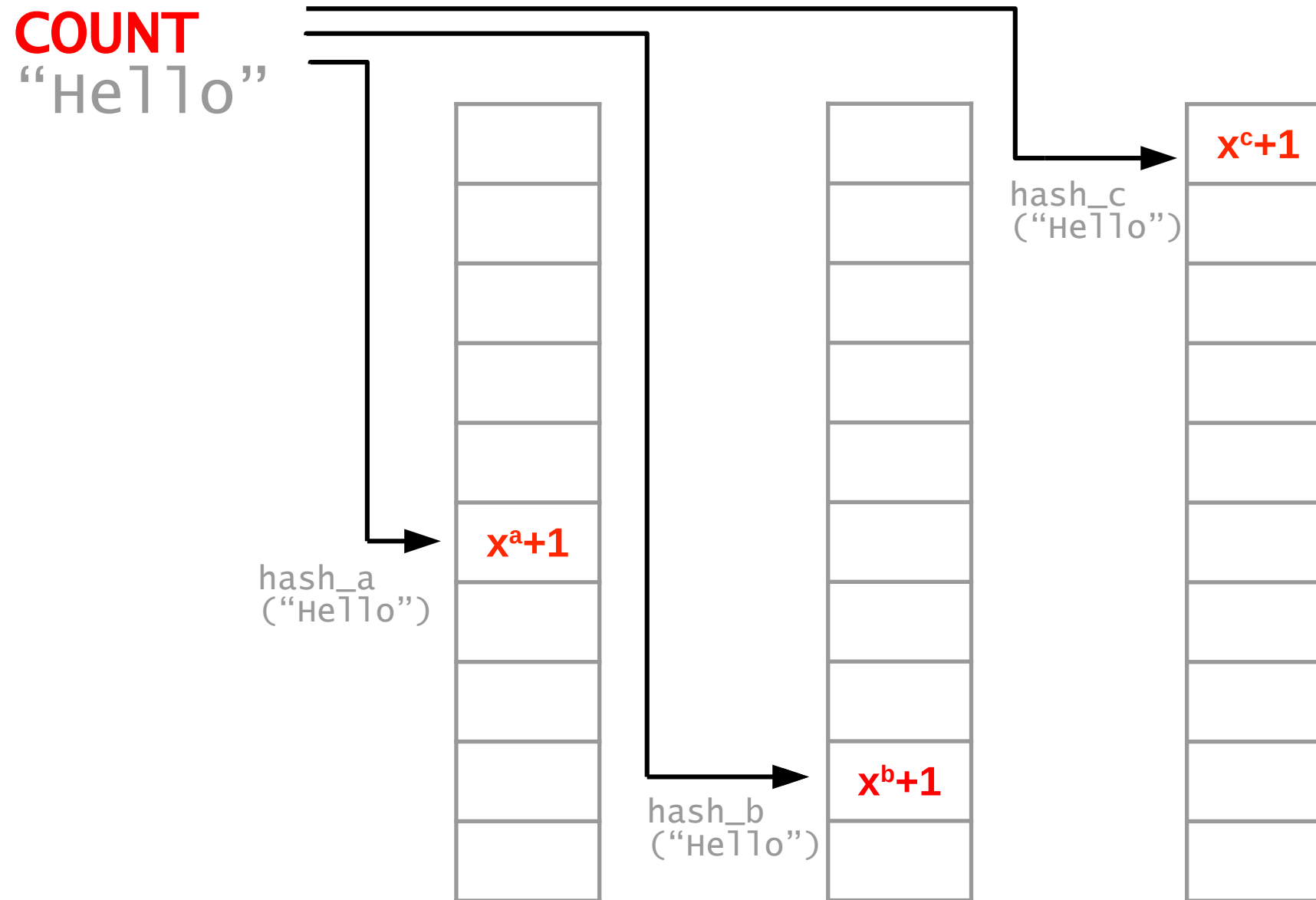
*Then, after counting **10000 elements in total**,  
the probability for **any estimate** to be  
off by **more than 100** is less than **5%**.*

A **CountMin sketch** uses the same principles as a counting bloom filter, but is **designed** to have **provable L1 error bounds** for frequency queries.

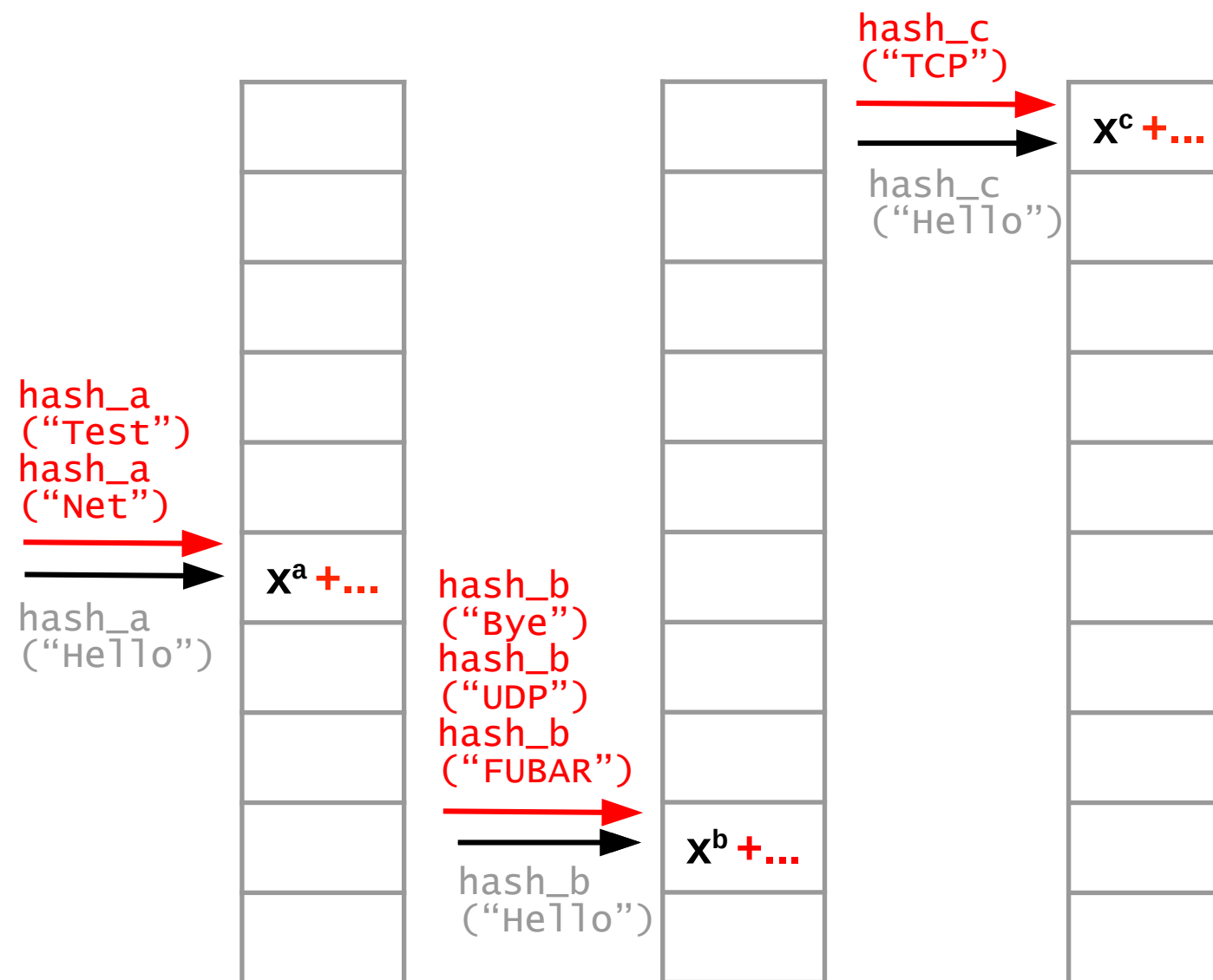
A CountMin Sketch uses multiple arrays and hashes.



