

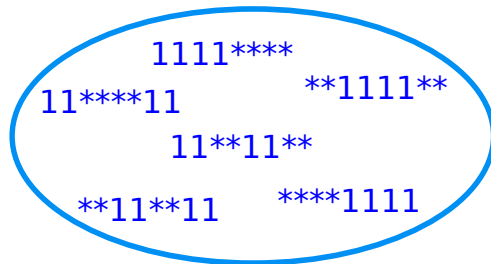


# Next generation sequencing

# Data structure in NGS mapping algorithm

- Hash
- Seed consecutive and not consecutive

# Hash and seed



ATGTCGATGCATAC

11 \*\* 11\* \*



AT--CG--

Hashes and store key + read sequence

READS

ACGCTTCGGGTTA

11 \*\* 11\* \*



AC--TT--

Looked up, if hit is found to a reads

Compute a score based on quality values and store position and score

REFERENCE

Full accuracy on the first 28 positions of the read





## Multiple Pattern Matching Problem

*given a set of patterns and a text, find all occurrences of any of the patterns in the text*

**input:** A set of  $k$  patterns  $P_1, P_2, \dots, P_k$  and text  $t=t_1\dots t_m$

**output:** All positions  $1 \leq i \leq m$  such that a substring of  $t$  starting at position  $i$  coincides with a pattern  $P_j$  for  $1 \leq j \leq k$

$t = \text{ATGGTCGGT}$

$P_1 = \text{GGT}$

$P_2 = \text{GGG}$

$P_3 = \text{CG}$

$P_4 = \text{ATG}$

Positions =  $\{1, 3, 6, 7\}$

This algorithm is resolved in  **$O(knm)$**  time where  $n$  is the length of the longest of the  $k$  patterns, by  $k$  applications of the PATTERN MATCHING algorithm

# Brute Force Is Too Slow

- The runtime of the brute force approach is too high!
  - Single *Pattern*:  $O(|Genome| * |Pattern|)$
  - Multiple *Patterns*:  $O(|Genome| * |Patterns|)$
  - $|Patterns|$  = combined length of *Patterns*

# Multiple Pattern Matching Problem



Genome

Pattern 1



Pattern 2



Pattern N





# Packing Patterns onto a Bus



# Tries

# Patterns

banana

pan

nab

antenna

bandana

ananas

nana

- **Trie:** a data structure for representing a collection of strings.
- Tries support fast pattern matching.

## Idea:

- Combine patterns into a rooted-tree with branches labeled by letters in the alphabet.
- Every string in Patterns is presented as a root-to-leaf path.

## Keyword Tree

The keyword tree for a set of patterns  $P_1, P_2, P_3, \dots, P_k$  is a rooted labeled tree satisfying the following condition:

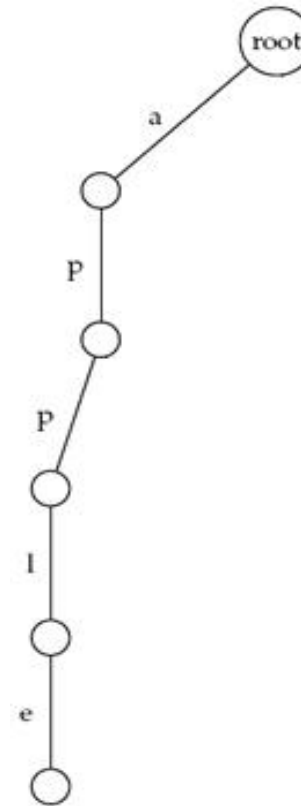
- Each edge of the tree is labeled with a letter of the alphabet
- Any two edges out of the same vertex have distinct labels
- Every pattern  $P_i$  from the set of patterns is spelled on some path from the root to the leaf

The time to build a keyword tree is  $O(N)$  where  $N$  is the total length of patterns  $P_1, P_2, P_3, \dots, P_k$

The time to solve the *Multiple Pattern Matching Problem* is  $O(N + nm)$

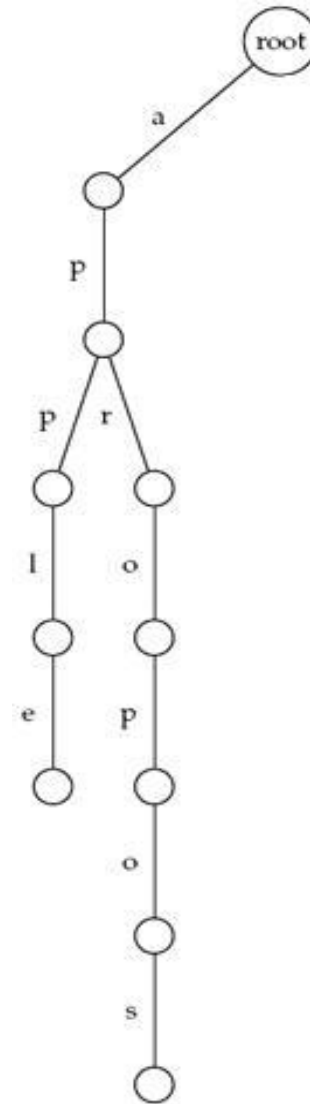
## Keyword Trees: Example

- **Keyword tree:**
  - Apple



## Keyword Trees: Example

- **Keyword tree:**
  - Apple
  - Apropos

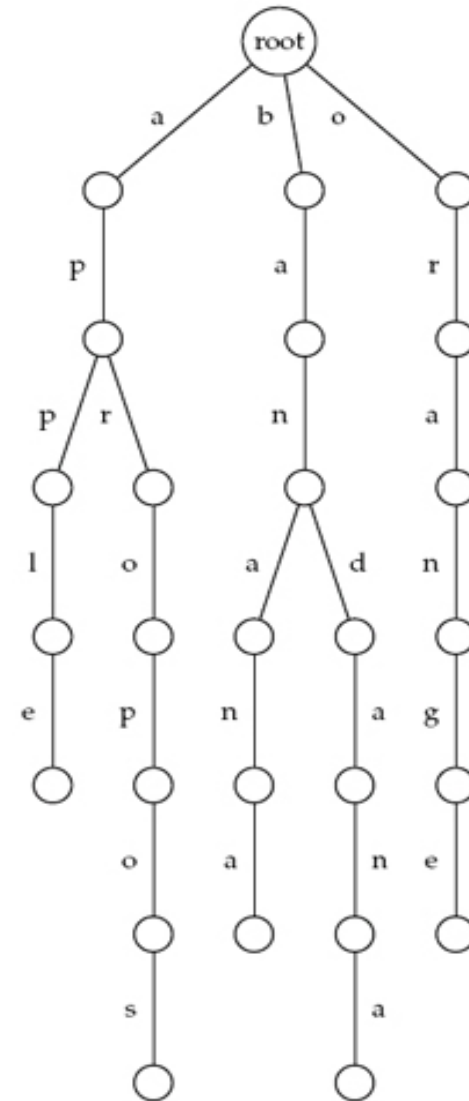






## Keyword Trees: Example

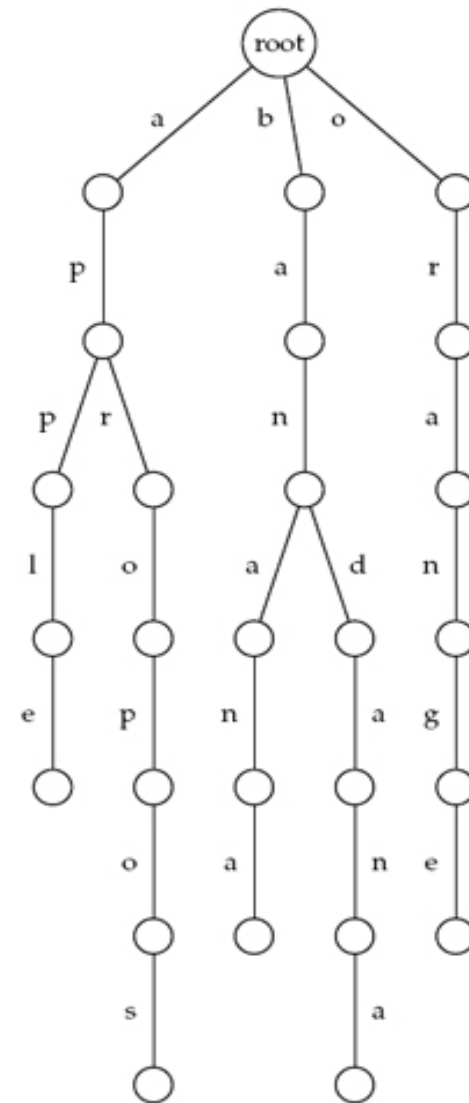
- **Keyword tree:**
  - Apple
  - Apropos
  - Banana
  - Bandana
  - Orange





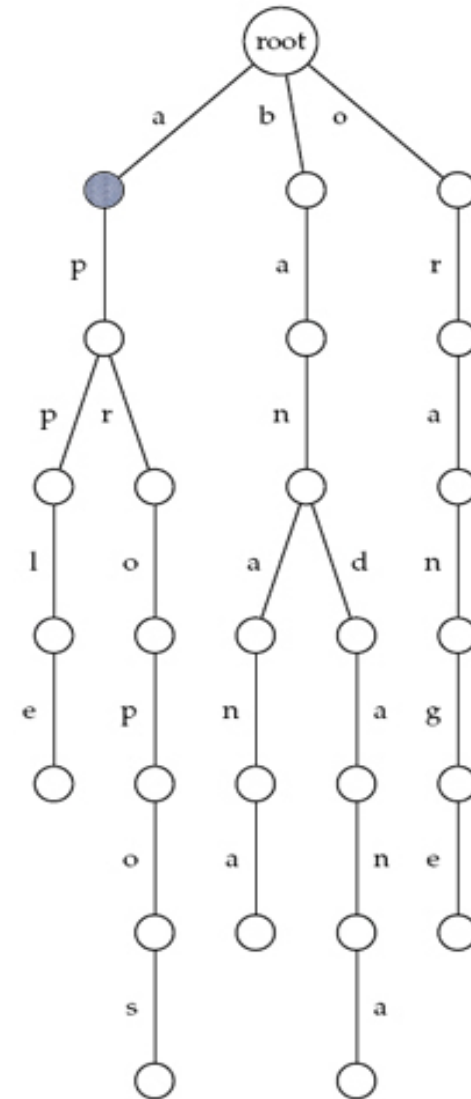
## Keyword Trees: Properties

- Stores a set of keywords in a rooted labeled tree.
- Each edge is labeled with a letter from an alphabet.
- Any two edges coming out of the same vertex have distinct labels.
- Every keyword stored can be spelled on a path from root to some leaf.
- Furthermore, every path from root to leaf gives a keyword.



