

TODAY: Strings

- tries & trays
- compressed tries
- suffix trees & arrays
- document retrieval
- linear-time construction

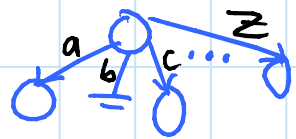
String matching: given text  $T$  & pattern  $P$ ,  
 here both strings over alphabet  $\Sigma$ ,  
 find some/all occurrences of  $P$  in  $T$   
 as substrings

- one-shot:  $O(T)$  time [Knuth, Morris, Pratt - SICOMP, 1977;  
 Boyer & Moore - CACM, 1977; Karp & Rabin - IBM JRD, 1987]
- static DS: preprocess  $T$ , query =  $P$
- goal:  $O(P)$  query  
 $O(T)$  space

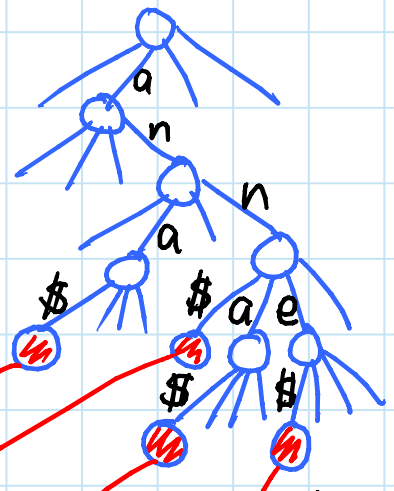
- other data structures consider when  $P$   
 has wildcards, or when  $P$  need not match  
 as an exact substring (Hamming/edit distance)  
 ~ see e.g. [Cole, Gottlieb, Lewenstein - STOC 2004]  
 [Maab & Novak - CPM 2005]

Warmup: predecessor among strings  $T_1, \dots, T_k$   
 (e.g. library search)

Trie = rooted tree with child branches labeled with letters in  $\Sigma$



- to represent strings as root-to-leaf paths in a trie, terminate them with a new letter \$ (otherwise can't distinguish prefixes as absent or present)



- e.g.:  $\{ana, ann, anna, anne\}$

- in-order traversal of leaves = sorted strings

Trie representation:  $T = \# \text{ nodes in trie} \leq \sum_{i=1}^k |T_i|$

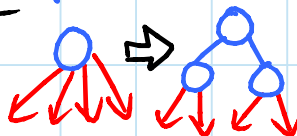
node stores children:

① as array

↳ blank cells store predecessor/successor

query  $O(P)$  space  $O(T|\Sigma|)$

② as balanced BST



$O(P \lg |\Sigma|)$   $O(T)$

③ as hash table

$O(P)$   $O(T)$

↳ doesn't support predecessor queries/sorting

③.5 as van Emde Boas/y-fast

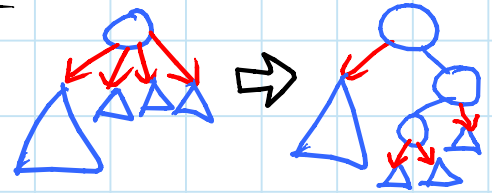
$O(P \lg \lg |\Sigma|)$   $O(T)$

③.75 = ③ + ③.5 (only need vEB when fall off)  $O(P + \lg \lg |\Sigma|)$   $O(T)$

[Farach-Colton - personal communication, 2012]:

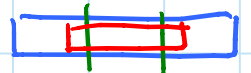
④ node stores children: as weight-balanced BST      query  $O(P + \lg k)$       space  $O(T)$   
     ↳ # descendant leaves in  $T$       ↳ # leaves

- split children in left & right halves to optimally balance sum of weights



⇒ every 2 edges followed either advances P letter or reduces # candidate T strings to 2/3

⇒ charge to  $O(P)$  or  $O(\lg k)$



⑤ leaf trimming (indirection)       $O(P + \lg \Sigma)$        $O(T)$

- cut below maximally deep nodes with  $\geq |\Sigma|$  descendant (leaves)

⇒ # leaves in top trie  $\leq |T|/|\Sigma|$

⇒ # branching top nodes  $\leq |T|/|\Sigma|$

- use ① on branching top nodes

& ① on top leaves (to find right bottom trie)

