# Alignment of Short Reads: Suffix Trees

## Ahmet Sacan

# The Newspaper Problem

# The Newspaper Problem as an Overlapping Puzzle

hoodie, appr... lie, appr... 2

...e have not yet named ...yet named any suspects, ...lt

...mation is wel... is welc... ...e ca

# Whole-Genome Sequencing



Hierarchical shotgun sequencing

Genomic DNA

BAC library

Organized mapped large clone contigs

BAC to be sequenced

Shotgun clones

Shotgun sequence   . . . ACCGTAAATGGGCTGATCATGCTTAAA
                              TGATCATGCTTAAACCCTGTGCATCCTACTG . . .

Assembly   . . . ACCGTAAATGGGCTGATCATGCTTAAACCCTGTGCATCCTACTG . . .

cow
2009

horse
2007

opossum
2007

macaque
2006

dog
2005

chimpanzee
2005

rat
2004

mouse
2002

human
2001

# From Reference Genome to Personal Genomes

- **Reference human genome** assembled in 2000.

CTGAGGATGGACTACGCTACTACTGATAGCTGTTT    reference genome

GAGGA             C**C**ACG             TGA–AG

CTGA       GGAC**C**        ACTAC    A–AGCT     reads

     GATGG     ACGCT               TGTTT

CTGAGGATGGAC**C**ACGCTACTACTGA–AGCTGTTT    personal genome

# Reference-based sequencing

- Map short reads to the reference genome.
- Exact "pattern matching"
  - Locate the occurrence of short read exactly.
- Approximate "pattern matching"
  - Allow mismatches, insertions, deletions.

# Searching short read (pattern) in genome (text)

*Pattern* drives along *Text*

```
p a n a m a b a n a n a s        Text
n a n a                          Pattern
```

# Searching short read (pattern) in genome (text)

*Pattern* drives along *Text*

p a n a m a b a n a n a s          *Text*
n a n a                            *Pattern*

# Searching for multiple patterns

Genome

CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTAC**C**ACATCGTAGCTAC

Pattern₁

Pattern₂

Pattern₃

Pattern₄

Pattern₅

# Searching for multiple patterns

# Computational Complexity

- How long does it take to find a single pattern?
  $O( |text| * |pattern| )$
- Multiple patterns?
  $O( |pattern1| * |text| ) + O( |pattern2| * |text| ) + ...$

  $= O( |text| * |patterns| )$
- Human genome:
  - $|Text| \approx 10^9$
  - $|Patterns| \approx 10^{12}$

# Pack Patterns onto a bus

CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTAC**C**ACATCGTAGCTAC

( Root )

## Patterns

banana

pan

nab

antenna

bandana

ananas

nana

( Root )

## Patterns

**banana**

pan

nab

antenna
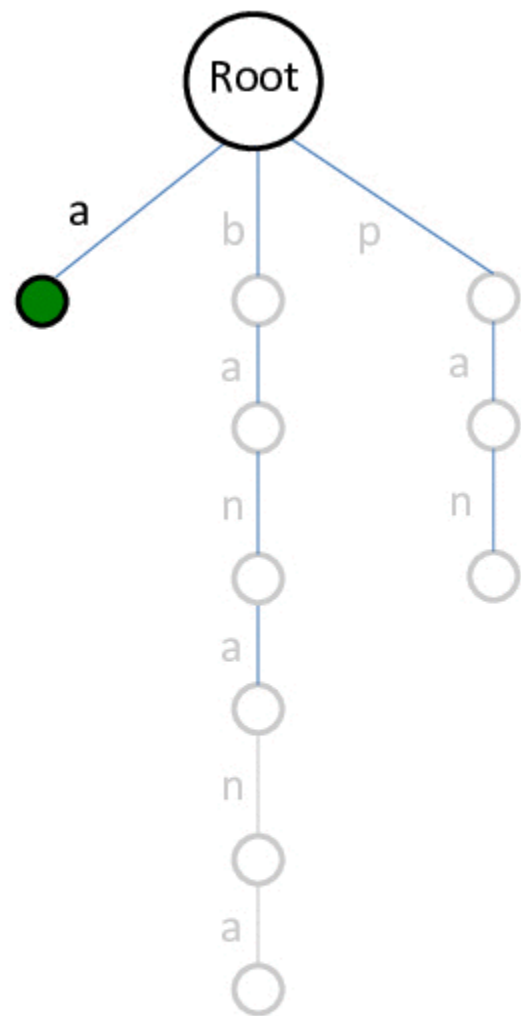
bandana

ananas

nana

## Patterns

banana
**pan**
and
nab
antenna
bandana
ananas
nana

# Patterns

banana
pan
**and**
nab
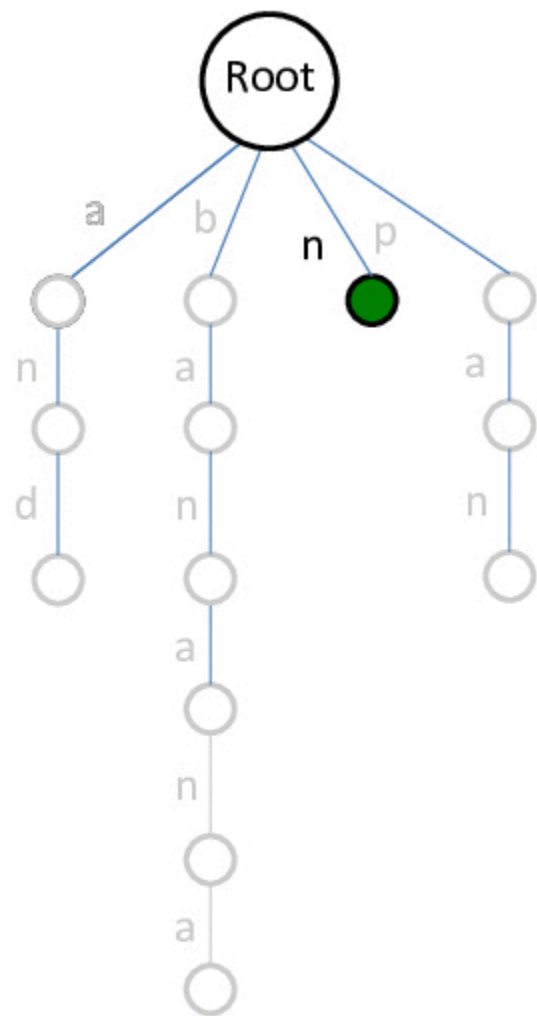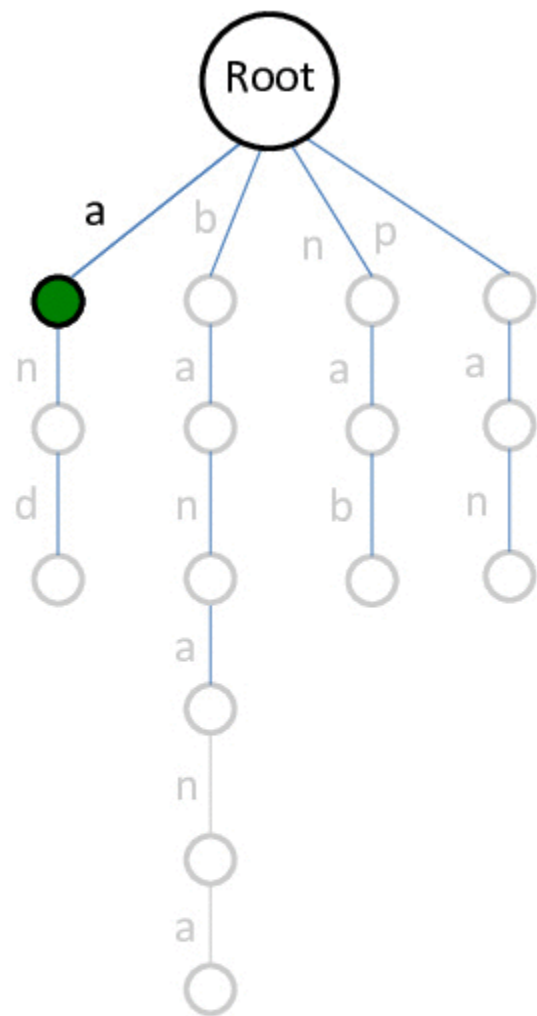antenna
bandana
ananas
nana

