# Data Structures and Algorithms
# Bloom Filters

CS 225
Brad Solomon

UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN

Department of Computer Science

Why store

K,v
K,v

when you
can store

1
0

# Announcements

MP_mosaic survey EC reached

MP_traversal survey EC **not reached** (Have until 11/20 to submit!)

MP_puzzle released, due after break. **Break doesn't count as a week**

↳ lab this week

# Learning Objectives

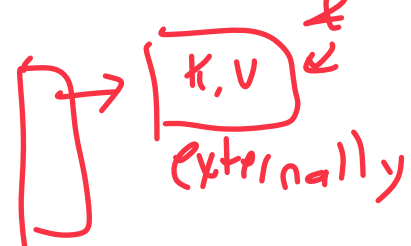Review when you would prefer different data structures

Build a conceptual understanding of a bloom filter

Review probabilistic data structures and one-sided error

Formalize the math behind the bloom filter

*(handwritten: ← review)*

*(handwritten: ↑ Wednesday)*

# Which collision resolution strategy is better?

- Big Records: **Open hashing (Seperate chaining)**
  - $\hookrightarrow$ Can pass by ref   $\hookrightarrow$ closed hashing
  - $\hookrightarrow$ $\boxed{K,V}$ externally
  - at this scale cant allocate

- Structure Speed: **Closed hashing (Double hashing)**

# What structure do hash tables implement?   Dictionaries

# What constraint exists on hashing that doesn't exist with BSTs?

① Probabilistic!

② Simple uniform hashing assumption (SUHA)

③ Pseudo-amortized

# Why talk about BSTs at all?   $\hookrightarrow$ B/c resize when $\alpha < 1$

$\hookrightarrow$ Some data is not hashable   $\hookrightarrow$ ordered dataset useful (Nearest neighbor)

# Running Times



Handwritten: → Exact lookup only

Handwritten: Trees in general → Approx match (KD Tree)

| | Hash Table | AVL | Linked List |
|---|---|---|---|
| **Find** | Expectation*: O(1)*** <br><br> Worst Case: O(n) <br><br> *(Expectation)* | O(log n) | O(n) |
| **Insert** | Expectation*: O(1)*** <br><br> Worst Case: O(n) | *(Guaranteed)* O(log n) | O(1) |
| **Storage Space** | O(n) | O(n) | O(n) |

# Memory-Constrained Data Structures

What method would you use to build a search index on a collection of objects *in a memory-constrained environment*?
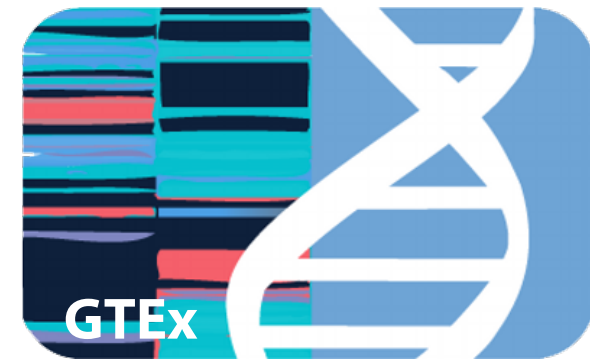
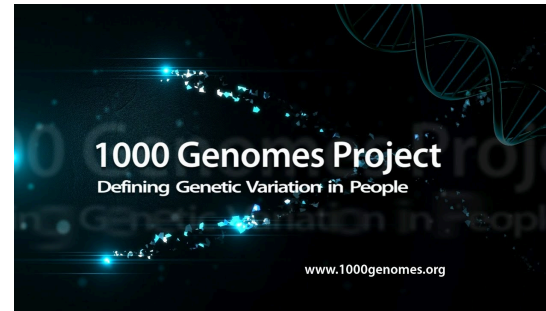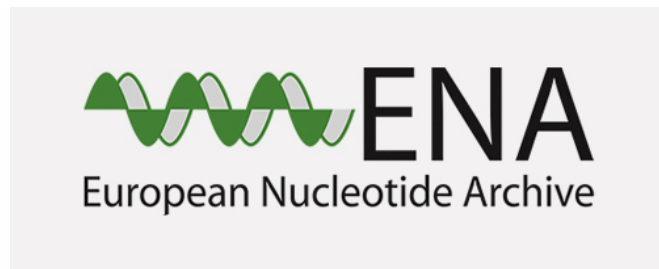**Constrained by Big Data (Large $N$)**



Google Index Estimate: >60 billion webpages

Google Universe Estimate (2013): >130 trillion webpages

# Memory-Constrained Data Structures

What method would you use to build a search index on a collection of objects *in a memory-constrained environment*?

**Constrained by Big Data (Large $N$)**



Sequence Read Archive Size: >60 petabases ($10^{15}$)

# Memory-Constrained Data Structures

What method would you use to build a search index on a collection of objects *in a memory-constrained environment*?
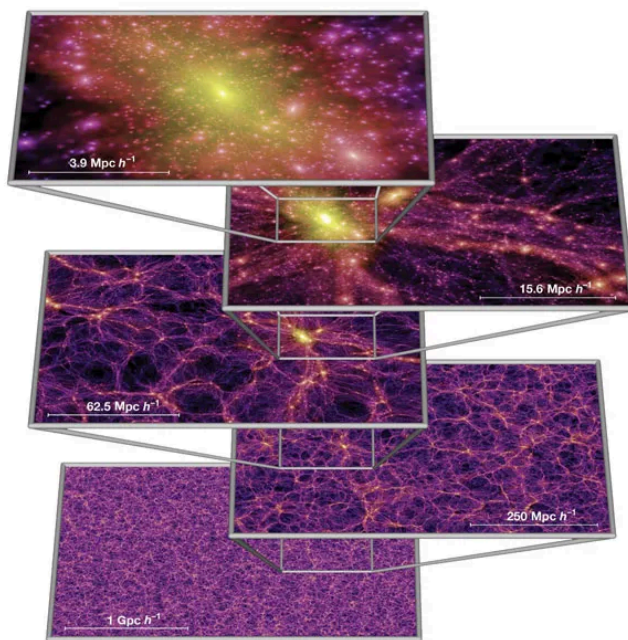
## Constrained by Big Data (Large $N$)