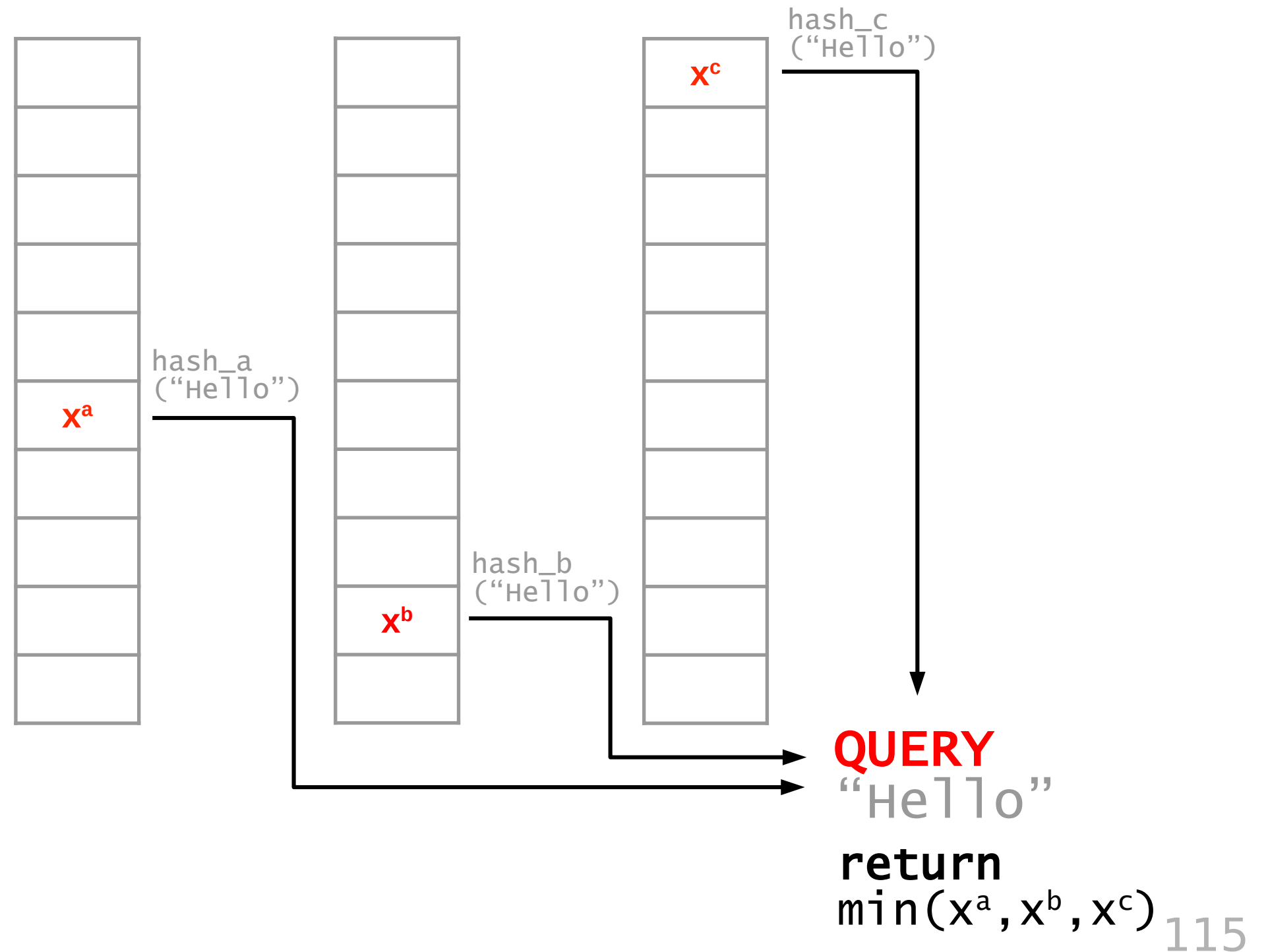
To count, increment **all** hash-indexed fields by 1.



# Hash collisions cause over-counting.



Returning the **minimum value** minimizes the error.

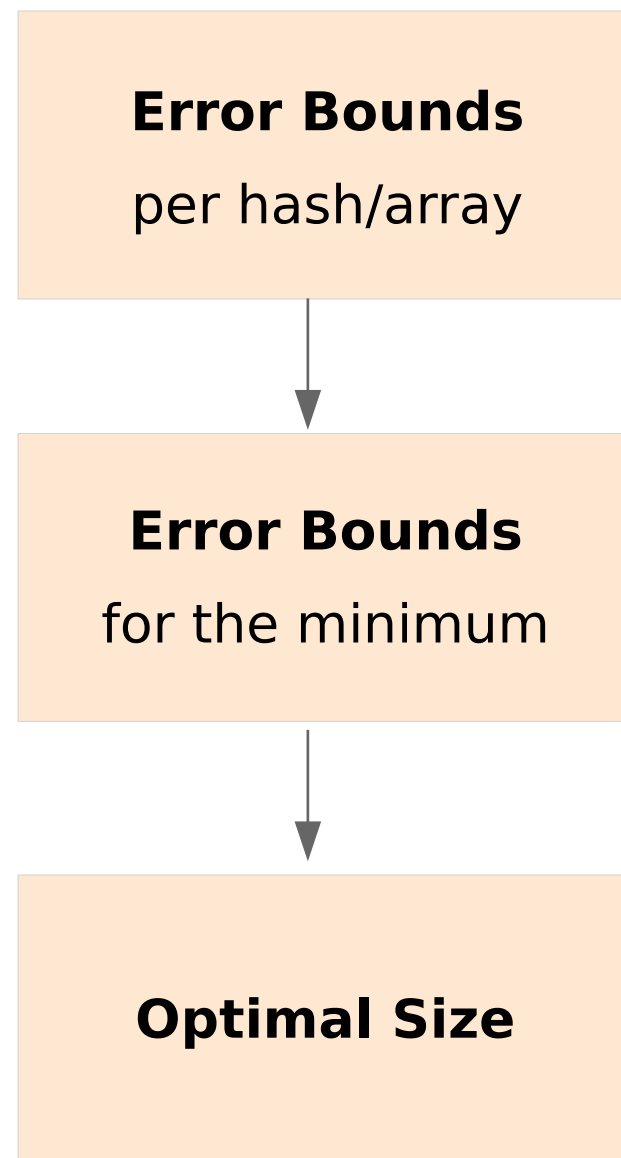


A CountMin sketch uses the same principles as a counting bloom filter, but is designed to have **provable L1 error bounds** for frequency queries.

$$Pr \left[ \underbrace{\hat{x}_i}_{\text{estimated frequency}} - \underbrace{x_i}_{\text{true frequency}} \geq \underbrace{\varepsilon \|\mathbf{x}\|_1}_{\text{sum of frequencies}} \right] \leq \delta$$



Understanding the error bounds allows  
**dimensioning** the sketch optimally.



**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$\hat{x}_i = \min_{h \in h_1 \dots h_d} \hat{x}_i^h$$

*estimated frequency* *estimate for specific hash*

The error bounds can be derived  
with **Markov's Inequality**

$$\Pr [X \geq c \cdot E[X]] \leq \frac{1}{c}$$

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

# The error bounds can be derived with **Markov's Inequality**

$$\Pr [\hat{x}_i^h - x_i \geq c \cdot E [\hat{x}_i^h - x_i]] \leq \frac{1}{c}$$

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$\Pr [\hat{x}_i^h - x_i \geq c \cdot E [\hat{x}_i^h - x_i]] \leq \frac{1}{c}$$

$$\hat{x}_i^h = x_i + \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

*true  
frequency*

*over-counting  
from hash collisions*

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



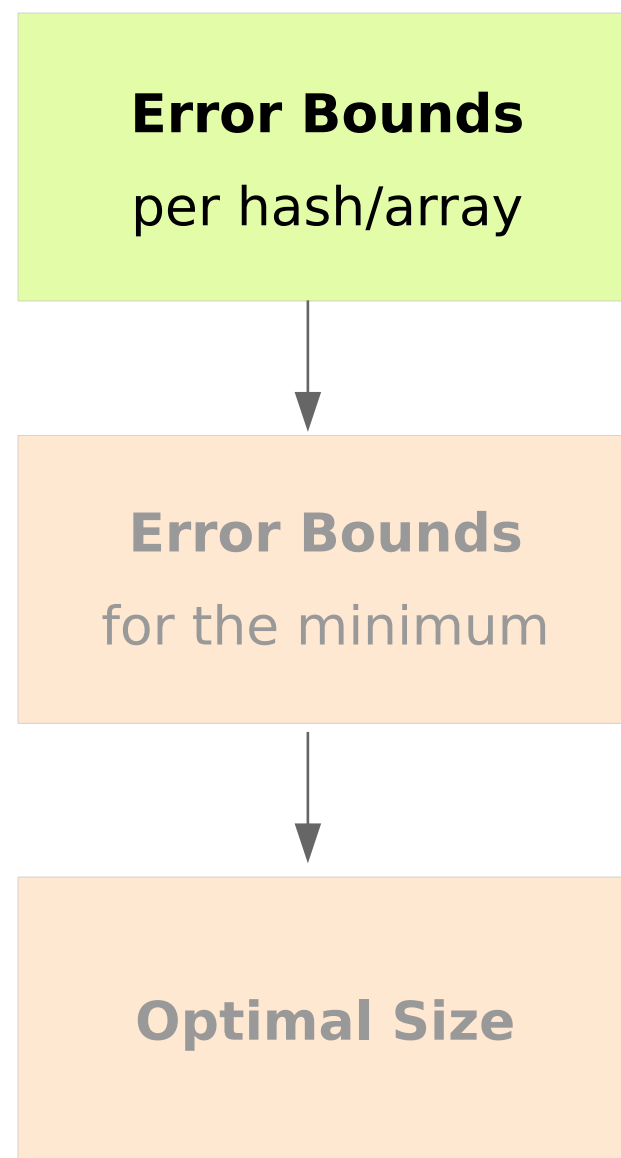
**Optimal Size**

$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h = x_i + \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

