

Tries

Trees &

kmp / 2 algorithms

string matching

Binary numbers

String Hashing

100

string

problems

30

hashing

10

tries

68-69

- Priyansh Agarwal

1 or 2

100 kmp / 2 algorithm problems

↓ 99% of these

can also be solved using

Hashing

String matching

$O(1)$

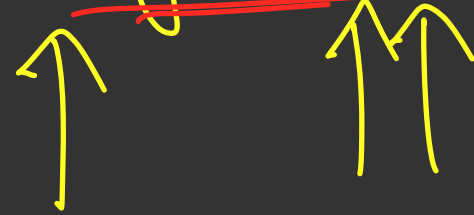
"priyansh"



$O(n)$

\times

"priyansh"



$O(m)$

$O(\min(n, m))$

"abcabcabcabc" $\rightarrow S \rightarrow \underline{O(n)}$

how many times does the string

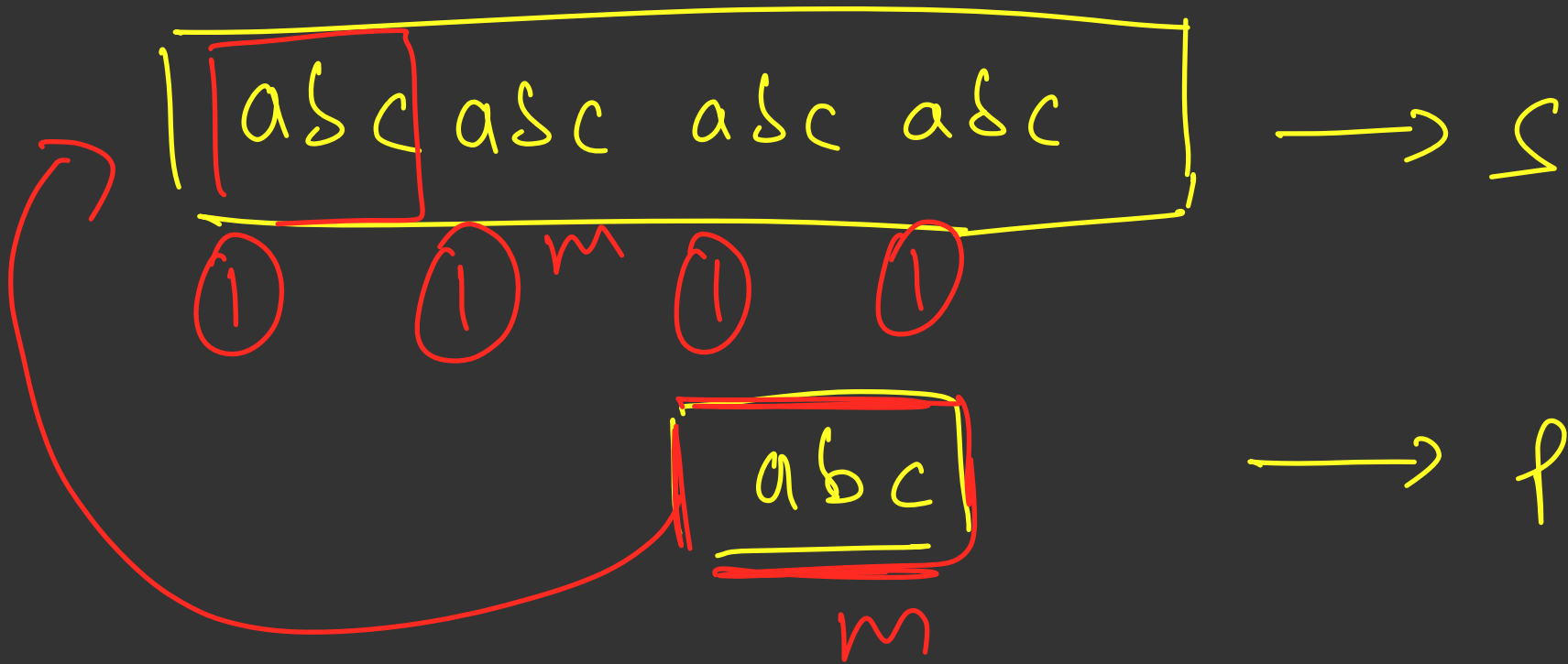
"abc" come as a

substring

\rightarrow

$\rightarrow \underline{O(m)}$

$O(n \times m)$



Substring $(i, i+m-1)$ in S
compare with P

$$O(n \times m)$$



① iterate over all indices $\rightarrow O(n)$

② compare the current substring with the pattern $\rightarrow O(m)$

$O(n \times m)$

$O(\log m)$

$O(1)$ time

① $kmp \geq$ ①

$O(n + m)$
time

worst case T.C

10^5

② Hashing

$O(n \log m)$
time



$O(n)$ time

100

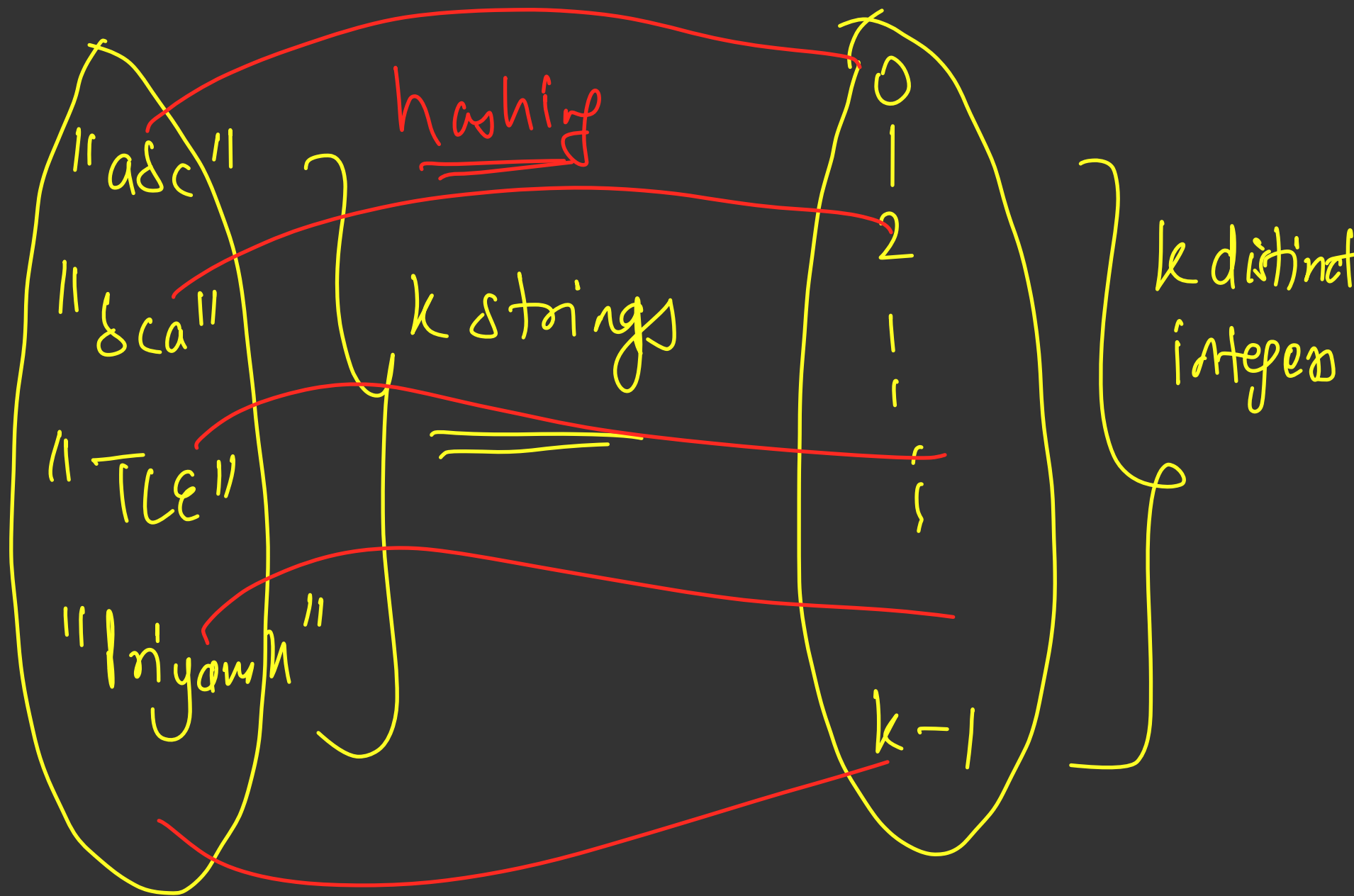
Why String Hashing?

