# COL 351: Analysis and Design of Algorithms

## Lecture 22

# Pattern Matching

**Given:** String $T = (t_{n-1}, \ldots, t_1, t_0)$ and a pattern $X = (x_{k-1}, \ldots, x_1, x_0)$, both binary.

**Find:** If there exists a sub-string of T that is identical to X.

Pattern

$X = (x_{k-1}, \ldots, x_1, x_0)$

$N_X = 2^{k-1}x_{k-1} + \cdots + 2^1 x_1 + 2^0 x_0$

(decimal form of $X$)

$X = 11001$

$N_X = 16 + 8 + 1 = 25$

$T = (t_{n-1}, \ldots, t_1, t_0)$

$N_T(j) = 2^{k-1}t_{j+k-1} + \cdots + 2^1 t_{j+1} + 2^0 t_j$

(decimal form of $(t_{j+k-1}, \cdots, t_{j+1}, t_j)$)

# Algorithm

Flag= False

**For** $j = 0$ to $(n - k)$:

      **If** $N_X = N_T(j)$ **then**

            Flag $=$ True

Return Flag

Time $= \ \mathrm{O}(nk)$

# Algorithm

$p =$ random prime in range $[2, n^4]$.

**Hash Function** $H : z \rightarrow z \mod p$

Flag= False

**For** $j = 0$ to $(n - k)$:

      **If** $H(N_X) = H(N_T(j))$ **then**

           Flag = True

Return Flag

Show:

- Answer returned is correct with probability $(1 - 1/n)$.
- Implementation in $O(n)$ time.

# Computing random prime in range $[2, n^4]$

> **Prime Number Theorem:** Number of primes in the range [2,L] is $\Theta\left(\dfrac{L}{\log L}\right)$.

- Probability $\left(\begin{array}{l}\text{a random number in} \\ \text{range } [2, L] \text{ is prime}\end{array}\right) \geq \dfrac{c}{\log L}$

## AKS Primality Test (By Agarwal, Kayal, Saxena)

- Checks if a number $n$ is prime or not in $O(\log^c(n))$ time, for some fixed $c \geq 1$.

# Observations

**Claim 1:** For any integer $z \leqslant 2^k$, the number of distinct prime factors of $z$ is at most $k$.

$$z = \underbrace{p_1}_{} \cdot \underbrace{p_2}_{} \circ \cdots \bullet \underbrace{p_\alpha}^{i_1 \quad i_2 \quad i_\alpha} \leq 2^k$$

$$\geqslant 2$$

$$\Rightarrow \alpha \leq k$$

**Claim 2:** For any $j$, the number of distinct prime factors of $(N_T(j) - N_X)$ is at most $n$.

$$N_T(j) \leq 2^n$$

$$N_X \qquad \leq 2^k \leq 2^n$$

$$\therefore N_T(j) - N_X \leq 2^n$$

$$\Rightarrow \text{Distinct Prime factors} \leq n$$

**Claim 3:** For any $j \in [0, n-k]$ with $N_T(j) \neq N_X$.

$$\mathbf{Prob}\left( p \text{ divides } \underbrace{(N_T(j) - N_X)}_{} \right) \leqslant \frac{1}{n^2}.$$

Probability of error
at location $j$

$$\frac{\text{No of Prime factors of } N_T(j) - N_X}{\text{No of choices for 'p'}}$$

$$\leq \frac{n}{\Theta(n^4 / \log n^4)} \leq \frac{1}{n^2}$$

**Claim 4:** If $X$ is not a substring of $T$. Then

$$\mathbf{Prob}\left( \exists \, j, \text{ such that } H(N_T(j)) = H(N_X) \right) \leqslant \frac{1}{n}.$$

By union bound,

$$\text{Prob}\left( \text{Error in entire algo} \right) \leq \sum_{j=0}^{n-k} \text{Prob}\left( \text{Error at location } j \right) \leq \frac{n-k+1}{n^2} \leq \frac{1}{n}$$

# How to recursively compute $N_T(j)$?

$$N_T(j) \;=\; \text{Decimal form of } (t_{j+k-1}, \cdots, t_{j+1}, t_j)$$

Removed

$$N_T(j+1) \;=\; \text{Decimal form of } (t_{j+k}, t_{j+k-1}, \cdots, t_{j+1})$$

Added

$$N_T(j+1) = \frac{N_T(j) - t_j}{2} \;+\; 2^{k-1}\, t_{j+k}$$

# Recursively Computing Hash of $N_T(j)$

$$N_T(j+1) = \frac{N_T(j) - t_j}{2} + 2^{k-1} t_{j+k}$$

$$H(N_T(j+1)) = N_T(j+1)$$

$$= \left( (N_T(j) - t_j) \cdot 2^{-1} - 2^{k-1} \cdot t_{j+k} \right) \bmod p$$

$$= \left( \underbrace{N_T(j) \bmod p}_{H(N_T(j))} - \underbrace{t_j \bmod p}_{O(1) \text{ time}} \right) \cdot \underbrace{2^{p-2} \bmod p}_{A} - \underbrace{2^{k-1} \bmod p}_{B} \cdot \overbrace{t_{j+k} \bmod p}^{O(1) \text{ time}} \bmod p$$

Claim : $H(N_T(j+1))$ can be obtained from $H(N_T(j))$ in

$O(1)$ space, $O(1)$ time if we know $A, B$.

# Efficiently Computing A and B

$$A = 2^{p-2} \mod p \qquad\qquad B = 2^{k-1} \mod p$$

$$2^z \pmod{p} = \begin{cases} \left(2^{\lfloor z/2 \rfloor} \pmod{p}\right)^2 \mod p & \text{if } z \text{ is even} \\[2em] \left(2^{\lfloor z/2 \rfloor} \pmod{p}\right)^2 \cdot 2 \mod p & \text{if } z \text{ is odd} \end{cases}$$

Claim: If we know $2^{\lfloor z/2 \rfloor} \pmod{p}$ then we can compute $2^z \pmod{p}$ in $O(1)$ space, time.

Computing $A$ — $O(\log p) = O(\log n)$ Time, $O(1)$ Space

Computing $B$ — $O(\log k) = O(\log n)$ Time, $O(1)$ Space

# Base Case

$$H\left(N_T(0)\right) = \text{Decimal}\left(t_{k-1} \quad \cdots \quad t_1 \quad t_0\right) \bmod p$$

with $t_{k-1} \to 2^{k-1}$, $t_1 \to 2^1$, $t_0 \to 2^0$

We can compute $2^{i+1} \bmod p$ from $2^i \bmod p$ in $O(1)$ space, $O(1)$ query time.

Claim: We can compute $H(N_T(0))$ and $H(N_x)$ in $O(k)$ time and $O(1)$ space.