

# VE281

## Data Structures and Algorithms

### **Universal Hashing and Bloom Filter**

#### **Learning Objectives:**

- Know what is Universal Hashing
- Know what Bloom filter is and how it works
- Know the advantages and disadvantages of Bloom filter

# Universal Hashing

- ❖ Collision is bad!
  - ❖ For  $x \neq y$ ,  $h(x) = h(y)$
- ❖ Given any fixed hashing scheme, an **adversary** can create a sequence of inputs that maximizes collisions
- ❖ Remedy?
- ❖ The idea of randomization
  - ❖ Similar to quickSort and randomSelect

# Universal Hashing

- ❖ A scheme to produce hashing functions
  - ❖  $h = u(p)$
  - ❖ We talked about by creating hash function by taking the mod of prime numbers ( $p$ )
- ❖ Foil **adversaries** by randomly picking  $p$ !

# Definition of Universal Hashing

- ❖ A randomized algorithm  $\mathbf{H}$  for constructing hash functions  $h$
- ❖  $h : U \rightarrow \{1, \dots, M\}$
- ❖  $\mathbf{H}$  is universal if:
  - ❖ for all  $x \neq y$  in  $U$
  - ❖  $Pr_{h \leftarrow H}[h(x) = h(y)] \leq 1/M$
- ❖  $\mathbf{H}$  is also called as a **universal hash function family**

# Other Universal Hash Function Families

- ❖ The Matrix method
- ❖ Keys:  $u$ -bits long
- ❖ Table size:  $M=2^b$
- ❖  $h$ :  $b$ -by- $u$  0/1 matrix
- ❖  $H$ :  $h(x) = h \cdot x$

$$\begin{array}{c} h \quad x \quad h(x) \\ \begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 \end{array} \quad \begin{array}{c} 1 \\ 0 \\ 1 \\ 0 \end{array} = \begin{array}{c} 1 \\ 1 \\ 0 \end{array} \end{array}$$

# Proof of Universal

**Claim 10.4** For  $x \neq y$ ,  $\Pr_h[h(x) = h(y)] = 1/M = 1/2^b$ .

	$h$	$x$
1	0	0
0	1	1
1	1	0

# Short break

# Bloom Filter

- ❖ Invented by Burton Bloom in 1970
- ❖ Supports **fast insert** and **find**
- ❖ Comparison to hash tables:
  - ❖ Pros: more space efficient
  - ❖ Cons:
    1. Can't store an associated object
    2. No deletion (There are variations support deletion, but this operation is complicated)
    3. Small **false positive** probability: may say x has been inserted even if it hasn't been
      - ❖ But no false negative (x is inserted, but says not inserted)

# Bloom Filter Applications

- ❖ When to use bloom filter?
  - ❖ If the false positive is not a concern, no associated objects, no deletion, and you look for space efficiency
- ❖ Original application: spell checker
  - ❖ 40 years ago, space is a big concern, it's OK to tolerate some error
- ❖ Canonical application: list of forbidden passwords
  - ❖ Don't care about the false positive issue
- ❖ Modern applications: network routers
  - ❖ Limited memory, need to be fast
  - ❖ Applications include keeping track of blocked IP address, keeping track of contents of caches, etc.

# Bloom Filter Implementation: Components

- ❖ An array of  $n$  **bits**. Each bit 0 or 1
  - ❖  $n = b|S|$ , where  $b$  is small real number. For example,  $b \approx 8$  for 32-bit IP address (That's why it is space efficient)
- ❖  $k$  hash functions  $h_1, \dots, h_k$ , each mapping inside  $\{0,1, \dots, n - 1\}$ .
  - ❖  $k$  usually small.

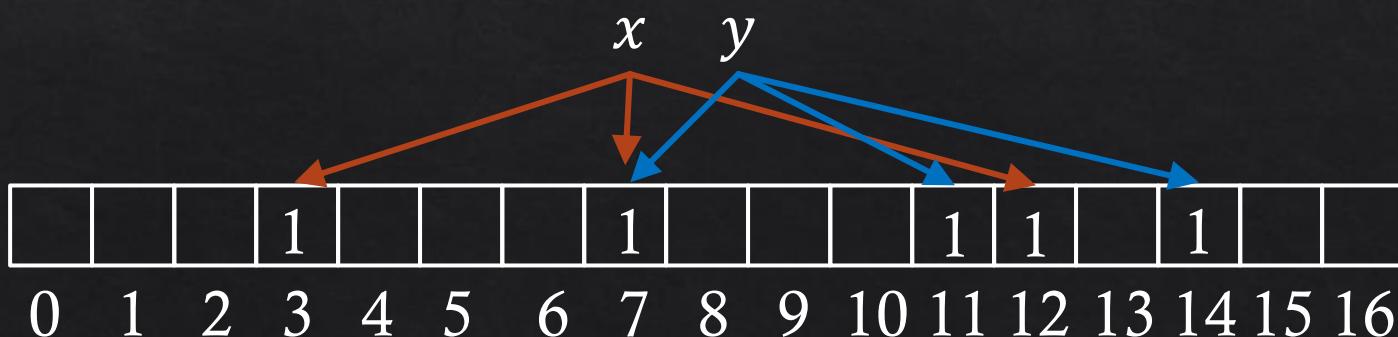
# Bloom Filter Insert

- ◆ Initially, the array is all-zero.
- ◆ Insert  $x$ : For  $i = 1, 2, \dots, k$ , set  $A[h_i(x)] = 1$ 
  - ◆ No matter whether the bit is 0 or 1 before

