

# String Algorithms

HKOI 2017 training

Anson Ho

# Content

- introduction
- trie
- KMP
- hashing
- suffix array
- sth else..

# What is String?

a string is traditionally a sequence of characters (Wiki)

or more familiarly

- `char s[SIZE];`
- `“Hello World!”`
- `{ ‘a’, ‘b’, ‘c’ }`

# ASCII code

- American Standard Code for Information Interchange
- 8 bits = 1 byte = sizeof(char)
- 256 different states (0~255)
- '0' = 48      '9' = 48+9 = 57
- 'A' = 65      'Z' = 65+26-1 = 90
- 'a' = 97      'z' = 97+26-1 = 122
- ' ' = 32      '+' = 43      '-' = 45
- null = 0      '\n' = 10      '\r' = 13

# Glossary

concatenation

- "addition" in string
- $1+10 = 11$
- "1"+"10" = "110"
- C++: "1""10"
- Haskell: "1"++"10"

# Glossary

lexicographical order

- comparison in string
- dictionary order
- alphabetical order
- numerical order of ASCII code
- “1” < “2” < “A” < “B” < “a” < “b”

# Glossary

lexicographical order

- same lengths -> numerical order
- “123” < “132”
- “ab” < “ac”
  
- different lengths -> add zeros
- “a” < “ab”      {97,0} < {97,98}
- “ab” < “z”      {97,98} < {122,0}

# Glossary

prefix

- bicycle
- triangle
- hyperlink

suffix

- girls\_
- beautifull



# Glossary

substring

- prefix of suffix, suffix of prefix
- contiguous subarray

subsequence

- obtained by deleting characters
- order is kept

# Glossary

“Alice and Bob”

substring

- “”, “ce a”, “Alice and Bob”

subsequence

- all of the above
- “A B”, “AaBb”

# Glossary

palindrome

- same as the reversed string
- “a”, “”
- “radar”, “abccba”
- longest palindromic subsequence
- $O(N^2)$  with DP

# I/O

- Input and Output
- What will you use to read a string from `stdin`?
- (in C++)

# I/O

scanf printf

- read until space, line break or EOF
- ‘\0’ is automatically added
  
- `int n;`
- `char x[5]; // max 4 char.`
- `scanf(“%d%s”, &n, x);`
- `printf(“%s”, x);`

# I/O

gets puts

- read until line break or EOF
  - ‘\0’ is automatically added
  - input will not include ‘\n’
  - output with ‘\n’ at the end
- 
- `char x[5];`
  - `gets(x);`
  - `puts(x);`

# I/O

`cin cout`

- SLOW
- `sync_with_stdio(false)`

`getchar putchar`

- single character
- = `scanf printf` with `%c`
- unlocked

# I/O

sscanf sprintf

- read from a string
- write to a string
  
- `int n, m;`
- `char x[12];`
- `sprintf(x, "%d%d", n, m);`
- `sscanf(x, "%d", &n);`



# Container

`char[]`

- static size
- easy for I/O
- null terminator (`'\0'`)
  - reserve one space for it

`std::vector<char>`

- dynamic size
- compatible with other STL tools
  - `sort vector< vector<char> >`
  - `map, set`

# Container

`<string>`

- C++ object
- supports `+`, `==`, `<`, `<=`, ...
- `char x[5];`
- `string str(x);`
- `str[0]=str[1];`
- `printf(“%s\n”,str.c_str());`

# <cstring>

- memcpy          strcpy
  - memcmp         strcmp
  - memset
  - strlen
- 
- don't use if you forget how to use
  - same time complexity as brute force
  - maybe faster in terms of constant

# Trie

- not a formal word
- (not found in Oxford Dictionary)
- a tree
- storing multiple string
- searching or counting strings or prefixes

# Trie

```
struct node{  
    int count,child[26];  
    node *child[26];  
};
```

- each node represents a character
- stores the positions of its children which are the next character
- stores the count of strings ending with the node or in its subtree

# Trie

$N$ : total length

$M$ : length of the target string

- time complexity:

- insertion:  $O(M)$

- deletion:  $O(M)$

- searching/ counting:  $O(M)$

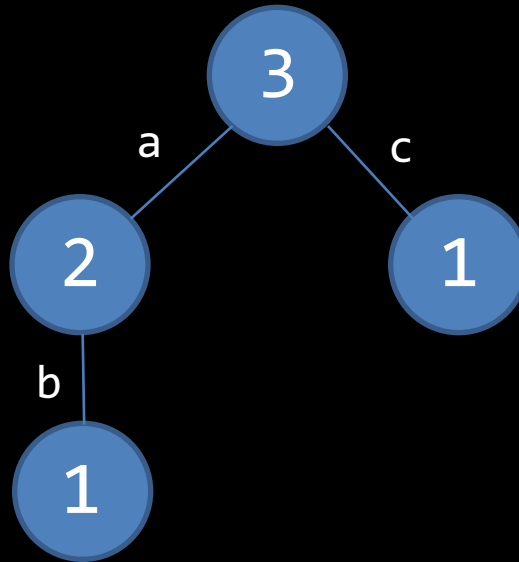
- sorting:  $O(26N)$

- space complexity:  $O(26N)$

actually all of them are DFS

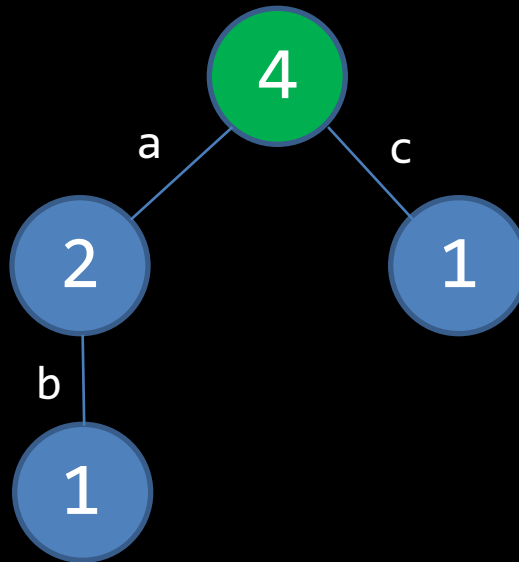
# Trie

insert "abc" to {"a", "ab", "c"}



# Trie

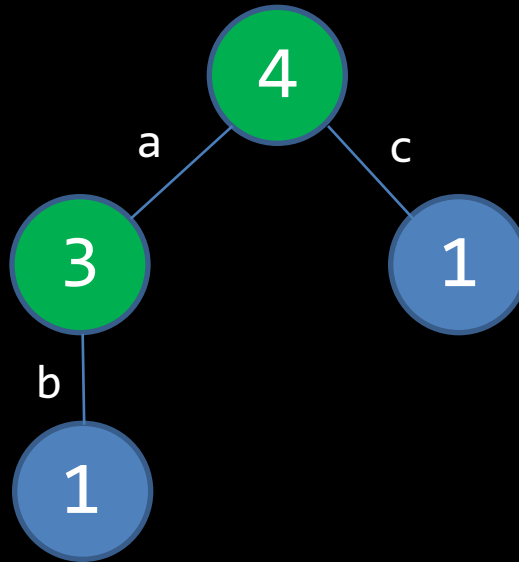
insert "abc" to {"a", "ab", "c"}





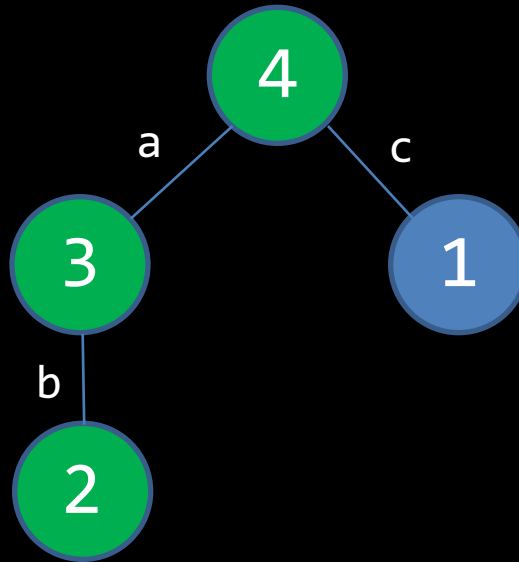