**Patterns**

banana
pan
and
**nab**
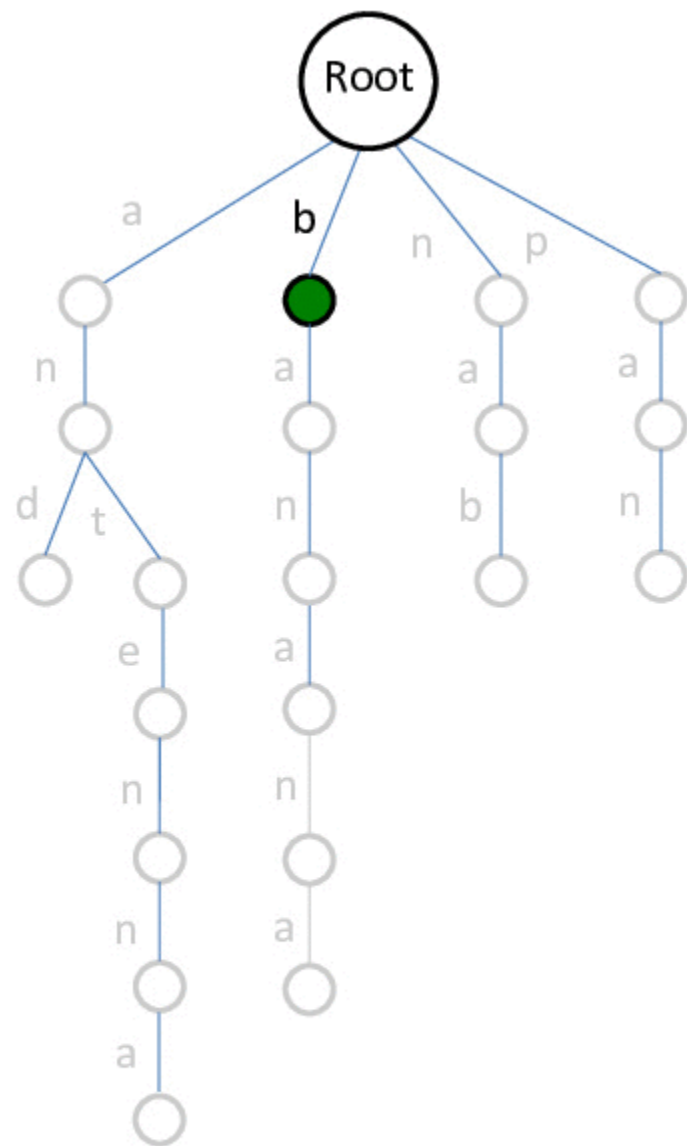antenna
bandana
ananas
nana

## Patterns

banana
pan
and
nab
**antenna**
bandana
ananas
nana

**Patterns**

banana
pan
and
nab
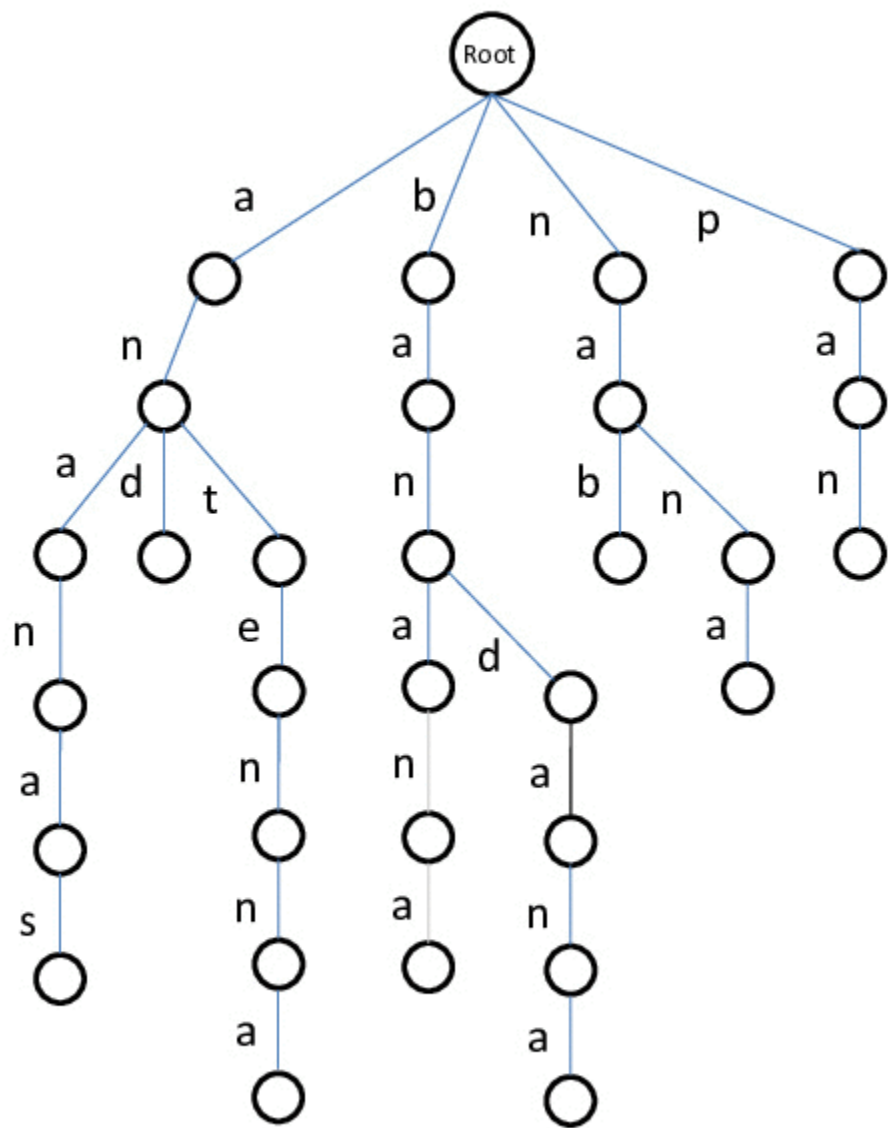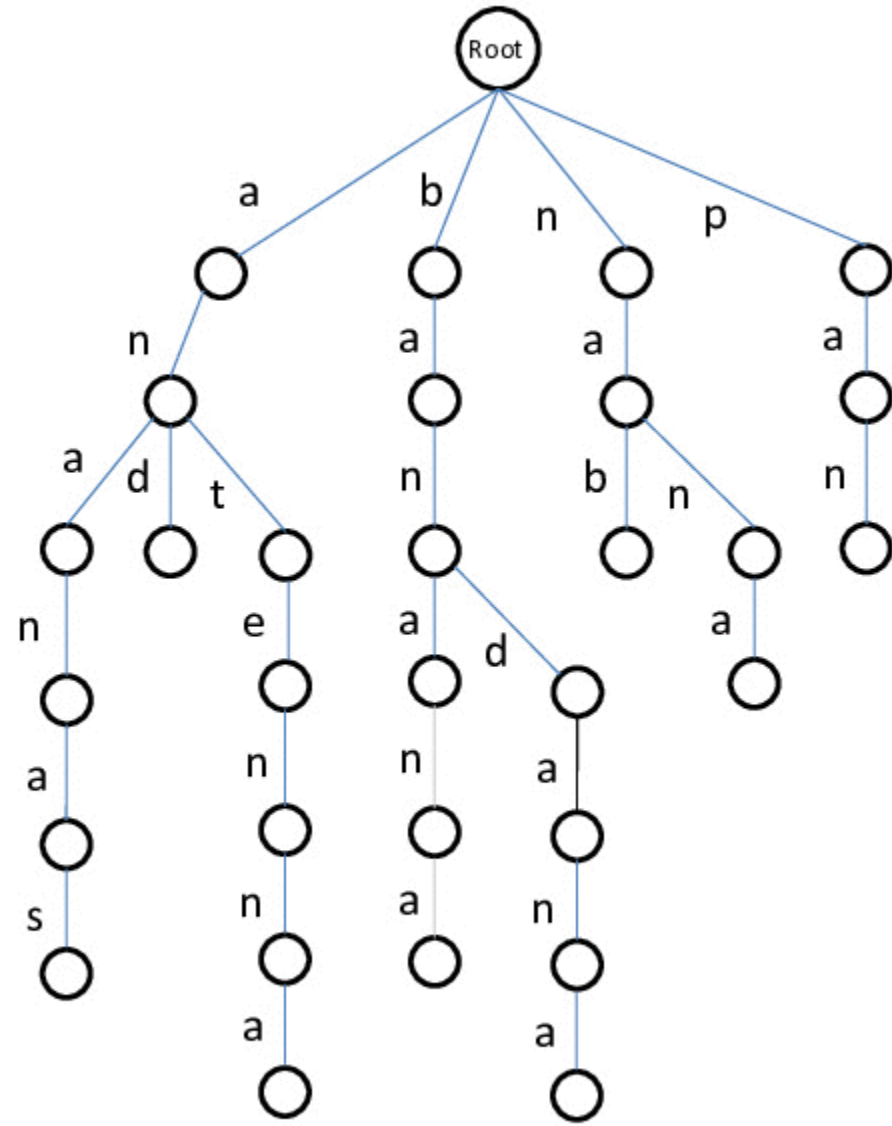antenna
**bandana**
ananas
nana

