## CS 473: Algorithms

Ruta Mehta

University of Illinois, Urbana-Champaign

Spring 2021

## CS 473: Algorithms, Spring 2021

# **Fingerprinting**

Lecture 11 March 2, 2021

Most slides are courtesy Prof. Chekuri

Ruta (UIUC) CS473 2 Spring 2021 2 / 30

### Hashing:

- **1** To insert x in dictionary store x in table in location h(x)
- ② To lookup y in dictionary check contents of location h(y)

Ruta (UIUC) CS473 3 Spring 2021 3 / 30

### Hashing:

- **1** To insert x in dictionary store x in table in location h(x)
- ② To lookup y in dictionary check contents of location h(y)

### **Bloom Filter:** tradeoff space for false positives

- What if elements (x) are unwieldy objects such a long strings, images, etc with non-uniform sizes.
- ② To insert x in dictionary, set *bit* at location h(x) to 1 (initially all bits are set to 0)
- **1** To lookup y if bit in location h(y) is 1 say yes, else no.

Ruta (UIUC) CS473 3 Spring 2021 3 / 30

### Hashing:

- **1** To insert x in dictionary store x in table in location h(x)
- 2 To lookup y in dictionary check contents of location h(y)

### **Bloom Filter:** tradeoff space for false positives

- What if elements (x) are unwieldy objects such a long strings, images, etc with non-uniform sizes.
- ② To insert x in dictionary, set *bit* at location h(x) to 1 (initially all bits are set to 0)
- **3** To lookup y if bit in location h(y) is 1 say yes, else no.

Issue: False positives due to collisions.

Ruta (UIUC) CS473 3 Spring 2021 3 / 30

**Bloom Filter:** tradeoff space for false positives

### Reducing false positives:

- **1** Pick k hash functions  $h_1, h_2, \ldots, h_k$  independently
- ② Insert x: for  $1 \le i \le k$  set bit in location  $h_i(x)$  in table i to 1

Ruta (UIUC) CS473 4 Spring 2021 4 / 30

**Bloom Filter:** tradeoff space for false positives

#### Reducing false positives:

- Pick k hash functions  $h_1, h_2, \ldots, h_k$  independently
- ② Insert x: for  $1 \leq i \leq k$  set bit in location  $h_i(x)$  in table i to 1
- **3** Lookup y: compute  $h_i(y)$  for  $1 \le i \le k$  and say yes only if each bit in the corresponding location is 1, otherwise say no. If probability of false positive for one hash function is  $\alpha < 1$  then with k independent hash function it is

Ruta (UIUC) CS473 4 Spring 2021 4 / 30

**Bloom Filter:** tradeoff space for false positives

### Reducing false positives:

- Pick k hash functions  $h_1, h_2, \ldots, h_k$  independently
- ② Insert x: for  $1 \le i \le k$  set bit in location  $h_i(x)$  in table i to 1
- **3** Lookup y: compute  $h_i(y)$  for  $1 \le i \le k$  and say yes only if each bit in the corresponding location is 1, otherwise say no. If probability of false positive for one hash function is  $\alpha < 1$  then with k independent hash function it is  $\alpha^k$ .

Ruta (UIUC) CS473 4 Spring 2021 4 / 30

## Take away points

- Hashing is a powerful and important technique for dictionaries.
   Many practical applications.
- Randomization fundamental to understand hashing.
- Good and efficient hashing possible in theory and practice with proper definitions (universal, perfect, etc).
- Related ideas of creating a compact fingerprint/sketch for objects is very powerful in theory and practice.

Ruta (UIUC) CS473 5 Spring 2021 5 / 30

## Fingerprinting

Source: Wikipedia

Process of mapping a large data item to a much shorter bit string, called its fingerprint.

Fingerprints uniquely identifies data "for all practical purposes".

Ruta (UIUC) CS473 6 Spring 2021 6 / 30

## Fingerprinting

Source: Wikipedia

Process of mapping a large data item to a much shorter bit string, called its fingerprint.

Fingerprints uniquely identifies data "for all practical purposes".

Typically used to avoid comparison and transmission of bulky data.

Eg: Web browser can store/fetch file fingerprints to check if it is changed.

Ruta (UIUC) CS473 6 Spring 2021 6 / 30

## Fingerprinting

Source: Wikipedia

Process of mapping a large data item to a much shorter bit string, called its fingerprint.

Fingerprints uniquely identifies data "for all practical purposes".

Typically used to avoid comparison and transmission of bulky data.

Eg: Web browser can store/fetch file fingerprints to check if it is changed.

As you may have guessed, fingerprint functions are hash functions.

Ruta (UIUC) CS473 6 Spring 2021 6 / 30

### Outline

### Use of hash functions for designing fast algorithms

### **Problem**

Given a text T of length m and pattern P of length n,  $m \gg n$ , find all occurrences of P in T.

Ruta (UIUC) CS473 7 Spring 2021 7 / 30

### Outline

### Use of hash functions for designing fast algorithms

#### Problem

Given a text T of length m and pattern P of length n,  $m \gg n$ , find all occurrences of P in T.

### Karp-Rabin Randomized Algorithm

Ruta (UIUC) CS473 7 Spring 2021 7 / 30

### Outline

### Use of hash functions for designing fast algorithms

#### **Problem**

Given a text T of length m and pattern P of length n,  $m \gg n$ , find all occurrences of P in T.

