

# Advanced Algorithms – COMS31900

---

## Pattern Matching part one

### Suffix Trees

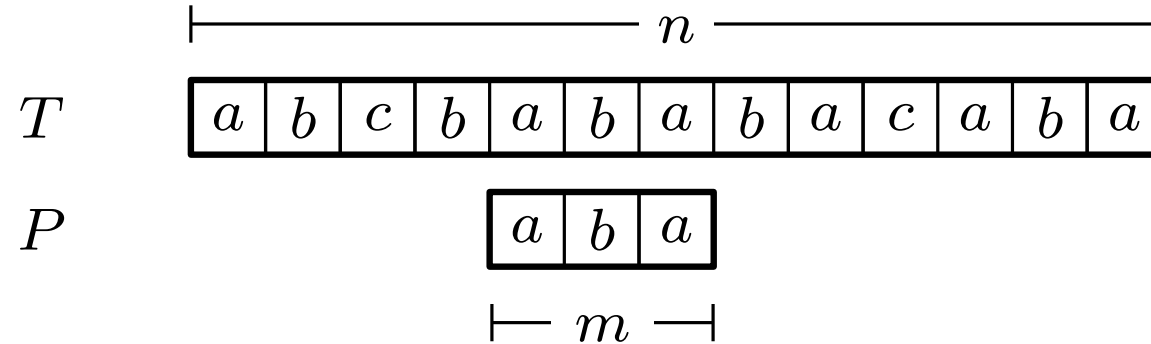
---

Raphaël Clifford

Slides by Benjamin Sach

# Exact pattern matching

**Input** A text string  $T$  (length  $n$ ) and a pattern string  $P$  (length  $m$ )



**Goal:** Find all the locations where  $P$  **matches** in  $T$

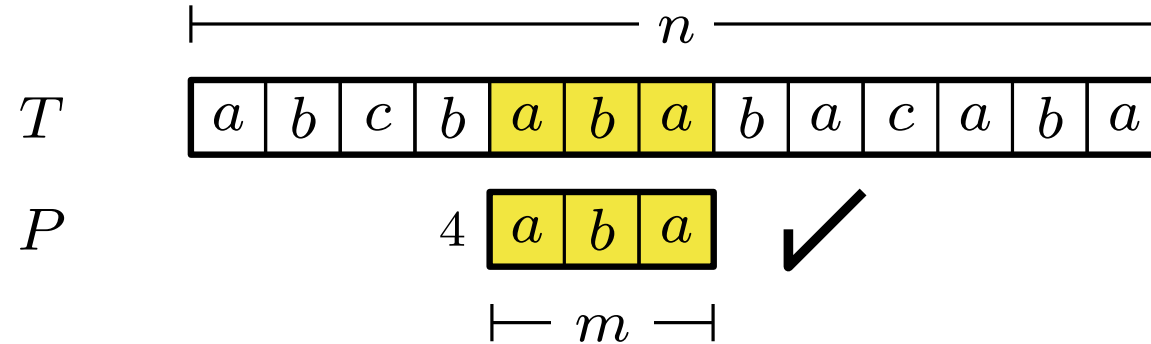
$P$  matches at location  $i$  iff

for all  $0 \leq j \leq m$  we have that  $P[j] = T[i + j]$

*(our strings are zero-indexed)*

# Exact pattern matching

**Input** A text string  $T$  (length  $n$ ) and a pattern string  $P$  (length  $m$ )



**Goal:** Find all the locations where  $P$  **matches** in  $T$

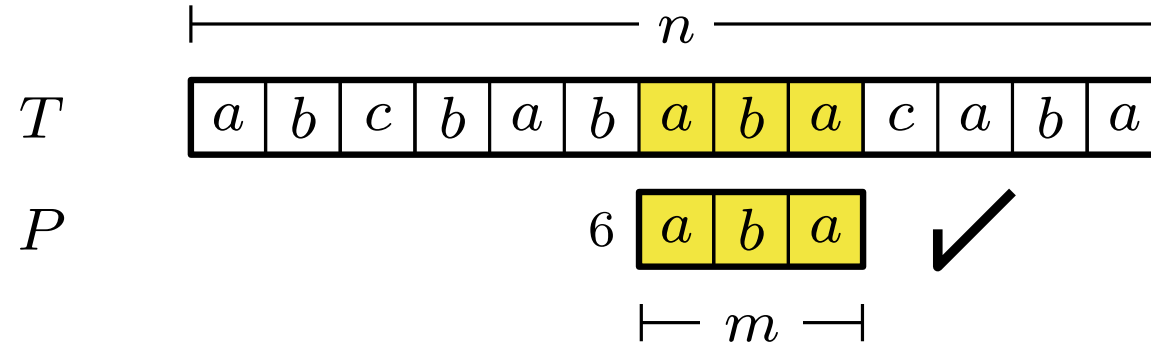
$P$  matches at location  $i$  iff

for all  $0 \leq j \leq m$  we have that  $P[j] = T[i + j]$

*(our strings are zero-indexed)*

# Exact pattern matching

**Input** A text string  $T$  (length  $n$ ) and a pattern string  $P$  (length  $m$ )



**Goal:** Find all the locations where  $P$  **matches** in  $T$

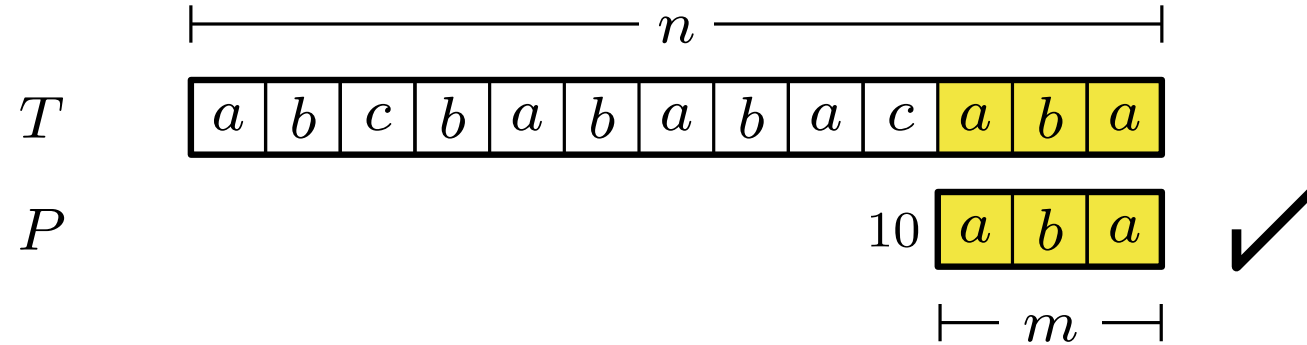
$P$  matches at location  $i$  iff

for all  $0 \leq j \leq m$  we have that  $P[j] = T[i + j]$

*(our strings are zero-indexed)*

# Exact pattern matching

**Input** A text string  $T$  (length  $n$ ) and a pattern string  $P$  (length  $m$ )



**Goal:** Find all the locations where  $P$  **matches** in  $T$

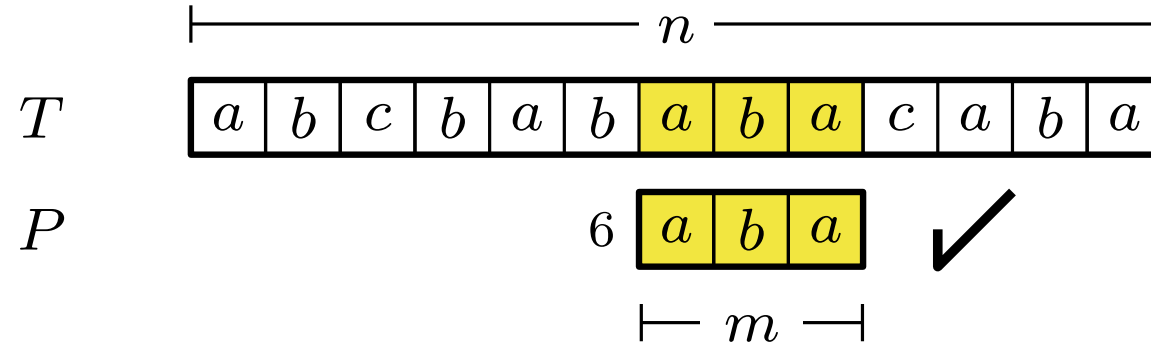
$P$  matches at location  $i$  iff

for all  $0 \leq j \leq m$  we have that  $P[j] = T[i + j]$

*(our strings are zero-indexed)*

# Exact pattern matching

**Input** A text string  $T$  (length  $n$ ) and a pattern string  $P$  (length  $m$ )



**Goal:** Find all the locations where  $P$  **matches** in  $T$

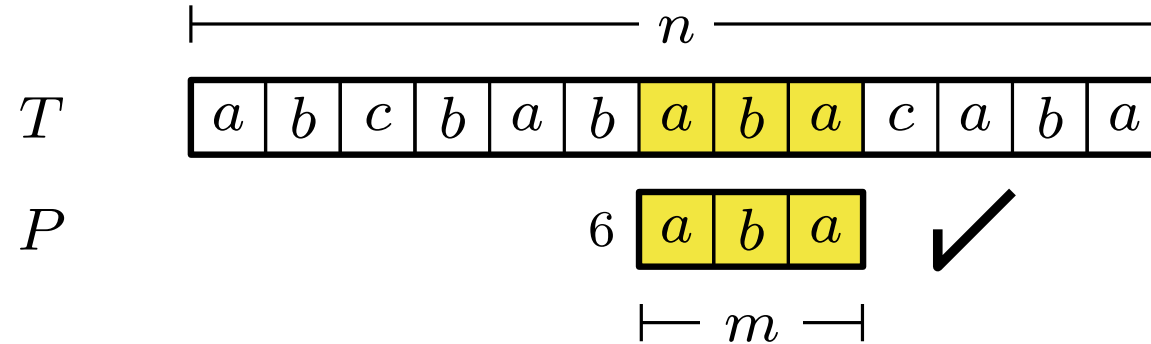
$P$  matches at location  $i$  iff

for all  $0 \leq j \leq m$  we have that  $P[j] = T[i + j]$

*(our strings are zero-indexed)*

# Exact pattern matching

**Input** A text string  $T$  (length  $n$ ) and a pattern string  $P$  (length  $m$ )



**Goal:** Find all the locations where  $P$  **matches** in  $T$

$P$  matches at location  $i$  iff

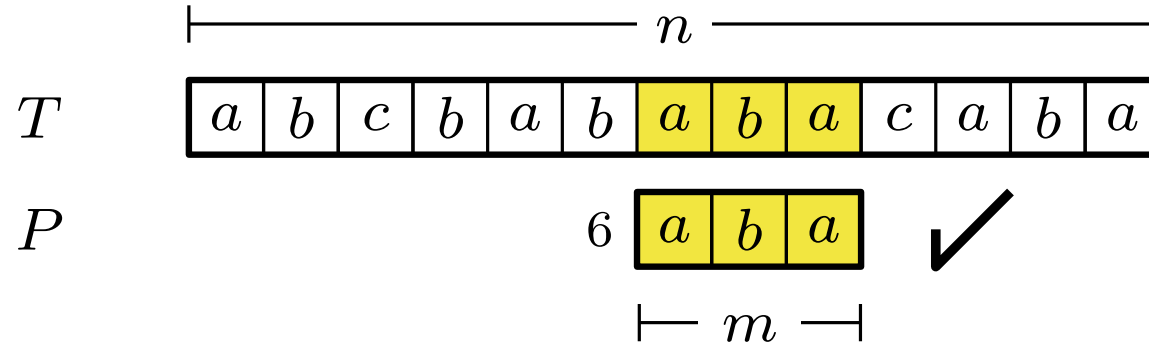
for all  $0 \leq j \leq m$  we have that  $P[j] = T[i + j]$

*(our strings are zero-indexed)*

- A naive algorithm takes  $O(nm)$  time

# Exact pattern matching

**Input** A text string  $T$  (length  $n$ ) and a pattern string  $P$  (length  $m$ )



**Goal:** Find all the locations where  $P$  **matches** in  $T$

$P$  matches at location  $i$  iff

for all  $0 \leq j \leq m$  we have that  $P[j] = T[i + j]$

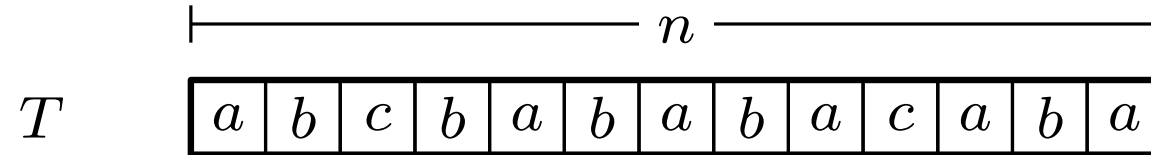
*(our strings are zero-indexed)*

- A naive algorithm takes  $O(nm)$  time
- Many  $O(n)$  time algorithms are known (for example KMP)



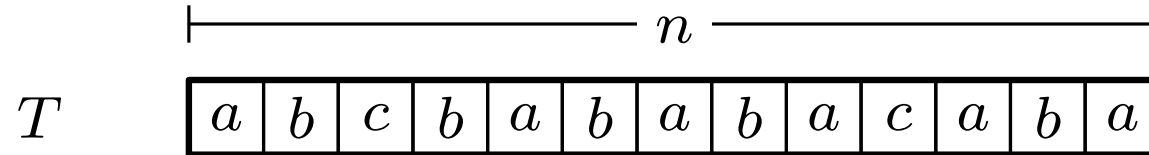
# Text indexing

Preprocess a text string  $T$  (length  $n$ ) to answer pattern matching queries...

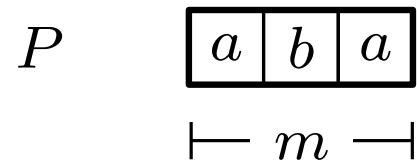


# Text indexing

Preprocess a text string  $T$  (length  $n$ ) to answer pattern matching queries...

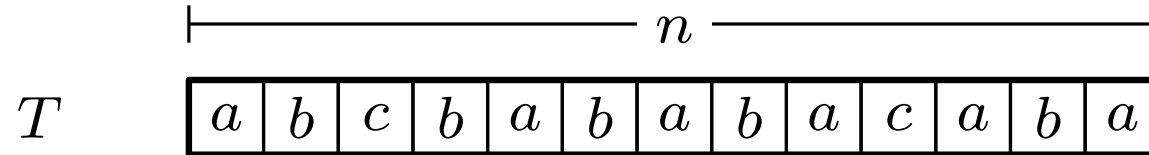


After preprocessing, a **query** is a pattern  $P$  (length  $m$ ),

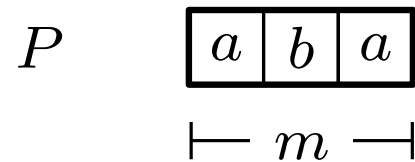


# Text indexing

Preprocess a text string  $T$  (length  $n$ ) to answer pattern matching queries...



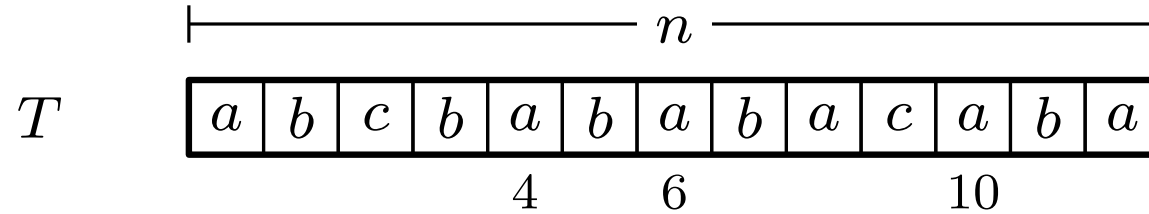
After preprocessing, a **query** is a pattern  $P$  (length  $m$ ),



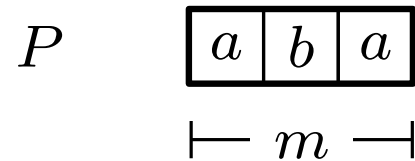
the output is a list of all matches in  $T$ .

# Text indexing

Preprocess a text string  $T$  (length  $n$ ) to answer pattern matching queries...



After preprocessing, a **query** is a pattern  $P$  (length  $m$ ),

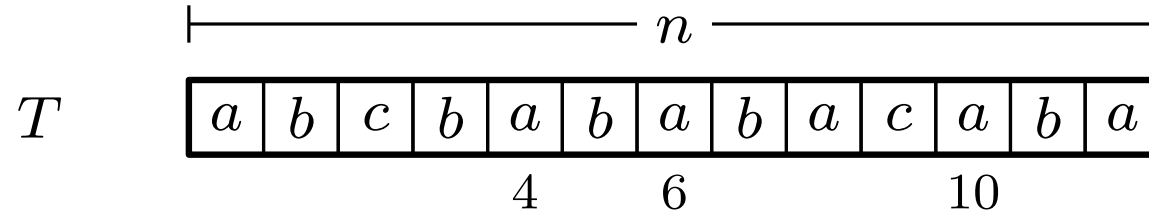


the output is a list of all matches in  $T$ .

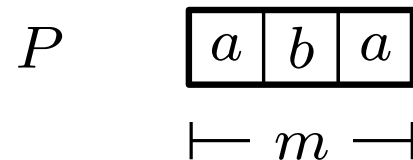
e.g. 4, 6, 10

# Text indexing

Preprocess a text string  $T$  (length  $n$ ) to answer pattern matching queries...



After preprocessing, a **query** is a pattern  $P$  (length  $m$ ),



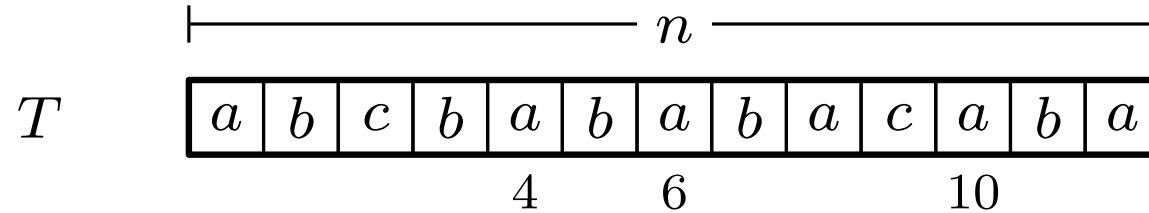
the output is a list of all matches in  $T$ .

e.g. 4, 6, 10

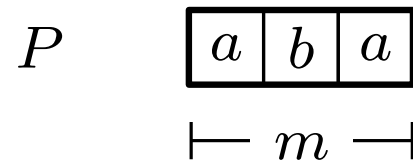
- A naive algorithm takes  $O(n)$  query time (using KMP)

# Text indexing

Preprocess a text string  $T$  (length  $n$ ) to answer pattern matching queries...



After preprocessing, a **query** is a pattern  $P$  (length  $m$ ),



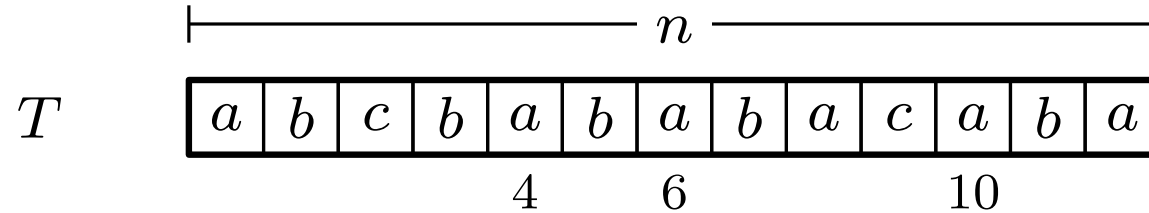
the output is a list of all matches in  $T$ .

e.g. 4, 6, 10

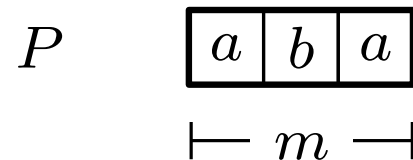
- A naive algorithm takes  $O(n)$  query time (using KMP)
- We want a query time which depends only on  $m$  and  $occ$ 
  - $occ$  is the number of occurrences (matches)

# Text indexing

Preprocess a text string  $T$  (length  $n$ ) to answer pattern matching queries...



After preprocessing, a **query** is a pattern  $P$  (length  $m$ ),



the output is a list of all matches in  $T$ .

e.g. 4, 6, 10

- A naive algorithm takes  $O(n)$  query time (using KMP)
- We want a query time which depends only on  $m$  and  $occ$ 
  - $occ$  is the number of occurrences (matches)
- We also want  $O(n)$  space and fast preprocessing (prep.) time

# The atomic suffix tree

$T$ 

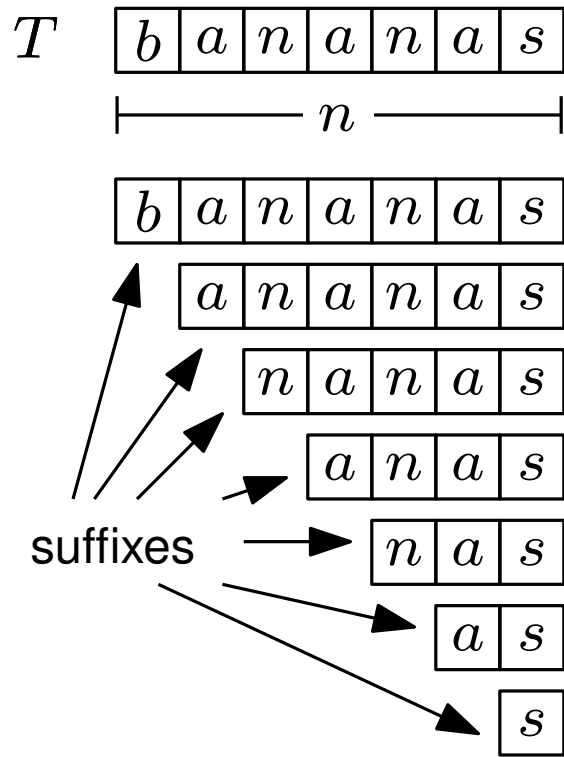
$b$	$a$	$n$	$a$	$n$	$a$	$s$
-----	-----	-----	-----	-----	-----	-----

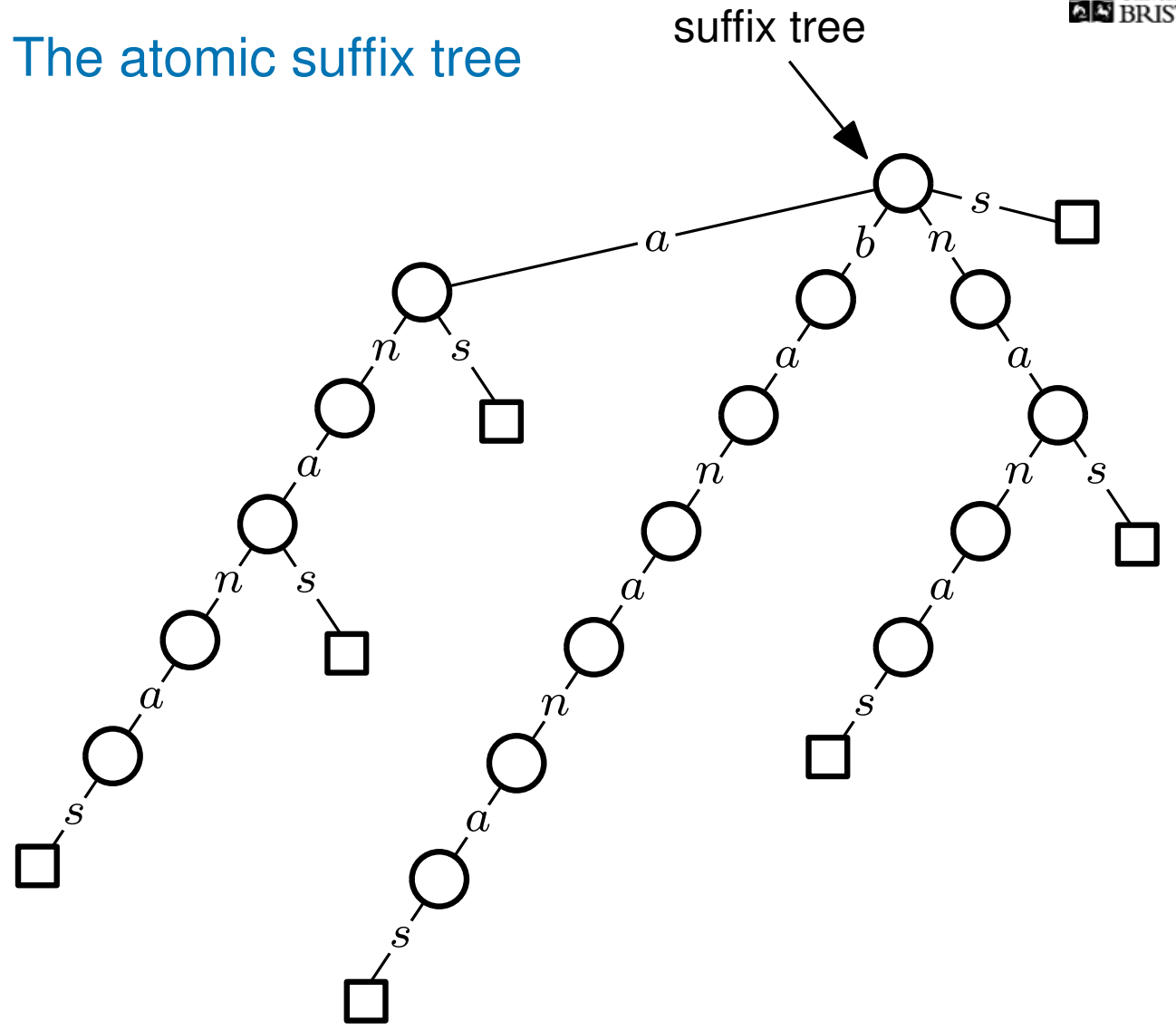
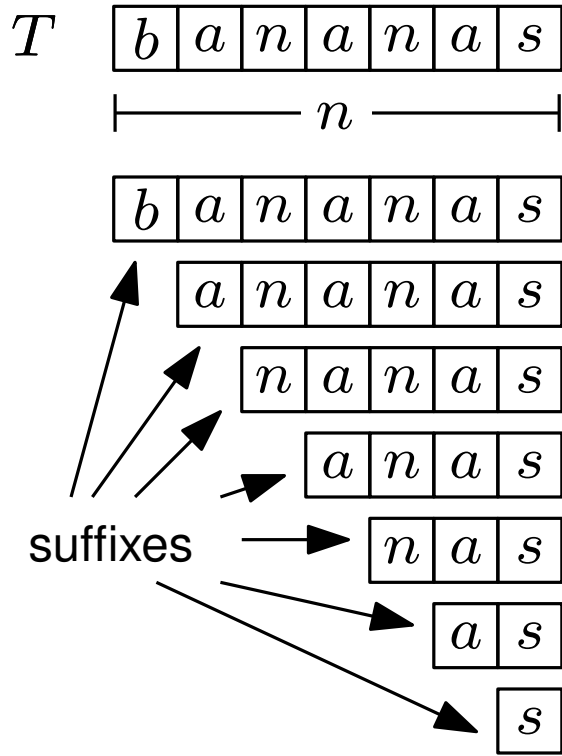
		—	$n$	—		
--	--	---	-----	---	--	--



# The atomic suffix tree



# The atomic suffix tree

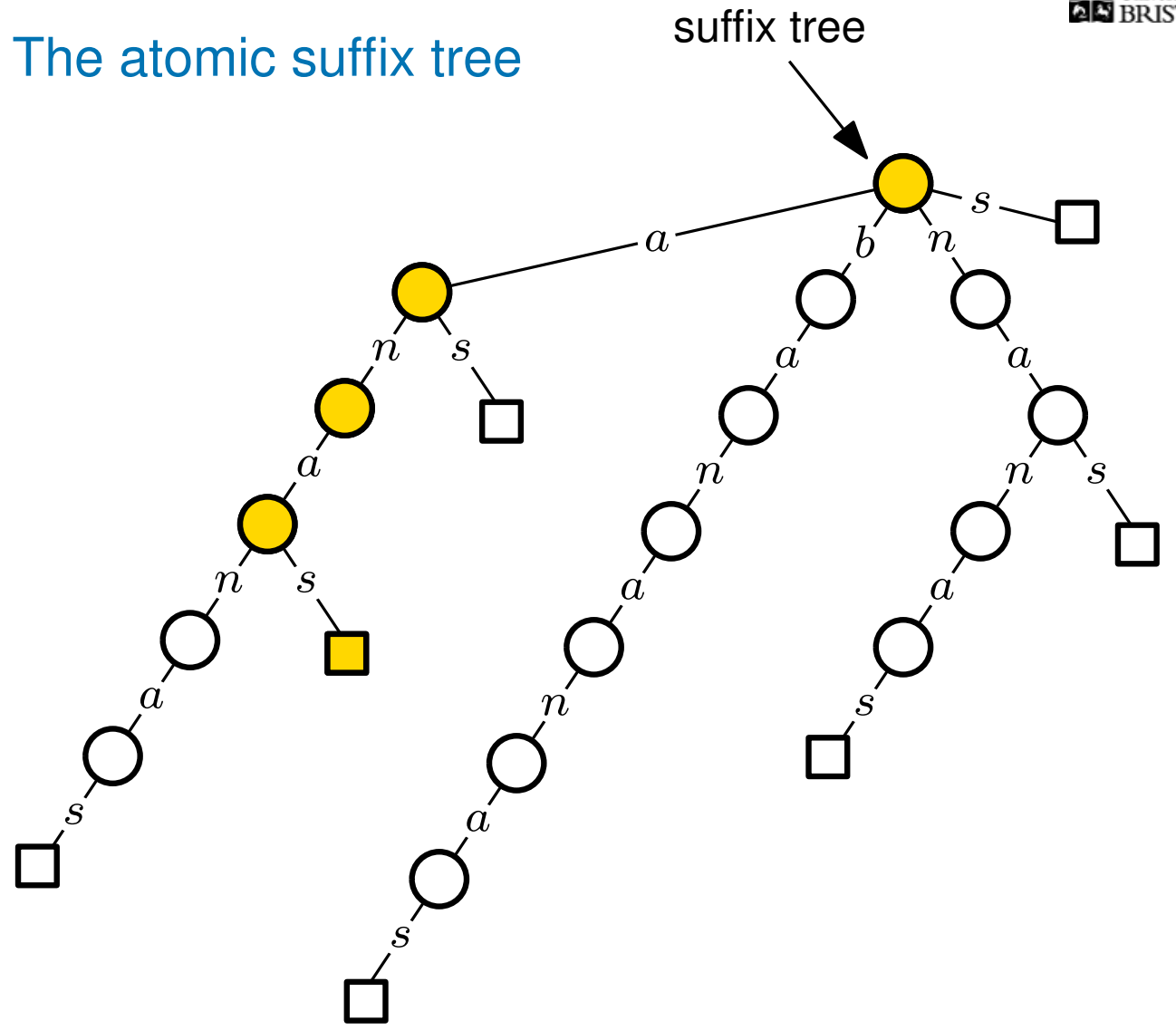
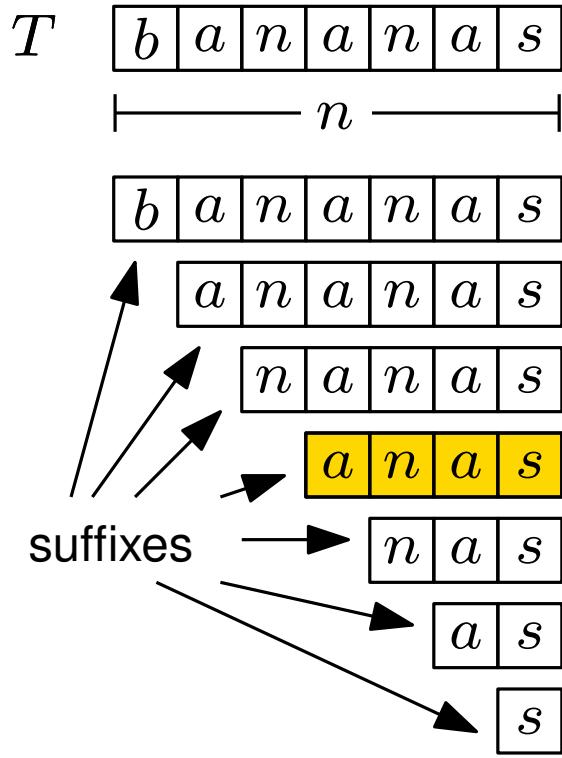








# The atomic suffix tree







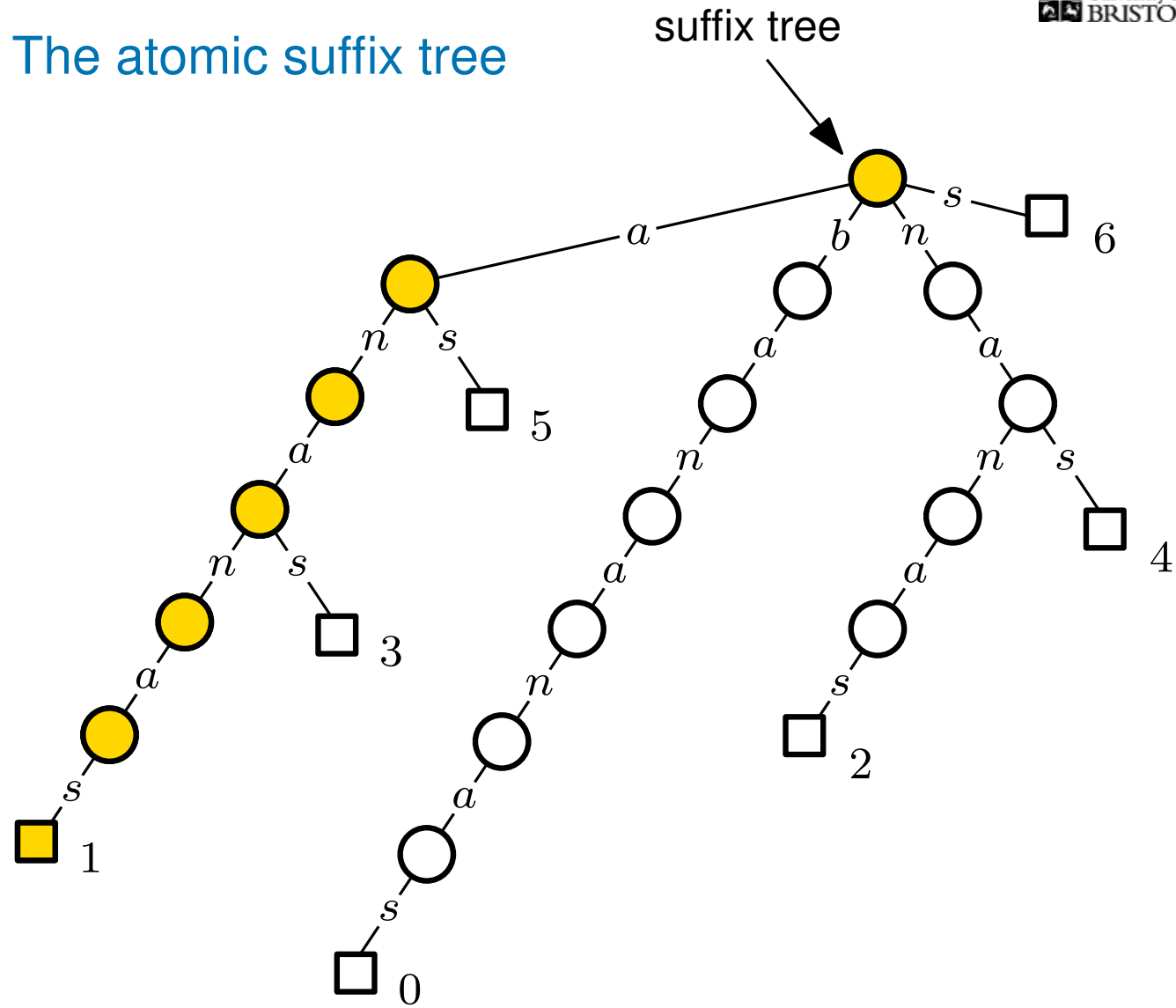
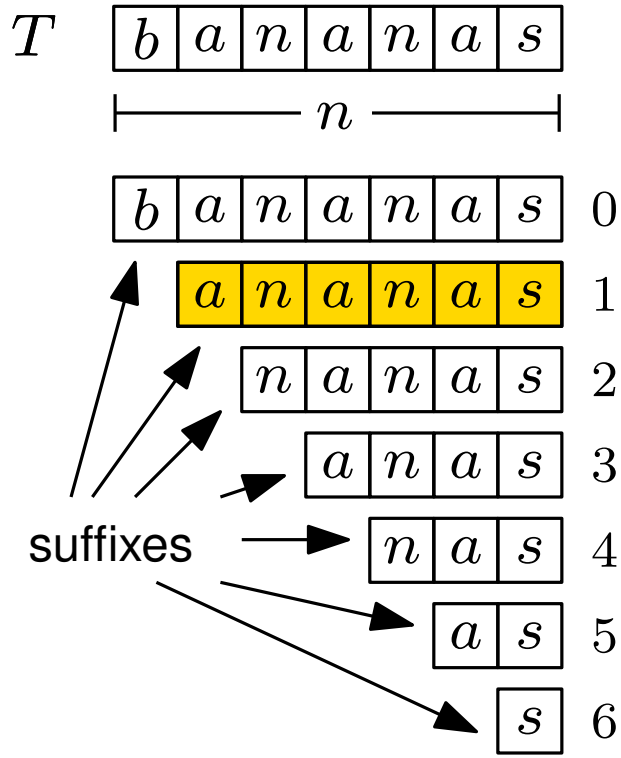






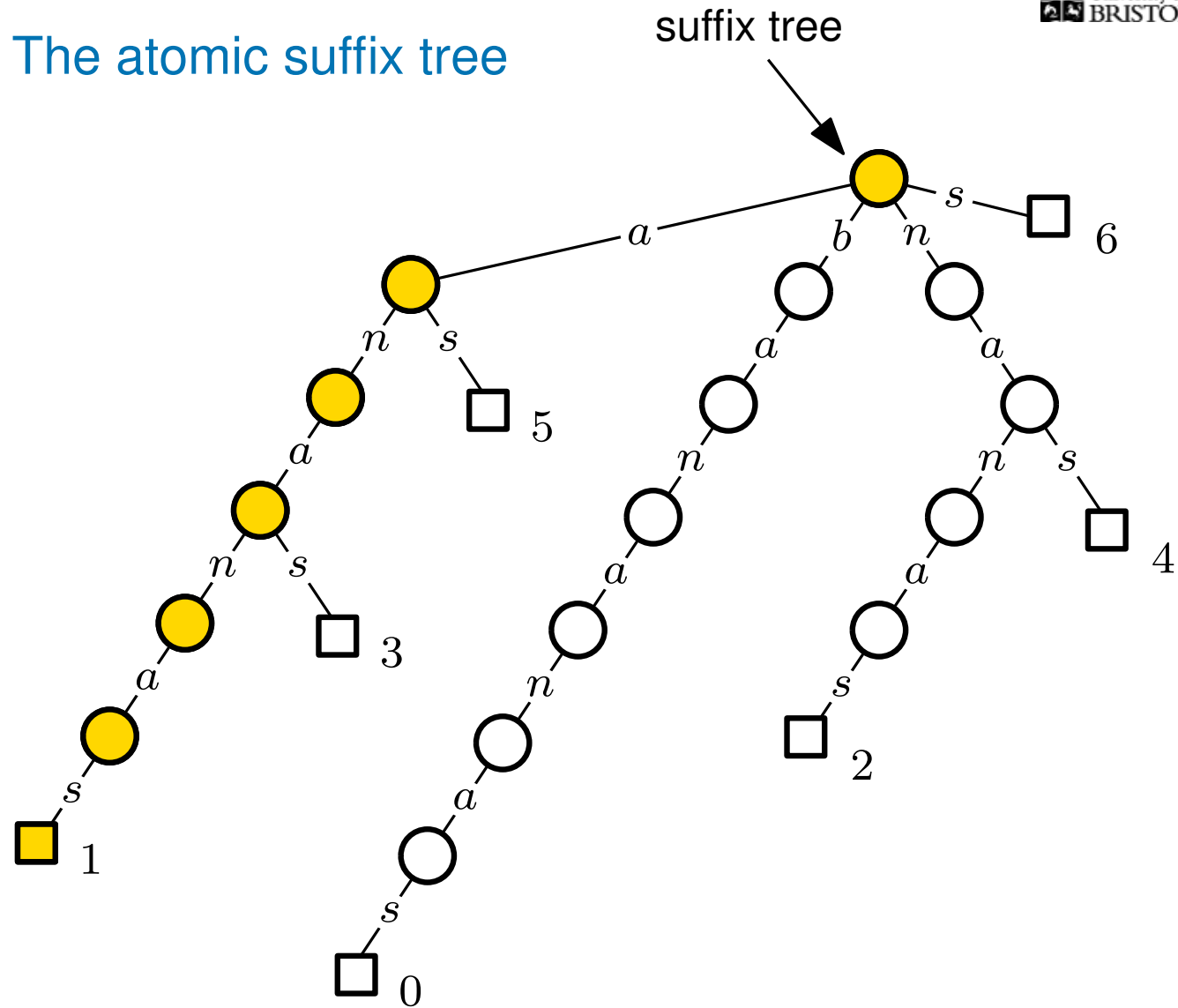
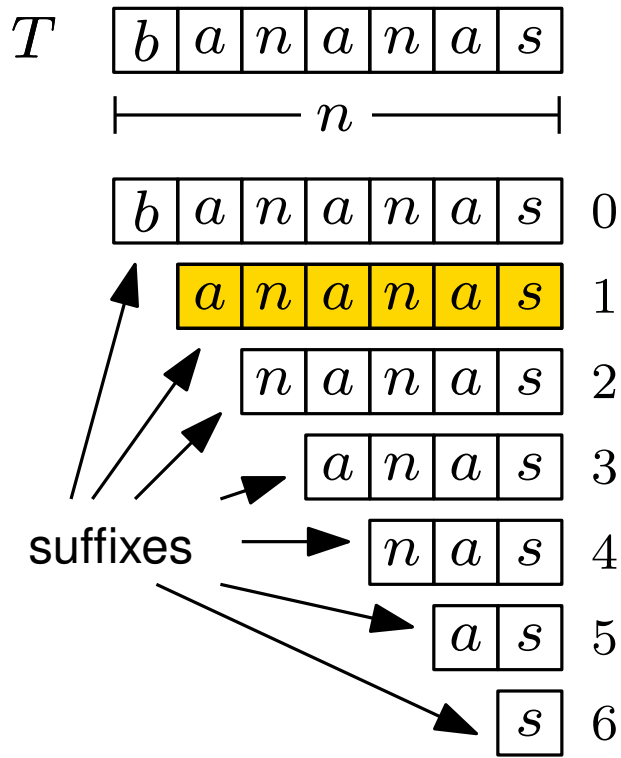


# The atomic suffix tree



- The suffix tree contains every suffix of  $T$  as a root to leaf path

# The atomic suffix tree

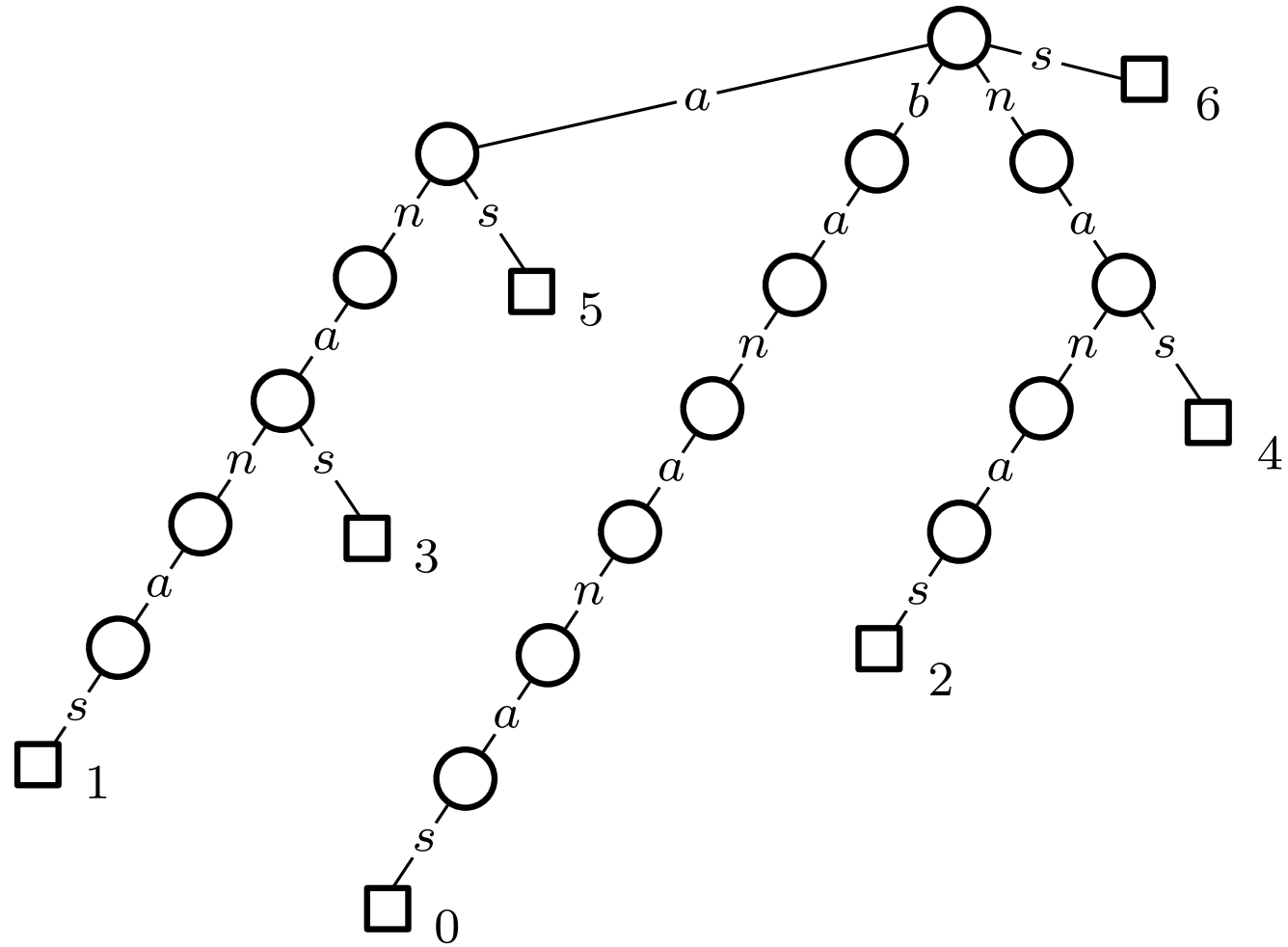
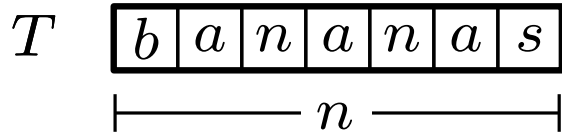


- The suffix tree contains every suffix of  $T$  as a root to leaf path
- Every edge is labelled with a character from  $T$



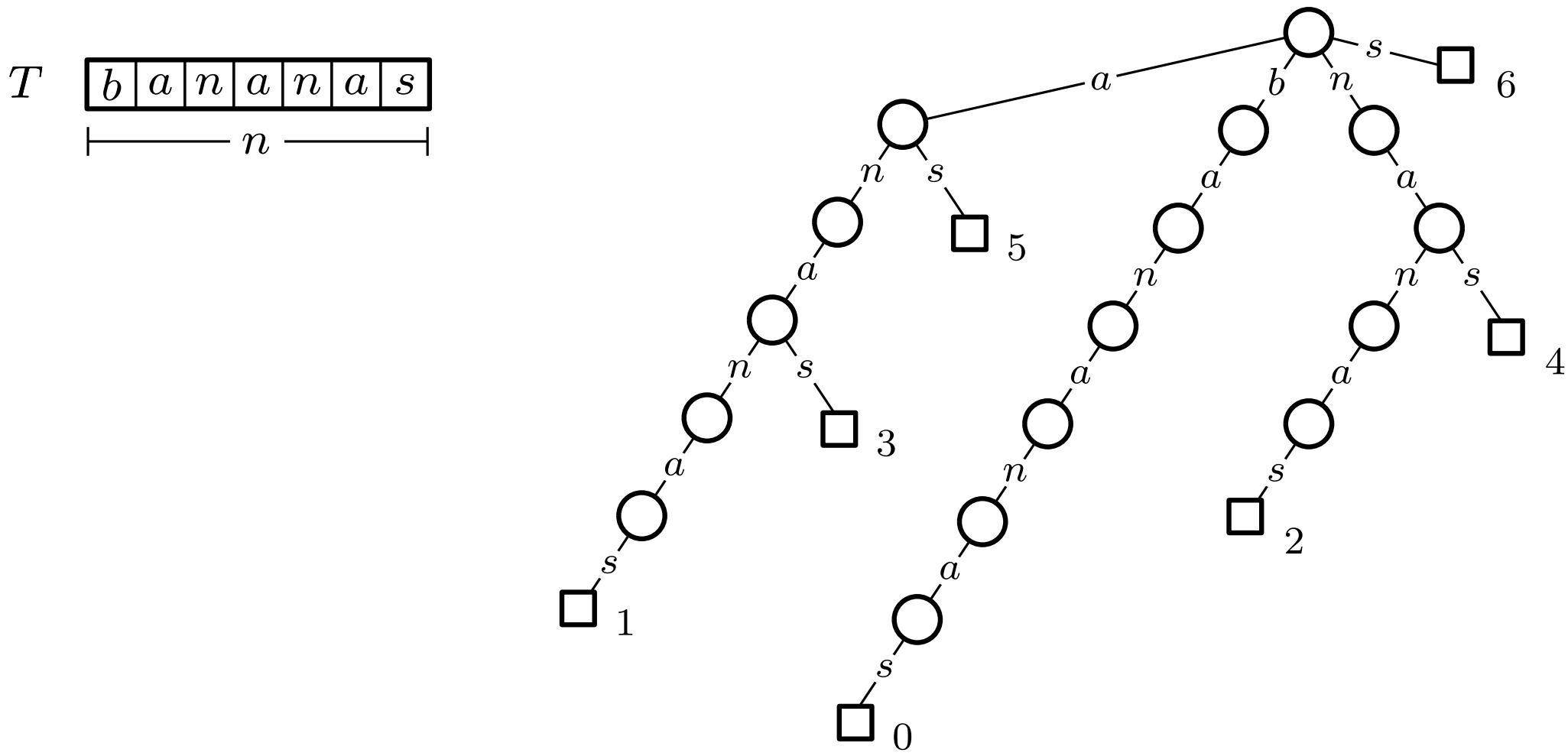


# Searching in an atomic suffix tree



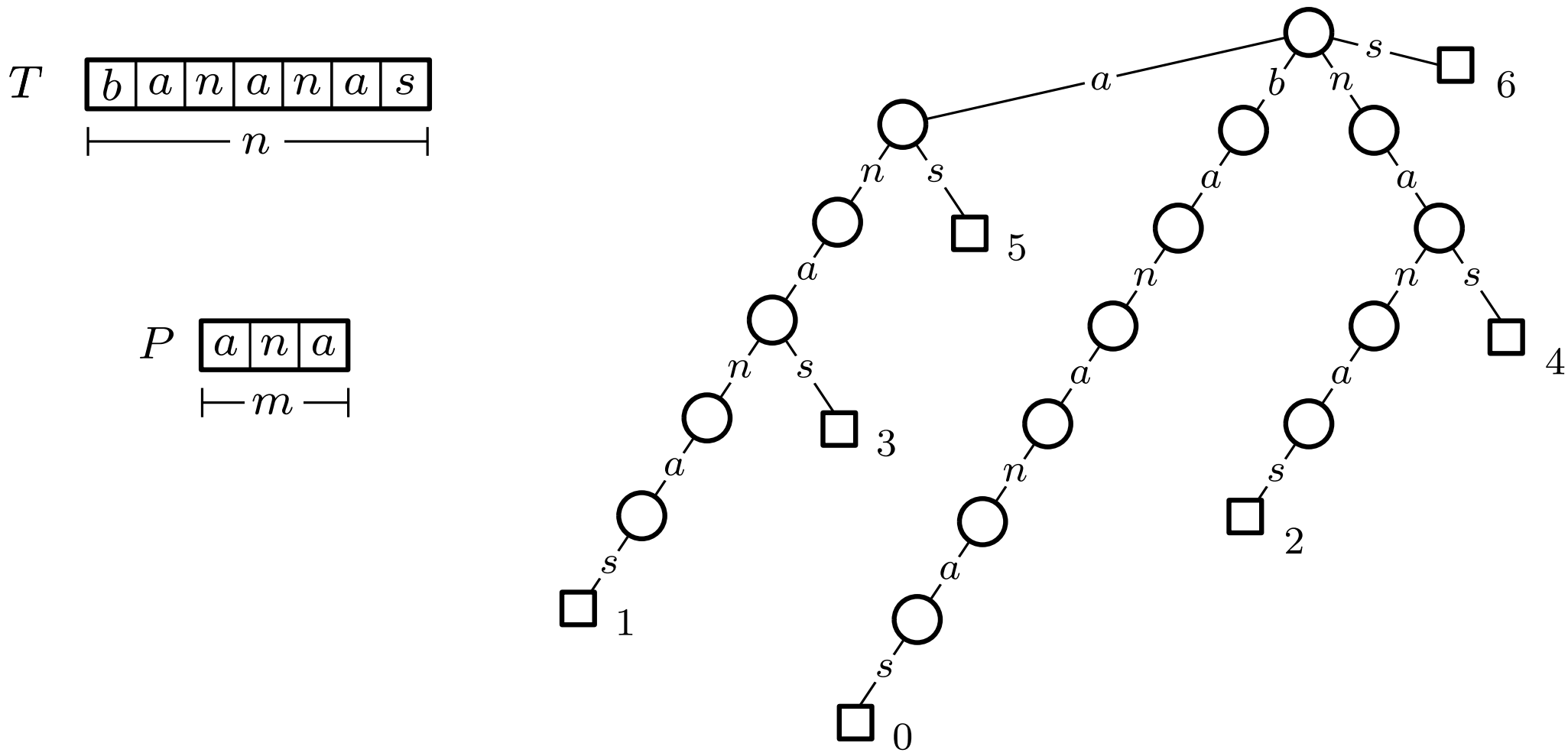


# Searching in an atomic suffix tree



*How do you find a pattern?*

# Searching in an atomic suffix tree

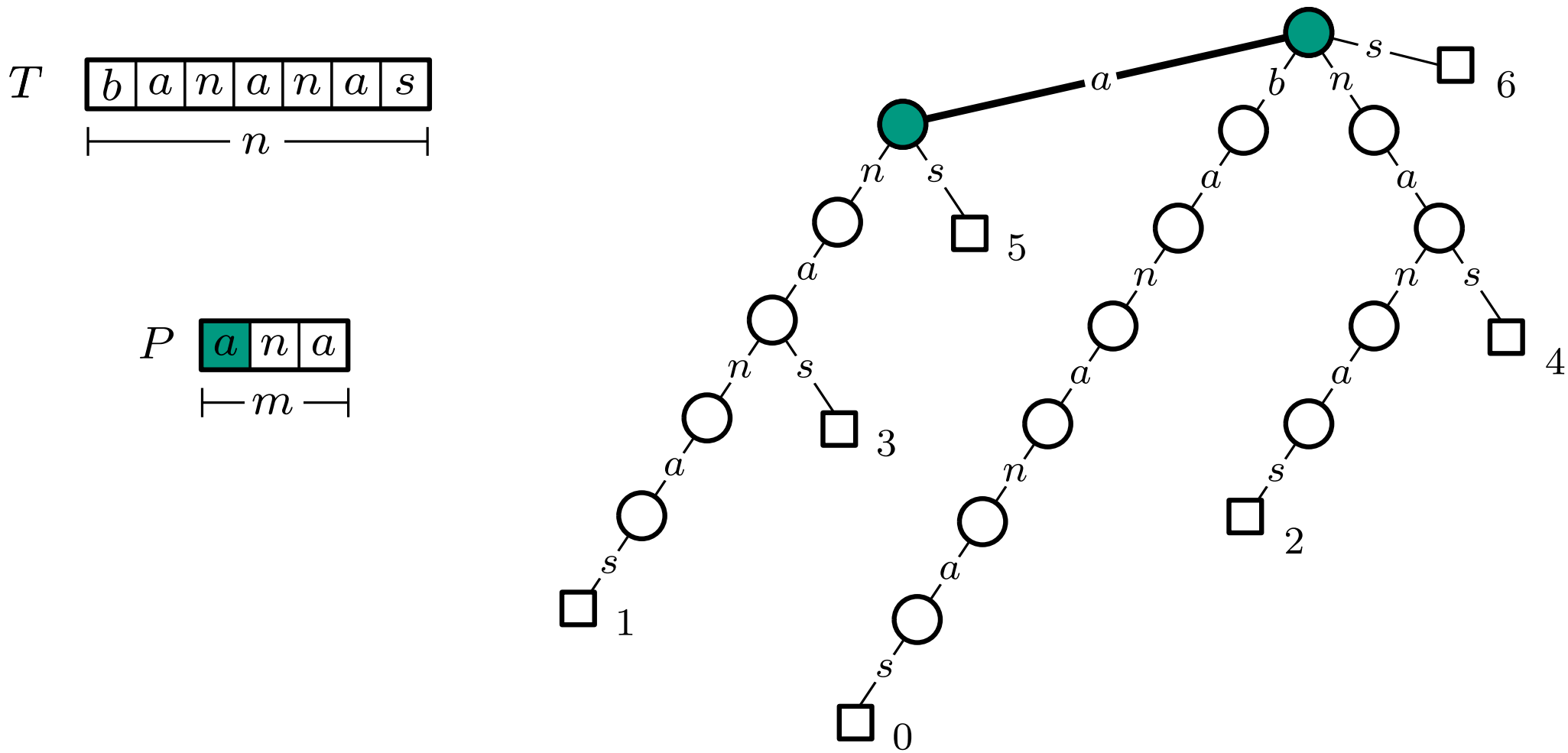


*How do you find a pattern?*





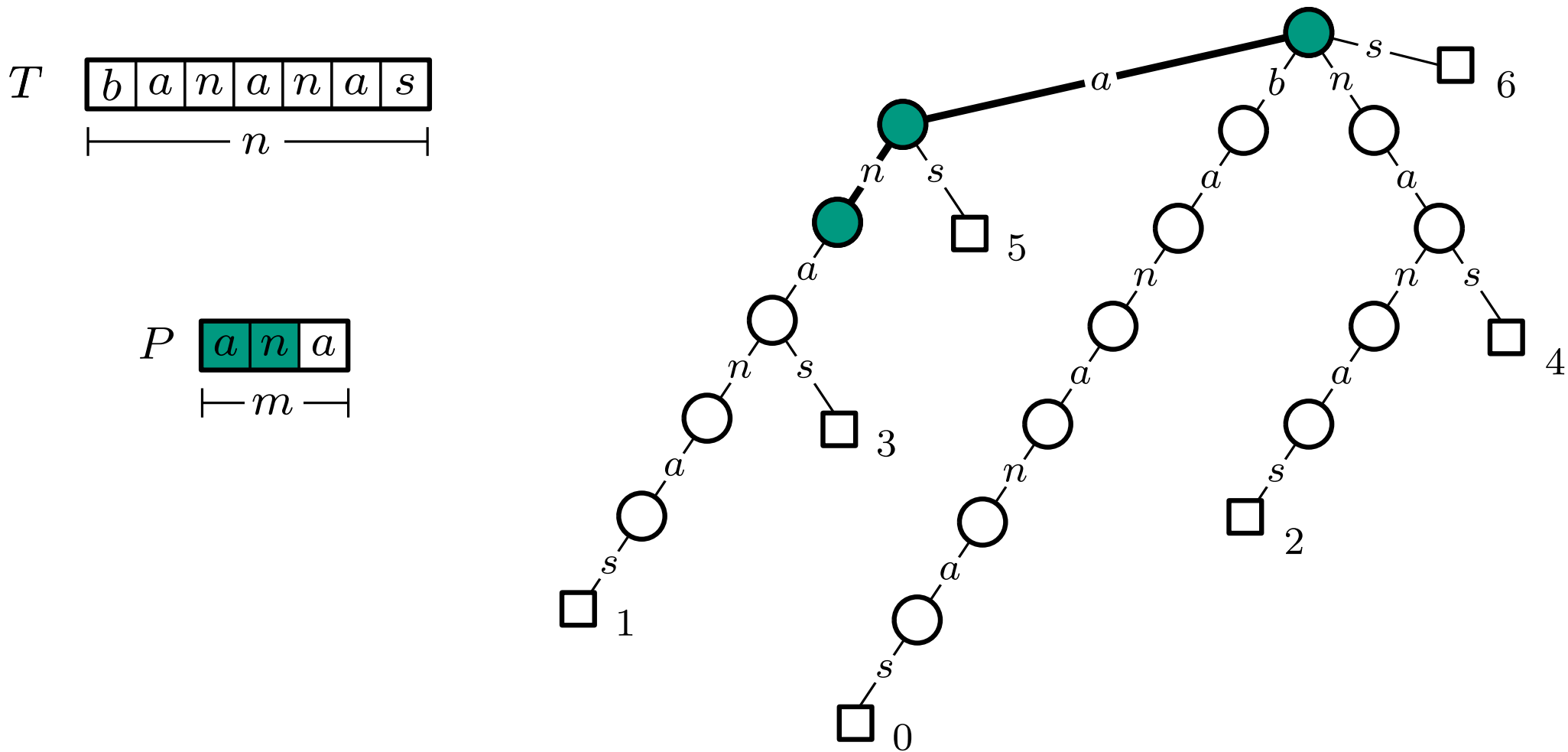
# Searching in an atomic suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