# p a n a m a b a n a n a s



**TrieMatching**(*Text, Patterns*): drive Trie(*Patterns*) along *Text* at each position of *Text*
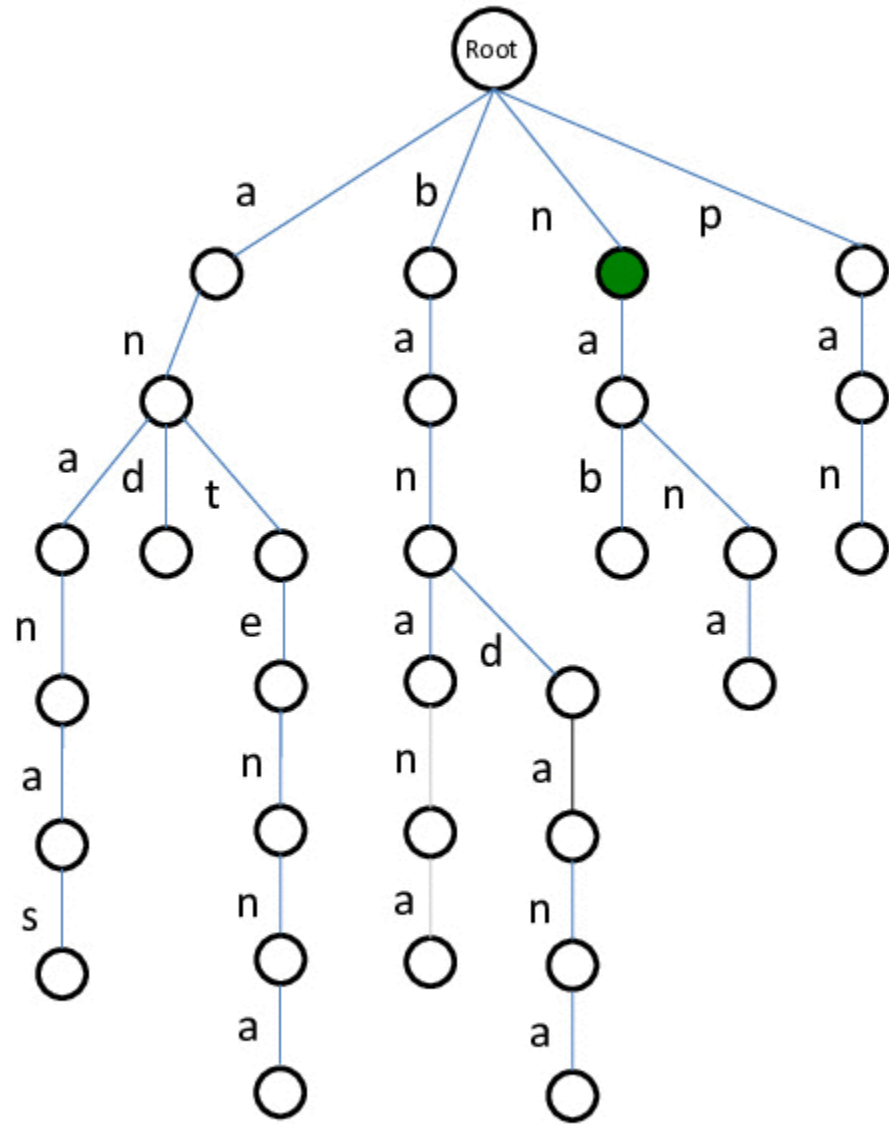
- walk down Trie(*Patterns*) by spelling symbols of *Text*

- a pattern from *Patterns* matches *Text* each time you reach a leaf!

For simplicity, we assume that no pattern is a substring of another pattern
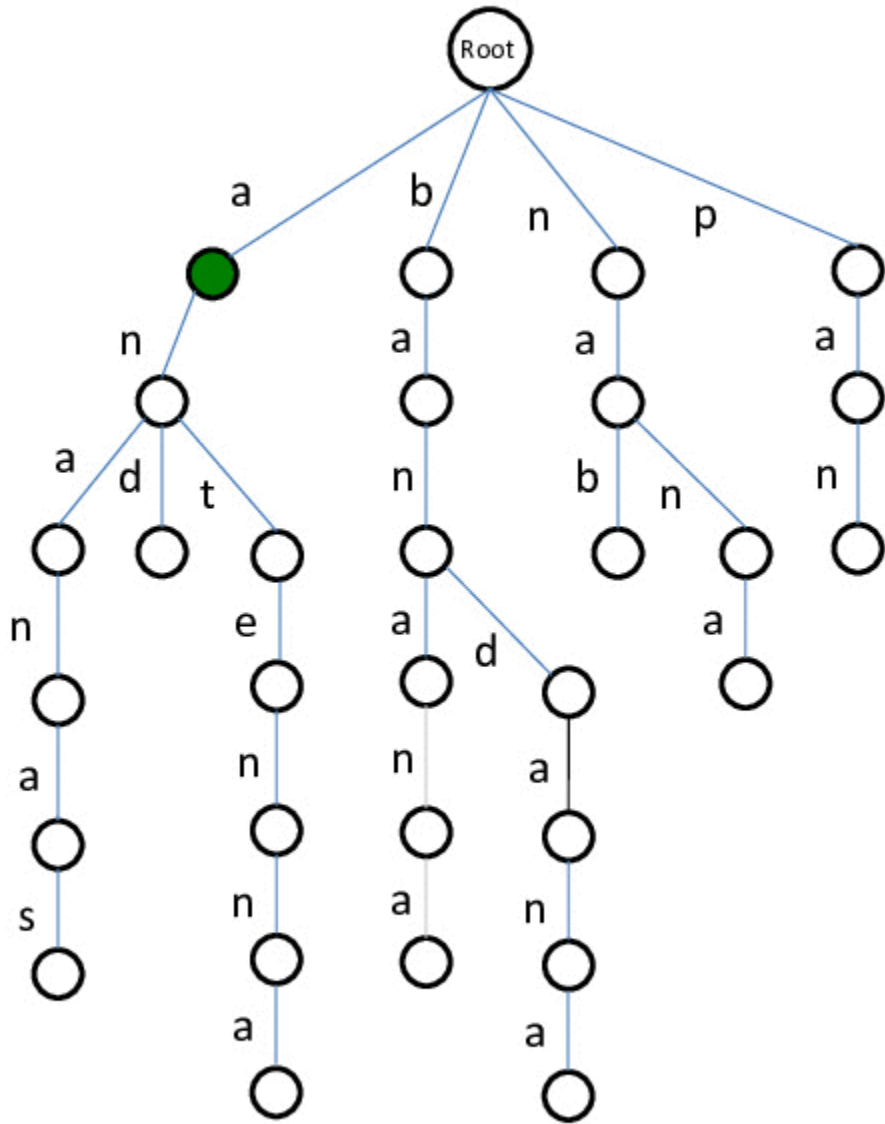
p a n a m a b a n a n a s

# pa**n**amabananas

# pan**a**mabananas

# Computational Complexity

- How long does it take to find multiple patterns?
  - Brute force: O(|text| * |patterns| )

  - Pattern trie: O(|text| * |LongestPattern| )
- Space complexity
  - Number of edges in pattern trie: O(|Patterns|)
  - Human genome: |Patterns| ≈ 10^12