### Karp-Rabin Randomized Algorithm

#### It involves:

- Sampling a prime
- String equality via mod p arithmetic
- Rabin's fingerprinting scheme rolling hash
- Karp-Rabin pattern matching algorithm: O(m + n) time.

Ruta (UIUC) CS473 7 Spring 2021 7 / 30

## Part I

# Sampling a Prime

## Sampling a prime

### Problem

Given an integer x > 0, sample a prime uniformly at random from all the primes between 1 and x.

Ruta (UIUC) CS473 9 Spring 2021 9 / 30

## Sampling a prime

#### **Problem**

Given an integer x > 0, sample a prime uniformly at random from all the primes between 1 and x.

#### Procedure

- Sample a number p uniformly at random from  $\{1, \ldots, x\}$ .
- 2 If p is a prime, then output p. Else go to Step (1).

Ruta (UIUC) CS473 9 Spring 2021 9 / 30

## Sampling a prime

#### **Problem**

Given an integer x > 0, sample a prime uniformly at random from all the primes between 1 and x.

#### Procedure

- Sample a number p uniformly at random from  $\{1, \ldots, x\}$ .
- ② If p is a prime, then output p. Else go to Step (1).

## Checking if p is prime

- Agrawal-Kayal-Saxena primality test: deterministic but slow
- Miller-Rabin randomized primality test: fast but randomized outputs 'prime' when it is not with very low probability.

Ruta (UIUC) CS473 9 Spring 2021 9 / 30

Is the returned prime sampled uniformly at random?

Ruta (UIUC) CS473 10 Spring 2021 10 / 30

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[algorithm outputs p^*] = 1/\pi(x)$ .

Ruta (UIUC) CS473 10 Spring 2021 10 / 30

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[algorithm outputs p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round. Pr[A] =

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[algorithm outputs p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round.  $Pr[A] = \pi(x)/x$ .

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[algorithm outputs p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round.  $Pr[A] = \pi(x)/x$ .

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[algorithm outputs p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round.  $Pr[A] = \pi(x)/x$ .

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[\text{algorithm outputs } p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round.  $Pr[A] = \pi(x)/x$ .

$$Pr[A \cap B] =$$

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[algorithm outputs p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round.  $Pr[A] = \pi(x)/x$ .

$$Pr[A \cap B] = Pr[B] = 1/x$$
. Why?

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[algorithm outputs p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round.  $Pr[A] = \pi(x)/x$ .

Event B: number (prime)  $p^*$  is picked. Pr[B] = 1/x.

 $Pr[A \cap B] = Pr[B] = 1/x$ . Why? Because  $B \subset A$ .

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[\text{algorithm outputs } p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round.  $Pr[A] = \pi(x)/x$ .

$$Pr[A \cap B] = Pr[B] = 1/x$$
. Why? Because  $B \subset A$ .

$$Pr[B|A] =$$

Is the returned prime sampled uniformly at random?  $\pi(x)$ : number of primes in  $\{1,\ldots,x\}$ ,

#### Lemma

For a fixed prime  $p^* \le x$ ,  $\Pr[algorithm outputs p^*] = 1/\pi(x)$ .

### Proof.

Event A: a prime is picked in a round.  $Pr[A] = \pi(x)/x$ .

Event B: number (prime)  $p^*$  is picked. Pr[B] = 1/x.

 $Pr[A \cap B] = Pr[B] = 1/x$ . Why? Because  $B \subset A$ .

$$\Pr[B|A] = \frac{\Pr[A \cap B]}{\Pr[A]} = \frac{\Pr[B]}{\Pr[A]} = \frac{1/x}{\pi(x)/x} = \frac{1}{\pi(x)}$$

Ruta (UIUC) CS473 10 Spring 2021 10 / 30

## Sampling a prime: Expected number of samples

#### Procedure

- **1** Sample a number p uniformly at random from  $\{1, \ldots, x\}$ .
- ② If p is a prime, then output p. Else go to Step (1).

## Running time in expectation

**Q:** How many samples in expectation before termination?

**A:**  $x/\pi(x)$ . Exercise.

 $\pi(x)$ : Number of primes between **0** and x.

J. Hadamard and C. J. de la Vallée-Poussin (1896)

Prime Number Theorem:  $\lim_{x\to\infty} \frac{\pi(x)}{x/\ln x} = 1$ 

 $\pi(x)$ : Number of primes between **0** and x.

## J. Hadamard and C. J. de la Vallée-Poussin (1896)

Prime Number Theorem:  $\lim_{x\to\infty} \frac{\pi(x)}{x/\ln x} = 1$ 

## Chebyshev (from 1848)

$$\pi(x) \ge \frac{7}{8} \frac{x}{\ln x} = (1.262..) \frac{x}{\lg x} > \frac{x}{\lg x}$$

 $\pi(x)$ : Number of primes between **0** and x.

## J. Hadamard and C. J. de la Vallée-Poussin (1896)

Prime Number Theorem:  $\lim_{x\to\infty} \frac{\pi(x)}{x/\ln x} = 1$ 

## Chebyshev (from 1848)

$$\pi(x) \ge \frac{7}{8} \frac{x}{\ln x} = (1.262..) \frac{x}{\lg x} > \frac{x}{\lg x}$$

•  $y \sim \{1, \dots, x\}$  u.a.r., then y is a prime w.p.  $\frac{\pi(x)}{x} > \frac{1}{\lg x}$ .

 $\pi(x)$ : Number of primes between **0** and x.