# Trie

insert "abc" to {"a", "ab", "c"}



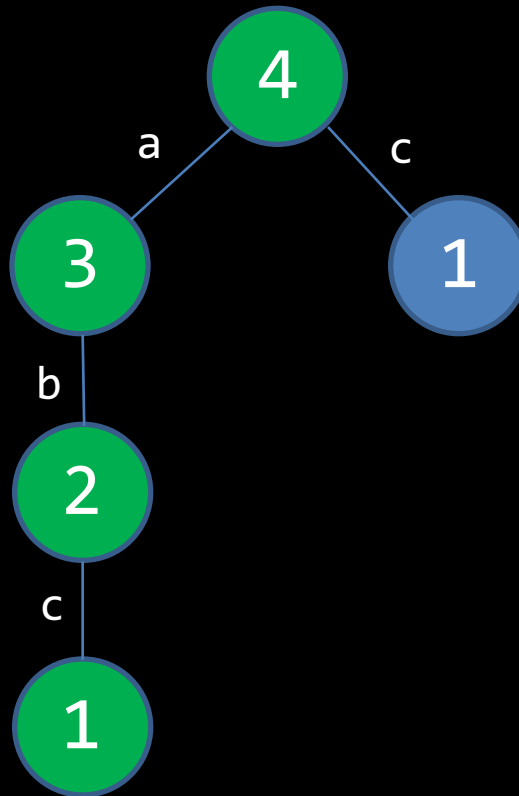
# Trie

insert "abc" to {"a", "ab", "c"}



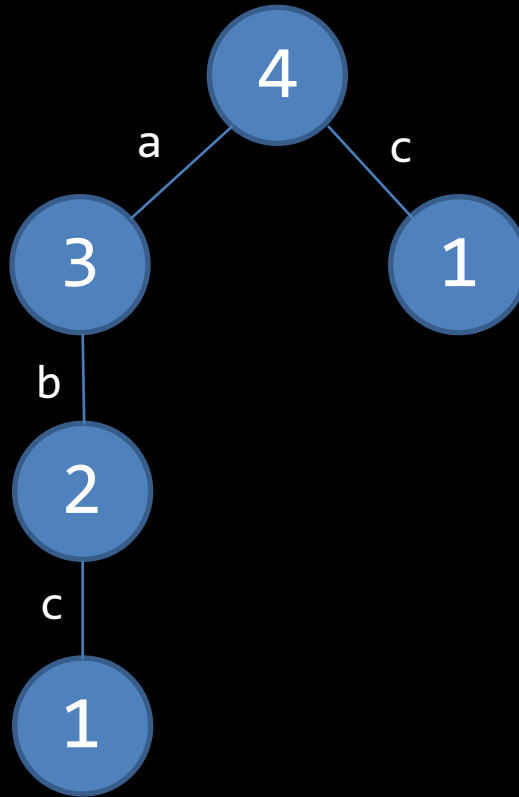
# Trie

insert "abc" to {"a", "ab", "c"}



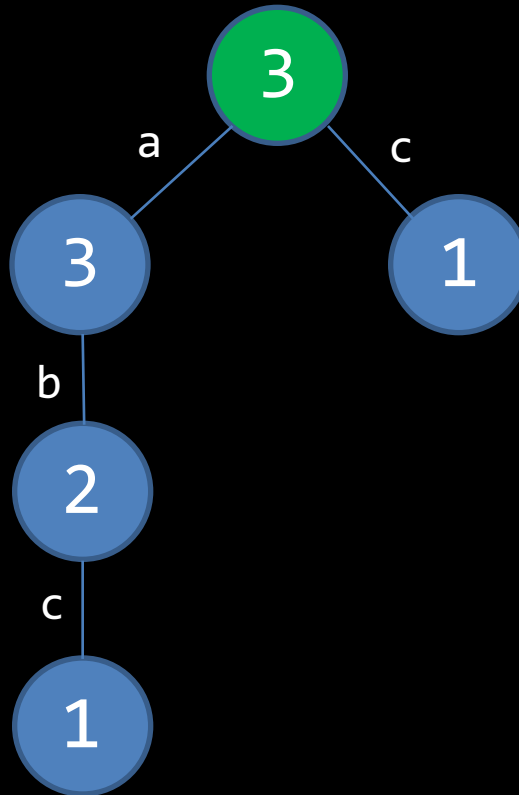
# Trie

delete "ab" from {"a", "ab", "abc", "c"}



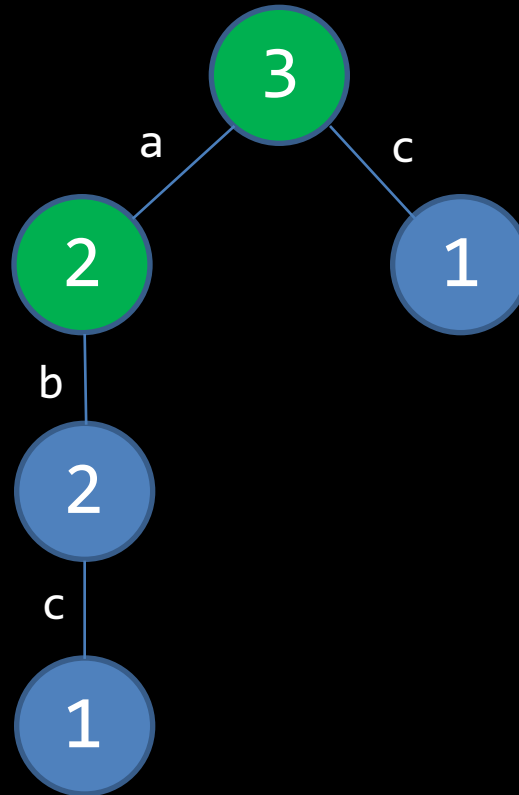
# Trie

delete "ab" from {"a", "ab", "abc", "c"}



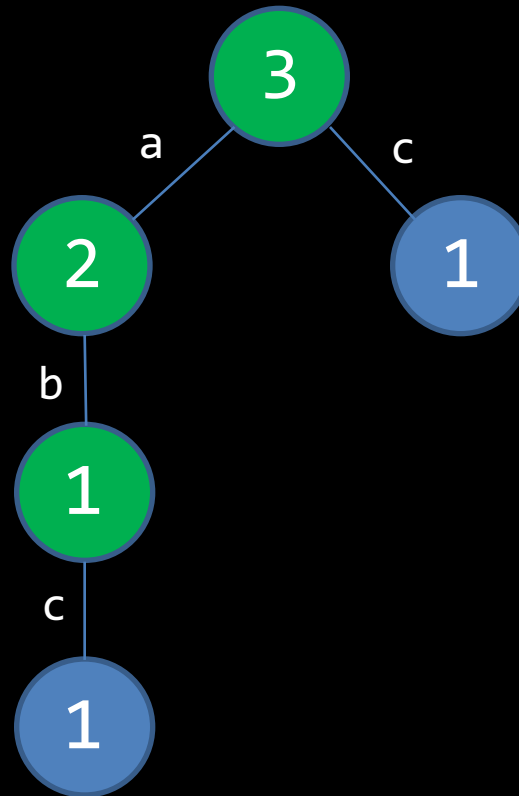