# Alternative: Pack *Text* onto a bus

- Generate all suffixes of Text
  - The suffixes represent all possible places a pattern can match
- Form a trie from these suffixes (**suffix trie**)
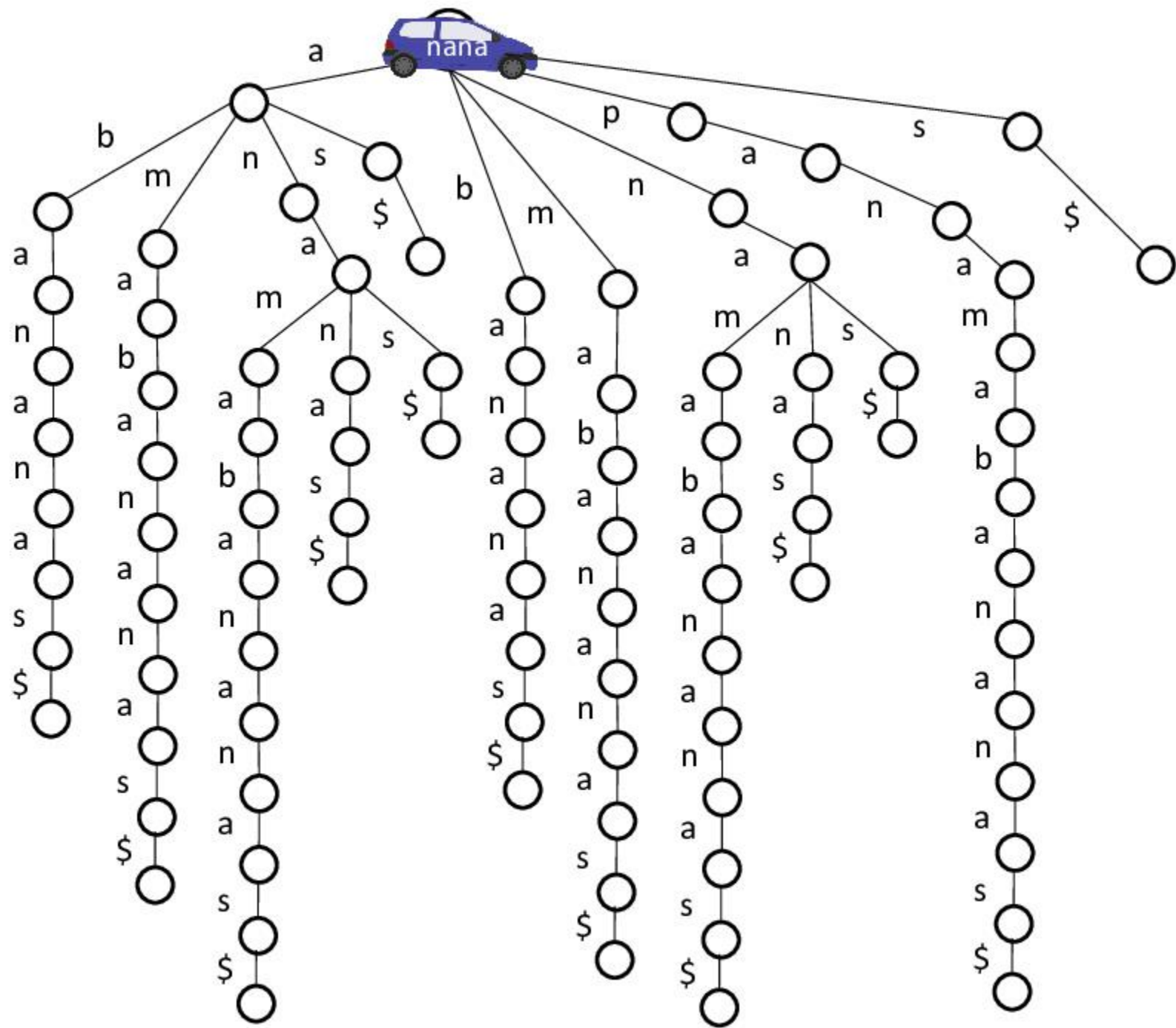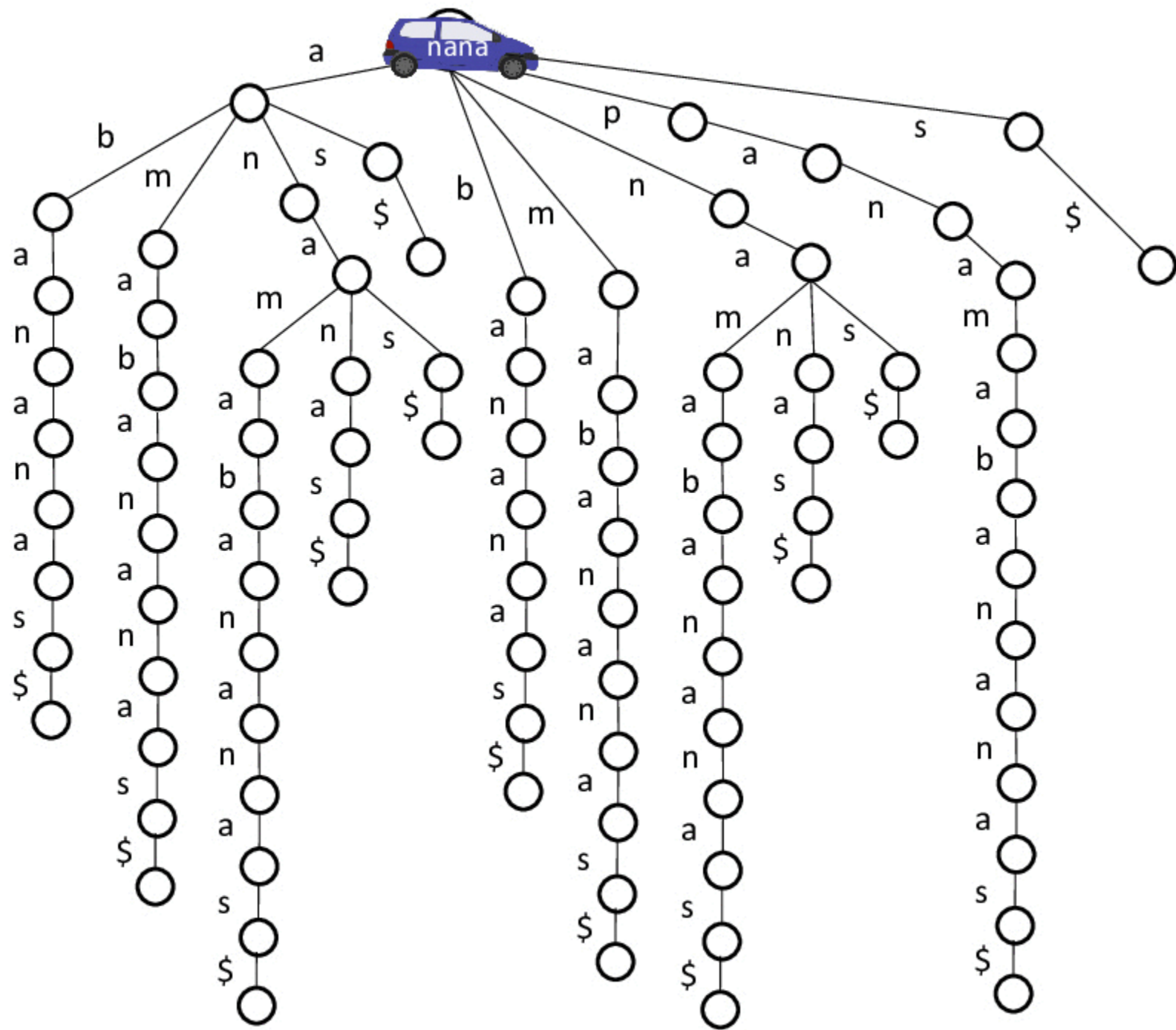- For each pattern, walk down from the root of the trie to see if there is a match.

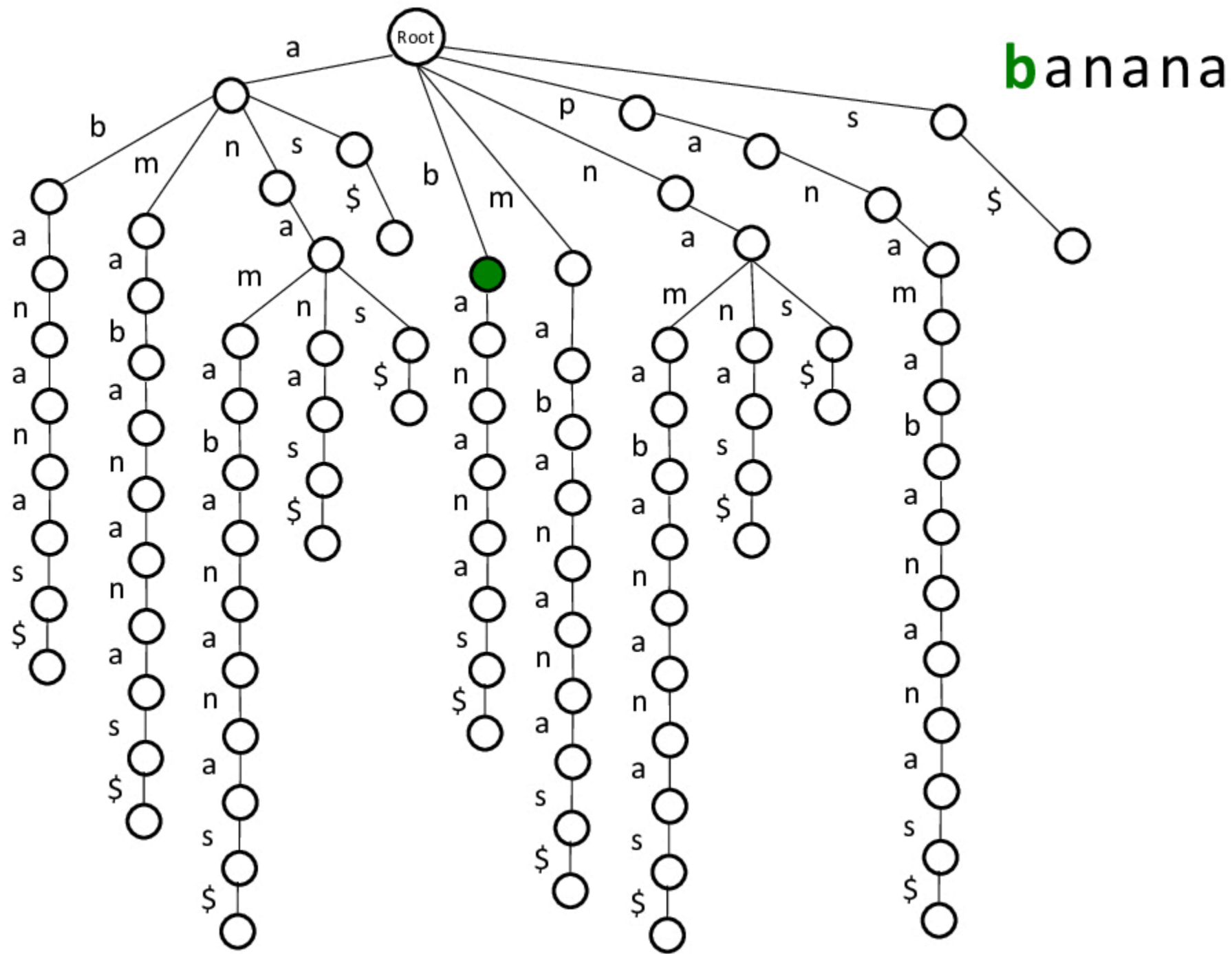( Root )                    p a n a m a b a n a n a s