## J. Hadamard and C. J. de la Vallée-Poussin (1896)

Prime Number Theorem:  $\lim_{x\to\infty} \frac{\pi(x)}{x/\ln x} = 1$ 

## Chebyshev (from 1848)

$$\pi(x) \ge \frac{7}{8} \frac{x}{\ln x} = (1.262..) \frac{x}{\lg x} > \frac{x}{\lg x}$$

- $y \sim \{1, \dots, x\}$  u.a.r., then y is a prime w.p.  $\frac{\pi(x)}{x} > \frac{1}{\lg x}$ .
- If we want  $k \ge 4$  primes then  $x \ge 2k \lg k$  suffices.

$$\pi(x) \ge \pi(2k \lg k) = \frac{2k \lg k}{\lg 2 + \lg k + \lg \lg k} \ge \frac{k(2 \lg k)}{2 \lg k} = k$$

Ruta (UIUC) CS473 12 Spring 2021 12 / 30

## Part II

# String Equality

#### **Problem**

Alice, the captain of a Mars lander, receives an N-bit string x, and Bob, back at mission control, receives a string y. They know nothing about each others strings, but want to check if x = y.

#### **Problem**

Alice, the captain of a Mars lander, receives an N-bit string x, and Bob, back at mission control, receives a string y. They know nothing about each others strings, but want to check if x = y.

Alice sends Bob x, and Bob confirms if x = y. But sending N bits is costly! Can they share less communication and check equality?

#### **Problem**

Alice, the captain of a Mars lander, receives an N-bit string x, and Bob, back at mission control, receives a string y. They know nothing about each others strings, but want to check if x = y.

Alice sends Bob x, and Bob confirms if x = y. But sending N bits is costly! Can they share less communication and check equality?

#### Possibilities:

- If want 100% surety then NO.
- If OK with 99.99% surety then  $O(\lg N)$  may suffice!!!

#### **Problem**

Alice, the captain of a Mars lander, receives an N-bit string x, and Bob, back at mission control, receives a string y. They know nothing about each others strings, but want to check if x = y.

Alice sends Bob x, and Bob confirms if x = y. But sending N bits is costly! Can they share less communication and check equality?

#### Possibilities:

- If want 100% surety then NO.
- If OK with 99.99% surety then  $O(\lg N)$  may suffice!!!
  - If x = y, then Pr[Bob says equal] = 1.
  - If  $x \neq y$ , then Pr[Bob says un-equal] = 0.9999.

Ruta (UIUC) CS473 14 Spring 2021 14 / 30

#### **Problem**

Alice, the captain of a Mars lander, receives an N-bit string x, and Bob, back at mission control, receives a string y. They know nothing about each others strings, but want to check if x = y.

Alice sends Bob x, and Bob confirms if x = y. But sending N bits is costly! Can they share less communication and check equality?

#### Possibilities:

- If want 100% surety then NO.
- If OK with 99.99% surety then  $O(\lg N)$  may suffice!!!
  - If x = y, then Pr[Bob says equal] = 1.
  - If  $x \neq y$ , then Pr[Bob says un-equal] = 0.9999.

#### HOW?

x, y: N-bit strings.

x, y: N-bit strings.

(Recall) If  $M = \lceil 2(5N) \lg 5N \rceil$ , then 5N primes in  $\{1, \dots, M\}$ .

Ruta (UIUC) CS473 15 Spring 2021 15 / 30

x, y: N-bit strings.

(Recall) If  $M = \lceil 2(5N) \lg 5N \rceil$ , then 5N primes in  $\{1, \ldots, M\}$ .

#### Procedure

Define  $h_p(x) = x \mod p$ 

• Alice picks a random prime p from  $\{1, \ldots M\}$ .

x, y: N-bit strings.

(Recall) If  $M = [2(5N) \lg 5N]$ , then 5N primes in  $\{1, \ldots, M\}$ .

#### Procedure

Define  $h_p(x) = x \mod p$ 

- Alice picks a random prime p from  $\{1, \ldots, M\}$ .
- ② She sends Bob prime p, and also  $h_p(x) = x \mod p$ .
- **3** Bob checks if  $h_p(y) = h_p(x)$ . If so, he says equal else un-equal.

Ruta (UIUC) CS473 Spring 2021

x, y: N-bit strings.

(Recall) If  $M = \lceil 2(5N) \lg 5N \rceil$ , then 5N primes in  $\{1, \ldots, M\}$ .

#### Procedure

Define  $h_p(x) = x \mod p$ 

- Alice picks a random prime p from  $\{1, \ldots M\}$ .
- ② She sends Bob prime p, and also  $h_p(x) = x \mod p$ .
- **3** Bob checks if  $h_p(y) = h_p(x)$ . If so, he says *equal* else *un-equal*.

#### Lemma

If x = y then Bob always says equal.

x, y: N-bit strings.

(Recall) If  $M = \lceil 2(5N) \lg 5N \rceil$ , then 5N primes in  $\{1, \ldots, M\}$ .

#### Procedure

Define  $h_p(x) = x \mod p$ 

- Alice picks a random prime p from  $\{1, \ldots M\}$ .
- ② She sends Bob prime p, and also  $h_p(x) = x \mod p$ .
- **3** Bob checks if  $h_p(y) = h_p(x)$ . If so, he says equal else un-equal.

#### Lemma

If  $x \neq y$  then,  $\Pr[Bob \ says \ equal] \leq 1/5$  (error probability).

Ruta (UIUC) CS473 16 Spring 2021 16 / 30

x, y: N-bit strings.

