

CS11921 - Fall 2024

# Algorithm Design & Analysis

Bloom Filters

Ibrahim Albluwi

# Did I see you before?

**Problem.** Given  $n$  elements (where  $n$  is large), answer queries of the form:  
Was element  $k$  seen before?

**Requirement.** Be as memory efficient as possible.

**Assumption.** A small possibility of incorrect YES answers is acceptable.

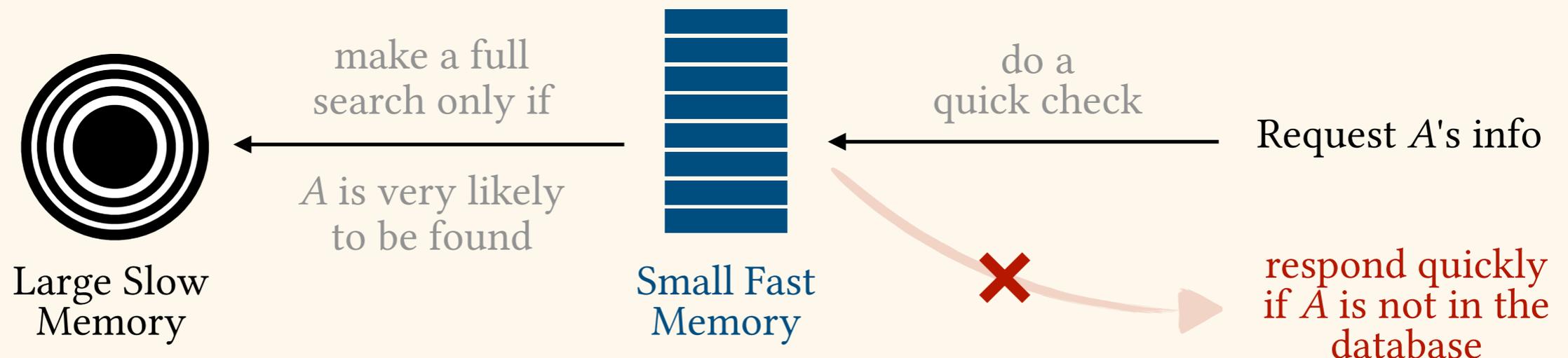


**Problem.** Design a *memory efficient* data structure to perform the following operations:

**INSERT**( $x$ ): Adds  $x$  to the set.

**CONTAINS**( $x$ ): Returns FALSE only if  $x$  is not in the set.  
Returns TRUE if  $x$  is *very likely* to be in the set.

## Application.



# Did I see you before?

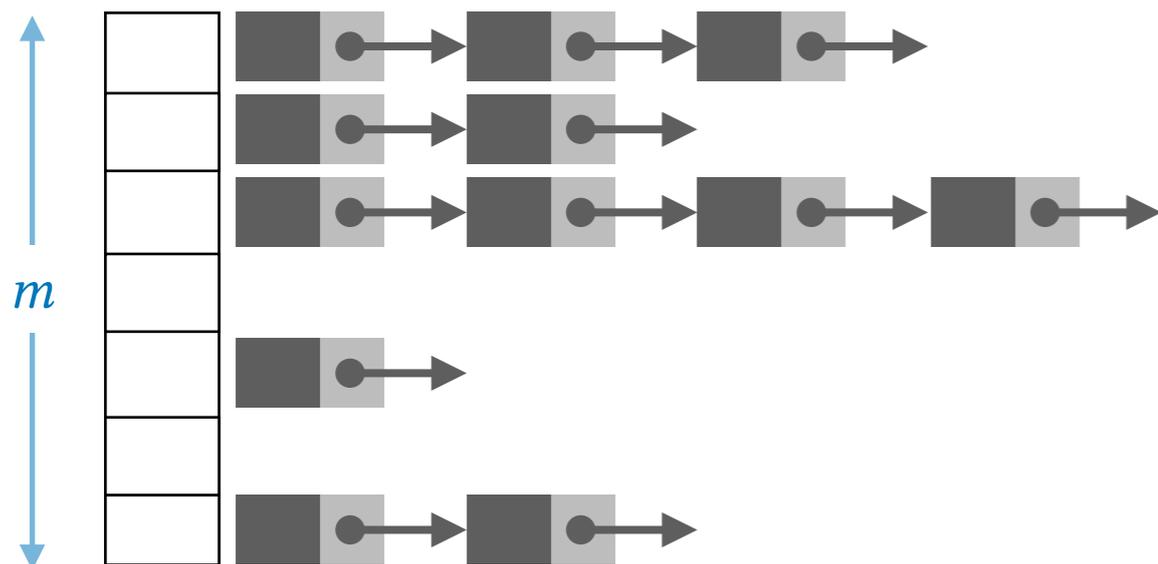
**Problem.** Design a *memory efficient* data structure to perform the following operations:

**INSERT**( $x$ ): Adds  $x$  to the set.

**CONTAINS**( $x$ ): Returns FALSE only if  $x$  is not in the set.  
Returns TRUE if  $x$  is *very likely* to be in the set.

