# Data Structures and Algorithms
# Bloom Filters

CS 225
G Carl Evans

**UNIVERSITY OF ILLINOIS**
**URBANA-CHAMPAIGN**

Department of Computer Science

# Memory-Constrained Data Structures

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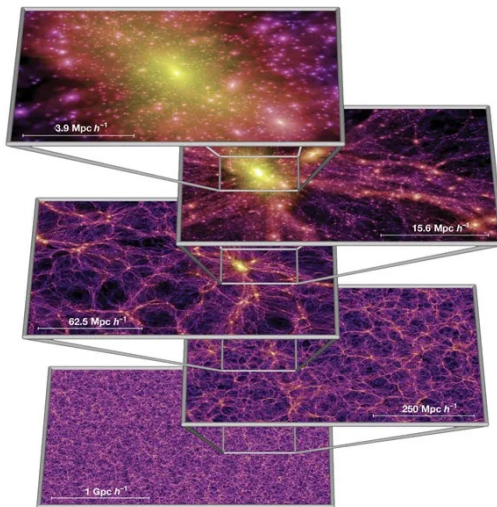
## Constrained by Big Data (Large $N$)



| 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

# Bloom Filter: Insertion

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

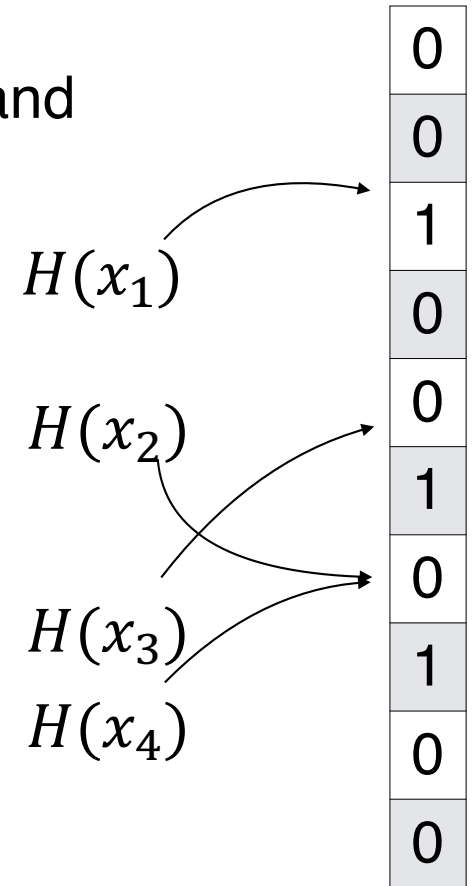**h(k) = k % 7**

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

# 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

$H(x_1)$

$H(x_2)$

$H(x_3)$

$H(x_4)$

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

# Bloom Filter: Search

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

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

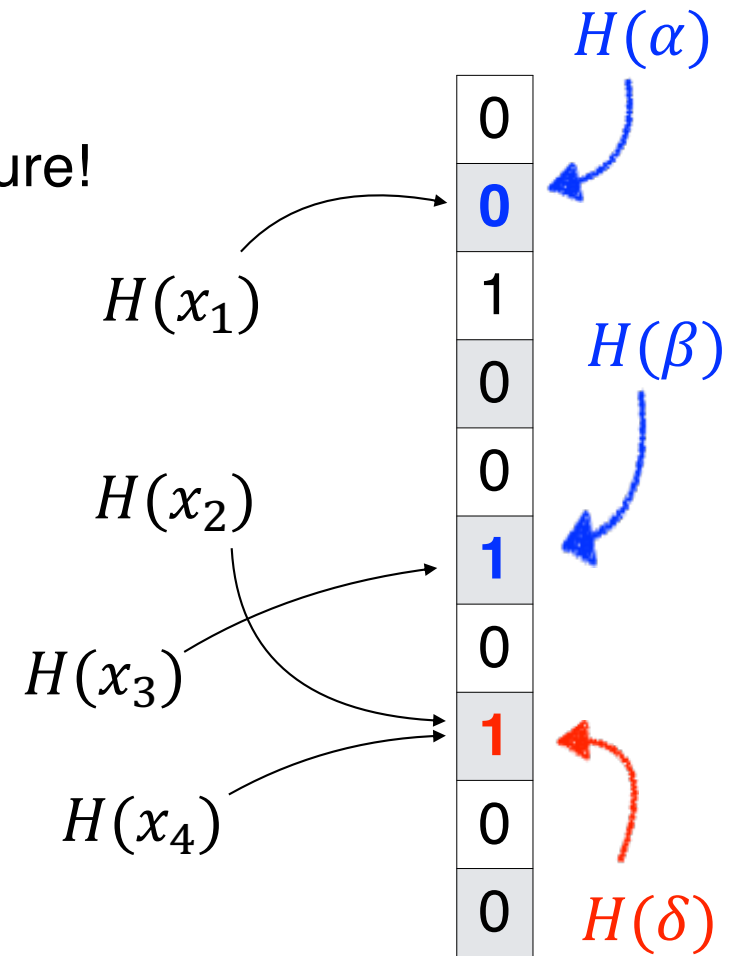| | |
|---|---|
| 0 | 0 |
| 1 | 1 |
| 2 | 1 |
| 3 | 0 |
| 4 | 1 |
| 5 | 0 |
| 6 | 1 |