(Recall) If  $M = \lceil 2(sN) \lg sN \rceil$ , then sN primes in  $\{1, \ldots, M\}$ .

#### Procedure

Define  $h_p(x) = x \mod p$ 

- Alice picks a random prime p from  $\{1, \ldots M\}$ .
- ② She sends Bob prime p, and also  $h_p(x) = x \mod p$ .
- **3** Bob checks if  $h_p(y) = h_p(x)$ . If so, he says *equal* else *un-equal*.

#### Lemma

If  $x \neq y$  then,  $\Pr[Bob \ says \ equal] \leq 1/s$  (error probability).

Ruta (UIUC) CS473 17 Spring 2021 17 / 30

## Question.

Let x = 6 = 2 \* 3. If we draw a p u.a.r. from  $\{2, 3, 5, 7\}$ , then what is the probability that  $x \mod p = 0$ ?

- **(A)** 0.
- **(B)** 1.
- **(C)** 1/4.
- **(D)** 1/2.
- (E) none of the above.

## Question.

Let x = 6 = 2 \* 3. If we draw a p u.a.r. from  $\{2, 3, 5, 7\}$ , then what is the probability that  $x \mod p = 0$ ?

- (A) 0.
- **(B)** 1.
- **(C)** 1/4.
- **(D)** 1/2.
- (E) none of the above.

Now, let y = 21. What is the probability that  $(y - x) \mod p = 15 \mod p = 0$ ?

- (A) 0.
- **(B)** 1.
- (C) 1/4.
- (D) 1/2.

Error probability

$$x, y$$
 N-bit string,  $M = \lceil 2(sN) \lg sN \rceil$ , and  $h_p(x) = x \mod p$ 

#### Lemma

If 
$$x \neq y$$
 then,  $\Pr[Bob \ says \ equal] = \Pr[h_p(x) = h_p(y)] \leq 1/s$ 

#### Proof.

Given 
$$x \neq y$$
,  $h_p(x) = h_p(y) \Rightarrow x \mod p = y \mod p$ .

CS473 19 Spring 2021 19 / 30

Error probability

$$x, y$$
 N-bit string,  $M = \lceil 2(sN) \lg sN \rceil$ , and  $h_p(x) = x \mod p$ 

#### Lemma

If 
$$x \neq y$$
 then,  $\Pr[Bob \ says \ equal] = \Pr[h_p(x) = h_p(y)] \leq 1/s$ 

#### Proof.

Given 
$$x \neq y$$
,  $h_p(x) = h_p(y) \Rightarrow x \mod p = y \mod p$ .

• 
$$D = |x - y|$$
, then  $D \mod p = 0$ , and  $D \le 2^N$ .

Ruta (UIUC) CS473 19 Spring 2021 19 / 30

Error probability

$$x, y$$
 N-bit string,  $M = \lceil 2(sN) \lg sN \rceil$ , and  $h_p(x) = x \mod p$ 

#### Lemma

If 
$$x \neq y$$
 then,  $\Pr[Bob \ says \ equal] = \Pr[h_p(x) = h_p(y)] \leq 1/s$ 

#### Proof.

Given  $x \neq y$ ,  $h_p(x) = h_p(y) \Rightarrow x \mod p = y \mod p$ .

- D = |x y|, then  $D \mod p = 0$ , and  $D \le 2^N$ .
- $D = p_1 \dots p_k$  prime factorization.

Error probability

$$x, y$$
 N-bit string,  $M = \lceil 2(sN) \lg sN \rceil$ , and  $h_p(x) = x \mod p$ 

#### Lemma

If 
$$x \neq y$$
 then,  $\Pr[Bob \ says \ equal] = \Pr[h_p(x) = h_p(y)] \leq 1/s$ 

#### Proof.

Given  $x \neq y$ ,  $h_p(x) = h_p(y) \Rightarrow x \mod p = y \mod p$ .

- D = |x y|, then  $D \mod p = 0$ , and  $D \le 2^N$ .
- $D = p_1 \dots p_k$  prime factorization. All  $p_i \geq 2 \Rightarrow D \geq 2^k$ .

Ruta (UIUC) CS473 19 Spring 2021 19 / 30

Error probability

$$x, y$$
 N-bit string,  $M = \lceil 2(sN) \lg sN \rceil$ , and  $h_p(x) = x \mod p$ 

#### Lemma

If 
$$x \neq y$$
 then,  $\Pr[Bob \ says \ equal] = \Pr[h_p(x) = h_p(y)] \leq 1/s$ 

#### Proof.

Given  $x \neq y$ ,  $h_p(x) = h_p(y) \Rightarrow x \mod p = y \mod p$ .

- D = |x y|, then  $D \mod p = 0$ , and  $D \le 2^N$ .
- $D = p_1 \dots p_k$  prime factorization. All  $p_i \ge 2 \Rightarrow D \ge 2^k$ .
- $2^k \le D \le 2^N \Rightarrow k \le N$ . D has at most N divisors.

Ruta (UIUC) CS473 19 Spring 2021 19 / 30

Error probability

$$x, y$$
 N-bit string,  $M = \lceil 2(sN) \lg sN \rceil$ , and  $h_p(x) = x \mod p$ 

#### Lemma

If 
$$x \neq y$$
 then,  $\Pr[Bob \ says \ equal] = \Pr[h_p(x) = h_p(y)] \leq 1/s$ 

#### Proof.

Given  $x \neq y$ ,  $h_p(x) = h_p(y) \Rightarrow x \mod p = y \mod p$ .