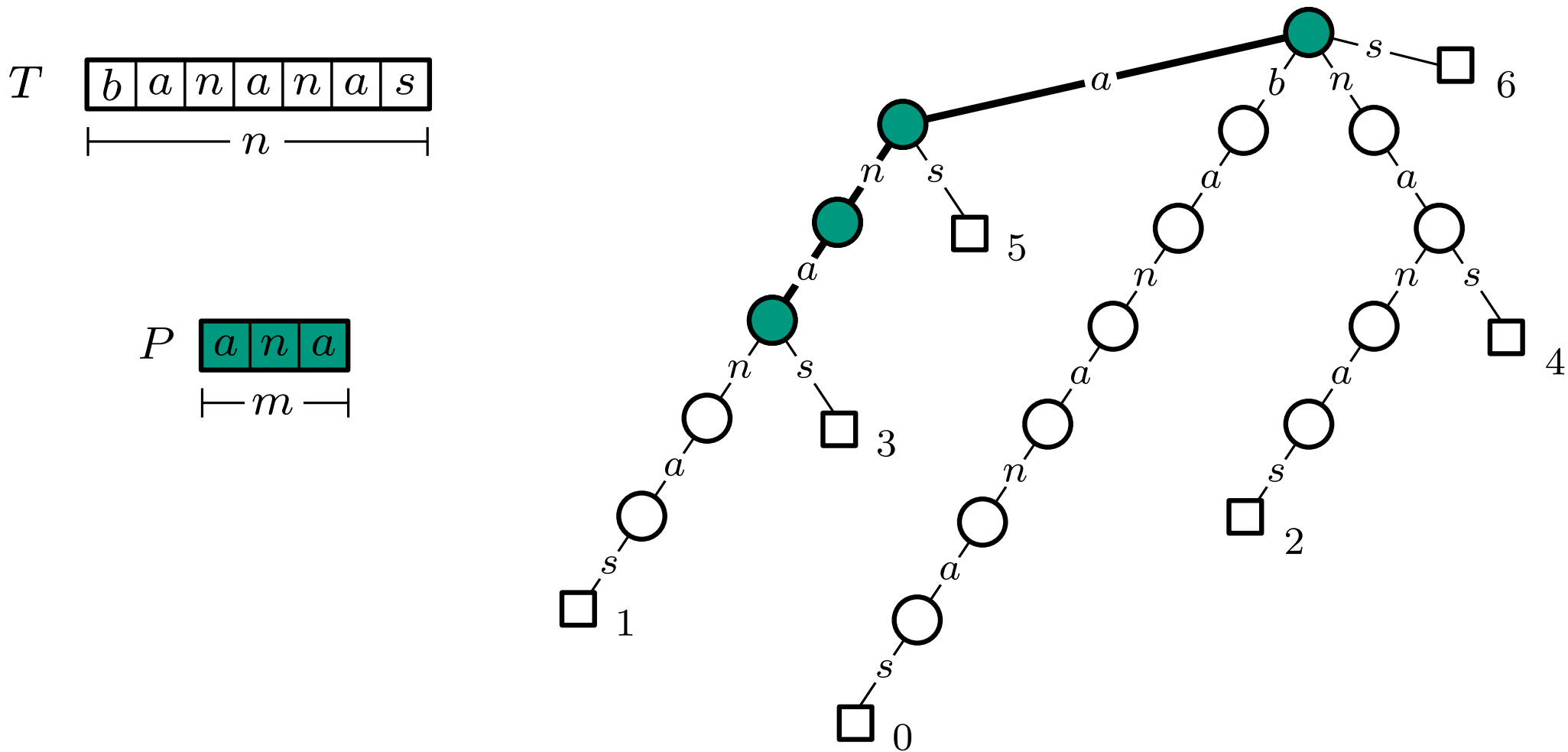
# Searching in an atomic suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

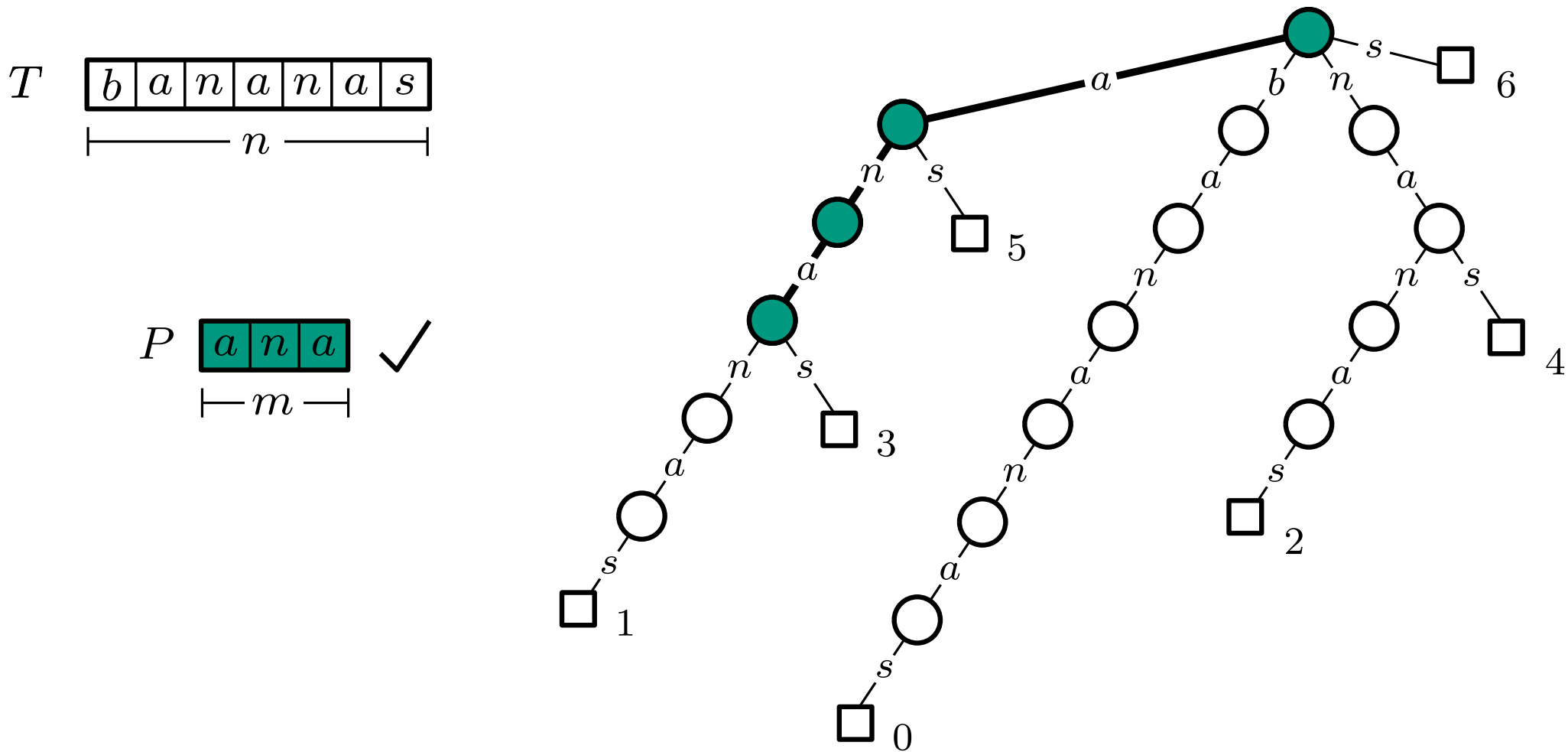
# Searching in an atomic suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

# Searching in an atomic suffix tree



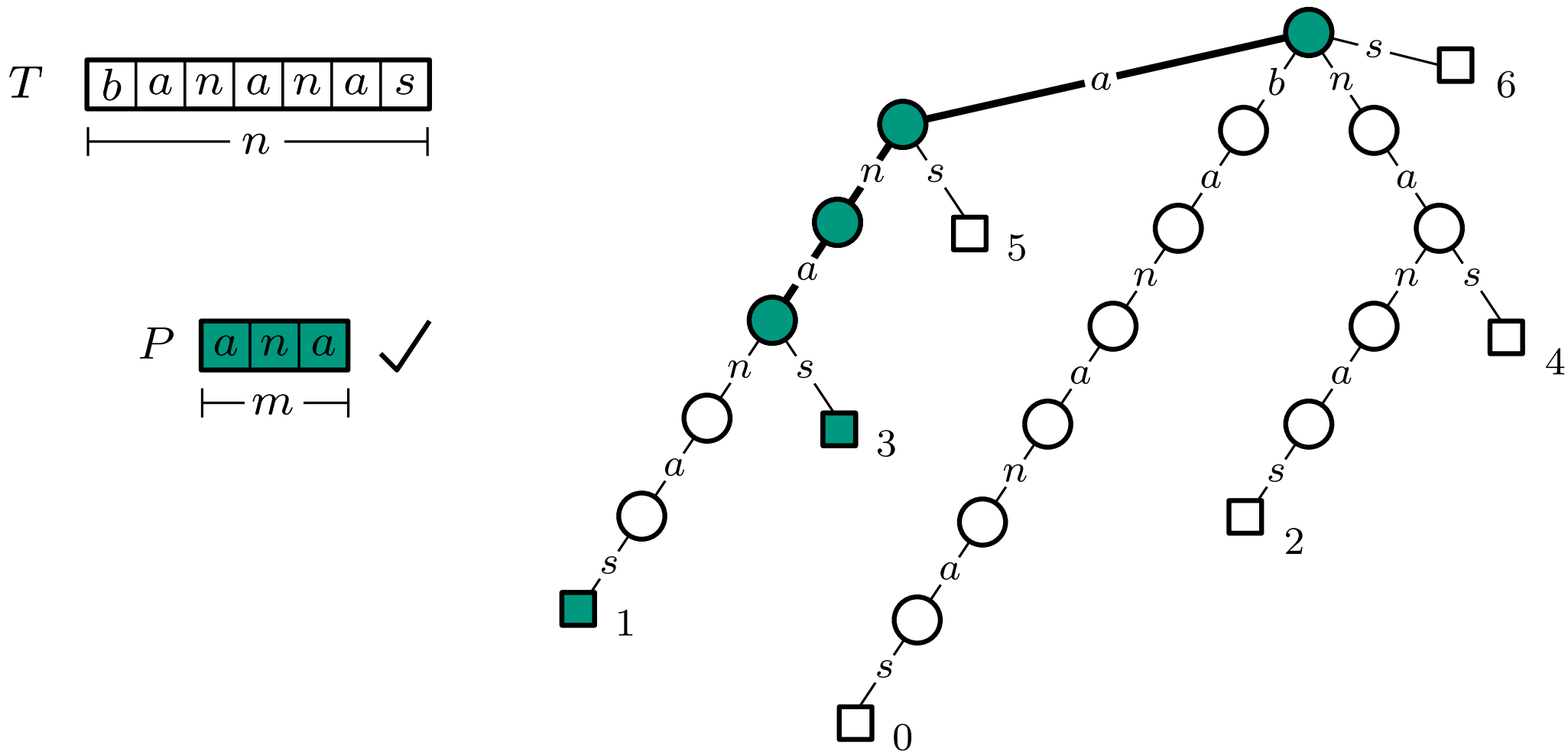
*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree



# Searching in an atomic suffix tree

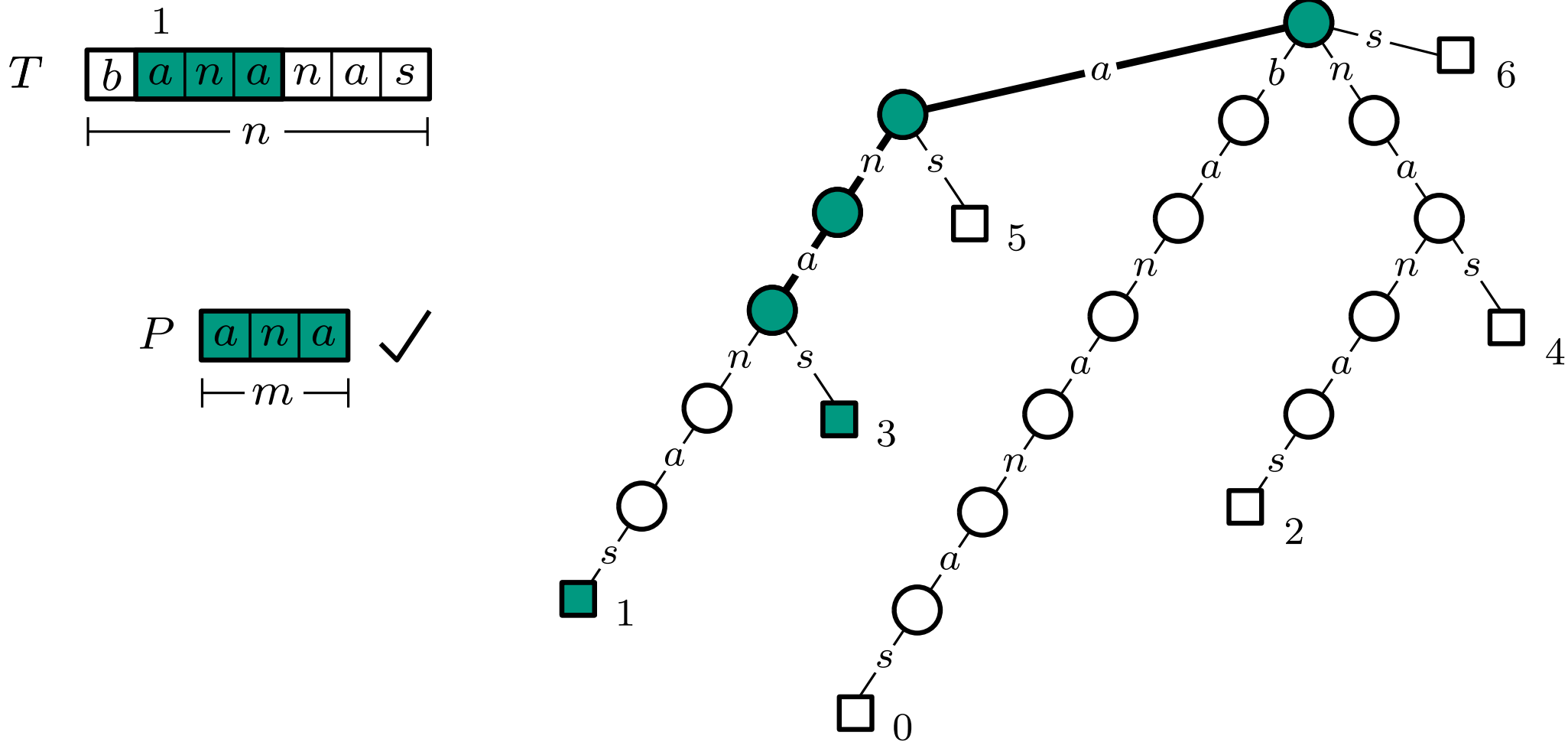


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in an atomic suffix tree



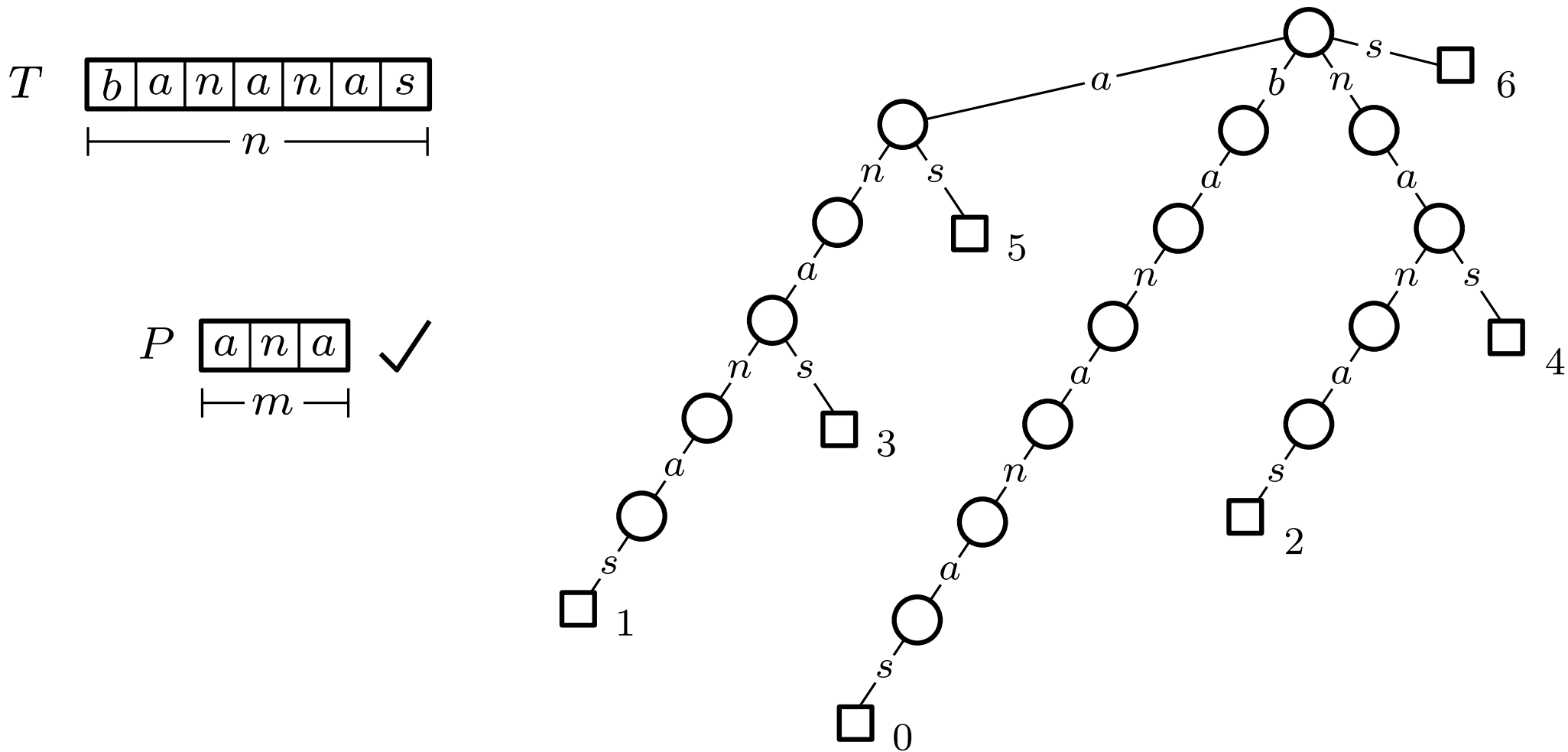
*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree



# Searching in an atomic suffix tree

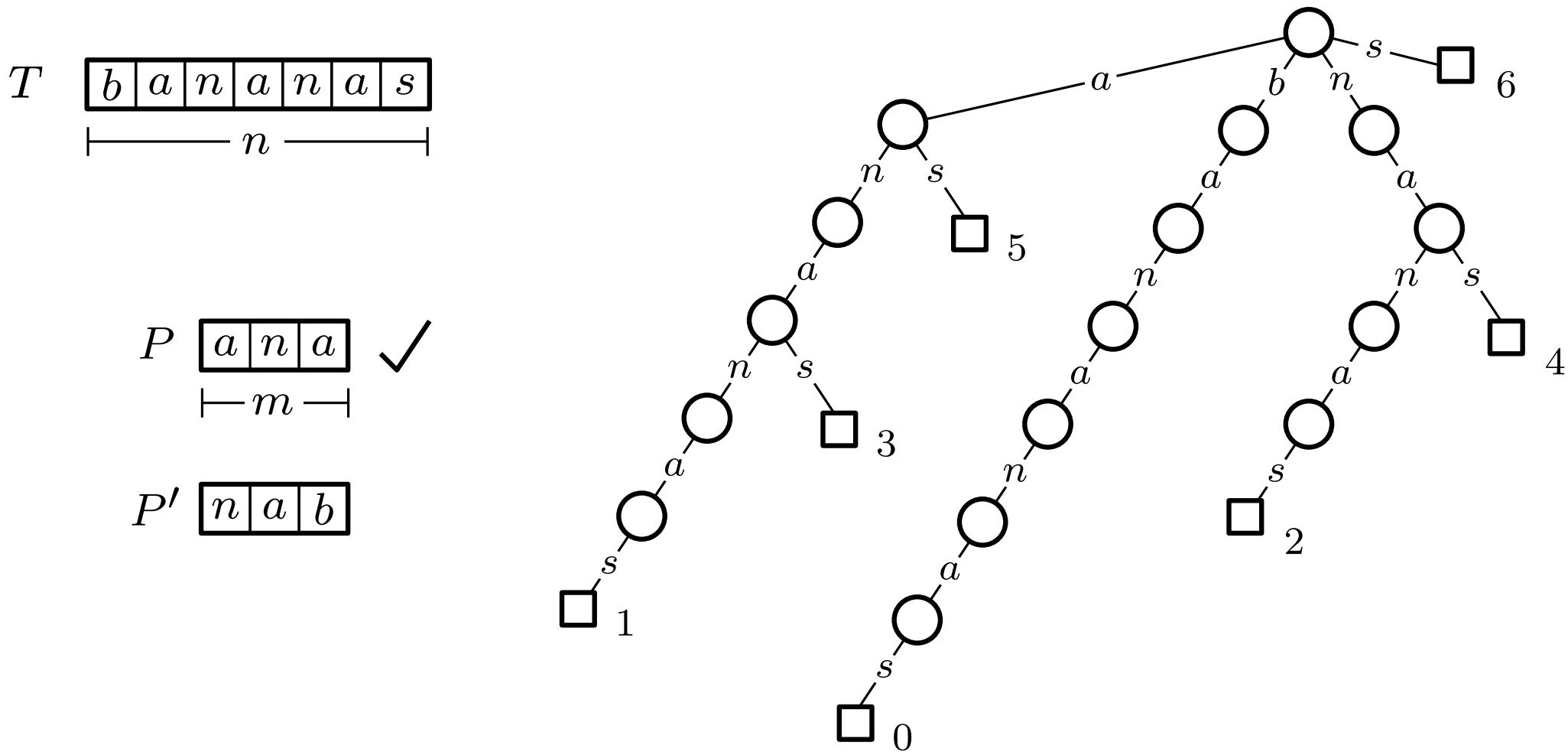


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in an atomic suffix tree



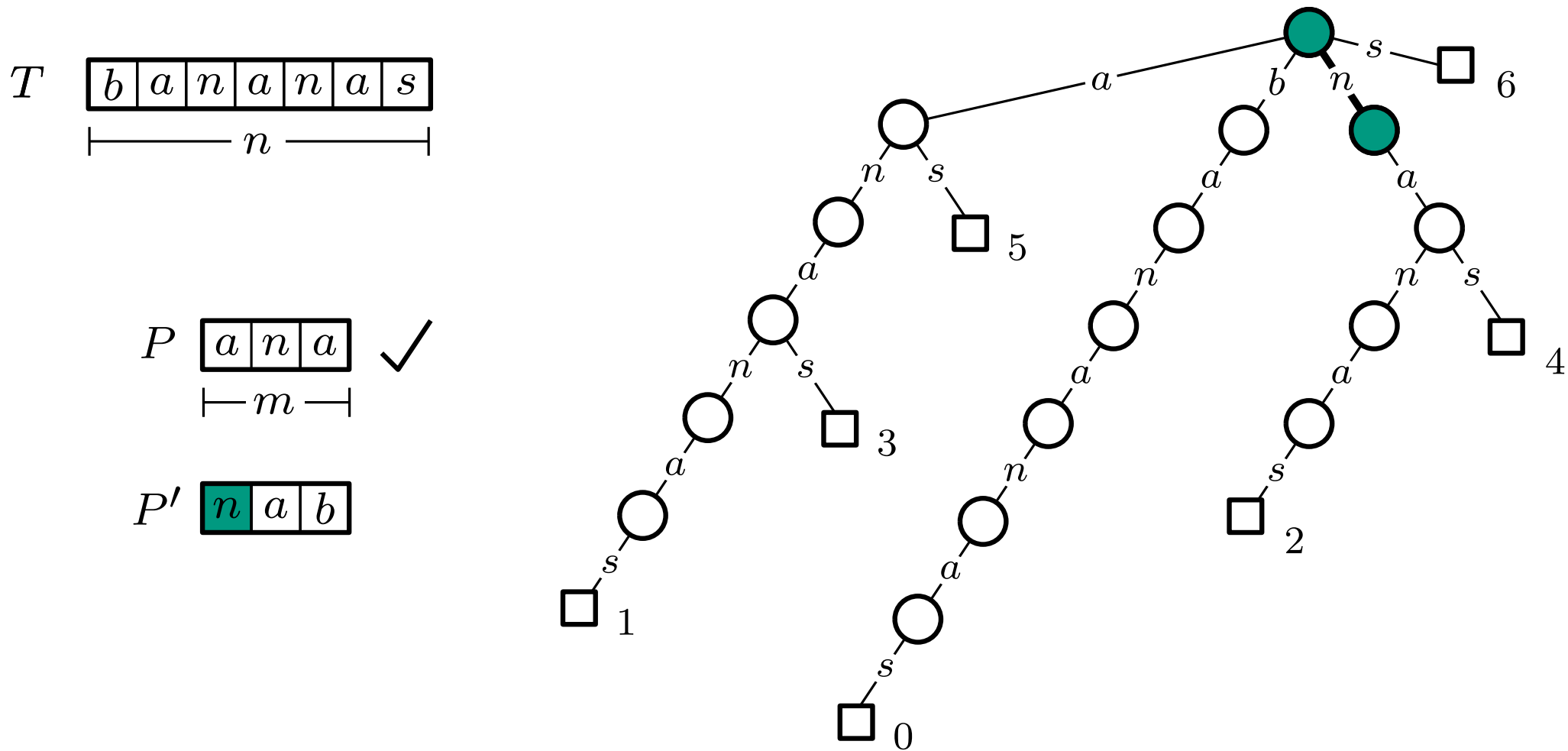
*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree



# Searching in an atomic suffix tree



*How do you find a pattern?*

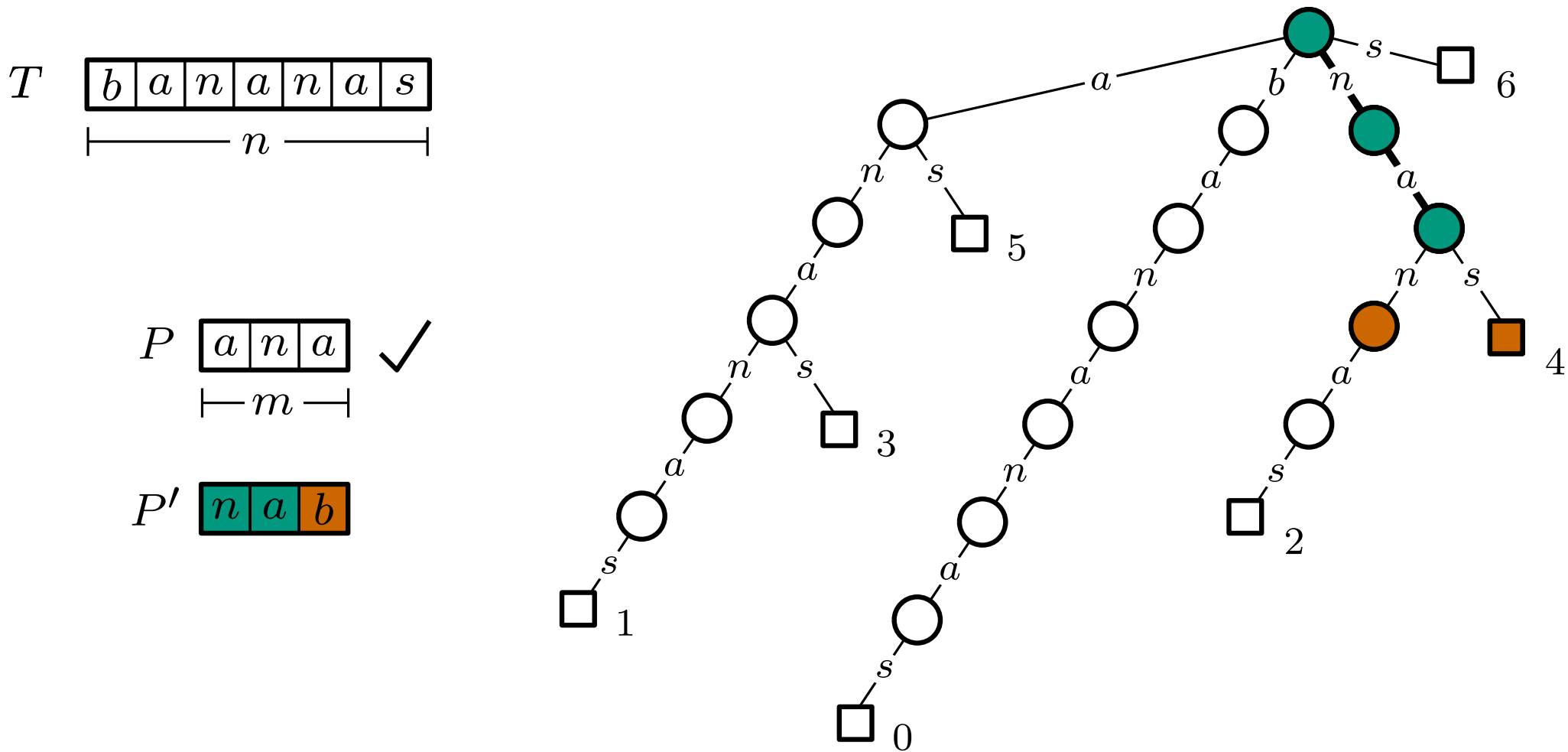
start at the root and walk down the tree

... matches occur at the leaves of the subtree





# Searching in an atomic suffix tree

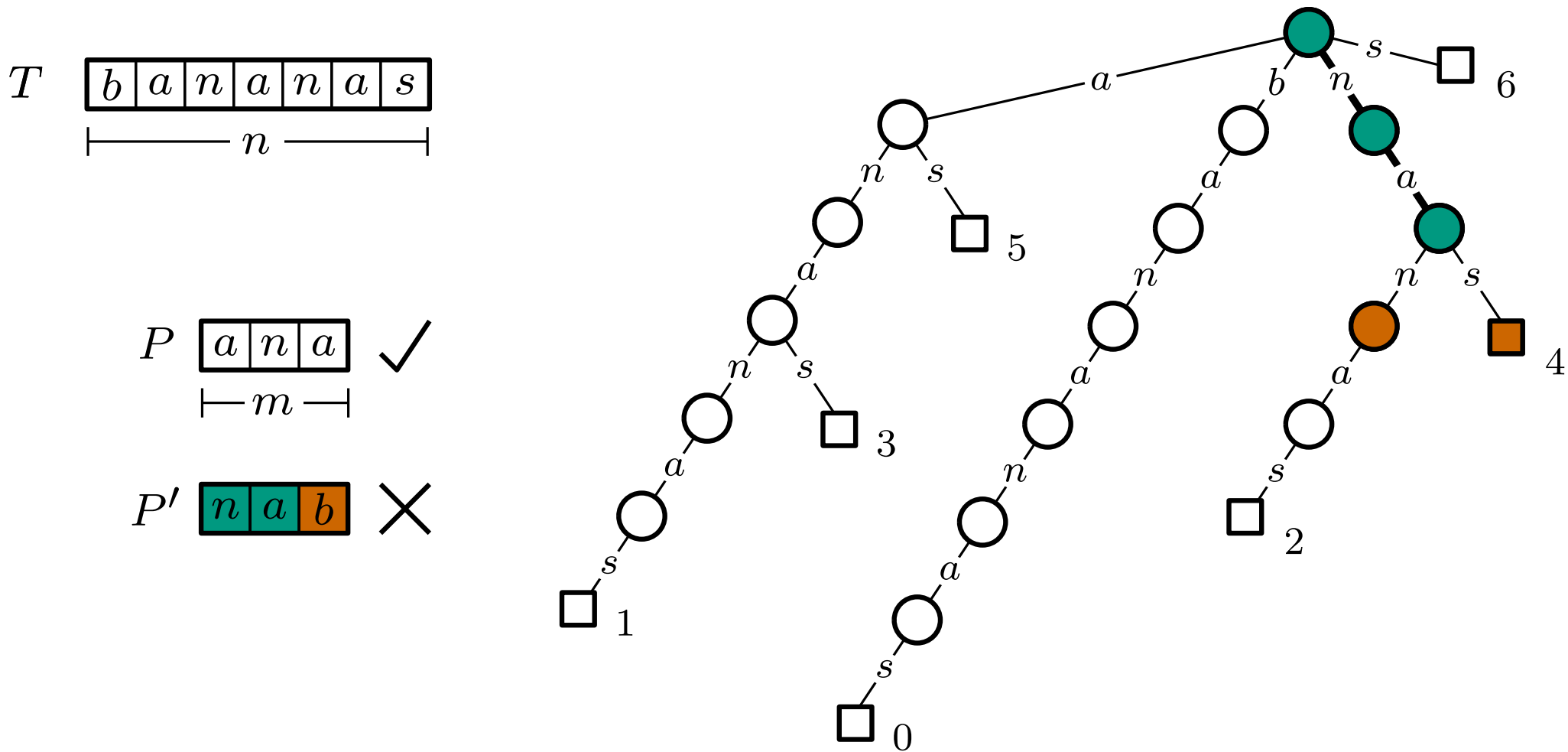


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in an atomic suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

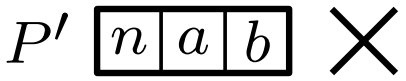
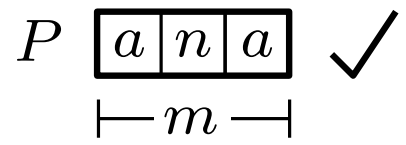
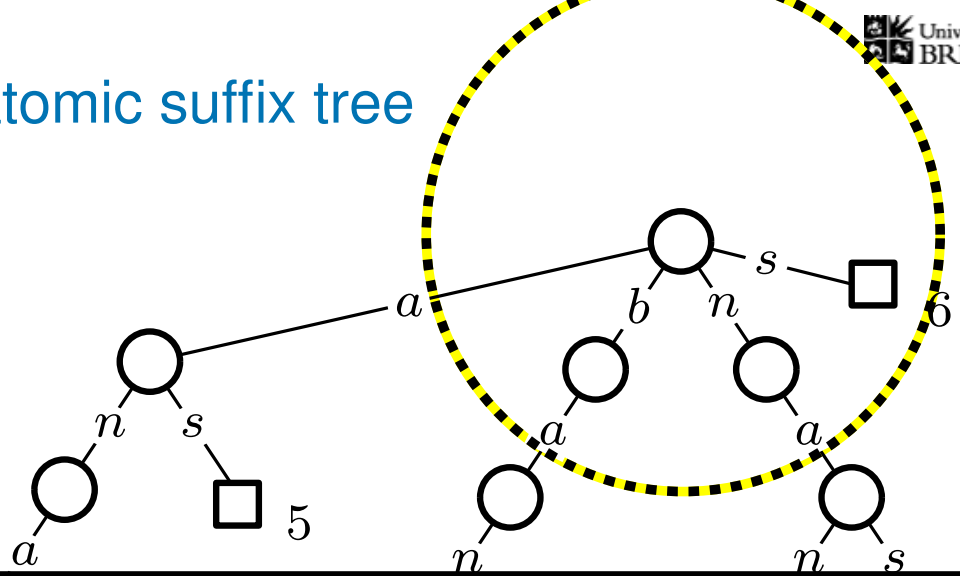
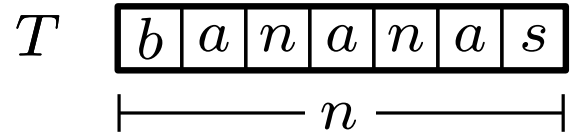
... matches occur at the leaves of the subtree



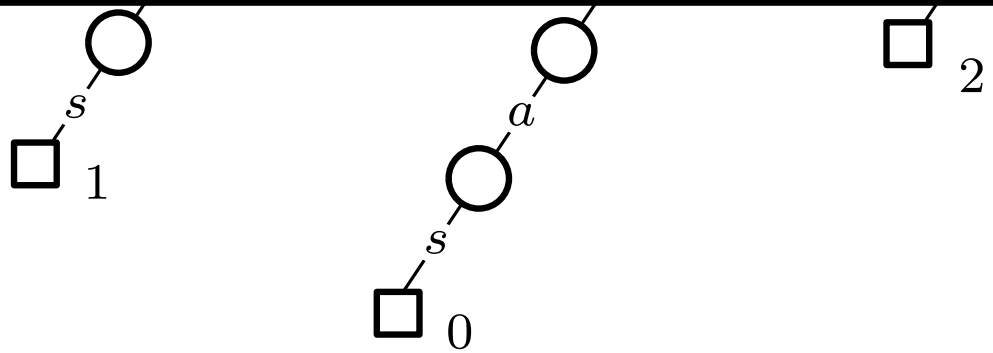




# Searching in an atomic suffix tree



**WARNING!** How long does it take to find the correct child?  
*There could be  $n$  edges here!*  
**In this lecture we assume the alphabet size is a constant**



How do you find a pattern?

start at the root and walk down the tree  
 ... matches occur at the leaves of the subtree

We can decide whether  $P$  matches *somewhere* in  $O(m)$  time

*(we'll worry about outputting the matches later)*

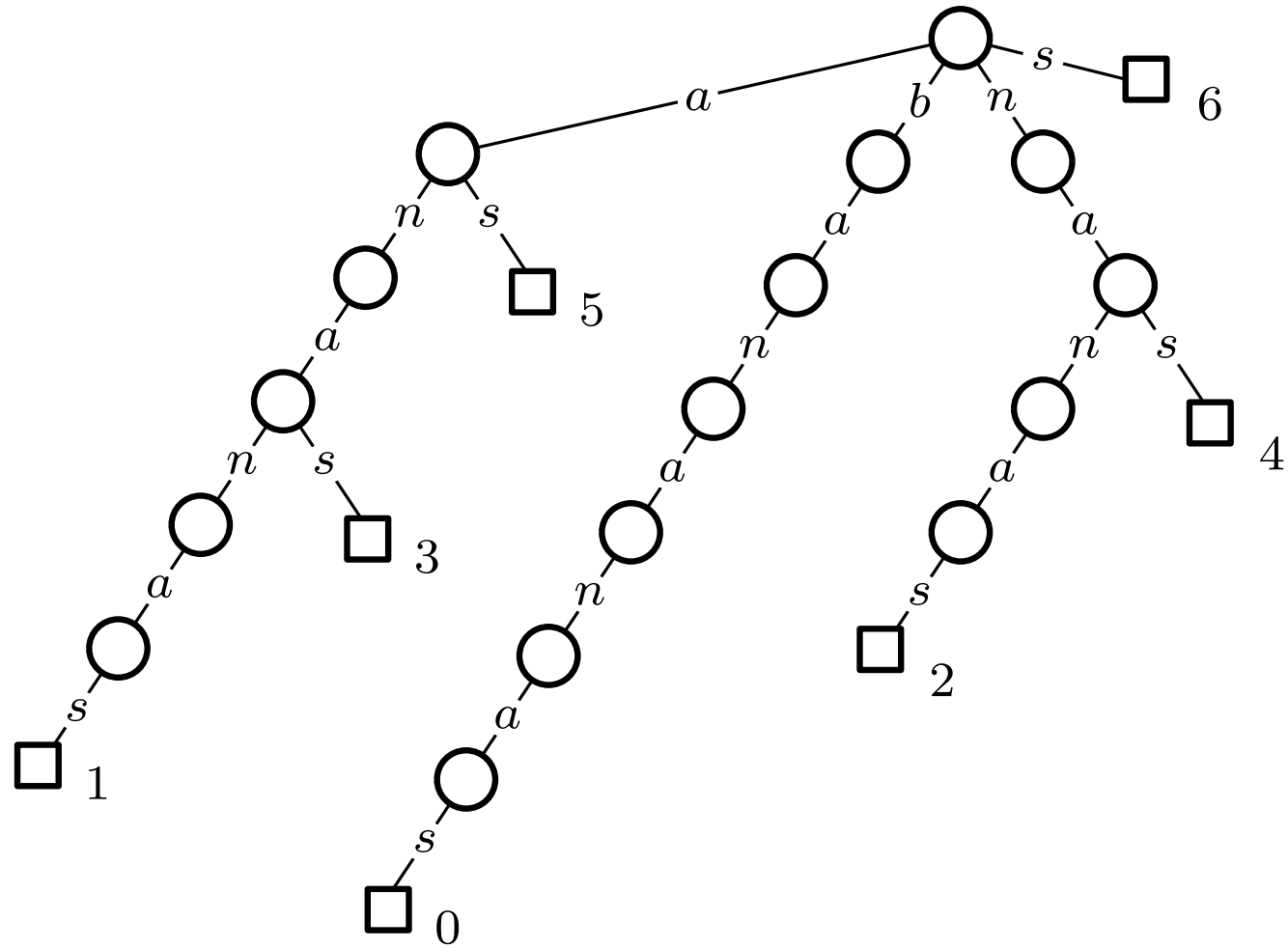
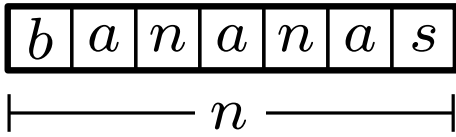






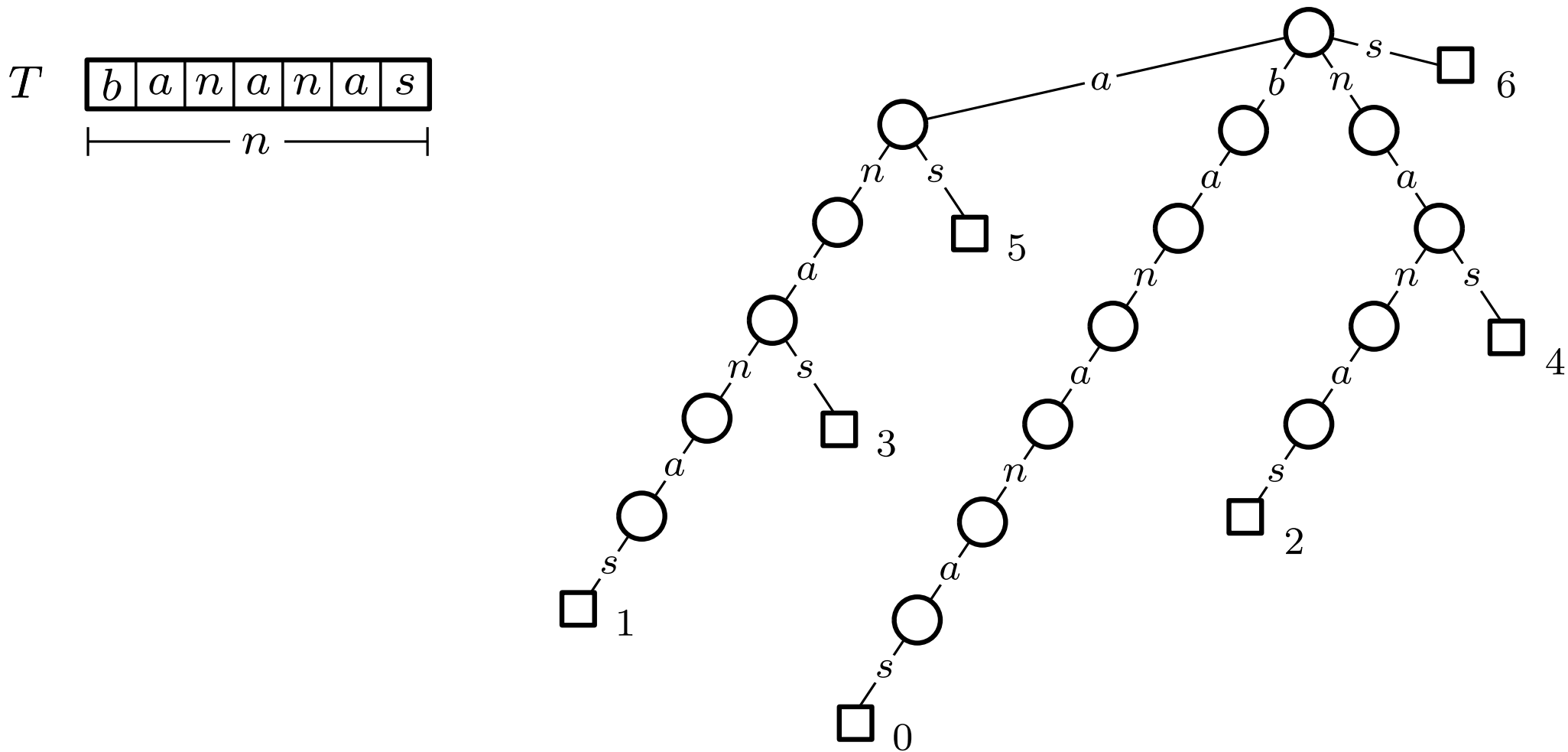
# how large is the atomic suffix tree?

$T$



There are at most  $n$  leaves

# how large is the atomic suffix tree?



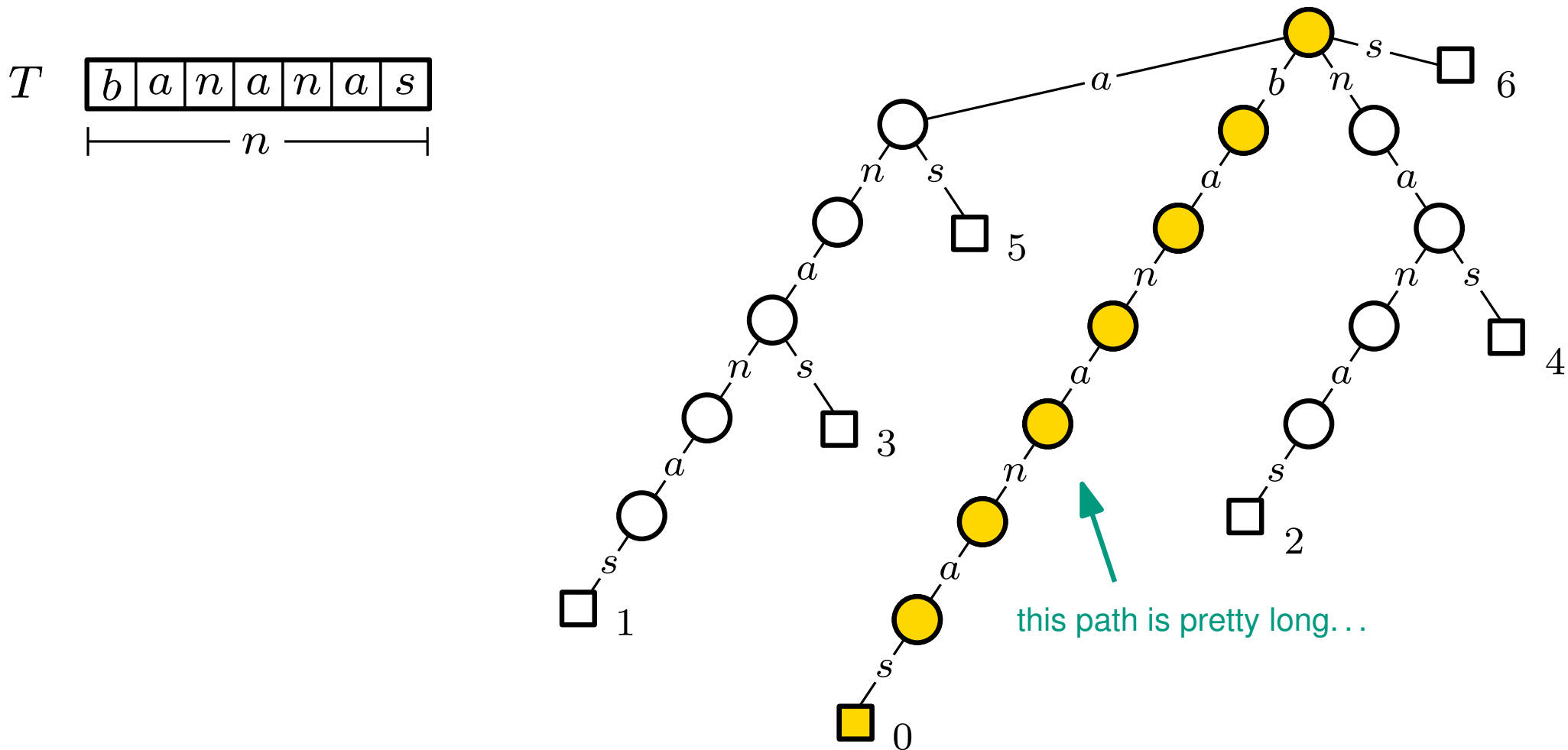
There are at most  $n$  leaves  
 that's good right?







# how large is the atomic suffix tree?

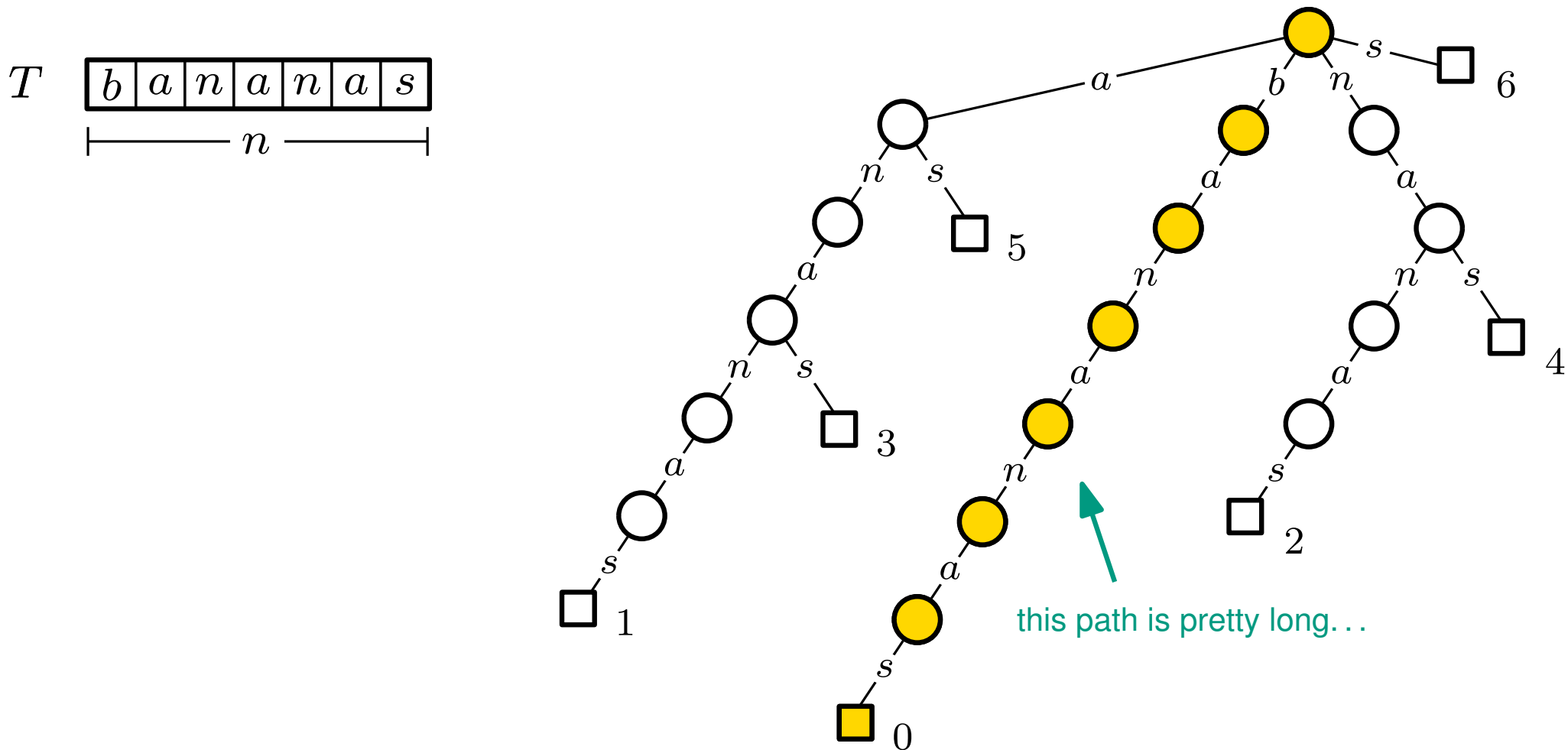


There are at most  $n$  leaves  
that's good right?