- Optimizing Brute Force solutions
- Comparing 2 strings
- Daddy of String Algos

A _ _ _ _ _ n characters

B _ _ _ _ _ n $O(n) \propto$

whether $A == B$ or not



"abc" → 0

"bca" → 1

"aaa" → 2

"aac" → 3

"abc"

hashing function

0

a

b

c

⋮

z



~ size 100



10^{18} integers

log log

Given a string \rightarrow map it to an
integer

"niyauh _ _ _ _" \rightarrow 6

" _ _ _ _ " \rightarrow 7

Requirements

always

every possible string cannot be mapped to a unique integer

- If $a == b$, $\text{Hash}(a) == \text{Hash}(b)$

- If $a \neq b$, $\text{Hash}(a) \neq \text{Hash}(b)$

a and b are strings

- If this is true with a very high probability we are good to go
- Calculating Hash function should be fast enough
- Hash of a string shouldn't change in the code
- Hash value should be itself $O(1)$

$$||abc|| \rightarrow 1$$

$$||\delta aa|| \rightarrow 2$$

$$||caa|| \rightarrow 2$$

$$||a\delta a|| \rightarrow 3$$

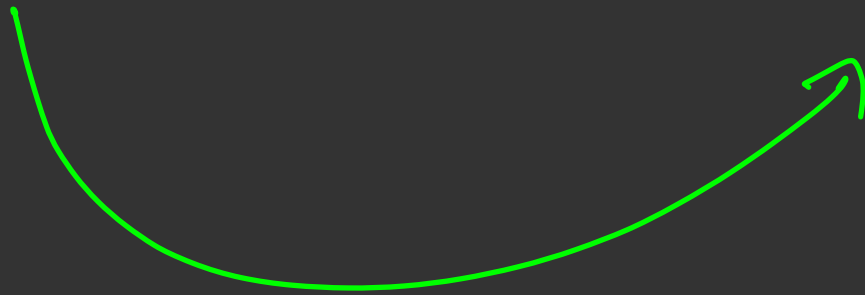
$$\left\{ \begin{array}{l} 1 \\ 2 \\ 3 \end{array} \right.$$

n

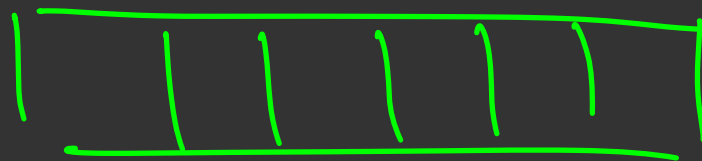
Strings

m

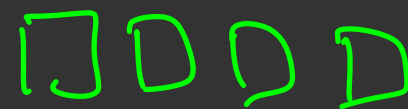
integers



$$n > m$$



m boxes



n balls

$n > m$

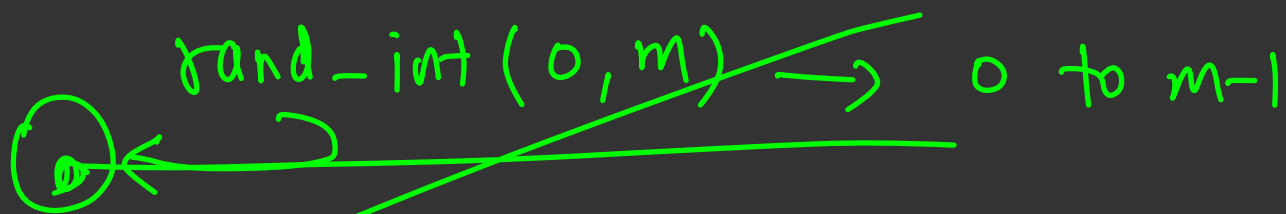
$n > m$

if we want to put all balls in m
boxes then atleast 1 box must have
more than 1 ball

n strings

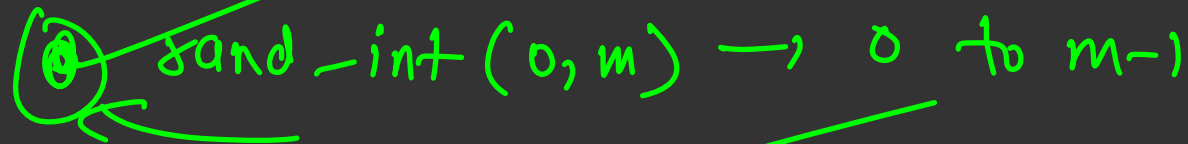
m integers

~~$\text{rand_int}(0, m) \rightarrow 0 \text{ to } m-1$~~



A horizontal line connects the first string (0) to the first integer (0). This line and the text above it are crossed out with a diagonal line.