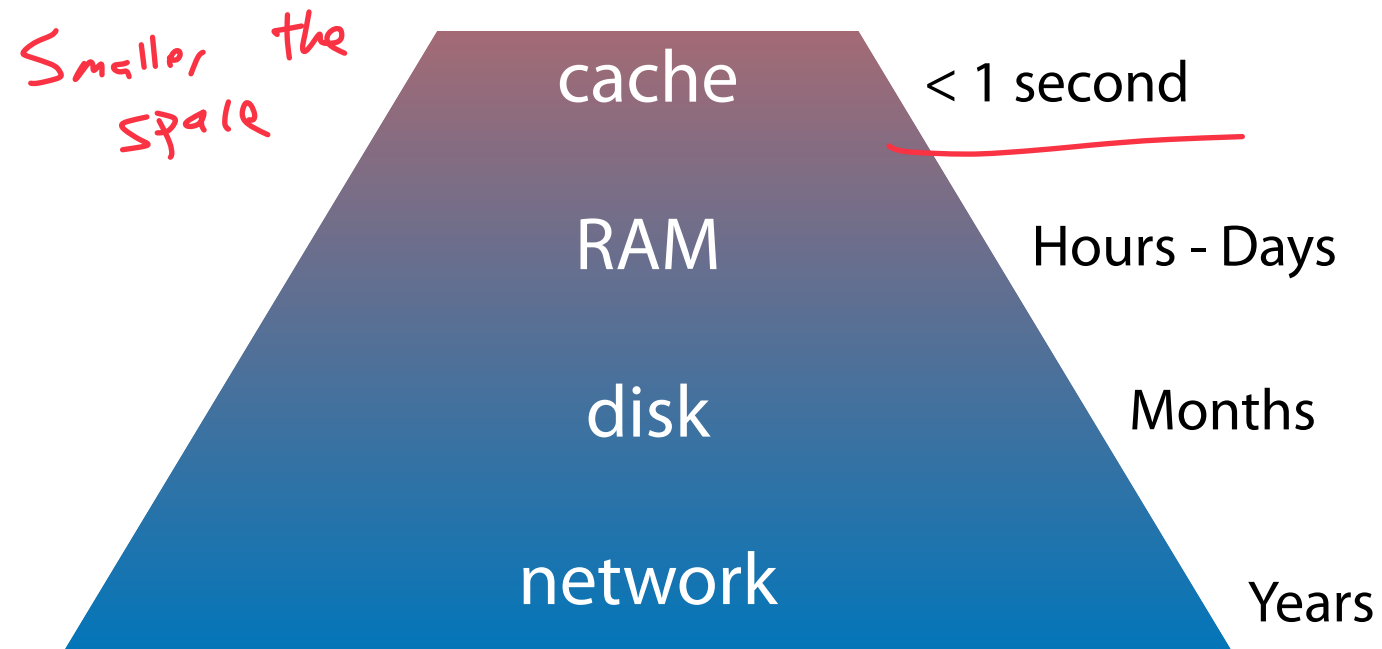| Sky Survey Projects | Data Volume |
| --- | --- |
| DPOSS (The Palomar Digital Sky Survey) | 3 TB |
| 2MASS (The Two Micron All-Sky Survey) | 10 TB |
| GBT (Green Bank Telescope) | 20 PB |
| GALEX (The Galaxy Evolution Explorer ) | 30 TB |
| SDSS (The Sloan Digital Sky Survey) | 40 TB |
| SkyMapper Southern Sky Survey | 500 TB |
| PanSTARRS (The Panoramic Survey Telescope and Rapid Response System) | ~ 40 PB expected |
| LSST (The Large Synoptic Survey Telescope) | ~ 200 PB expected |
| SKA (The Square Kilometer Array) | ~ 4.6 EB expected |

Table: http://doi.org/10.5334/dsj-2015-011

Estimated total volume of one array: 4.6 EB

Image: https://doi.org/10.1038/nature03597

# Memory-Constrained Data Structures

What method would you use to build a search index on a collection of objects *in a memory-constrained environment*?

**Constrained by resource limitations**

Smaller the space

| | |
|---|---|
| cache | < 1 second |
| RAM | Hours - Days |
| disk | Months |
| network | Years |

(Estimates are Time x 1 billion courtesy of https://gist.github.com/hellerbarde/2843375)

# Memory-Constrained Data Structures

What method would you use to build a search index on a collection of objects *in a memory-constrained environment*?

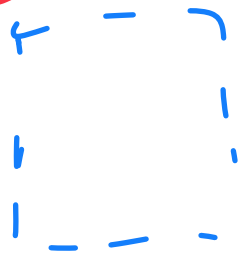1) Make a space-efficient encoding (Compress the data)

2) Throw away data we dont care about

3) Make a distributed network

Make a bloom filter

# Reducing storage costs

100 billion

1) Throw out information that isn't needed

# of sides

Length of perimeter :  $\times$

| 1 | 2 2 | 2 1 |
|---|-----|-----|
| square | diamond | triangle |
| 4 | 4 | 3 |

$\times$

Use data stream to process each item individually

2) Compress the dataset " A A A G G G "

A3   G3   $\equiv$   run length encoding

# Reducing a hash table

What can we remove from a hash table?

$$H(k_1) = i_1$$

# Reducing a hash table

What can we remove from a hash table?

Take away values

$$H(k_1) = i_1$$

# Reducing a hash table

What can we remove from a hash table?

Take away values and keys

$$H(k_1) = i_1$$

$m$

Something hashed here at some time

# Reducing a hash table

What can we remove from a hash table?

