

AKAHLLO

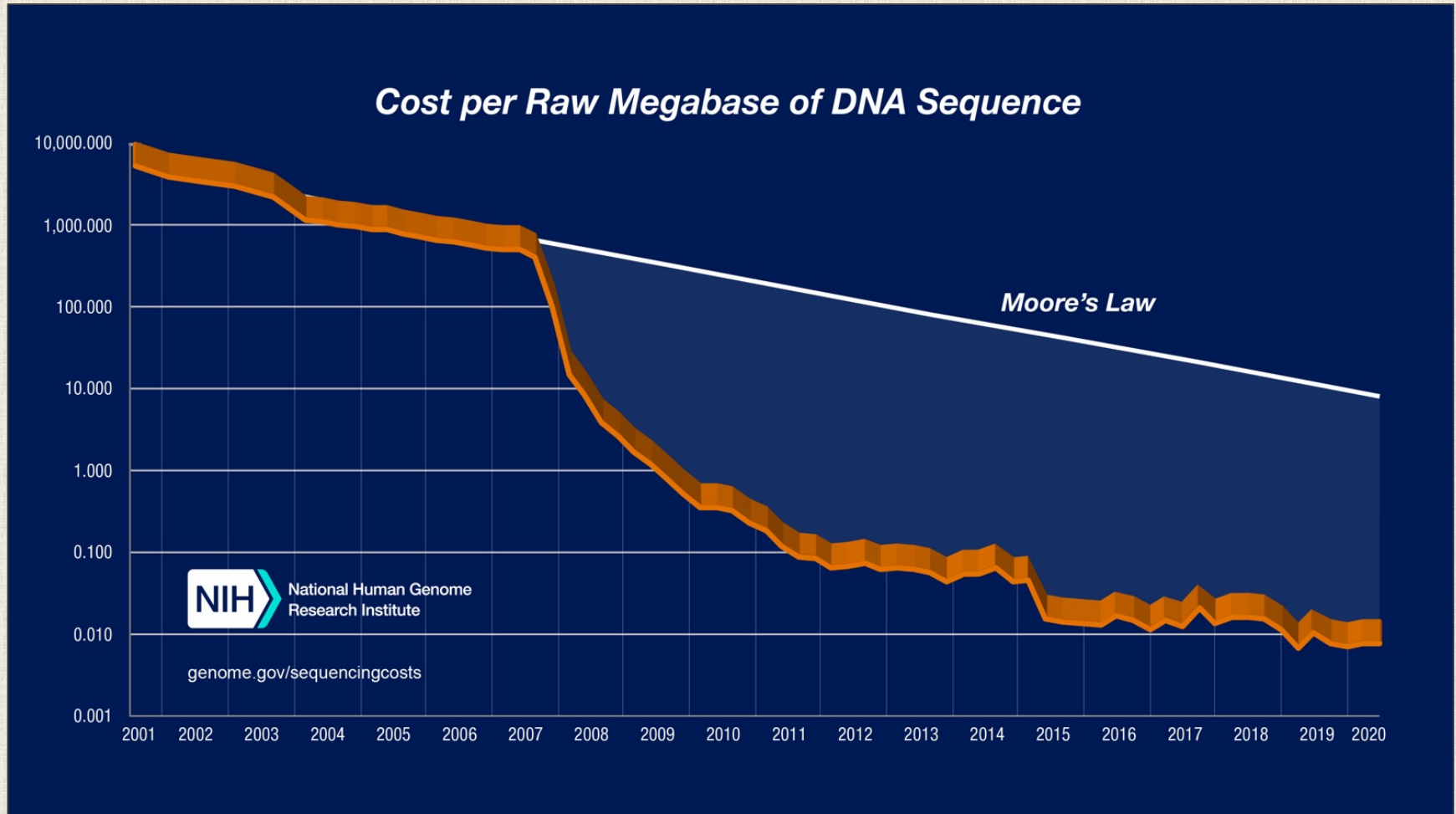
puzzle by Laura Callaghan



Read Mapping

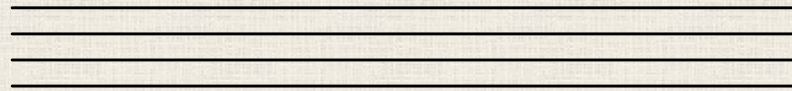


Plummeting Sequencing Costs Mean Millions of Human Genomes

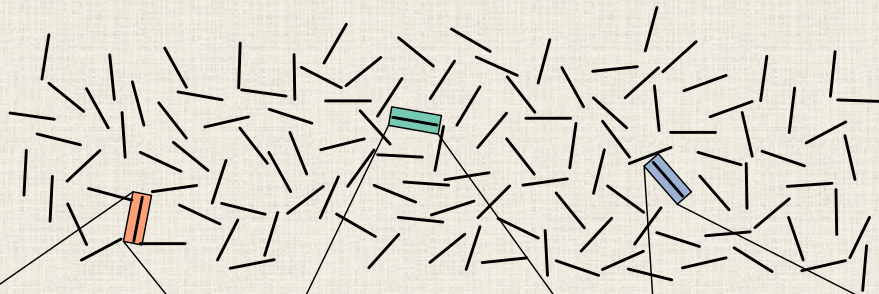


Recall: Overview of Genome Sequencing

Multiple identical copies of a genome



Shatter the genome into reads



Sequence the reads
(Lab)

AGAATATCA

TGAGAATAT

GAGAATATC

Assemble the genome using overlapping reads
(Computational)

AGAATATCA

GAGAATATC

TGAGAATAT

...TGAGAATATCA...

Recall: Overview of Genome Sequencing

But why would we do this if we already have a good sense of the human genome?

Assemble the
genome using
overlapping reads
(Computational)

AGAATATCA
GAGAATATC
TGAGAATAT
. . . TGAGAATATCA . . .

Recall: Overview of Genome Sequencing

But why would we do this if we already have a good sense of the human genome?

Analogy: Would you assemble a jigsaw puzzle without looking at the photo on the box?

Assemble the genome using overlapping reads
(Computational)

AGAATATCA
GAGAATATC
TGAGAATAT
...TGAGAATATCA...

As different as we may seem, we are genetically very similar

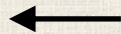
Two humans share ~99.9% of their genomes. The rest is typically single-nucleotide variations.



```
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAG
CTACCACATCGTAGCTACGATGCATTAGCAAGCTATCGATCGAT
CGATCGATTATCTACGATCGATCGATCGATCACTATACGAGCTA
CTACGTACGTACGATCGCGGGACTATTATCGACTACAGATAAAA
CATGCTAGTACAACAGTATACATAGCTGCGGGATACGATTAGCT
AATAGCTGACGATATCCGAT
```



```
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAG
CTACAAACATCGTAGCTACGATGCATTAGCAAGCTATCGATCGAT
CGATCGATTATCTACGATCGATCGATCGATCACTATACGAGCTA
CTACGTACGTACGATCGCGTGGACTATTATCGACTACAGATGAAA
CATGCTAGTACAACAGTATACATAGCTGCGGGATACGATTAGCT
AATAGCTGACGATATCCGAT
```



As different as we may seem, we are genetically very similar

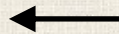
Single-nucleotide polymorphism (SNP): One-nucleotide mutation present in at least 1% of humans.



CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAG
CTAC**C**ACATCGTAGCTACGATGCATTAGCAAGCTATCGATCGAT
CGATCGATTATCTACGATCGATCGATCGATCACTATACGAGCTA
CTACGTACGTACGATCGCG**G**GACTATTATCGACTACAGAT**A**AAA
CATGCTAGTACAACAGTATACATAGCTGCGGGATACGATTAGCT
AATAGCTGACGATATCCGAT



CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAG
CTAC**A**ACATCGTAGCTACGATGCATTAGCAAGCTATCGATCGAT
CGATCGATTATCTACGATCGATCGATCGATCACTATACGAGCTA
CTACGTACGTACGATCGCG**T**GACTATTATCGACTACAGAT**G**AAA
CATGCTAGTACAACAGTATACATAGCTGCGGGATACGATTAGCT
AATAGCTGACGATATCCGAT



A First Attempt at Read Mapping

Read-mapping: Comparison of sequencing reads *Patterns* from an individual against a **reference human genome** *Text* stored in a database.

CTGAGGATGGACTACGCTACTACTGATAGCTGTTT
GAGGA CCACG TGA-A

Text

Patterns

A First Attempt at Read Mapping

Read-mapping: Comparison of sequencing reads *Patterns* from an individual against a **reference human genome** *Text* stored in a database.

CTGAGGATGGACTACGCTACTACTGATAGCTGTTT
GAGGA CCACG TGA-A

Text

Patterns

STOP: How could we map reads against a reference genome?

A First Attempt at Read Mapping

Read-mapping: Comparison of sequencing reads *Patterns* from an individual against a **reference human genome** *Text* stored in a database.

CTGAGGATGGACTACGCTACTACTGATAGCTGTTT
GAGGA CCACG TGA-A

Text

Patterns

Idea: We could *align* each read against the reference genome. But what problem are we solving?

Recall these Two Problems from Text

A **fitting alignment** of v and w is an alignment of a substring of v against all of w .

Fitting Alignment Problem:

- **Input:** Two strings and a scoring matrix.
- **Output:** A fitting alignment of the strings with maximum alignment score according to the scoring matrix.

Global

CGTAGGCTTAAGGTTA
A-TAG----A---T-A

Local

CGTAGGCTTAAGGTTA
ATAGATA

Fitting

CGTAGGCTTAAGGTTA
ATAGA--TA

A First Attempt at Read Mapping

Read-mapping: Comparison of sequencing reads *Patterns* from an individual against a **reference human genome** *Text* stored in a database.

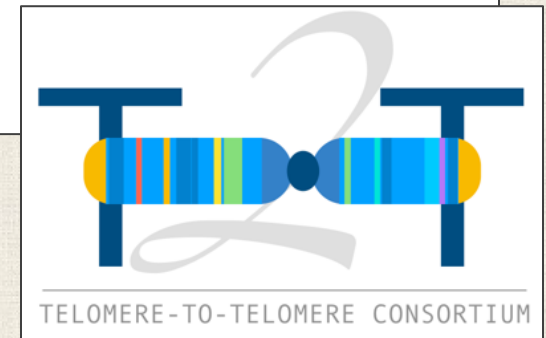
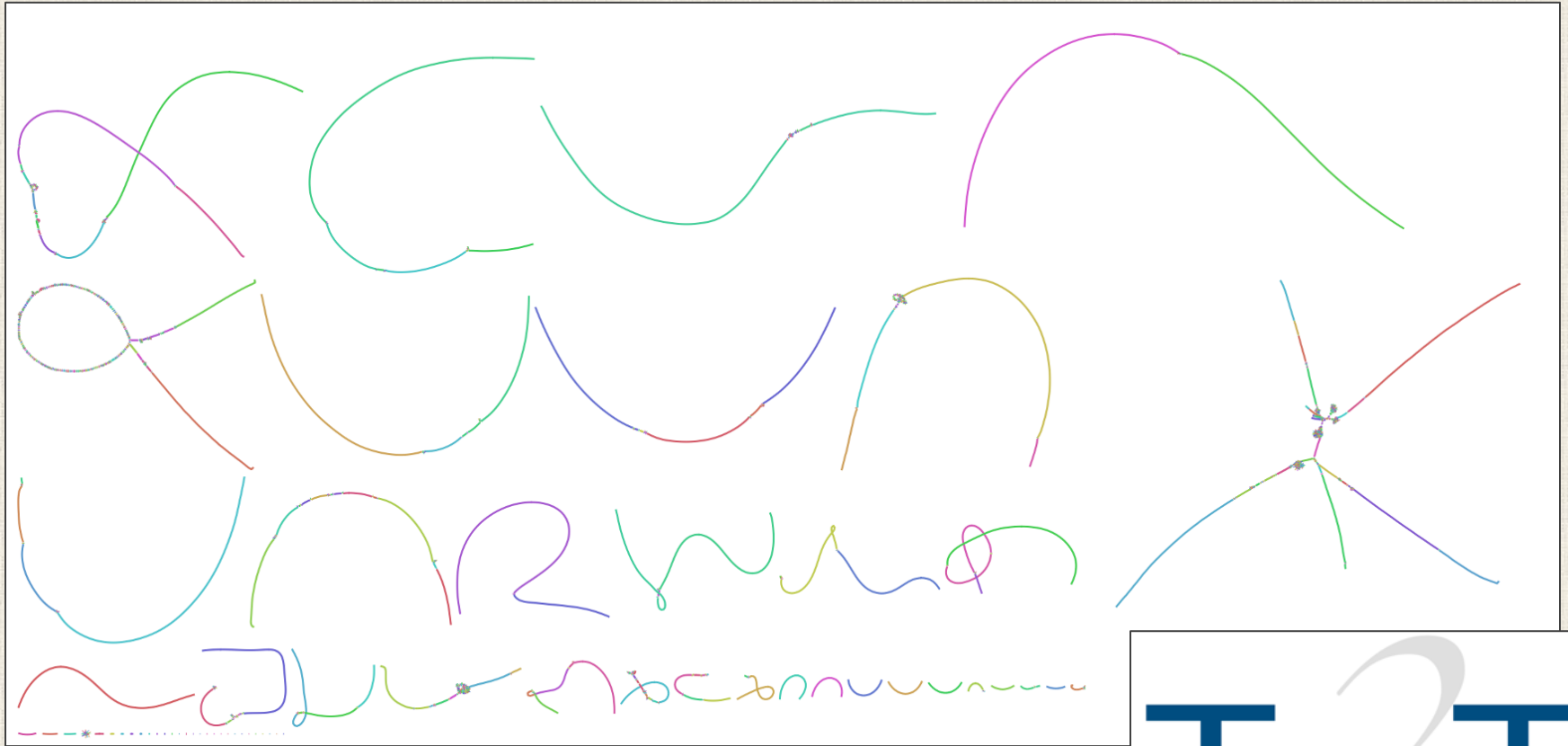
CTGAGGATGGACTACGCTACTACTGATAGCTGTTT
GAGGA CCACG TGA-A

Text

Patterns

Problem: Fitting alignment has runtime $O(|Pattern| * |Text|)$ for each string *Pattern* in our reads – too long!

Remember: The First *Full* Human Genome Wasn't Sequenced Until 2020!



Reference Genomes are Not Diverse

GRCh38.p13

Description: Genome Reference Consortium Human Build 38 patch release 13 (GRCh38.p13)

Organism name: [Homo sapiens \(human\)](#)

BioProject: [PRJNA31257](#)

Submitter: Genome Reference Consortium

Date: 2019/02/28

Assembly type: haploid-with-alt-loci

Release type: patch

Assembly level: Chromosome

Genome representation: full

RefSeq category: reference genome

GenBank assembly accession: GCA_000001405.28 (latest)

RefSeq assembly accession: GCF_000001405.39 (latest)

RefSeq assembly and GenBank assembly identical: no ([hide details](#))

- Only in GenBank: 1 unplaced scaffold (in primary assembly-unit)
- Data displayed for RefSeq version

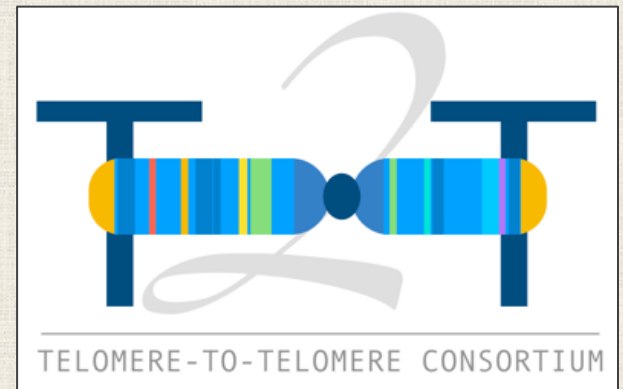
https://www.ncbi.nlm.nih.gov/assembly/GCF_000001405.39

When we “map” reads against a reference human genome, the most commonly used reference is 70% made up of RP11, or “some guy from Buffalo”.

The first complete human reference is taken from a single sperm

This reference is quickly being replaced by “CHM13”, the first ever complete human genome.

CHM13 is a cell line grown from a “molar pregnancy”, an empty egg fertilized by a sperm. Both copies of chromosomes come from this sperm.



Long, Cheap Reads Give Hope of a "Pan Genome" Era

TOWARDS A
COMPLETE
REFERENCE OF
HUMAN GENOME
DIVERSITY



STOP: Why might mapping reads from one ethnicity against a genome from a different ethnicity lead to biased conclusions?

INTRODUCTION TO MULTIPLE PATTERN MATCHING

Simplifying Assumption: No Mutations

We will first find all the *exact* matches of reads in the reference. If we are looking for variants, then we can eliminate from consideration areas of the reference genome where exact matches occur.

Simplifying Assumption: No Mutations

We will first find all the *exact* matches of reads in the reference. If we are looking for variants, then we can eliminate from consideration areas of the reference genome where exact matches occur.

Key Point: This lecture will assume that we are working with highly accurate reads rather than messy reads.

Simplifying Assumption: No Mutations

We will first find all the *exact* matches of reads in the reference. If we are looking for variants, then we can eliminate from consideration areas of the reference genome where exact matches occur.

Multiple Pattern Matching Problem: *Find all occurrences of a collection of patterns in a text.*

- **Input:** A string *Text* and a collection *Patterns* containing (shorter) strings.
- **Output:** All starting positions in *Text* where a string from *Patterns* appears as a substring.

Wait: Isn't This Too Simple?

STOP: What is the probability of a randomly selected 250-mer matching in two genomes with 99.9% similarity? (Assume mutations are uniform.)

Multiple Pattern Matching Problem: *Find all occurrences of a collection of patterns in a text.*

- **Input:** A string *Text* and a collection *Patterns* containing (shorter) strings.
- **Output:** All starting positions in *Text* where a string from *Patterns* appears as a substring.

Wait: Isn't This Too Simple?

Answer: We need all 250 positions to match, so the probability is $(1-0.001)^{250} = 0.779$. So nearly 80% of 250-mers will match in two 99.9% similar genomes.

Multiple Pattern Matching Problem: *Find all occurrences of a collection of patterns in a text.*

- **Input:** A string *Text* and a collection *Patterns* containing (shorter) strings.
- **Output:** All starting positions in *Text* where a string from *Patterns* appears as a substring.

Wait: Isn't This Too Simple?

BruteForcePatternMatching: slide each string down *Text*, one at a time.

Multiple Pattern Matching Problem: *Find all occurrences of a collection of patterns in a text.*

- **Input:** A string *Text* and a collection *Patterns* containing (shorter) strings.
- **Output:** All starting positions in *Text* where a string from *Patterns* appears as a substring.

Wait: Isn't This Too Simple?

STOP: What is the Big-O runtime of **BruteForcePatternMatching?**

Multiple Pattern Matching Problem: *Find all occurrences of a collection of patterns in a text.*

- **Input:** A string *Text* and a collection *Patterns* containing (shorter) strings.
- **Output:** All starting positions in *Text* where a string from *Patterns* appears as a substring.

So Things Are Much Harder Than They Seem ...

Answer: For one read, it is $O(|Text| * |Pattern|)$. For multiple reads, it is $O(|Text| * |Patterns|)$, where $|Patterns|$ is the total length of *Patterns*.

Multiple Pattern Matching Problem: *Find all occurrences of a collection of patterns in a text.*

- **Input:** A string *Text* and a collection *Patterns* containing (shorter) strings.
- **Output:** All starting positions in *Text* where a string from *Patterns* appears as a substring.

HERDING PATTERNS INTO A TRIE

In Brute Force, Patterns Travel One at a Time

Text

CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACCACATCGTAGCTAC

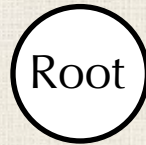


Idea: Organize Patterns and Make a Single Pass Down the Reference Genome

Text

CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACCACATCGTAGCTAC





Trie: Graph constructed from *Patterns*.

- There is a single root node with indegree 0, denoted *root*; all other nodes have indegree 1.
- Each edge of $Trie(Patterns)$ is labeled with a letter.
- Edges leading out of a given node have distinct labels.
- Every string in *Patterns* is spelled out by concatenating the letters along some path from the root downward.
- Every path from the root to a **leaf** (node with outdegree 0) spells a string from *Patterns*.

Patterns

banana

pan

and

nab

antenna

bandana

ananas

nana

Root

Patterns

banana

pan

and

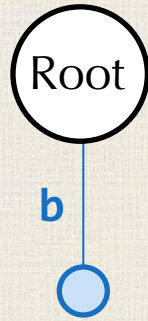
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

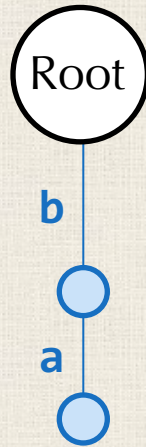
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

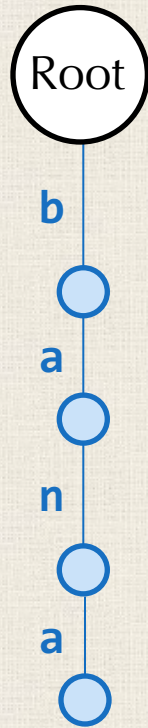
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

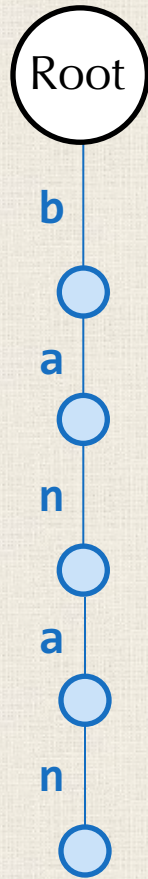
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

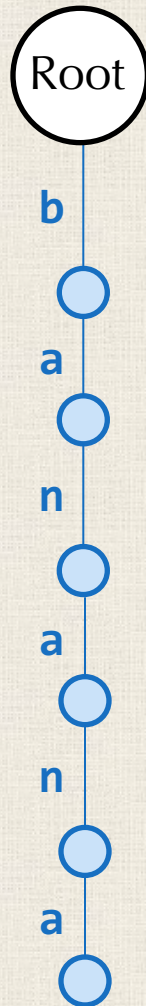
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

nab

antenna

bandana

ananas

nana

Patterns

banana

pan

and

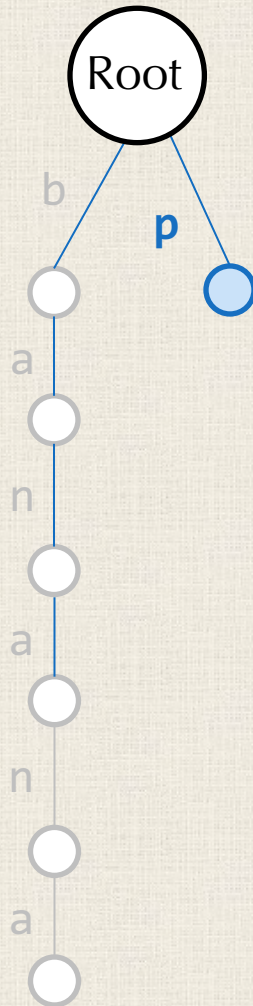
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

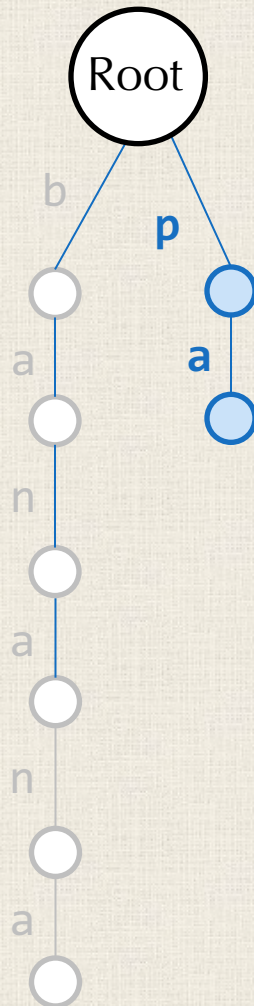
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

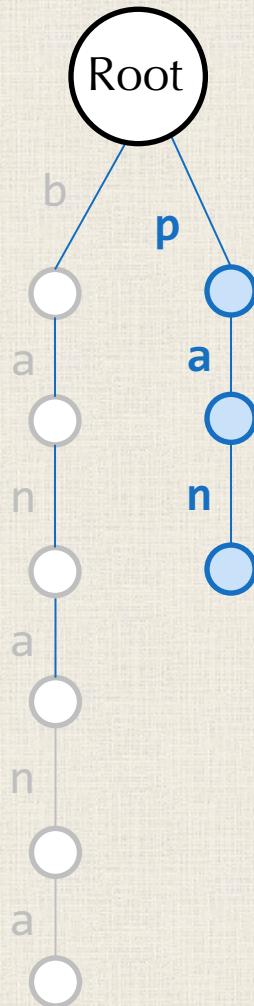
nab

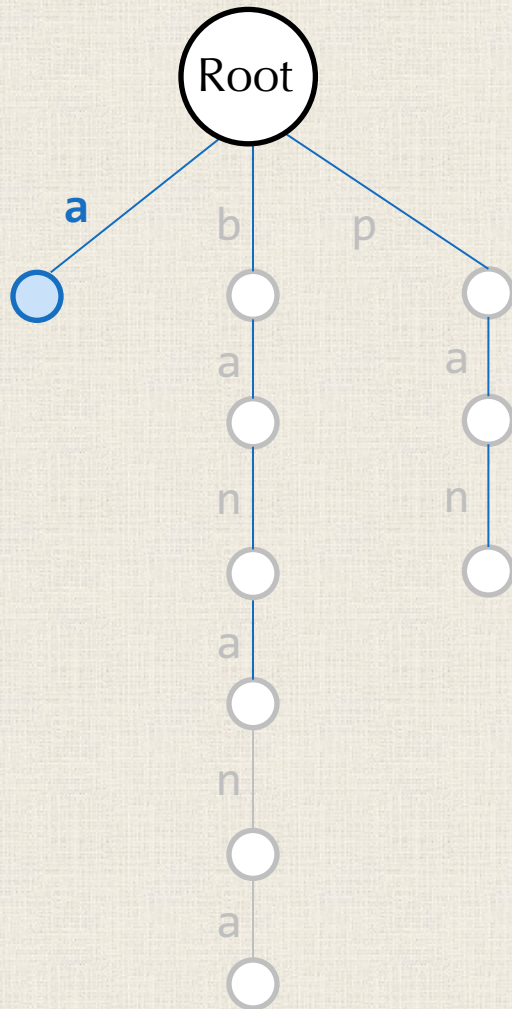
antenna

bandana

ananas

nana





Patterns

banana

pan

and

nab

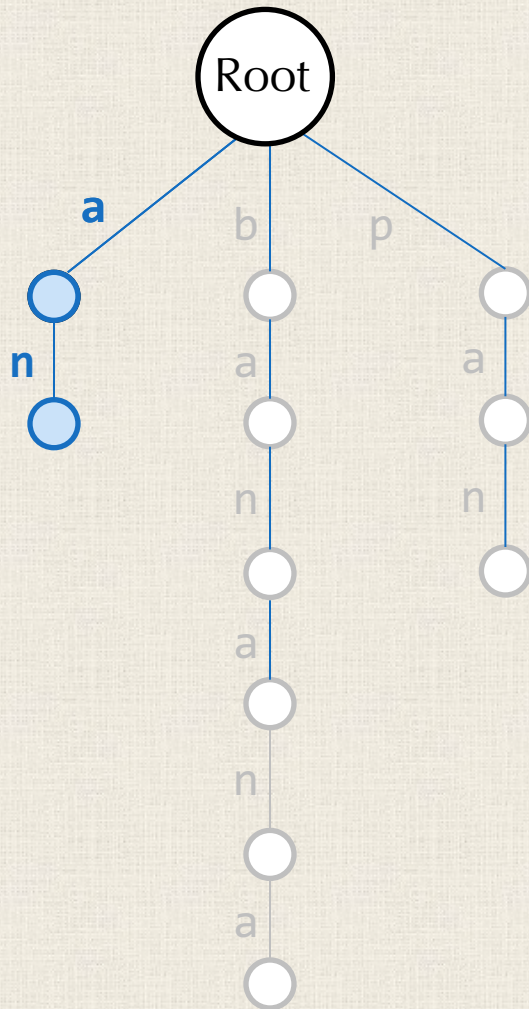
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

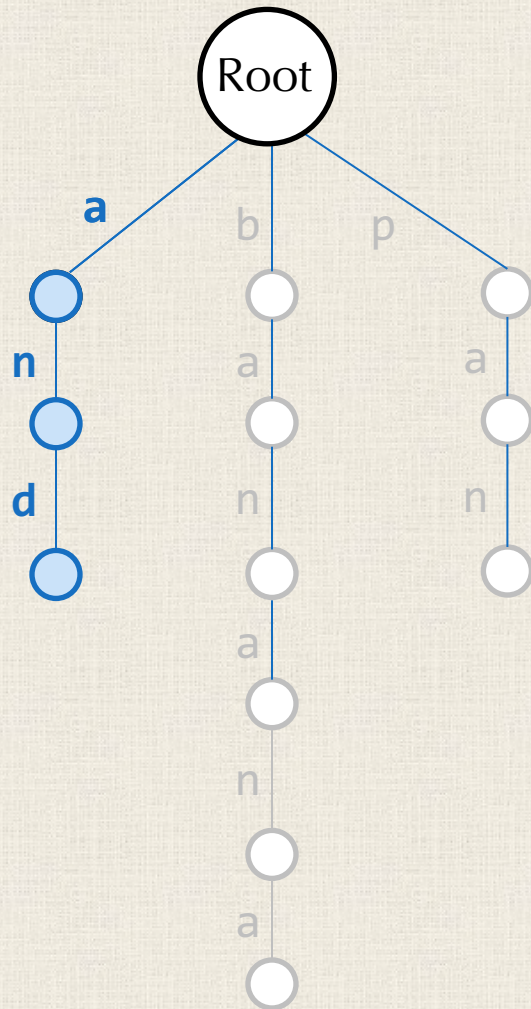
antenna

bandana

ananas

nana

Patterns



banana

pan

and

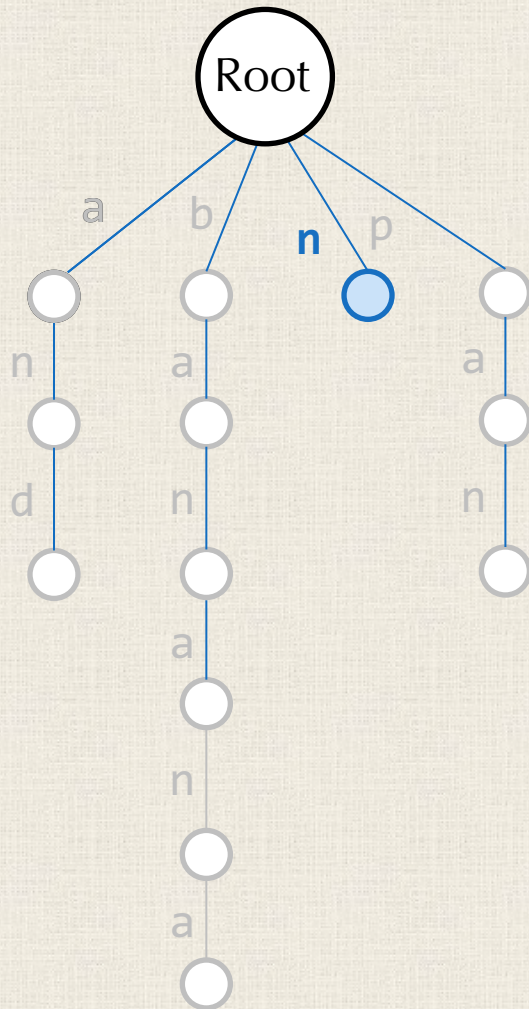
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

nab

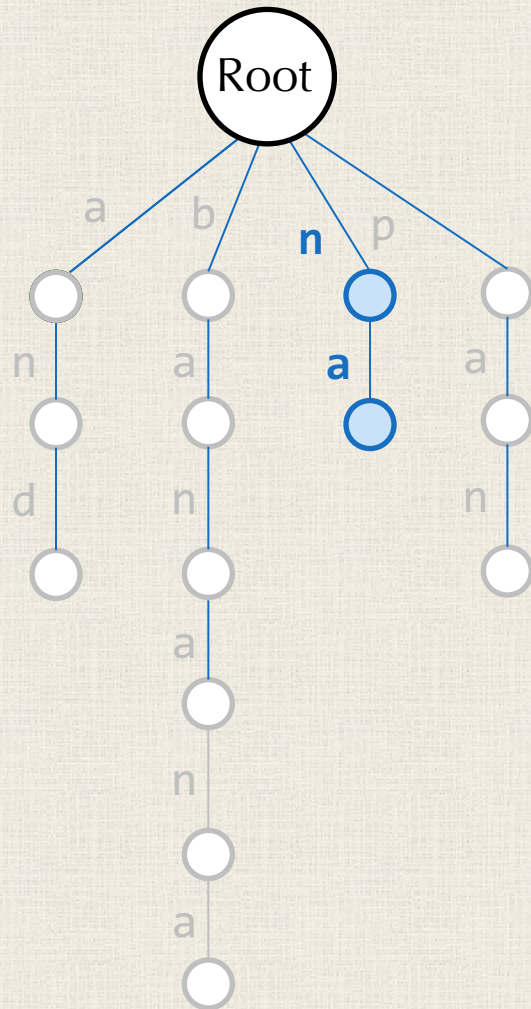
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

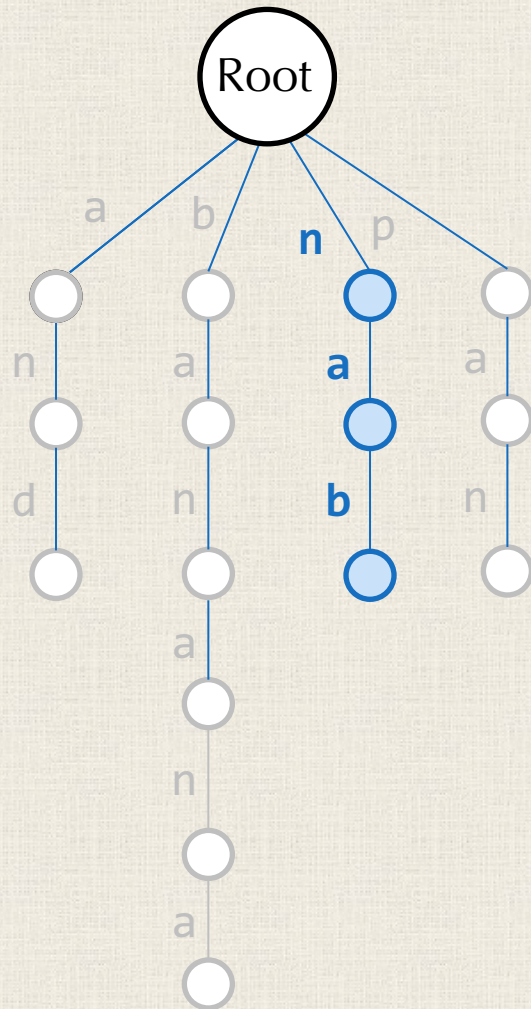
antenna

bandana

ananas

nana

Patterns



banana

pan

and

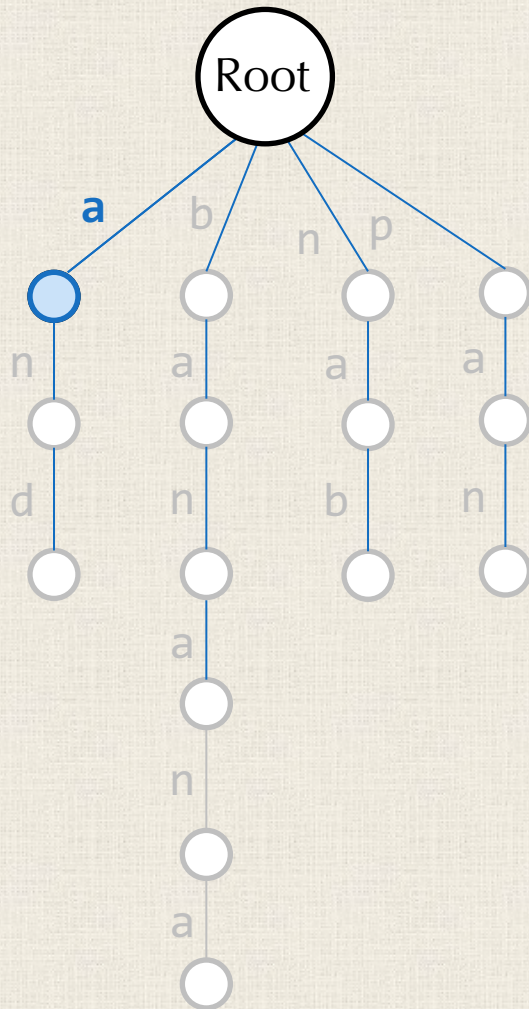
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

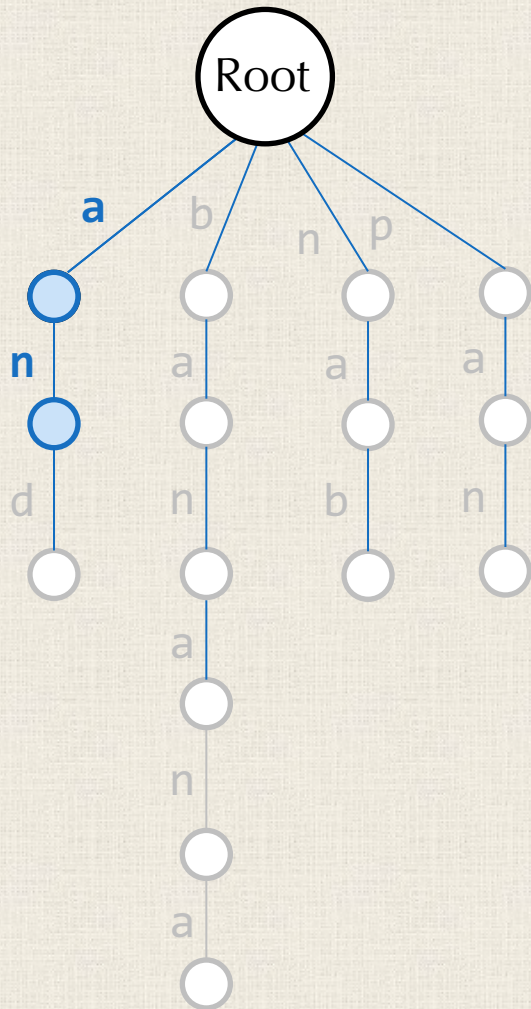
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

nab

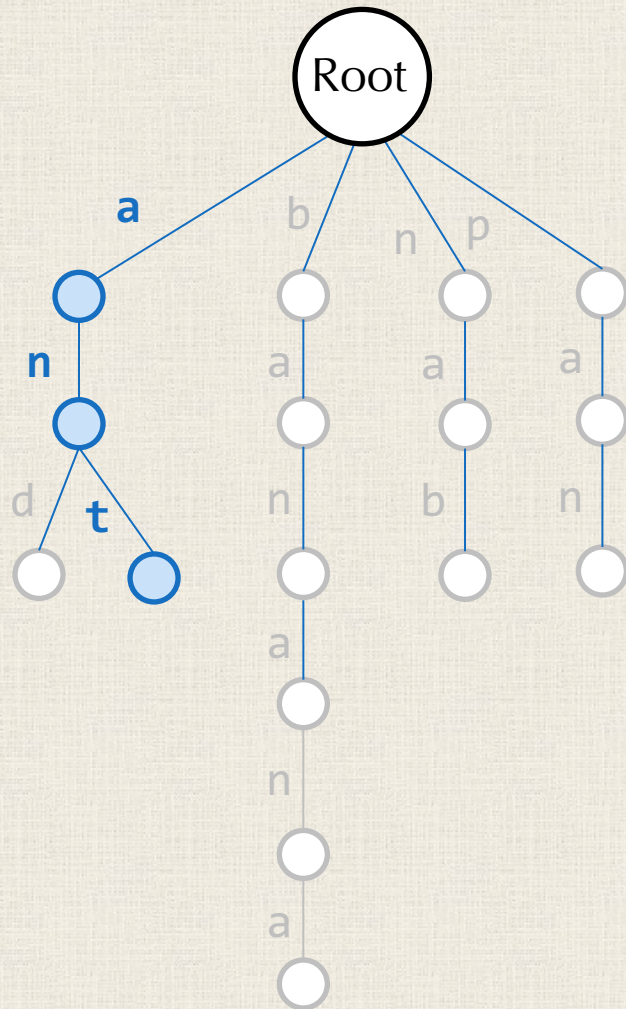
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

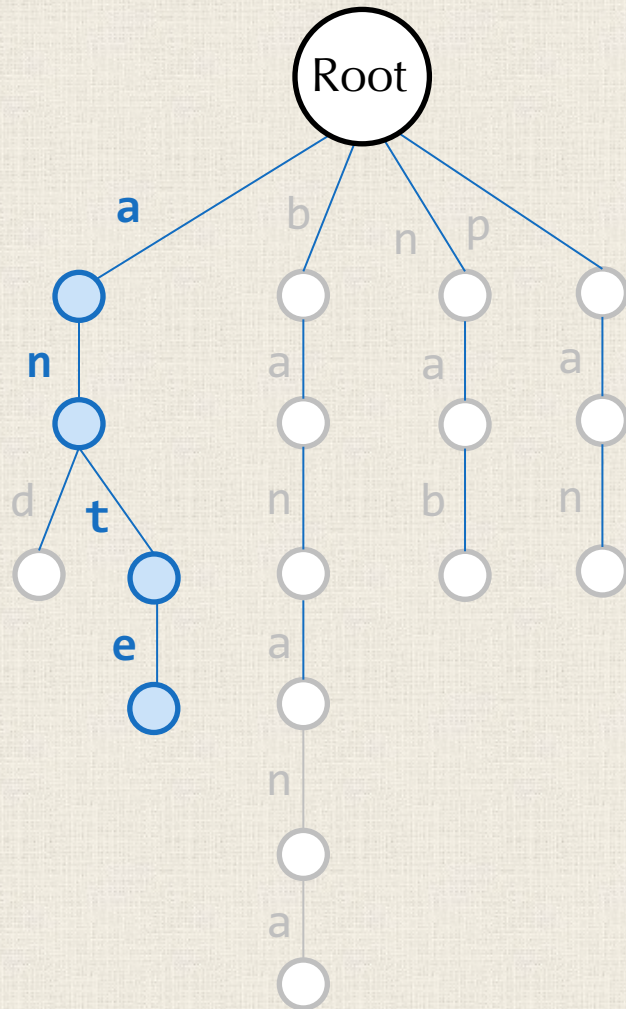
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

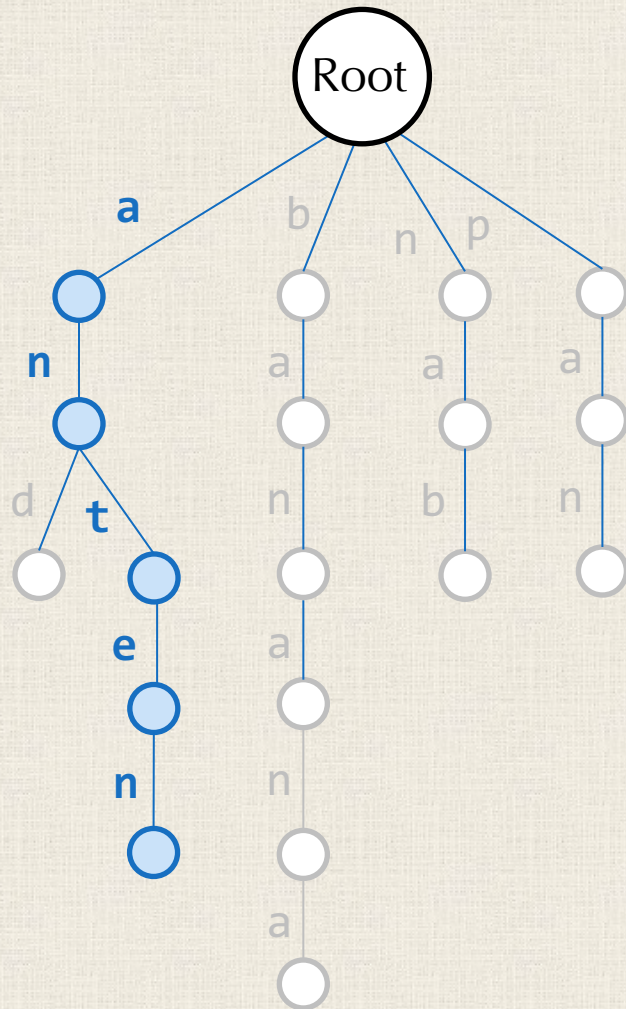
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

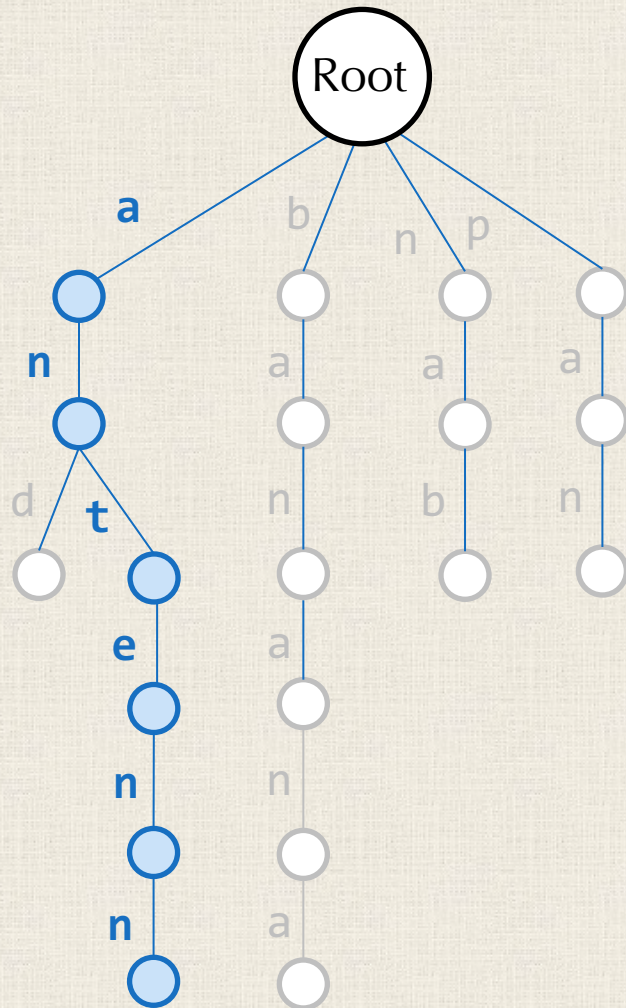
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

antenna

bandana

ananas

nana

Patterns

banana

pan

and

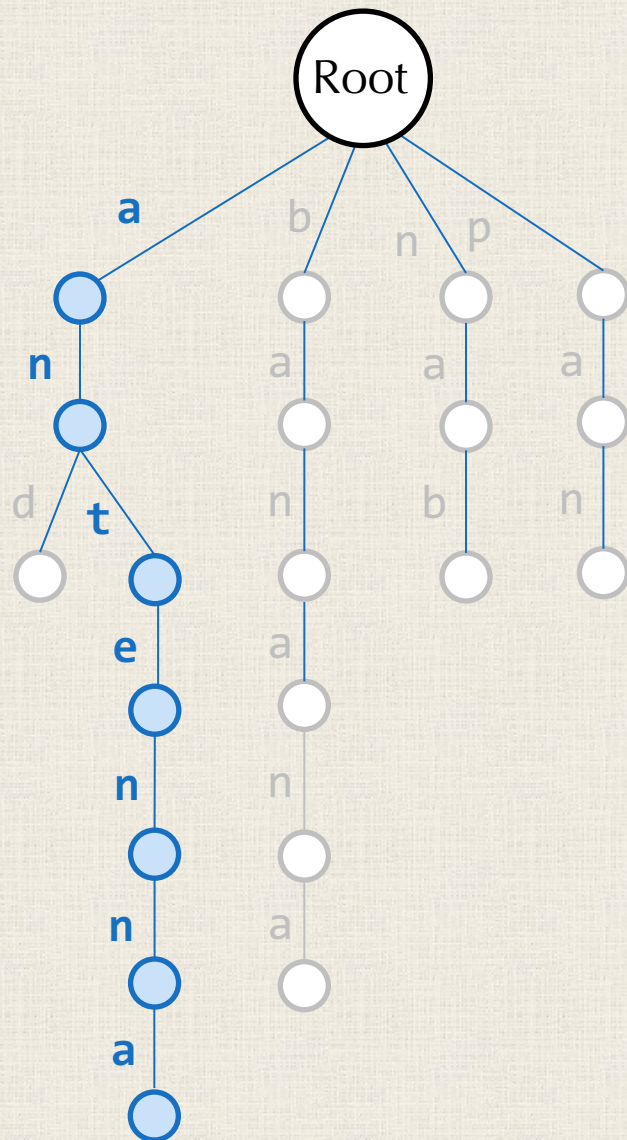
nab

antenna

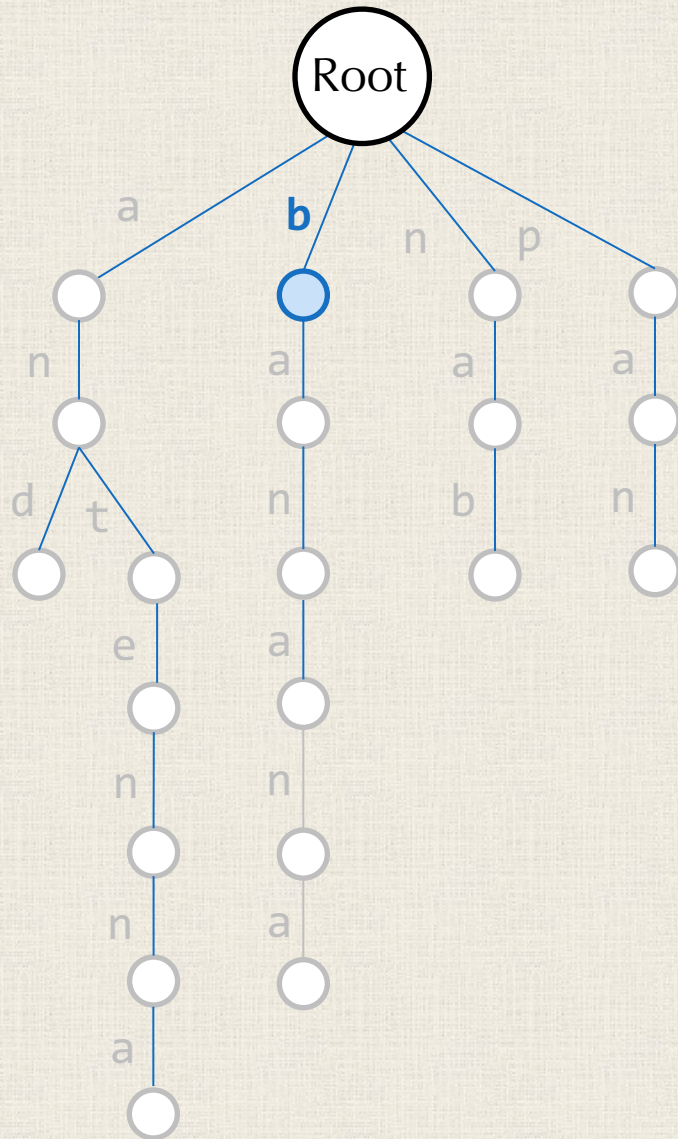
bandana

ananas

nana



Patterns



banana

pan

and

nab

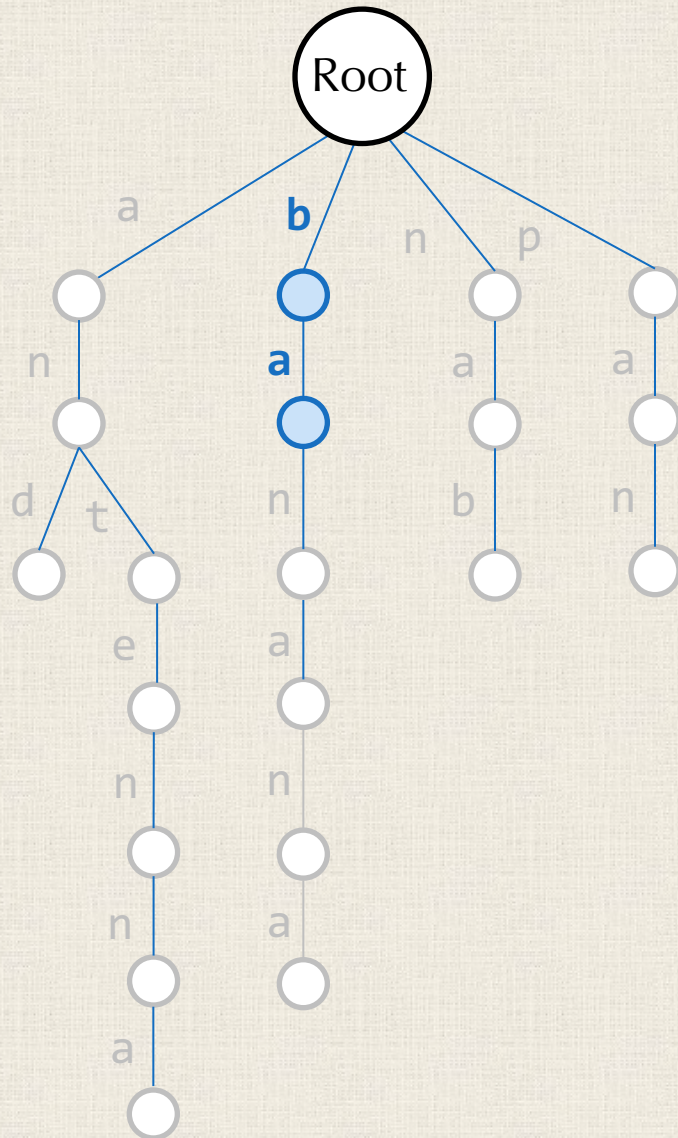
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

antenna

bandana

ananas

nana

Patterns

banana

pan

and

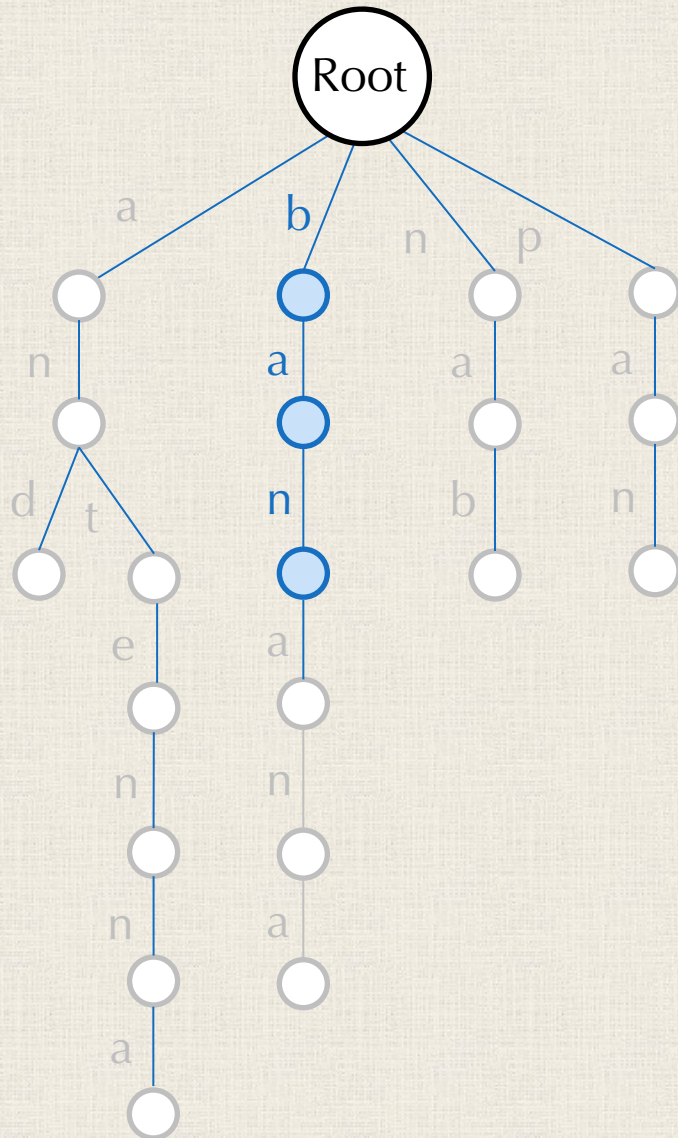
nab

antenna

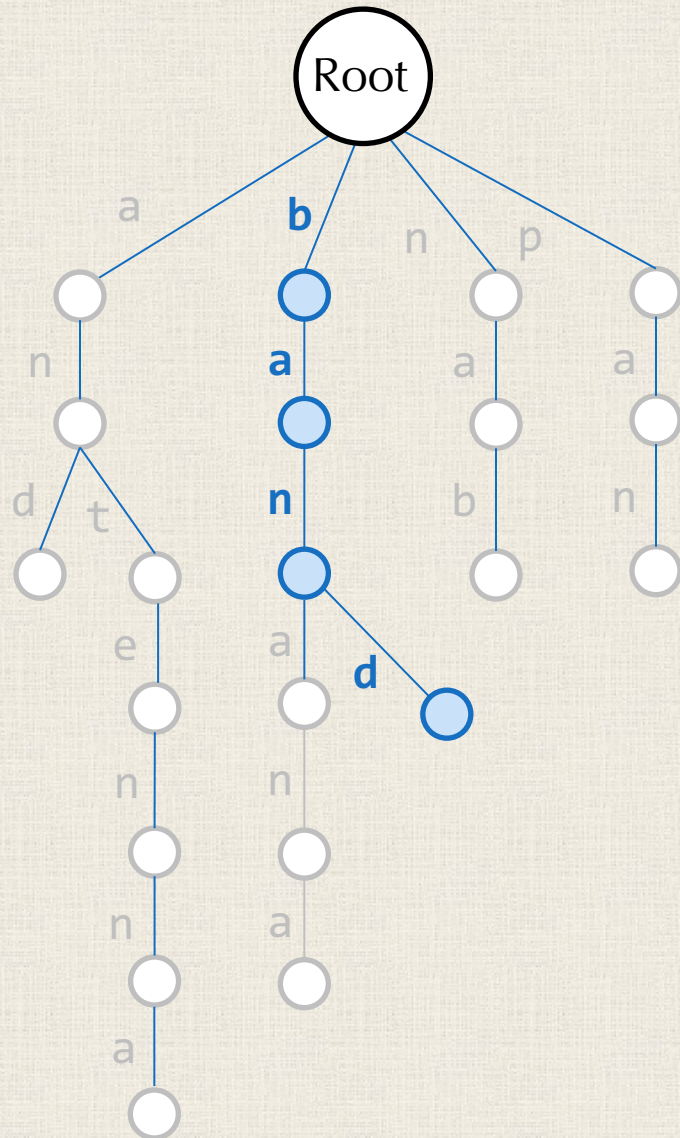
bandana

ananas

nana



Patterns



banana

pan

and

nab

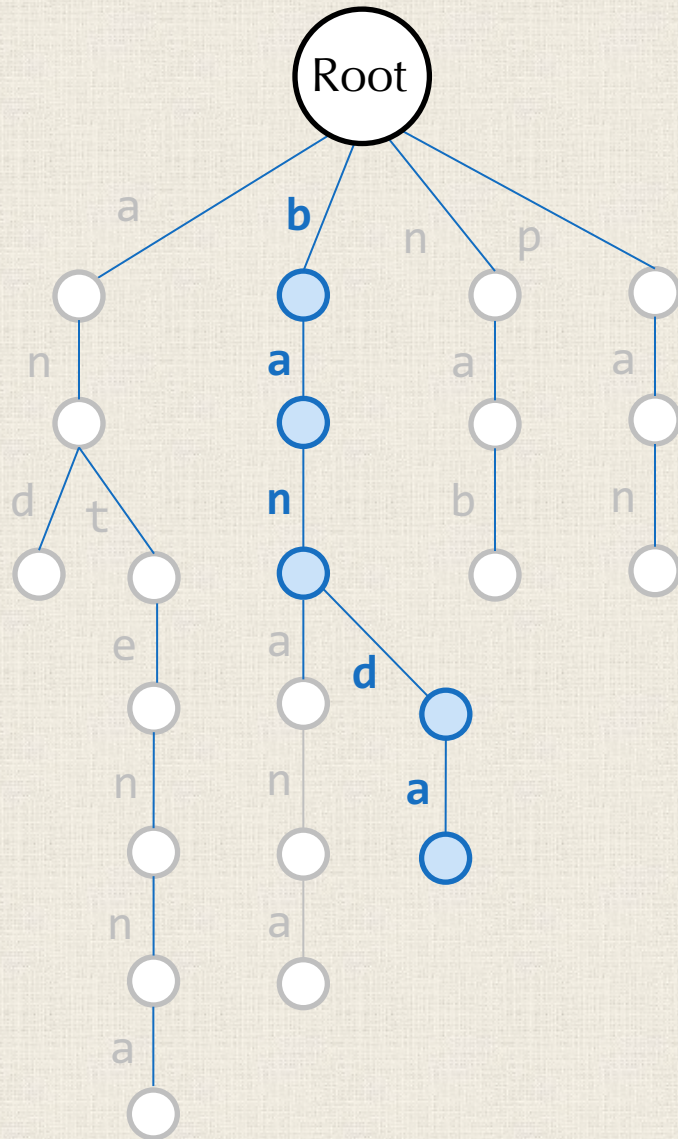
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

antenna

bandana

ananas

nana

Patterns

banana

pan

and

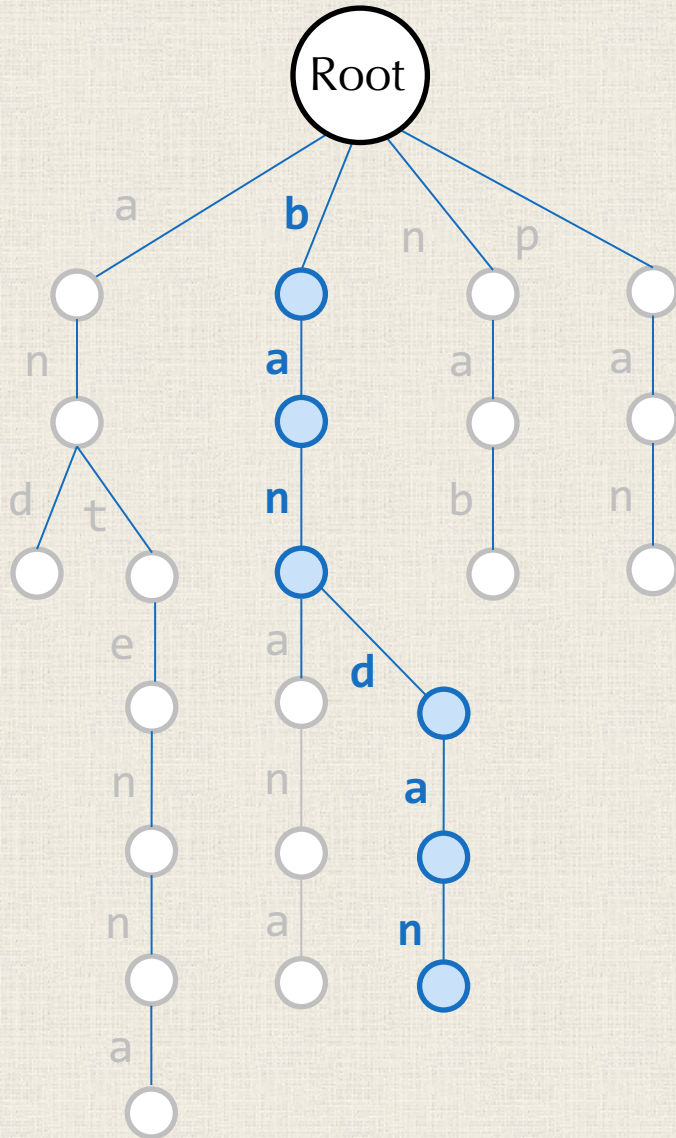
nab

antenna

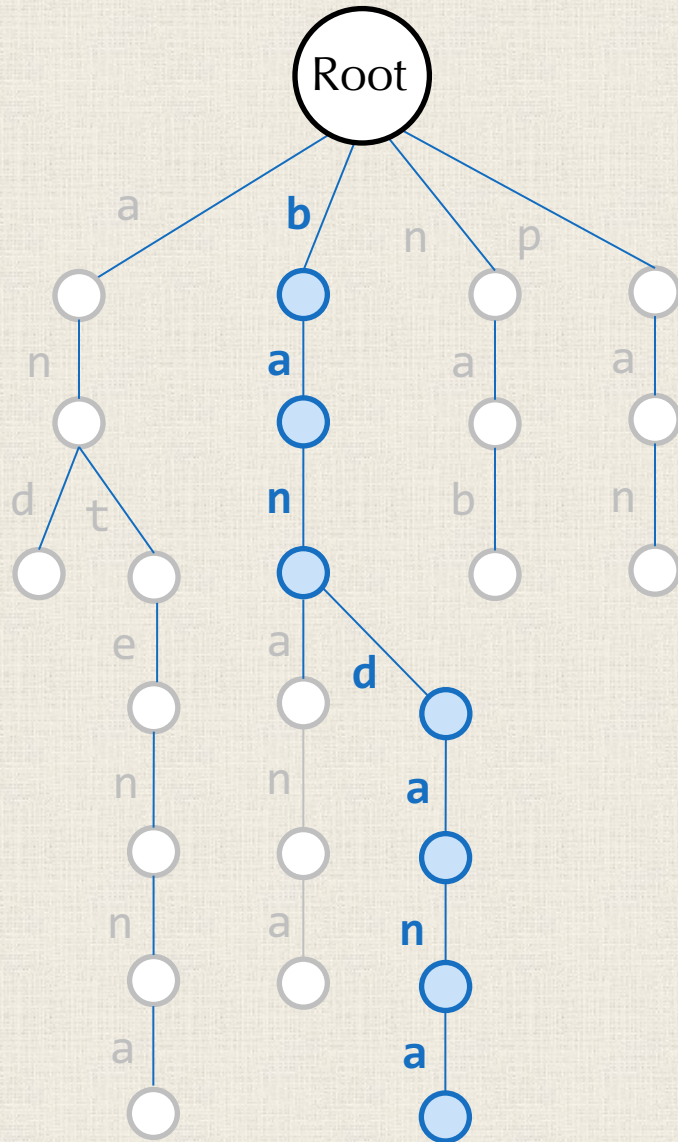
bandana

ananas

nana



Patterns



banana

pan

and

nab

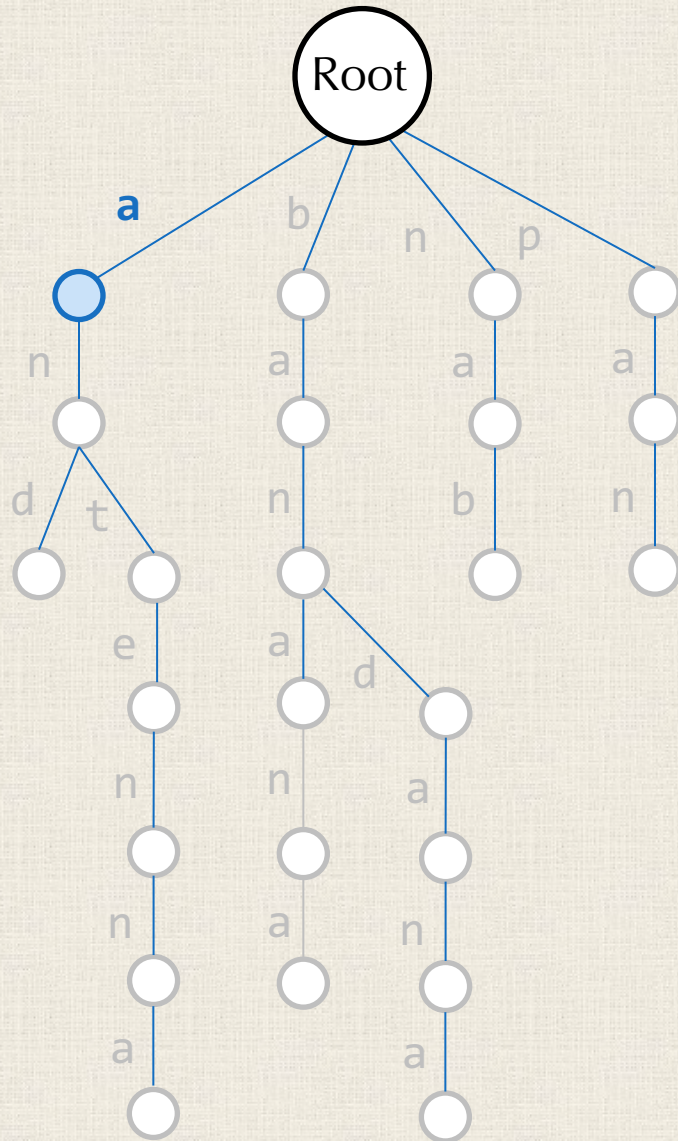
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

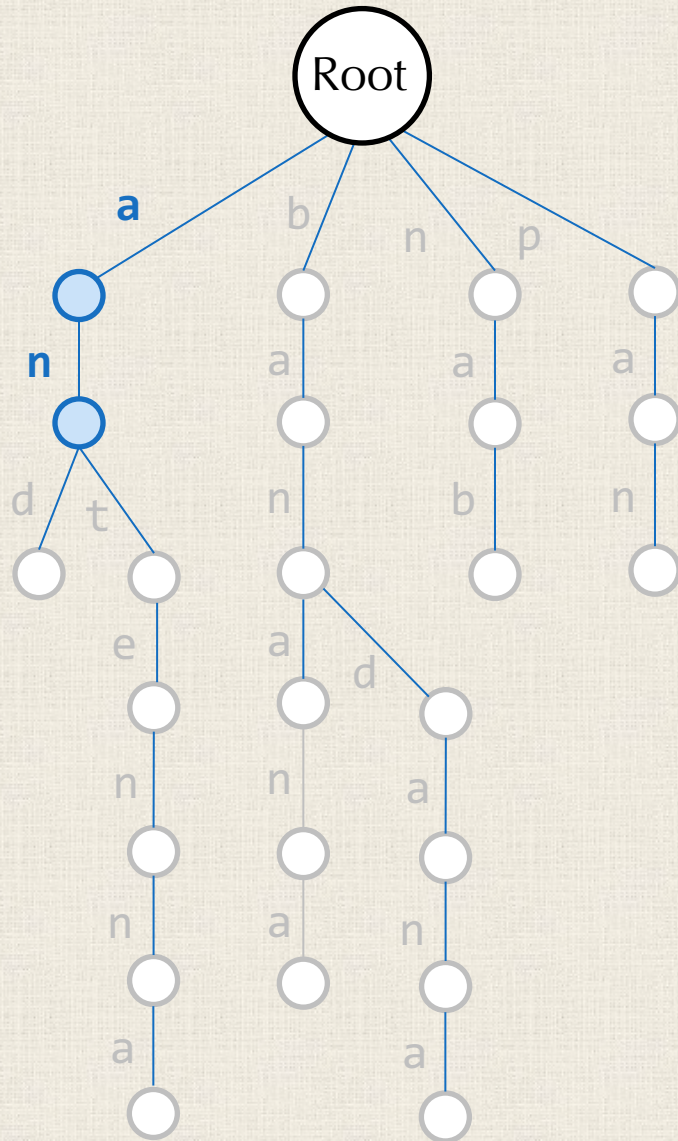
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

