

# String Hashing

- Sidhant Bansal

Credits: CP-Algorithms, Codeforces Blogs

# Agenda

- 
1. Motivation
  2. What is String Hashing
  3. How to do String Hashing
  4. Analysis ←
  5. Applications
  6. Pitfalls

# Motivation

Why?

$$O(N^2) / O(N^3) \rightarrow \text{TLE}.$$

String Hashing lets you most often than not “optimise” a brute-force approach to a string related question.

How ?

Well at a very high-level, think of it like it speeds-up the “comparison of any two strings”.

Earlier comparing two strings (are equal or not or which one is smaller/larger) used to be **O(N)** now it is **O(1) (NOTE: correct with high probability, NOT 100% of the times)** or **O(logN)** for smaller/larger lexical comparison.

$$n, m \rightarrow O(n+m)$$



$O(1)$  or  $O(\log N)$

## Most Powerful String DS

There are many string DS:

- KMP, Suffix Array, Suffix Trees, Palindromic Trees, Aho-Cohrasik, 2, --

However all problems  $\leq$  Medium-Hard difficulty of these Tags are most often than not solvable by String Hashing as well. (As an alternate solution)

# So why doesn't everyone use it?

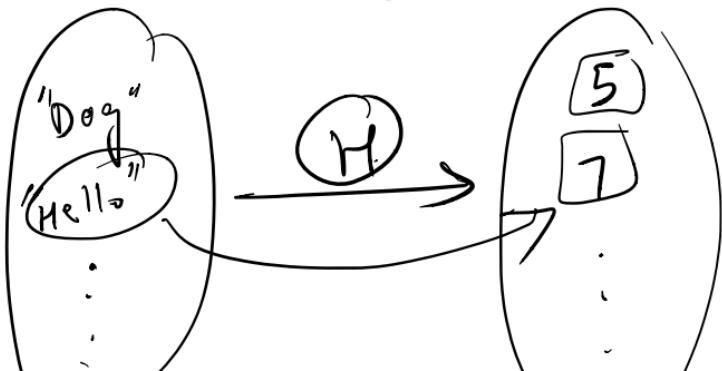
- It can be slightly tedious to implement. And you have to be smart about the modulos, primes, etc you pick.
- Basically your implementation needs to be clean for it to pass (i.e not give WA neither give TLE)

# What is String Hashing

# What is String Hashing?

The goal is to be able to have a “fast” method to map any string to an integer,  
such that

- If  $s_1 \neq s_2$ , then  $H(s_1) \neq H(s_2)$
- If  $s_1 == s_2$ , then  $H(s_1) == H(s_2)$
- Given  $s$ , we should be able to compute  $H(s)$  fast (i.e. in  $O(1)$  or  $O(\log N)$ )
- $H(s)$  should be an integer/long long/pair of ints/pair of long long  $\rightarrow O(1)$  in  
“size”



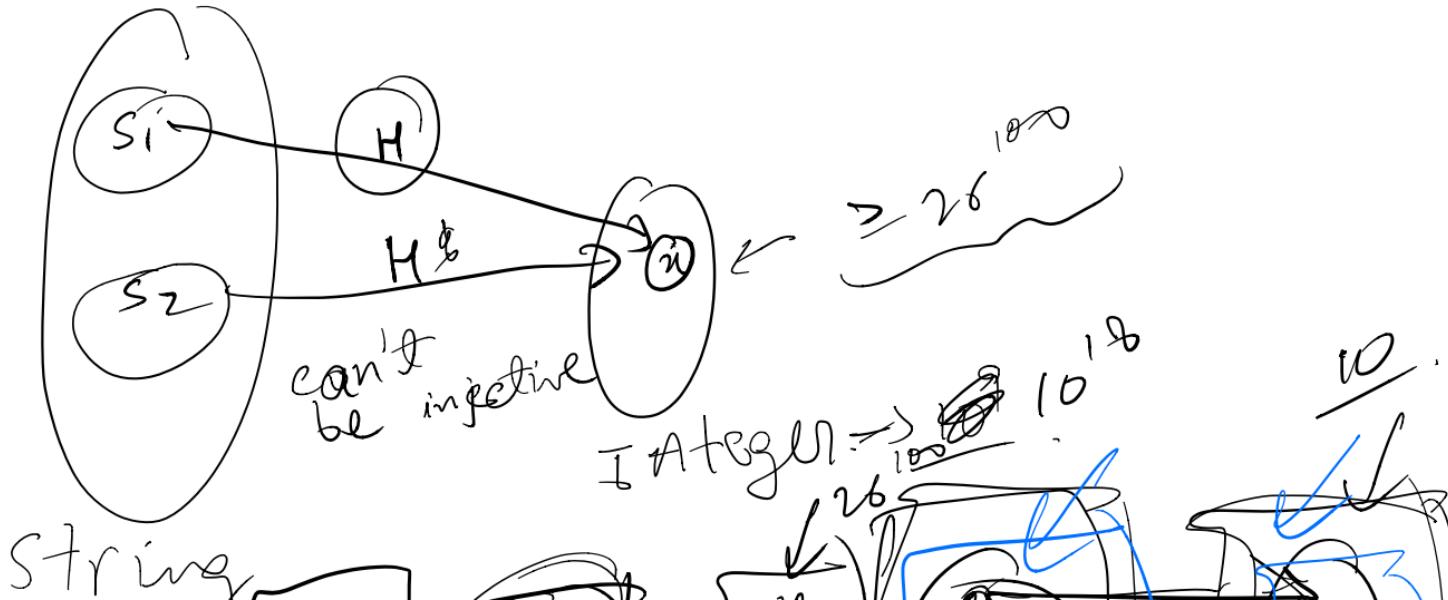
$5 \neq 7 \rightarrow "Dog" \neq "Hello"$