Take away values and keys

This is a **bloom filter**

$$H(k_1) = i_1$$



← Something hashed here

← Nothing hashed here

$m$

$O(m)$ Storage (Small constant 1 bit per block)

# Bloom Filter ADT

**Constructor**

**Insert**

**Find**

# Bloom Filter: Insertion

S = { 16, 8, 4, 13, 29, 11, 22 }

h(k) = k % 7

| | |
|---|---|
| 0 | 0 |
| 1 | 0 1 |
| 2 | 0 1 |
| 3 | 0 |
| 4 | 0 1 |
| 5 | 0 |
| 6 | 0 1 |

← h(29) = 1

Empty BF is vector/array of 0s

1) Hash key to hash value (address)

2) Set bit to 1 at address

No collision possible

1 says "Something hashed at some point"

# Bloom Filter: Insertion

An item is inserted into a bloom filter by hashing and then setting the hash-valued bit to 1

If the bit was already one, it stays 1

↳Each 1 is One or more

inserts

$H(x_1)$

$H(x_2)$

$H(x_3)$

$H(x_4)$

| 0 |
|---|
| 0 |
| 1 |
| 0 |
| 0 |
| 1 |
| 0 |
| 1 |
| 0 |
| 0 |

# Bloom Filter: Deletion

**S = { 16, 8, 4, 13, 29, 11, 22 }**

**h(k) = k % 7**

| | |
|---|---|
| 0 | 0 |
| 1 | ~~1~~ 0 |
| 2 | 1 |
| 3 | 0 |
| 4 | 1 |
| 5 | 0 |
| 6 | ~~1~~ 0 |

Oh no!

**_delete(13)**

1) Hash key

2) Set bit to 0

**_delete(29)**

_find(8)

# Bloom Filter: Deletion

Due to hash collisions and lack of information, items cannot be deleted!
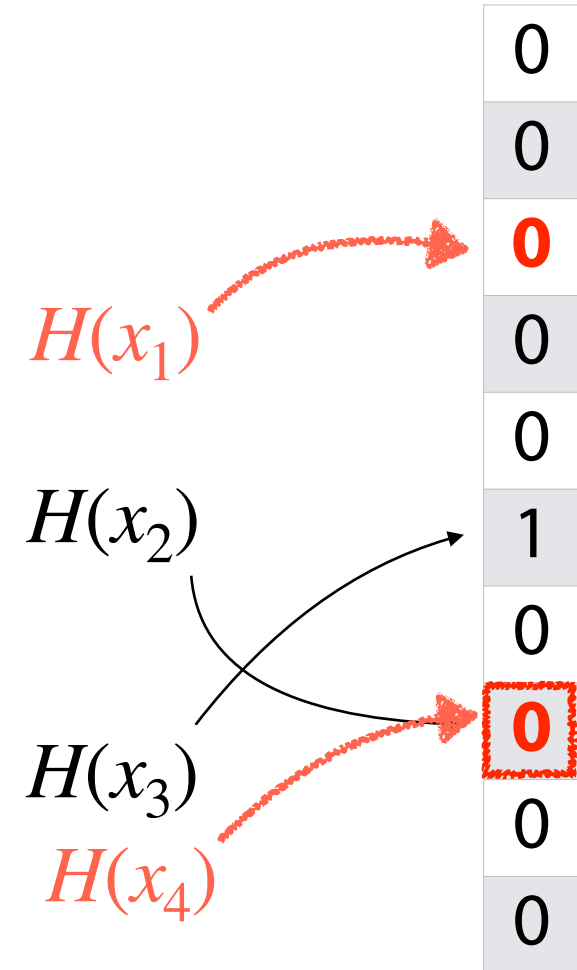
① We never know if it is safe to set 1 to 0.

② We need the property that 1 stays 1 forever

$H(x_1)$

$H(x_2)$

$H(x_3)$

$H(x_4)$

| 0 |
| 0 |
| **0** |
| 0 |
| 0 |
| 1 |
| 0 |
| **0** |
| 0 |
| 0 |

# Bloom Filter: Search

**S = { 16, 8, 4, 13, 29, 11, 22 }**

↑
**h(k) = k % 7**   6

| 0 | 0 |
|---|---|
| 1 | 1 |
| 2 | 1 |
| 3 | 0 |
| 4 | 1 |
| 5 | 0 |
| 6 | 1 |

_find(16) → true! ✓ correct

1) Hash key
2) Look up value

_find(20) → True ✗ False
  ↳ 20 % 7 = 6

_find(3) → False ✓ correct

Probabilistic accuracy!

# Bloom Filter: Search

The bloom filter is a *probabilistic* data structure!

If the value in the BF is 0:

100% of the time item not in set

If the value in the BF is 1:

item might be in set

or

is a hash collision

$H(\alpha)$

| |
|---|
| 0 |
| **0** |
| 1 |
| 0 |
| 0 |
| **1** |
| 0 |
| **1** |
| 0 |
| 0 |

$H(x_1)$

$H(\beta)$

$H(x_2)$

$H(x_3)$

$H(x_4)$

$H(\delta)$

# Probabilistic Accuracy: Malicious Websites

Imagine we have a detection oracle that identifies if a site is malicious



"Not malicious"



"Malicious"

# Probabilistic Accuracy: Malicious Websites

Imagine we have a detection oracle that identifies if a site is malicious
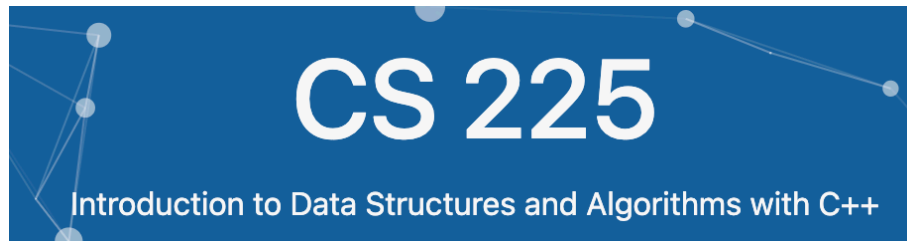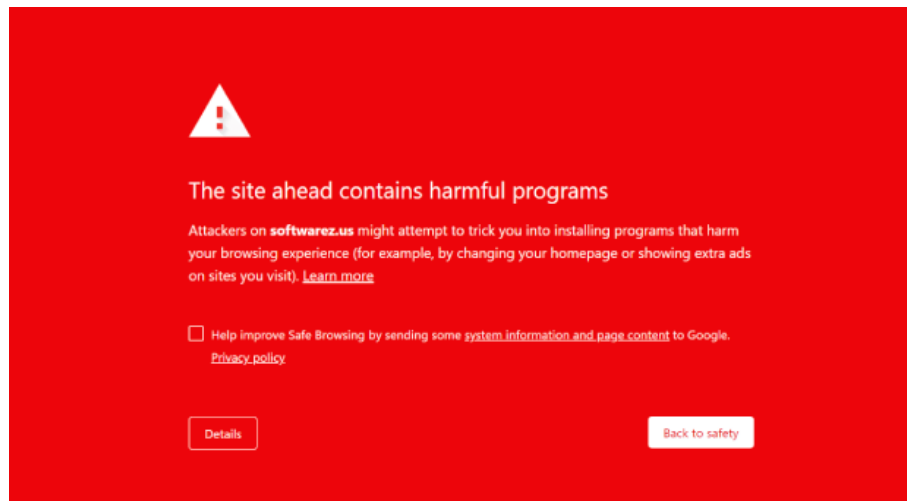
True Positive: Oracle says M | actual website is M