**Solution # 0.** Use a hash table.

- Quick and always returns a correct answer!
- Not memory efficient!



A **C++ implementation** requires at least:

- $m \times 64$  bits for each linked list object.
- $n \times 64$  bits for each next pointer.
- $n \times 32$  bits for each element.  
(assuming the elements are integers)

Assuming  $m = n/4$ , the total is at least:

$$\frac{64}{4}n + 96n = 112n \text{ bits to store } n \text{ integers.}$$

(much more than that is needed  
if the elements are strings, for example)

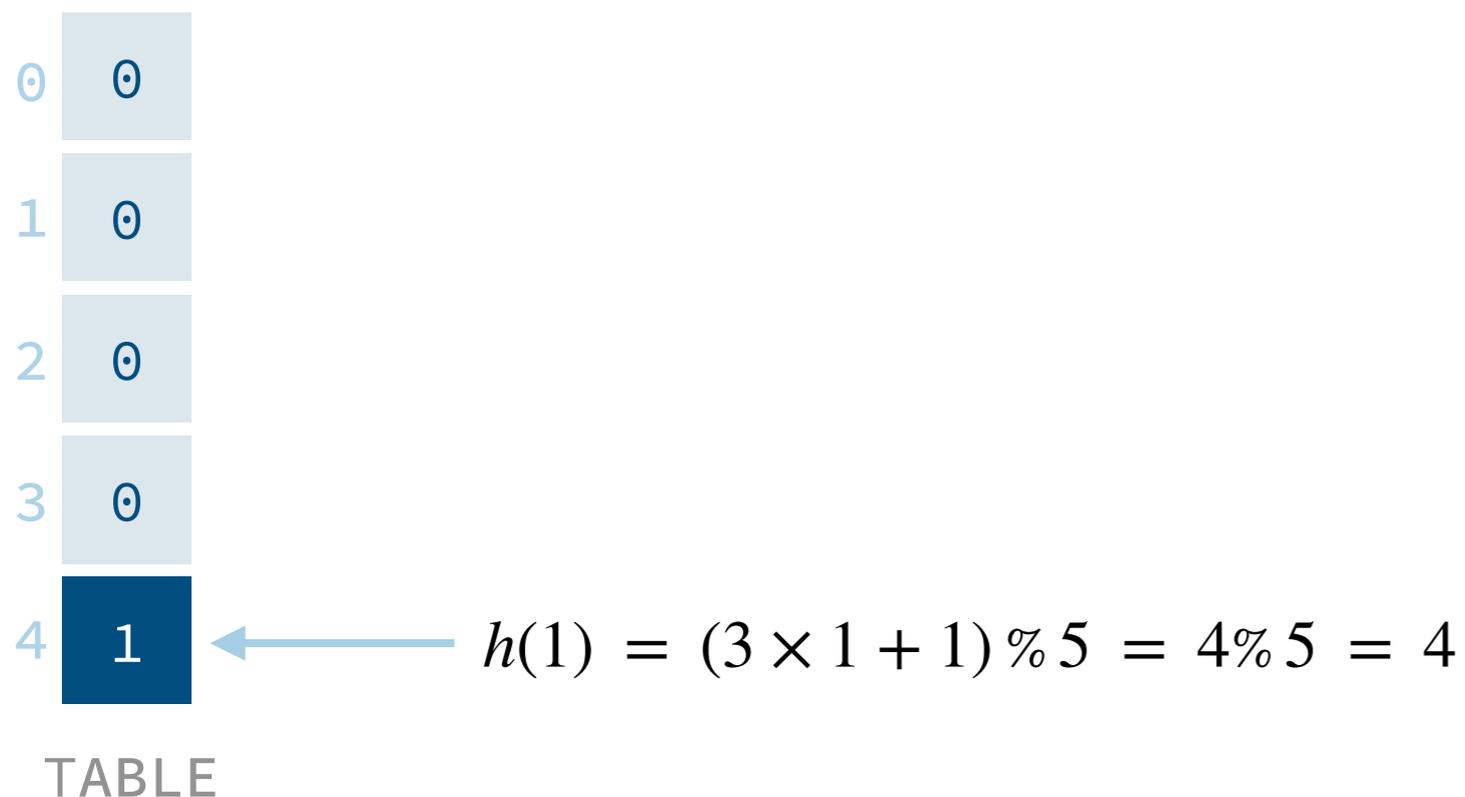
# Did I see you before?