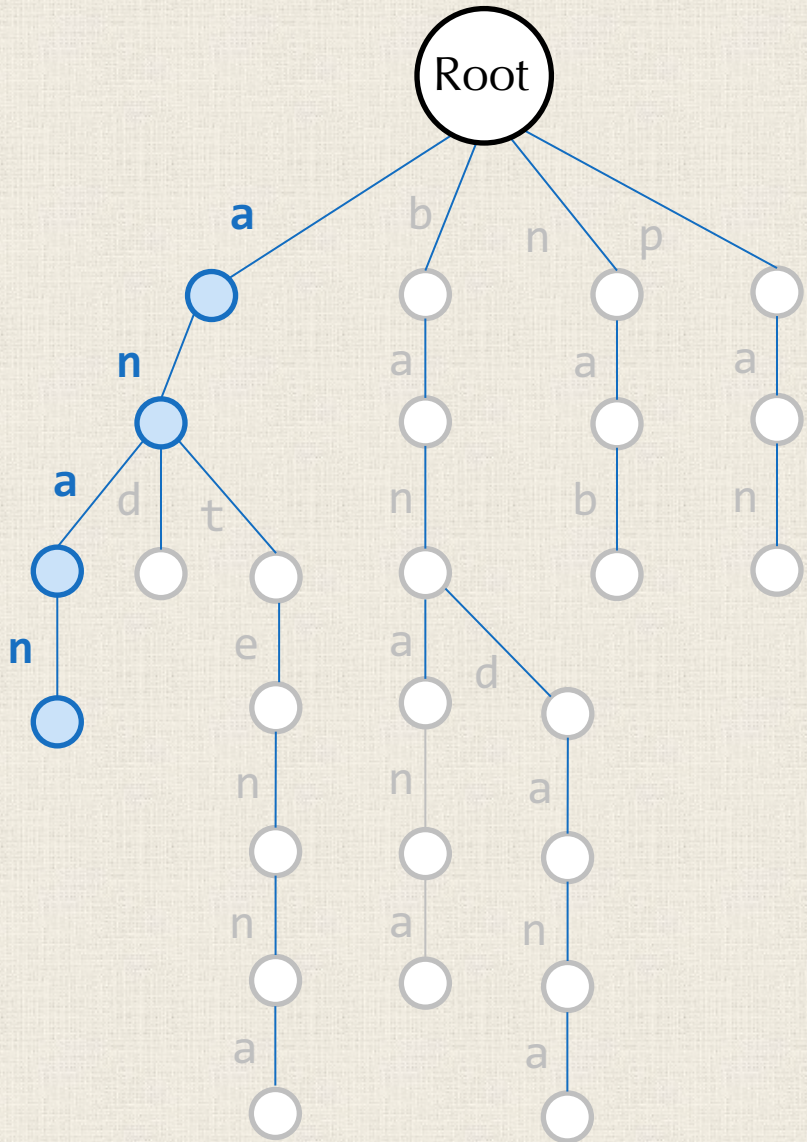
antenna

bandana

ananas

nana

Patterns



banana

pan

and

nab

antenna

bandana

ananas

nana

Patterns

banana

pan

and

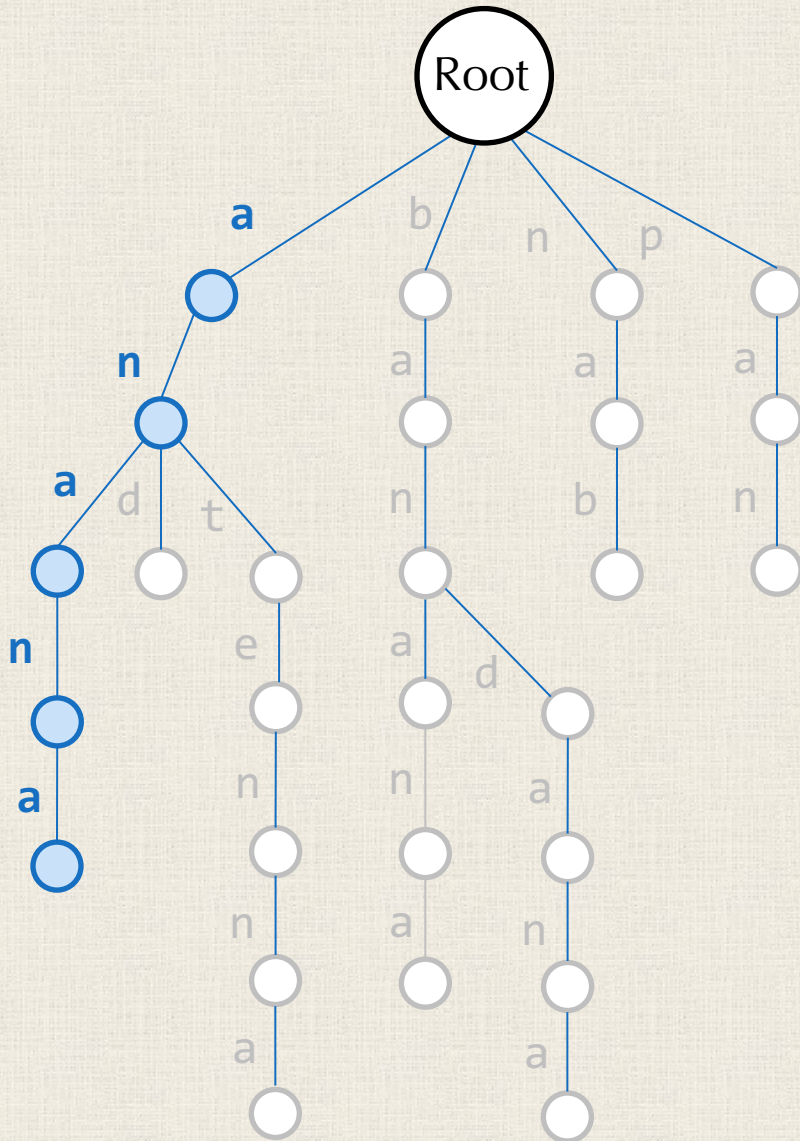
nab

antenna

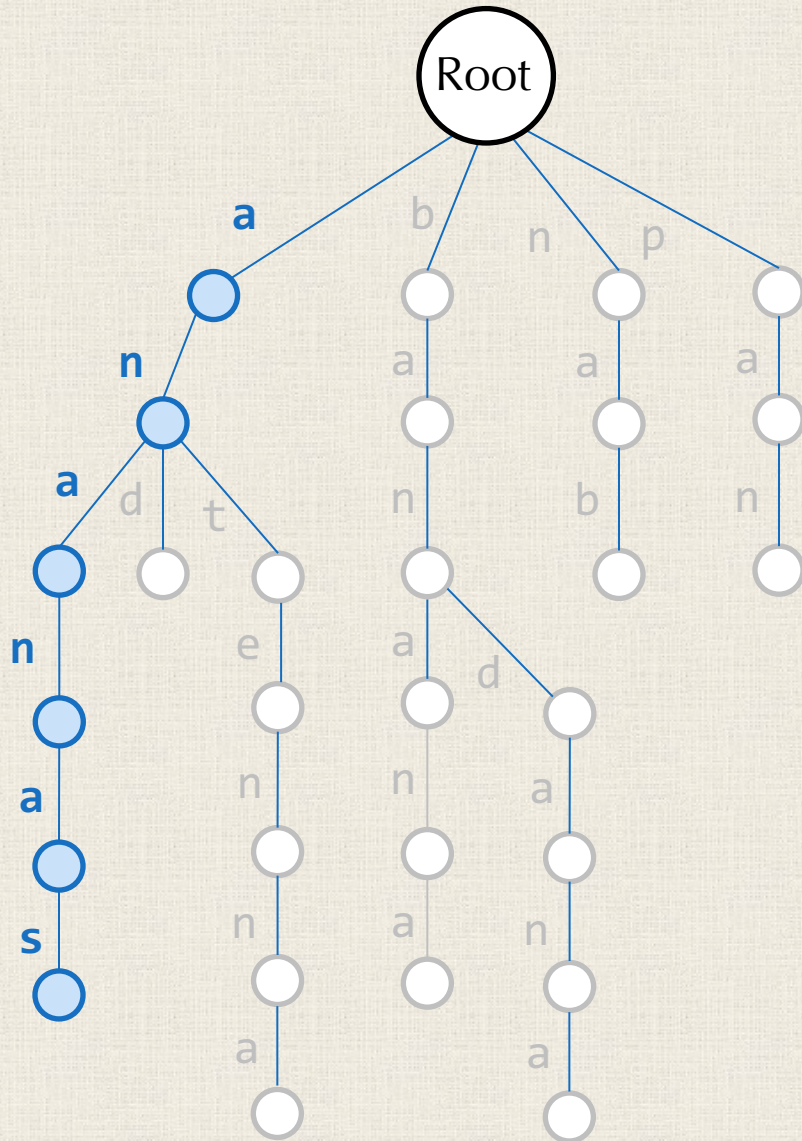
bandana

ananas

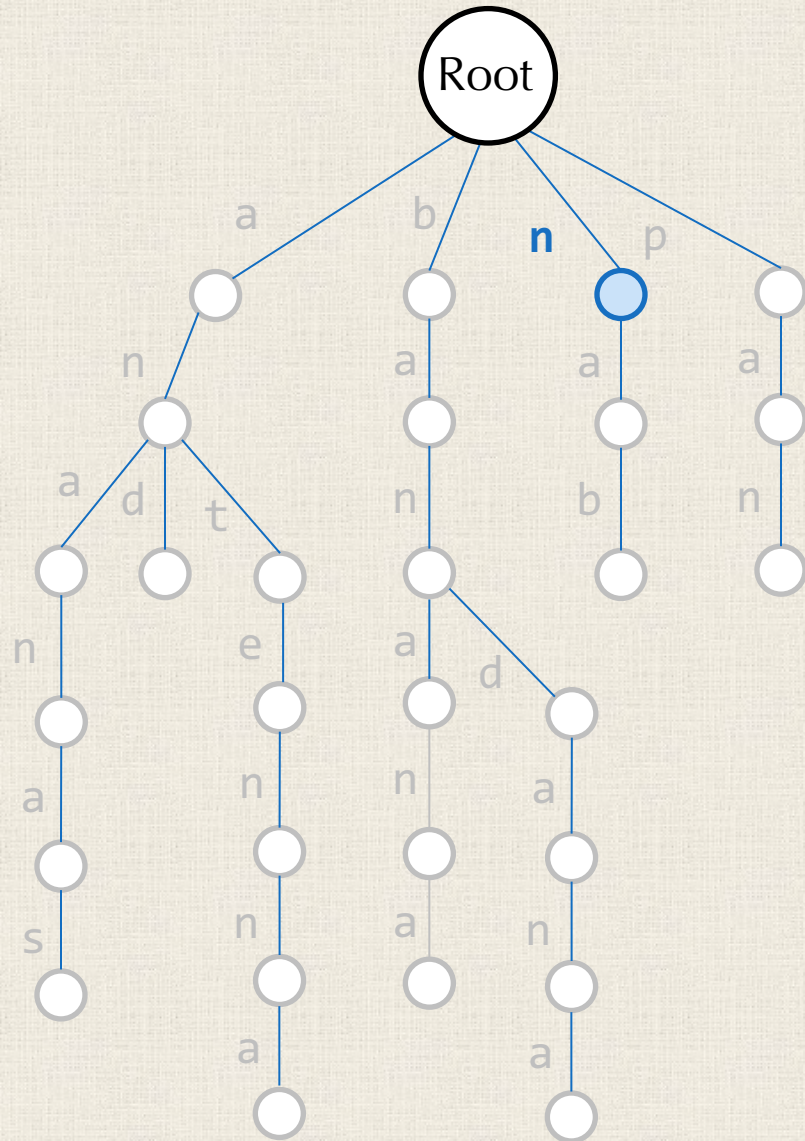
nana



Patterns



banana
pan
and
nab
antenna
bandana
ananas
nana



Patterns

banana

pan

and

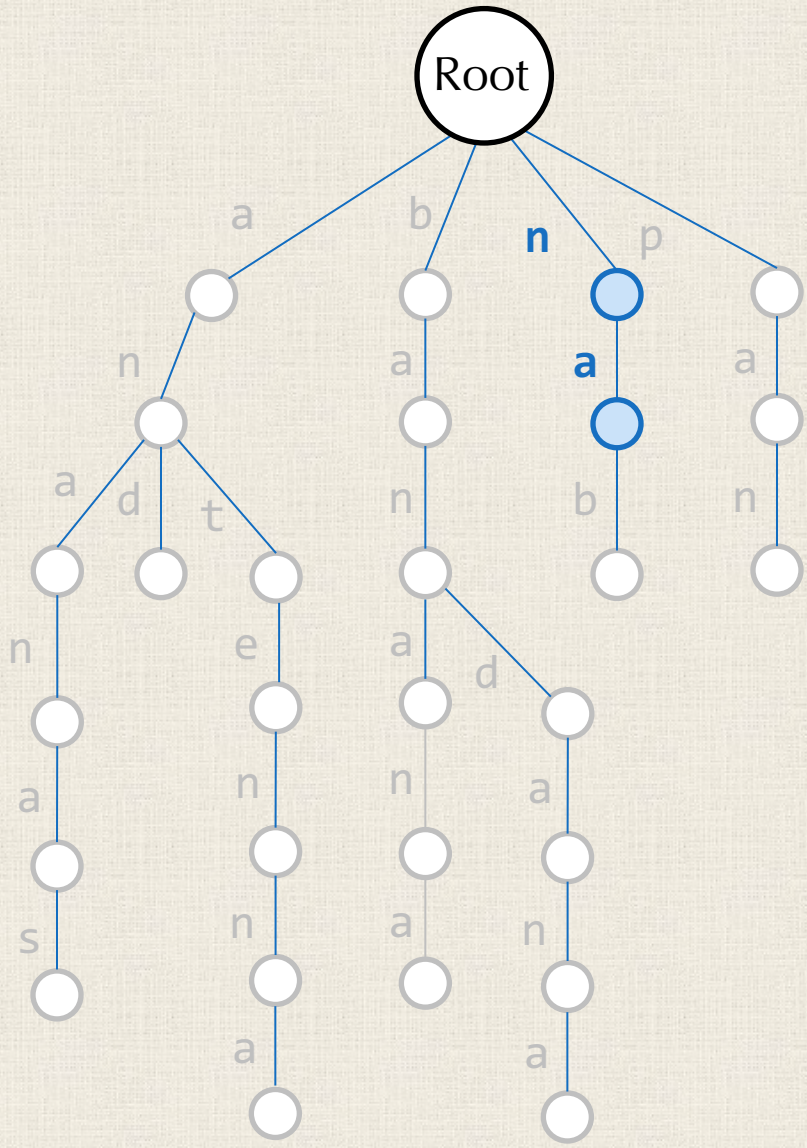
nab

antenna

bandana

ananas

nana



Patterns

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

Patterns

banana

pan

and

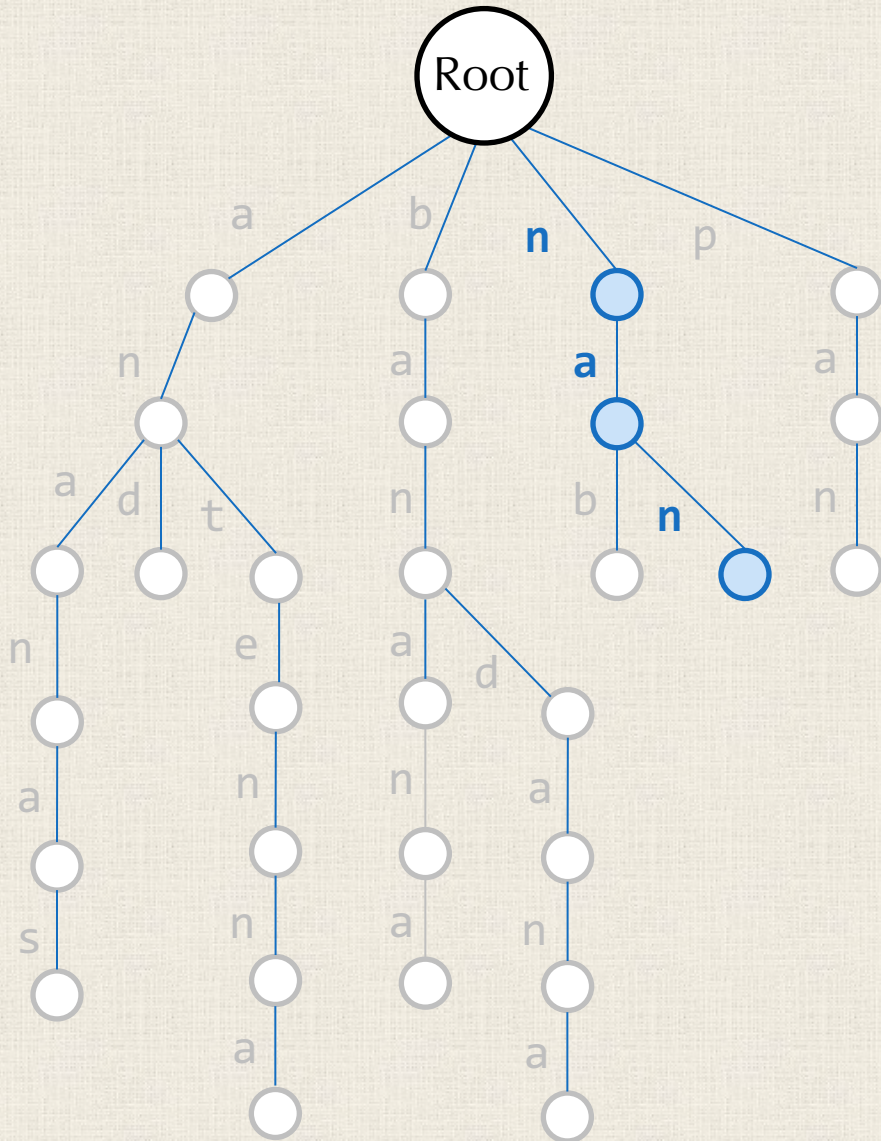
nab

antenna

bandana

ananas

nana



Patterns

banana

pan

and

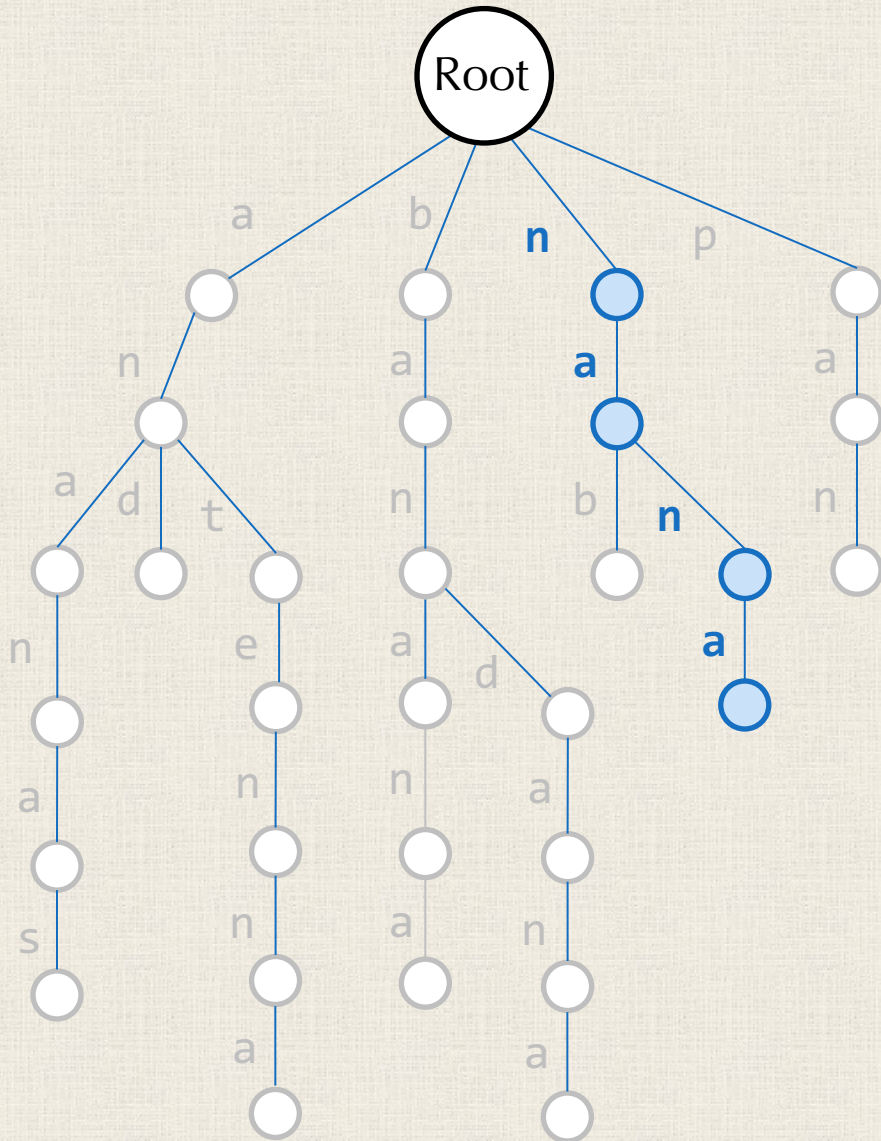
nab

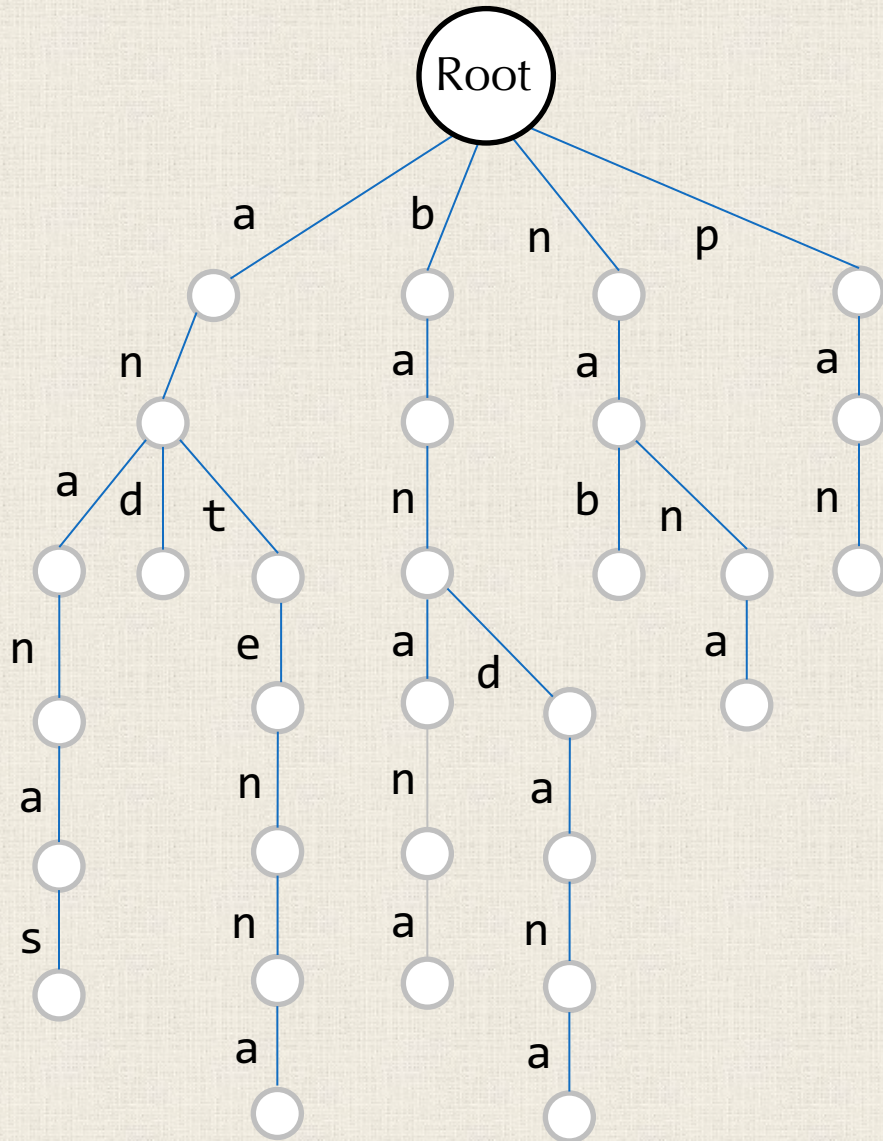
antenna

bandana

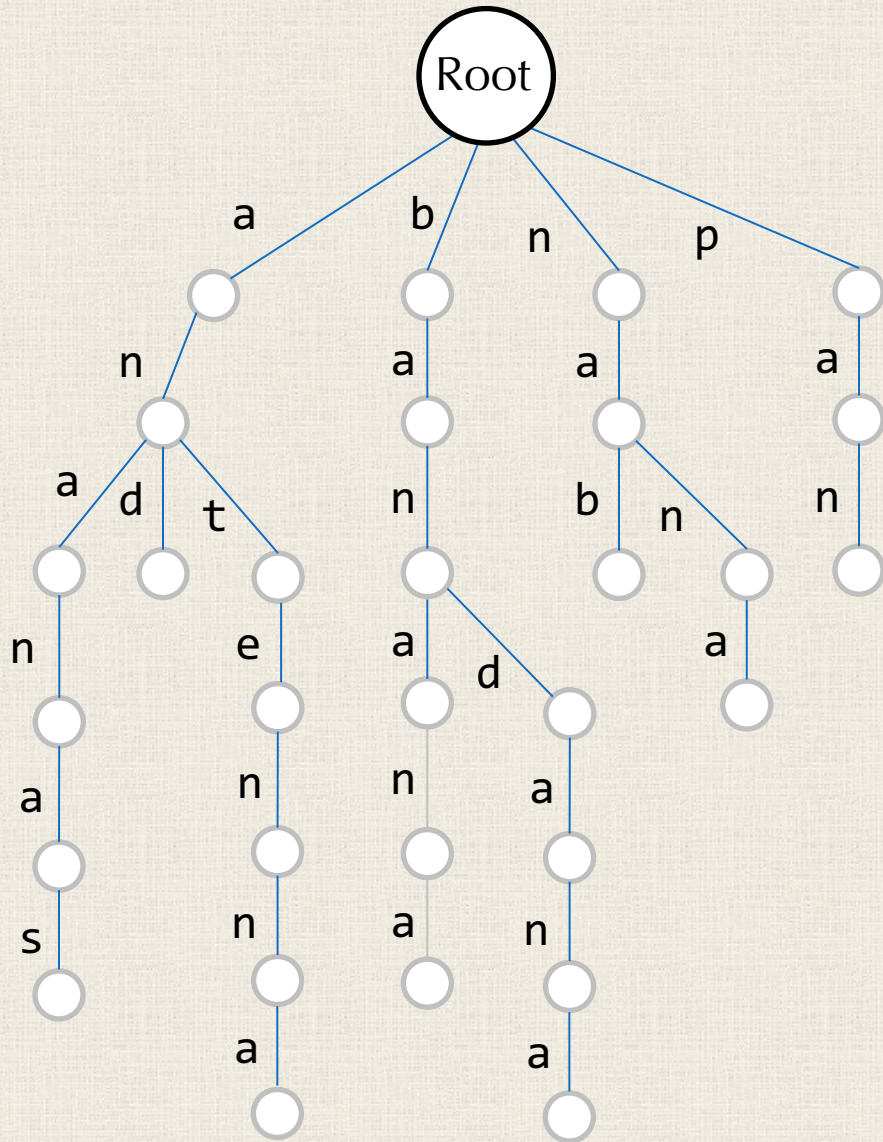
ananas

nana





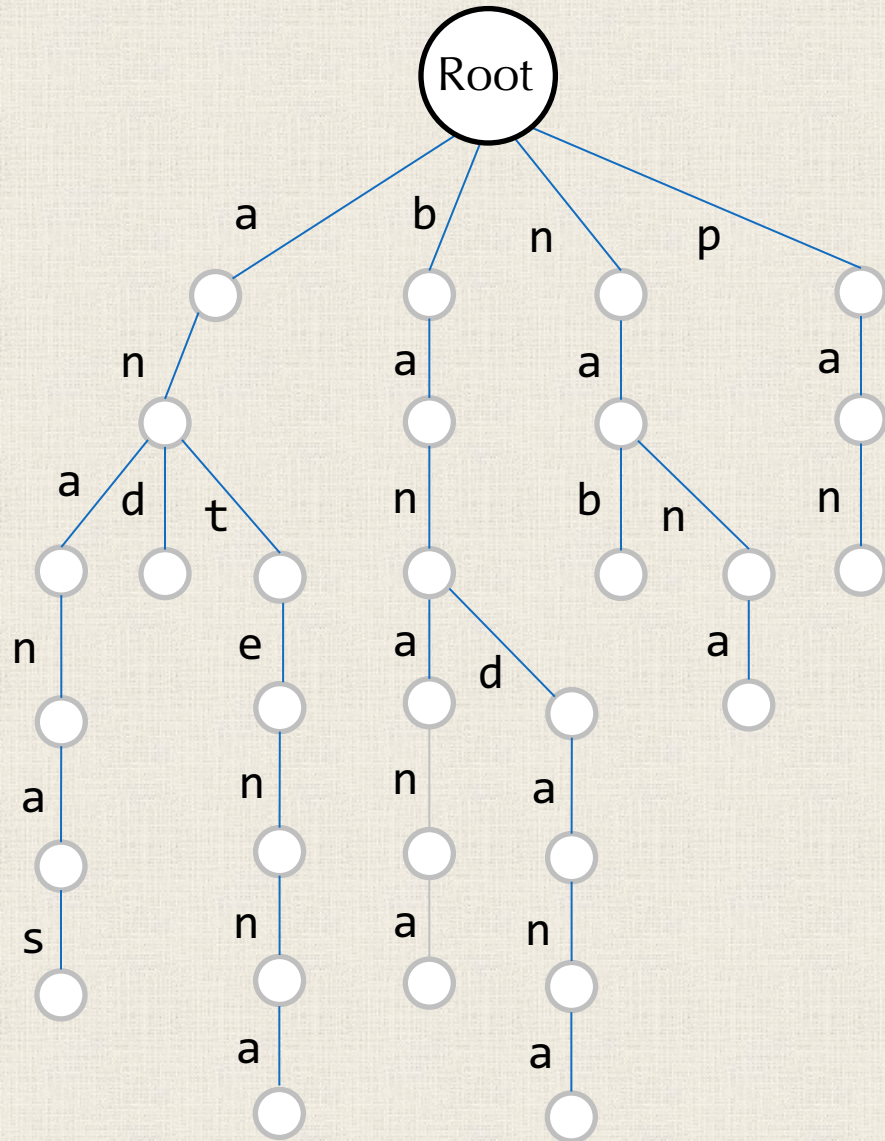
STOP: How can we use the trie for pattern matching? (e.g, say *Text* = “panamabananas”).



Answer: We can test whether any string in *Patterns* matches a prefix of *Text* by "sliding it" down the *Text*.

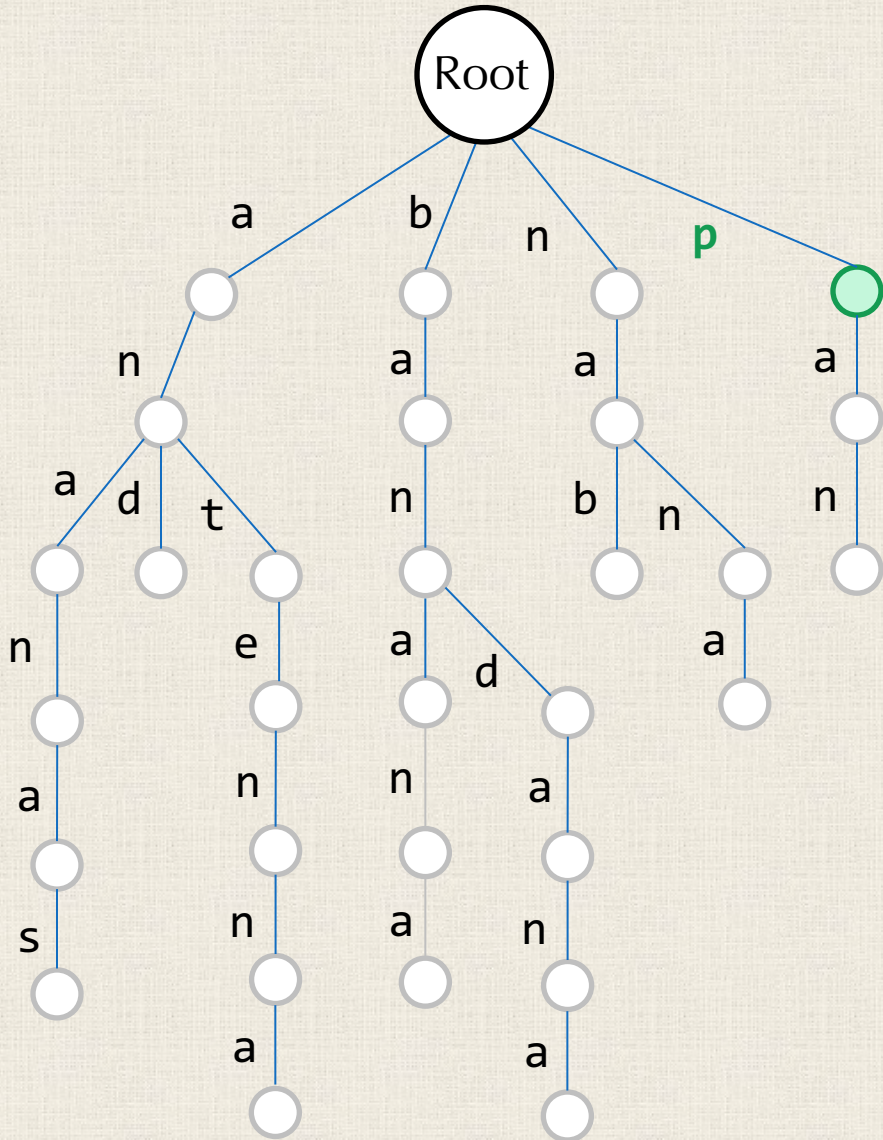


panamabanas



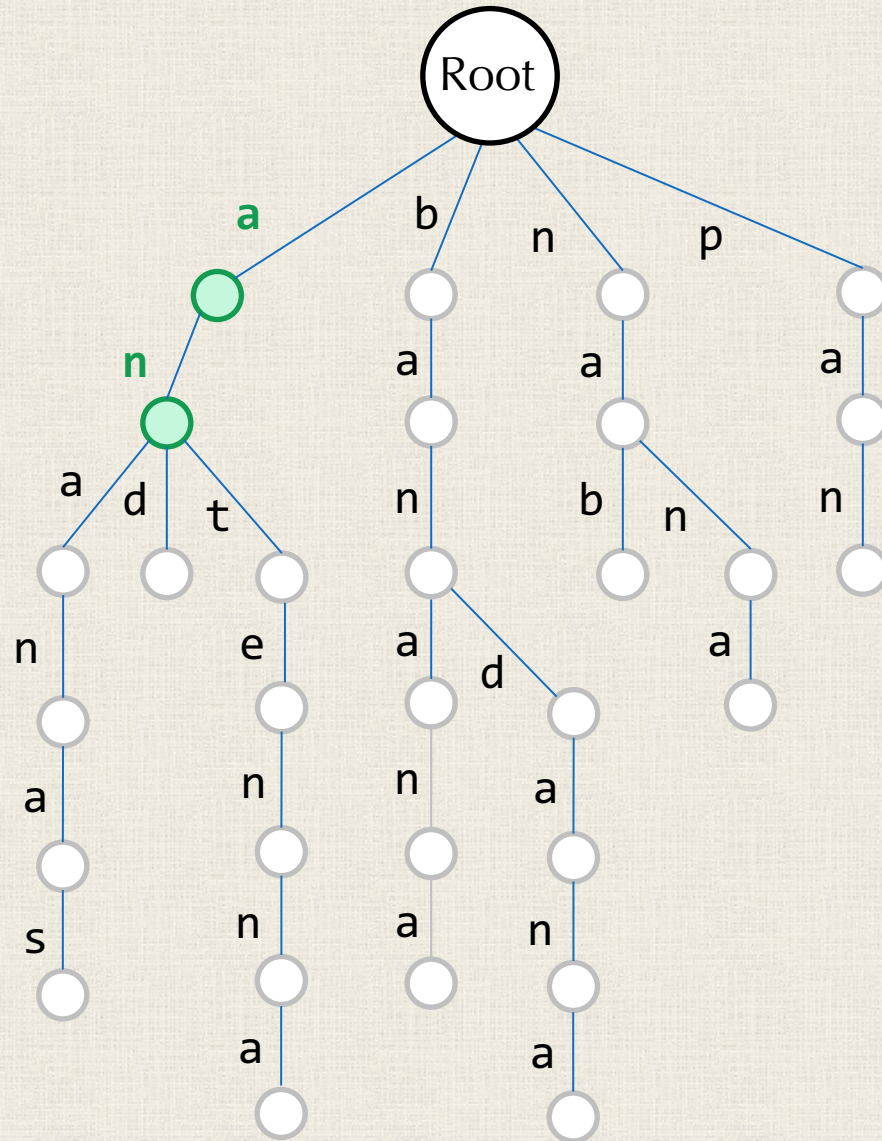
Pattern	Positions
---------	-----------

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



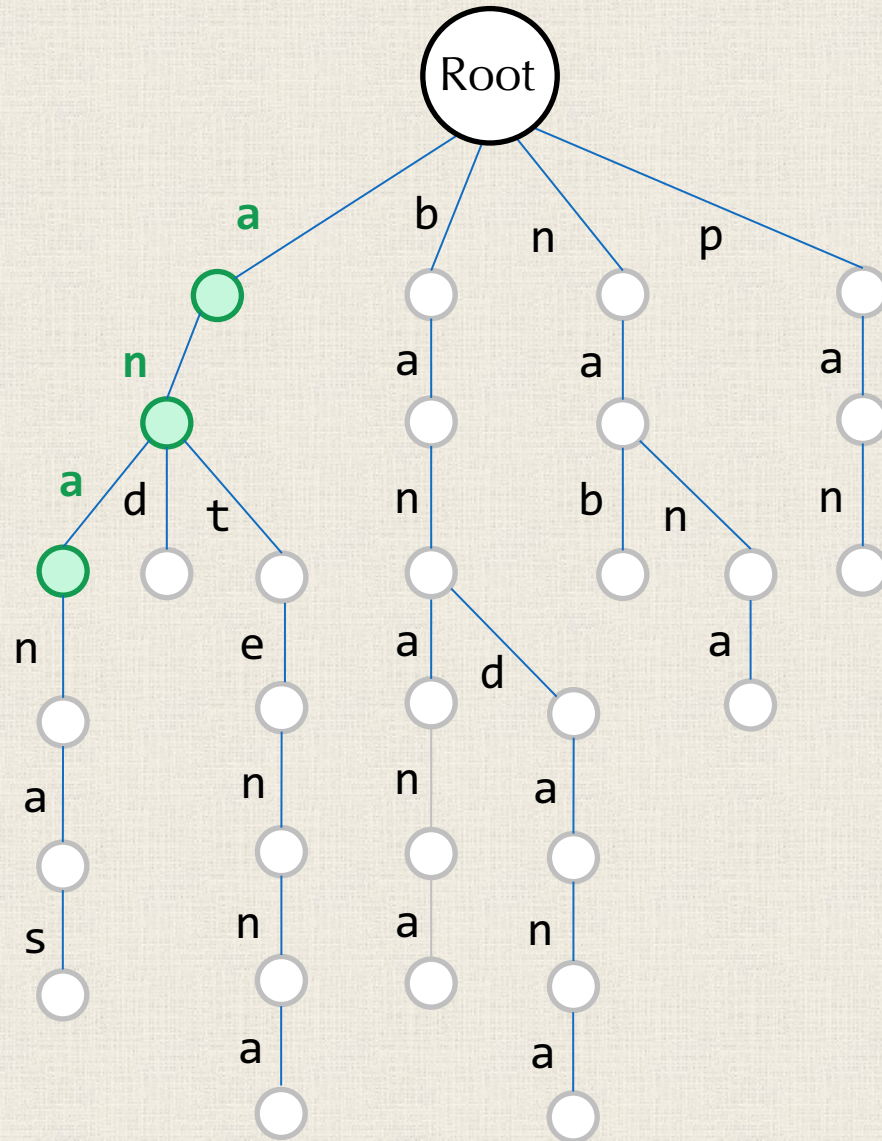
Pattern	Positions
---------	-----------

panamabanas



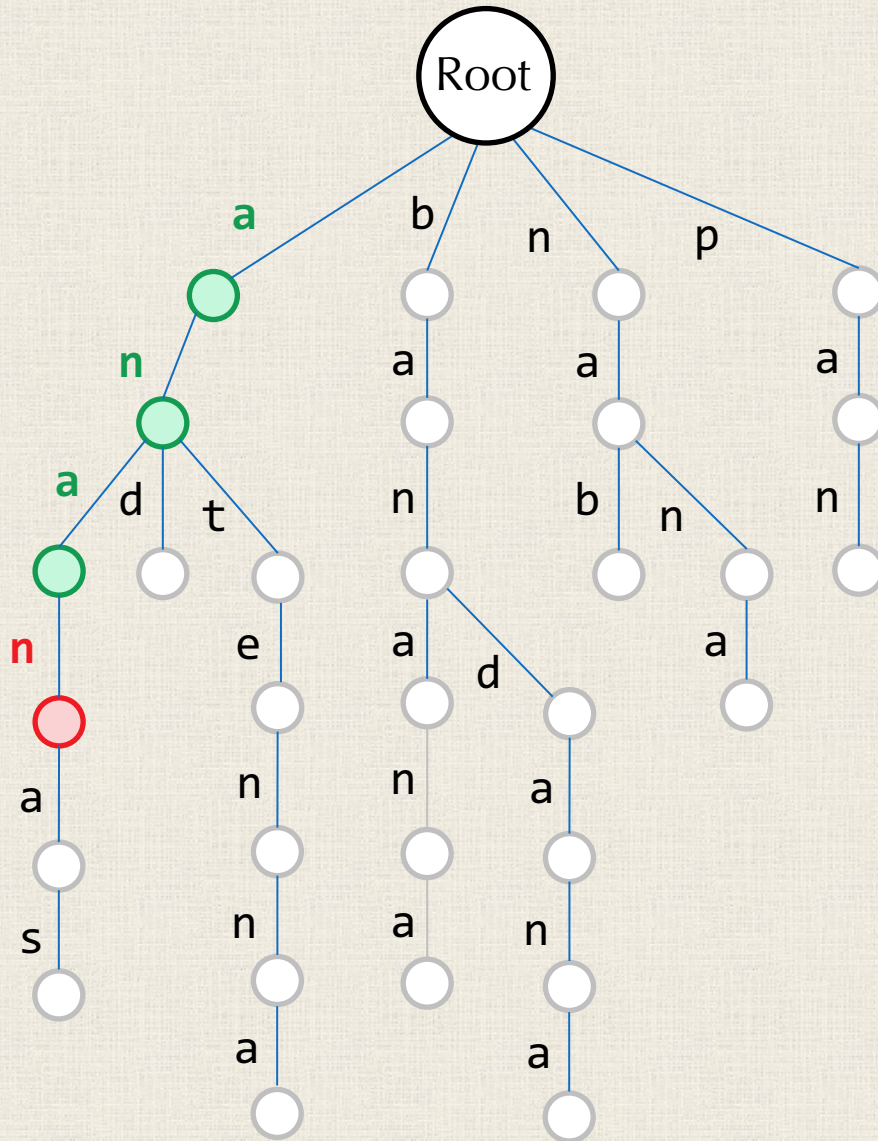
Pattern	Positions
pan	\emptyset

panamabananas



Pattern	Positions
pan	\emptyset

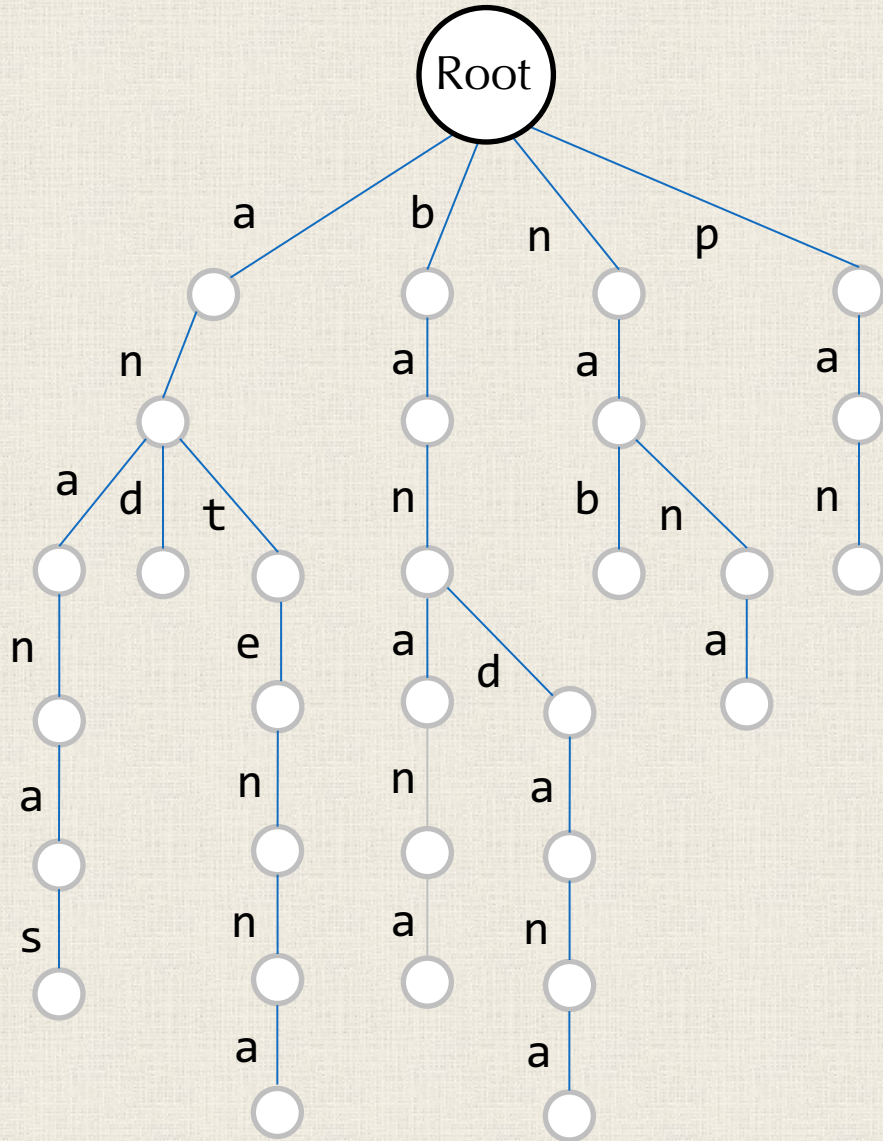
p a n a m a b a n a n a s



By how the trie is organized, we can be sure that there are no patterns matching this position of *Text*.