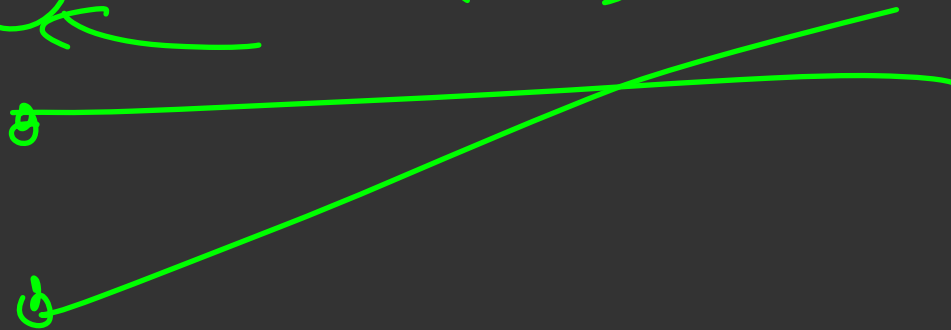
~~$\text{rand_int}(0, m) \rightarrow 0 \text{ to } m-1$~~



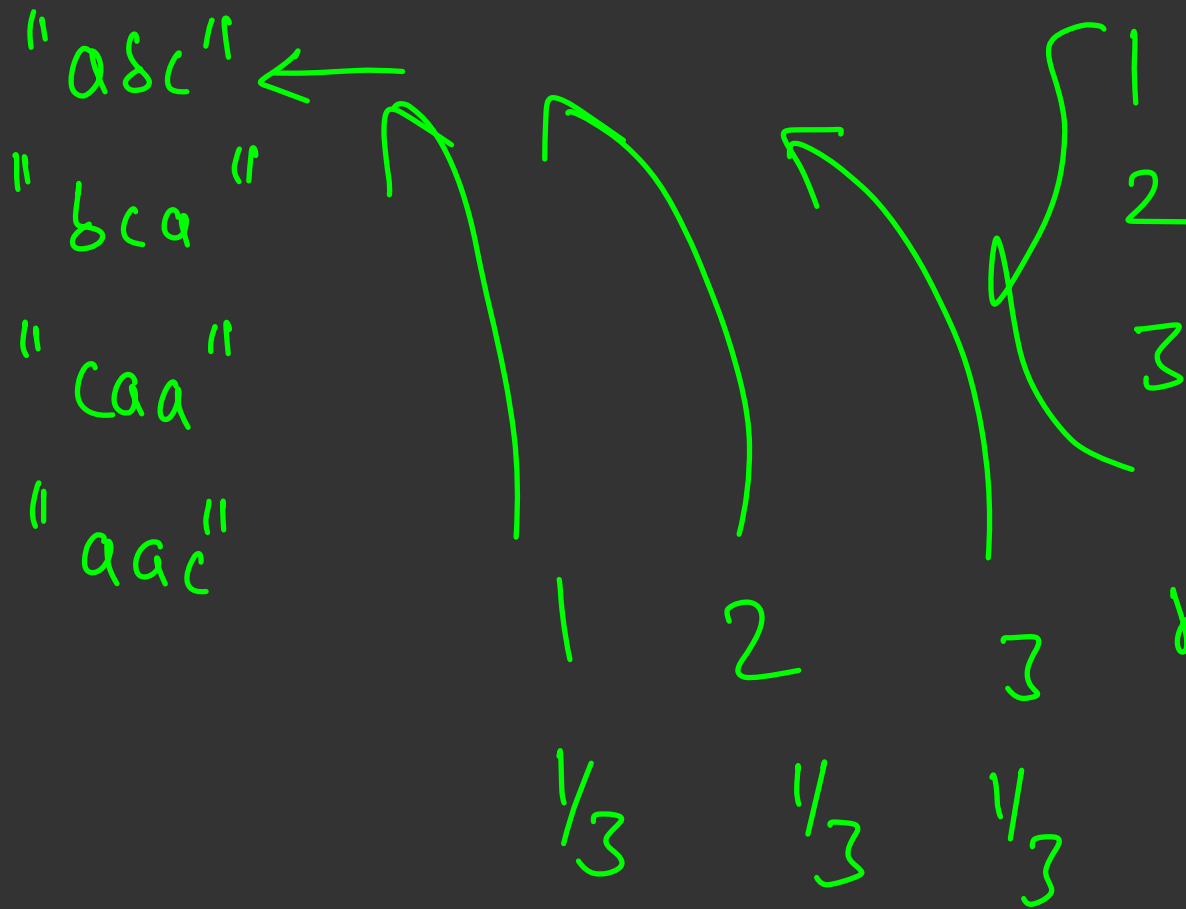
A horizontal line connects the second string (1) to the first integer (0). This line and the text above it are crossed out with a diagonal line.

3

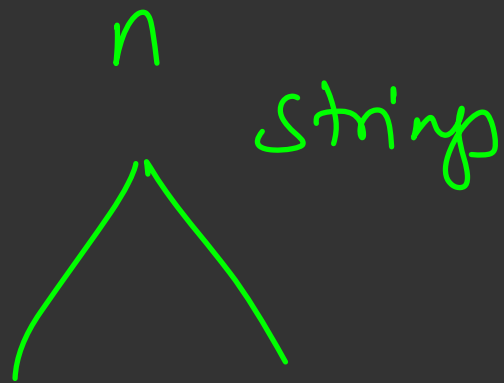
4



A horizontal line connects the third string (3) to the second integer (1). This line is not crossed out.



rand - int(1,3)