**Solution # 1.** Use a hash table, but don't store the elements, and don't resolve collisions!

**INSERT**( $x$ ): Set  $\text{TABLE}[h(x)]$  to 1

**CONTAINS**( $x$ ): Return  $\text{TABLE}[h(x)] == 1$

**Example.** Insert  $\{1, 3, 4\}$  using  $h(x) = (3x + 1) \% m$



# Did I see you before?

**Solution # 1.** Use a hash table, but don't store the elements, and don't resolve collisions!

**INSERT**( $x$ ): Set  $\text{TABLE}[h(x)]$  to 1

**CONTAINS**( $x$ ): Return  $\text{TABLE}[h(x)] == 1$

**Example.** Insert  $\{1, 3, 4\}$  using  $h(x) = (3x + 1) \% m$

0	1	← $h(3) = (3 \times 3 + 1) \% 5 = 10 \% 5 = 0$
1	0	
2	0	
3	0	
4	1	

TABLE

# Did I see you before?

**Solution # 1.** Use a hash table, but don't store the elements, and don't resolve collisions!

**INSERT**( $x$ ): Set  $\text{TABLE}[h(x)]$  to 1

**CONTAINS**( $x$ ): Return  $\text{TABLE}[h(x)] == 1$

**Example.** Insert  $\{1, 3, 4\}$  using  $h(x) = (3x + 1) \% m$

0	1	
1	0	
2	0	
3	1	← $h(4) = (3 \times 4 + 1) \% 5 = 13 \% 5 = 3$
4	1	

TABLE

# Did I see you before?

**Solution # 1.** Use a hash table, but don't store the elements, and don't resolve collisions!

**INSERT**( $x$ ): Set  $\text{TABLE}[h(x)]$  to 1

**CONTAINS**( $x$ ): Return  $\text{TABLE}[h(x)] == 1$

**Example.** Insert  $\{1, 3, 4\}$  using  $h(x) = (3x + 1) \% m$

0	1	← <b>CONTAINS</b> (3) = TRUE $\text{TABLE}[h(3)] == 1$
1	0	
2	0	
3	1	
4	1	

TABLE

# Did I see you before?

**Solution # 1.** Use a hash table, but don't store the elements, and don't resolve collisions!

**INSERT**( $x$ ): Set  $\text{TABLE}[h(x)]$  to 1

**CONTAINS**( $x$ ): Return  $\text{TABLE}[h(x)] == 1$

**Example.** Insert  $\{1, 3, 4\}$  using  $h(x) = (3x + 1) \% m$

0	1
1	0
2	0
3	1
4	1

TABLE

← **CONTAINS**(5) = FALSE

$\text{TABLE}[h(5) = (15 + 1) \% 5 = 1] == 0$

# Did I see you before?

**Solution # 1.** Use a hash table, but don't store the elements, and don't resolve collisions!

**INSERT**( $x$ ): Set  $\text{TABLE}[h(x)]$  to 1

**CONTAINS**( $x$ ): Return  $\text{TABLE}[h(x)] == 1$

**Example.** Insert  $\{1, 3, 4\}$  using  $h(x) = (3x + 1) \% m$

0	1
1	0
2	0
3	1
4	1

TABLE



**CONTAINS**(6) = TRUE

$\text{TABLE}[h(6) = (18 + 1) \% 5 = 4] == 1$

**INCORRECT.** We did not add 6 to the table.

# Did I see you before?

**Solution # 1.** Use a hash table, but don't store the elements, and don't resolve collisions!

**INSERT**( $x$ ): Set  $\text{TABLE}[h(x)]$  to 1

**CONTAINS**( $x$ ): Return  $\text{TABLE}[h(x)] == 1$

**Example.** Insert  $\{1, 3, 4\}$  using  $h(x) = (3x + 1) \% m$

0	1
1	0
2	0
3	1
4	1

TABLE

**IF**  $\text{TABLE}[h(x)] == 0$ :

$x$  definitely was **not** seen before  
(otherwise, it would have been set to 1).

**IF**  $\text{TABLE}[h(x)] == 1$ :

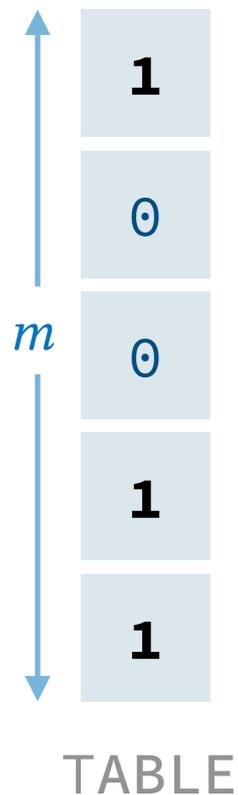
$x$  was seen before or another element  $y$   
was seen before, where  $h(x) = h(y)$ .



After  $n$  insertions, what is the probability that **CONTAINS**( $x$ ) finds  $\text{TABLE}[h(x)]$  to be 1?

# Solution # 1 Analysis

**Probability of Error.** After  $n$  insertions, what is the probability that **CONTAINS**( $x$ ) finds  $\text{TABLE}[h(x)]$  to be 1?