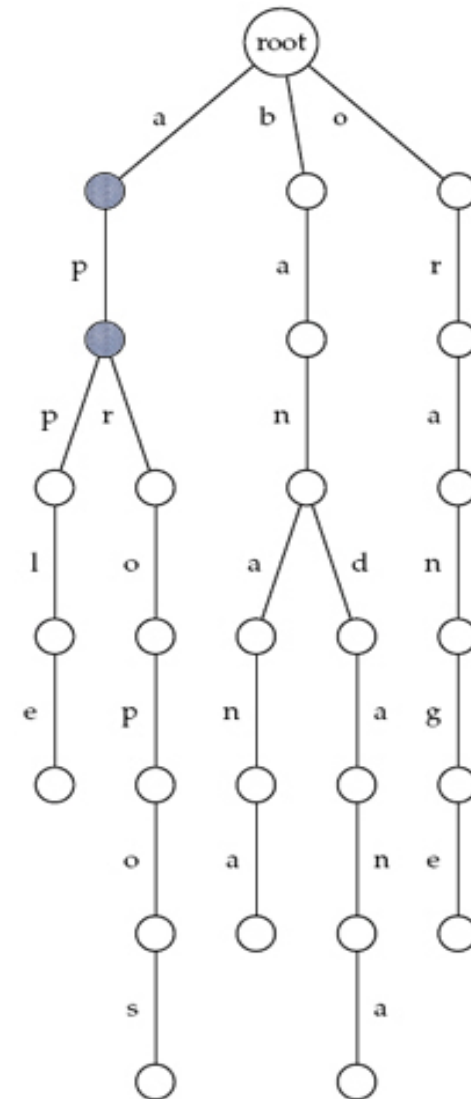
## Keyword Trees: Threading

- Thread “appeal”
  - appeal



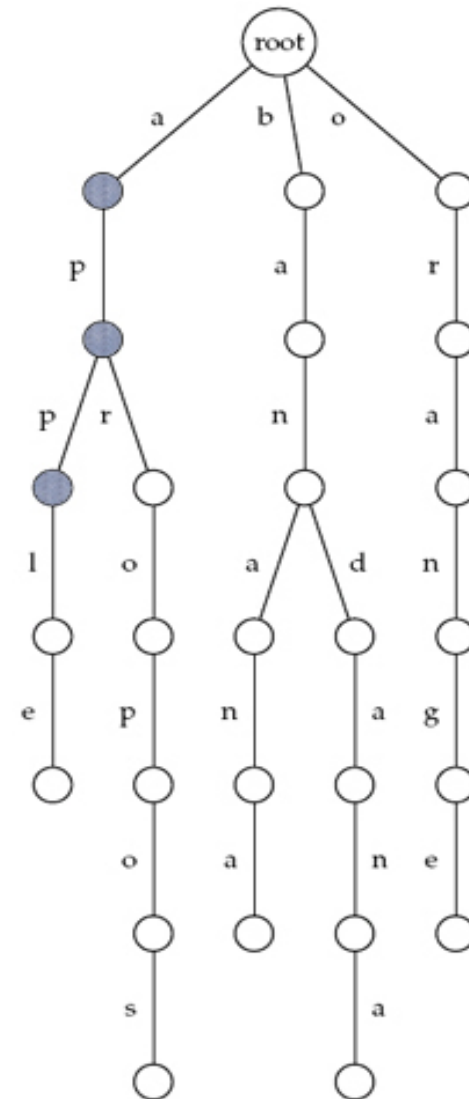
## Keyword Trees: Threading

- Thread “appeal”
  - appeal



# Keyword Trees: Threading

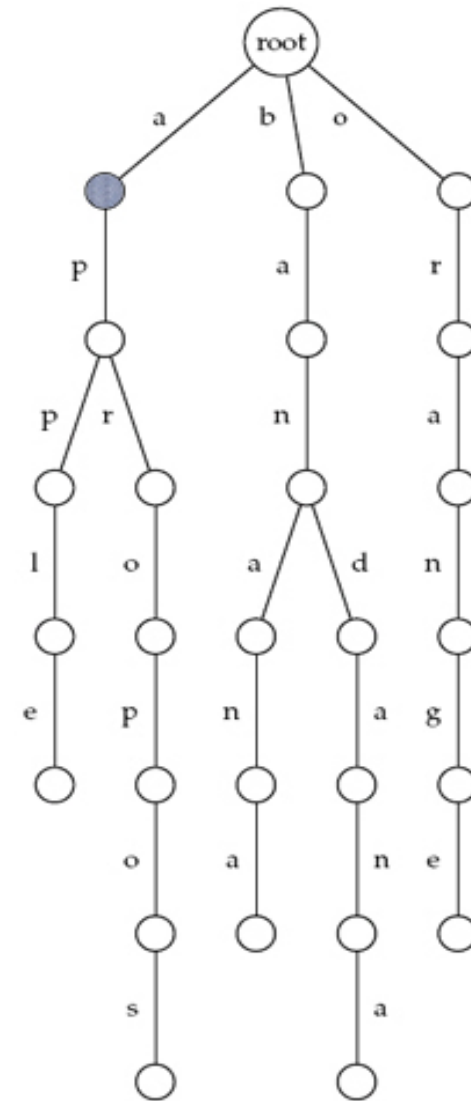
- Thread “appeal”
  - appeal





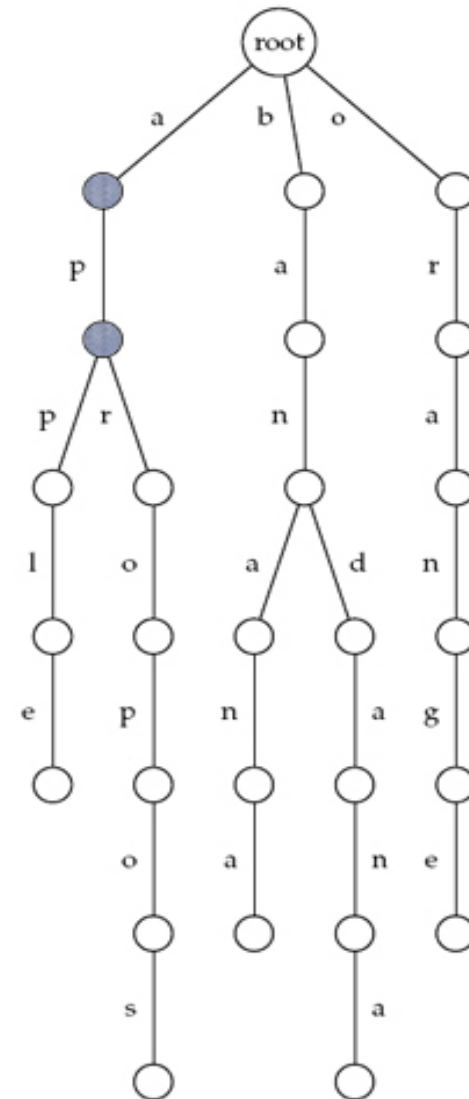
## Keyword Trees: Threading

- Thread “apple”
  - apple



## Keyword Trees: Threading

- Thread “apple”
  - apple

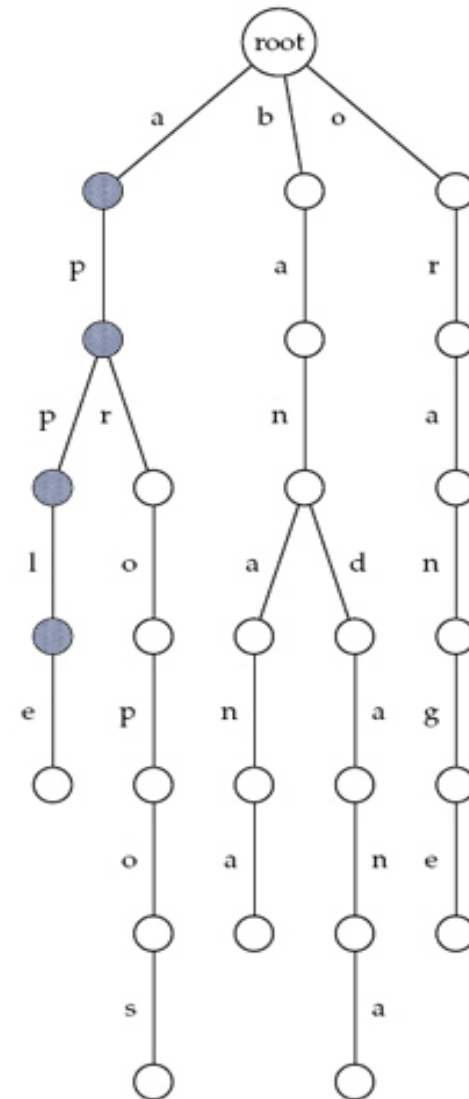






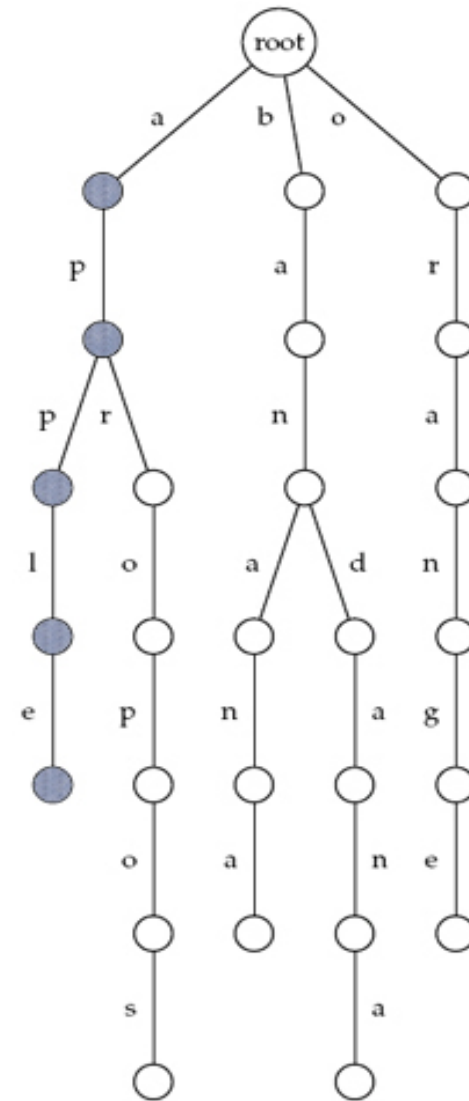
## Keyword Trees: Threading

- Thread “apple”
  - apple



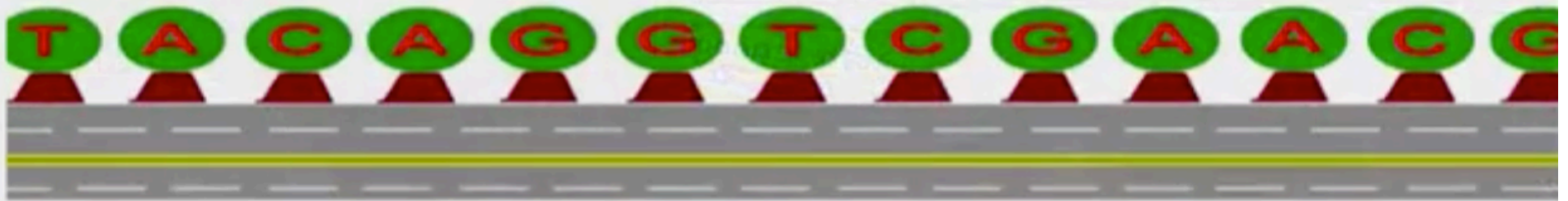
# Keyword Trees: Threading

- Thread “apple”
  - apple



# Using the Trie for Pattern Matching

- **TrieMatching:** Slide the trie down the genome.