**_find(20)**

**_find(3)**

# Bloom Filter: Search

The bloom filter is a *probabilistic* data structure!

If the value in the BF is 0:

If the value in the BF is 1:

$H(x_1)$

$H(x_2)$

$H(x_3)$

$H(x_4)$

$H(\alpha)$

$H(\beta)$

$H(\delta)$

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

# Probabilistic Accuracy: Malicious Websites

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



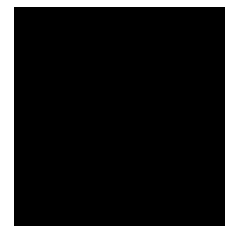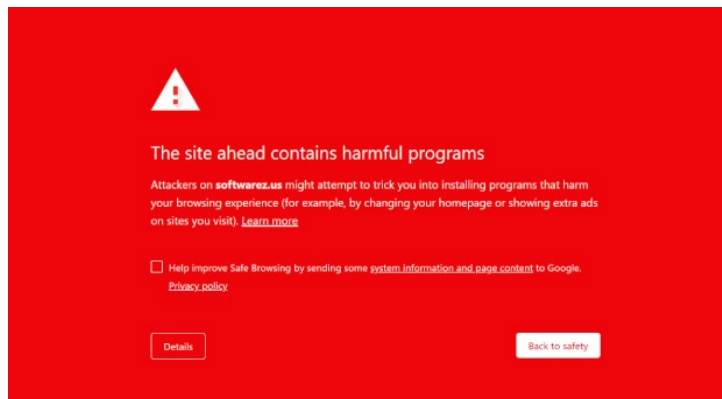"Not malicious"

"Malicious"