After 1 insertion:

$$P(\text{TABLE}[i] \text{ is } 1) = \frac{1}{m}$$

$$P(\text{TABLE}[i] \text{ is not } 1) = 1 - \frac{1}{m}$$

After  $n$  insertions:

$$P(\text{TABLE}[i] \text{ is still not } 1) = \left(1 - \frac{1}{m}\right)^n = \left(\left(1 - \frac{1}{m}\right)^m\right)^{\frac{n}{m}}$$

knowing that  $\lim_{x \rightarrow \infty} \left(1 - \frac{1}{x}\right)^x = e^{-1}$

$$P(\text{TABLE}[i] \text{ is still not } 1) = e^{-\frac{n}{m}}$$

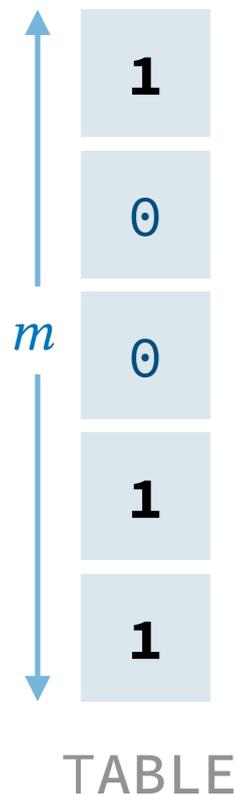
$$P(\text{TABLE}[i] \text{ is } 1) = 1 - e^{-\frac{n}{m}}$$



Is this low enough?

# Solution # 1 Analysis

**Probability of Error.** After  $n$  insertions, what is the probability that **CONTAINS**( $x$ ) finds  $\text{TABLE}[h(x)]$  to be 1?



After 1 insertion:

$$P(\text{TABLE}[i] \text{ is } 1) = \frac{1}{m}$$

$$P(\text{TABLE}[i] \text{ is not } 1) = 1 - \frac{1}{m}$$

After  $n$  insertions:

$$P(\text{TABLE}[i] \text{ is still not } 1) = \left(1 - \frac{1}{m}\right)^n = \left(\left(1 - \frac{1}{m}\right)^m\right)^{\frac{n}{m}}$$

knowing that  $\lim_{x \rightarrow \infty} \left(1 - \frac{1}{x}\right)^x = e^{-1}$

$$P(\text{TABLE}[i] \text{ is still not } 1) = e^{-\frac{n}{m}}$$

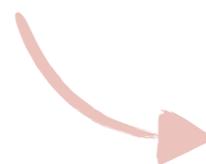
$$P(\text{TABLE}[i] \text{ is } 1) = 1 - e^{-\frac{n}{m}}$$

$m$	probability of error
100	$\approx 1$
1000	$\approx 1$
2000	$\approx 0.993$
5000	$\approx 0.865$
10,000	$\approx 0.632$
50,000	$\approx 0.181$
100,000	$\approx 0.095$
500,000	$\approx 0.02$
<b>1,000,000</b>	<b><math>\approx 0.01</math></b>

Assuming  $n = 10,000$

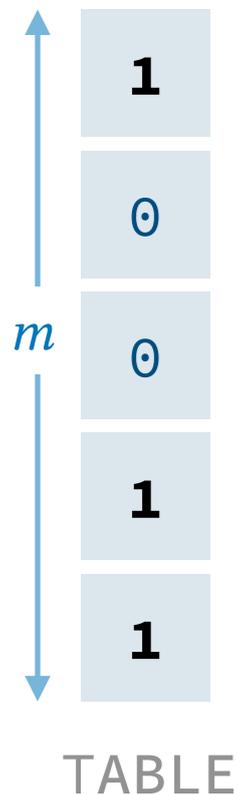
the probability of a wrong YES is

$$1 - e^{-\frac{10000}{1000000}} = 1 - e^{-\frac{1}{100}} \approx 0.01$$



# Solution # 1 Analysis

**Probability of Error.** After  $n$  insertions, what is the probability that **CONTAINS**( $x$ ) finds  $\text{TABLE}[h(x)]$  to be 1?



After 1 insertion:

$$P(\text{TABLE}[i] \text{ is } 1) = \frac{1}{m}$$

$$P(\text{TABLE}[i] \text{ is not } 1) = 1 - \frac{1}{m}$$

After  $n$  insertions:

$$P(\text{TABLE}[i] \text{ is still not } 1) = \left(1 - \frac{1}{m}\right)^n = \left(\left(1 - \frac{1}{m}\right)^m\right)^{\frac{n}{m}}$$

knowing that  $\lim_{x \rightarrow \infty} \left(1 - \frac{1}{x}\right)^x = e^{-1}$

$$P(\text{TABLE}[i] \text{ is still not } 1) = e^{-\frac{n}{m}}$$

$$P(\text{TABLE}[i] \text{ is } 1) = 1 - e^{-\frac{n}{m}}$$

Assuming  $n = 10,000$



We need 1 million bits ( $100n$  bits) to get a 1% error rate.