Unfortunately there can be *lots* of internal nodes

7 characters

# how large is the atomic suffix tree?



There are at most  $n$  leaves  
that's good right?

Unfortunately there can be *lots* of internal nodes

7 characters

23 nodes



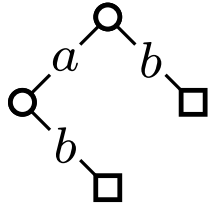
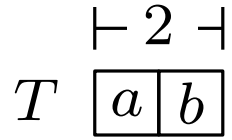


how large is the atomic suffix tree?

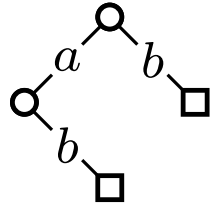
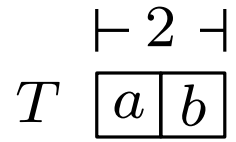
how large is the atomic suffix tree?

$T$   $\begin{array}{|c|c|} \hline a & b \\ \hline \end{array}$

how large is the atomic suffix tree?

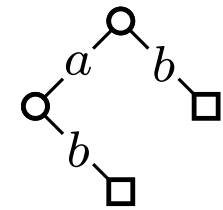
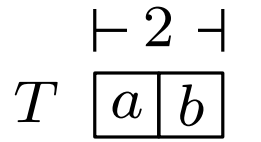


how large is the atomic suffix tree?

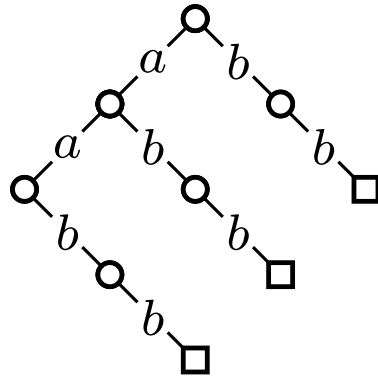
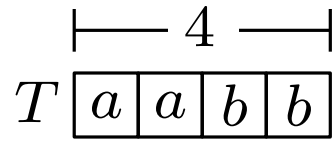


4 nodes

# how large is the atomic suffix tree?

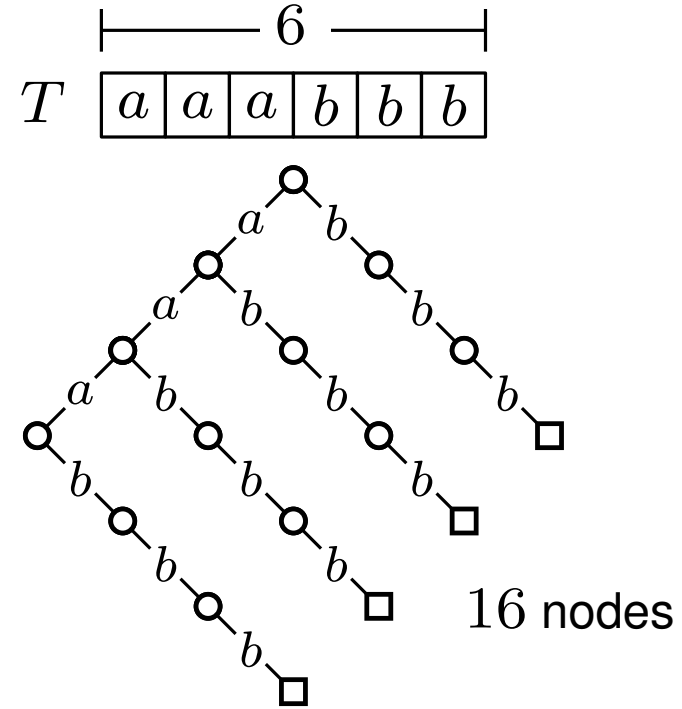
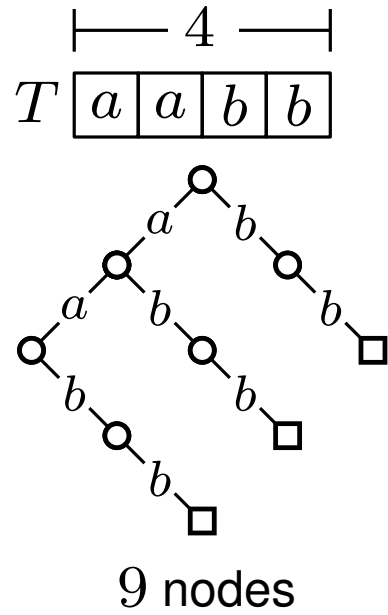
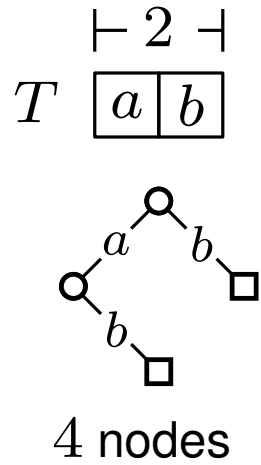


4 nodes

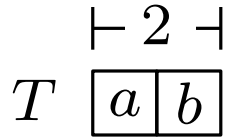


9 nodes

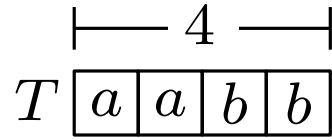
## how large is the atomic suffix tree?



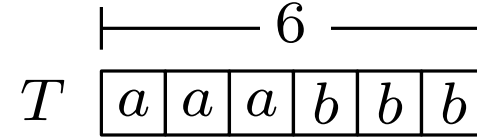
# how large is the atomic suffix tree?



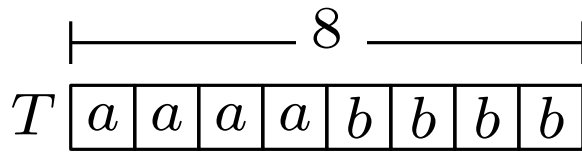
4 nodes



9 nodes

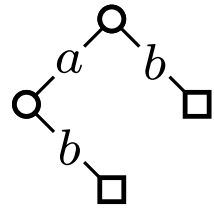
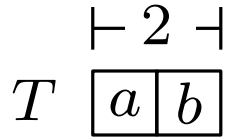


16 nodes

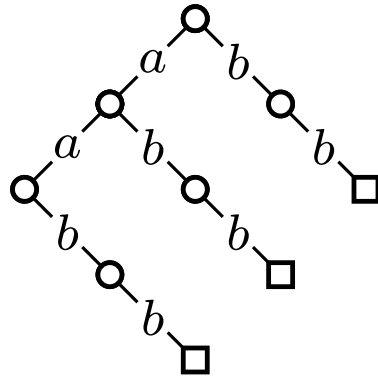
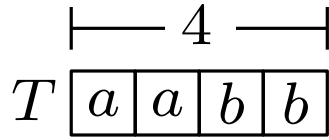


25 nodes

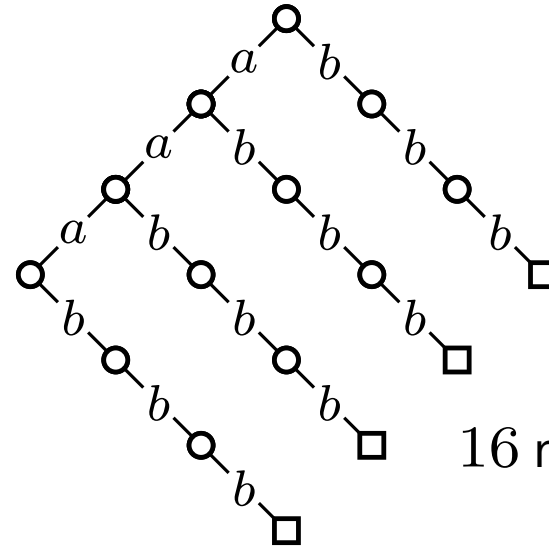
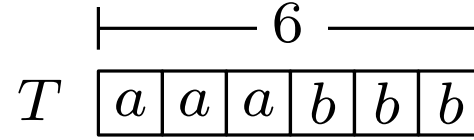
# how large is the atomic suffix tree?



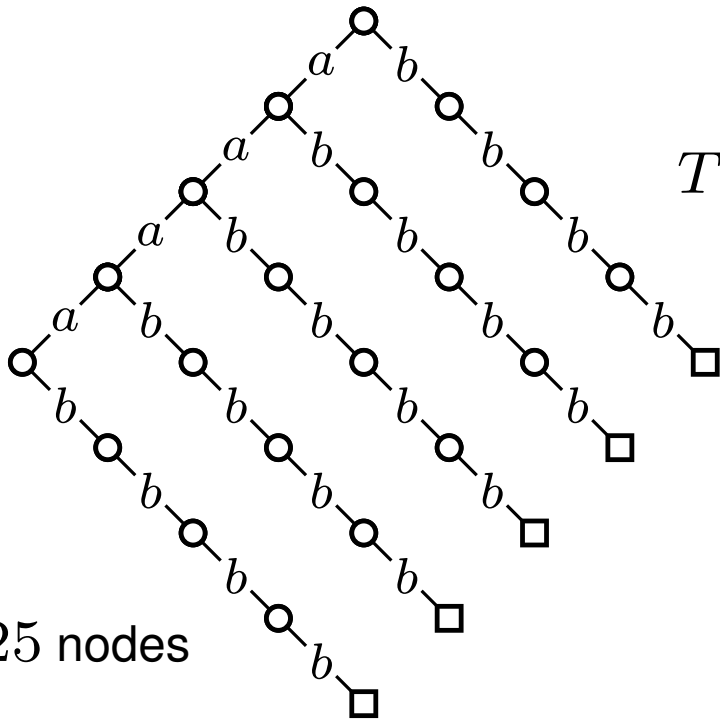
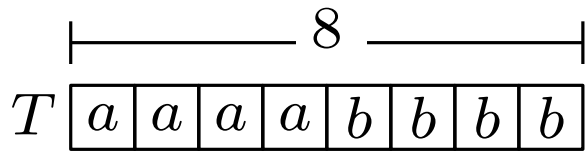
4 nodes



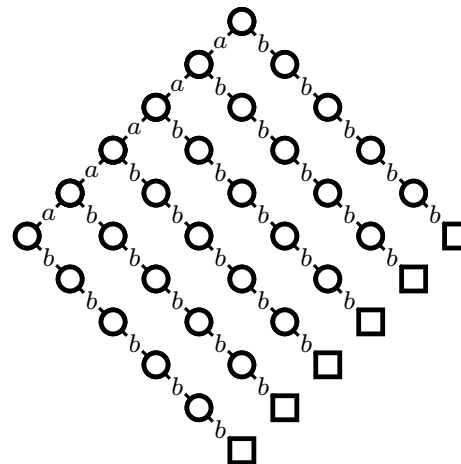
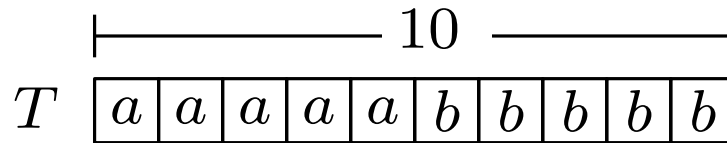
9 nodes



16 nodes



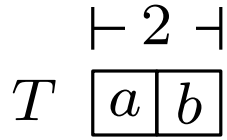
25 nodes



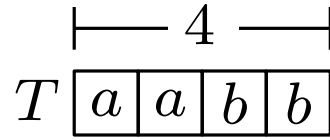
36 nodes



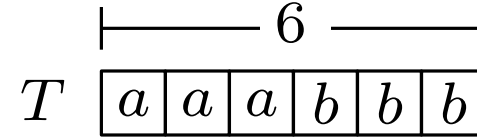
# how large is the atomic suffix tree?



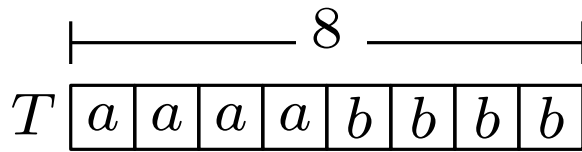
4 nodes



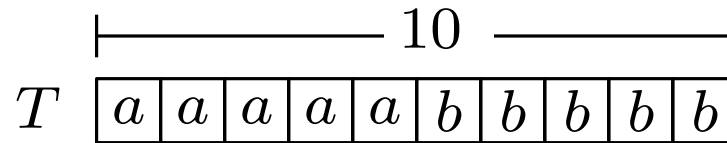
9 nodes



16 nodes



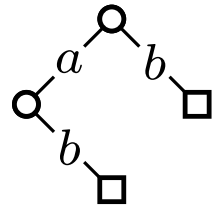
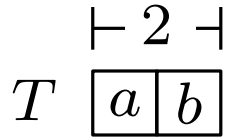
25 nodes



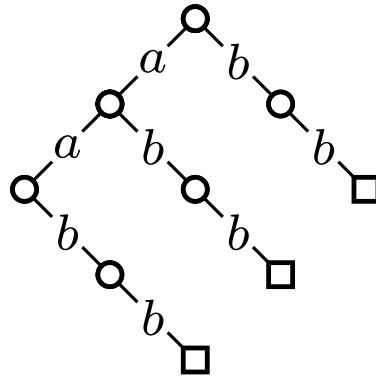
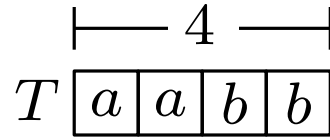
36 nodes

An atomic suffix tree can have  $((n/2) + 1)^2$  nodes

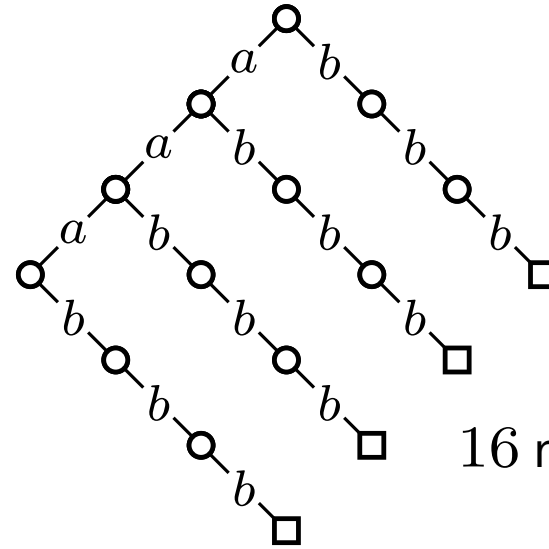
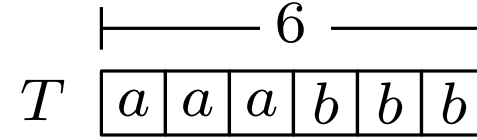
# how large is the atomic suffix tree?



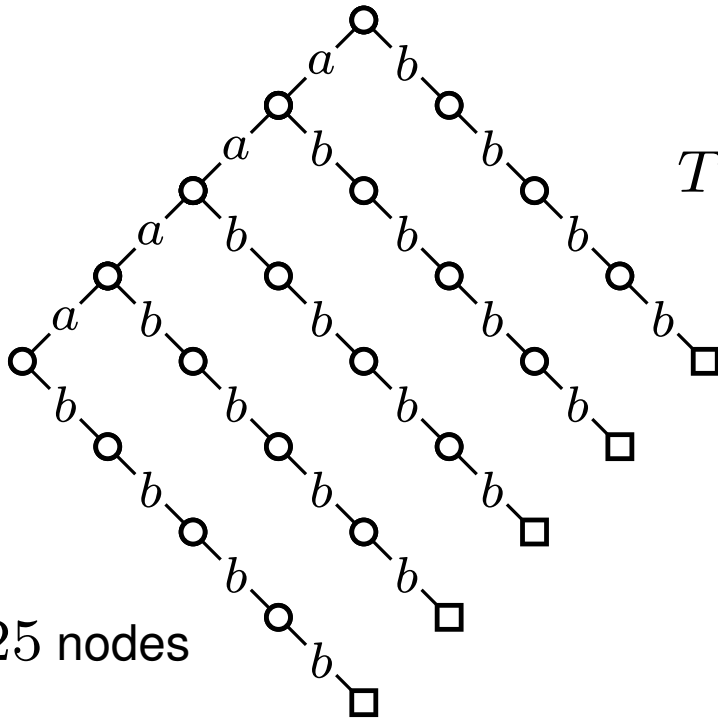
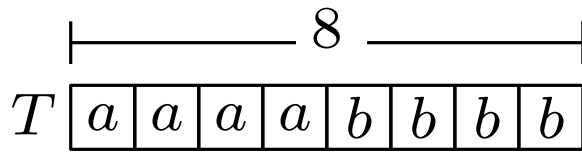
4 nodes



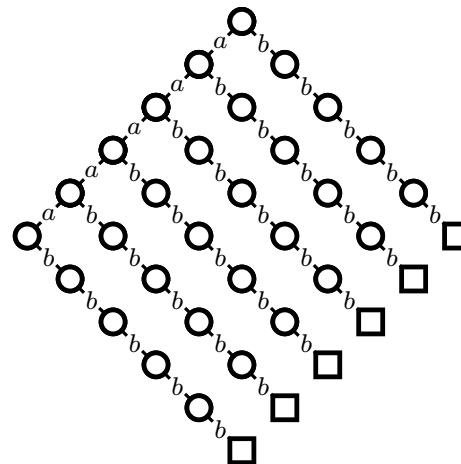
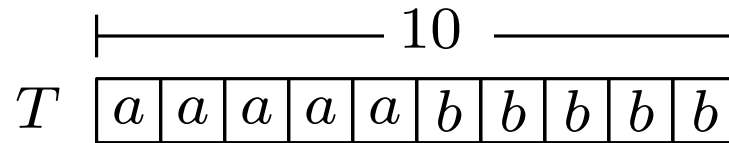
9 nodes



16 nodes



25 nodes

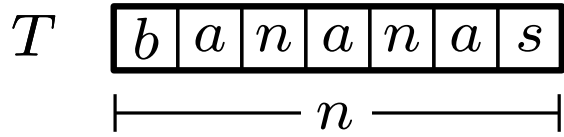


36 nodes

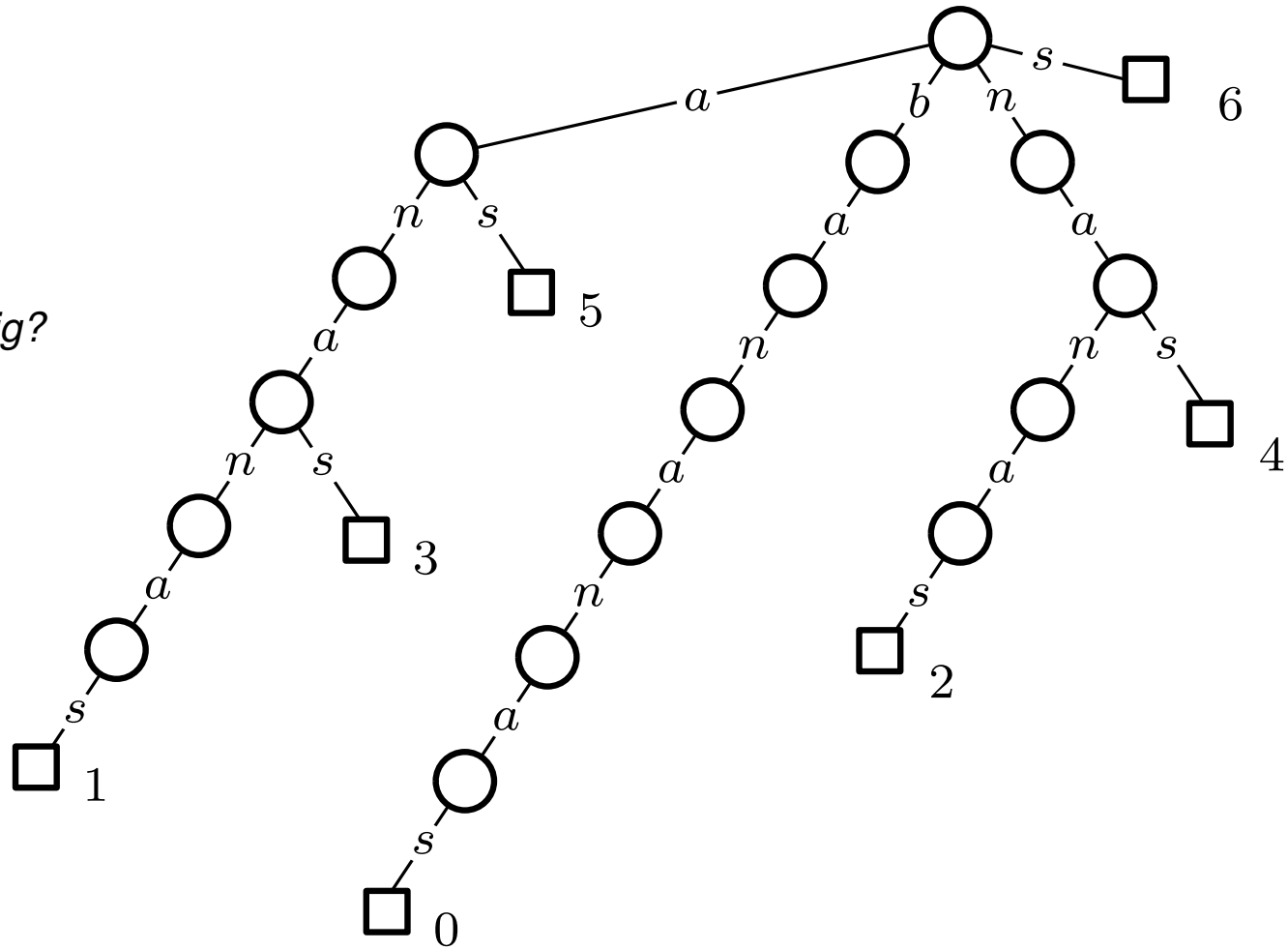
An atomic suffix tree can have  $((n/2) + 1)^2$  nodes

*this is far too big!*

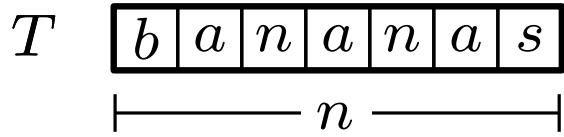
# Compacted suffix trees



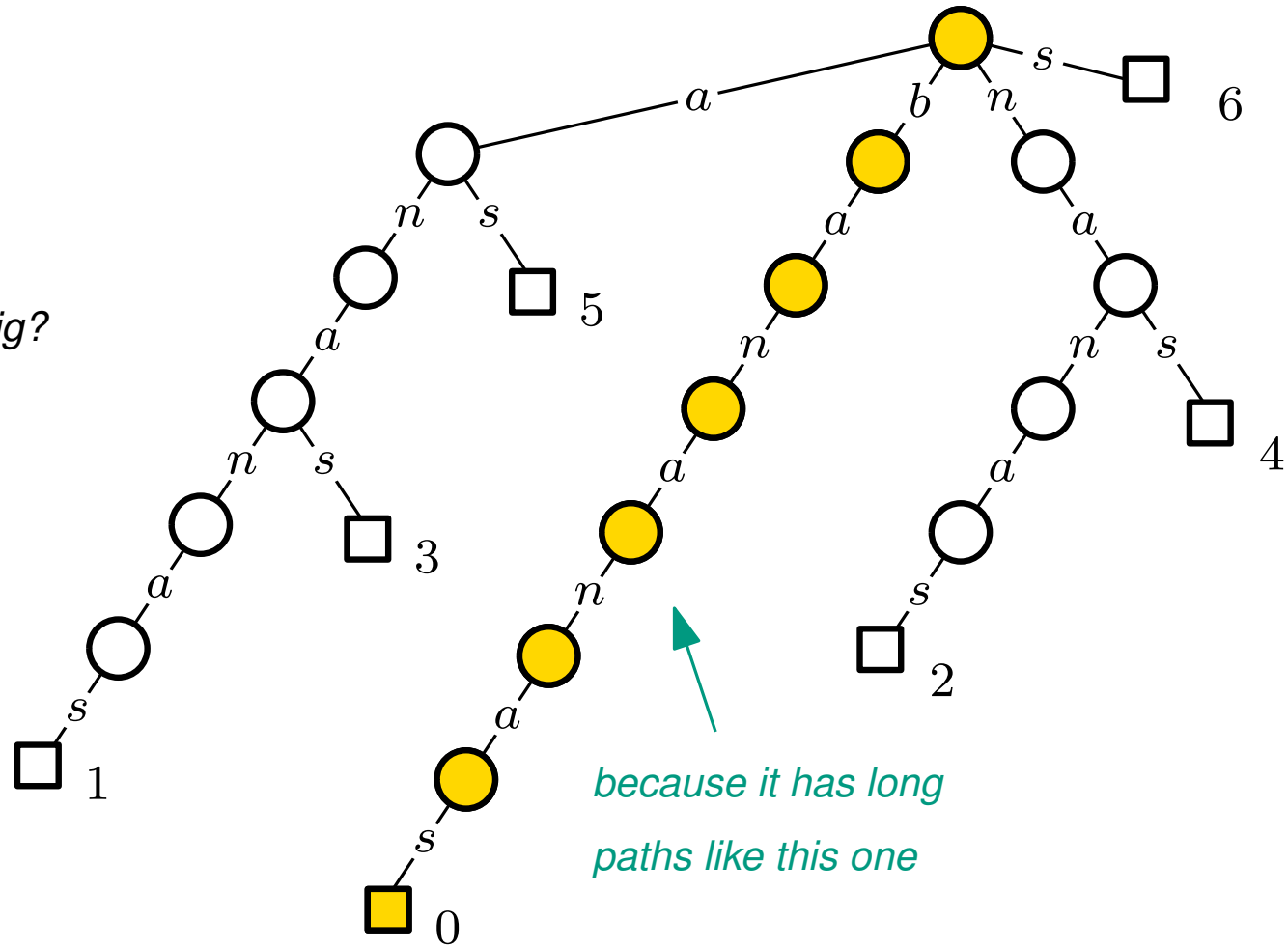
Why is the atomic suffix tree so big?



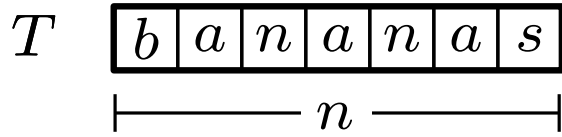
# Compacted suffix trees



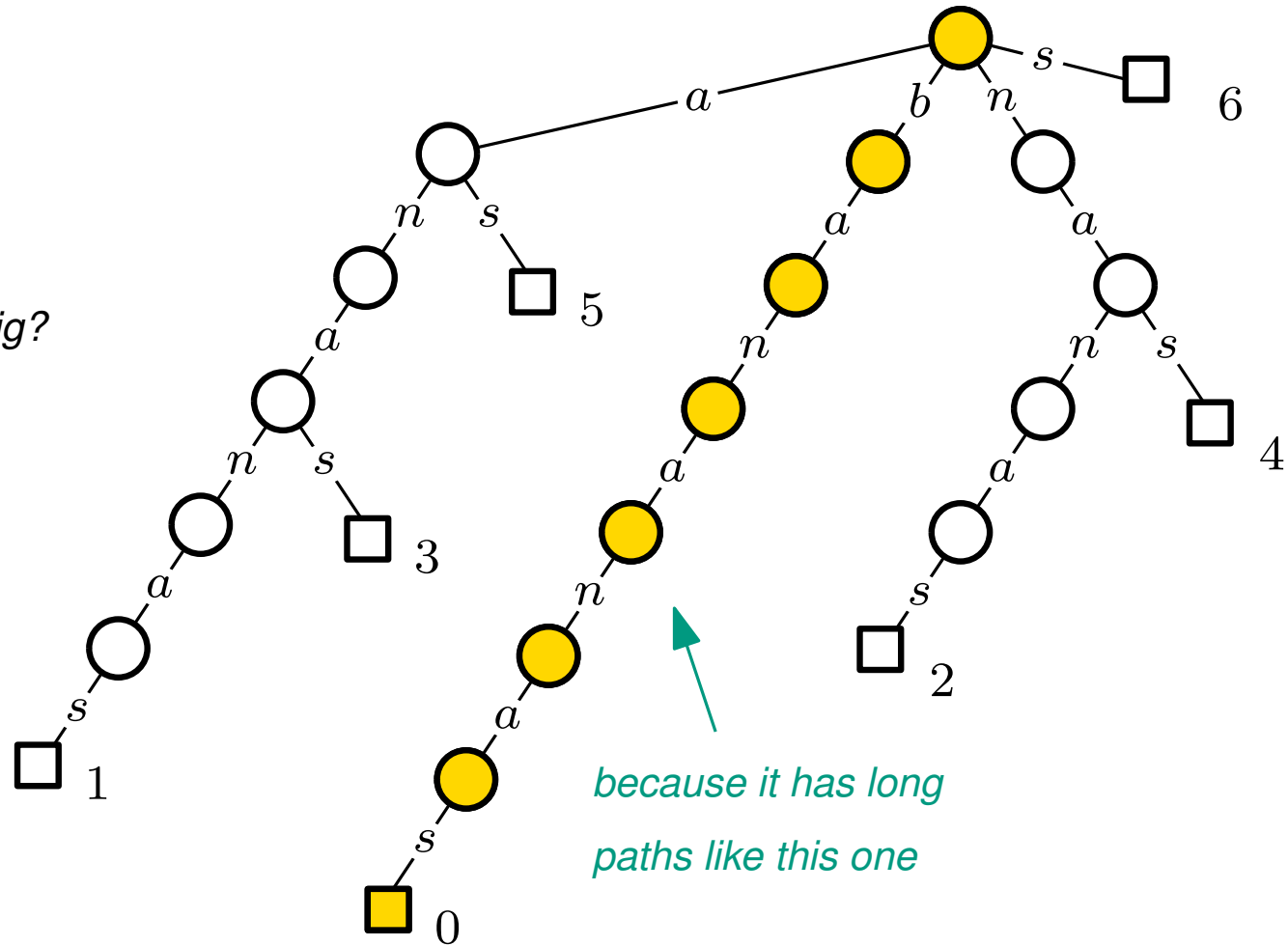
Why is the atomic suffix tree so big?



# Compacted suffix trees

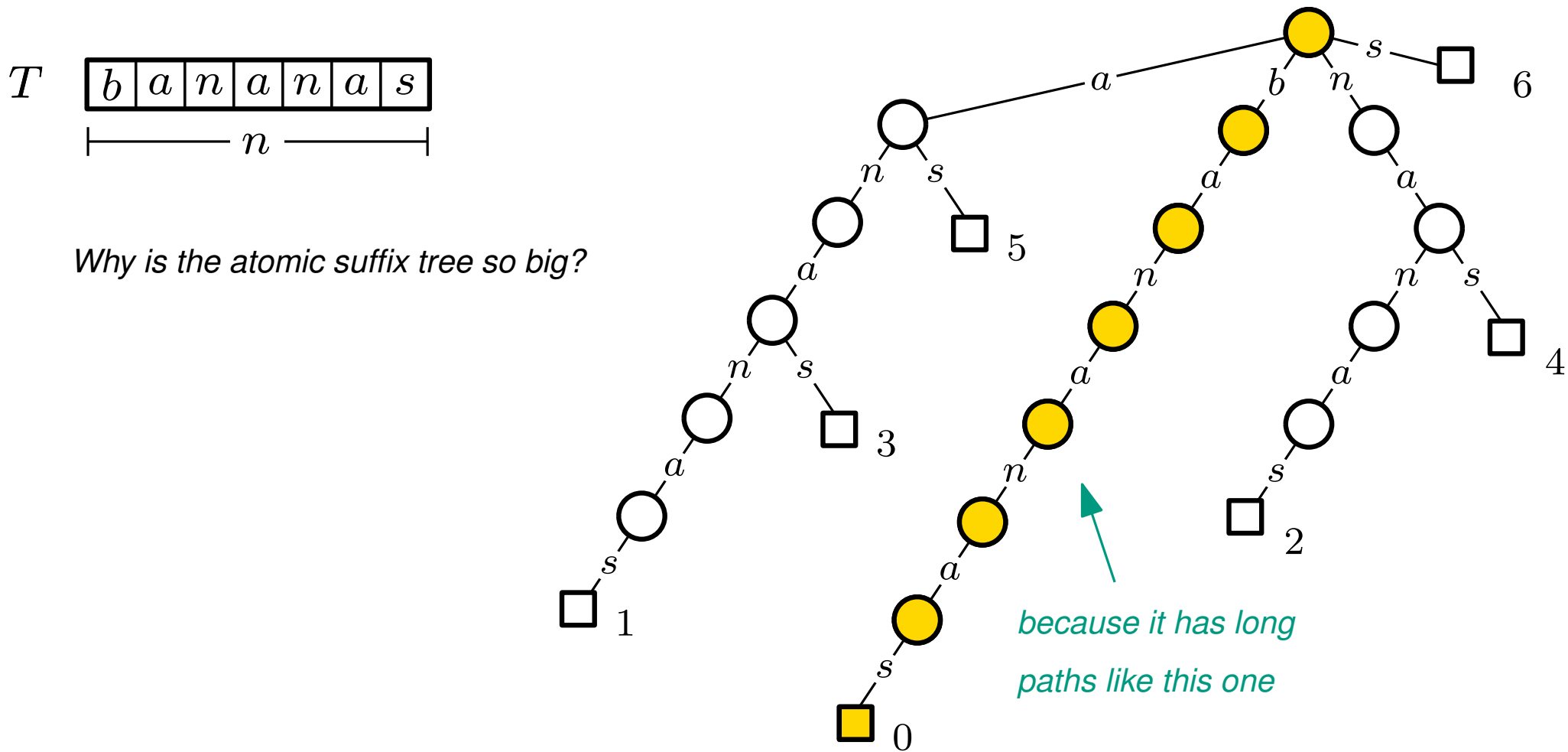


Why is the atomic suffix tree so big?



**Main Idea** replace each non-branching path with a single edge

# Compacted suffix trees

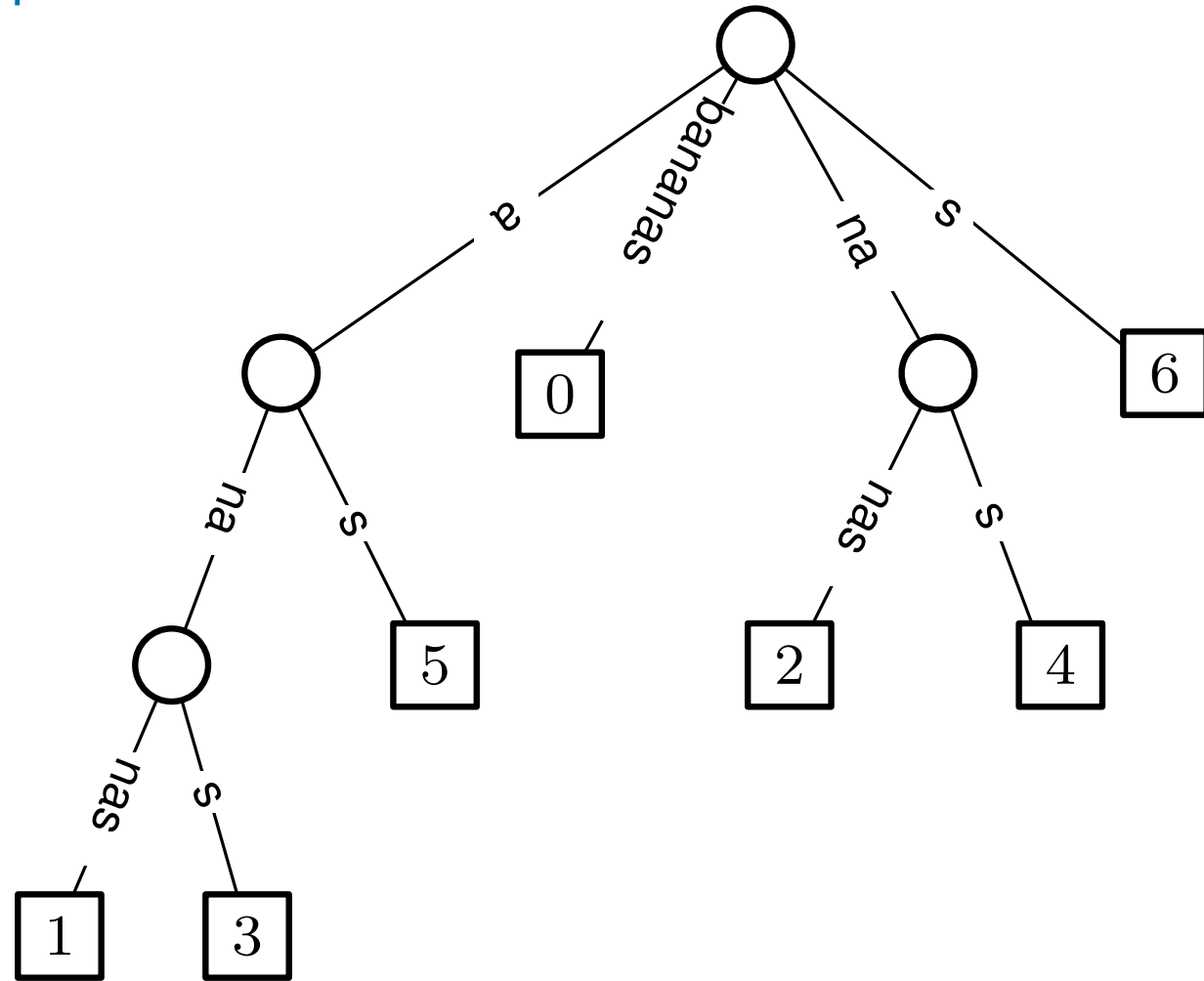
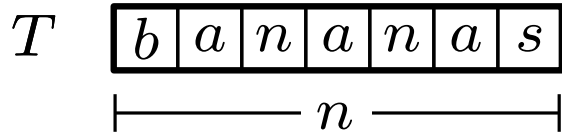


**Main Idea** replace each non-branching path with a single edge

- edges are now labelled with substrings



# Compacted suffix trees



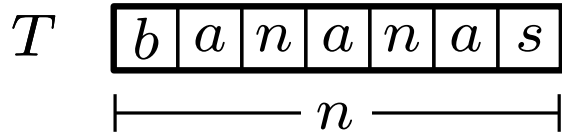
**Main Idea** replace each non-branching path with a single edge

- edges are now labelled with substrings

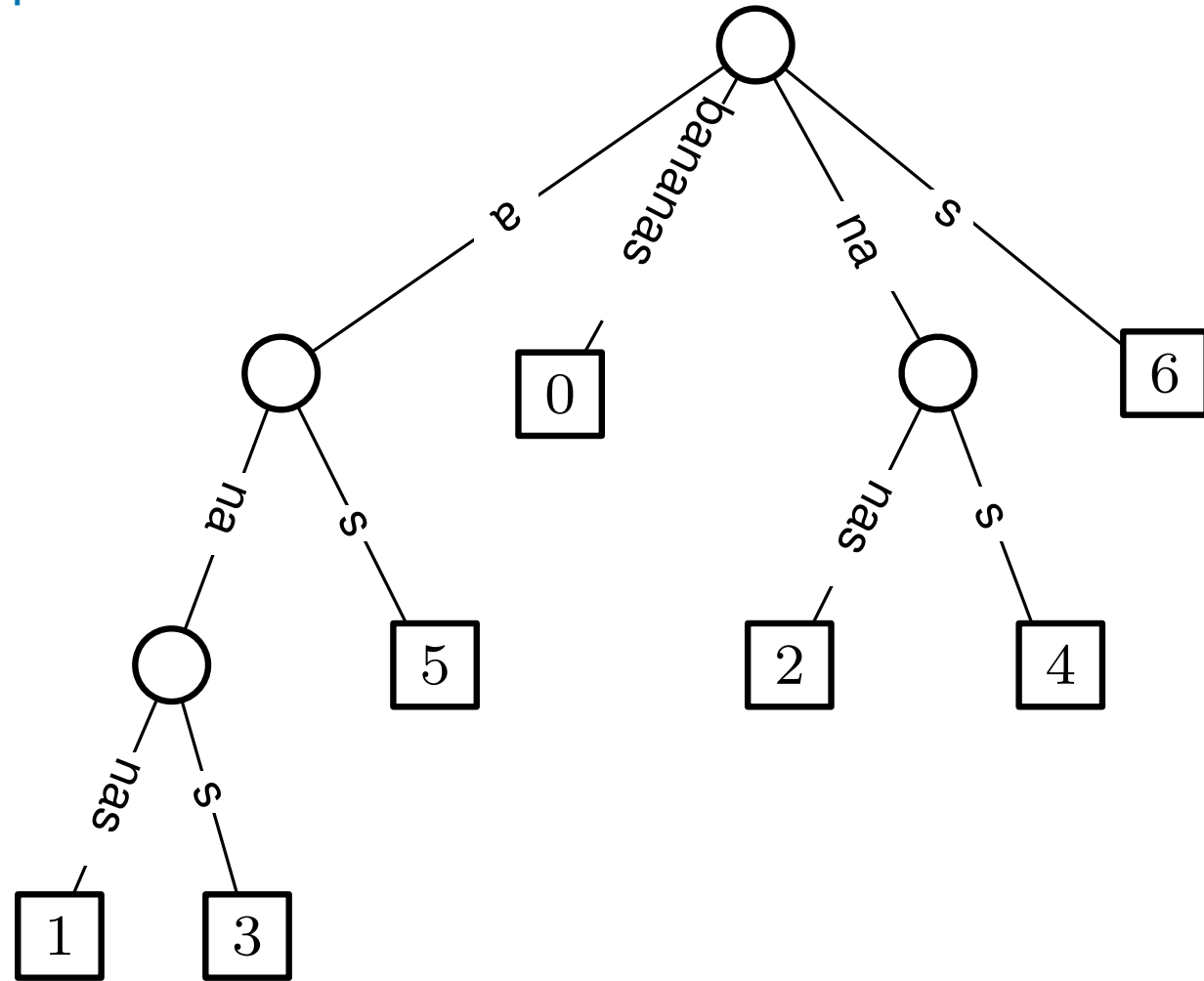
*(instead of single characters)*



# Compacted suffix trees



- There are at most  $n$  leaves

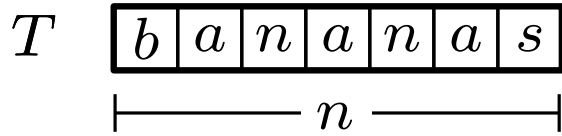


**Main Idea** replace each non-branching path with a single edge

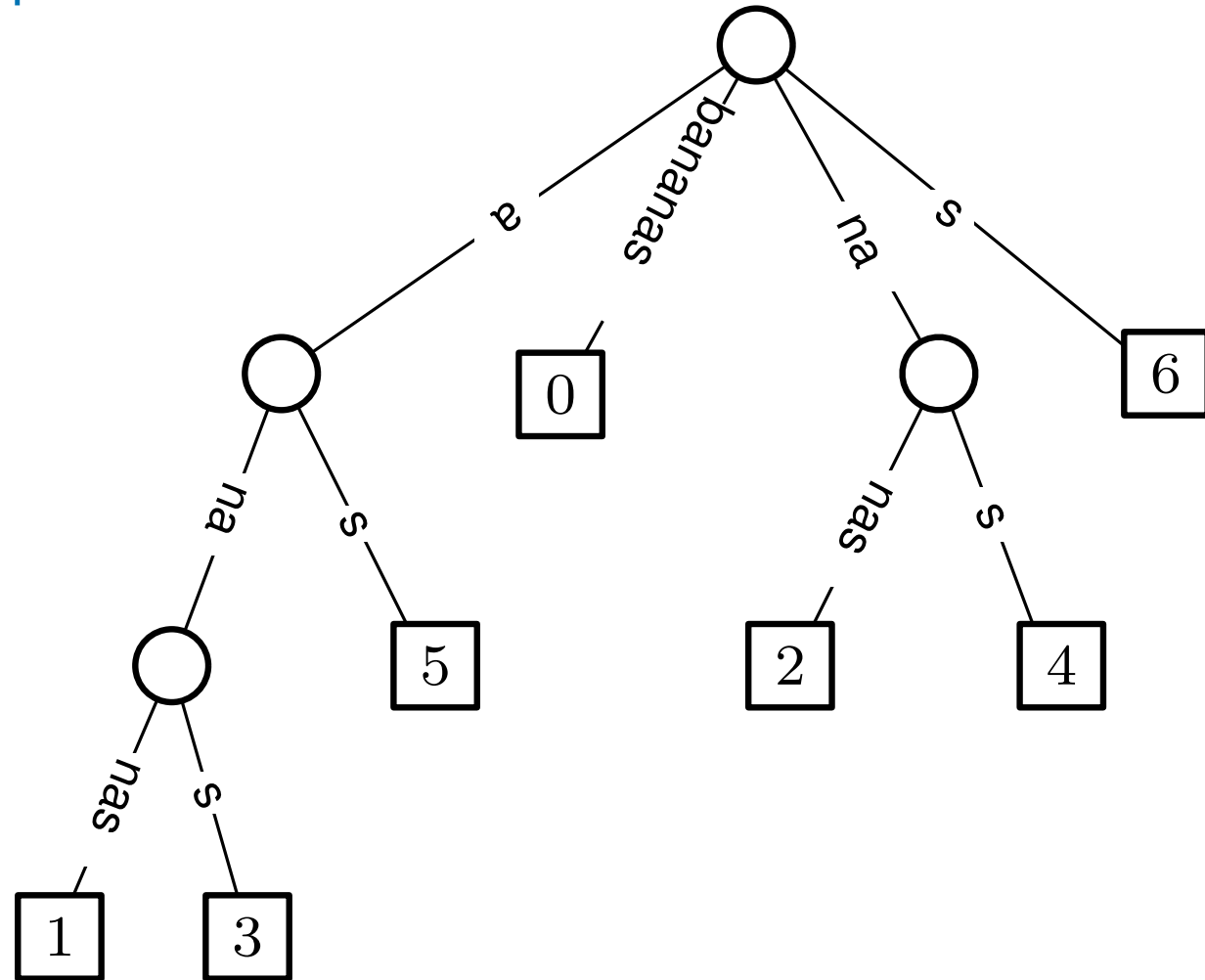
- edges are now labelled with substrings

*(instead of single characters)*

# Compacted suffix trees



- There are at most  $n$  leaves
- Every internal node has two or more children

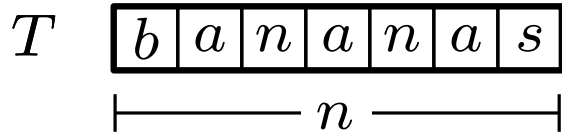


**Main Idea** replace each non-branching path with a single edge

- edges are now labelled with substrings

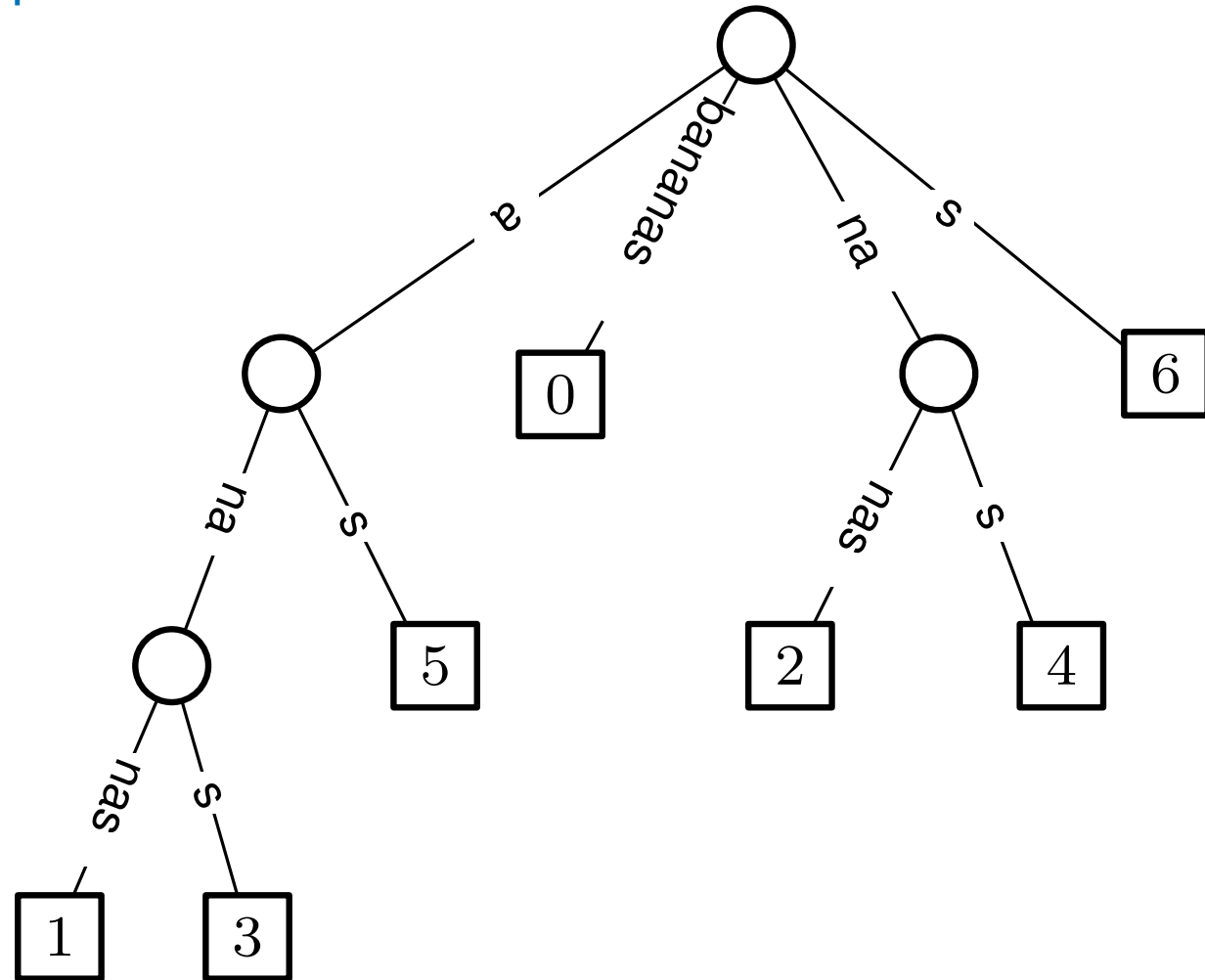
*(instead of single characters)*

# Compacted suffix trees



- There are at most  $n$  leaves
- Every internal node has two or more children

so there are  $O(n)$  edges

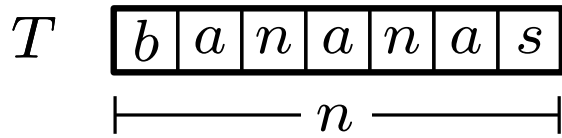


**Main Idea** replace each non-branching path with a single edge

- edges are now labelled with substrings

*(instead of single characters)*

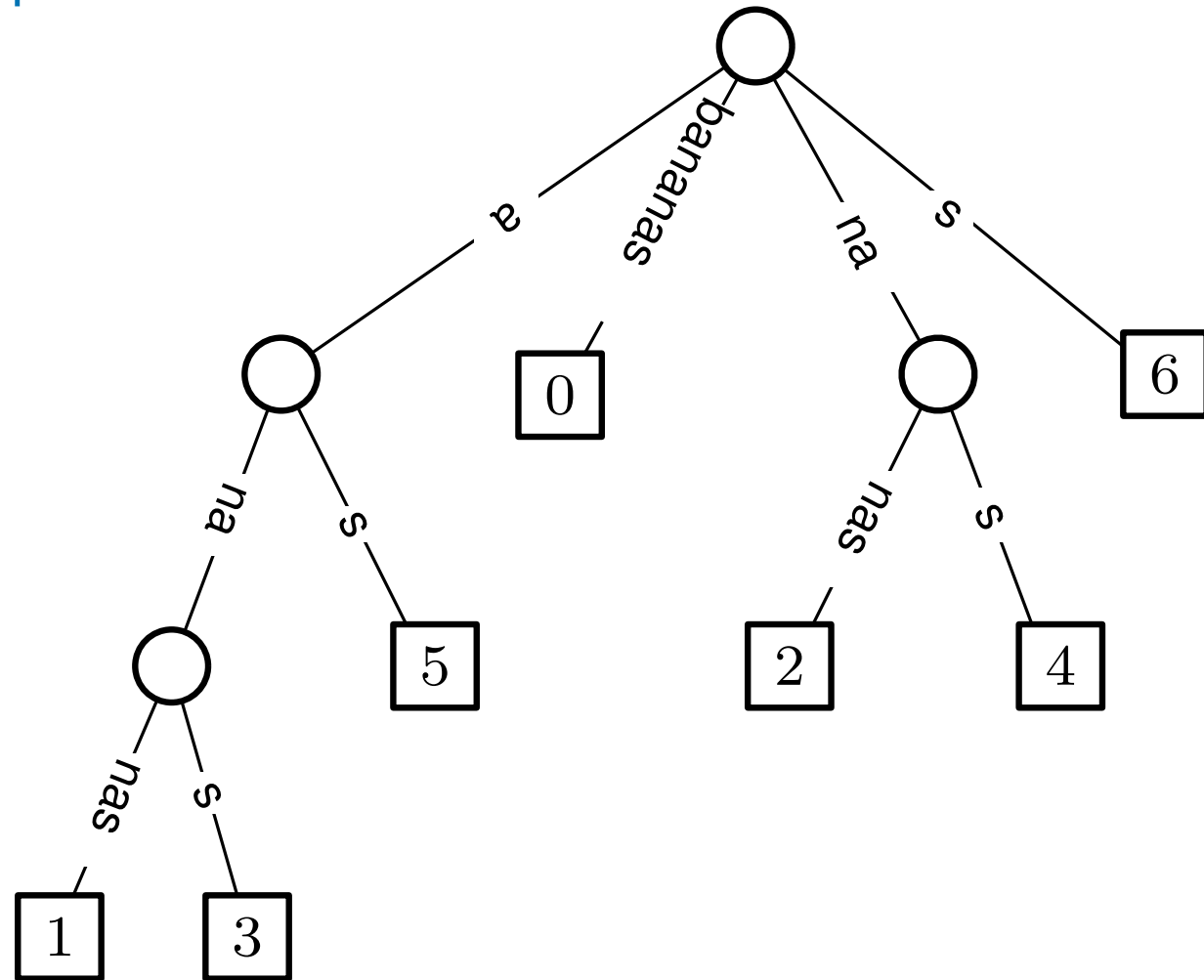
# Compacted suffix trees



- There are at most  $n$  leaves
- Every internal node has two or more children

so there are  $O(n)$  edges

don't the edges take up lots of space?