# Probabilistic Accuracy: Malicious Websites

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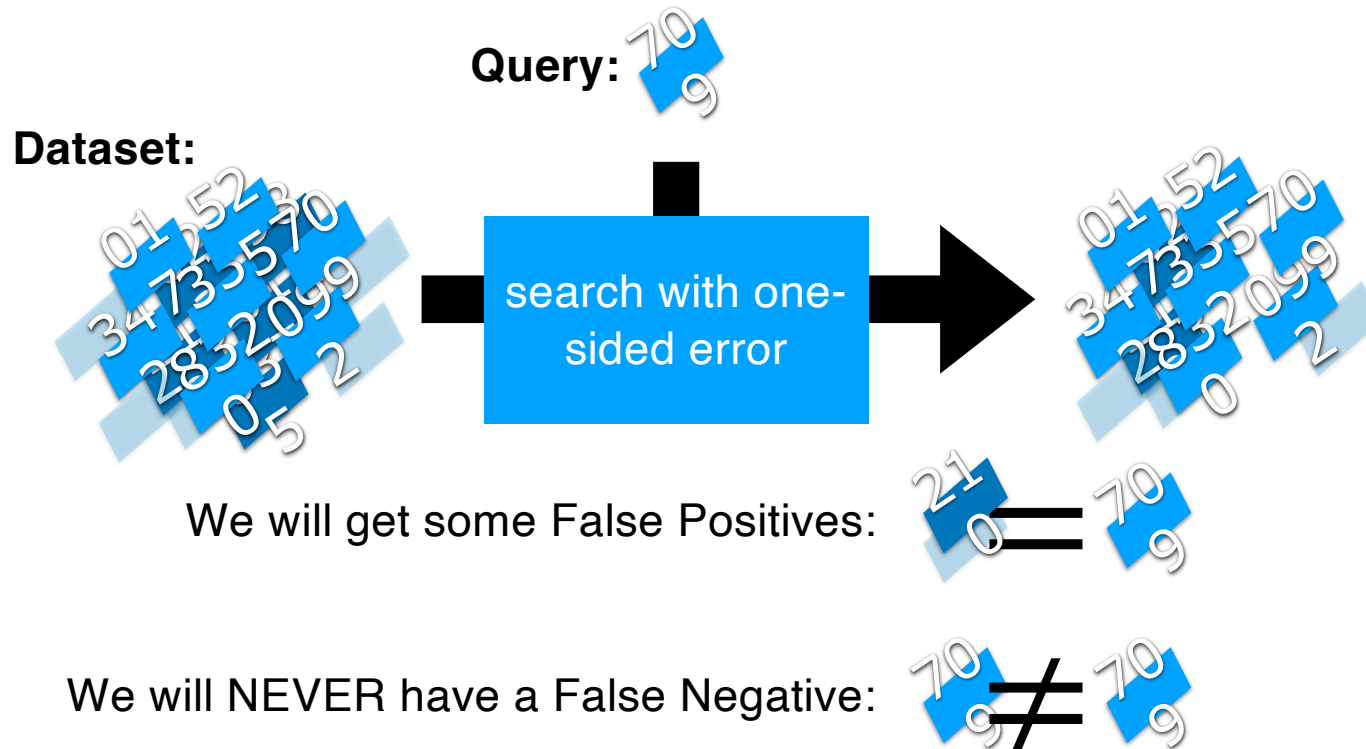
True Positive:

False Positive:
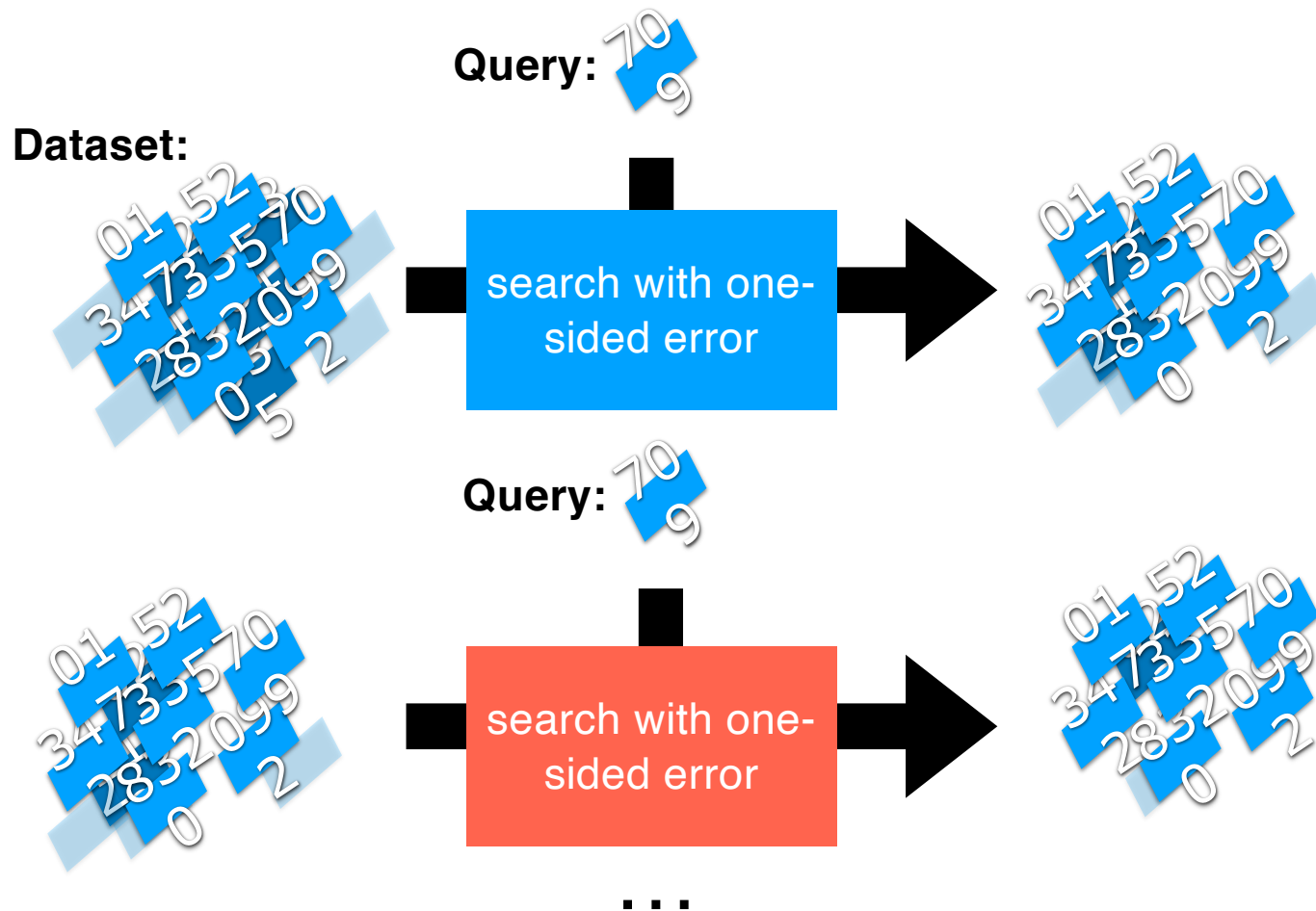
False Negative:

True Negative:

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

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

# Probabilistic Accuracy: One-sided error

# 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 |

# 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 |

$$\begin{bmatrix} 01 & 52 & 3 \\ 34 & 735 & 570 \\ 28 & 32 & 099 \\ 0 & 5 & 2 \end{bmatrix}$$

# 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

$h_1$ $h_2$ $h_3$ $\cdots$ $h_k$

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

# Bloom Filter: Repeated Trials



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

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

# Bloom Filter: Repeated Trials



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

If *all* queries yield 1, item *may* be in the set; or we might have collided *k* times
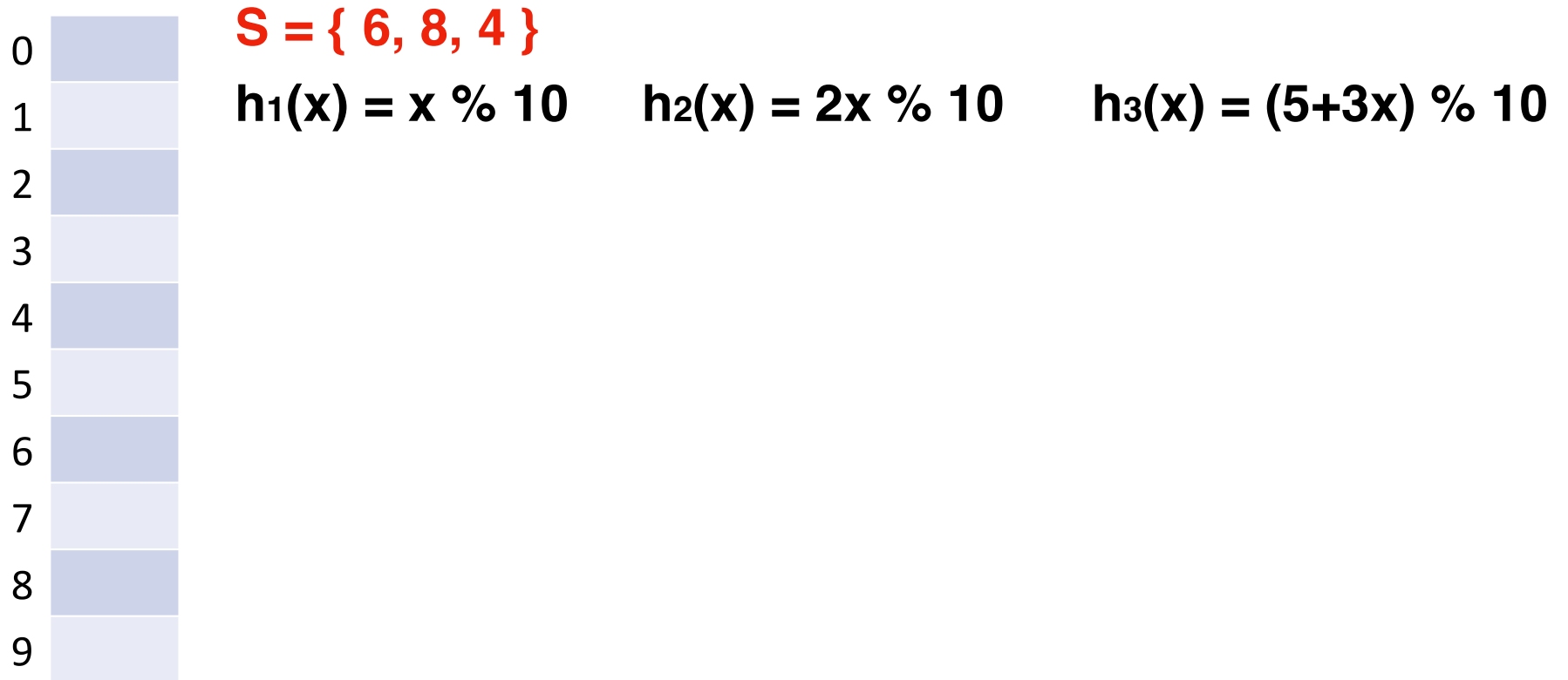
# Bloom Filter: Repeated Trials