Not much better than normal hash tables (require  $\geq 112n$  bits)

$m$	probability of error
100	$\approx 1$
1000	$\approx 1$
2000	$\approx 0.993$
5000	$\approx 0.865$
10,000	$\approx 0.632$
50,000	$\approx 0.181$
100,000	$\approx 0.095$
500,000	$\approx 0.02$
<b>1,000,000</b>	<b><math>\approx 0.01</math></b>

# Solution # 2 : Bloom Filters

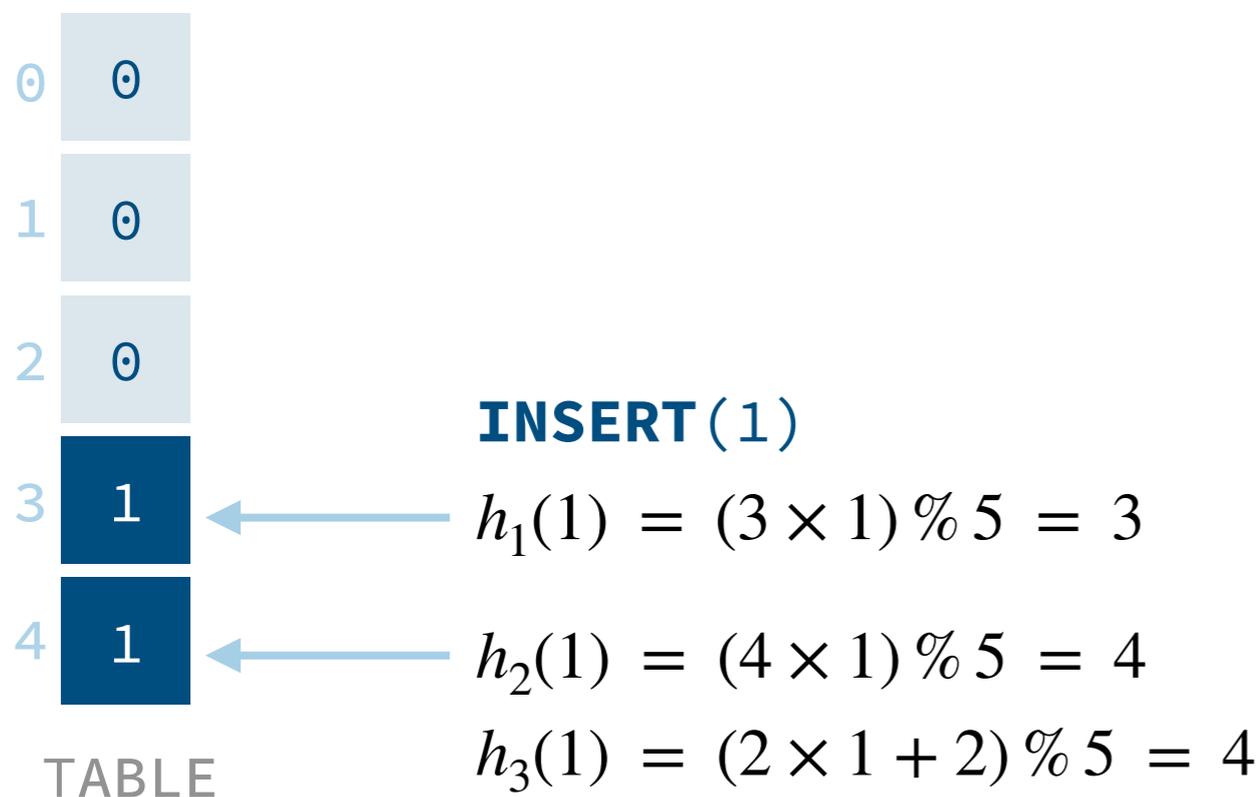
**Solution # 2.** Use multiple hash functions:  $h_1(x), h_2(x), \dots, h_k(x)$

**INSERT**( $x$ ): Set  $T[h_1(x)]=1$        $T[h_2(x)]=1$       ...  $T[h_k(x)]=1$

**CONTAINS**( $x$ ): Return  $T[h_1(x)]=1$  **AND**  $T[h_2(x)]=1$  **AND** ...  $T[h_k(x)]=1$

**Example.**  $k = 3$ ,     $h_1(x) = 3x \% m$ ,     $h_2(x) = 4x \% m$ ,     $h_3(x) = (2x + 2) \% m$ .

Assume that  $m = 5$ , and we insert the elements  $\{1, 3, 5\}$ .