# String Poll Question

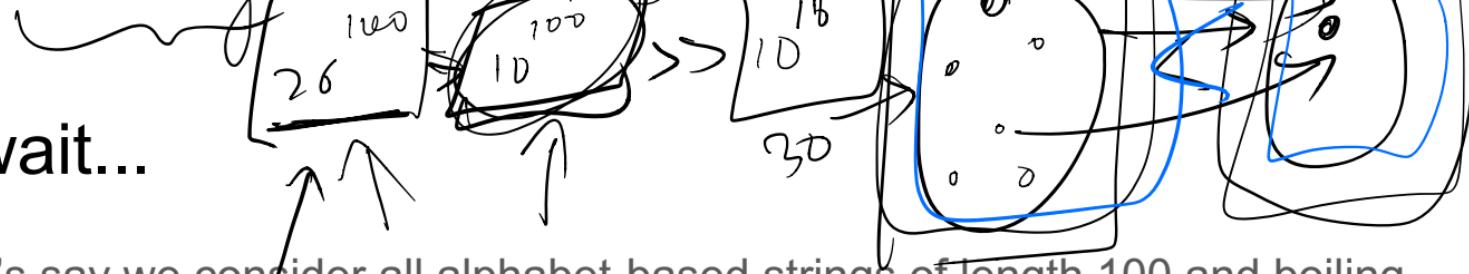
Integer

$\rightarrow H$ ?

- Do you think there exists a string hashing which can provide us with all that power mentioned in the previous slide? (True / False)



But wait...



- Let's say we consider all alphabet-based strings of length 100 and boiling them down to an integer.
- How many unique strings can we have?  $(26^{100})$
- And how many different integers do we have?  $(2^{32} \sim 10^9)$

But you can see  $26^{100} >> 10^9$

4 bytes = 32 bits  $\Rightarrow 2^{32}$

So how can we ensure that "if  $s_1 \neq s_2$ , then  $H(s_1) \neq H(s_2)$ "

$$[A-Z] \Rightarrow 26$$

$$\downarrow$$

---

$$26 \times 26 \times 26 \dots$$

$$H(s)^D \% MDP$$

$$H(n+S)$$

$$\times 26 = 26^{100}$$

$$H(A B C D) = \sum_{i=1}^{n-1} S_i \cdot p^i + S_n \cdot p^n \mod$$

$$[n + \phi(H(s))] \approx F(s)$$

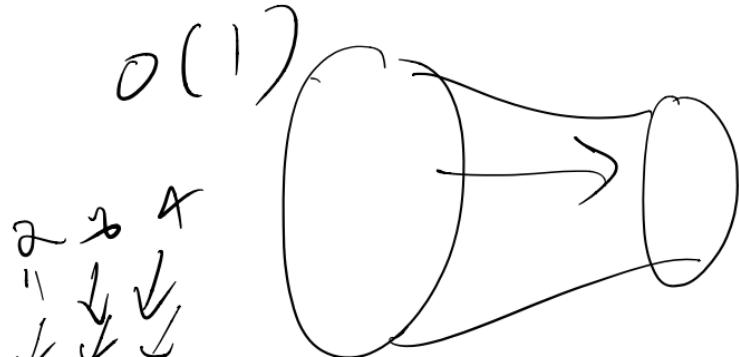
of MOD

$$= n + p \left( \sum_{i=0}^{n-1} s_i \cdot p^i \right)$$

## Fixed definition of String Hashing

The goal is to be able to have a “fast” method to map any string to an integer, such that

- If  $s_1 \neq s_2$ , then  $H(s_1) \neq H(s_2)$  (with high probability)  $(H(s_1) = H(s_2) \sim \text{low probability})$
- If  $s_1 == s_2$ , then  $H(s_1) == H(s_2)$  (with 100% probability)
- Given  $s$ , we should be able to compute  $H(s)$  fast (i.e. in  $O(1)$  or  $O(\log N)$ )
- $H(s)$  should be an integer/long long/pair of ints/pair of long long



$$(2^{14})_{2^6} = (2 \times 2^6 + 1 \times 2^6)^2$$

BAD

in base = 26

# How to do String Hashing

$$P(p=30) \left\{ \begin{array}{l} c_0 \times 3^0 + \\ c_1 \times 3^1 + \\ c_2 \times 3^2 + \\ \dots \end{array} \right.$$

"DAB"

$$\begin{aligned} H("DAB") &= (D \times 3^0 + A \times 3^1 + B \times 3^2) \% (10^9 + 7) \\ &= (4 \times 3^0 + 1 \times 3^1 + 2 \times 3^2) \% (10^9 + 7) \end{aligned}$$