- D = |x y|, then  $D \mod p = 0$ , and  $D \le 2^N$ .
- $D = p_1 \dots p_k$  prime factorization. All  $p_i \ge 2 \Rightarrow D \ge 2^k$ .
- $2^k \le D \le 2^N \Rightarrow k \le N$ . D has at most N divisors.
- Probability that a random prime p from  $\{1, \ldots, M\}$  is a divisor  $= \frac{k}{\pi(M)} \le \frac{N}{\pi(M)}$

Error probability

$$x, y$$
 N-bit string,  $M = \lceil 2(sN) \lg sN \rceil$ , and  $h_p(x) = x \mod p$ 

#### Lemma

If 
$$x \neq y$$
 then,  $\Pr[Bob \ says \ equal] = \Pr[h_p(x) = h_p(y)] \leq 1/s$ 

#### Proof.

Given  $x \neq y$ ,  $h_p(x) = h_p(y) \Rightarrow x \mod p = y \mod p$ .

- D = |x y|, then  $D \mod p = 0$ , and  $D \le 2^N$ .
- $D = p_1 \dots p_k$  prime factorization. All  $p_i \ge 2 \Rightarrow D \ge 2^k$ .
- $2^k \le D \le 2^N \Rightarrow k \le N$ . D has at most N divisors.
- Probability that a random prime p from  $\{1, \ldots, M\}$  is a divisor  $= \frac{k}{\pi(M)} \le \frac{N}{\pi(M)} \le \frac{N}{M/\lg M} = \frac{N}{2(sN)\lg sN} \lg M \le \frac{1}{s}$

Ruta (UIUC) CS473 19 Spring 2021 19 / 30

## Low Error Probability

**1** Choose large enough s. Error prob: 1/s.

## Low Error Probability

- **1** Choose large enough s. Error prob: 1/s.
- Alice repeats the process R times, and Bob says equal only if he gets equal all R times.

## Low Error Probability

- **1** Choose large enough s. Error prob: 1/s.
- Alice repeats the process R times, and Bob says equal only if he gets equal all R times.

Error probability:  $\frac{1}{s^R}$ .

## Low Error Probability

- **1** Choose large enough s. Error prob: 1/s.
- Alice repeats the process R times, and Bob says equal only if he gets equal all R times.

Error probability:  $\frac{1}{s^R}$ . For s = 5, R = 10,  $\frac{1}{5^{10}} \le 0.000001$ .

## Low Error Probability

- Choose large enough s. Error prob: 1/s.
- Alice repeats the process R times, and Bob says equal only if he gets equal all R times.

Error probability:  $\frac{1}{sR}$ . For s = 5, R = 10,  $\frac{1}{510} \le 0.000001$ .

$$M = \lceil 2(sN) \lg sN \rceil$$

#### Amount of Communication

Each round sends 2 integers < M. # bits:  $2 \lg M < 4(\lg s + \lg N)$ .

Ruta (UIUC) CS473 Spring 2021 20 / 30

## Low Error Probability

- **1** Choose large enough s. Error prob: 1/s.
- Alice repeats the process R times, and Bob says equal only if he gets equal all R times.

Error probability:  $\frac{1}{s^R}$ . For  $s=5, R=10, \frac{1}{5^{10}} \leq 0.000001$ .

$$M = \lceil 2(sN) \lg sN \rceil$$

#### Amount of Communication

Each round sends 2 integers  $\leq M$ . # bits:  $2 \lg M \leq 4(\lg s + \lg N)$ .

If x and y are copies of Wikipedia, about 25 billion characters. If 8 bits per character, then  $N \approx 2^{38}$  bits.

Ruta (UIUC) CS473 20 Spring 2021 20 / 30

## Low Error Probability

- **1** Choose large enough s. Error prob: 1/s.
- Alice repeats the process R times, and Bob says equal only if he gets equal all R times.

Error probability:  $\frac{1}{s^R}$ . For  $s=5, R=10, \frac{1}{5^{10}} \leq 0.000001$ .

$$M = \lceil 2(sN) \lg sN \rceil$$

#### Amount of Communication

Each round sends 2 integers  $\leq M$ . # bits:  $2 \lg M \leq 4(\lg s + \lg N)$ .

If x and y are copies of Wikipedia, about 25 billion characters. If 8 bits per character, then  $N \approx 2^{38}$  bits.

Second approach will send  $10(2 \lg (10N \lg 5N)) \le 1280$  bits.

Ruta (UIUC) CS473 20 Spring 2021 20 / 30

## Part III

# Karp-Rabin Pattern Matching Algorithm

Given a string T of length m and pattern P of length n, s.t.  $m \gg n$ , find all occurrences of P in T.

## Example

**T**=abracadabra, **P**=ab.

Given a string T of length m and pattern P of length n, s.t.  $m \gg n$ , find all occurrences of P in T.

## Example

**T**=abracadabra, **P**=ab.

Solution  $S = \{1, 8\}$ .

Given a string T of length m and pattern P of length n, s.t.  $m \gg n$ , find all occurrences of P in T.

## Example

**T**=abracadabra, **P**=ab.

Solution  $S = \{1, 8\}$ .

For j > i, let  $T_{i...j} = T[i]T[i+1]...T[j]$ .

Given a string T of length m and pattern P of length n, s.t.  $m \gg n$ , find all occurrences of P in T.

## Example

**T**=abracadabra, **P**=ab.

Solution  $S = \{1, 8\}$ .