*hash collision*

$$= \begin{cases} 1, & \text{if } h(x_i) = h(x_j) \\ 0, & \text{otherwise} \end{cases}$$



$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h - x_i = \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

*estimation  
error*

*over-counting  
from hash collisions*

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

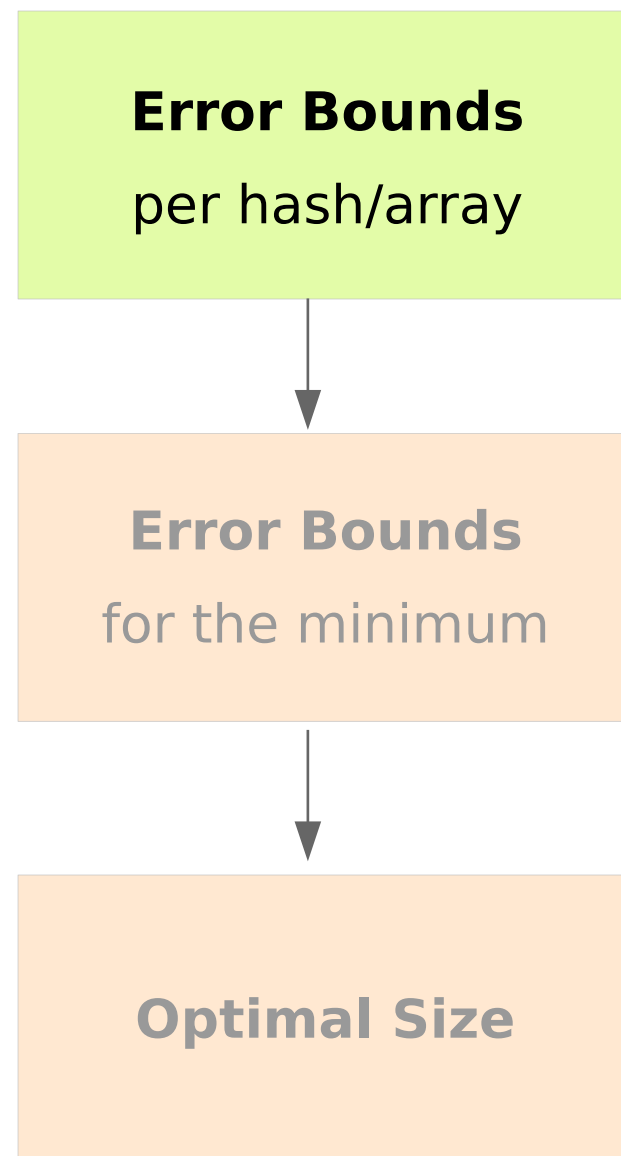
$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h - x_i = \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

$$E \left[ \hat{x}_i^h - x_i \right] = E \left[ \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j) \right]$$



We treat the **data as a constant** and the **hash as a random function** with certain properties.

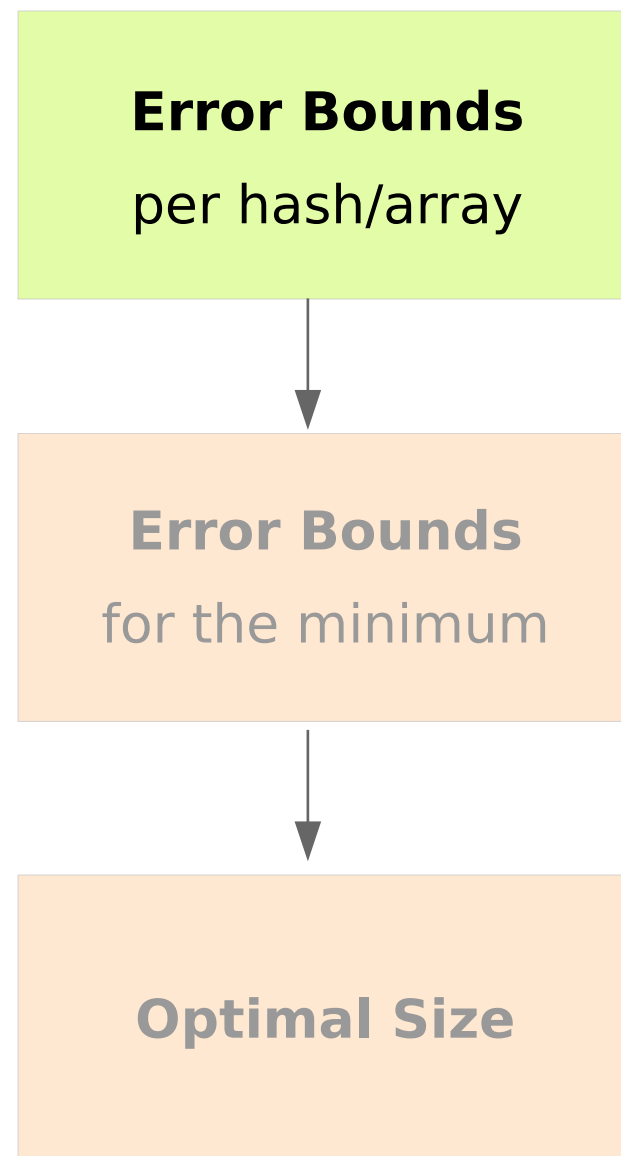


$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h - x_i = \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

$$E \left[ \hat{x}_i^h - x_i \right] = E \left[ \sum_{x_j \neq x_i} \underbrace{x_j}_{\text{constant}} \underbrace{1_h(x_i, x_j)}_{\text{random}} \right]$$

We treat the **data as a constant** and the **hash as a random function** with certain properties.

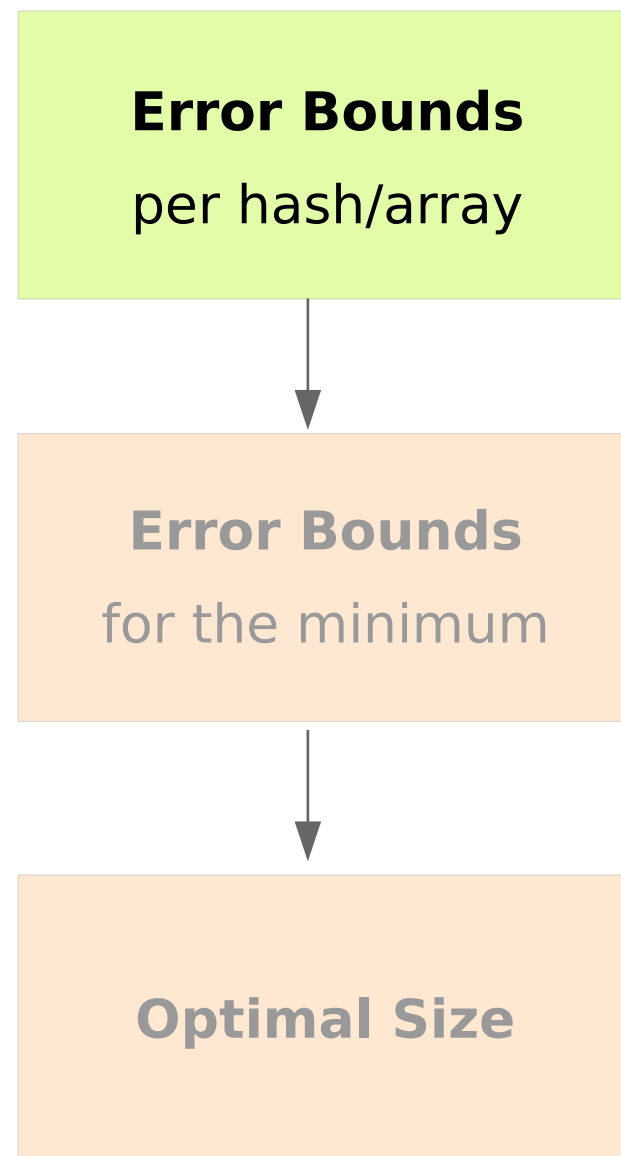


$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h - x_i = \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

$$E \left[ \hat{x}_i^h - x_i \right] = \sum_{x_j \neq x_i} x_j E \left[ 1_h(x_i, x_j) \right]$$

We treat the **data as a constant** and the **hash as a random function** with certain properties.

