

# Fingerprinting II and string comparison

## Schwartz-Zippel theorem

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

$$r_i \in R \text{ S}$$

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

Deterministically  $O(n)$

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

$$x_i, b_i \in \{0,1\}$$

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

If comparison is an expensive operation, then SZT gives us the following solution:

→ Interpret  $X$  and  $Y$  as polynomials

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

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

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

(choose  $\vec{r} \in \{0,1\}^n$ )

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

## Context: Database comparisons

→ Two distant databases  $X$  and  $Y$ . Are they the same?

→ Expensive operation: transmitting a bit

Is the method above efficient?

NO. Random  $r$  needs to be sent and it is as large as the database

## Solution 1

Interpret  $X$  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

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

for a suitably randomly chosen prime  $p$ . Specifically  $p < k$ .

If  $p$  is small, fingerprints are also small. However there is a trade-off between the size of  $p$  and the probability of an error.

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

$$x - y \equiv 0 \pmod{p}$$

( $x - y$  is divisible by  $p$ )

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

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

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

$$P_r(x - y = 0 \pmod{p} \mid x \neq y) = \frac{\# \text{ bad primes}}{\# \text{ primes we check from}} \leq \frac{4}{\pi(k)} \leq \frac{\ln k \cdot 4}{1.2 \cdot k}$$

# Bad primes: How many primes can divide  $x - y$ ?

$$x - y < 2^n$$

What is the smallest number with  $n$  distinct prime divisors?

$$\prod_{i=1}^n p_i > 2^n = \frac{1}{4} 2^{2n}$$

$p_i$  -  $i$ th smallest prime

$$\text{for } k = t \cdot n \ln(t \cdot n)$$

$$P_r < \frac{\ln(t \cdot n \ln(t \cdot n)) \cdot 4}{t \cdot n \ln(t \cdot n)} \in O\left(\frac{1}{t}\right)$$

How many bits  $x$  needs to send to  $y$ ?

$$\text{for } t = n \quad \text{a prime of } \log_2(2) \text{ bits} \quad k = n^2 \cdot 2 \ln(n)$$

$$\in O(\log n)$$

and the fingerprint of the same size.

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

Solution 1

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

Solution 2

Choose  $p$  and randomize over  $z$

$$v \in_{\mathbb{R}} \mathbb{S} \subseteq \mathbb{Z}_p$$

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

To compare with solution 1, we would like the probability of success to be roughly  $\frac{1}{n}$ .

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

What needs to be sent?

$$v < p \approx O(\log n) \text{ bits}$$

$$\text{and } X(v) \pmod p \quad O(\log n) \text{ bits}$$

3<sup>rd</sup> method

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

$$P(\text{num}(X)) \text{ and } P(\text{num}(Y)) \text{ and compare}$$

⇒ UNIVERSAL HASHING

$$\mathcal{H}: X \rightarrow Y \quad c = \Pr(f(M) = c) = \frac{|\mathcal{H}|}{|Y|}$$