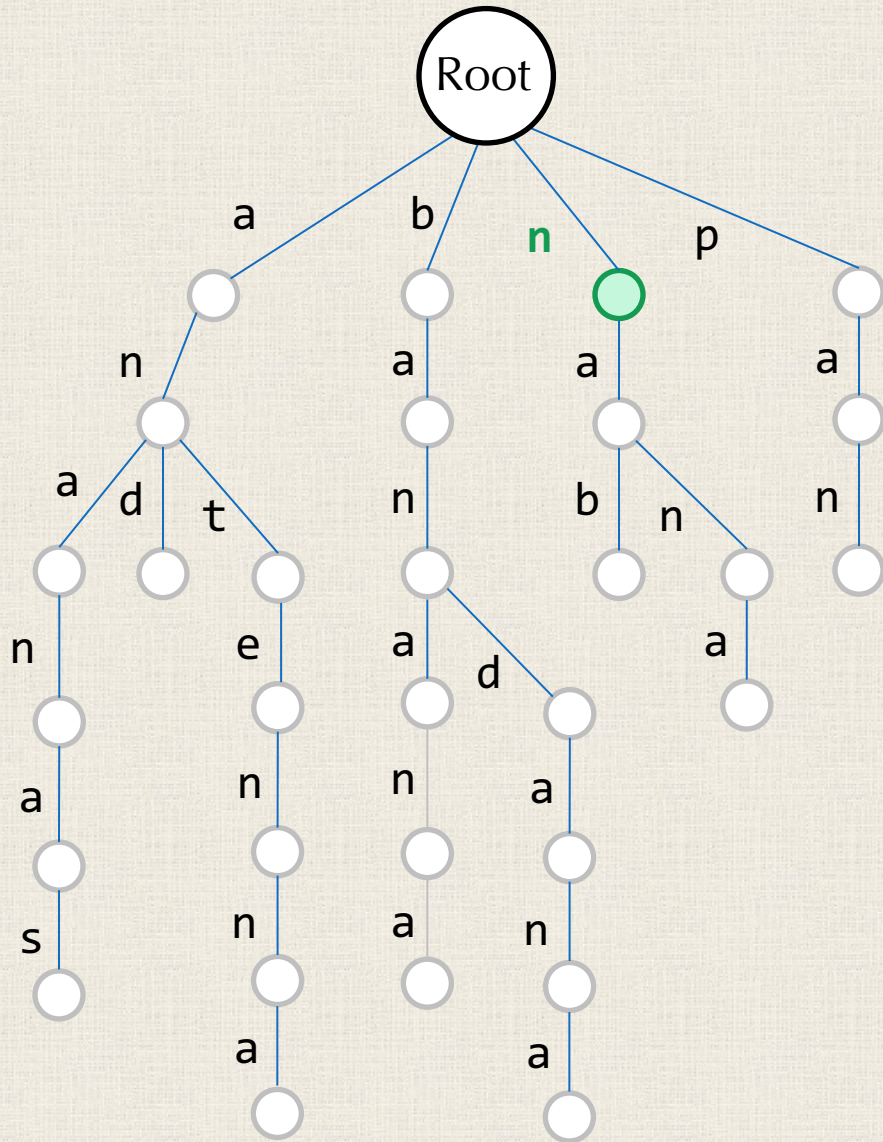
Pattern	Positions
pan	\emptyset

panamabanas



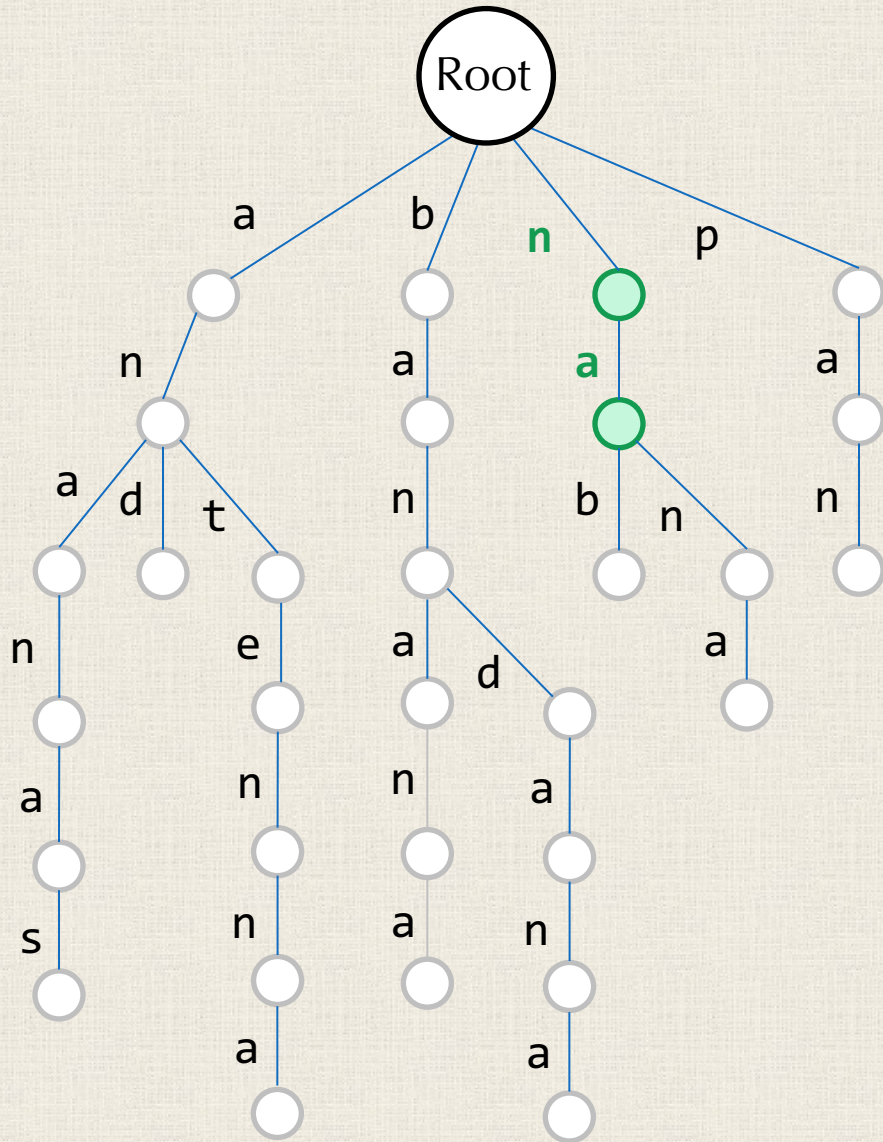
Pattern	Positions
pan	\emptyset

panabanas



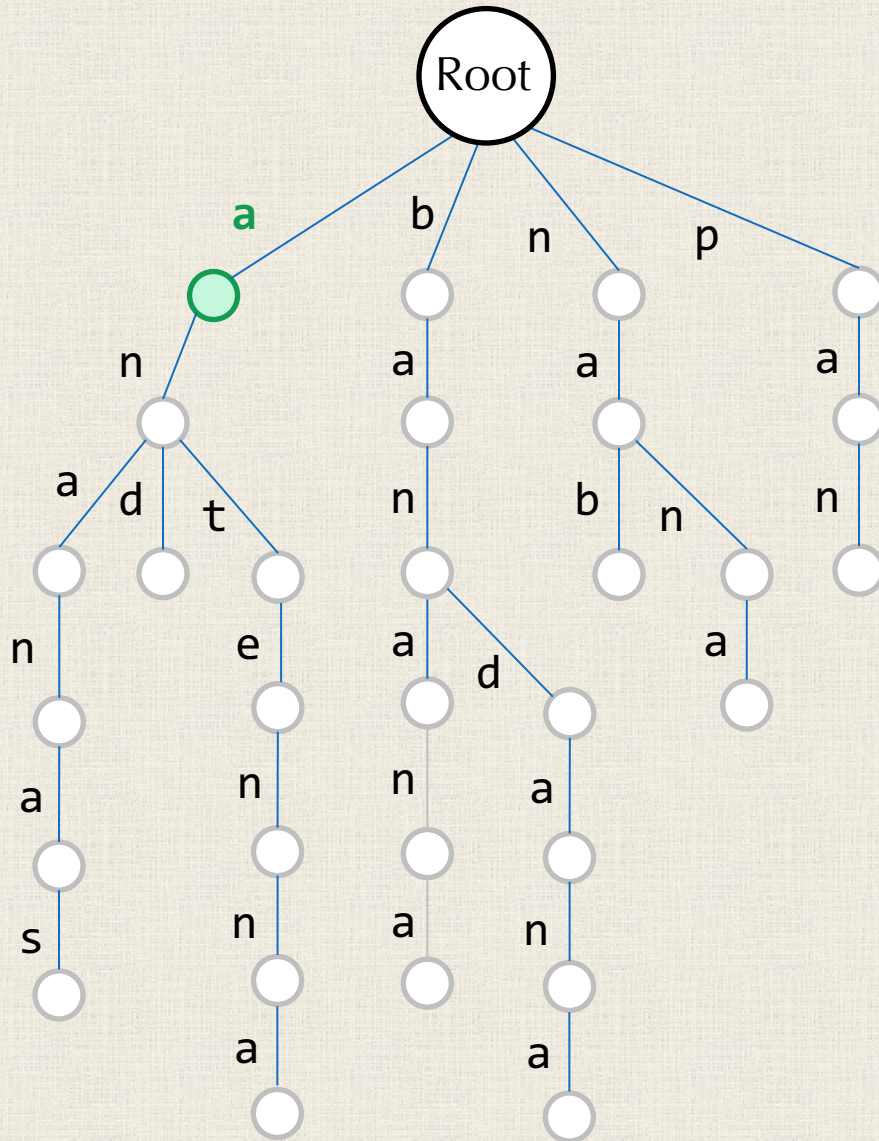
Pattern	Positions
pan	\emptyset

panabanas



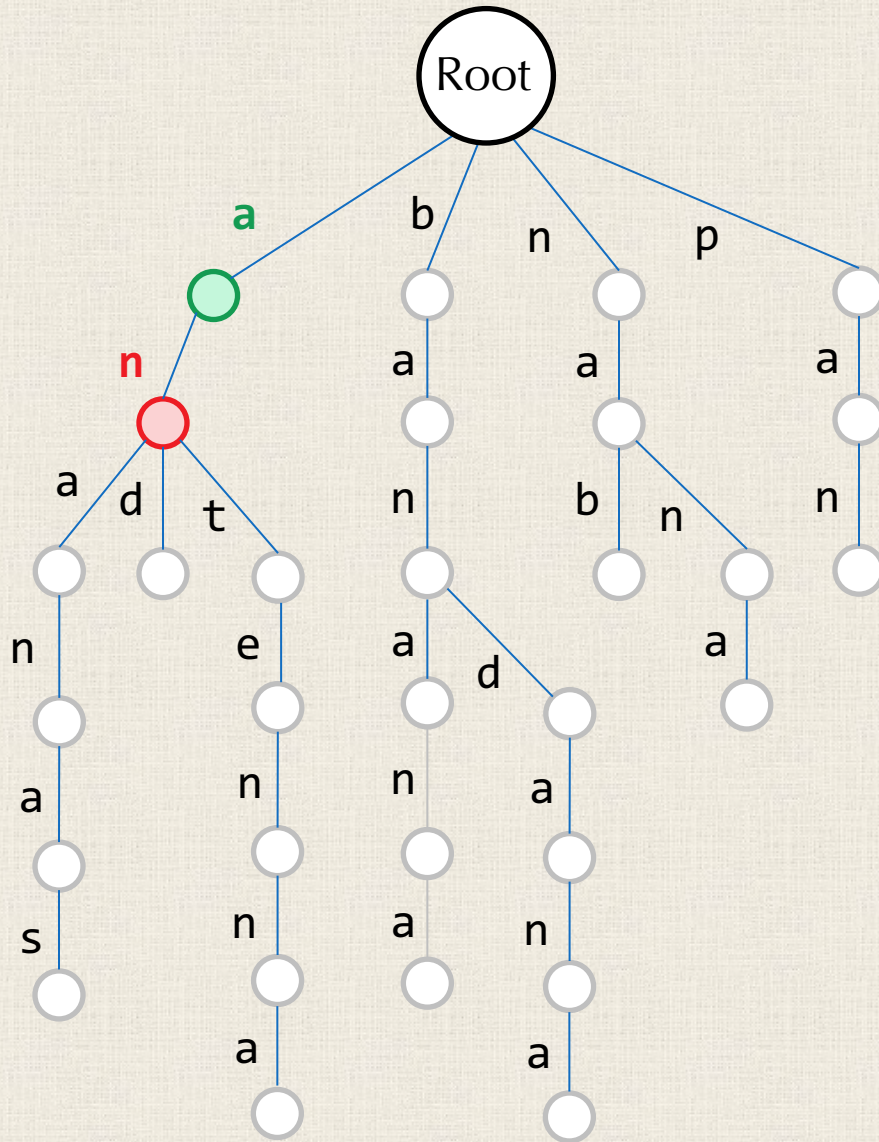
Pattern	Positions
pan	\emptyset

pan**a**mabanas



Pattern	Positions
pan	∅

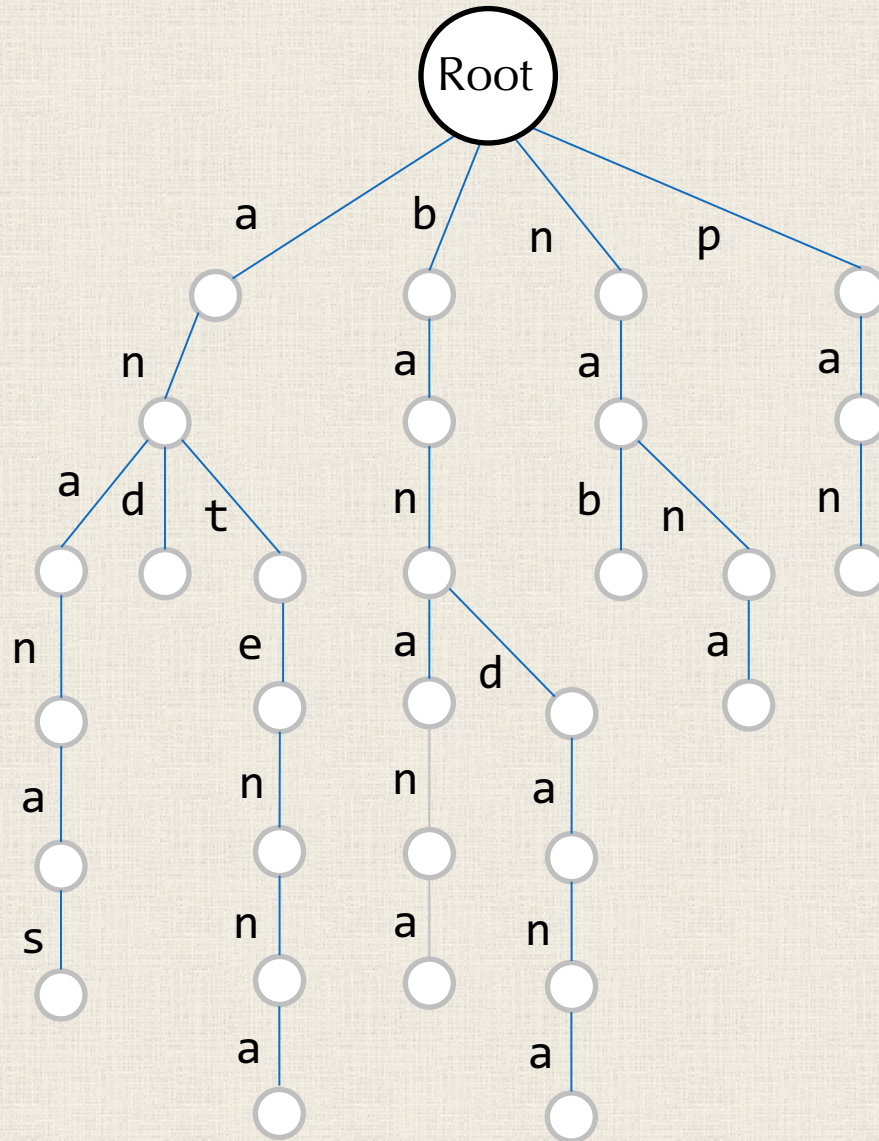
panamabananas



Pattern	Positions
---------	-----------

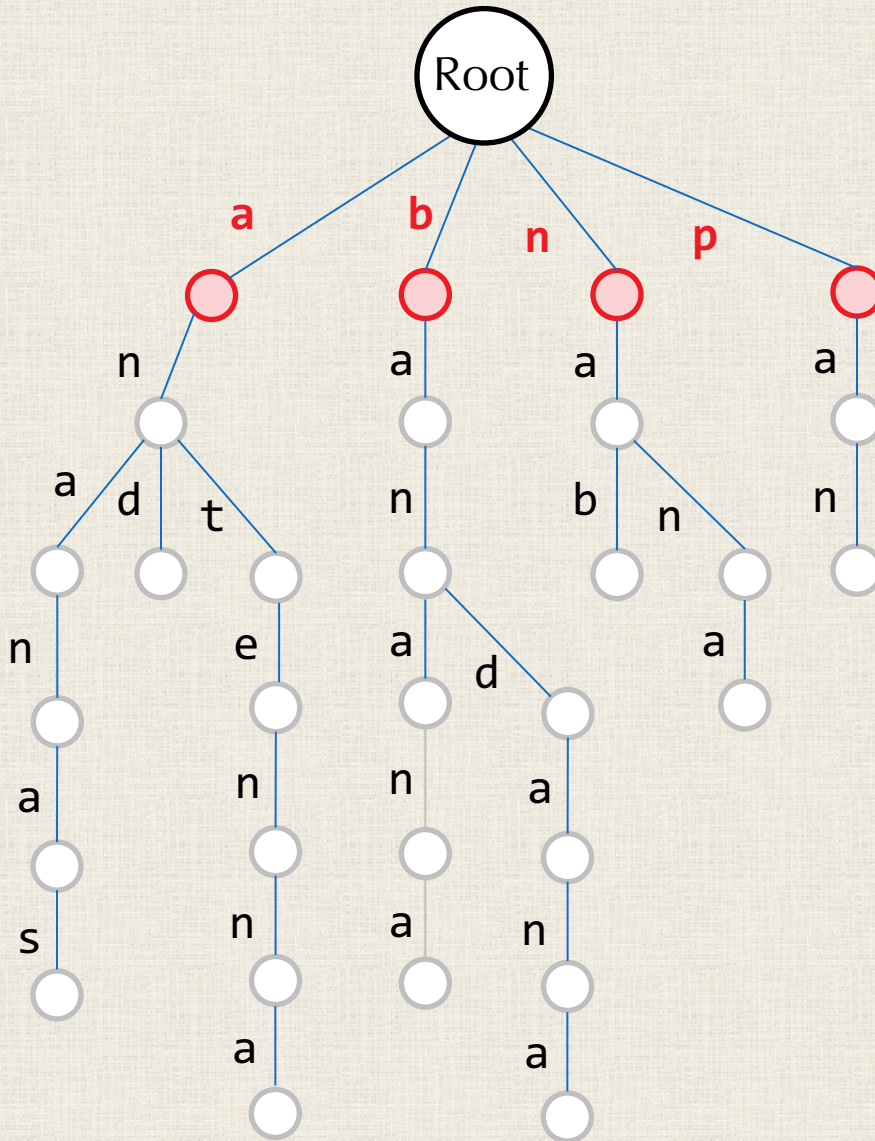
pan	∅
-----	---

panamabananas



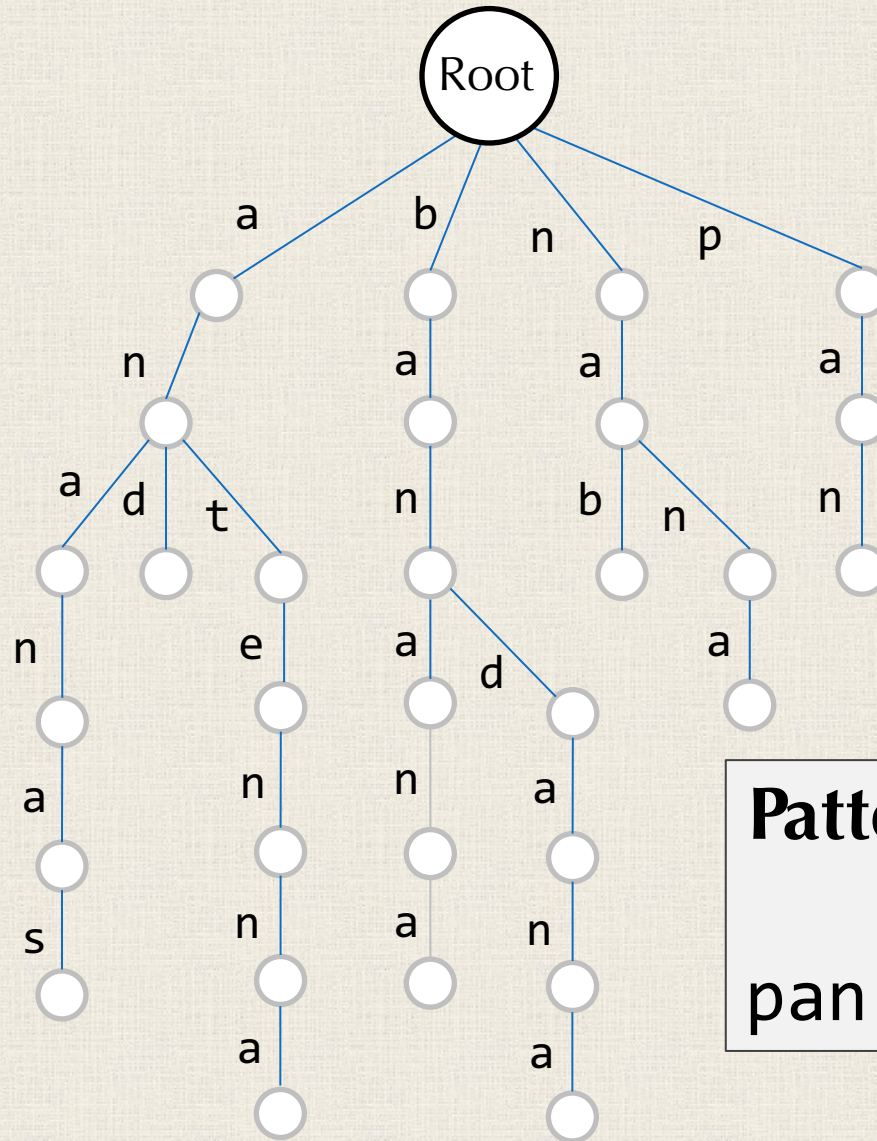
Pattern	Positions
pan	∅

panamabananas



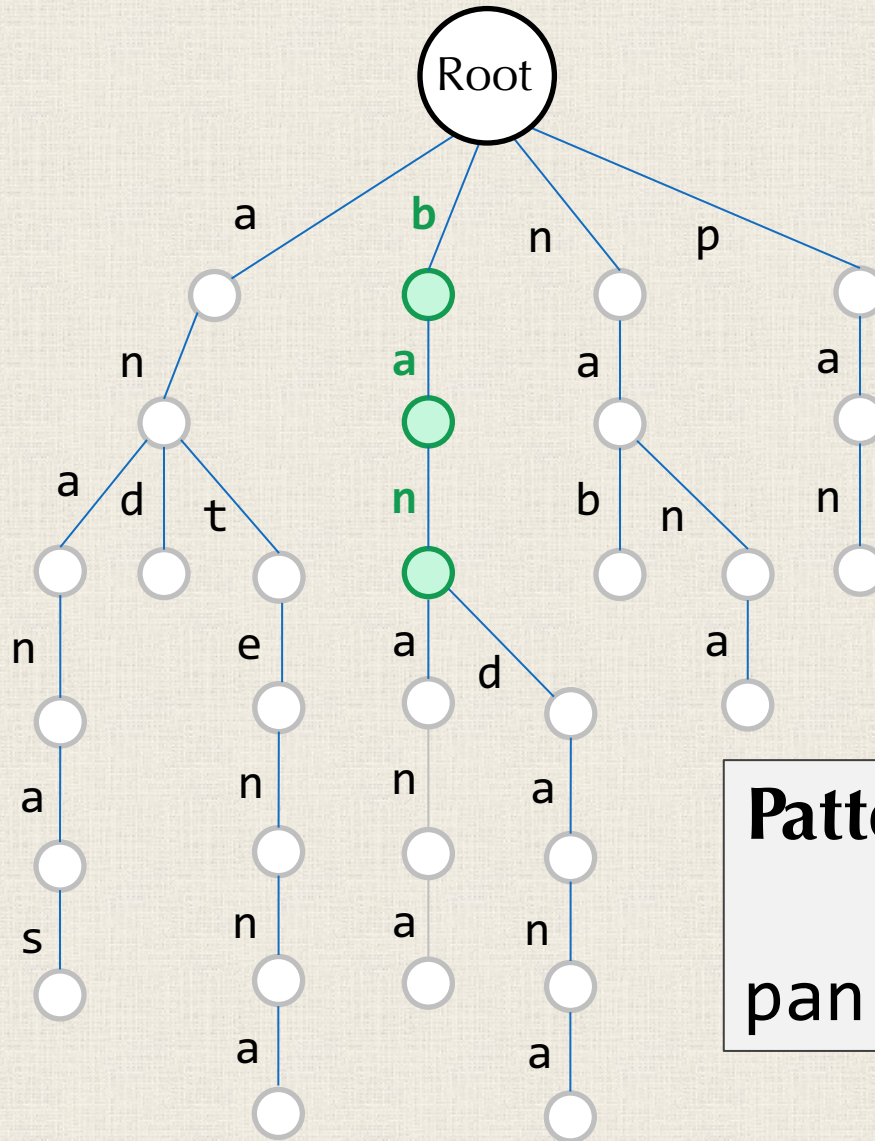
Pattern	Positions
pan	∅

panama bananas



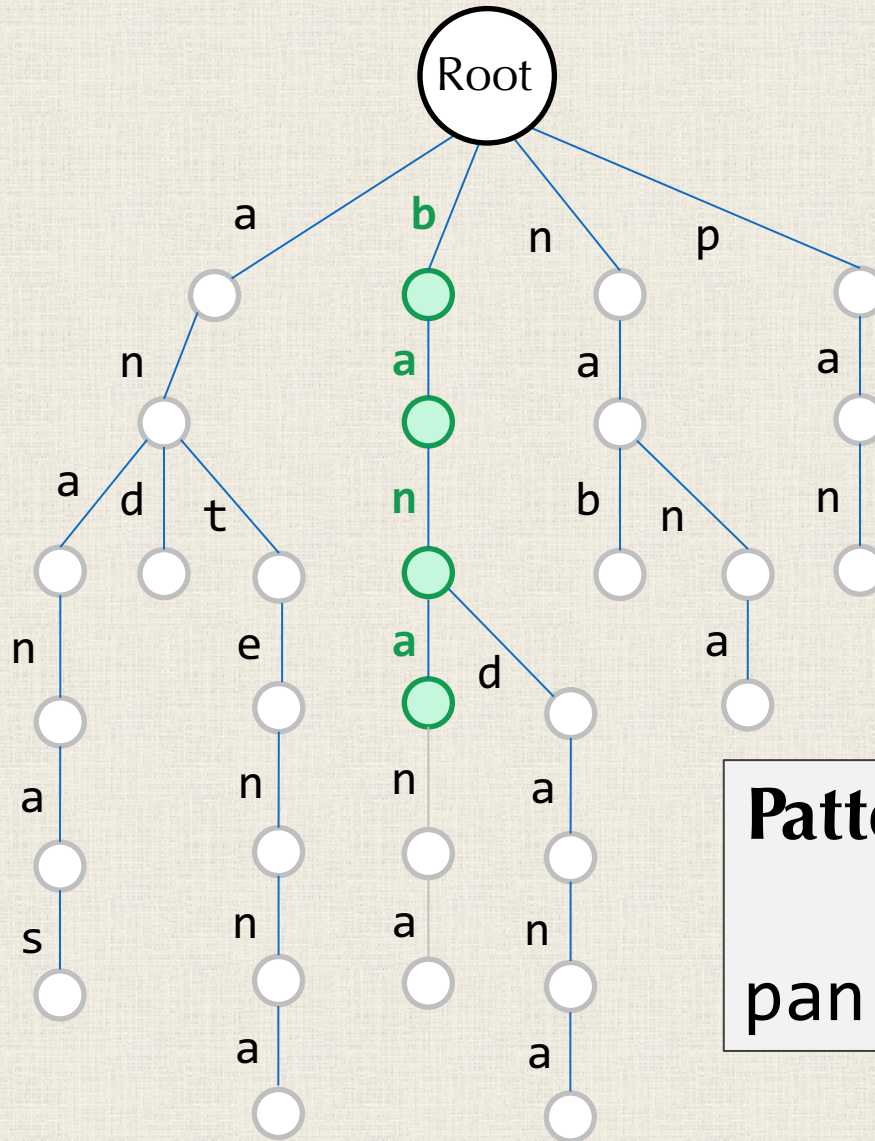
Pattern	Positions
pan	∅

panama bananas



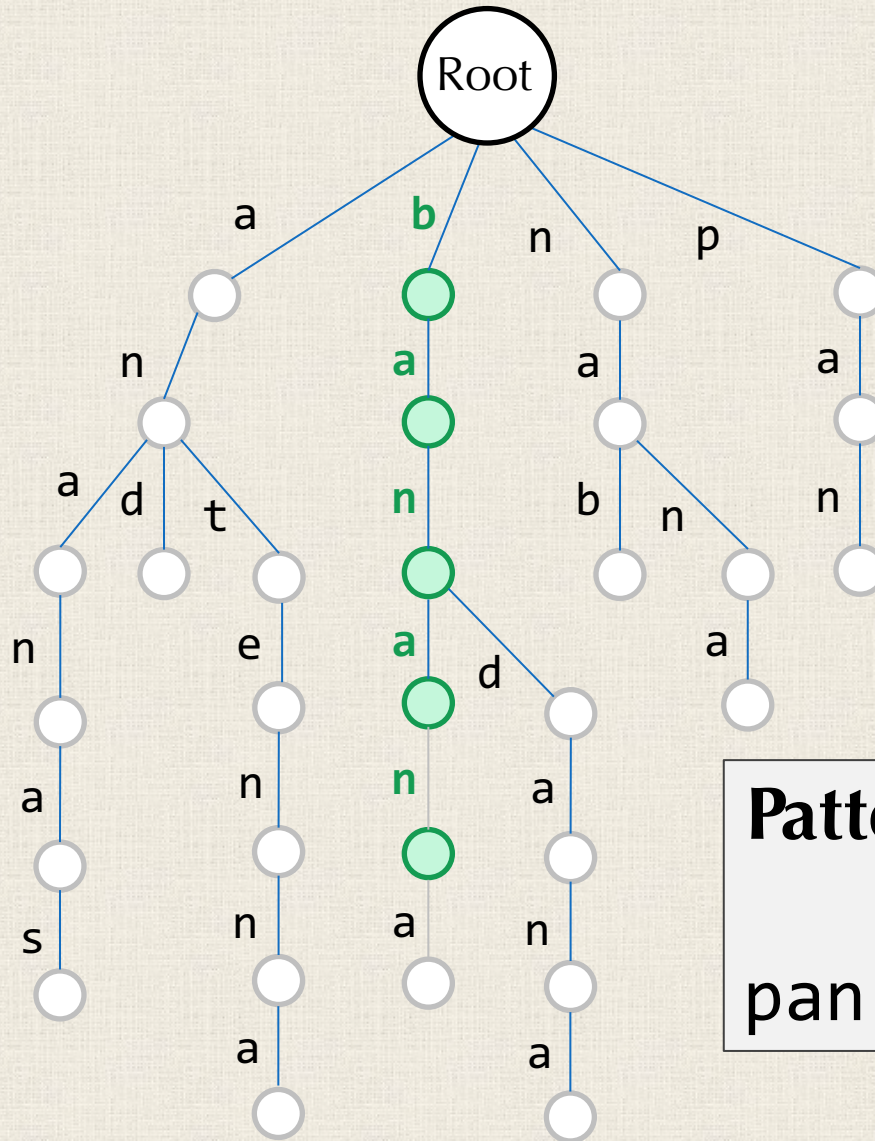
Pattern	Positions
pan	\emptyset

panama bananas



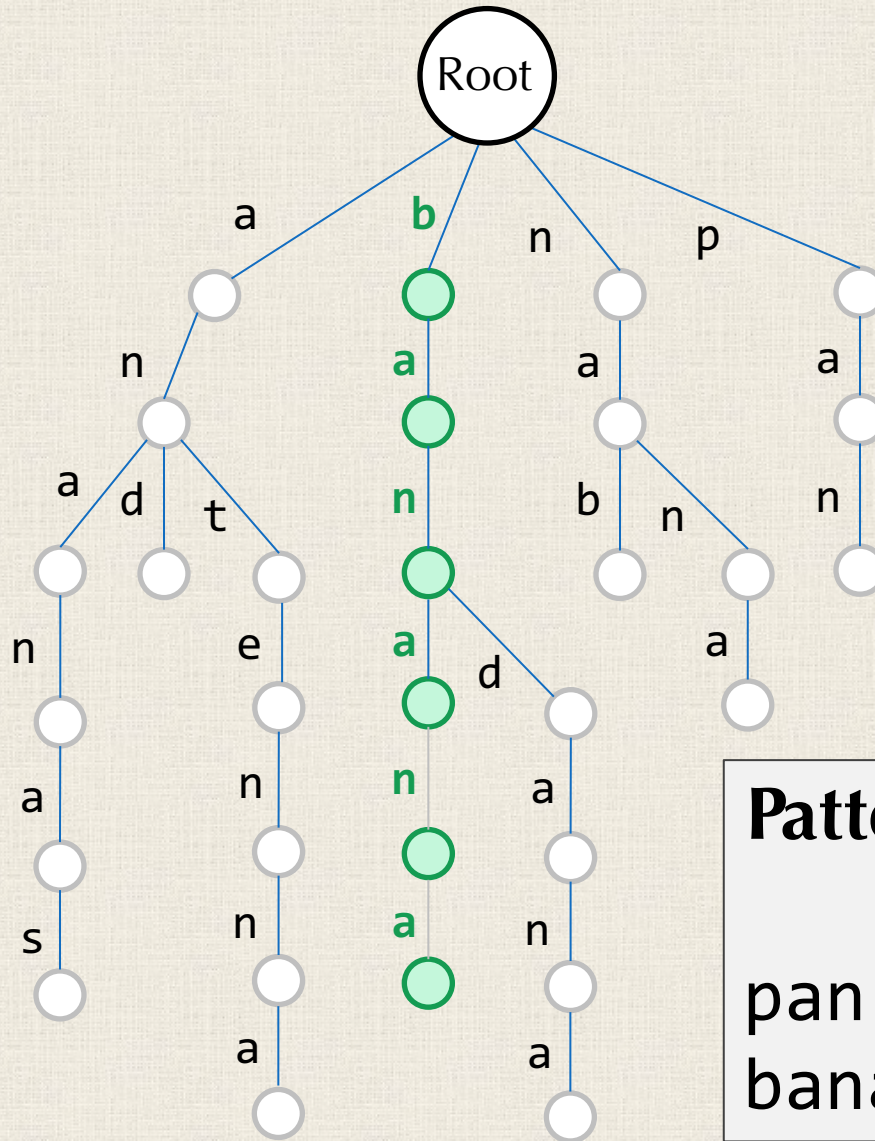
Pattern	Positions
pan	\emptyset

panama bananas



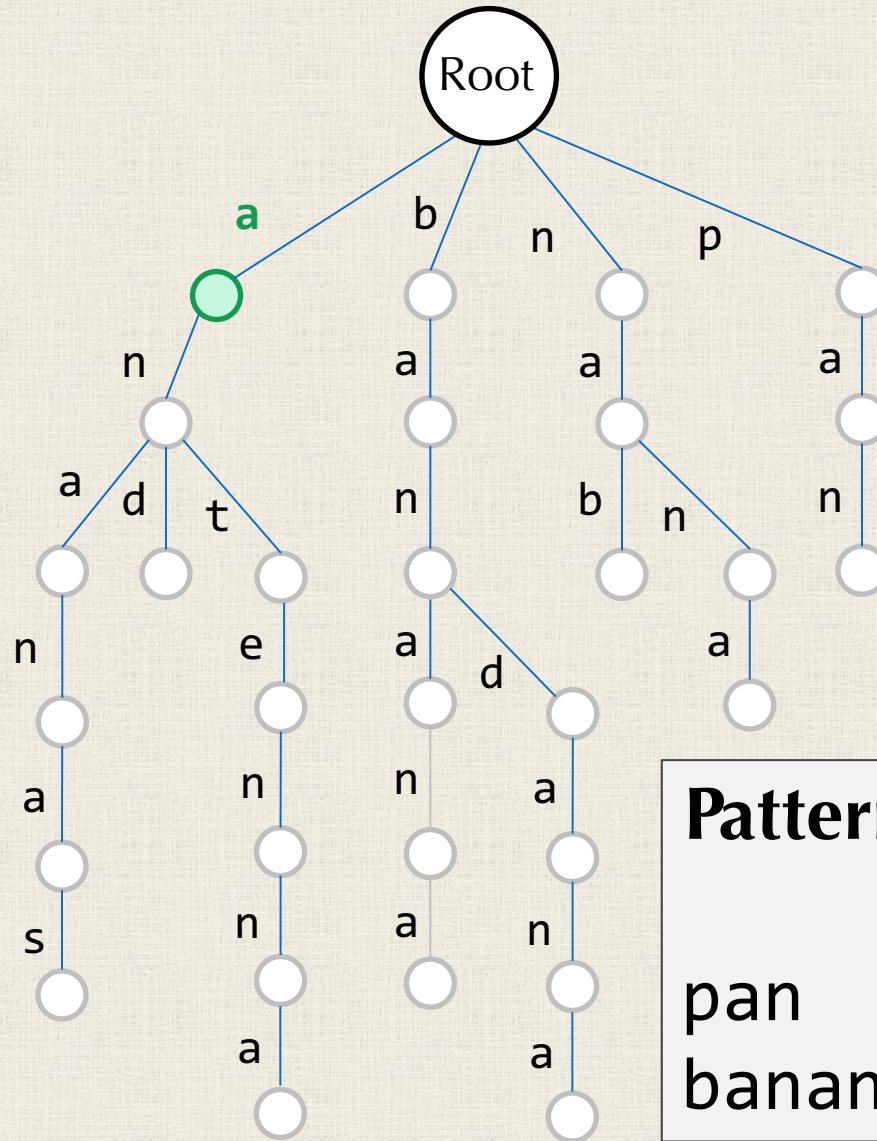
Pattern	Positions
pan	\emptyset

panama bananas



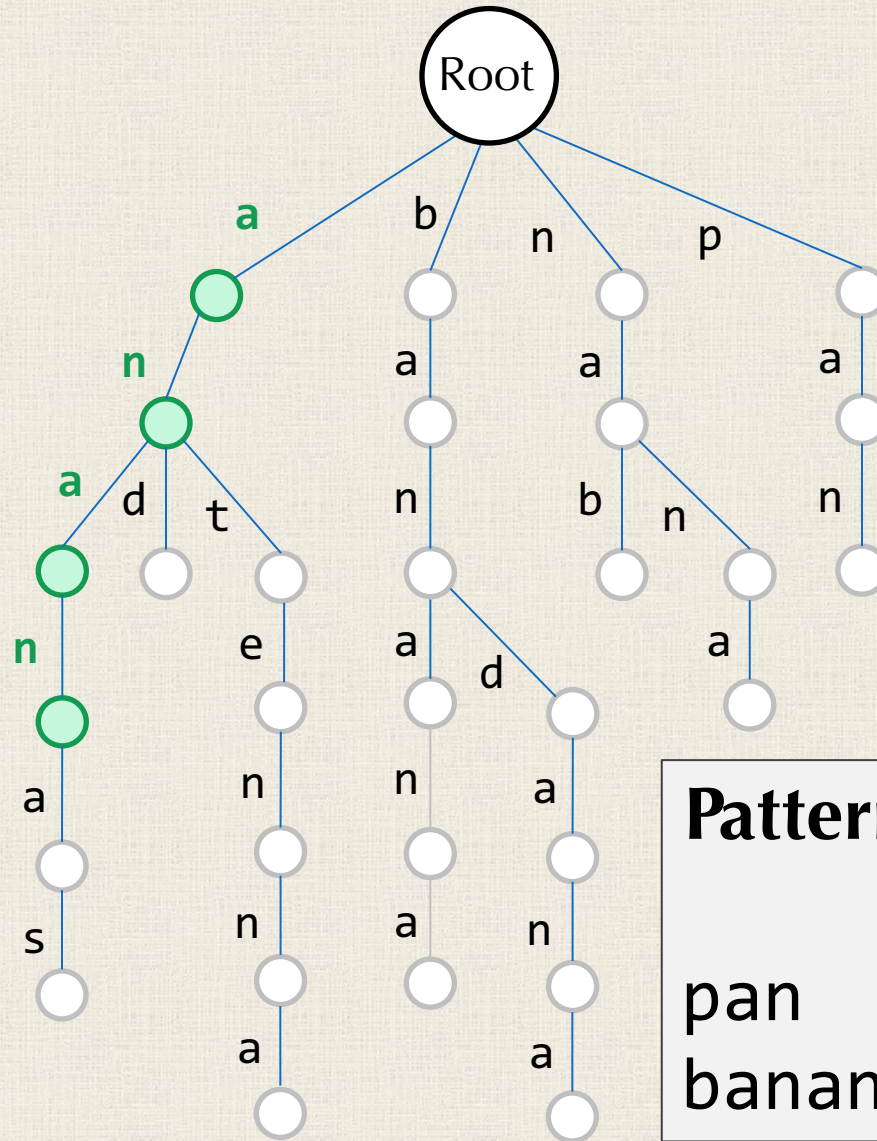
Pattern	Positions
pan	0
banana	6

panama banana



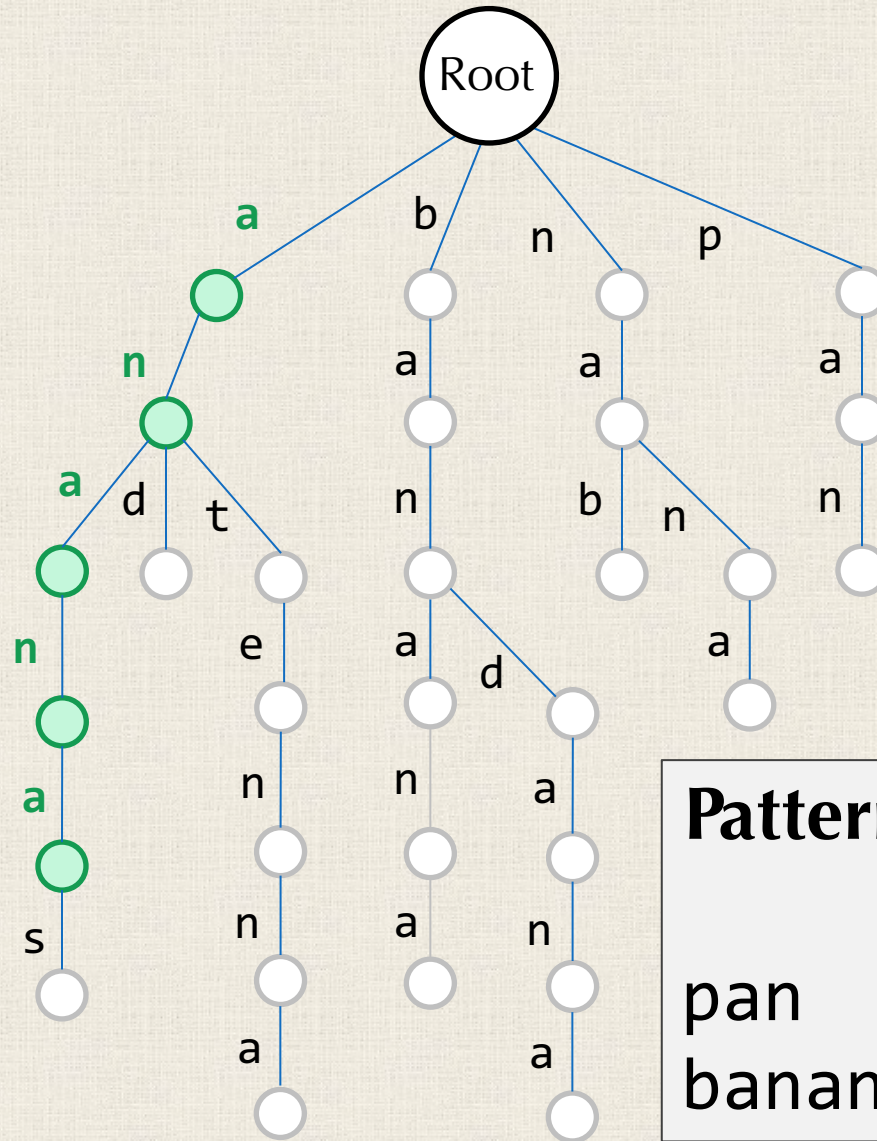
Pattern	Positions
pan	0
banana	6

panama bananas



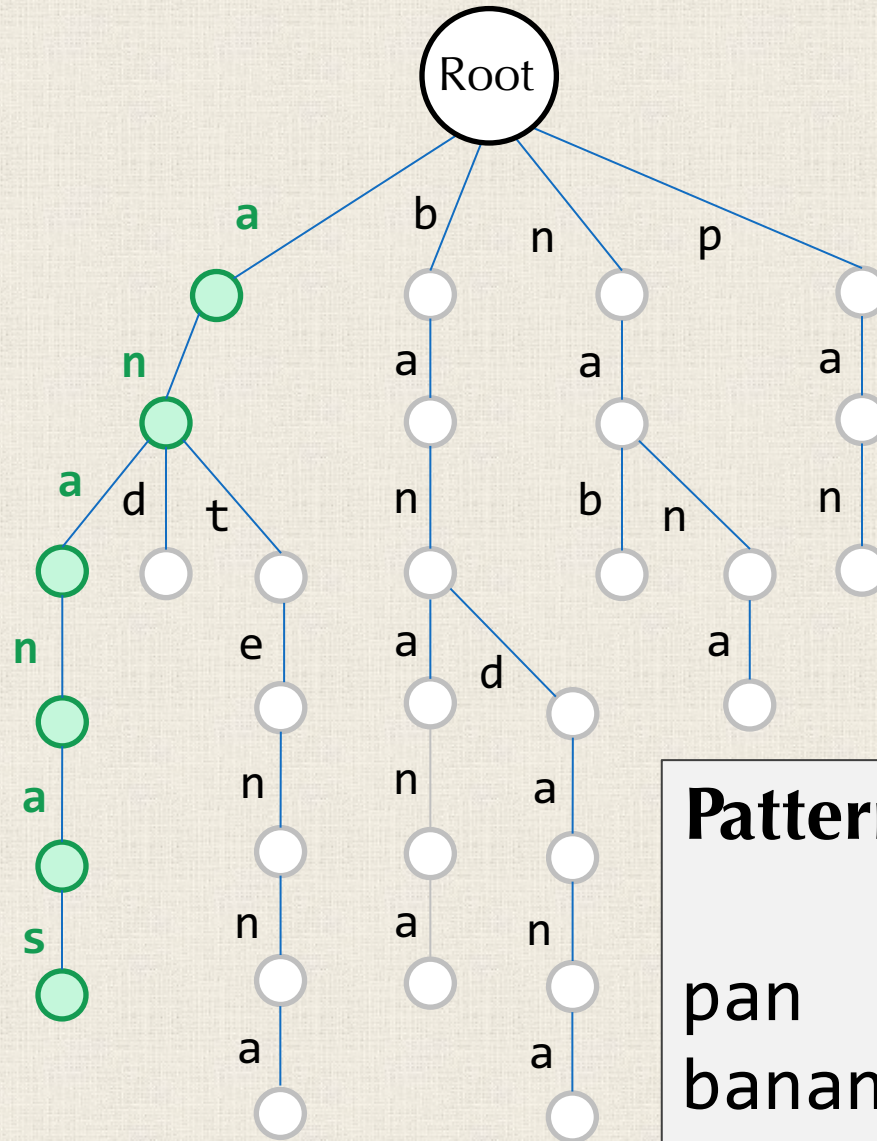
Pattern	Positions
pan	0
banana	6

panama bananas



Pattern	Positions
pan	0
banana	6

panama bananas



Pattern	Positions
pan	0
banana	6
ananas	7

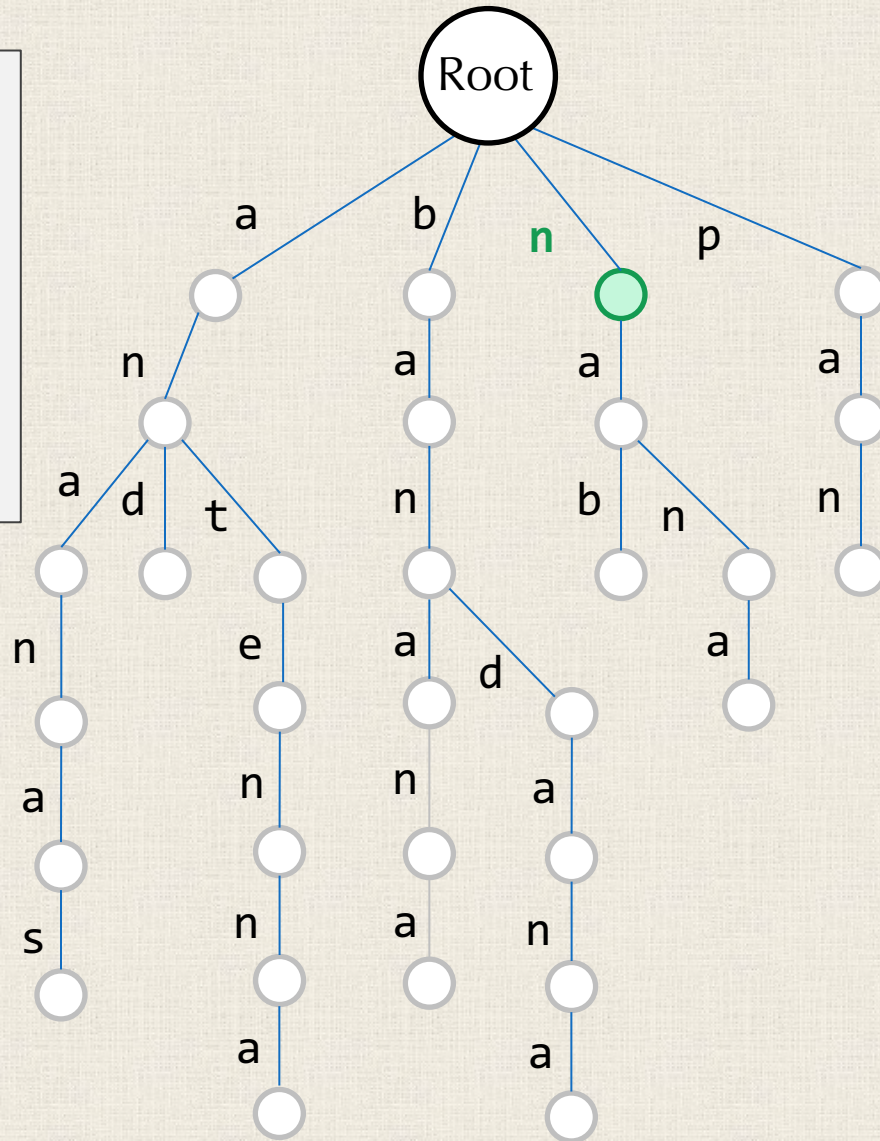
panamabananas

Pattern	Positions
---------	-----------

pan	0
-----	---

banana	6
--------	---

ananas	7
--------	---



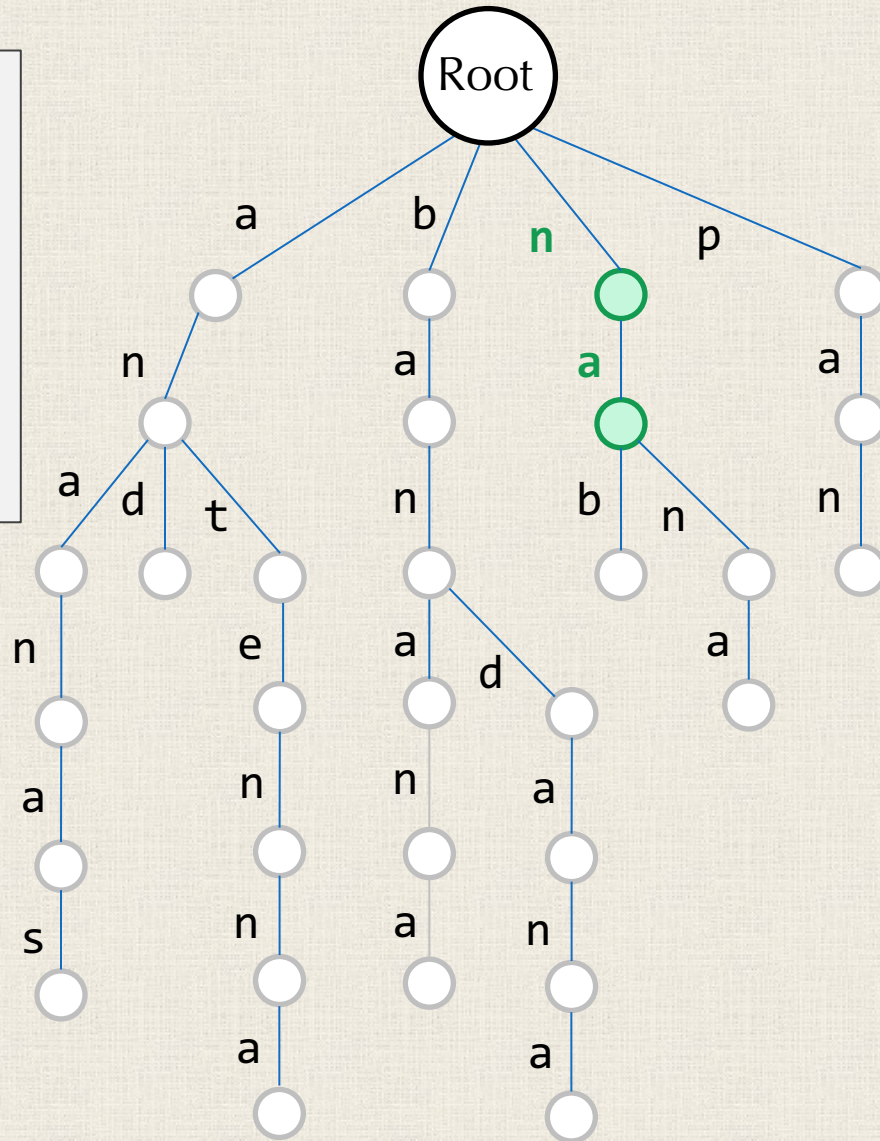
panamabananas

Pattern	Positions
---------	-----------

pan	0
-----	---

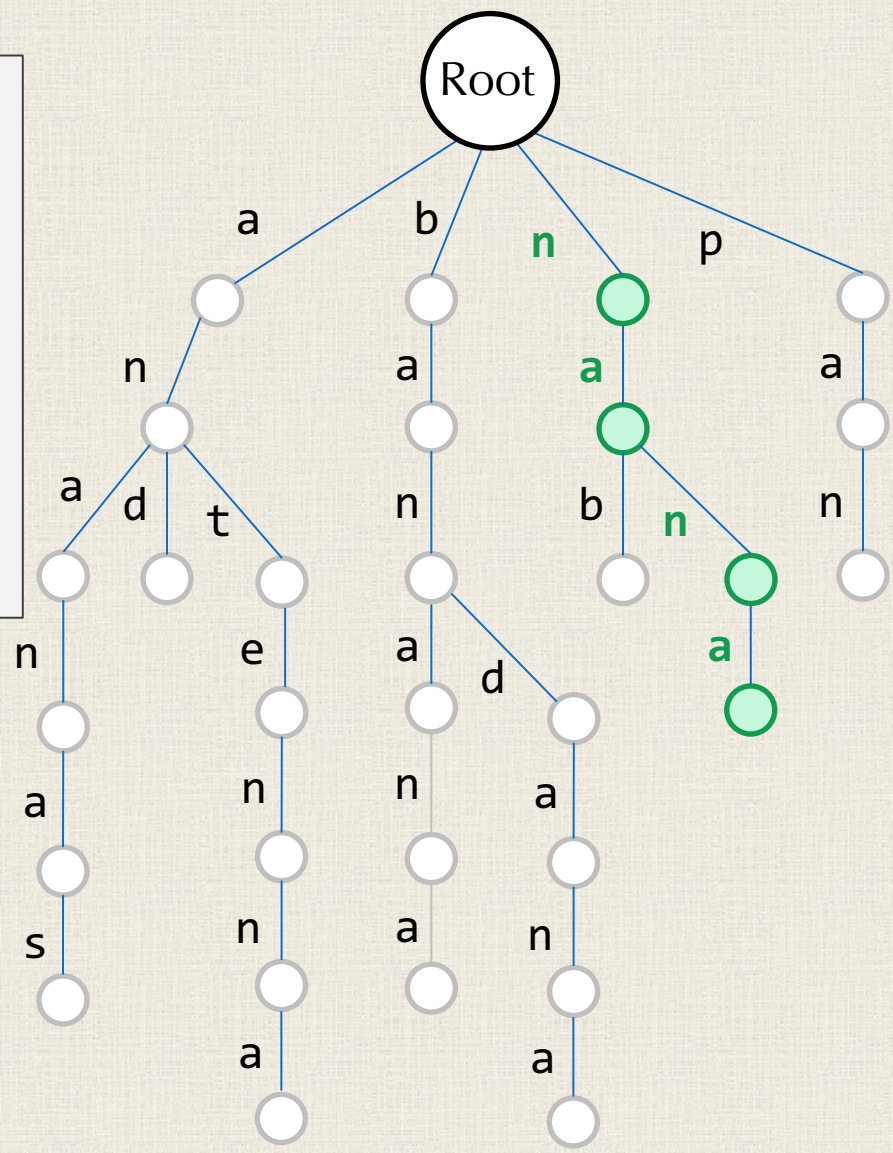
banana	6
--------	---

ananas	7
--------	---



panamabannanas

Pattern	Positions
pan	0
banana	6
ananas	7
nana	8



panamabananas

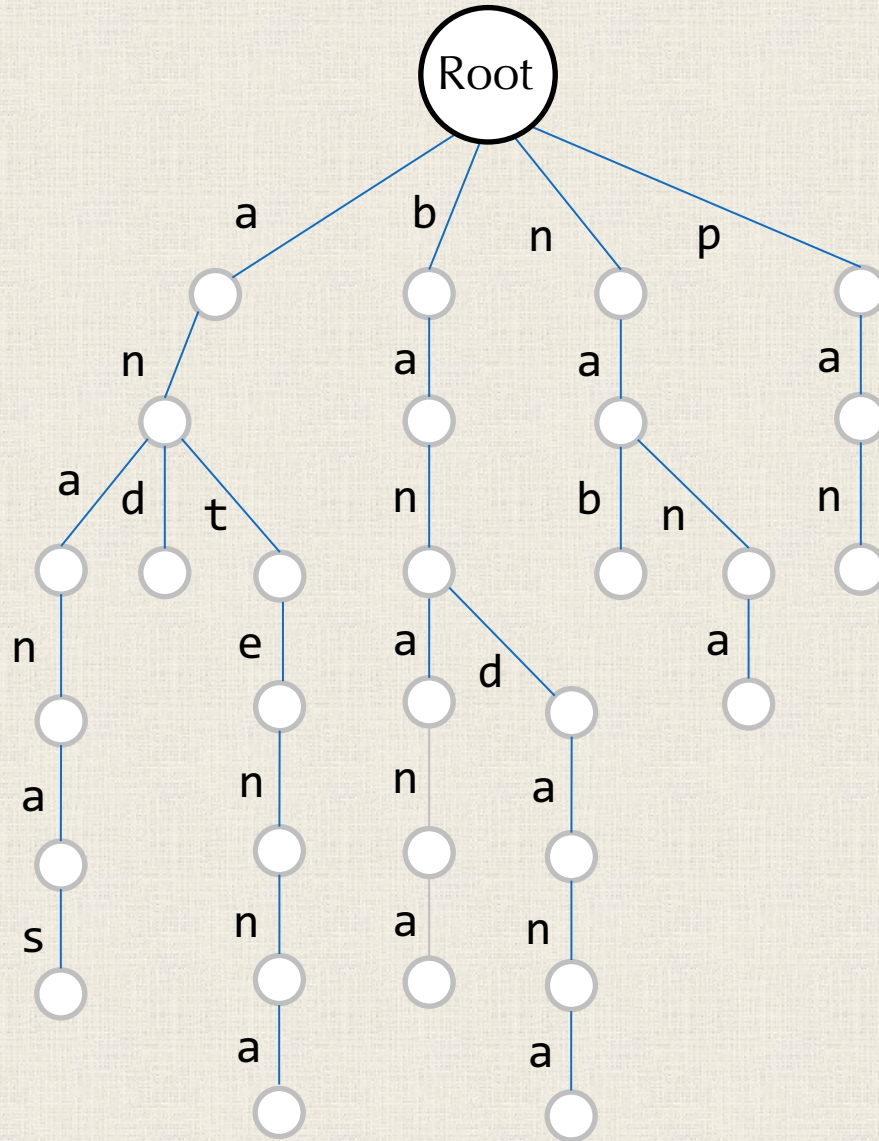
Pattern	Positions
---------	-----------

pan	0
-----	---

banana	6
--------	---

ananas	7
--------	---

nana	8
------	---



panamabanas

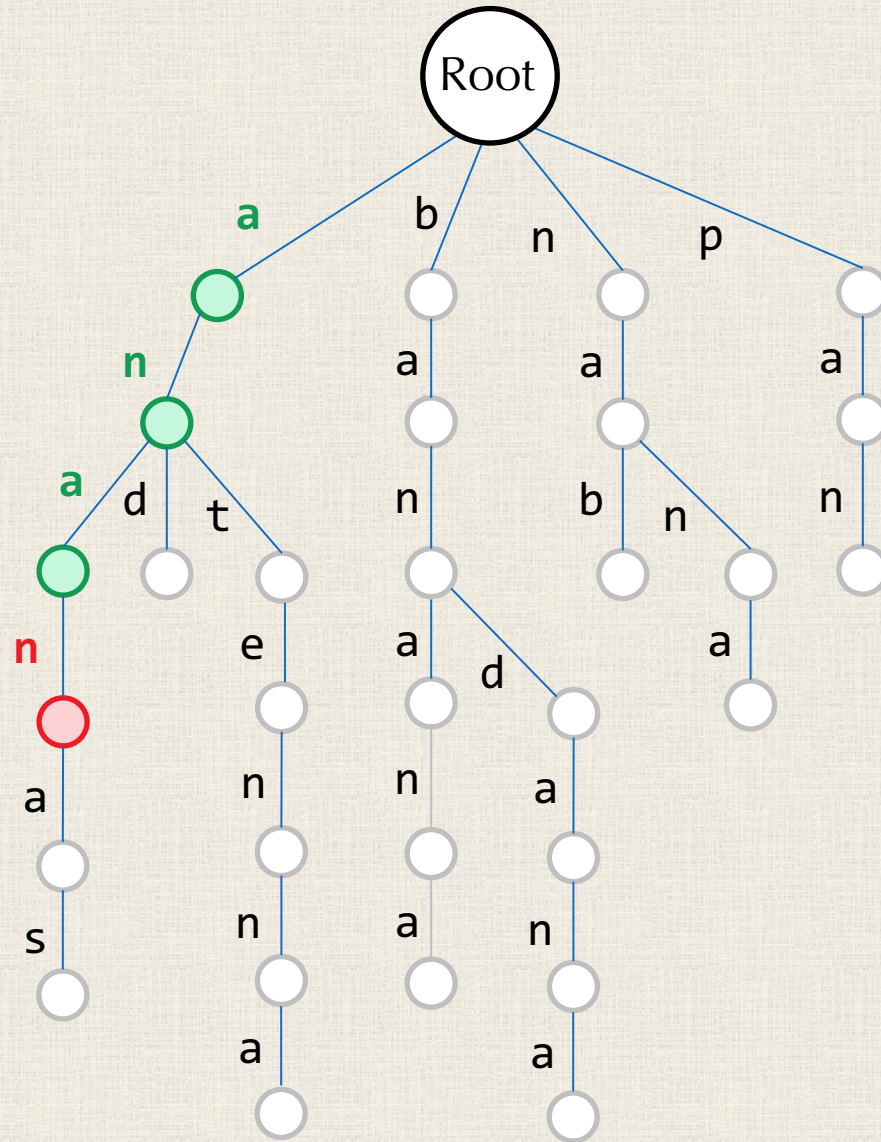
Pattern	Positions
---------	-----------

pan	0
-----	---

banana	6
--------	---

ananas	7
--------	---

nana	8
------	---



panamabanas

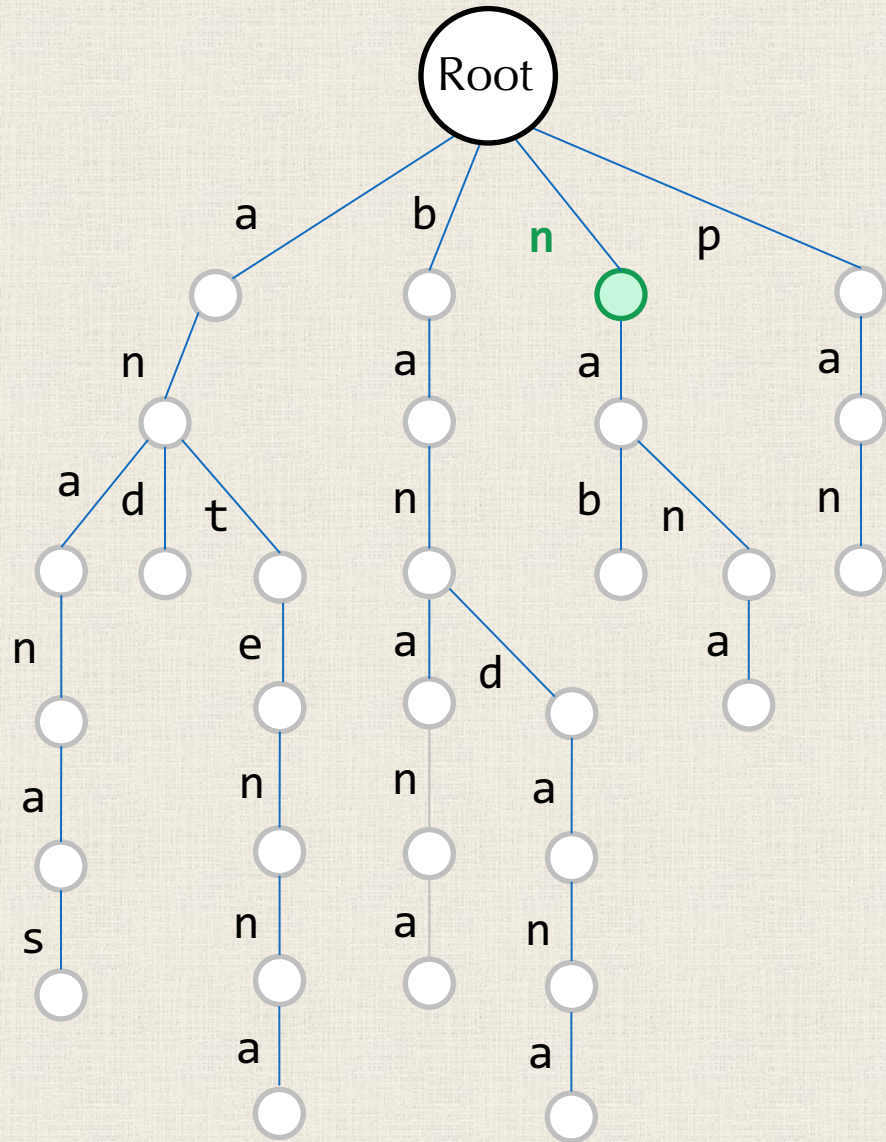
Pattern	Positions
---------	-----------

pan	0
-----	---

banana	6
--------	---

ananas	7
--------	---

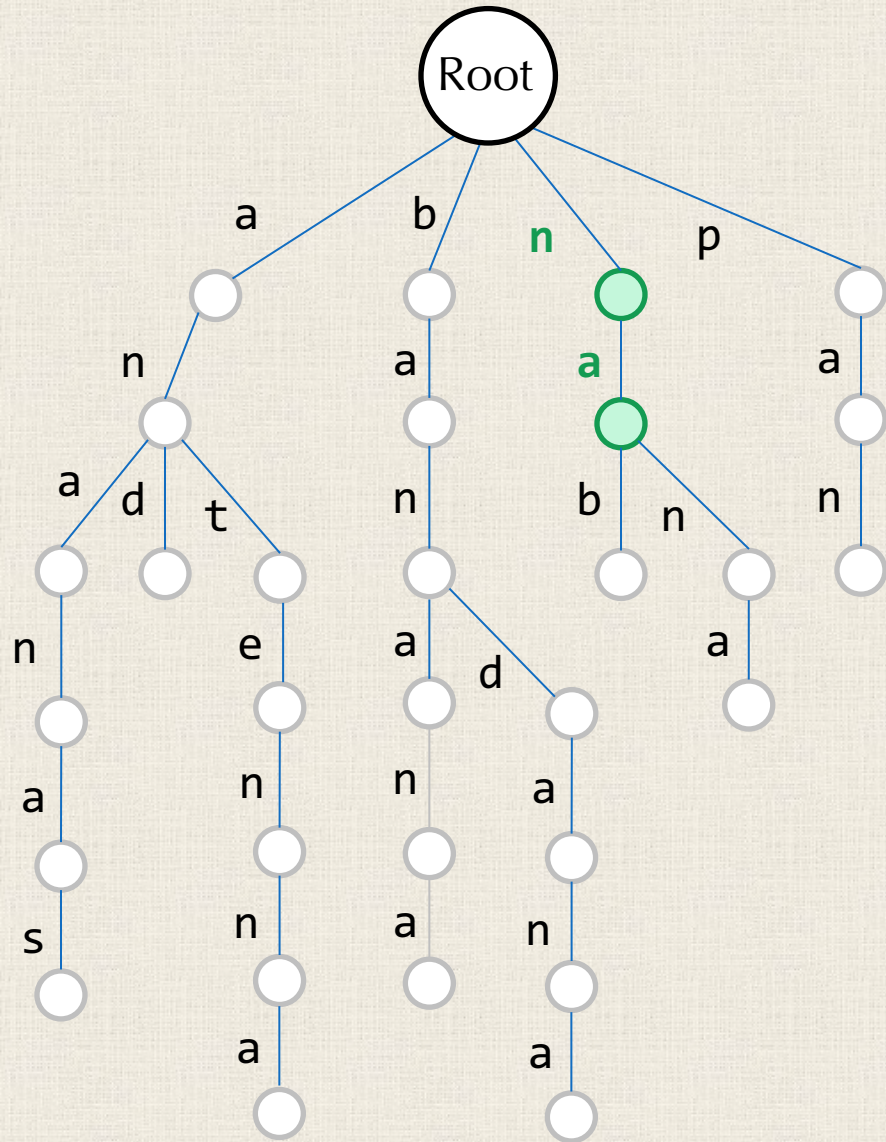
nana	8
------	---



panamabanas

Pattern	Positions
---------	-----------

pan	0
banana	6
ananas	7
nana	8



panamabananas

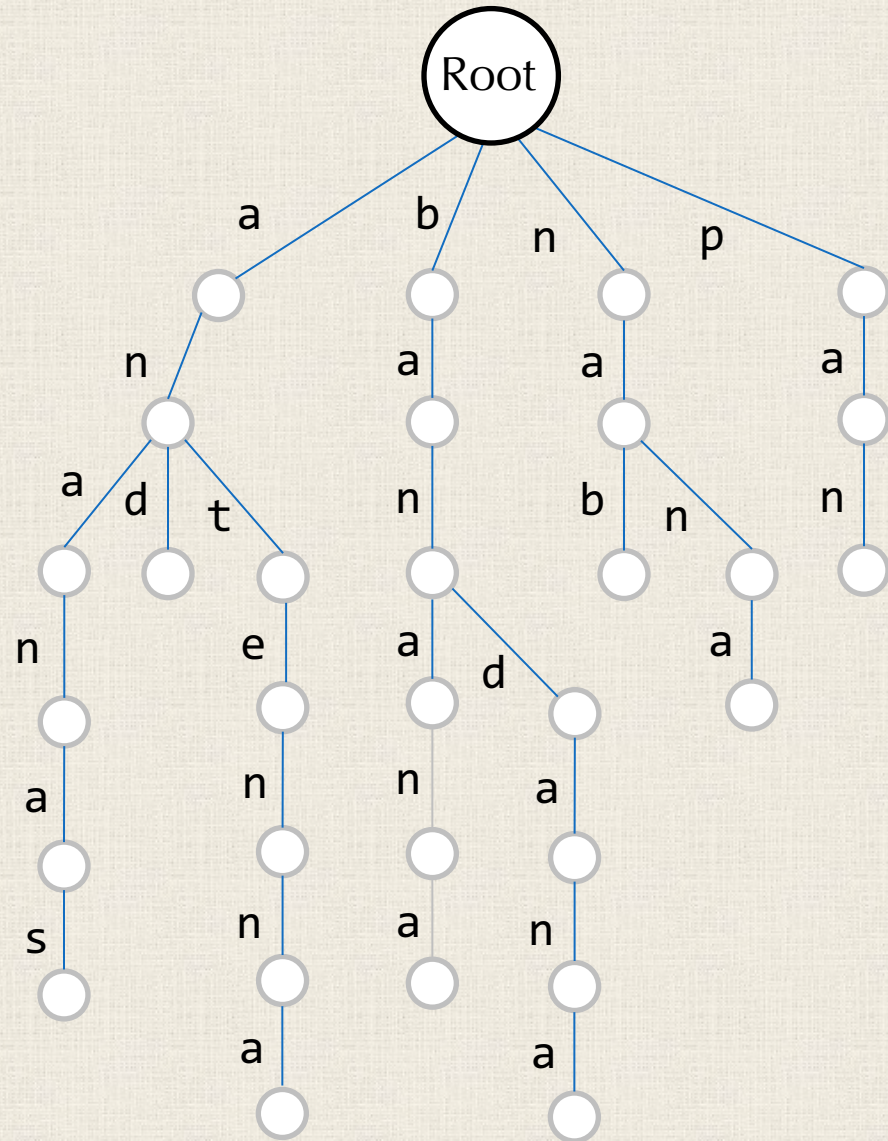
Pattern	Positions
---------	-----------

pan	0
-----	---

banana	6
--------	---

ananas	7
--------	---

nana	8
------	---



panamabananas

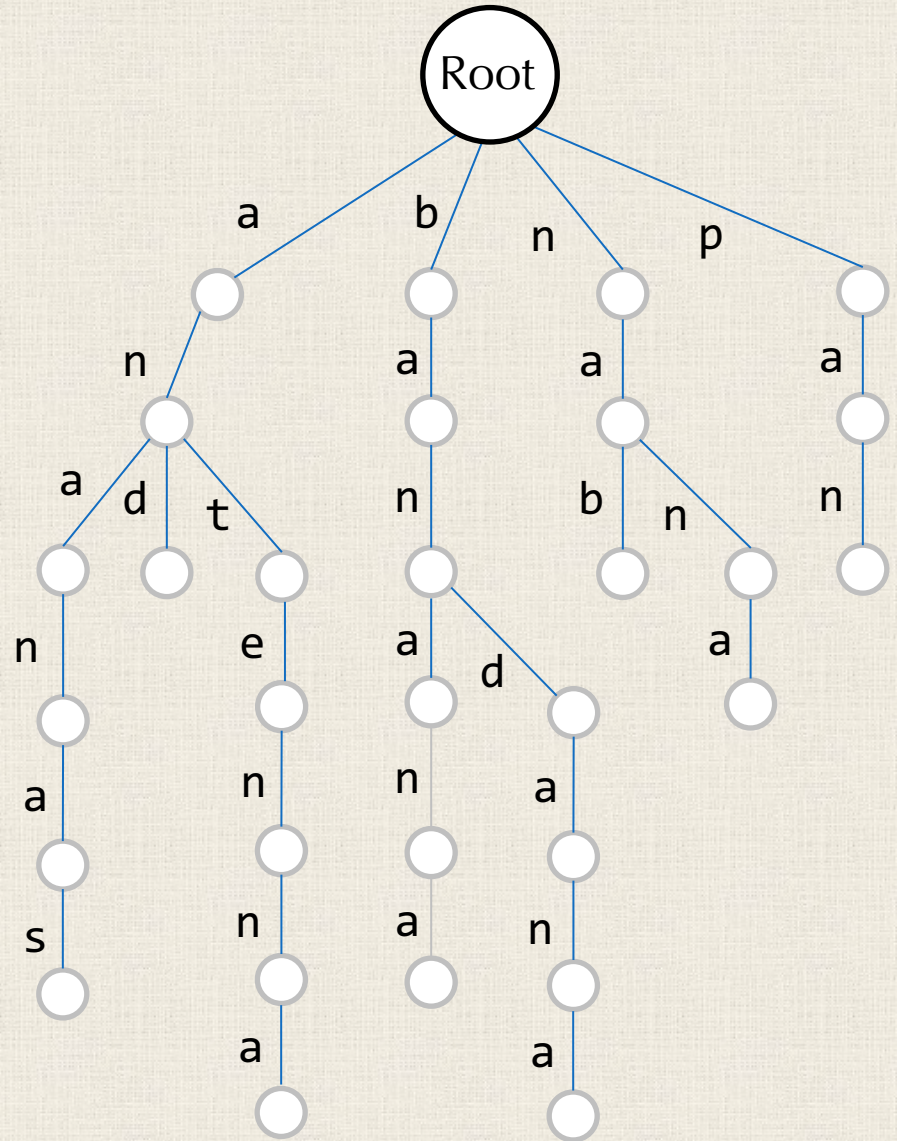
Pattern	Positions
---------	-----------

pan	0
-----	---

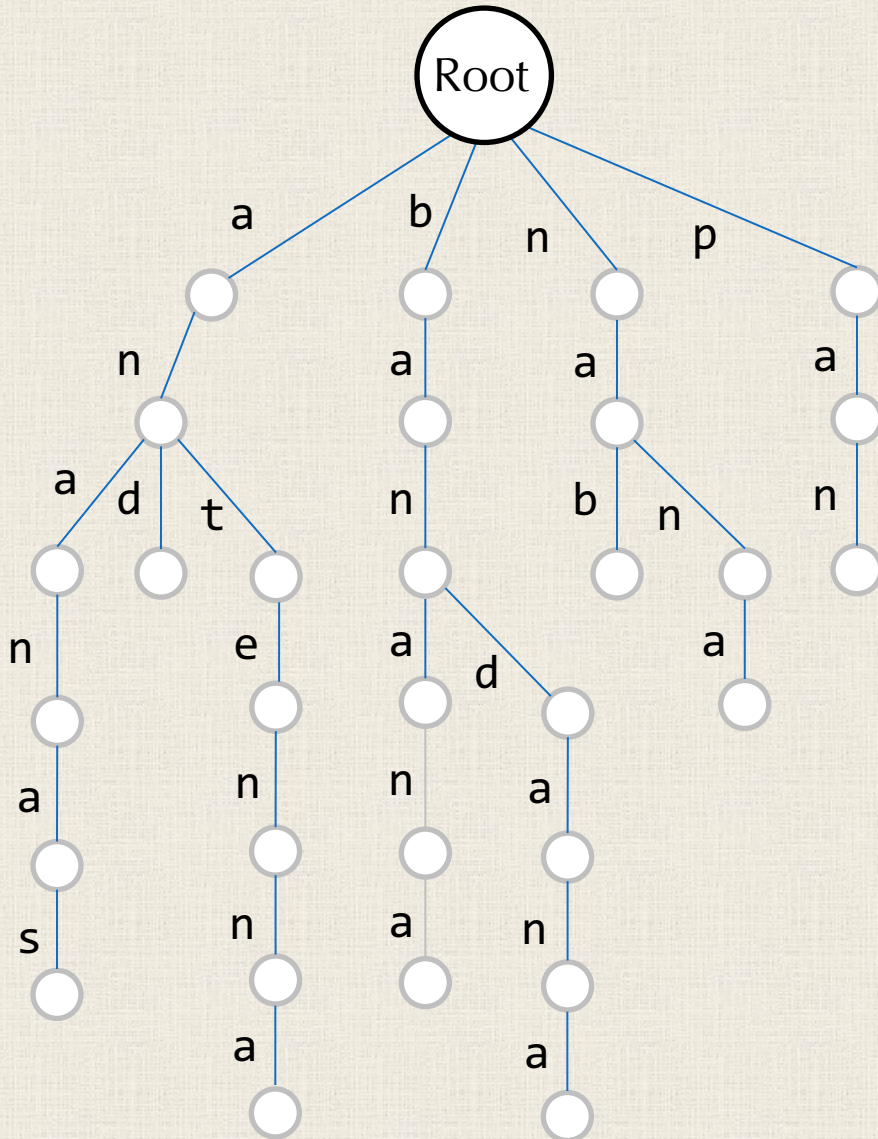
banana	6
--------	---

ananas	7
--------	---

nana	8
------	---



What if one pattern is a prefix of another?



STOP: If we add the pattern “band” to our dataset, why would it cause problems?

Success?

Runtime of Brute Force: $O(|Text|^*|Patterns|)$

Success?

Runtime of Brute Force: $O(|Text|^*|Patterns|)$

STOP: What is the running time of multiple pattern matching with a trie?

Success?

Runtime of Brute Force: $O(|Text|^*|Patterns|)$

STOP: What is the running time of multiple pattern matching with a trie?

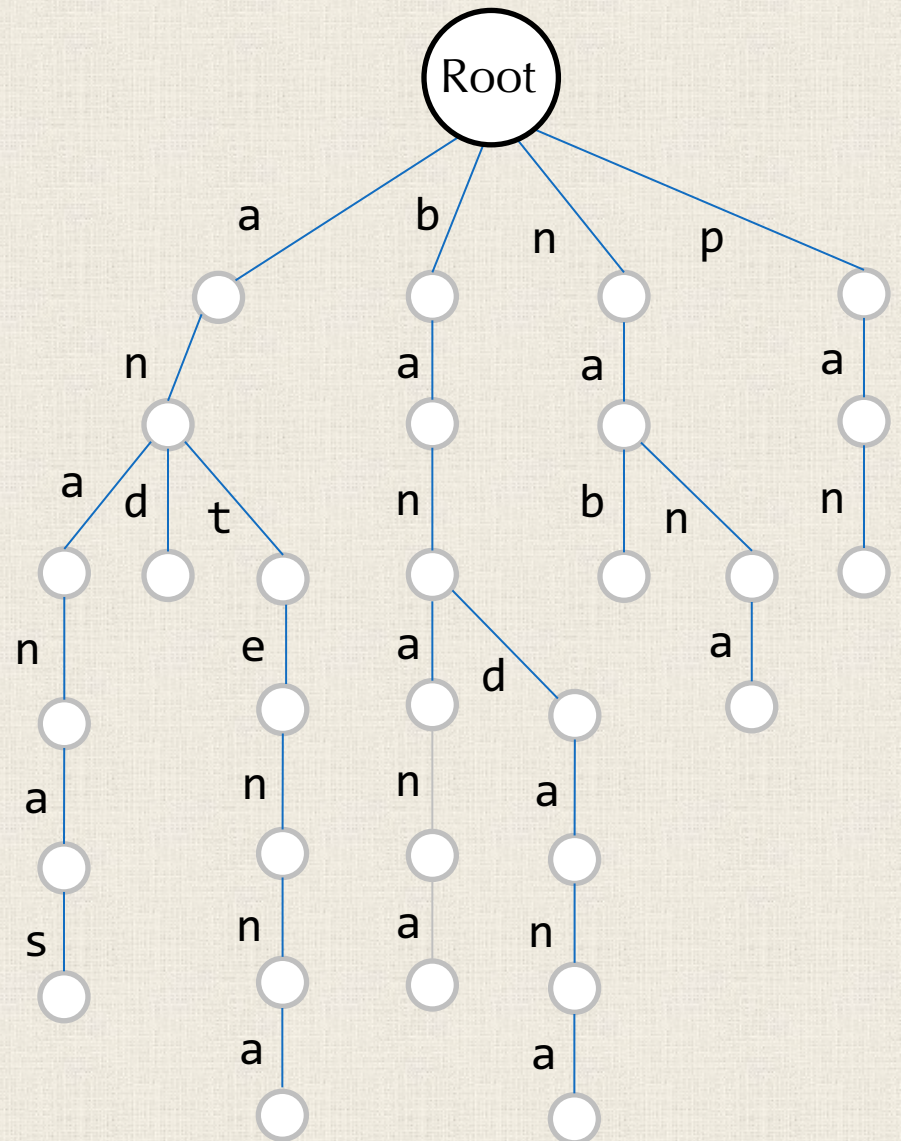
Runtime of Pattern Matching with Trie:

- Trie Construction: $O(|Patterns|)$
- Pattern Matching: $O(|Text| * |LongestPattern|)$

But We Neglected Memory ...

Our trie has 30 edges
(38 when we add the \$
symbol), and $|Patterns| = 39$.

Worst case: Trie takes up
 $O(|Patterns|)$, which for
a human sequencing
project may be ~1 TB!



PREPROCESSING THE GENOME INSTEAD

Preprocessing the Genome

If $|Patterns| > |Text|$, then why not try to construct a data structure from *Text*?

This would take less memory and allow us to “teleport” each string *Pattern* to its occurrences in *Text*.

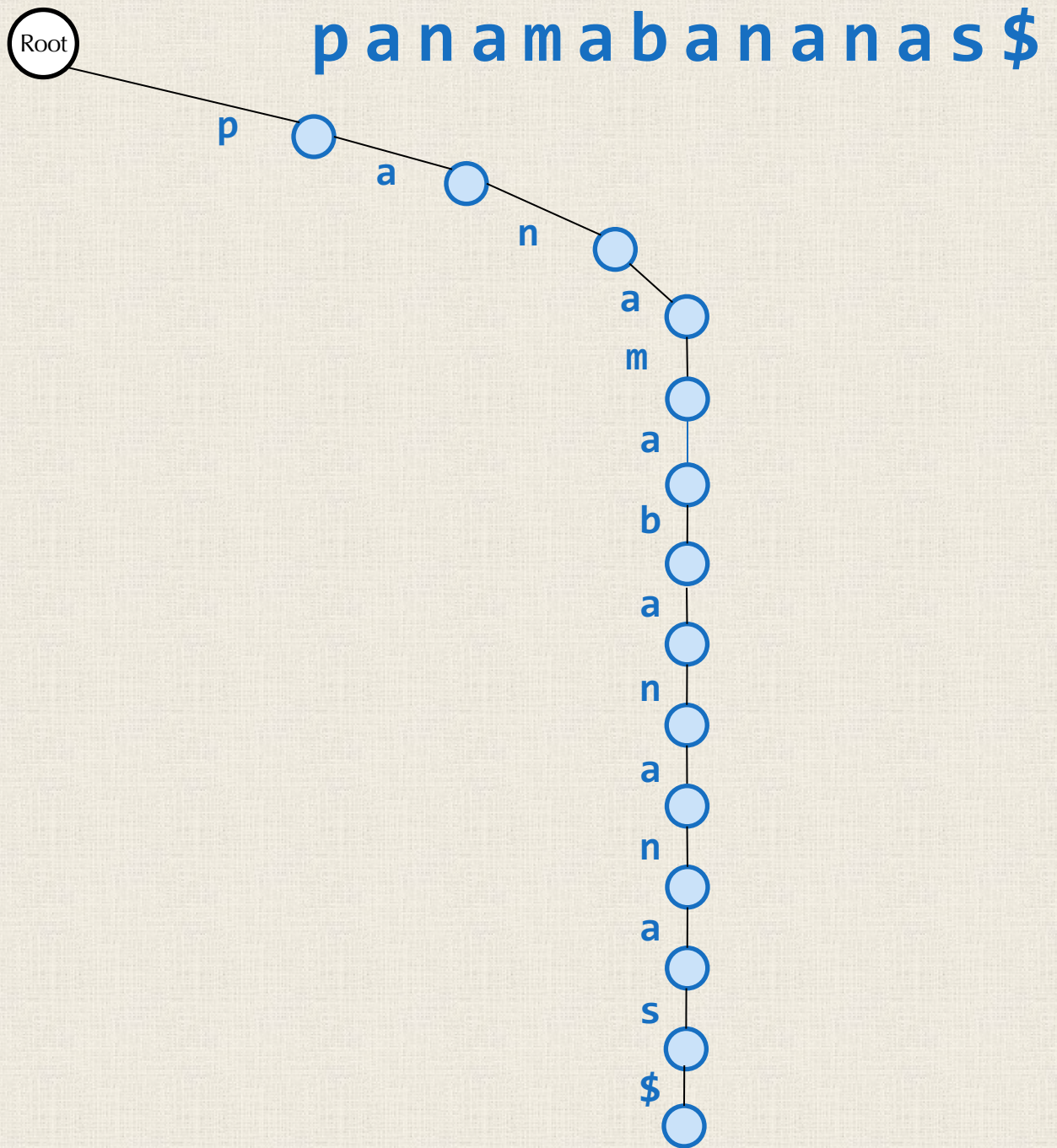




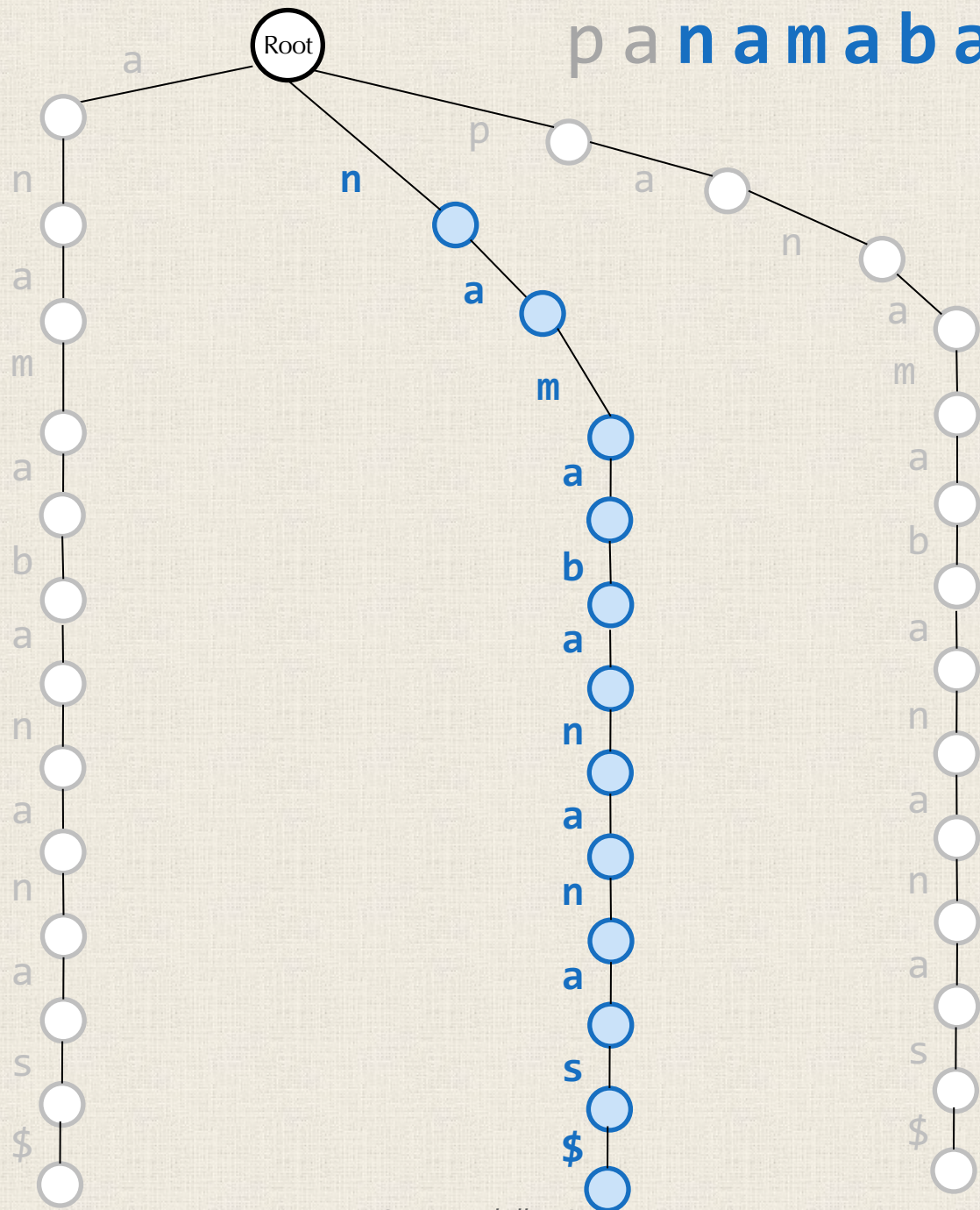
panamabananas\$

Suffix trie: Trie formed from the *suffixes* of *Text*.

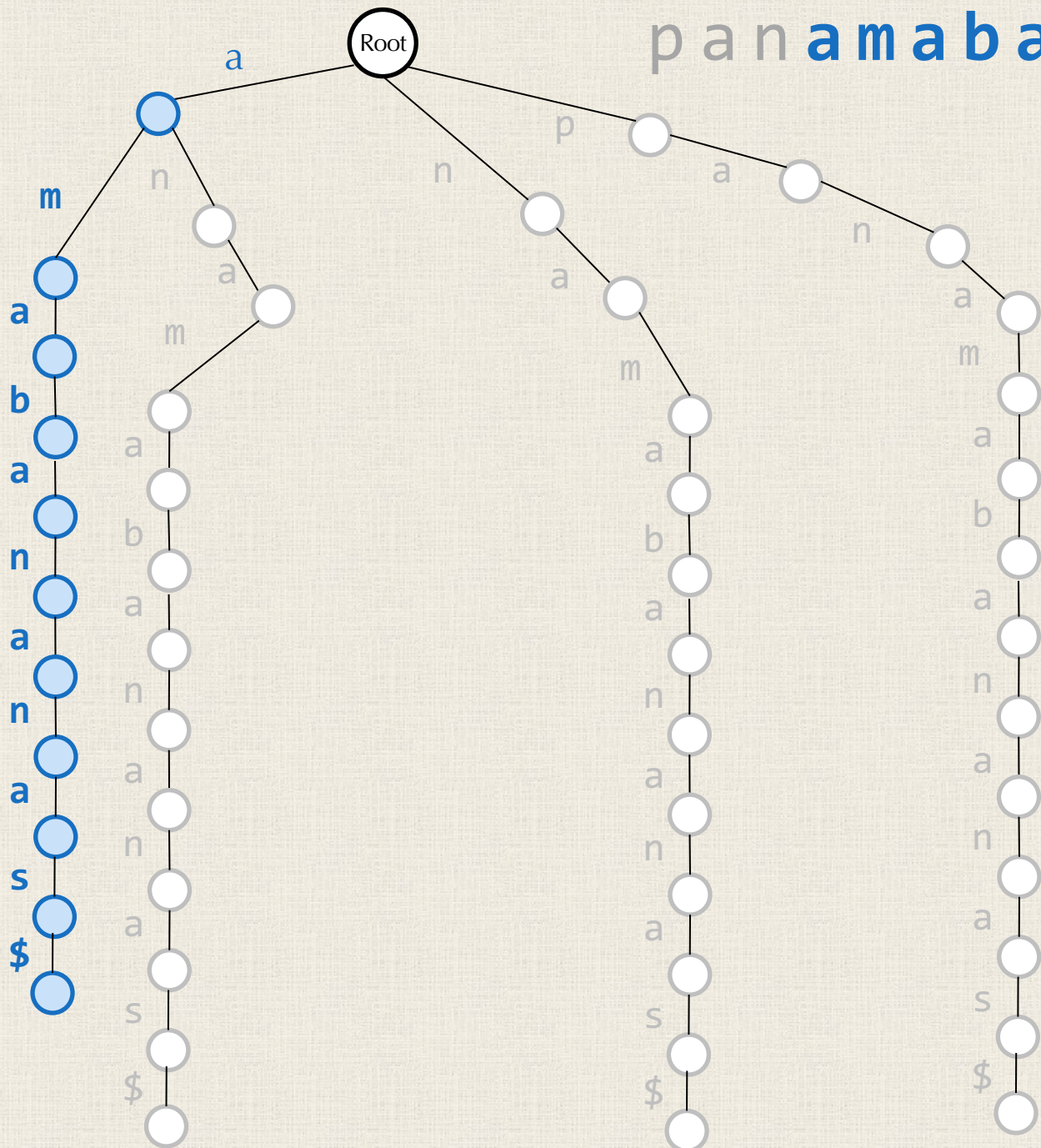
Note: We augment *Text* with a \$ to mark the end of the string (why?)