# Trie

delete "ab" from {"a", "ab", "abc", "c"}



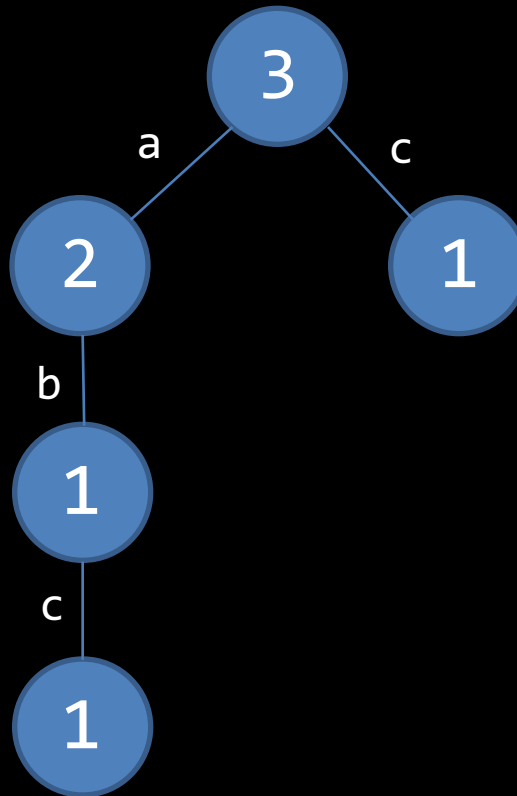
# Trie

delete "ab" from {"a", "ab", "abc", "c"}



# Trie

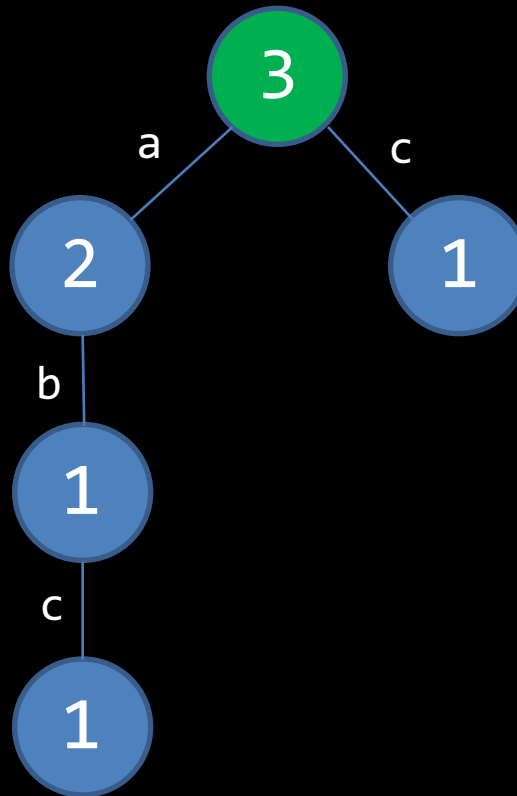
count string with the prefix “a”





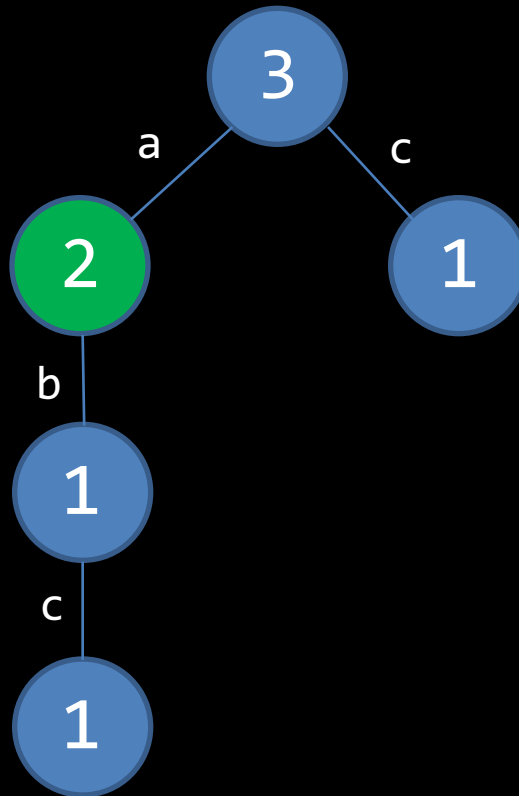
# Trie

count string with the prefix “a”



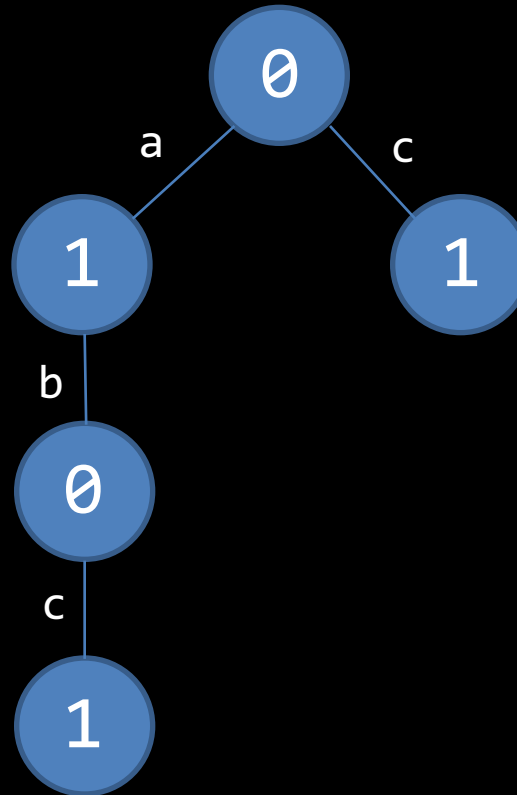
# Trie

count string with the prefix “a”