# Solution # 2 : Bloom Filters

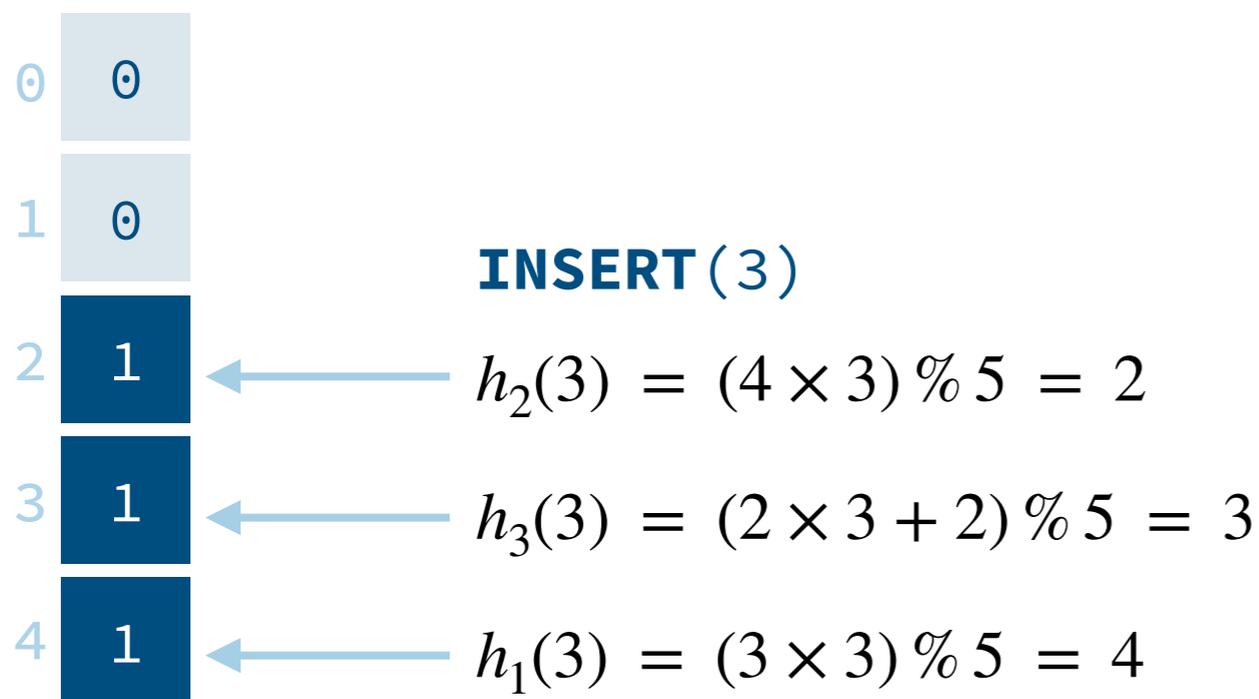
**Solution # 2.** Use multiple hash functions:  $h_1(x)$ ,  $h_2(x)$ , ...,  $h_k(x)$

**INSERT**( $x$ ): Set  $T[h_1(x)]=1$        $T[h_2(x)]=1$       ...  $T[h_k(x)]=1$

**CONTAINS**( $x$ ): Return  $T[h_1(x)]=1$  **AND**  $T[h_2(x)]=1$  **AND** ...  $T[h_k(x)]=1$

**Example.**  $k = 3$ ,  $h_1(x) = 3x \% m$ ,  $h_2(x) = 4x \% m$ ,  $h_3(x) = (2x + 2) \% m$ .

Assume that  $m = 5$ , and we insert the elements  $\{1, 3, 5\}$ .



TABLE

# Solution # 2 : Bloom Filters

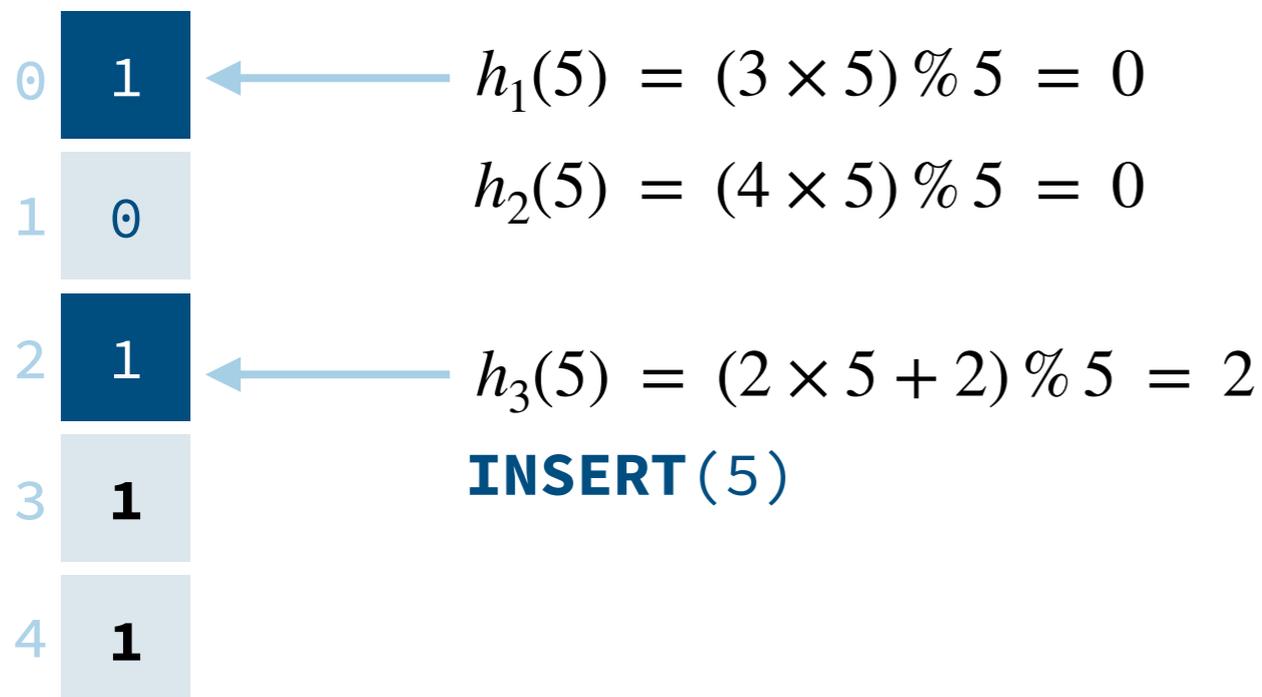
**Solution # 2.** Use multiple hash functions:  $h_1(x), h_2(x), \dots, h_k(x)$