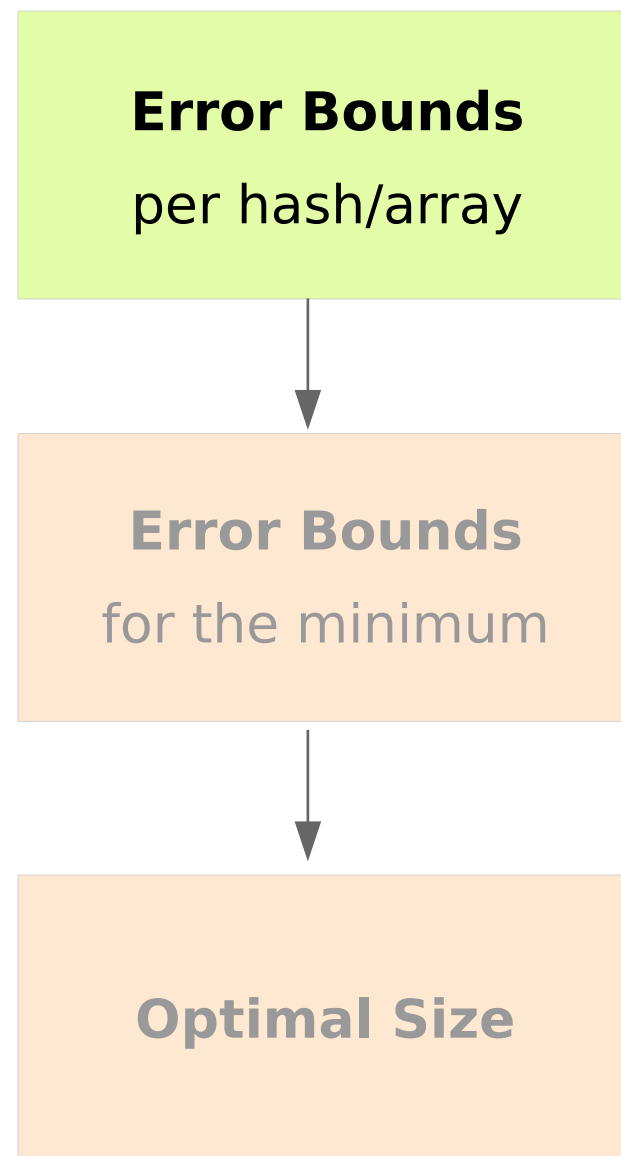


$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h - x_i = \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

$$E \left[ \hat{x}_i^h - x_i \right] = \sum_{x_j \neq x_i} x_j \underbrace{E \left[ 1_h(x_i, x_j) \right]}_{\leq \frac{1}{w}}$$

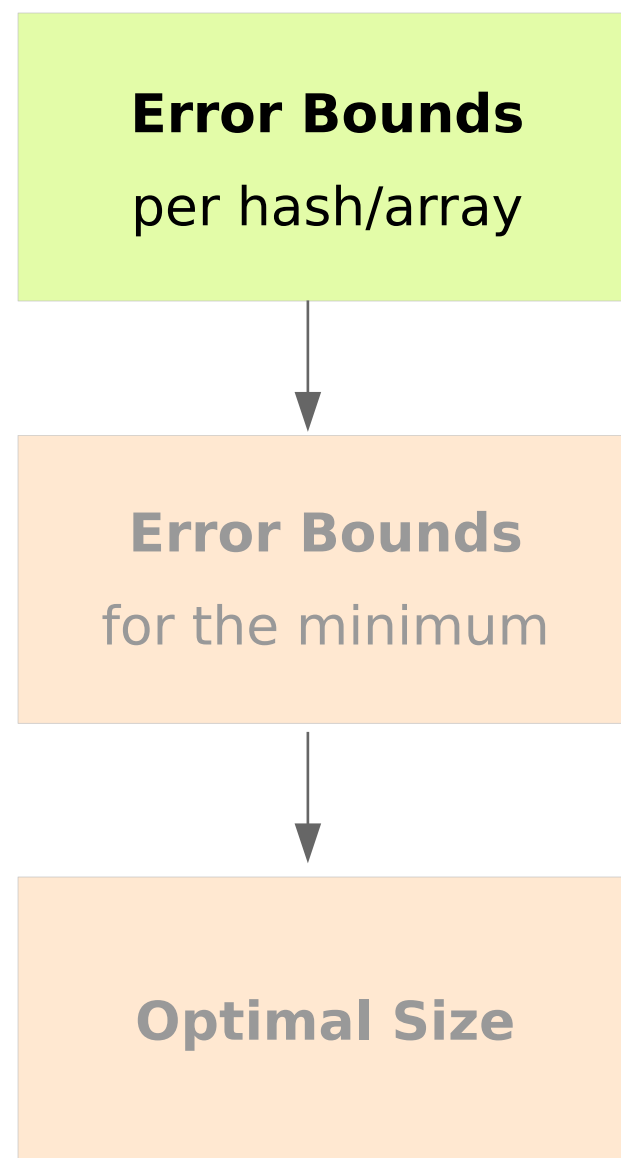
We treat the **data as a constant** and the **hash as a random function** with certain properties.



$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h - x_i = \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

$$E \left[ \hat{x}_i^h - x_i \right] \leq \sum_{x_j \neq x_i} x_j \frac{1}{w}$$



$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h - x_i = \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

$$E \left[ \hat{x}_i^h - x_i \right] \leq \sum_{x_j \neq x_i} x_j \frac{1}{w} \leq \sum_{x_j} x_j \frac{1}{w}$$

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot E \left[ \hat{x}_i^h - x_i \right] \right] \leq \frac{1}{c}$$

$$\hat{x}_i^h - x_i = \sum_{x_j \neq x_i} x_j 1_h(x_i, x_j)$$

$$E \left[ \hat{x}_i^h - x_i \right] \leq \sum_{x_j \neq x_i} x_j \frac{1}{w} \leq \| \mathbf{x} \|_1 \frac{1}{w}$$

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$\Pr \left[ \hat{x}_i^h - x_i \geq c \cdot \underbrace{E \left[ \hat{x}_i^h - x_i \right]}_{\leq \frac{1}{w} \|\mathbf{x}\|_1} \right] \leq \frac{1}{c}$$

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$\Pr \left[ \hat{x}_i^h - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right] \leq \frac{1}{c}$$



**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$\Pr \left[ \hat{x}_i^h - x_i \geq \underbrace{\epsilon^h}_{\frac{c}{w}} \mid \|\mathbf{x}\|_1 \leq \underbrace{\delta^h}_{\frac{1}{c}} \right]$$

**Error Bounds**  
per hash/array



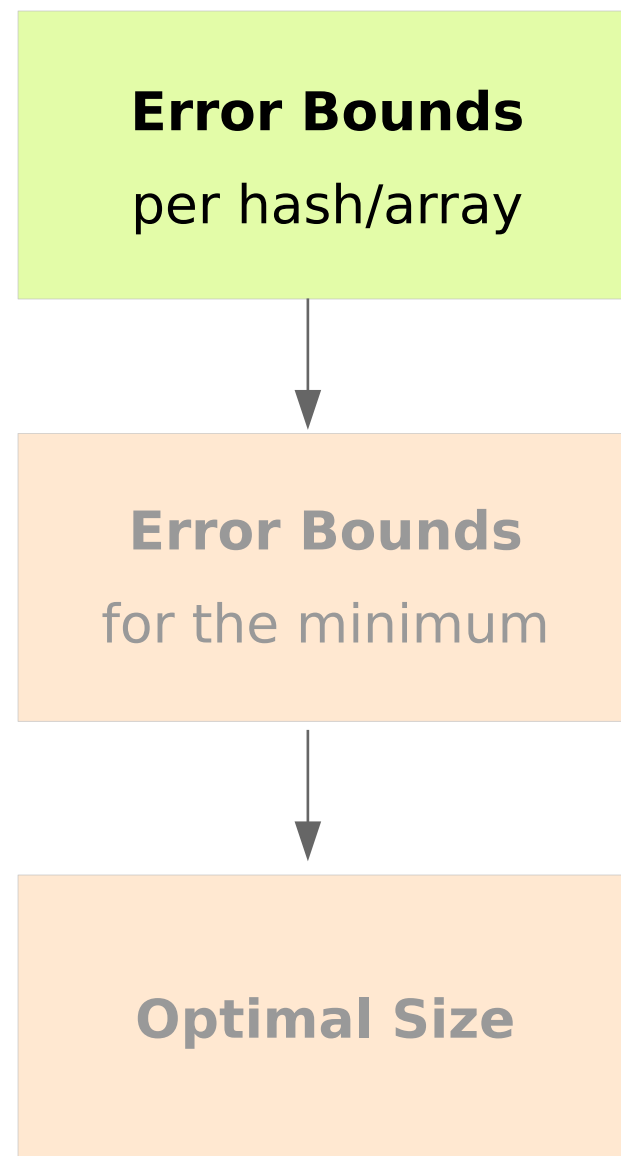
**Error Bounds**  
for the minimum



**Optimal Size**

$$\Pr \left[ \hat{x}_i^h - x_i \geq \underbrace{\epsilon^h}_{\frac{c}{w}} \mid \|\mathbf{x}\|_1 \right] \leq \underbrace{\delta^h}_{\frac{1}{c}}$$