For 
$$j > i$$
, let  $T_{i...j} = T[i]T[i+1]...T[j]$ .

## Brute force algorithm

 $S = \emptyset$ . For each  $i = 1 \dots m - n + 1$ 

• If  $T_{i...i+n-1} = P$  then  $S = S \cup \{i\}$ .

Given a string T of length m and pattern P of length n, s.t.  $m \gg n$ , find all occurrences of P in T.

## Example

**T**=abracadabra, **P**=ab.

Solution  $S = \{1, 8\}$ .

For 
$$j > i$$
, let  $T_{i...j} = T[i]T[i+1]...T[j]$ .

### Brute force algorithm

 $S = \emptyset$ . For each  $i = 1 \dots m - n + 1$ 

• If  $T_{i...i+n-1} = P$  then  $S = S \cup \{i\}$ .

O(mn) run-time.

## Using Hash Function

Pick a prime p u.a.r. from  $\{1, \ldots, M\}$ .  $h_p(x) = x \mod p$ .

## Brute force algorithm using hash function

$$S = \emptyset$$
. For each  $i = 1 \dots m - n + 1$ 

• If 
$$h_p(T_{i...i+n-1}) = h_p(P)$$
 then  $S = S \cup \{i\}$ .

Ruta (UIUC) CS473 23 Spring 2021 23 / 30

## Using Hash Function

Pick a prime p u.a.r. from  $\{1, \ldots, M\}$ .  $h_p(x) = x \mod p$ .

## Brute force algorithm using hash function

$$S = \emptyset$$
. For each  $i = 1 \dots m - n + 1$ 

• If 
$$h_p(T_{i...i+n-1}) = h_p(P)$$
 then  $S = S \cup \{i\}$ .

If x is of length n, then computing  $h_p(x)$  takes O(n) running time.

Overall O(mn) running time.

Ruta (UIUC) CS473 23 Spring 2021 23 / 30

# Using Hash Function

Pick a prime p u.a.r. from  $\{1, \ldots, M\}$ .  $h_p(x) = x \mod p$ .

## Brute force algorithm using hash function

$$S = \emptyset$$
. For each  $i = 1 \dots m - n + 1$ 

• If 
$$h_p(T_{i...i+n-1}) = h_p(P)$$
 then  $S = S \cup \{i\}$ .

If x is of length n, then computing  $h_p(x)$  takes O(n) running time.

Overall O(mn) running time.

Can we compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  fast?

## mod p math

Let a and b be (non-negative) integers.

$$(a+b) \bmod p = ((a \bmod p) + (b \bmod p)) \bmod p$$

## mod p math

Let a and b be (non-negative) integers.

$$(a+b) \bmod p = ((a \bmod p) + (b \bmod p)) \bmod p$$

$$(a \cdot b) \mod p = ((a \mod p) \cdot (b \mod p)) \mod p$$

$$x = T_{i...i+n-1}$$
 and  $x' = T_{i+1...i+n}$ 

## Example

x = 1011001, and x' = 0110010 (or x' = 0110011).

$$x = T_{i...i+n-1}$$
 and  $x' = T_{i+1...i+n}$ .

## Example

x = 1011001, and x' = 0110010 (or x' = 0110011).

$$x' = 2(x - x_{hb}2^{n-1}) + x'_{lb}$$

$$x = T_{i...i+n-1}$$
 and  $x' = T_{i+1...i+n}$ 

## Example

x = 1011001, and x' = 0110010 (or x' = 0110011).

$$x' = 2(x - x_{hb}2^{n-1}) + x'_{lb}$$
  
=  $2x - x_{hb}2^{n} + x'_{lb}$ 

$$x = T_{i...i+n-1}$$
 and  $x' = T_{i+1...i+n}$ .

## Example

x = 1011001, and x' = 0110010 (or x' = 0110011).

$$x' = 2(x - x_{hb}2^{n-1}) + x'_{lb}$$
  
=  $2x - x_{hb}2^{n} + x'_{lb}$ 

$$h_{p}(x') = x' \mod p$$

$$= (2(x \mod p) - x_{hb}(2^{n} \mod p) + x'_{lb}) \mod p$$

$$= (2h_{p}(x) - x_{hb}h_{p}(2^{n}) + x'_{lb}) \mod p$$

$$x = T_{i...i+n-1}$$
 and  $x' = T_{i+1...i+n}$ .

## Example

x = 1011001, and x' = 0110010 (or x' = 0110011).

$$x' = 2(x - x_{hb}2^{n-1}) + x'_{lb}$$
  
=  $2x - x_{hb}2^{n} + x'_{lb}$ 

$$h_{p}(x') = x' \mod p$$

$$= (2(x \mod p) - x_{hb}(2^{n} \mod p) + x'_{lb}) \mod p$$

$$= (2h_{p}(x) - x_{hb}h_{p}(2^{n}) + x'_{lb}) \mod p$$

$$= (2h_{p}(T_{i...i+n-1}) - T_{i}h_{p}(2^{n}) + T_{i+n}) \mod p$$

p: a random prime from  $\{1,\ldots,M\}$ . Rolling hash:  $h_p(T_{i+1\ldots i+n})=(2h_p(T_{i\ldots i+n-1})-T_ih_p(2^n)+T_{i+n}) \mod p$ .