Example:  $n = 17$ , 3 hash functions

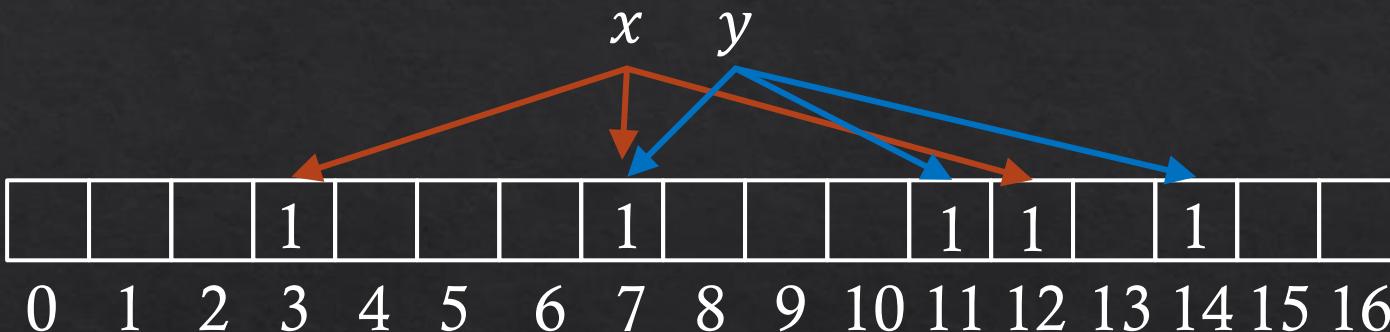
$$h_1(x) = 7, h_2(x) = 3, h_3(x) = 12$$

$$h_1(y) = 11, h_2(y) = 14, h_3(y) = 7$$



# Bloom Filter Find

- ◊ Find  $x$ : return true if and only if  $A[h_i(x)] = 1, \forall i = 1, \dots, k$



Suppose  $h_1(x) = 7, h_2(x) = 3, h_3(x) = 12$ . Find  $x$ ? Yes!

Suppose  $h_1(z) = 3, h_2(z) = 11, h_3(z) = 5$ . Find  $z$ ? No!

- ◊ No false negative: if  $x$  was inserted,  $\text{find}(x)$  guaranteed to return true
- ◊ False positive possible: consider  $h_1(w) = 11, h_2(w) = 12, h_3(w) = 7$  in the above example

# Heuristic Analysis of Error Probability

- ❖ Intuition: should be a trade-off between space (array size) and false positive probability
  - ❖ Array size decreases, more reuse of bits, false positive probability increases
- ❖ Goal: analyze the false positive probability
- ❖ Setup: Insert data set  $S$  into the Bloom filter, use  $k$  hash functions, array has  $n$  bits
- ❖ Assumption: All  $k$  hash functions map keys uniformly random and these hash functions are independent

# Probability of a Slot Being 1

- ❖ For an arbitrary slot  $j$  in the array, what's the probability that the slot is 1?
- ❖ Consider when slot  $j$  is 0
  - ❖ Happens when  $h_i(x) \neq j$  for all  $i = 1, \dots, k$  and  $x \in S$
  - ❖  $\Pr(h_i(x) \neq j) = 1 - \frac{1}{n}$
  - ❖  $\Pr(A[j] = 0) = \left(1 - \frac{1}{n}\right)^{k|S|} \approx e^{-\frac{k|S|}{n}} = e^{-\frac{k}{b}}$ 
    - ❖  $b = \frac{n}{|S|}$  denotes # of bits per object
  - ❖  $\Pr(A[j] = 1) \approx 1 - e^{-\frac{k}{b}}$

# False Positive Probability

- ❖ For  $x$  not in  $S$ , the false positive probability happens when all  $A[h_i(x)] = 1$  for all  $i = 1, \dots, k$ 
  - ❖ The probability is  $\epsilon \approx (1 - e^{-\frac{k}{b}})^k$
- ❖ For a fixed  $b$ ,  $\epsilon$  is minimized when  $k = (\ln 2) \cdot b$
- ❖ The minimal error probability is  $\epsilon \approx \left(\frac{1}{2}\right)^{\ln 2 \cdot b} \approx 0.6185^b$ 
  - ❖ Error probability decreases exponentially with  $b$
- ❖ Example:  $b = 8$ , could choose  $k$  as 5 or 6. Min error probability  $\approx 2\%$