(Root)

p a n a m a b a n a n a s **$**
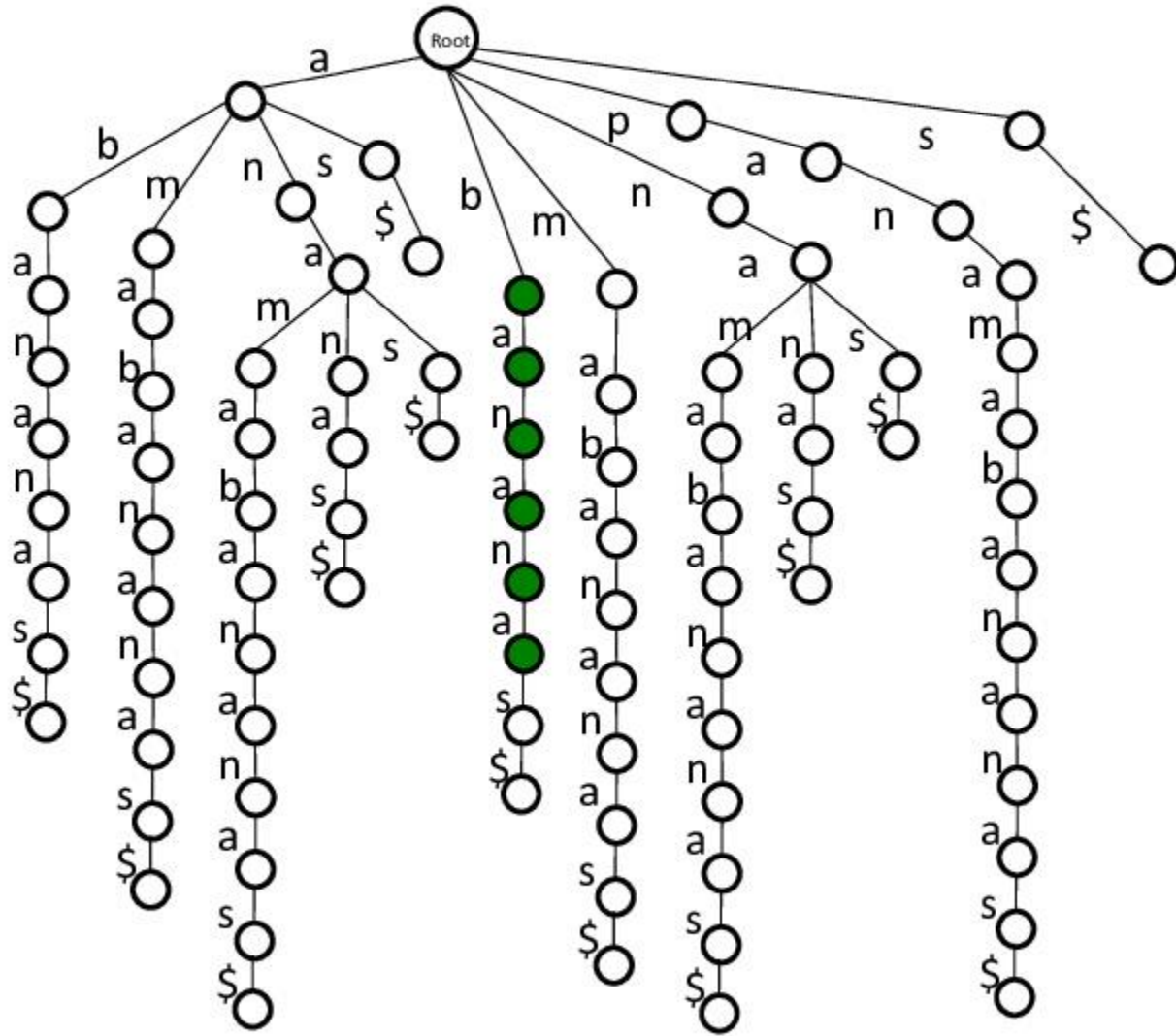
Adding "$" sign in the end  (we'll explain later why)
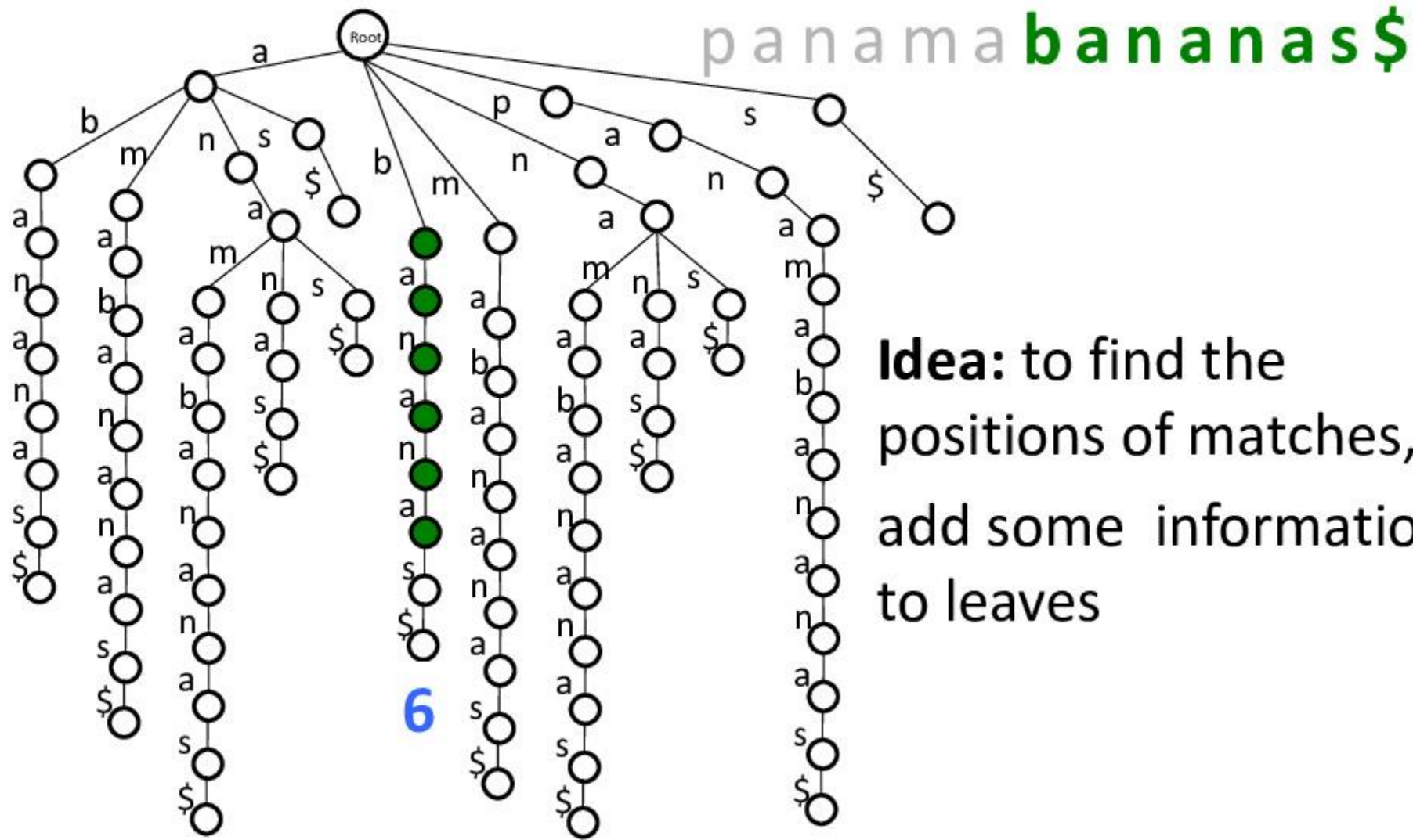
banana