- p: a random prime from  $\{1,\ldots,M\}$ . Rolling hash:  $h_p(T_{i+1,\ldots,i+p}) = (2h_p(T_{i,\ldots,i+p-1}) - T_ih_p(2^n) + T_{i+p}) \mod p$ .
  - Set  $S = \emptyset$ . Compute  $h_p(T_{1...n})$ ,  $h_p(2^n)$ , and  $h_p(P)$ .
  - ② For each i = 1, ..., m n + 1
    - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
    - **2** Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$  by applying rolling hash.

- p: a random prime from  $\{1, \ldots, M\}$ .
- Rolling hash:  $h_p(T_{i+1...i+n}) = (2h_p(T_{i...i+n-1}) T_ih_p(2^n) + T_{i+n}) \mod p$ .
  - ① Set  $S = \emptyset$ . Compute  $h_p(T_{1...n})$ ,  $h_p(2^n)$ , and  $h_p(P)$ .
  - ② For each i = 1, ..., m n + 1
    - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
    - **2** Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$  by applying rolling hash.

## Running Time

• In Step 1, computing  $h_p(x)$  for an n bit x is in O(n) time.

- p: a random prime from  $\{1, \ldots, M\}$ .
- Rolling hash:  $h_p(T_{i+1...i+n}) = (2h_p(T_{i...i+n-1}) T_ih_p(2^n) + T_{i+n}) \mod p$ .
  - ① Set  $S = \emptyset$ . Compute  $h_p(T_{1...n})$ ,  $h_p(2^n)$ , and  $h_p(P)$ .
  - ② For each  $i = 1, \ldots, m n + 1$ 
    - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
    - **2** Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$  by applying rolling hash.

## Running Time

- In Step 1, computing  $h_p(x)$  for an n bit x is in O(n) time.
- Assuming  $O(\lg M)$  bit arithmetic can be done in O(1) time,
  - Since  $h_p(.)$  produces  $\lg M$  bit numbers, both steps inside for loop can be done in O(1) time.

- p: a random prime from  $\{1, \ldots, M\}$ .
- Rolling hash:  $h_p(T_{i+1...i+n}) = (2h_p(T_{i...i+n-1}) T_ih_p(2^n) + T_{i+n}) \mod p$ .
  - ① Set  $S = \emptyset$ . Compute  $h_p(T_{1...n})$ ,  $h_p(2^n)$ , and  $h_p(P)$ .
  - ② For each i = 1, ..., m n + 1
    - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
    - **2** Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$  by applying rolling hash.

## Running Time

- In Step 1, computing  $h_p(x)$  for an n bit x is in O(n) time.
- Assuming  $O(\lg M)$  bit arithmetic can be done in O(1) time,
  - Since  $h_p(.)$  produces  $\lg M$  bit numbers, both steps inside for loop can be done in O(1) time.
  - Overall O(m+n) time.

- p: a random prime from  $\{1, \ldots, M\}$ .
- Rolling hash:  $h_p(T_{i+1...i+n}) = (2h_p(T_{i...i+n-1}) T_ih_p(2^n) + T_{i+n}) \mod p$ .
  - ① Set  $S = \emptyset$ . Compute  $h_p(T_{1...n})$ ,  $h_p(2^n)$ , and  $h_p(P)$ .
  - ② For each i = 1, ..., m n + 1
    - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
    - 2 Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$  by applying rolling hash.

## Running Time

• In Step 1, computing  $h_n(x)$  for an *n* bit x is in O(n) time.

Assuming  $O(\lg M)$  bit arithmetic can be done in O(1) time,

- Since  $h_p(.)$  produces  $\lg M$  bit numbers, both steps inside for **loop** can be done in O(1) time.
- Overall O(m+n) time. Can't do better. Ruta (UIUC)

26 / 30

- **1** For each i = 1, ..., m n + 1
  - **1** If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - Ompute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

#### Lemma

If match at any position i then  $i \in S$ . In otherwords if  $T_{i...i+n-1} = P$ , then  $i \in S$ .

All matched positions are in **S**.

- **1** For each i = 1, ..., m n + 1
  - **1** If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - ② Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

#### Lemma

If match at any position i then  $i \in S$ . In otherwords if  $T_{i...i+n-1} = P$ , then  $i \in S$ .

All matched positions are in S.

Can it contain unmatched positions?

- **1** For each i = 1, ..., m n + 1
  - **1** If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - ② Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

#### Lemma

If match at any position i then  $i \in S$ . In otherwords if  $T_{i...i+n-1} = P$ , then  $i \in S$ .

All matched positions are in S.

Can it contain unmatched positions? YES!

- **1** For each i = 1, ..., m n + 1
  - **1** If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - ② Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

#### Lemma

If match at any position i then  $i \in S$ . In otherwords if  $T_{i...i+n-1} = P$ , then  $i \in S$ .

All matched positions are in **S**.

Can it contain unmatched positions? YES! With what probability?

Pr[S contains an index i, while there is no match at i]

- **1** For each i = 1, ..., m n + 1
  - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - $oldsymbol{o}$  Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

Pr[S contains an index i, while there is no match at i]

- For each i = 1, ..., m n + 1
  - **1** If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - $oldsymbol{o}$  Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

Set 
$$M = \lceil 2(sn) \lg sn \rceil$$
. Given  $x \neq y$ ,  $\Pr[h_p(x) = h_p(y)] \leq 1/s$ .

Pr[S contains an index i, while there is no match at i]

- **1** For each i = 1, ..., m n + 1
  - **1** If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - lacktriangle Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

Set 
$$M = \lceil 2(sn) \lg sn \rceil$$
. Given  $x \neq y$ ,  $\Pr[h_p(x) = h_p(y)] \leq 1/s$ .