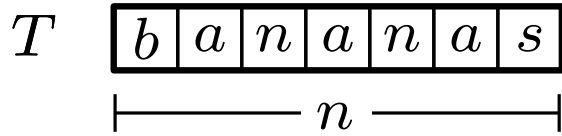


**Main Idea** replace each non-branching path with a single edge

- edges are now labelled with substrings

*(instead of single characters)*

# Compacted suffix trees

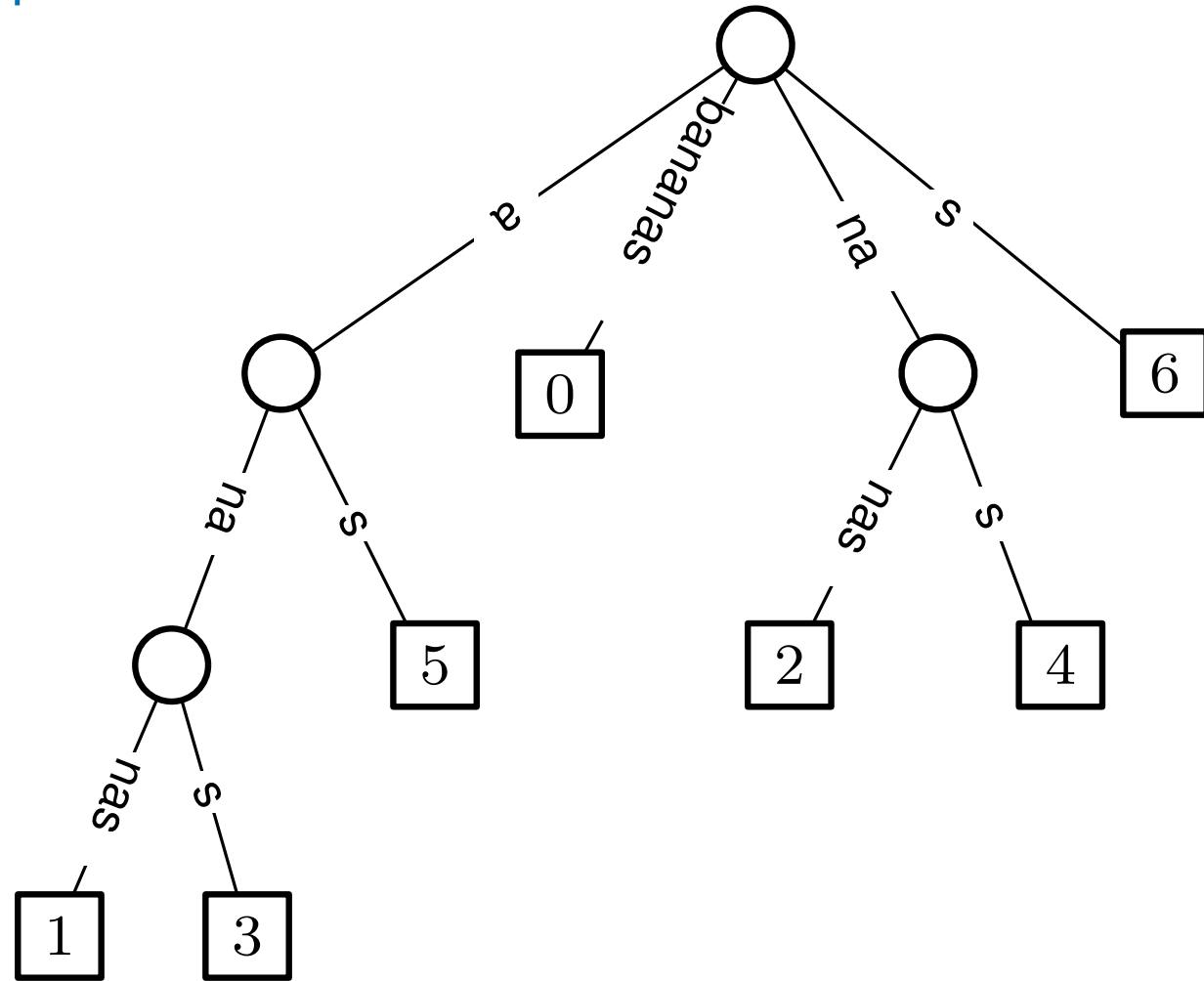


- There are at most  $n$  leaves
- Every internal node has two or more children

so there are  $O(n)$  edges

don't the edges take up lots of space?

we only store the end points

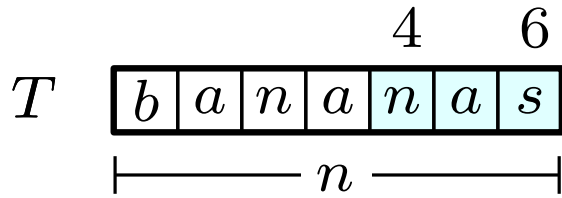


**Main Idea** replace each non-branching path with a single edge

- edges are now labelled with substrings

*(instead of single characters)*

# Compacted suffix trees

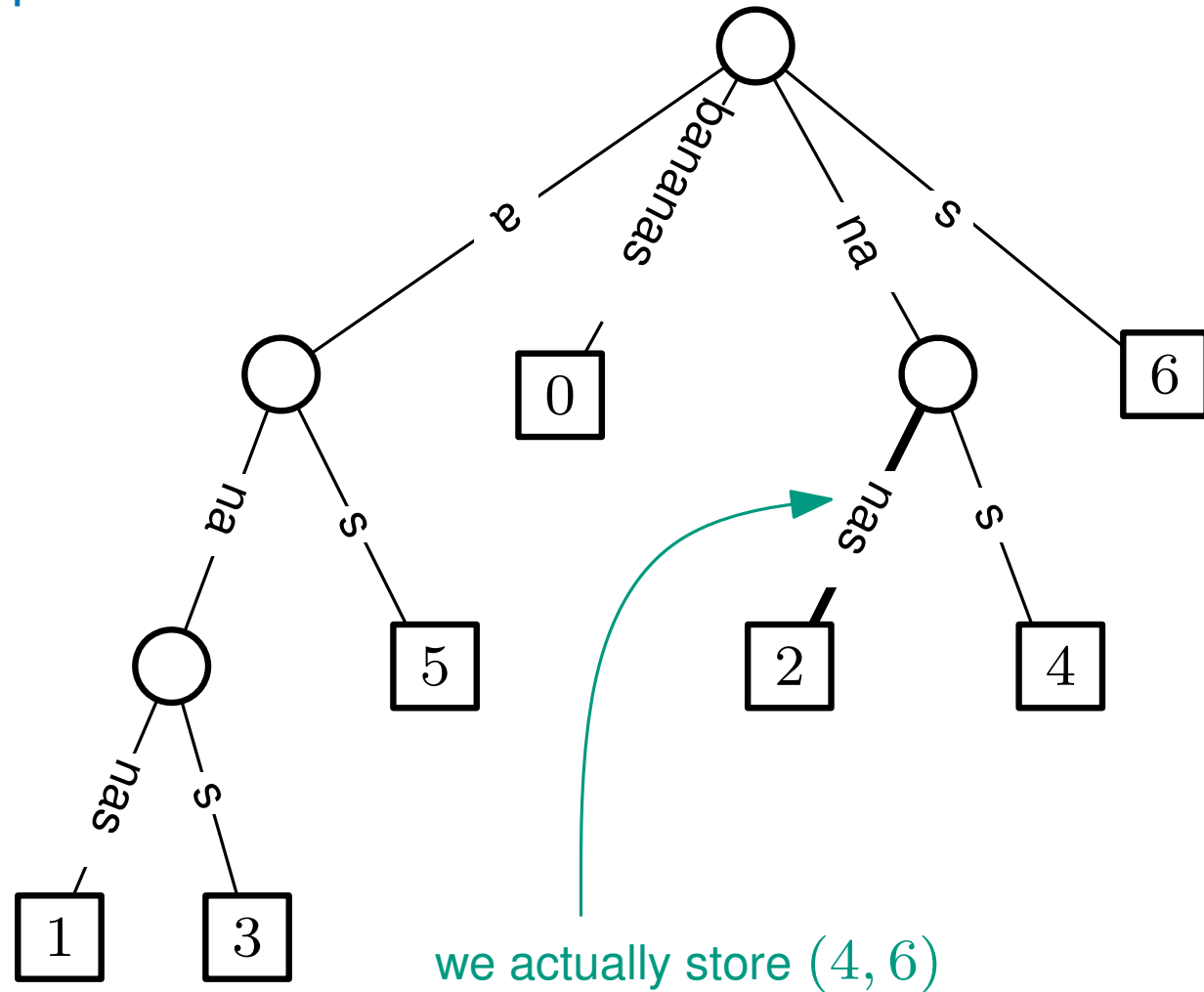


- There are at most  $n$  leaves
- Every internal node has two or more children

so there are  $O(n)$  edges

don't the edges take up lots of space?

we only store the end points

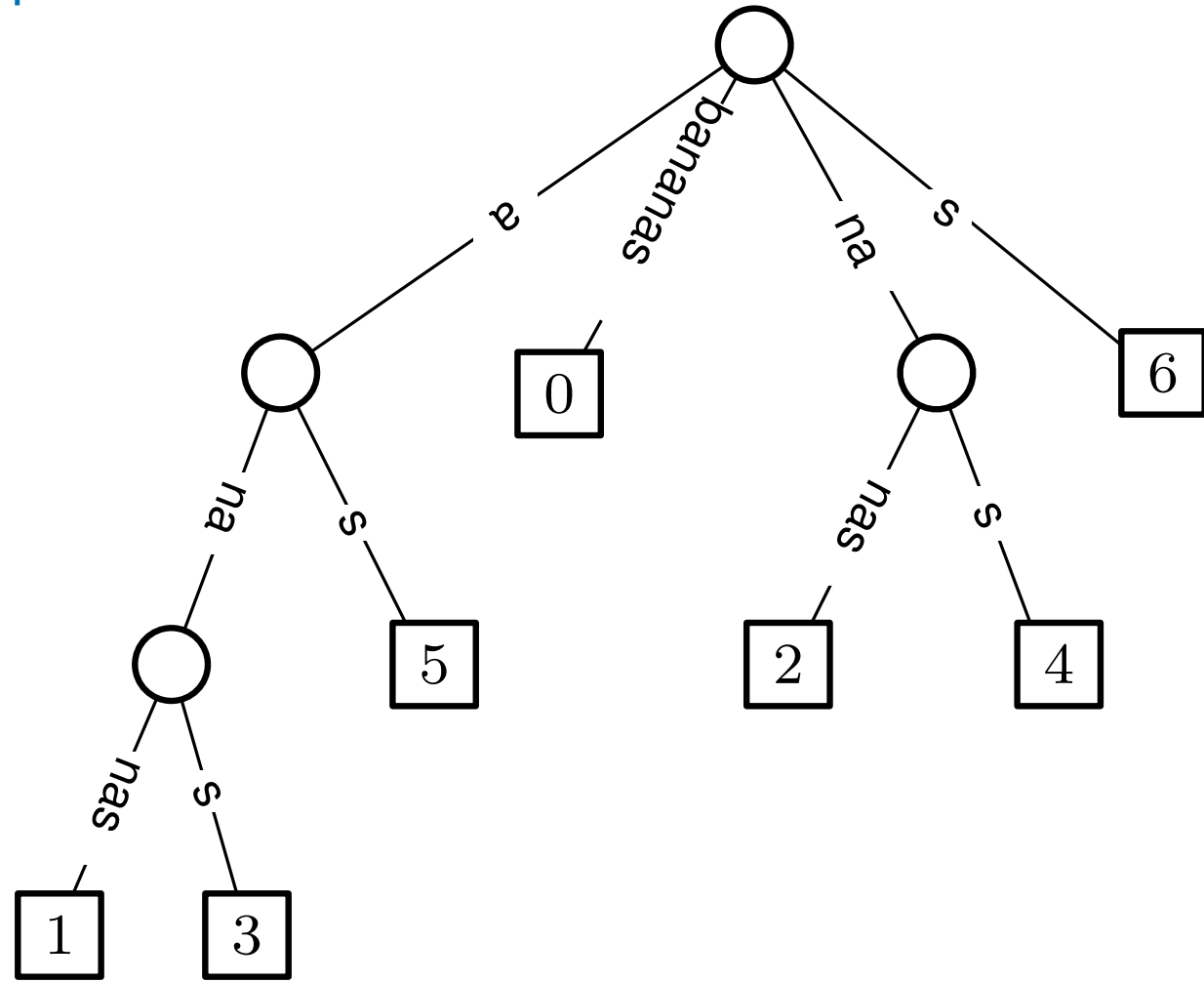
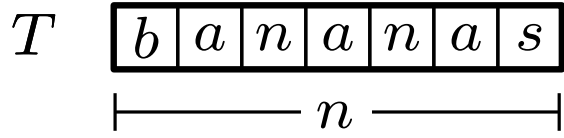


**Main Idea** replace each non-branching path with a single edge

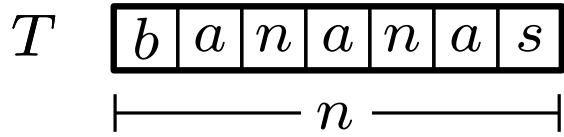
- edges are now labelled with substrings

*(instead of single characters)*

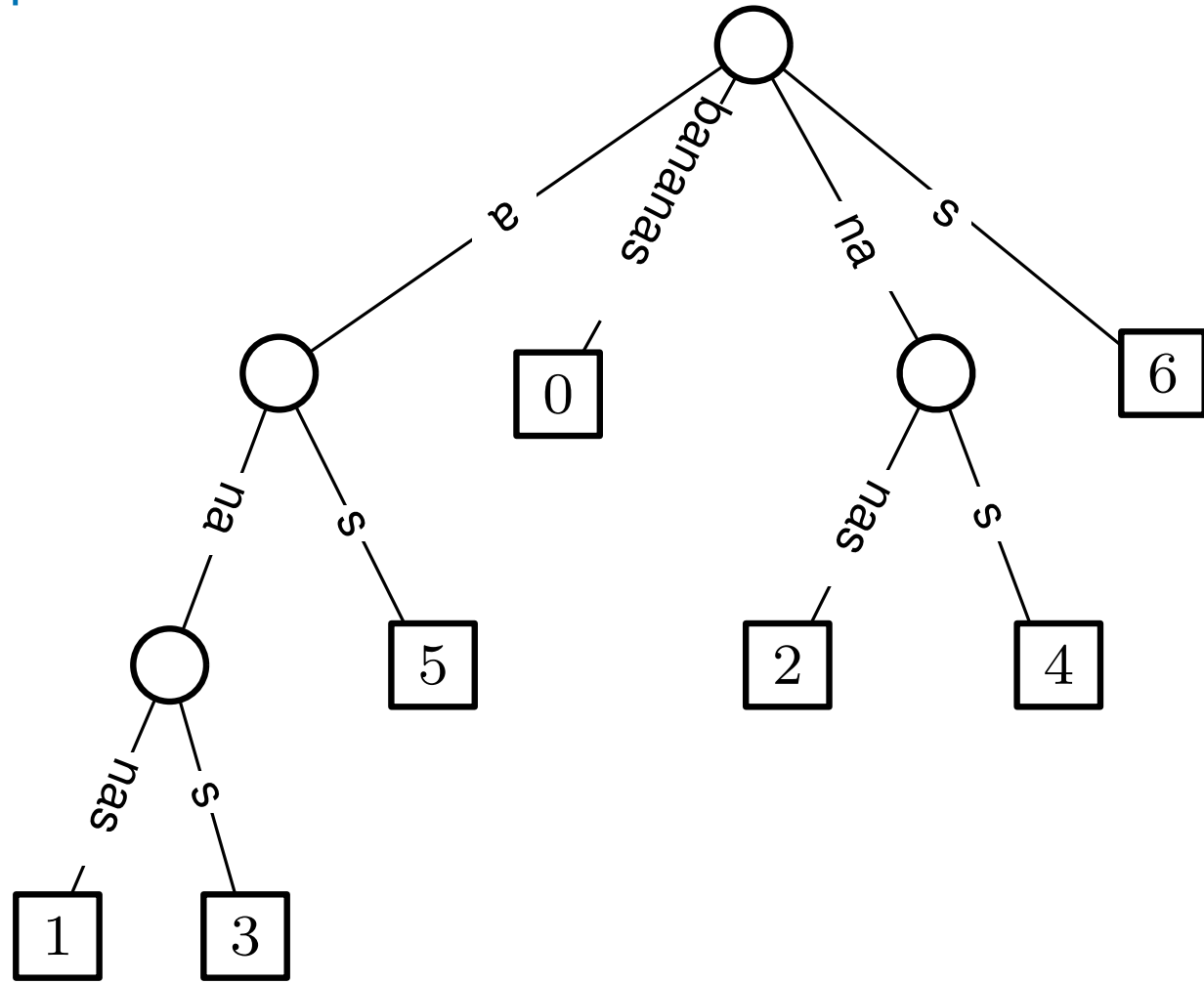
# Compacted suffix trees



# Compacted suffix trees

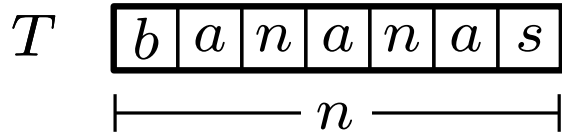


**Compacted Suffix Tree of  $T$**



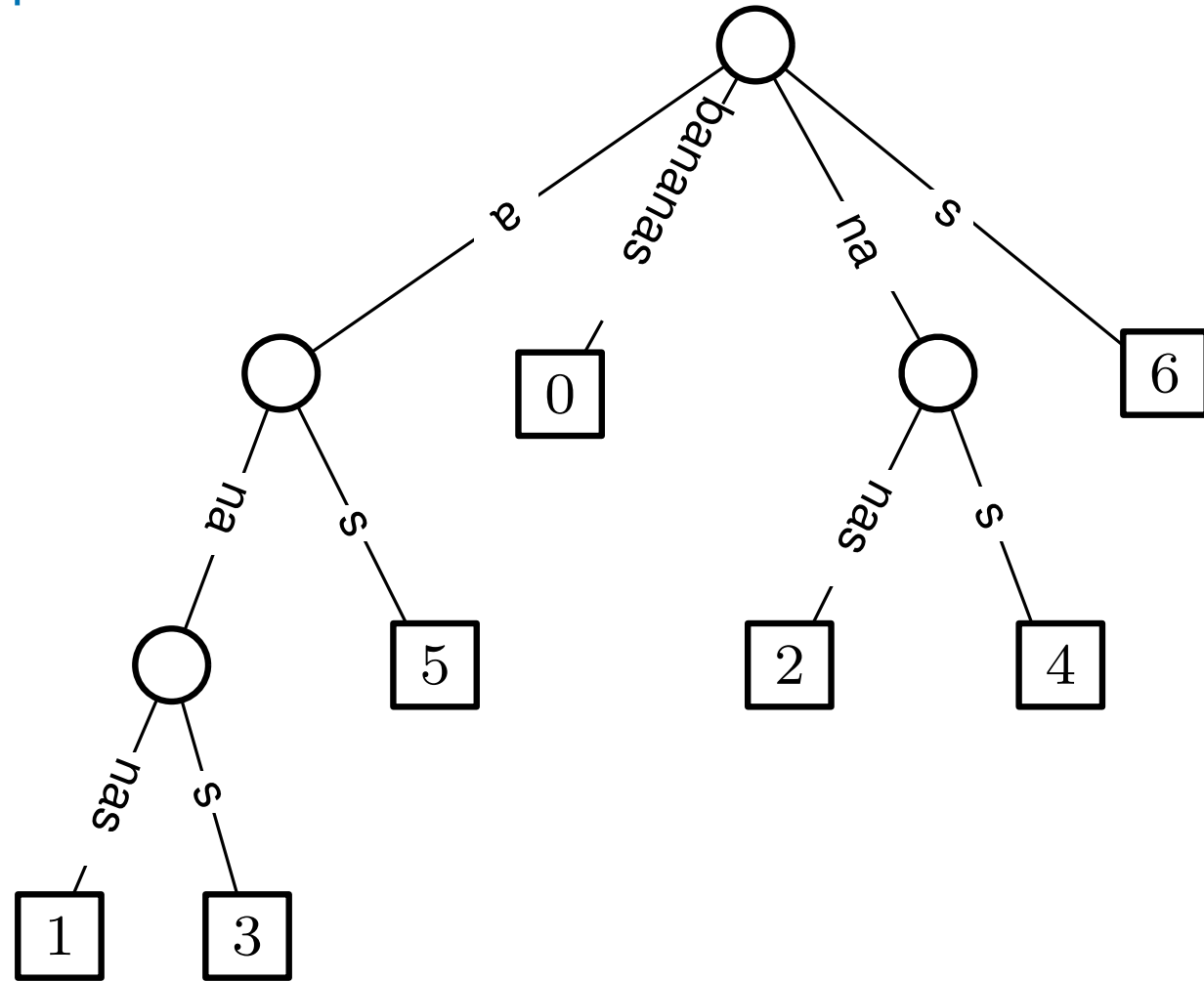


# Compacted suffix trees

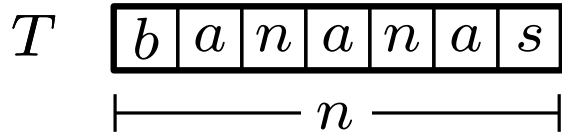


**Compacted Suffix Tree of  $T$**

- A rooted tree with  $n$  leaves

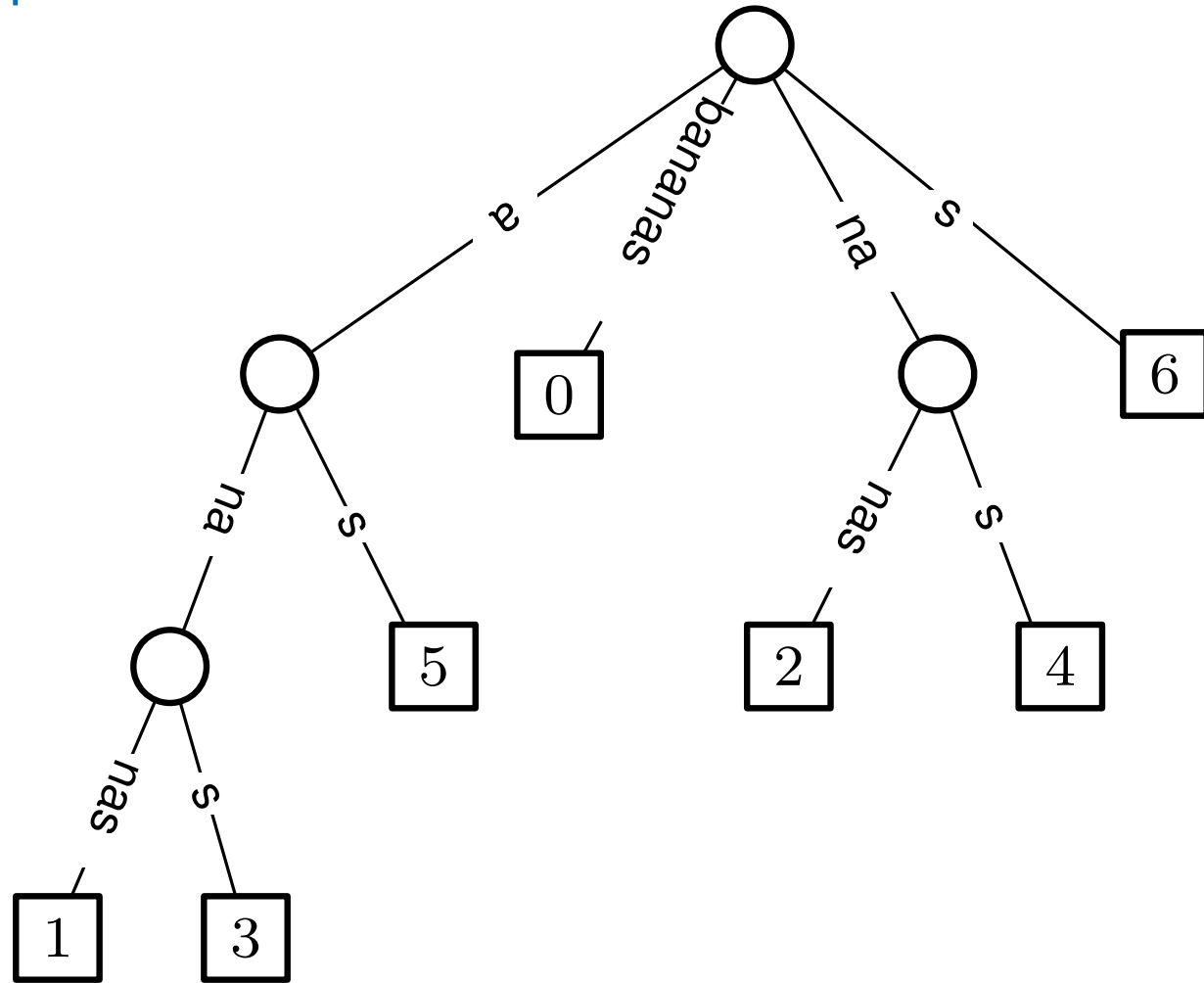


# Compacted suffix trees

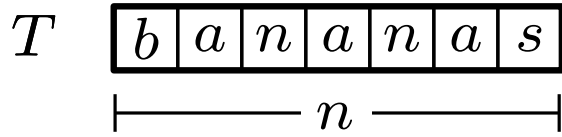


## Compacted Suffix Tree of $T$

- A rooted tree with  $n$  leaves
- Every internal node has two or more children

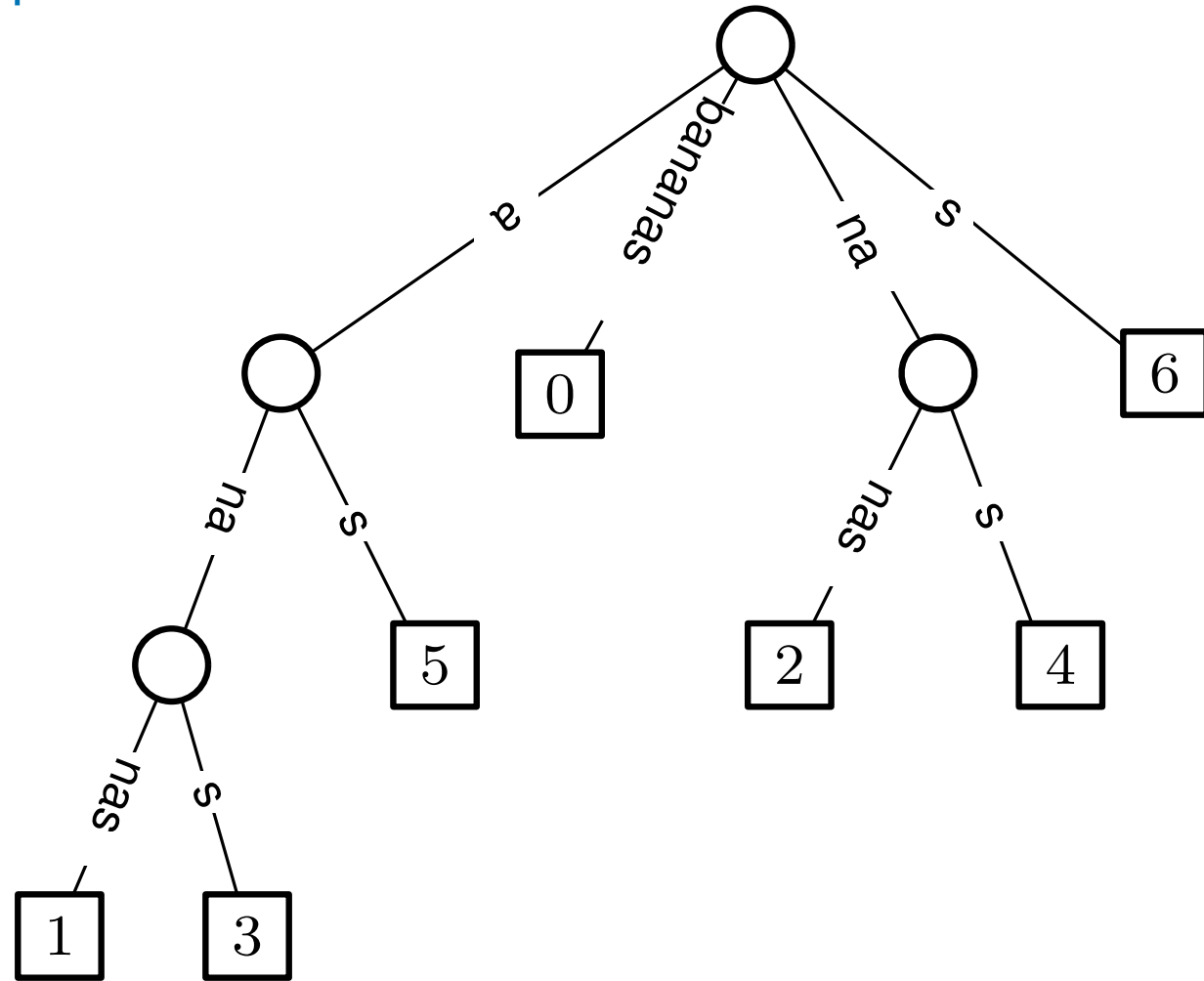


# Compacted suffix trees

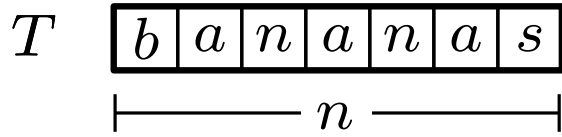


## Compacted Suffix Tree of $T$

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring

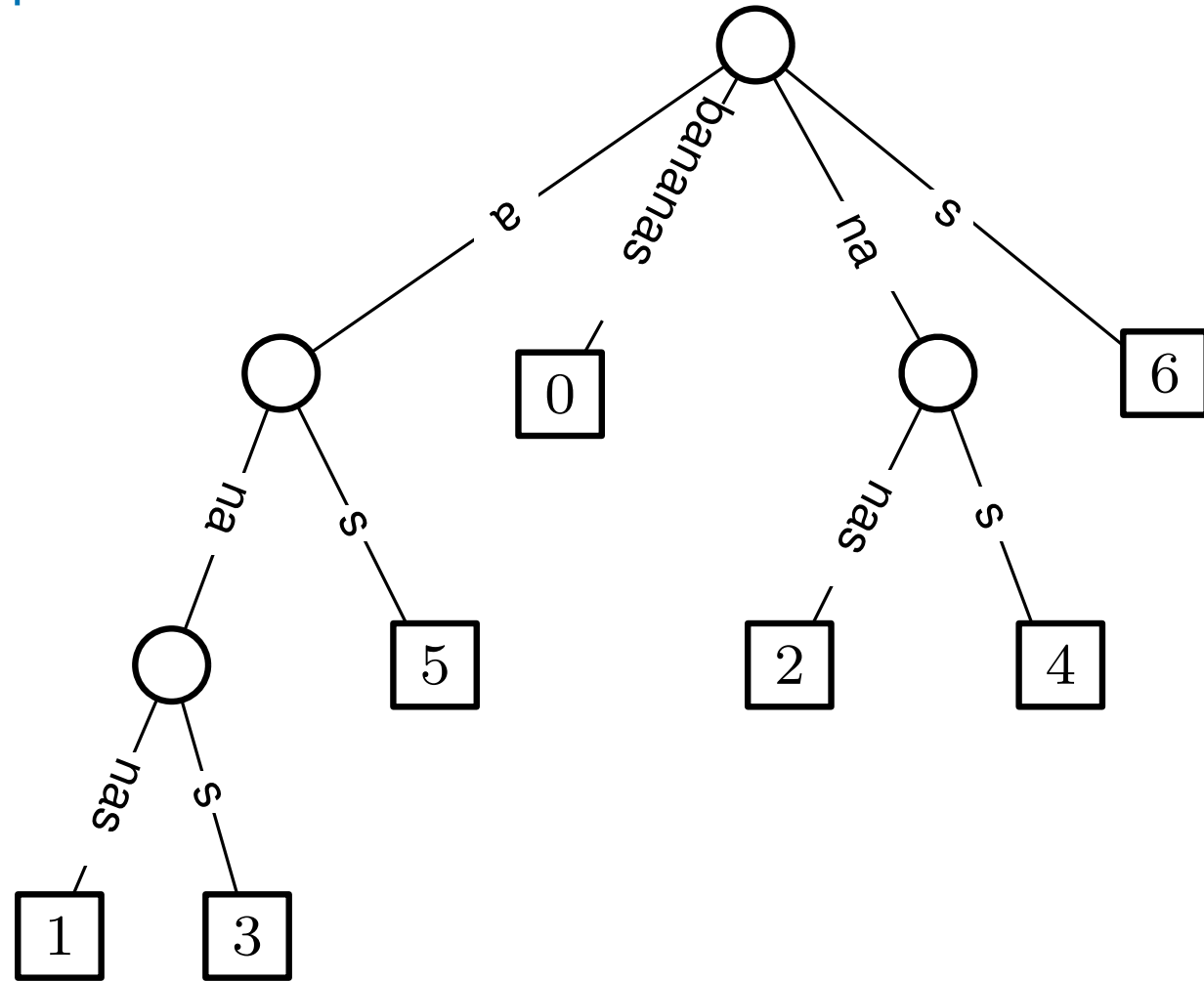


# Compacted suffix trees

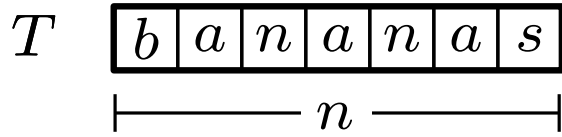


## Compacted Suffix Tree of $T$

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring
- No two edges leaving the same node have the same first character

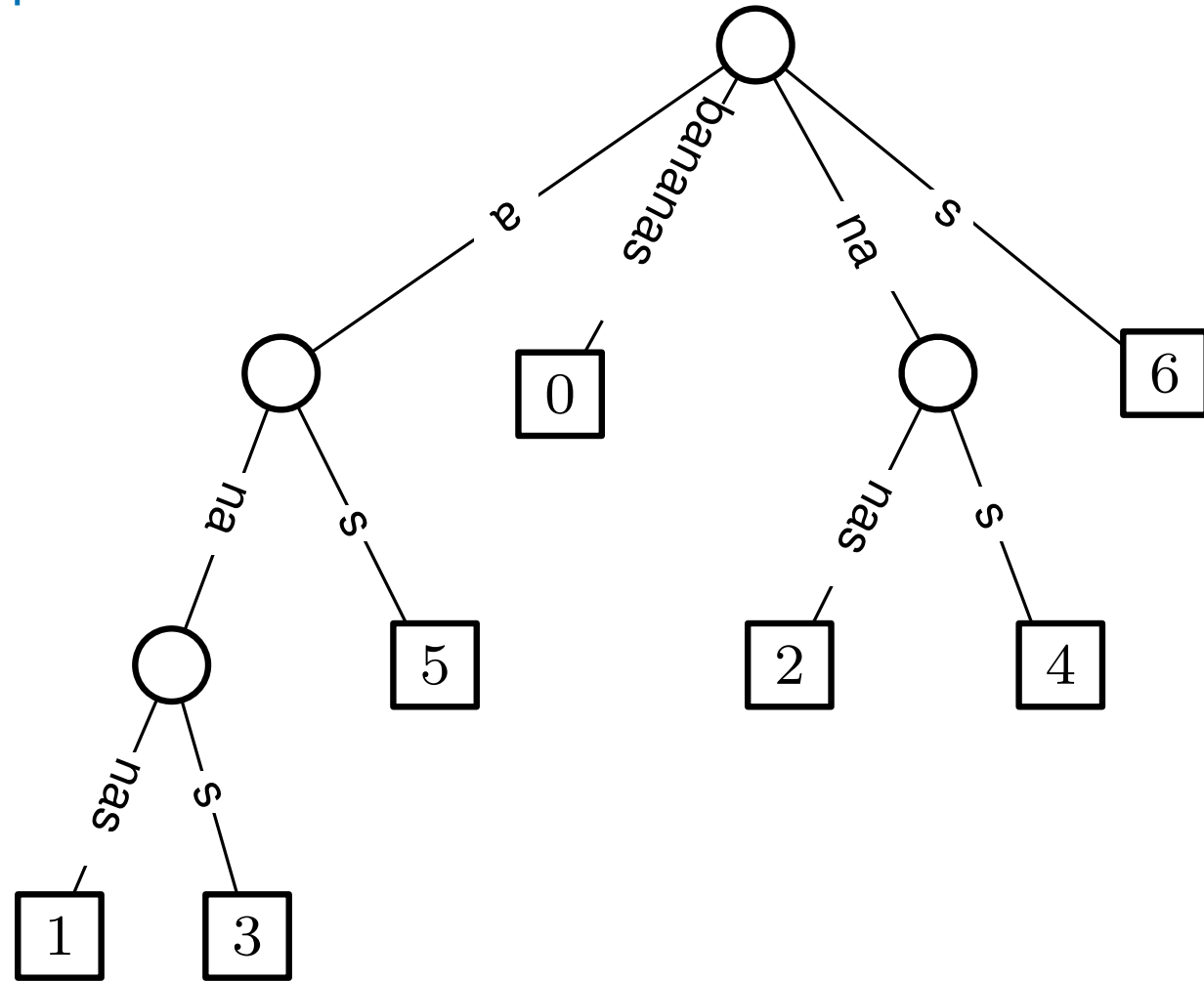


# Compacted suffix trees



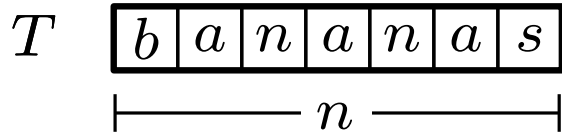
**Compacted Suffix Tree of  $T$**

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring



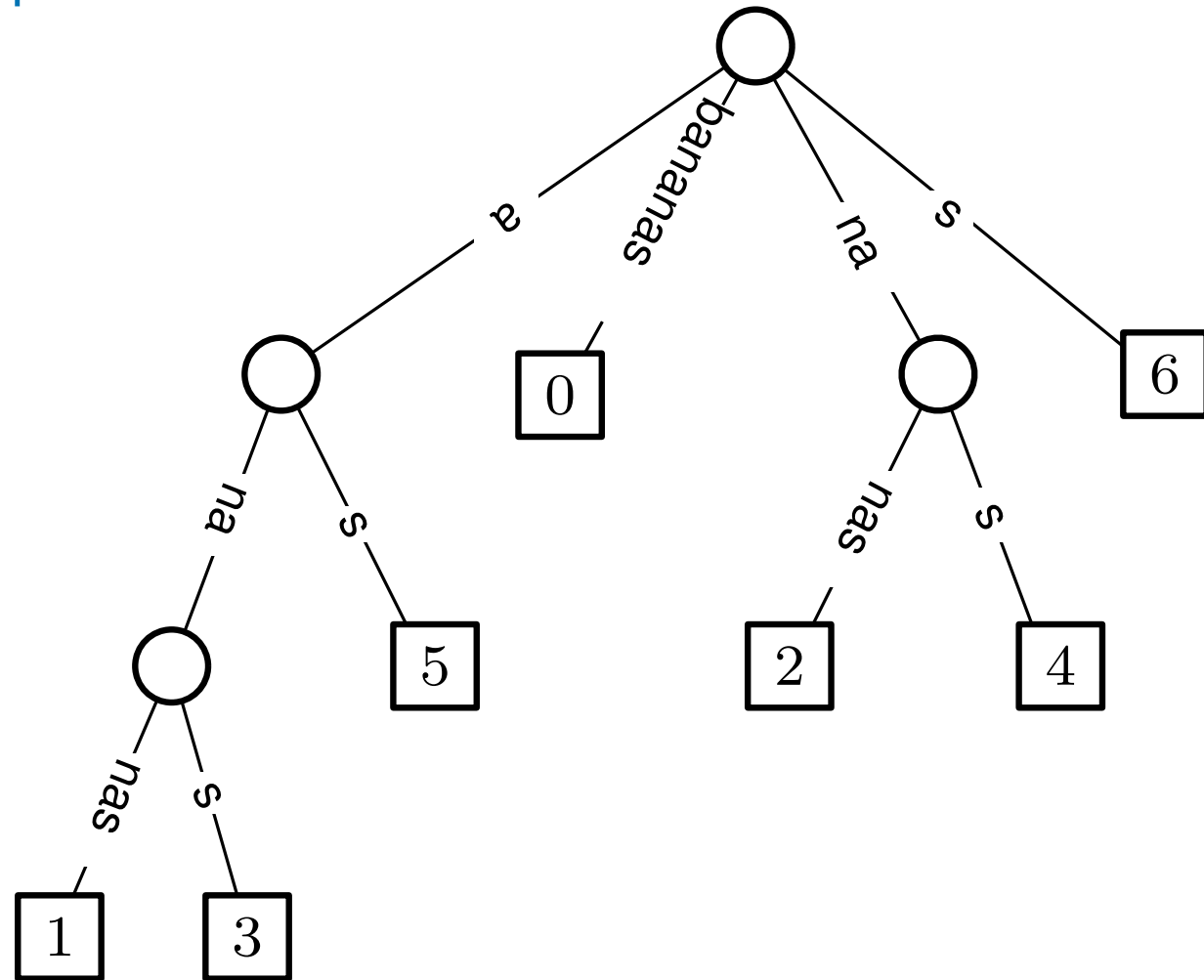
- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$

# Compacted suffix trees



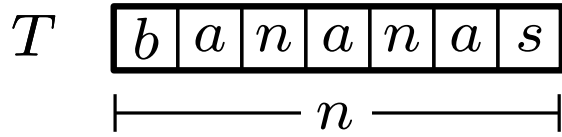
**Compacted Suffix Tree of  $T$**

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring



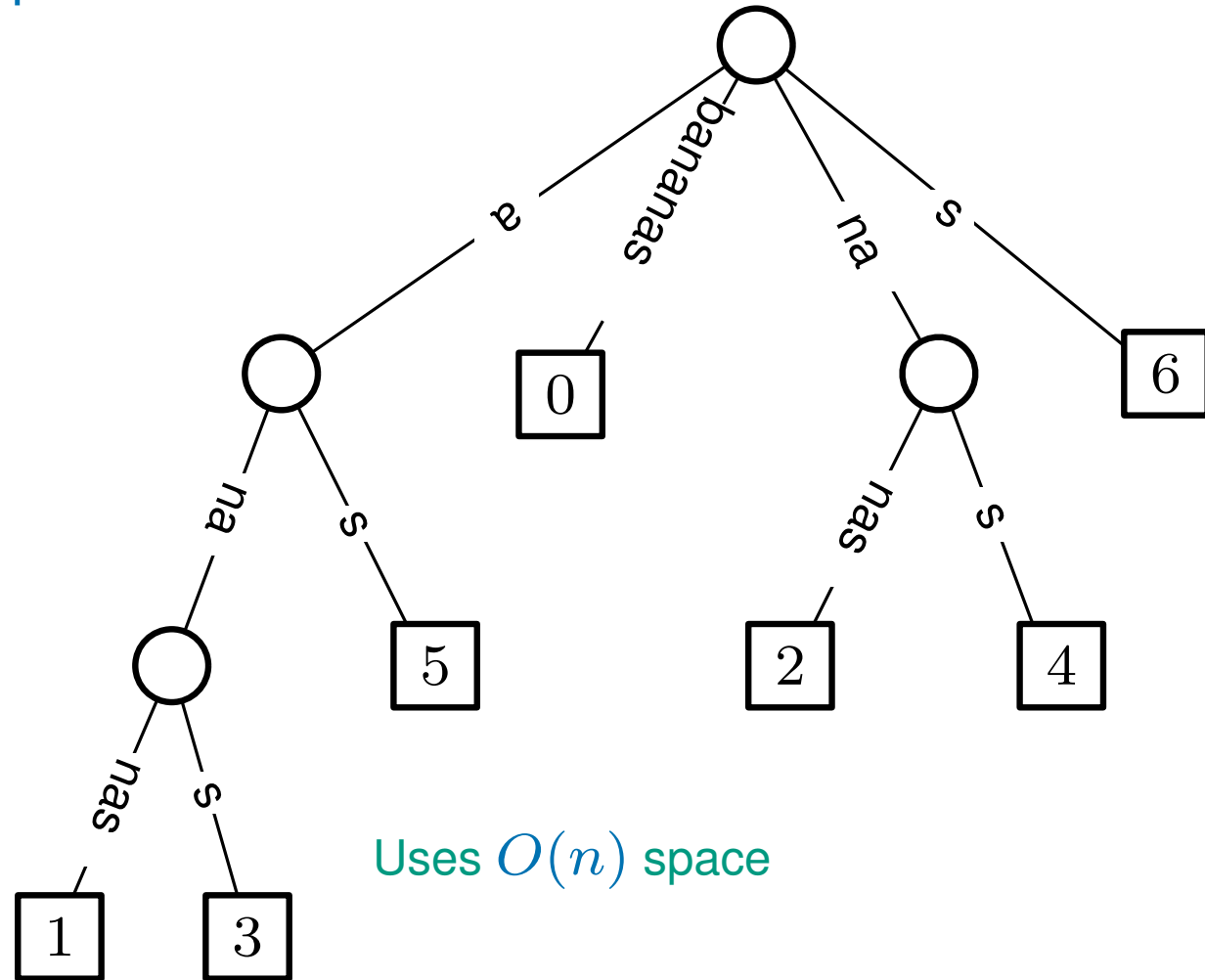
- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix

# Compacted suffix trees



## Compacted Suffix Tree of $T$

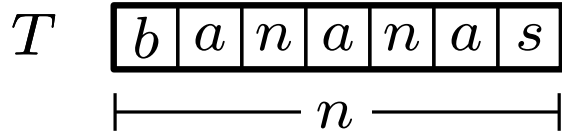
- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring



Uses  $O(n)$  space

- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix

# Compacted suffix trees

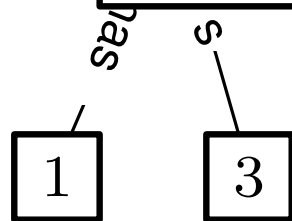


**Sanity Check**

Does the compacted suffix tree always exist?

## Compacted Suffix Tree of $T$

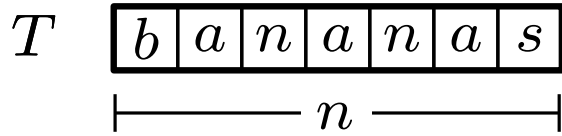
- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring
- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix



Uses  $O(n)$  space



# Compacted suffix trees



## Compacted Suffix Tree of $T$

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring
- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix

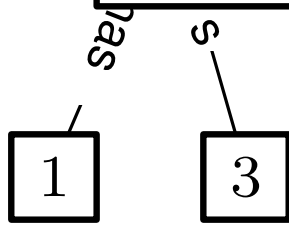
**Sanity Check**

Does the compacted suffix tree always exist?

$T$ 

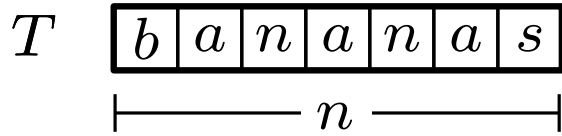
$b$	$b$
-----	-----

*this doesn't have  $n$  leaves*



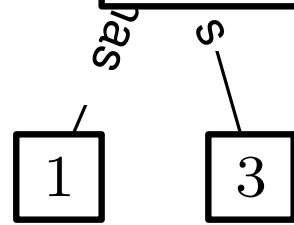
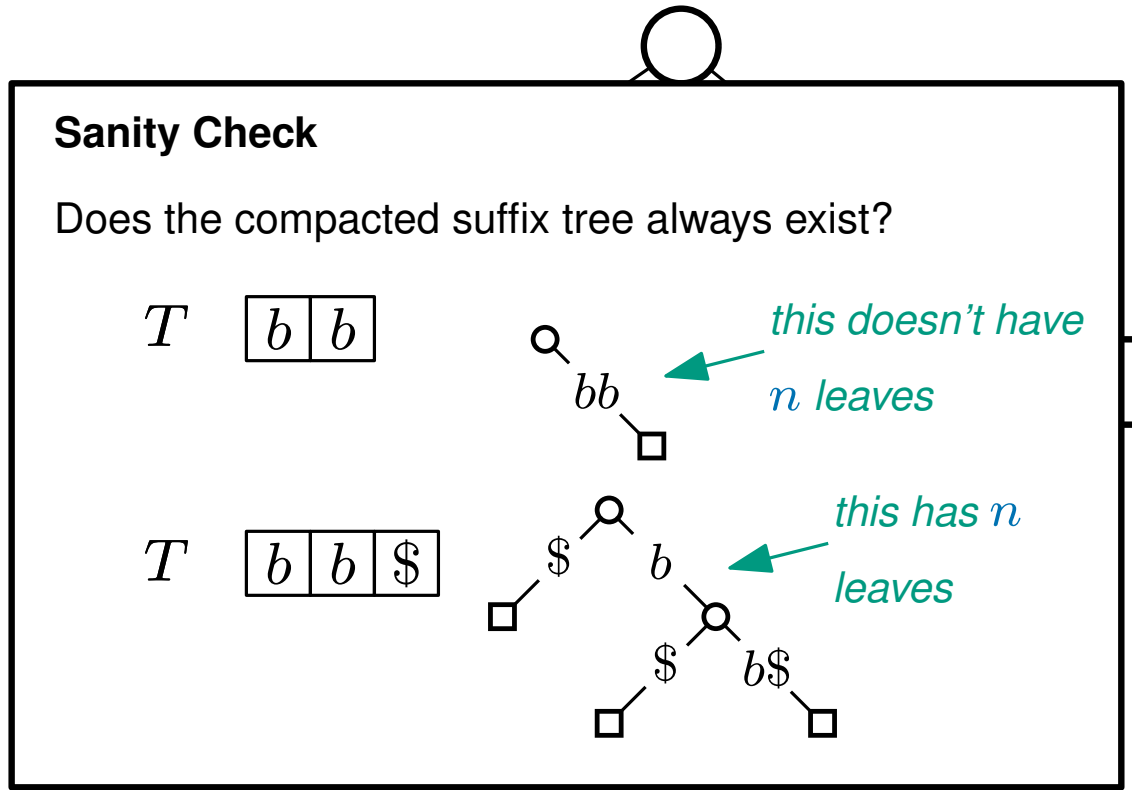
Uses  $O(n)$  space

# Compacted suffix trees



## Compacted Suffix Tree of $T$

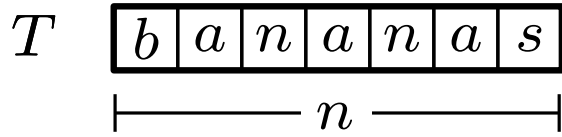
- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring



Uses  $O(n)$  space

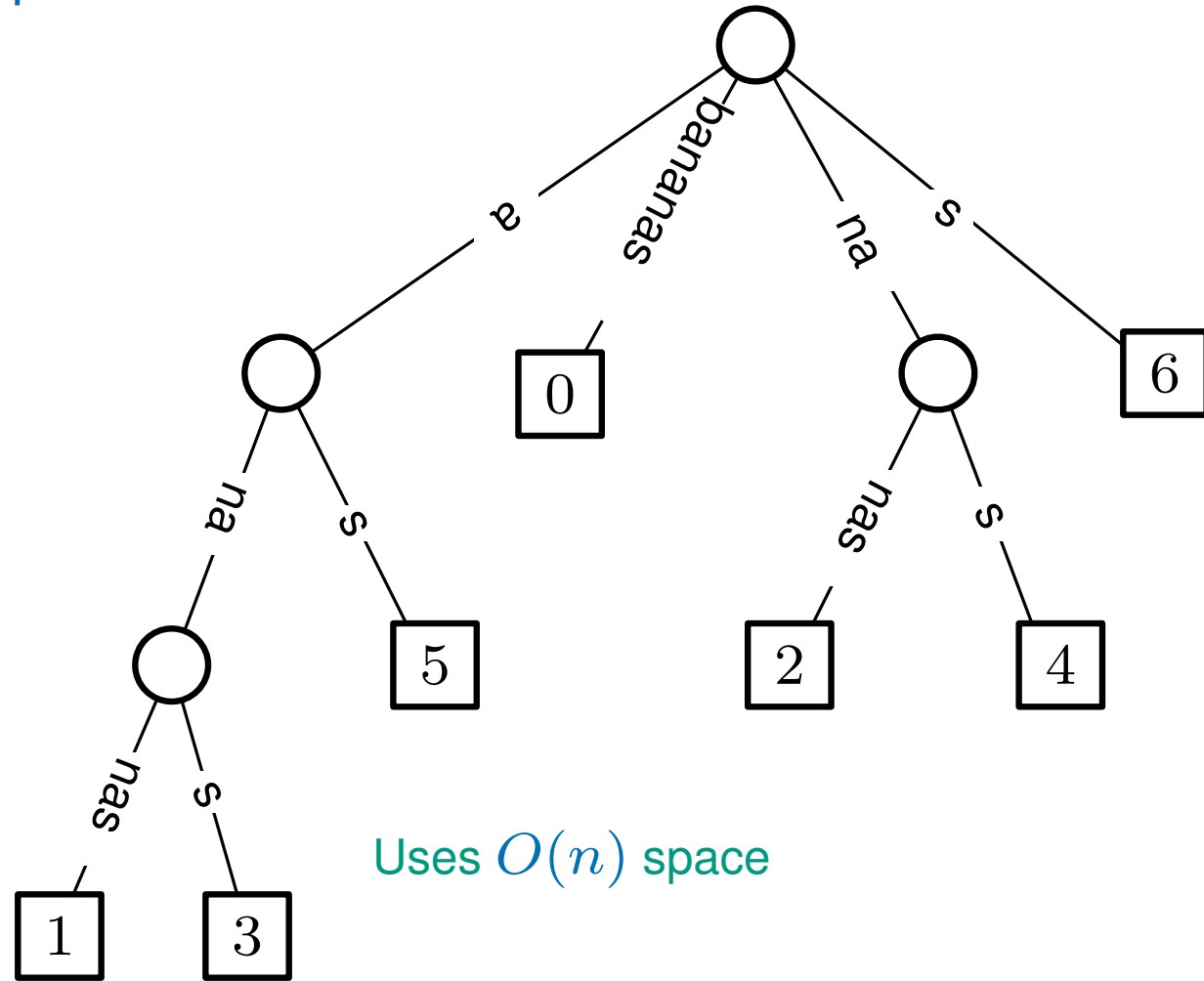
- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix

# Compacted suffix trees



**Compacted Suffix Tree of  $T$**

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring



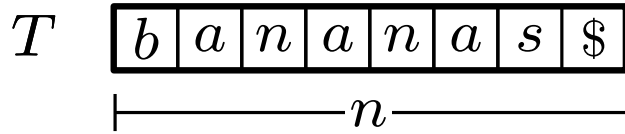
Uses  $O(n)$  space

- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix



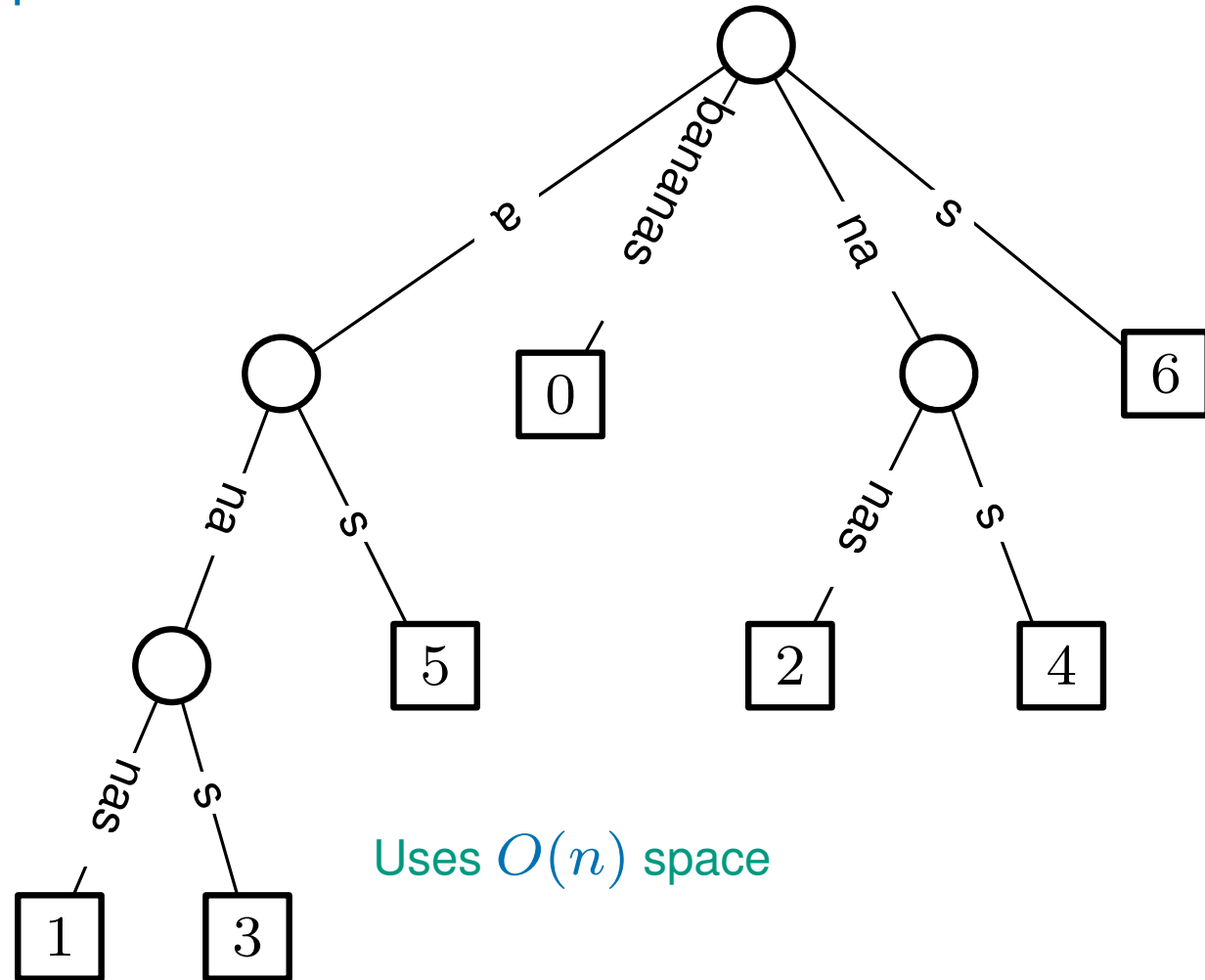
# Compacted suffix trees

Step one: Add a \$ (unique symbol) to  $T$



Compacted Suffix Tree of  $T$

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring

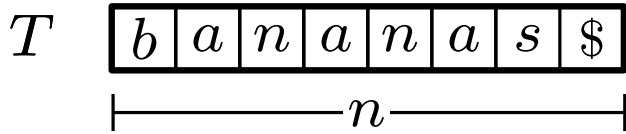


Uses  $O(n)$  space

- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix

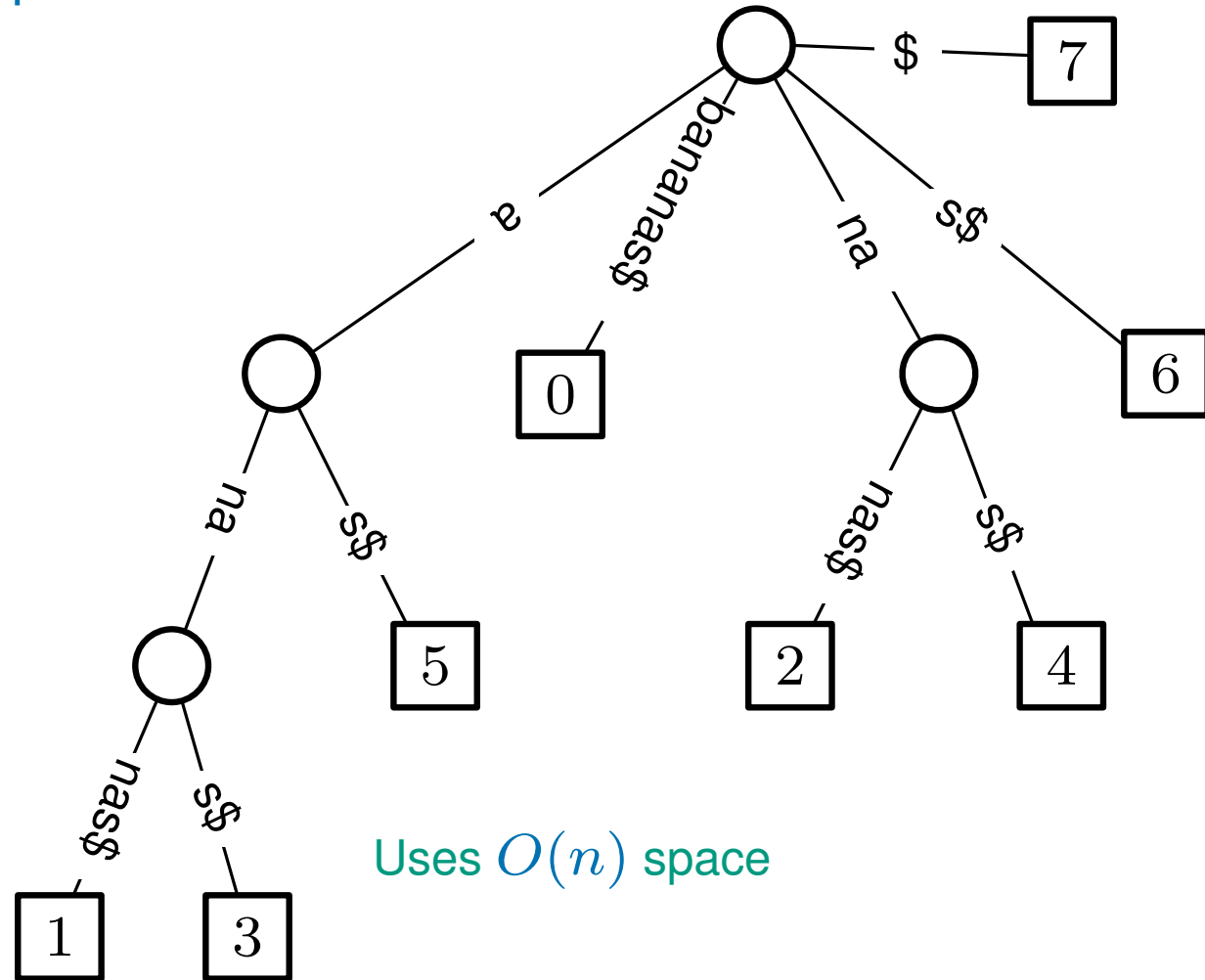
# Compacted suffix trees

**Step one:** Add a \$ (unique symbol) to  $T$



**Compacted Suffix Tree of  $T$**

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring

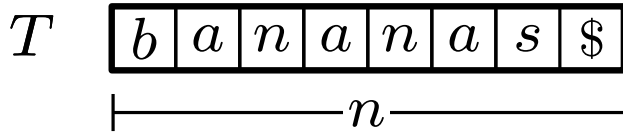


Uses  $O(n)$  space

- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix

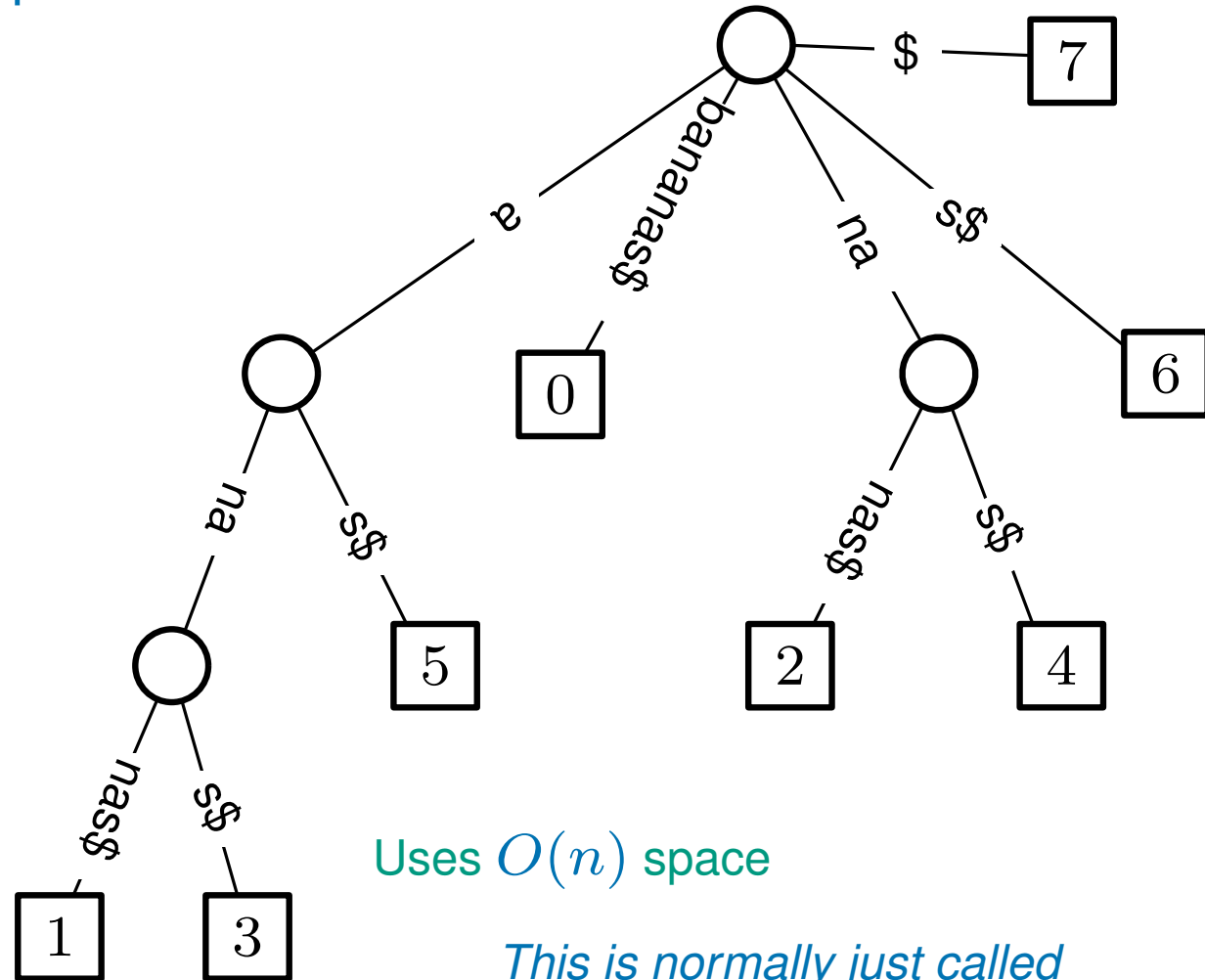
# Compacted suffix trees

**Step one:** Add a \$ (unique symbol) to  $T$



**Compacted Suffix Tree of  $T$**

- A rooted tree with  $n$  leaves
- Every internal node has two or more children
- Every edge is labelled with a substring

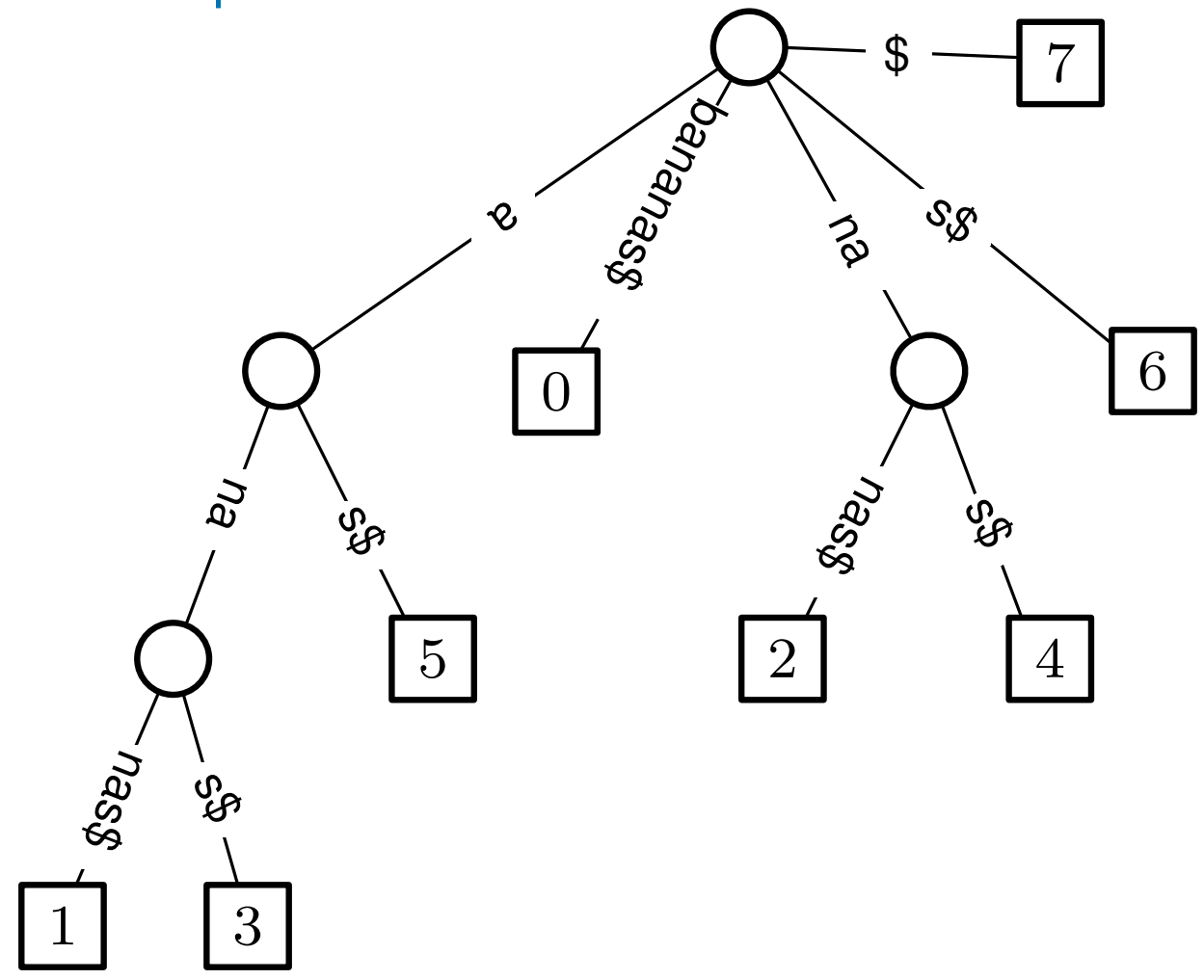
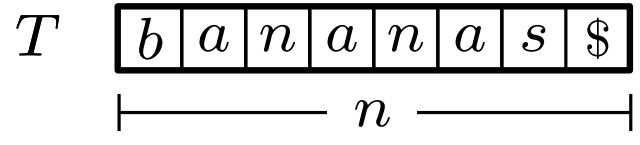


Uses  $O(n)$  space

*This is normally just called a suffix tree*

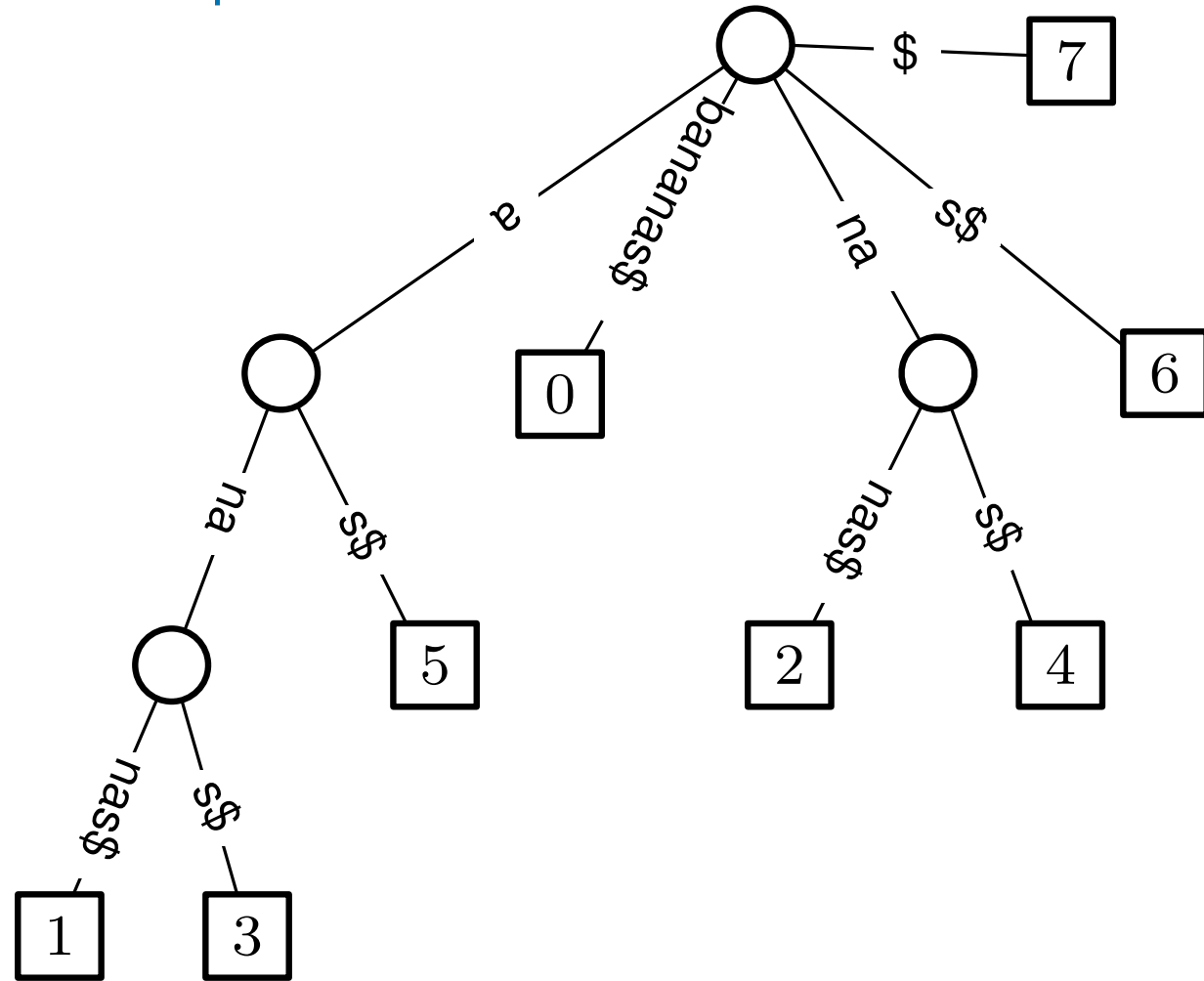
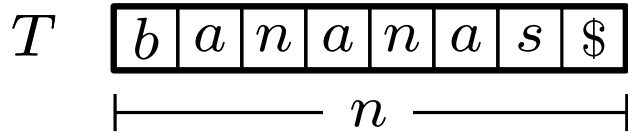
- No two edges leaving the same node have the same first character
- Each leaf is labelled with a location in  $T$
- Any root-to-leaf path spells out the corresponding suffix

# Searching in a compacted suffix tree



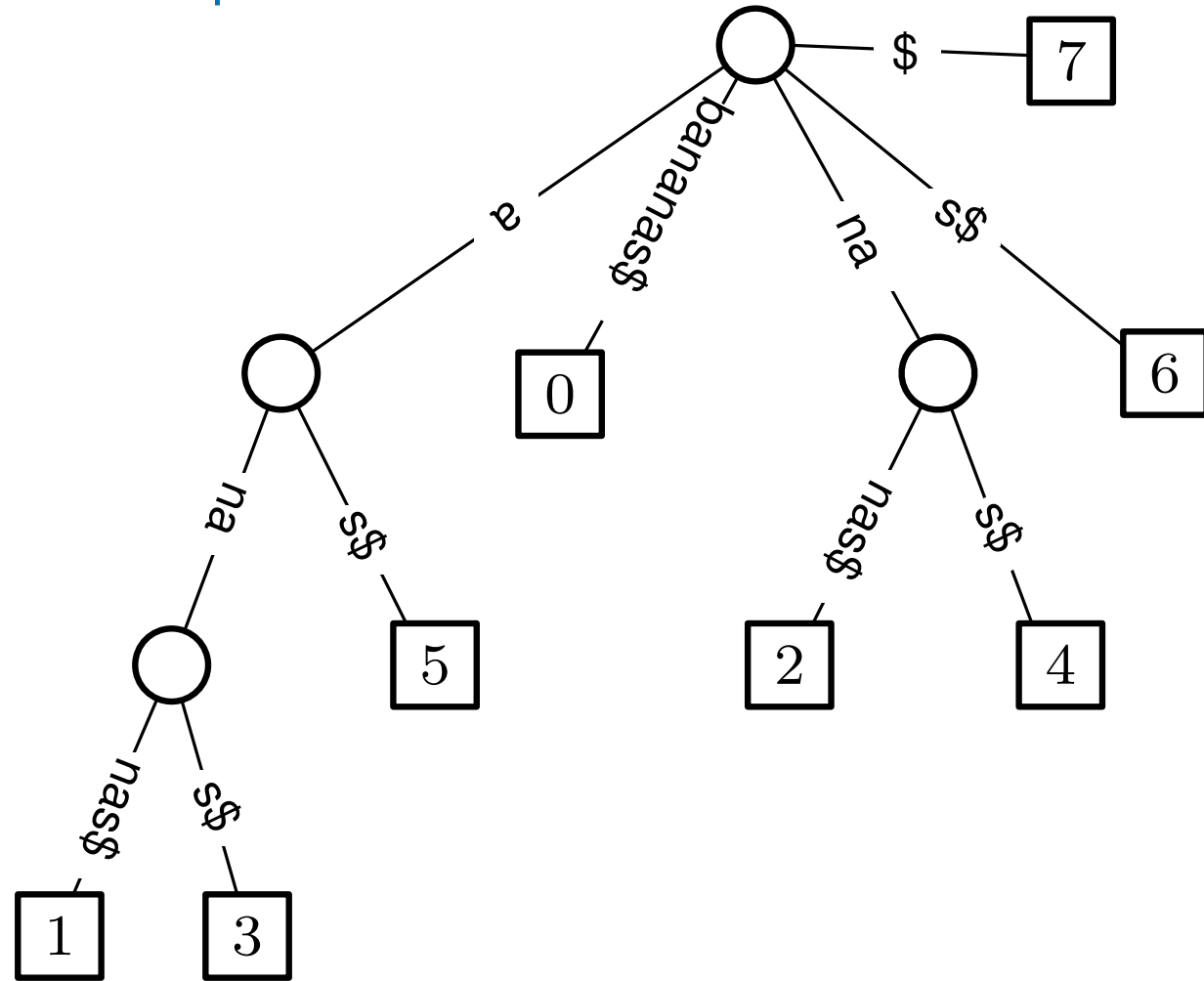
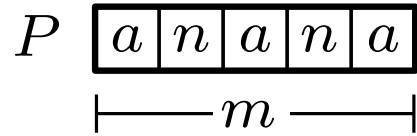
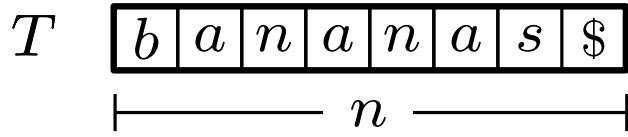


# Searching in a compacted suffix tree



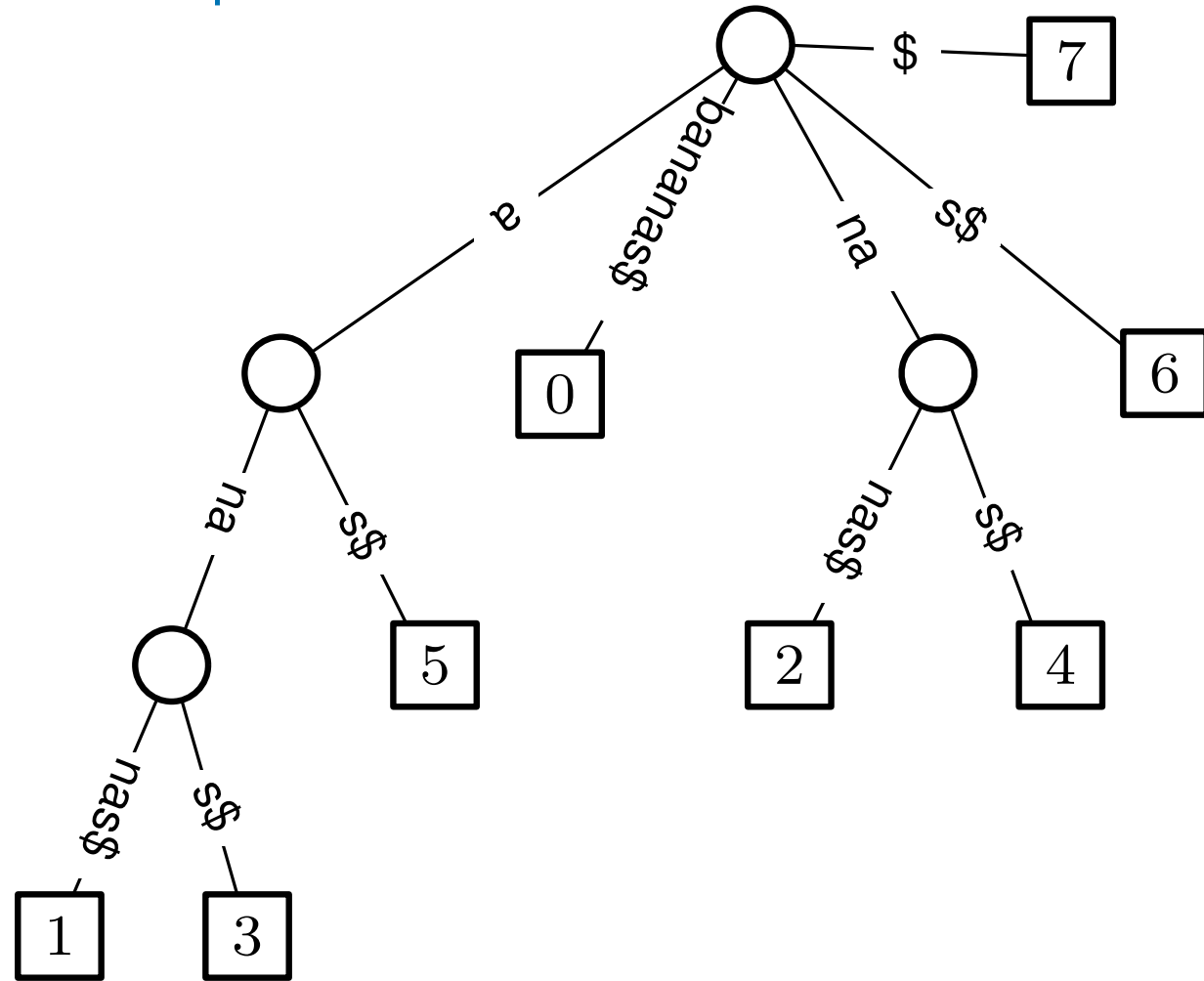
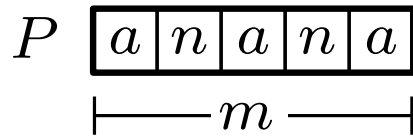
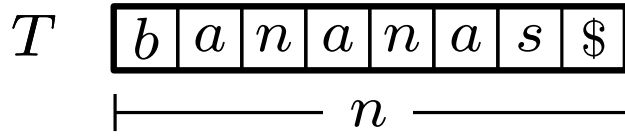
*How do you find a pattern?*

# Searching in a compacted suffix tree



How do you find a pattern?

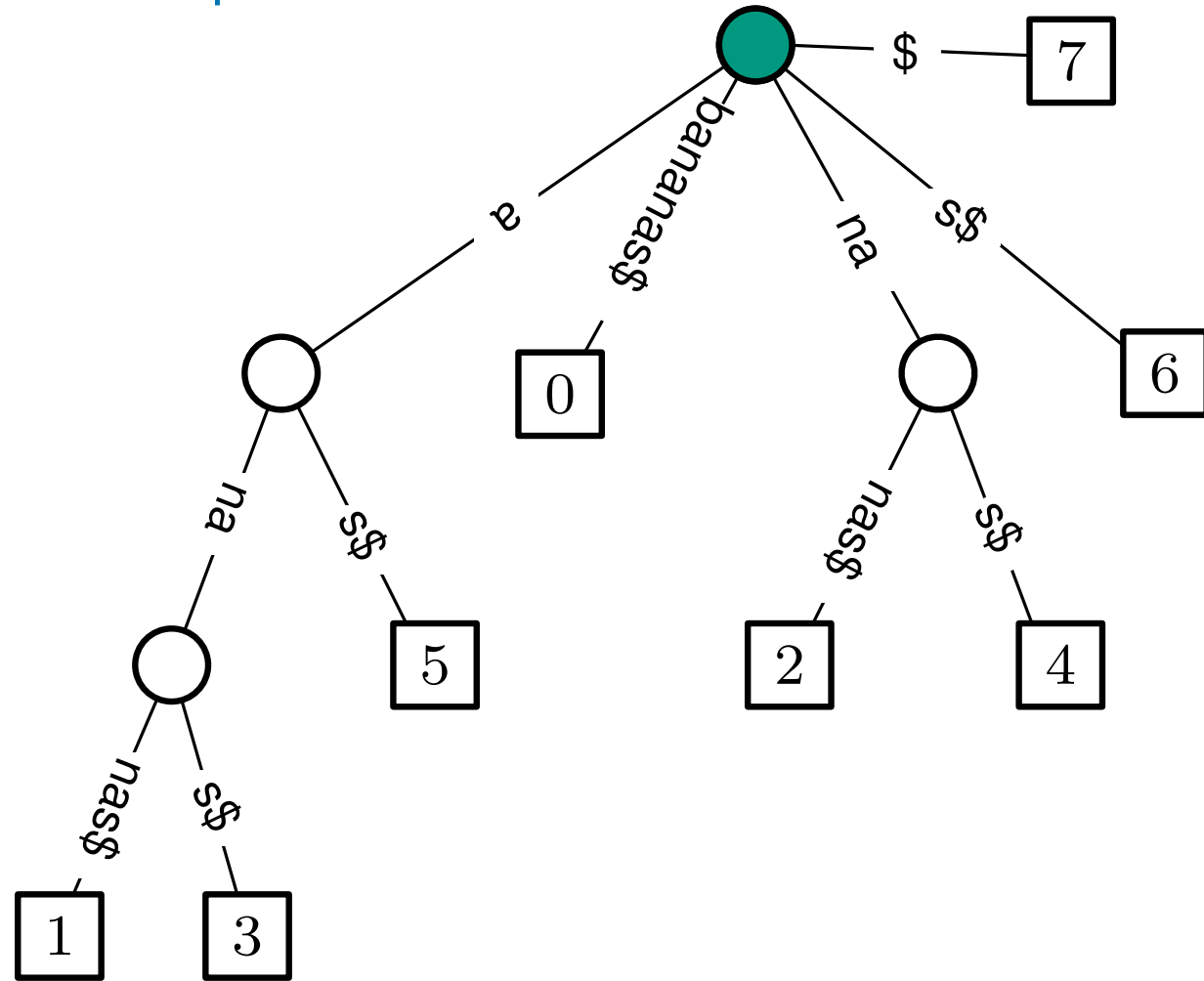
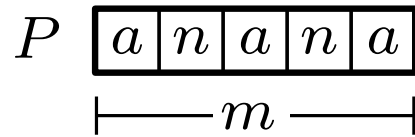
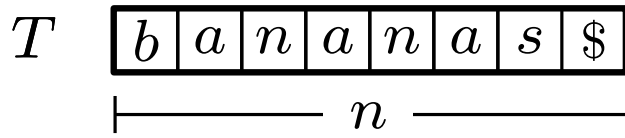
# Searching in a compacted suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

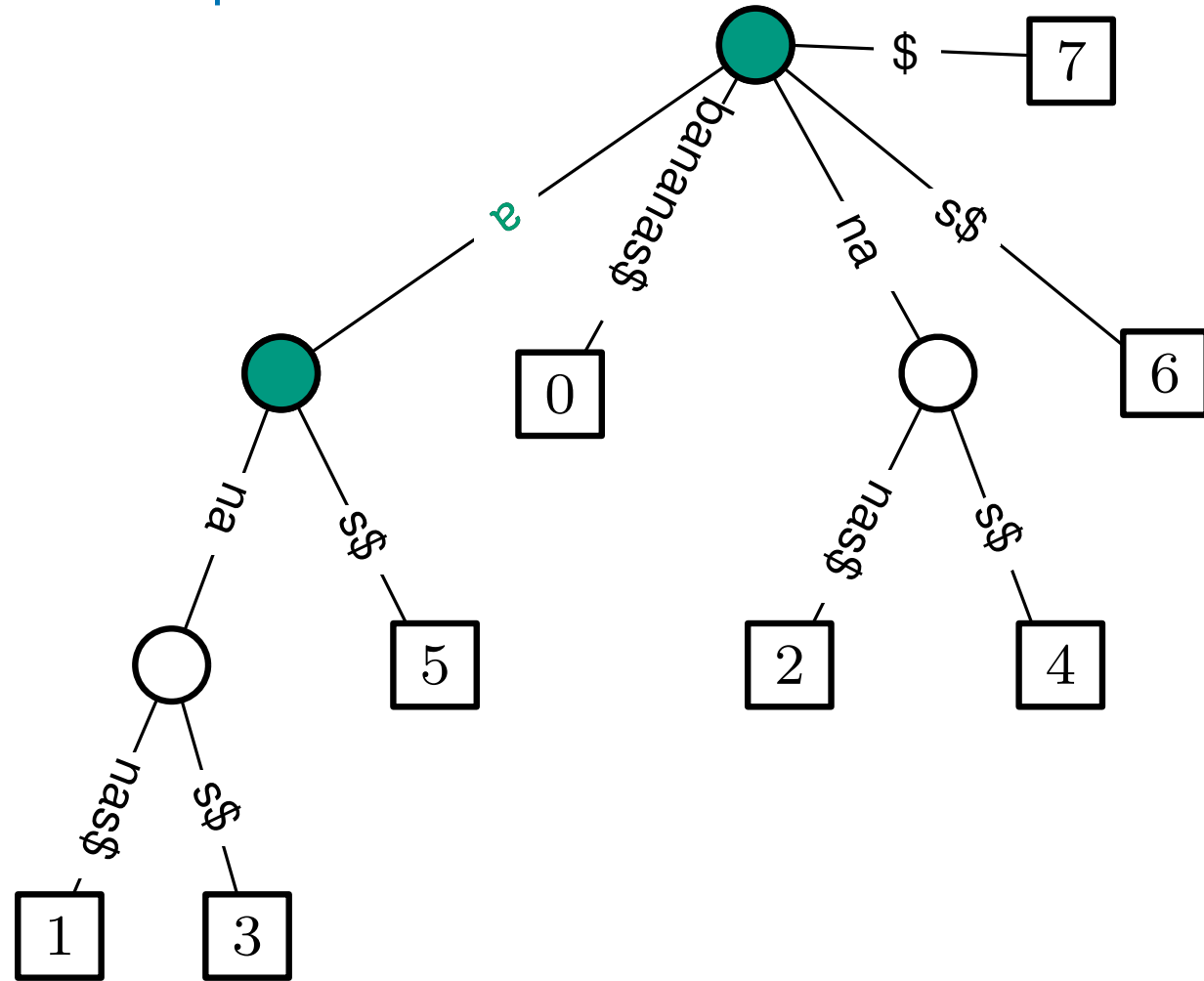
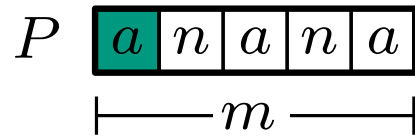
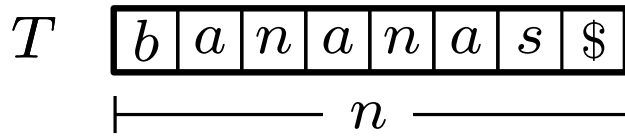
# Searching in a compacted suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

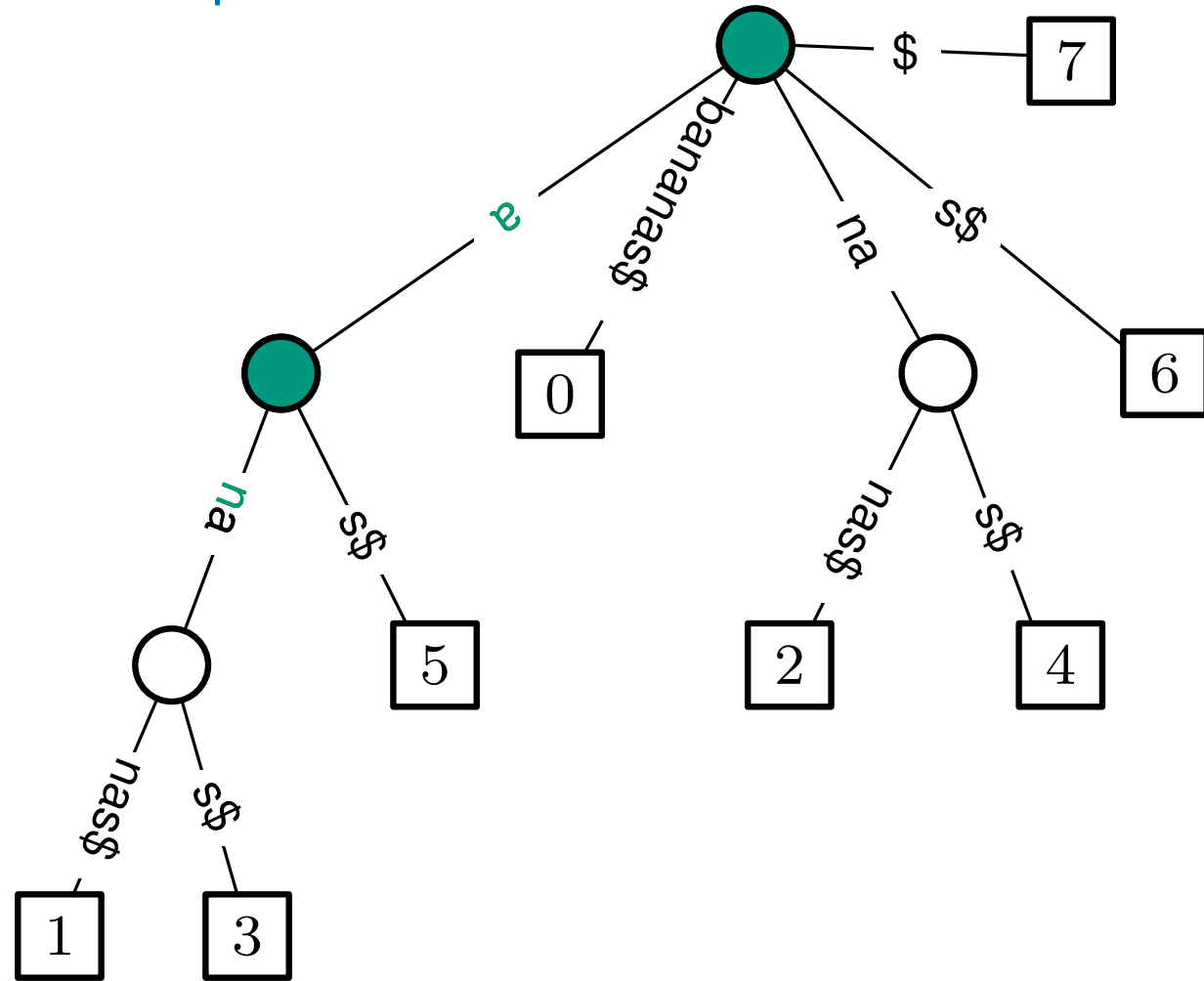
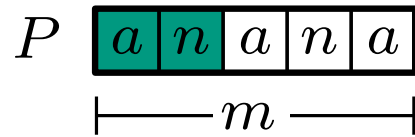
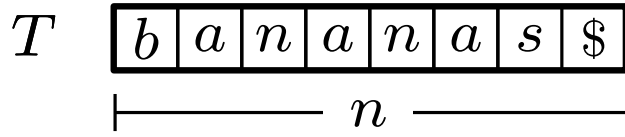
# Searching in a compacted suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

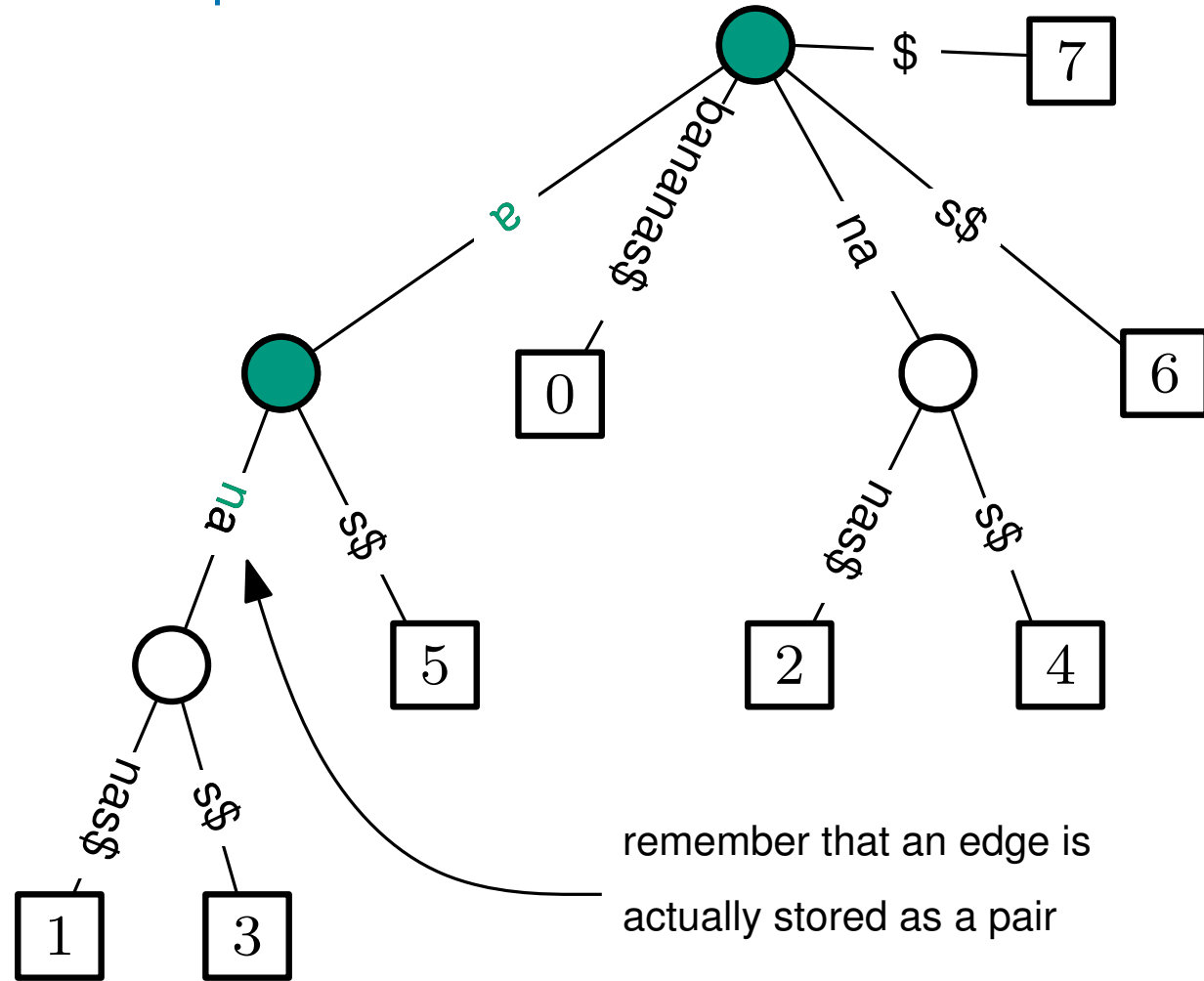
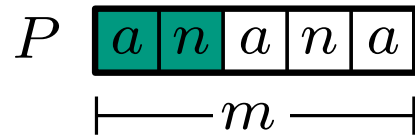
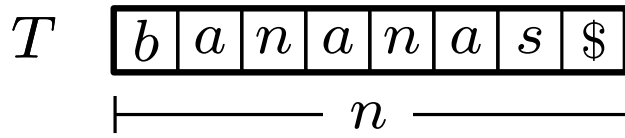
# Searching in a compacted suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

# Searching in a compacted suffix tree

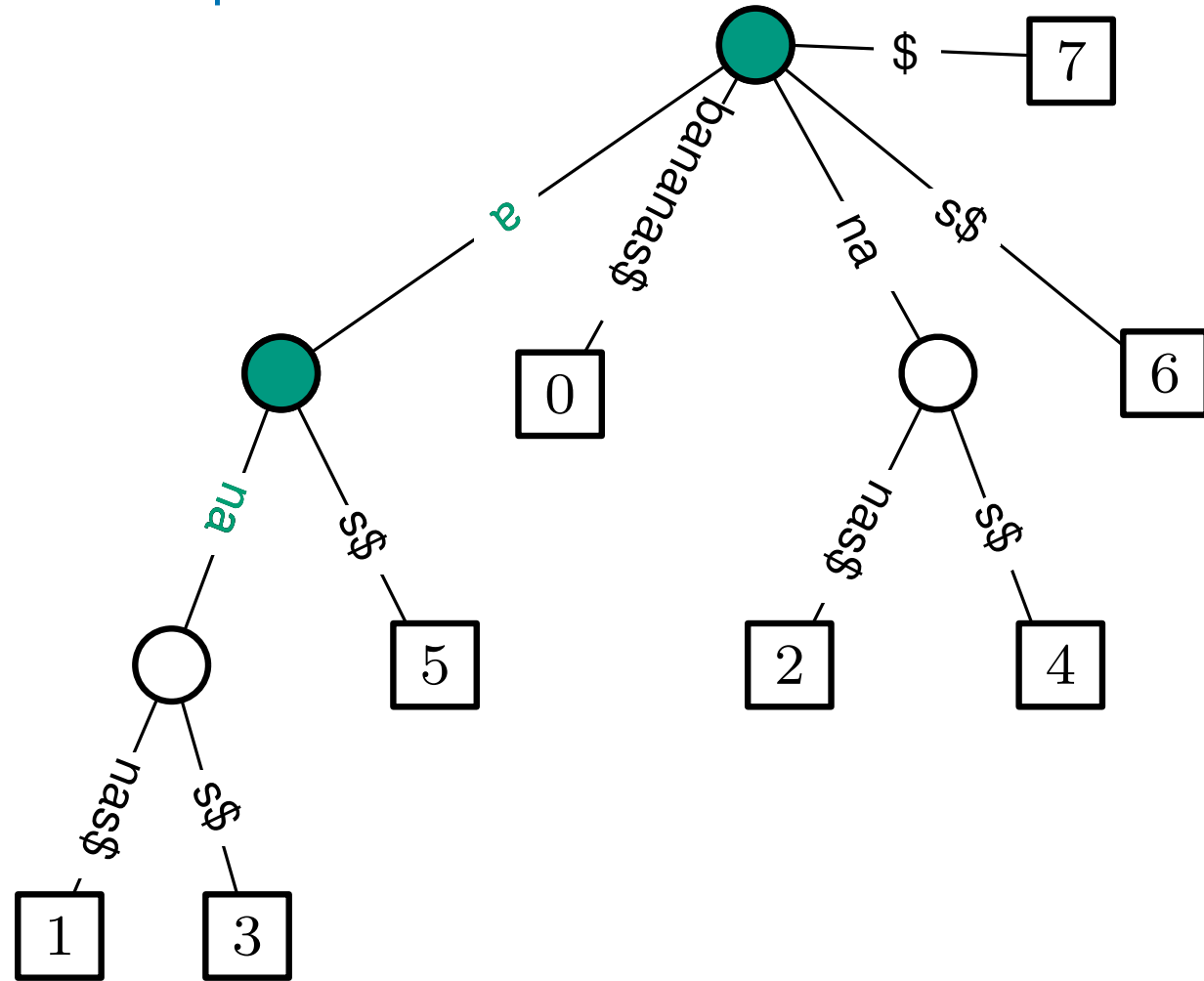
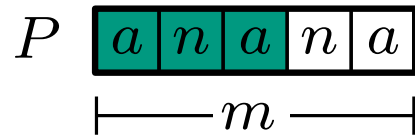
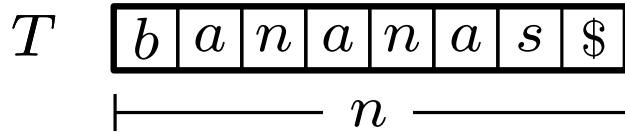


remember that an edge is actually stored as a pair  
*we're actually looking in  $T$*

*How do you find a pattern?*

start at the root and walk down the tree

# Searching in a compacted suffix tree

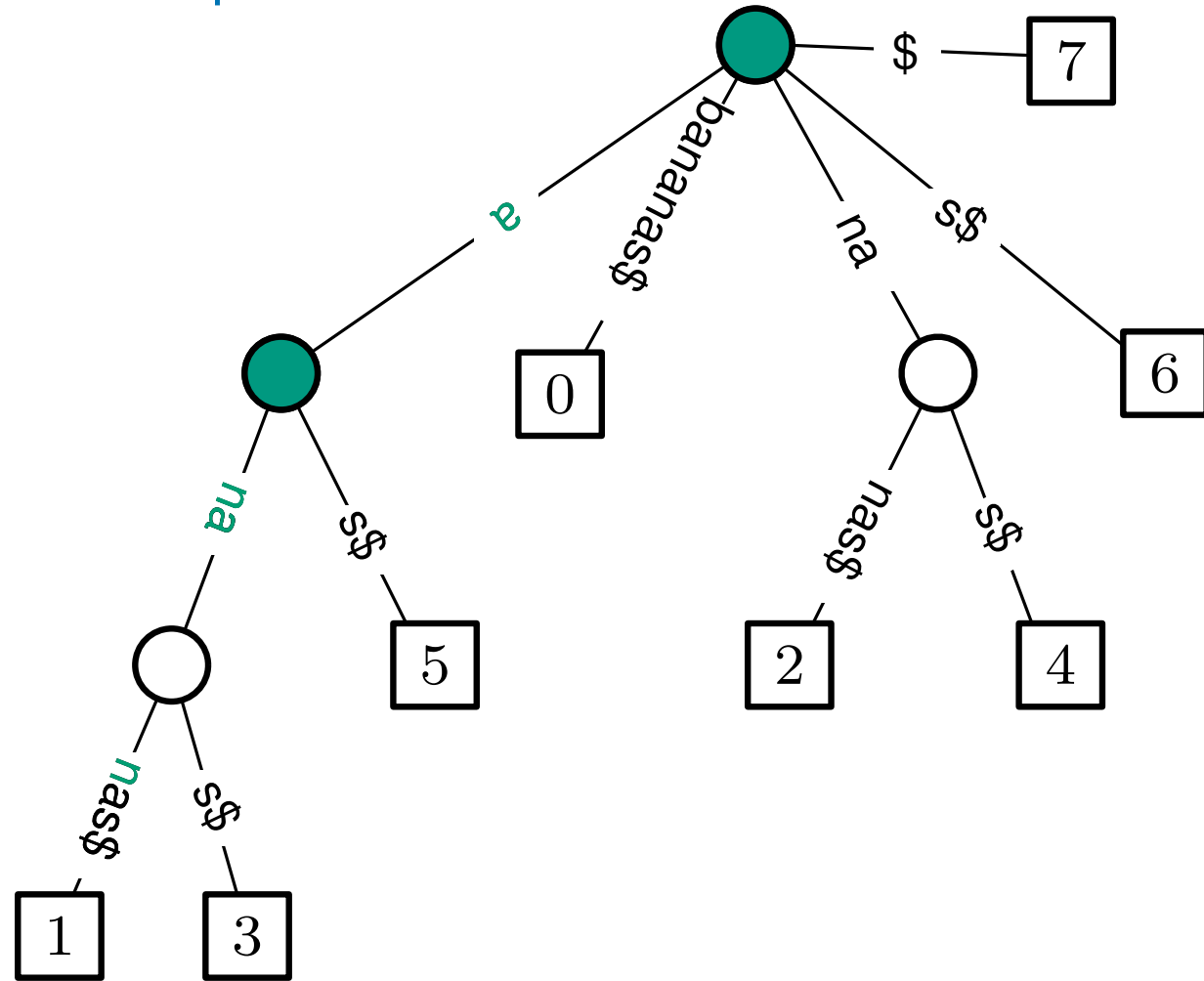
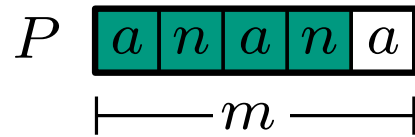
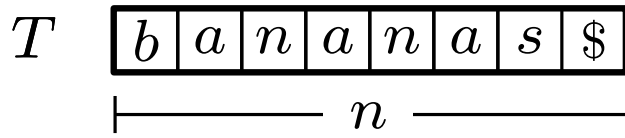


*How do you find a pattern?*

start at the root and walk down the tree



# Searching in a compacted suffix tree

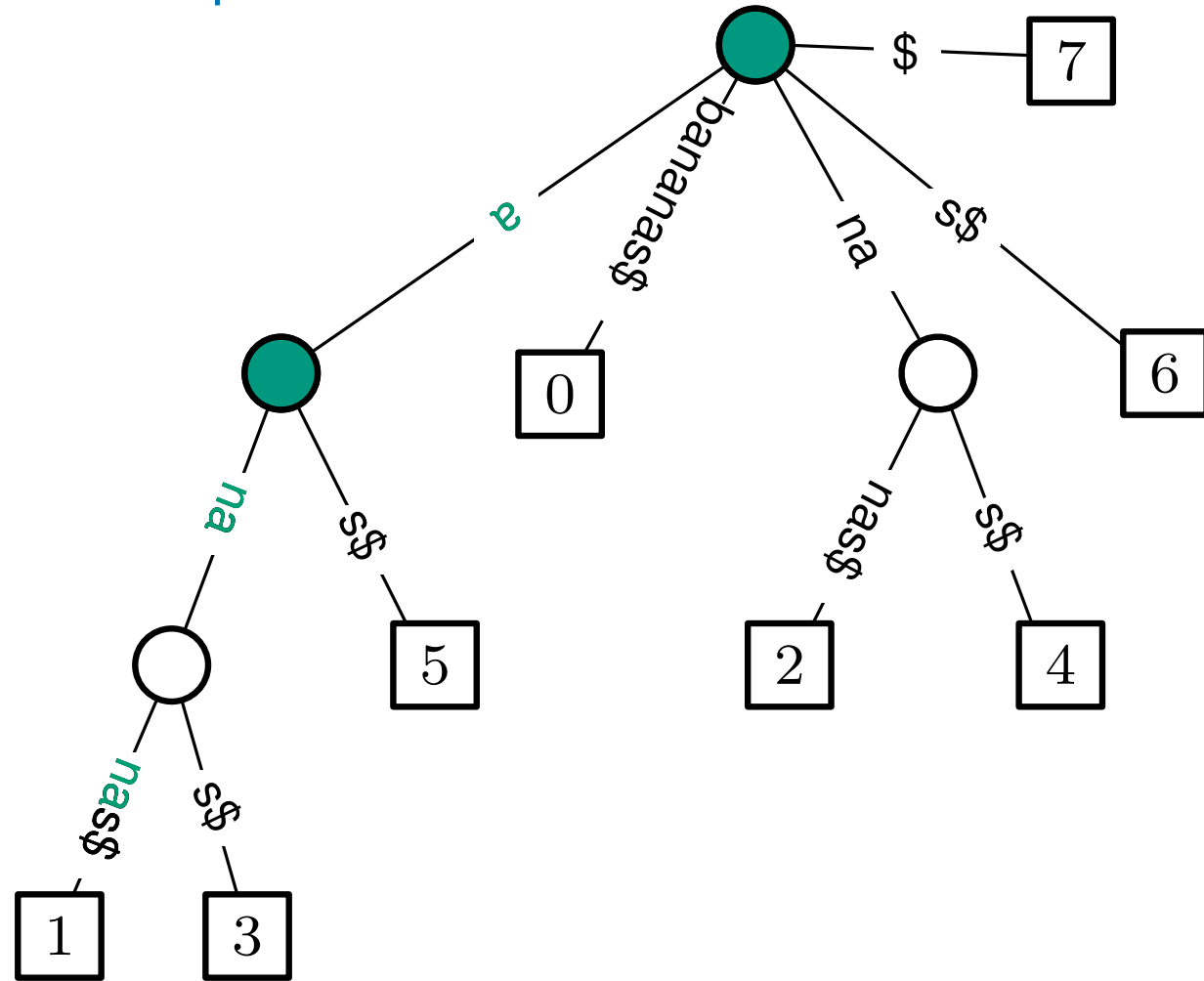
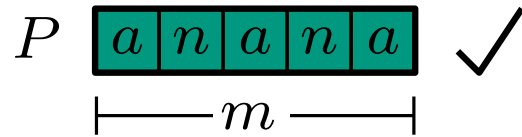
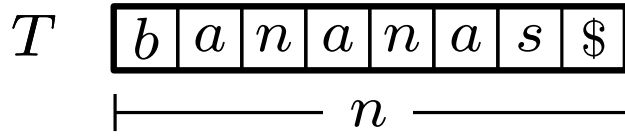