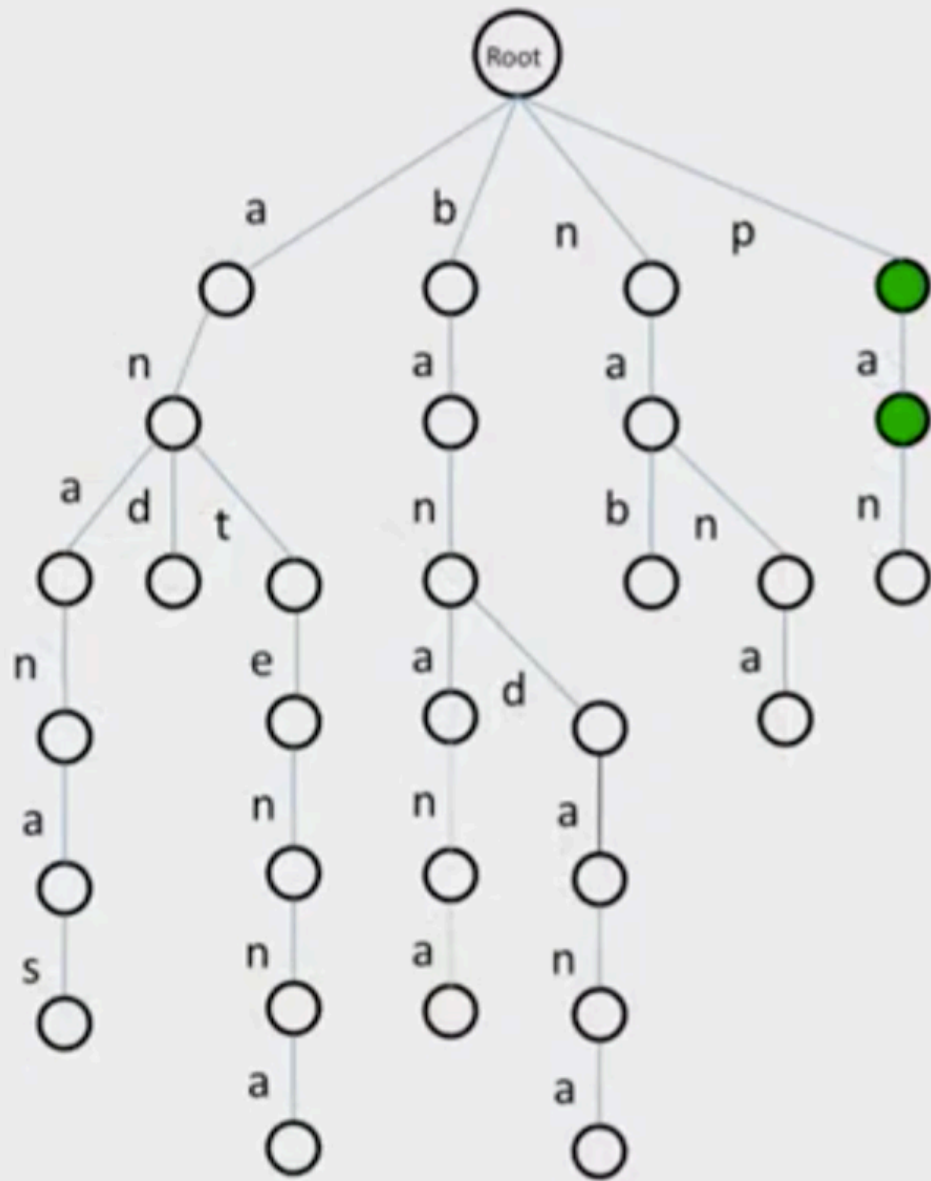


- At each position, walk down the trie and see if we can reach a leaf by matching symbols.