False Positive: Oracle says M | actual not M

False Negative: not M | actual M

True Negative: not M | not M

# Imagine we have a **bloom filter** that **stores malicious sites…**



|  | **Bit Value = 1** | **Bit Value = 0** |
|---|---|---|
| **Item Inserted** | $H(z)$ → 0 1 0 0 1 'Yes' True Positive | $H(z)$ → 0 0 0 0 1 'No' False Negative — Not possible (was inserted) |
| **Item NOT inserted** | 0 1 0 0 1 'Yes' False Positive | 0 0 0 0 1 'No' True Negative |

# Probabilistic Accuracy: One-sided error

**Query:**

**Dataset:**

search with one-sided error

We will get some False Positives: =

We will NEVER have a False Negative: ≠

Will remove some bad but no good items

# Probabilistic Accuracy: One-sided error

# Bloom Filter: Repeated Trials

Use many hashes/filters; add each item to each filter

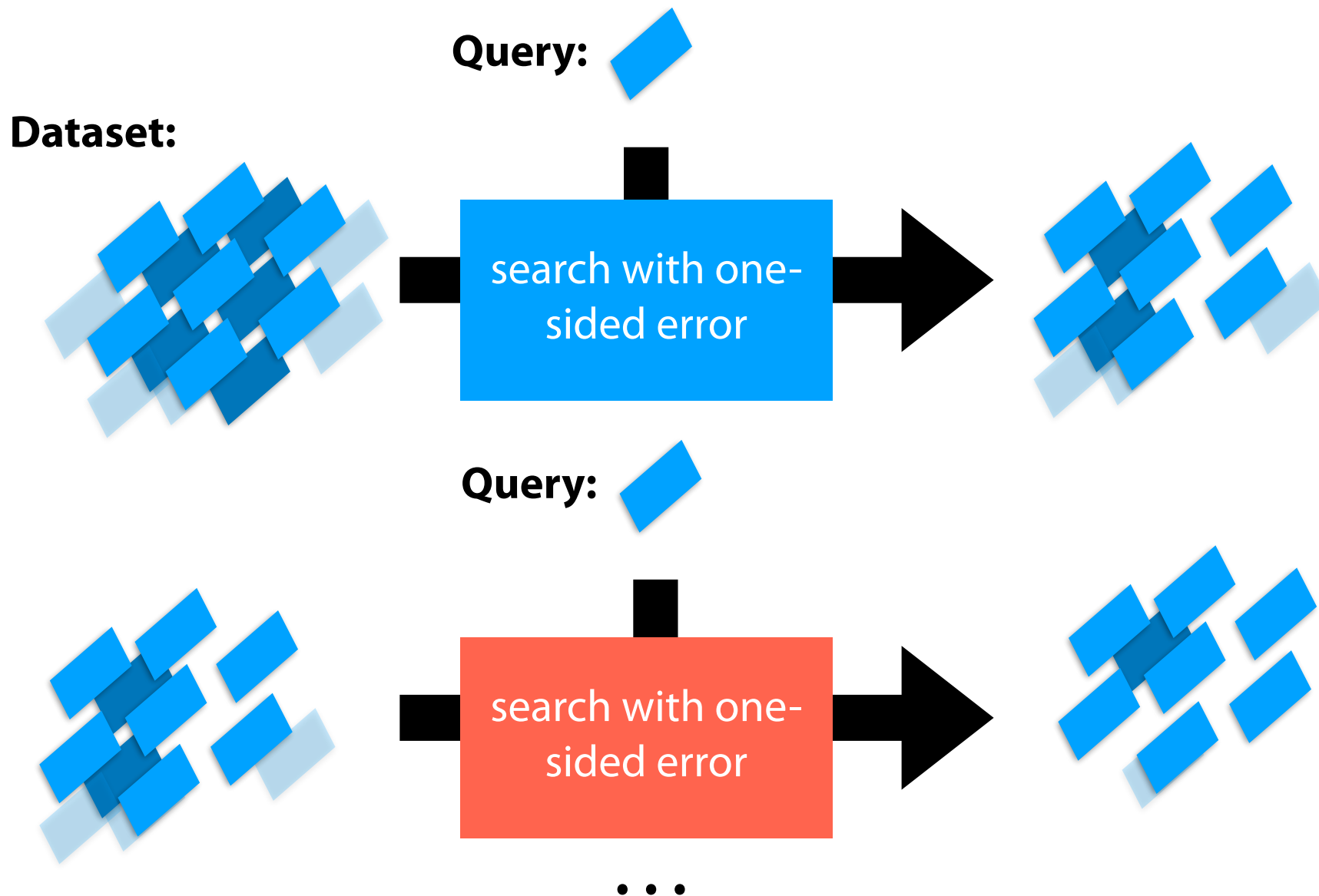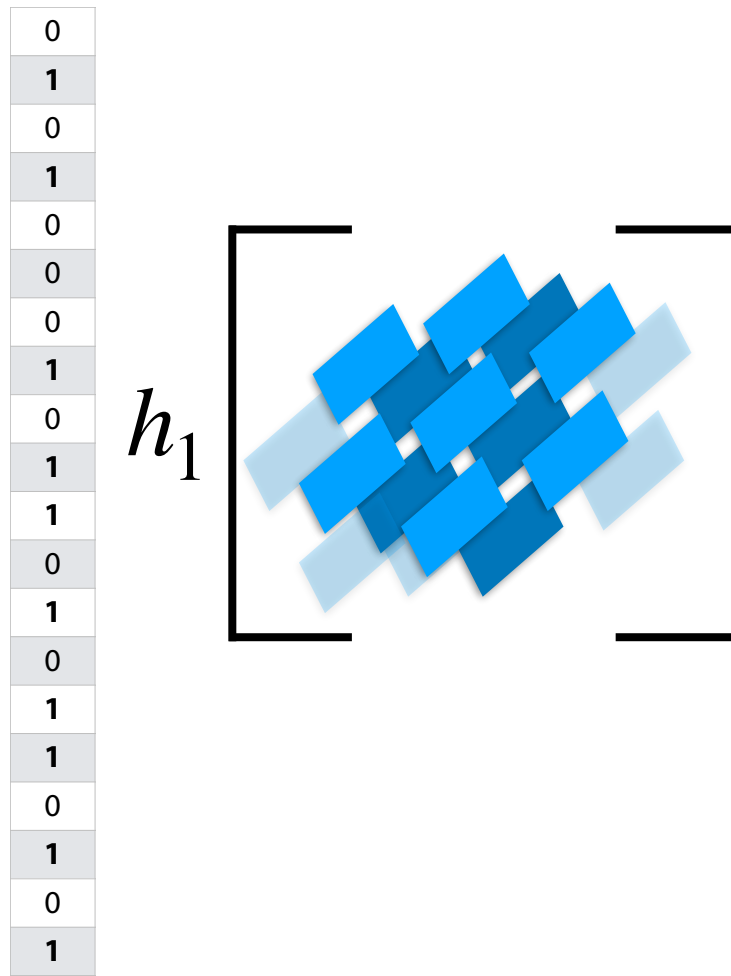| |
|---|
| 0 |
| 1 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 1 |
| 0 |
| 1 |

