*hashtables.*

# Bloom Filters

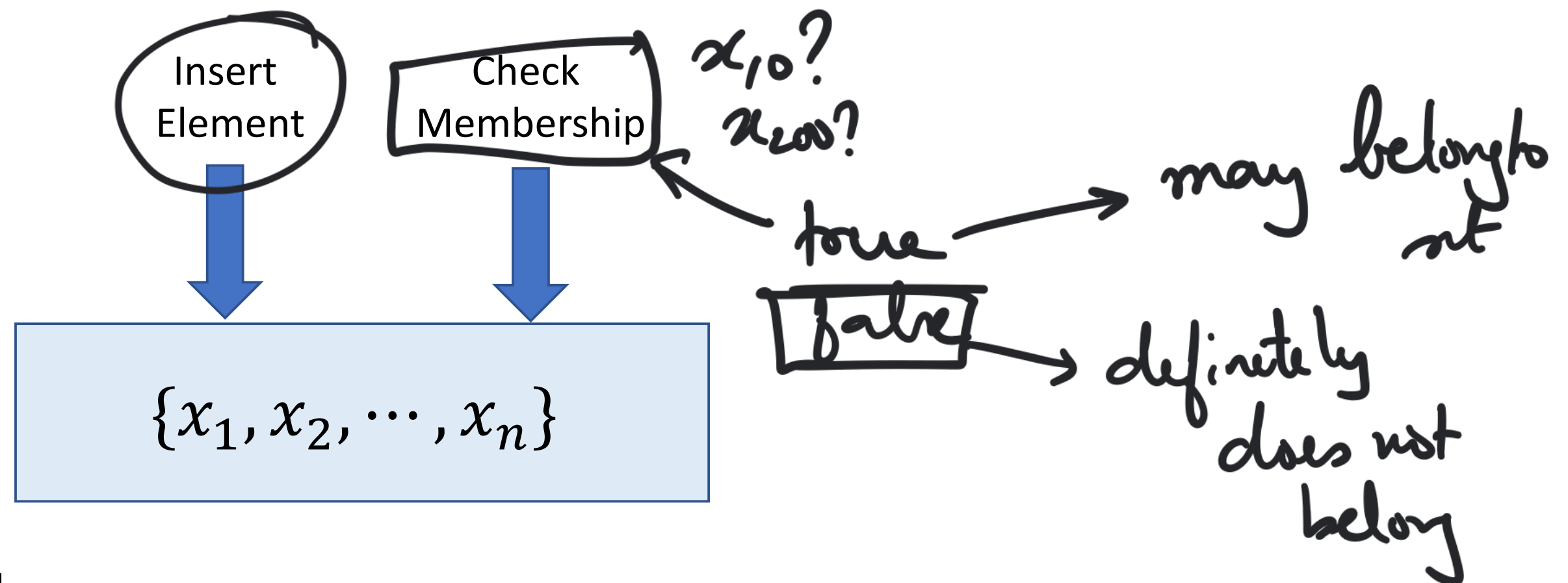Sriram Sankaranarayanan

Data Structures and Algorithms

# What is a Bloom Filter?

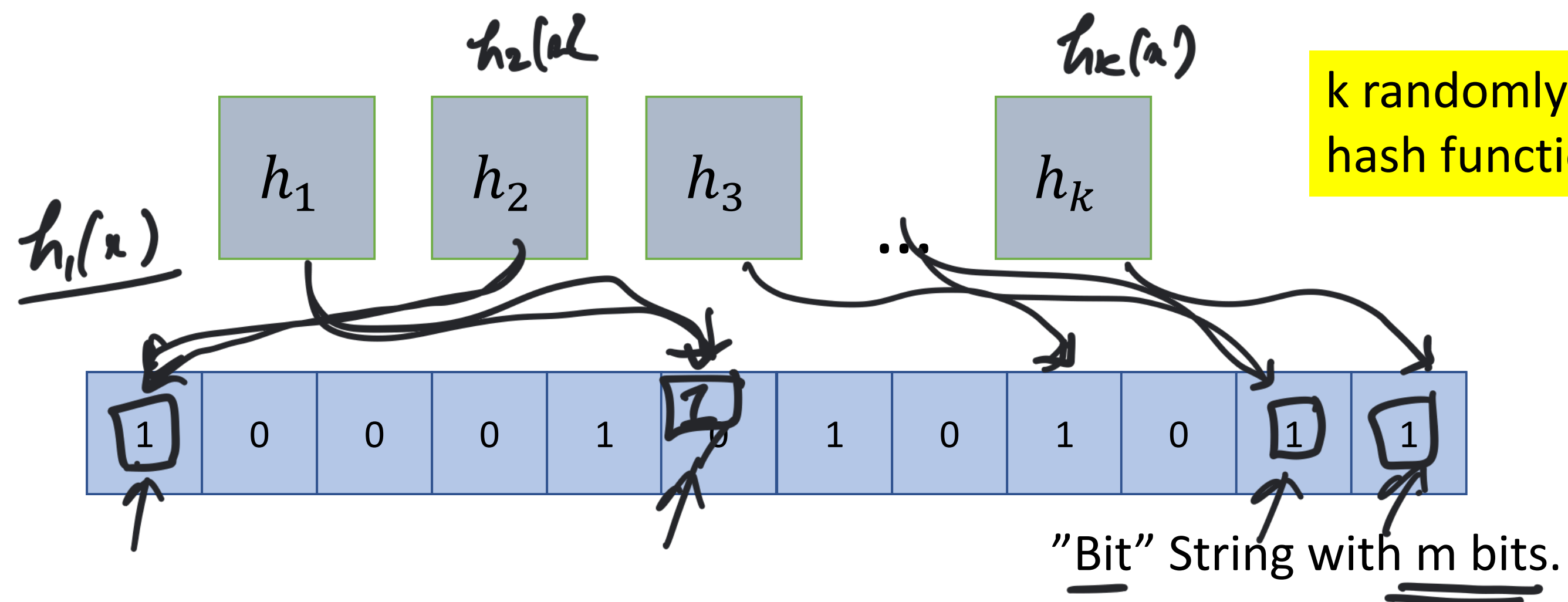$$\{\underline{x}_1, .., \underline{x}_n\}$$