$\rightarrow \text{MOD} \leftarrow \text{INTEGER}$

$$H("DAB") =$$

$D \times 3^0 \checkmark$

# Polynomial Rolling Hash Technique

- Given a string  $S = s_0 s_1 \dots s_{N-1}$

$$H(S) = \left( \sum_{i=0}^{N-1} (s_i * p^i) \% MOD \right)$$

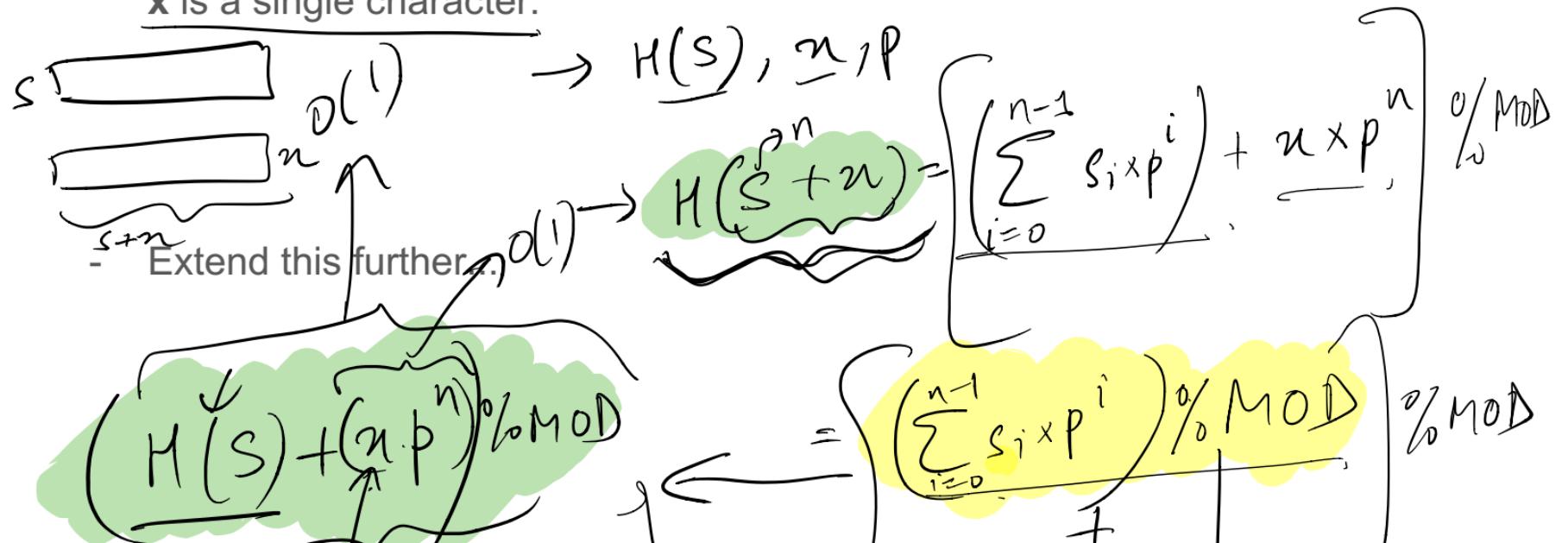
where  $p = 30 \rightarrow |\Sigma| = 2^6$   
 $\{A, Z\} \rightarrow 2^6$   
 $\{0 - 9\} \rightarrow 10$

- The hash is a polynomial in  $p$
- This hash depends on:
  - The length of the string
  - The order of the characters
  - And is contained within an integer (because of modulo arithmetic)

$$= (a \% m \% D) + (b \% p \% D)$$

## Why did we say "Rolling"

- If we know the hash of  $S$ , i.e  $H(S)$
- Then we can easily calculate the hash of  $S + x$  (i.e.  $H(S + x)$ ) or  $x + S$ , where  $x$  is a single character.



# Why did we say "Rolling"

- If we know the hash of every prefix of string  $S$ .  $O(N)$
- Let the prefixes of string  $S$  be  $P_0, P_1, P_2, \dots, P_{N-1}$
- Example:  $S = "ABCDEF"$ 
  - $P_0 = "A"$ ,  $P_1 = "AB"$ ,  $P_2 = "ABC"$ , ..
- Then if we know  $H(P_{i-1})$  then we can calculate  $H(P_i)$  in  $O(1)$
- So given  $H(P_i)$  for all  $i$ .
- Now we can calculate hash of any substring of  $S$ , i.e  $S[i\dots j]$ .
  - But how ?

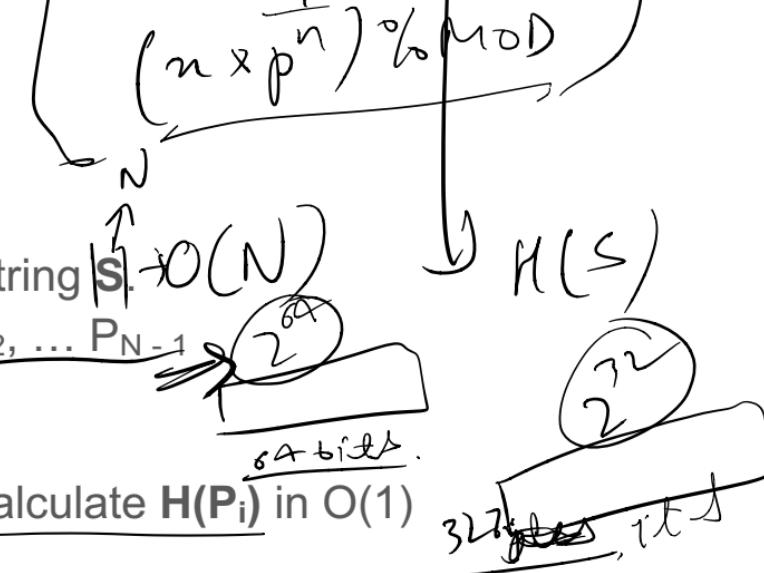
$S = \text{APPLE}$

$P_0 = A$

$P_1 = AP$

$P_2 = APP$

$$\begin{aligned} H(P_2), L &\Rightarrow H(P_2 + L) \\ &\Rightarrow H(APPL) \Rightarrow H(P_3) \end{aligned}$$



$$\begin{aligned} P_3 &= \text{A PPL} \\ P_4 &= \text{A PPLE} \end{aligned}$$

Explain how to get  $H(S[i \dots j])$  ?