$h_1$

# Bloom Filter: Repeated Trials

Use many hashes/filters; add each item to each filter

$h_1$

| |
|---|
| 0 |
| 1 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 1 |
| 0 |
| 1 |

$h_2$

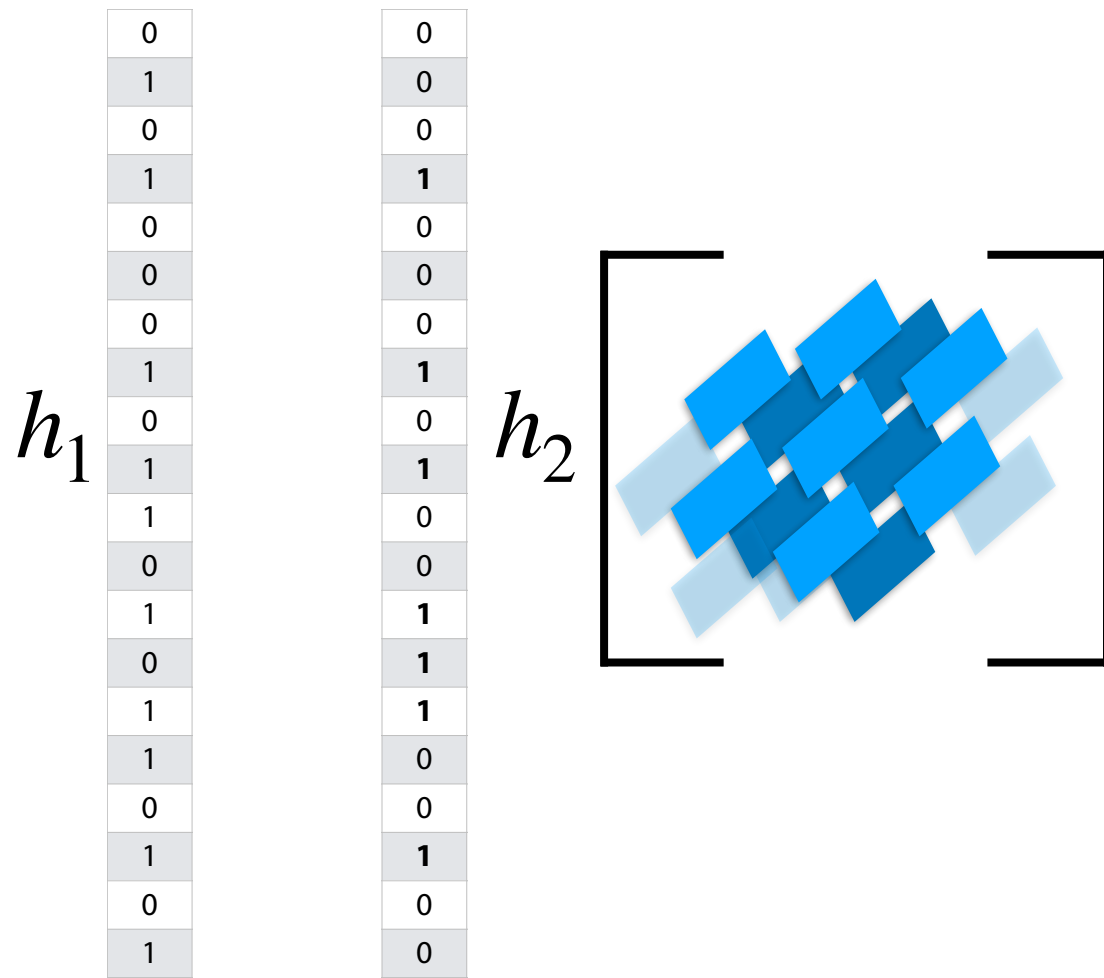