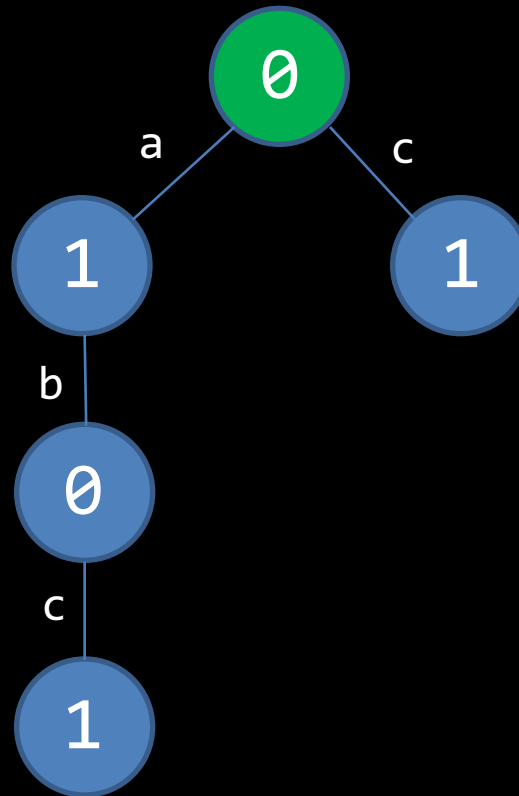
# Trie

count string "a"



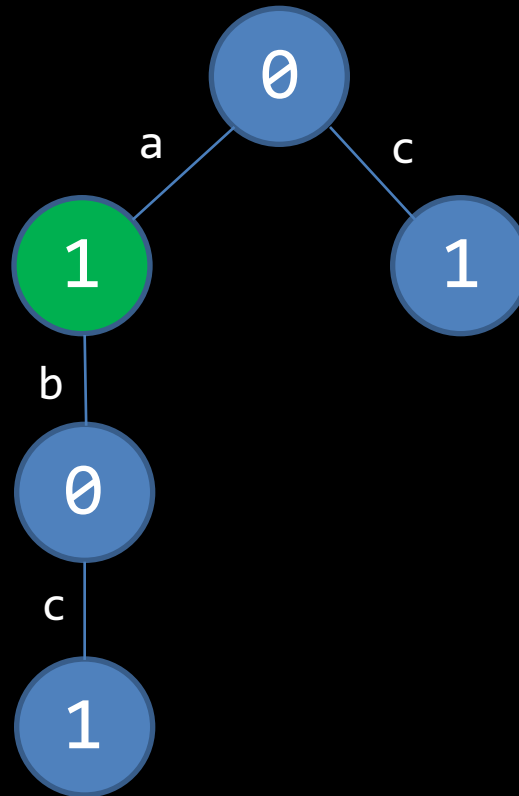
# Trie

count string "a"



# Trie

count string "a"