$$H(S[i \dots j]) = \left( \sum_{k=i}^j s_k * p^{k-i} \right) \% MOD$$

$\left( p \times 30^0 + p \times 30^1 + L \times 30^2 \right) \% MOD$

$$= \left( \frac{H(P_j) - H(P_{i-1})}{p^i} \right) \% MOD = ((H(P_j) - H(P_{i-1})) p^{-i}) \% MOD$$

$$[H(P_i), H_i]$$

$$\left( \sum_{k=j}^n s_k * p^{k-n} \right) \% MOD$$

$\Rightarrow$   $\sum_{k=j}^n s_k * p^{k-n}$

$a = -10$

$b = 3$

$mod = 5$

$(-10 + 3) \% 5$

$\Downarrow$

-2 or 2

(A)

(B)

```

1 const int N = 1e6 + 6;
2 #define int long long;
3 const int mod = 1e9 + 7;
4 const int base = 33;
5
6 int add(int a, int b, int mod){  $\rightarrow (a+b) \% mod$ 
7     int res = (a + b) % mod;
8     if(res < 0) -
9         res += mod;
10    return res;
11 }
12
13 int mult(int a, int b, int mod){  $\rightarrow (axb) \% mod$ 
14     int res = (a * 1LL * b) % mod;
15     if(res < mod) -
16         res += mod;
17     return res;
18 }
19
20 int power(int a, int b, int mod){  $\rightarrow (a^b) \% mod$  in  $O(\log b)$ 
21     int res = 1;
22     while(b){
23         if((b % 2) == 1)
24             res = mult(res, a, mod);
25         a = mult(a, a, mod);
26         b /= 2;
27     }
28     return res;
29 }
```

$\rightarrow \text{int}(2^b)$

Russian

Multiplication

```
31 int pw[N];
32 int inv[N];
33 int hash[N];
34
35 void precalc() {
36     pw[0] = 1;
37     for(int i = 1; i < N; i++)
38         pw[i] = mult(pw[i - 1], base, mod);
39
40     int pw_inv = power(base, mod - 2, mod);
41     inv[0] = 1;
42     for(int i = 1; i < N; i++)
43         inv[i] = mult(inv[i - 1], pw_inv, mod);
44 }
45
46 void build(string s){
47     int n = s.length();
48     for(int i = 0; i < n; ++i){
49         hash[i] = add((i == 0) ? 0 : hash[i - 1], mult(pw[i], s[i] - 'a' + 1, mod), mod);
50     }
51 }
52
53 int getHash(int x, int y){
54     int res = add(hash[y], (x == 0) ? 0 : -hash[x - 1], mod);
55     res = mult(res, (x == 0) ? 1 : inv[x], mod);
56     return res;
57 }
```

$$H(S[0 \dots i]) = H(S[1 \dots i]) + \cancel{1 + \dots + 30} \pmod{M}$$

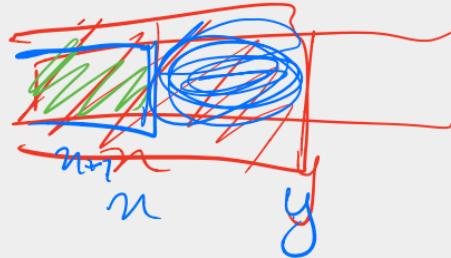
$a' - r'$

# Analysis

$$(A x^2 + B x^3 + C x^4 + \dots)$$

11

0.



$$p = 2$$

$\frac{\partial u}{\partial M D}$

$30(n) = 1/30$

## Analysis

$$n = \left(30^{-1}\right) \in \{0, \text{MOD}\} \text{ s.t. } \left[30 \cdot n = 1\right] \% \text{ MOD}$$

Fix the values first.

- Firstly, your  $p = \text{base} > \# \text{alphabets}$ . So generally keep it  $> 26$  or  $> 36$ .
- Your MOD should be a prime, so that calculation of  $(1/\text{base})$  under MOD is easy.

When mod is prime:  
(Using Euler's Theorem)

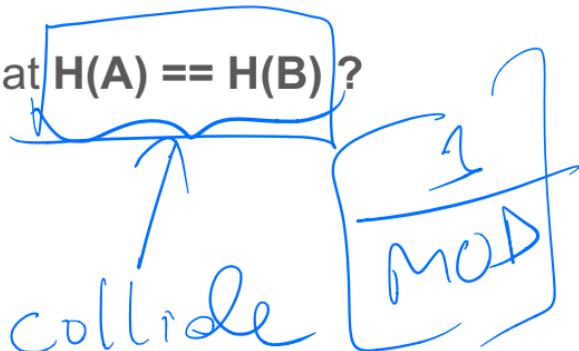
$$\left(\frac{1}{\text{base}}\right) \% \text{MOD} = (\text{base}^{\text{MOD}-2}) \% \text{MOD}$$

```
34
35 void precalc() {
36     pw[0] = 1;
37     for(int i = 1; i < N; i++)
38         pw[i] = mult(pw[i - 1], base, mod);
39
40     int pw_inv = power(base, mod - 2, mod);
41     inv[0] = 1;
42     for(int i = 1; i < N; i++)
43         inv[i] = mult(inv[i - 1], pw_inv, mod);
44 }
```

# Why high probability? (Poll Question)

Given two random string  $A \neq B$ , what is the probability that  $H(A) == H(B)$ ?

- A. 0
- B.  $1 / (\sqrt{MOD})$
- C.  $1 / MOD$
- D.  $1 / MOD^2$



collide  
"Bord"

$A = "Apple"$

$0, \frac{1}{\sqrt{MOD}}, \frac{1}{MOD}, \frac{1}{MOD^2}$

$MOD$

$Apple$   
 $orange$   
 $n$

But wait...it amplifies.