False positive: Pr[S contains an i, while no match at i]

Pr[S contains an index i, while there is no match at i]

- **1** For each i = 1, ..., m n + 1
  - **1** If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - ② Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

Set 
$$M = \lceil 2(sn) \lg sn \rceil$$
. Given  $x \neq y$ ,  $\Pr[h_p(x) = h_p(y)] \leq 1/s$ .

## False positive: Pr[S contains an i, while no match at i]

• Given  $T_{i...i+n-1} \neq P$ ,  $\Pr[i \in S] \leq 1/s$ .

Pr[S contains an index i, while there is no match at i]

- **1** For each i = 1, ..., m n + 1
  - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - lacktriangle Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

Set 
$$M = \lceil 2(sn) \lg sn \rceil$$
. Given  $x \neq y$ ,  $\Pr[h_p(x) = h_p(y)] \leq 1/s$ .

## False positive: Pr[S contains an i, while no match at i]

- Given  $T_{i...i+n-1} \neq P$ ,  $\Pr[i \in S] \leq 1/s$ .
- **Pr**[Any index in *S* is wrong]

Pr[S contains an index i, while there is no match at i]

- **1** For each i = 1, ..., m n + 1
  - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - 2 Compute  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

Set 
$$M = \lceil 2(sn) \lg sn \rceil$$
. Given  $x \neq y$ ,  $\Pr[h_p(x) = h_p(y)] \leq 1/s$ .

## False positive: Pr[S contains an i, while no match at i]

- Given  $T_{i...i+n-1} \neq P$ ,  $\Pr[i \in S] < 1/s$ .
- Pr[Any index in S is wrong] < m/s (Union bound).

Pr[S contains an index i, while there is no match at i]

- **1** For each i = 1, ..., m n + 1
  - If  $h_p(T_{i...i+n-1}) = h_p(P)$ , then  $S = S \cup \{i\}$ .
  - Occupate  $h_p(T_{i+1...i+n})$  using  $h_p(T_{i...i+n-1})$  and  $h_p(2^n)$ .

Set 
$$M = \lceil 2(sn) \lg sn \rceil$$
. Given  $x \neq y$ ,  $\Pr[h_p(x) = h_p(y)] \leq 1/s$ .

## False positive: Pr[S contains an i, while no match at i]

- Given  $T_{i...i+n-1} \neq P$ ,  $\Pr[i \in S] \leq 1/s$ .
- $Pr[Any index in S is wrong] \le m/s$  (Union bound).
- To ensure S is correct with at least 0.99 probability, we need

$$1 - \frac{m}{s} = 0.99 \Leftrightarrow \frac{m}{s} = \frac{1}{100} \Leftrightarrow s = 100m$$

Back to running time

## Running Time

• In Step 1, computing  $h_p(x)$  for an n bit x is in O(n) time.

Assuming  $O(\lg M)$  bit arithmetic can be done in O(1) time,

- Since  $h_p(.)$  produces  $\lg M$  bit numbers, both steps inside for loop can be done in O(1) time.
- Overall O(m+n) time. Can't do better.

$$M = \lceil 200mn \lg 100mn \rceil \Rightarrow \lg M = O(\lg m)$$

Back to running time

## Running Time

• In Step 1, computing  $h_p(x)$  for an n bit x is in O(n) time.

Assuming  $O(\lg M)$  bit arithmetic can be done in O(1) time,

- Since  $h_p(.)$  produces  $\lg M$  bit numbers, both steps inside for loop can be done in O(1) time.
- Overall O(m+n) time. Can't do better.

$$M = \lceil 200mn \lg 100mn \rceil \Rightarrow \lg M = O(\lg m)$$

Even if T is entire Wikipedia, with bit length  $m \approx 2^{38}$ ,

Back to running time

## Running Time

• In Step 1, computing  $h_p(x)$  for an n bit x is in O(n) time.

Assuming  $O(\lg M)$  bit arithmetic can be done in O(1) time,

- Since  $h_p(.)$  produces  $\lg M$  bit numbers, both steps inside for loop can be done in O(1) time.
- Overall O(m+n) time. Can't do better.

$$M = \lceil 200mn \lg 100mn \rceil \Rightarrow \lg M = O(\lg m)$$

Even if T is entire Wikipedia, with bit length  $m \approx 2^{38}$ ,

 $\lg M pprox 64$  (assuming bit-length of  $n \leq 2^{16}$ )

Back to running time

## Running Time

• In Step 1, computing  $h_p(x)$  for an n bit x is in O(n) time.

Assuming  $O(\lg M)$  bit arithmetic can be done in O(1) time,

- Since  $h_p(.)$  produces  $\lg M$  bit numbers, both steps inside for loop can be done in O(1) time.
- Overall O(m+n) time. Can't do better.

$$M = \lceil 200mn \lg 100mn \rceil \Rightarrow \lg M = O(\lg m)$$

Even if T is entire Wikipedia, with bit length  $m \approx 2^{38}$ ,

$$\lg M pprox 64$$
 (assuming bit-length of  $n \leq 2^{16}$ )

64-bit arithmetic is doable on laptops!

## Take away points

- Hashing is a powerful and important technique. Many practical applications.
- Randomization fundamental to understand hashing.
- Good and efficient hashing possible in theory and practice with proper definitions (universal, perfect, etc).
- Related ideas of creating a compact fingerprint/sketch for objects is very powerful in theory and practice.

Ruta (UIUC) CS473 30 Spring 2021 30 / 30