**INSERT**( $x$ ): Set  $T[h_1(x)]=1$        $T[h_2(x)]=1$       ...  $T[h_k(x)]=1$

**CONTAINS**( $x$ ): Return  $T[h_1(x)]=1$  **AND**  $T[h_2(x)]=1$  **AND** ...  $T[h_k(x)]=1$

**Example.**  $k = 3$ ,     $h_1(x) = 3x \% m$ ,     $h_2(x) = 4x \% m$ ,     $h_3(x) = (2x + 2) \% m$ .

Assume that  $m = 5$ , and we insert the elements  $\{1, 3, 5\}$ .



TABLE

# Solution # 2 : Bloom Filters

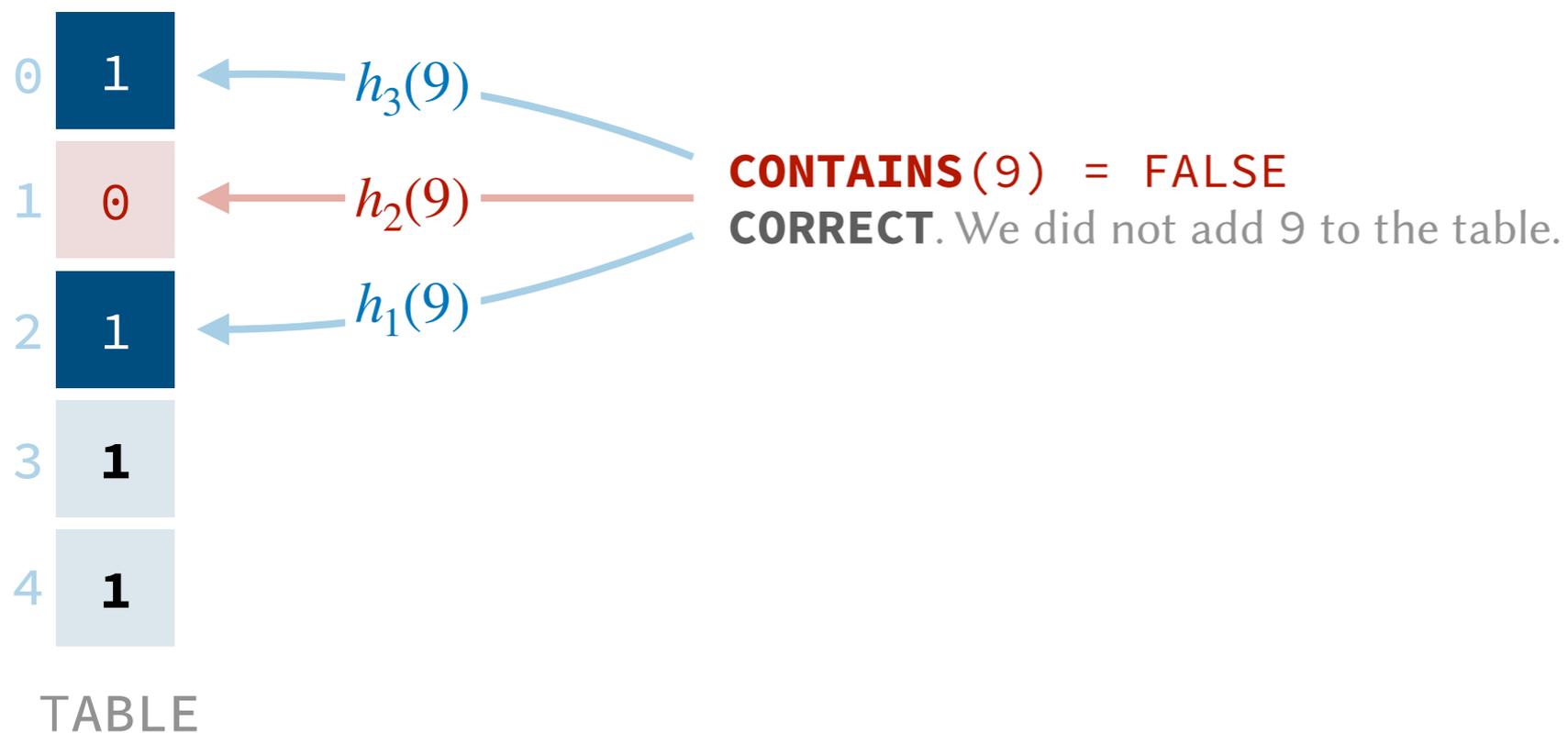
**Solution # 2.** Use multiple hash functions:  $h_1(x), h_2(x), \dots, h_k(x)$

**INSERT**( $x$ ): Set  $T[h_1(x)]=1$        $T[h_2(x)]=1$       ...  $T[h_k(x)]=1$

**CONTAINS**( $x$ ): Return  $T[h_1(x)]=1$  **AND**  $T[h_2(x)]=1$  **AND** ...  $T[h_k(x)]=1$

**Example.**  $k = 3$ ,  $h_1(x) = 3x \% m$ ,  $h_2(x) = 4x \% m$ ,  $h_3(x) = (2x + 2) \% m$ .

Assume that  $m = 5$ , and we insert the elements  $\{1, 3, 5\}$ .