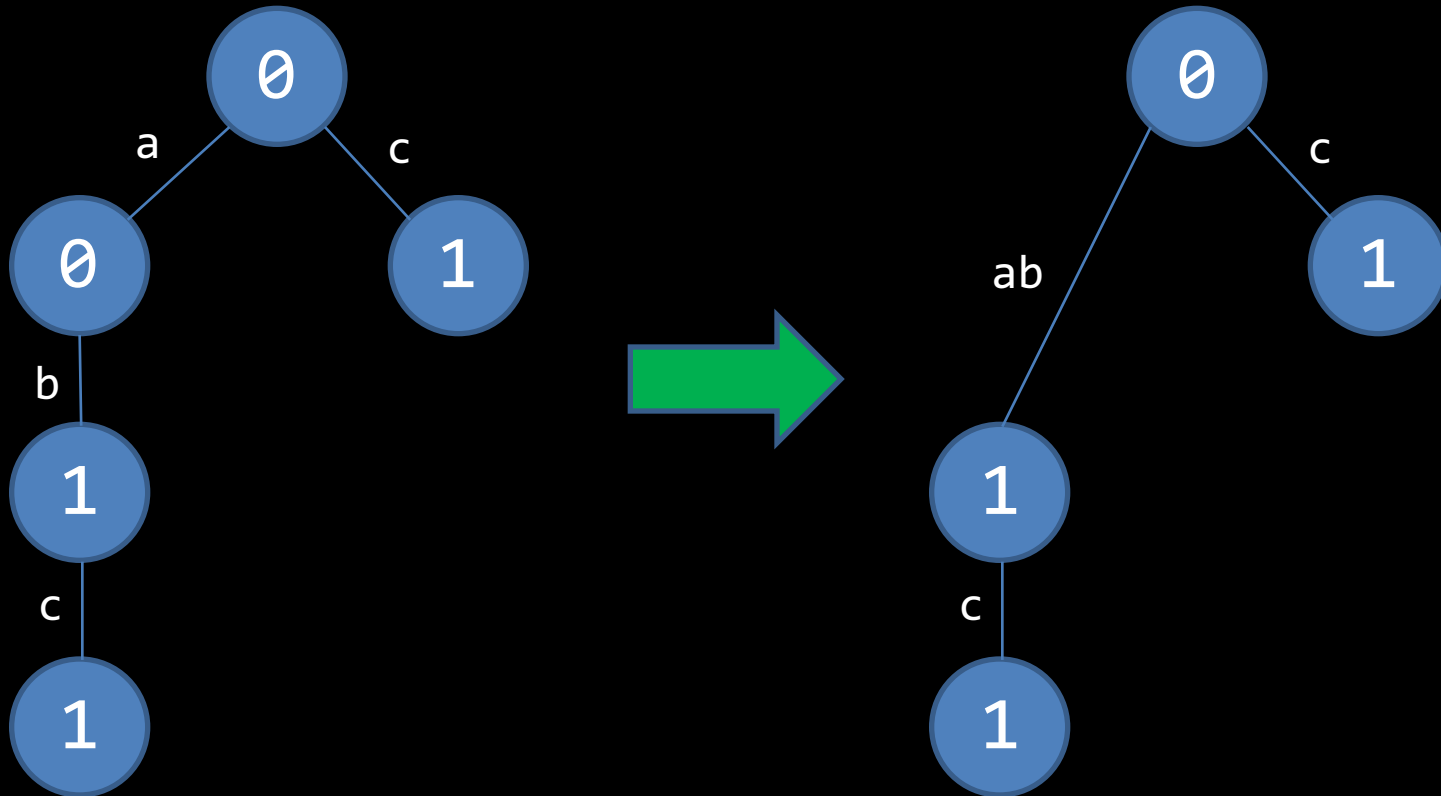
# Trie

radix tree/ compressed trie

- a variant (or an optimization)
- some useless nodes in trie are merged with its parent
- harder implementation

# Trie

radix tree/ compressed trie



# KMP

## HKOJ 01002 A Counting Problem

- given 2 string  $x$ ,  $y$
- length  $\leq 1000$
- count the occurrence of  $y$  in  $x$



# KMP

HKOJ 01002 A Counting Problem

for each starting position in  $x$   
if the  $|y|$  characters follow  
exactly match with  $y$

$ans++$

- $O(NM)$

# KMP

## HKOJ 01002 A Counting Problem

- given 2 string  $x, y$
- length  $\leq 1000$   $1e5$
- count the occurrences of  $y$  in  $x$
- a faster algorithm is needed

# KMP

- KMP is a solution
  - Knuth-Morris-Pratt Algorithm
- $O(N+M)$
- How it works?

KMP

brute force

AHKOIHHKOHKABC  
HKOHK

# KMP

brute force

AHKOIHHK OIHK ABC

HK OIHK

KMP

brute force

AHKOIHHKOIHKABC  
HKOIHK

# KMP

brute force

AHKOIHHK OIHK ABC  
HK OIHK

# KMP

brute force

AHKOIHHK OIHK ABC  
HK OIHK



# KMP

brute force

AHKOIHHKIOIHKABC  
HKOIHK

# KMP

brute force

AHKOIHHK OIHK ABC  
HK OIHK

# KMP

brute force

AHKOI**H**HKOIHKABC  
HKOI**H**K

# KMP

brute force

AHKOIHKOIHKABC  
HKOIHK

KMP

brute force

AHKOIHHKOHKABC  
HKOHK



# KMP

brute force

AHKOIHHK OIHK ABC  
HK OIHK

# KMP

optimization 1

AHKOIHKOIHKABC  
HKOIHK

KMP

optimization 1

-HKOIHH-----  
HKOIHK



# KMP

optimization 1

-HKOIHH-----  
HKOIHK

# KMP

optimization 1

-HKOIH~~H~~-----

HKOIHK

# KMP

optimization 1

-HKOIH<sup>I</sup>H-----  
HKOIH<sup>I</sup>HK

# KMP

optimization 1

-HKOIHH-----

HKOIHK

# KMP

## optimization 1

HKOIHK -> move 4 units

HKOIHK -> move 4 units

HKOIHK -> move 1 unit

HKOIHK -> move 2 units

HKOIHK -> move 1 unit

# KMP

## optimization 1

HKOIHK -> move 4 units

HKOIHK -> move 4 units

HKOIHK -> move 1 unit

HKOIHK -> move 2 units

HKOIHK -> move 1 unit

# KMP

optimization 1

HKOIHK -> move 4 units

length of move = length of green string - length of underline / 2

# KMP

## optimization 1

- precompute an array of length of underline in  $O(M)$
- simple dynamic programming



# KMP

optimization 1

$O(\text{number of operations in the searching stage}) =$

$O(\text{number of increments of the green part} + \text{number of decrements of the green part}) = O(N)$

overall time complexity =  $O(N+M)$

# KMP

optimization 2

HKOIHK -> move 4 units

length of move = length of green  
string - length of underline / 2

# KMP

optimization 2

---HKOI#-----  
HKOIHK

KMP

optimization 2

---HKOI#-----  
HKOIHK

# KMP

optimization 2

---HKOI#-----  
1HKOIHK

# KMP

optimization 2

---HKOI#-----  
12HKOIHK

# KMP

optimization 2

---HKOI#-----  
123HKOIHK

# KMP

optimization 2

---HKOI#-----  
1234HKOIHK



# KMP

optimization 2

*HKOIHK* -> move 4+1 units

*HKOIHKOI* -> move 4+4 units

HKOIHHKI -> move 5-1 units

HKOIHKKI -> move 5-1+1 units

# KMP

optimization 2

HKOIH*KKI* -> move 5-1+1 units

length of move = length of green  
string - length of underline / 2 +  
length of italic / 2