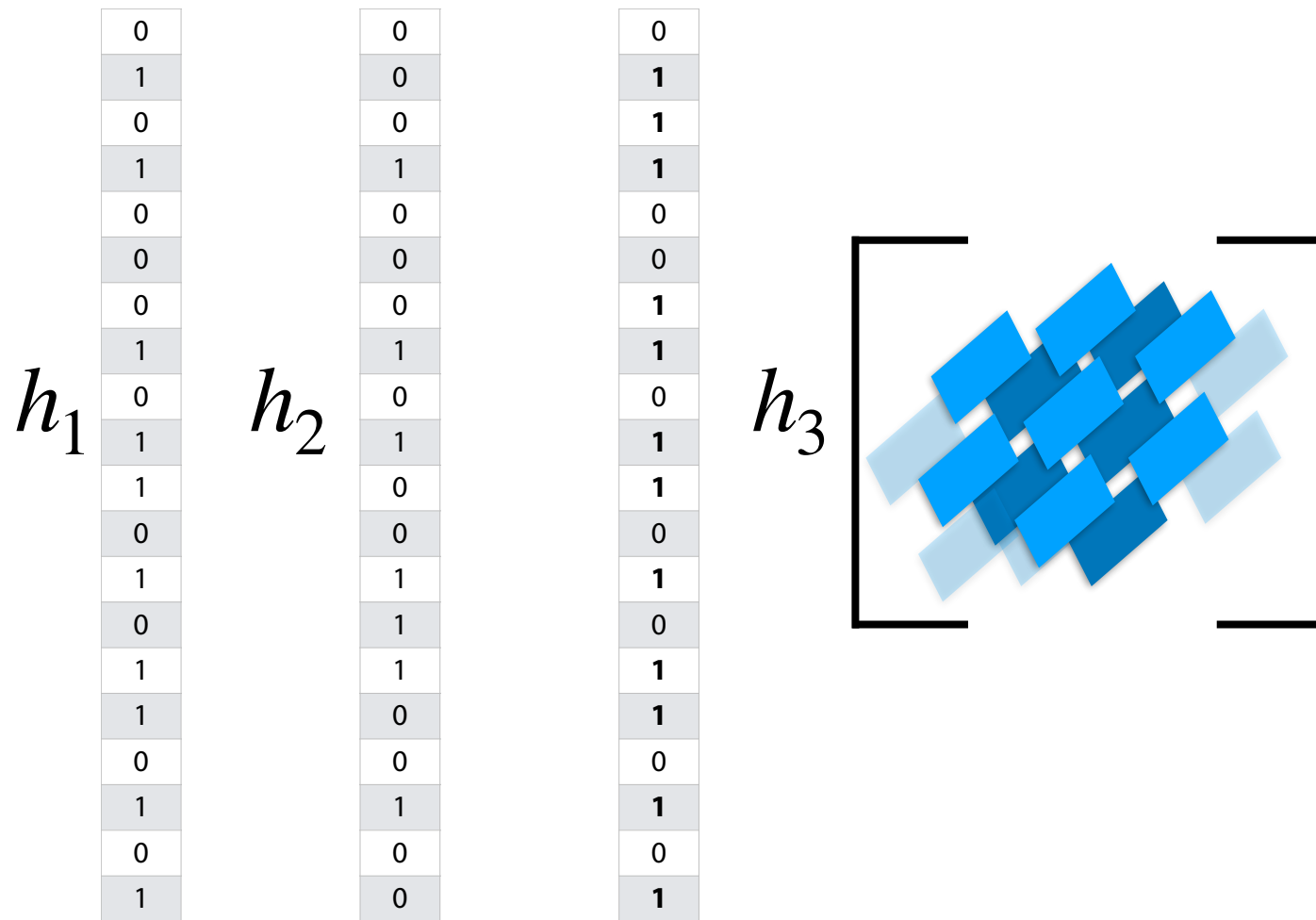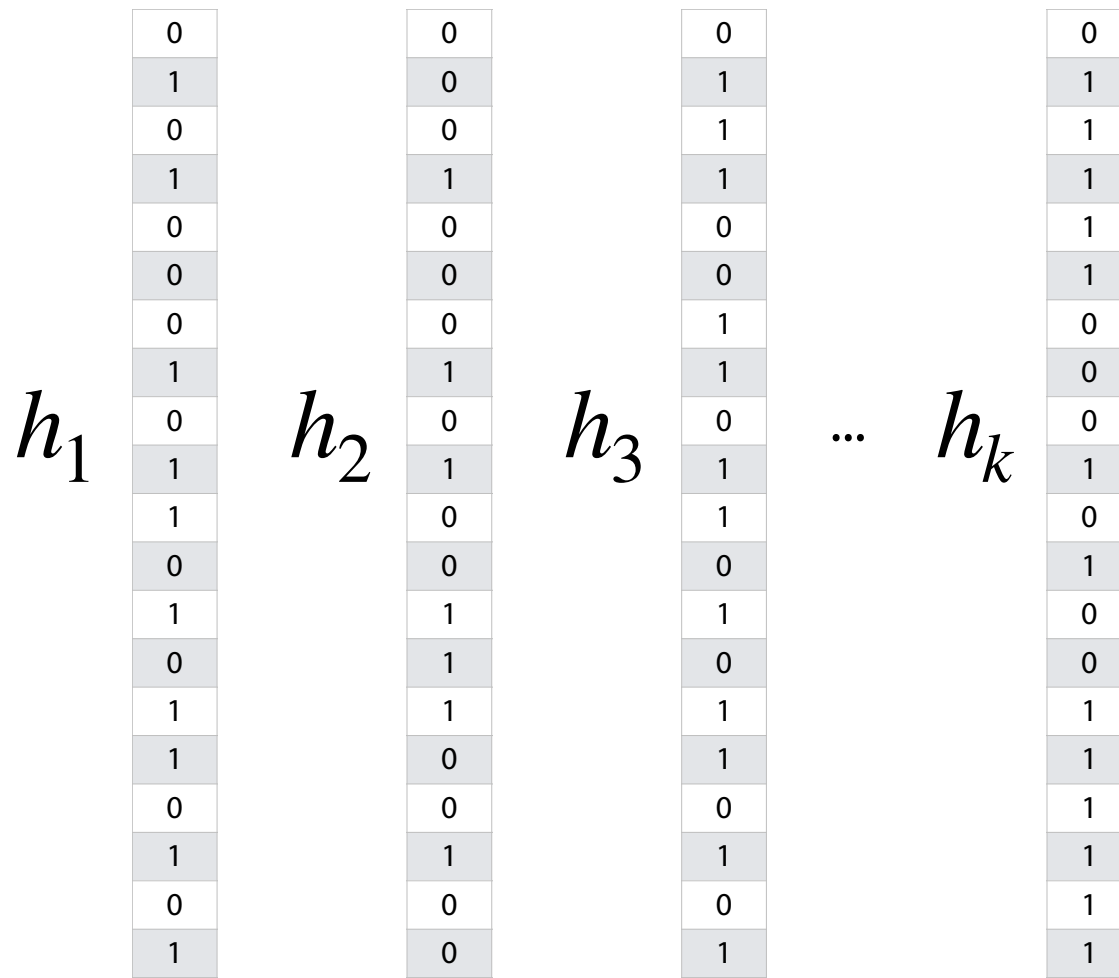| |
|---|
| 0 |
| 0 |
| 0 |
| **1** |
| 0 |
| 0 |
| 0 |
| **1** |
| 0 |
| **1** |
| 0 |
| 0 |
| **1** |
| **1** |
| **1** |
| 0 |
| 0 |
| **1** |
| 0 |
| 0 |