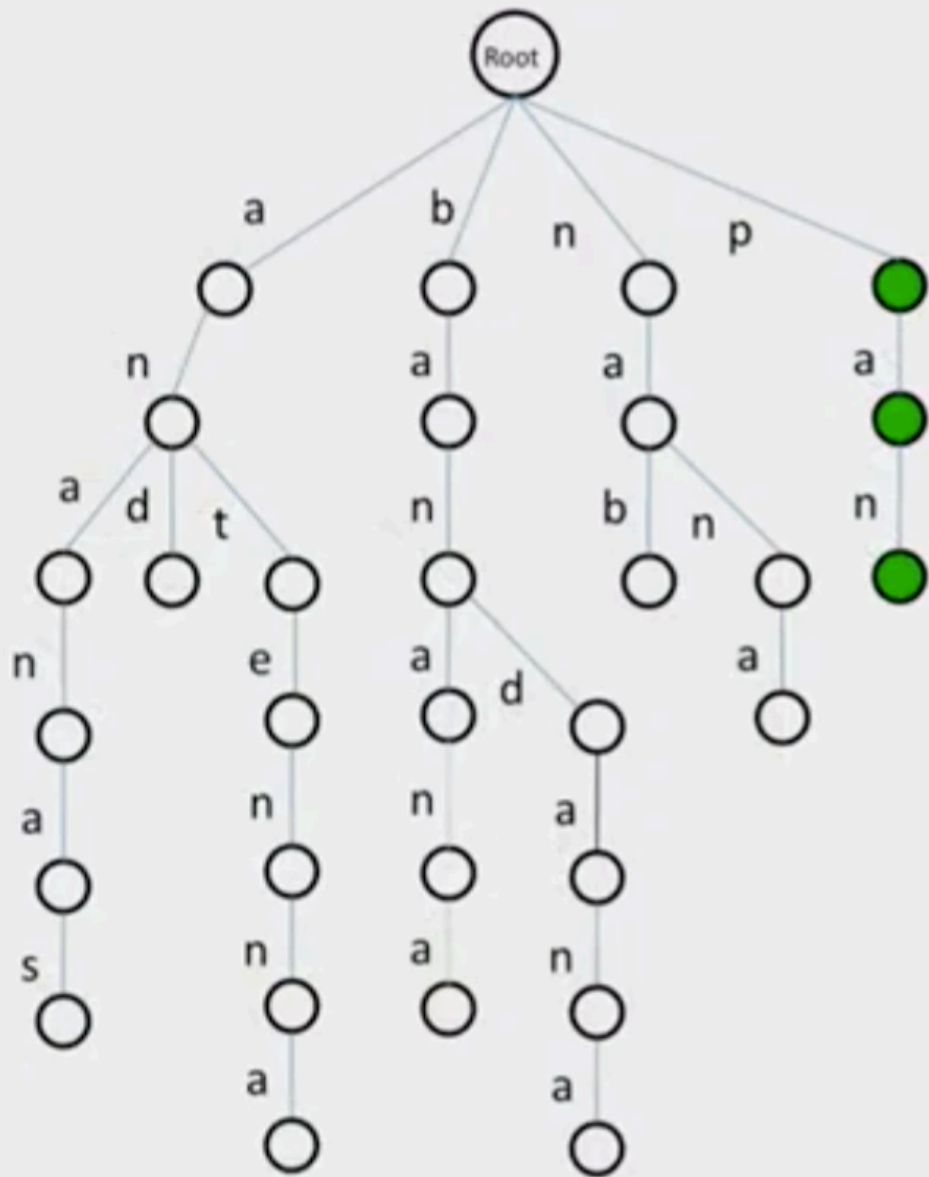




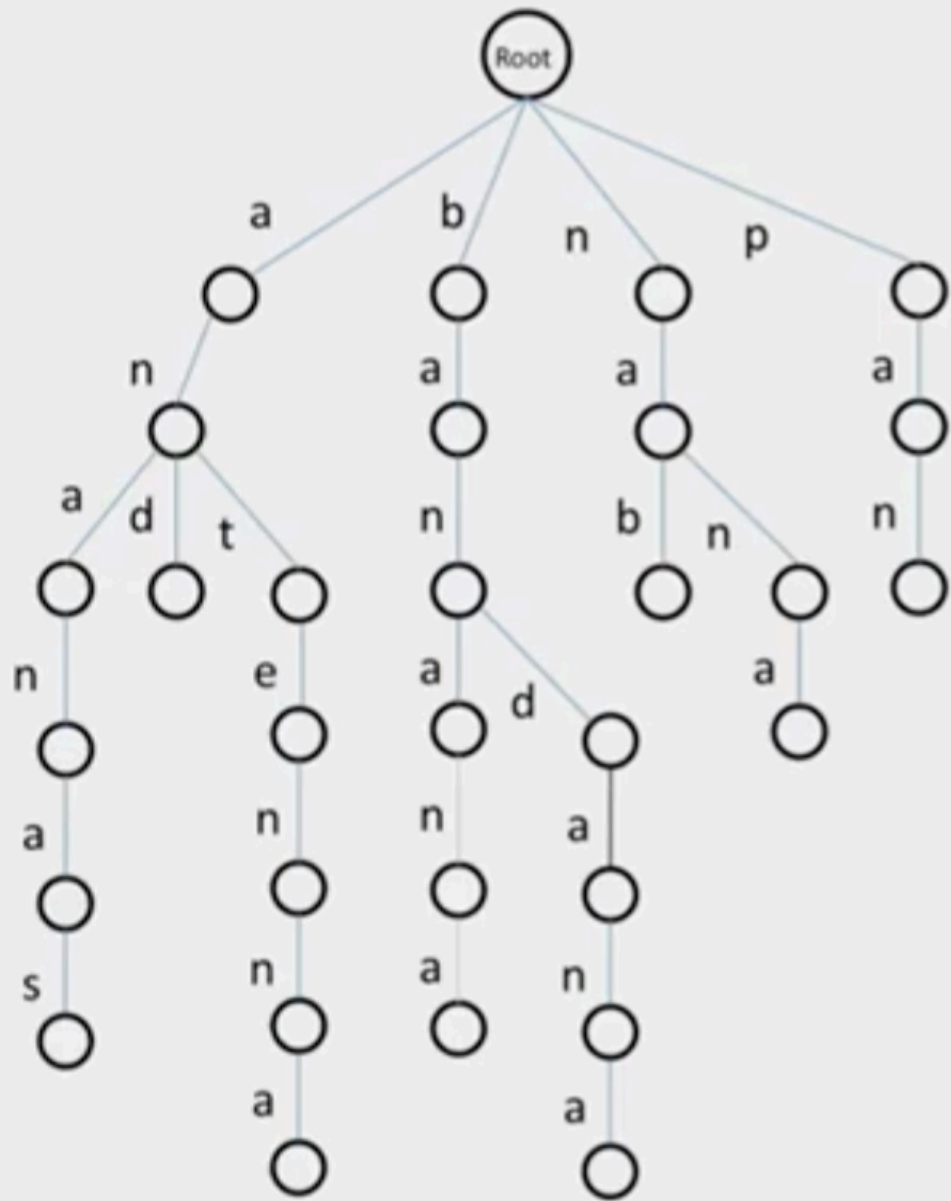
panamabanas



panamabananas



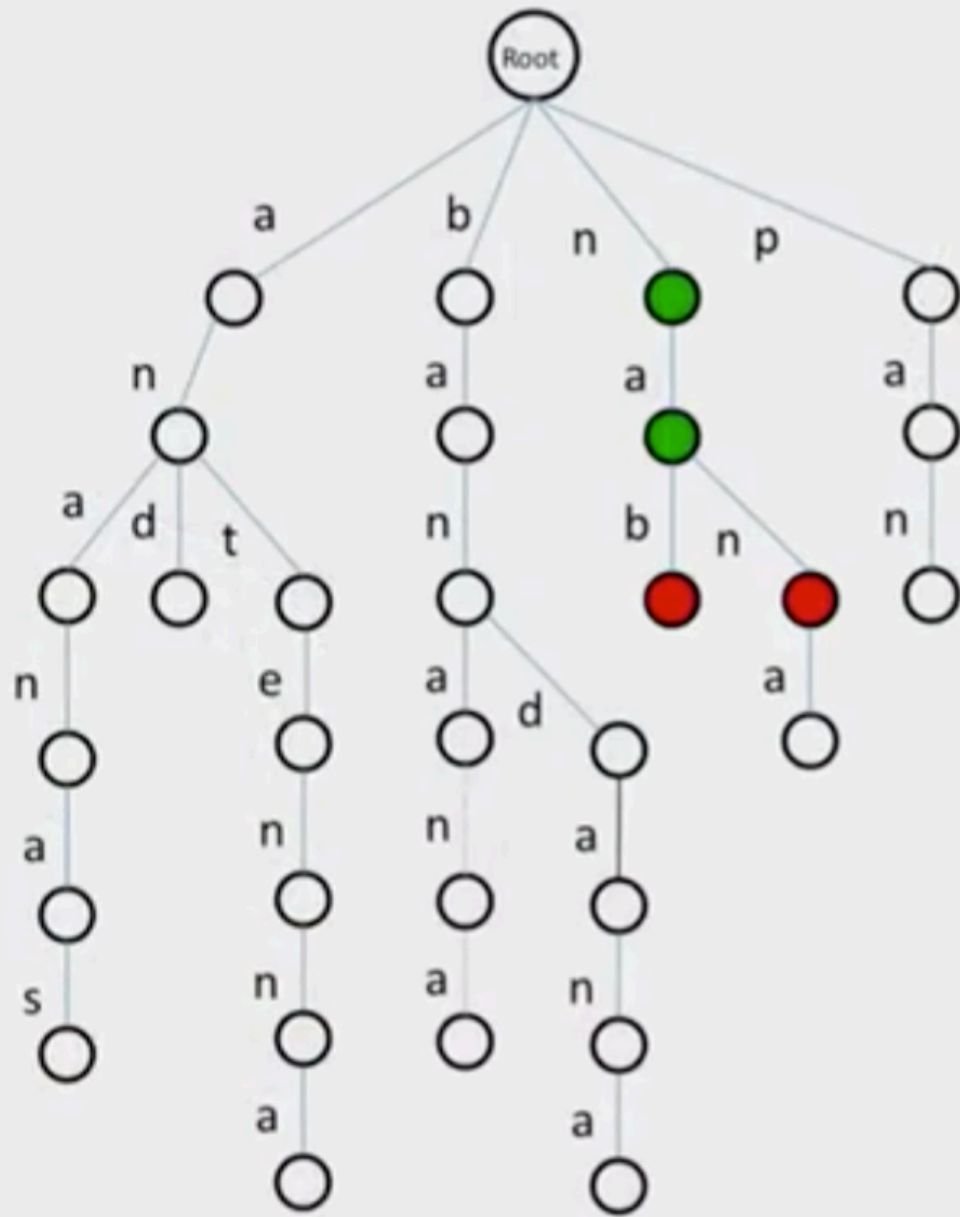
panamabananas



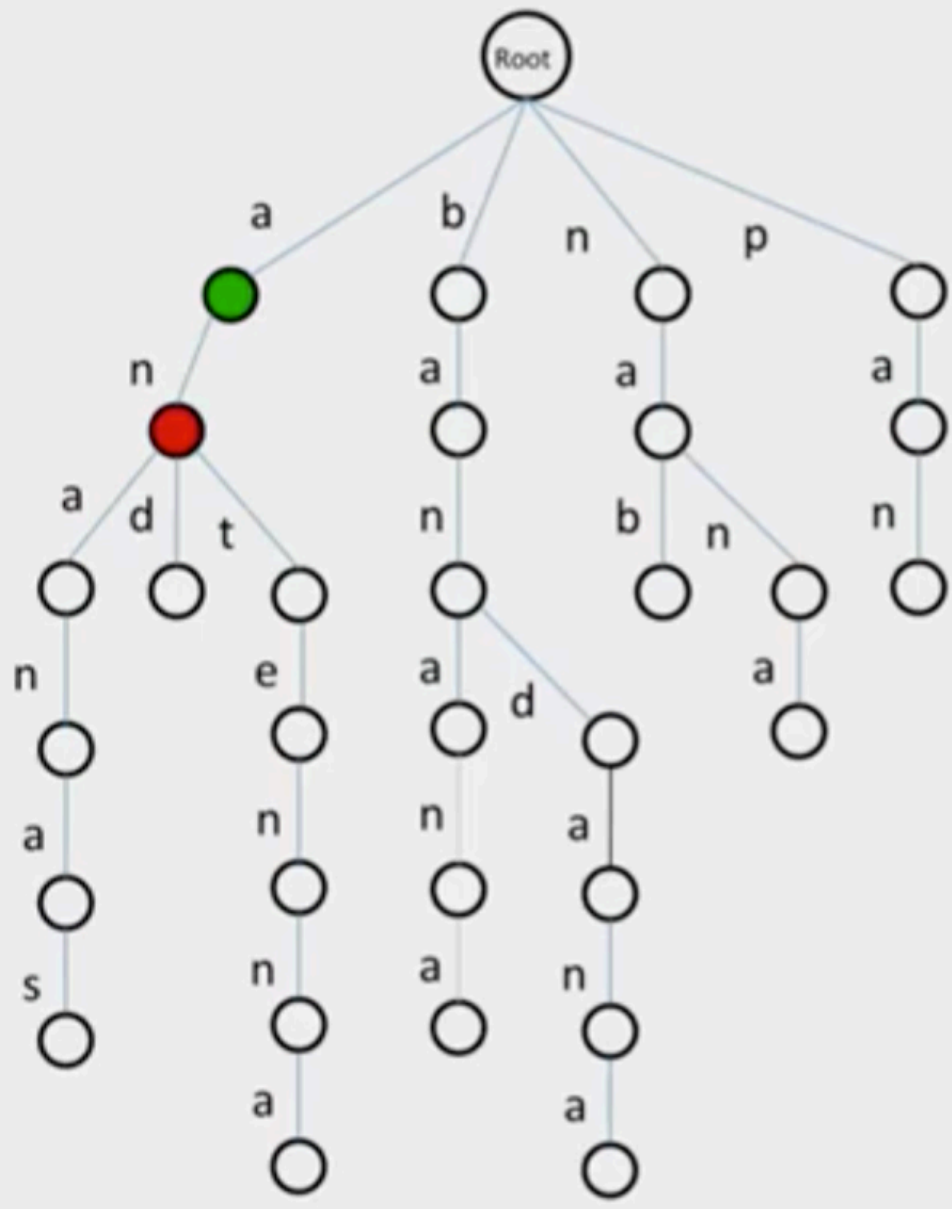




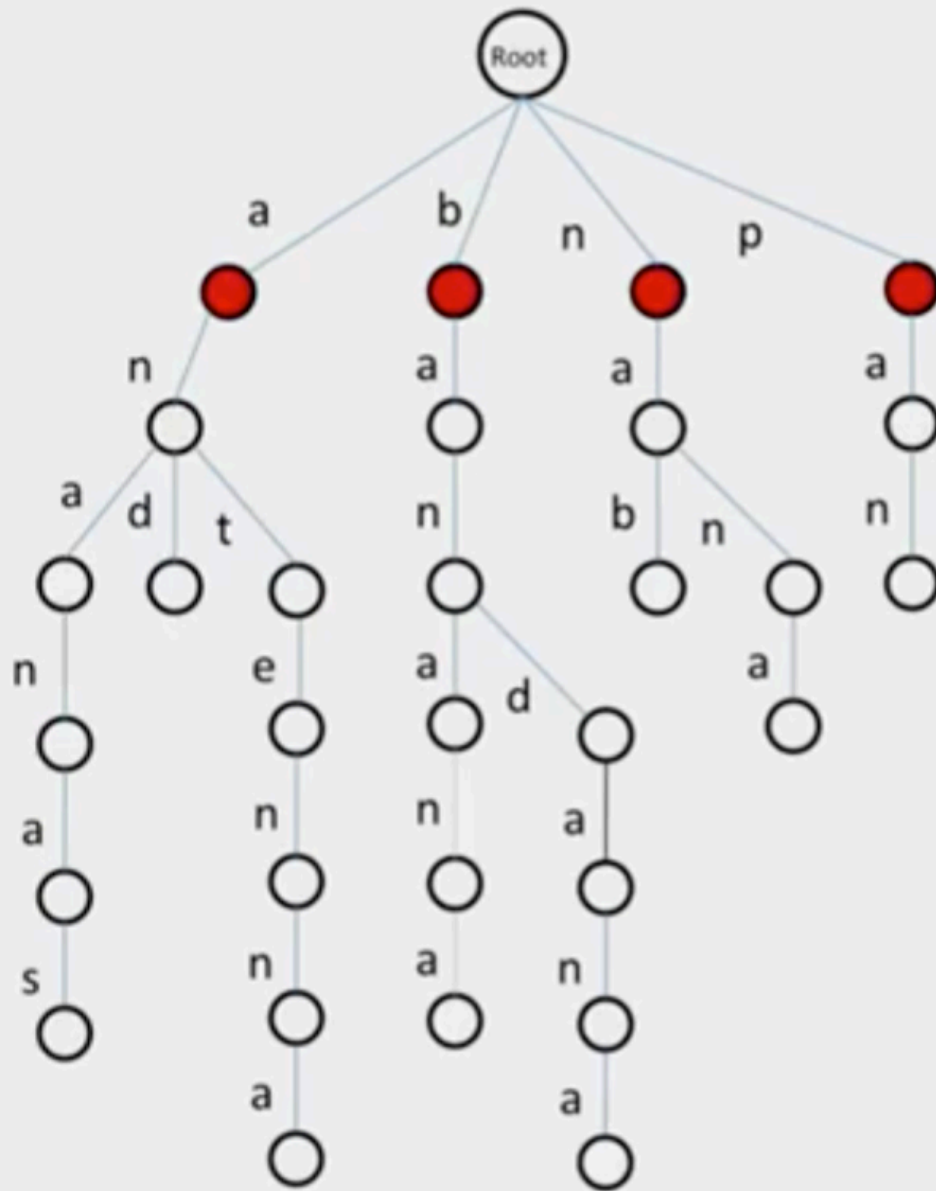
pa**n**a**m**abanas



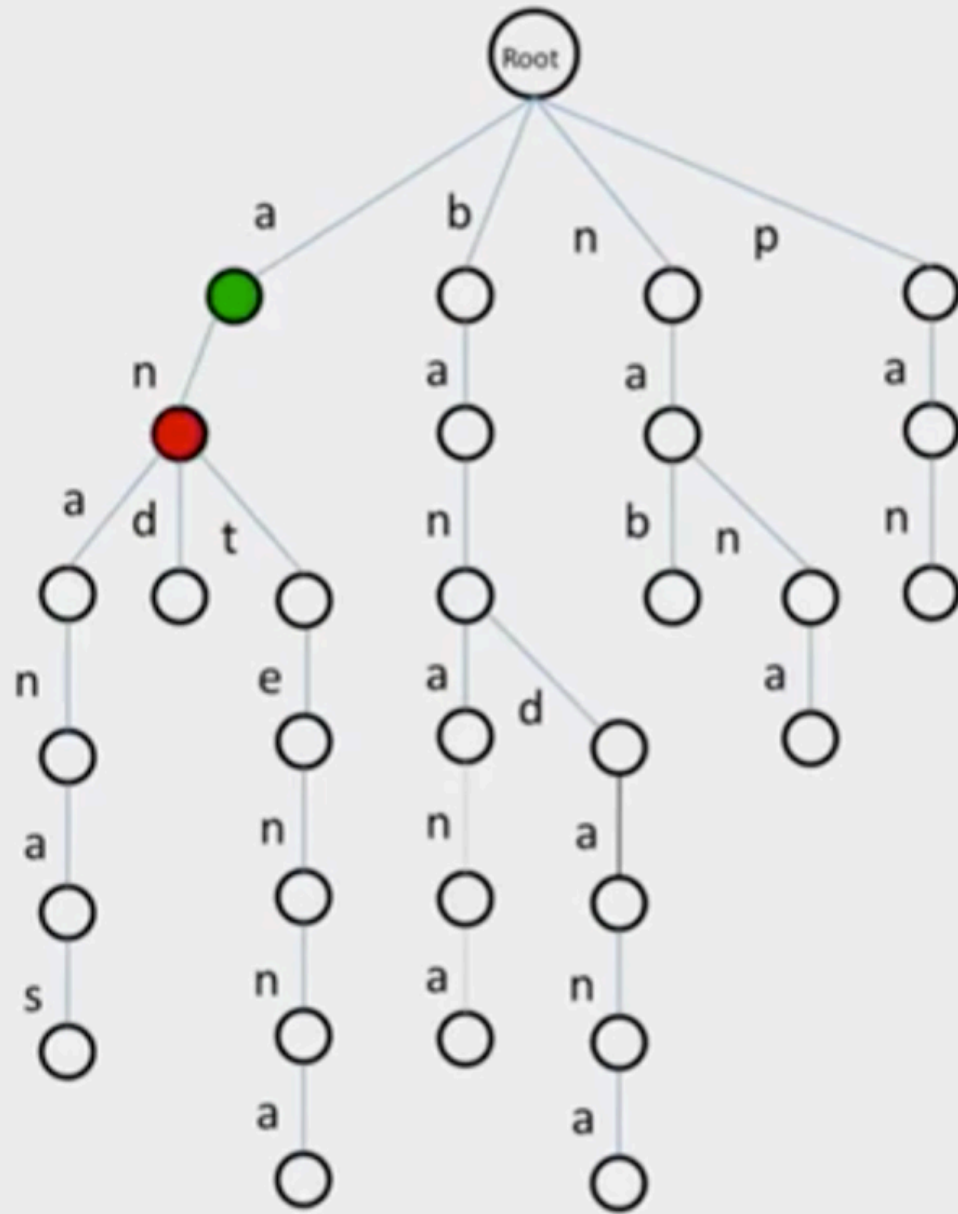
panamabananas



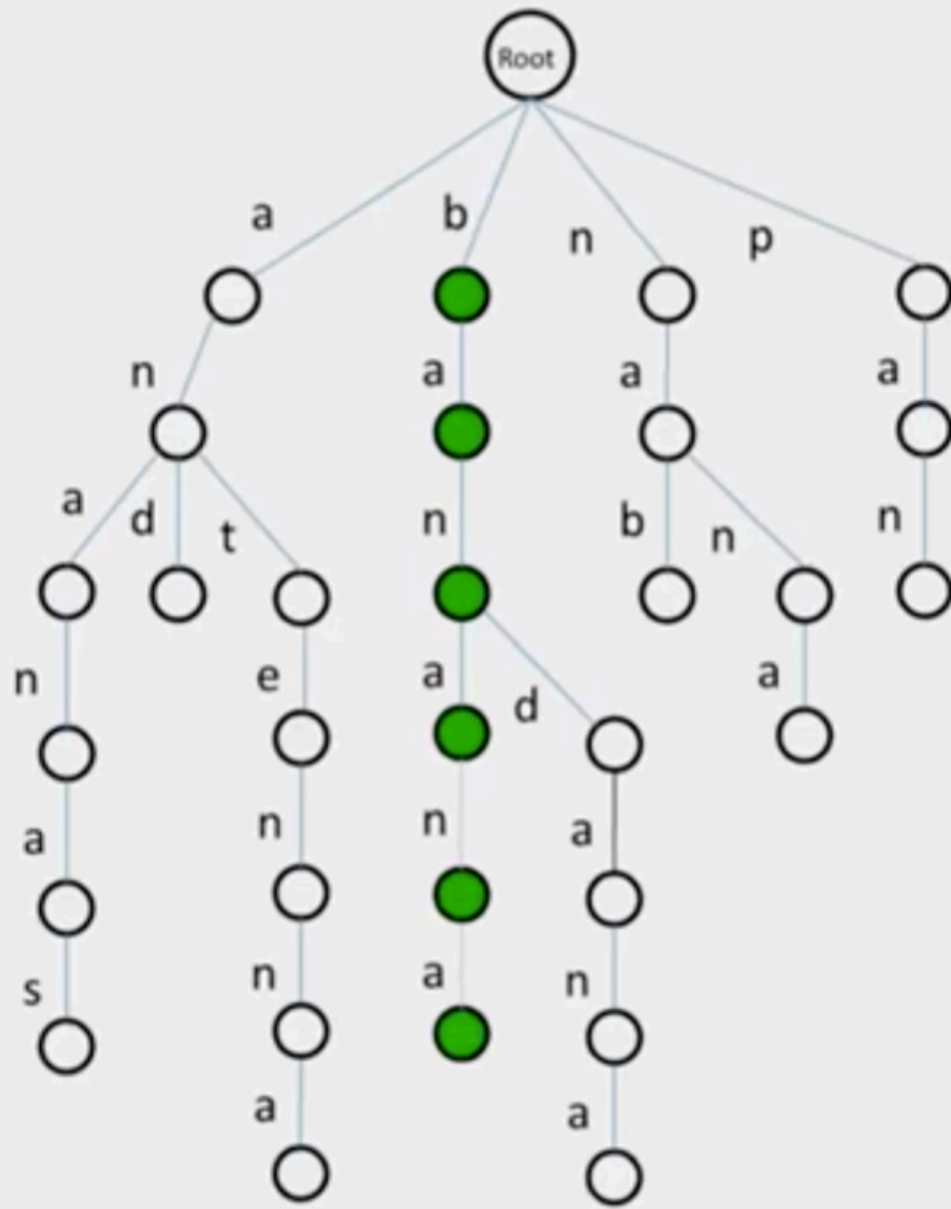
pana**m**abanas



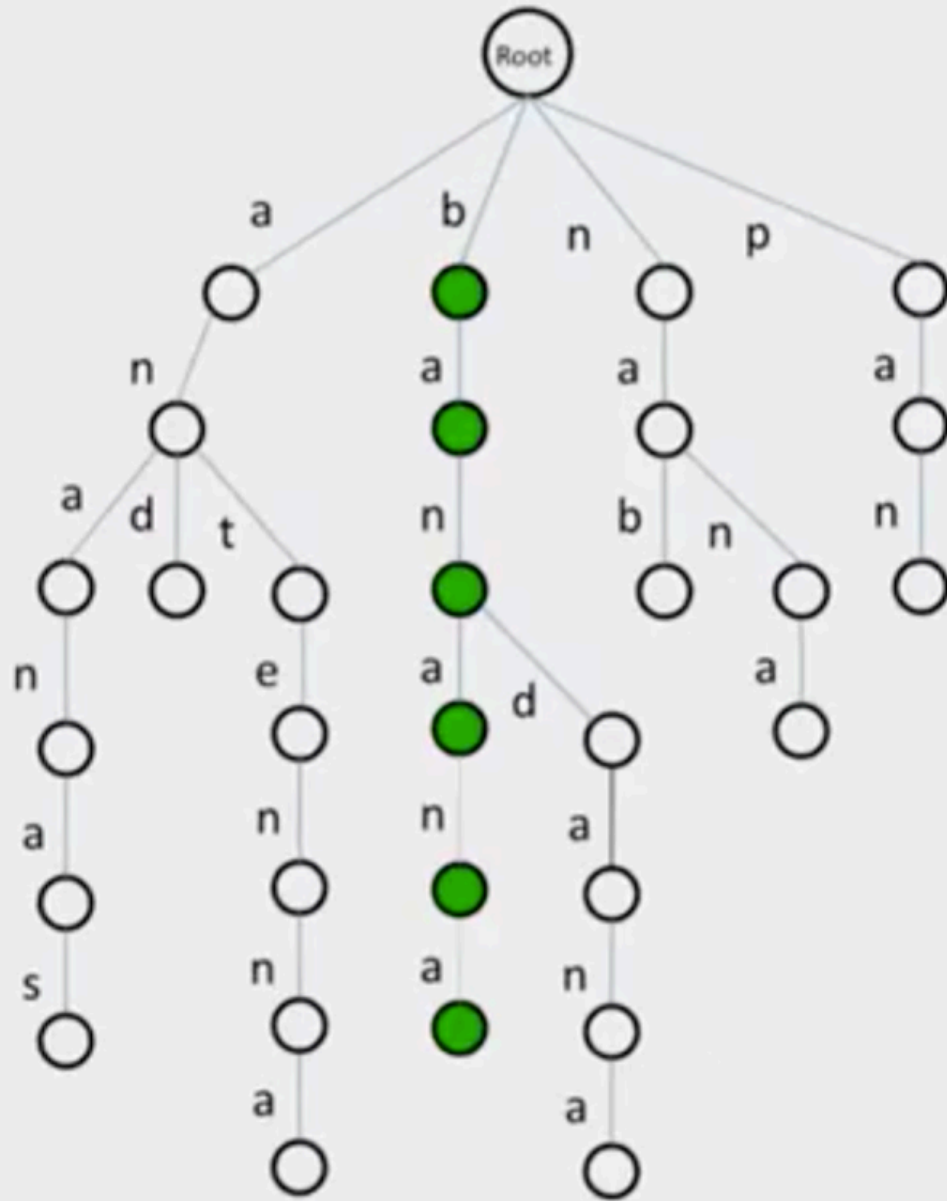