# Where Are the Matches???



bananas$

# Where Are the Matches???



p a n a m a **b a n a n a s $**
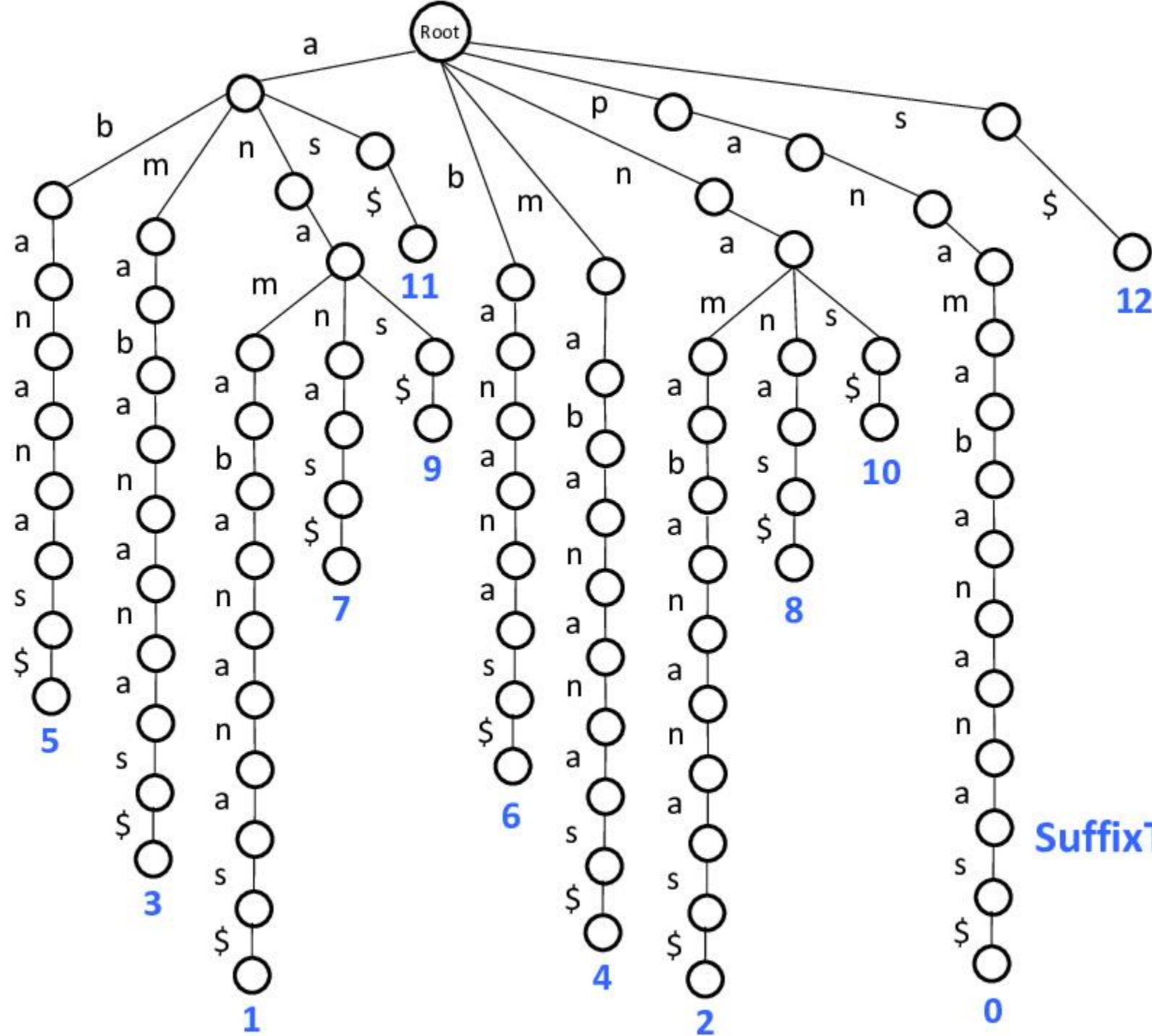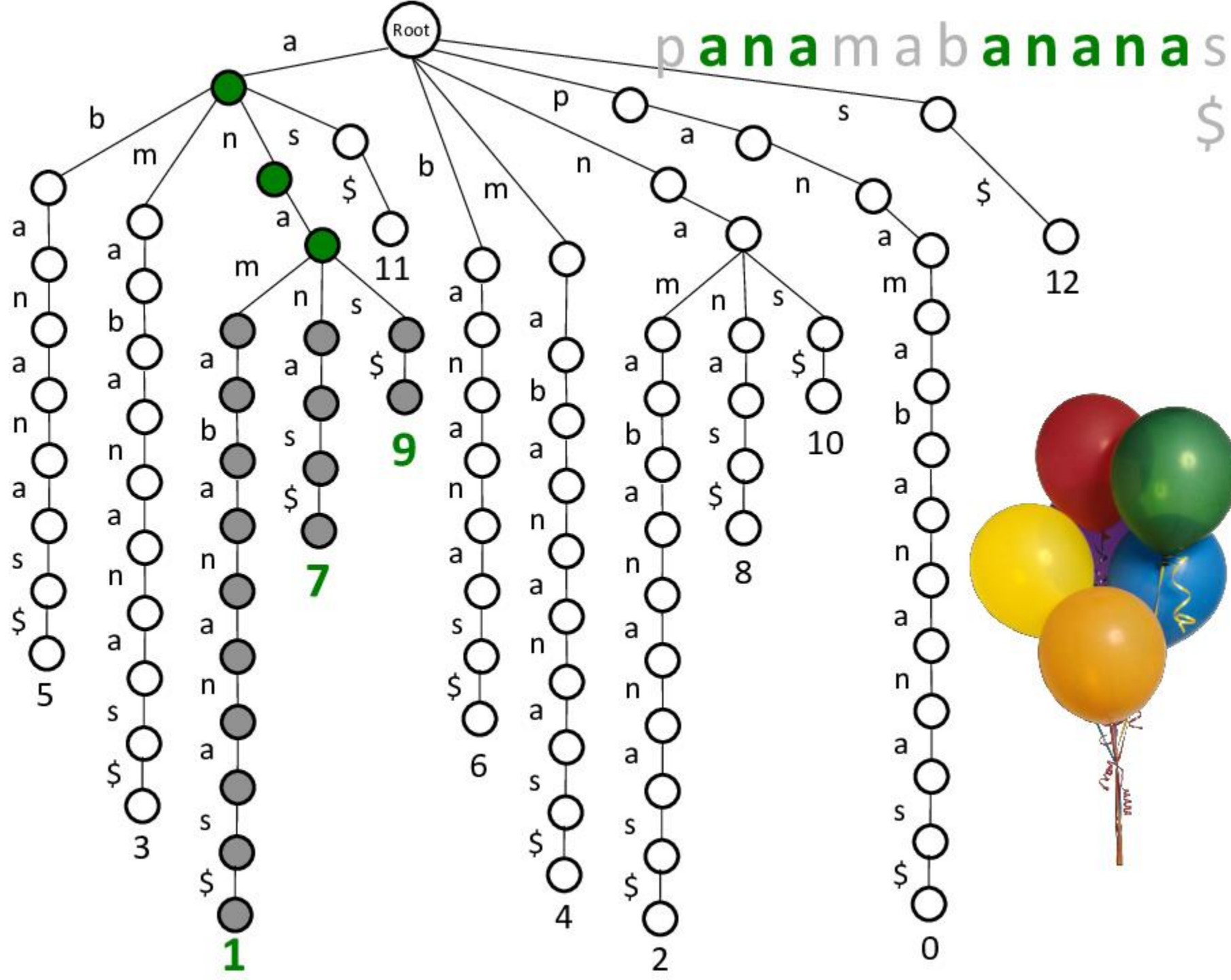
**Idea:** to find the positions of matches, add some information to leaves

SuffixTrie(*Text*)

# Identifying position of match in text

- Once we find a match, walk down to the leaves to get the position(s) of the matches.

