

# Polynomial hashes

$$\mathbb{Z}_p^d \rightarrow \mathbb{Z}_p \rightsquigarrow X = (x_1, x_2, x_3, \dots, x_d)$$

$$\text{Def: } h_a(x) = \sum_{i=1}^d x_i \cdot a^{i-1} \\ = x_1 + x_2 \cdot a + x_3 \cdot a^2 + \dots + x_d \cdot a^{d-1}$$

$\approx$  polynomial with coeffs.  $x_i$   
evaluated at point  $a$

⊙  $h_a(x)$  can be calculated in  $O(d)$  time

Def family of polynomial hash functions  
 $\mathcal{R} = \{h_a \mid a \in \mathbb{Z}_p\}$

⊙ Only  $O(1)$  time to pick  $h_a$   
(vs.  $\Theta(d)$  for scalar product family)

Thm  $\mathcal{R}$  is  $d$ -universal

Proof

Fact polynomial  $p$  of degree  $d$   
has  $\leq d$  roots unless  $p \equiv 0$

• let  $x, y \in \mathbb{Z}_p^d$

$$\underbrace{h_a(x)}_{\text{poly}(a)} = \underbrace{h_a(y)}_{\text{poly}(a)} \Leftrightarrow \underbrace{h_a(x) - h_a(y)}_{\text{poly}(a)} = 0$$

$$\text{def } p(t) = h_t(x) - h_t(y) = \\ = \sum_{i=1}^d (x_i - y_i) t^{i-1}$$

!  $t$  is variable !  
 $x, y$  constants

•  $x \neq y \Rightarrow \exists i : x_i \neq y_i \Rightarrow p(t) \neq 0$

•  $h_a(x) = h_a(y) \Leftrightarrow a$  is a root of  $p(t)$

$\Rightarrow p(t)$  has  $\leq d$  roots

$$\Rightarrow \Pr[h_a(x) = h_a(y)] \leq \frac{d}{p} \quad \square$$

• Compose with linear functions to get 2-independence  $\leftarrow$  lemma 6

Corollary for prime  $p$  and #buckets  $m$  s.t.  $p \geq 4 \cdot d \cdot m$

family  $\mathcal{R} = \mathcal{L}$  is  $(2, 5/2)$ -independent

# Variable length strings

- $x$  is a string of length  $\leq L$
- pad  $x$  with some character  $\#$  to get length exactly  $L$
- then use  $\mathcal{R}$
- $\#$  must not appear elsewhere in  $x$   
 $\rightsquigarrow \#$  is "new" character
- good choice is  $\# = 0$   
 (saves us work)

# Application: Rabin-Karp algorithm

input: long text  $\sigma$  (haystack)  
 short string  $r$  (needle)  
 $|r| = d, |\sigma| = n$

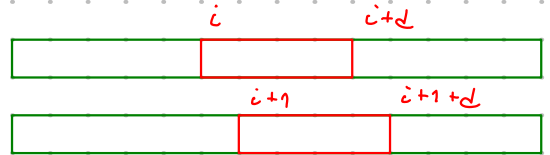
output: all occurrences of  $r$  in  $\sigma$

## Naive algorithm:

for  $i = 1 \dots (n-d)$ :  $\Rightarrow$  time  $O(n \cdot d)$

check  $\sigma[i:i+d] \stackrel{?}{=} r$   
 $O(d)$

→ sliding window



## Rabin-Karp:

- first compare hash of  $\sigma[i:i+d]$  with hash of  $r$
  - $h(r) \neq h(k) \Rightarrow r \neq k$
  - $h(r) = h(k) \Rightarrow$  check  $r \stackrel{?}{=} k$
- $O(1)$  on average for reasonable  $h$
- $\Rightarrow O(n+d + d \cdot (\# \text{ occurrences} + \# \text{ collisions}))$   
 (Also works well for multiple needles)

## Problem:

we need  $h(\sigma[i:i+d]) \rightarrow h(\sigma[i+1:i+1+d])$   
 in  $O(1)$  time

Solution: use polynomial (rolling) hash  $\mathcal{R}$

$$\begin{aligned}
 H_i &= h_a(\sigma[i:i+d]) = \sigma_i \cdot a^{d-1} + \sigma_{i+1} \cdot a^{d-2} + \dots + \sigma_{i+d-1} \cdot a^0 \\
 H_{i+1} &= h_a(\sigma[i+1:i+1+d]) = \sigma_{i+1} \cdot a^{d-1} + \dots + \sigma_{i+d-1} \cdot a^1 + \sigma_{i+d} \cdot a^0
 \end{aligned}$$

$$\left. \begin{aligned}
 & \dots + \sigma_{i+d-1} \cdot a^0 \\
 & \dots + \sigma_{i+d-1} \cdot a^1 + \sigma_{i+d} \cdot a^0
 \end{aligned} \right\} \begin{aligned}
 H_{i+1} &= a \cdot H_i \\
 & - \sigma_i \cdot a^d \\
 & + \sigma_{i+d}
 \end{aligned}$$