panam**a**b**b**ananas



panama **bananas**



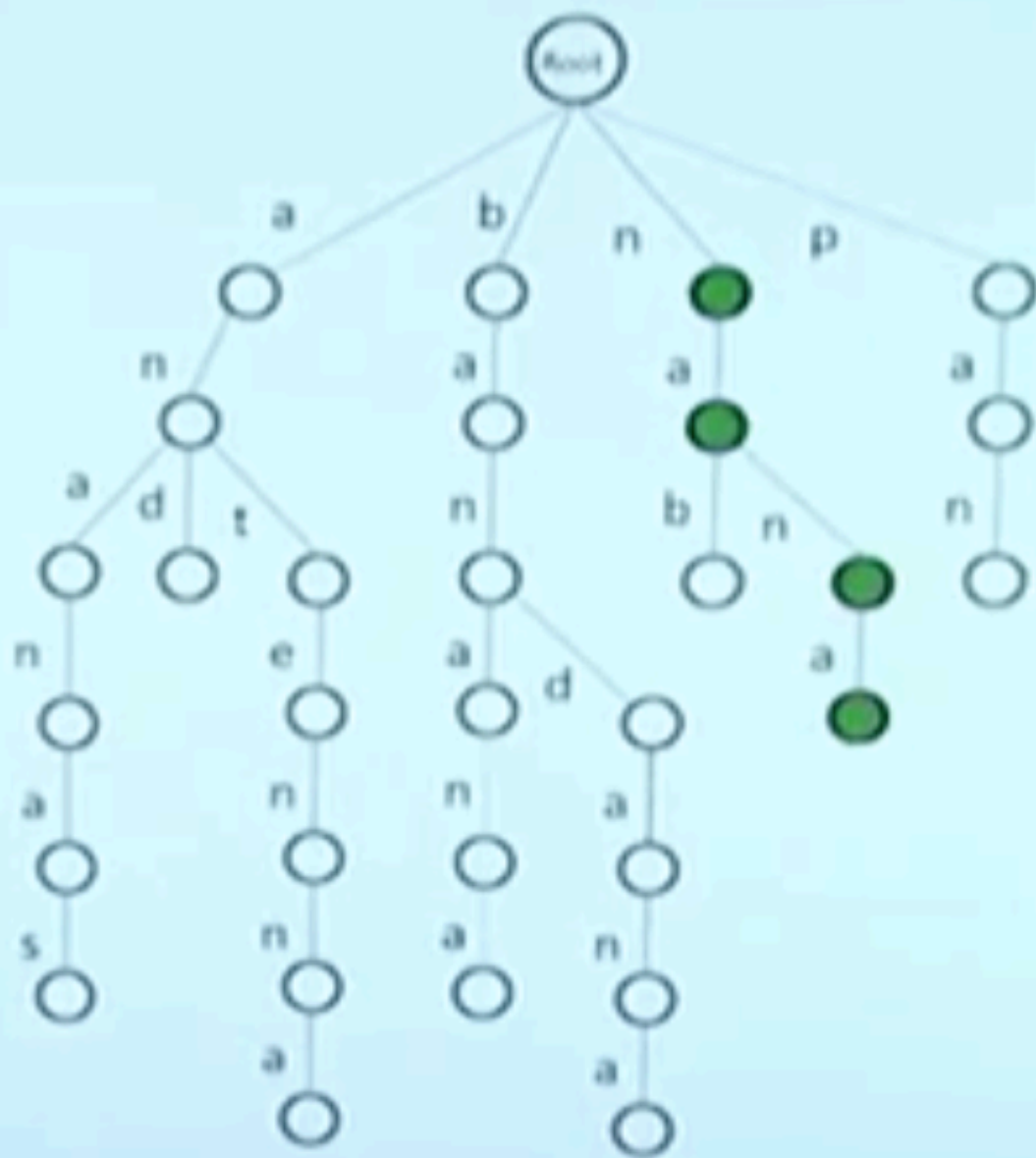
panama **bananas**





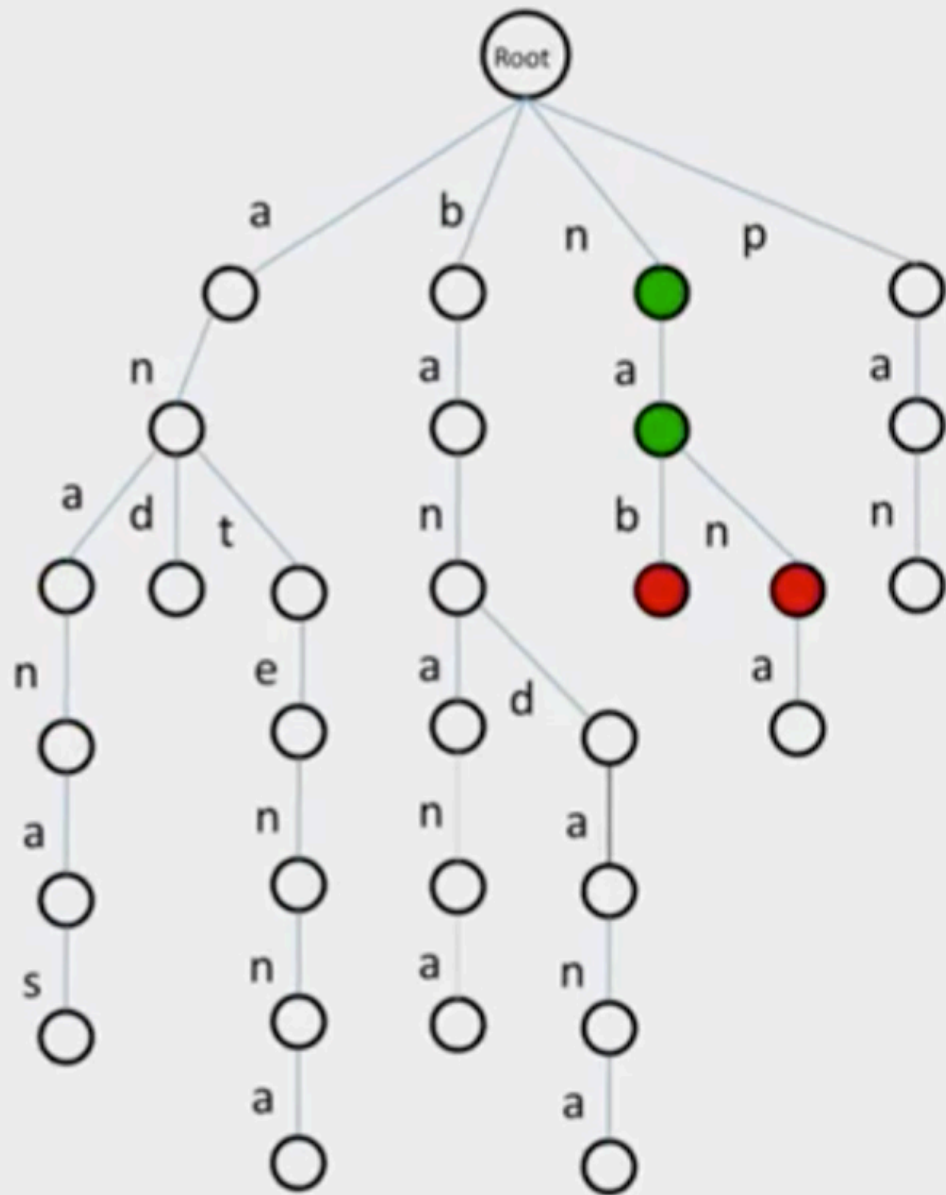


panamabananas

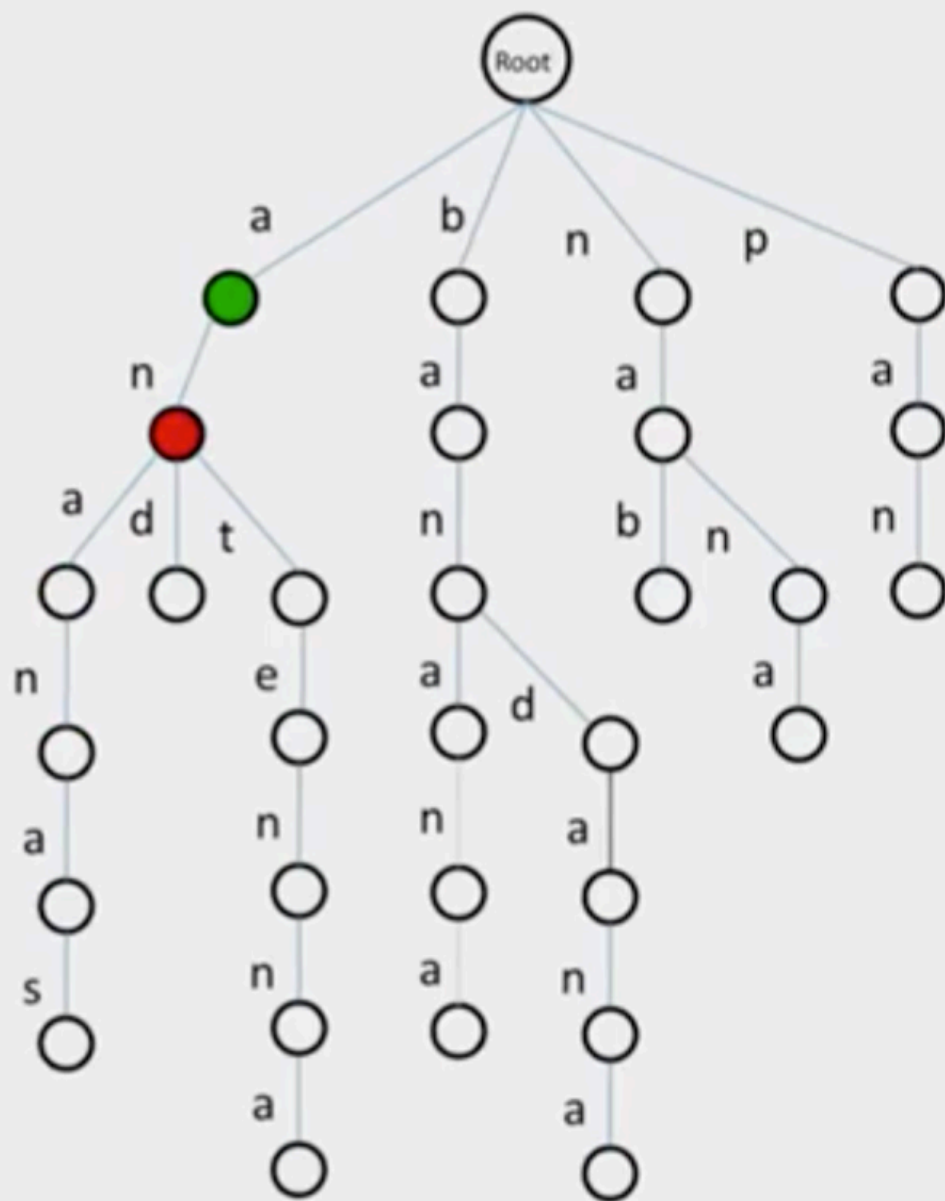




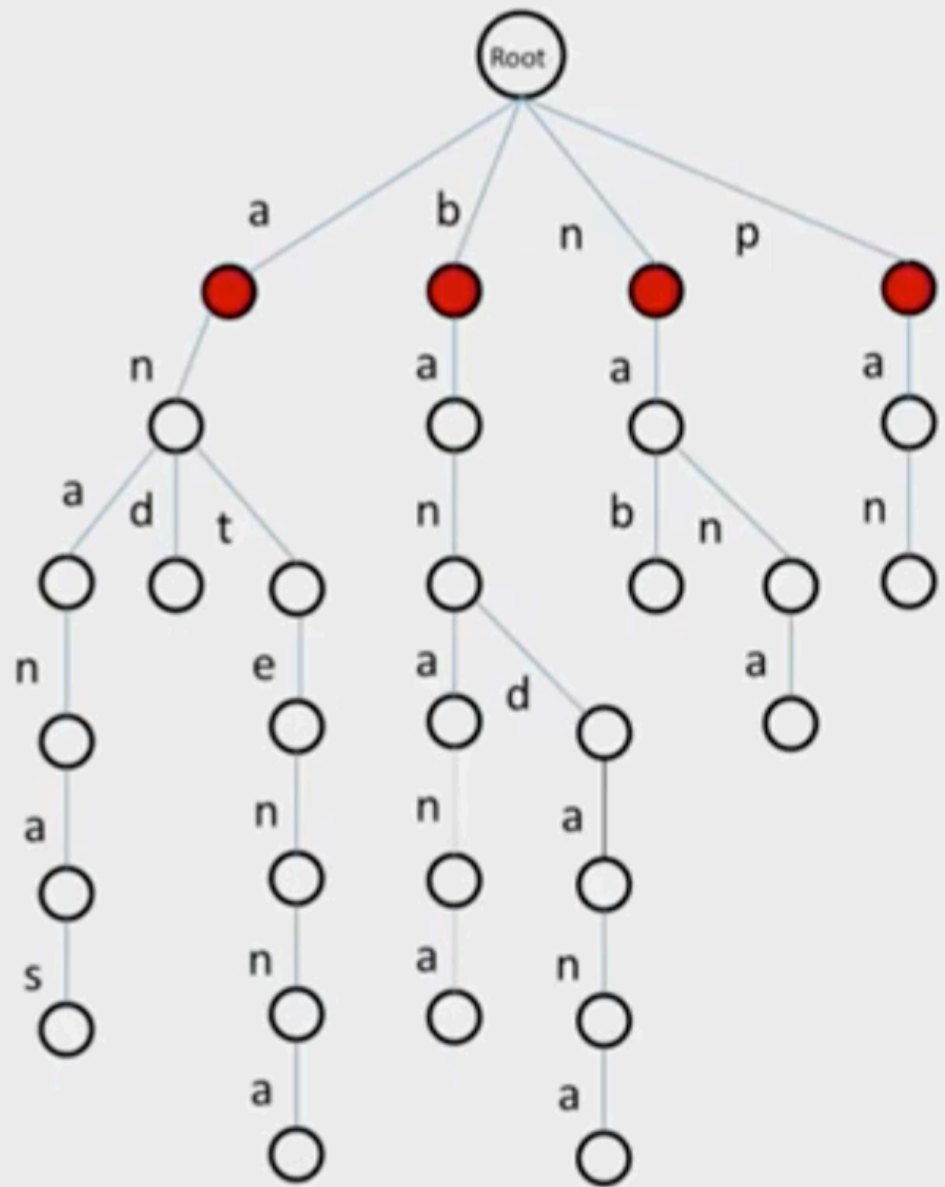
panamabanas



panamabanas



panamabanas



# Success!

- Runtime of Brute Force:
  - Total:  $O(|Genome| * |Patterns|)$
- Runtime of Trie Matching:
  - Trie Construction:  $O(|Patterns|)$
  - Pattern Matching:  $O(|Genome| * |LongestPattern|)$

# Bowtie

Bowtie is based on **Burrows Wheeler Transformation** (BWT) paired with a **Compressed Suffix Array** (CSA)

## Suffix trie

- Fact:  $y$  occurs in  $x$  if  $y$  is a prefix of a suffix of  $x$
- We build a trie with all the suffixes of the text  $x$
- Example:  $x = \mathbf{GATTACA}$  we build the tree on:
  - **GATTACA**
  - **ATTACA**
  - **TTACA**
  - **TACA**
  - **ACA**
  - **CA**
  - **A**





