

1. Fingerprinting and string comparison + pattern matching
2. primality testing

Schwartz-Zippel thm.

$$\Pr(Q(v_1, \dots, v_n) = 0 \mid Q \neq 0) \leq \frac{\deg(Q)}{|S|}$$

for  $v_i \in S$ .

**Problem.** Verify whether two strings  $X, Y \in \{0,1\}^n$  are equal.

deterministically  $O(n)$

$$X = (x_1, \dots, x_n)$$

$$Y = (y_1, \dots, y_n)$$

USE S-Z theorem

interpret  $X$  and  $Y$  as multivariable polynomials:

$$X(z_1, \dots, z_n) = \sum_{i=1}^n x_i z_i \pmod{p}$$

$$Y(z_1, \dots, z_n) = \sum_{i=1}^n y_i z_i \pmod{p}$$

$$X(z_1, \dots, z_n) - Y(z_1, \dots, z_n) \stackrel{?}{=} 0$$

choose  $r \in \{0,1\}^n$  ←

by Schwartz-Zippel theorem

$$\Pr(X(v_1, \dots, v_n) - Y(v_1, \dots, v_n) \mid [X-Y](\vec{z}) \neq 0) \leq \frac{\deg[X-Y]}{2} = \frac{1}{2}$$

## Context: Database comparison

- > two distant databases  $X$  and  $Y$  are they the same?
- > is the above method efficient in number of transmitted bits?  
NO! random  $r$  needed to calculate the fingerprint is as long as the database.

## Solution 1

Interpret both  $X$  and  $Y$  as numbers

$$\text{num}(X) = \sum_{i=1}^n x_i 2^{i-1}$$

$$\text{num}(Y) = \sum_{i=1}^n y_i 2^{i-1}$$

Compare:

$X \bmod p$  with  $Y \bmod p$

if  $p$  is well chosen, fingerprints are small and probability of error is also small.

When does this give a wrong answer? ( $X \neq Y$ , but  $X \equiv Y \pmod{p}$ )

$X - Y \equiv 0 \pmod{p}$  (read as  $X - Y$  is divisible by  $p$ ).

$\pi(k)$  - all primes smaller than  $k$ .  $\pi(k) \approx \frac{k}{\ln k}$

$$\text{for } k \geq 29 \quad \pi(k) \leq (1.2 \dots) \frac{k}{\ln k}$$

$$\frac{1}{\pi(z)} \approx \frac{n}{k}$$

$$Pr(X - Y \equiv 0 \pmod{p} \mid X \neq Y) = \frac{\# \text{ bad primes}}{\# \text{ primes we choose from}} = \frac{n}{\pi(z)} \leq \frac{\ln k \cdot n}{k}$$

# bad primes: How many divisors can  $X - Y$  have at most?

What is the largest value of  $X - Y$ ?

$$X - Y \leq 2^n$$

Smallest number with  $n$  prime divisors

$$= \prod_{i=1}^n p_i > 2^n = \prod_{i=1}^n 2$$

$$p_i > p_{i-1} > 2$$

$p_i$  -  $i$ th smallest prime

$$\# \text{ bad primes} < n$$

for  $k = \lceil t \cdot n \cdot \log(t \cdot n) \rceil$

$$Pr < \frac{\ln(t \cdot n \cdot \log(t \cdot n)) \cdot n}{t \cdot n \cdot \log(t \cdot n)} \in O\left(\frac{1}{t}\right)$$

How many bits do we need to send? for  $t = n$  we need to send prime  $p$   $O(\log n)$  bits and the hash  $O(\log n)$

What we did:

$$X = \sum_{i=0}^{n-1} x_i \cdot z^{i-1} \pmod{p}$$

method 1:

choose  $z = 2$  and randomize over  $p$

method 2:

choose  $p$  and randomize over  $z$

Method 2 can be analysed using polynomial comparison.

$X(z)$  and  $Y(z)$  are polynomials, by  $\mathbb{F}_z$

$$Pr[(X-Y)(z) = 0 \mid (X-Y)(z) \neq 0] \leq \frac{\deg(X-Y)}{|S|} = \frac{n-1}{|S|}$$

$\forall z \in S$  to match method 1 we want this to be roughly  $\frac{1}{n}$

$\Rightarrow |S|$  is roughly  $n^2$ .  $\therefore p$  needs to be larger than  $n^2$ .

How many bits do we need to send?

Prime  $p \sim O(\log n)$

number  $t < p \sim O(\log n)$

3rd method:

choose a random polynomial  $P \pmod p$  and evaluate  $P(\text{num}(X))$

and  $P(\text{num}(Y))$  and compare.  $\leadsto$  this is an idea behind universal hashing.

## Pattern matching

$X$  - a string of length  $n$

$Y$  - a string of length  $m$  with  $m < n$   $\mid$  wlog  $X, Y \in \{0,1\}^n$

Is  $Y$  a substring of  $X$ ?

Is  $T$  a substring of  $A$  :

Naive algorithm  $\approx O(m \cdot n)$  comparisons

Better solutions (Knuth-Morris)  $\approx O(m+n)$  comparisons

Rabin-Karp in  $O(m+n)$  time  $\leftarrow$  Monte Carlo

Main idea: Compare fingerprints.

Imagine calculating fingerprints is for free. How many comparisons?

Each fingerprint is  $O(\log m)$  bits long

$$\approx O(n \log m)$$

↓  
this is not what is meant  
in the slides.

Strings are arrays of objects :

both  $X$  and  $Y$  are in memory of a computer in an indexed array,  
that is each  $x_i$  and  $y_i$  need to be addressed separately in the  
memory. This is the expensive operation.

So in analysis of Rabin Karp algorithm, the expensive operation  
is calculating the hash  $\text{num}(Y) = \sum_{i=1}^m x_i 2^{m-i} \bmod p$  because  
need to access each bit of array  $Y$ .  $\approx O(m)$

Comparison of fingerprints is comparison of two integers, therefore  
in  $O(1)$ .

Naive calculation of hashes results in  $O(m \cdot n)$

What we need is cheaper fingerprint calculation:

$$\text{let } X_j = (x_j, \dots, x_{j+m-1})$$

$$F(x_j) = x_j \cdot 2^{m-1} + x_{j+1} \cdot 2^{m-2} + \dots + x_{j+m-1}$$

$$F(x_{j+n}) = x_{j+n} \cdot 2^{m-1} + x_{j+n+1} \cdot 2^{m-2} + \dots + x_{j+n+m-1}$$

$$F(x_{j+n}) = 2 \cdot [F(x_j) - 2^{m-1} \cdot x_j] + x_{j+n}$$

Time analysis

The hash of $\gamma$ : $F(\gamma)$	$m$ steps	} $O(m+n)$
First hash of $X_n$ : $F(X_n)$	$m$ steps	
$n$ following hashes each in $O(1)$	$n$ steps	