# Bloom Filter: Repeated Trials

Use many hashes/filters; add each item to each filter

# Bloom Filter: Repeated Trials

Use many hashes/filters; add each item to each filter

# Bloom Filter: Repeated Trials

| | | | | | | |
|---|---|---|---|---|---|---|
| 0 | | 0 | | 0 | | 0 |
| 1 | | 0 | | 1 | | 1 |
| 0 | | 0 | | 1 | | 1 |
| 1 | | 1 | | 1 | | 1 |
| 0 | | 0 | | 0 | | 1 |
| 0 | | 0 | | 0 | | 1 |
| 0 | | 0 | | 1 | | 0 |
| 1 | | 1 | | 1 | | 0 |
| 0 | | 0 | | 0 | | 0 |
| 1 | | 1 | | 1 | | 1 |
| 1 | | 0 | | 1 | | 0 |
| 0 | | 0 | | 0 | | 1 |
| 1 | | 1 | | 1 | | 0 |
| 0 | | 1 | | 0 | | 0 |
| 1 | | 1 | | 1 | | 1 |
| 1 | | 0 | | 1 | | 1 |
| 0 | | 0 | | 0 | | 1 |
| 1 | | 1 | | 1 | | 1 |
| 0 | | 0 | | 0 | | 1 |
| 1 | | 0 | | 1 | | 1 |

$$h_{\{1,2,3,...,k\}}(y)$$

# Bloom Filter: Repeated Trials



Do k find()
ops

$$h_{\{1,2,3,\ldots,k\}}(y)$$

If *any* query yields 0,
item is not in the set

random  collision

# Bloom Filter: Repeated Trials

$$h_{\{1,2,3,\ldots,k\}}(z)$$

If *all* queries yield 1, item *may* be in the set; or we might have collided *k* times

# Bloom Filter: Repeated Trials

Using repeated trials, even a very bad filter can still have a very low FPR!

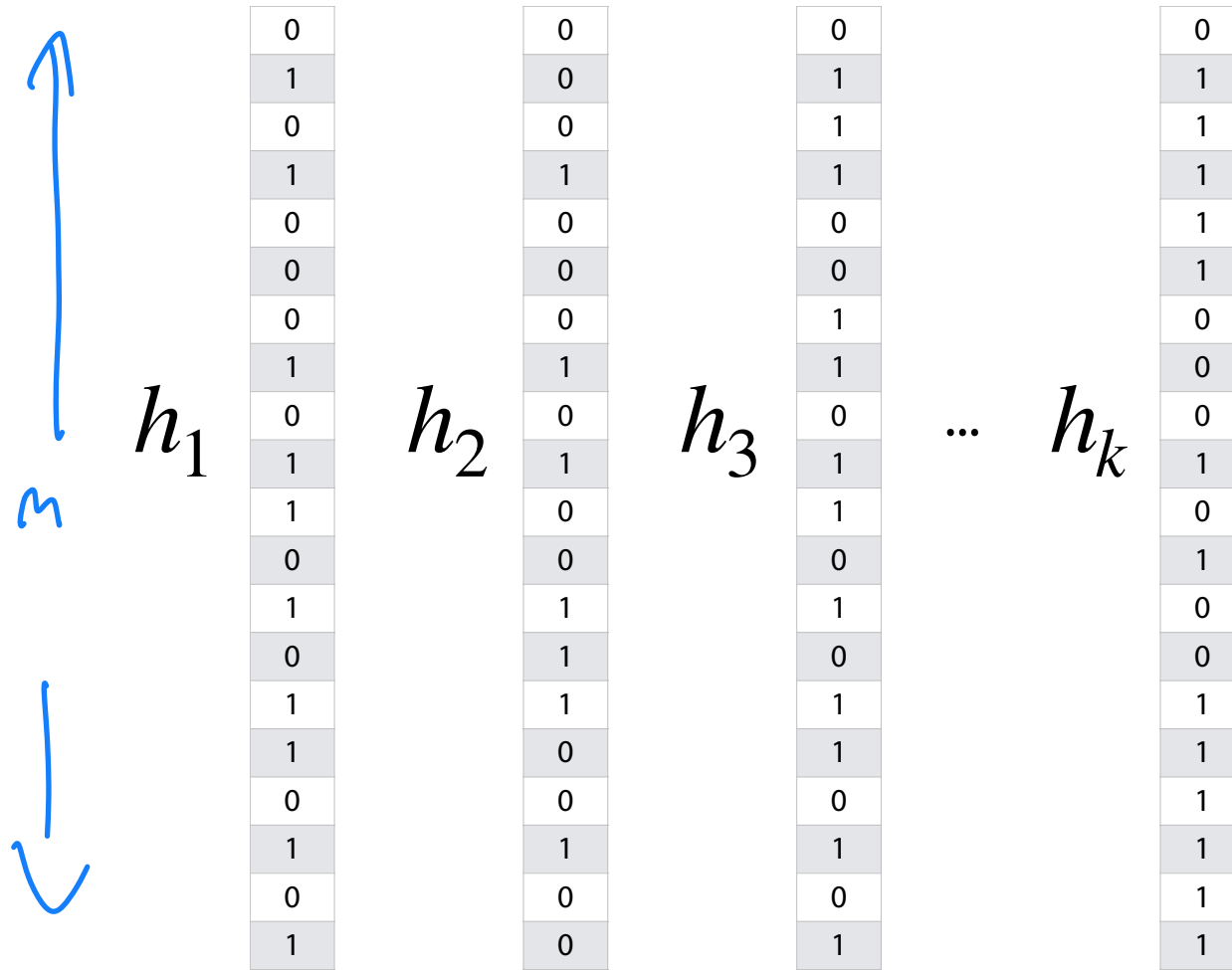If we have $k$ bloom filter, each with a FPR $p$, what is the likelihood that **all** filters return the value '1' for an item we didn't insert?