# Solution # 2 : Bloom Filters

**Solution # 2.** Use multiple hash functions:  $h_1(x), h_2(x), \dots, h_k(x)$

**INSERT**( $x$ ): Set  $T[h_1(x)] = 1$        $T[h_2(x)] = 1$       ...  $T[h_k(x)] = 1$

**CONTAINS**( $x$ ): Return  $T[h_1(x)] == 1$  **AND**  $T[h_2(x)] == 1$  **AND** ...  $T[h_k(x)] == 1$

**Example.**  $k = 3$ ,  $h_1(x) = 3x \% m$ ,  $h_2(x) = 4x \% m$ ,  $h_3(x) = (2x + 2) \% m$ .

Assume that  $m = 5$ , and we insert the elements  $\{1, 3, 5\}$ .

0	1
1	0
2	1
3	1
4	1

TABLE

$h_1(6)$

$h_2(6)$

$h_3(6)$

**CONTAINS**(6) = TRUE

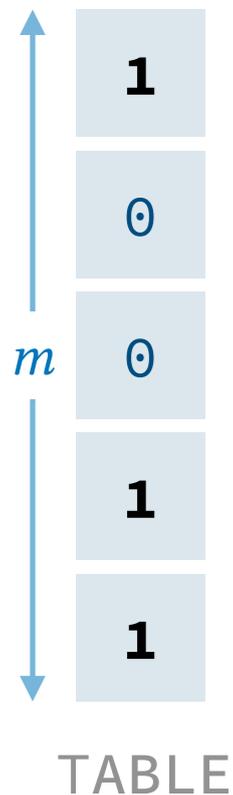
**WRONG.** We did not add 6 to the table.



After  $n$  insertions, what is the probability that **CONTAINS**( $x$ ) finds  $T[h_1(x)]$  and  $T[h_2(x)]$  ... and  $T[h_k(x)]$  to be 1?

# Solution # 2 : Analysis

**Probability of Error.** After  $n$  insertions, what is the probability that **CONTAINS**( $x$ ) finds  $\tau[h_1(x)]$  and  $\tau[h_2(x)] \dots$  and  $\tau[h_k(x)]$  to be 1?



After 1 insertion:

$$P(\tau[i] \text{ is not } 1) = \left(1 - \frac{1}{m}\right)^k$$

After  $n$  insertions:

$$P(\tau[i] \text{ is still not } 1) = \left(1 - \frac{1}{m}\right)^{kn} = \left(\left(1 - \frac{1}{m}\right)^m\right)^{\frac{kn}{m}}$$

knowing that  $\lim_{x \rightarrow \infty} \left(1 - \frac{1}{x}\right)^x = e^{-1}$

$$P(\tau[i] \text{ is still not } 1) = e^{-\frac{kn}{m}}$$

$$P(\tau[i] \text{ is } 1) = 1 - e^{-\frac{kn}{m}}$$

$$P(\tau[h_1(x)] \text{ and } \tau[h_2(x)] \dots \text{ and } \tau[h_k(x)] = 1) =$$

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

Fixing  $n$  and  $m$ , we find (using differentiation) that the optimal value of  $k$  that minimizes the probability of error is:

$$k = \frac{m}{n} \ln 2$$



Is this low enough?

# Solution # 2 : Analysis

**Probability of Error.** After  $n$  insertions, what is the probability that **CONTAINS**( $x$ ) finds  $\tau[h_1(x)]$  and  $\tau[h_2(x)]$  ... and  $\tau[h_k(x)]$  to be 1?

$$\epsilon = (1 - e^{-\frac{kn}{m}})^k$$

$$k = \frac{m}{n} \ln 2$$

Assuming  $n = 10,000$

$m$	optimal $k$	probability of error
10,000	$< 1$	$\approx 0.67$
20,000	$\approx 1.386$	$\approx 0.48$
30,000	$\approx 2.079$	$\approx 0.33$
40,000	$\approx 2.773$	$\approx 0.23$
50,000	$\approx 3.466$	$\approx 0.16$
100,000	$\approx 6.931$	$\approx 0.026$
125,000	$\approx 8.664$	$\approx 0.01$

We need 125,000 bits ( $12.5n$  bits) to get a 1% error rate (using 8-9 hash functions)



Much better than normal hash tables  
(require  $\geq 112n$  bits)