6.3.17

Topic: Strings

Terminology:
text $T$ , $|T| = n$

pattern $P$ , $|P| = m$

alphabet $\Sigma$
↑ e.g binary , small (good), large (poly(n))

# Pattern matching

Find occurence of $P$ in $T$ (as a substring)

There exist solutions in $O(n+m)$:

- Knuth, Mons, Pratt alg.
- Boyer - Moore

etc...

Idea: preprocess pattern

Data structure perspective:
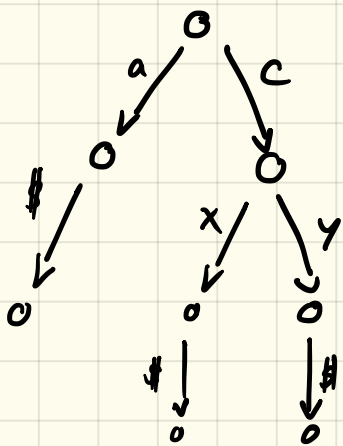preprocess $T$ in $O(n)$
query $P$ in $O(m)$

## Problem 1

Preprocess $T_1, \ldots, T_k$ $\quad n = \sum_i |T_i|$

Query $P$ is $P = T_i$,

or find pred./succ. of $P$ in $\{T_i\}$
w.r.t. lex. order

Trie: Rooted tree, where strings correspond
to root-leaf paths.

If the tree is in-order $\leadsto$ sorted $T$.



$T = \{a, cx, cy\}$

| nodes store children | query | space |
|---|---|---|
| ① array + blank cells store [1] pointers | $O(m)$ | $O(n\|\Sigma\|)$ |
| ② balanced search tree (BST) (set/map) | $O(m \log \|\Sigma\|)$ | $O(n)$ |
| ③ hash table [2] | $O(m)$ | $O(n)$ |
| ③.5 van Emde-Boas tree | $O(m \lg\lg \|\Sigma\|)$ | $O(n)$ |
| ③ + ③.5 hash + vEB [3] | $O(m + \lg\lg \|\Sigma\|)$ | $O(n)$ |
| ④ weight balanced BST | $O(m + \log k)$ [4] | $O(n)$ |
| ⑤ indirection | $O(m + \log \|\Sigma\|)$ | $O(n)$ |



1
2 no PRED/SUCC
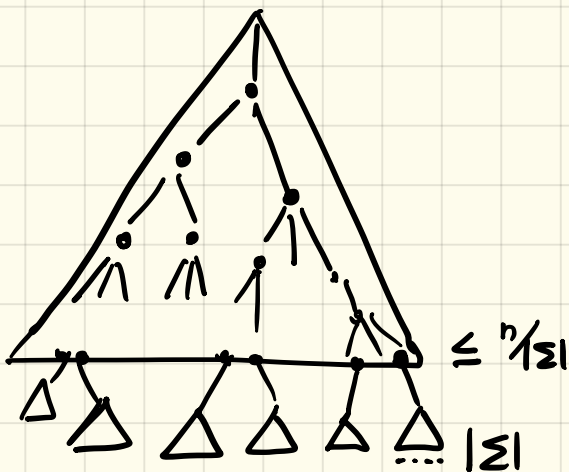3 works w.h.p. / vEB only used once
4 cf. next page

④



[a-b]   [c-e]
a   b   [c-d]   e
c   d

Claim: Going down twice:
- advances P
- reduces # candidates to ⅔

⑤



$\leq {}^n/_{|\Sigma|}$

$|\Sigma|$

① on branching nodes (BN)
   in top part  ($\leq {}^n/_\Sigma$ BN)