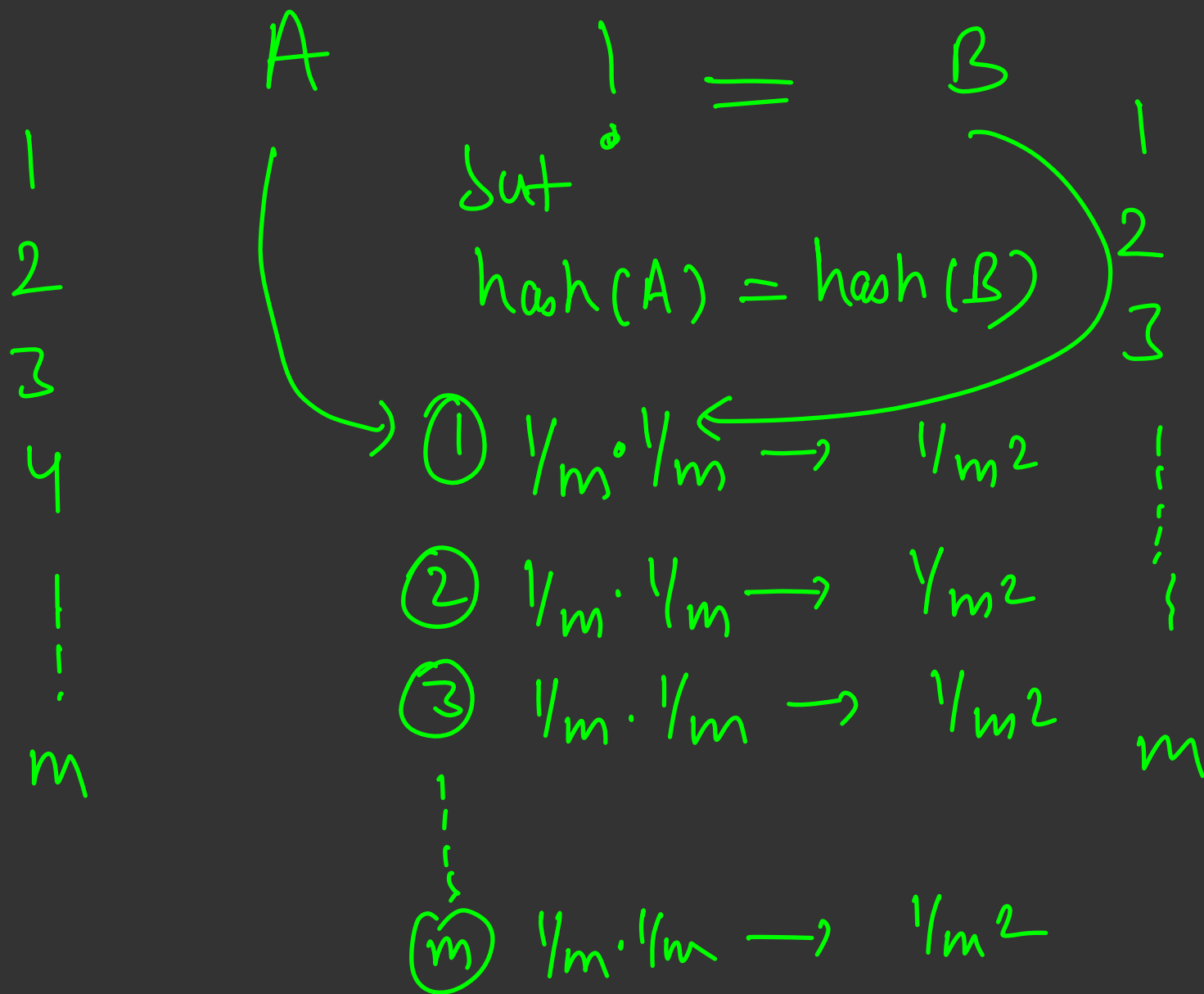


m integer

$[a] \neq [b]$

what is the probability that

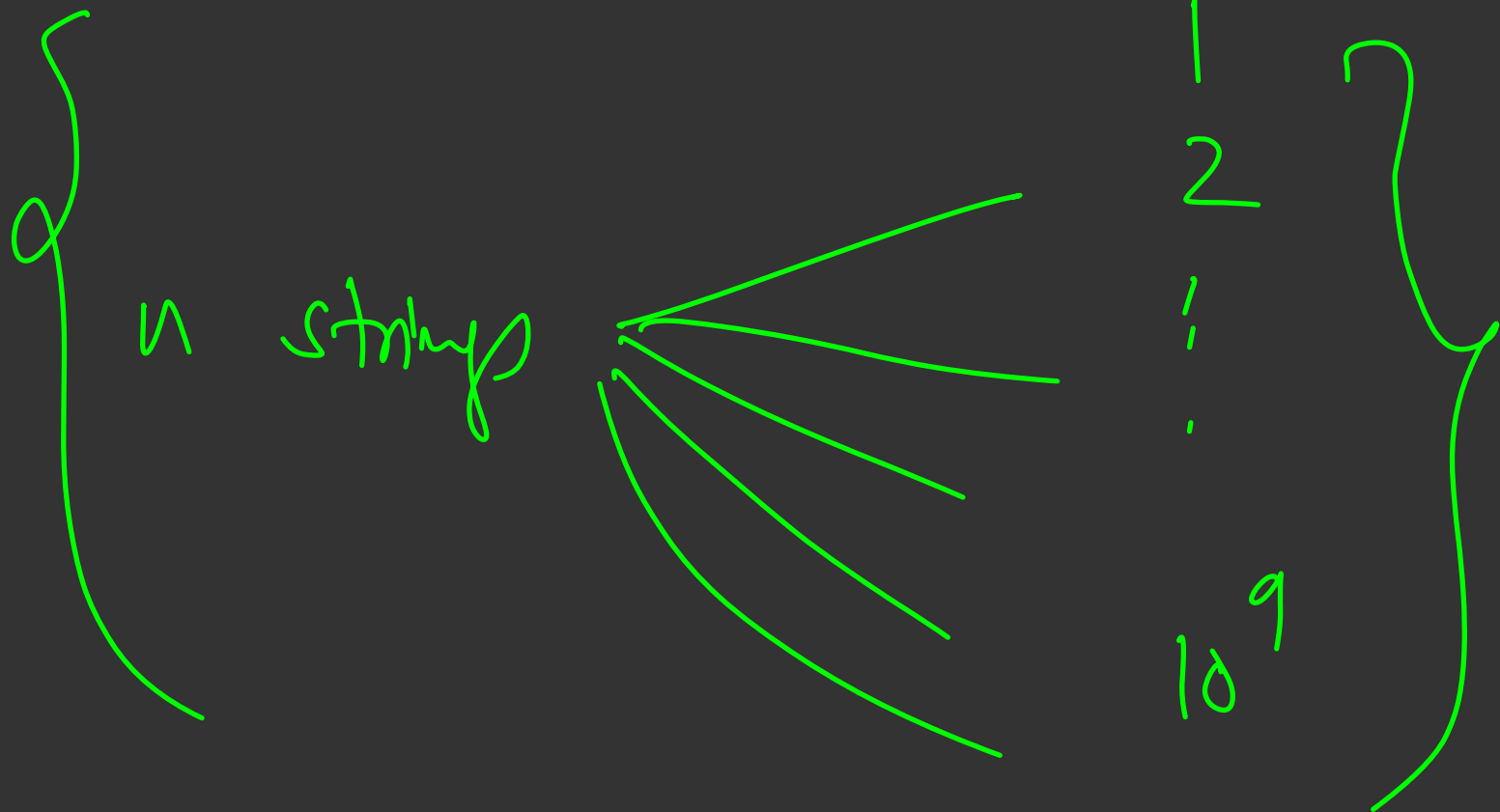
$$\text{hash}(a) == \text{hash}(b)$$



$$\left(\frac{1}{m^2} + \frac{1}{m^2} + \frac{1}{m^2} + \dots \right) m \text{ times} \rightarrow \underline{\underline{\frac{1}{m}}}$$