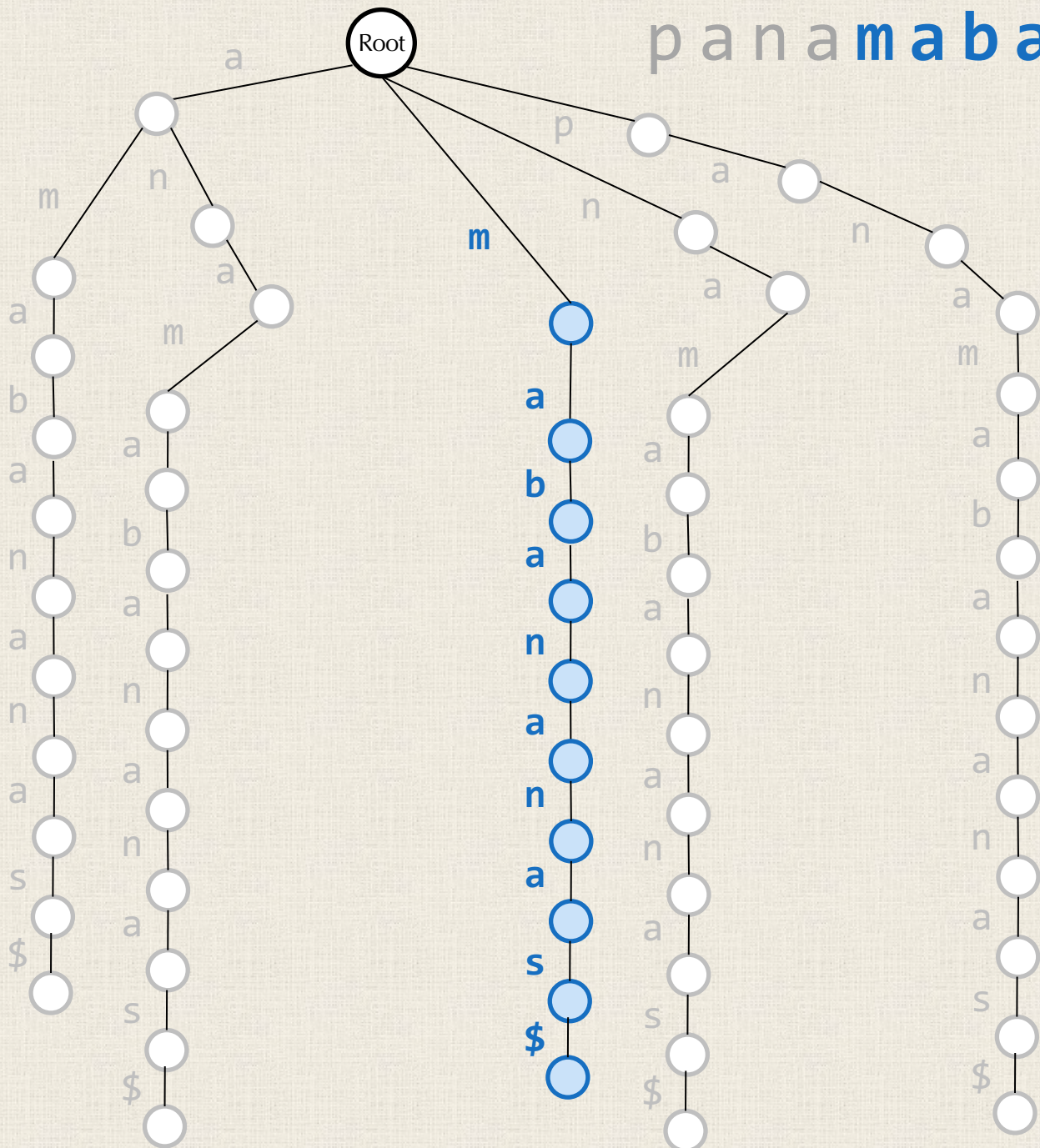
panamabananas\$



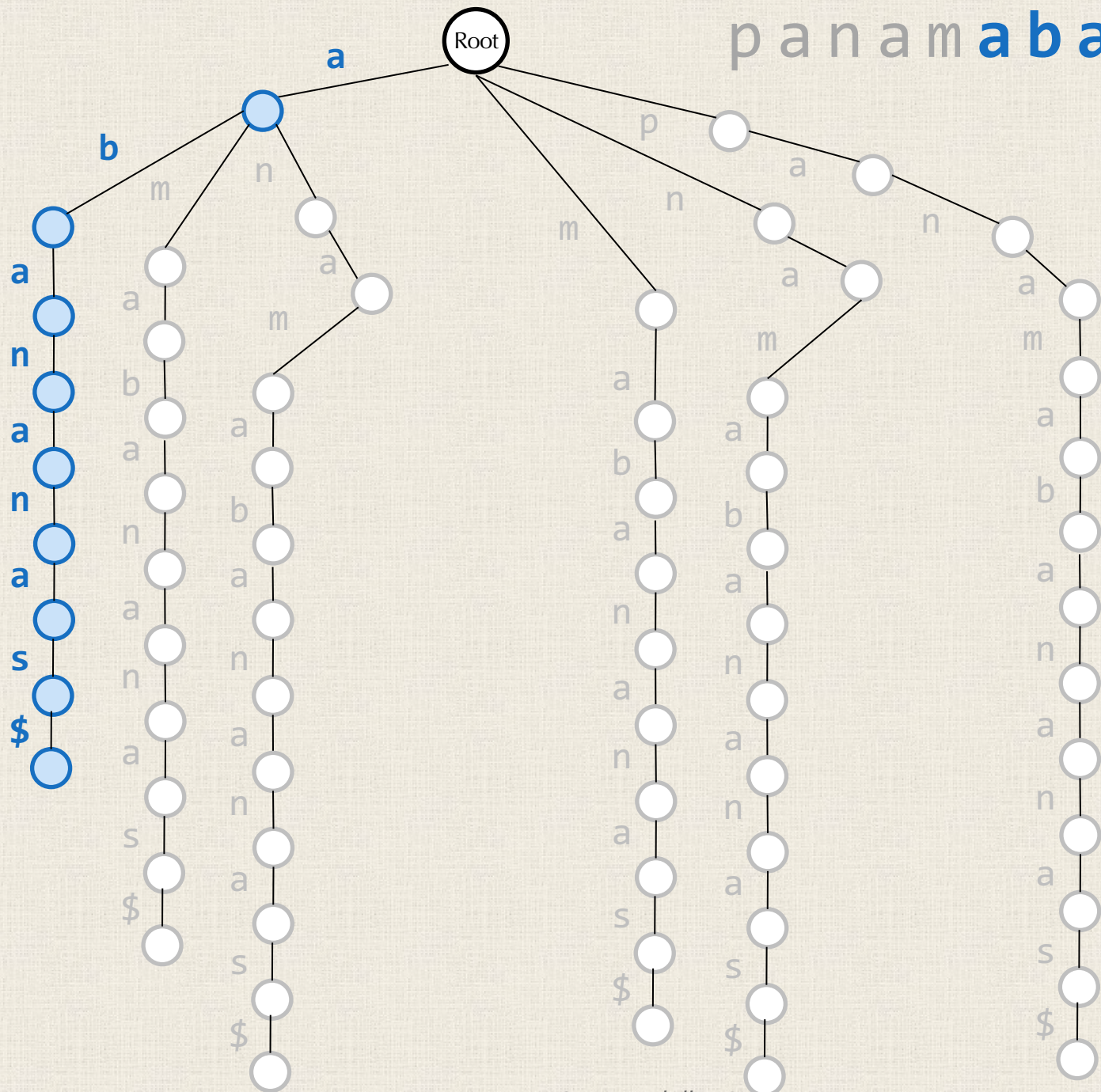
panamabananas\$



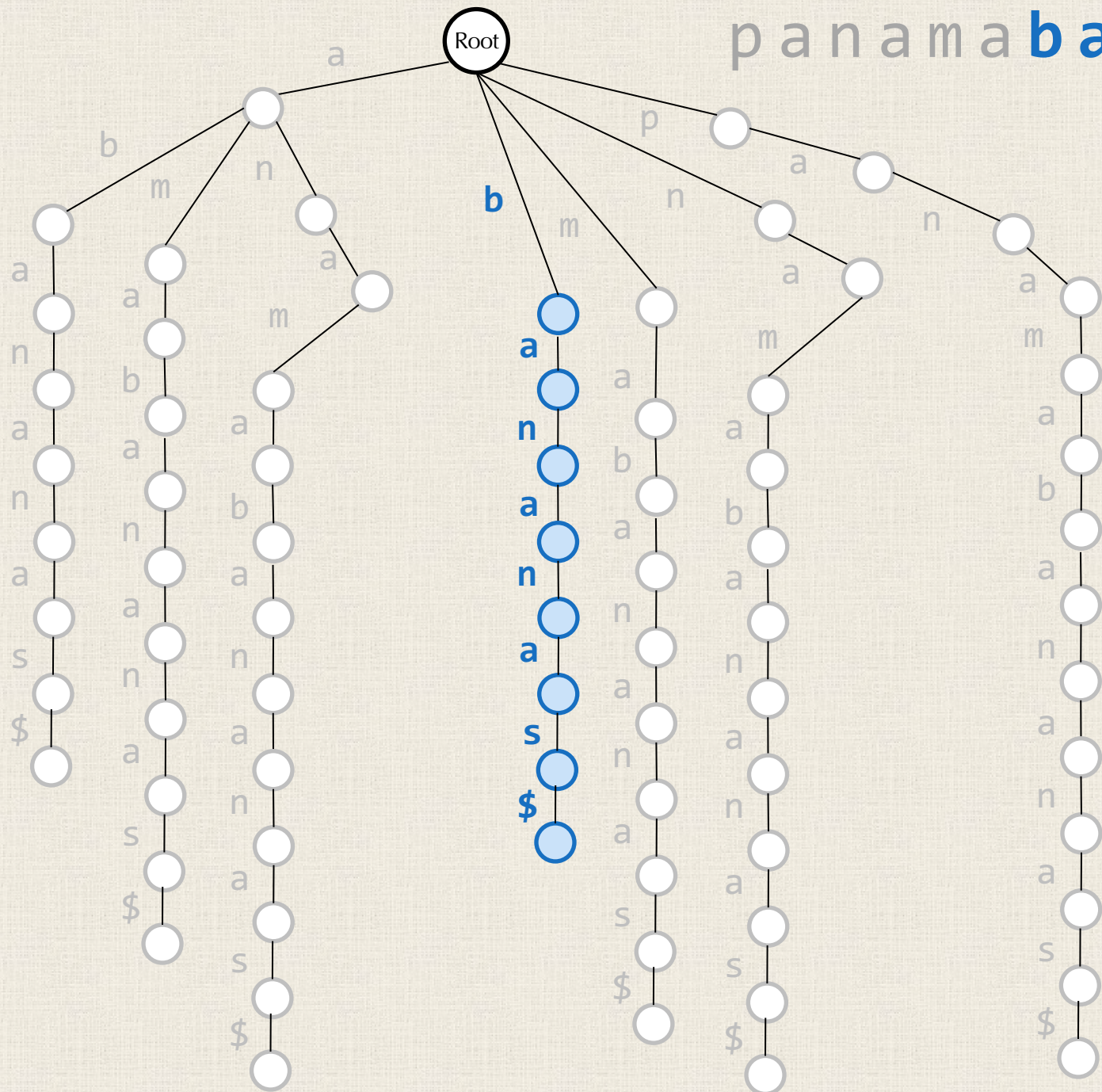
panamabananas\$



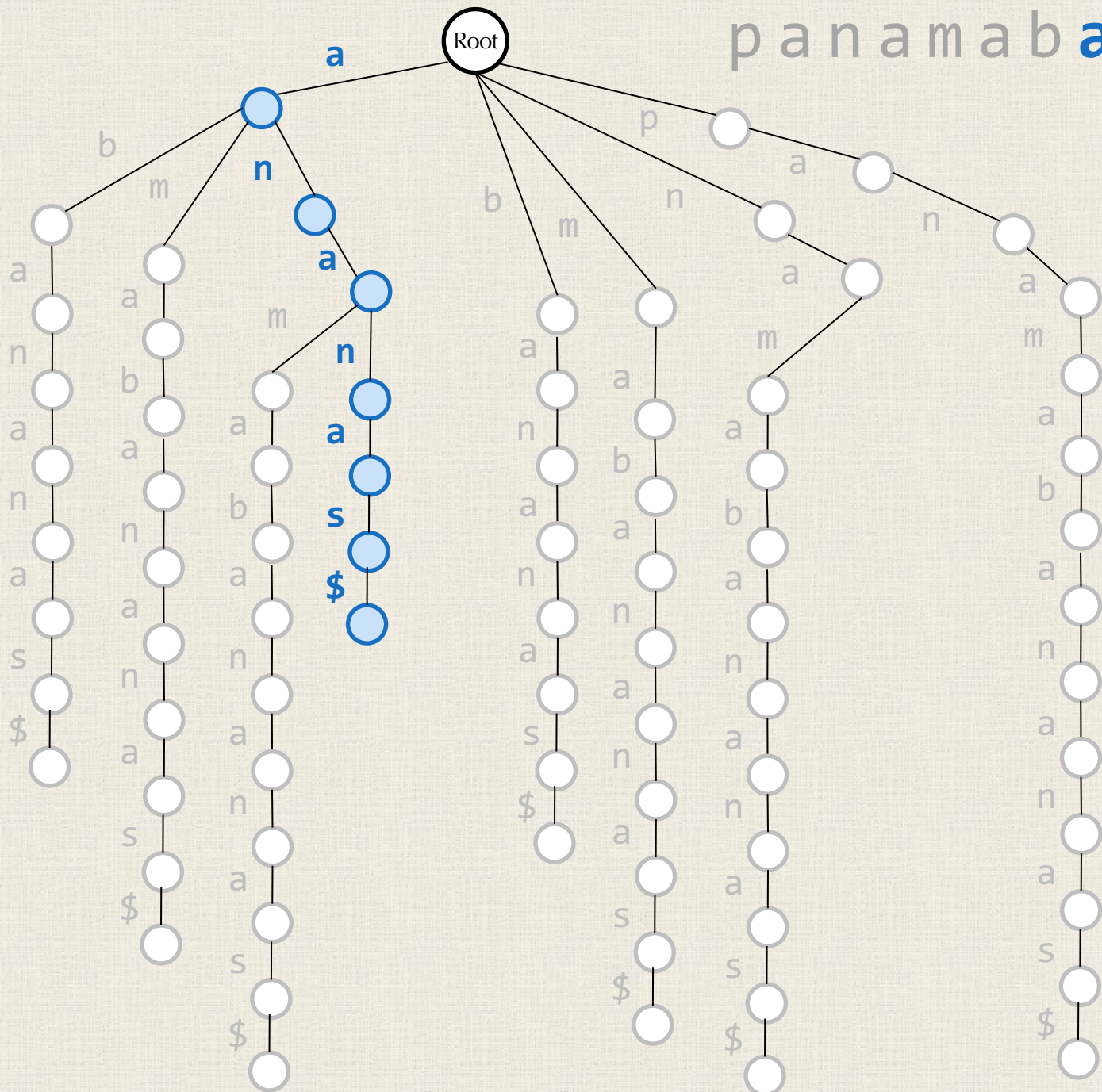
panamabanas\$



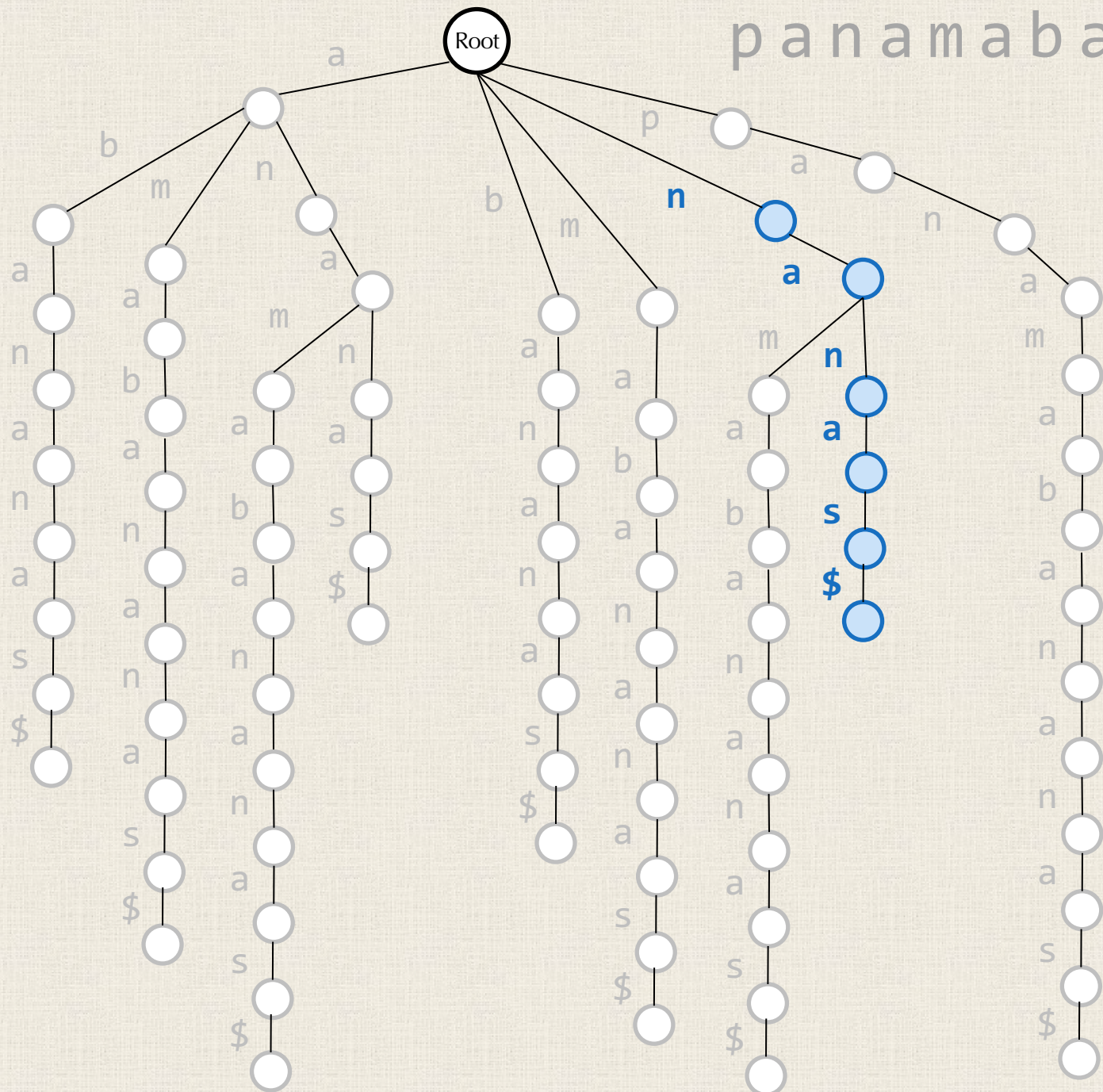
panama **bananas** \$



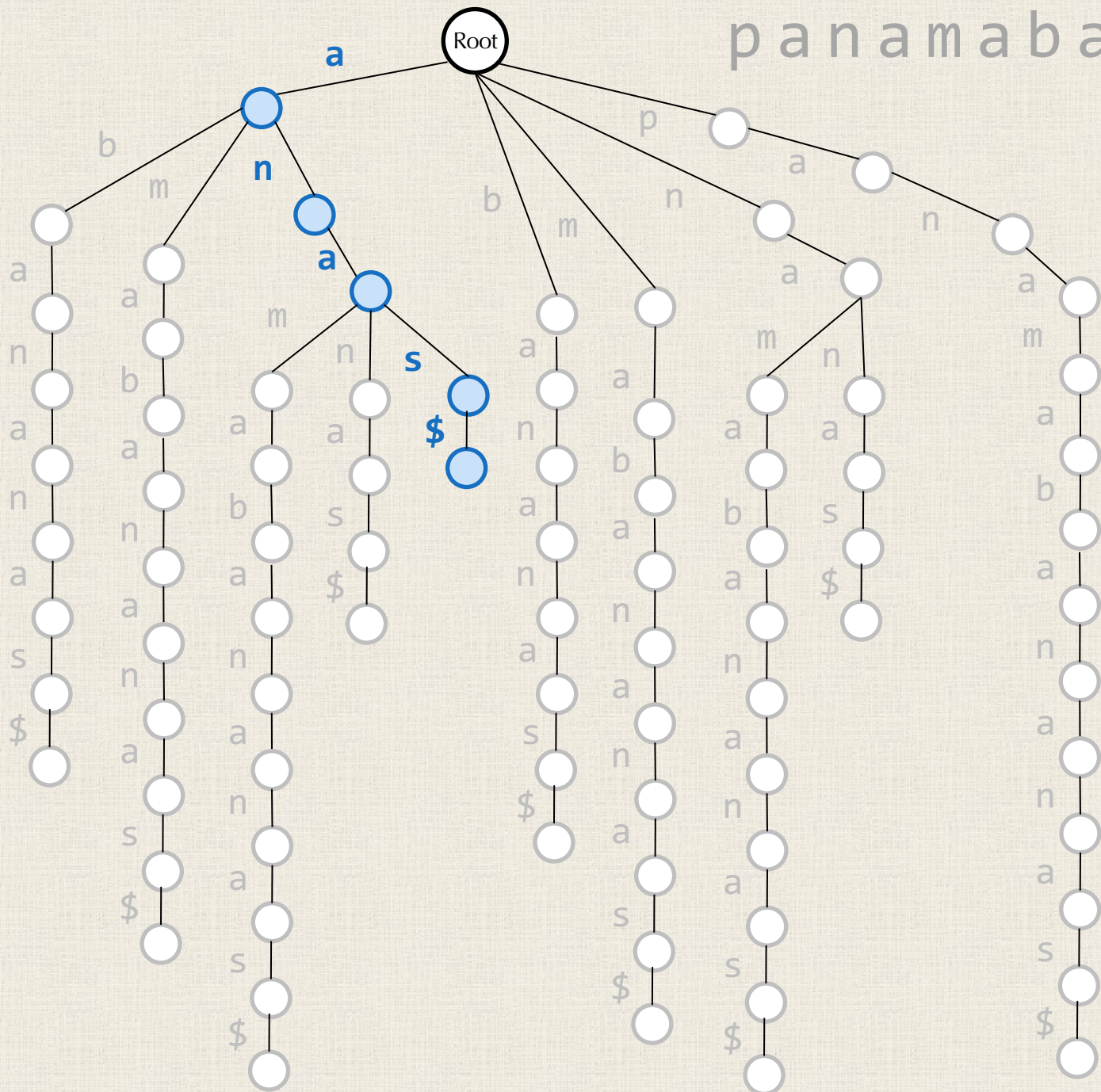
panamabananas\$



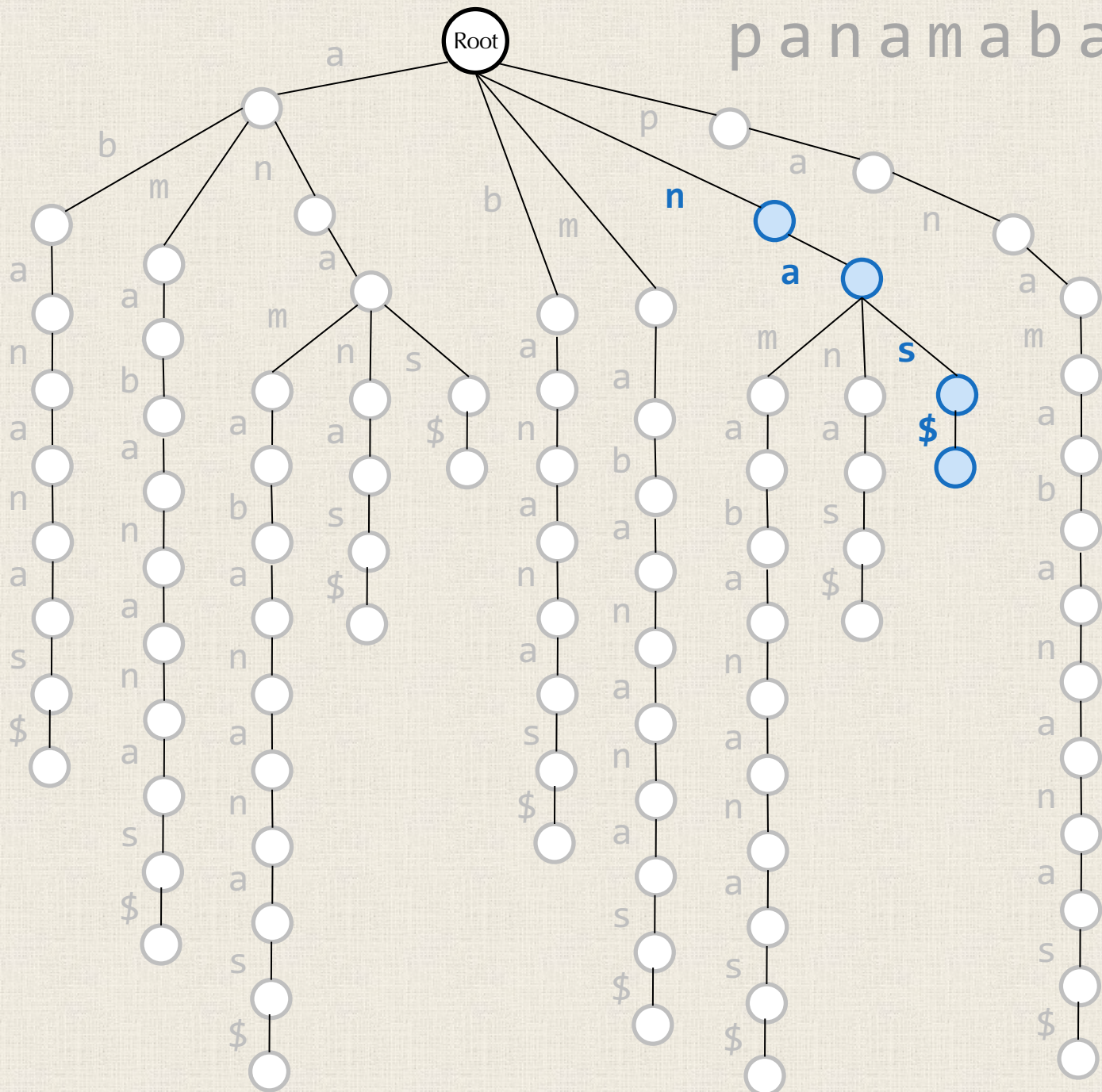
panamabananas\$



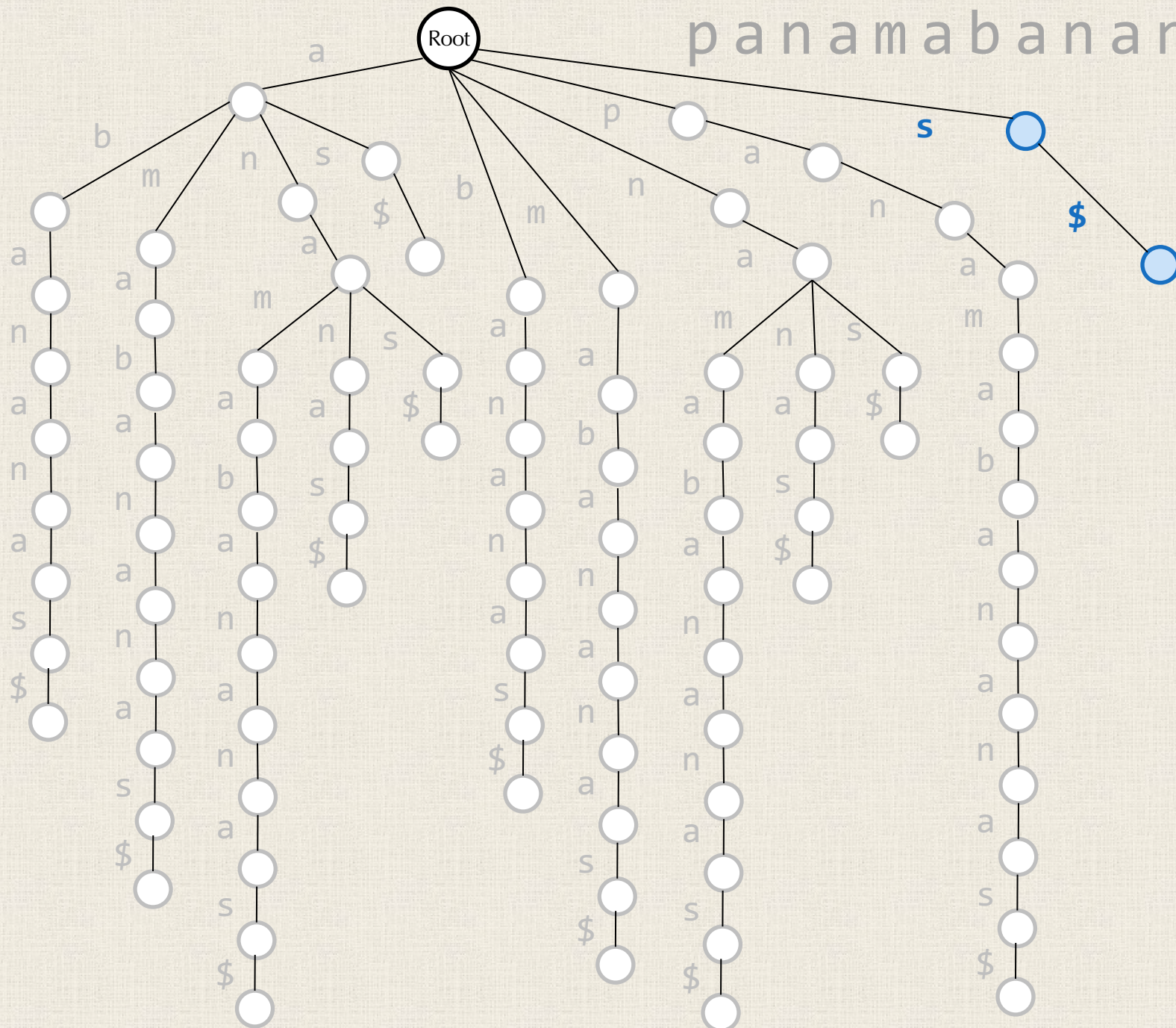
panamabanas\$



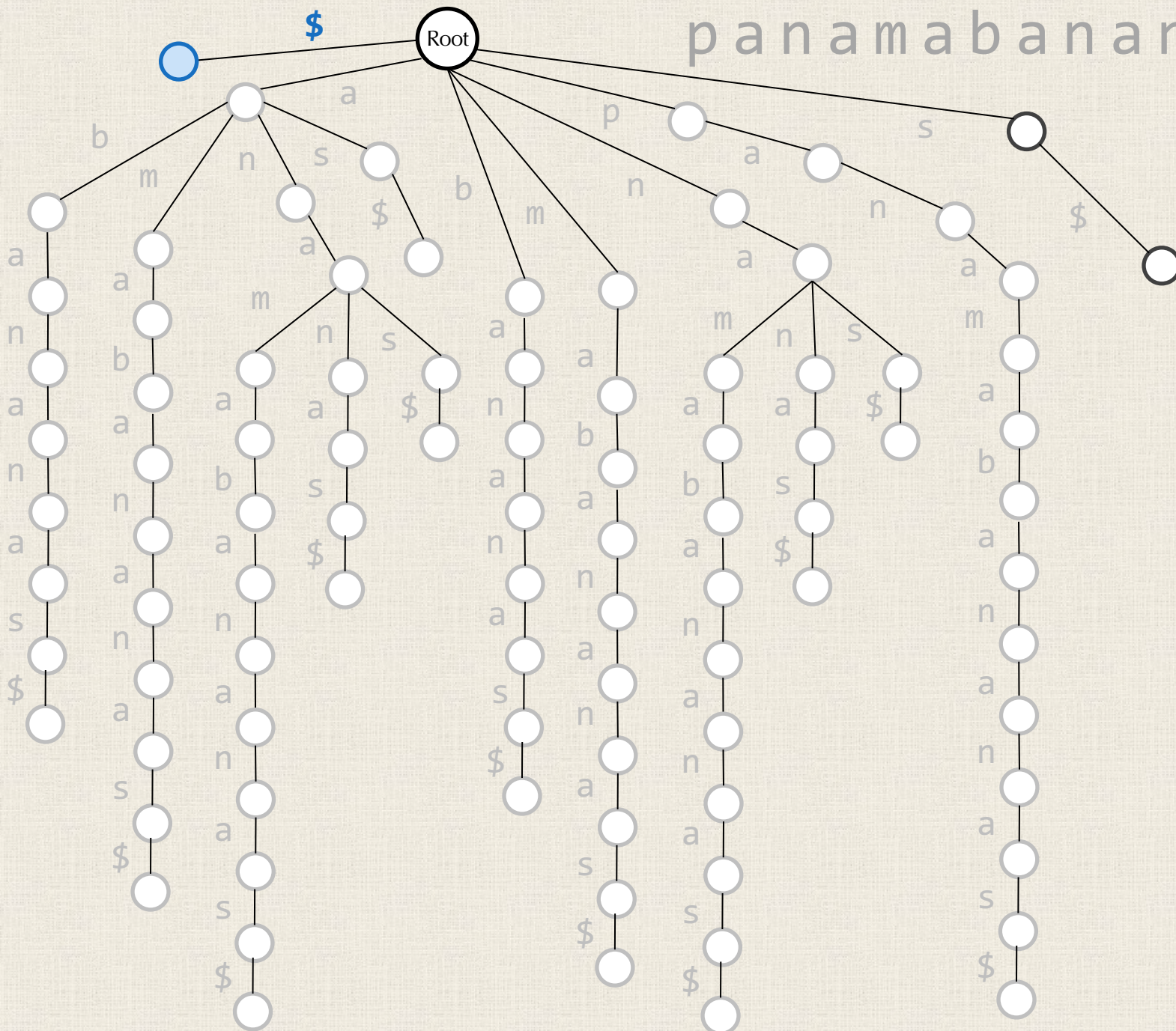
panamabanas\$

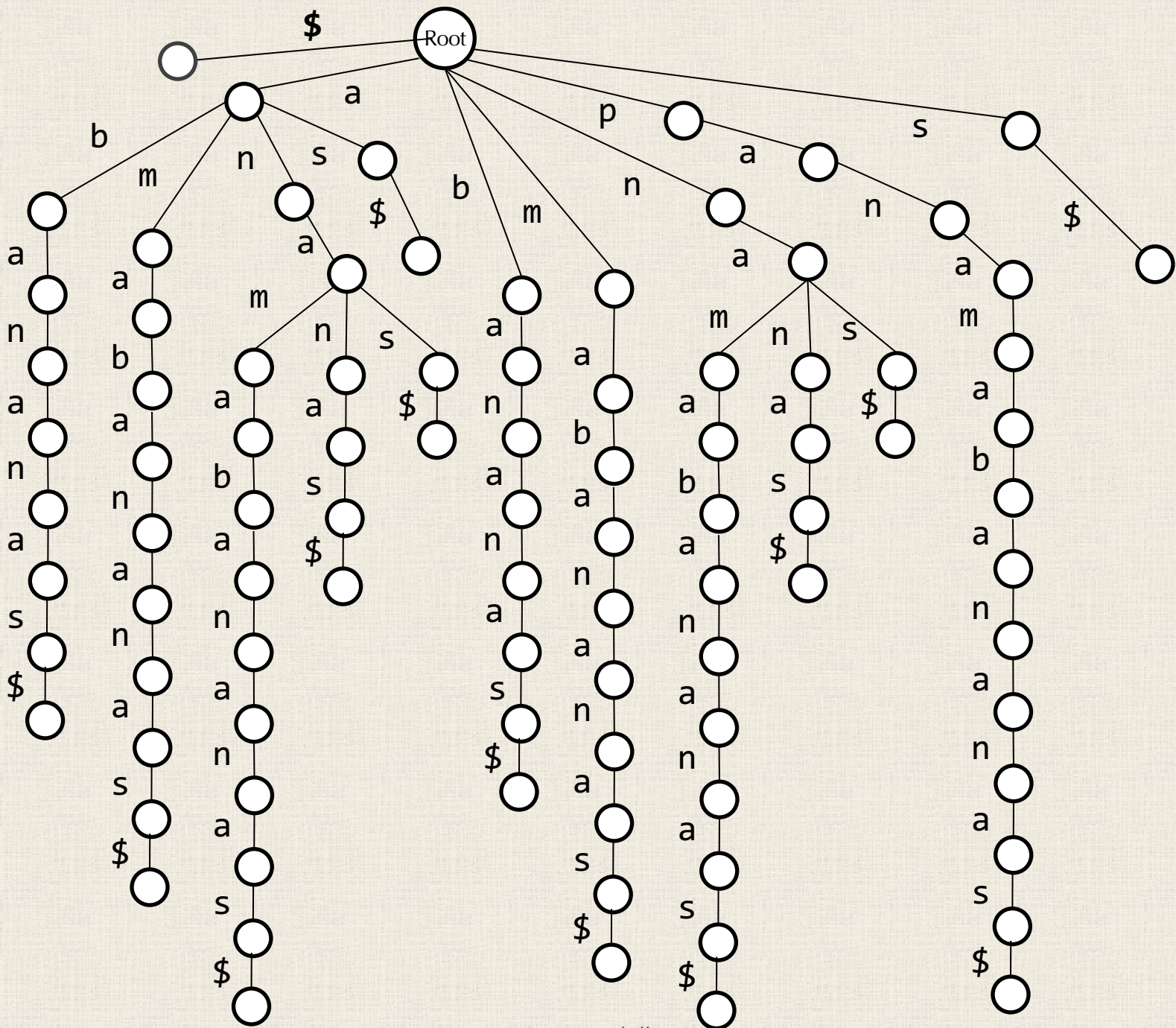


panamabanana s\$

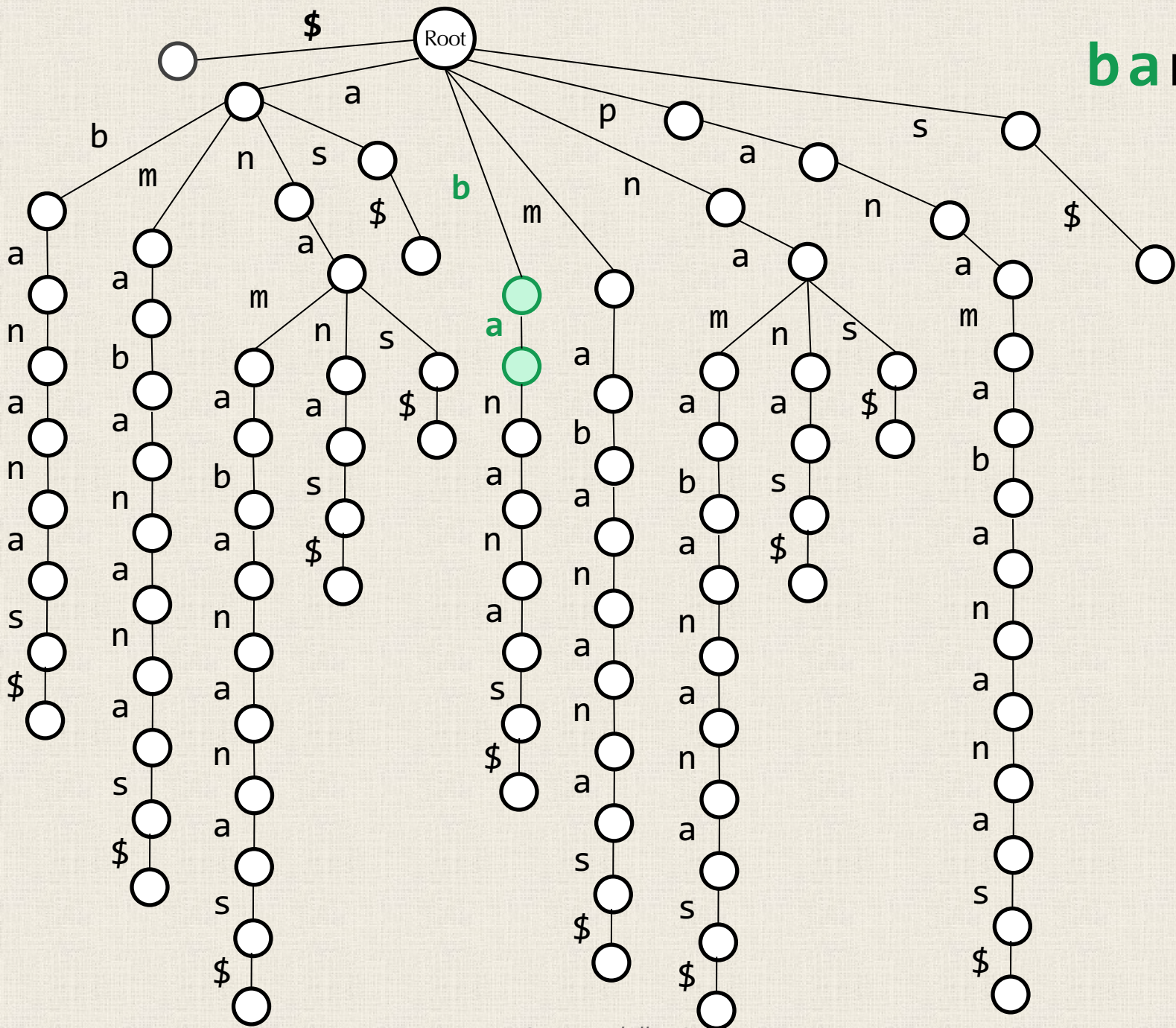


panamabananas\$

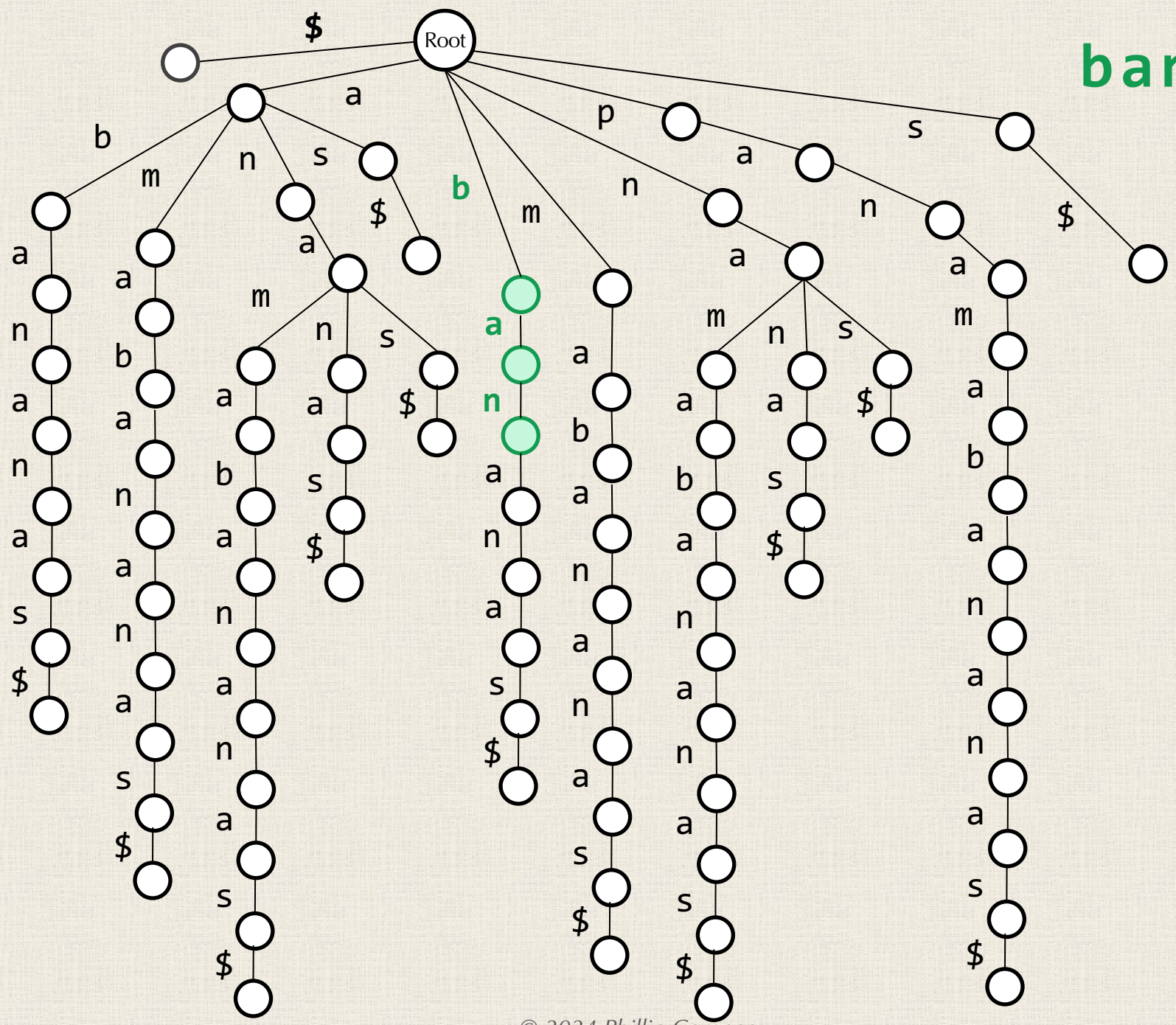




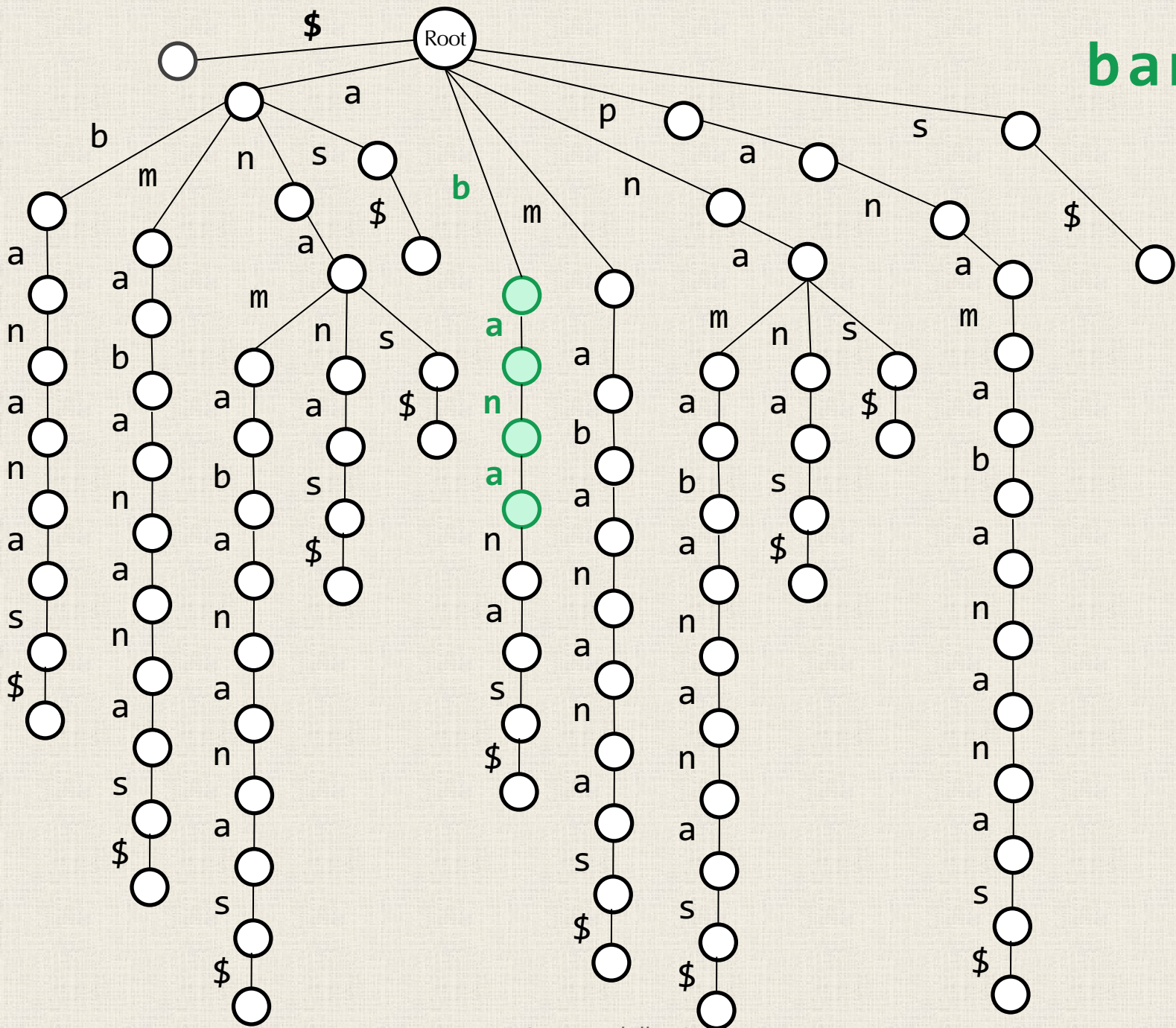
banana



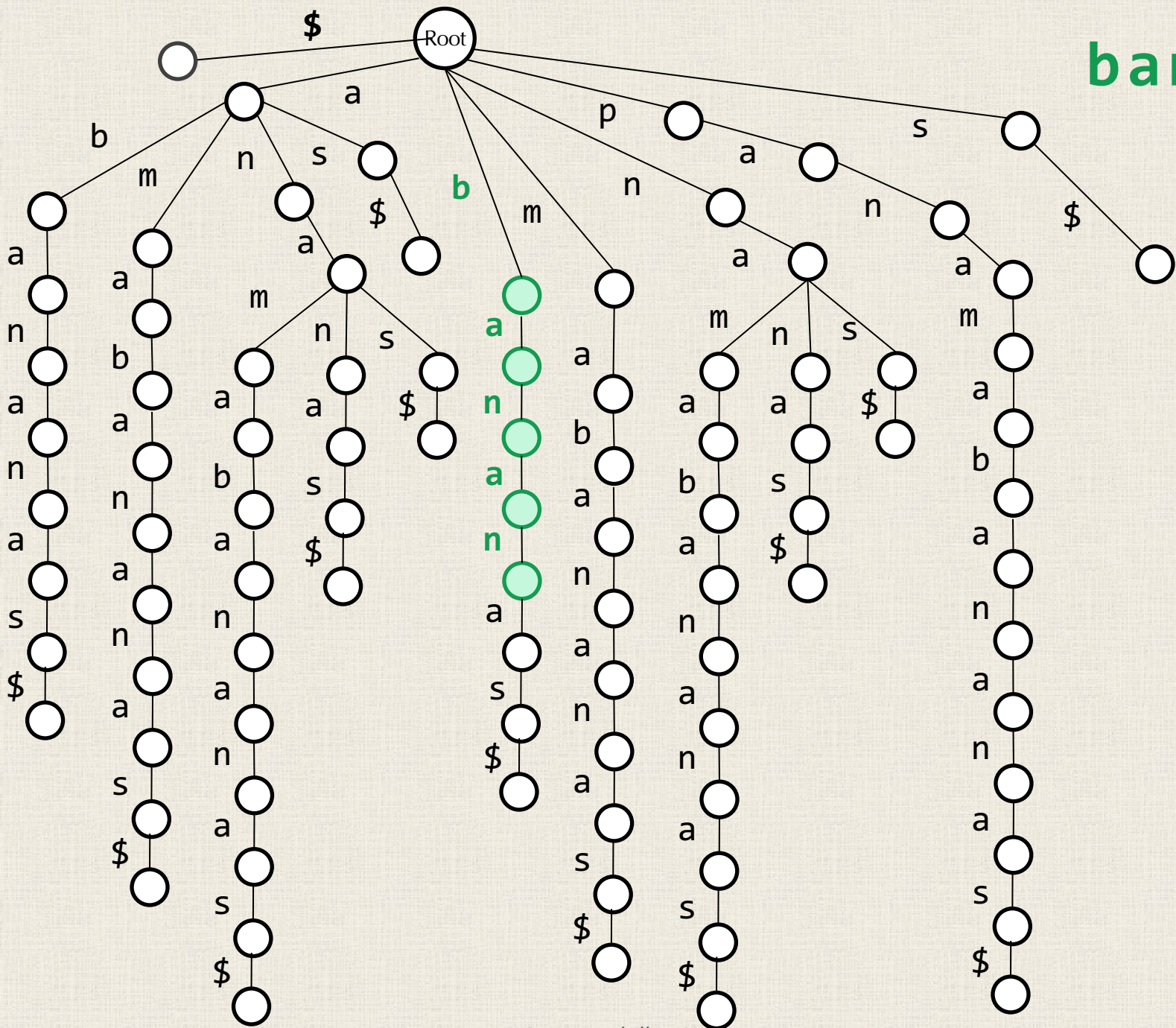
banana

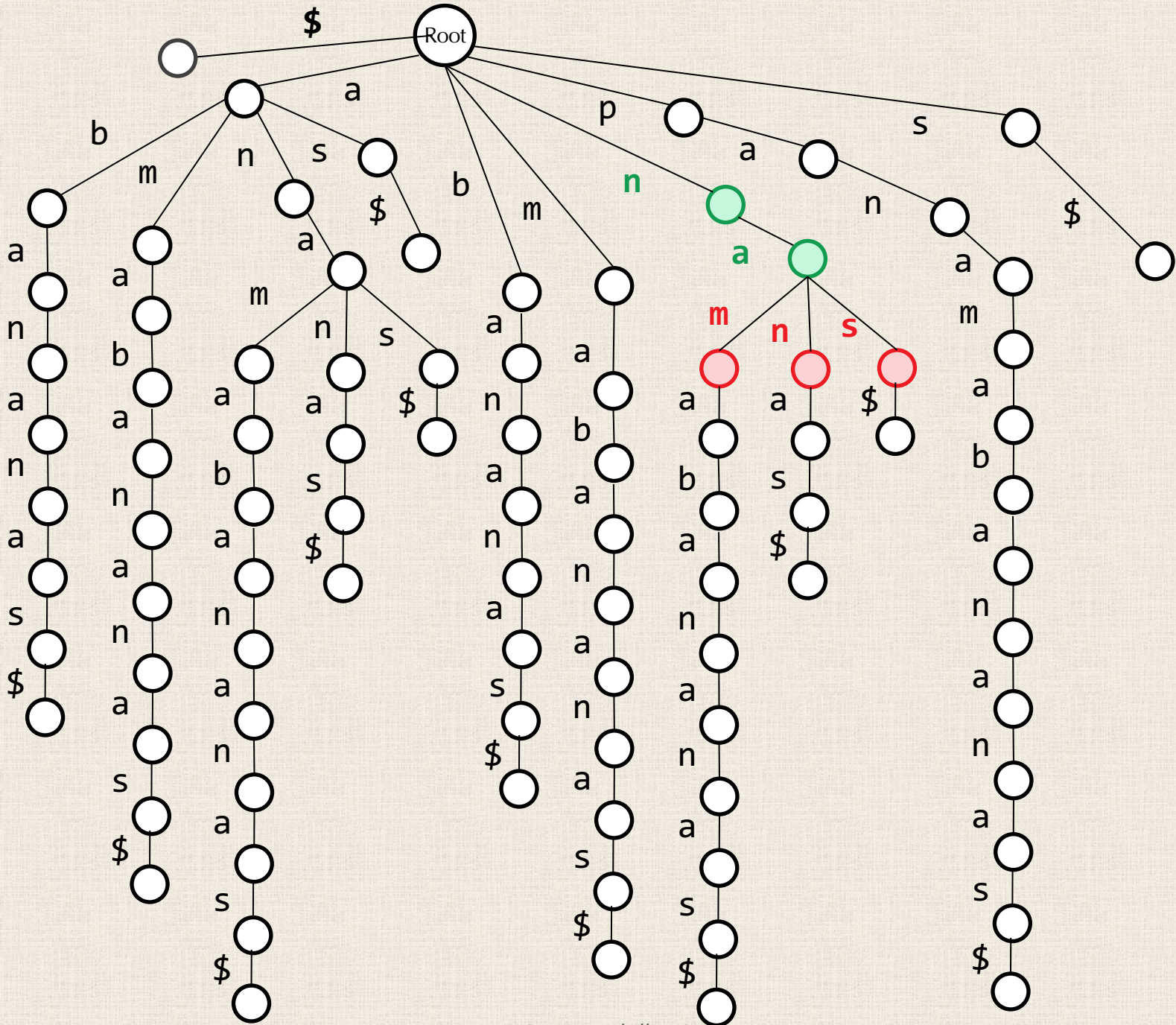


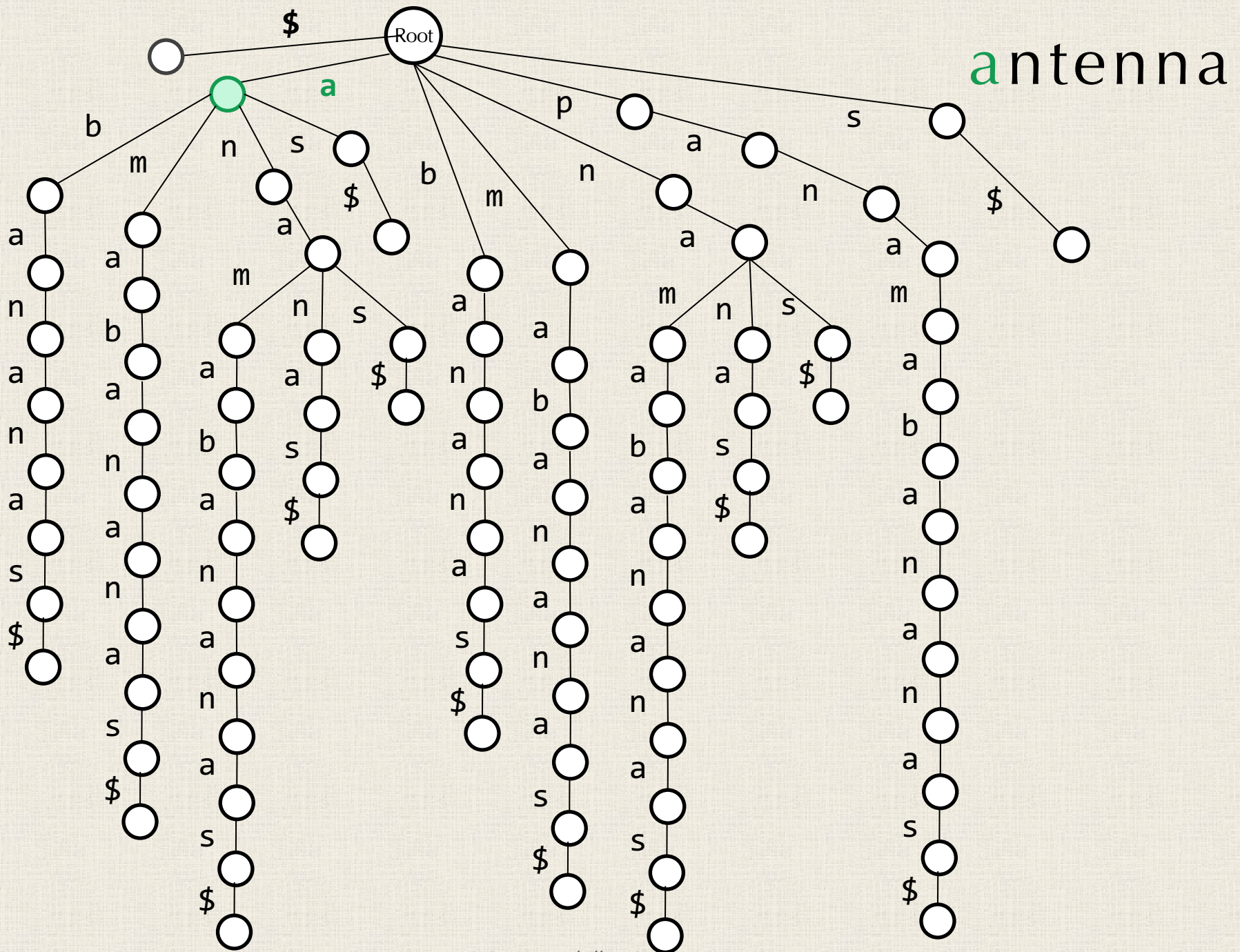
banana

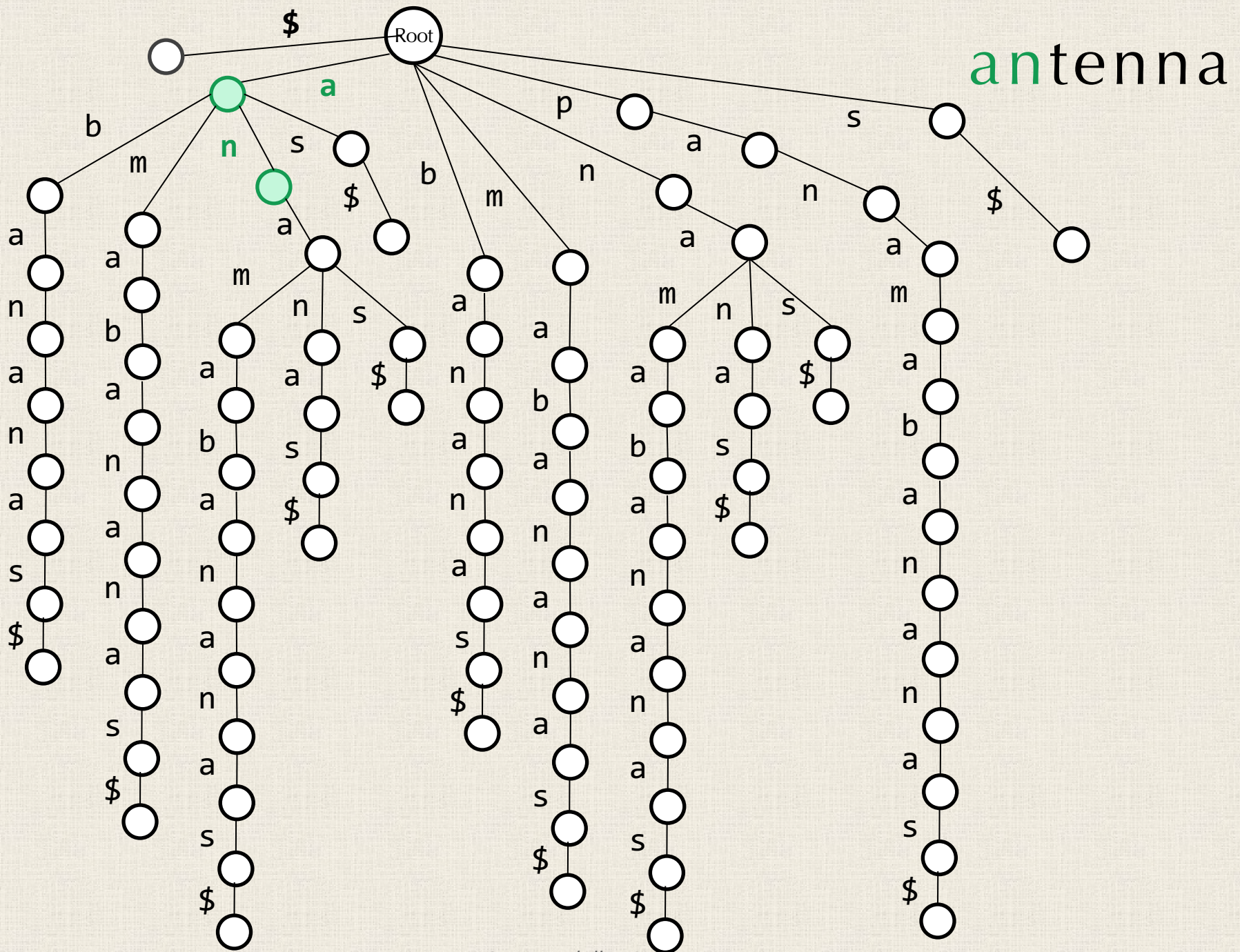


banana

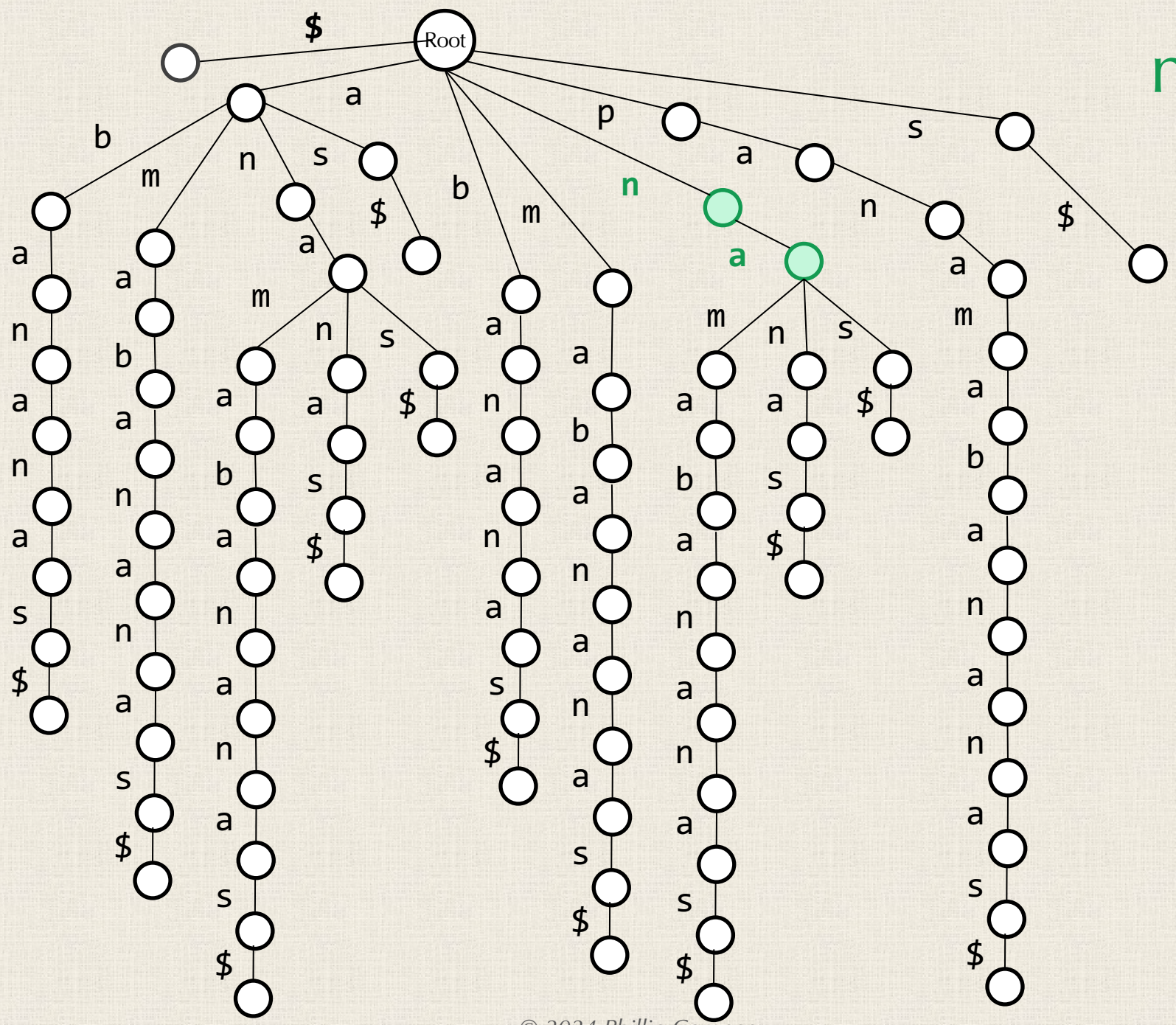




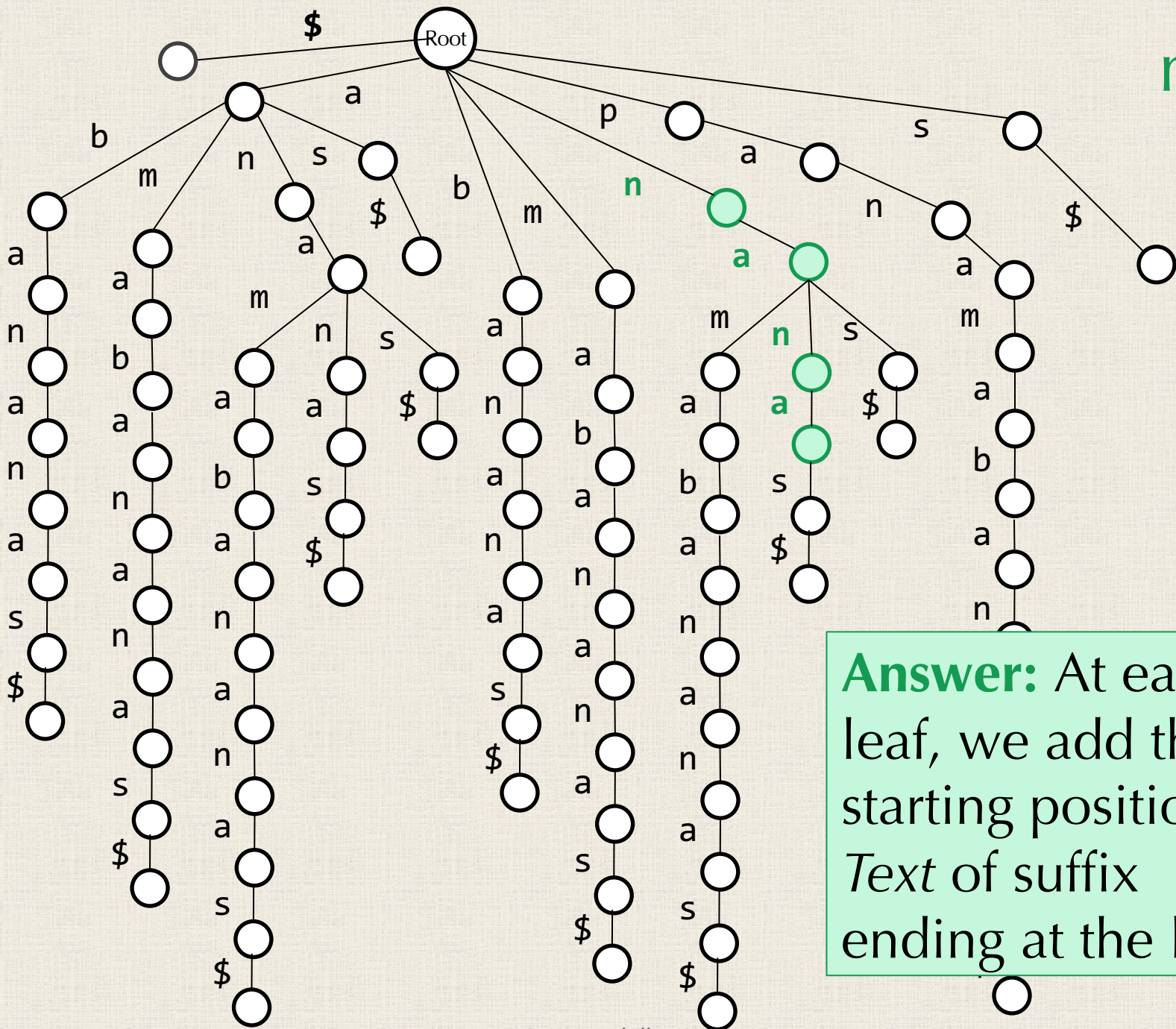




nana

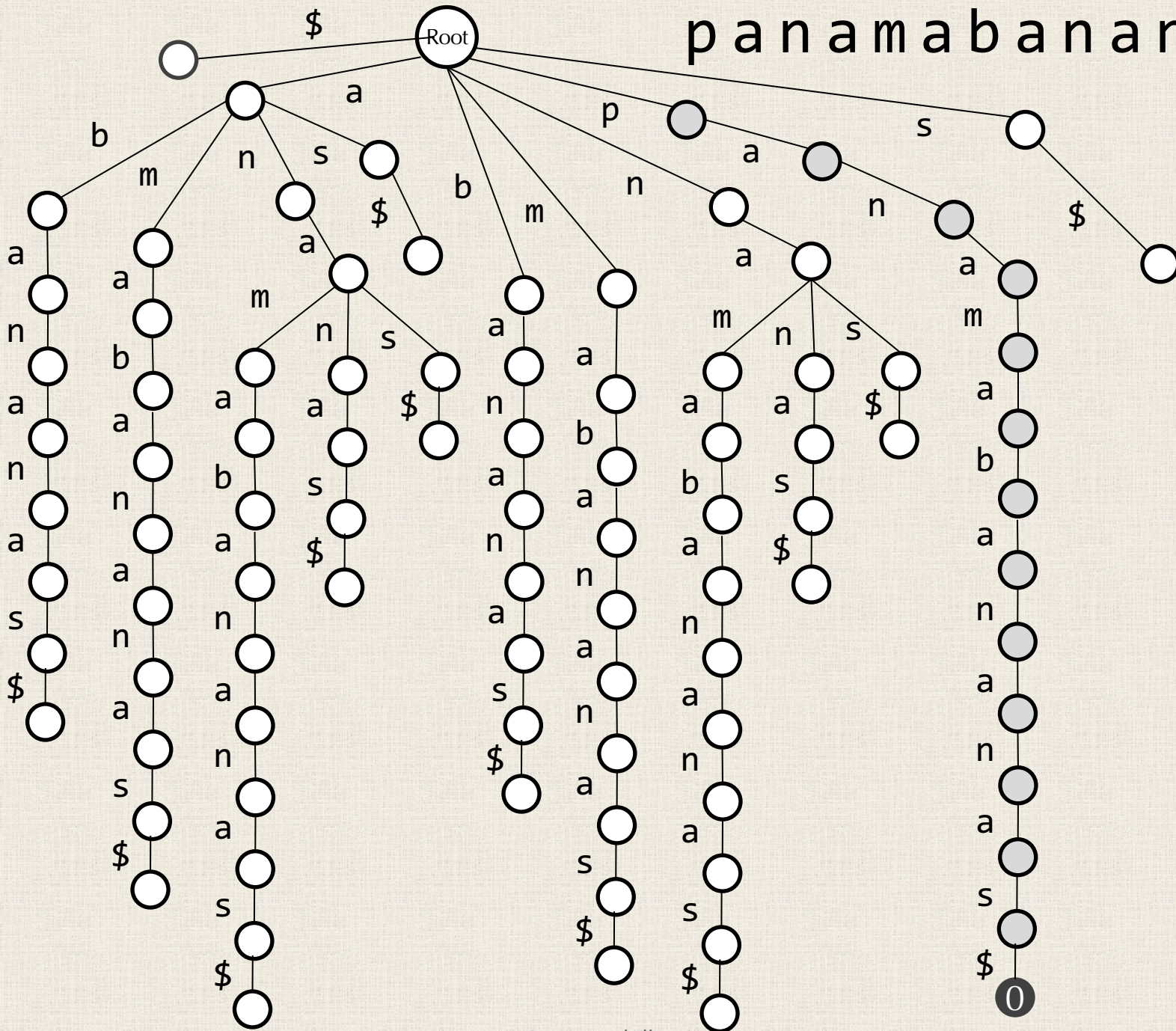


nana

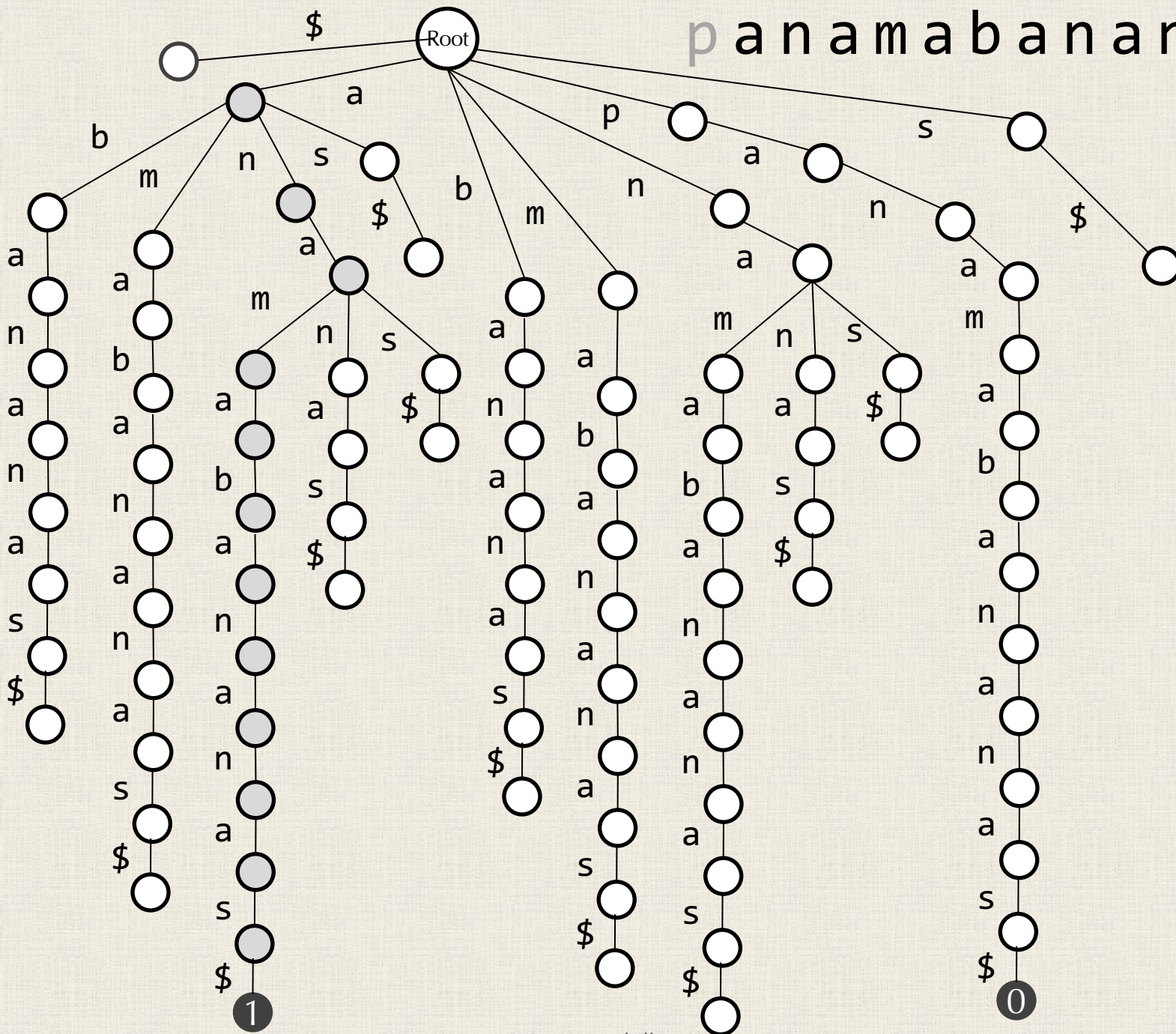


Answer: At each leaf, we add the starting position in *Text* of suffix ending at the leaf.

