

# Finger printing and string comparison

## Schwartz-Zippel theorem

$$\Pr(Q(r_1, \dots, r_n) = 0 \mid Q \neq 0) \leq \frac{\deg Q}{|S|}$$

$$r_i \in_R S$$

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

Deterministically $O(n)$	$X = (x_1, \dots, x_n) \quad x_i, y_i \in \{0, 1\}$ $Y = (y_1, \dots, y_n)$
--------------------------	--

If comparison is an expensive operation, then S-Z theorem gives us solution:

→ Interpret  $X$  and  $Y$  as polynomials:

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

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

$$X(\vec{z}) - Y(\vec{z}) \stackrel{?}{=} 0$$

Choose  $\vec{v} \in \{0, 1\}^n$

by S-Z theorem

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

## Context: Database comparison

- Two distant databases  $X$  and  $Y$ . Are they the same?
- Expensive operation: transmitting a bit (sending messages between the databases).

Is the method above efficient? ↑

NO! Random  $r$  needs to be distributed and it is as long as the database!

## Solution 1

Interpret both  $X$  and  $Y$  as numbers:

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

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

Compare

$X \bmod p$  and  $Y \bmod p$  fingerprints

for suitably <sup>randomly</sup> chosen prime  $p$ . (smaller than  $\epsilon$ )  
If  $p$  is small, fingerprints are small (in bits). However there is a tradeoff between the size of  $p$  and the probability of an error.

Error can happen if  $X \neq Y$  but  $X \equiv Y \pmod p$

$$\boxed{X - Y \equiv 0 \pmod p} \quad (\text{read } X - Y \text{ is divisible by } p)$$

$\pi(k)$  - number of primes smaller than  $k$ .

$$\pi(k) = O\left(\frac{k}{\ln k}\right) \leftarrow$$

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

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

$$\Pr(X - Y \equiv 0 \pmod{p} \mid X \neq Y) = \frac{\# \text{ bad primes}}{\# \text{ primes we chose from}} \stackrel{?}{=} \frac{n}{\pi(k)} < \frac{\ln k \cdot \ln}{k \cdot (1.2)}$$

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

What is the largest value of  $X - Y$ ?  $X - Y < 2^h$

What is smallest number with  $n$  prime divisors?

$$\prod_{i=1}^n p_i > 2^n = \prod_{i=1}^n 2 \Rightarrow \# \text{ bad primes} < n$$

$p_i$  -  $i^{\text{th}}$  smallest prime

for  $k = t \cdot n \ln(tn)$

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

How many bits does  $X$  need to send to  $Y$ ?

for  $t = n$  a prime of  $O(\log_2(n))$  bits and the fingerprint  $O(\log_2 n)$

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

Solution 1:

Choose  $z = 2$  and randomize over  $p$ .

Solution 2:

Choose  $p$  and randomize over  $z$

**Method 2** analysis using S-Z theorem

$X(z)$  and  $Y(z)$  are polynomials mod  $p$

$$r \in_R S \subseteq \mathbb{Z}/p$$

$$\Pr \left\{ (X-Y)(r) = 0 \pmod{p} \mid (X-Y)(z) \neq 0 \right\} \leq \frac{\deg(X-Y)}{|S|} = \frac{n-1}{|S|}$$

To match the method 1, we would like this probability to be roughly  $1/n$

$$\Rightarrow |S| = n^2 \Rightarrow p \text{ needs to be larger than } n^2$$

What needs to be sent?

$r < p \approx O(\log(n))$  bits  
and  $X(r) \pmod{p} \approx O(\log(n))$  bits } twice as large

**3rd method:**

Choose a random polynomial  $P \pmod{p}$  and evaluate

$P(\text{num}(x))$  and  $P(\text{num}(y))$  and compare.

⇒ UNIVERSAL HASHING