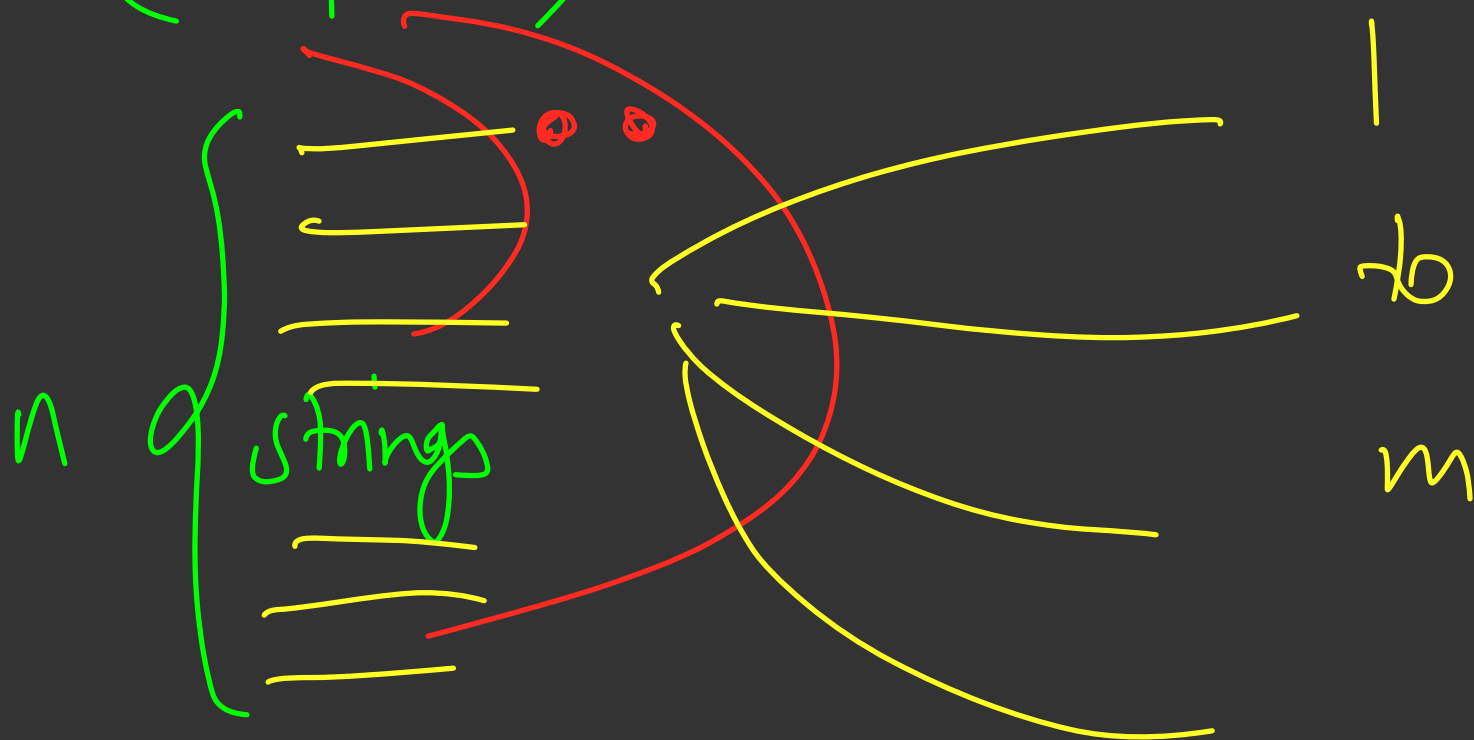
$$\underline{m = 10^9}$$

$$\boxed{1/10^9}$$



(A, B)

if $(A \neq B)$ but $h(A) = h(B)$



$(A == B)$

$hash(A) == hash(B)$

"asc" →

"sac" → 2

"aac" → 3

"sao" → 1

1
2
3

is "asc" == "sac"
↑ 1 ↑ 1

What is the probability that this fails

if $(A)_0 \geq B$ but $\text{hash}(A) == \text{hash}(B)$

$\text{hash}(A) \rightarrow$	1	$1/m$	$\text{hash}(B) \rightarrow$	1	$\rightarrow 1/m$
	2	$1/m$		2	$1/m$
	3			3	
	\vdots				
	m				

if my algorithm has a probability
of failing as $1/m$

how many times can I run my
algorithm on average before it fails

m comparisons



Probability of event happening = p

then

average no. of tries required for event to happen is $\frac{1}{p}$

A and B of my

we are comparing A and B



average

$1/m$

10^8

$$m \rightarrow \underline{10^9}$$

$$1/10^9$$

$$\underline{10^8}$$

$$\boxed{\underline{10^9}}$$

Codeforces problem \rightarrow 100 test cases

every test case require 10^8 comparisons

10^{10} comparisons

10^9

$S \longrightarrow$ randomly generating an integer
from 1 to m

Polynomial Rolling Hash

$$\begin{aligned}\text{hash}(s) &= s[0] + s[1] \cdot p + s[2] \cdot p^2 + \dots + s[n-1] \cdot p^{n-1} \pmod{m} \\ &= \sum_{i=0}^{n-1} s[i] \cdot p^i \pmod{m},\end{aligned}$$

(p, m) if $(A == B)$ then
 $\text{hash}(A) = \text{hash}(B)$

$$\text{hash}(s) = (s_0 \cdot p^0 + s_1 \cdot p^1 + s_2 \cdot p^2 + \dots + s_{n-1} \cdot p^{n-1}) \% m$$

p and m \rightarrow range of integers allowed

(0 to $m-1$)

① if $A = B$ then

$\text{hash}(A) = \text{hash}(B)$ with
probability 1

② if $A \neq B$ then $\text{hash}(A) \neq$
 $\text{hash}(B)$ with a very high probability

③ $\text{hash}(\text{any string}) \rightarrow \underline{0 \text{ to } m-1}$
 m distinct intes

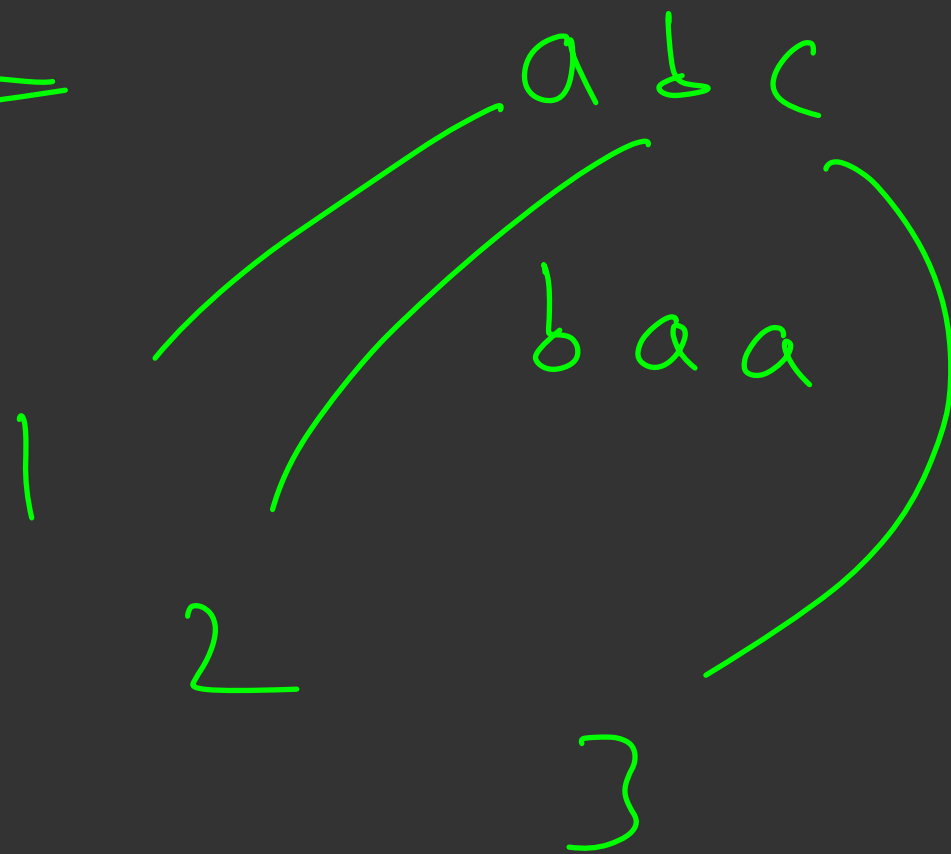
Assume that mod m did not
exist

'asc'
—
'baa'
—

292
300
————

(a (—————)
b (—————)