---

## Pattern Matching with Suffix Trees

- To find any pattern in a text:
    1. Build suffix tree for text of length  $m$ — $O(m)$  time
    2. Thread the pattern of length  $n$  through the suffix tree of the text— $O(n)$  time.
  - Therefore the runtime of the Pattern Matching Problem is only  $O(m + n)$ !
-

---

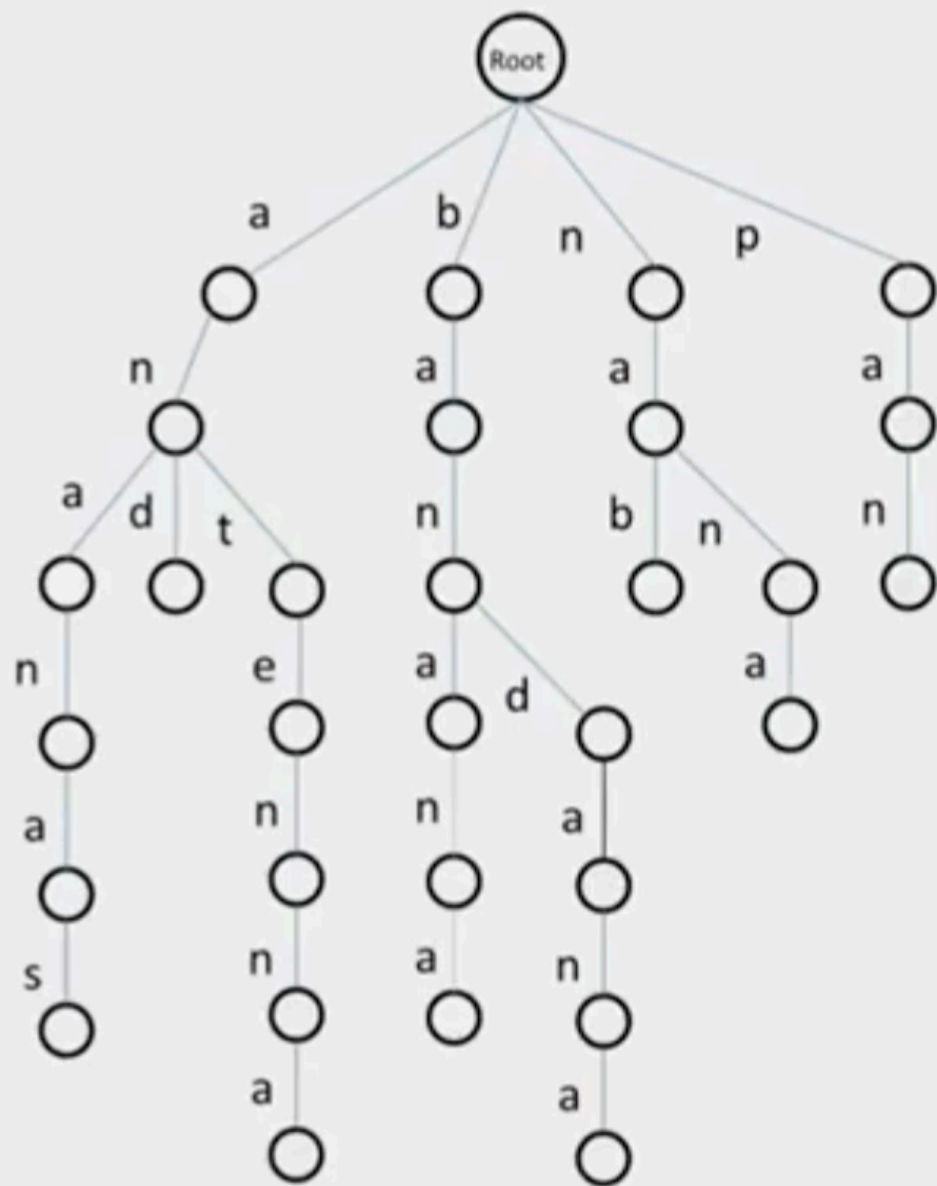
## Pattern Matching with Suffix Trees

### SuffixTreePatternMatching( $p, t$ )

1. Build **suffix tree** for text  $t$
  2. Thread pattern  $p$  through **suffix tree**
  3. **if** threading is complete
  4.     **output** positions of all  $p$ -matching leaves in the tree
  5. **else**
  6.     **output** "Pattern does not appear in text"
-

# Memory Analysis of TrieMatching

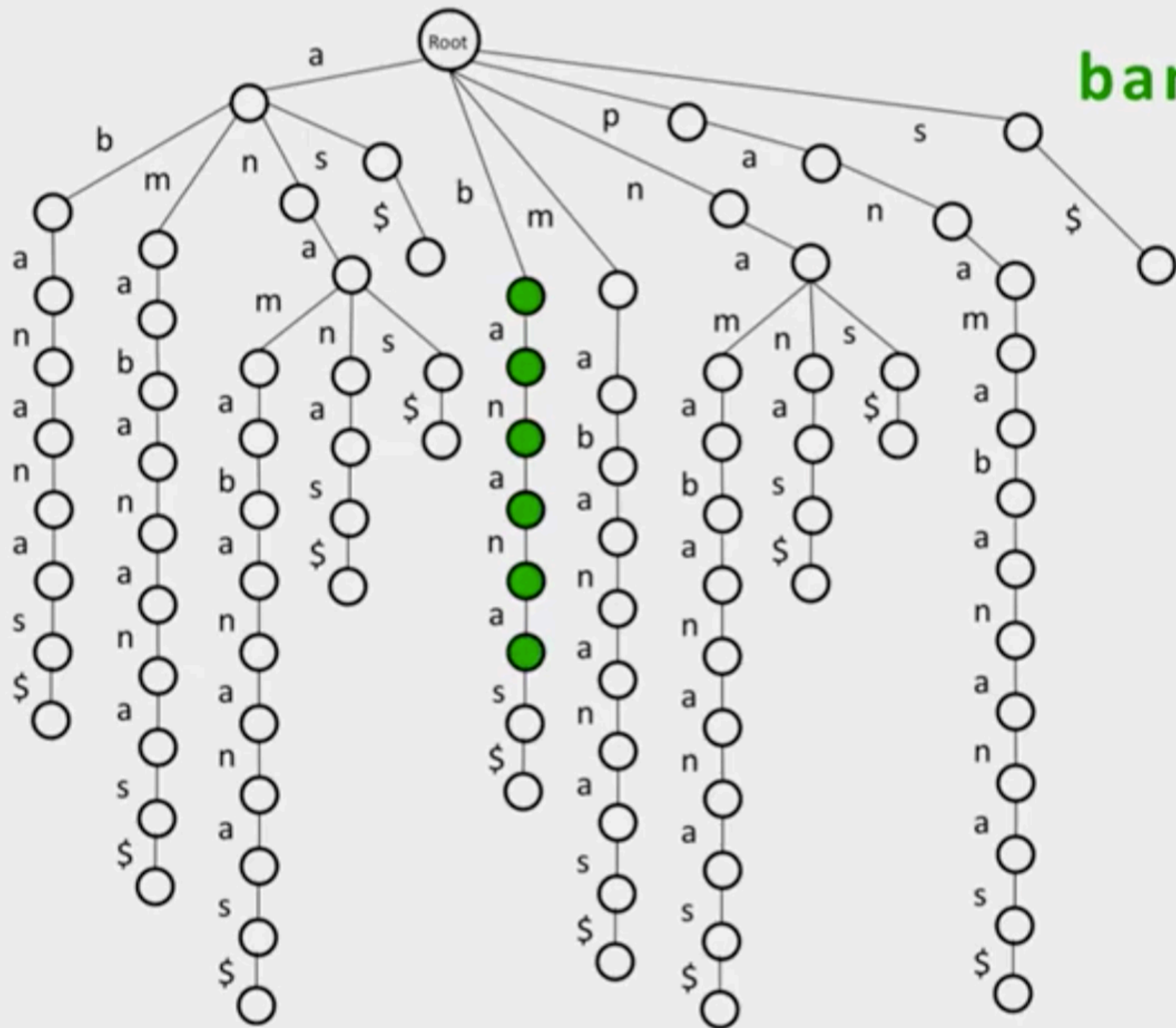
- Son completely forgot about memory!
- Worst case: # edges =  $O(|Patterns|)$



# Preprocessing the Genome

- Form all suffixes of *Genome*.
- How can we combine these suffixes into a data structure?
- Let's use a trie!

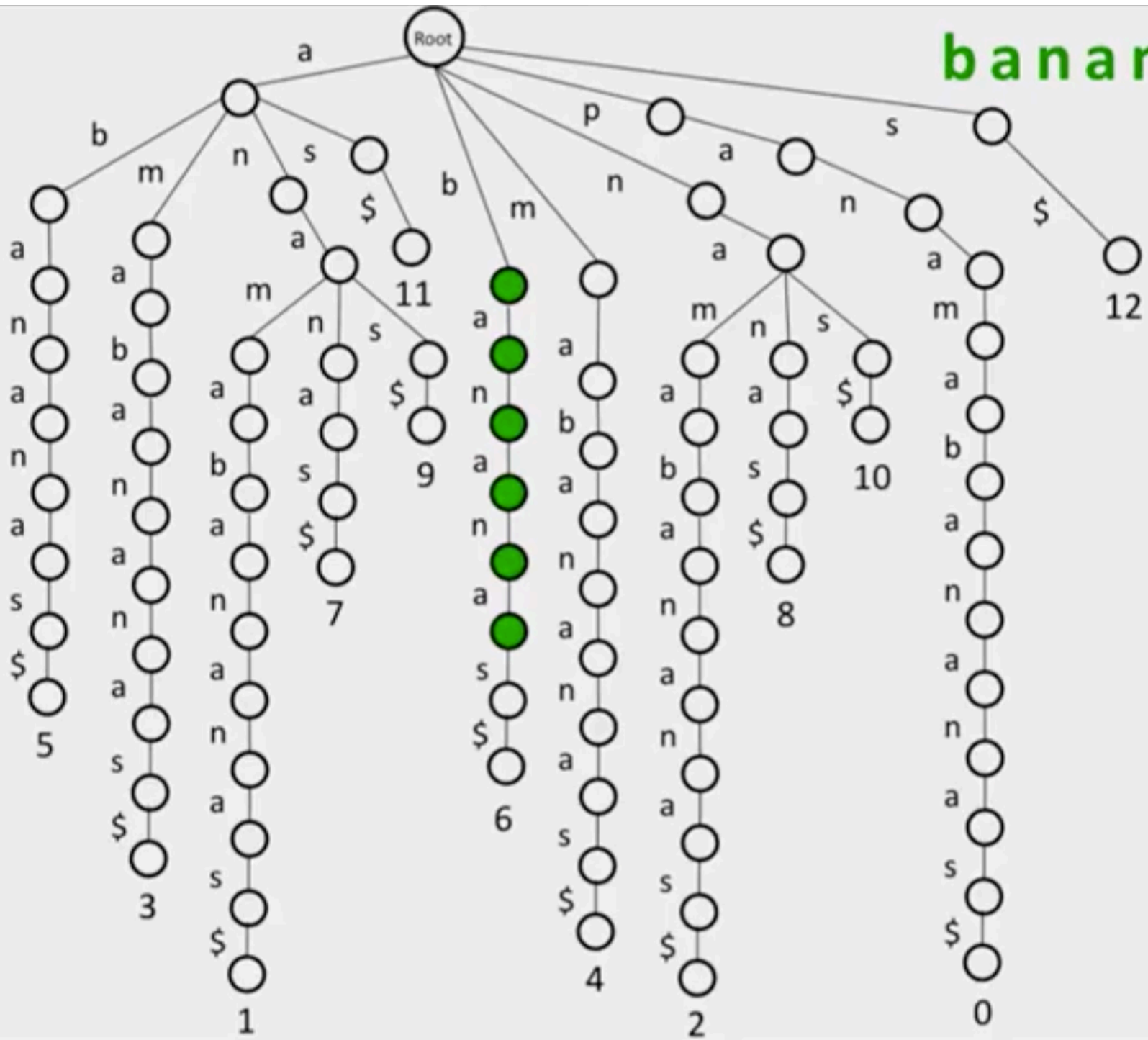
banana



## Where Are the Matches?

- To find where the pattern matches are, we need to add a little more information to the suffix trie.
- At each leaf (\$), we add the starting position in *Genome* of the suffix ending at that leaf.