*How do you find a pattern?*

start at the root and walk down the tree



# Searching in a compacted suffix tree

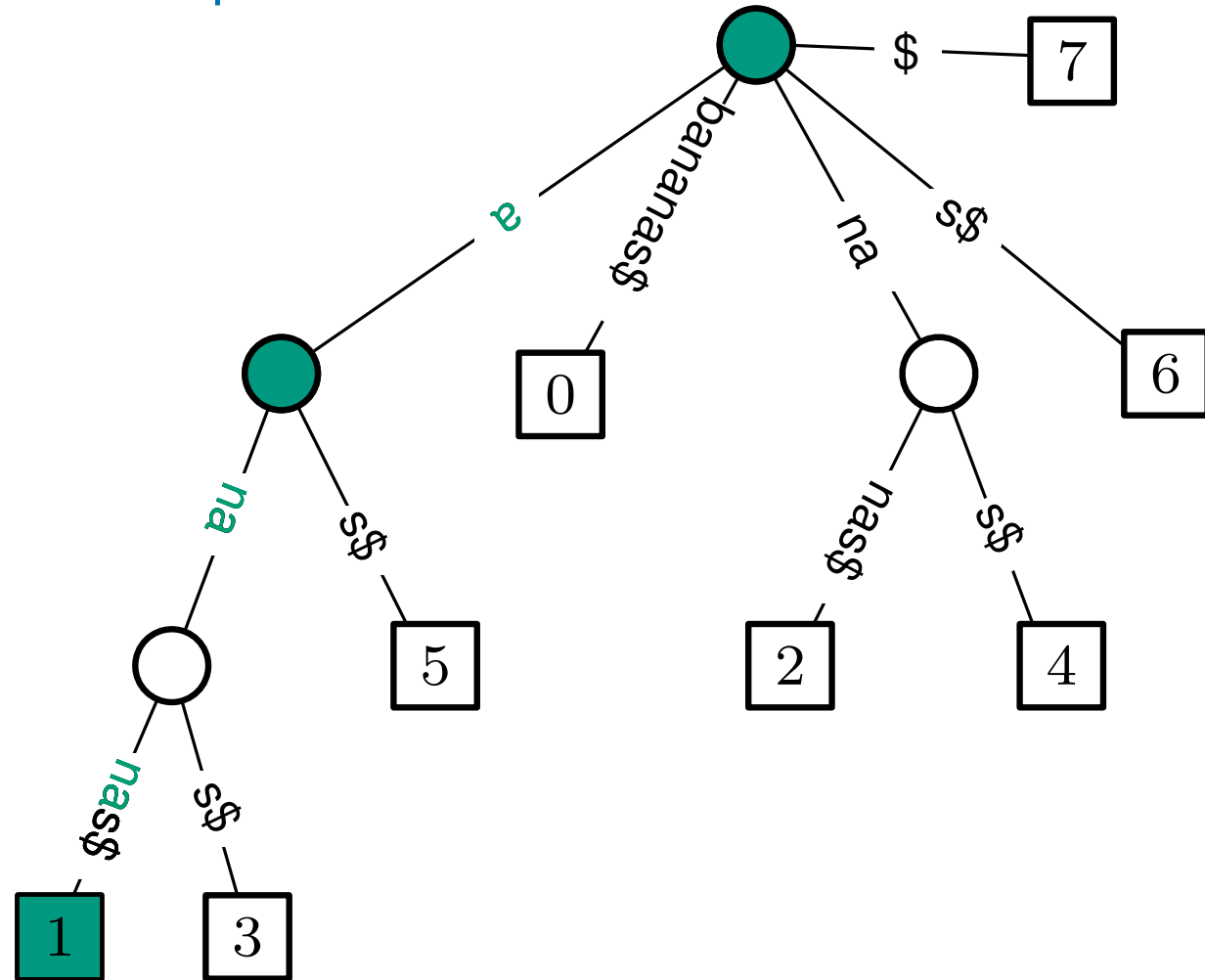
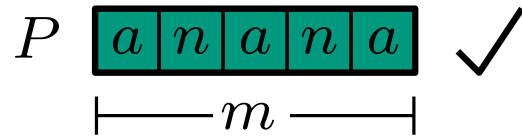
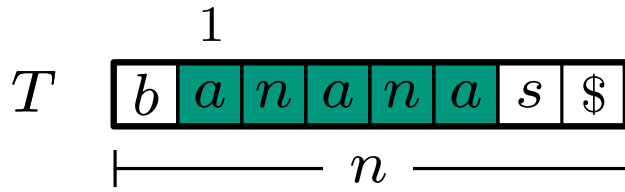


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in a compacted suffix tree

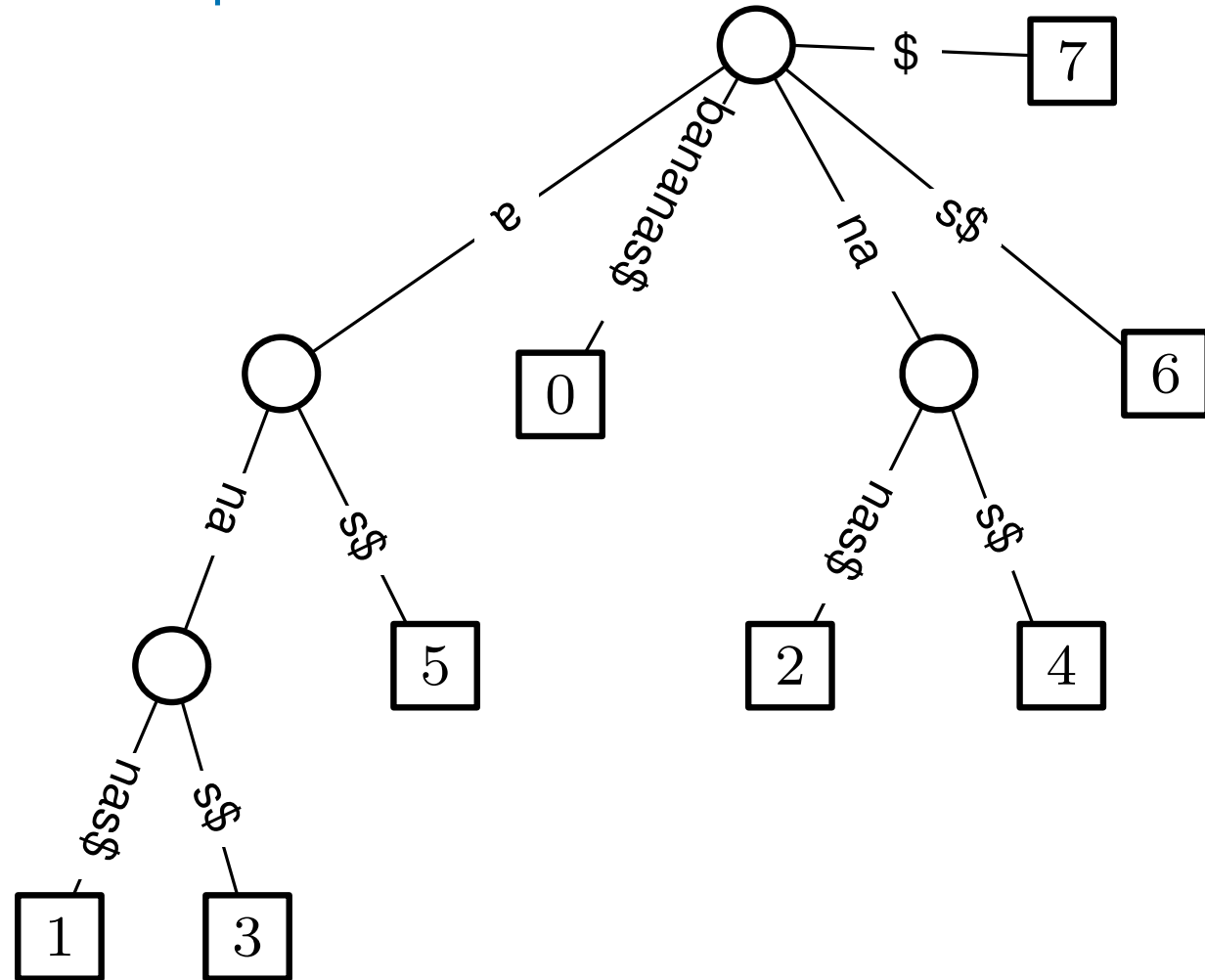
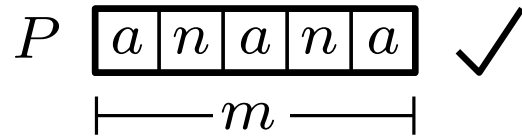
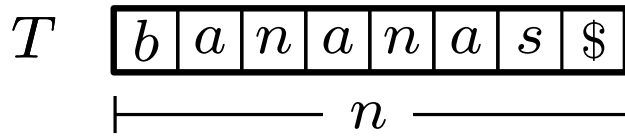


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in a compacted suffix tree

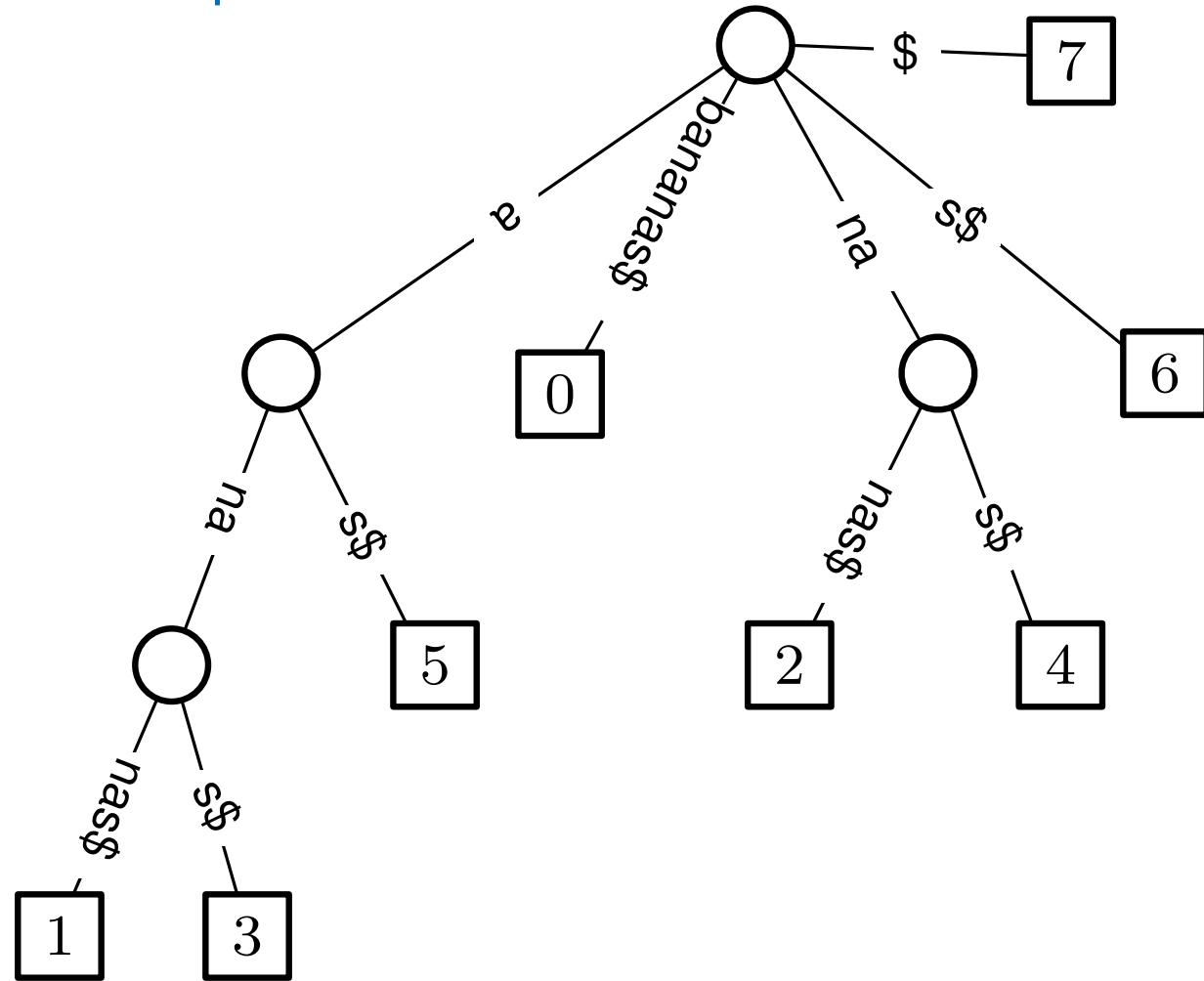
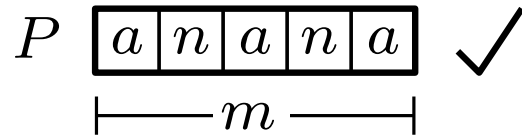
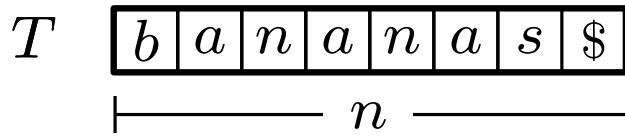


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in a compacted suffix tree



*How do you find a pattern?*

start at the root and walk down the tree

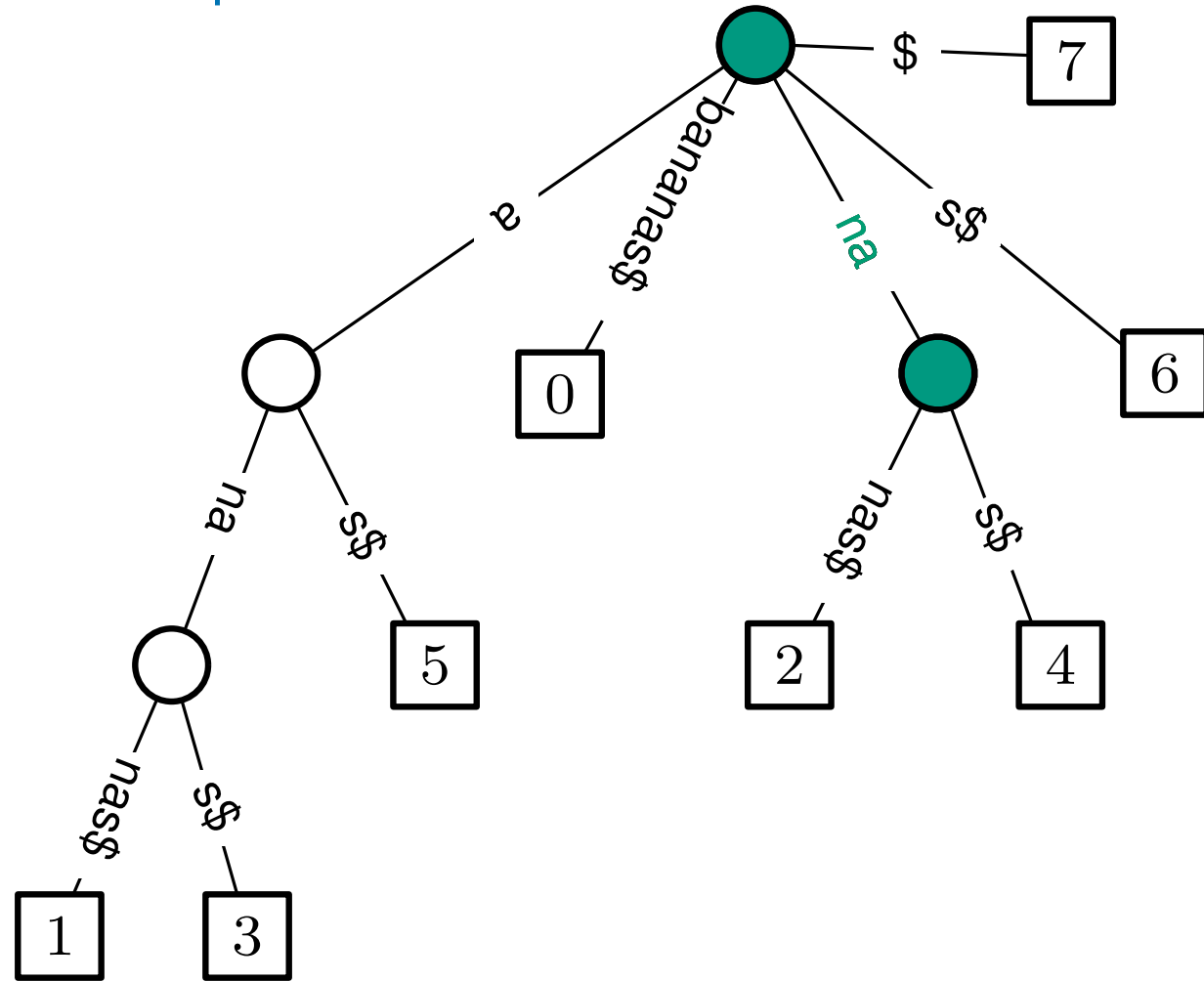
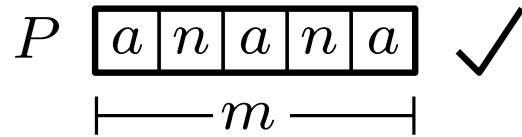
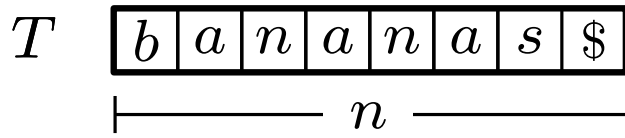
... matches occur at the leaves of the subtree







# Searching in a compacted suffix tree

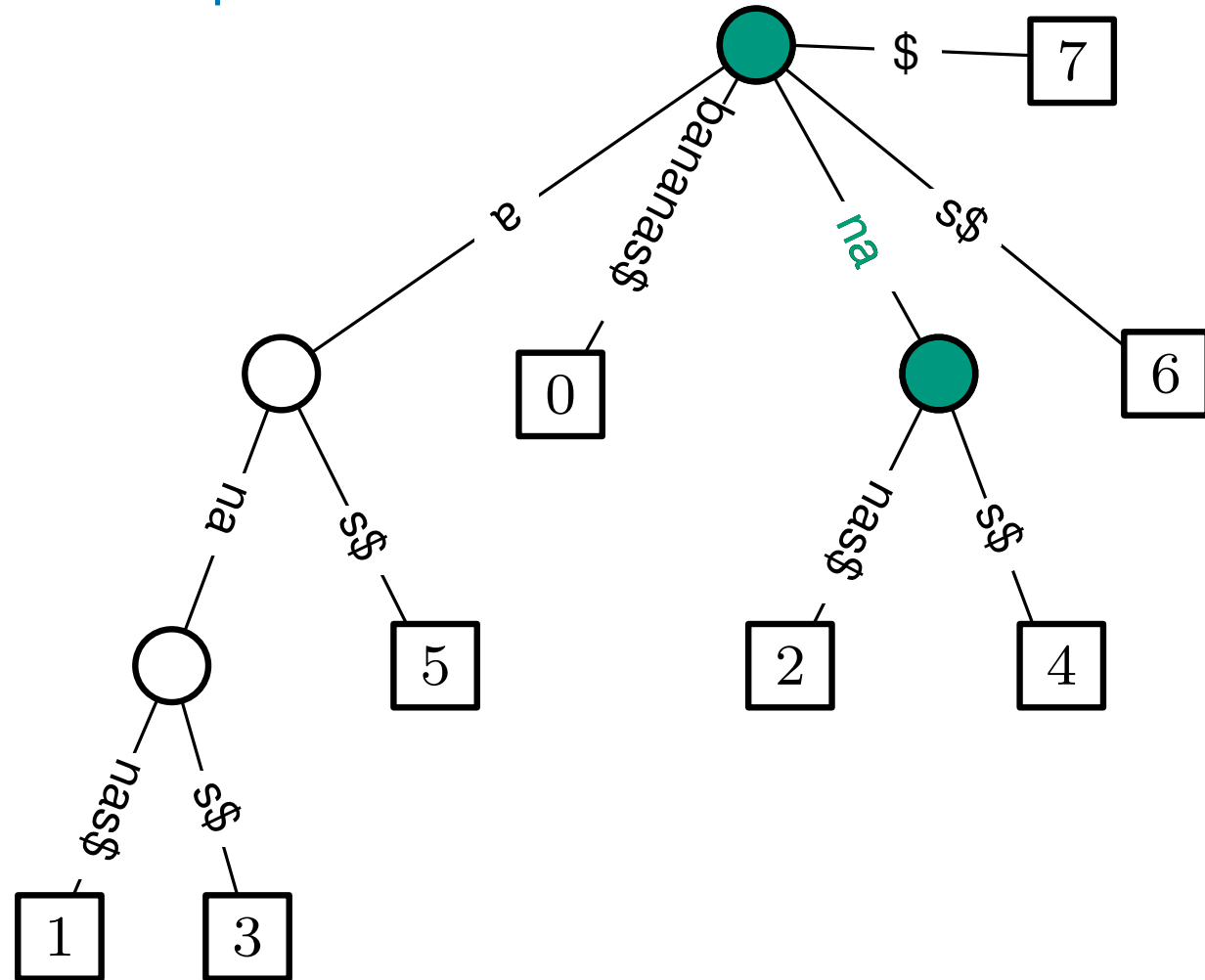
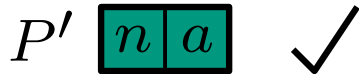
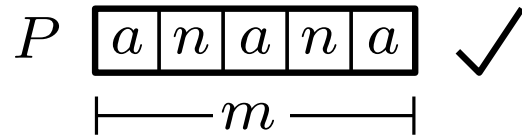
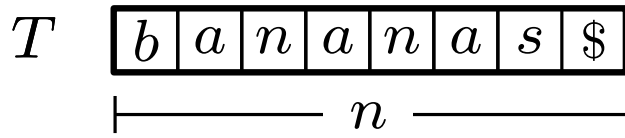


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in a compacted suffix tree

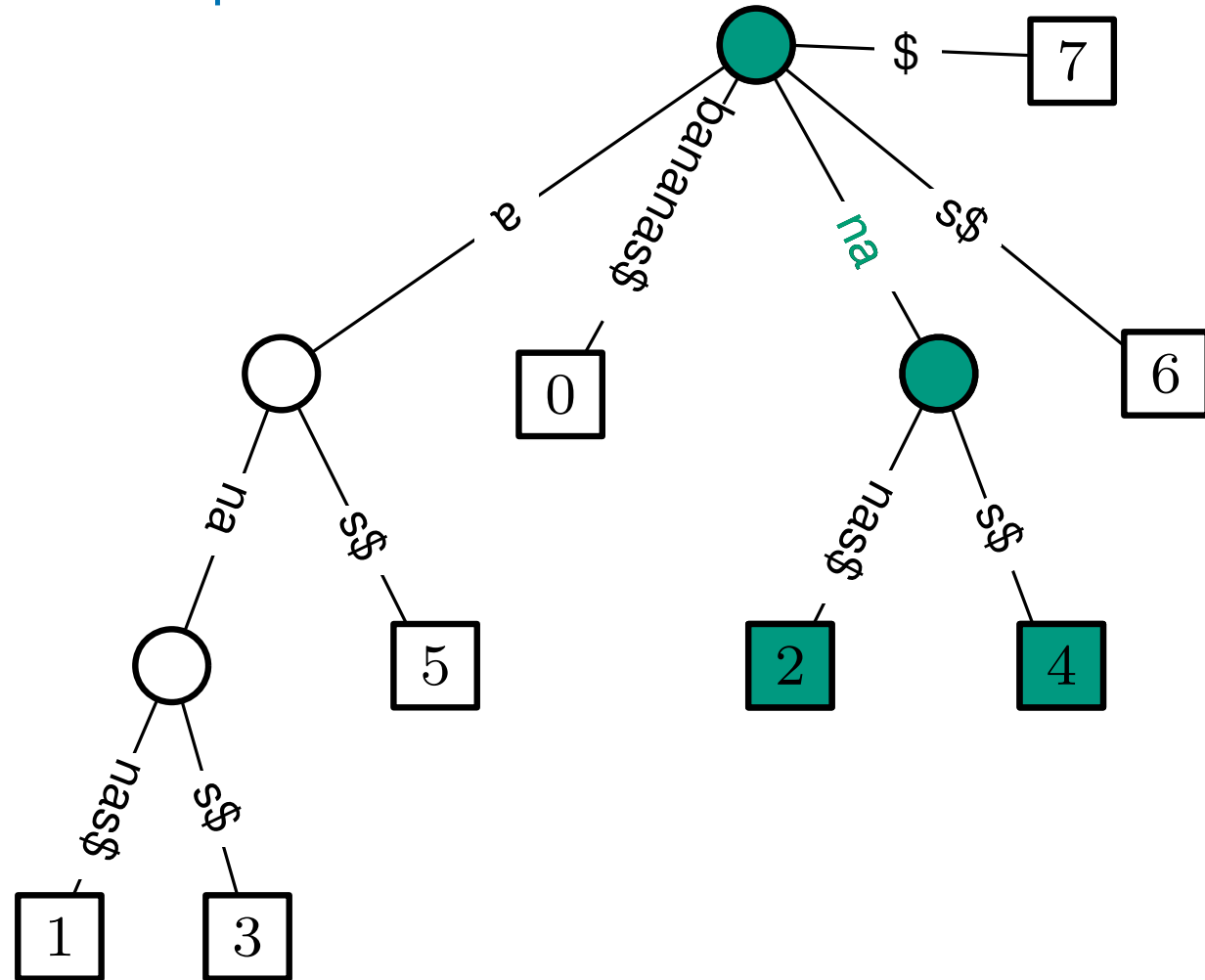
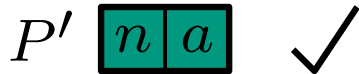
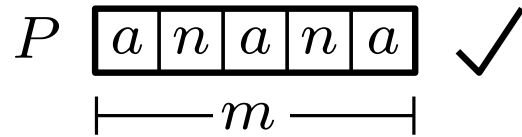
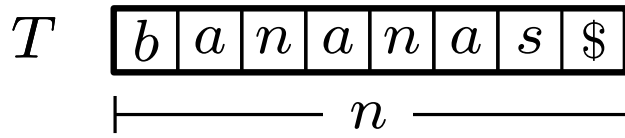


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in a compacted suffix tree

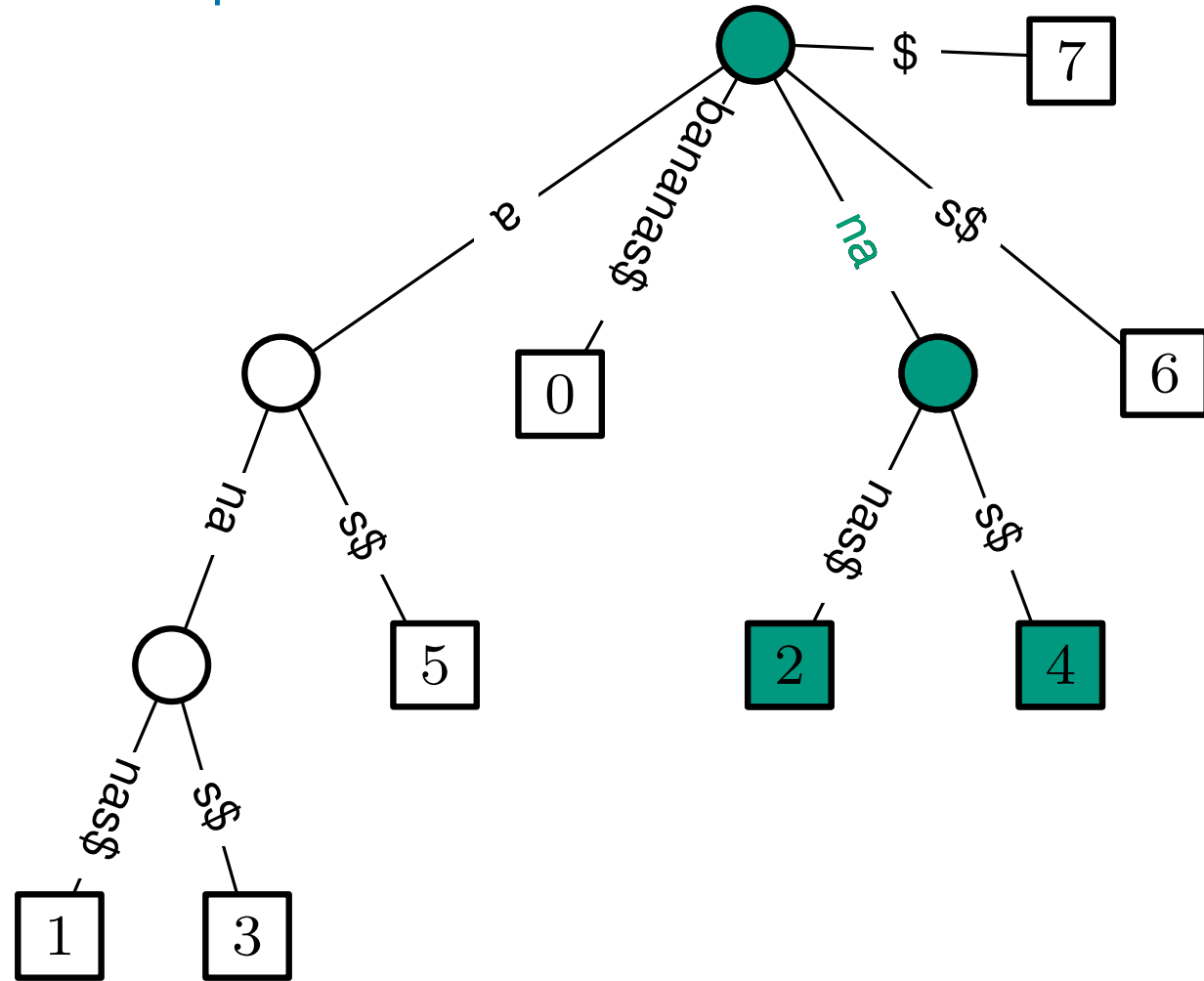
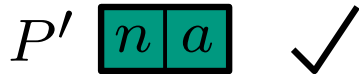
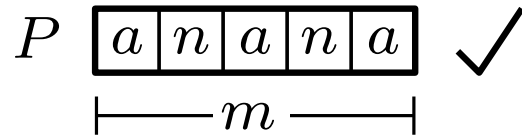
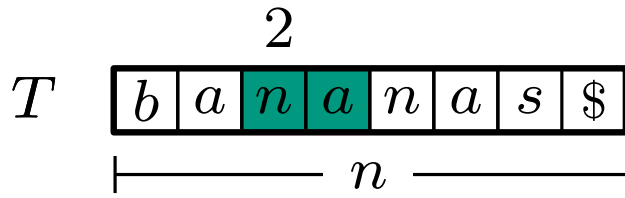


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in a compacted suffix tree



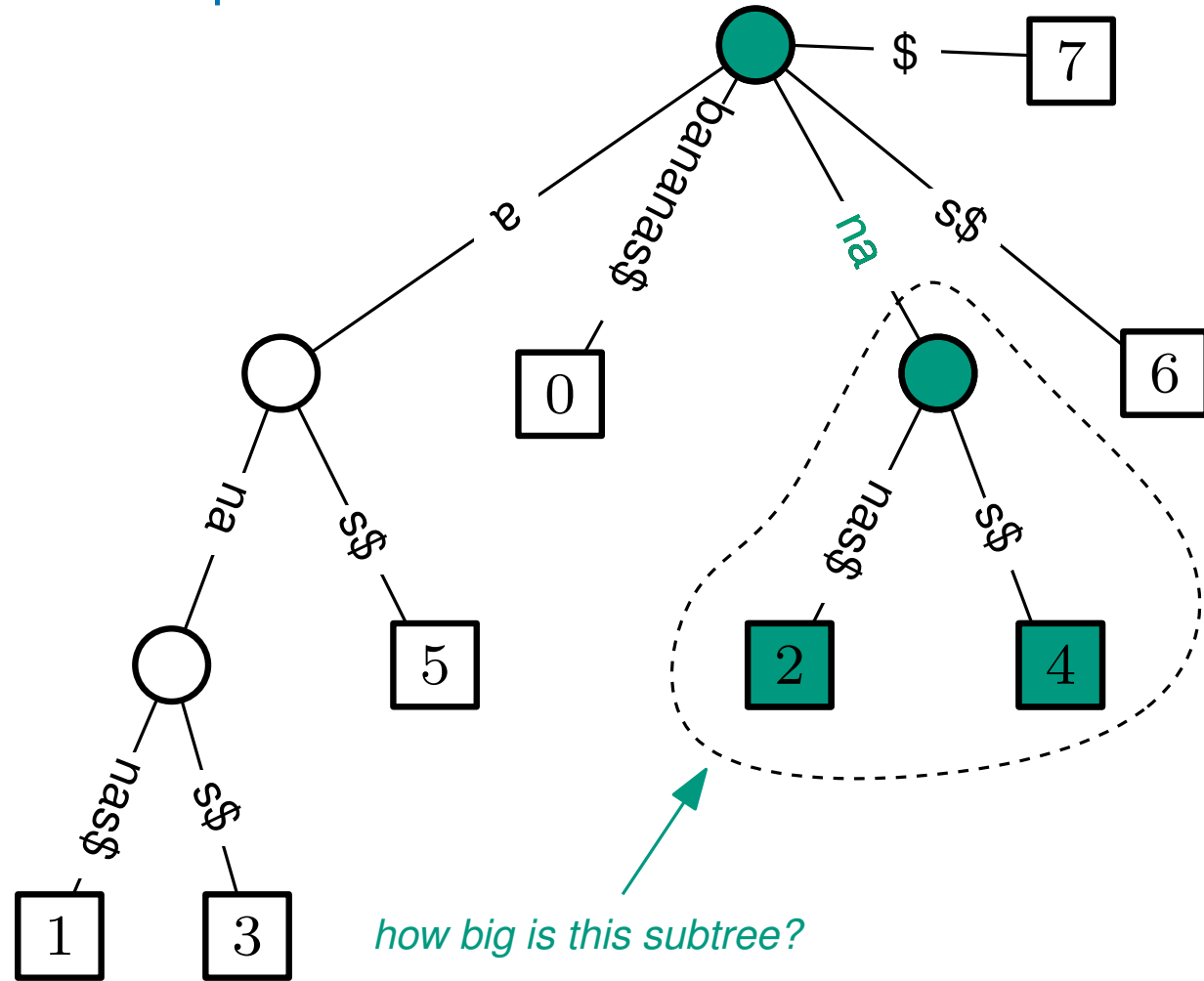
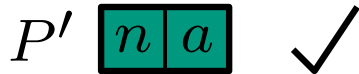
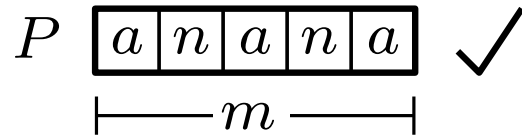
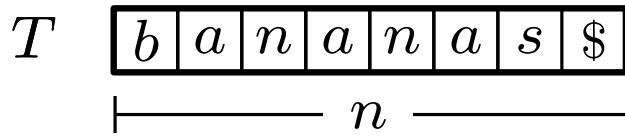
*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree



# Searching in a compacted suffix tree

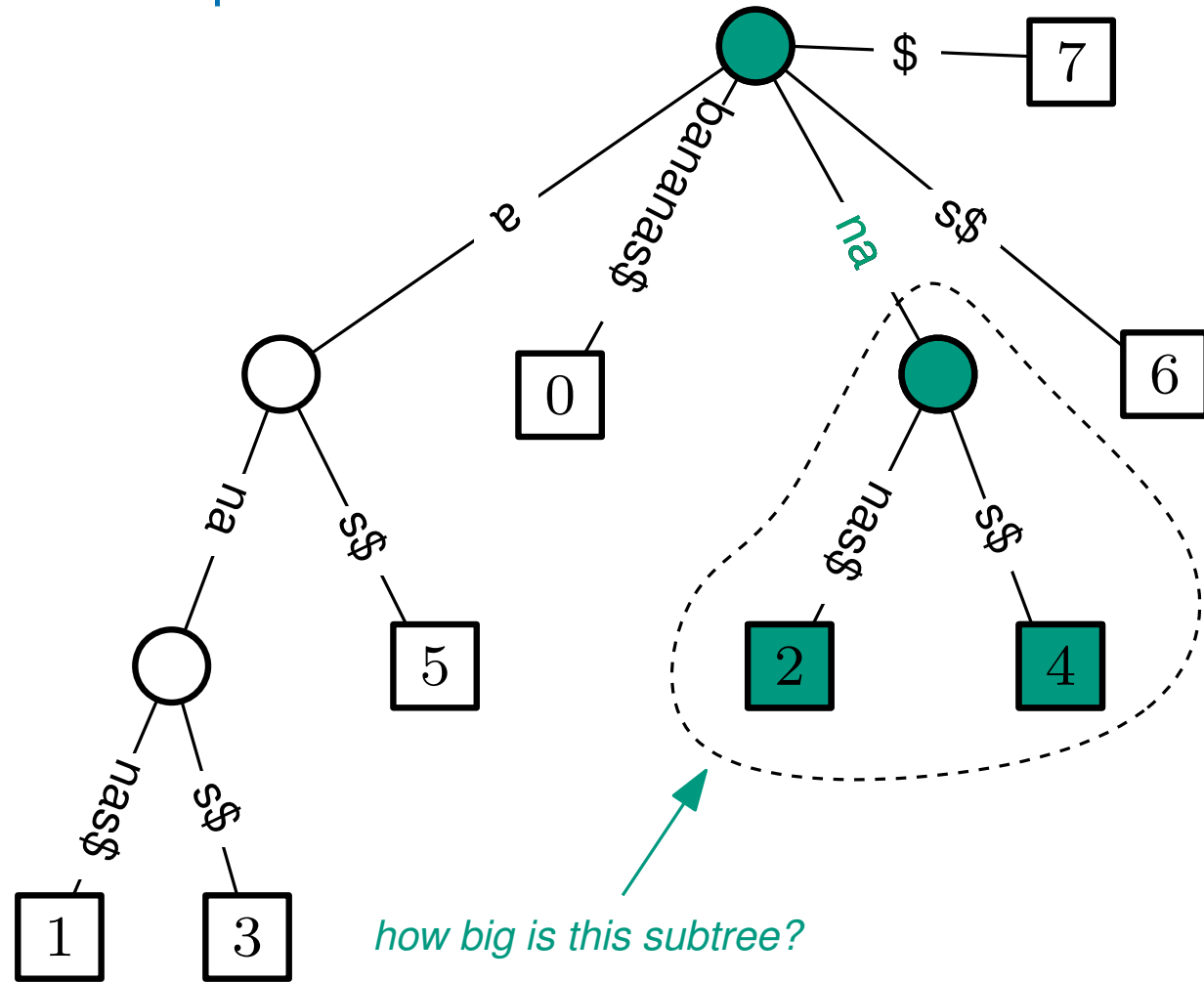
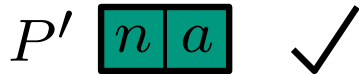
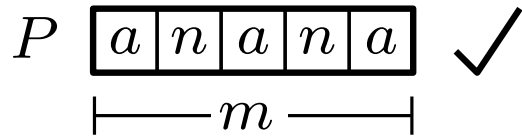
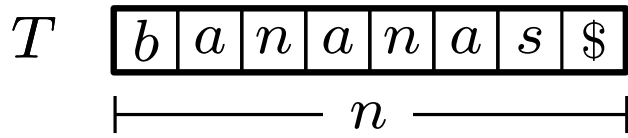


*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in a compacted suffix tree



*how big is this subtree?*

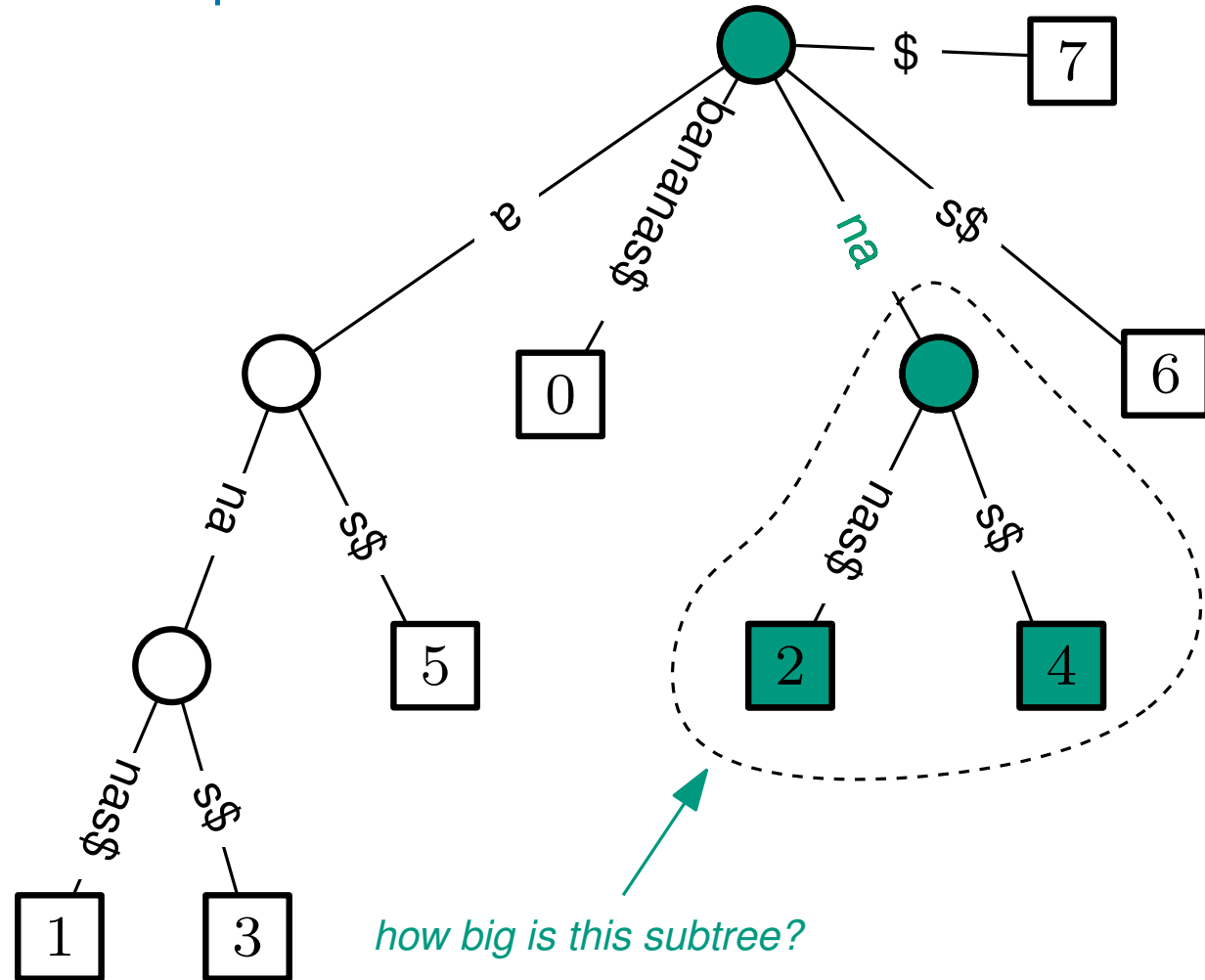
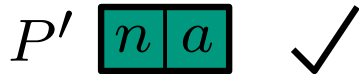
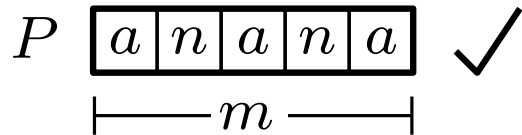
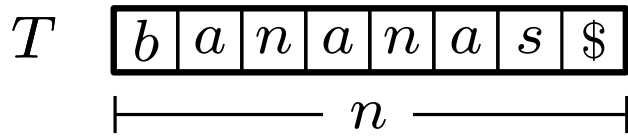
$O(\text{occ})$  because it has  $\text{occ}$  leaves  
 (and each internal node has at least two children)

*How do you find a pattern?*

start at the root and walk down the tree

... matches occur at the leaves of the subtree

# Searching in a compacted suffix tree



*how big is this subtree?*  
 $O(occ)$  because it has  $occ$  leaves  
 (and each internal node has at least two children)

*How do you find a pattern?*

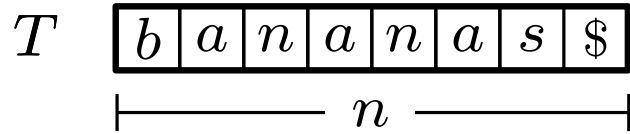
start at the root and walk down the tree  
 ... matches occur at the leaves of the subtree

We can find all the matches in  $O(m + occ)$  time (by looking at the whole subtree)



you should  
*never actually*  
 do it like this

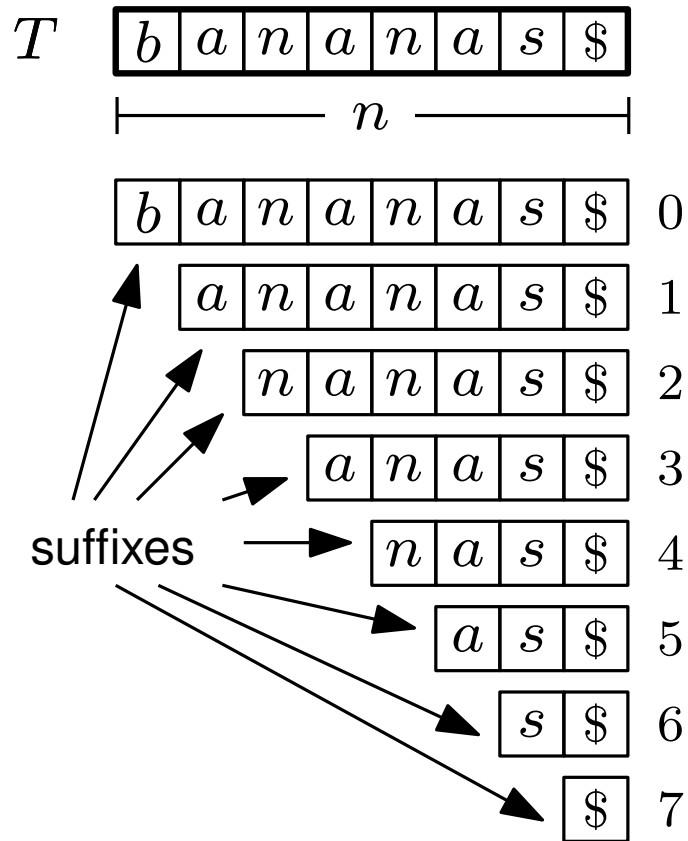
**Naively** constructing a compacted suffix tree

Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this **Naively** constructing a compacted suffix tree

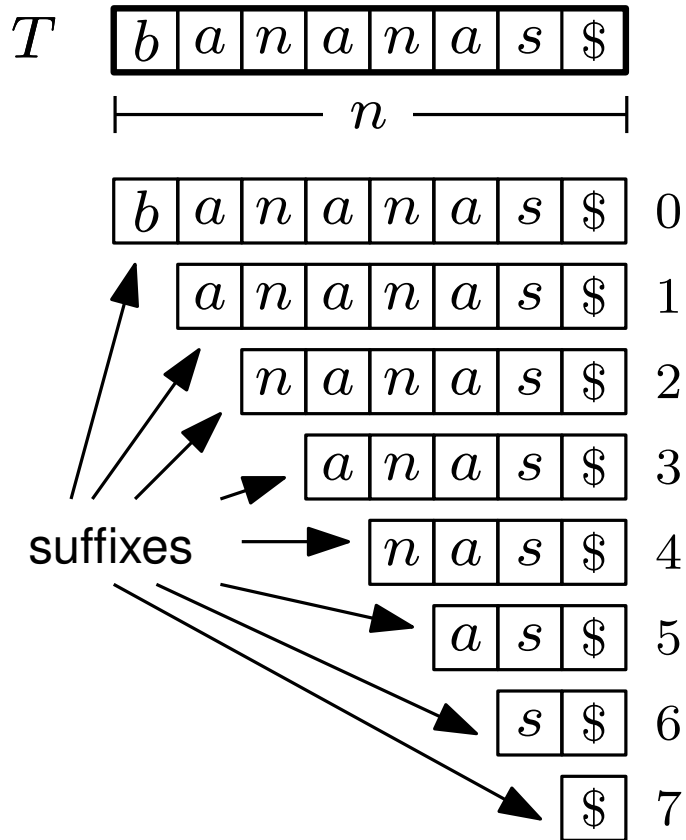


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should  
never actually  
do it like this

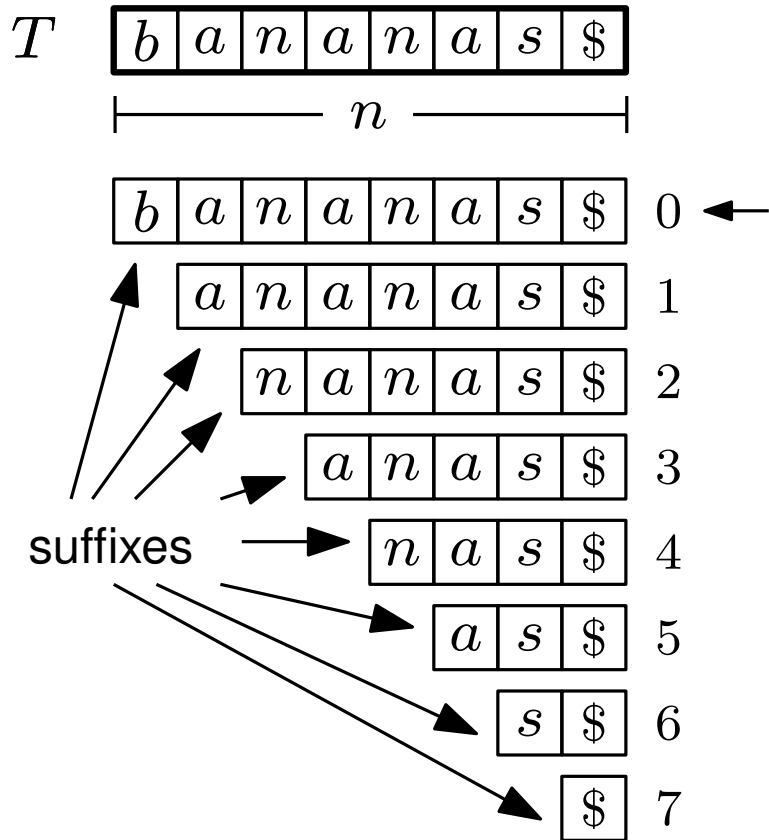
# Naively constructing a compacted suffix tree



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

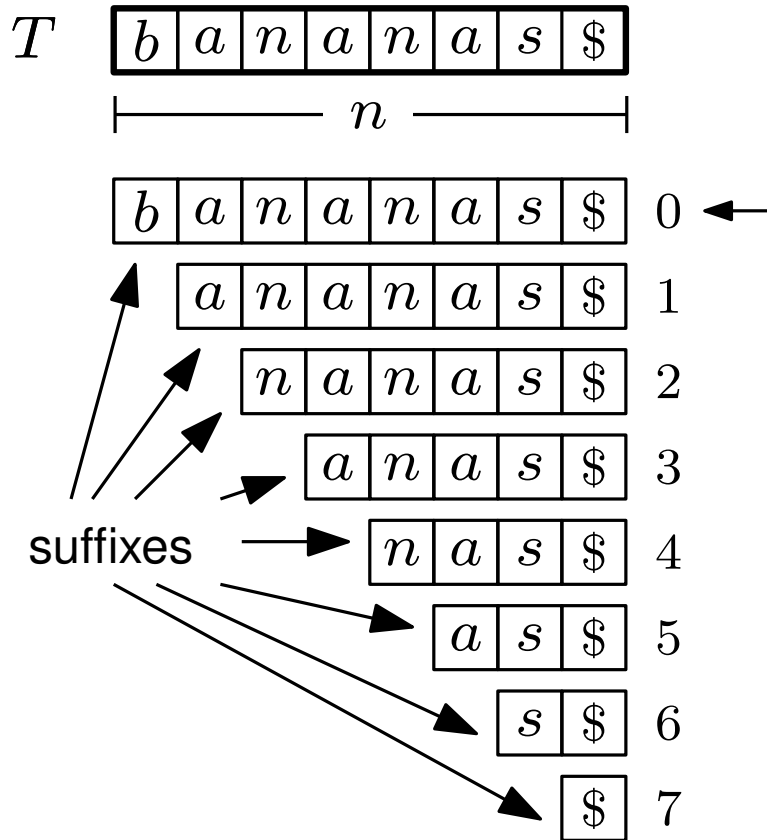
you should never actually do it like this **Naively** constructing a compacted suffix tree



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this **Naively** constructing a compacted suffix tree

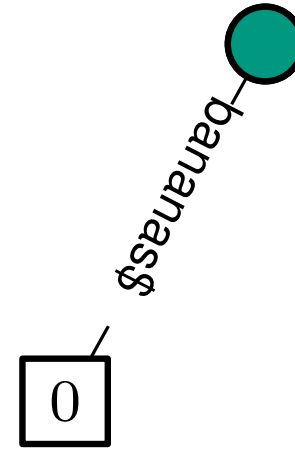
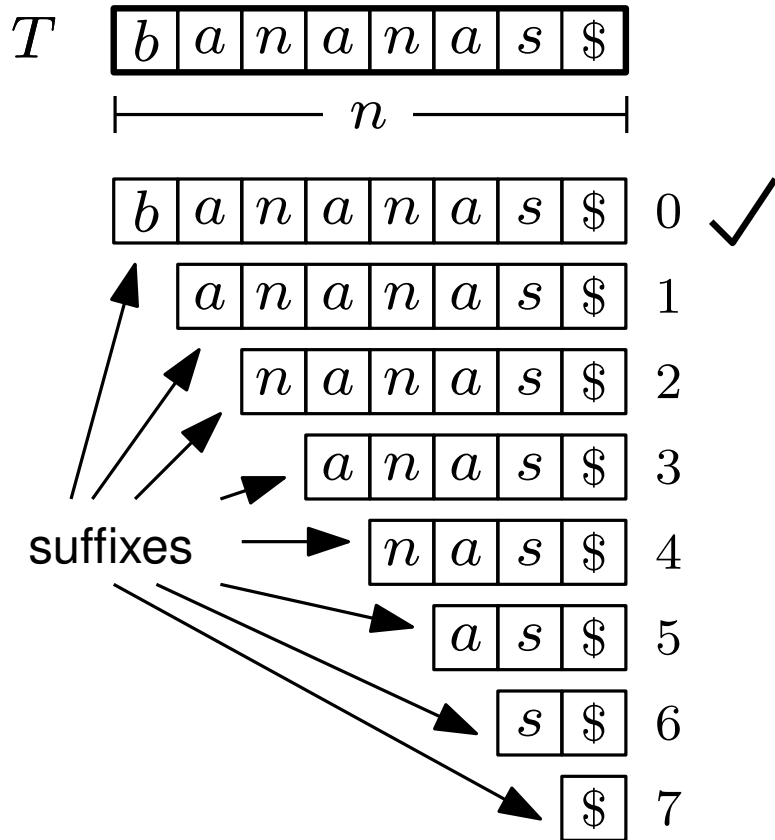


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

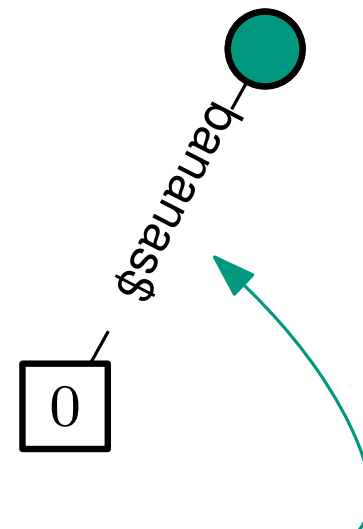
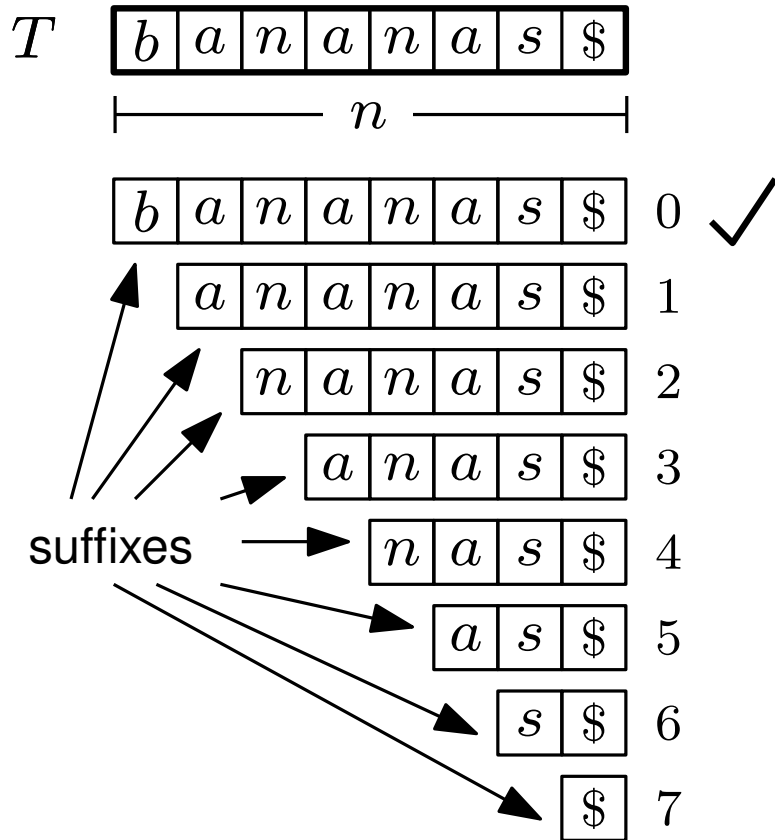


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree



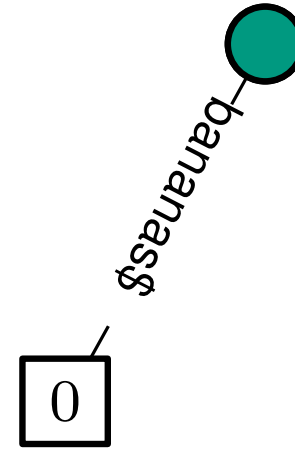
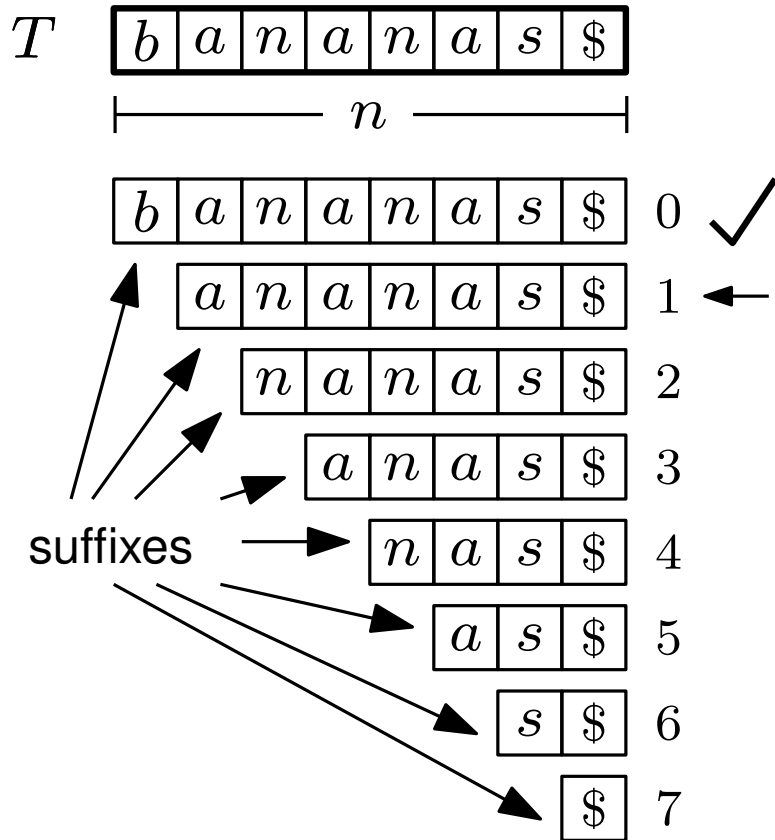
we actually store this as (0, 7)

Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree



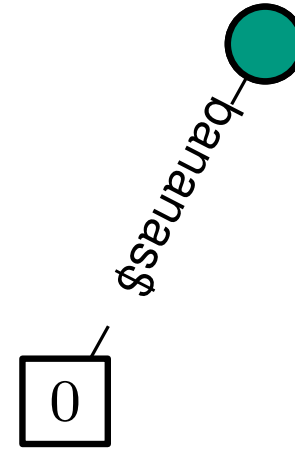
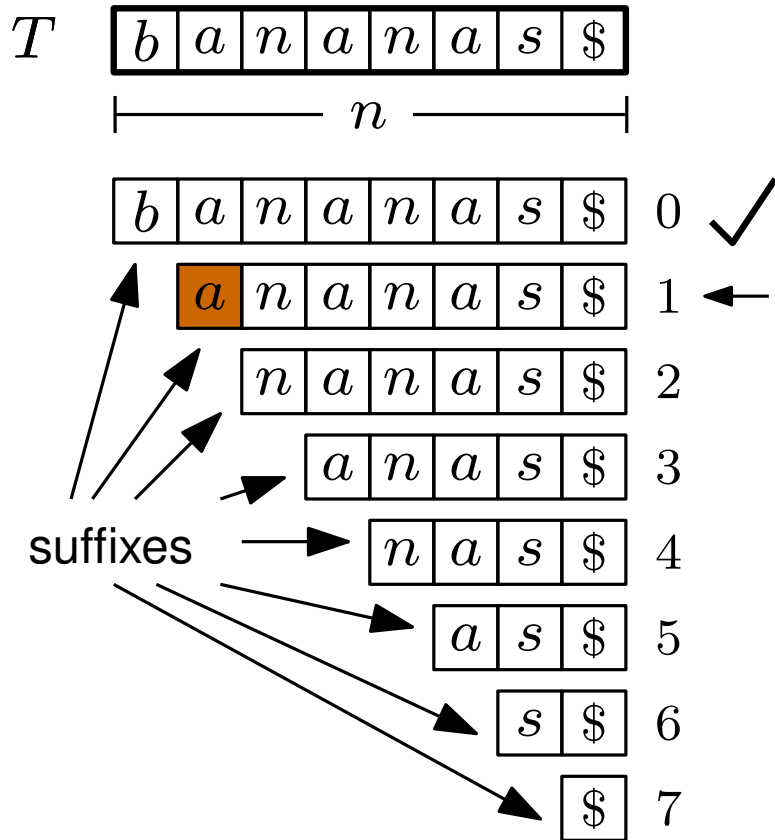
Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*



you should never actually do it like this

# Naively constructing a compacted suffix tree

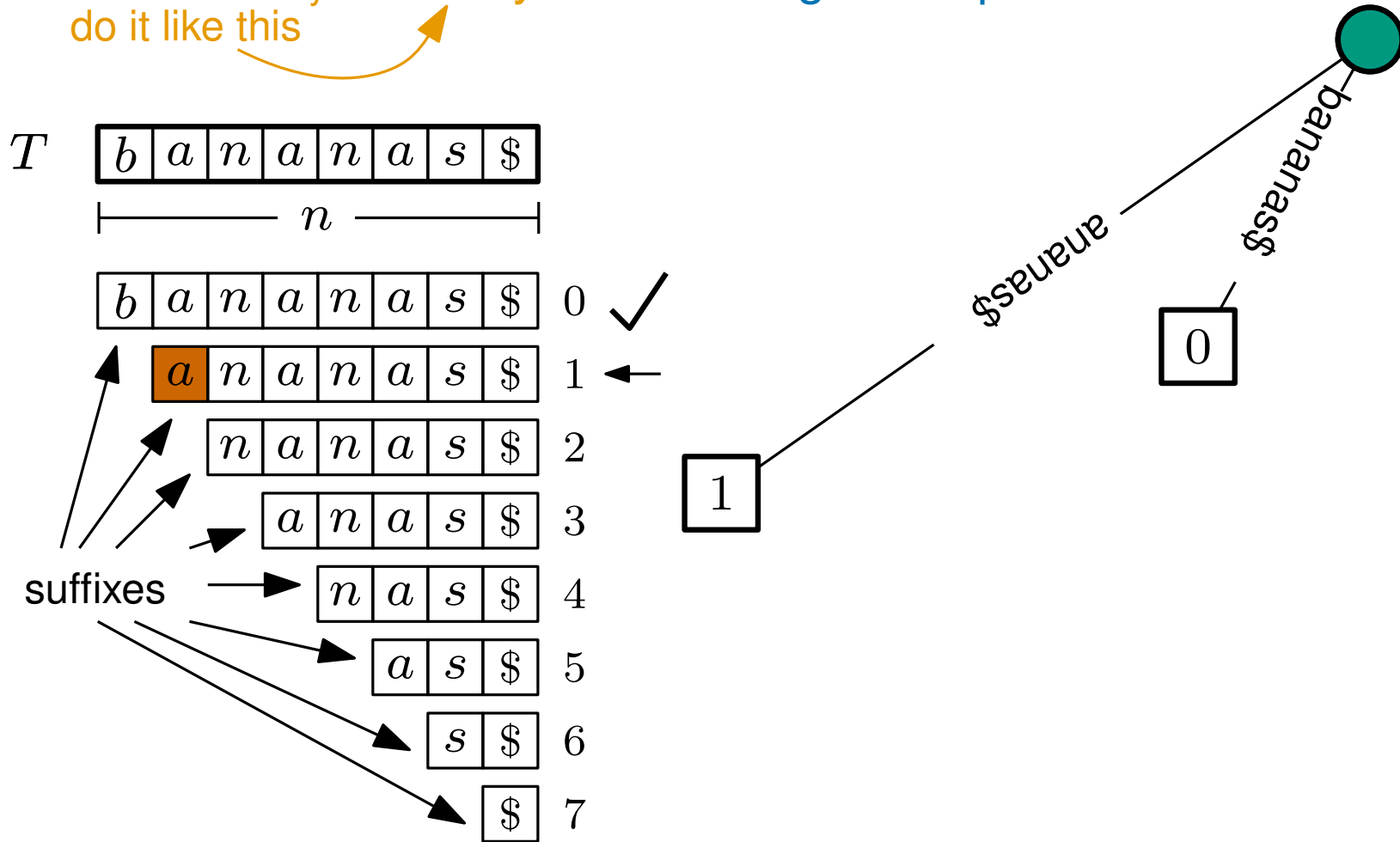


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

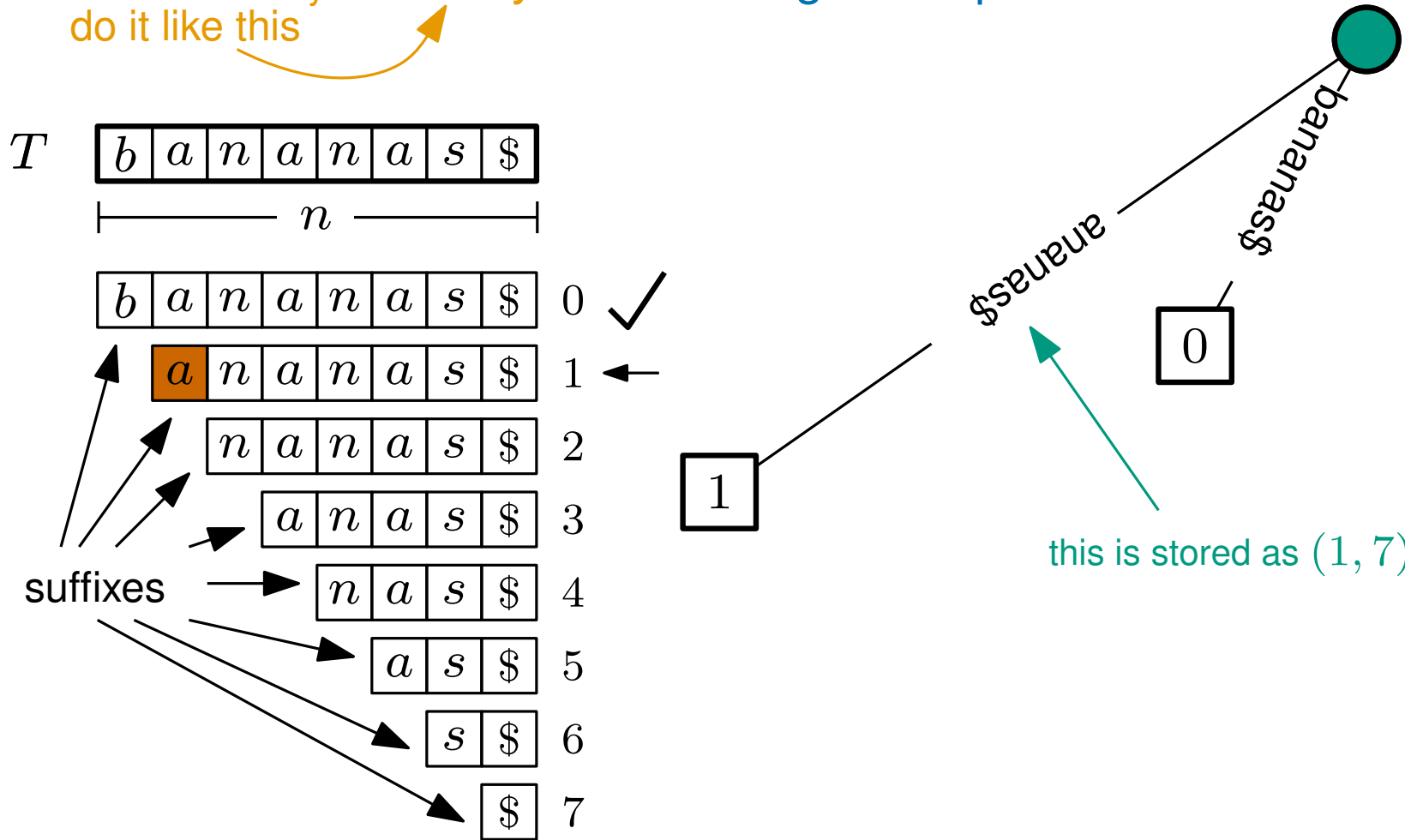


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

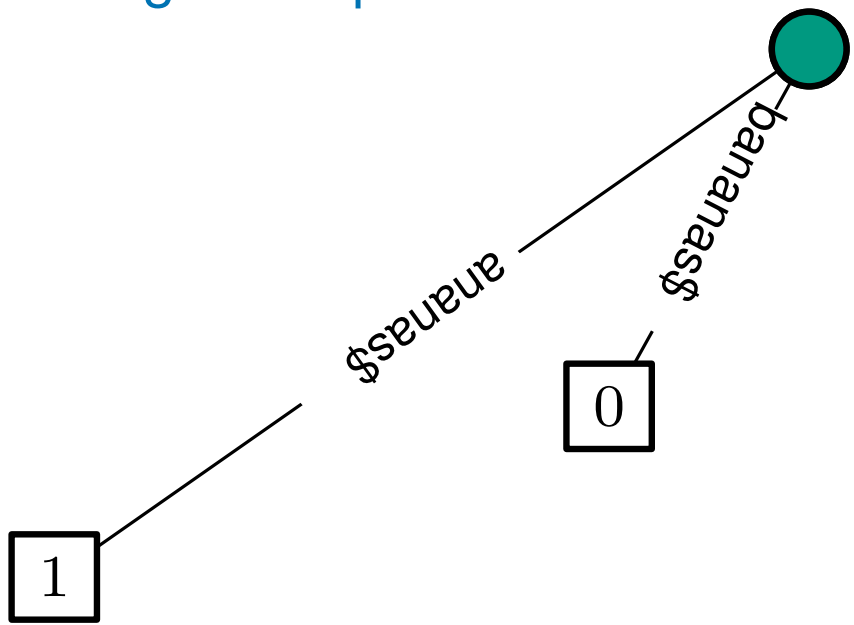
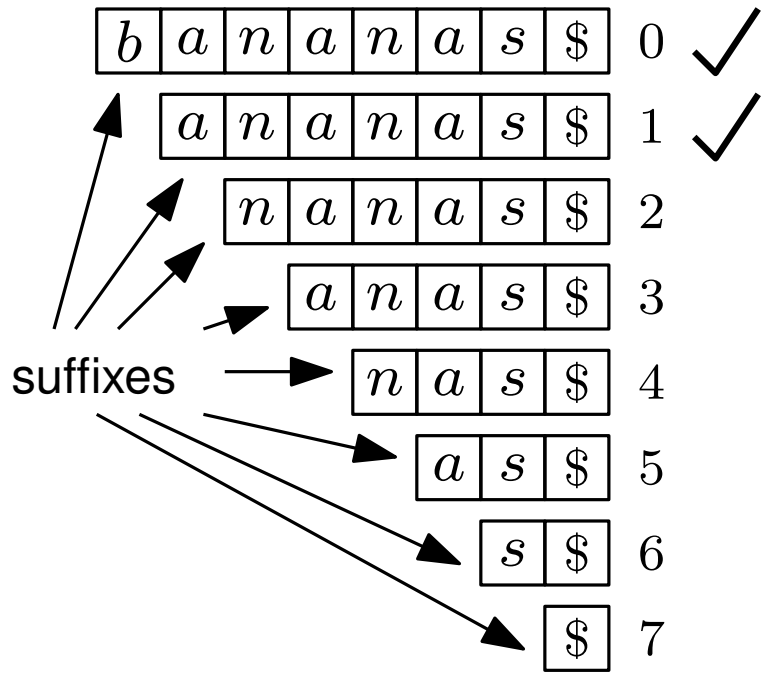
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|

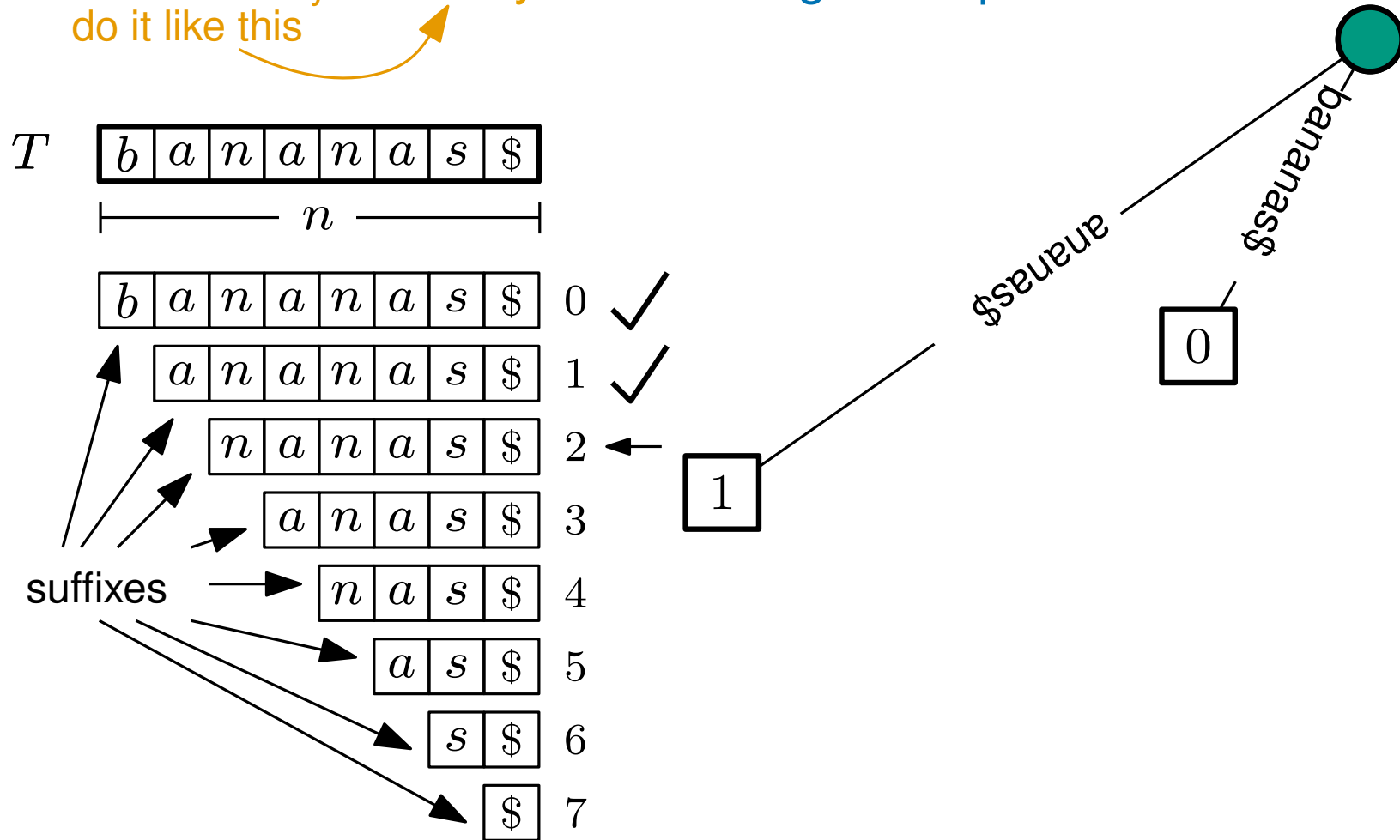


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

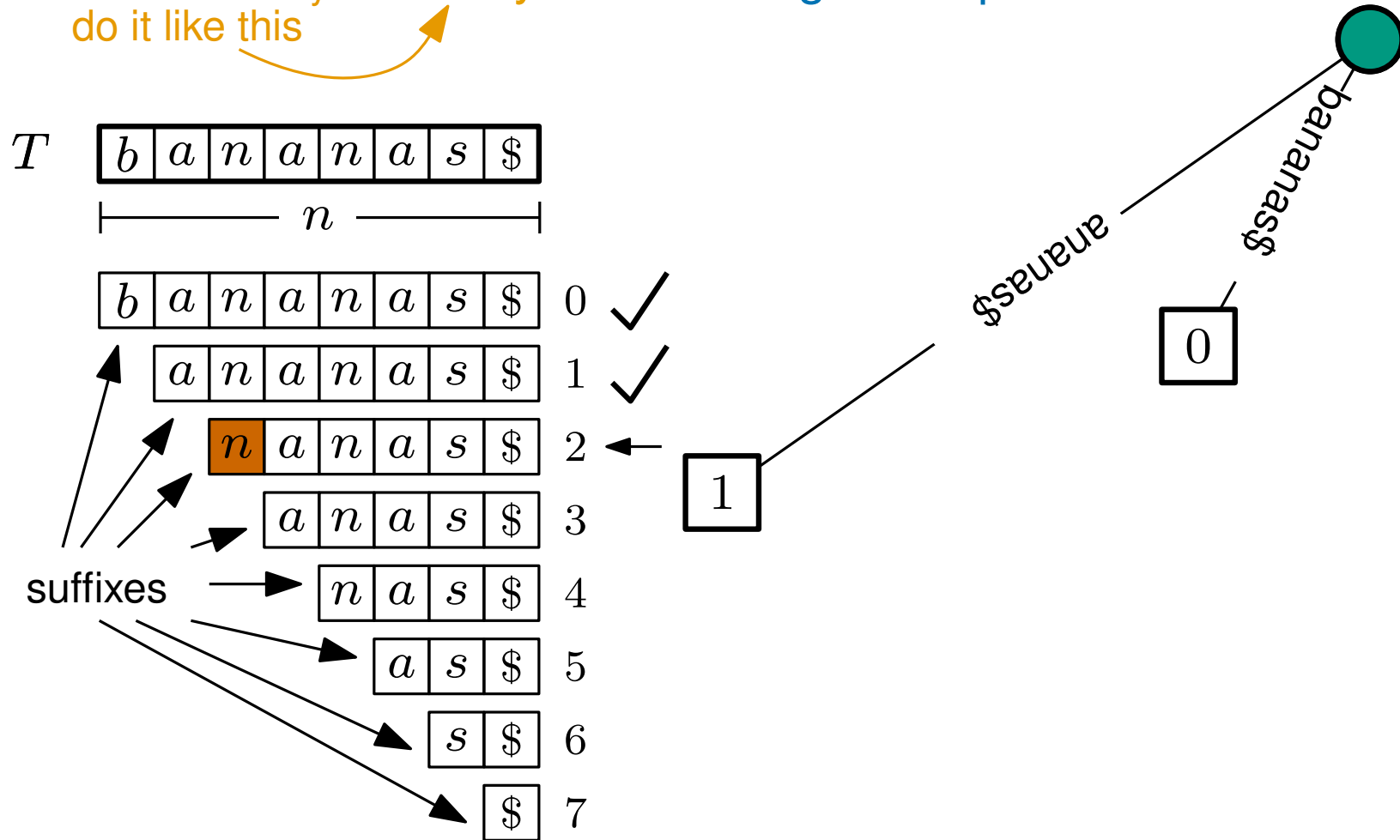


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

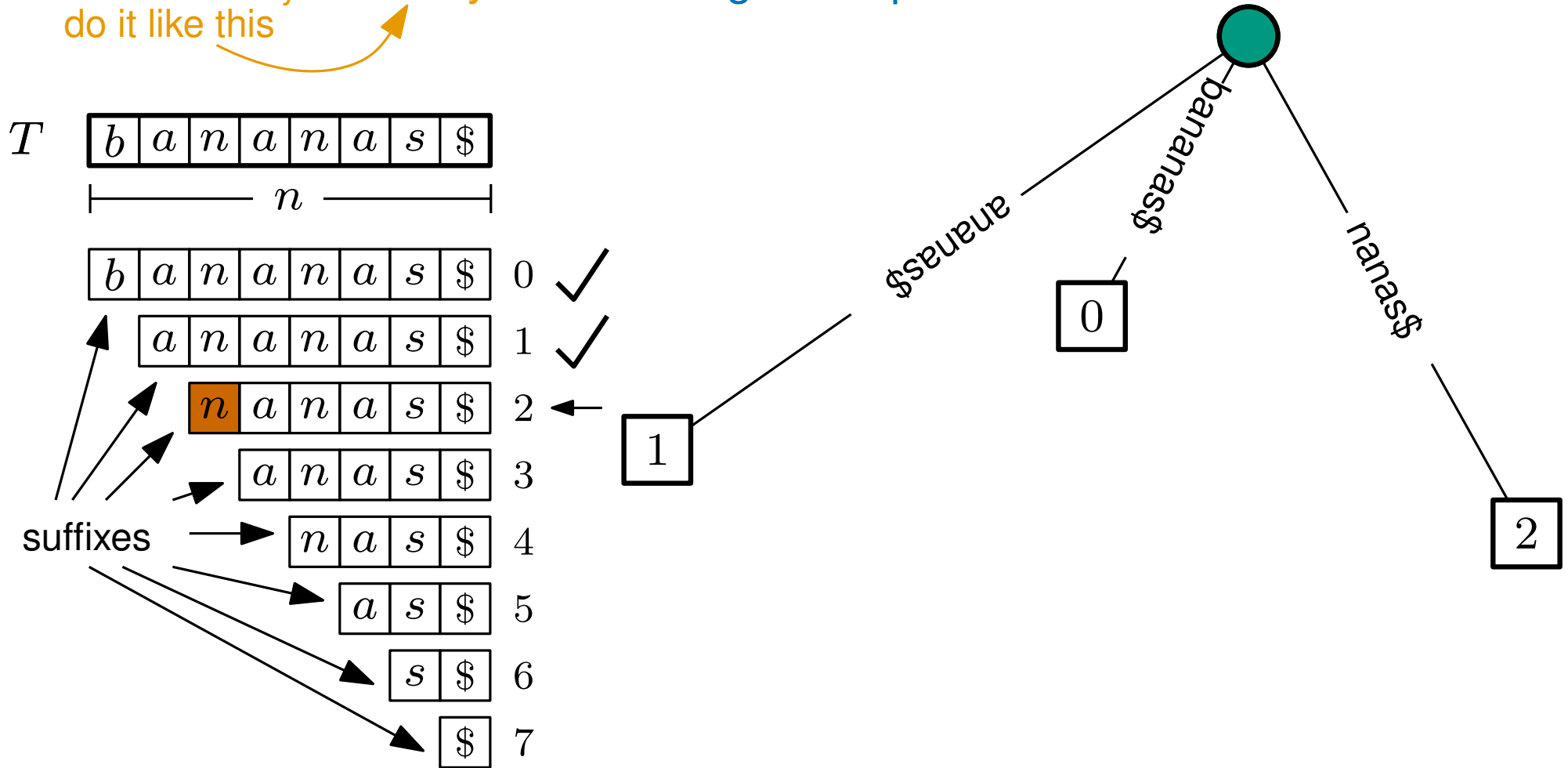


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

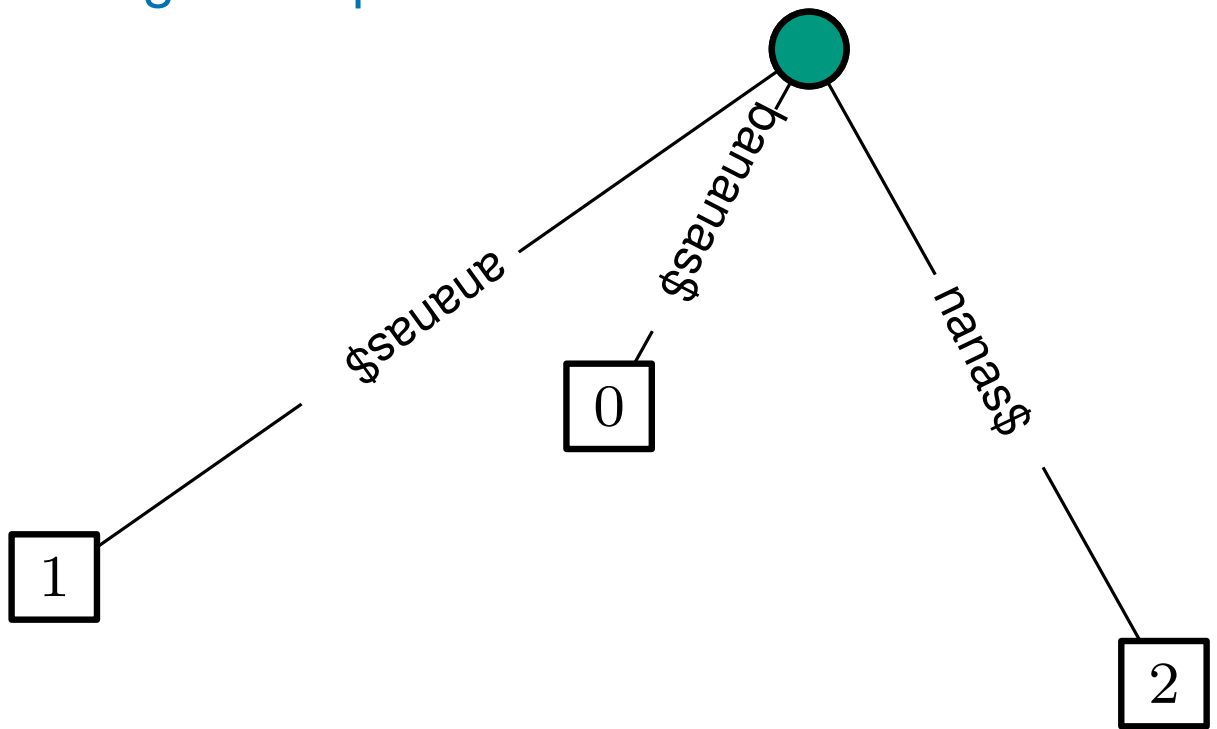
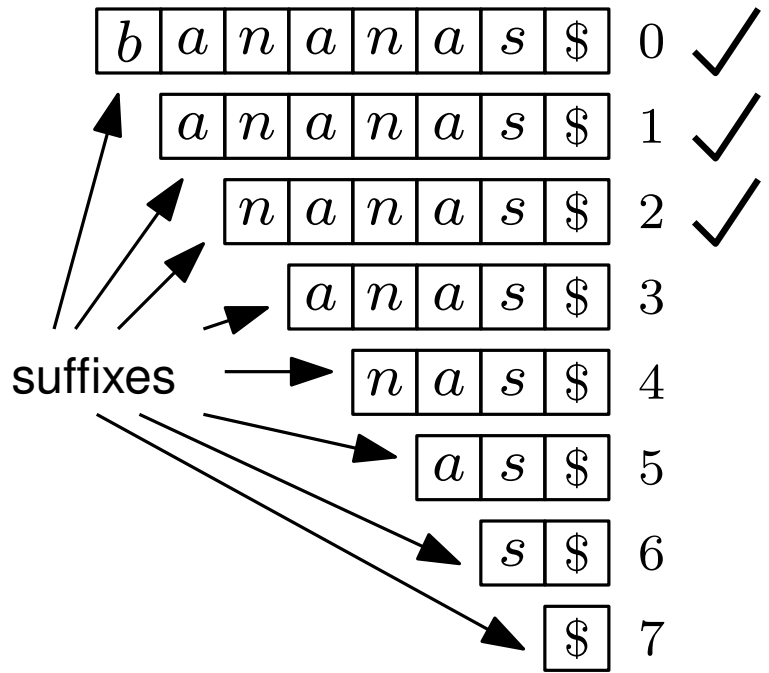
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*



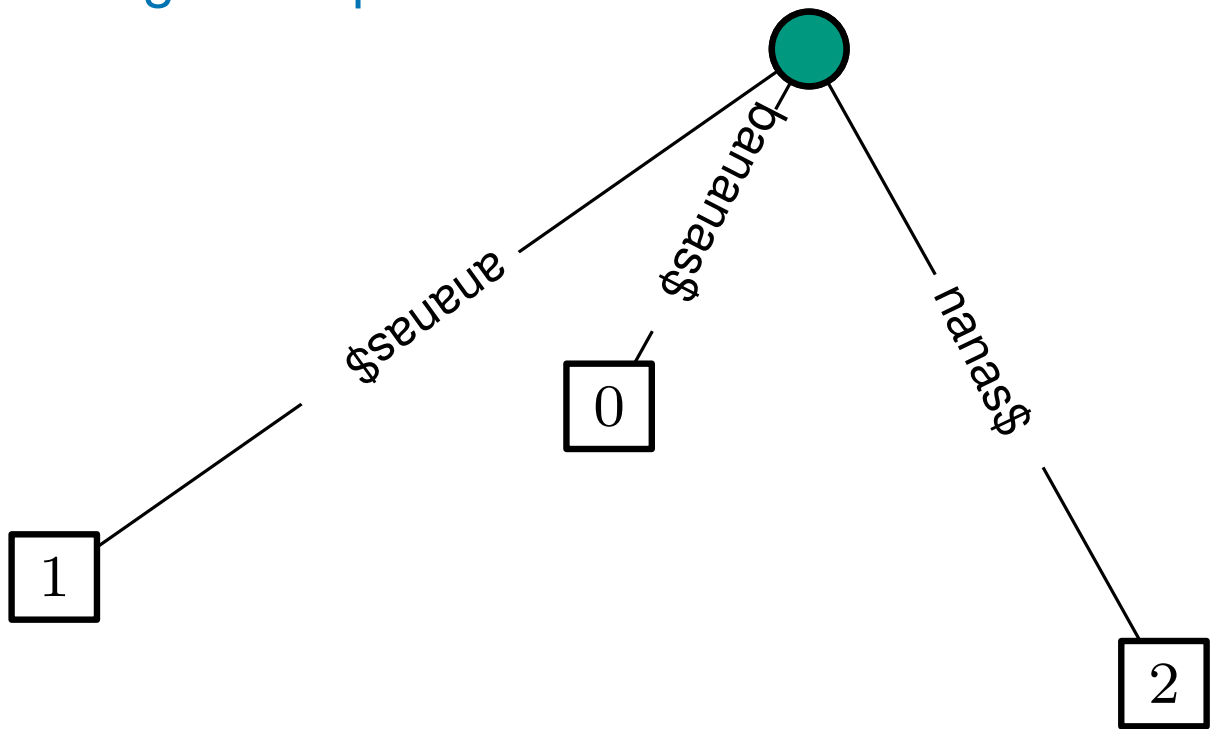
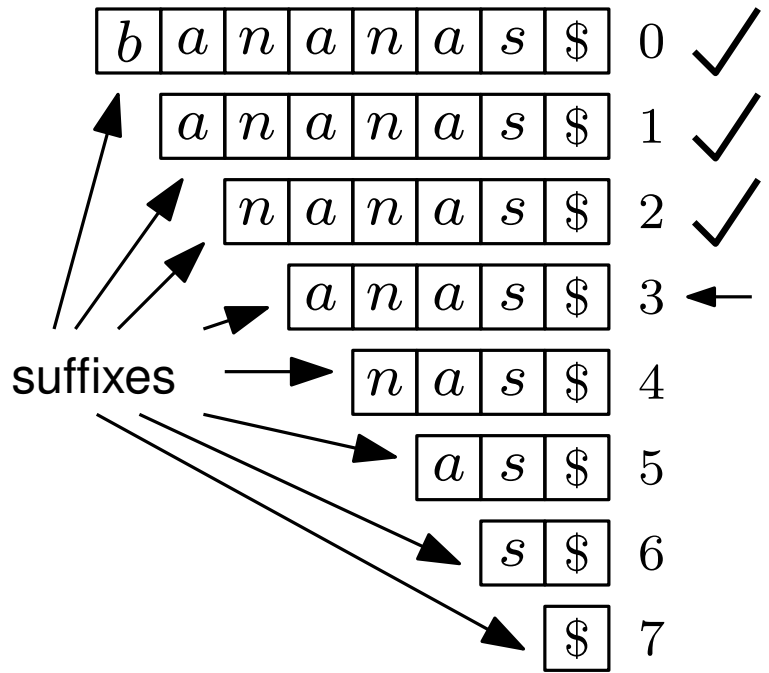
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|

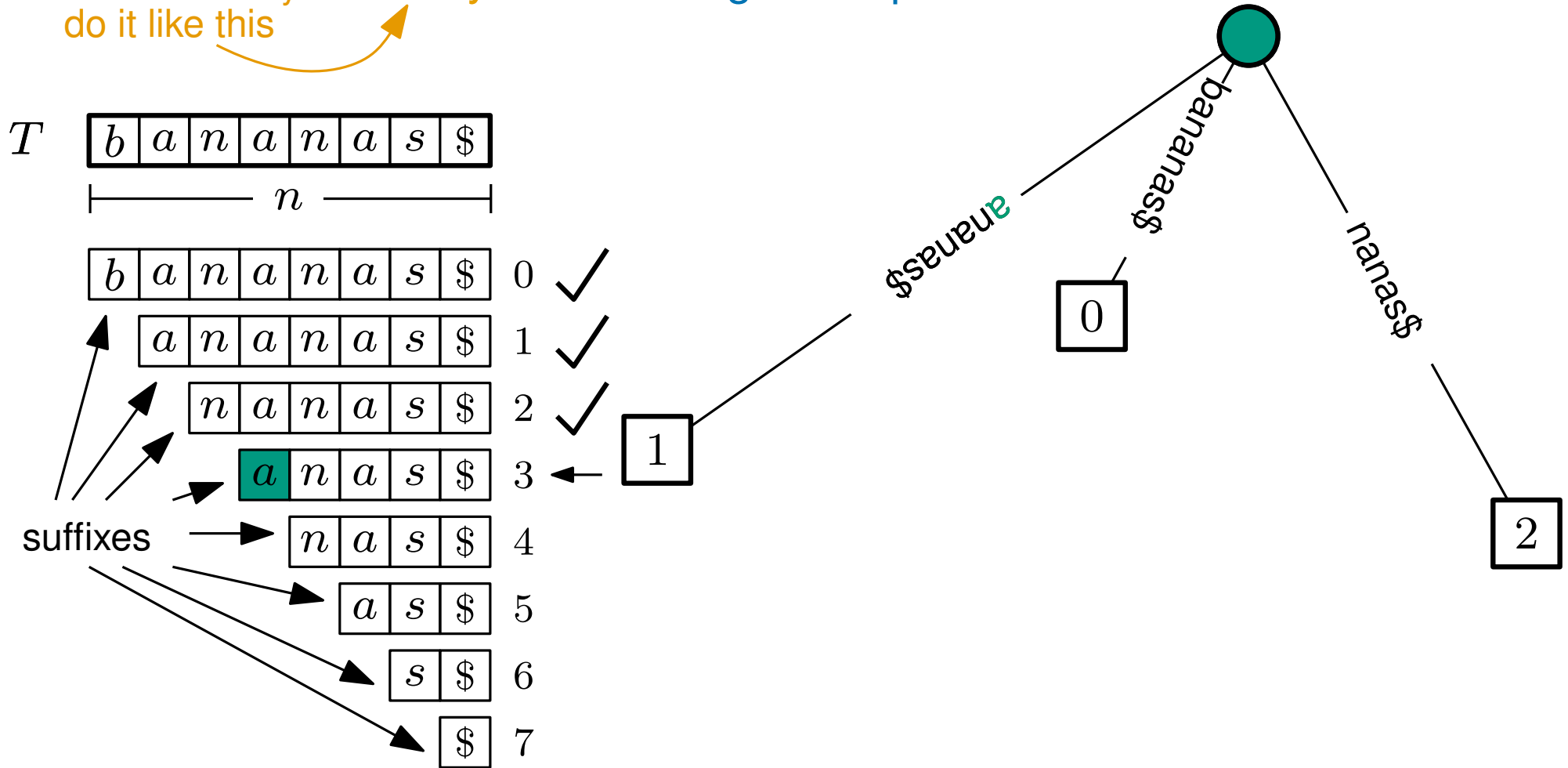


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

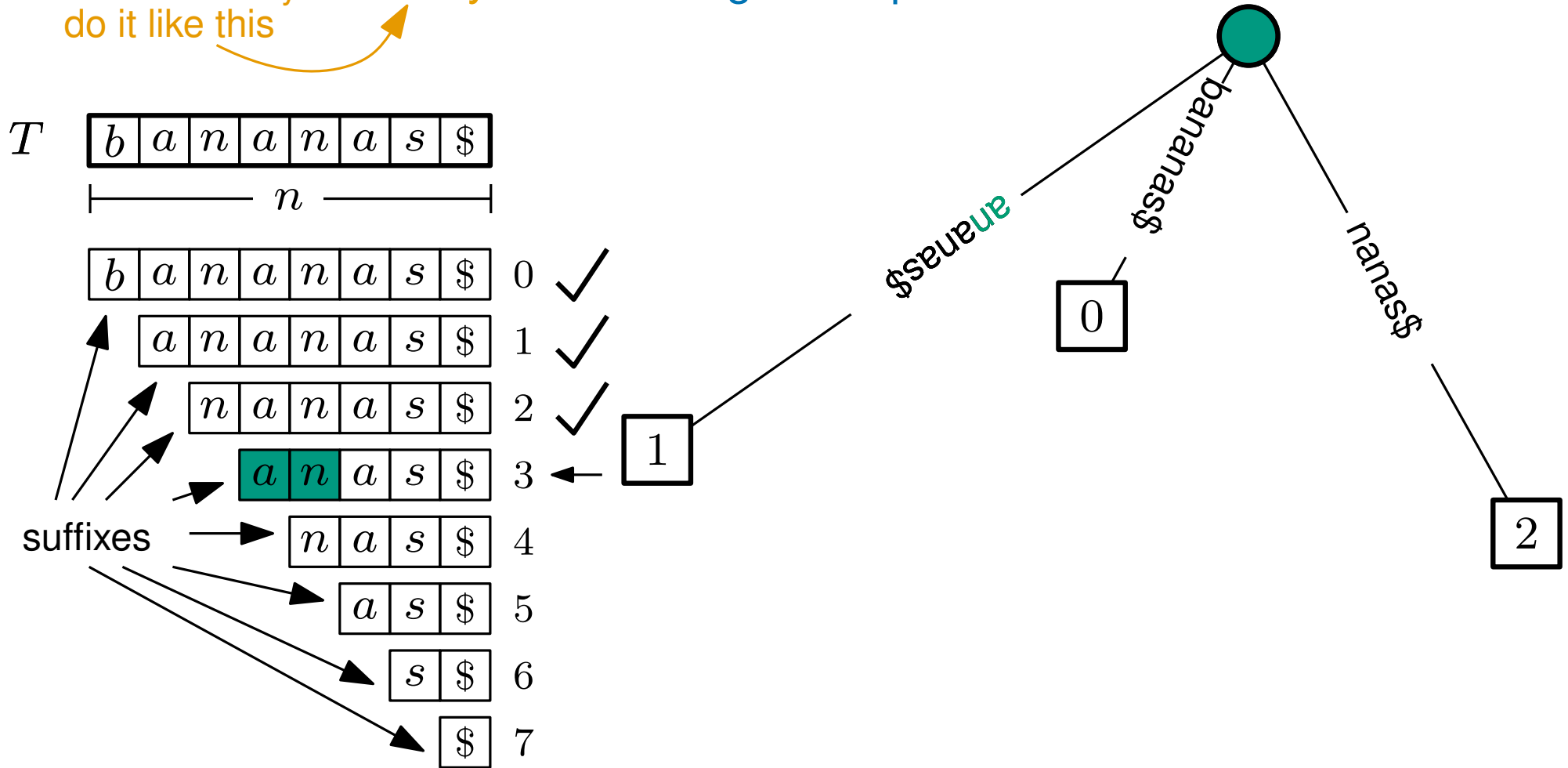


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

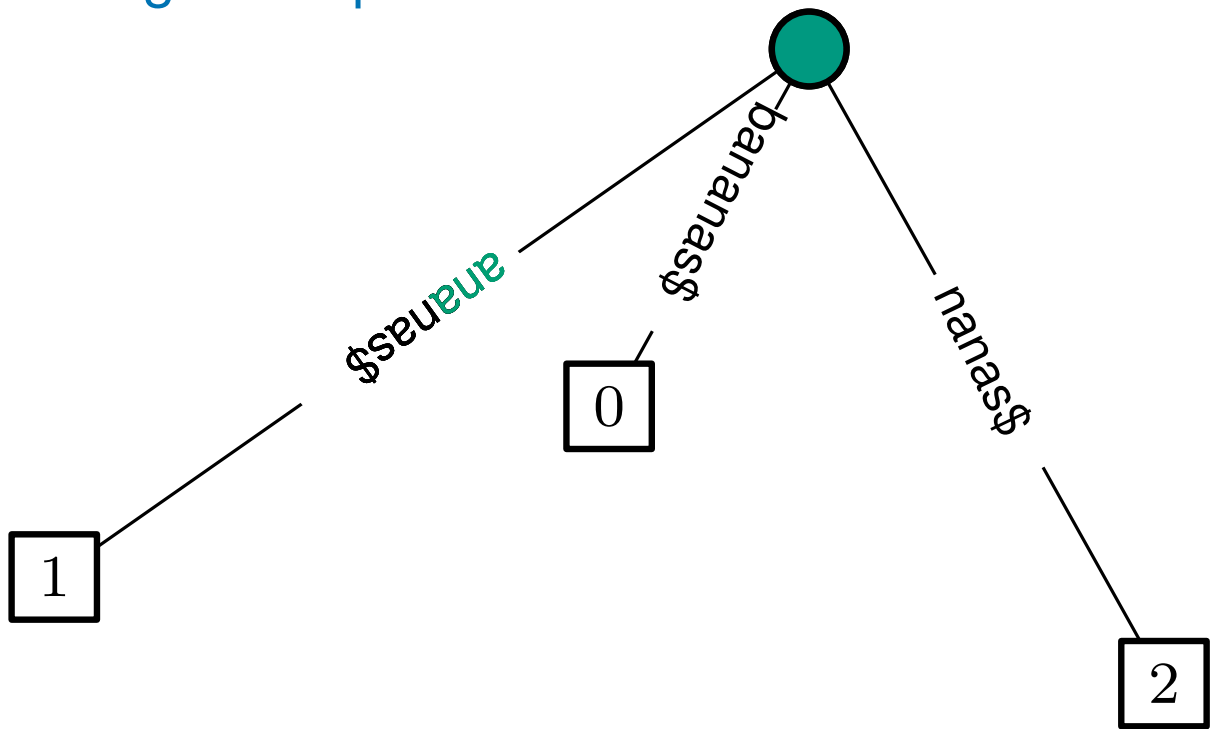
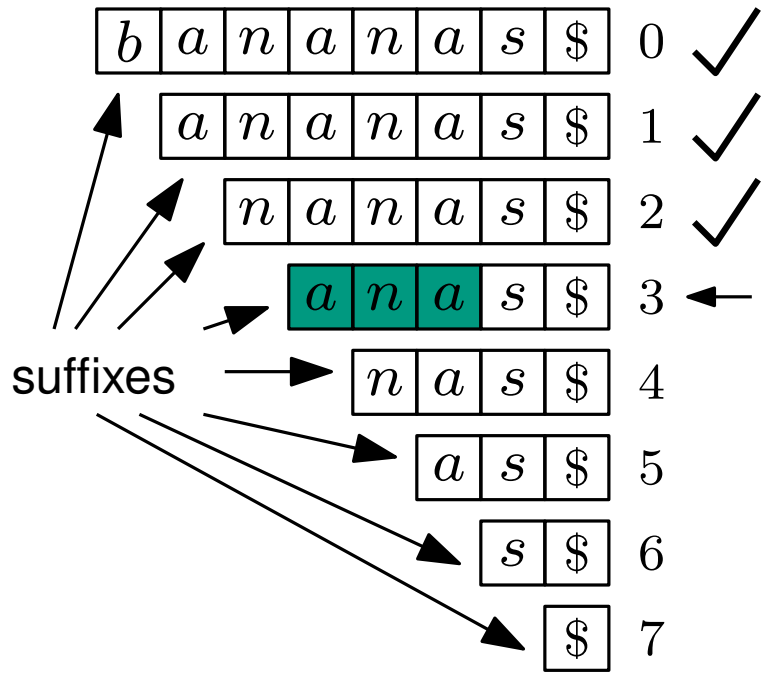
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|

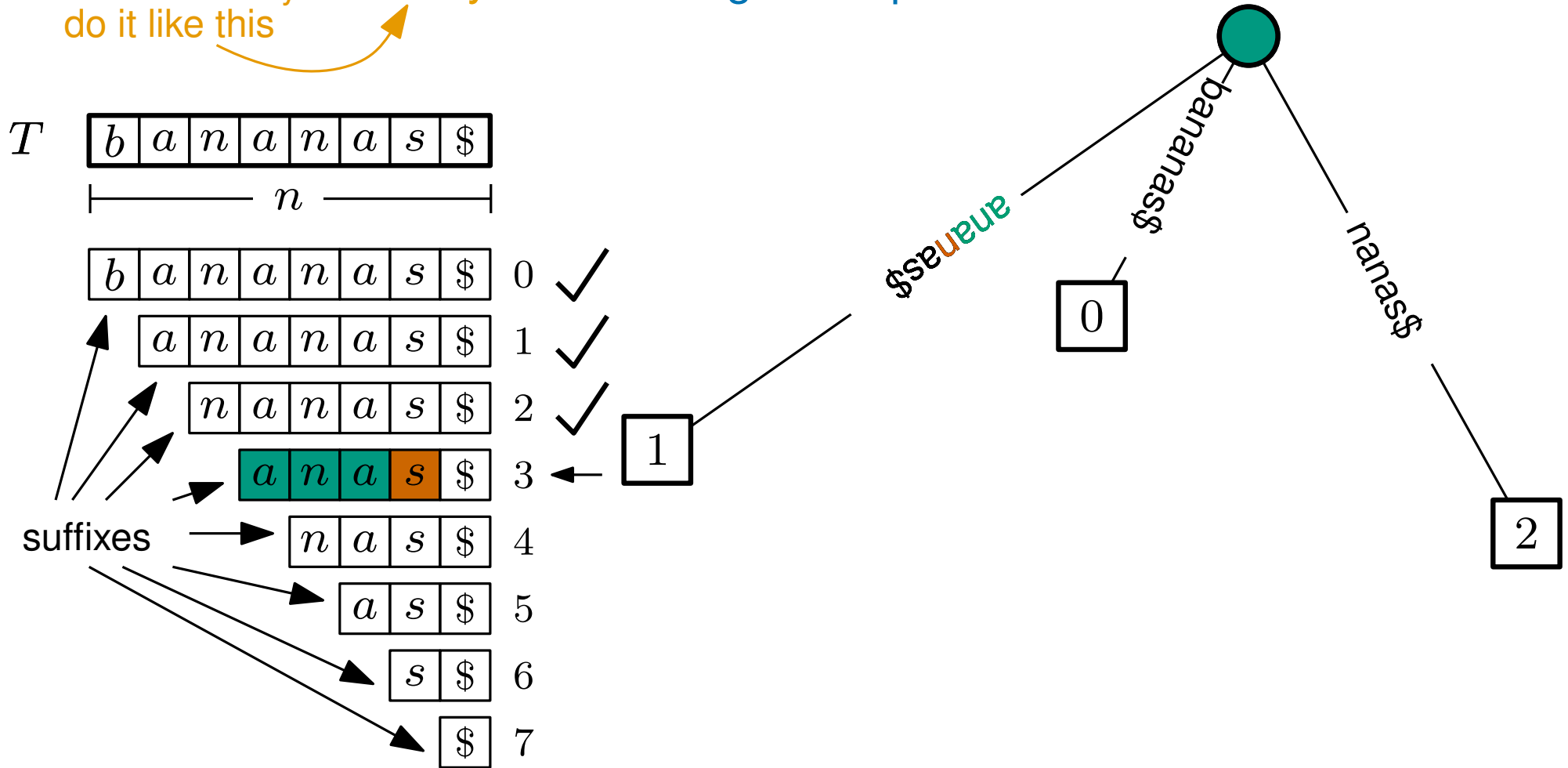


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

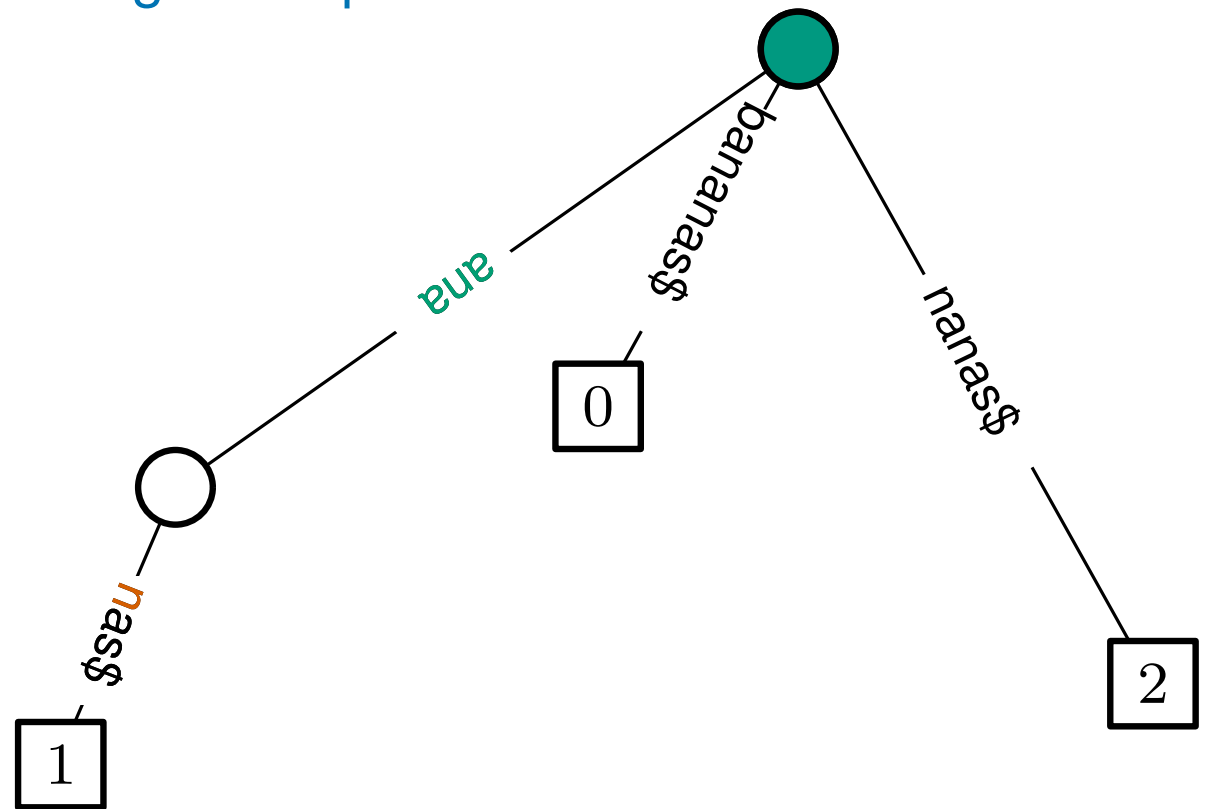
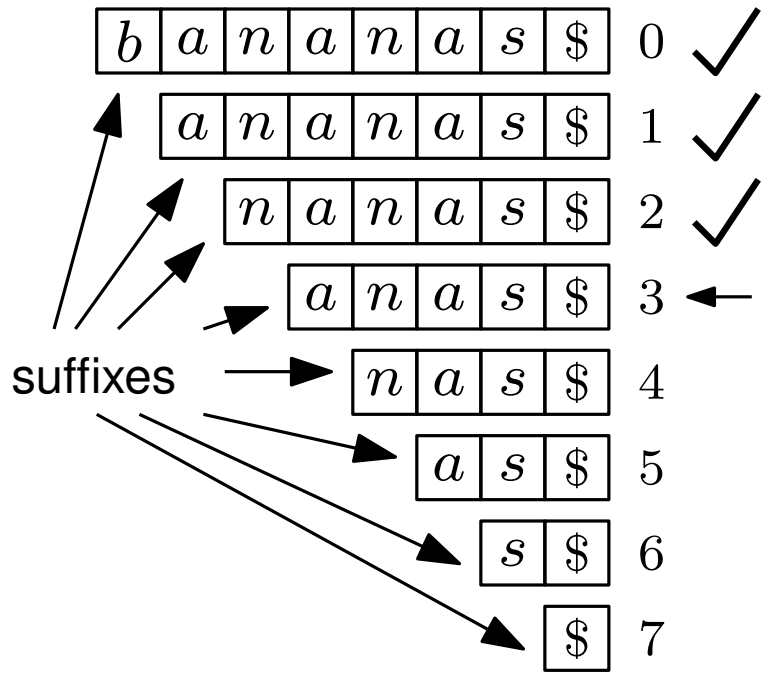
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