27



d → 4

e → 5

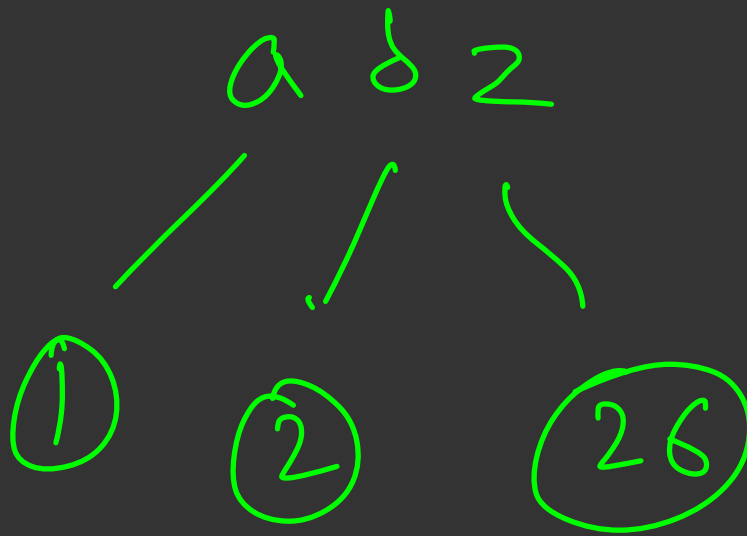
a \rightarrow 1

b \rightarrow 2

c \rightarrow 3

⋮

z \rightarrow 26



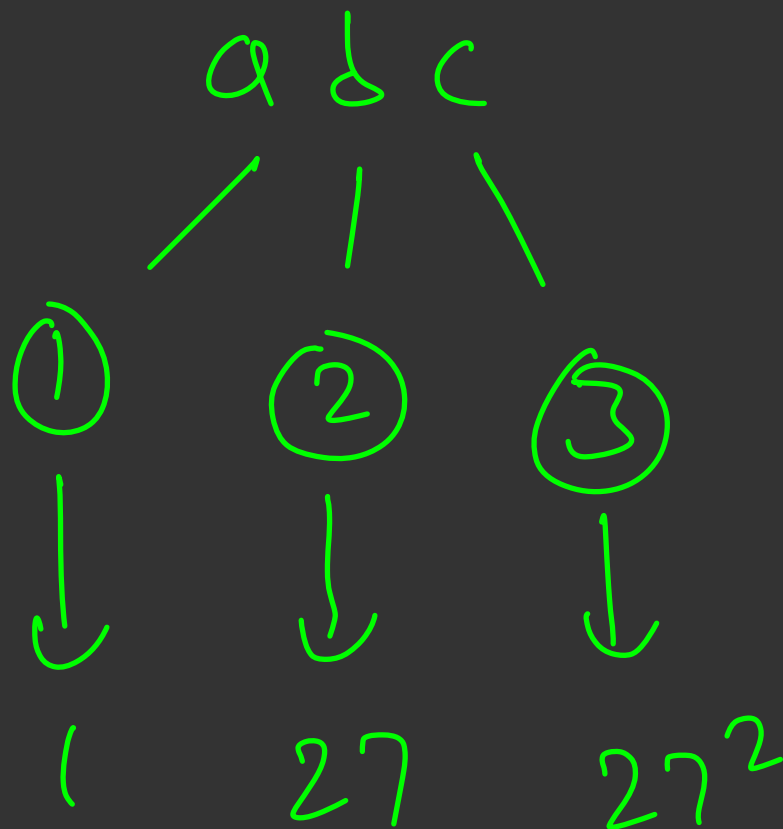
$$\boxed{① + ② \cdot 27 + (26) 27^2}$$

A

B

integer

$$\begin{array}{lcl}
 292 & \longrightarrow & \boxed{\frac{2 \cdot 10^2}{3 \cdot 10^2} + 9 \cdot 10^1 + 2 \cdot 1} \\
 300 & \longrightarrow & \boxed{3 \cdot 10^2 + 0 \cdot 10^1 + 0 \cdot 1}
 \end{array}$$



1 0 1 0 1
0 0 1 0 1

base 2

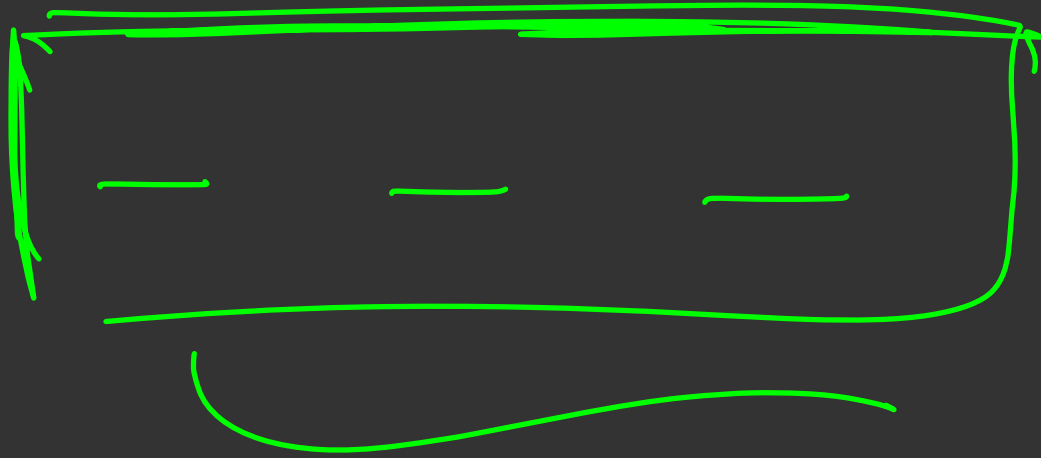
$$1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0$$