So for any given comparison the probability that we will get a false positive is  $1/\text{MOD}$ .

$$A, B \ (A \neq B) \rightarrow \Pr[\text{accidentally giving } H(A) = H(B)] = \frac{1}{\text{MOD}}$$

If we do K such comparisons in our algorithm, then on average (i.e. in expectation)  $K/\text{MOD}$  of them will give us false positive.

$$\begin{aligned} & H(x_1) \neq H(x_2) \quad x_1 \neq x_2 \\ & H(x_1) = H(x_2) \quad \text{m. f. n.} \quad K \cdot \frac{1}{\text{MOD}} \\ & H(x_1) = H(x_3) \quad O(K) \\ & \dots \\ & H(x_1) \quad H(x_2) \quad (A) O \\ & \quad \quad \quad (B) T \\ & \quad \quad \quad (C) N \end{aligned}$$

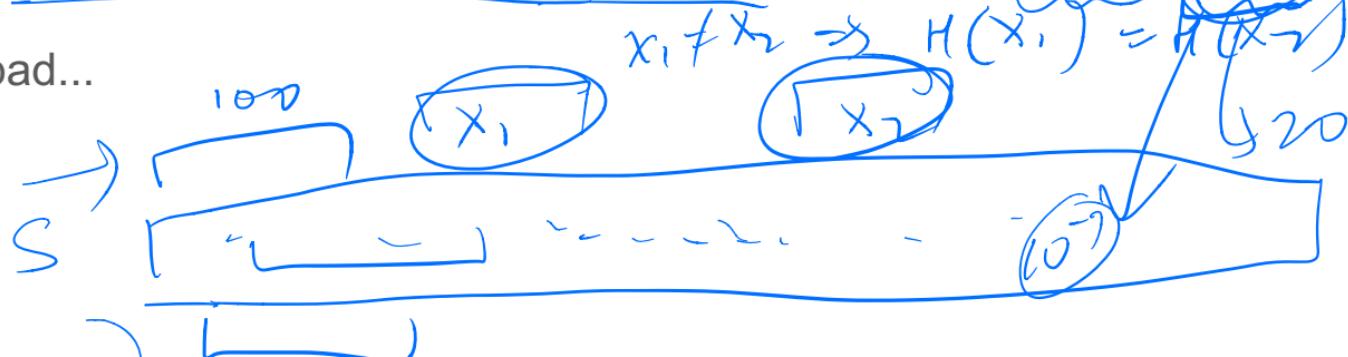
$$h(x_1) = h(x_k)$$

$$\begin{array}{c} 70 \\ \cancel{70} \\ h(x_1) \quad h(x_{71}) \\ 0 \quad 35 \end{array}$$

## Example where it may break

- Let's say your MOD =  $10^9 + 7$
- Maybe you hashed a string S of length  $10^6$  and put all the 100 length substrings of S in a set to find the number of distinct 100 length substrings
- How many comparisons did you do?
- Umm, roughly  $(10^6 - 100 + 1)^2 \sim 10^{12}$
- And each comparison had a fail probability of  $\sim 10^{-9}$
- So roughly how many failed comparisons did you do?  $10^{12} \cdot 10^{-9} = 1000 > 0$ .

This is bad...



$(10-100)$

$10^9$

## The Fix

$$H(1 \dots 10^9) = 3D$$

$$H(2 \dots 10^9) = 35 \dots$$

⋮

→ set →  $\#$

$\frac{1}{MOD^m}$

Use two (or more) mods.

Let  $H(S) = (H_1(S), H_2(S))$  using two different modulus

Ex.  $MOD_1 = 10^9 + 7$ ,  $MOD_2 = 10^9 + 9$

Now probability of accidentally matching for two strings is roughly  $10^{-18}$ .

So # of false positive  $\sim 10^{12 - 18} \sim 10^{-6}$  (super low)

$$H("APPLE") = \{ H_1(APPLE, 10^9 + 7), H_2(APPLE, 10^9 + 9) \}$$

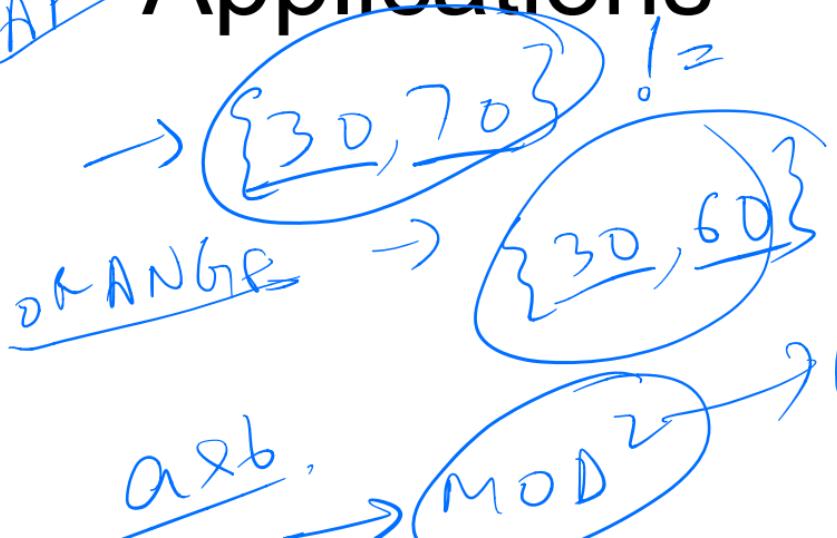
$$H_1(\text{APPLE}) \approx H_1(\text{ORANGE}) \Rightarrow \frac{1}{\text{MOD}}$$