& ② on rest of top (⇒ nonbranching in  $T$ )

⇒  $O(T)$  space on top

- bottom trees have  $< |\Sigma|$  descendant (leaves)

⇒ ④ achieves  $O(P + \lg \Sigma)$  query time

↳ simplification by Farach-Colton of:

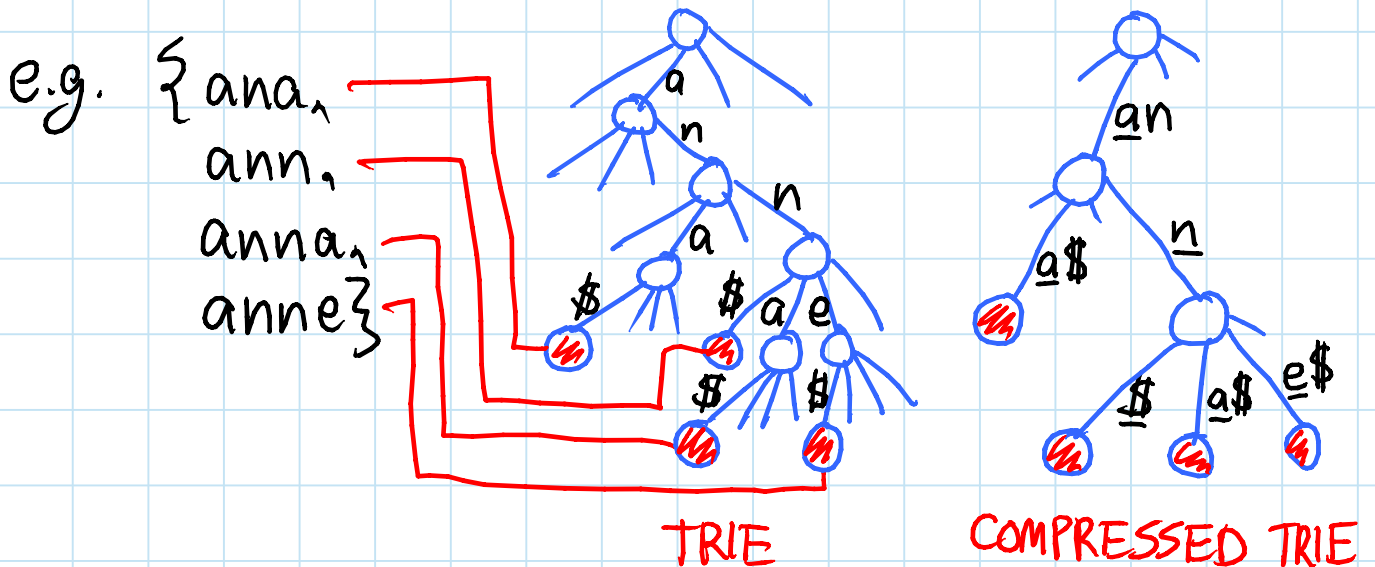
⑥ suffix trays

$O(P + \lg \Sigma)$        $O(T)$

[Cole, Kopelowitz, Lewenstein - ICALP 2006]

Application: sorting strings  $T_1, \dots, T_k$   
 - repeatedly insert into trie/tray  
 $\Rightarrow O(T + k \lg \Sigma)$   
 - typically  $O(T)$  &  $\ll O(Tk \lg k)$  via comparison

Compressed trie: contract nonbranching paths to single edge, keyed by first letter of path



- same representations apply,  
 with  $T = \#$  compressed nodes

## Suffix tree (trie):

Compressed trie of all  
|T| suffixes  $T[i:]$  of T  
(with \$ appended)