$$0 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0$$

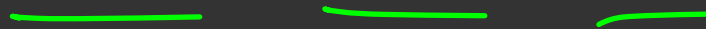
3 9 2 4
2 9 2 4

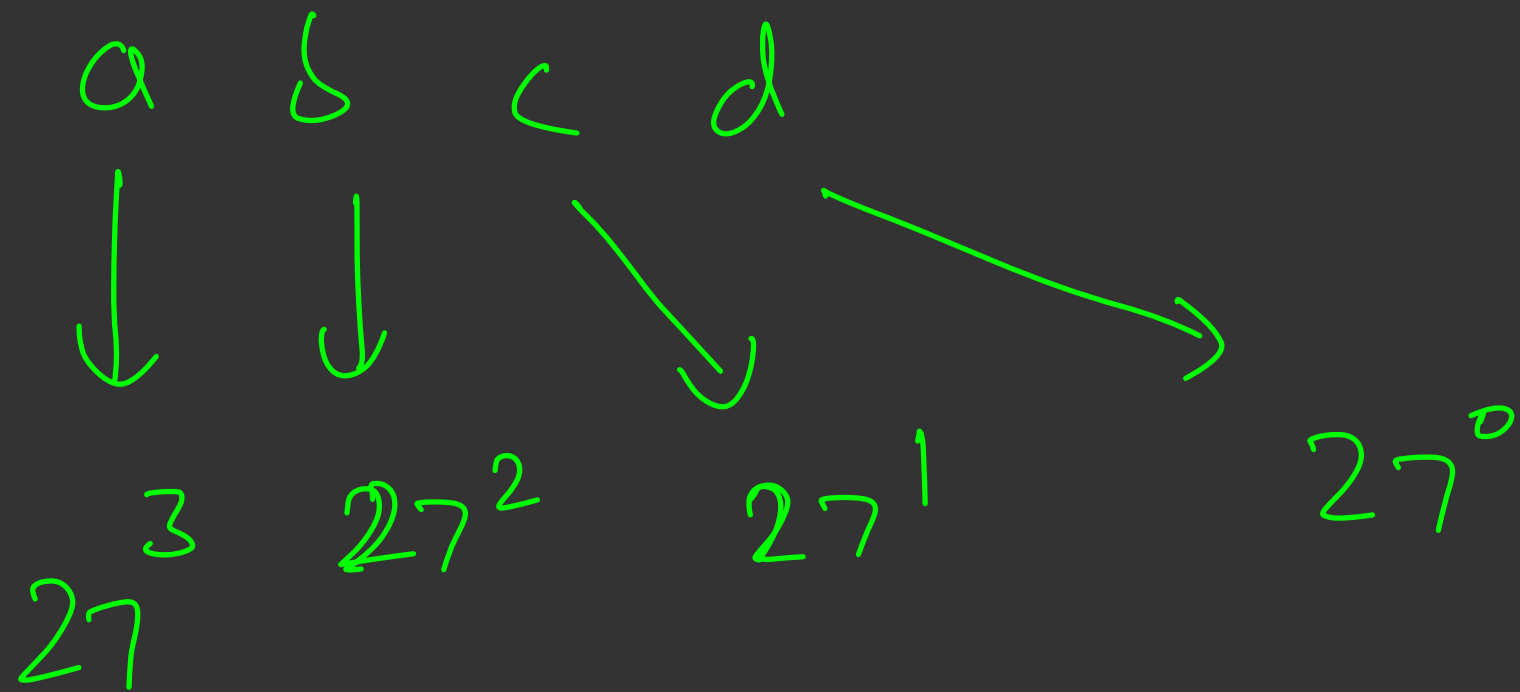
$$3 \cdot 10^3 + 9 \cdot 10^2 + 2 \cdot 10^1 + 4 \cdot 10^0$$

$$3 \cdot 10^3$$



3000





26^3

— — —

$$\text{hash}(s) = [s_0 \cdot p^0 + s_1 \cdot p^1 + s_2 \cdot p^2 + \dots + s_{n-1} \cdot p^{n-1}] \% m$$

$p > \text{no. of unique characters}$

$$(p, m)$$

string can contain character a, b, c

3 unique character

if $p = 10^5$ and $m \leq 10^9$

$$\text{hash}(s) = [s_0 \cdot p^0 + s_1 \cdot p^1 + s_2 \cdot p^2 + \dots + s_{n-1} \cdot p^{n-1}] \% m$$

p should not divide m
 m should not divide p } $p > \text{unique character}$
 $m \leq 10^9$

fermat's theorem

$$\rightarrow a^{m-1} = 1 \pmod{m}$$

for $m = \text{prime}$

$$\rightarrow a^{m-2} = a^{-1} \pmod{m}$$

Strings only contain lower case letters

$$l = 27, \quad m = 10^9 + 7$$

Strings with uppercase & lowercase

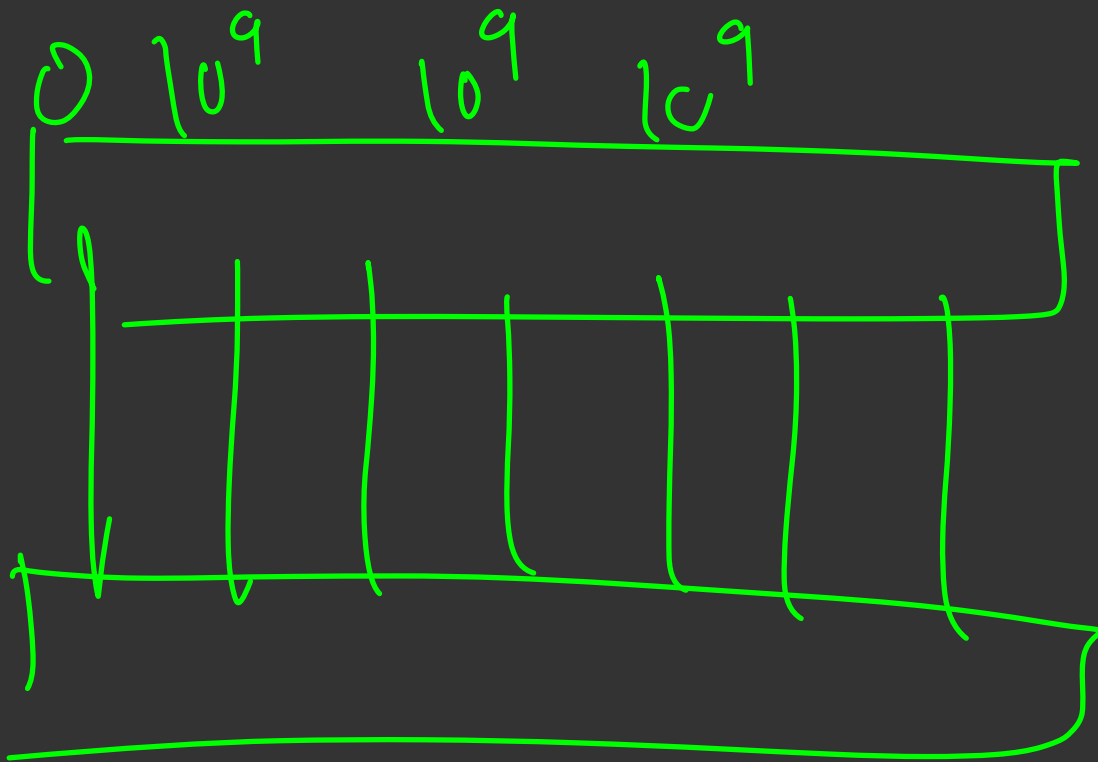
$$l = 60, \quad m = 10^9 + 7$$

$f = 1001$

$$1 \quad w = 10^9 + \gamma$$

$$p = 10^9 + 5$$

$$m = (0^9 + 7)$$



Hashing a string

$$\text{hash}(s) = (s_0 \cdot p^0 + s_1 \cdot p^1 + \dots + s_{n-1} \cdot p^{n-1}) \% m$$

$p >$ unique character

0 to 10^9

$m \leq 10^9$ and prime

$$\text{hash(arr)} = (a[0] \cdot p^0 + a[1] \cdot p^1 + \dots + a[n-1] \cdot p^{n-1}) \% m$$

$p >$ unique element $p \rightarrow 10^9 + 3$

Hash of a string $S = H$
 $\equiv \textcircled{n}$

find hash of string $S + ch$

$$(H + ch \cdot p^n) \% m$$

$$\text{hash}(s) = (s_0 \cdot p^0 + s_1 \cdot p^1 + \dots + s_{n-1} \cdot p^{n-1}) \% m$$

$$\text{hash}(s + ch)$$

$$\rightarrow (H + ch \cdot p^n) \% m$$

hash of a string S is H

what is the hash of

$(ch + S)$

$$\text{hash}(s) = \left(\underbrace{s_0 \cdot p^0 + s_1 \cdot p^1 + \dots + s_{n-1} \cdot p^{n-1}}_{H} \right) \% m$$

$$\text{hash}(ch + s)$$

$$= \left(ch \cdot p^0 + \underbrace{s_0 \cdot p^1 + s_1 \cdot p^2 + \dots + s_{n-1} \cdot p^n}_{H \cdot p} \right) \% m$$

$$= \left(ch + \underline{\underline{H \cdot p}} \right) \% m$$

Why Rolling?