- A fast set data structure based on hashing.

Insert Element      Check Membership      $x_{10}$?  $x_{200}$?

true → may belong to set

false → definitely does not belong

$$\{x_1, x_2, \cdots, x_n\}$$

- Based on hash-tables.
- Approximate in nature: false positives possible.

true

# Basic Idea

Falve Positives.

$h_2(x)$     $h_k(x)$

$h_1$     $h_2$     $h_3$     $h_k$

$h_1(x)$

| 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 |

"Bit" String with m bits.

**Insert element** $x$: Set the bits $h_1(x), h_2(x), \ldots, h_k(x)$

signature

Membership of element $x$: Are the bits $h_1(x), h_2(x), \ldots h_k(x)$ all set to 1?

# Bloom Filter: Properties

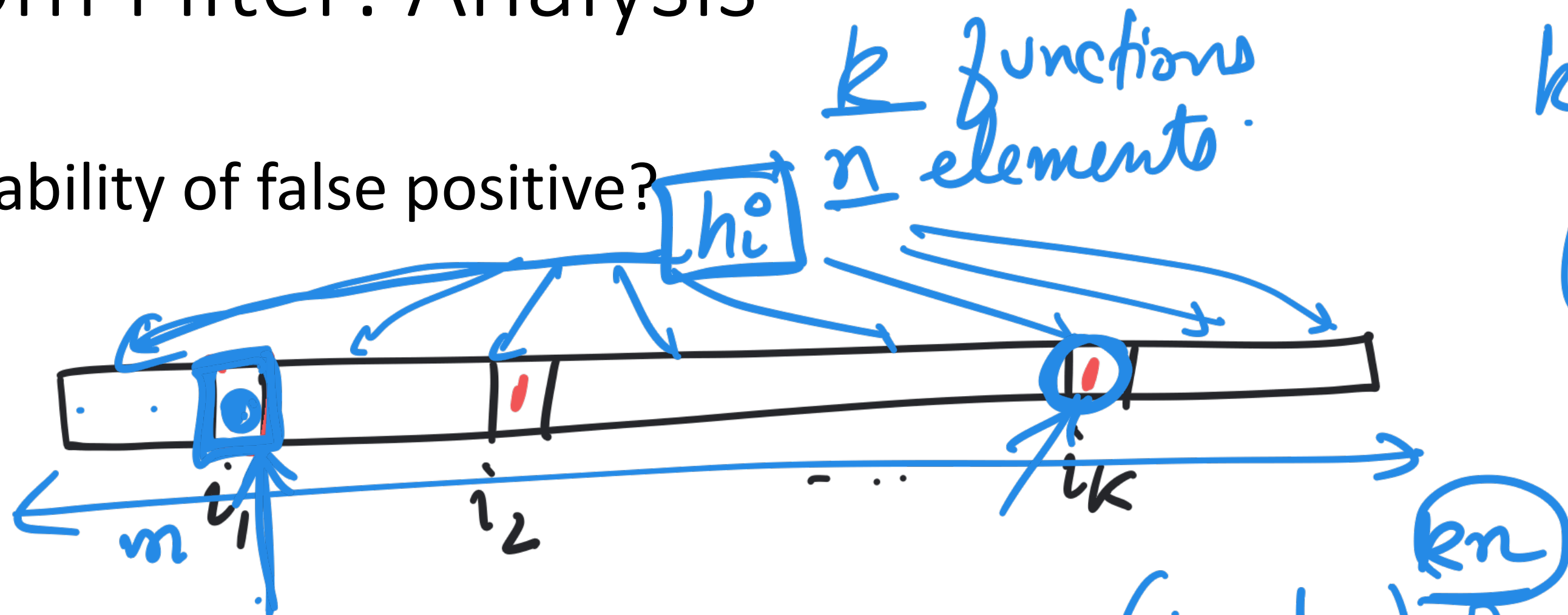- Constant time insertion and membership check
  - More precisely $\Theta(k)$

- If element was inserted, membership query will return true.