$$\frac{1}{\text{MOD}_1}$$

$$\frac{1}{\text{MOD}_2} \sim$$

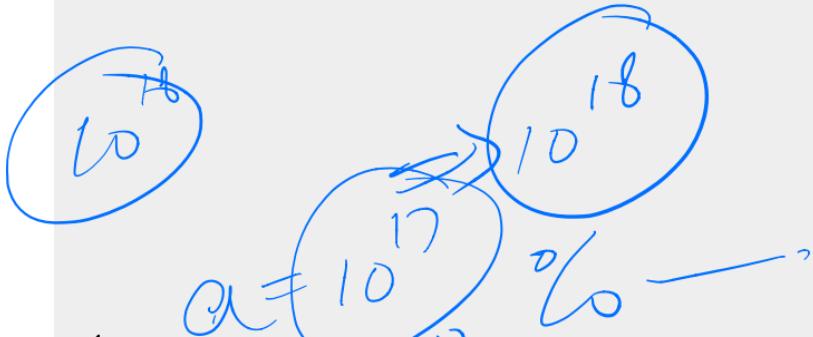
$$\frac{1}{\text{MOD}} \sim 10^{-10}$$

# APPLICATIONS

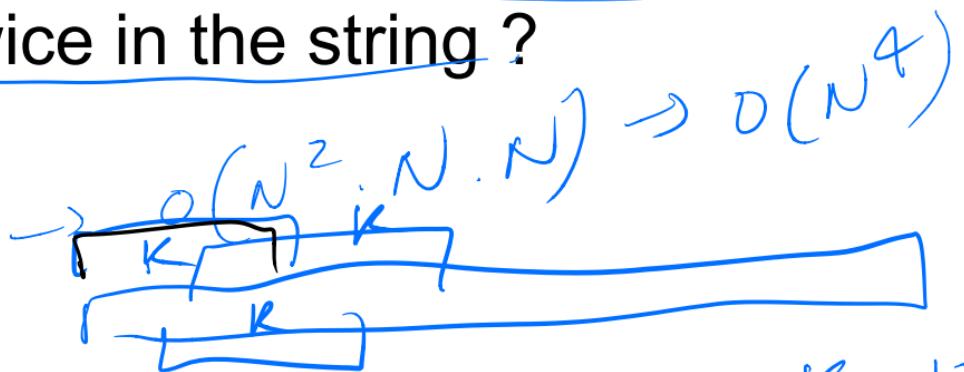
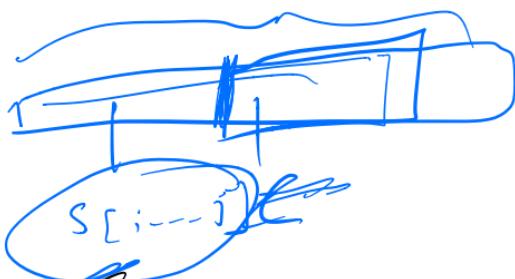


$$H(\rightarrow) \approx 30$$

$$\{q_{30}, 30\}$$



Given a string  $S$ , find the largest substring which occurs at-least twice in the string ?



Is there a substring of length  $K$ , that occurs  $\geq 2$  times.

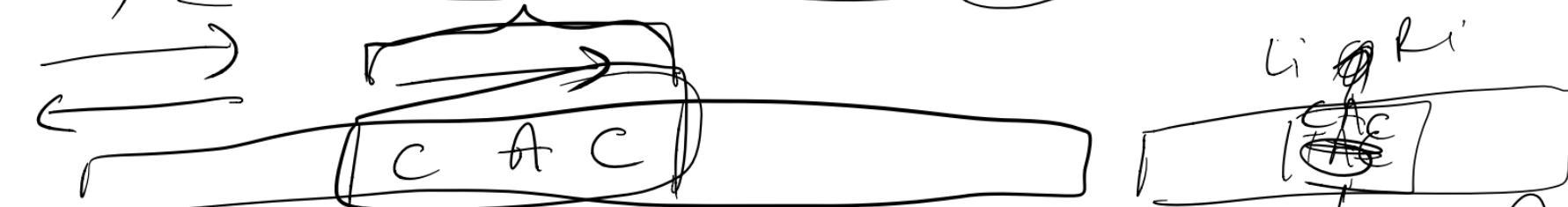
Yes/No.  $\rightarrow O(N \log N) \leftarrow O(N \log N)$

Binary searching.

Given a string S and Q queries of the form (L<sub>i</sub>, R<sub>i</sub>), tell if S[L<sub>i</sub>...R<sub>i</sub>] is a palindrome or not?

Segment Tree  $O(1)$  Change  $S[n] = y$

- What if we character updates also?



$T \rightarrow 2T$   $L_i \leftarrow$   $R_i$

$H(S) \Rightarrow A B B A$

$H(S^R) \Rightarrow$

$T \times 3^0 + A \times 3^0 + C \times 3^0$

$(i, i+)$   
 Given a string **S** find the length of the longest  
 palindromic substring in **S**

$\rightarrow$  Palindromic Tree  
 $\rightarrow$  Manacher

$\rightarrow$  Suffix Array

$\rightarrow$  Hashing

$i+1$

$R$

$\frac{n-1}{2}$

$\frac{k-1}{2}$

$S^R$

$S_{n-d}$

$\frac{n-1}{2}$

$S_i \quad S_j \quad S_1$

$O(N)$

$A$

$S_{i+1}^o$

$S_R$

$L = B$     $R =$   
 $\downarrow$     $\downarrow$   
 $\downarrow$     $\downarrow$   
 $i+1$     $S_R$

Largest common prefix of two strings **S** and **T**. If we are allowed to swap two characters in one of the strings ?

-- B v

S 1

=



# Pitfalls

# Pitfalls

*using ll (ints)*

- Slow implementation (unnecessary modulo operations)
- Off-by-one errors
- Always picking  $10^9 + 7$  and  $10^9 + 9$  as your mods and 26/30/36 as your base.  
Well then you can be hacked by the author.

100%.

Our randomisation argument wasn't completely correct (It was a simplified version and had some loopholes).