- ✓ Hash(s) = H, what is Hash(s + x), x = character

- $(H + xp^n) \% m$

→ log n time

- ✓ Hash(s) = H, what is Hash(x + s), x = character

- $(x + Hp) \% m$

→ O(1) time

- How about calculating the Hash of a substring quickly?

- How about just comparing two strings?

hash of a substring:

Given a string calculate its hash
of size n

→ Time = $O(n)$ time

$O(1)$ time to add

character at back

$O(1)$ time

$$\text{hash}(s) = (s_0 \cdot p^0 + s_1 \cdot p^1 + \dots + s_{n-1} \cdot p^{n-1}) \% m$$

hash of substring from l to r

$$(s_l \cdot p^0 + s_{l+1} \cdot p^1 + \dots + s_r \cdot p^{r-l})$$

hash(s)

$$= \left[s_0 \cdot p^0 + s_1 \cdot p^1 + \dots + s_{l-1} \cdot p^{l-1} + \underbrace{s_l \cdot p^l + \dots + s_{\sigma} \cdot p^{\sigma} + \dots + s_{n-1} \cdot p^{n-1}}_{p^l} \right] \% m$$

$$\left[s_l \cdot p^0 + s_{l+1} \cdot p^1 + \dots + s_{\sigma} \cdot p^{\sigma-l} \right] \% m$$

hash[i]

$$\rightarrow s_0 \cdot p^0 + s_1 \cdot p^1 + \dots + s_i \cdot p^i$$

$$\text{hash}[x] - \text{hash}[l-1] = \left[\underbrace{s_l \cdot p^l + s_{l+1} \cdot p^{l+1} + \dots + s_x \cdot p^x}_{p^l} \right]$$

$$= (s_0 \cdot p^0 + s_1 \cdot p^1 + \dots + s_{l-1} \cdot p^{l-1})$$

hash(substring from l to r)

$$= \left[\frac{\text{hash}(r) - \text{hash}(l-1)}{p^d} \right] \% m$$

$$= \left(\text{hash}(r) - \text{hash}(l-1) \right) \cdot \underline{\text{modinv}(p^d)} \% m$$

$$\underline{O(\log m)}$$

$$a \rightarrow \left(\frac{1}{a} \right) \% m \quad \hookrightarrow \text{prime}$$

$$(a^{m-2})$$

$$(p^d) \rightarrow \left(\frac{1}{p^d} \right) \rightarrow (p^d)^{m-2}$$

"abcabcabc" $\rightarrow O(n)$

find out how many abc occurs $O(m)$

$O(n) + O(m)$

Pitfalls?

- $A \neq B$ but $\text{Hash}(A) == \text{Hash}(b)$
 - Probability = $1 / m$
- Comparing 50 such strings, probability of a collision = $50 / m$
- Let's look at an example problem to see how it fails
 - Substring comparison problem
- Solution????

$$\underline{10^{10}}$$

$$\underline{10^9}$$

What if you're doing 10^7

comparisons for $n = 10^9$

$$\frac{1}{10^9}$$

→

$$10^7$$

$$\frac{10^7}{10^9}$$

$$\underline{\underline{1/100}}$$

$$\text{hash1} \rightarrow p = 30, \quad m = 10^9 + 7$$

$$\text{hash2} \rightarrow p = 30, \quad m = 10^9 + 9$$

$$\text{Single Integer} \rightarrow \frac{1}{m} \rightarrow \underline{\underline{\frac{1}{m^2}}}$$

$$\begin{array}{ccc}
 10^{12} & 10^{12} & m \\
 \uparrow & \uparrow & \\
 (a \times b) \% m & \xrightarrow{\quad} & \text{overflow}
 \end{array}$$

$$(a + b) \% m$$

Problems

- ✓ Number of different substrings
- Palindrome queries
- ✓ Largest string which repeats twice
- ✓ Longest palindromic substring

$O(n^2)$

5 level

3 low

1 —

5
3
1

6
4
2

aaabab

5
4
3
2
1

1

Binary search on length of string

$O(\log n)$

k

for(int i=0; i<=n-k; i++)

Binary search on length $O(\log n)$

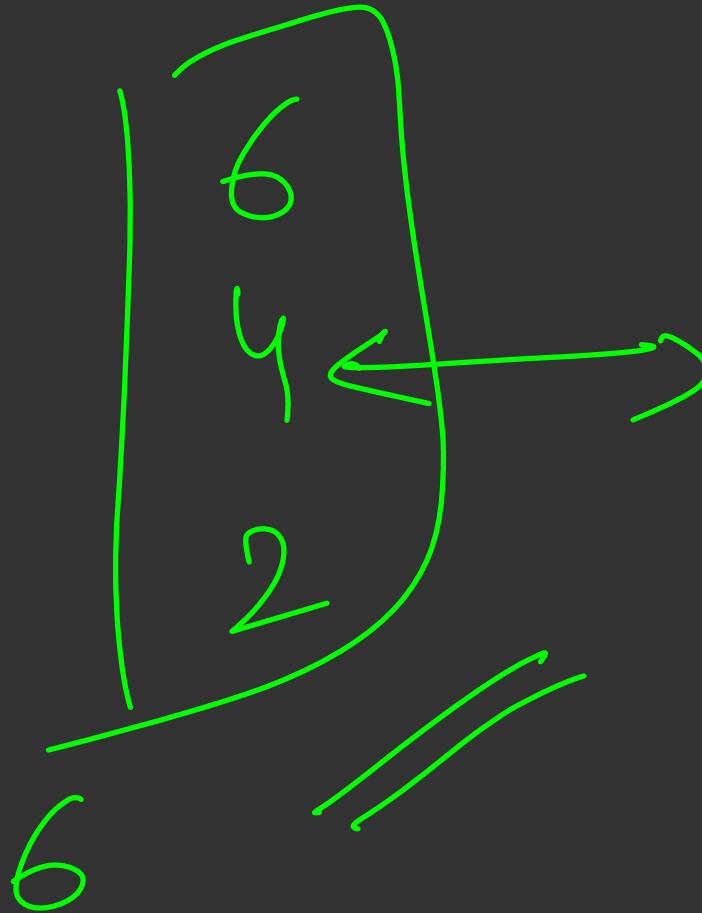
→ $map < int, int >$ occure

```
for (int i = 0 ; i <= n - k ; i++)
```

```
    int hash = substringHash(i, i+k-1)
```

```
    if (occure.find(hash) != occure.end())
```


Longest palindromic substring



```
for (int i = 0; i ≤ n - k; i++)
```

```
    int first_half =
```


a b c c b a

original string \longleftrightarrow
reverse string \longleftrightarrow

→ x a b c c b a y z

z y a b c a b c x

S → xypr abc cbc mn

Srev → nm abc cba rqp y