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?

# Bloom Filter: Repeated Trials

But doesn't this hurt our storage costs by storing $k$ separate filters?

$$h_1 \quad h_2 \quad h_3 \quad \cdots \quad h_k$$

# Bloom Filter: Repeated Trials

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

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

S = { 6, 8, 4 }

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

# 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)

_find(16)

# Bloom Filter

A probabilistic data structure storing a set of values

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

Built from a bit vector of length $m$ and $k$ hash functions
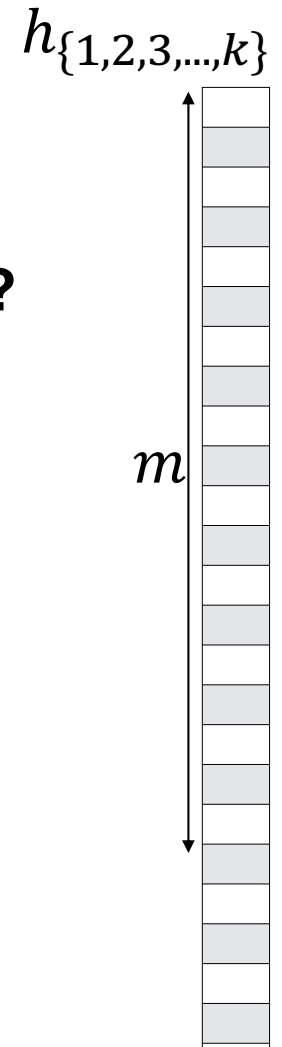
Insert / Find runs in: _____

Delete is not possible (yet)!

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

# Bloom Filter: Error Rate

Given bit vector of size $m$ and $k$ SUHA hash function

**What is our expected FPR after $n$ objects are inserted?**
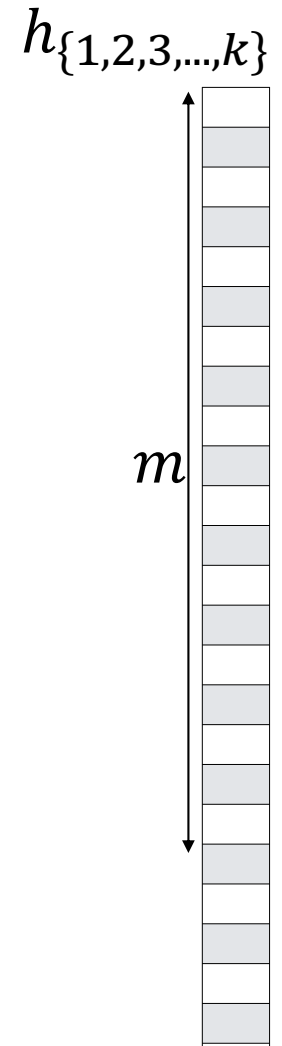
$h_{\{1,2,3,...,k\}}$