Ideally you want to randomly pick the base as well as the mod. See this [Codeforces Blog Post](#) for precise details.

3

So rule of thumb is: Submit the code once with BASE = 33/36 (some typical value) and MOD<sub>1</sub> = 10<sup>9</sup> + 7 and MOD<sub>2</sub> = 10<sup>9</sup> + 9. If it WAs then either pick some other primes from your codebook or during runtime randomise the picking of base and mod from a hard-coded list.

1 10 : 22 → 40

33

33

# Further Readings

- <https://codeforces.com/blog/entry/60445> (More in-depth)
- <https://www.quora.com/q/threadssiiithyderabad/String-Hashing-for-competitive-programming>  
(Quora Post)
- <https://cp-algorithms.com/string/string-hashing.html> (CP Algorithms)

⑩

- <https://codeforces.com/blog/entry/4898> (Anti-Hash tests)
- <https://codeforces.com/blog/entry/60442> (Extensive discussion on anti-hash testing)  
<http://rng-58.blogspot.com/2017/02/hashing-and-probability-of-collision.html>

(Hashing for unordered elements)

MOD  $\Rightarrow$   $[10^9 + 7, 10^9 + 9]$

20  
prime

codebook  $\Rightarrow$  10-15 big parts  
 $\hookrightarrow$  Rng 58. ~~Path HS~~  $\hookrightarrow \sim 10^9$   
integer

# Thank You

$a, b \sim 10^{11}$

$\times 10^{11} \nleq 10^{22}$

Q&A

mod

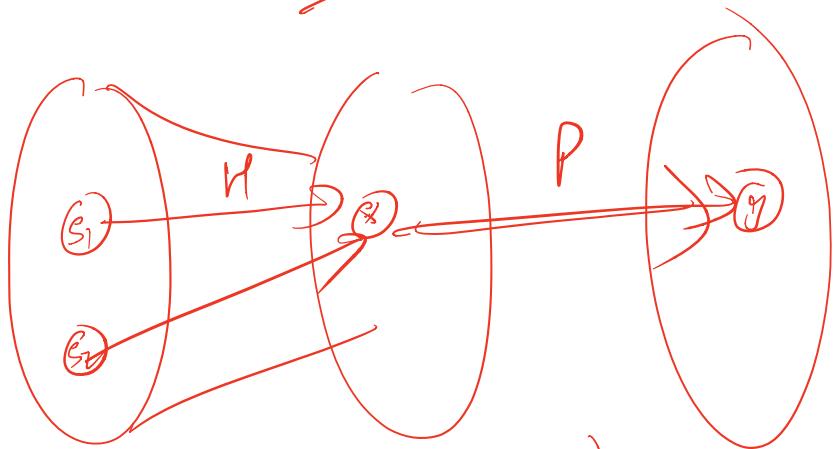
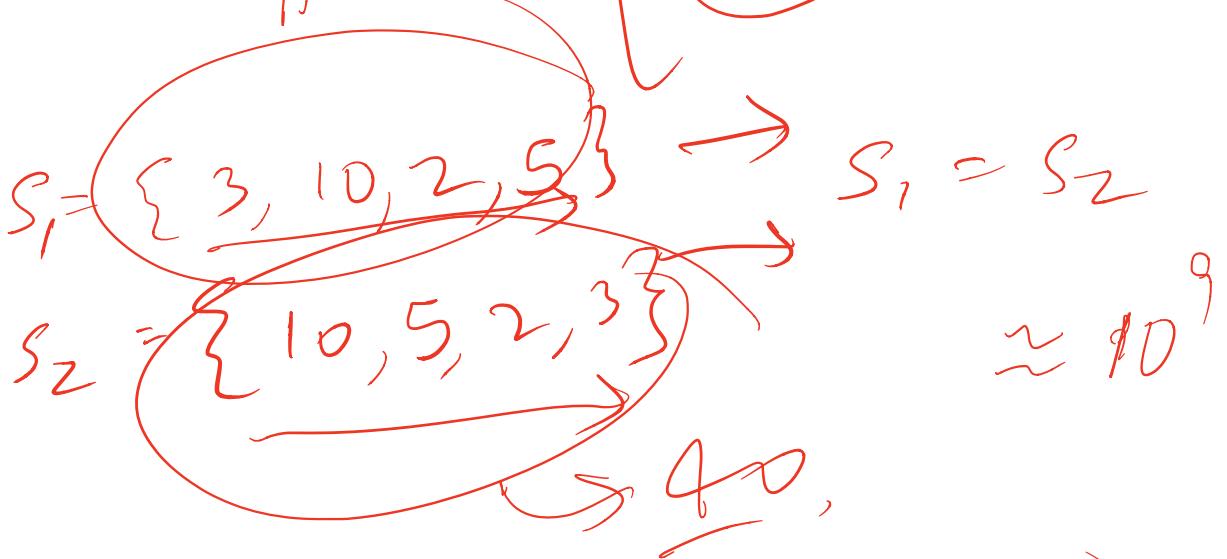
$a + b = \text{mod}$