① on top leaves

② on non-branching

④ on bottom

⟶  string sorting  $O(n + k \log|\Sigma|)$
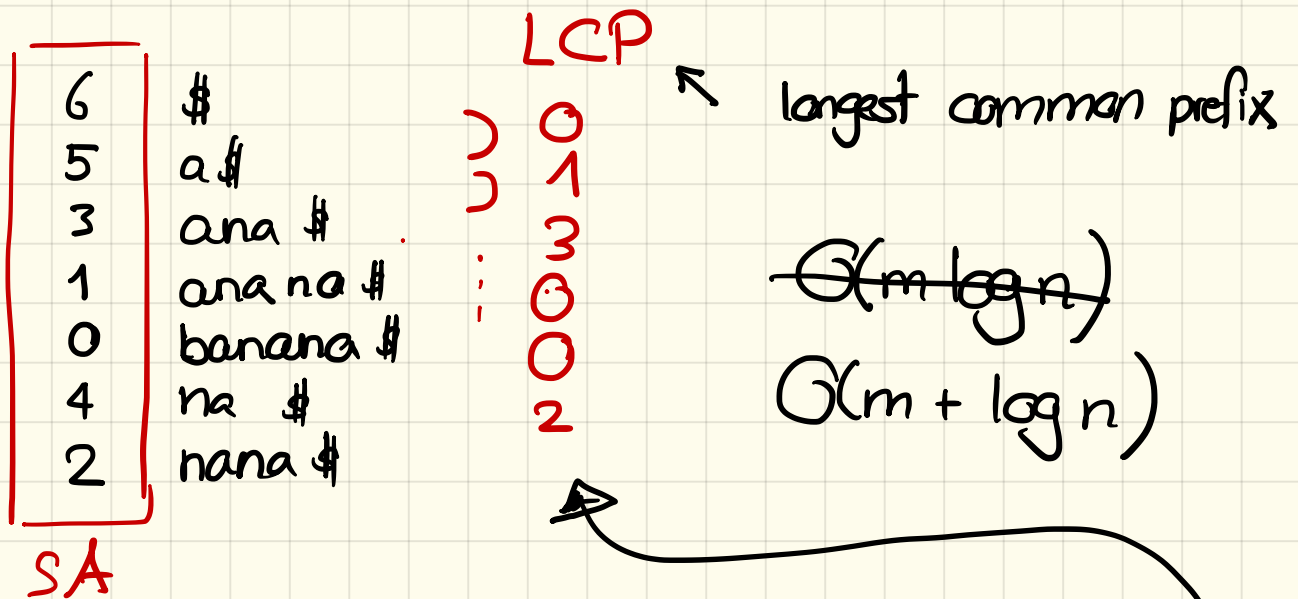
# Compressed trie



Assume long input (e.g. war and piece)
and short queries

## Suffix tree : compressed tie of all $T[i]$



banana #

$n+1$ leaves

$O(n)$ space

$O(m + \lg |\Sigma| + k)$ all occurences of $P$

└ # of occurences

hoshing: $O(m + k)$    all occ. of $P$

E.g.   Longest Match of $T[i:], [j:] = LCA(i,j)$

$\rightarrow$  How do we construct a suffix tree?
        In $O(n)$ .... complicated ....

$\rightarrow$  Consider instead  <u>suffix array</u>.

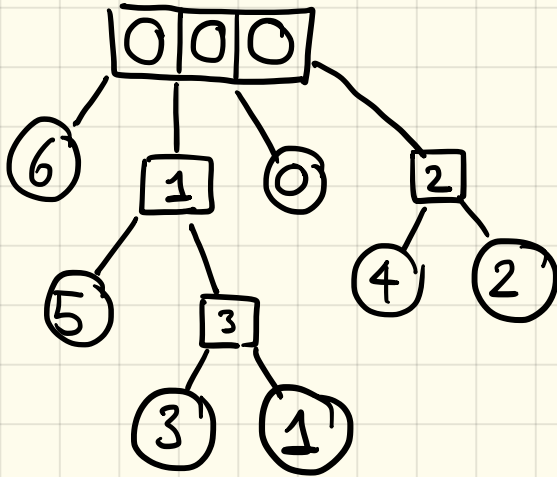| SA | | LCP | |
|---|---|---|---|
| 6 | # | 0 | $\nwarrow$ longest common prefix |
| 5 | a# | 1 | |
| 3 | ana # | 3 | |
| 1 | anana # | 0 | |
| 0 | banana # | 0 | ~~$O(m \log n)$~~ |
| 4 | na # | 2 | $O(m + \log n)$ |
| 2 | nana # | | |

$LCP(T[i:], [j:]) = RMQ$ in  of corrsp. pos. in

<u>Claim</u>  Construction of SA equiv. to ST.

## $ST \rightarrow SA$: in order traversal
→ euler tour gives LCP

## $SA \rightarrow ST$: Cartesian tree $O(n)$ on LCP (mostly)



## SA construction in $O(n + sort(\Sigma))$:

1) Sort $\Sigma$,

2) replace $\Sigma$ with $[1, \ldots, |\Sigma|]$

3) $T_0 = [\ (T[3i], T[3i+1], T[3i+2])\ $ for $i = 0, \ldots]$

$T_1 = [\ ''\ 3i+1 \qquad '' +2 \quad '' \ +3 \qquad '' \qquad ]$

$T_2 = [\ ''\ 3i+2 \qquad '' +3 \quad +4 \qquad '' \qquad ]$

4) recurse on $T_0 \circ T_1$ = relative order of $T[3i], T[3i+1]$

5) radix sort of $T_2$

$$T_2[i:] \approx T[3i+2:] \approx (T[3i+2], T[3i+3])$$
$$\approx (T[3i+2], T_0[i+1:])$$

6) merge $T_0 \cdot T_1$ with $T_2$

| $T_0 \cdot T_1$ |
|---|

| $T_2$ |
|---|

$$T_0[i:] \text{ vs } T_2[j:]$$
$$\approx (T[3i], T_1[i:]) \text{ vs } (T[3j+2], T_0[j+1:])$$

similarly $T_2$ vs $T_1$

e.g.   banana $\#$

$T_0$    ban    ana    $\$$

$T_1$    ana    na$\#$

$T_2$    nan    a$\#$