$$p^k \longrightarrow P = 0.5 \qquad (0.5)^{10} \approx 0.00097$$

$k = 10$

# Bloom Filter: Repeated Trials

But doesn't this hurt our storage costs by storing $k$ separate filters?
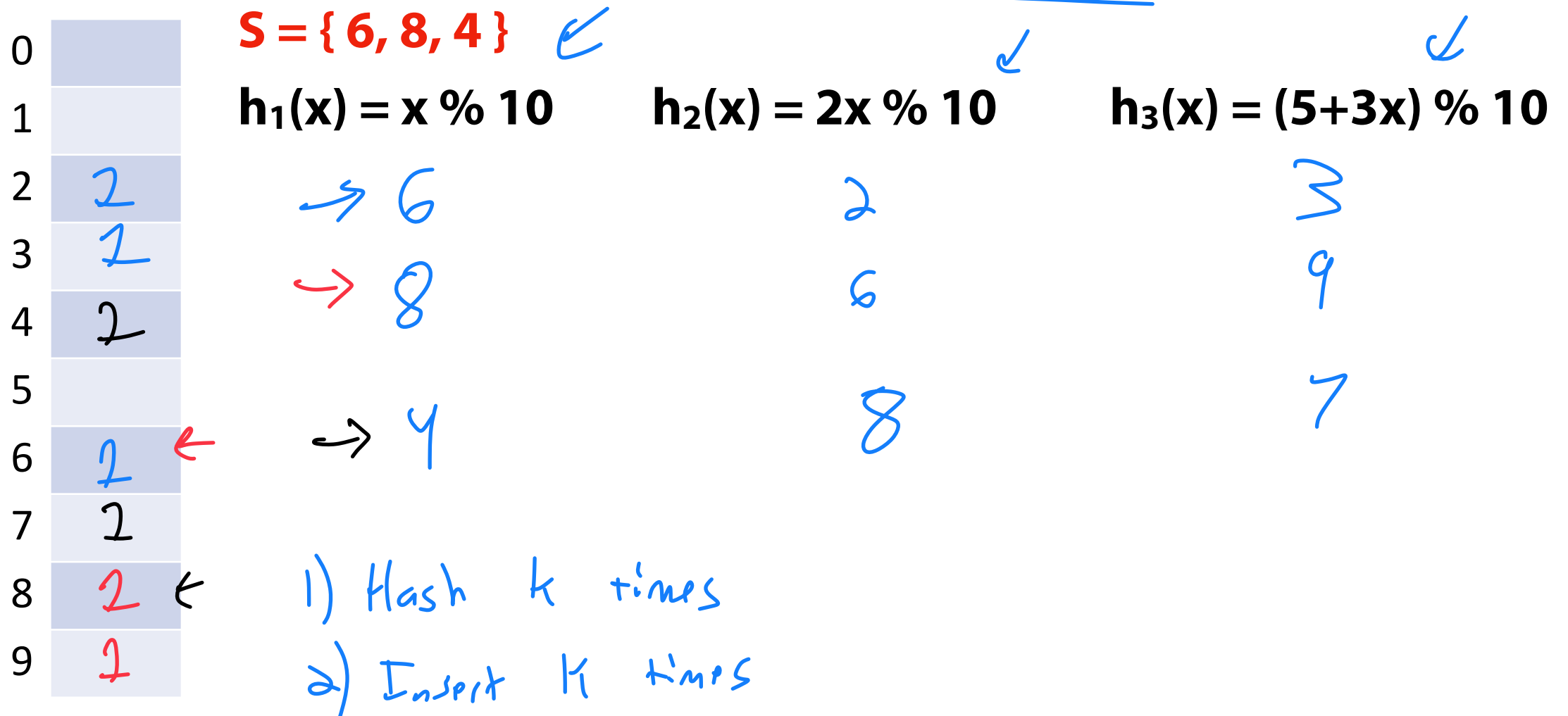
$h_1$  $h_2$  $h_3$  ...  $h_k$

$m$

current storage

$k \cdot m$

can be good

vs $(n)$

# Bloom Filter: Repeated Trials

Rather than use a new filter for each hash, one filter can use $k$ hashes

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | 2 |
| 3 | 1 |
| 4 | 2 |
| 5 | |
| 6 | 1 |
| 7 | 2 |
| 8 | 2 |
| 9 | 1 |

$S = \{ 6, 8, 4 \}$

$h_1(x) = x \% 10$   $h_2(x) = 2x \% 10$   $h_3(x) = (5+3x) \% 10$

$\rightarrow 6$         2           3
$\rightarrow 8$         6           9
                                    7
$\rightarrow 4$         8

1) Hash $k$ times

2) Insert $k$ times

# Bloom Filter: Repeated Trials

Rather than use a new filter for each hash, one filter can use $k$ hashes

| | |
|---|---|
| 0 | 0 |
| 1 | 0 |
| 2 | 1 |
| 3 | 1 |
| 4 | 1 |
| 5 | 0 |
| 6 | 1 |
| 7 | 1 |
| 8 | 1 |
| 9 | 1 |

$h_1(x) = x \% 10$　　$h_2(x) = 2x \% 10$　　$h_3(x) = (5+3x) \% 10$

_find(1)　　　2　　2　　8

1) Hash $k$ times
2) Do $k$ lookups

$$\left[ find(1) \Rightarrow False \right]$$

_find(16)　6　　　2　　　3

↳ True

(this is false positive!)

# Bloom Filter

$$H = \{h_1, h_2, \ldots, h_k\}$$

A probabilistic data structure storing a set of values

Built from a bit vector of length $m$ and $k$ hash functions

Insert / Find runs in: _____

Delete is not possible (yet)!

| |
|---|
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 1 |
| 0 |
| 1 |
| 0 |
| 0 |