$10^{12}$

40

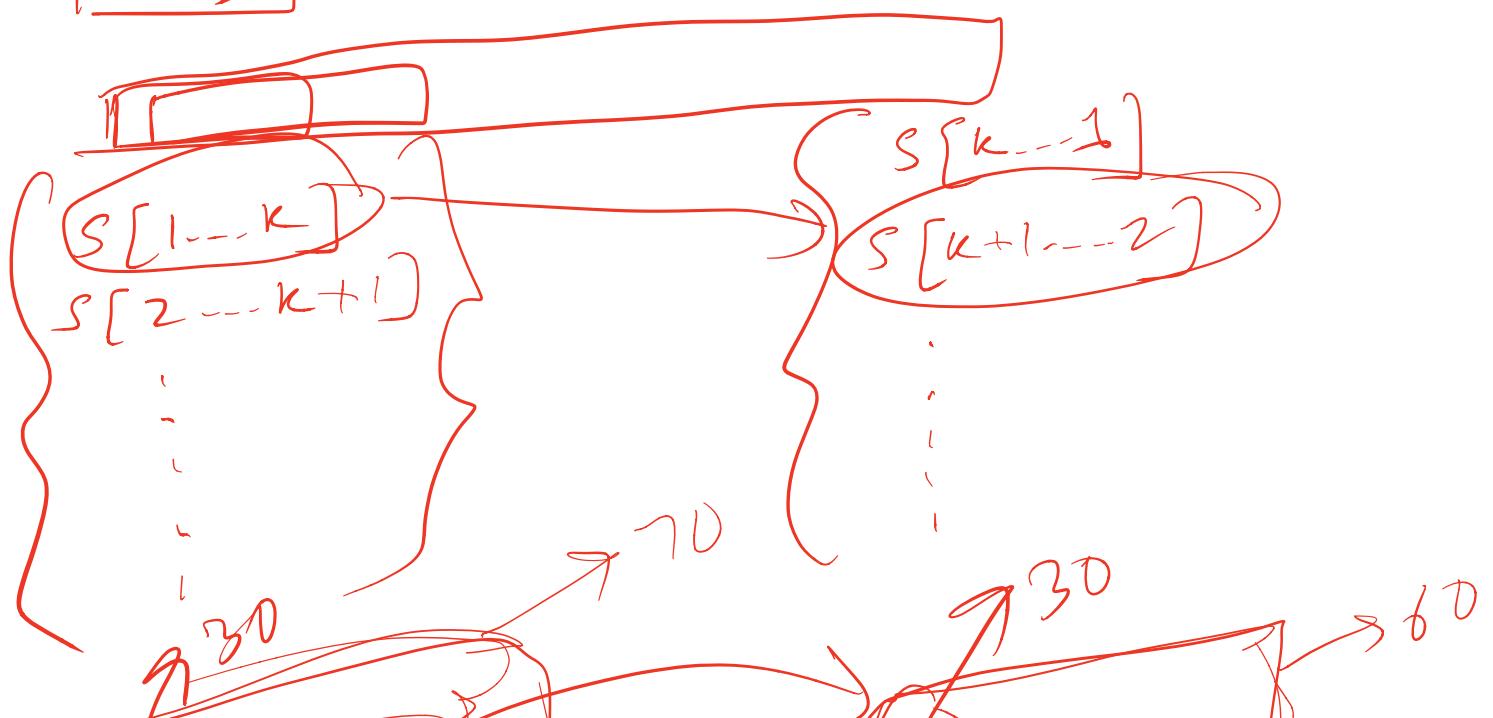
11

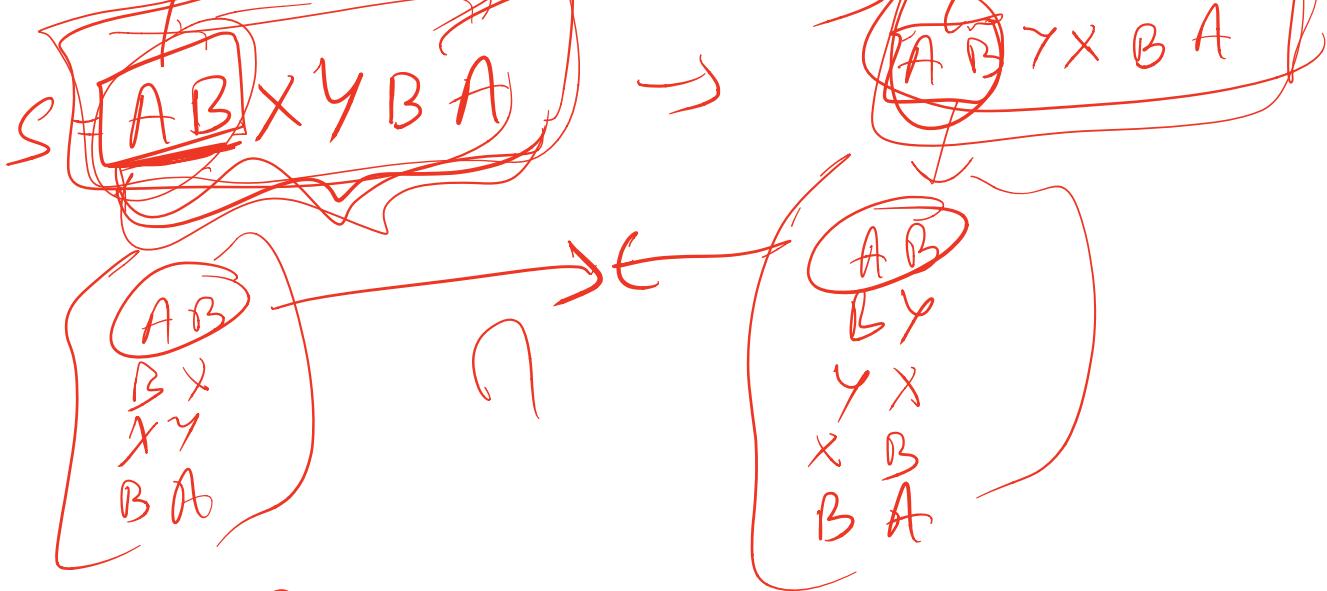
35



$$P(H(S_1)) = y = P(H(S_2))$$

→





$$R = 2 \cdot H$$

↓

$$S + n \rightarrow f(H(S), n)$$

$$H(S+n) \Rightarrow \boxed{?}$$