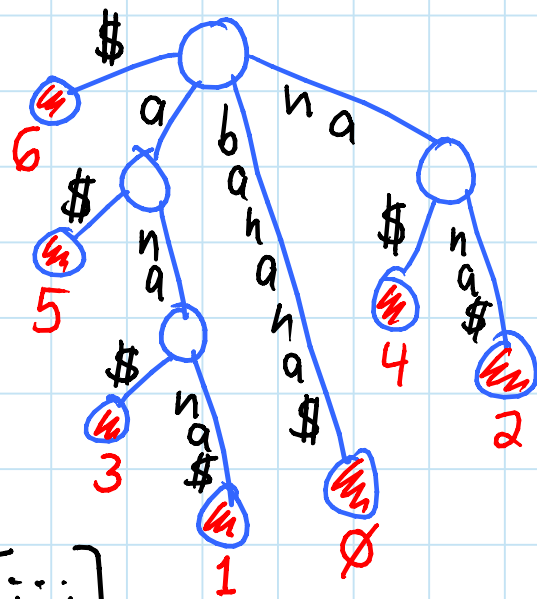
- e.g.: b a n a n a \$  
          0 1 2 3 4 5 6

- |T|+1 leaves

- edge label = substring  $T[i:j]$

⇒ store as two indices (i,j)

⇒  $O(T)$  space



## Applications:

- search for P gives subtree whose leaves correspond to all occurrences of P

-  $O(P)$  time via hash (+VEB) → leaves can still be sorted in T

-  $O(P + \lg \Sigma)$  via trays ⇒ leaves sorted in T

- list first k occurrences in  $O(k)$  more time

- every node points to leftmost descend. leaf

- leaves connected via linked list

- # occurrences in  $O(1)$  more time (subtree sizes)

- longest repeated substring in T:  $O(T)$  time

= branching node of maximum "letter depth"

- longest substring match of  $T[i:]$  vs.  $T[j:]$ :

$O(1)$  via LCA query

- all occurrences of  $T[i:j] = (|T|-j)$ th "weighted" level ancestor of leaf for  $T[i:]$  for compression
  - store nodes in long path/ladder of  $L_{15}$  in van Emde Boas predecessor DS  $\Rightarrow O(\lg \lg T)$
  - can't afford lookup tables at the bottom...
  - use ladder decomposition on bottom trees  $\Rightarrow$  jump to top of  $O(\lg \lg n)$  ladders (to reach height  $O(\lg n)$ )
  - only need predecessor query on last ladder  $\Rightarrow O(\lg \lg T)$  query &  $O(T)$  space

[Abbott, Baran, Demaine, ... - 6.897, Spr. 2005, L19.5]

- multiple documents via mult.  $\$s$ :  $T = T_1 \$_1 \dots T_k \$_k$
- count # distinct documents containing  $P$ 
  - store # distinct  $\$s$  below each node
- longest common substring in  $O(T)$ 
  - = branching node with  $\geq 2$  distinct  $\$s$  below
- find  $d$  distinct documents containing  $P$  in  $O(d)$  more "document retrieval problem" [Muthukrishnan - SODA 2002]
  - each  $\$i$  stores leaf # of previous  $\$i$
  - in interval  $[l, n]$  of leaves below a node, want first  $\$i$ , i.e.  $\$i$  storing  $< l$ , for each occ.  $i$ 

make these leaves (trim below)
  - so find  $m = \text{RMQ}(l, n)$  on array of stored values
  - if stored value at leaf  $m$  is  $< i$ : [L15]
    - found desired  $\$i$  ~ output it
    - recurse in intervals  $[l, m-1]$  &  $[m+1, n]$
- $\Rightarrow O(1)$  time per output (& can stop anytime)

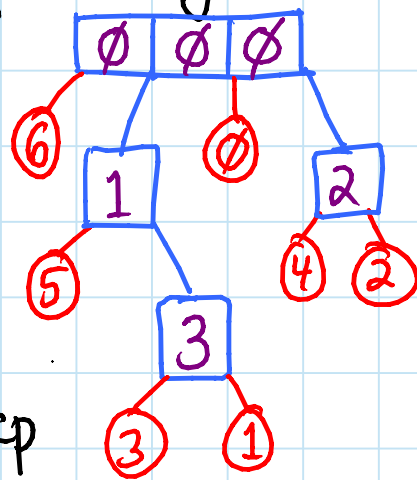


Suffix arrays: sort the suffixes of  $T$   
 just store the indices  $\Rightarrow O(T)$  space