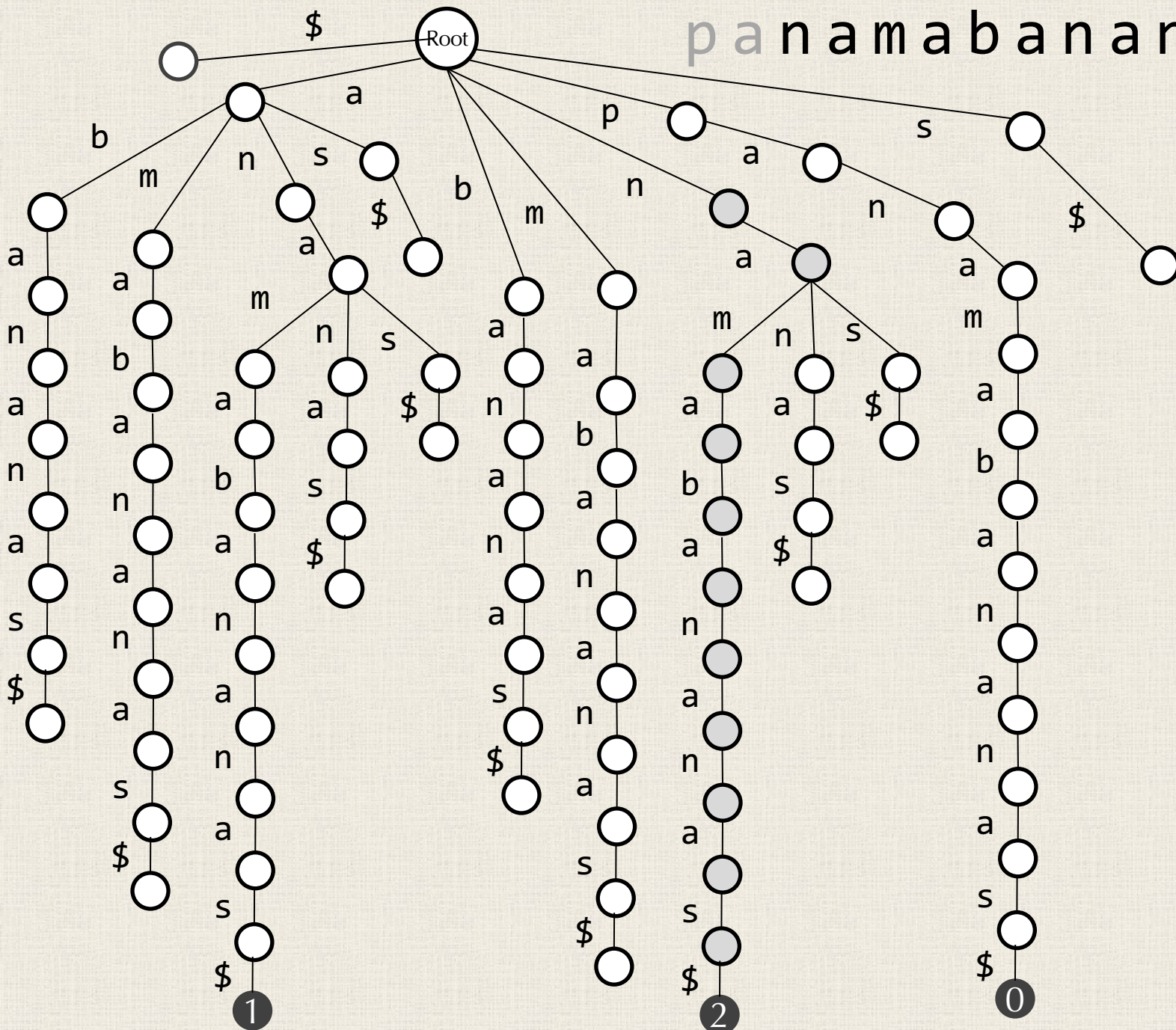
panamabananas\$



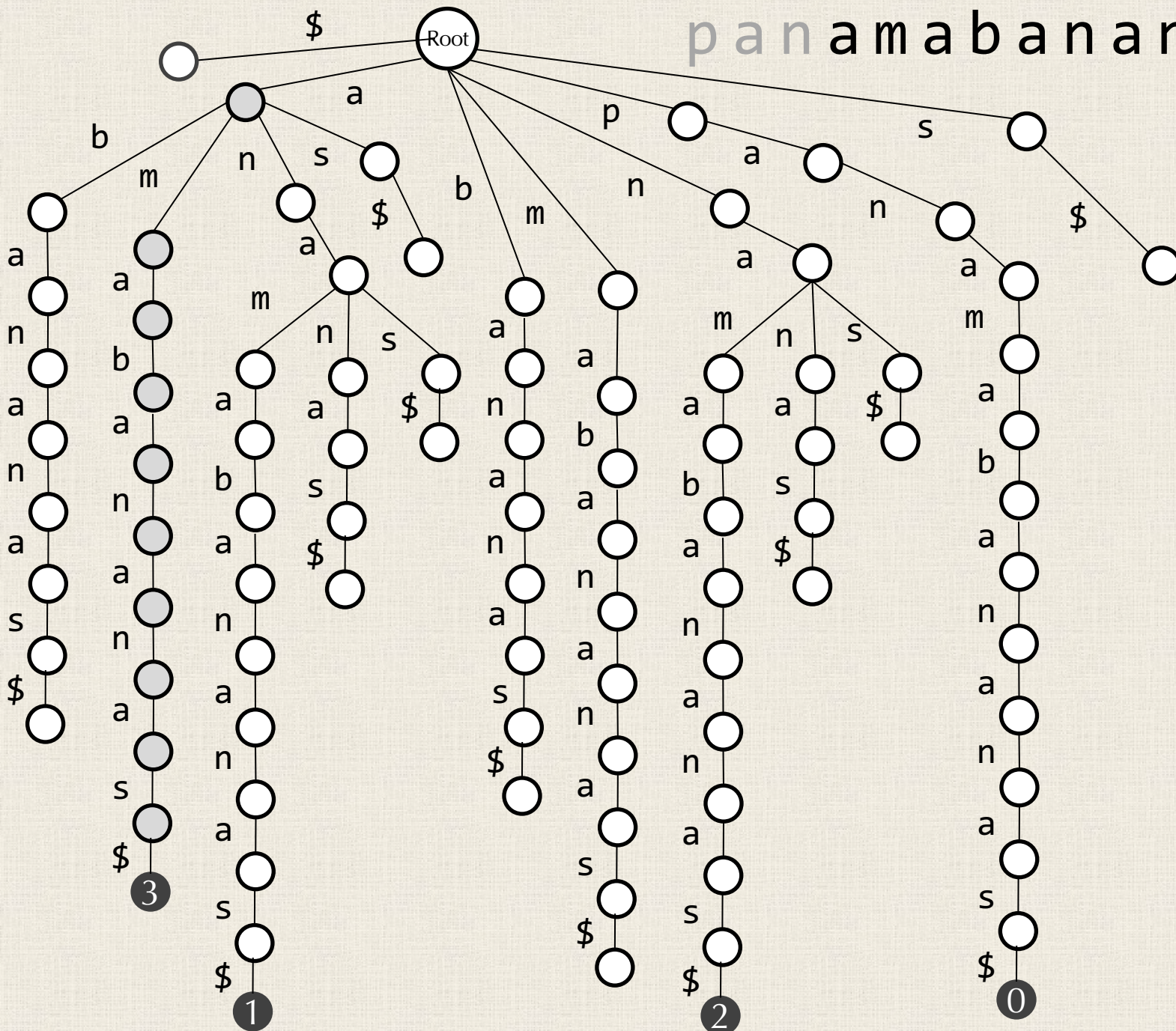
panamabananas\$



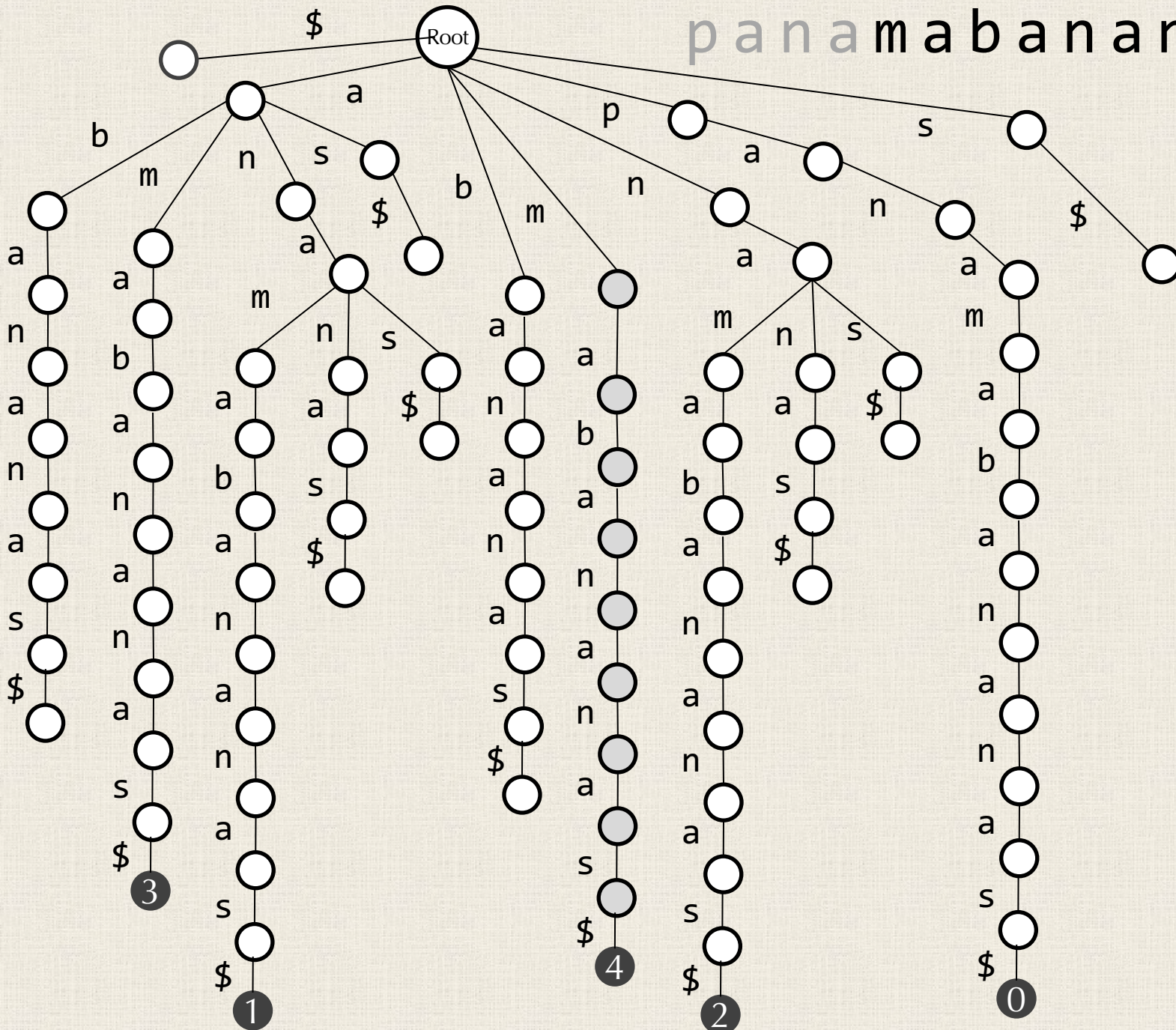
panamabananas\$



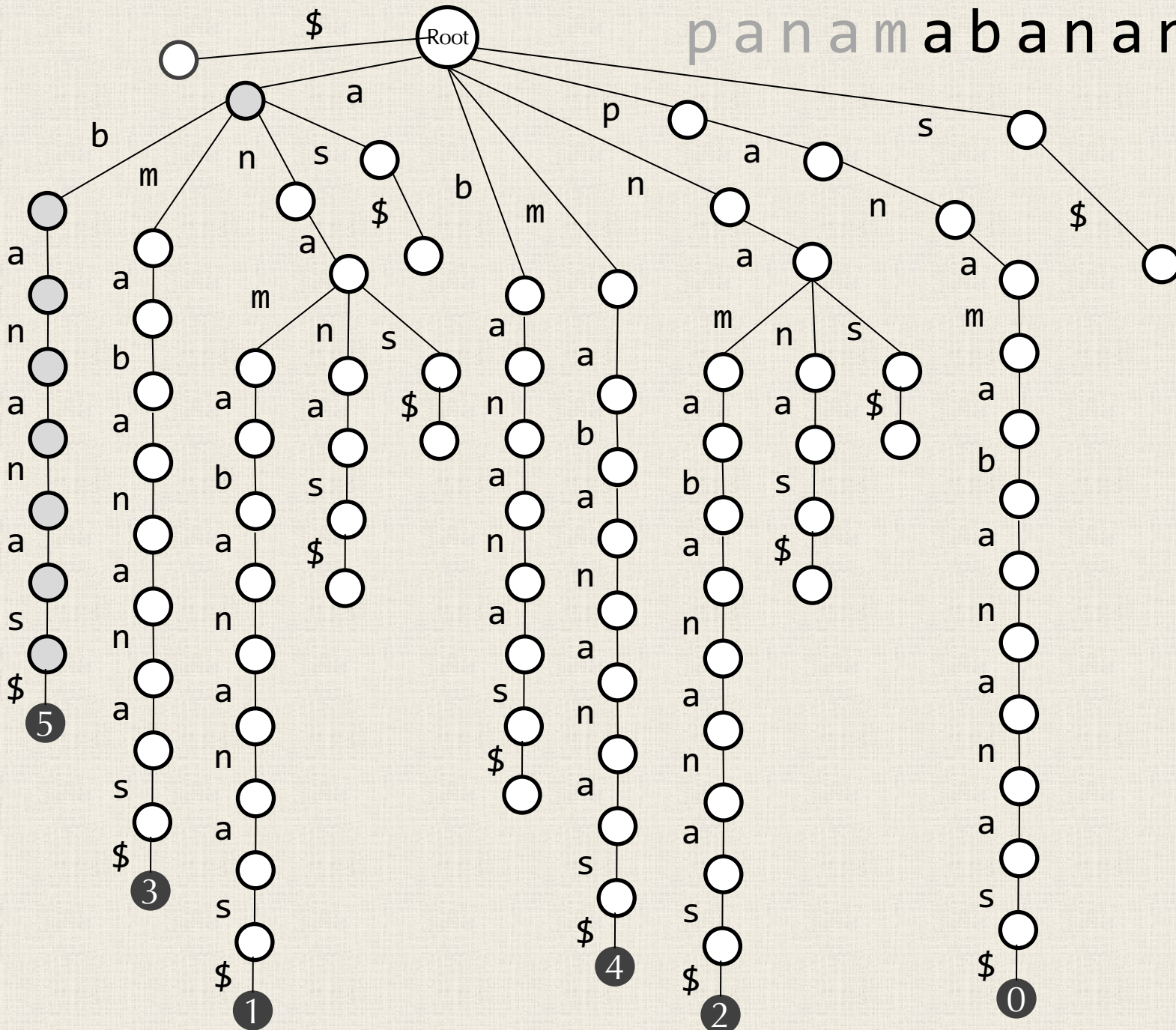
panamabananas\$



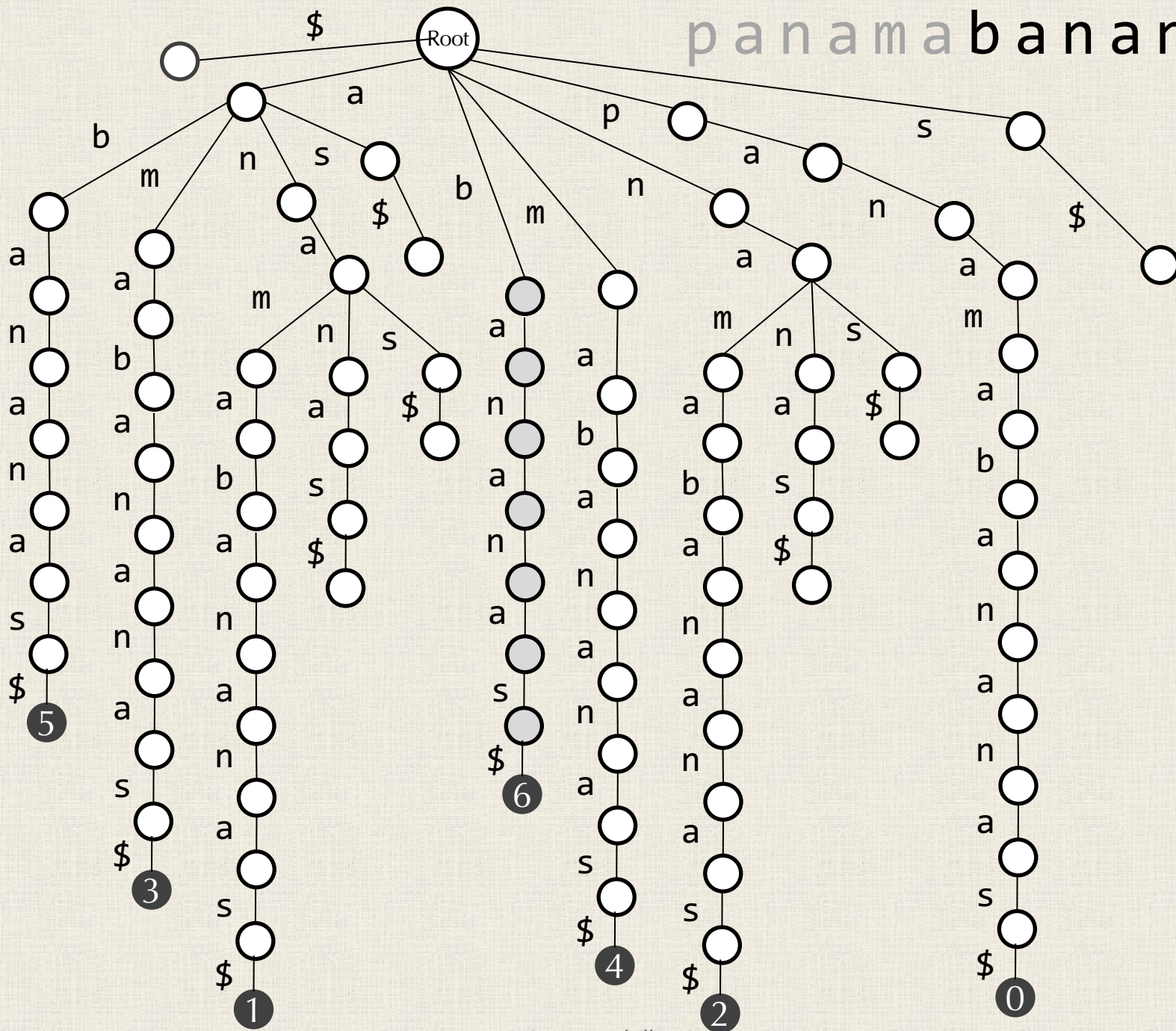
panamabananas\$



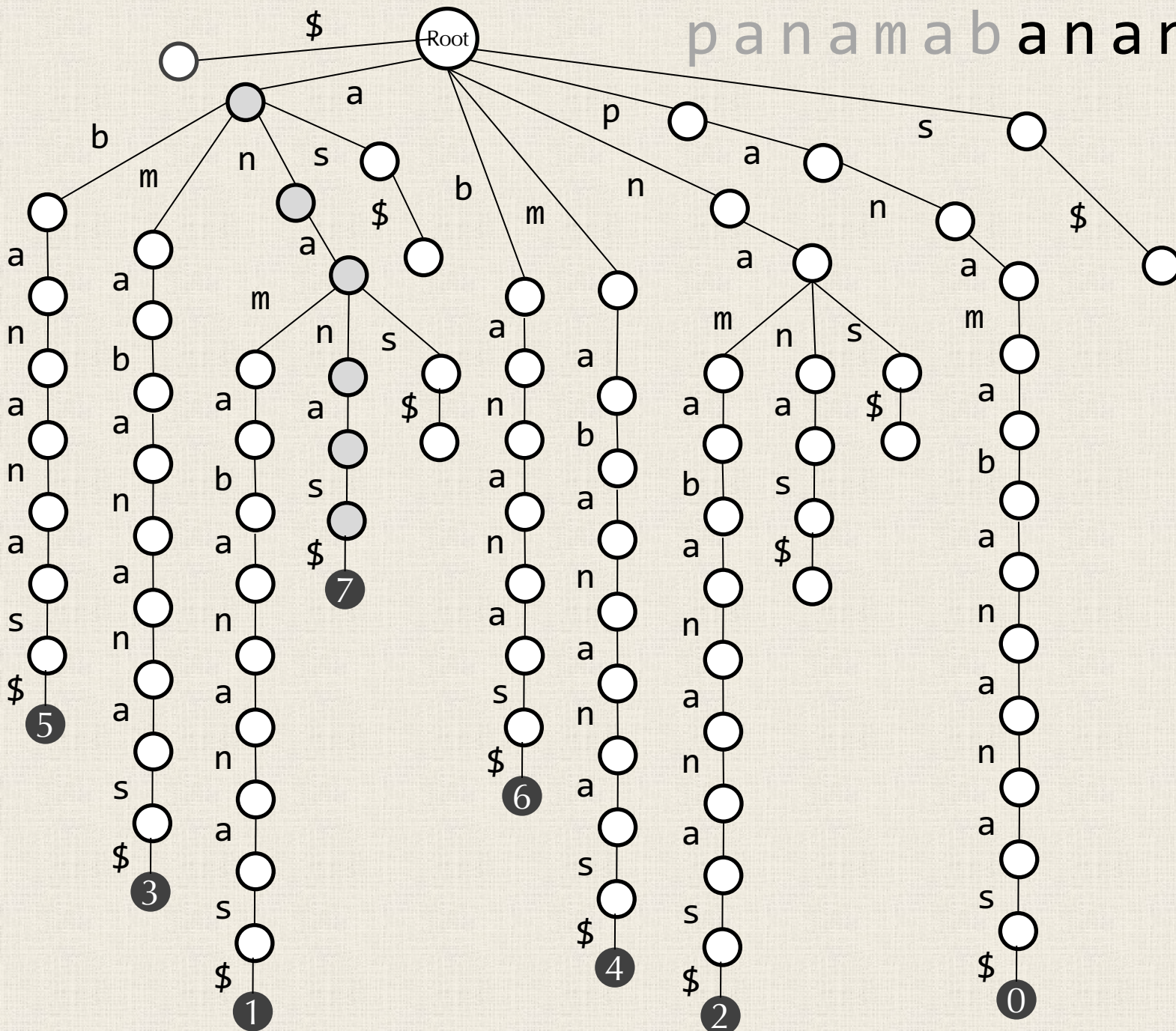
panamabananas\$



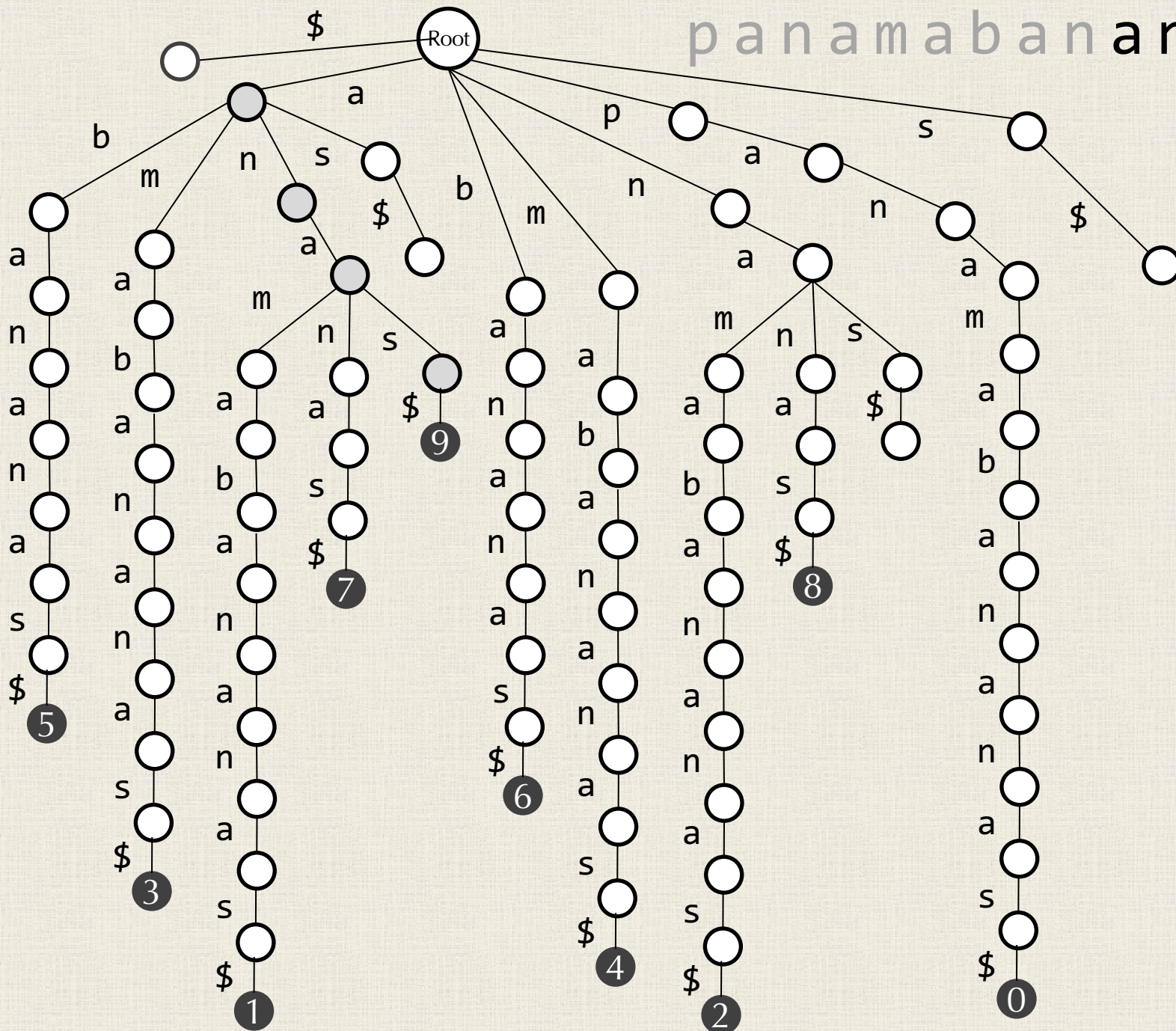
panamabananas\$



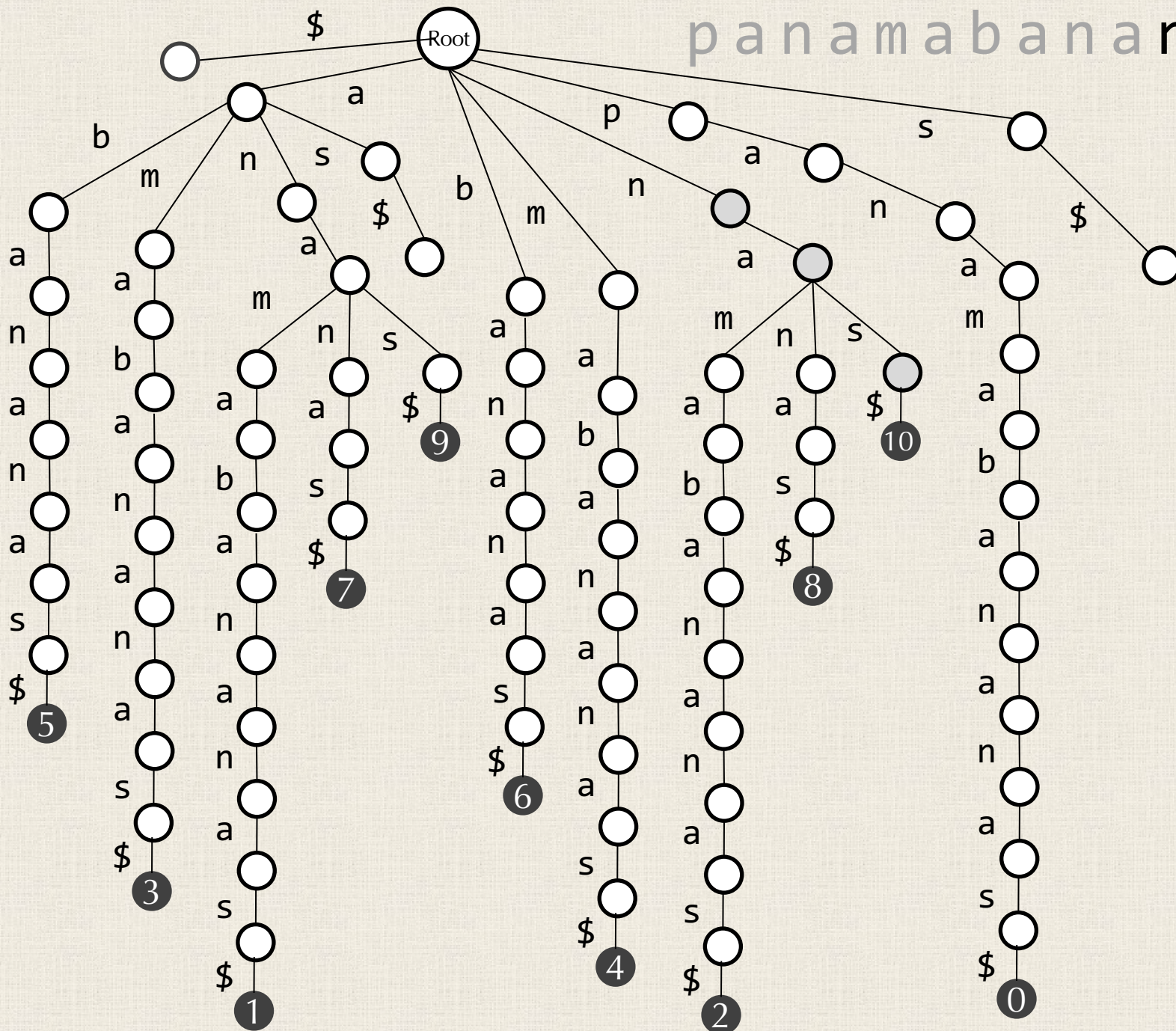
panamabananas\$



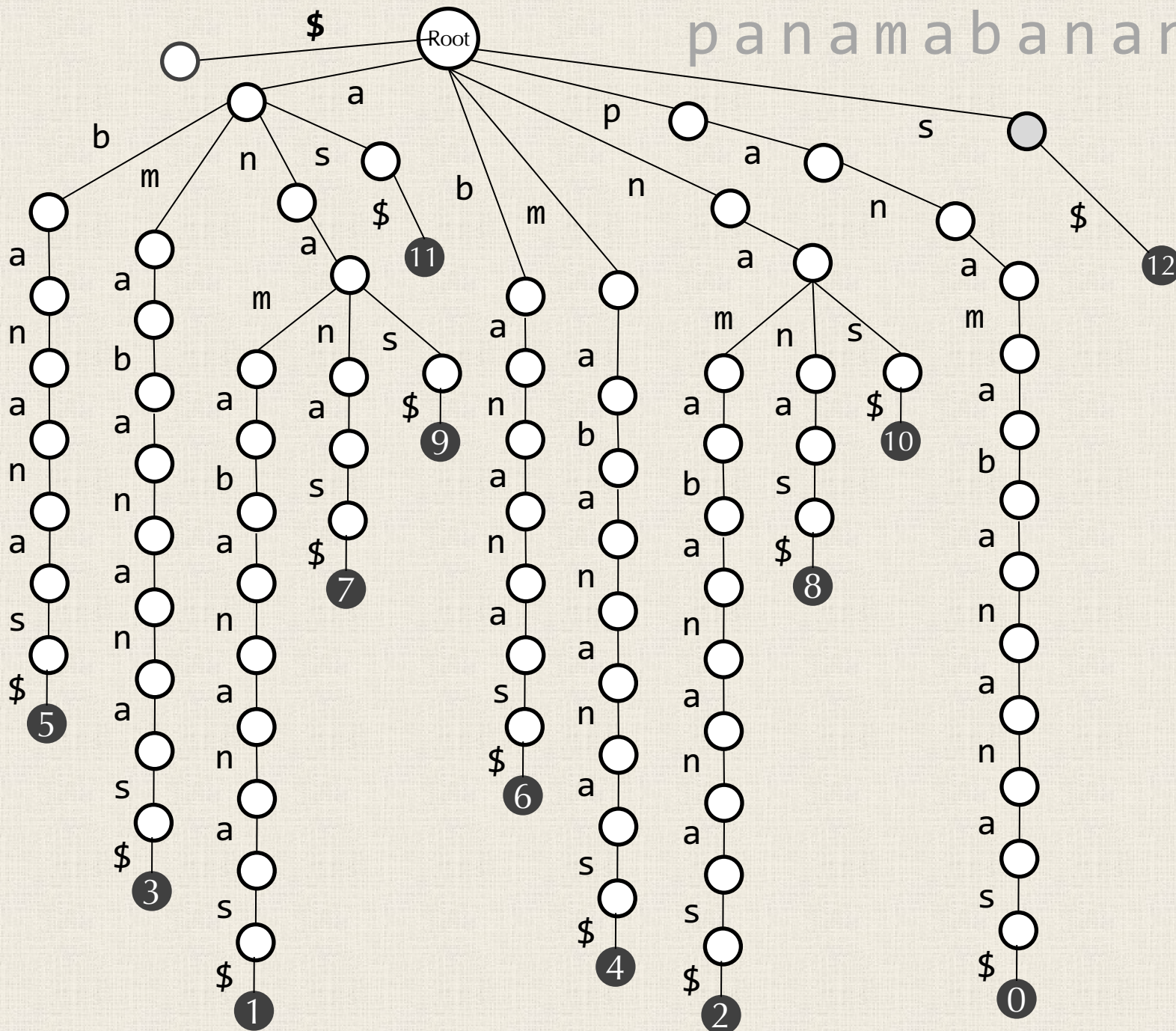
panamabananas\$



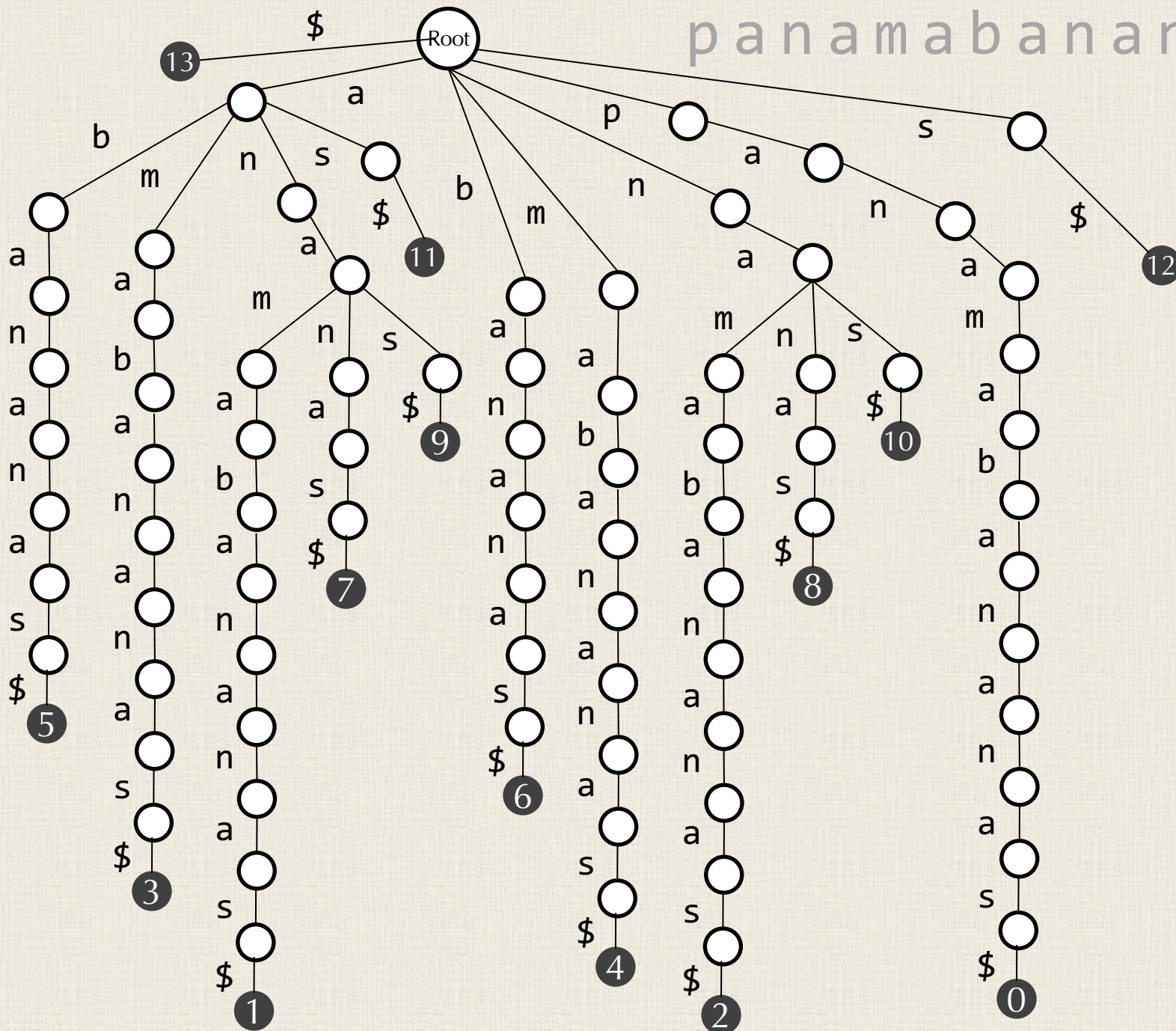
panamabananas\$



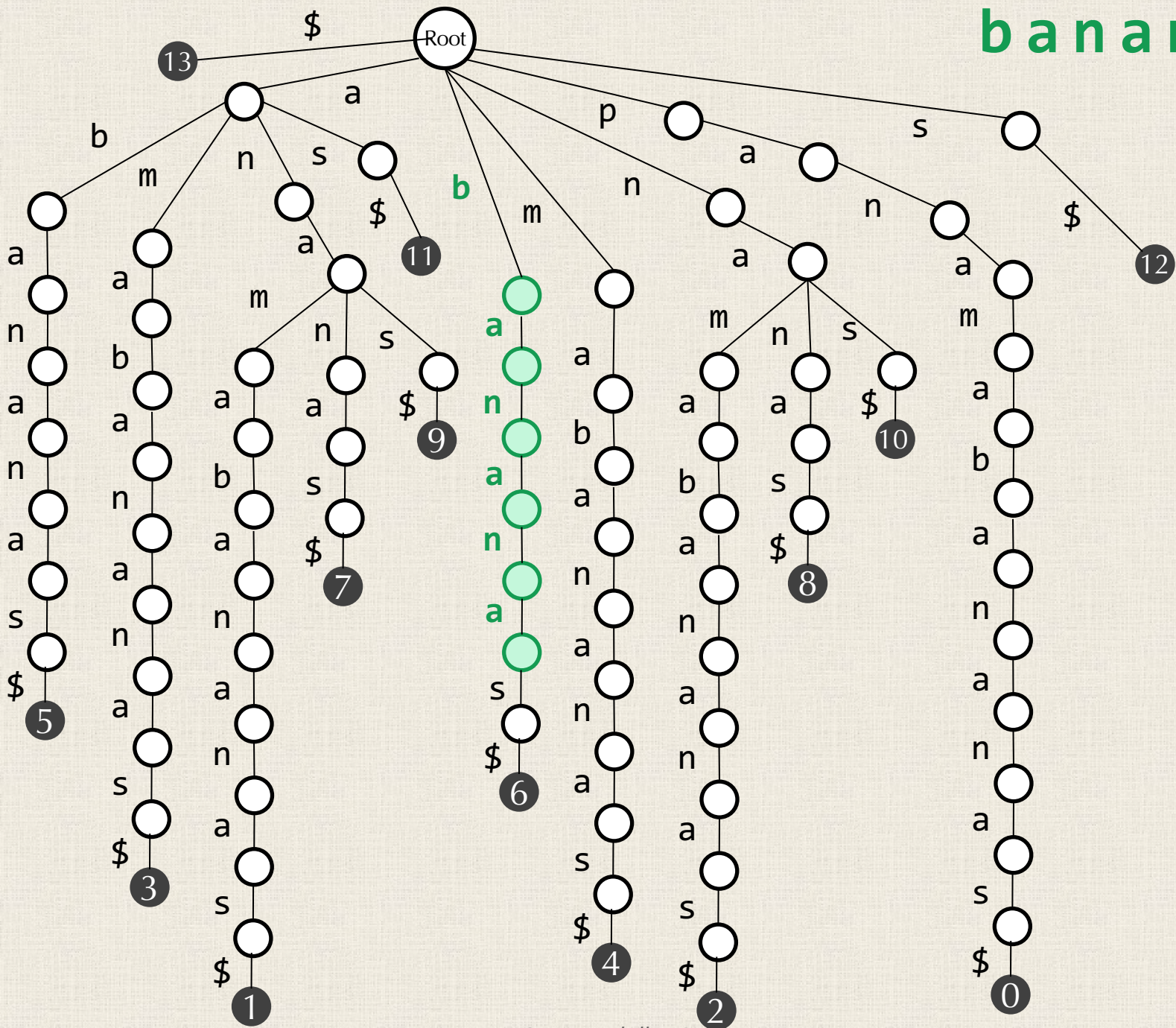
panamabananas\$



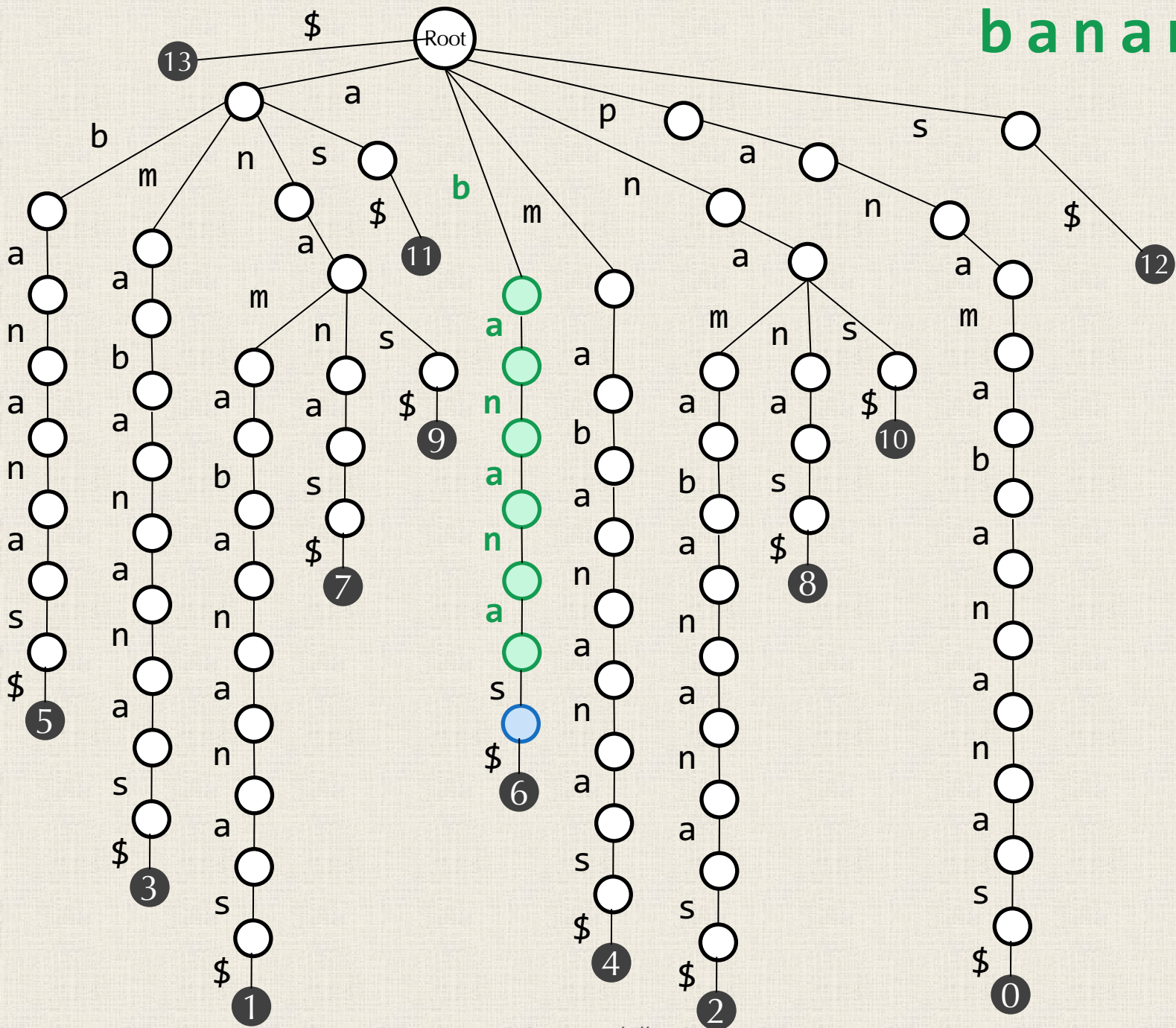
panamabananas\$



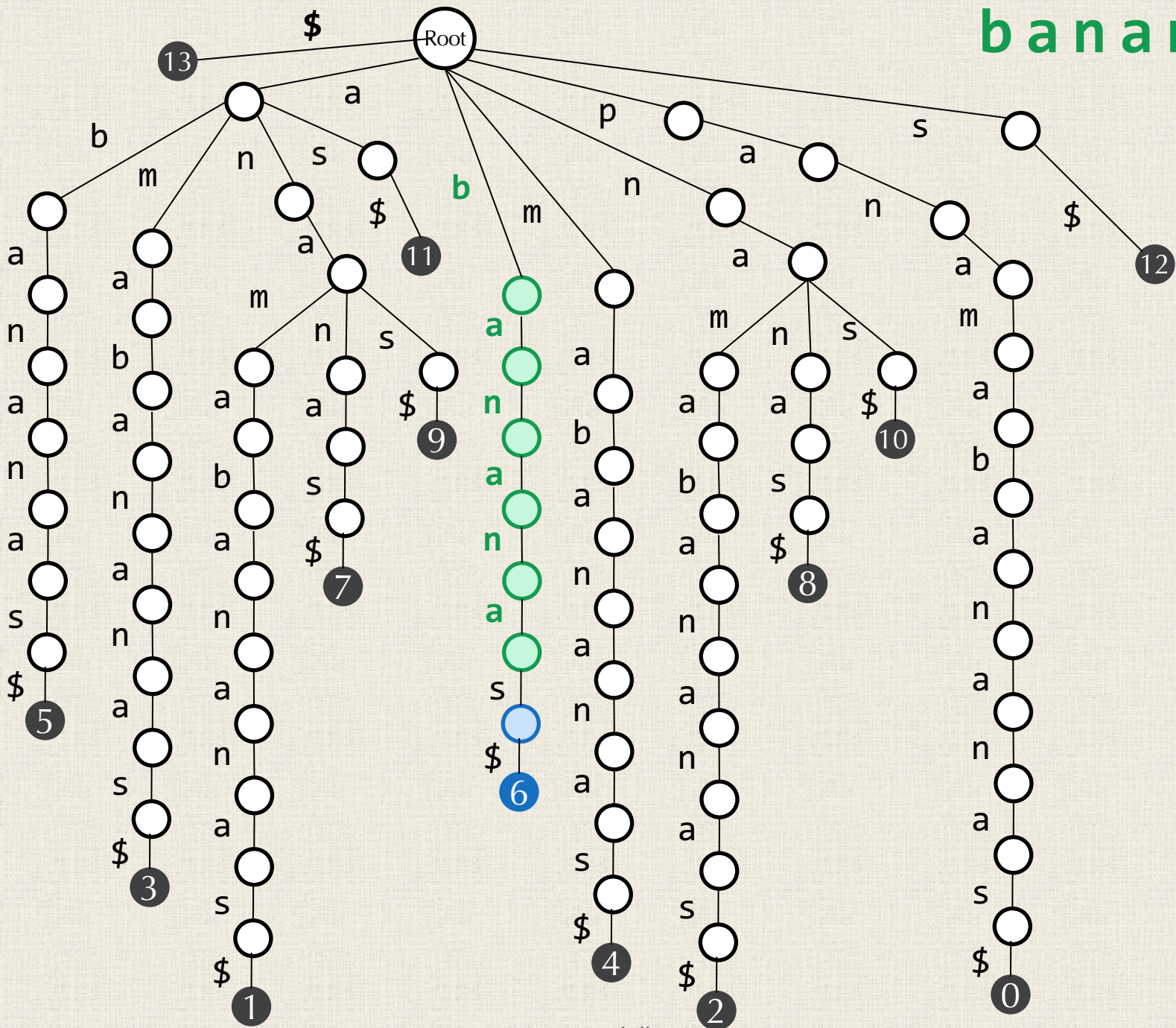
banana



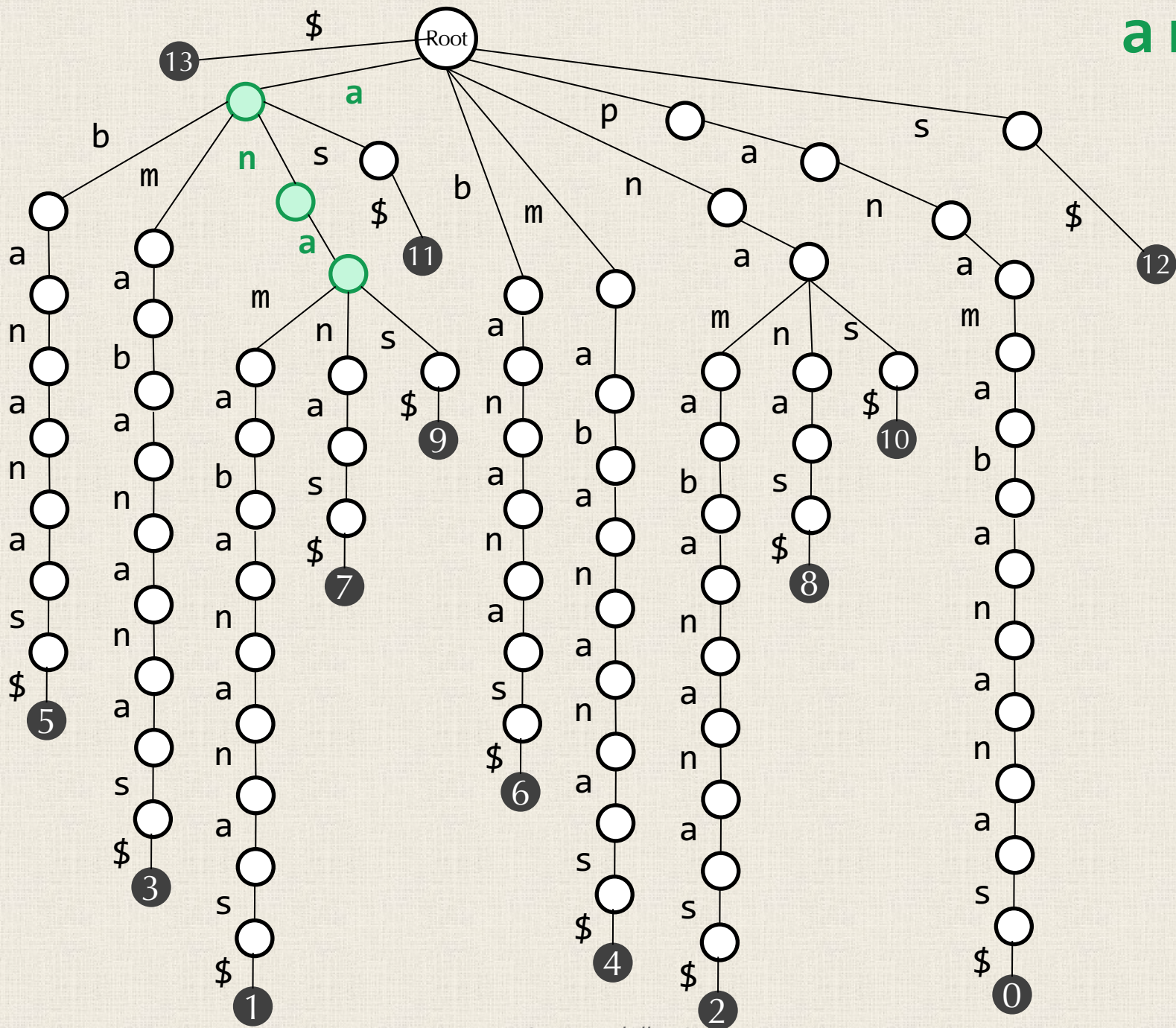
banana



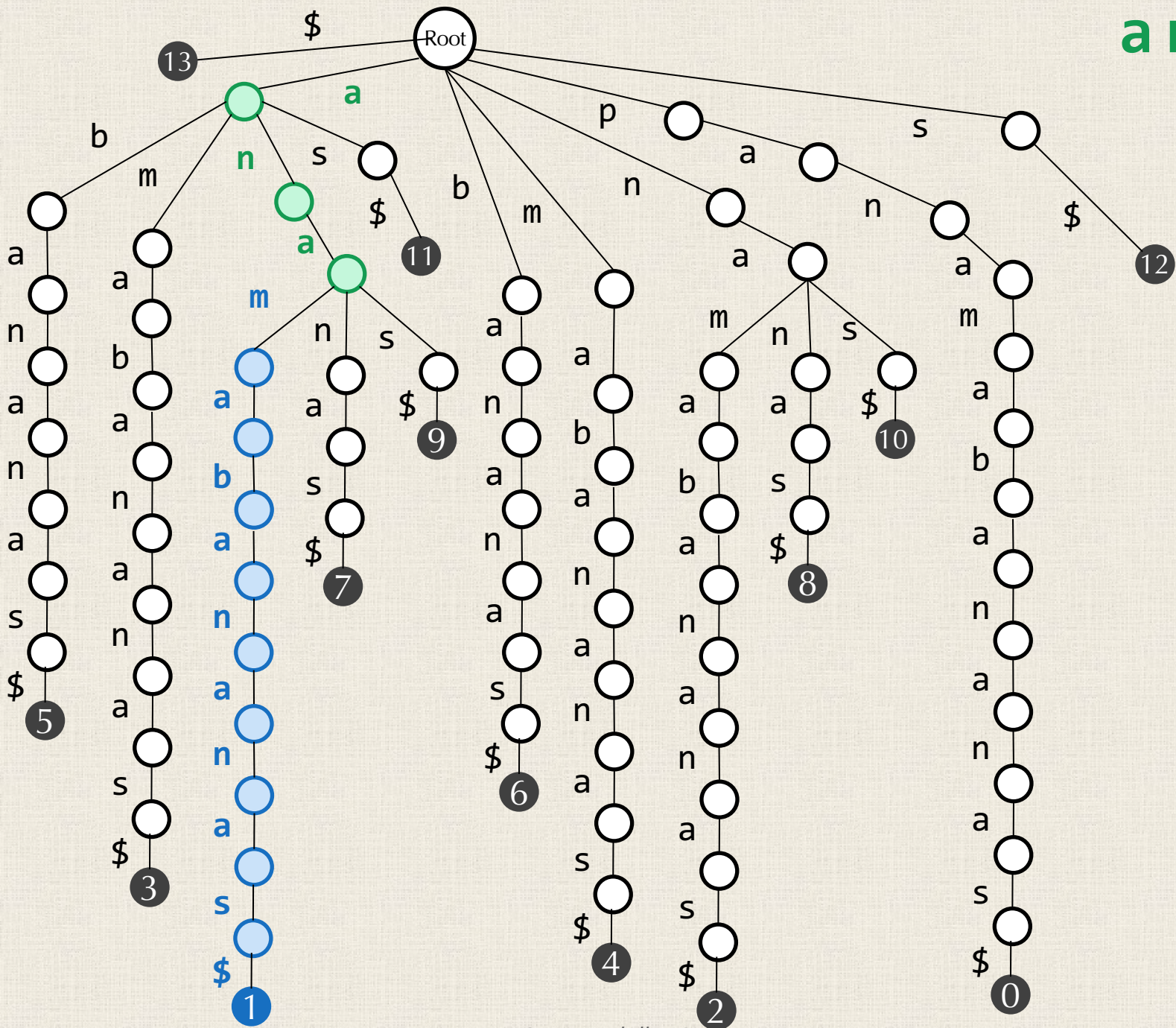
banana



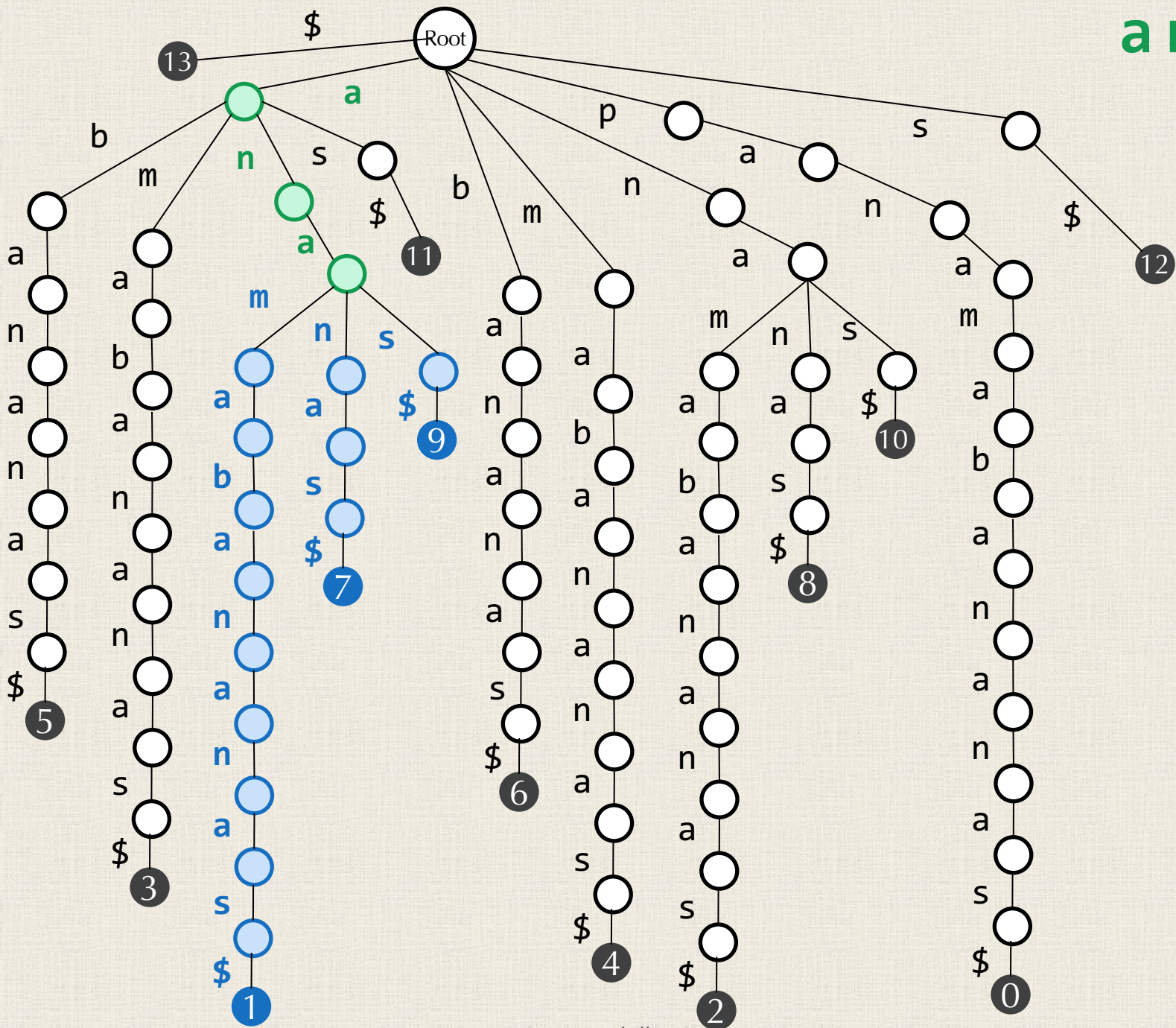
ana



ana



ana



We Have Memory Troubles Again

Suffixes

Constructing a trie for a collection of *Patterns* takes $O(|Patterns|)$ runtime and memory.

Constructing a suffix trie takes $O(Suffixes(Text)) = O(|Text|^2)$ runtime and memory.

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


SUFFIX TREES

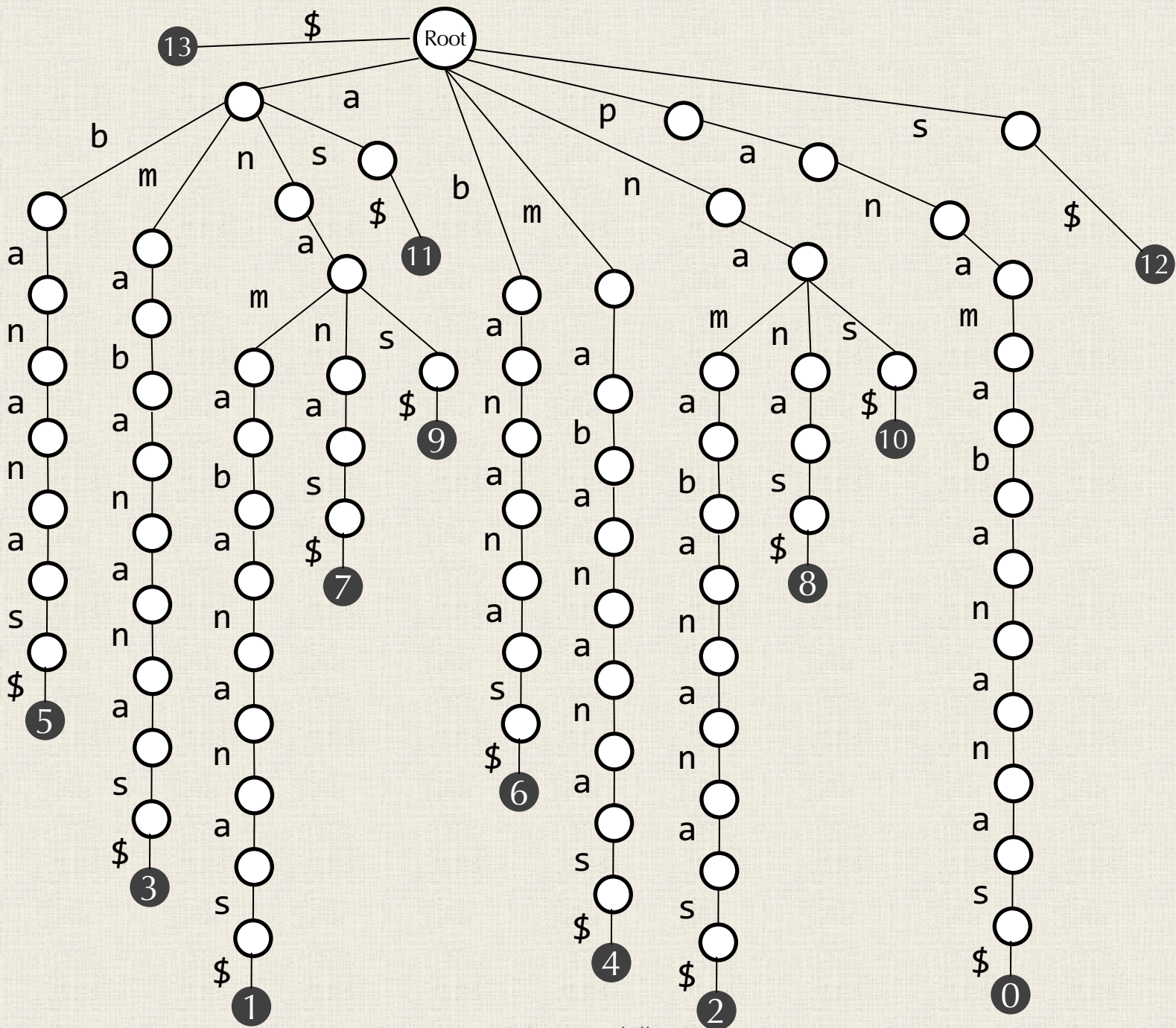
Compressing the Suffix Trie

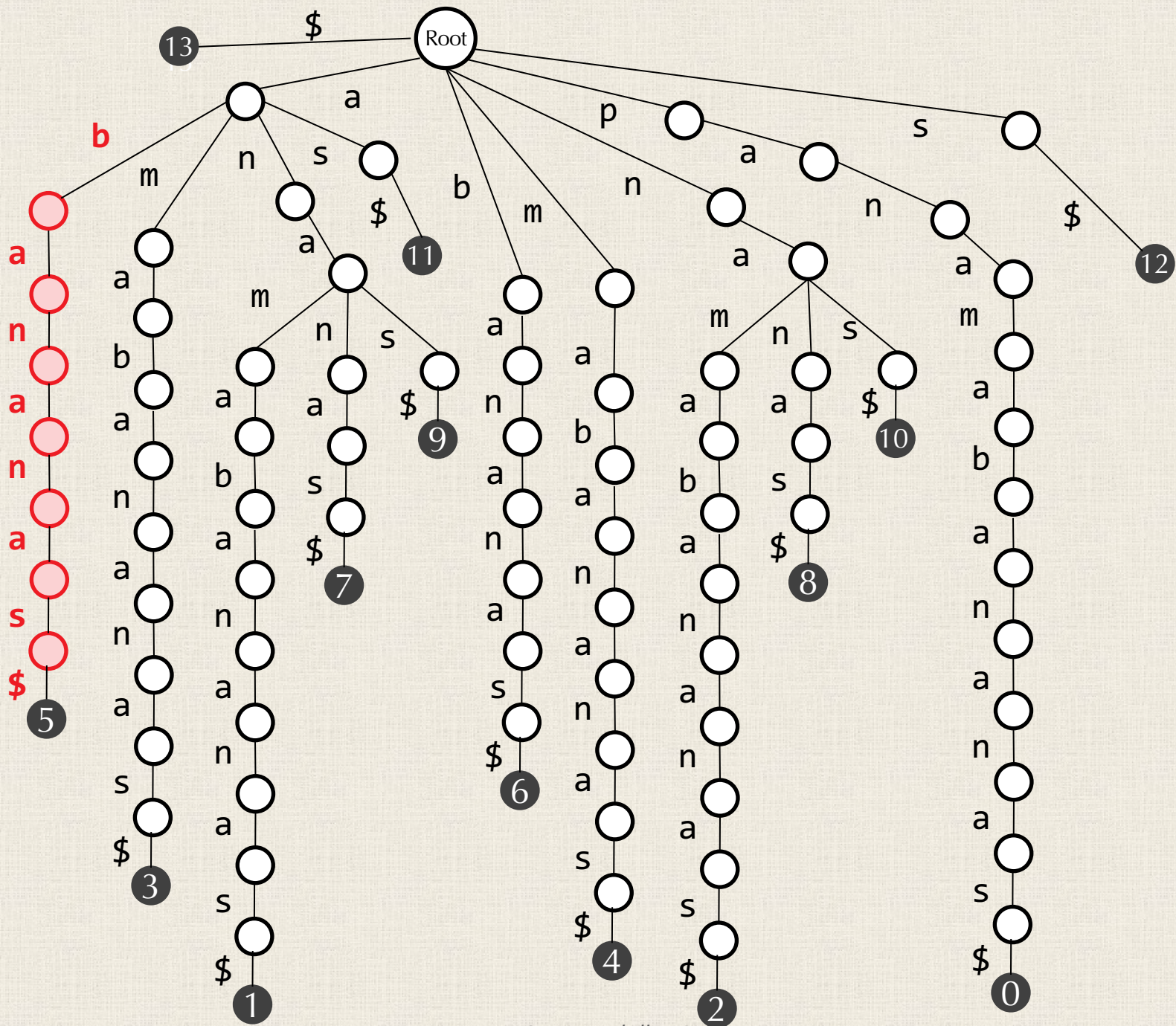
The suffix trie will only offer a memory improvement if $|Text|^2 < |Patterns|$ (i.e., for very short genomes and many reads). But we can compress the suffix trie into a smaller structure.

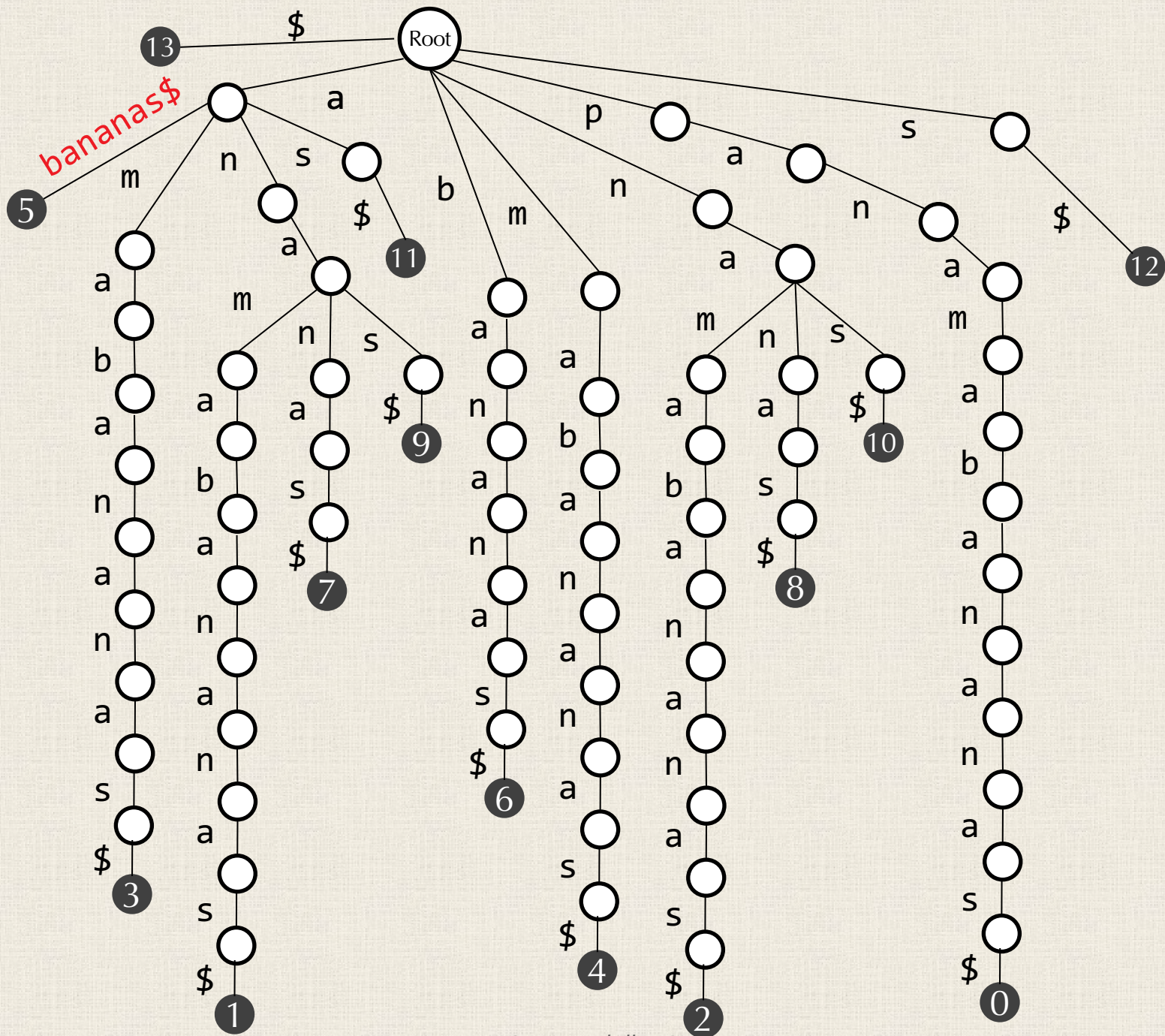
Compressing the Suffix Trie

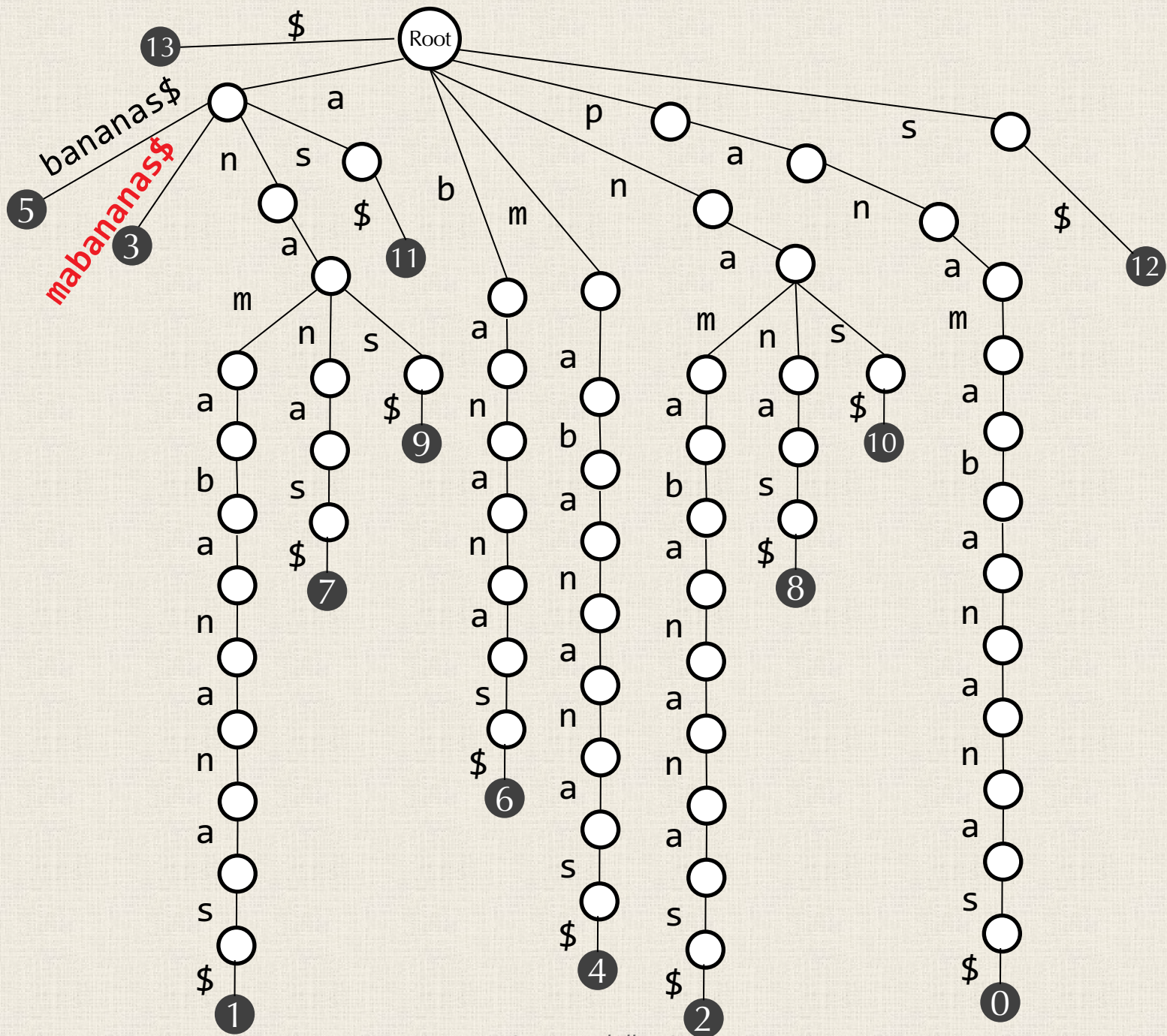
The suffix trie will only offer a memory improvement if $|Text|^2 < |Patterns|$ (i.e., for very short genomes and many reads). But we can compress the suffix trie into a smaller structure.

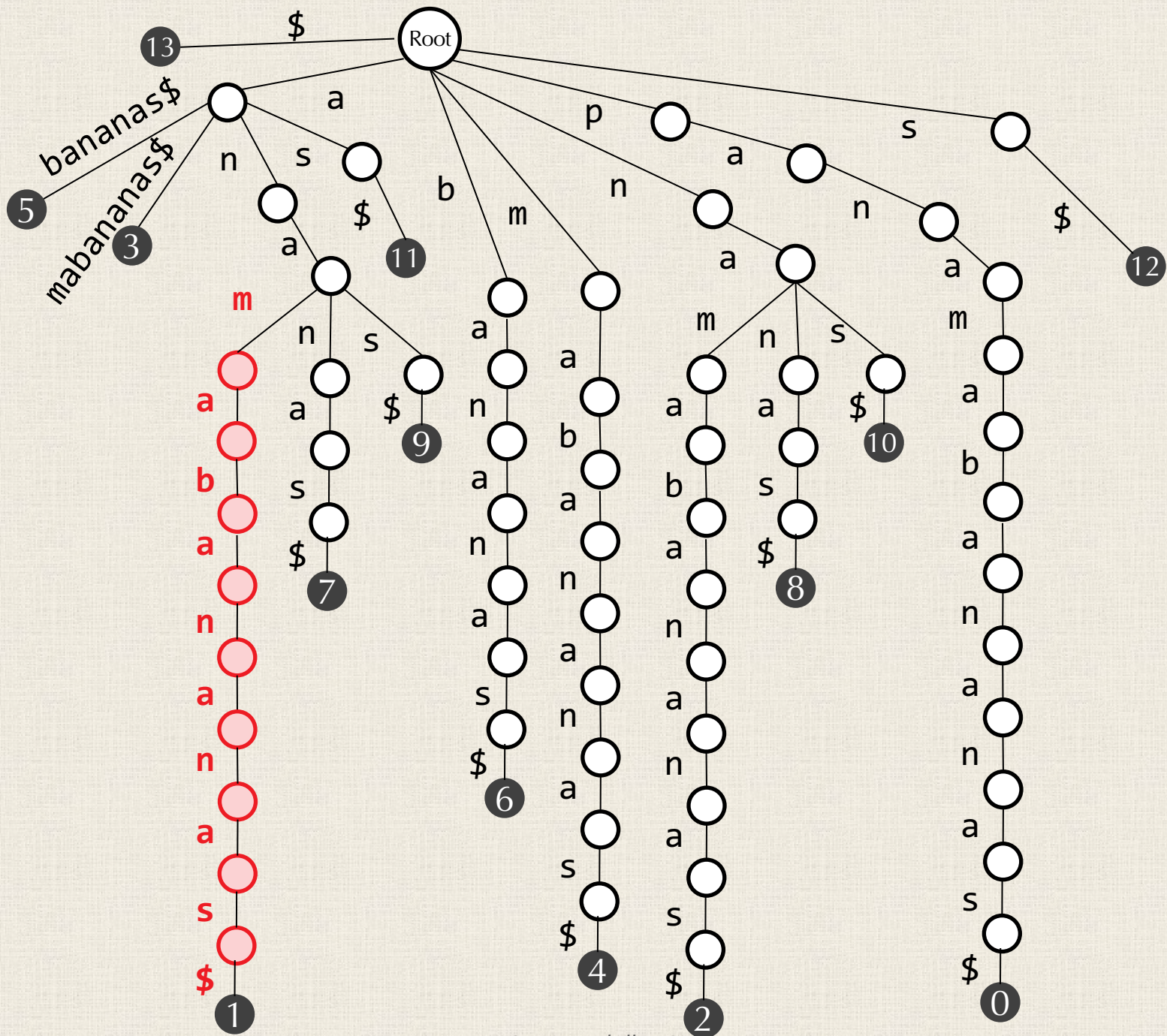
Suffix Tree: Formed by combining the edges on any *maximal non-branching path* of the suffix trie into a single edge.