# Memory Footprint of Suffix Trie

The suffix trie is formed from |*Text*| suffixes with total length:

$$|Text| * (|Text| - 1)/2$$
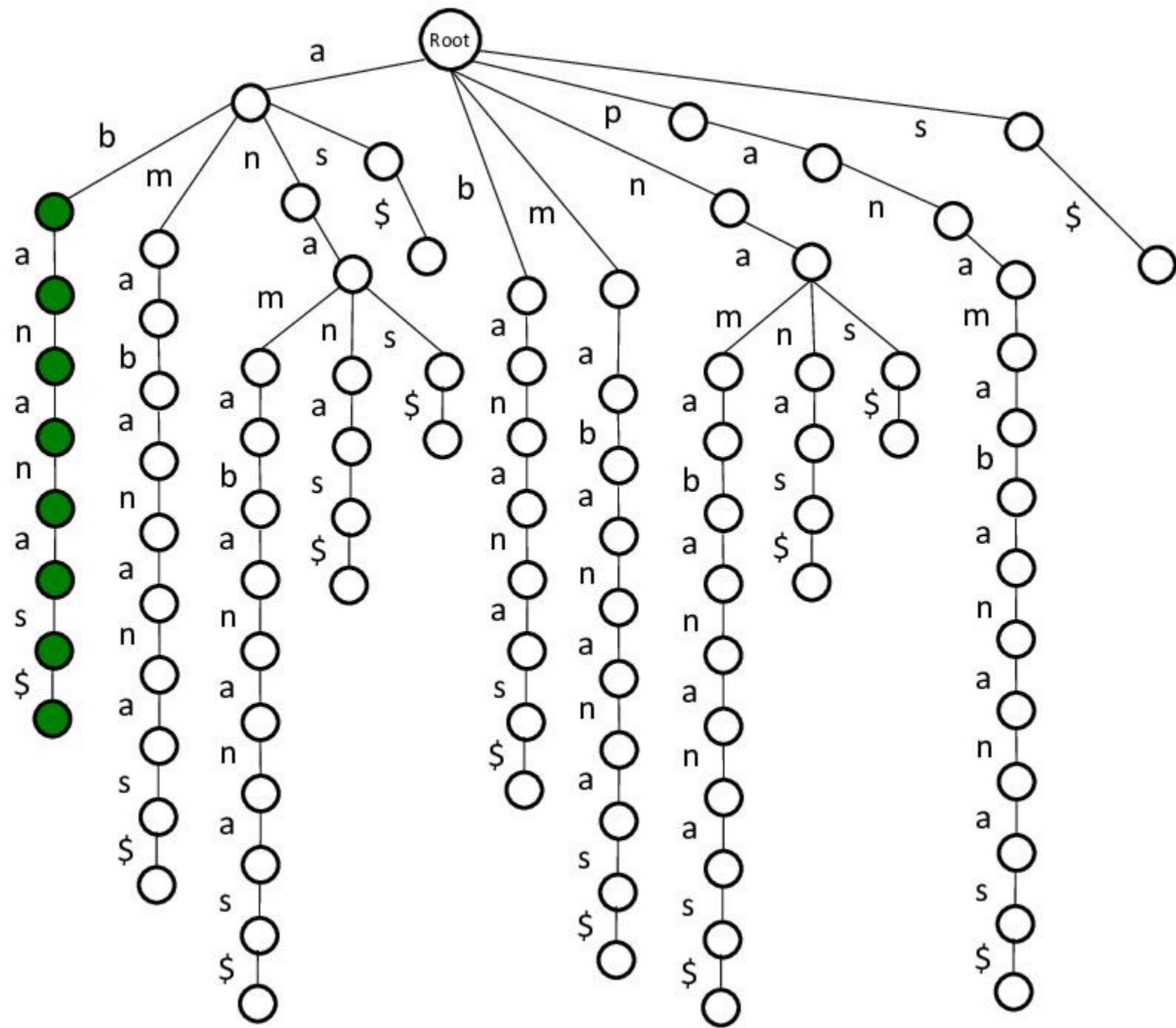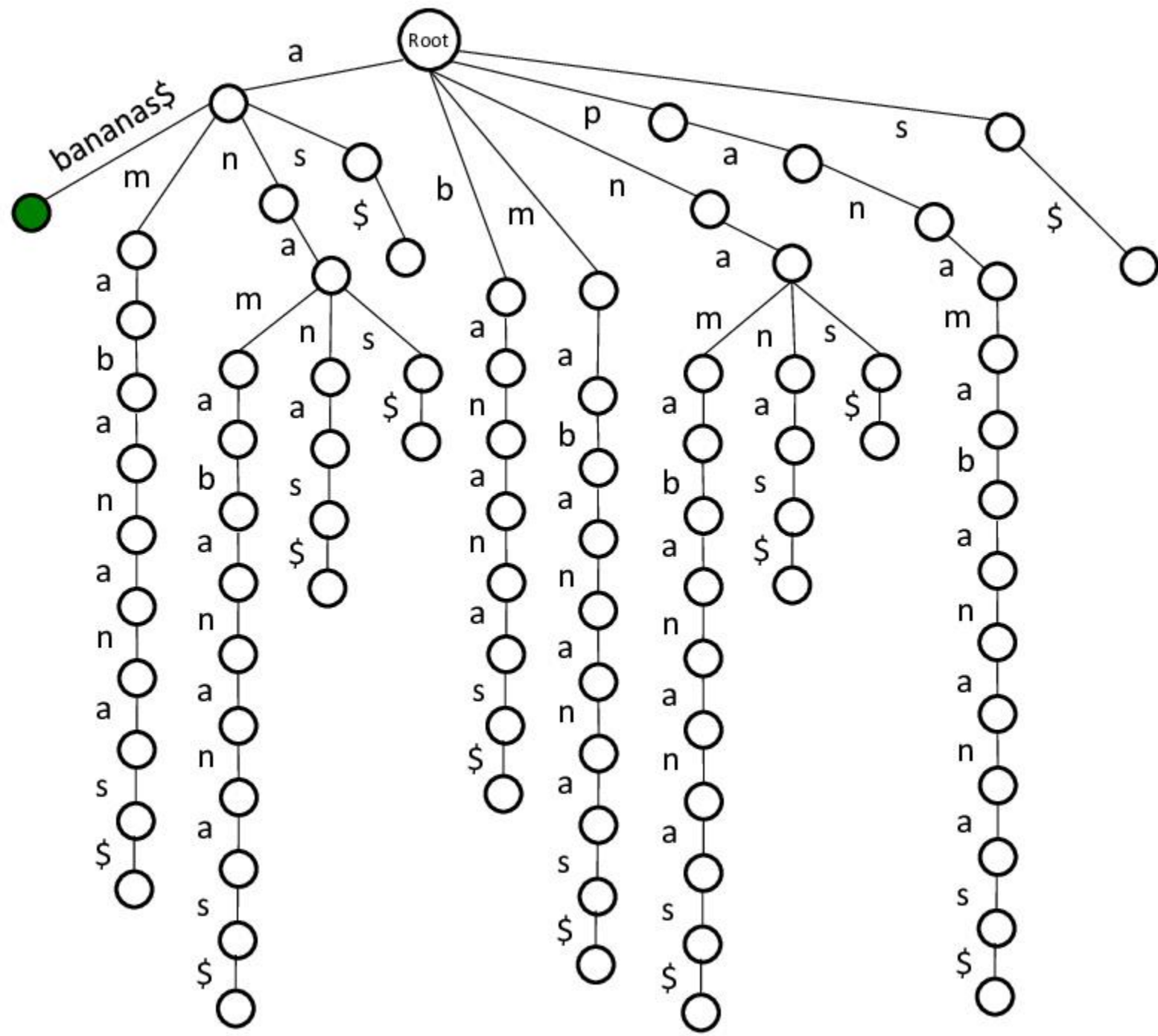
For human genome:
- $|Text| \approx 3 * 10^9$
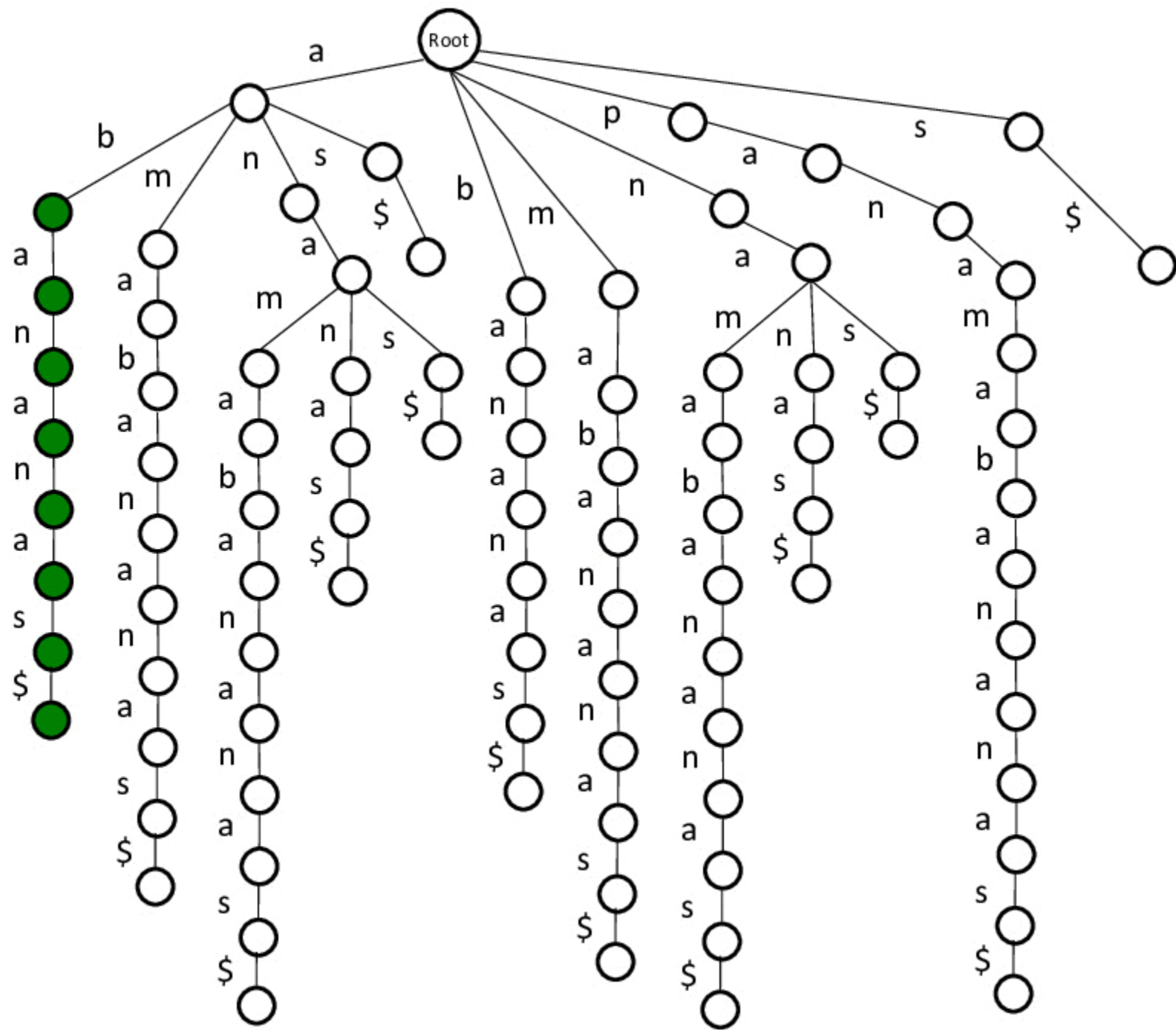
|*Text*| symbols

```
p a n a m a b a n a n a s $
  a n a m a b a n a n a s $
    n a m a b a n a n a s $
      a m a b a n a n a s $
        m a b a n a n a s $
          a b a n a n a s $
            b a n a n a s $
              a n a n a s $
                n a n a s $
                  a n a s $
                    n a s $
                      a s $
                        s $
                          $
```
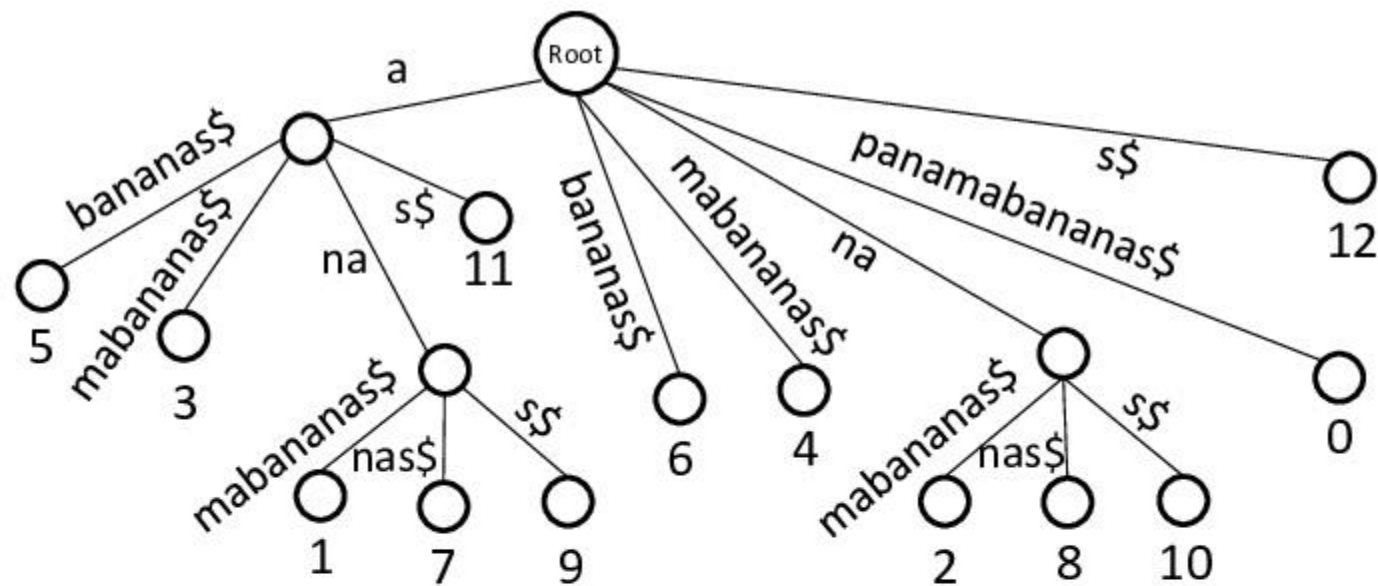
|*Text*| suffixes