*The **estimate** for each hash has a well defined **L1 error bound**.*



$$\Pr \left[ \hat{x}_i^h - x_i \geq \underbrace{\varepsilon^h}_{\frac{c}{w}} \|\mathbf{x}\|_1 \right] \leq \underbrace{\delta^h}_{\frac{1}{c}}$$

*The **estimate** for each hash has a well defined **L1 error bound**.*

***What about the minimum?***

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$\Pr [\hat{x}_i - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1] \leq ?$$

**Error Bounds**  
per hash/array



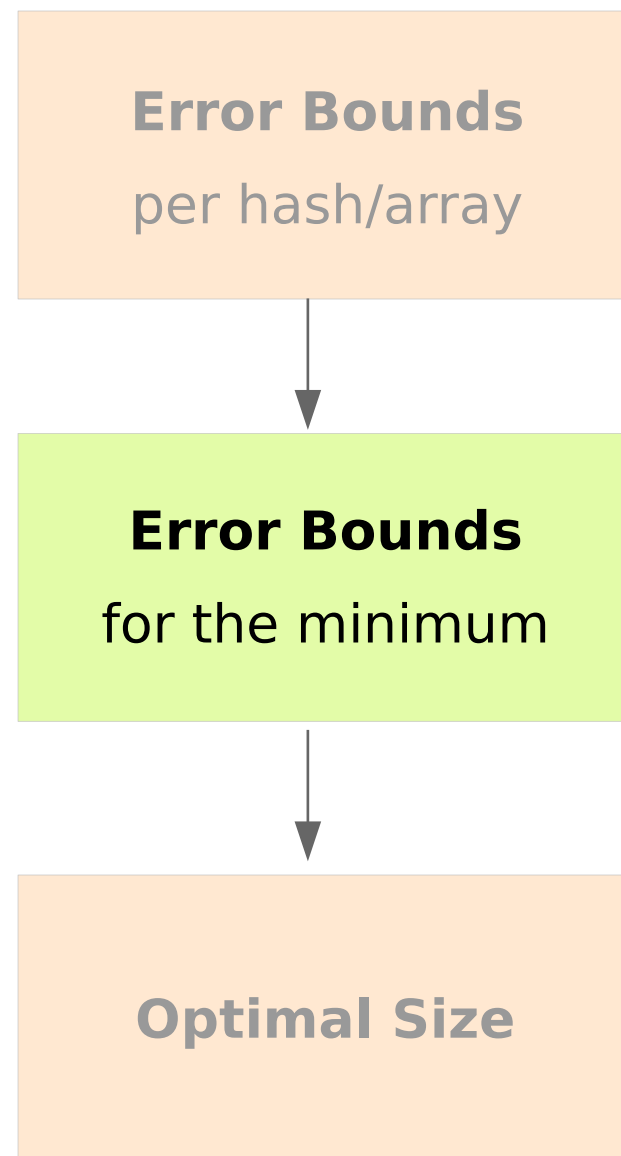
**Error Bounds**  
for the minimum



**Optimal Size**

$$Pr \left[ \underbrace{\min_{h \in h_1 \dots h_d} \hat{x}_i^h}_{\hat{x}_i} - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right] \leq ?$$

Multiple hash functions work like **independent trials**.



$$Pr \left[ \underbrace{\min_{h \in h_1 \dots h_d} \hat{x}_i^h}_{\hat{x}_i} - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right] \leq ?$$

$\Leftrightarrow$

$$\prod_{h \in h_1 \dots h_d} Pr \left[ \hat{x}_i^h - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right] \leq ?$$

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$Pr \left[ \underbrace{\min_{h \in h_1 \dots h_d} \hat{x}_i^h}_{\hat{x}_i} - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right] \leq ?$$

$\Leftrightarrow$

$$\prod_{h \in h_1 \dots h_d} \underbrace{Pr \left[ \hat{x}_i^h - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right]}_{\leq \frac{1}{c}} \leq ?$$

*error bound per hash*

**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$Pr \left[ \underbrace{\min_{h \in h_1 \dots h_d} \hat{x}_i^h}_{\hat{x}_i} - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right] \leq ?$$

$\Leftrightarrow$

$$\prod_{h \in h_1 \dots h_d} \underbrace{Pr \left[ \hat{x}_i^h - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right]}_{\leq \frac{1}{c}} \leq \frac{1}{c^d}$$



**Error Bounds**  
per hash/array



**Error Bounds**  
for the minimum



**Optimal Size**

$$Pr \left[ \underbrace{\min_{h \in h_1 \dots h_d} \hat{x}_i^h}_{\hat{x}_i} - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right] \leq \frac{1}{c^d}$$

**Error Bounds**  
per hash/array

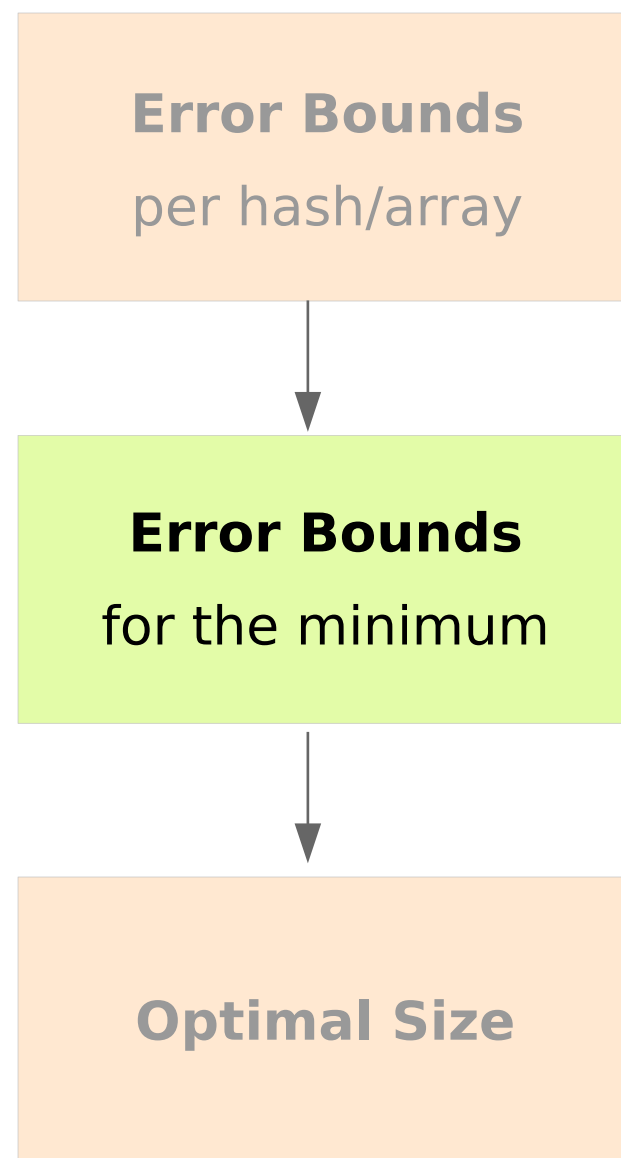


**Error Bounds**  
for the minimum



**Optimal Size**

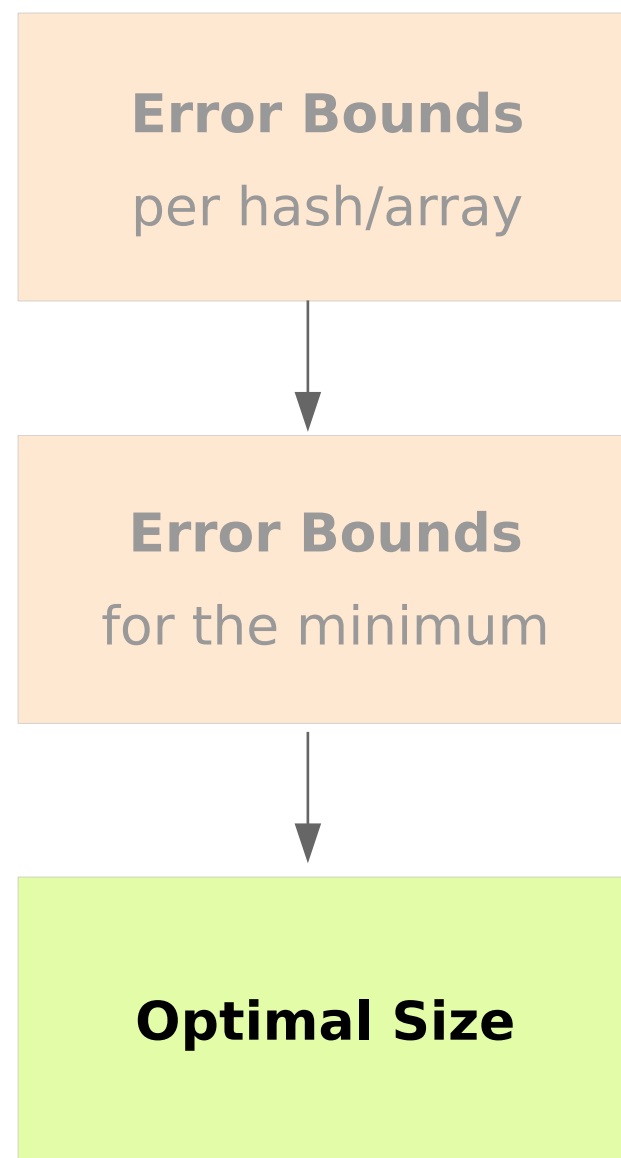
$$\Pr \left[ \hat{x}_i - x_i \geq \frac{c}{w} \|\mathbf{x}\|_1 \right] \leq \frac{1}{c^d}$$



$$\Pr \left[ \hat{x}_i - x_i \geq \underbrace{\epsilon}_{\frac{c}{w}} \mid \|\mathbf{x}\|_1 \leq \underbrace{\delta}_{\frac{1}{c^d}} \right] \leq \frac{1}{c^d}$$