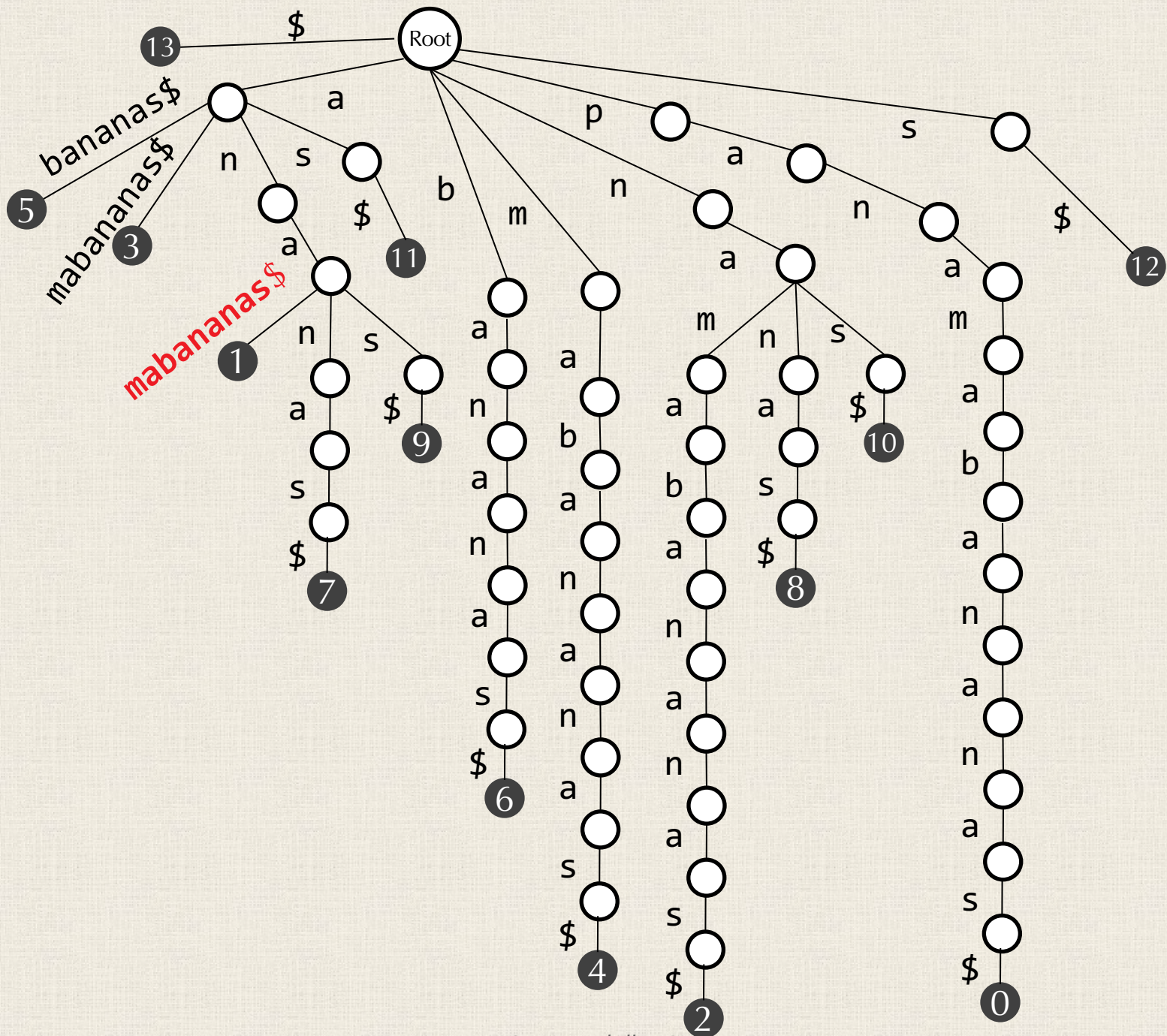


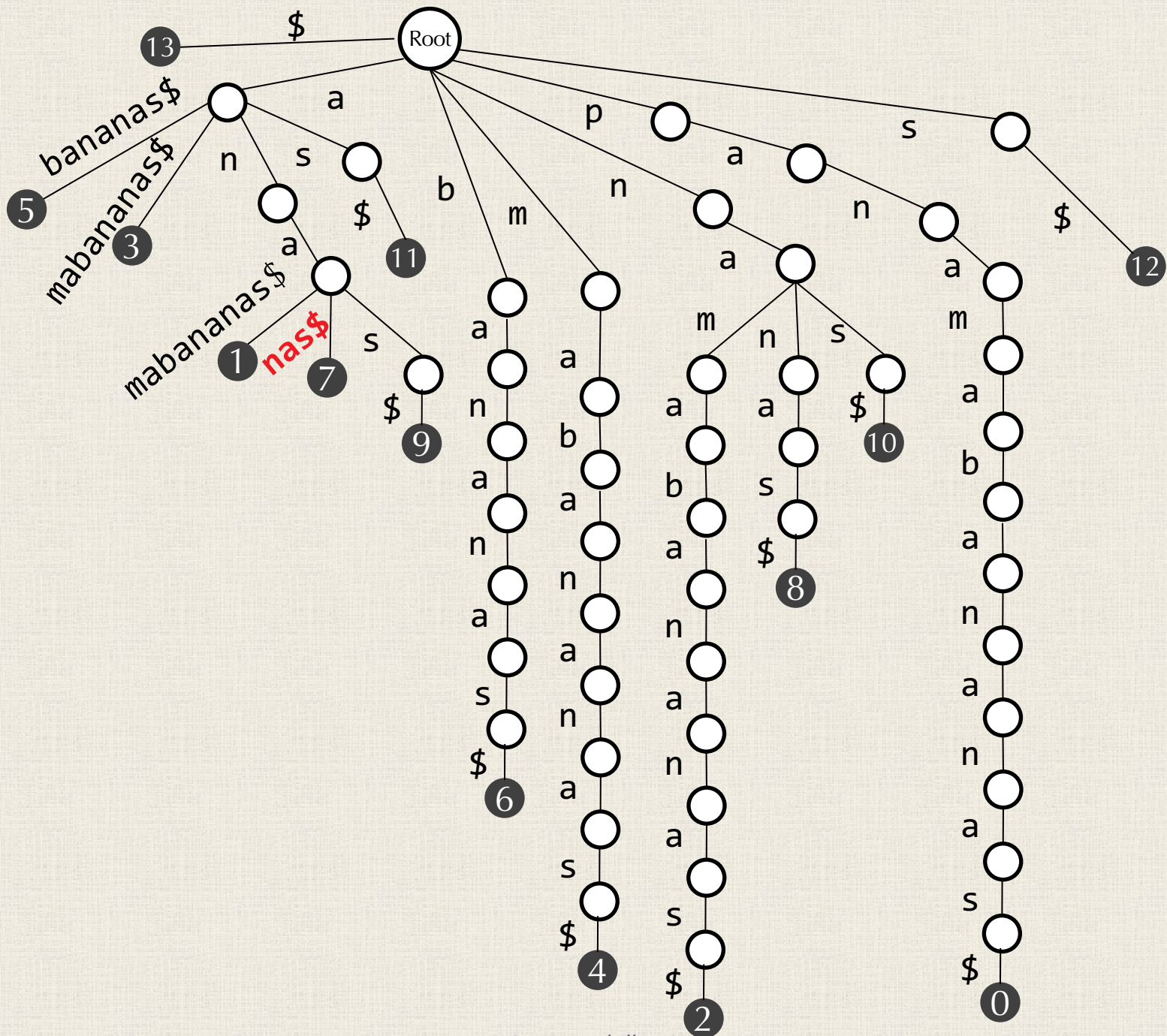


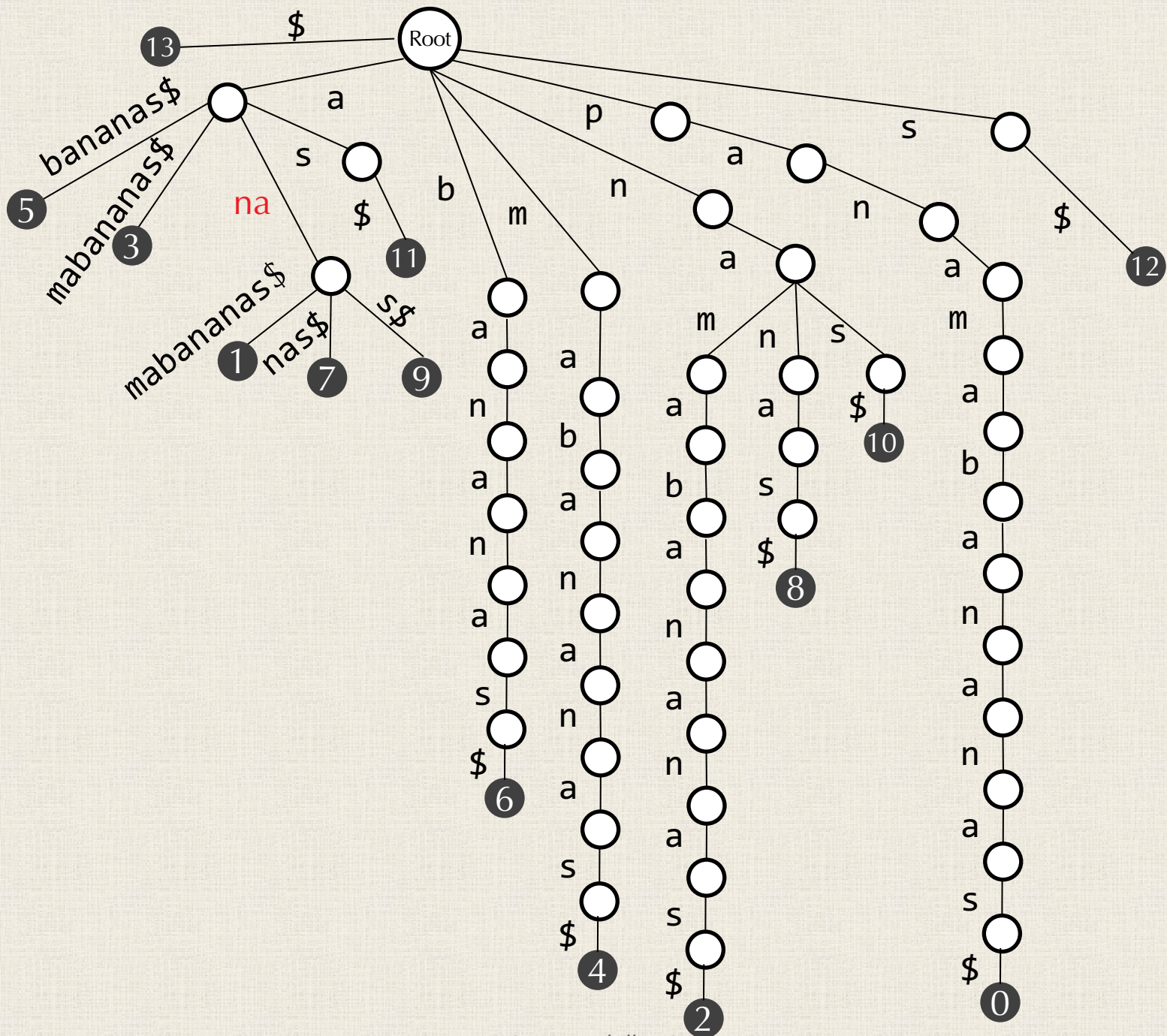


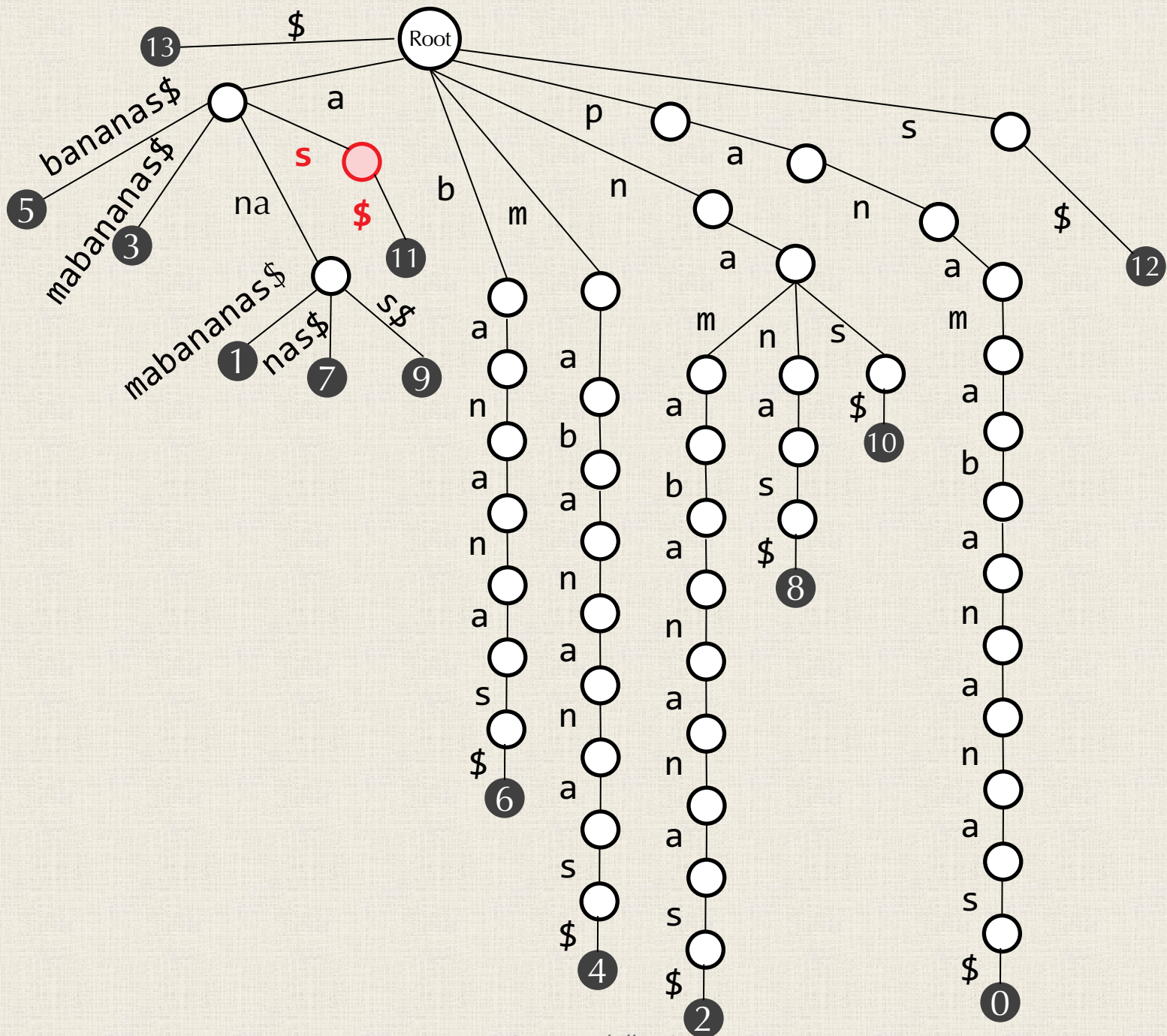


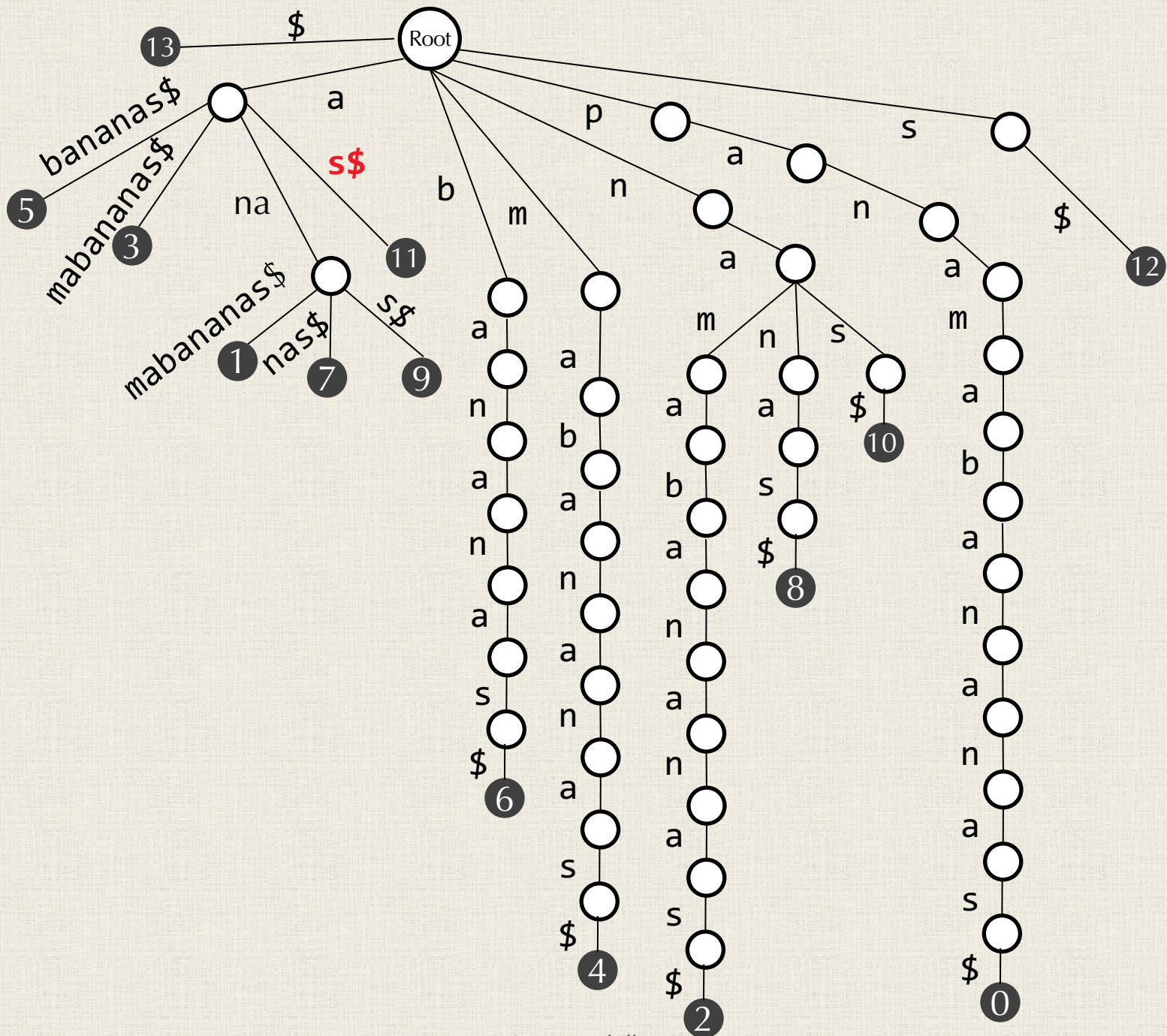


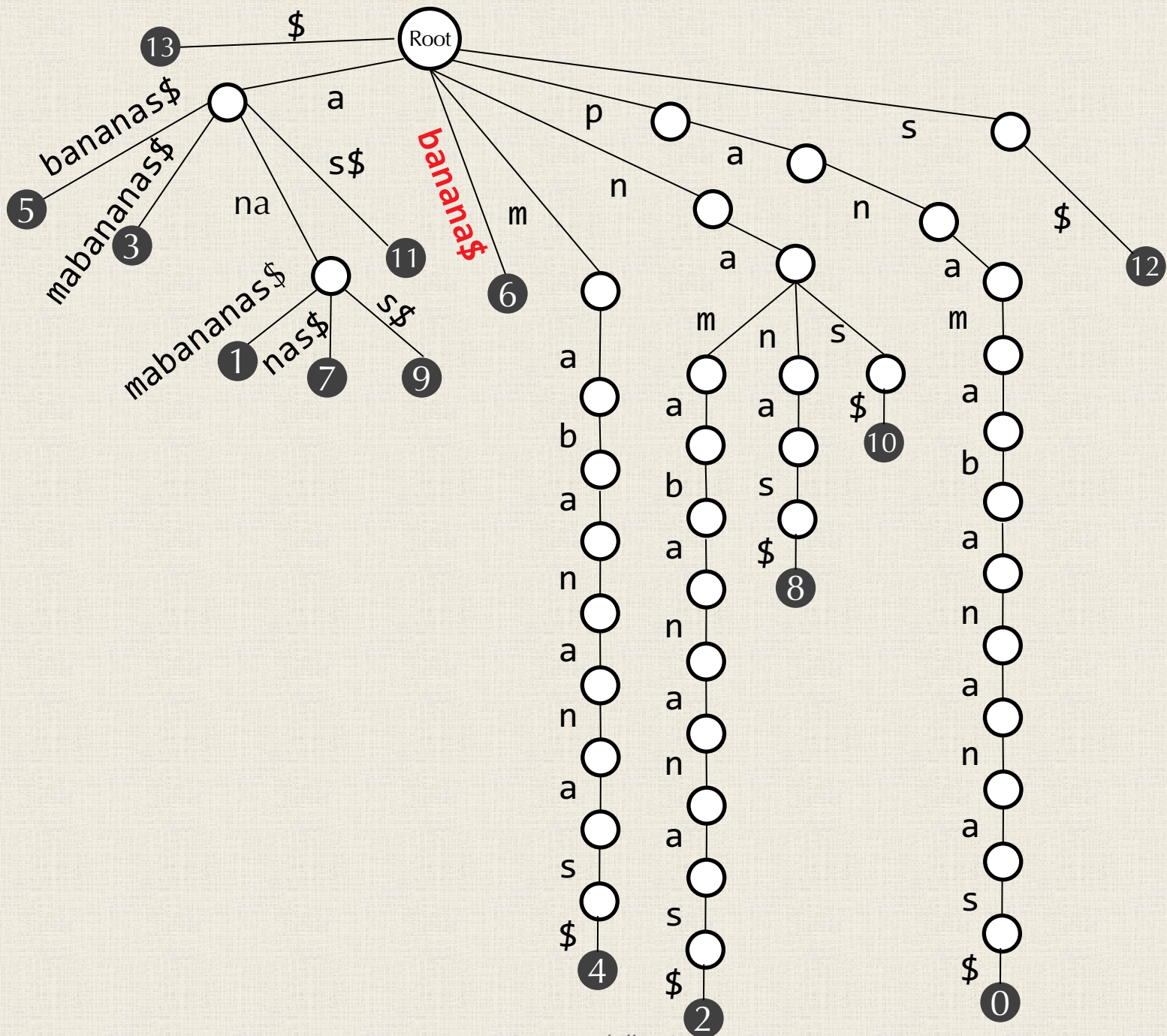


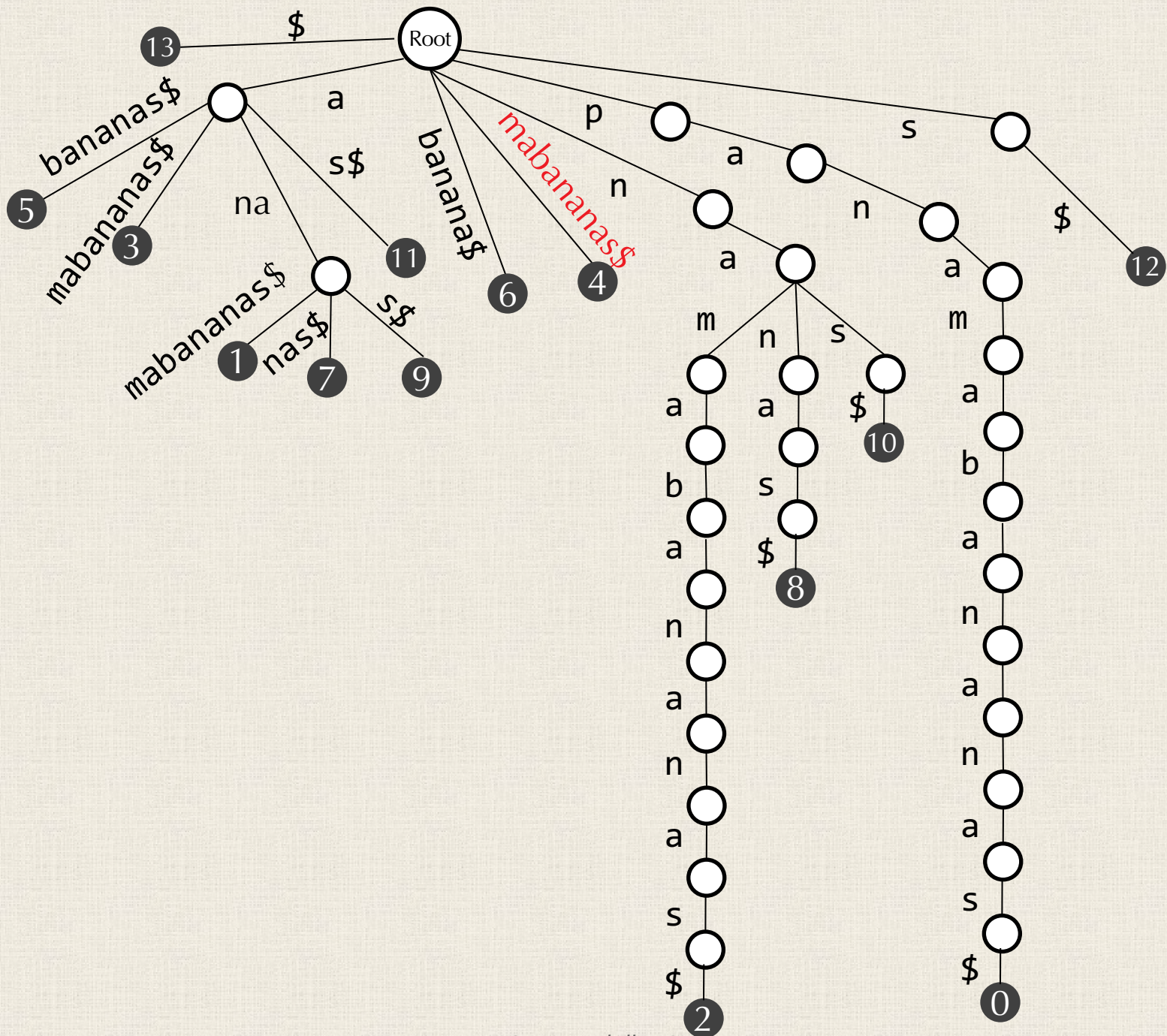


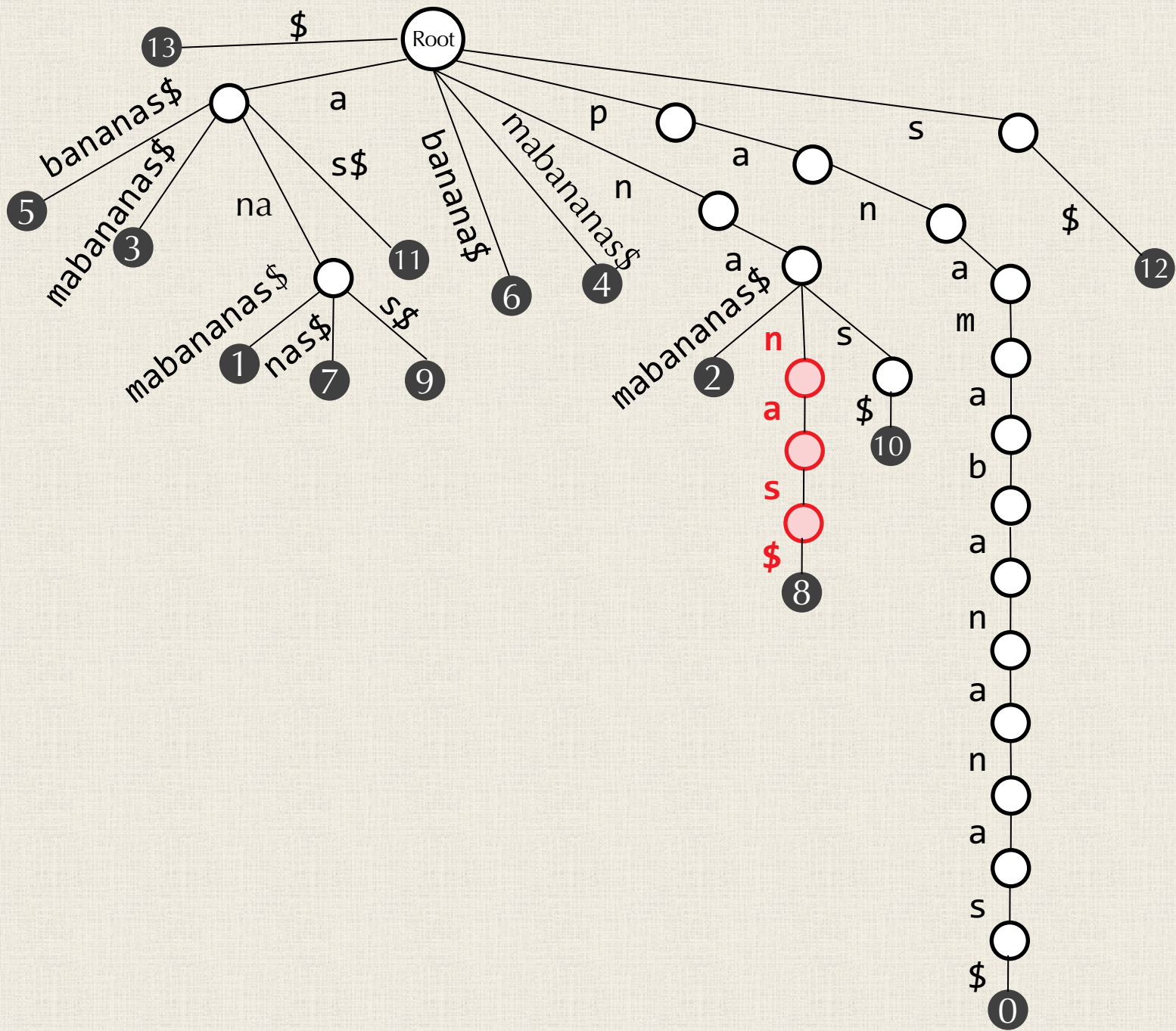


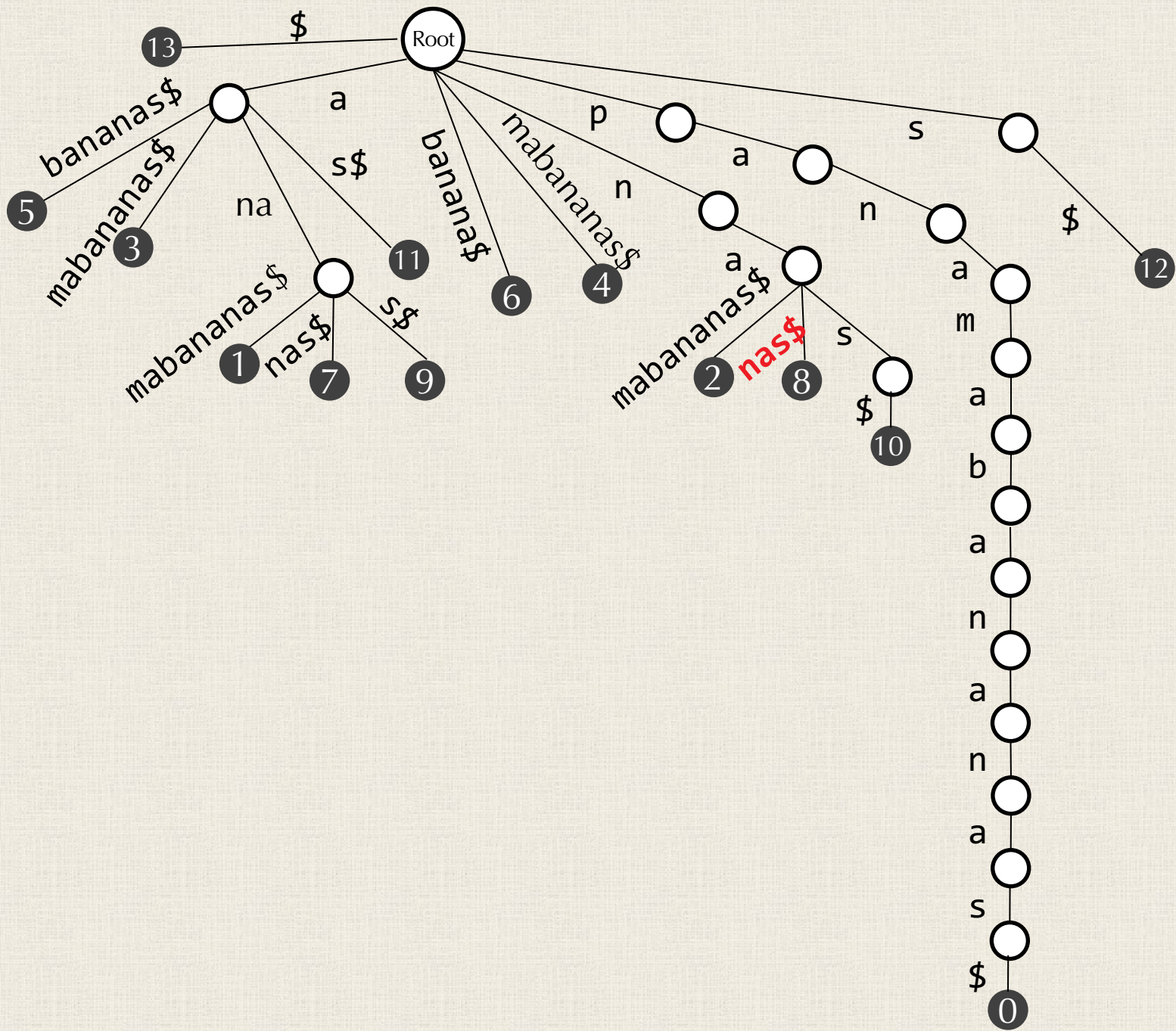


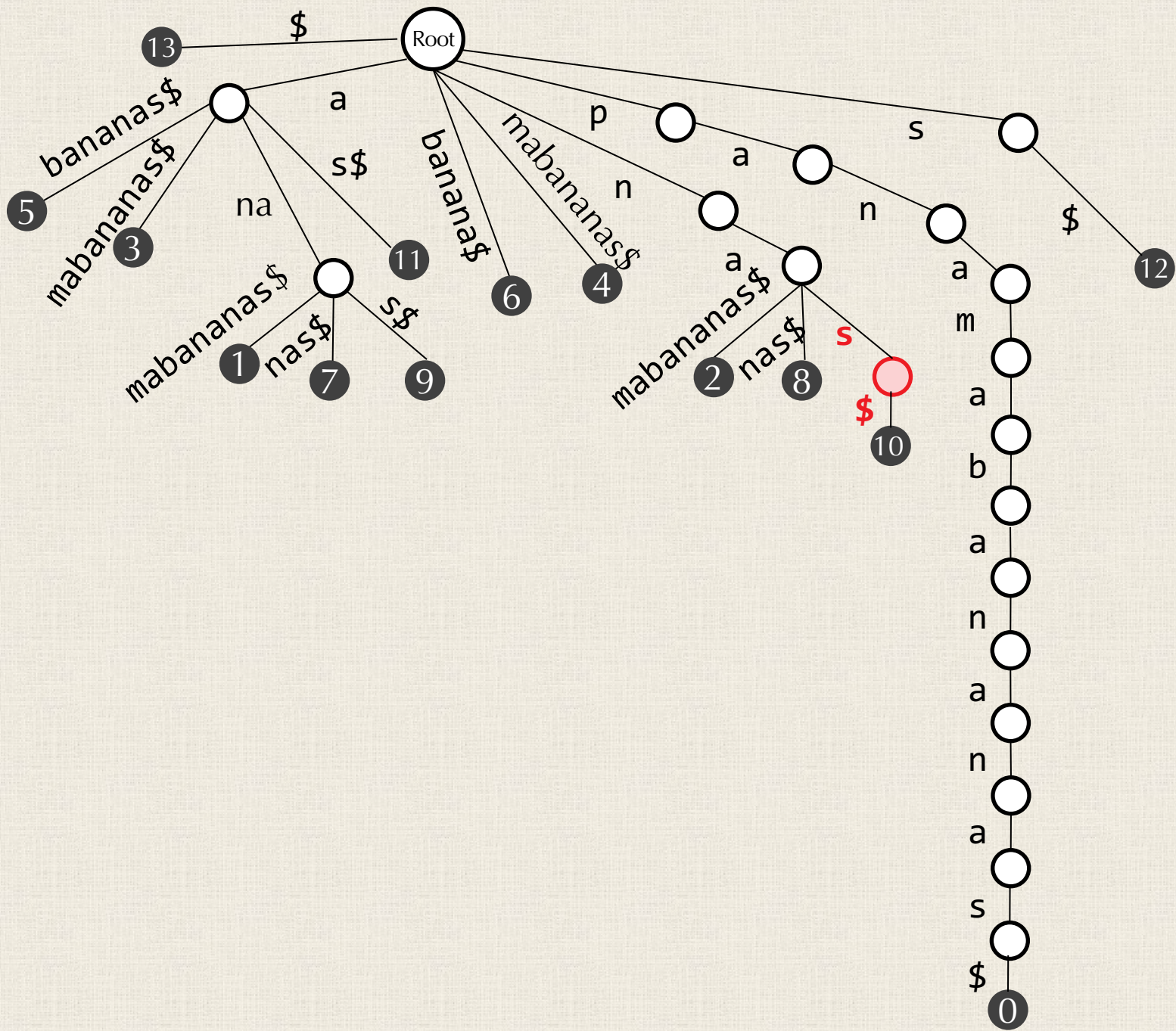


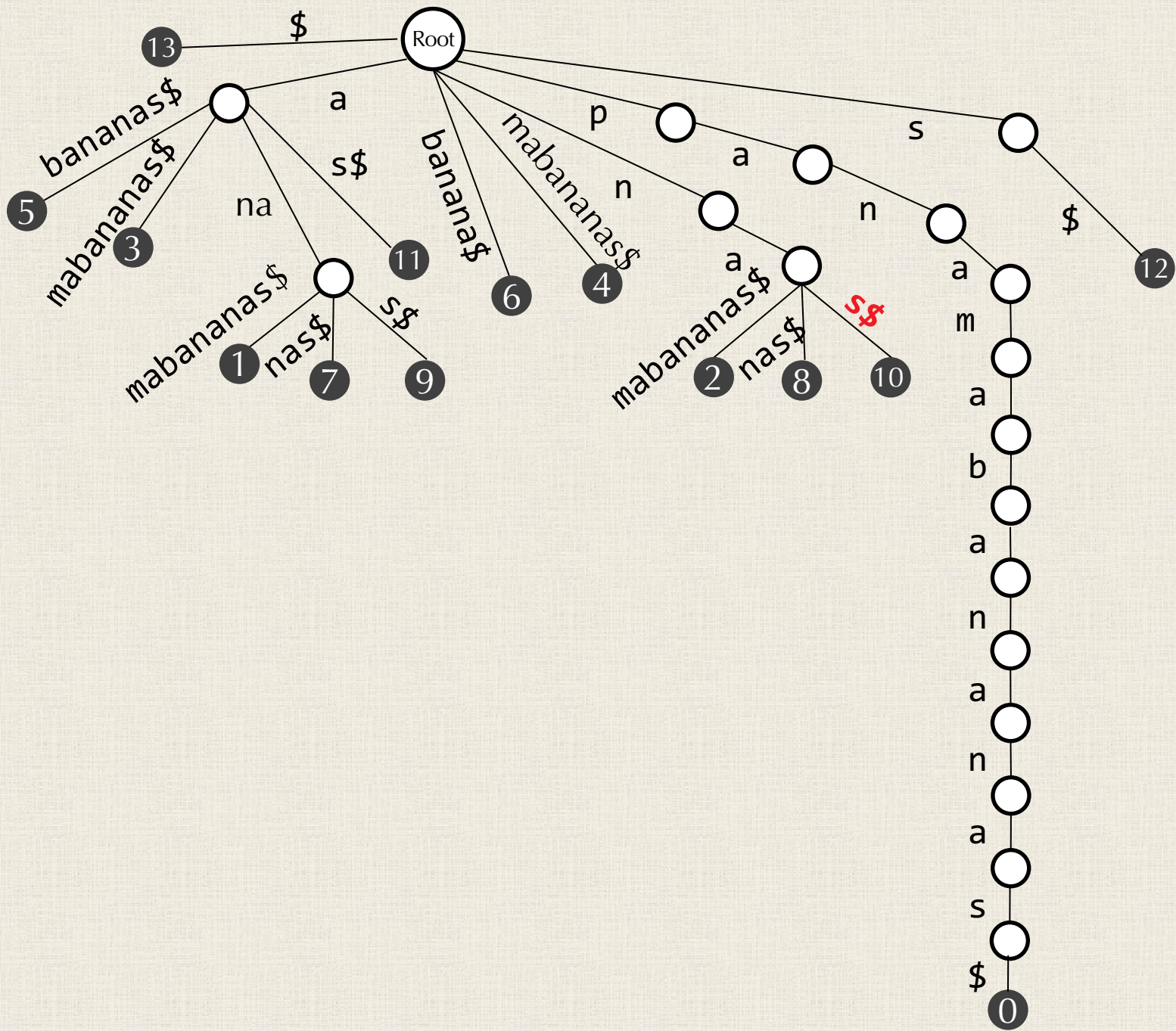


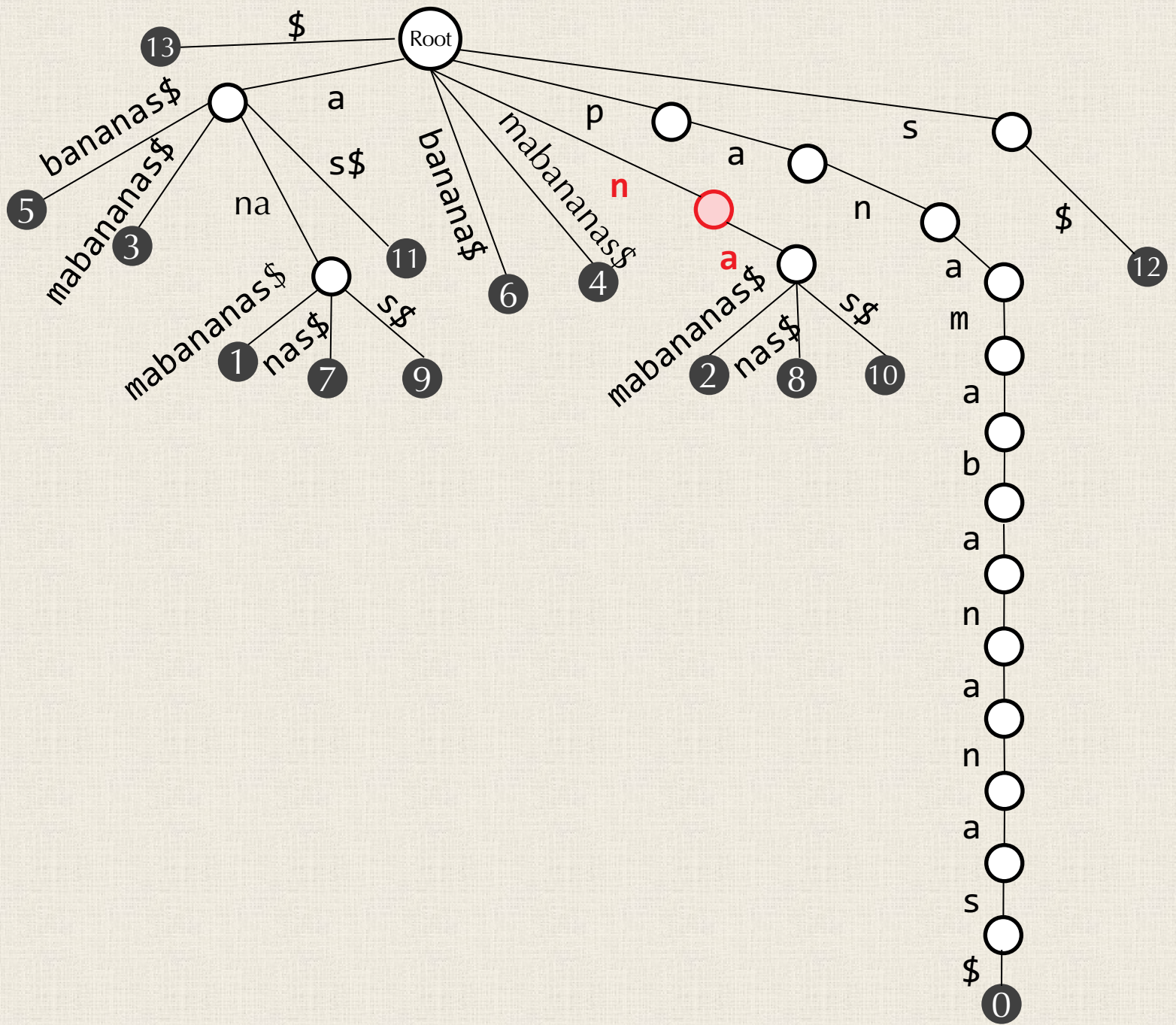


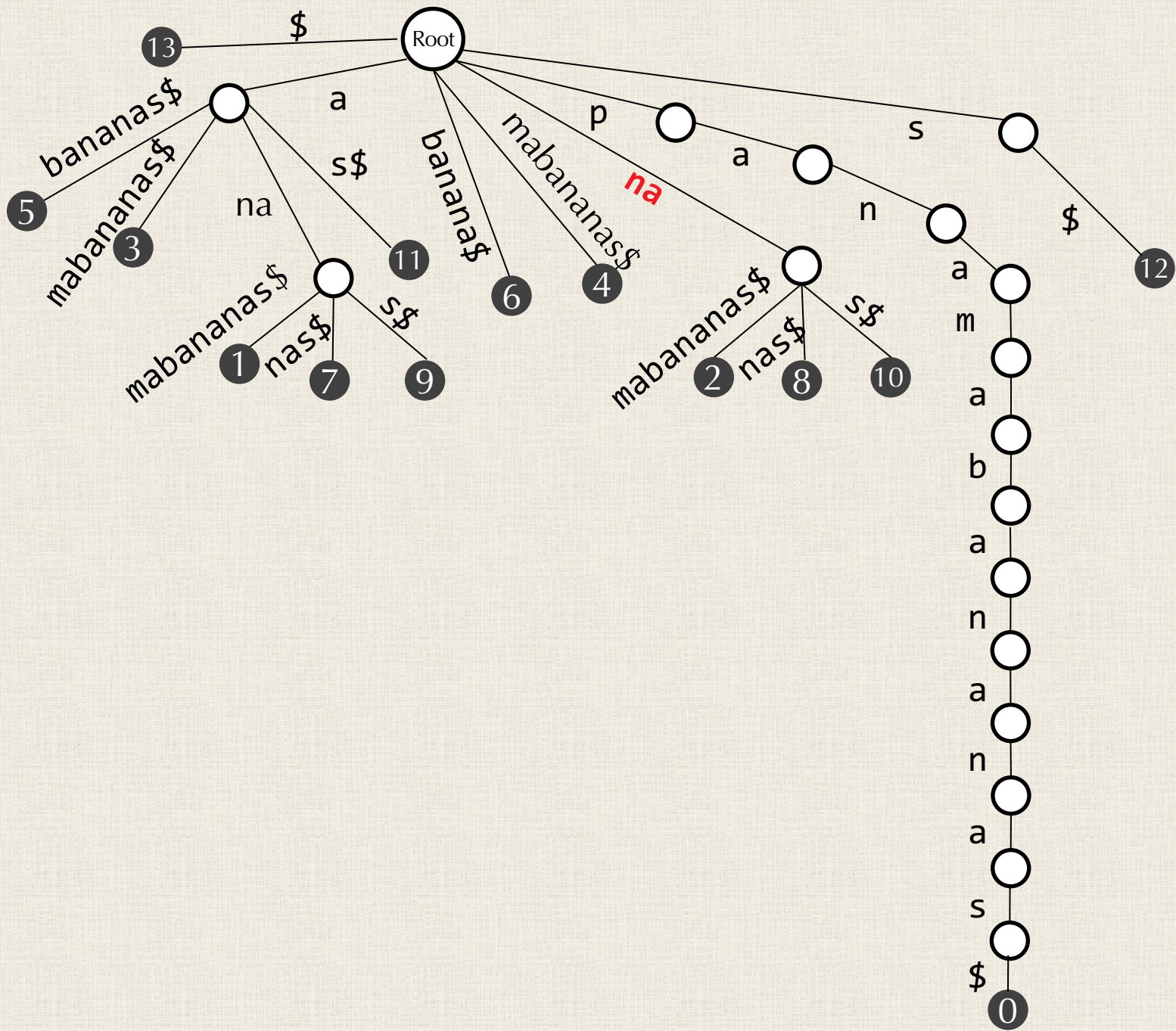


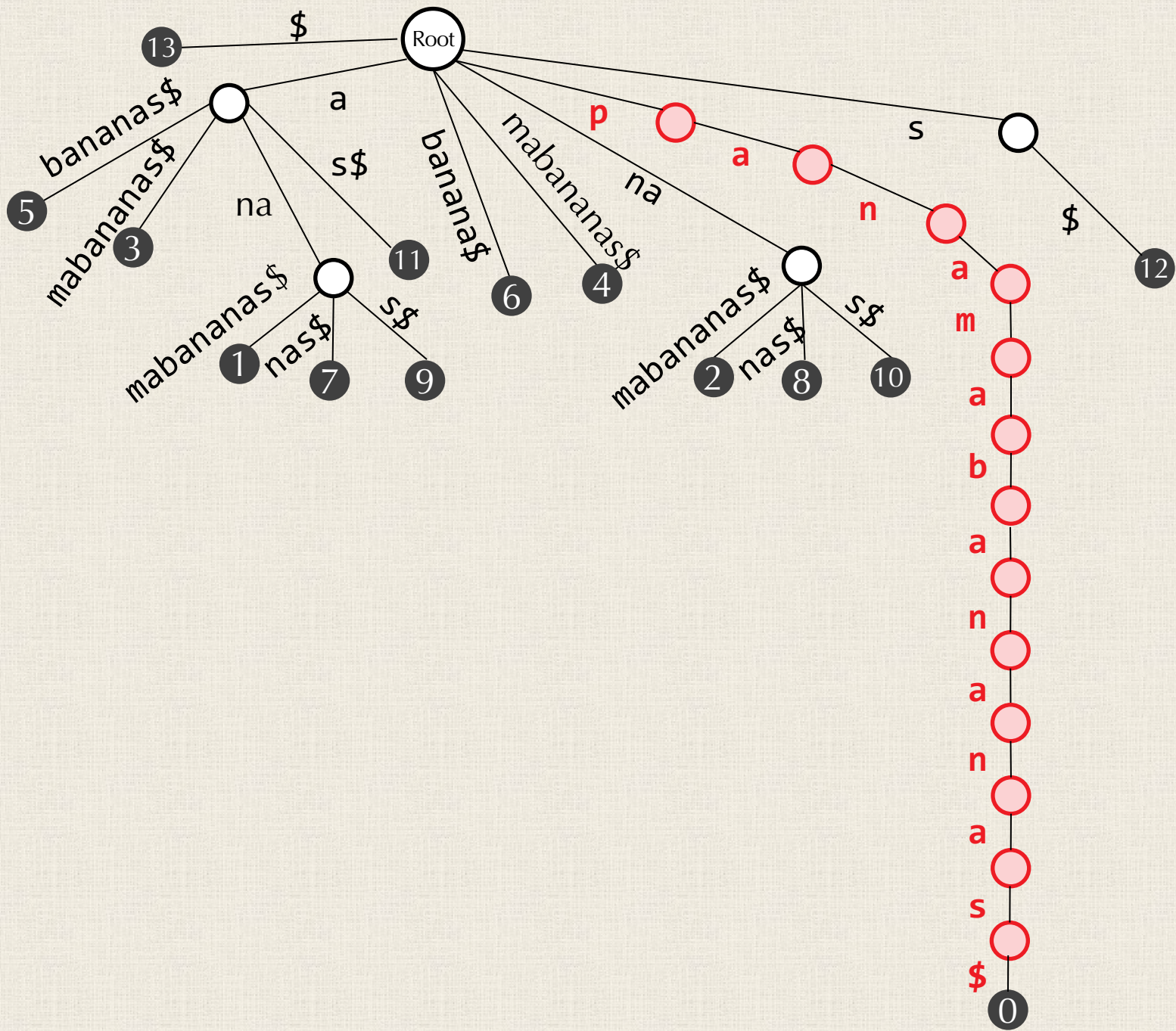


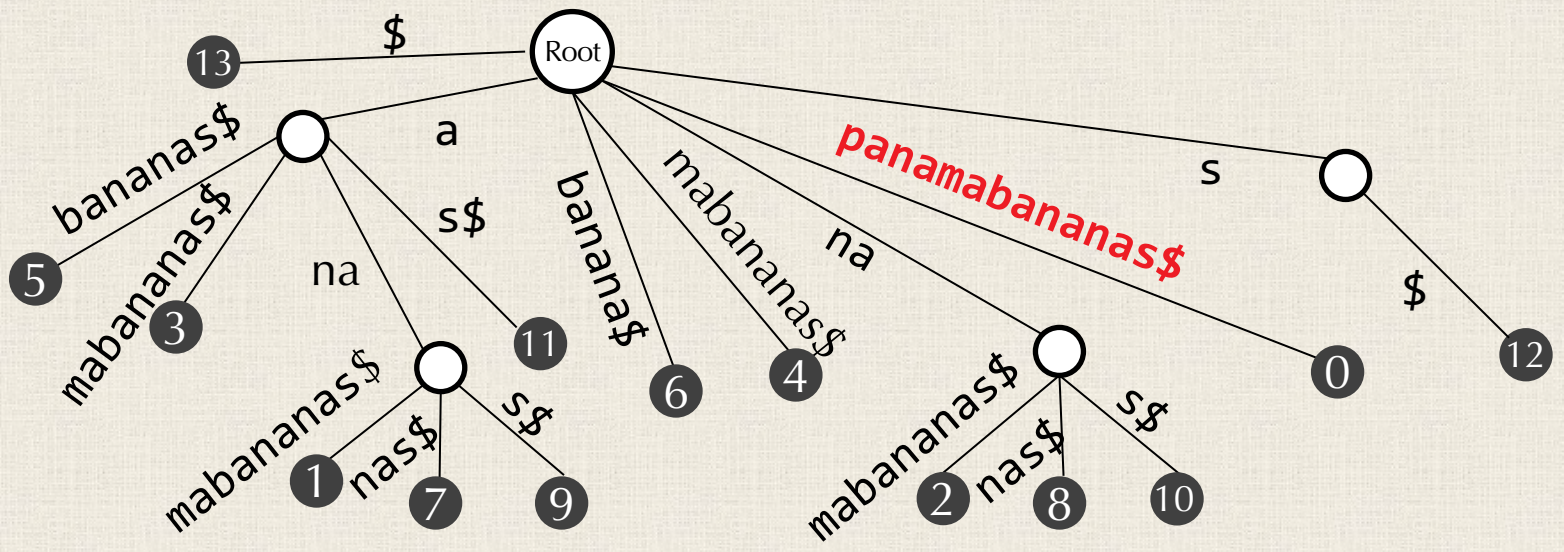


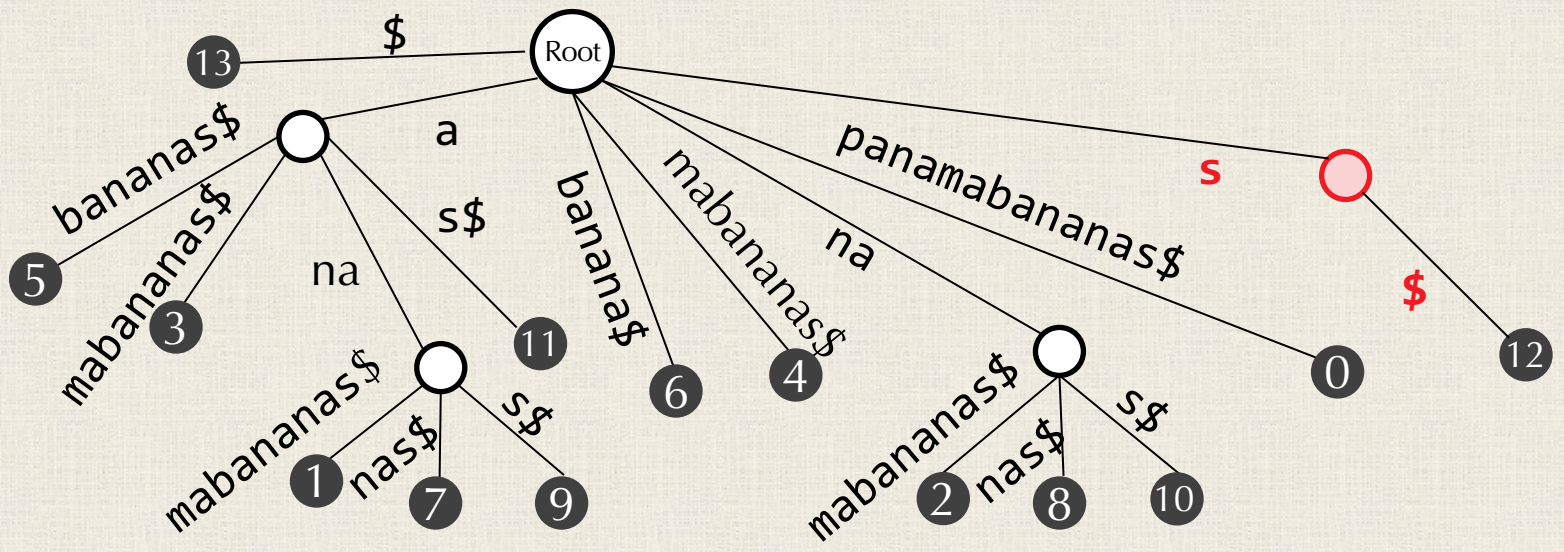


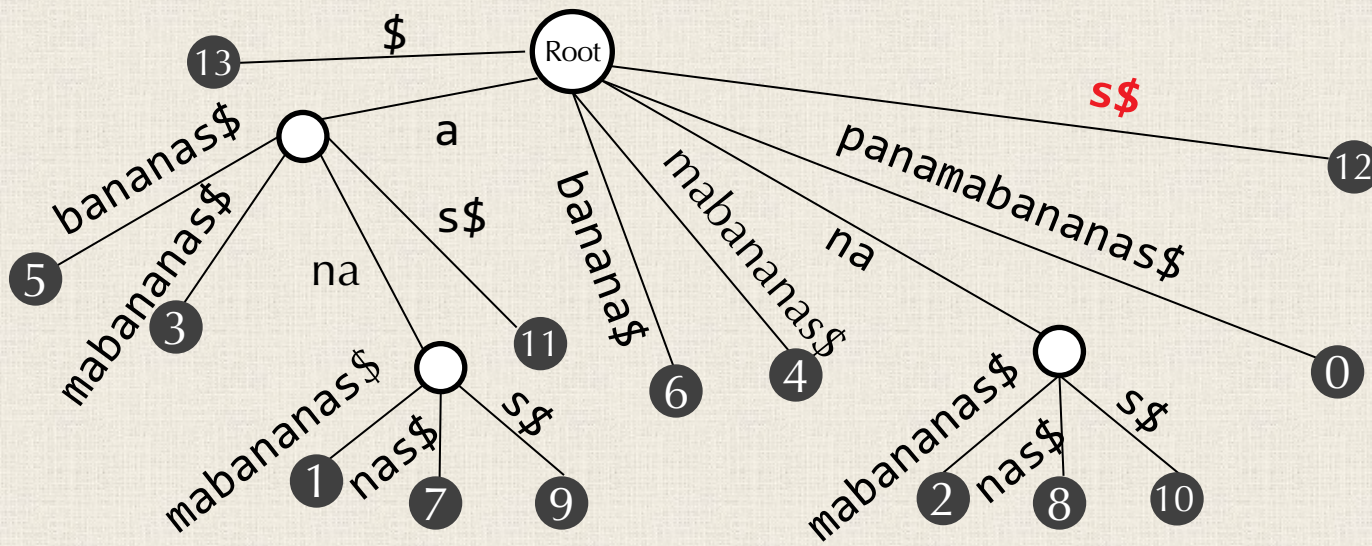


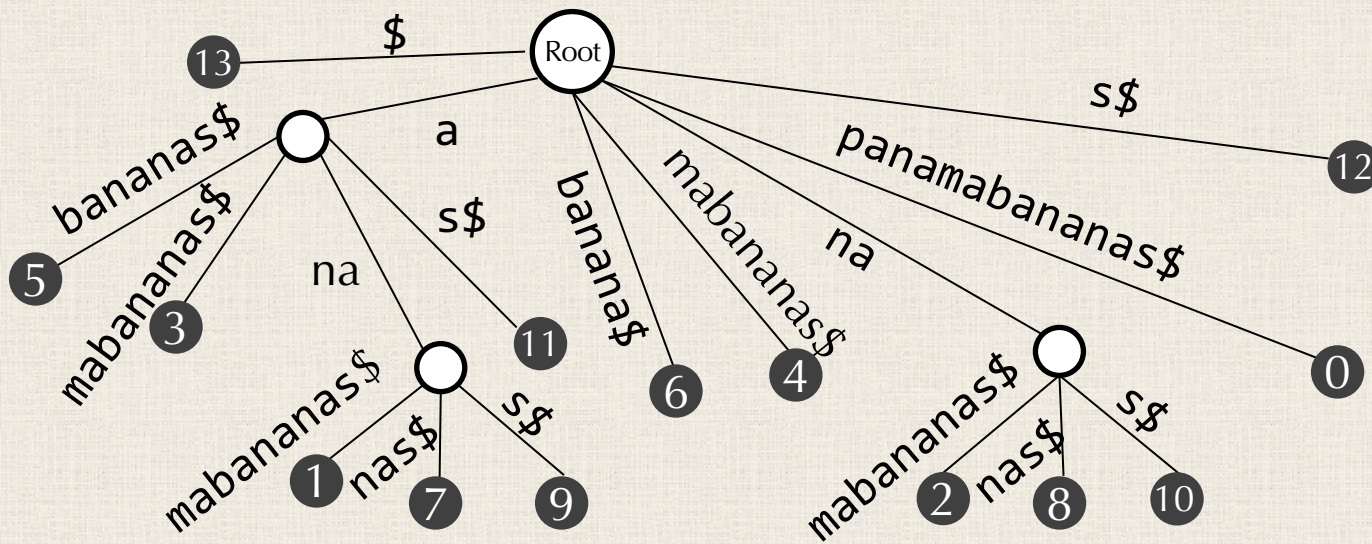










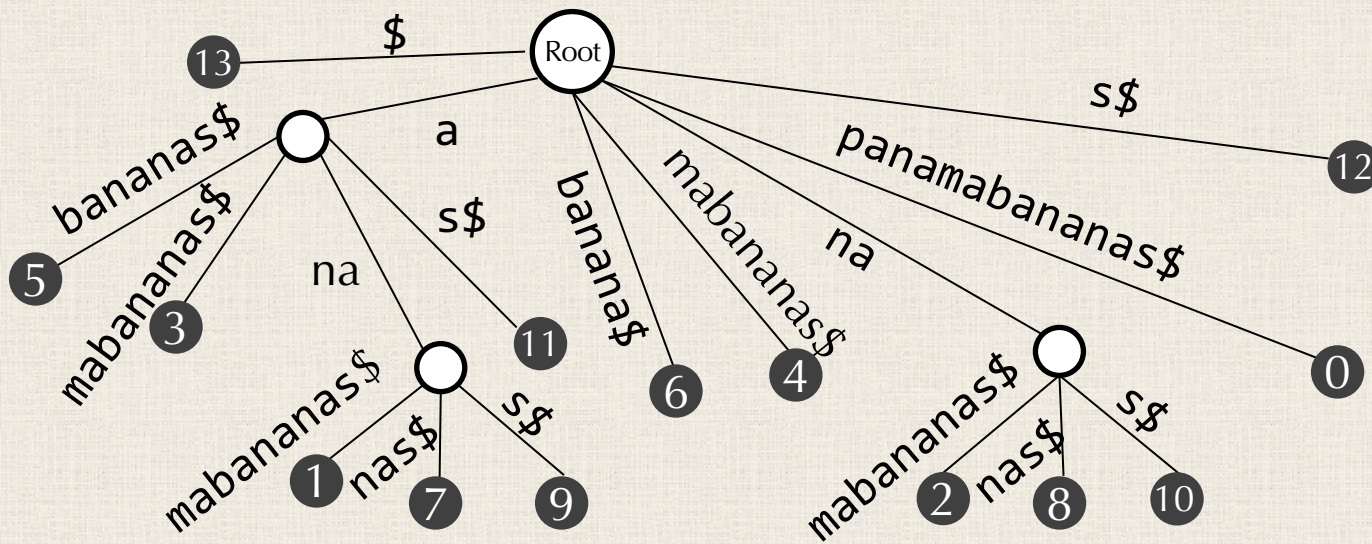


Linear pattern matching algorithms - IEEE Conference ...

by P Weiner · 1973 · Cited by 2270 · Related articles

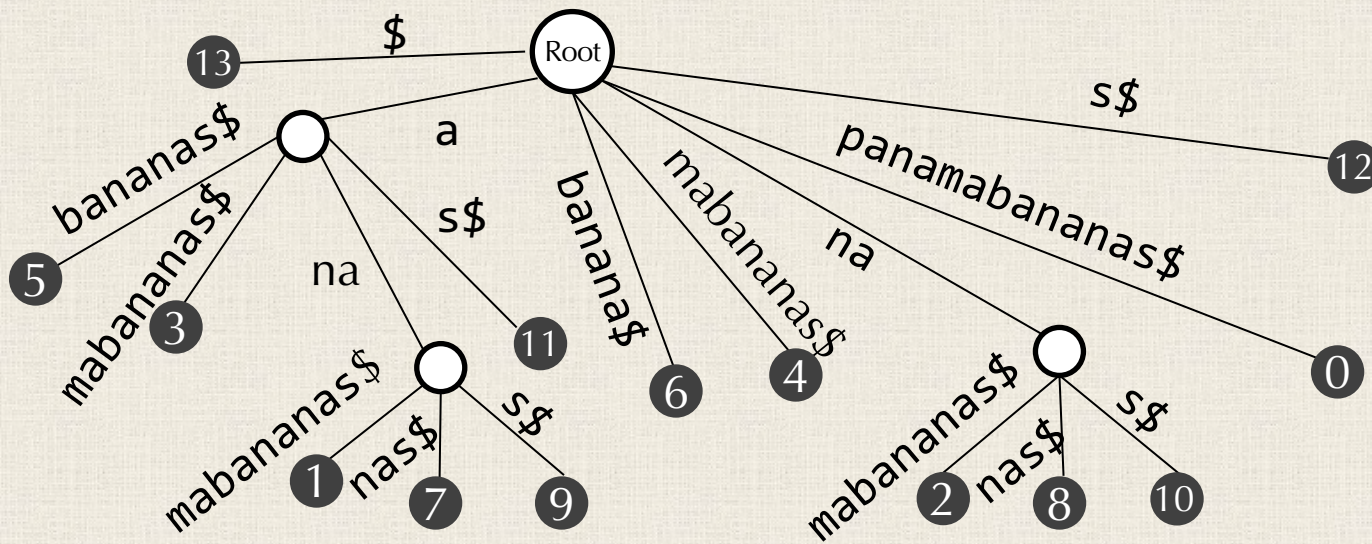
A **linear time algorithm** for obtaining a compacted version of a bi-tree associated with a given **string** is presented. ... With this construction as the basic tool, we indicate how to solve several **pattern matching** problems, including some from [4] in **linear** time.





Theorem: The suffix tree of *Text* has:

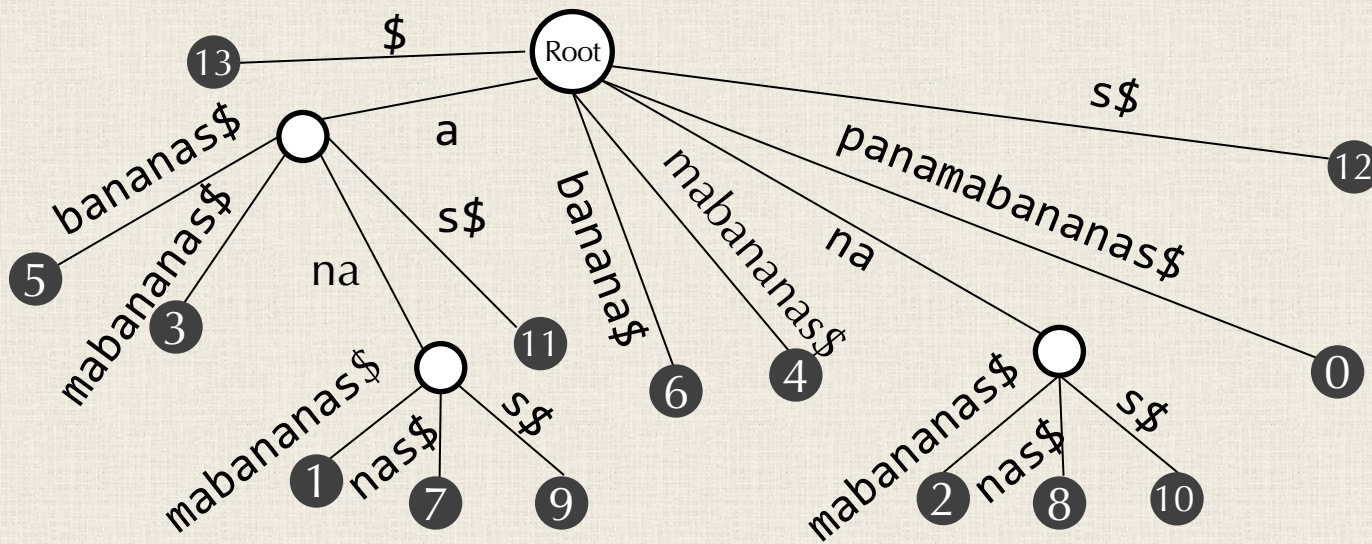
1. exactly $|Text| + 1$ leaves (degree 1)
2. at most $|Text|$ internal nodes (larger degree).



Theorem: The suffix tree of *Text* has:

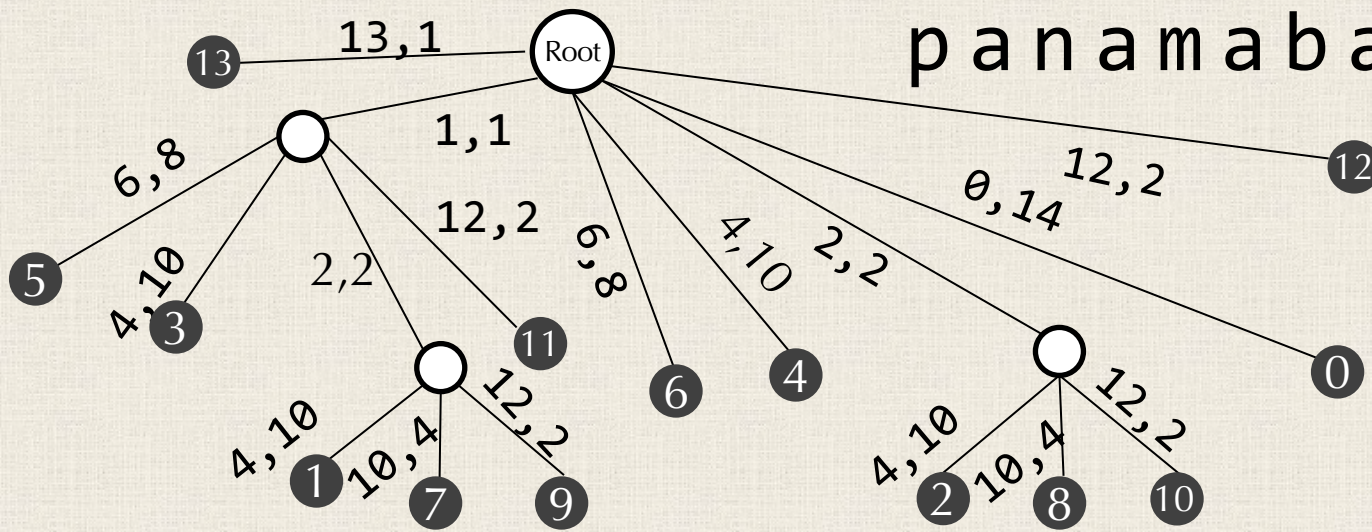
1. exactly $|Text| + 1$ leaves (degree 1)
2. at most $|Text|$ internal nodes (larger degree).

Thus, the suffix tree only takes $O(|Text|)$ space. (It is also possible to construct the suffix tree in just $O(|Text|)$ time; that's a story for another day.)



STOP: But wait ... we are still storing all the suffixes! So why are suffix trees more memory efficient than suffix tries?

panamabananas\$



STOP: But wait ... we are still storing all the suffixes! So why are suffix trees more memory efficient than suffix tries?

Answer: We can replace each edge with two integers: the starting position of the string, and its length (constant total memory).

Runtime and Memory Analysis

Runtime:

- $O(|Text|)$ to construct the suffix tree *directly*.
- $O(|Patterns|)$ to find pattern matches.
- Total: $O(|Text| + |Patterns|)$

Memory:

- $O(|Text|)$ to construct the suffix tree *directly*.
- $O(|Text|)$ to store the suffix tree.
- Total: $O(|Text|)$

From Theory to Practice

We are trained to see an algorithm with $O(|Text|)$ runtime and memory and celebrate.

From Theory to Practice

We are trained to see an algorithm with $O(|Text|)$ runtime and memory and celebrate.

But a data structure with $2 * |Text|$ memory is very different from one with $1000 * |Text|$ memory.

From Theory to Practice

We are trained to see an algorithm with $O(|Text|)$ runtime and memory and celebrate.

But a data structure with $2 * |Text|$ memory is very different from one with $1000 * |Text|$ memory.

It takes ~60 GB to store a suffix tree for a 3 GB human genome (4 bytes for each edge of tree). Can we do better?

SUFFIX ARRAYS

Fun Time! *The Clock Game*



<https://www.youtube.com/watch?v=oc9H8bo8yg0>

Suffix Array of “panamabananas\$”

20 years go by...

Suffix array: Starting positions of sorted suffixes of *Text*.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

20 years go by...

Suffix array: Starting positions of sorted suffixes of *Text*.

Can be constructed in linear time, and takes only about 12 GB for human genome.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

STOP: How can we use the suffix array for pattern matching?

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

STOP: How can we use the suffix array for pattern matching?

Answer: Like the Clock Game, we can pick a value, obtain “higher” or “lower”, and move the next value to the middle of the unsearched suffixes.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

Answer: Like the Clock Game, we can pick a value, obtain “higher” or “lower”, and move the next value to the middle of the unsearched suffixes.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
anas\$	9
ananas\$	7
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
anas\$	9
ananas\$	7
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
ab ananas\$	5
a m a bananas\$	3
a n a m a bananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Next, we find the *last* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Next, we find the *last* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Next, we find the *last* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

First, we find the *first* occurrence of “ana”.

Next, we find the *last* occurrence of “ana”.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

STOP: How do we know the starting locations of these three occurrences?

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

Answer: The suffix array tells us!

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

STOP: What is the runtime of the binary search approach using the suffix array?

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Suffix Array of “panamabananas\$”

This approach, shown for *Pattern* = “ana”, is called **binary search**.

Answer: Each pattern takes $O(\log(|Text|))$.

Sorted Suffixes	Suffix Array
\$	13
abananas\$	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

The Binary Search Pseudocode that 90% of Programmers Get Wrong

```
PatternMatchingSuffixArray(Text, Pattern)
```

```
  s ← SuffixArray(Text)
```

```
  minIndex ← 0
```

```
  maxIndex ← |Text| - 1
```

```
  while minIndex ≤ maxIndex
```

```
    midIndex ← floor((minIndex + maxIndex)/2)
```

```
    if Pattern > suffix of Text starting at s(midIndex)
```

```
      minIndex ← midIndex + 1
```

```
    else
```

```
      maxIndex ← midIndex - 1
```

```
  if Pattern matches suffix of Text starting at s(midIndex)
```

```
    first ← minIndex
```

```
  else
```

```
    return no occurrences
```

```
  minIndex ← first
```

```
  maxIndex ← |Text| - 1
```

```
  while minIndex ≤ maxIndex
```

```
    midIndex ← floor((minIndex + maxIndex)/2)
```

```
    if Pattern matches suffix of Text starting at s(midIndex)
```

```
      minIndex ← midIndex + 1
```

```
    else
```

```
      maxIndex ← midIndex - 1
```

```
  last ← maxIndex
```

```
  return (first, last)
```

Suffix Arrays are Another Example of Computational Biologists Influencing CS

Suffix arrays: a new method for on-line string searches

U Manber, G Myers - *siam Journal on Computing*, 1993 - SIAM

A new and conceptually simple data structure, called a suffix array, for on-line string searches is introduced in this paper. Constructing and querying suffix arrays is reduced to a sort and search paradigm that employs novel algorithms. The main advantage of suffix arrays over suffix trees is that, in practice, they use three to five times less space. From a complexity standpoint, suffix arrays permit on-line string searches of the type, "Is W a substring of A?" to be answered in time $O(P+\log N)$, where P is the length of W and N is the ...

☆ 🔖 Cited by 2870 Related articles All 43 versions ⇨

Note: This is a paper produced by a computational biologist (and future head of search for Google) in a non-biological context seven years before the publication of the first draft human genome.

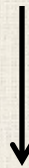
THE BURROWS-WHEELER TRANSFORM

Idea to Reduce Memory: Compress the Genome

Run-length encoding: compresses every run of n identical symbols X into the substring “ nX ”.

Text

GGGGGGGGGGCCCCCCCCAAAATTTTTTTTTTTTTTTTCCCCCG



10G11C7A15T5C1G

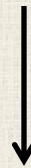
Run-length encoding

Idea to Reduce Memory: Compress the Genome

Run-length encoding: compresses every run of n identical symbols X into the substring “ nX ”.

Text

GGGGGGGGGGCCCCCCCCAAAATTTTTTTTTTTTTTTTCCCCCG



10G11C7A15T5C1G

Run-length encoding

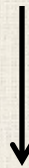
Problem: Genomes don't have lots of runs...

Idea to Reduce Memory: Compress the Genome

Run-length encoding: compresses every run of n identical symbols X into the substring “ nX ”.

Text

GGGGGGGGGGCCCCCCCCAAAATTTTTTTTTTTTTTTTCCCCCG

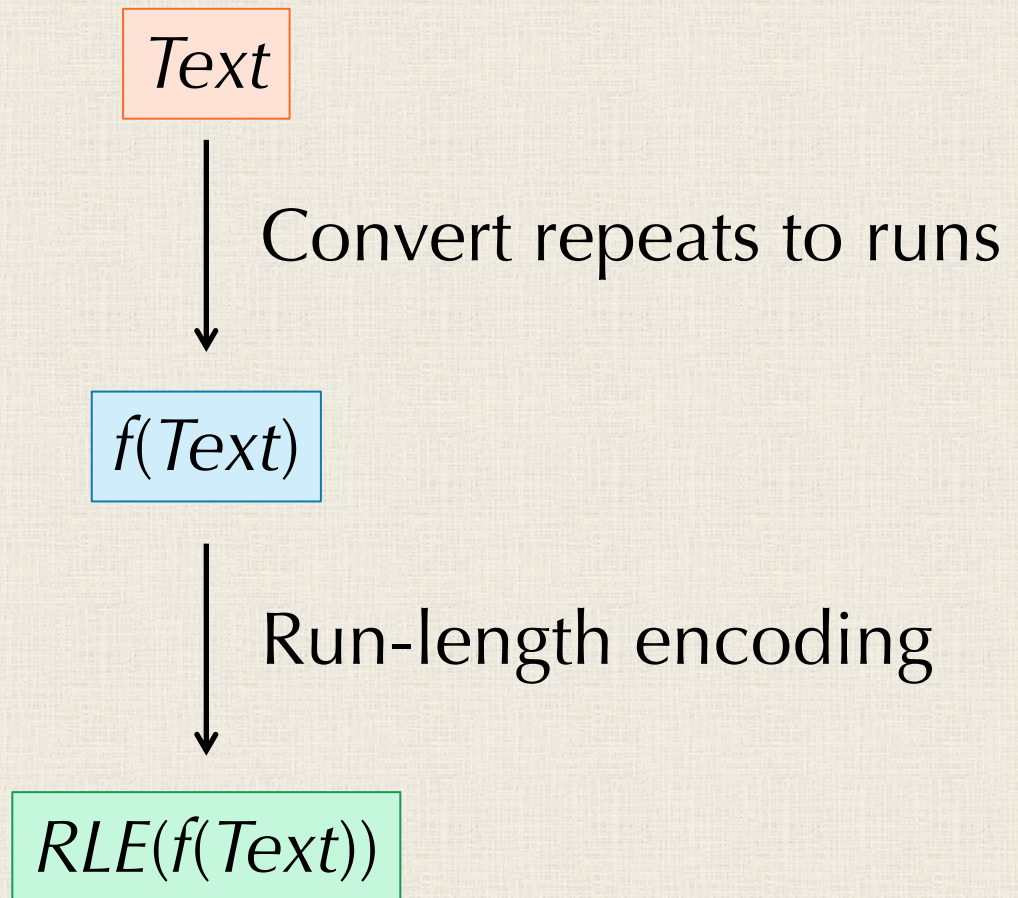


10G11C7A15T5C1G

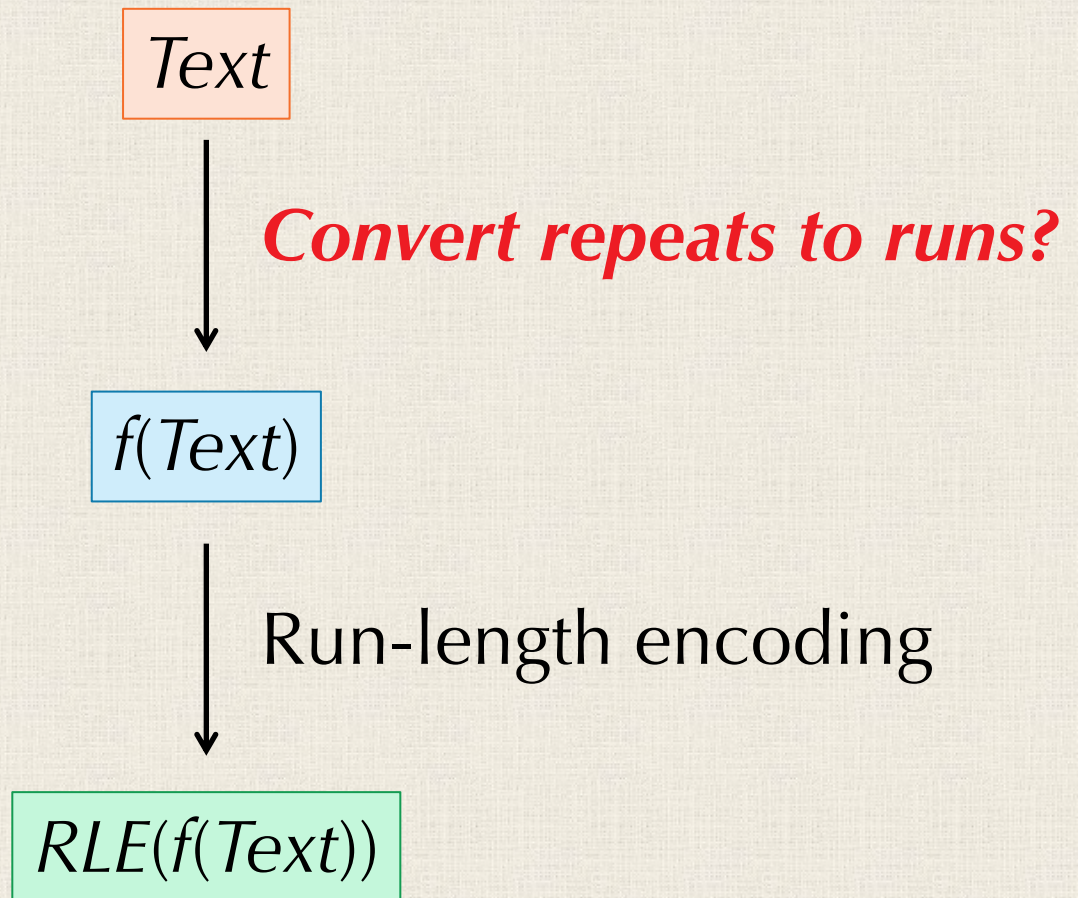
Run-length encoding

...but they do have a lot of repeats!

Converting Repeats to Runs?



Converting Repeats to Runs?



Converting Repeats to Runs?



One way we could convert repeats into runs would be to simply sort *Text* lexicographically.

Converting Repeats to Runs?



STOP: What is wrong with sorting the string?

Converting Repeats to Runs?



STOP: What is wrong with sorting the string?

Answer: There is no way to “invert” the compression to recover *Text*. That is, f should be an injection!

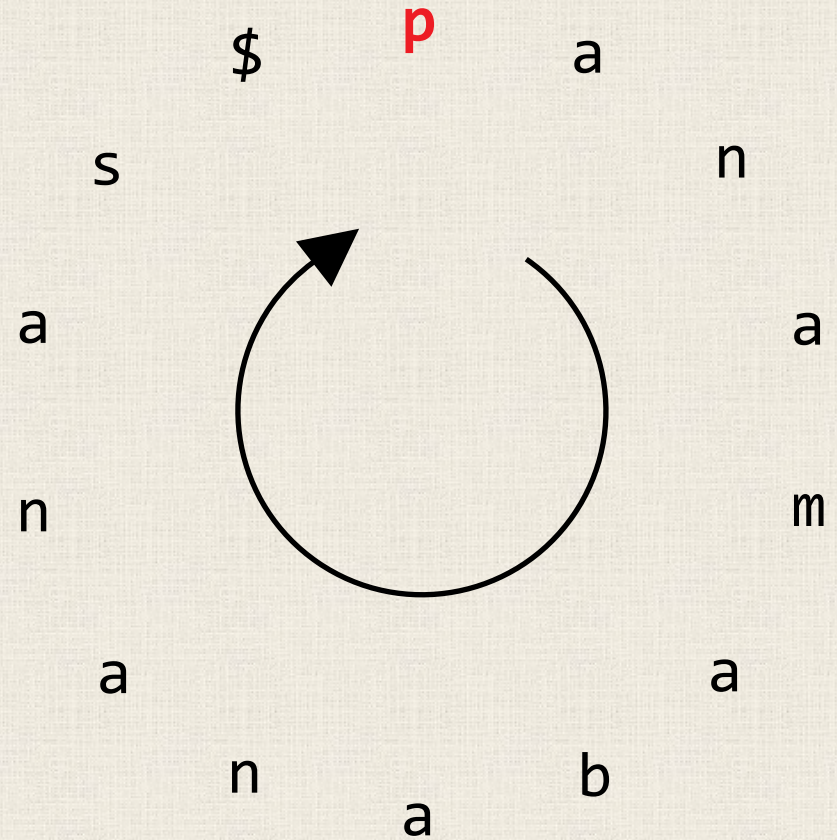
Converting Repeats to Runs?



Formally, define $S_{A, n}$ as the set of all strings of length $n+1$, with a “\$” and n symbols taken from an alphabet A . We want a function $f: S_{A, n} \rightarrow S_{A, n}$ that is **invertible**: if $f(x) = f(y)$, then $x = y$.