Since each suffix adds one leaf and at most one internal vertex to the suffix tree:

- \# vertices < 2|*Text*|

- memory footprint of the suffix tree: $O(|\textit{Text}|)$
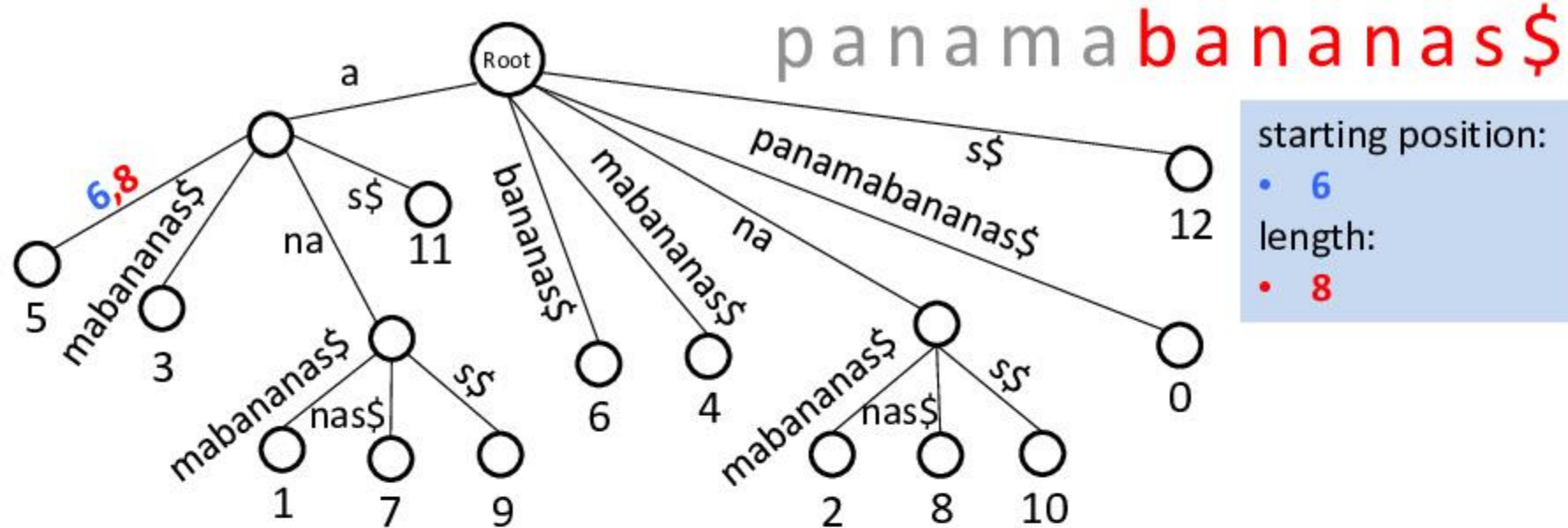
Since each suffix adds one leaf and at most one internal vertex to the suffix tree:

- # vertices < 2|*Text*|
- memory footprint of the suffix tree: O(|*Text*|)
- storing edge labels

# Overview of Suffix Tree

- Fast **Exact** Multiple Pattern Matching
  - Time: O( |text| + |patterns| )
  - Memory: O( |text| )
    - Actual implementation still too demanding in memory requirements. ~20 * |text|
- Construction: O( |text| )
- Need better method that:
  - Can handle mutations (approximate matching)
  - Has smaller memory footprint