- False positives possible.
  - Membership query may return true but element may not have been inserted.

# Bloom Filter: Analysis

• Probability of false positive?



$k$ functions

$n$ elements

$h_i$

$kn$

$\lim_{m \to \infty} \left(1 - \frac{1}{m}\right)^{kn} \approx e$

$\Pr.\left( \text{bit } i_1 \underline{\text{ was not set}} \right) = \left(1 - \frac{1}{m}\right)^{kn}$

$\approx \boxed{e^{-kn/m}}$

$\Pr.\left( i_1 \text{ was set} \right) = \left(1 - e^{-kn/m}\right)$

$\Pr.\left( \underline{i_1 \text{ was set}} \text{ and } \underline{i_2} \ldots \underline{i_k} \right) = \left(1 - e^{-kn/m}\right)^k \leftarrow \text{False positive}$
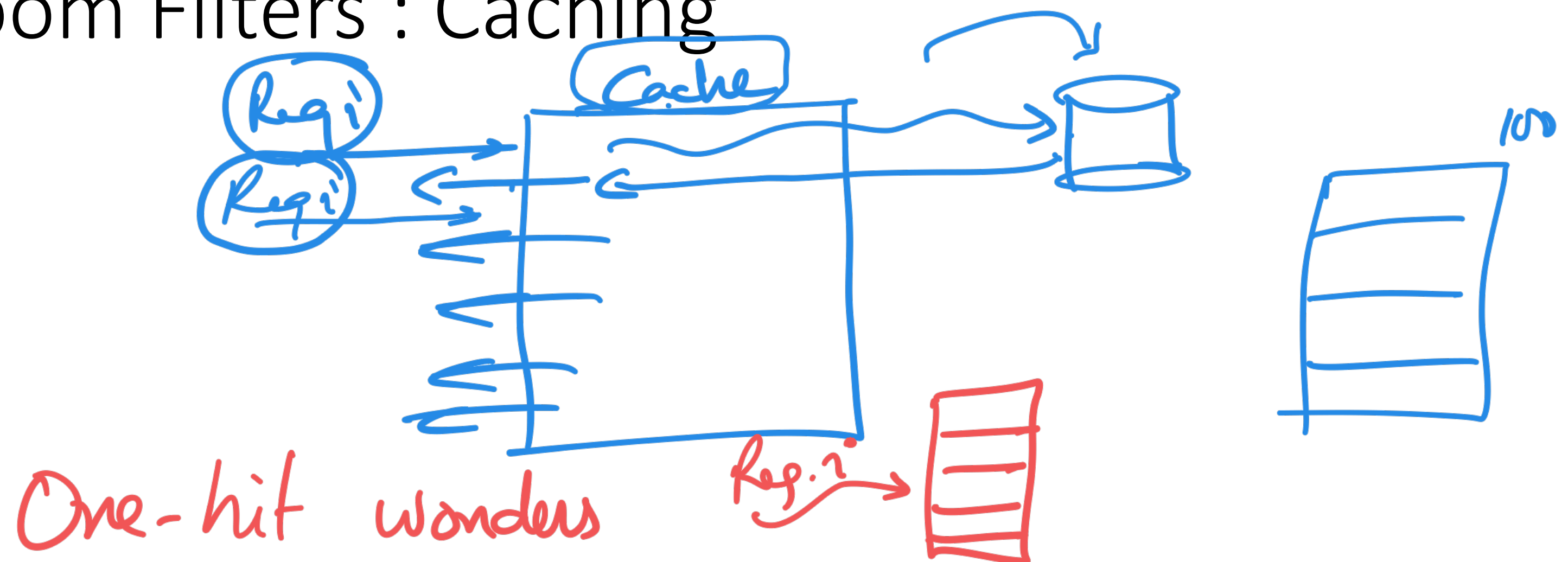
# Bloom Filter By Numbers

- n = 5,000 strings (these could be long strings) inserted
- m = 25,000 bit vector size (5 bits/element)
- k = 3 hash functions.

- Probability of false positives is

$$\left(1 - e^{-\frac{kn}{m}}\right)^k = (1 - e^{-0.6})^3 = 0.09$$

# Bloom Filters : Caching



One-hit wonders

*Maggs, Bruce M.; Sitaraman, Ramesh K. (July 2015), "Algorithmic nuggets in content delivery" (PDF), SIGCOMM Computer Communication Review, **45** (3): 52–66.*