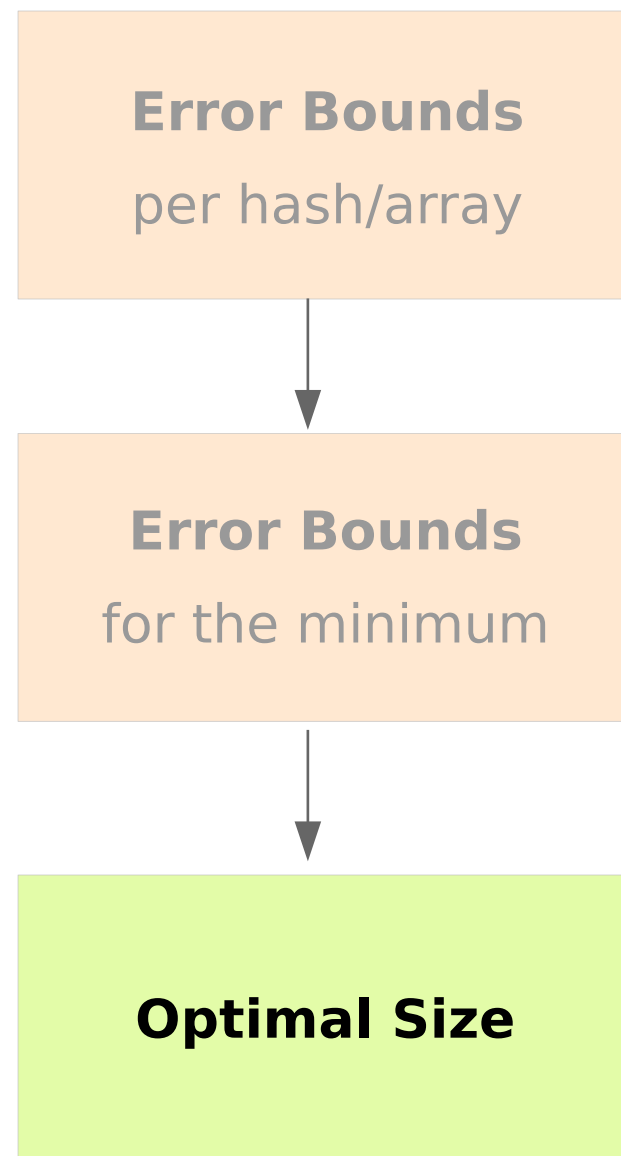
*We have proven the error bounds!*  
***But what about the constant  $c$ ?***

For **every**  $c$ , there is a pair  $(d, w)$  achieving the error bound and confidence  $(\varepsilon, \delta)$ .



$$\varepsilon = \frac{c}{w} \Rightarrow w = \left\lceil \frac{c}{\varepsilon} \right\rceil \quad (\text{hash range})$$
$$\delta = \frac{1}{c^d} \Rightarrow d = \left\lceil \log_c \frac{1}{\delta} \right\rceil \quad (\text{\#hashes})$$

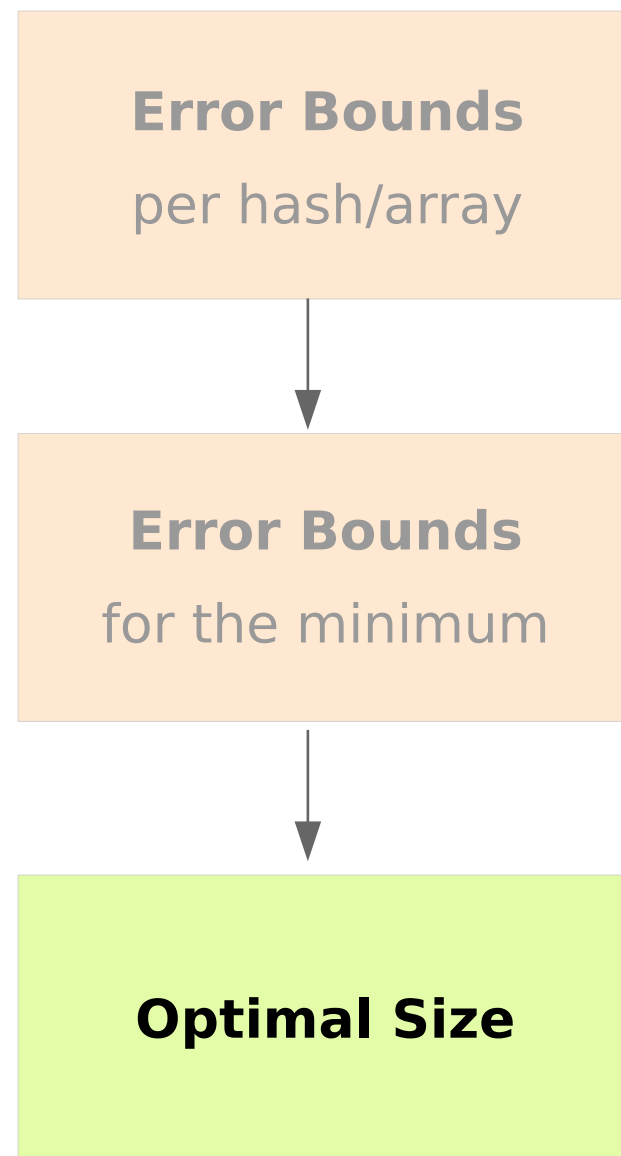
Choosing  $c=e$  **minimizes** the total **number of counters**.



$$\varepsilon = \frac{e}{w} \Rightarrow w = \left\lceil \frac{e}{\varepsilon} \right\rceil \quad (\text{hash range})$$
$$\delta = \frac{1}{e^d} \Rightarrow d = \left\lceil \ln \frac{1}{\delta} \right\rceil \quad (\text{\#hashes})$$

$$d \cdot w = \frac{c}{\varepsilon} \log_c \frac{1}{\delta} \stackrel{\text{minimize}}{=} \frac{e}{\varepsilon} \ln \frac{1}{\delta}$$

# A **CountMin** sketch recipe



**Given  $\varepsilon, \delta$  , choosing**

$$w = \left\lceil \frac{e}{\varepsilon} \right\rceil \quad (\text{hash range})$$

$$d = \left\lceil \ln \frac{1}{\delta} \right\rceil \quad (\text{\#hashes})$$

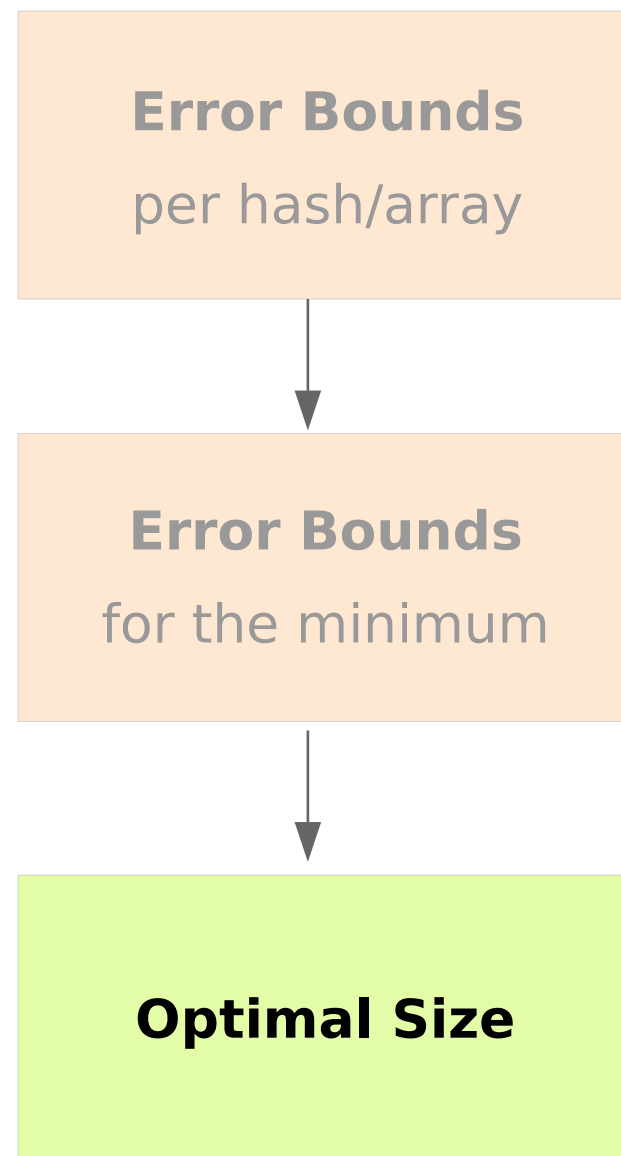
*requires the **minimum number of counters** s.t. the CountMin Sketch can guarantee that*