The Burrows-Wheeler Transform

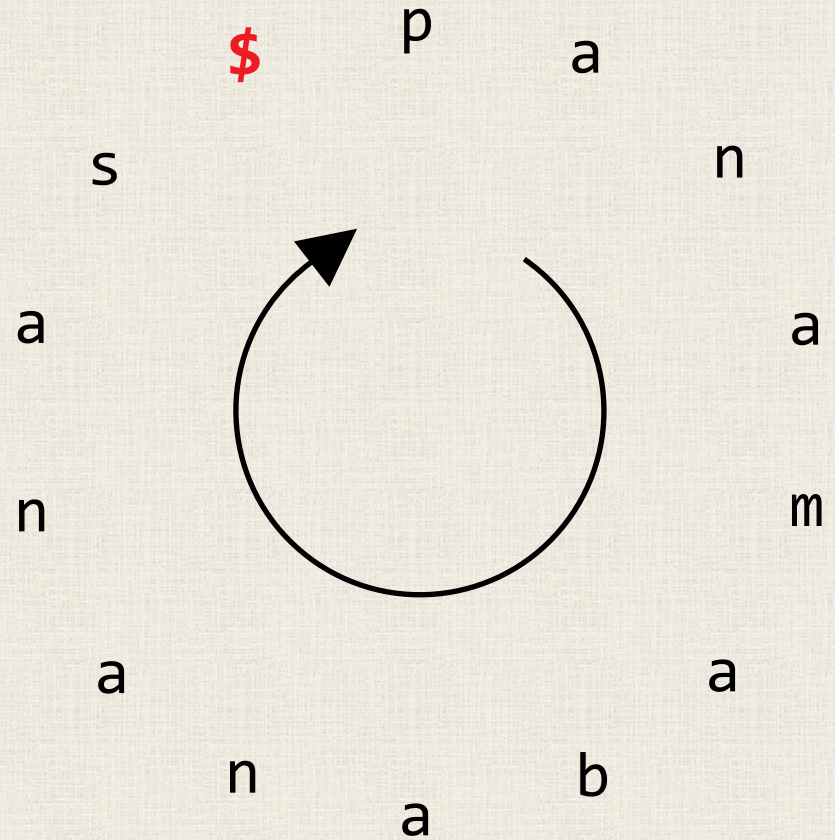
panamabananas\$



Form all cyclic rotations of
"panamabananas\$"

The Burrows-Wheeler Transform

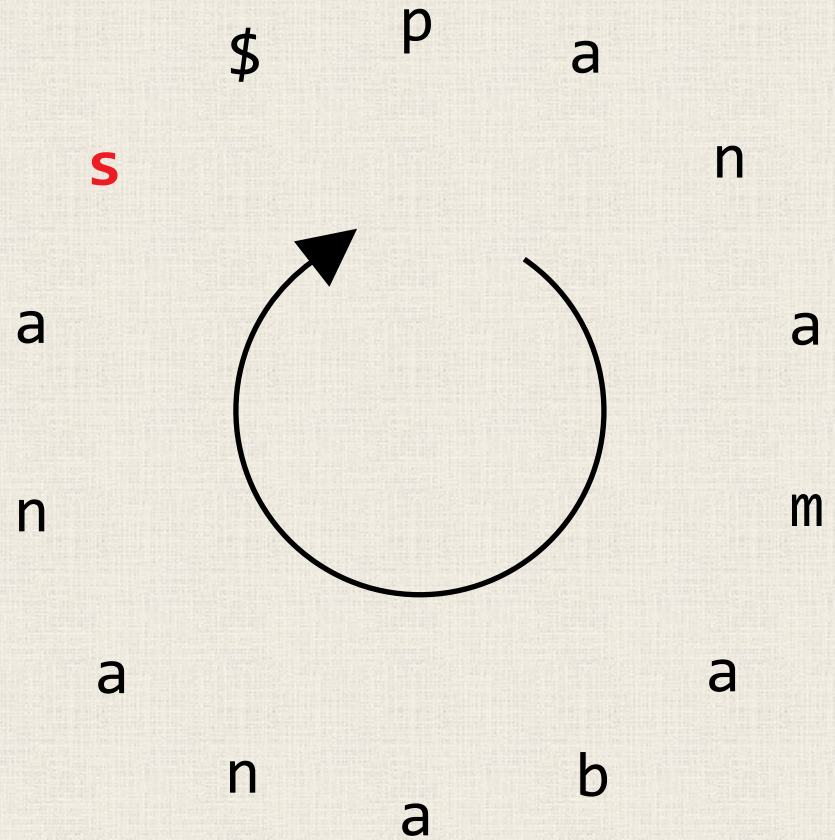
```
panamabananas$  
$panamabananas
```



Form all cyclic rotations of
"panamabananas\$"

The Burrows-Wheeler Transform

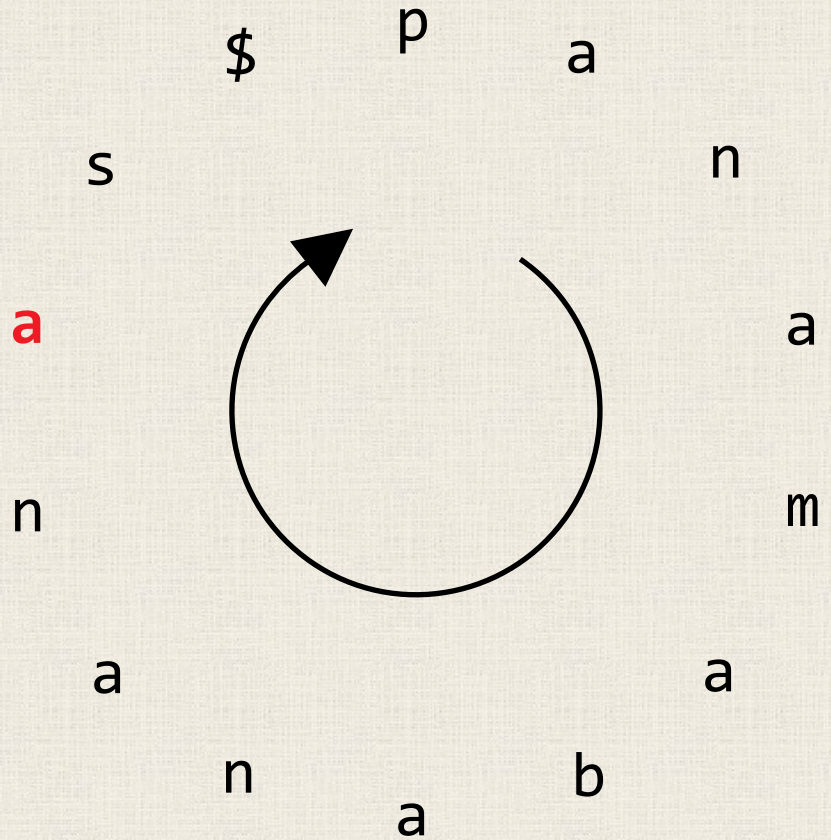
```
panamabananas$  
$panamabananas  
s$panamabanana
```



Form all cyclic rotations of
"panamabananas\$"

The Burrows-Wheeler Transform

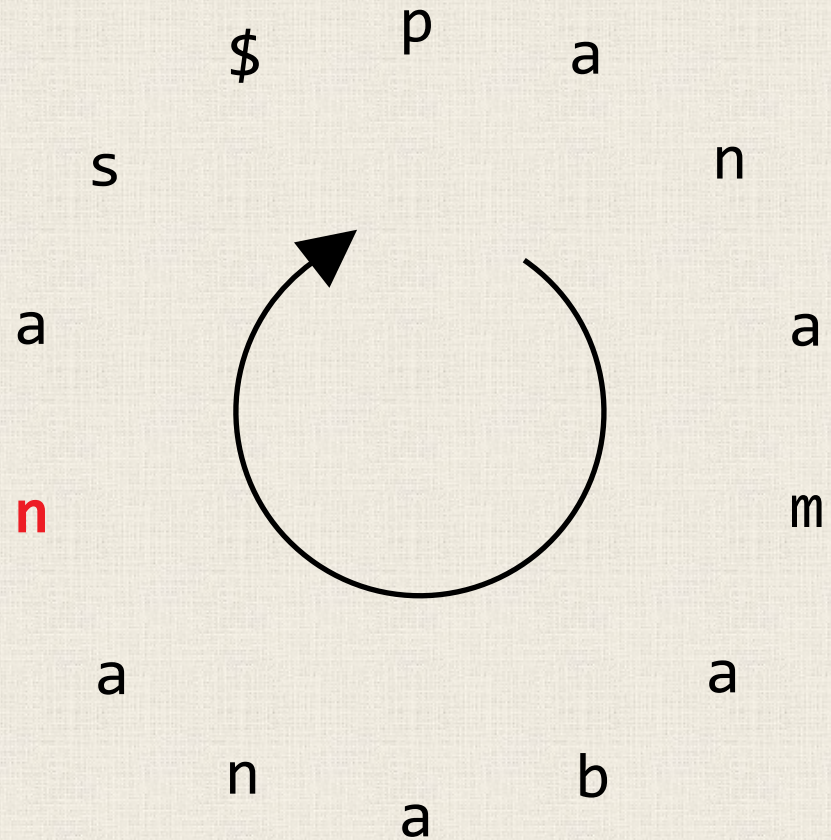
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanan
```



Form all cyclic rotations of
"panamabananas\$"

The Burrows-Wheeler Transform

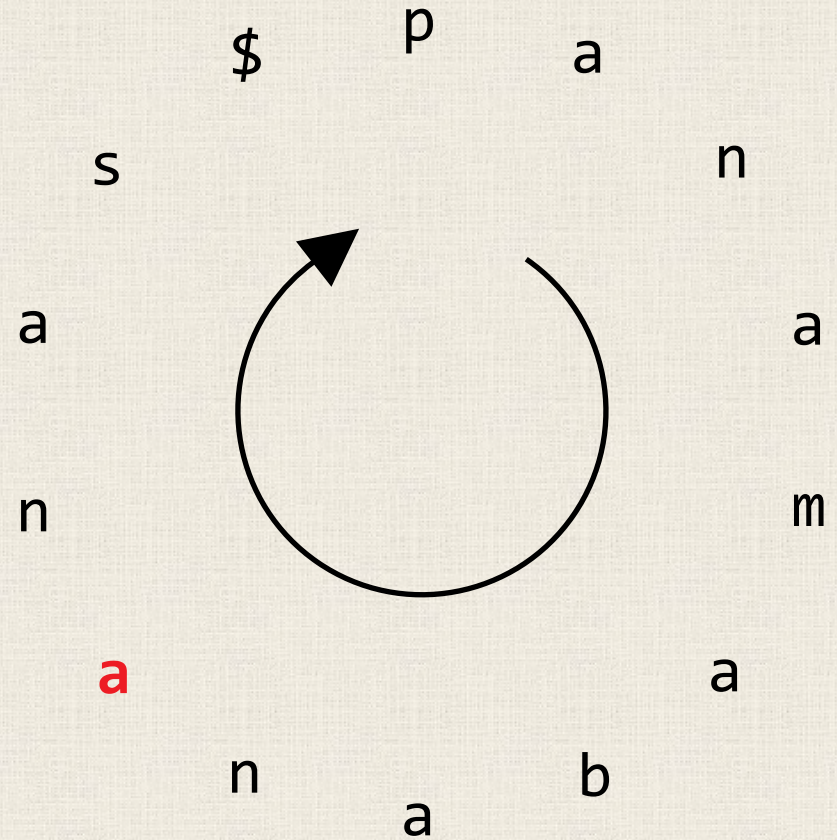
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanana  
nas$panamabana
```



Form all cyclic rotations of
"panamabananas\$"

The Burrows-Wheeler Transform

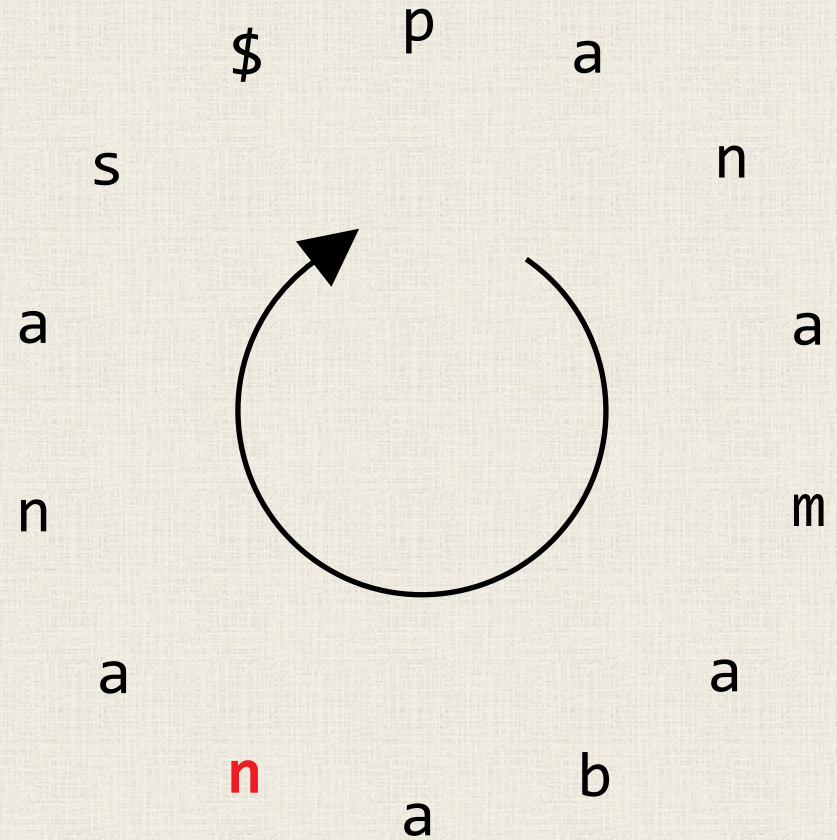
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanana  
nas$panamabana  
anas$panamaban
```



Form all cyclic rotations of
"panamabananas\$"

The Burrows-Wheeler Transform

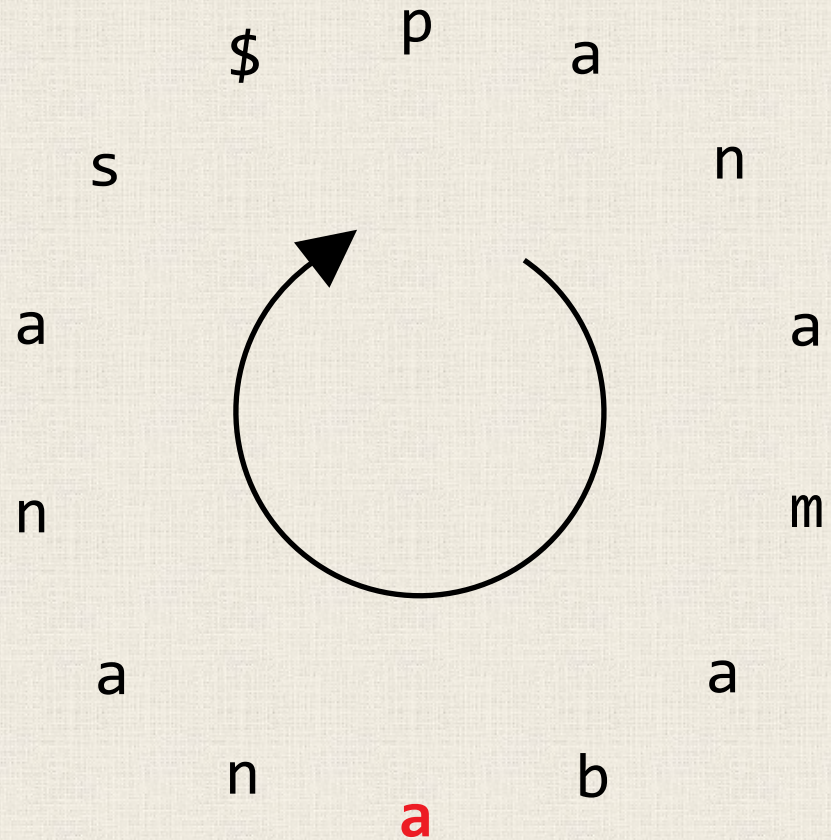
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanana  
nas$panamabana  
anas$panamaban  
nanas$panamaba
```



Form all cyclic rotations of
“panamabananas\$”

The Burrows-Wheeler Transform

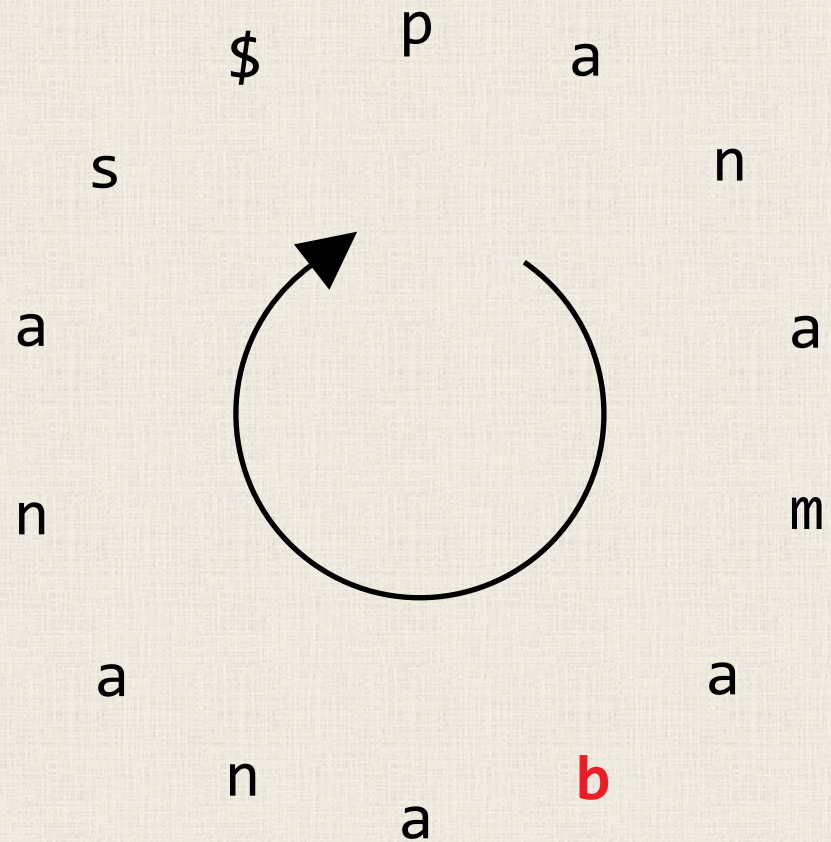
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanana  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
ananas$panamab
```



Form all cyclic rotations of
"panamabananas\$"

The Burrows-Wheeler Transform

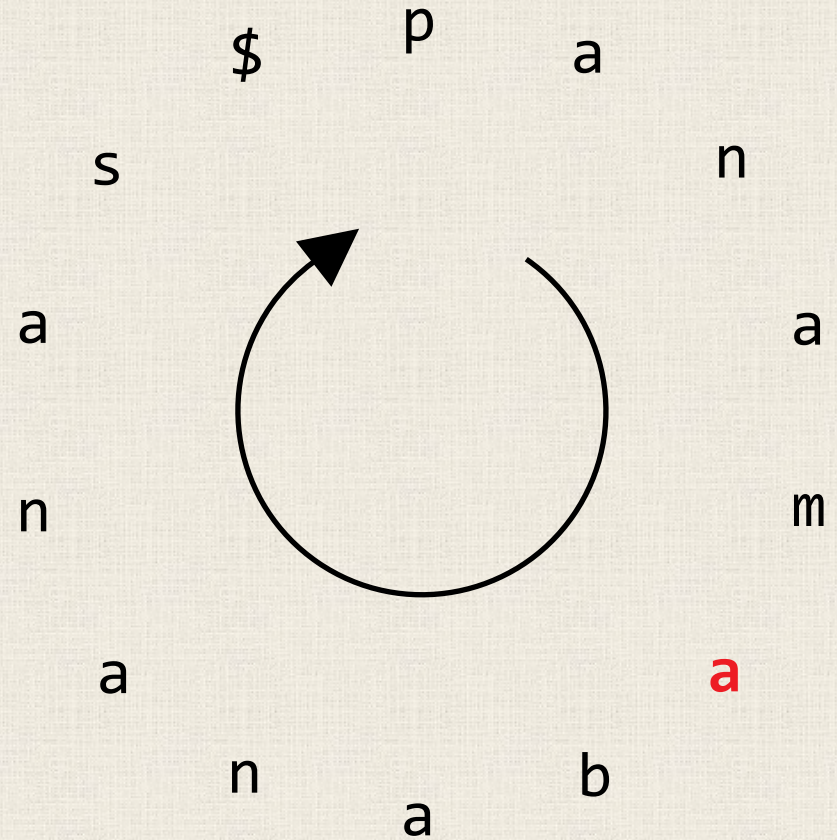
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamaban  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
ananas$panamab  
bananas$panama
```



Form all cyclic rotations of
“panamabananas\$”

The Burrows-Wheeler Transform

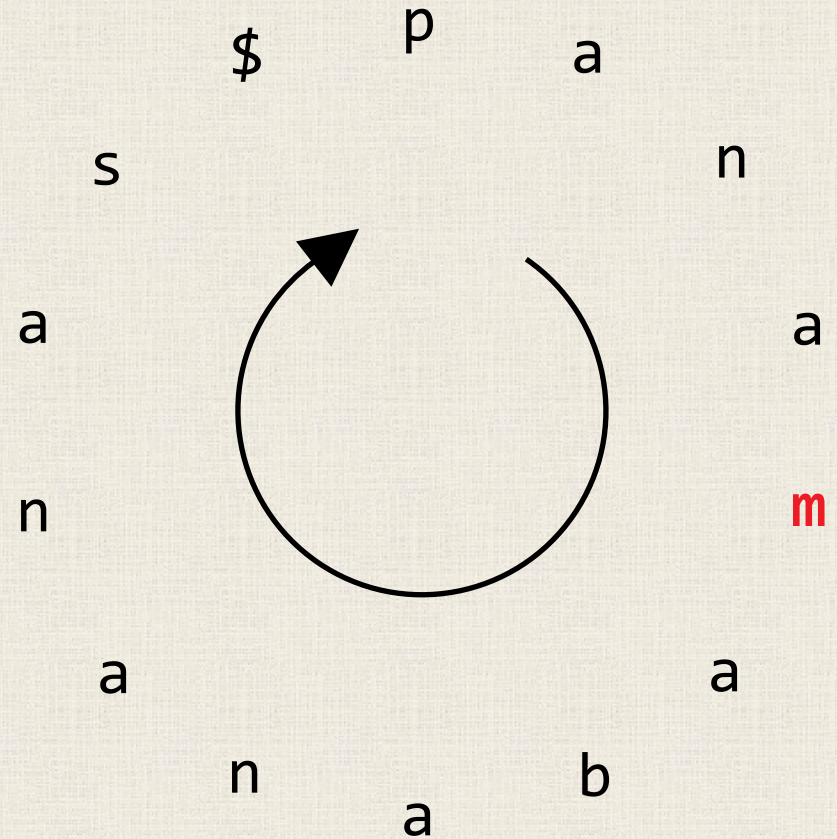
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanana  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
ananas$panamab  
bananas$panama  
abananas$panam
```



Form all cyclic rotations of
"panamabananas\$"

The Burrows-Wheeler Transform

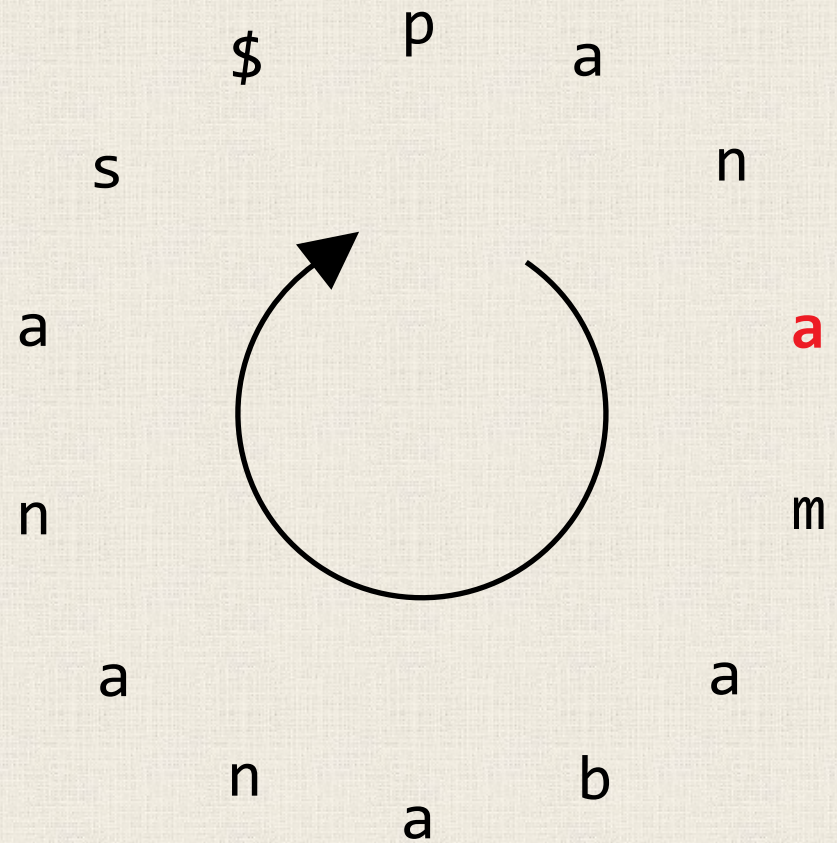
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanana  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
ananas$panamab  
bananas$panama  
abananas$panam  
mabananas$pana
```



Form all cyclic rotations of
“panamabananas\$”

The Burrows-Wheeler Transform

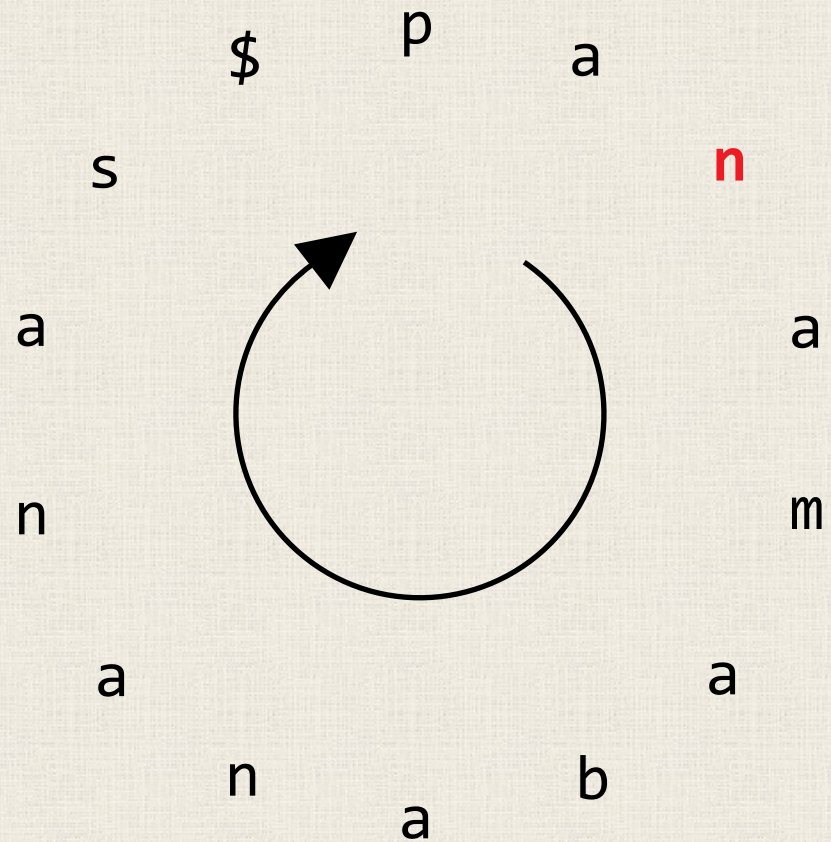
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanan  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
ananas$panamab  
bananas$panama  
abananas$panam  
mabananas$pana  
amabananas$pan
```



Form all cyclic rotations of
“panamabananas\$”

The Burrows-Wheeler Transform

```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanan  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
ananas$panamab  
bananas$panama  
abananas$panam  
mabananas$pana  
amabananas$pan  
namabananas$pa
```

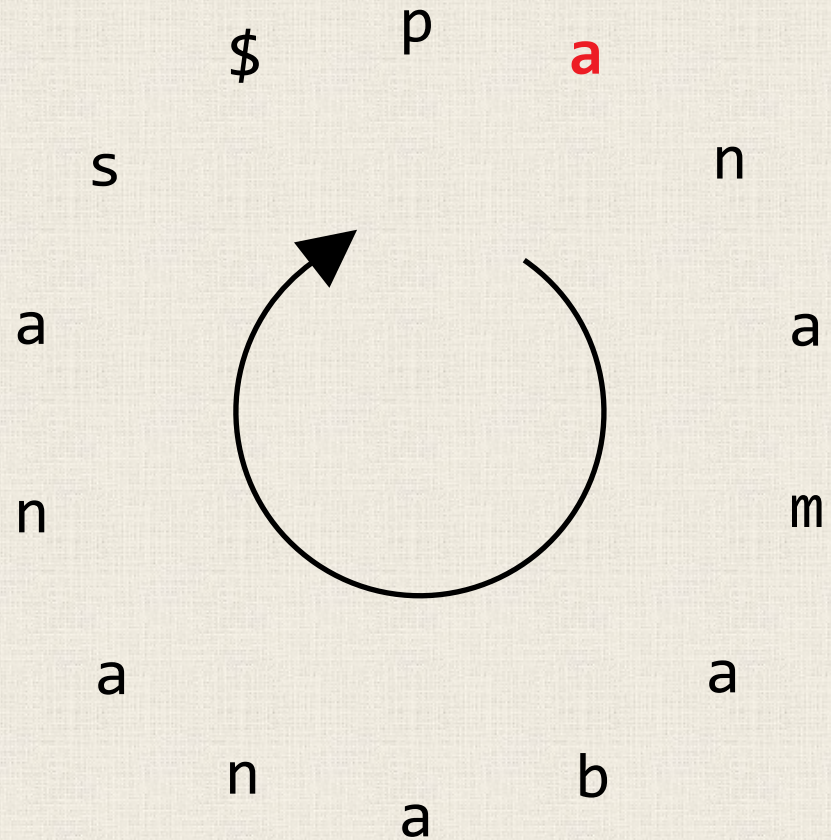


Form all cyclic rotations of
“panamabananas\$”

The Burrows-Wheeler Transform

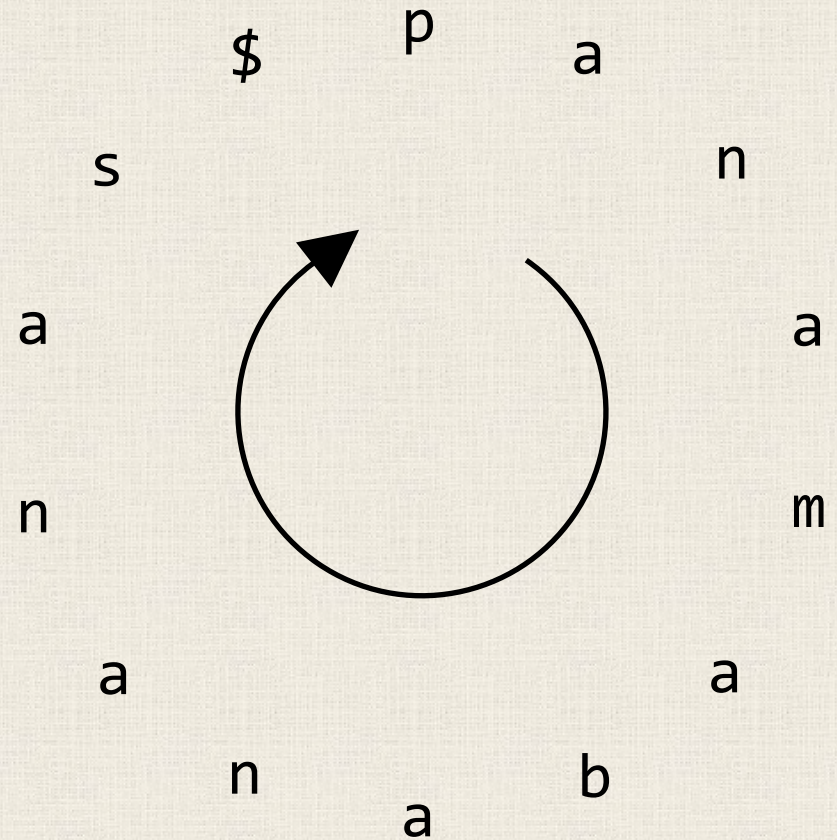
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamaban  
nas$panamabana  
anas$panamaban  
anas$panamaba  
ananas$panamab  
ananas$panamab  
bananas$panama  
abananas$panam  
mabananas$pana  
amabananas$pan  
namabananas$pa  
anamabananas$p
```

Form all cyclic rotations of
“panamabananas\$”



The Burrows-Wheeler Transform

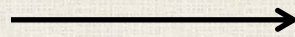
```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamabanan  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
ananas$panamab  
bananas$panama  
abananas$panam  
mabananas$pana  
amabananas$pan  
namabananas$pa  
anamabananas$p
```



Form all cyclic rotations of
“panamabananas\$”

The Burrows-Wheeler Transform

```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamaban  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
ananas$panamab  
bananas$panama  
abananas$panam  
mabananas$pana  
amabananas$pan  
namabananas$pa  
anamabananas$p
```



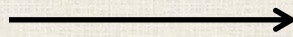
```
$panamabananas  
abananas$panam  
amabananas$pan  
anamabananas$p  
ananas$panamab  
anas$panamaban  
as$panamaban  
bananas$panama  
mabananas$pana  
namabananas$pa  
nanas$panamaba  
nas$panamabana  
panamabananas$  
s$panamabanana
```

Form all cyclic rotations of
“panamabananas\$”

Sort the strings
lexicographically
(\$ comes first)

The Burrows-Wheeler Transform

```
panamabananas$  
$panamabananas  
s$panamabanana  
as$panamaban  
nas$panamabana  
anas$panamaban  
nanas$panamaba  
nanas$panamab  
bananas$panama  
abananas$panam  
mabananas$pana  
amabananas$pan  
namabananas$pa  
anamabananas$p
```



```
$panamabananas  
abananas$panam  
amabananas$pan  
anamabananas$p  
ananas$panamab  
anas$panamaban  
as$panamaban  
bananas$panama  
mabananas$pana  
namabananas$pa  
nanas$panamaba  
nas$panamabana  
panamabananas$  
s$panamabanana
```

Form all cyclic rotations of
“panamabananas\$”

We call this matrix of
symbols $M(\text{Text})$

The Burrows-Wheeler Transform

panamabananas\$
 \$panamabananas
 s\$panamabanana
 as\$panamabanan
 nas\$panamabana
 anas\$panamaban
 nanas\$panamaba
 ananas\$panamab
 bananas\$panama
 abananas\$panam
 mabananas\$pana
 amabananas\$pan
 namabananas\$pa
 anamabananas\$p



\$panamabananas**s**
 abananas\$panam
 amabananas\$pan
 anamabananas\$**p**
 ananas\$panama**b**
 anas\$panamaban
 as\$panamabanan
 bananas\$panama
 mabananas\$pana
 namabananas\$pa
 nanas\$panamaba
 nas\$panamabana
 panamabananas\$**s**
 s\$panamabanana**a**

Form all cyclic rotations of
Text = "panamabananas\$"

$BWT(Text)$ = last column of
 $M(text)$ = "smn**p**bnnaaaaa**s**a".

Let's Examine BWT(*Text*) for *Text* = Watson & Crick, 1953

nd Corey (1). They kindly made their manuscript availa a
nd criticism, especially on interatomic distances. We a
nd cytosine. The sequence of bases on a single chain d a
nd experimentally (3,4) that the ratio of the amounts o u
nd for this reason we shall not comment on it. We wish a
nd guanine (purine) with cytosine (pyrimidine). In oth a
nd ideas of Dr. M. H. F. Wilkins, Dr. R. E. Franklin a
nd its water content is rather high. At lower water co a
nd pyrimidine bases. The planes of the bases are perpe a
nd stereochemical arguments. It has not escaped our no a
nd that only specific pairs of bases can bond together u
nd the atoms near it is close to Furberg's 'standard co a
nd the bases on the inside, linked together by hydrogen a
nd the bases on the outside. In our opinion, this stru a
nd the other a pyrimidine for bonding to occur. The hy a
nd the phosphates on the outside. The configuration of a
nd the ration of guanine to cytosine, are always very c a
nd the same axis (see diagram). We have made the usual u
nd their co-workers at King's College, London. One of a

Where do you think the “u”s come from?

nd Corey (1). They kindly made their manuscript available a
nd criticism, especially on interatomic distances. We a
nd cytosine. The sequence of bases on a single chain a
nd experimentally (3,4) that the ratio of the amounts of u
nd for this reason we shall not comment on it. We wish a
nd guanine (purine) with cytosine (pyrimidine). In other a
nd ideas of Dr. M. H. F. Wilkins, Dr. R. E. Franklin a
nd its water content is rather high. At lower water content a
nd pyrimidine bases. The planes of the bases are perpendicular a
nd stereochemical arguments. It has not escaped our notice a
nd that only specific pairs of bases can bond together u
nd the atoms near it is close to Furberg's 'standard conformation a
nd the bases on the inside, linked together by hydrogen bonds a
nd the bases on the outside. In our opinion, this structure a
nd the other a pyrimidine for bonding to occur. The hydrogen a
nd the phosphates on the outside. The configuration of a
nd the ratios of guanine to cytosine, are always very constant a
nd the same axis (see diagram). We have made the usual u
nd their co-workers at King's College, London. One of a

BWT Converts Repeats to Runs

Text

Burrows-Wheeler Transform!

Convert repeats to runs

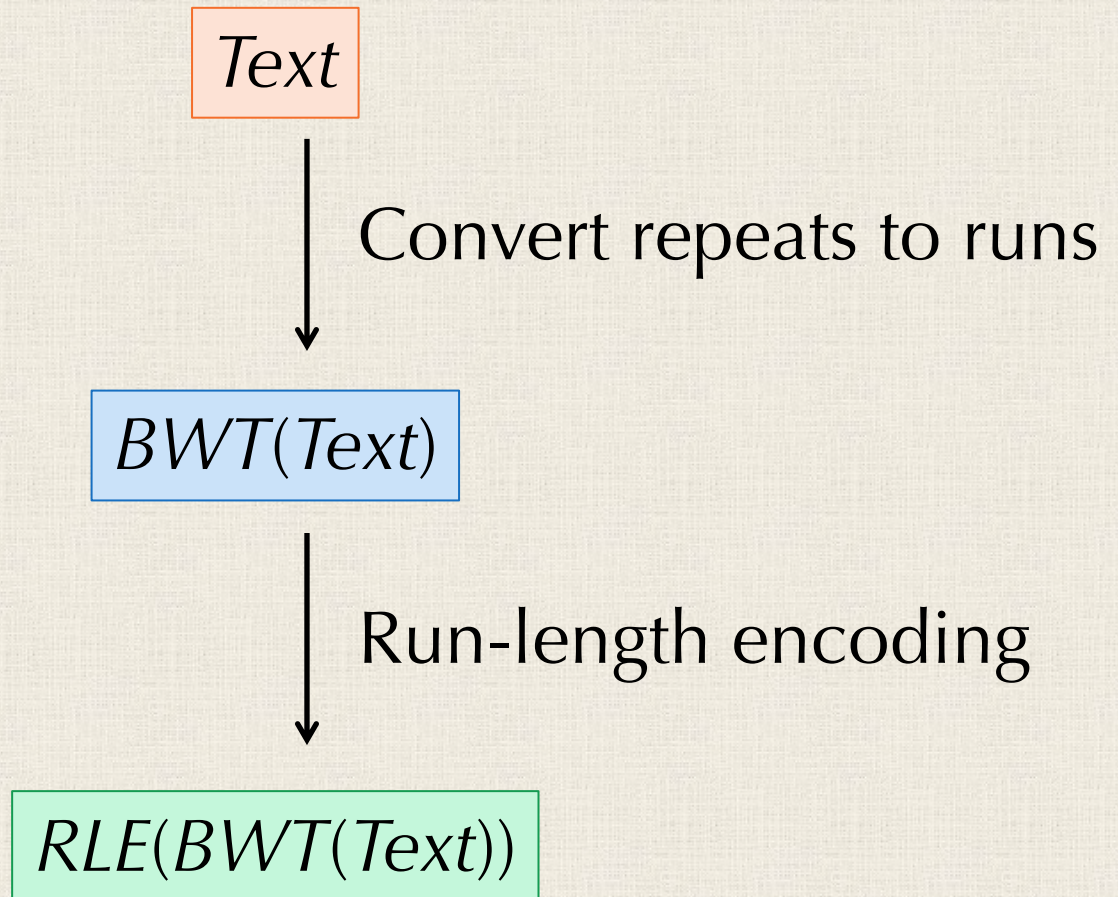
BWT(Text)

Run-length encoding

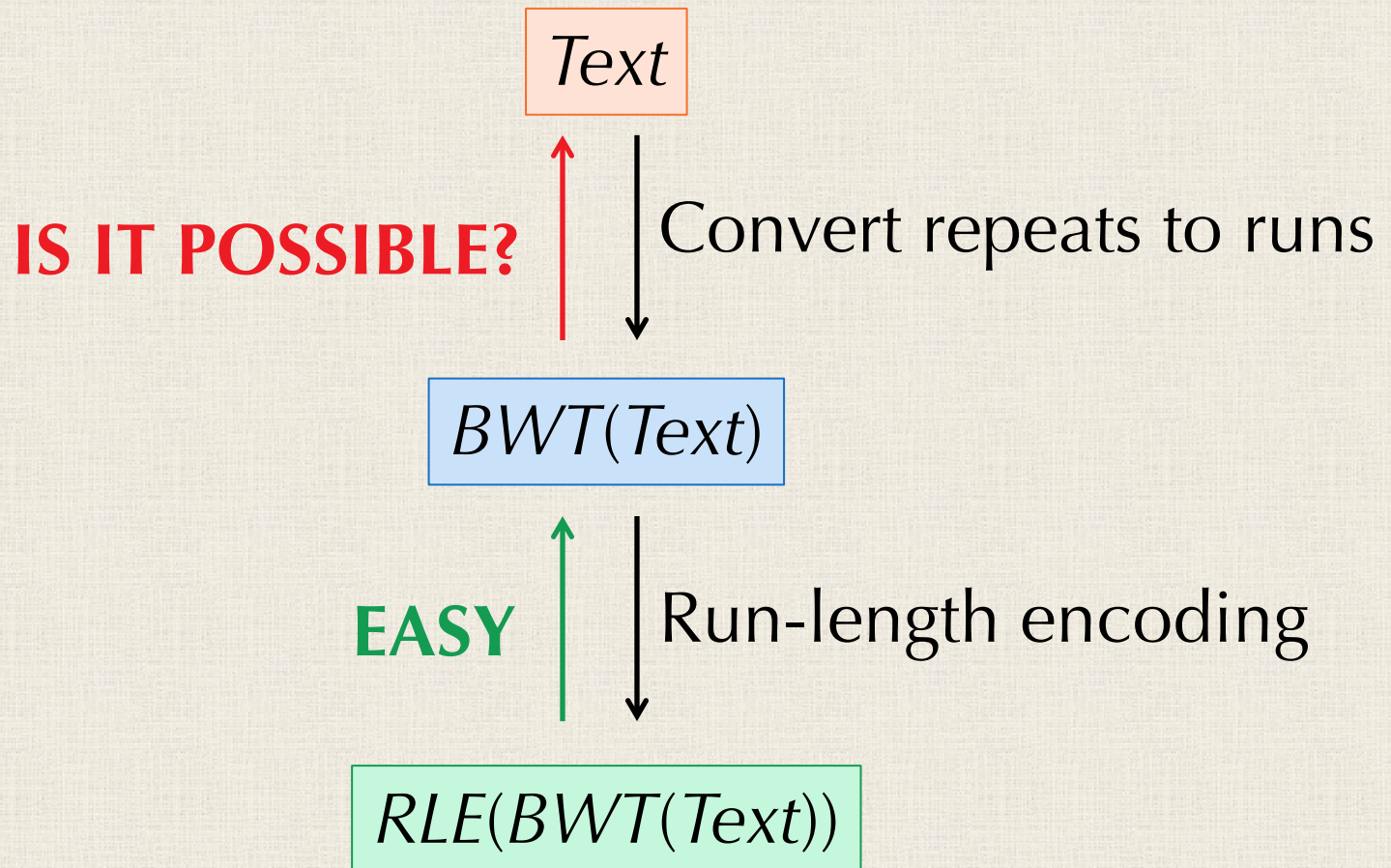
RLE(BWT(Text))

INVERTING THE BURROWS- WHEELER TRANSFORM

How Can We Decompress?

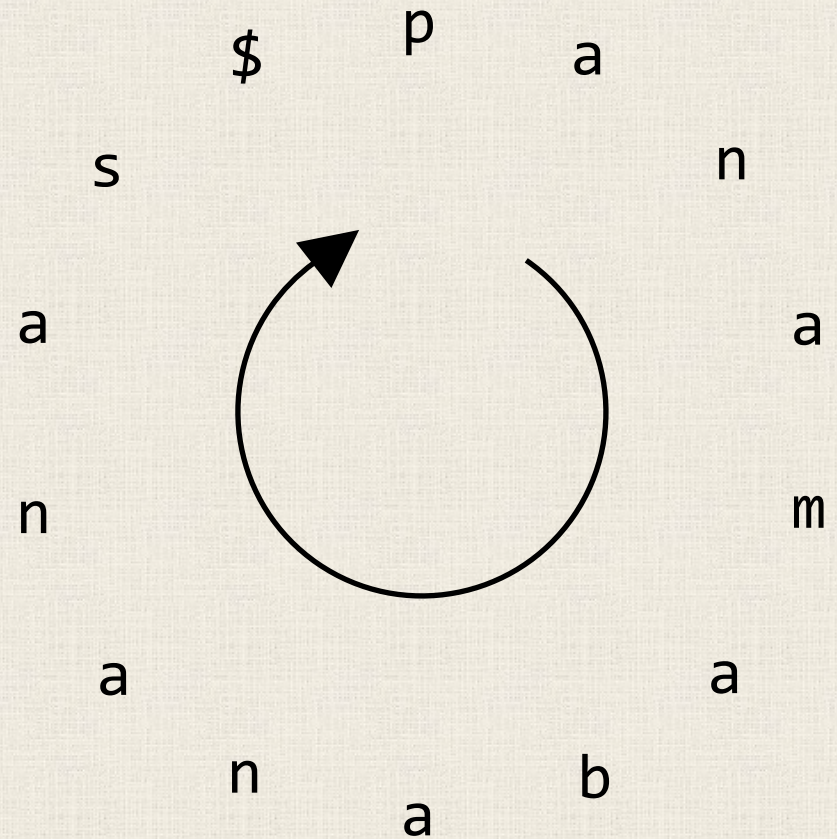


How Can We Decompress?



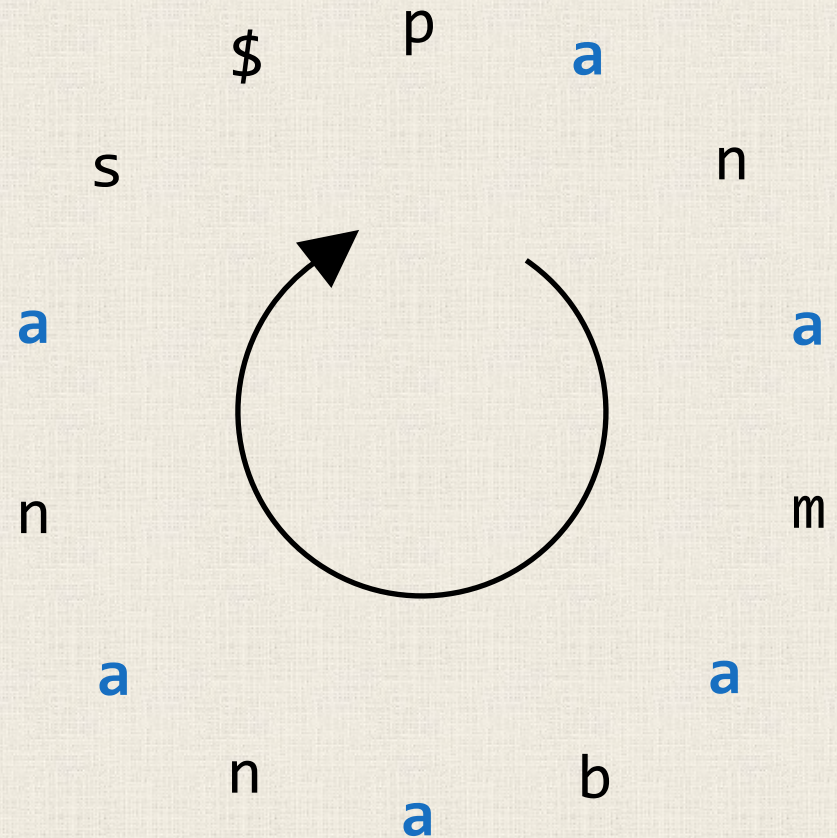
A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana



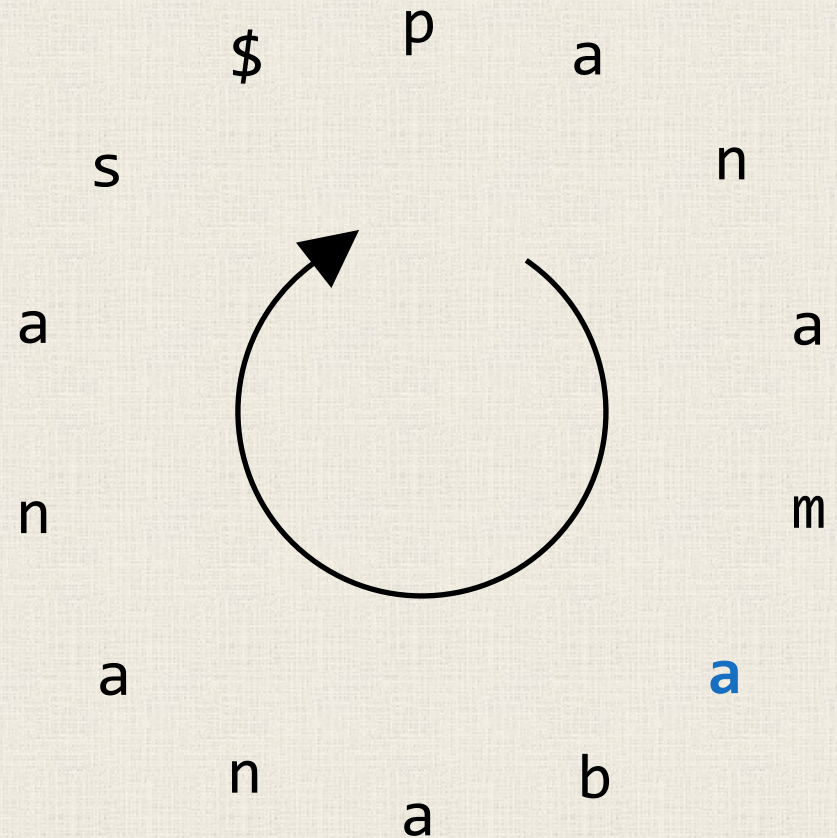
A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$p
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabanana



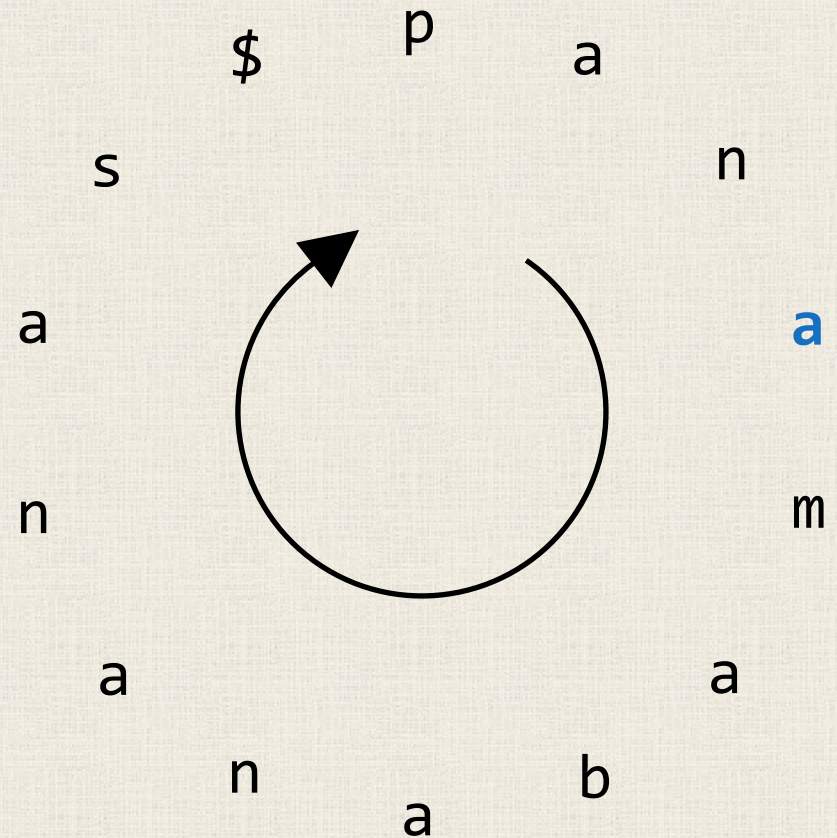
A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana



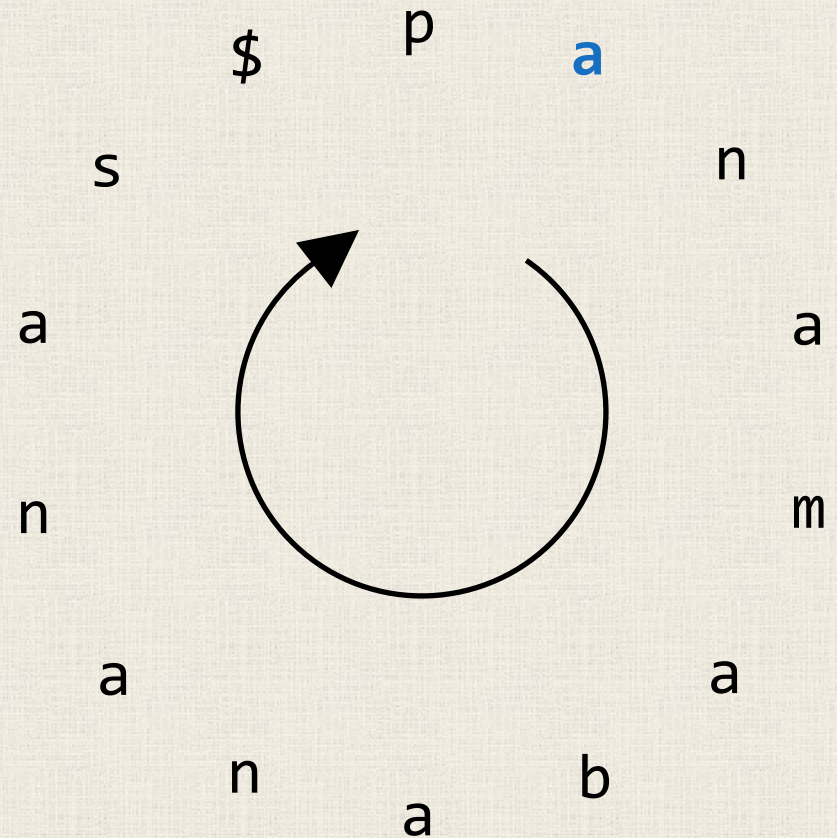
A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamabanan
bananas\$panama
mabananas\$pana**a**
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana



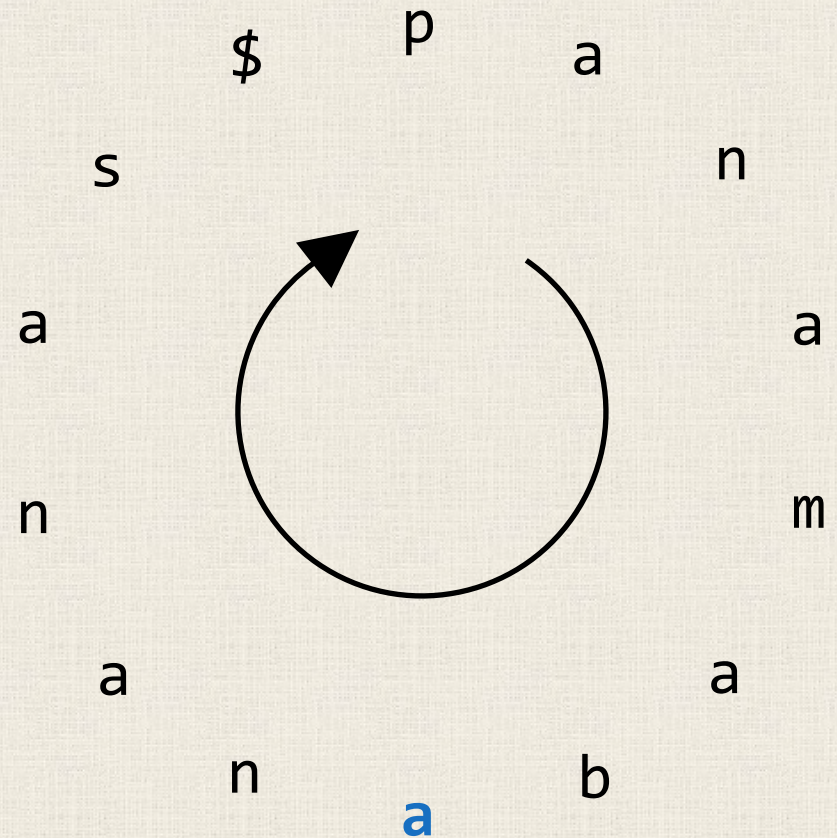
A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$p**a**
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana



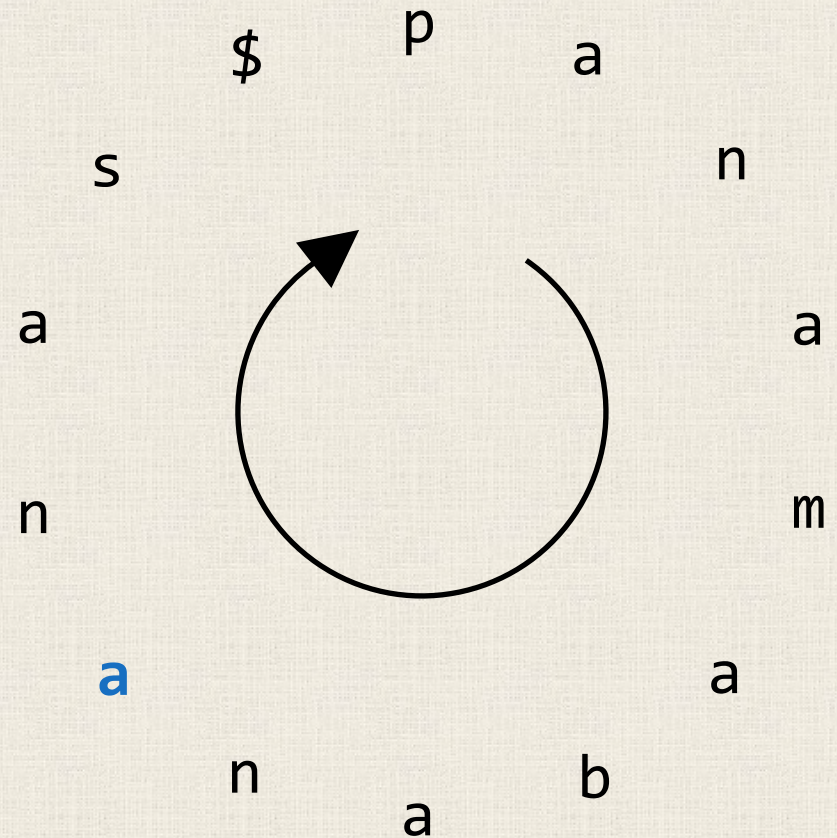
A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamabanan
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamab**a**
nas\$panamabana
panamabananas\$
s\$panamabanana



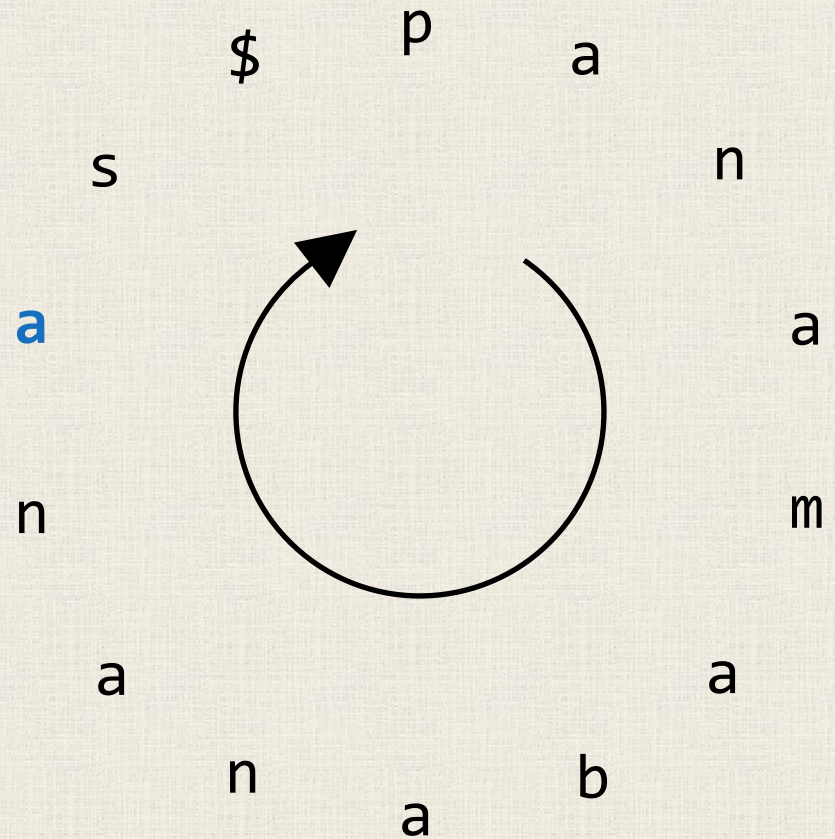
A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabanana



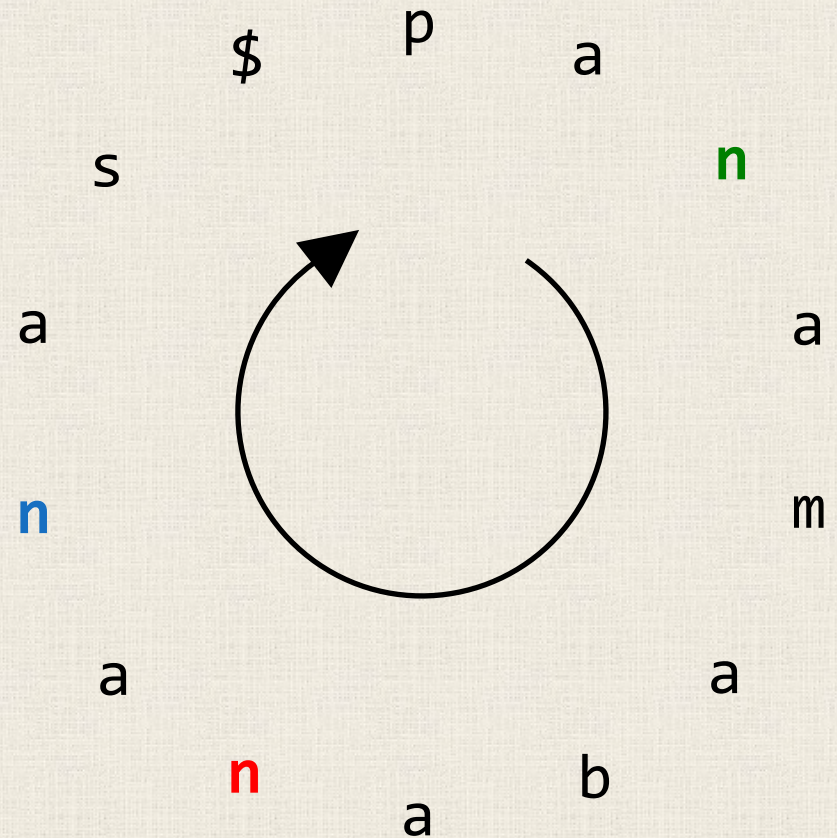
A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabanana**a**



A Strange Observation

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamabana
as\$panamabana
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabanana



Is It True in General?

- 1
- 2
- 3
- 4
- 5
- 6

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana

Is It True in General?

	\$panamabananas
1	abananas\$panam
2	amabananas\$pan
3	anamabananas\$p
4	ananas\$panamab
5	anas\$panamaban
6	as\$panamabanan
	bananas\$panama
	mabananas\$pana
	namabananas\$pa
	nanas\$panamaba
	nas\$panamabana
	panamabananas\$
	s\$panamabana

These strings are sorted

Is It True in General?

1 \$panamabananas
2 abananas\$panam
3 amabananas\$pan
4 anamabananas\$p
5 ananas\$panamab
6 anas\$panamaban
as\$panamabanan
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana

Chop off **a**



bananas\$panam
mabananas\$pan
namabananas\$p
nanas\$panamab
nas\$panamaban
s\$panamabanan

These strings are sorted

Is It True in General?

1 \$panamabananas
2 abananas\$panam
3 amabananas\$pan
4 anamabananas\$p
5 ananas\$panamab
6 anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana

Chop off **a**



bananas\$panam
mabananas\$pan
namabananas\$p
nanas\$panamab
nas\$panamaban
s\$panamaban

Still
sorted

These strings are sorted

Is It True in General?

1 \$panamabananas
2 abananas\$panam
3 amabananas\$pan
4 anamabananas\$p
5 ananas\$panamab
6 anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabanana

These strings are sorted

Chop off **a**

bananas\$panam
mabananas\$pan
namabananas\$p
nanas\$panamab
nas\$panamaban
s\$panamaban

Still
sorted

Add **a**
to end

bananas\$panama**a**
mabananas\$pana**a**
namabananas\$pa**a**
nanas\$panamaba**a**
nas\$panamabana**a**
s\$panamabanana**a**

Is It True in General?

1 \$panamabananas
2 abananas\$panam
3 amabananas\$pan
4 anamabananas\$p
5 ananas\$panamab
6 anas\$panamaban
as\$panamaban
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabanana

These strings are sorted

Chop off **a**

bananas\$panam
mabananas\$pan
namabananas\$p
nanas\$panamab
nas\$panamaban
s\$panamaban

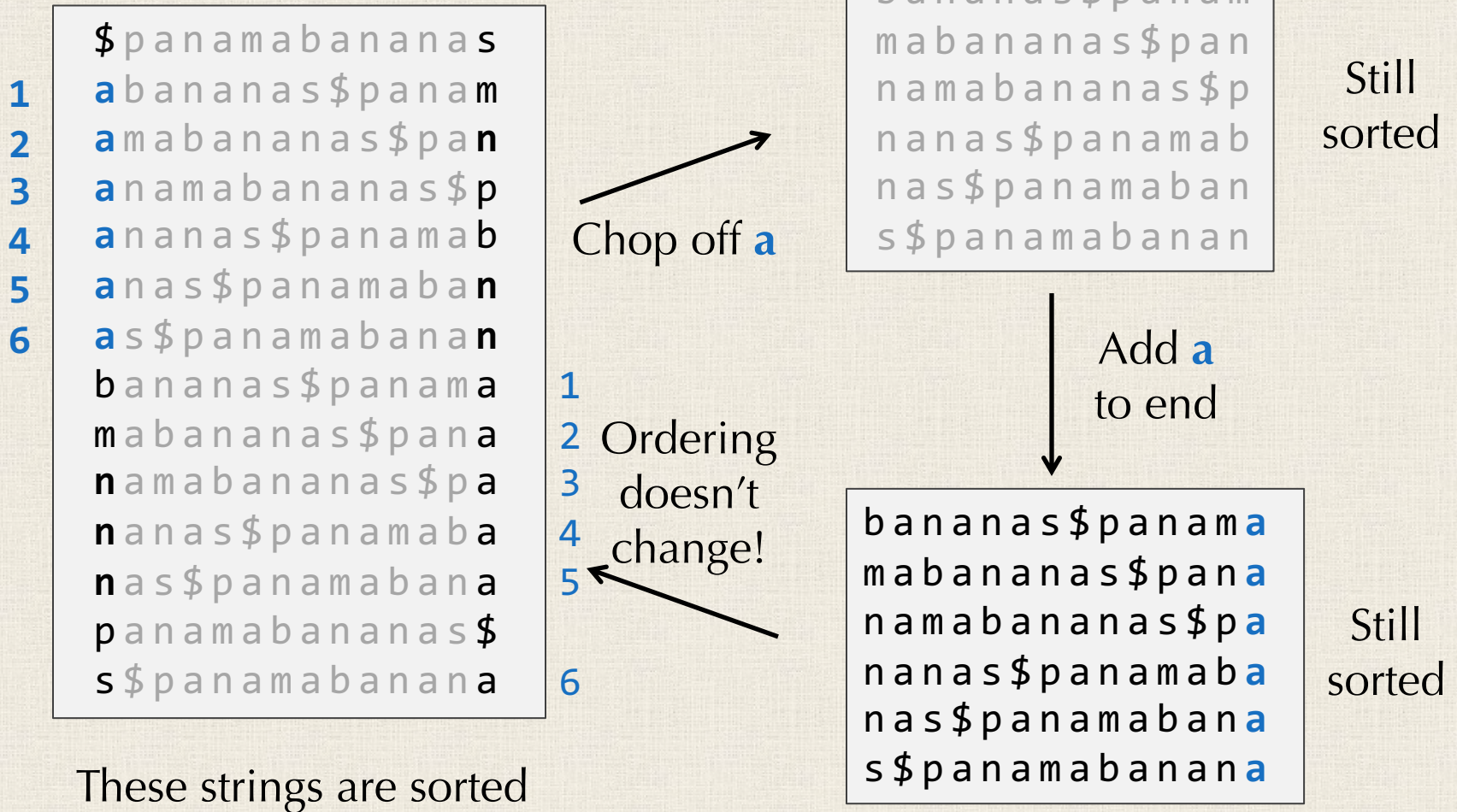
Still
sorted

Add **a**
to end

bananas\$panama**a**
mabananas\$pana**a**
namabananas\$pa**a**
nanas\$panamaba**a**
nas\$panamabana**a**
s\$panamabanana**a**

Still
sorted

Is It True in General?



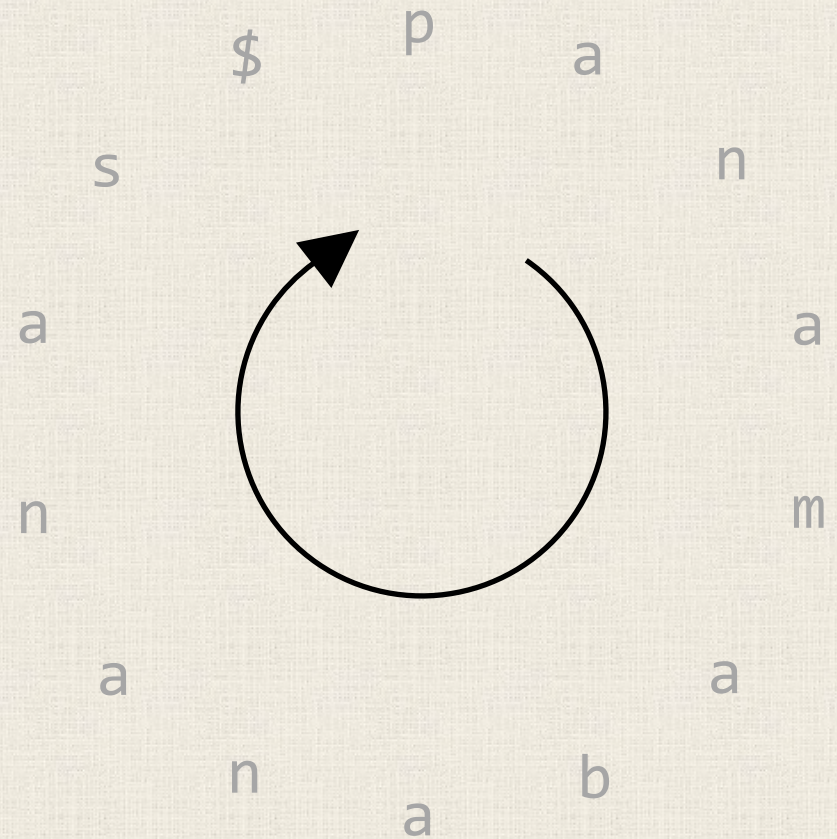
Is It True in General?

$\$$ ₁ panamabananas s ₁
 a ₁ bananas\$panam m ₁
 a ₂ mabananas\$pan n ₁
 a ₃ namabananas\$ p ₁
 a ₄ nanas\$panama b ₁
 a ₅ nas\$panamaba n ₂
 a ₆ s\$panamabana n ₃
 b ₁ ananas\$panama a ₁
 m ₁ abananas\$pana a ₂
 n ₁ amabananas\$pa a ₃
 n ₂ anas\$panamaba a ₄
 n ₃ as\$panamabana a ₅
 p ₁ anamabananas\$ s ₁
 s ₁ \$panamabana a ₆

First-Last Property: The k -th occurrence of *symbol* in *FirstColumn* and the k -th occurrence of *symbol* in *LastColumn* correspond to the same position of *symbol* in *Text*.

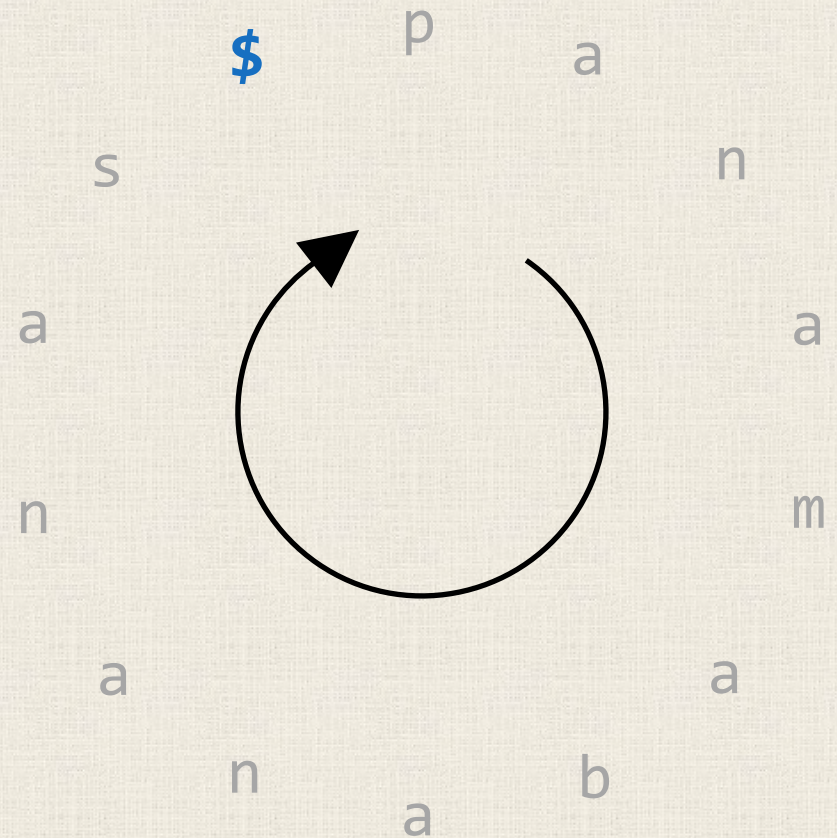
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$p3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



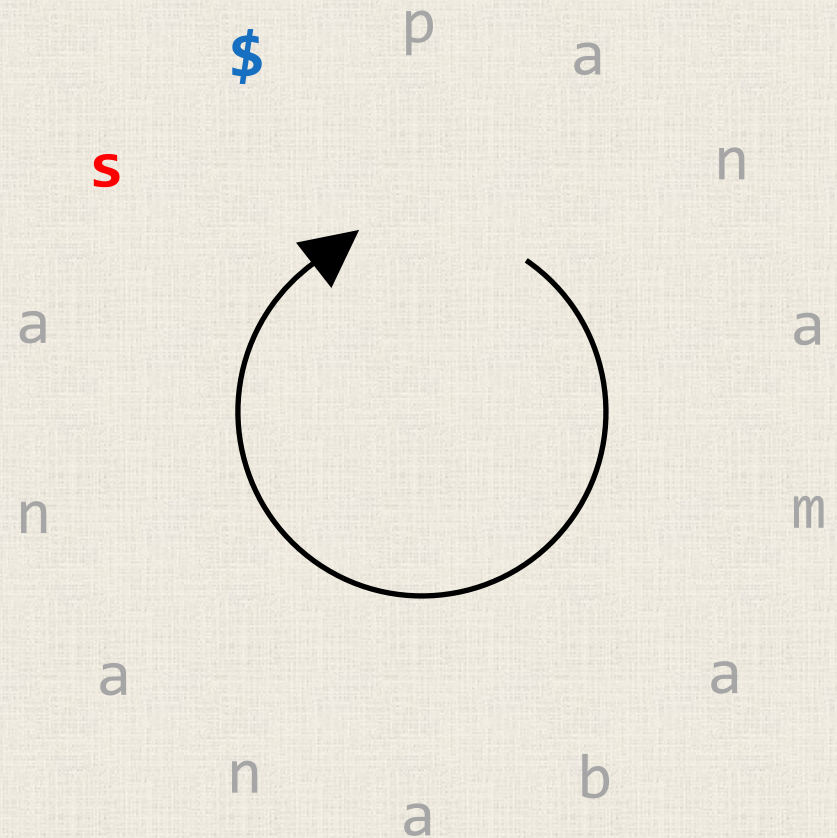
Efficient BWT Decompression

\$ ₁	p	a	n	a	m	a	b	a	n	a	n	a	s	\$ ₁
a ₁	b	a	n	a	n	a	s	\$	p	a	n	a	m	\$ ₁
a ₂	m	a	b	a	n	a	n	a	s	\$	p	a	n	\$ ₁
a ₃	n	a	m	a	b	a	n	a	n	a	s	\$	p	\$ ₁
a ₄	n	a	n	a	s	\$	p	a	n	a	m	a	b	\$ ₁
a ₅	n	a	s	\$	p	a	n	a	m	a	b	a	n	\$ ₂
a ₆	s	\$	p	a	n	a	m	a	b	a	n	a	n	\$ ₃
b ₁	a	n	a	n	a	s	\$	p	a	n	a	m	a	\$ ₁
m ₁	a	b	a	n	a	n	a	s	\$	p	a	n	a	\$ ₂
n ₁	a	m	a	b	a	n	a	n	a	s	\$	p	a	\$ ₃
n ₂	a	n	a	s	\$	p	a	n	a	m	a	b	a	\$ ₄
n ₃	a	s	\$	p	a	n	a	m	a	b	a	n	a	\$ ₅
p ₁	a	n	a	m	a	b	a	n	a	n	a	s	\$	\$ ₁
s ₁	\$	p	a	n	a	m	a	b	a	n	a	n	a	\$ ₆



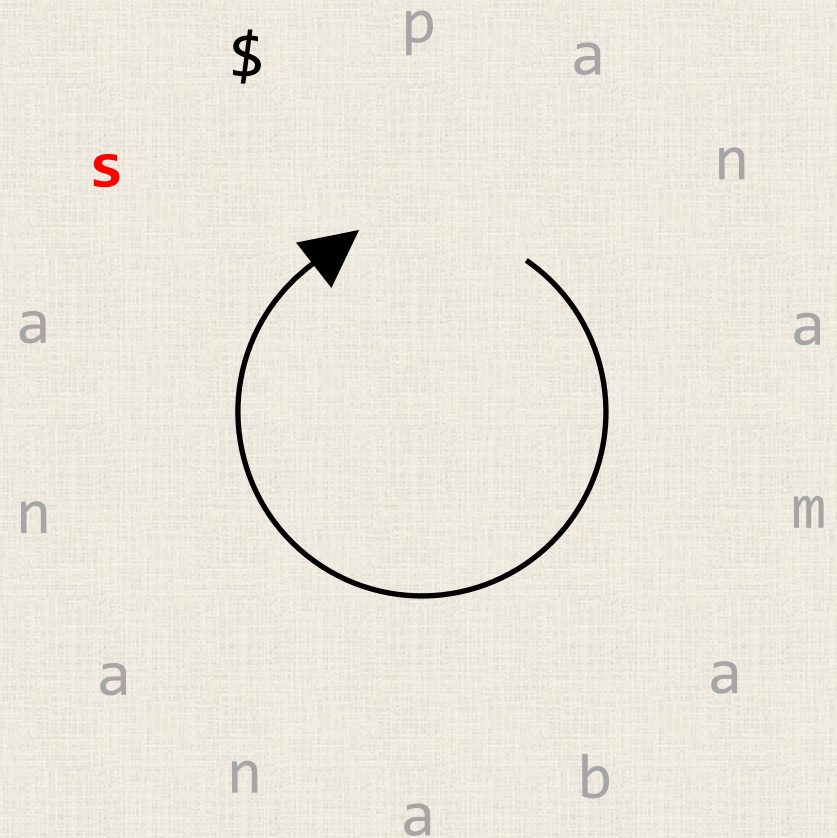
Efficient BWT Decompression

\$ ₁	p	a	n	a	m	a	b	a	n	a	n	a	s	\$ ₁
a ₁	b	a	n	a	n	a	s	\$	p	a	n	a	m	\$ ₁
a ₂	m	a	b	a	n	a	n	a	s	\$	p	a	n	\$ ₁
a ₃	n	a	m	a	b	a	n	a	n	a	s	\$	p	\$ ₁
a ₄	n	a	n	a	s	\$	p	a	n	a	m	a	b	\$ ₁
a ₅	n	a	s	\$	p	a	n	a	m	a	b	a	n	\$ ₂
a ₆	s	\$	p	a	n	a	m	a	b	a	n	a	n	\$ ₃
b ₁	a	n	a	n	a	s	\$	p	a	n	a	m	a	\$ ₁
m ₁	a	b	a	n	a	n	a	s	\$	p	a	n	a	\$ ₂
n ₁	a	m	a	b	a	n	a	n	a	s	\$	p	a	\$ ₃
n ₂	a	n	a	s	\$	p	a	n	a	m	a	b	a	\$ ₄
n ₃	a	s	\$	p	a	n	a	m	a	b	a	n	a	\$ ₅
p ₁	a	n	a	m	a	b	a	n	a	n	a	s	\$	\$ ₁
s ₁	\$	p	a	n	a	m	a	b	a	n	a	n	a	\$ ₆



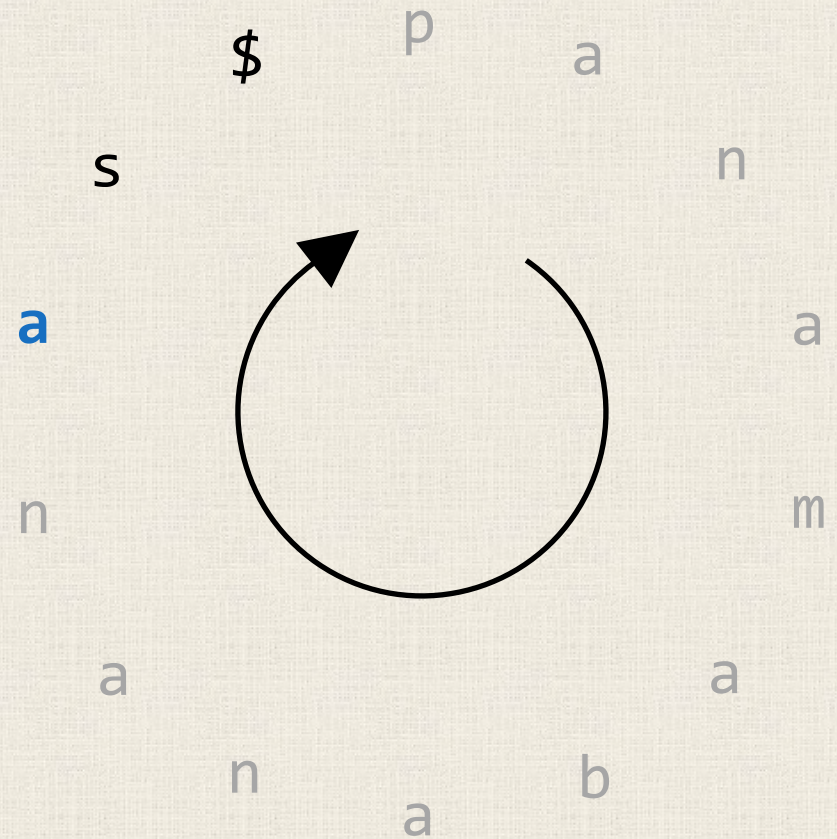
Efficient BWT Decompression

```
$1panamabananaS1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
S1$panamabanana6
```