$m$

# Bloom Filter: Error Rate

Given bit vector of size $m$ and 1 SUHA hash function

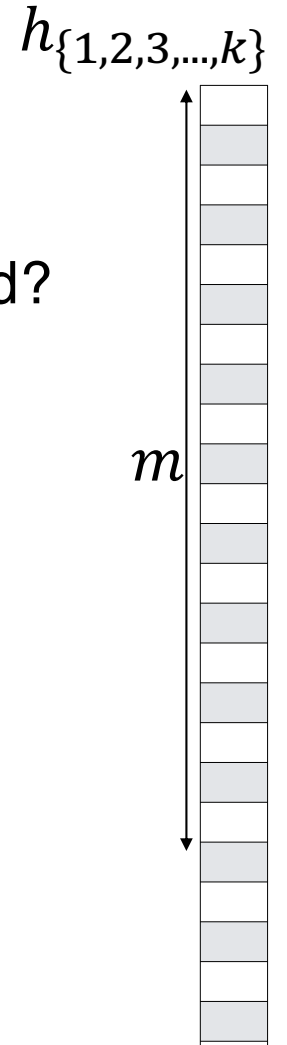What's the probability a specific bucket is 1 after one object is inserted?

Same probability given $k$ SUHA hash function?

$h_{\{1,2,3,...,k\}}$

$m$

# Bloom Filter: Error Rate

$h_{\{1,2,3,...,k\}}$

Given bit vector of size $m$ and $k$ SUHA hash function

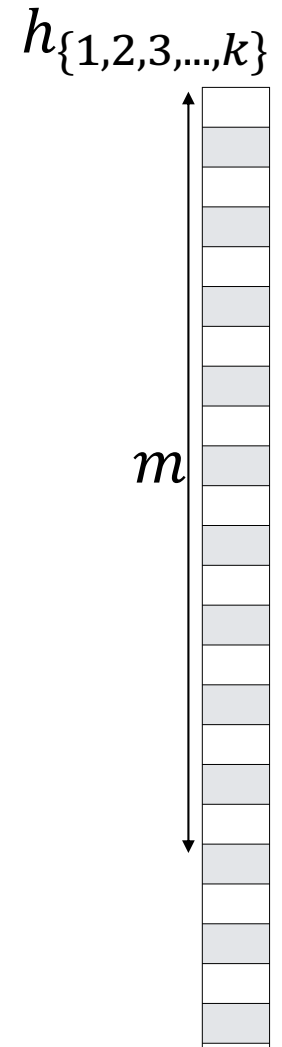Probability a specific bucket is 0 after one object is inserted?

$m$

After $n$ objects are inserted?

# Bloom Filter: Error Rate

Given bit vector of size $m$ and $k$ SUHA hash function

What's the probability a specific bucket is 1 after $n$ objects are inserted?

$h_{\{1,2,3,...,k\}}$

$m$

# Bloom Filter: Error Rate
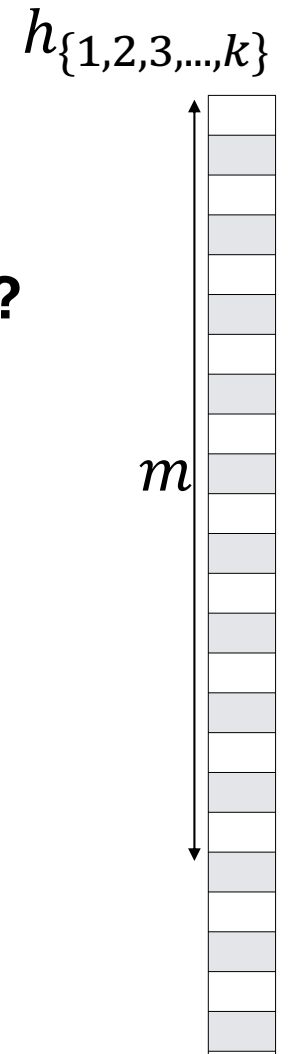
$h_{\{1,2,3,...,k\}}$

Given bit vector of size $m$ and $k$ SUHA hash function

**What is our expected FPR after $n$ objects are inserted?**

$m$

The probability my bit is 1 after $n$ objects inserted

$$\left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^{k}$$

The number of [assumed independent] trials

# Bloom Filter: Error Rate

Vector of size $m$, $k$ SUHA hash function, and $n$ objects

**To minimize the FPR, do we prefer…**

**(A) large** $k$        (**B) small** $k$

$$\left( 1 - \left( 1 - \frac{1}{m} \right)^{nk} \right)^{k}$$

$h_{\{1,2,3,\dots,k\}}$

$m$