$$\hat{x}_i - x_i \geq \varepsilon \|\mathbf{x}\|_1$$

*with a probability less than  $\delta$*

# A CountMin sketch recipe

(see extended slides for derivation)



**Given  $\varepsilon, \delta$  , choosing**

$$w = \left\lceil \frac{e}{\varepsilon} \right\rceil \quad (\text{hash range})$$

$$d = \left\lceil \ln \frac{1}{\delta} \right\rceil \quad (\text{\#hashes})$$

requires the **minimum number of counters** s.t. the CountMin Sketch can guarantee that

$$\hat{x}_i - x_i \geq \varepsilon \|\mathbf{x}\|_1$$

*with a probability less than  $\delta$*

A **CountMin sketch** uses the same principles as a counting bloom filter, but is **designed** to have **provable L1 error bounds** for frequency queries.

***CountMin sketch recipe***

**Choose**  $d = \left\lceil \ln \frac{1}{\delta} \right\rceil$ ,  $w = \left\lceil \frac{e}{\epsilon} \right\rceil$

**Then**  $\hat{x}_i - x_i \geq \epsilon \|\mathbf{x}\|_1$  with a probability less than  $\delta$



A **CountMin sketch** uses the same principles as a counting bloom filter, but is **designed** to have **provable L1 error bounds** for frequency queries.

→ only one design out of many!

A **Count sketch** uses the same principles as a counting bloom filter, but is **designed** to have **provable L2 error bounds** for frequency queries.

The Count sketch uses **additional hashing** to give **L2 error bounds**, but requires more **resources**.

### CountMin sketch

$h_1, \dots, h_d: U \rightarrow \{1, \dots, w\}$

**COUNT  $x_i$ :**

for  $h$  in  $h_1, \dots, h_d$ :

$\text{Reg}_h[h(x_i)] + 1$

**QUERY  $x_i$ :**

return  $\min_{h \text{ in } h_1, \dots, h_d} ($

$\text{Reg}_h[h(x_i)]$

$)$

The Count sketch uses **additional hashing** to give **L2 error bounds**, but requires more **resources**.

### CountMin sketch

$h_1, \dots, h_d: U \rightarrow \{1, \dots, w\}$

**COUNT  $x_i$ :**

for  $h$  in  $h_1, \dots, h_d$ :

$\text{Reg}_h[h(x_i)] + 1$

**QUERY  $x_i$ :**

return  $\min_{h \text{ in } h_1, \dots, h_d} ($

$\text{Reg}_h[h(x_i)]$

$)$

### Count sketch

$h_1, \dots, h_d: U \rightarrow \{1, \dots, w\}$

$g: U \rightarrow \{+1, -1\}$

**COUNT  $x_i$ :**

for  $h$  in  $h_1, \dots, h_d$ :

$\text{Reg}_h[h(x_i)] + g(x_i)$

**QUERY  $x_i$ :**

return  $\text{median}_{h \text{ in } h_1, \dots, h_d} ($

$\text{Reg}_h[h(x_i)] * g(x_i)$

$)$

The Count sketch uses **additional hashing** to give **L2 error bounds**, but requires more **resources**.

***CountMin sketch recipe***

***Choose***  $d = \left\lceil \ln \frac{1}{\delta} \right\rceil, w = \left\lceil \frac{e}{\epsilon} \right\rceil$

***Then***  $\hat{x}_i - x_i \geq \epsilon \|\mathbf{x}\|_1$  with a probability less than  $\delta$

The Count sketch uses **additional hashing** to give **L2 error bounds**, but requires more **resources**.

***CountMin sketch recipe***

**Choose**  $d = \left\lceil \ln \frac{1}{\delta} \right\rceil, w = \left\lceil \frac{e}{\epsilon} \right\rceil$

**Then**  $\hat{x}_i - x_i \geq \epsilon \|\mathbf{x}\|_1$  with a probability less than  $\delta$

***Count sketch recipe***

**Choose**  $d = \left\lceil \ln \frac{1}{\delta} \right\rceil, w = \left\lceil \frac{e}{\epsilon^2} \right\rceil$

**Then**  $\hat{x}_i - x_i \geq \epsilon \|\mathbf{x}\|_2$  with a probability less than  $\delta$

# Sketches are the new black

## OpenSketch

NSDI '13

[source]



## UnivMon

SIGCOMM '16

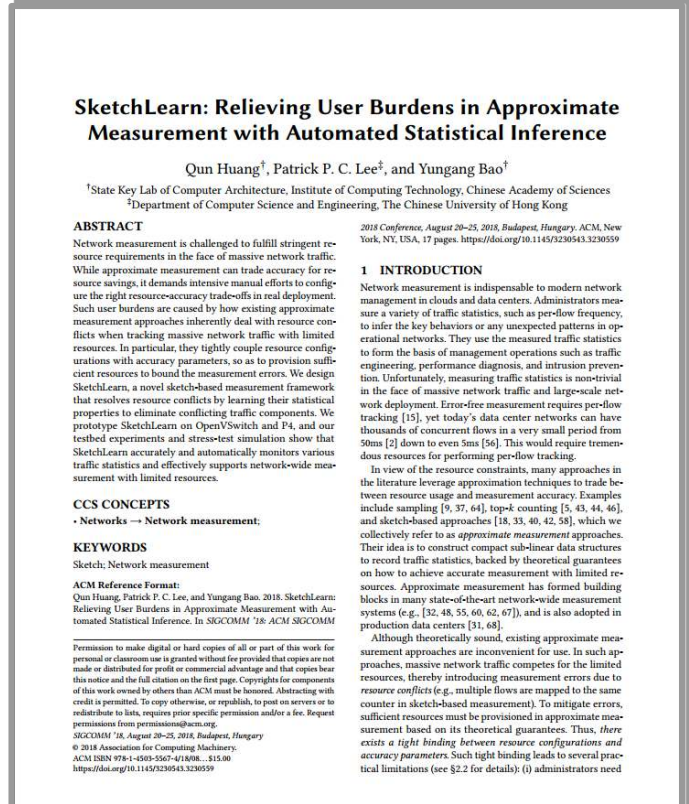
[source]



## SketchLearn

SIGCOMM '18

[source]





# Sketches are the new black

## LightGuardian

### NSDI '21

[source]

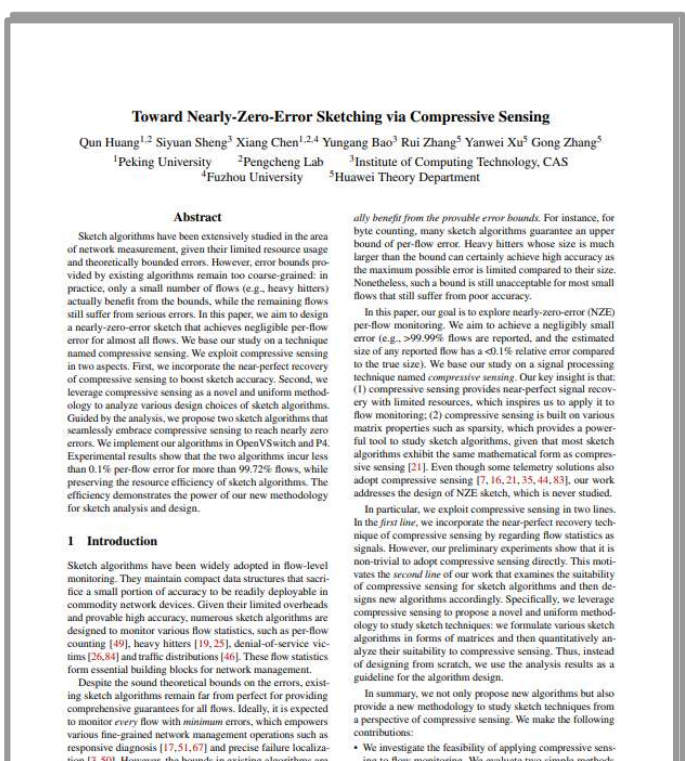


<sup>1</sup>Co-primary authors: Yikai Zhao, Kaicheng Yang, and Zirui Liu. Corresponding authors: Tong Yang (yangtong@gmail.com) and Yi Wang

## Nearly-Zero-Error

### NSDI '21

[source]



ally benefit from the provable error bounds. For instance, for byte counting, many sketch algorithms guarantee an upper bound of per-flow error. Heavy hitters whose size is much larger than the bound can certainly achieve high accuracy as the maximum possible error is limited compared to their size. Nonetheless, such a bound is still unacceptable for most small flows that still suffer from poor accuracy.