banana







# Memory Trouble Once Again

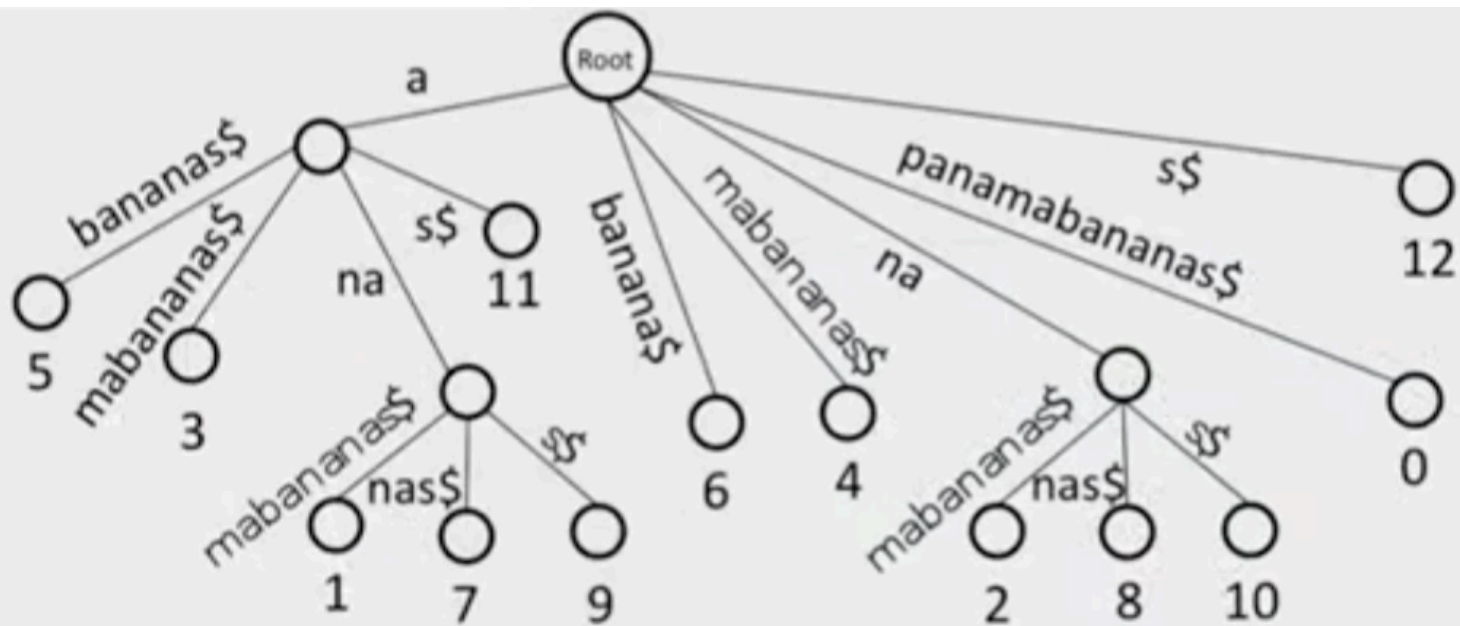
- Worst case: the suffix trie holds  $O(|\textit{Suffixes}|)$  nodes.
- For a *Genome* of length  $n$ ,  
 $|\textit{Suffixes}| = n(n - 1)/2 = O(n^2)$

## *Suffixes*

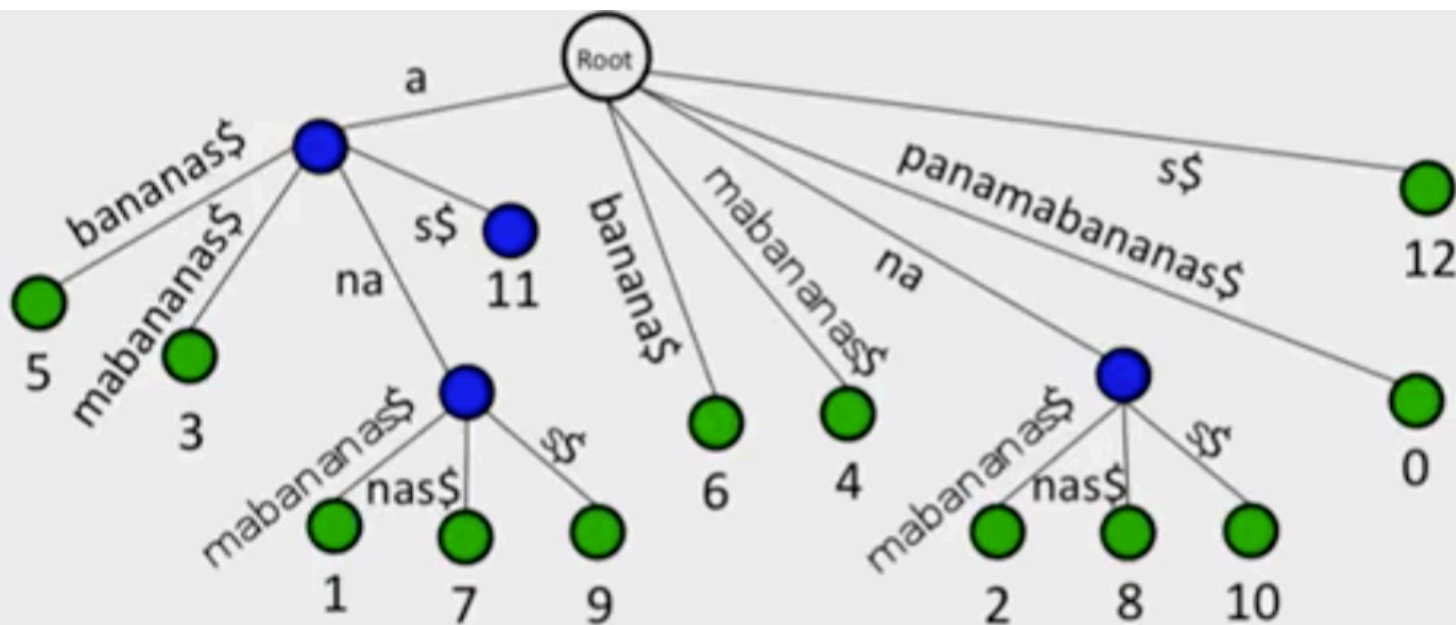
panamabananas\$  
anamabananas\$  
namabananas\$  
amabananas\$  
mabananas\$  
abananas\$  
bananas\$  
ananas\$  
nanas\$  
anas\$  
nas\$  
as\$  
s\$  
\$

## Compressing the Trie

- This doesn't mean that our idea was bad!
- To reduce memory, we can compress each “nonbranching path” of the tree into an edge.



- This data structure is called a **suffix tree**.

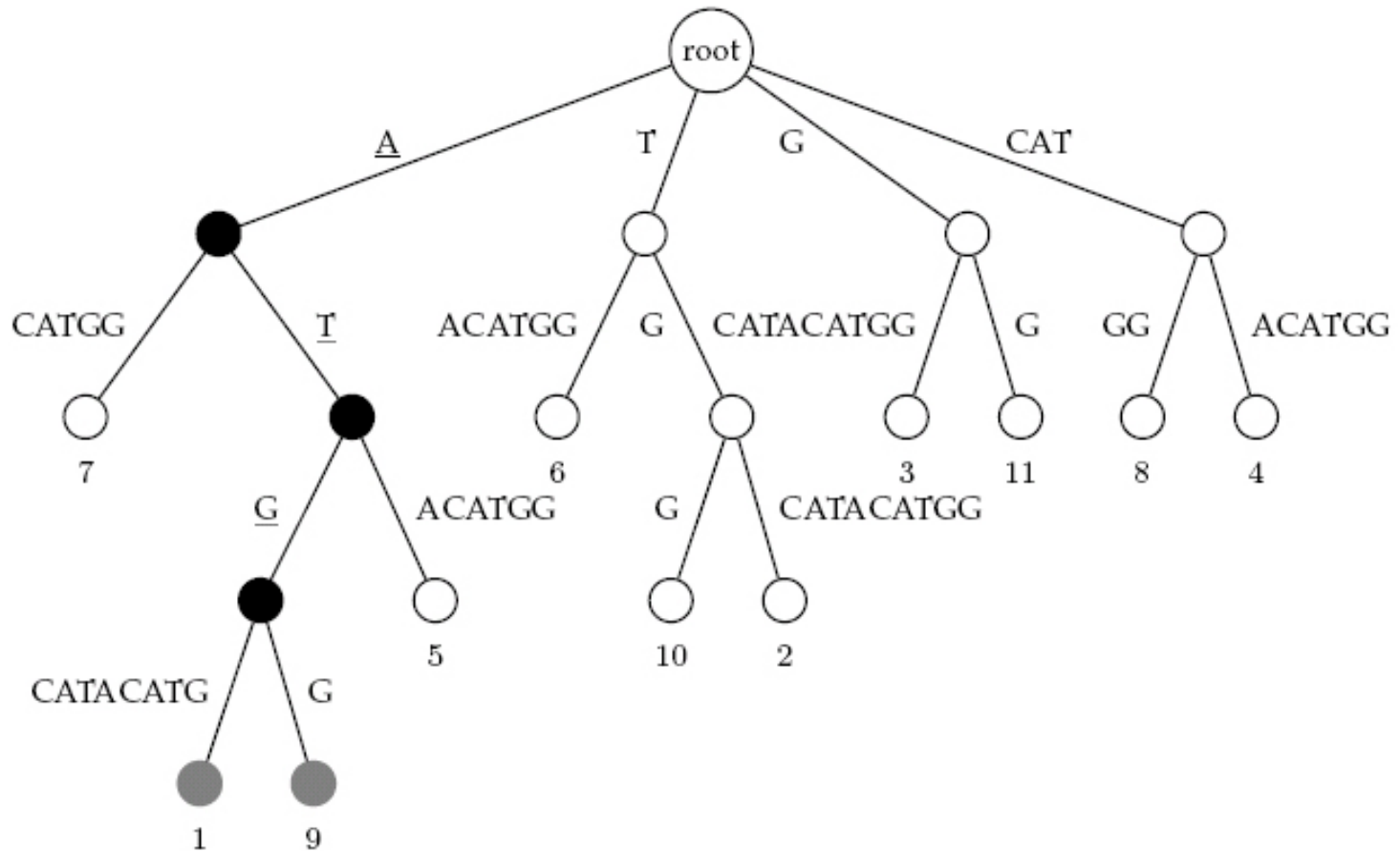


- This data structure is called a **suffix tree**.
- For any *Genome*, # nodes  $< 2 |Genome|$ .
  - # leaves =  $|Genome|$ ;
  - # internal nodes  $< |Genome| - 1$

# We are Not Finished Yet

- I am happy with the suffix tree, but I am not completely satisfied.
  - Runtime:  $O(|Genome| + |Patterns|)$
  - Memory:  $O(|Genome|)$
- However, big-O notation ignores constants!
  - The best known suffix tree implementations require  $\sim 20$  times the length of  $|Genome|$ .
  - Can we reduce this constant factor?

# Threading ATG through a Suffix Tree



**Suffix Trees.** We can reduce the number of edges in the **suffix trie** by combining the edges on any non-branching path into a single edge. The resulting data structure is called *suffix tree*.