- e.g. b a n a n a \$
- |   |    |              |
|---|----|--------------|
| 6 | \$ |              |
| 5 | a  | \$           |
| 3 | a  | n a \$       |
| 1 | a  | n a n a \$   |
| 0 | b  | a n a n a \$ |
| 4 | n  | a \$         |
| 2 | n  | a n a \$     |
- searchable in  $O(P \lg T)$  via binary search
  - $\text{lcp}[i] =$  length of longest common prefix of  $i$ th &  $(i+1)$ st suffix in order
  - when binary searching in interval  $\text{SA}[i:j]$ , only need to compare from letter  $\text{RMQ}_{\text{lcp}}(i, j-1)$
  - via RMQ of L15.  $O(P + \lg T)$  search [2007, PS4]

Suffix trees  $\leftrightarrow$  suffix arrays:

- ( $\rightarrow$ ) via in-order traversal of leaves
- ( $\leftarrow$ ) via Cartesian tree of lcp array
- put all mins at root (unlike L15)
- nonleaf child subtrees: recurse
- suffixes fit in between as leaves
- lcp value forming a node = letter depth of that node
- $\Rightarrow$  edge length = child lcp - parent lcp
- $\Rightarrow$  can reconstruct labels
- all doable in linear time [L15]
- lcp's computable in  $O(T)$  from SA [Kasai et al. - CPM 2001] or directly in suffix-array construction below



# Constructing suffix array ( $\Rightarrow$ tree) in $O(T + \text{sort}(\Sigma))$

[Kärkäinen & Sanders - IICALP 2003], inspired by  
[Farach - FOCSS 1997; Farach-Colton, Ferragina, Muthukrishnan - JACM 2000]

① sort  $\Sigma$  - initially in  $\text{sort}(\Sigma)$  time (or, if don't need children sorted, just number  $\Sigma$  arbitrarily)  
- later, radix sort in  $O(T)$  time

② replace each letter by its rank in  $\Sigma \Rightarrow \leq |\Sigma| \leq |T|$

③ form  $T_0 = \langle (T[3i], T[3i+1], T[3i+2]) \rangle$  for  $i = 0, 1, 2, \dots$

$T_1 = \langle (T[3i+1], T[3i+2], T[3i+3]) \rangle$  for  $i = 0, 1, 2, \dots$

$T_2 = \langle (T[3i+2], T[3i+3], T[3i+4]) \rangle$  for  $i = 0, 1, 2, \dots$

single "letter"

$\Rightarrow \text{suffixes}(T) \approx \bigcup_{i=0,1,2} \text{suffixes}(T_i)$

④ recurse on  $\langle T_0, T_1 \rangle \Rightarrow \frac{2}{3}|T|$  "letters"

$\rightarrow$  sorted order & lcp of  $\bigcup_{i=0,1} \text{suffixes}(T_i)$

⑤ radix sort suffixes( $T_2$ ) by writing

$T_2[i:] \approx T[3i+2:] = \langle T[3i+2], T[3i+3:] \rangle \approx \langle T[3i+2], T_0[i+1:] \rangle$

- also get lcp in suffixes( $T_2$ ): try to extend by 1

⑥ merge  $\bigcup_{i=0,1} \text{suffixes}(T_i)$  with suffixes( $T_2$ ) via:

-  $T_0[i:]$  vs.  $T_2[j:] = T[3i:]$  vs.  $T[3j+2:]$

$= \langle T[3i], T[3i+1:] \rangle$  vs.  $\langle T[3j+2], T[3j+3:] \rangle$

$T_1[i:]$

$T_0[j+1:]$

-  $T_1[i:]$  vs.  $T_2[j:] = T[3i+1:]$  vs.  $T[3j+2:]$

$= \langle T[3i+1], T[3i+2], T[3i+3:] \rangle \rightarrow T_0[i+1:]$

vs.  $\langle T[3j+2], T[3j+3], T[3j+4:] \rangle \rightarrow T_1[i+1:]$

- also get lcp: try to extend by 1 or 2

$\Rightarrow T(n) = T(\frac{2}{3}n) + O(n) = O(n)$  ( $n = |T|$ )