# KMP

## optimization 2

- by Knuth, the K in KMP
- no improvement on time complexity
- thus sometimes being ignored

# KMP

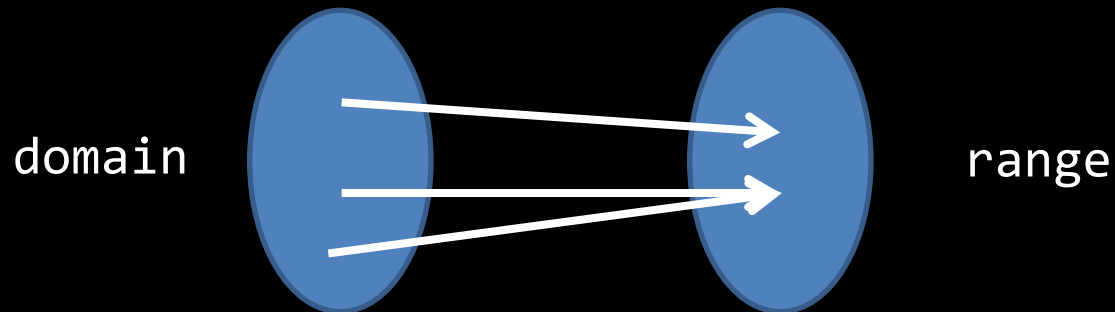
## Remarks

- KMP works with one target and multiple input
- e.g.: finding  $w$  from  $x, y, z$

Take a break

# Hashing

- hash table
  - taught in Data Structures (II)
- string hashing:  $\text{String} \rightarrow \text{Number}$ 
  - a function
  - not necessarily injective



# Hashing

- advantages
  - can be stored in array ~hash table
  - can be compared/searched in  $O(1)$
- disadvantages
  - inverse may not be easily found
  - not injective -> may crash
    - e.g.: output extra results in string matching

# Hashing

- rolling hash
  - “sliding window” fast (in  $O(1)$ )
- Rabin-Karp hash
  - common rolling hash
  - addition
    - can be replaced by xor
  - multiplication
  - mod



# Hashing

## Rabin-Karp hash

$$H = (c_0 a^{n-1} + c_1 a^{n-2} + \dots + c_{n-2} a + c_{n-1}) \bmod p$$

- $H$  = output number
- $c$  = ASCII code / A=0, B=1, ..., Z=25
- $a$  = 256 / 26
- $P$  = a large prime, e.g.  $1e9+7$

# Hashing

“HKOI”

$$c = 7, 10, 14, 8$$

$$a = 26$$

$$p = 64997$$

$$H = (7*26^3 + 10*26^2 + 14*26 + 8) \bmod 64997$$

$$H = 130164 \bmod 64997 = 170$$

# Hashing

“GO”

$$c = 6, 14$$

$$a = 26$$

$$p = 64997$$

$$H = (6 * 26 + 14) \bmod 64997$$

$$H = 170 \bmod 64997 = 170$$

# Hashing

$$H(\text{"HKOI"}) = 170 = H(\text{"GO"})$$

"HKOI" == "GO" ???

# Hashing

- assume range is any string
- surjective function
- all output with equal probability

If the input has 2 strings,

$$P(\text{crash}) = 1/64997 = 0.0000154$$

- negligible?

# Hashing

If the input has 700 strings,

$$P(\text{crash}) = 1 - (64996/64997)^{700}$$

$$P(\text{crash}) = 0.01$$

- NOT negligible

# Hashing

solution 1

- give up hashing
- use exact algorithm
- e.g.: KMP for string matching

# Hashing

solution 2a

- use a larger prime

$$p = 1e9+7$$

If the input has 700 strings,

$$P(\text{crash}) = 1 - [(1e9+6)/(1e9+7)]^{700}$$

$$P(\text{crash}) = 0.0000007$$



# Hashing

solution 2b

- use primes
- compute  $H$  more than once with different primes
- $P(\text{crash})$  decreases exponentially

# Hashing

## solution 2

- cannot be stored in array
- use map/set instead
- increase time complexity by  $\log N$
- not always need to store the result

# Hashing

“slide”

$$H = (c_0 a^{n-1} + c_1 a^{n-2} + \dots + c_{n-2} a + c_{n-1}) \bmod p$$

$$H' = (c_1 a^{n-1} + c_2 a^{n-2} + \dots + c_{n-1} a + c_n) \bmod p$$

$$H' = ((H - c_0 a^{n-1}) \times a + c_n) \bmod p$$

$$H = ((H' - c_n) \times a^{-1} + c_0 a^{n-1}) \bmod p$$

# Hashing

“slide”

- use long long if there is possibility of overflow
- subtraction
  - avoid negative numbers
  - $(a-b) \bmod p \rightarrow (a+p-b\%p)\%p$