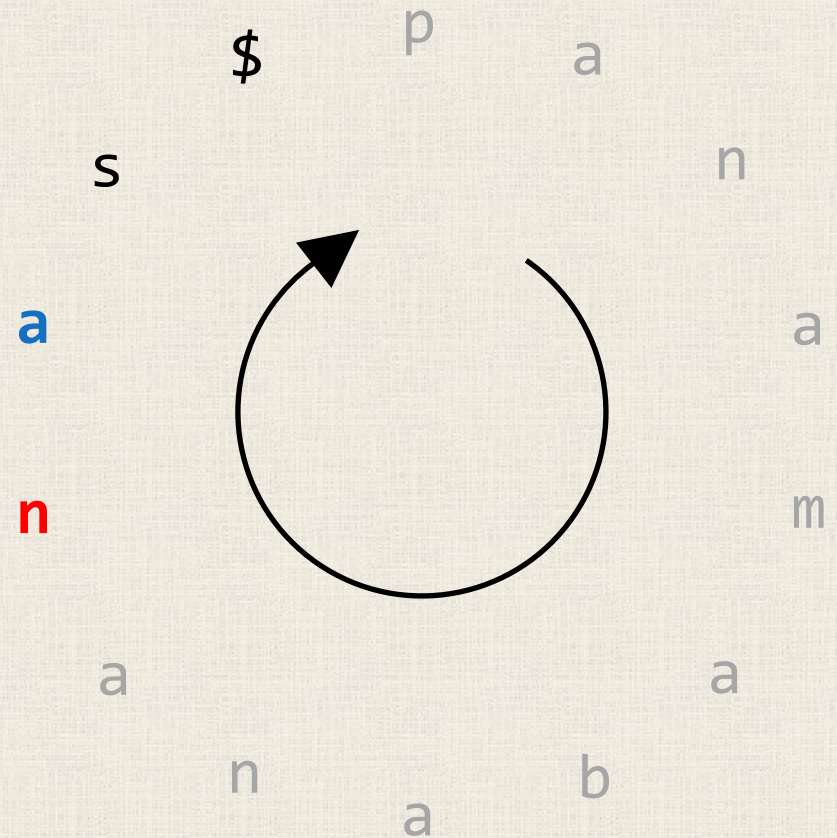
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$p3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabananaa6
```



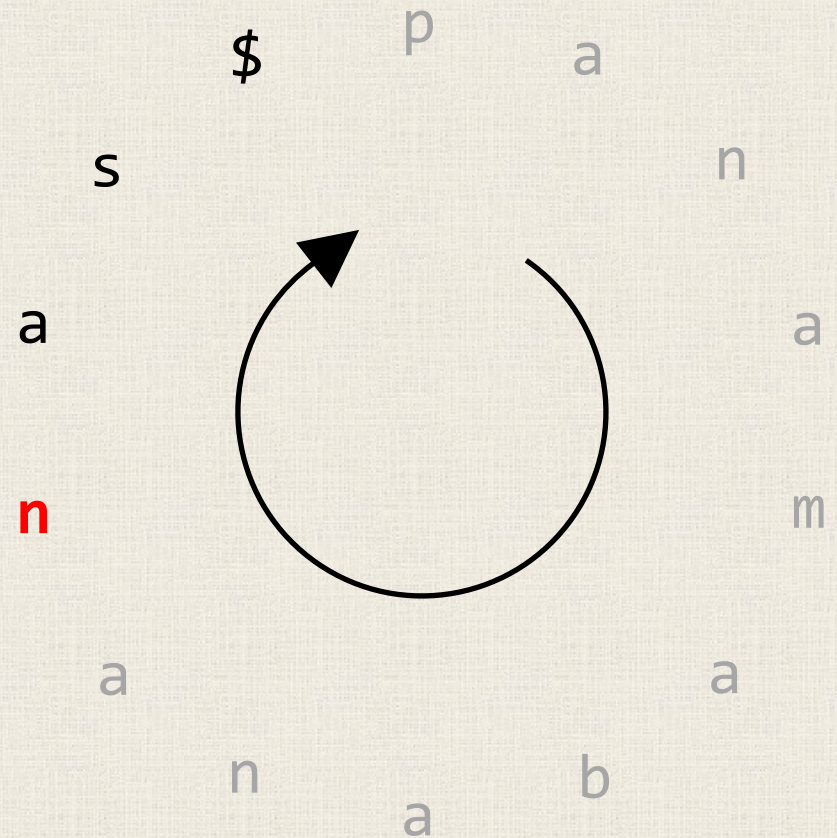
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamabanan3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanaa6
```



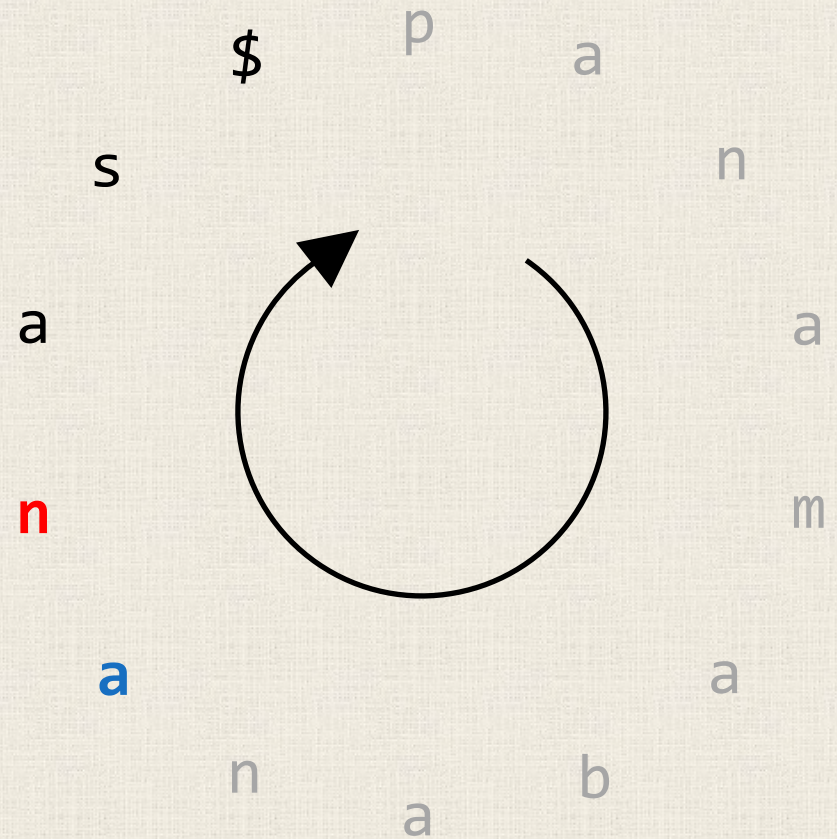
Efficient BWT Decompression

\$ ₁	p	a	n	a	m	a	b	a	n	a	n	a	s	\$ ₁
a ₁	b	a	n	a	n	a	s	\$	p	a	n	a	m	\$ ₁
a ₂	m	a	b	a	n	a	n	a	s	\$	p	a	n	\$ ₁
a ₃	n	a	m	a	b	a	n	a	n	a	s	\$	p	\$ ₁
a ₄	n	a	n	a	s	\$	p	a	n	a	m	a	b	\$ ₁
a ₅	n	a	s	\$	p	a	n	a	m	a	b	a	n	\$ ₂
a ₆	s	\$	p	a	n	a	m	a	b	a	n	a	n ₃	
b ₁	a	n	a	n	a	s	\$	p	a	n	a	m	a	\$ ₁
m ₁	a	b	a	n	a	n	a	s	\$	p	a	n	a	\$ ₂
n ₁	a	m	a	b	a	n	a	n	a	s	\$	p	a	\$ ₃
n ₂	a	n	a	s	\$	p	a	n	a	m	a	b	a	\$ ₄
n ₃	a	s	\$	p	a	n	a	m	a	b	a	n	a	\$ ₅
p ₁	a	n	a	m	a	b	a	n	a	n	a	s	\$	\$ ₁
s ₁	\$	p	a	n	a	m	a	b	a	n	a	n	a	\$ ₆



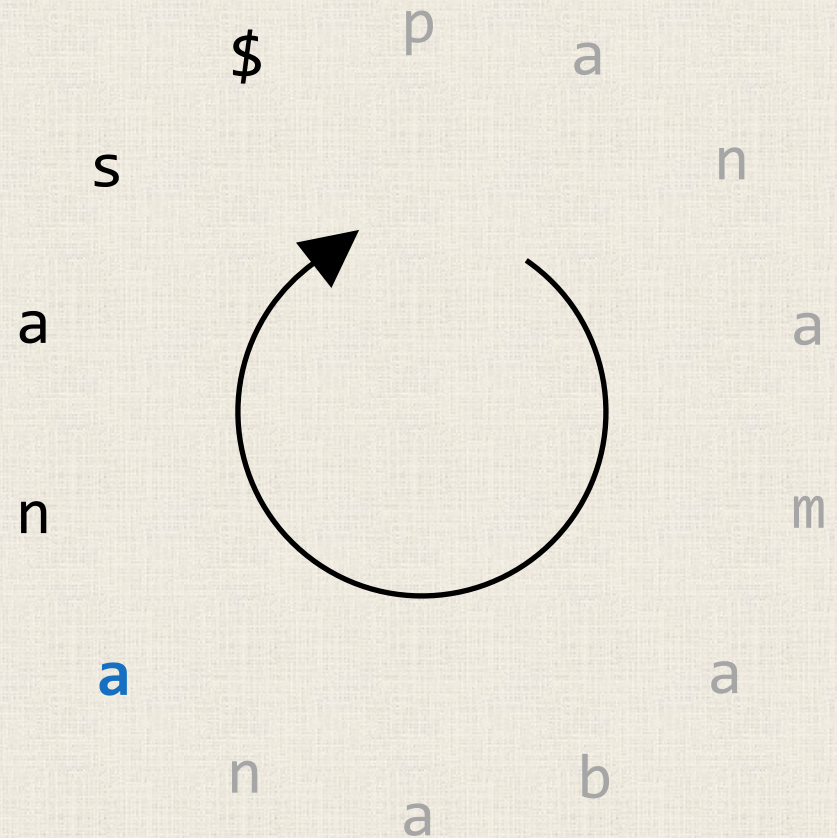
Efficient BWT Decompression

\$ ₁	p	a	n	a	m	a	b	a	n	a	n	a	s	\$ ₁
a ₁	b	a	n	a	n	a	s	\$	p	a	n	a	m	\$ ₁
a ₂	m	a	b	a	n	a	n	a	s	\$	p	a	n	\$ ₁
a ₃	n	a	m	a	b	a	n	a	n	a	s	\$	p	\$ ₁
a ₄	n	a	n	a	s	\$	p	a	n	a	m	a	b	\$ ₁
a ₅	n	a	s	\$	p	a	n	a	m	a	b	a	n	\$ ₂
a ₆	s	\$	p	a	n	a	m	a	b	a	n	a	n	\$ ₃
b ₁	a	n	a	n	a	s	\$	p	a	n	a	m	a	\$ ₁
m ₁	a	b	a	n	a	n	a	s	\$	p	a	n	a	\$ ₂
n ₁	a	m	a	b	a	n	a	n	a	s	\$	p	a	\$ ₃
n ₂	a	n	a	s	\$	p	a	n	a	m	a	b	a	\$ ₄
n₃	a	s	\$	p	a	n	a	m	a	b	a	n	a	\$₅
p ₁	a	n	a	m	a	b	a	n	a	n	a	s	\$	\$ ₁
s ₁	\$	p	a	n	a	m	a	b	a	n	a	n	a	\$ ₆



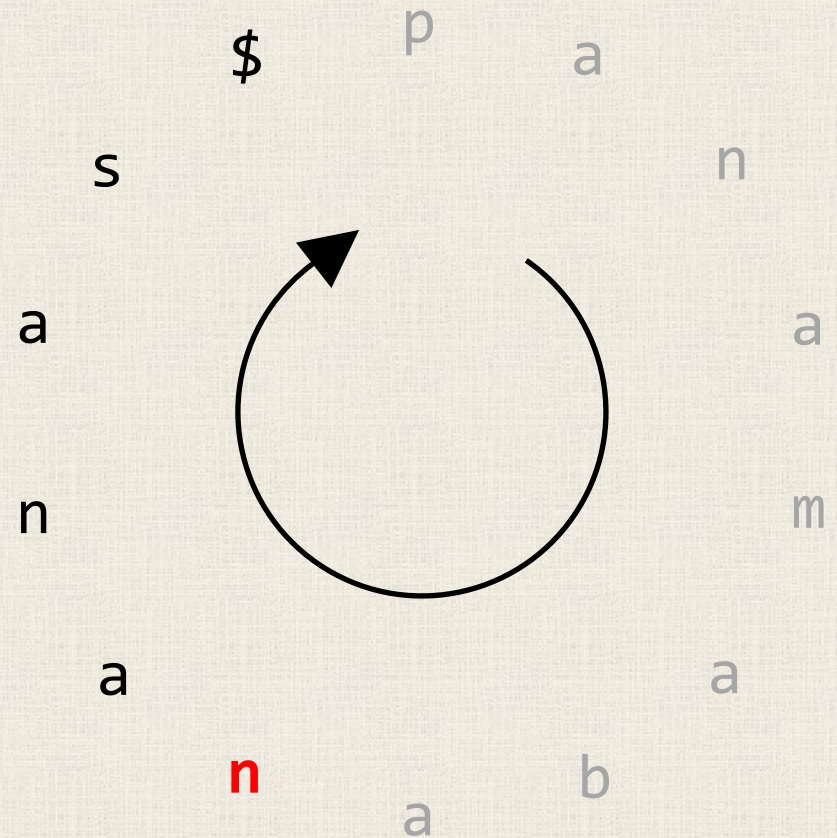
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$p3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



Efficient BWT Decompression

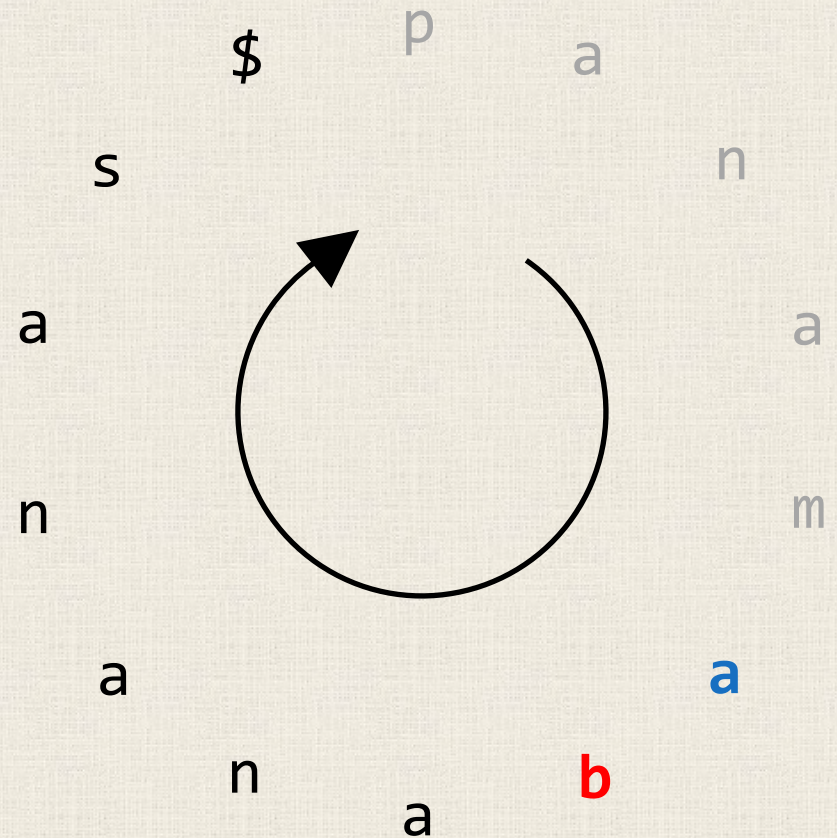
```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamabann3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



Efficient BWT Decompression

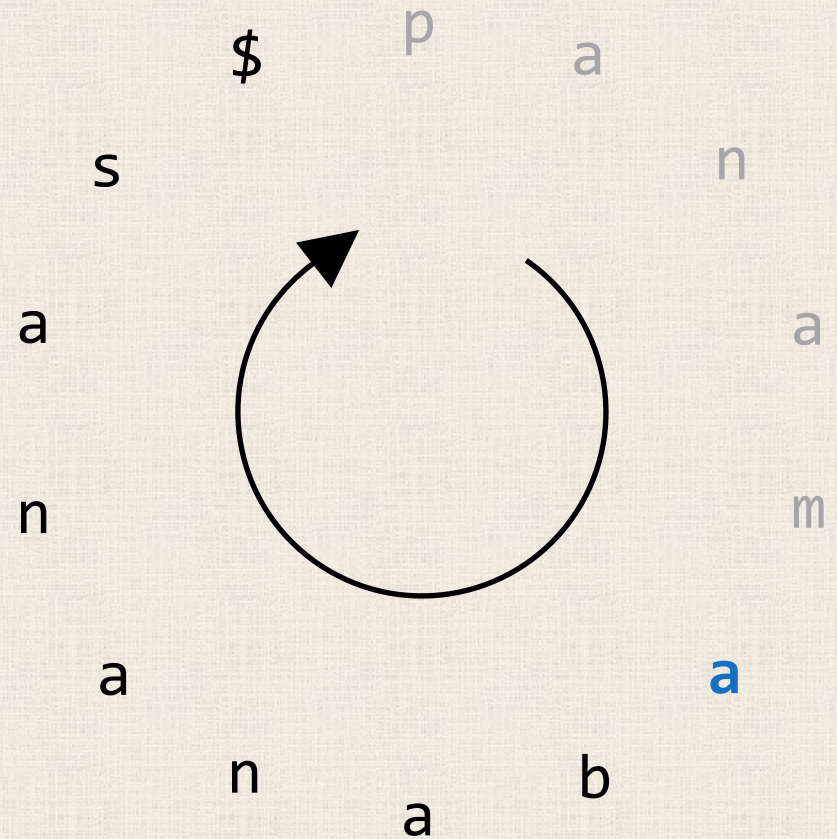
```

$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
    
```



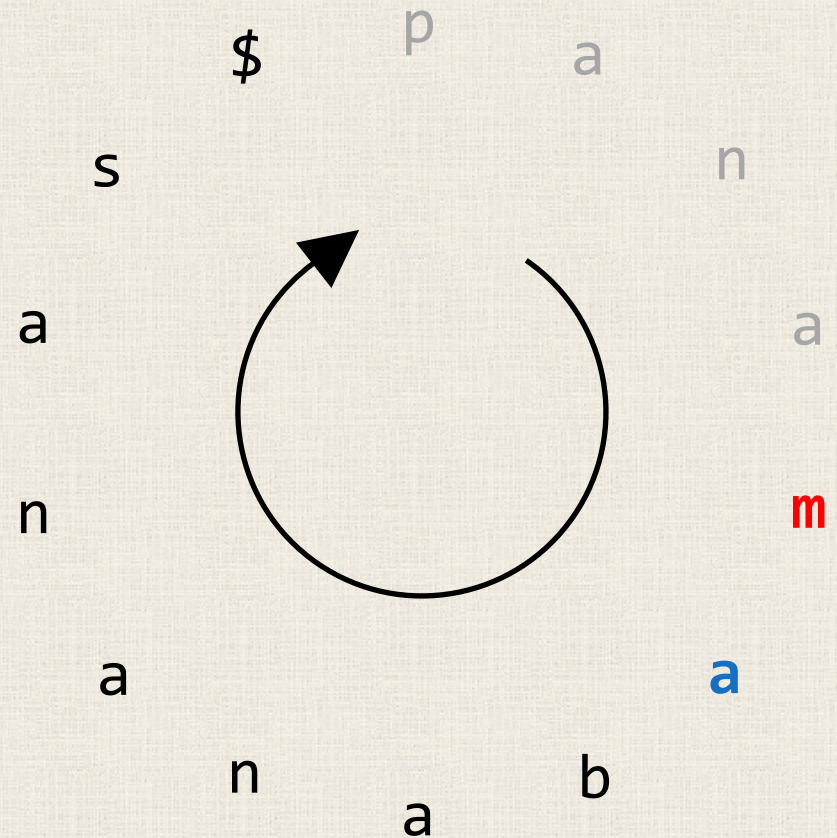
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



Efficient BWT Decompression

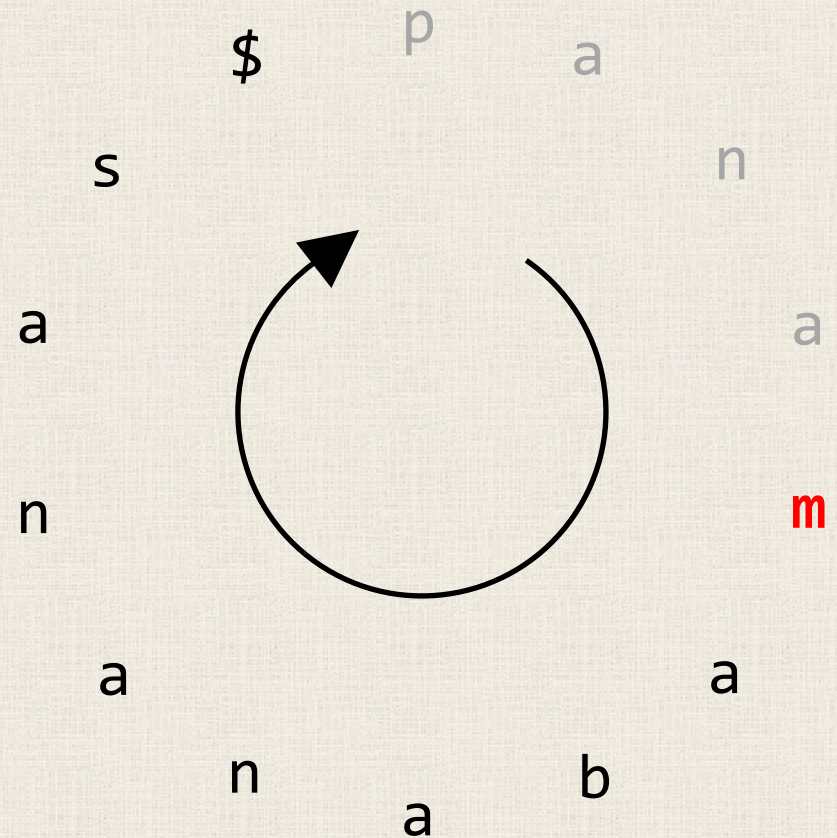
\$ ₁	p	a	n	a	m	a	b	a	n	a	s	\$	s ₁	
a ₁	b	a	n	a	n	a	s	\$	p	a	n	a	m ₁	
a ₂	m	a	b	a	n	a	n	a	s	\$	p	a	n ₁	
a ₃	n	a	m	a	b	a	n	a	n	a	s	\$	p ₁	
a ₄	n	a	n	a	s	\$	p	a	n	a	m	a	b ₁	
a ₅	n	a	s	\$	p	a	n	a	m	a	b	a	n ₂	
a ₆	s	\$	p	a	n	a	m	a	b	a	n	a	n ₃	
b ₁	a	n	a	n	a	s	\$	p	a	n	a	m	a ₁	
m ₁	a	b	a	n	a	n	a	s	\$	p	a	n	a ₂	
n ₁	a	m	a	b	a	n	a	n	a	s	\$	p	a ₃	
n ₂	a	n	a	s	\$	p	a	n	a	m	a	b	a ₄	
n ₃	a	s	\$	p	a	n	a	m	a	b	a	n	a ₅	
p ₁	a	n	a	m	a	b	a	n	a	n	a	s	\$	1
s ₁	\$	p	a	n	a	m	a	b	a	n	a	n	a	6



Efficient BWT Decompression

```

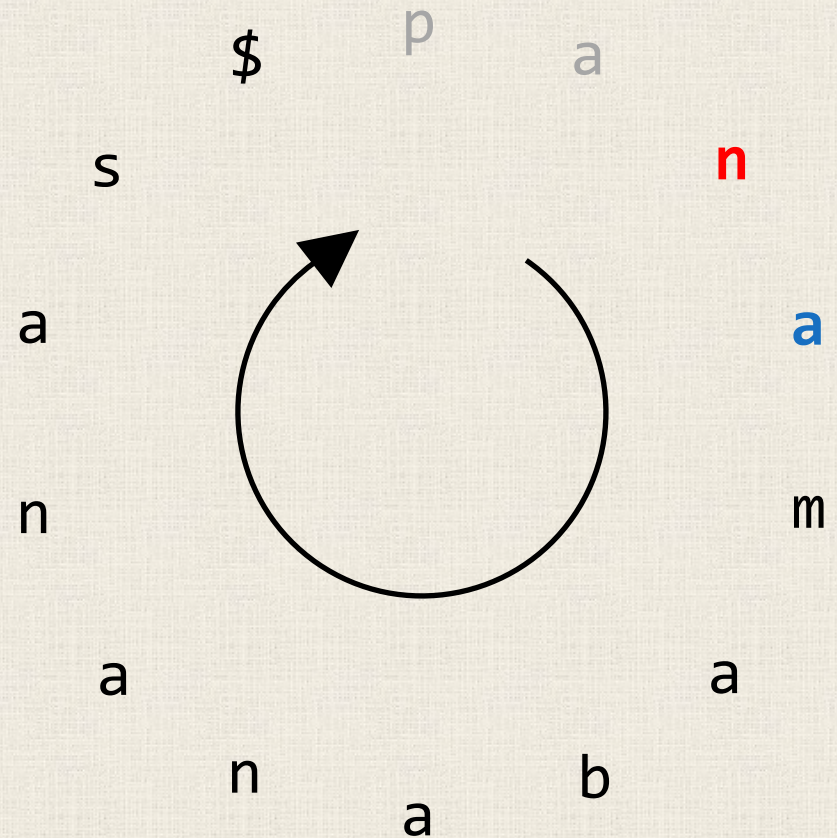
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
    
```



Efficient BWT Decompression

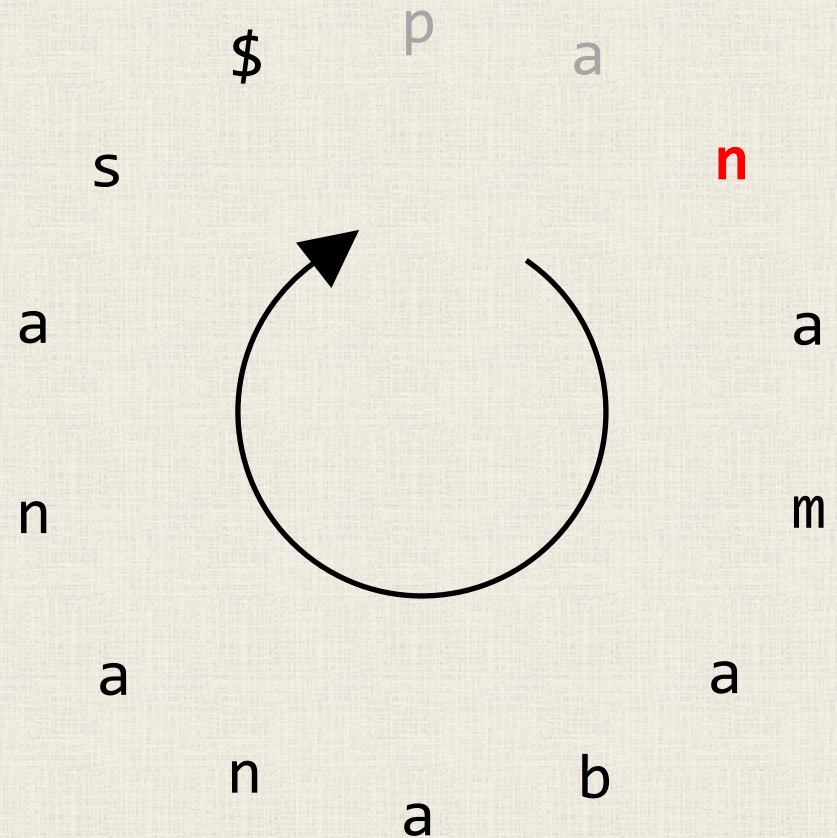
```

$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
    
```



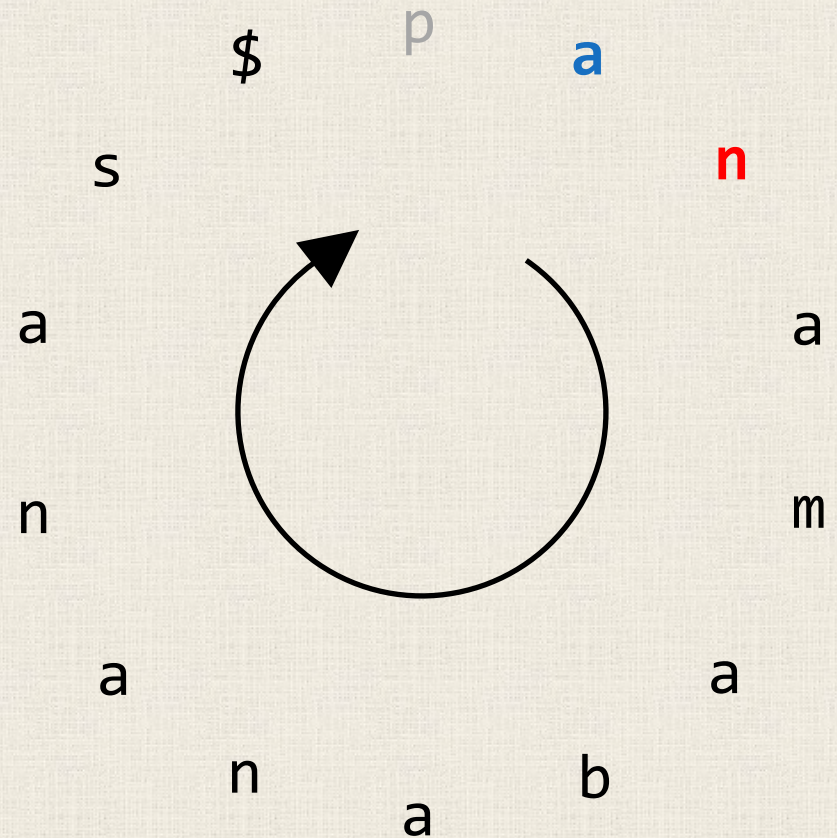
Efficient BWT Decompression

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1n1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



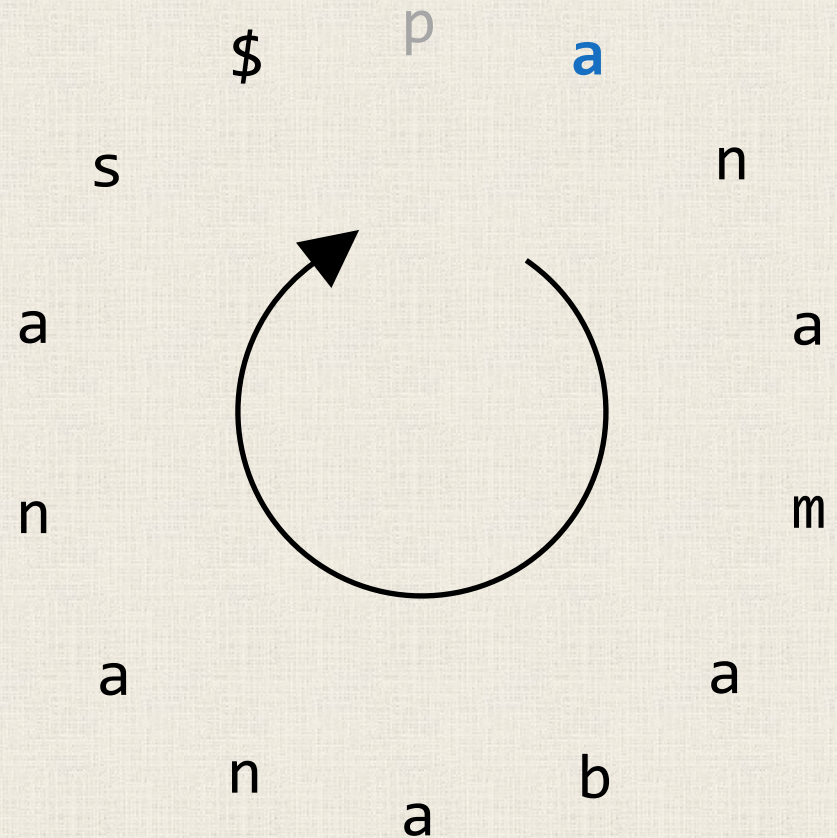
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



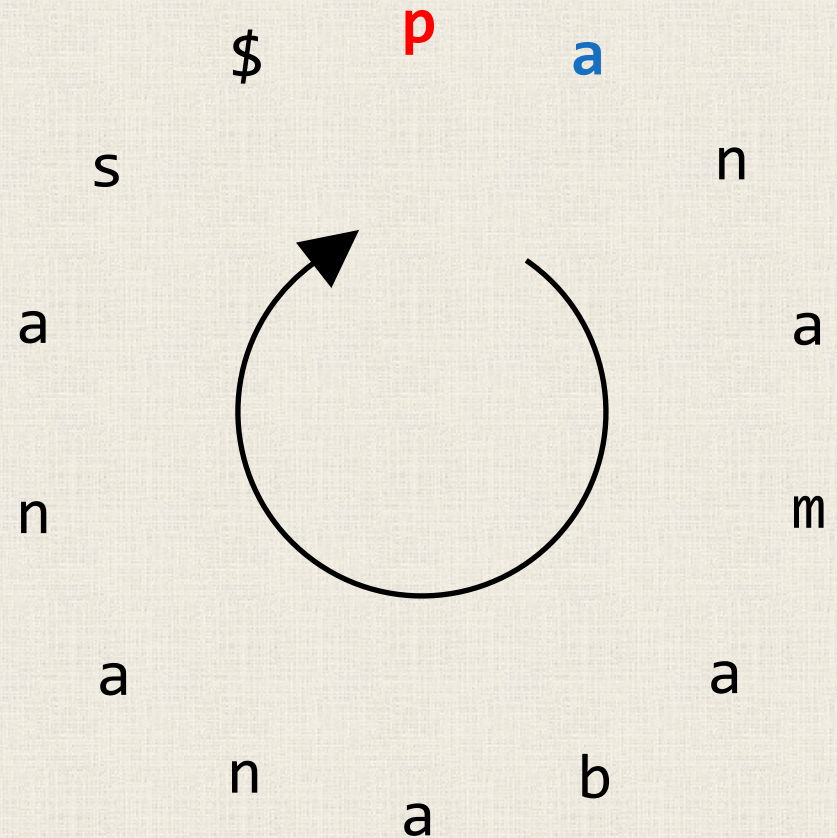
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



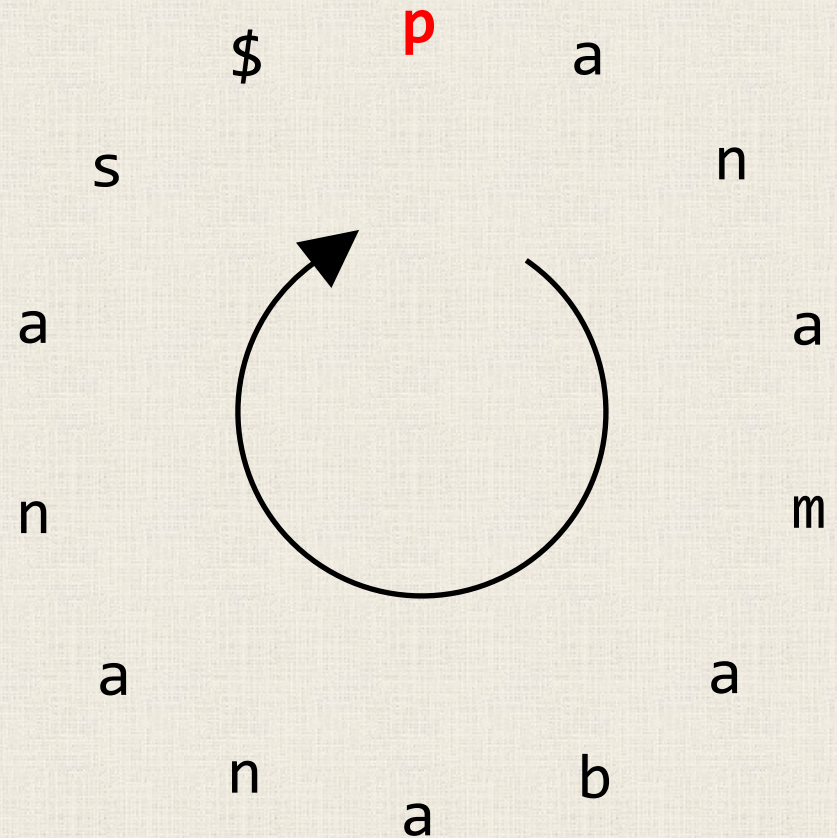
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



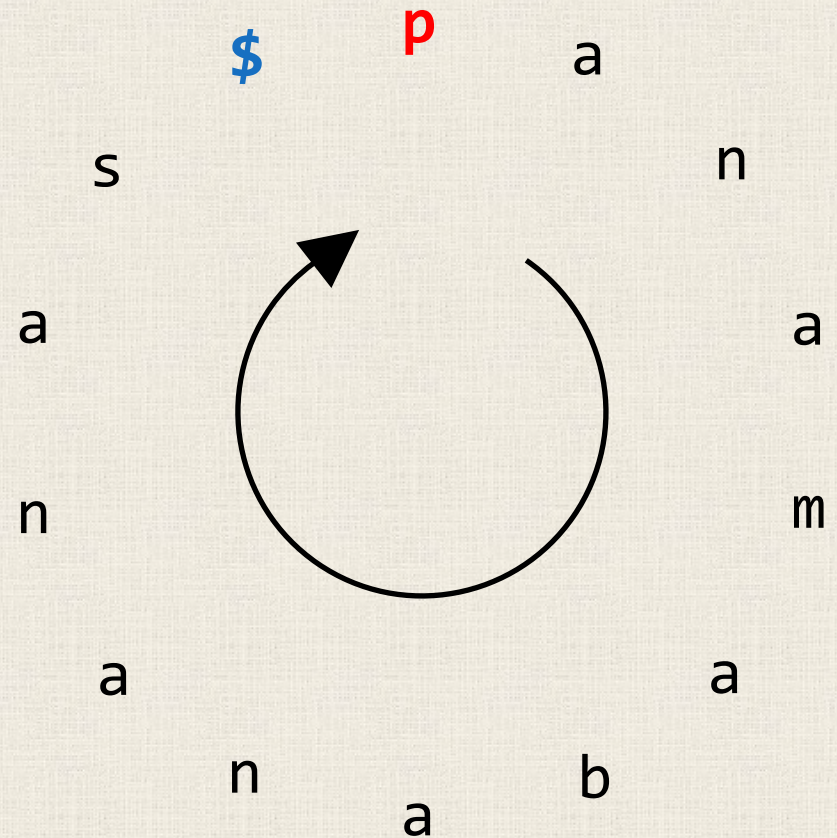
Efficient BWT Decompression

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



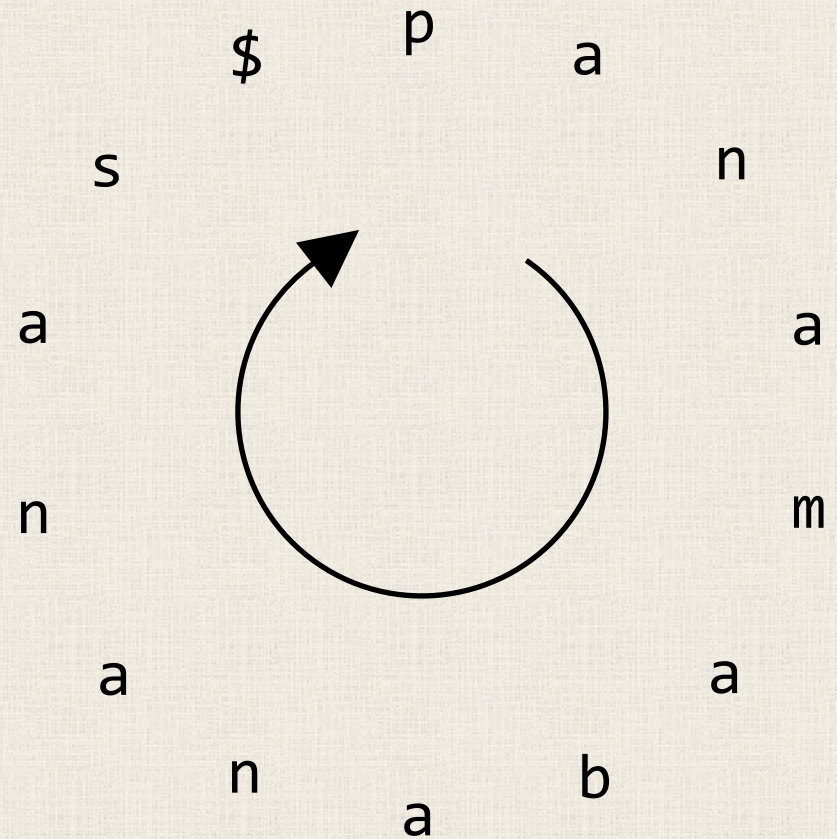
Efficient BWT Decompression

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```



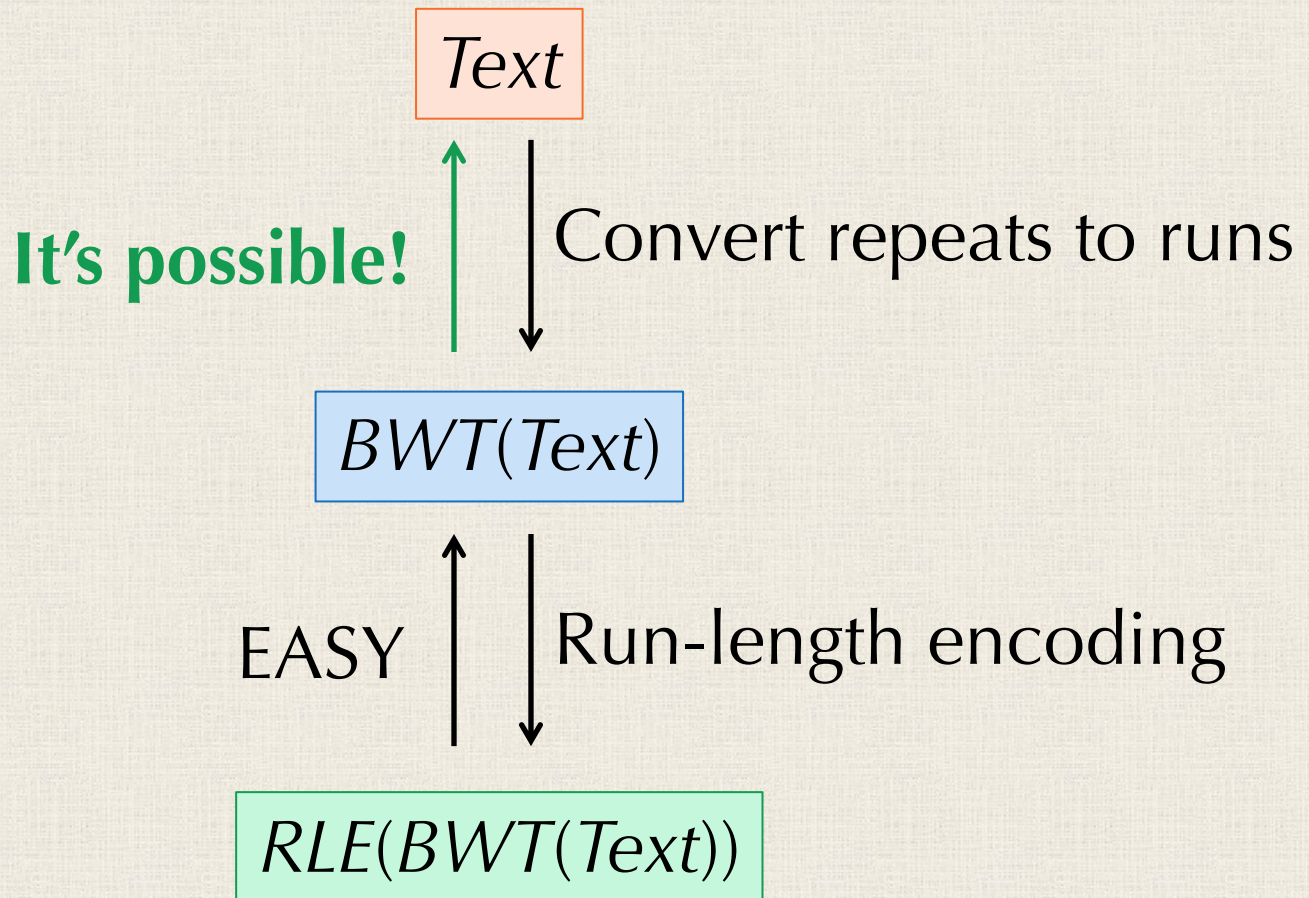
Efficient BWT Decompression

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$p3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabanana6
```



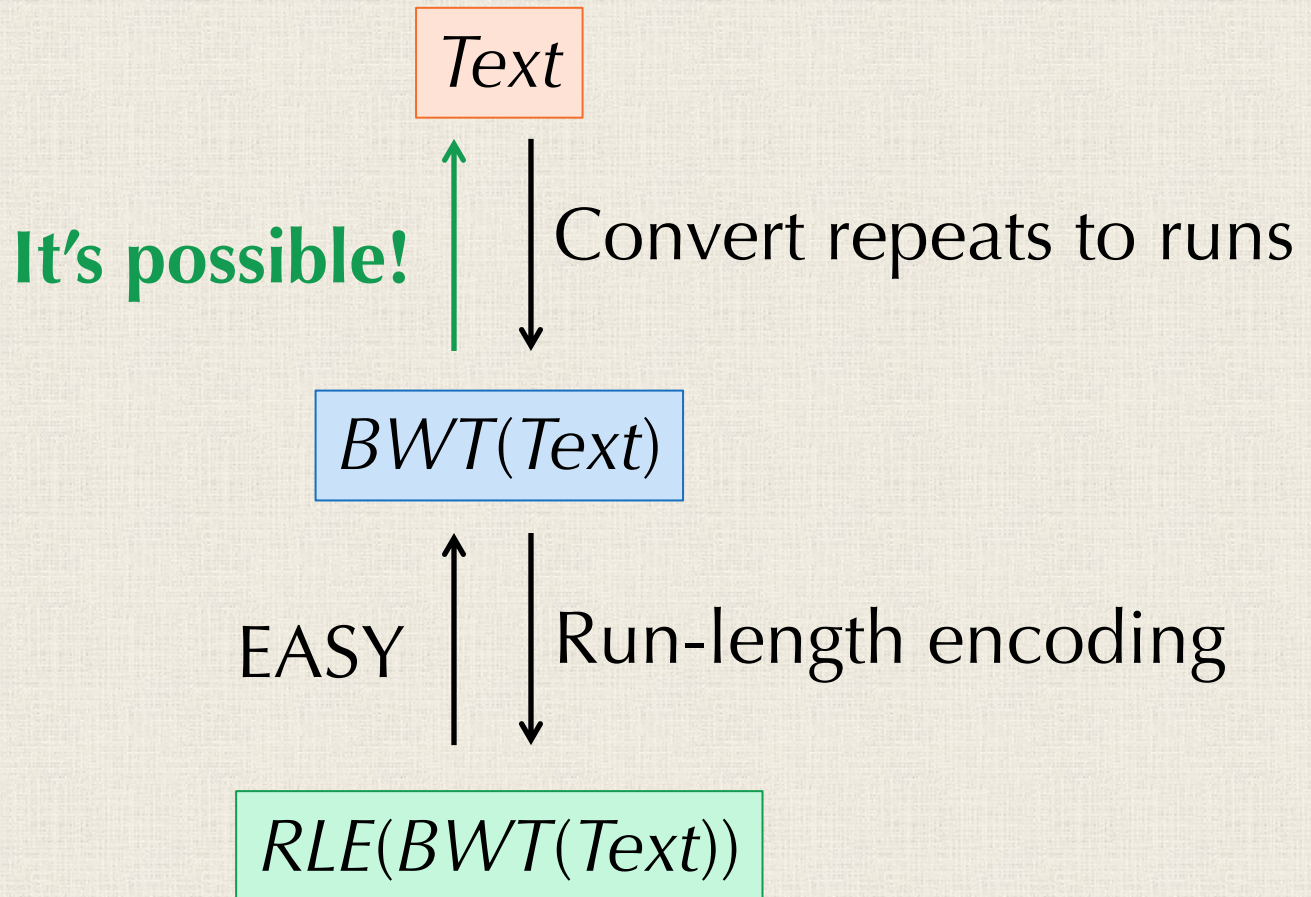
Memory: $2|Text| = O(|Text|)$.

So We Can Decompress BWT ...



So We Can Decompress BWT ...

But What About Pattern Matching?



PATTERN MATCHING WITH BURROWS-WHEELER

Recalling Our Goal

Pattern Matching with Suffix Array:

- Runtime: $O(|Text| + |Patterns|)$
- Memory: $O(|Text|)$
- Problem: suffix tree takes $\sim 20 \times |Text|$ space, but suffix array takes $\sim 4 \times |Text|$ space.

Recalling Our Goal

Pattern Matching with Suffix Array:

- Runtime: $O(|Text| + |Patterns|)$
- Memory: $O(|Text|)$
- Problem: suffix tree takes $\sim 20 \times |Text|$ space, but suffix array takes $\sim 4 \times |Text|$ space.

Can we use $BWT(Text)$ as our data structure instead?

Pattern Matches “Clump” at Start of $M(\text{Text})$

$\$$ ₁ panamabananas₁
a₁ bananas\$panam₁
a₂ mabananas\$pan₁
a₃namabananas\$p₁
a₄nanas\$panamab₁
a₅nas\$panamaban₂
a₆s\$panamaban₃
b₁ananas\$panama₁
m₁abananas\$pana₂
n₁amabananas\$pa₃
n₂anas\$panamaba₄
n₃as\$panamabana₅
p₁anamabananas\$₁
s₁\$panamabanana₆

Connecting $M(\text{Text})$ to Suffix Array

$\$$ ₁ panamabananas₁
a₁ bananas\$panam₁
a₂ mabananas\$pan₁
a₃ **n**amabananas\$p₁
a₄ **n**anas\$panamab₁
a₅ **n**as\$panamaban₂
a₆ s\$panamaban₃
b₁ ananas\$panama₁
m₁ abananas\$pana₂
n₁ amabananas\$pa₃
n₂ anas\$panamaba₄
n₃ as\$panamabana₅
p₁ anamabananas\$₁
s₁ \$panamabanana₆

Sorted Suffixes	Suffix Array
\$	13
abananas	5
amabananas\$	3
anamabananas\$	1
ananas\$	7
anas\$	9
as\$	11
bananas\$	6
mabananas\$	4
namabananas\$	2
nanas\$	8
nas\$	10
panamabananas\$	0
s\$	12

Connecting $M(\text{Text})$ to Suffix Array

$\$$ ₁ panamabananas₁
a₁ bananas\$panam₁
a₂ mabananas\$pan₁
a₃ namabananas\$p₁
a₄ nanas\$panamab₁
a₅ nas\$panamaban₂
a₆ s\$panamaban₃
b₁ ananas\$panama₁
m₁ abananas\$pana₂
n₁ amabananas\$pa₃
n₂ anas\$panamaba₄
n₃ as\$panamabana₅
p₁ anamabananas\$₁
s₁ \$panamabanana₆

We could find pattern matches easily if we had all the suffixes, but we would need $O(|\text{Text}|^2)$ space...

Connecting $M(\text{Text})$ to Suffix Array

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panamab1
a5nas$panamaban2
a6s$panamaban3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$_1
s1$panamabanana6
```

We could find pattern matches easily if we had all the suffixes, but we would need $O(|\text{Text}|^2)$ space...

We are going to pattern match using just two columns of $M(\text{Text})$:
FirstColumn and
LastColumn

We Match Patterns *Backward*

\$ ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s ₁
a ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m ₁
a ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₁
a ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p ₁
a ₄	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b ₁
a ₅	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₂
a ₆	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₃
b ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₁
m ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₂
n ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₃
n ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₄
n ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₅
p ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$ ₁
s ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₆

Searching for **ana** in *Text* =
panamabananas

We Match Patterns *Backward*

\$ ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s ₁
a ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m ₁
a ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₁
a ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p ₁
a ₄	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b ₁
a ₅	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₂
a ₆	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₃
b ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₁
m ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₂
n ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₃
n ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₄
n ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₅
p ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$ ₁
s ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₆

Searching for **ana** in *Text* =
panamabananas

We Match Patterns *Backward*

\$ ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s ₁
a ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m ₁
a ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₁
a ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p ₁
a ₄	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b ₁
a ₅	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₂
a ₆	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₃
b ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₁
m ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₂
n ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₃
n ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₄
n ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₅
p ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$ ₁
s ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₆

Searching for **ana** in *Text* =
panamabananas

We Match Patterns *Backward*

\$ ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s ₁
a ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m ₁
a ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₁
a ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p ₁
a ₄	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b ₁
a ₅	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₂
a ₆	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₃
b ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₁
m ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₂
n ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₃
n ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₄
n ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₅
p ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$ ₁
s ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₆

Searching for **ana** in *Text* =
panamabananas

We Match Patterns *Backward*

\$ ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s ₁
a ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m ₁
a ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₁
a ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p ₁
a ₄	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b ₁
a ₅	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₂
a ₆	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₃
b ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₁
m ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₂
n ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₃
n ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₄
n ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₅
p ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$ ₁
s ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₆

Searching for **ana** in *Text* = panamabananas

Now we can apply the First-Last Property and find where these three “n” are hiding in *FirstColumn*.

We Match Patterns *Backward*

\$ ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s ₁
a ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m ₁
a ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₁
a ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p ₁
a ₄	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b ₁
a ₅	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₂
a ₆	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₃
b ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₁
m ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₂
n ₁	a	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₃
n ₂	a	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₄
n ₃	a	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₅
p ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$ ₁
s ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₆

Searching for **ana** in *Text* =
panamabananas

We Match Patterns *Backward*

\$ ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s ₁
a ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m ₁
a ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₁
a ₃	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p ₁
a ₄	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b ₁
a ₅	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₂
a ₆	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₃
b ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₁
m ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₂
n ₁	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₃
n ₂	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₄
n ₃	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₅
p ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$ ₁
s ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₆

Searching for **ana** in *Text* = panamabananas

All three match, and again we apply the First-Last Property.

We Match Patterns *Backward*

\$ ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s ₁
a ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m ₁
a ₂	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₁
a ₃	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p ₁
a ₄	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b ₁
a ₅	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₂
a ₆	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n ₃
b ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₁
m ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₂
n ₁	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₃
n ₂	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₄
n ₃	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₅
p ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$ ₁
s ₁	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a ₆

Searching for **ana** in *Text* = panamabananas

All three match, and again we apply the First-Last Property.

We have found the occurrences of "ana"!

Matching “ana” backward with count arrays

\$?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	a
na	?	?	?	?	?	?	?	?	?	?	?	?	?	a
na	?	?	?	?	?	?	?	?	?	?	?	?	?	a
na	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	a

\$	a	b	m	n	p	s
0	0	0	0	0	0	1
0	0	0	1	0	0	1
0	0	0	1	1	0	1
0	0	0	1	1	1	1
0	0	1	1	1	1	1
0	0	1	1	2	1	1
0	0	1	1	3	1	1
0	1	1	1	3	1	1
0	2	1	1	3	1	1
0	3	1	1	3	1	1
0	4	1	1	3	1	1
0	5	1	1	3	1	1
1	5	1	1	3	1	1
1	6	1	1	3	1	1

We repeat this process and see that we must be looking for a₃ through a₅.

Matching “ana” backward with count arrays

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
ana	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
ana	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
ana	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

\$	a	b	m	n	p	s
0	0	0	0	0	0	1
0	0	0	1	0	0	1
0	0	0	1	1	0	1
0	0	0	1	1	1	1
0	0	1	1	1	1	1
0	0	1	1	2	1	1
0	0	1	1	3	1	1
0	1	1	1	3	1	1
0	2	1	1	3	1	1
0	3	1	1	3	1	1
0	4	1	1	3	1	1
0	5	1	1	3	1	1
1	5	1	1	3	1	1
1	6	1	1	3	1	1

Once we find a_3 through a_5 , we have now found that there are three occurrences of “ana”.

But where are they in *text*?

**WHERE ARE THE MATCHED
PATTERNS?**

Where are the Matches?

Multiple Pattern Matching Problem: *Find all occurrences of a collection of patterns in a text.*

- **Input:** A string *Text* and a collection *Patterns* containing (shorter) strings.
- **Output:** All ***starting positions*** in *Text* where a string from *Patterns* appears as a substring.

Where are the Matches?

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamabanan
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana

Example: We know that **ana** occurs 3 times, but where?

The Suffix Array Holds the Key

```
$panamabananas  
abananas$panam  
amabananas$pan  
anamabananas$p  
ananas$panamab  
anas$panamaban  
as$panamabanan  
bananas$panama  
mabananas$pana  
namabananas$pa  
nanas$panamaba  
nas$panamabana  
panamabananas$  
s$panamabana
```


The Suffix Array Holds the Key

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamabanan
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabana

Suffix Array

13
5
3
1
7
9
11
6
4
2
8
10
0
12

The Suffix Array Holds the Key

\$panamabananas
abananas\$panam
amabananas\$pan
anamabananas\$p
ananas\$panamab
anas\$panamaban
as\$panamabanan
bananas\$panama
mabananas\$pana
namabananas\$pa
nanas\$panamaba
nas\$panamabana
panamabananas\$
s\$panamabanna

Suffix Array

13

5

3

1

7

9

11

6

4

2

8

10

0

12

“But we’ve already used the suffix array for pattern matching! How is this useful?”

BURROWS AND WHEELER SET UP CHECKPOINTS

First: Dealing with the First Column

1	\$?	?	?	?	?	?	?	?	?	?	?	?	s
	a	?	?	?	?	?	?	?	?	?	?	?	?	m
	a	?	?	?	?	?	?	?	?	?	?	?	?	n
	a	?	?	?	?	?	?	?	?	?	?	?	?	p
	a	?	?	?	?	?	?	?	?	?	?	?	?	b
	a	?	?	?	?	?	?	?	?	?	?	?	?	n
	a	?	?	?	?	?	?	?	?	?	?	?	?	n
7	b	?	?	?	?	?	?	?	?	?	?	?	?	a
8	m	?	?	?	?	?	?	?	?	?	?	?	?	a
9	n	?	?	?	?	?	?	?	?	?	?	?	?	a
	n	?	?	?	?	?	?	?	?	?	?	?	?	a
	n	?	?	?	?	?	?	?	?	?	?	?	?	a
12	p	?	?	?	?	?	?	?	?	?	?	?	?	\$
13	s	?	?	?	?	?	?	?	?	?	?	?	?	a

STOP: Is there any way of representing the information in the first column using less data?

Answer: We only need to store the *first occurrence* of each symbol.

First: Dealing with the First Column

	\$?	?	?	?	?	?	?	?	?	?	?	?	?	s
1	a	?	?	?	?	?	?	?	?	?	?	?	?	?	m
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	p
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	b
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
7	b	?	?	?	?	?	?	?	?	?	?	?	?	?	a
8	m	?	?	?	?	?	?	?	?	?	?	?	?	?	a
9	n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
	n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
	n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
12	p	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
13	s	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Note: we don't even need to store that the first symbol occurs at position 1 of the first column.

First: Dealing with the First Column

	\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
1	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
7	b	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
8	m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
9	n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
	n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
	n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
12	p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
13	s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Note: we don't even need to store that the first symbol occurs at position 1 of the first column.

STOP: For a genome, how many integers will we need to store to hold the first occurrence of each symbol?

First: Dealing with the First Column

	\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
1	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
7	b	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
8	m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
9	n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
	n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
	n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
12	p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
13	s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Note: we don't even need to store that the first symbol occurs at position 1 of the first column.

Answer: Only three! The position of the first occurrence of C, G, and T. (\$ occurs at position 0; A occurs at position 1.)

Second: A “Partial” Suffix Array

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Suffix Array	
13	
5	
3	
1	
7	
9	
11	
6	
4	
2	
8	
10	
0	
12	

Partial suffix array: only stores values that are divisible by K for some integer K .

$K = 5$ at left.

Second: A “Partial” Suffix Array

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Suffix Array	
13	
5	
3	
1	
7	
9	
11	
6	
4	
2	
8	
10	
0	
12	

We can make at most $K - 1$ additional backtrack steps to determine where our pattern match is.

Second: A “Partial” Suffix Array

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Suffix Array	
13	
5	
3	
1	
7	
9	
11	
6	
4	
2	
8	
10	
0	
12	

We can make at most $K - 1$ additional backtrack steps to determine where our pattern match is.

Second: A “Partial” Suffix Array

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Suffix Array	
13	
5	
3	
1	
7	
9	
11	
6	
4	
2	
8	
10	
0	
12	

We can make at most $K - 1$ additional backtrack steps to determine where our pattern match is.

Second: A “Partial” Suffix Array

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Suffix Array	
13	
5	
3	
1	
7	
9	
11	
6	
4	
2	
8	
10	
0	
12	

We can make at most $K - 1$ additional backtrack steps to determine where our pattern match is.

Second: A “Partial” Suffix Array

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	b	a	n	a	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	a	n	a	?	?	?	?	?	?	?	?	?	?	a	
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Suffix Array	
13	
5	
3	
1	
7	
9	
11	
6	
4	
2	
8	
10	
0	
12	

We can make at most $K - 1$ additional backtrack steps to determine where our pattern match is.

Second: A “Partial” Suffix Array

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	b	a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	a	n	a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

Suffix Array	
13	
5	
3	
1	
7	
9	
11	
6	
4	
2	
8	
10	
0	
12	

We have a match! We know “abana” occurs at position 5, and we took two steps back, so this “ana” occurs at position 7.

Third: "Checkpoint" Count Arrays

\$?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	a

\$	a	b	m	n	p	s
0	0	0	0	0	0	1
0	0	0	1	0	0	1
0	0	0	1	1	0	1
0	0	0	1	1	1	1
0	0	1	1	1	1	1
0	0	1	1	2	1	1
0	0	1	1	3	1	1
0	1	1	1	3	1	1
0	2	1	1	3	1	1
0	3	1	1	3	1	1
0	4	1	1	3	1	1
0	5	1	1	3	1	1
1	5	1	1	3	1	1
1	6	1	1	3	1	1

"Count" arrays help us know what we've encountered up to given row.

Third: "Checkpoint" Count Arrays

\$?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	a

\$	a	b	m	n	p	s
0	0	0	0	0	0	1
0	0	0	1	0	0	1
0	0	0	1	1	0	1
0	0	0	1	1	1	1
0	0	1	1	1	1	1
0	0	1	1	2	1	1
0	0	1	1	3	1	1
0	1	1	1	3	1	1
0	2	1	1	3	1	1
0	3	1	1	3	1	1
0	4	1	1	3	1	1
0	5	1	1	3	1	1
1	5	1	1	3	1	1
1	6	1	1	3	1	1

"Count" arrays help us know what we've encountered up to given row.

Checkpoint arrays: only store every C arrays.

Third: "Checkpoint" Count Arrays

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

\$	a	b	m	n	p	s
0	0	0	0	0	0	1
0	0	0	1	0	0	1
0	0	0	1	1	0	1
0	0	0	1	1	1	1
0	0	1	1	1	1	1
0	0	1	1	2	1	1
0	0	1	1	3	1	1
0	1	1	1	3	1	1
0	2	1	1	3	1	1
0	3	1	1	3	1	1
0	4	1	1	3	1	1
0	5	1	1	3	1	1
1	5	1	1	3	1	1
1	6	1	1	3	1	1

STOP: Say we are at position 13 and see "a". How can we know how many "a"s have occurred so far?

Third: "Checkpoint" Count Arrays

\$?	?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	?	a

\$	a	b	m	n	p	s
0	0	0	0	0	0	1
0	0	0	1	0	0	1
0	0	0	1	1	0	1
0	0	0	1	1	1	1
0	0	1	1	1	1	1
0	0	1	1	2	1	1
0	0	1	1	3	1	1
0	1	1	1	3	1	1
0	2	1	1	3	1	1
0	3	1	1	3	1	1
0	4	1	1	3	1	1
0	5	1	1	3	1	1
1	5	1	1	3	1	1
1	6	1	1	3	1	1

Answer: Go to nearest checkpoint array, which has counted 3 "a"s up to position 10, and then count three more (total of 6).

“How are K and C chosen in practice?”

\$?	?	?	?	?	?	?	?	?	?	?	?	?	s
a	?	?	?	?	?	?	?	?	?	?	?	?	?	m
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	p
a	?	?	?	?	?	?	?	?	?	?	?	?	?	b
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
a	?	?	?	?	?	?	?	?	?	?	?	?	?	n
b	?	?	?	?	?	?	?	?	?	?	?	?	?	a
m	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
n	?	?	?	?	?	?	?	?	?	?	?	?	?	a
p	?	?	?	?	?	?	?	?	?	?	?	?	?	\$
s	?	?	?	?	?	?	?	?	?	?	?	?	?	a

\$	a	b	m	n	p	s
0	0	0	0	0	0	1
0	0	0	1	0	0	1
0	0	0	1	1	0	1
0	0	0	1	1	1	1
0	0	1	1	1	1	1
0	0	1	1	2	1	1
0	0	1	1	3	1	1
0	1	1	1	3	1	1
0	2	1	1	3	1	1
0	3	1	1	3	1	1
0	4	1	1	3	1	1
0	5	1	1	3	1	1
1	5	1	1	3	1	1
1	6	1	1	3	1	1

(Partial)
Suffix Array

13
5
3
1
7
9
11
6
4
2
8
10
0
12

So, We Need Just Four Items ... Total Memory ($K = C = 100$): $\sim 1.5|Text|!$

First Occurrence

BWT

Checkpoint Arrays

(Partial) Suffix Array

	s
1	m
	n
	p
	b
	n
	n
7	a
8	a
9	a
	a
	a
12	\$
13	a

\$	a	b	m	n	p	s
0	0	0	0	0	0	1
0	0	0	1	0	0	1
0	0	0	1	1	0	1
0	0	0	1	1	1	1
0	0	1	1	1	1	1
0	0	1	1	2	1	1
0	0	1	1	3	1	1
0	1	1	1	3	1	1
0	2	1	1	3	1	1
0	3	1	1	3	1	1
0	4	1	1	3	1	1
0	5	1	1	3	1	1
1	5	1	1	3	1	1
1	6	1	1	3	1	1

13
5
3
1
7
9
11
6
4
2
8
10
0
12

EPILOGUE: MISMATCH- TOLERANT READ MAPPING

Returning to Our Original Problem

Multiple Pattern Matching Problem: *Find all occurrences of a collection of patterns in a text.*

- **Input:** A string *Text* and a collection *Patterns* containing (shorter) strings.
- **Output:** All starting positions in *Text* where a string from *Patterns* appears as a substring.

Returning to Our Original Problem

Multiple *Approximate* Pattern Matching Problem:

Find all occurrences of a collection of patterns in a text.

- **Input:** A string *Text*, a collection *Patterns* containing (shorter) strings, **and an integer d .**
- **Output:** All starting positions in *Text* where a string from *Patterns* appears as a substring **with at most d mismatches.**

Method 1: Seeding

Say that *Pattern* appears in *Text* with just one mismatch.

Pattern

a c t t g g c t

Text

...g g c a c a c t a g g c t c c...

Method 1: Seeding

Say that *Pattern* appears in *Text* with just one mismatch.

<i>Pattern</i>	a c t t g g c t
<i>Text</i>	... g g c a c a c t a g g c t c c ...

Method 1: Seeding

Say that *Pattern* appears in *Text* with just one mismatch.

<i>Pattern</i>	a c t t g g c t
<i>Text</i>	... g g c a c a c t a g g c t c c ...

If we divide the strings in half, then one must match exactly!

Method 1: Seeding

Let's take another example of strings that have 3 mismatches.

Pattern

act**t**aggctcgggataa**t**cc

Text

...ggcacact**a**ag**t**ctcgggataa**g**cccc...

Method 1: Seeding

Let's take another example of strings that have 3 mismatches.

Pattern a c t **t** a | g **g** c t c | g g g a **t** | a a **t** c c

Text ... g g c a c a c t **a** a g **t** c t c | g g g a **t** | a a **g** c c c c ...

Now we can divide strings into four equal pieces and find at least one that matches.

Method 1: Seeding

Theorem: If *Pattern* occurs in *Text* with d mismatches, then we can divide *Pattern* into $d + 1$ “equal” pieces and find at least one exact match.



https://en.wikipedia.org/wiki/Pigeonhole_principle#/media/File:TooManyPigeons.jpg

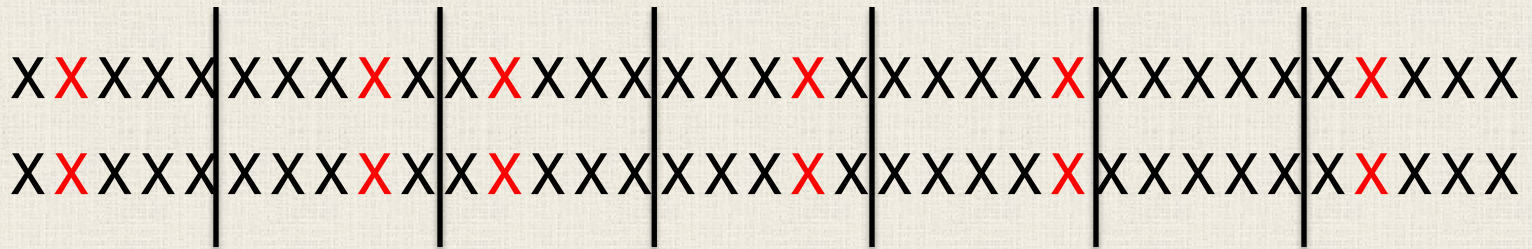
Method 1: Seeding

Theorem: If *Pattern* occurs in *Text* with d mismatches, then we can divide *Pattern* into $d + 1$ “equal” pieces and find at least one exact match.

XX

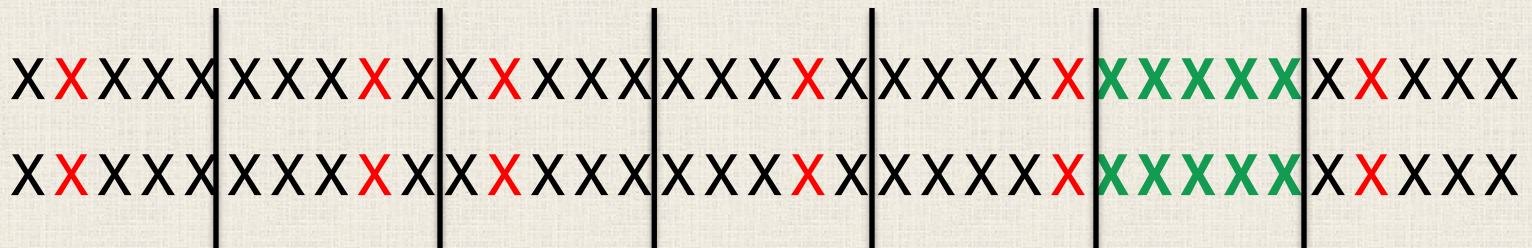
Method 1: Seeding

Theorem: If *Pattern* occurs in *Text* with d mismatches, then we can divide *Pattern* into $d + 1$ “equal” pieces and find at least one exact match.



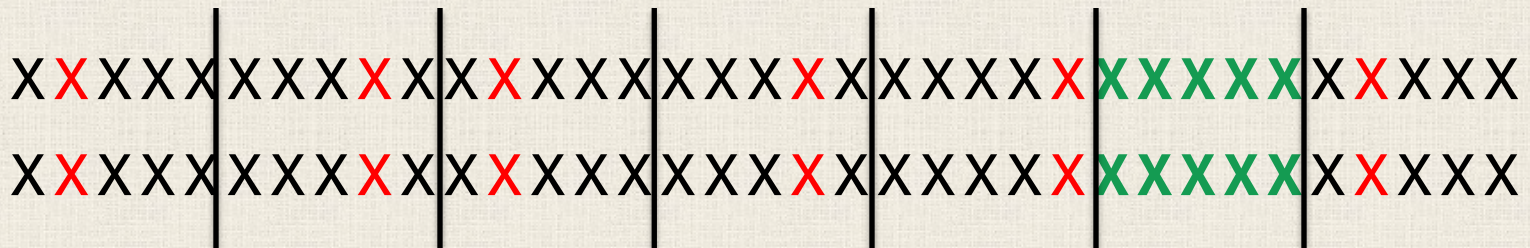
Method 1: Seeding

Theorem: If *Pattern* occurs in *Text* with d mismatches, then we can divide *Pattern* into $d + 1$ “equal” pieces and find at least one exact match.



Method 1: Seeding

Theorem (Formally): If two strings of length n match with at most d mismatches, then they must share a k -mer of length $k = \lfloor n/(d+1) \rfloor$.



Method 1: Seeding

Seed-and-extend for inexact pattern matching:

1. Divide *Pattern* into $d+1$ "equal" segments (called **seeds**).
2. Find which seeds match *Text* exactly, using an approach like BWT (**seed detection**).
3. Extend all seeds in both directions and align *Pattern* against this region of *Text* (**seed extension**) – can be generalized to include gaps.

Method 2: BWT Saves the Day Again

Recall: searching for ana in panamabananas

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```


Method 2: BWT Saves the Day Again

Recall: searching for **ana** in panamabananas

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panamab1  
a5nas$panamaban2  
a6s$panamaban3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabanana6
```

Method 2: BWT Saves the Day Again

Recall: searching for **ana** in panamabananas

```
$1panamabananas1  
a1bananas$panam1  
a2mabananas$pan1  
a3namabananas$p1  
a4nanas$panama1  
a5nas$panamaban2  
a6s$panamabana3  
b1ananas$panama1  
m1abananas$pana2  
n1amabananas$pa3  
n2anas$panamaba4  
n3as$panamabana5  
p1anamabananas$1  
s1$panamabana6
```

Method 2: BWT Saves the Day Again

Recall: searching for **ana** in panamabananas

If we allow 1 mismatch, then we need to keep the red letters around.

```
$1panamabananas1
a1bananas$panam1
a2mabananas$pan1
a3namabananas$p1
a4nanas$panama1
a5nas$panamaba2
a6s$panamabana3
b1ananas$panama1
m1abananas$pana2
n1amabananas$pa3
n2anas$panamaba4
n3as$panamabana5
p1anamabananas$1
s1$panamabana6
```


Method 2: BWT Saves the Day Again

Recall: searching for **ana** in panamabananas

If we allow 1 mismatch, then we need to keep the red letters around.

	# Mismatches
$\$$ ₁ panamabananas ₁	
a ₁ bananas\$pana m ₁	1
a ₂ mabananas\$pa n ₁	0
a ₃ namabananas\$ p ₁	1
a ₄ nanas\$panama b ₁	1
a ₅ nas\$panamaba n ₂	0
a ₆ s\$panamabana n ₃	0
b ₁ ananas\$panama ₁	
m ₁ abananas\$pana ₂	
n ₁ amabananas\$pa ₃	
n ₂ anas\$panamaba ₄	
n ₃ as\$panamabana ₅	
p ₁ anamabananas\$ ₁	
s ₁ \$panamabana ₆	

Method 2: BWT Saves the Day Again

Recall: searching for **ana** in panamabananas

Now we extend only strings that have at most 1 mismatch.

	# Mismatches
$\$$ ₁ panamabananas _s ₁	
a ₁ bananas\$pana m ₁	1
a ₂ mabananas\$pa n ₁	0
a ₃ namabananas\$ p ₁	1
a ₄ nanas\$panama b ₁	1
a ₅ nas\$panamaba n ₂	0
a ₆ s\$panamabana n ₃	0
b ₁ a _n anas\$panama _a ₁	
m ₁ a _b anas\$pana _a ₂	
n ₁ a _m abananas\$pa _a ₃	
n ₂ a _n as\$panamaba _a ₄	
n ₃ a _s \$panamabana _a ₅	
p ₁ a _n amabananas\$ _s ₁	
s ₁ \$panamabana _a ₆	

Method 2: BWT Saves the Day Again

Recall: searching for **ana** in panamabananas

Now we extend only strings that have at most 1 mismatch.

	# Mismatches
$\$$ ₁ panamabananas _s ₁	
a ₁ bananas\$panam ₁	
a ₂ mabananas\$pan ₁	
a ₃ namabananas\$p ₁	
a ₄ nanas\$panamab ₁	
a ₅ nas\$panamaban ₂	
a ₆ s\$panamaban ₃	
b ₁ a nanas\$panama ₁	1
m ₁ a bananas\$pana ₂	1
n ₁ a mabananas\$pa ₃	0
n ₂ a nas\$panamaba ₄	0
n ₃ a s\$panamabana ₅	0
p ₁ a namabananas\$ ₁	1
s ₁ \$panamabana _a ₆	

Method 2: BWT Saves the Day Again

Recall: searching for **ana** in panamabananas

One string produces a second mismatch (the \$), so we discard it.

	# Mismatches
$\$$ ₁ panamabananas s ₁	
a ₁ bananas \$ panam ₁	
a ₂ mabananas \$ pan ₁	
a ₃ namabananas \$ p ₁	
a ₄ nanas \$ panamab ₁	
a ₅ nas \$ panamaban ₂	
a ₆ s \$ panamaban ₃	
b ₁ a nanas \$ panam a ₁	1
m ₁ a bananas \$ pan a ₂	1
n ₁ a mabananas \$ pa a ₃	0
n ₂ a nas \$ panamab a ₄	0
n ₃ a s \$ panamabana a ₅	0
p ₁ a namabananas \$ s ₁	2
s ₁ \$ panamabanan a ₆	

Method 2: BWT Saves the Day Again

Recall: searching for **ana** in panamabananas

In the end, we have five 3-mers with at most 1 mismatch.

	# Mismatches
$\$$ ₁ panamabananas ₁	
a ₁ b ananas\$panam ₁	
a ₂ m abananas\$pan ₁	
a ₃ n amabananas\$p ₁	
a ₄ n anas\$panamab ₁	
a ₅ n as\$panamaban ₂	
a ₆ s\$panamaban ₃	
b ₁ ananas\$panama ₁	1
m ₁ abananas\$pana ₂	1
n ₁ amabananas\$pa ₃	0
n ₂ anas\$panamaba ₄	0
n ₃ as\$panamabana ₅	0
p ₁ anamabananas\$ ₁	2
s ₁ \$panamabanana ₆	

Method 2: BWT Saves the Day Again

Note: this approach faces complications in practice. We will have to consider a huge number of strings at the beginning, most of which may be frivolous.

In practice, we may require that a suffix of length m matches the string exactly to reduce the number of strings we have to consider initially.

Citations

Compression

[\[PDF\] A block-sorting lossless data compression algorithm](#)

www.hpl.hp.com/techreports/Compaq-DEC/SRC-RR-124.pdf ▼

by AB Lossless - 1994 - [Related articles](#)

May 10, 1994 - **A Block-sorting Lossless. Data Compression Algorithm.** M. Burrows and D.J. Wheeler digital. Systems Research Center. 130 Lytton Avenue.

Pattern
Matching

[Opportunistic data structures with applications - IEEE Conference ...](#)

<https://ieeexplore.ieee.org/document/892127>

by P Ferragina - 2000 - [Cited by 1016](#) - [Related articles](#)

Opportunistic data structures with applications. Abstract: We address the issue of compressing and **indexing** data. We devise a data structure whose space occupancy is a function of the entropy of the underlying data set.

Read
Mapping

[Fast and accurate short read alignment with Burrows–Wheeler transform](#)

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2705234/> ▼

by H Li - 2009 - [Cited by 17996](#) - [Related articles](#)

May 18, 2009 - Results: We implemented Burrows-Wheeler Alignment tool (**BWA**), a new read alignment package that is based on backward search with ...

[Ultrafast and memory-efficient alignment of short DNA sequences to ...](#)

<https://genomebiology.biomedcentral.com/articles/10.1186/gb-2009-10-3-r25> ▼

by B Langmead - 2009 - [Cited by 13797](#) - [Related articles](#)

Mar 4, 2009 - **Bowtie** is an ultrafast, memory-efficient alignment program for aligning short DNA sequence reads to large genomes. For the human genome, ...