In this paper, our goal is to explore nearly-zero-error (NZE) per-flow monitoring. We aim to achieve a negligibly small error (e.g., >99.99% flows are reported, and the estimated size of any reported flow has a <0.1% relative error compared to the true size). We base our study on a signal processing technique named *compressive sensing*. Our key insight is that: (1) compressive sensing provides near-perfect signal recovery with limited resources, which inspires us to apply it to flow monitoring; (2) compressive sensing is built on various matrix properties such as sparsity, which provides a powerful tool to study sketch algorithms, given that most sketch algorithms exhibit the same mathematical form as compressive sensing [21]. Even though some telemetry solutions also adopt compressive sensing [7, 16, 21, 35, 44, 83], our work addresses the design of NZE sketch, which is never studied.

In particular, we exploit compressive sensing in two lines. In the *first line*, we incorporate the near-perfect recovery technique of compressive sensing by regarding flow statistics as signals. However, our preliminary experiments show that it is non-trivial to adopt compressive sensing directly. This motivates the *second line* of our work that examines the suitability of compressive sensing for sketch algorithms and then designs new algorithms accordingly. Specifically, we leverage compressive sensing to propose a novel and uniform methodology to study sketch techniques: we formulate various sketch algorithms in forms of matrices and then quantitatively analyze their suitability to compressive sensing. Thus, instead of designing from scratch, we use the analysis results as a guideline for the algorithm design.

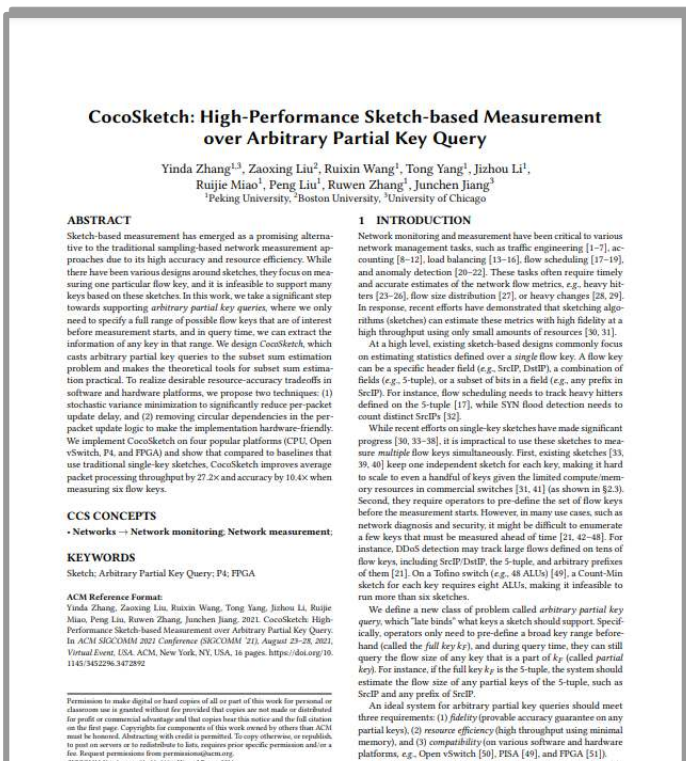
In summary, we not only propose new algorithms but also provide a new methodology to study sketch techniques from a perspective of compressive sensing. We make the following contributions:

- We investigate the feasibility of applying compressive sensing to flow monitoring. We evaluate two simple methods and show that simple utilization either suffers from poor scalability or fails to reach the expected accuracy level.

## CocoSketch

### SIGCOMM '21

[source]



Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
SIGCOMM '21, August 23–28, 2021, Virtual Event, USA.  
© 2021 Association for Computing Machinery.  
ACM ISBN 978-1-4503-8387-2/21/08...\$15.00.  
<https://doi.org/10.1145/342298.347292>



Sketches have limitations. They **compute statistical summaries** and favor elements with **high frequency**.

$$Pr \left[ \hat{x}_i - x_i \geq \varepsilon \|\mathbf{x}\|_1 \right] \leq \delta$$

*estimation error*      *relative to sum of all elements*

Sketches have limitations. They **compute statistical summaries** and favor elements with **high frequency**.

*Let*  $\varepsilon = 0.01$ ,  $\|\mathbf{x}\|_1 = 10000$  ( $\Rightarrow \varepsilon \cdot \|\mathbf{x}\|_1 = 100$ )

*Assume two flows*  $x_a$ ,  $x_b$ ,

*with*  $\|x_a\|_1 = 1000$ ,  $\|x_b\|_1 = 50$

*high frequency*

*low frequency*

Sketches have limitations. They **compute statistical summaries** and favor elements with **high frequency**.

*Let*  $\varepsilon = 0.01$ ,  $\|\mathbf{x}\|_1 = 10000$  ( $\Rightarrow \varepsilon \cdot \|\mathbf{x}\|_1 = 100$ )

*Assume two flows*  $x_a$ ,  $x_b$ ,

*with*  $\|x_a\|_1 = 1000$ ,  $\|x_b\|_1 = 50$

**Error relative to stream size: 1%**

Sketches have limitations. They **compute statistical summaries** and favor elements with **high frequency**.

*Let*  $\varepsilon = 0.01$  ,  $\|\mathbf{x}\|_1 = 10000$  ( $\Rightarrow \varepsilon \cdot \|\mathbf{x}\|_1 = 100$ )

*Assume two flows*  $x_a$  ,  $x_b$  ,

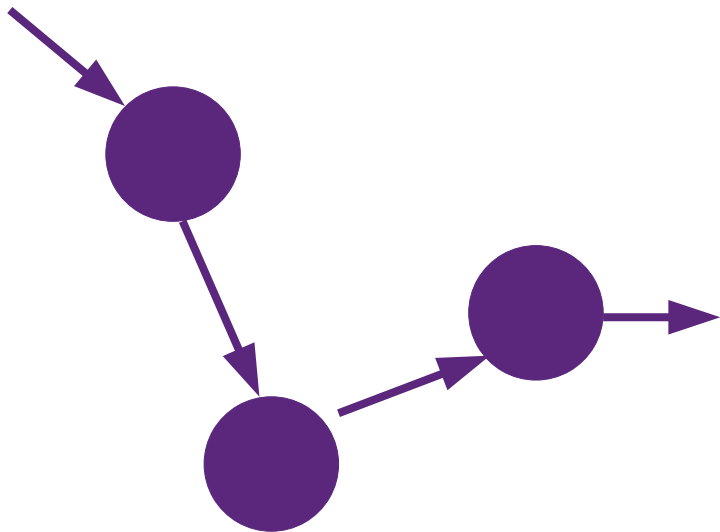
*with*  $\|x_a\|_1 = 1000$  ,  $\|x_b\|_1 = 50$

**Error relative to stream size: 1%**

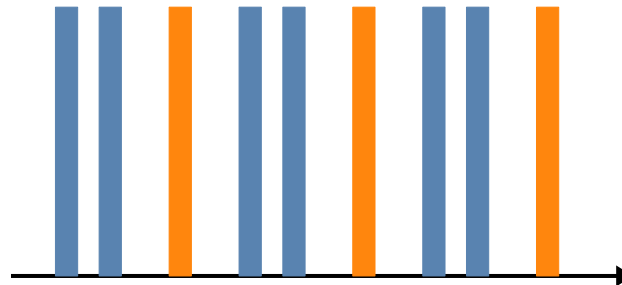
**flow size:**  $x_a$ : 10%,  $x_b$ : **200%**

## Other problems a sketch **can't** handle:

*causality*




*patterns*



*rare things*





← You are looking at a stream of data (packets).  
Today, I'll show you how set membership and frequency queries can be realized in P4.

**PART 1**

Is a certain element (e.g. ip address) in the stream?

→ Bloom filter

**PART 2**

How frequently does an element appear?

→ CountMin Sketch, Count Sketch, ...

**TAKEAWAY**

Probabilistic data structures provide  
**trade-offs between resources and error, and**  
**provable guarantees** to rely on.