# Hashing

“slide”

- inverse

- taught in Mathematics in OI

- Fermat's little theorem

- $p$  needs to be a prime first

- power mod

- $a^{-1} \bmod p \rightarrow a^{(p-2)} \bmod p$

# Hashing

## segment tree

- taught in Data Structures (III)
- find the hash of substring in  $O(\log N)$
- with precomputation in  $O(N \log N)$

## Arithmetic Sequence

Problem setter: Alex Poon

### Challenge

If the constraint is  $1 \leq N \leq 100,000$  and  $R = N$  for all subtasks

How can we write a program to verify if a sequence is valid in  $O(N \lg N)$ ?

tag: Data structure (III), (String algorithm/data structure (II))

# Suffix array

- Indices of all non-empty suffixes
- sorted
- “IOIHKGHKOI”
- “HKGHKOI”
- “HKOI”
- ...

# Suffix array

- IOIHKGHKOI • 5
- OIHKGHKOI • 3
- IHKGHKOI • 6
- HKGHKOI • 9
- KGHKOI • 2
- GHKOI • 0
- HKOI • 4
- KOI • 7
- OI • 8
- I • 1



# Suffix array

$O(N^2 \log N)$

- brute force
- user-defined compare function
- or sort  
`vector<pair<vector<char>,int> >`

# Suffix array

$O(N \log N \log N)$

- compare suffixes by the length of the power of two
- no. of sort:  $O(\log N)$
- no. of comparisons in each sort:  $O(N \log N)$

# Suffix array

$O(N \log N \log N)$

len=1,2,4,8...

rank[i]=the rank of the first len characters of the ith suffix

# Suffix array

- IOIHKGHKOI
- OIHKGHKOI
- IHKGHKOI
- HKGHKOI
- KGHKOI
- GHKOI
- HKOI
- KOI
- OI
- I

# Suffix array

- IO
  - OI
  - IH
  - HK
  - KG
  - GH
  - HK
  - KO
  - OI
  - I\_
- len=2

# Suffix array

- IO
  - OI
  - IH
  - HK
  - KG
  - GH
  - HK
  - KO
  - OI
  - I\_
- 4
  - 7
  - 3
  - 1
  - 5
  - 0
  - 1
  - 6
  - 7
  - 2

# Suffix array

$O(N \log N \log N)$

len = 1

- trivial

# Suffix array

$O(N \log N \log N)$

$len > 1$

- `rank[]` is describing the status of  $len/2$  till the end of the sorting of `len`



# Suffix array

$O(N \log N \log N)$

len > 1

compare(a,b) //overriding '<'

if r[a]!=r[b] return r[a]<r[b]

return r[a+len/2]<r[b+len/2]

- r[i]=-1 if i>=n

# Suffix array

$O(N \log N \log N)$

$len > 1$

- update `rank[]` at the end of sorting
- handle the case of equality

$r[a] == r[b]$  and  $r[a + len/2] == r[b + len/2]$

# Suffix array

$O(N \log N \log N)$

$len \geq N$

- the whole process is done after this sorting

# Suffix array

faster than  $O(N \log N \log N)$

- it exists
- harder to understand/implement
- not so important unless you aim for NOI

# Suffix array

longest common prefix lcp[]

- useful in many application of suffix array
- brute force needs  $O(N^2)$
- actually  $O(N)$  is enough

# Suffix array

longest common prefix lcp[]

$i$  : ABC...

$i+1$  : BCD...

$\text{lcp}[i] = 0$

$j$  : ABC...

$j+1$  : ABD...

$\text{lcp}[j] = 2$

# Suffix array

longest common prefix lcp[]

j : ABC...

j+1 : ABD...

$\text{lcp}[j] = 2$

k : BC...

k+1 : BD...

$\text{lcp}[k] \geq \text{lcp}[j] - 1$

# Suffix array

longest common prefix  $lcp[]$

- the 2 new suffixes obtained by deleting the first character may not be neighbors
- the property  $lcp[k] \geq lcp[j]-1$  still holds
- can be computed in  $O(N)$  by the original order of the suffixes



# Suffix array

string matching

- binary search
- find the suffix beginning with the target
- i.e. target is the prefix of it
- $O(T \log N)$
- can be improved by `lcp[]`

# Suffix array

other application

- longest repeated substring
- longest common substring
- longest palindromic substring

# Other string algorithms

- Z algorithm
  - string matching in  $O(N)$
- Manacher algorithm
  - longest palindromic substring in  $O(N)$

# Other string algorithms

- suffix tree
  - a tree storing all the suffixes of a string
  - compute suffix array in  $O(N)$
- Aho-Corasick Algorithm
  - multi-target version of KMP
  - linear time complexity

# Practice

[judge.hkoi.org/tag/strings](http://judge.hkoi.org/tag/strings)

[codeforces.com/problemset/tags/strings](http://codeforces.com/problemset/tags/strings)

- HKOJ 01002 A Counting Problem
- HKOJ M0932 String Rotation
  
- UOJ 35 後綴排序
- NOI 2015 品酒大會

# Reference

- 演算法筆記
- [www.csie.ntnu.edu.tw/~u91029/](http://www.csie.ntnu.edu.tw/~u91029/)

The end