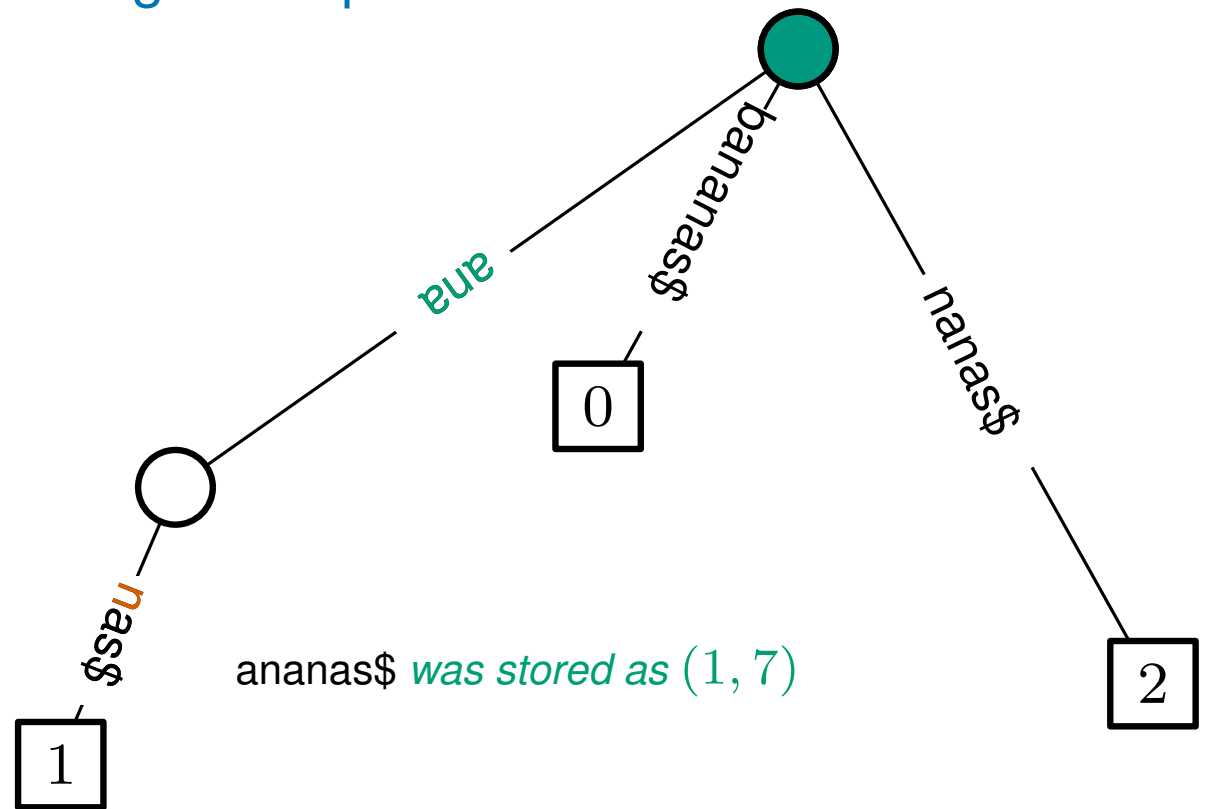
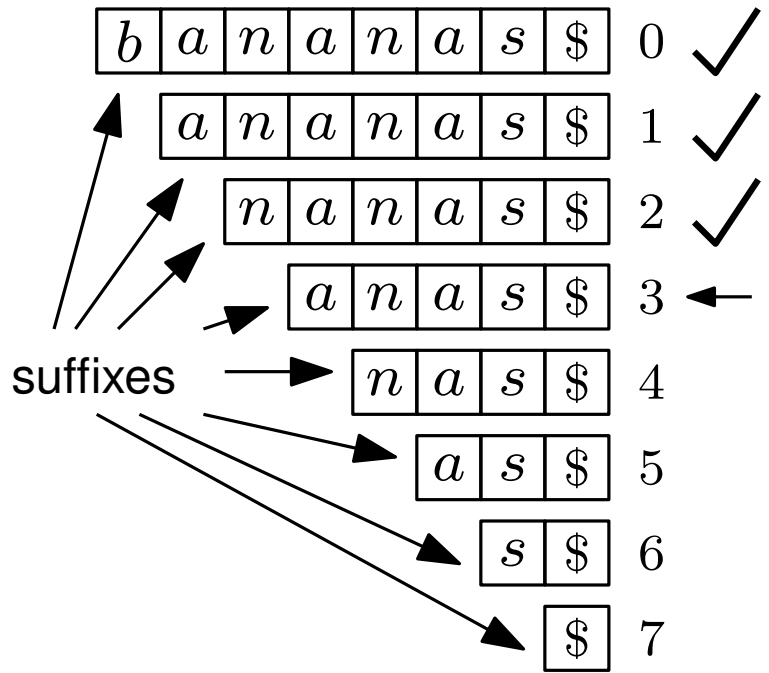
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

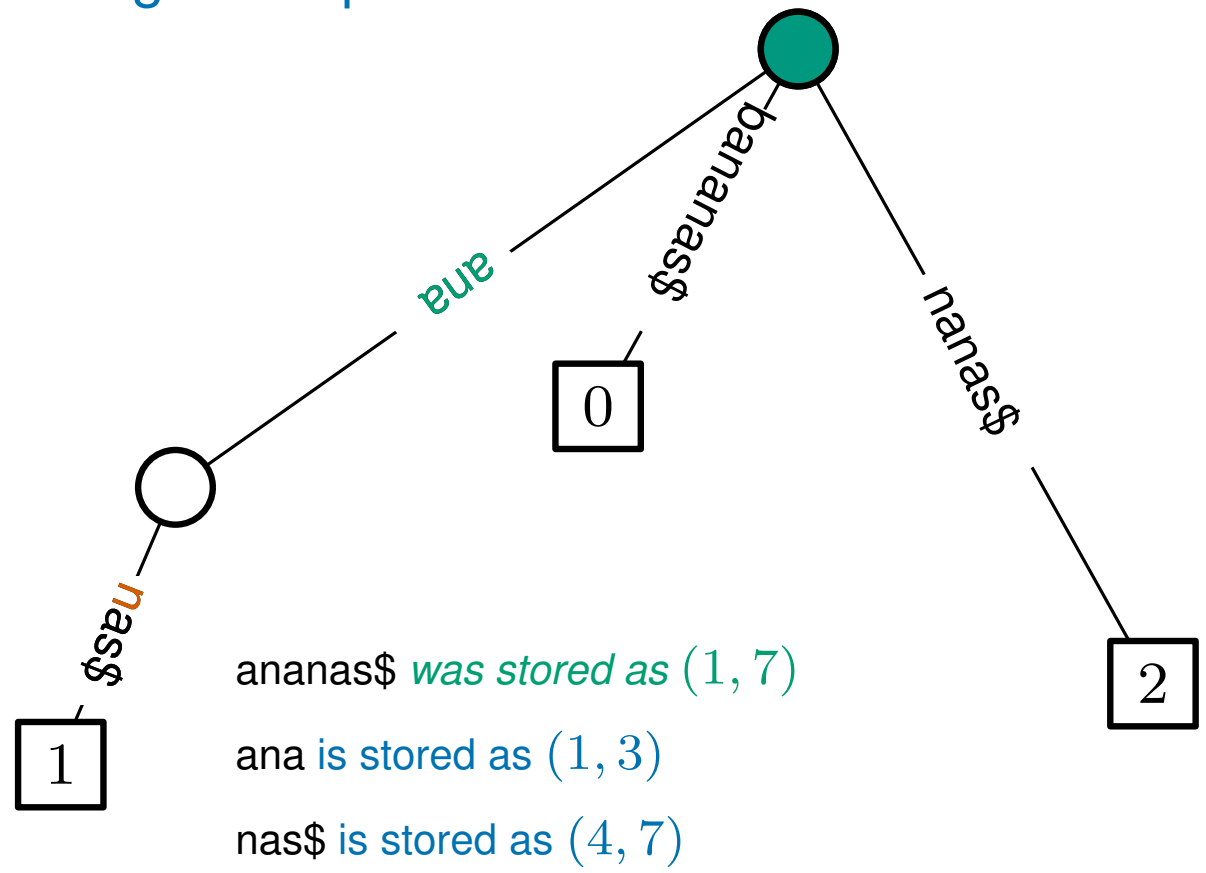
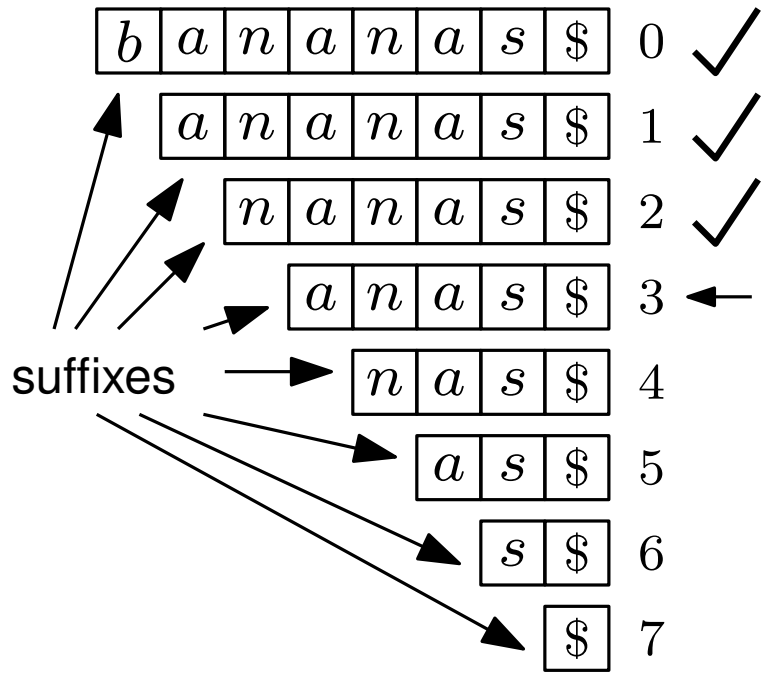
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*



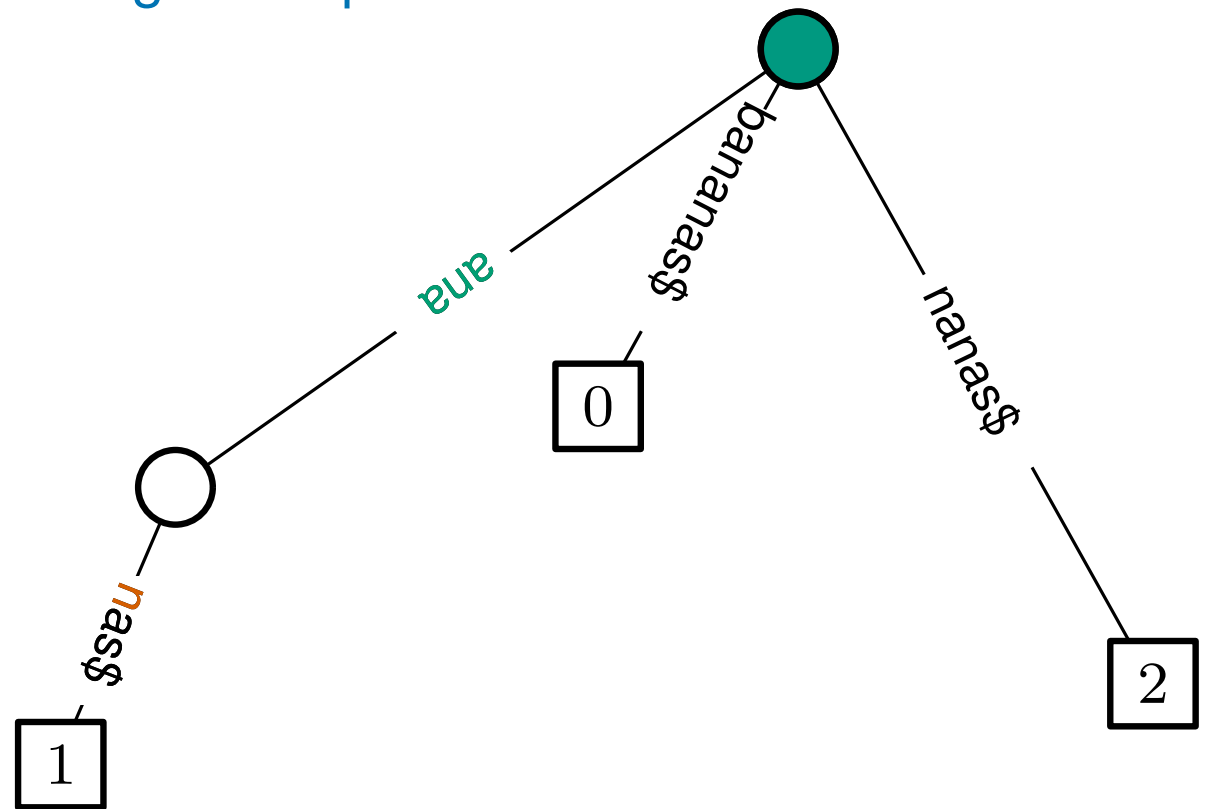
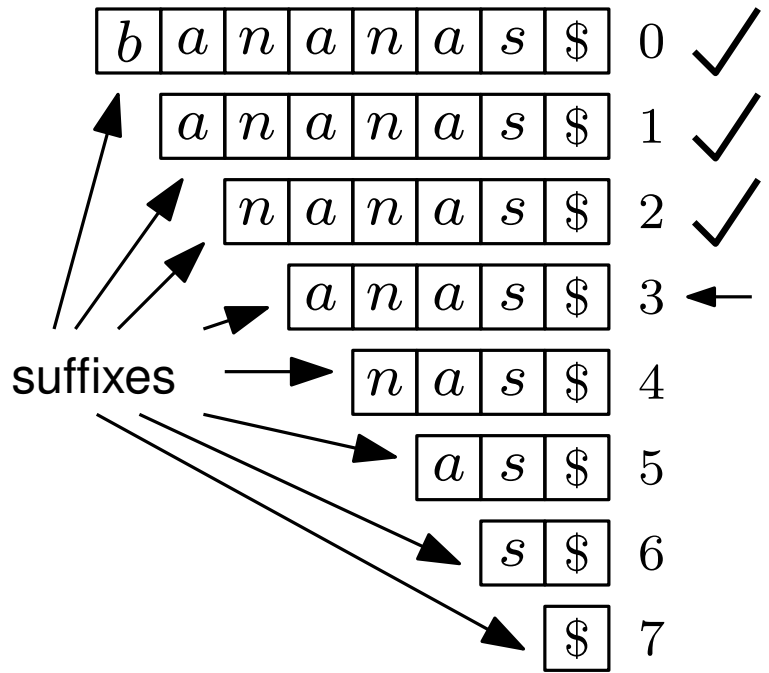
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

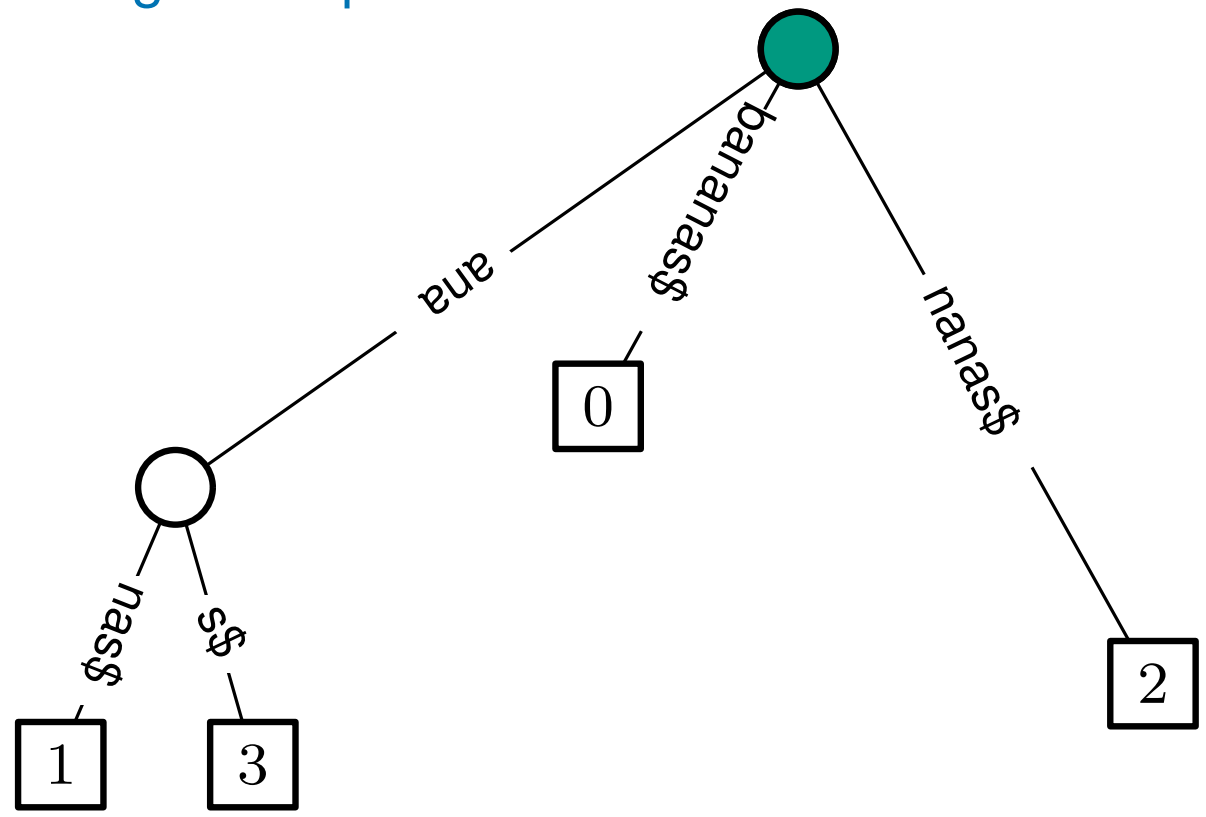
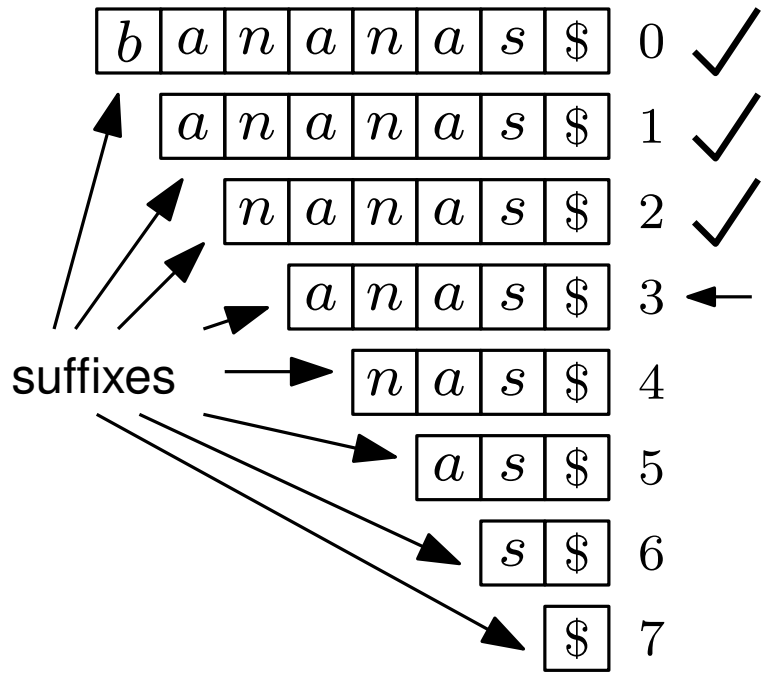
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

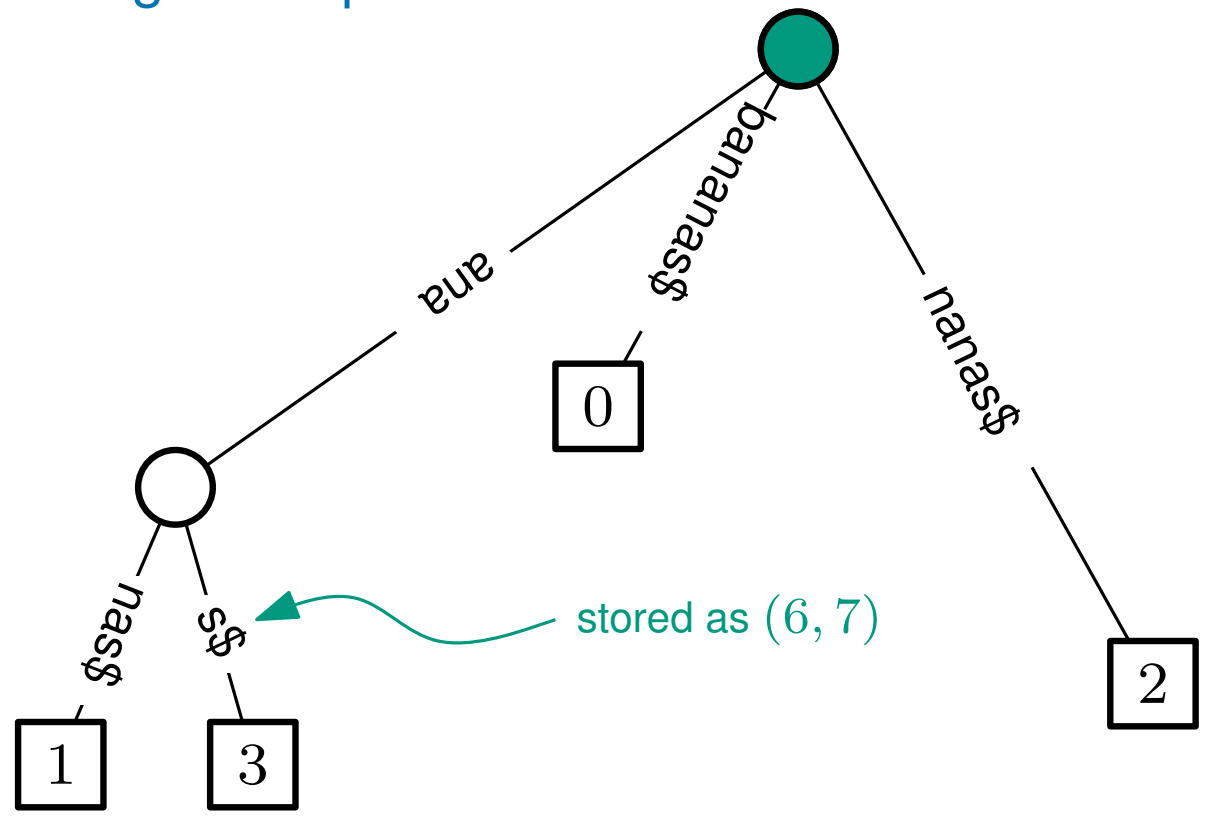
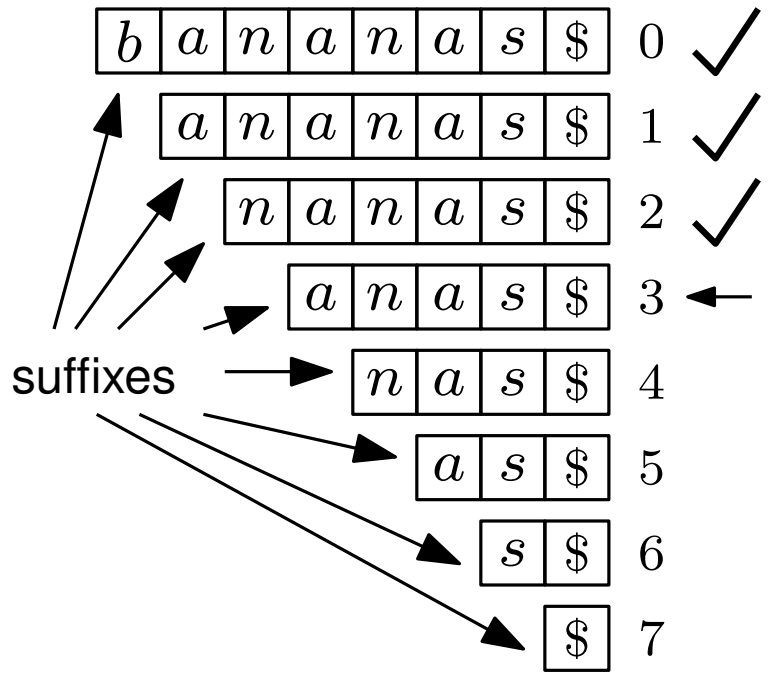
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

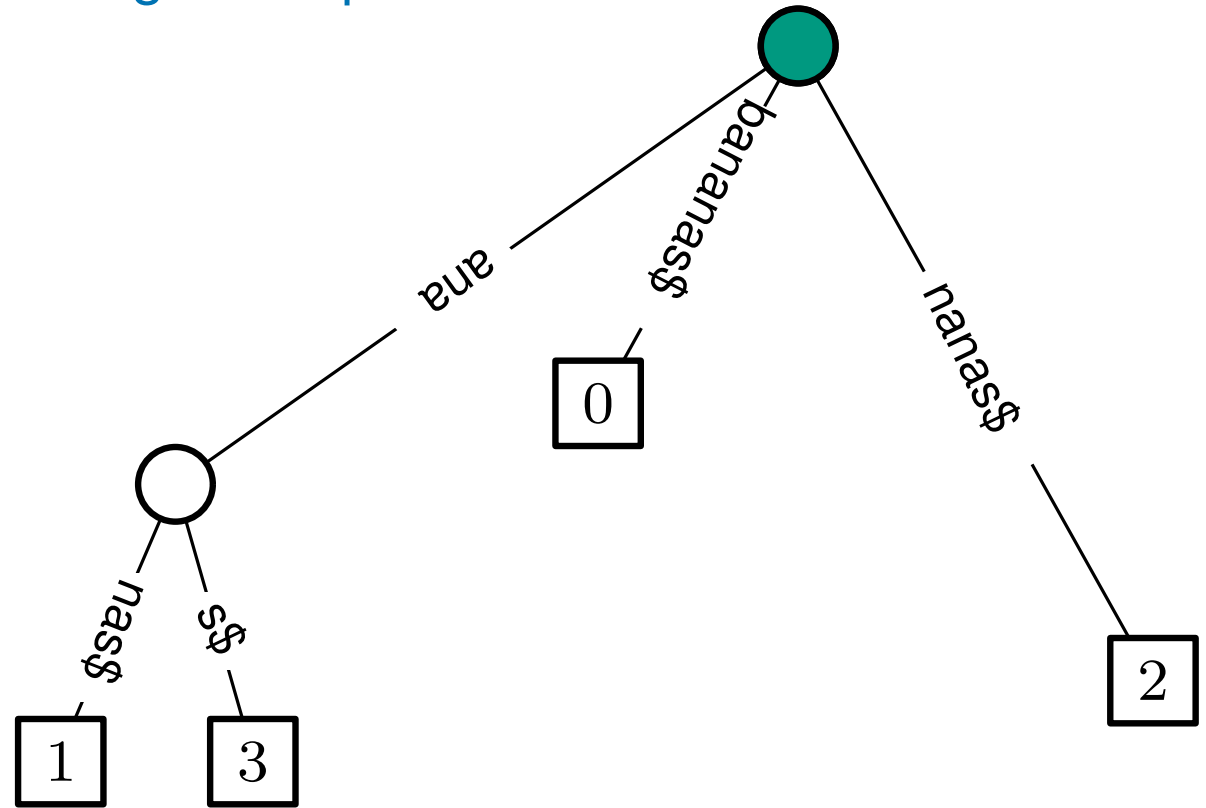
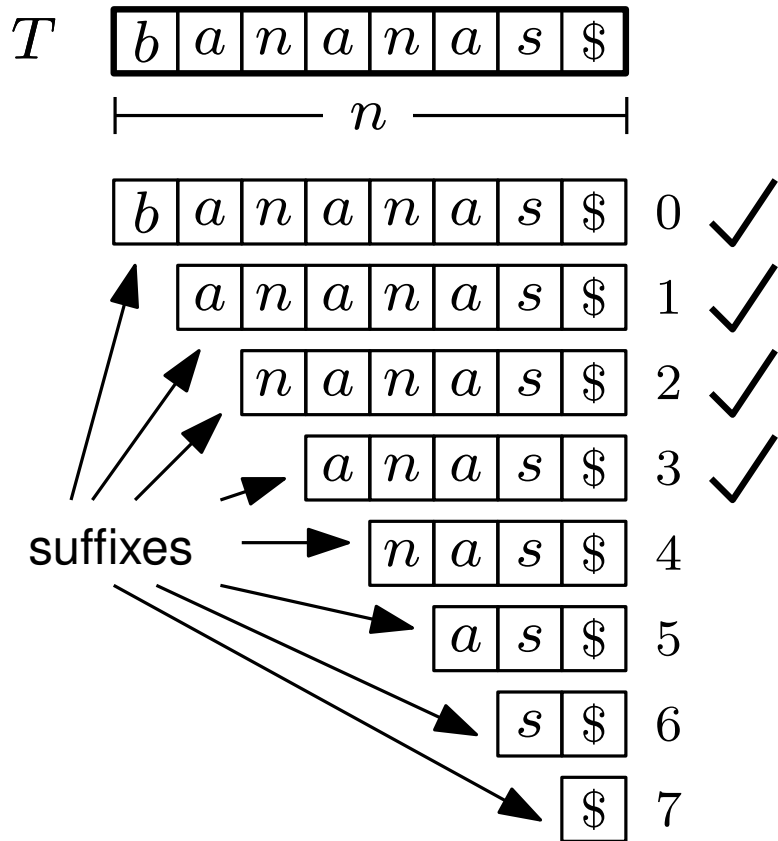
  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this **Naively** constructing a compacted suffix tree



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

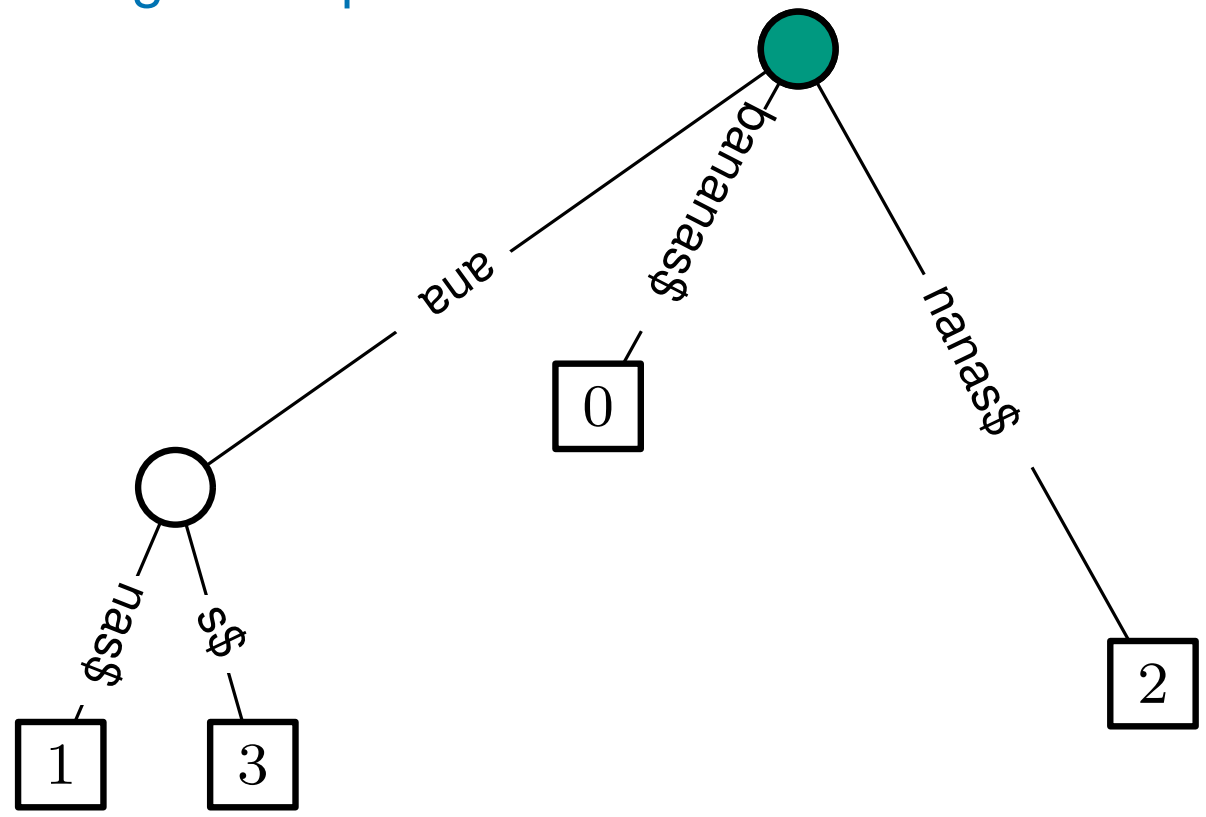
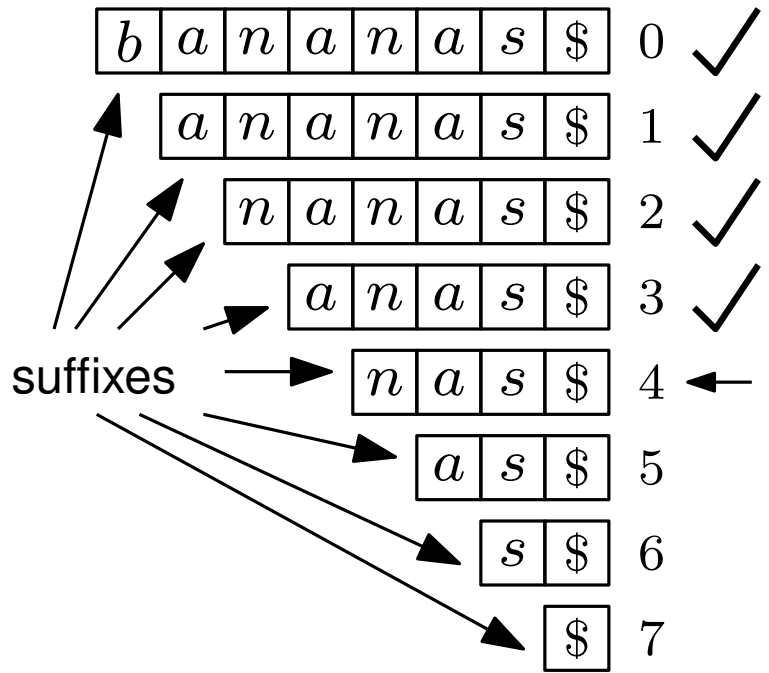
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

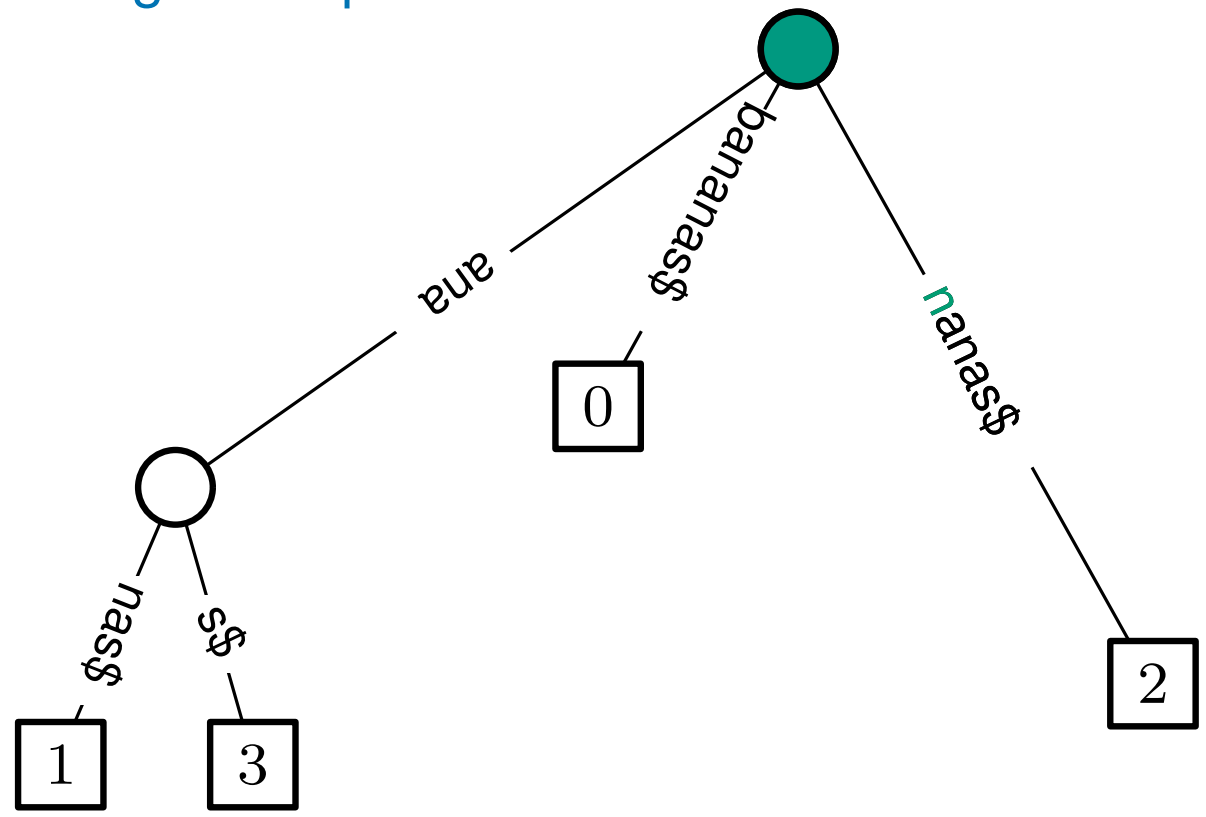
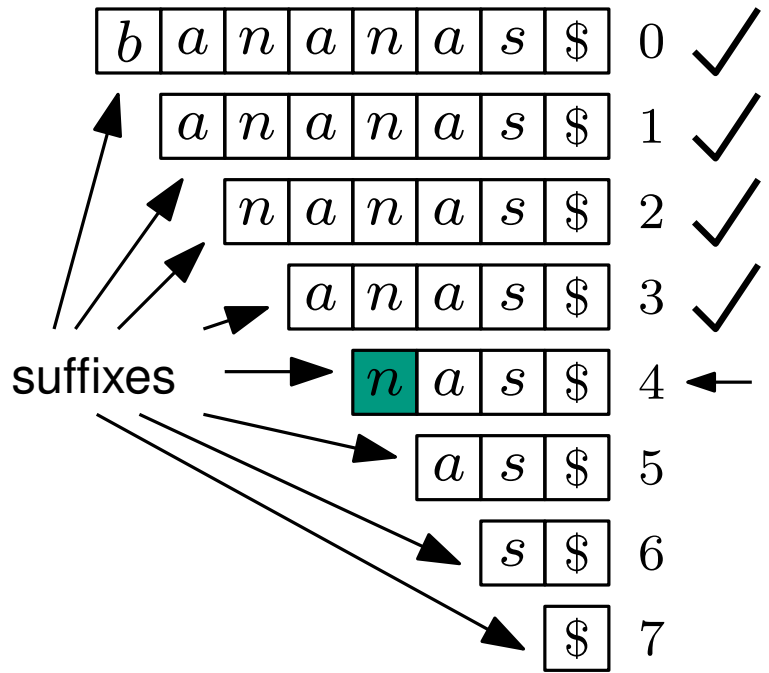
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

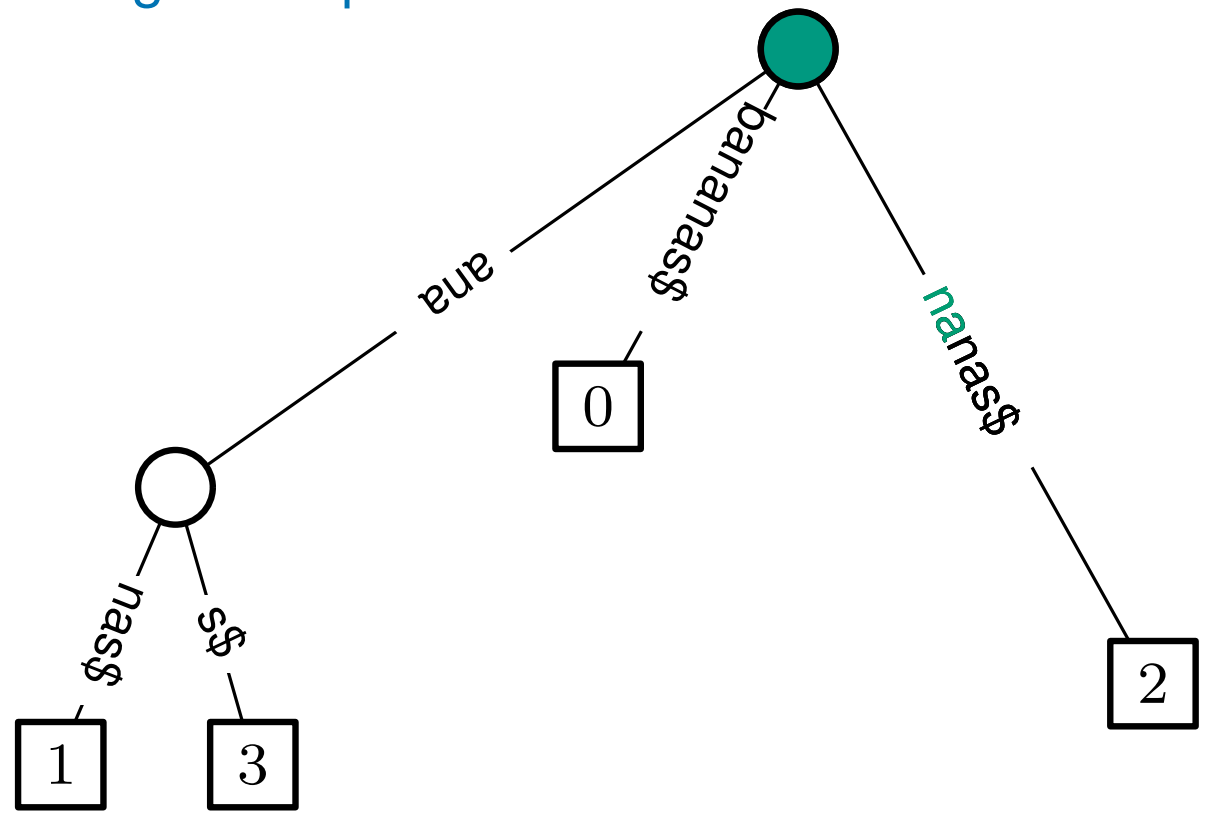
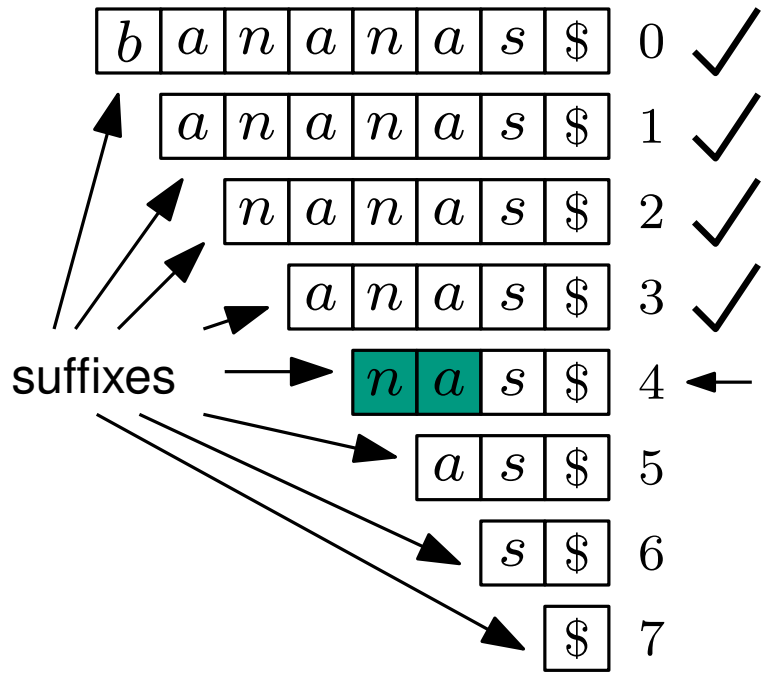
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

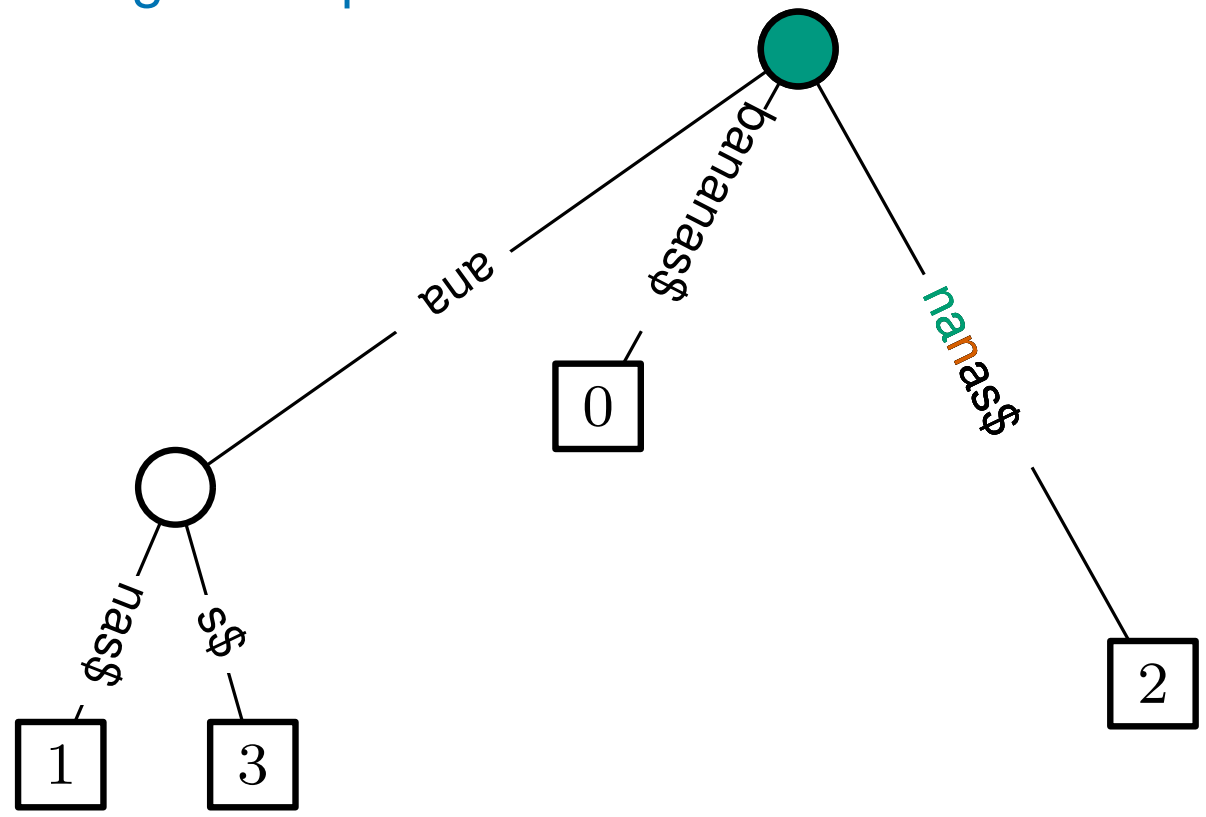
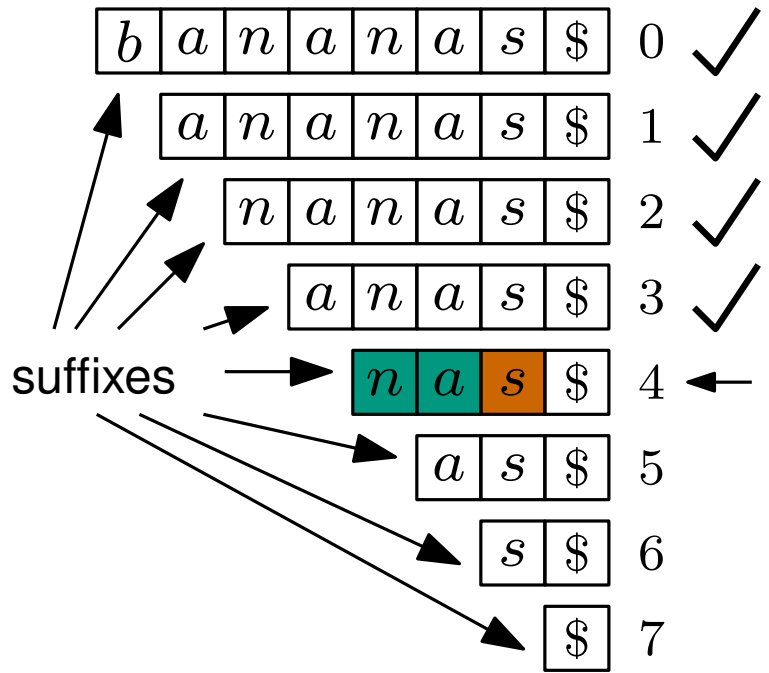
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*



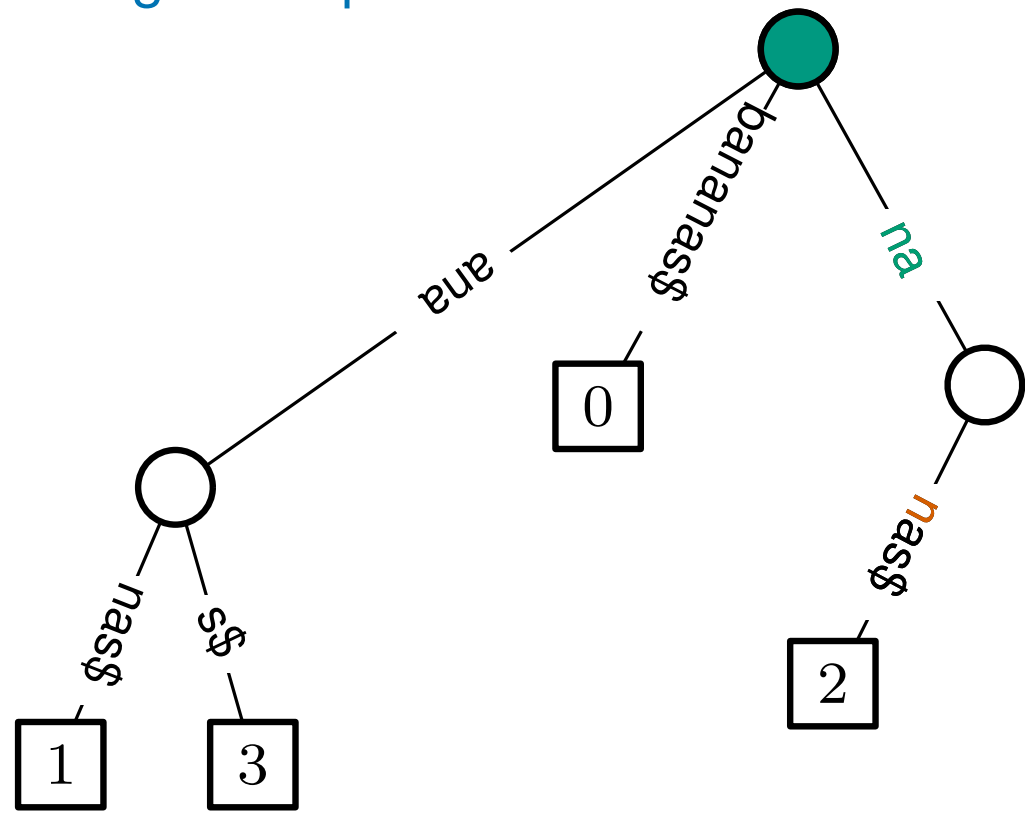
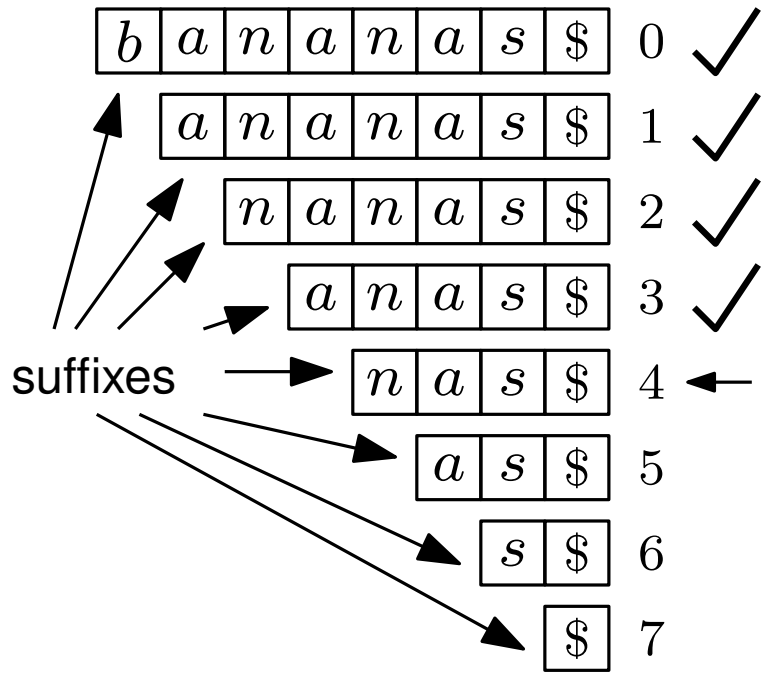
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

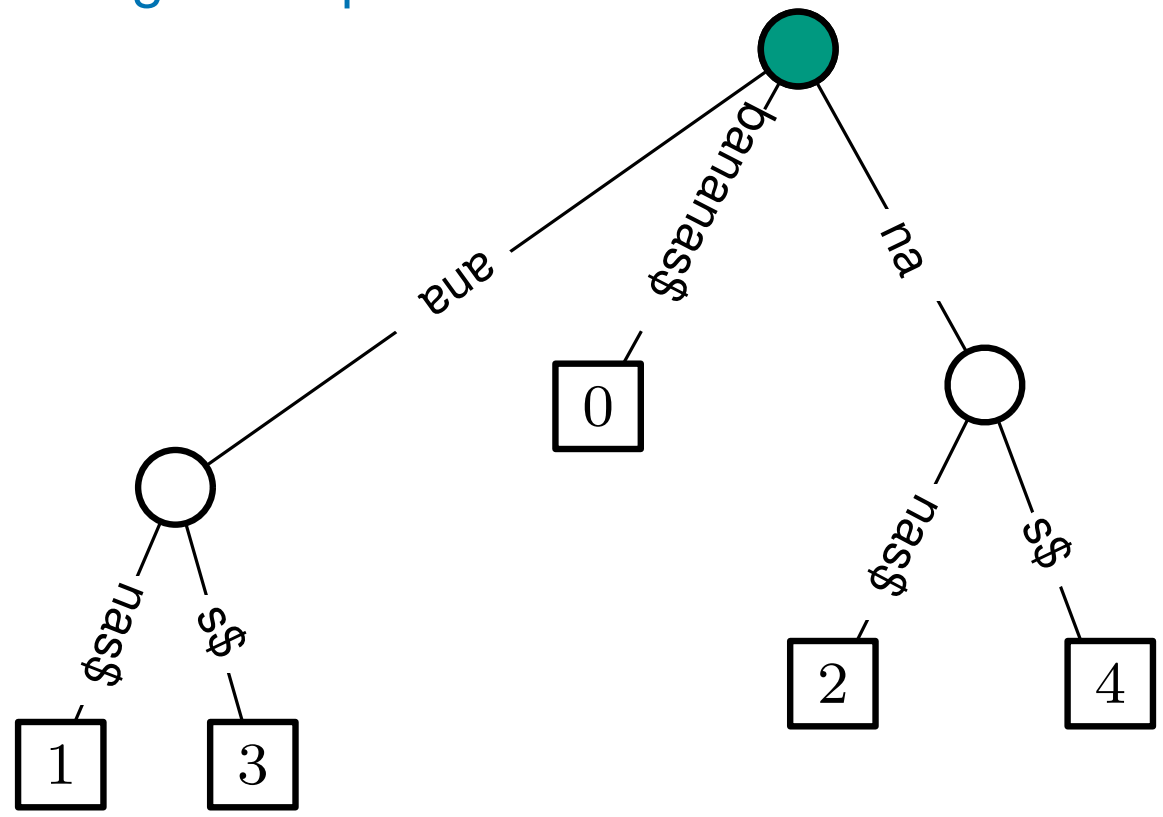
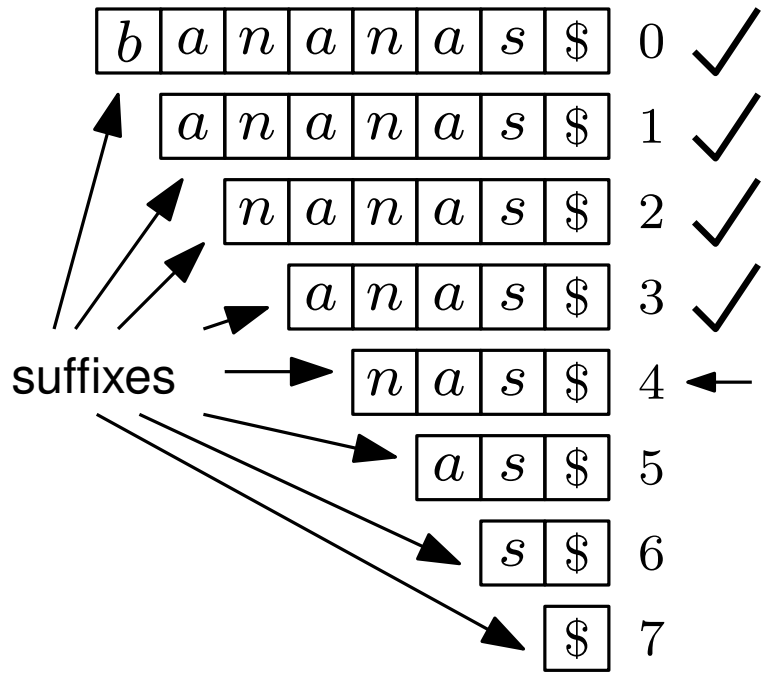
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

$b$	$a$	$n$	$a$	$n$	$a$	$s$	$\$$
-----	-----	-----	-----	-----	-----	-----	------

  
|-----  $n$  -----|

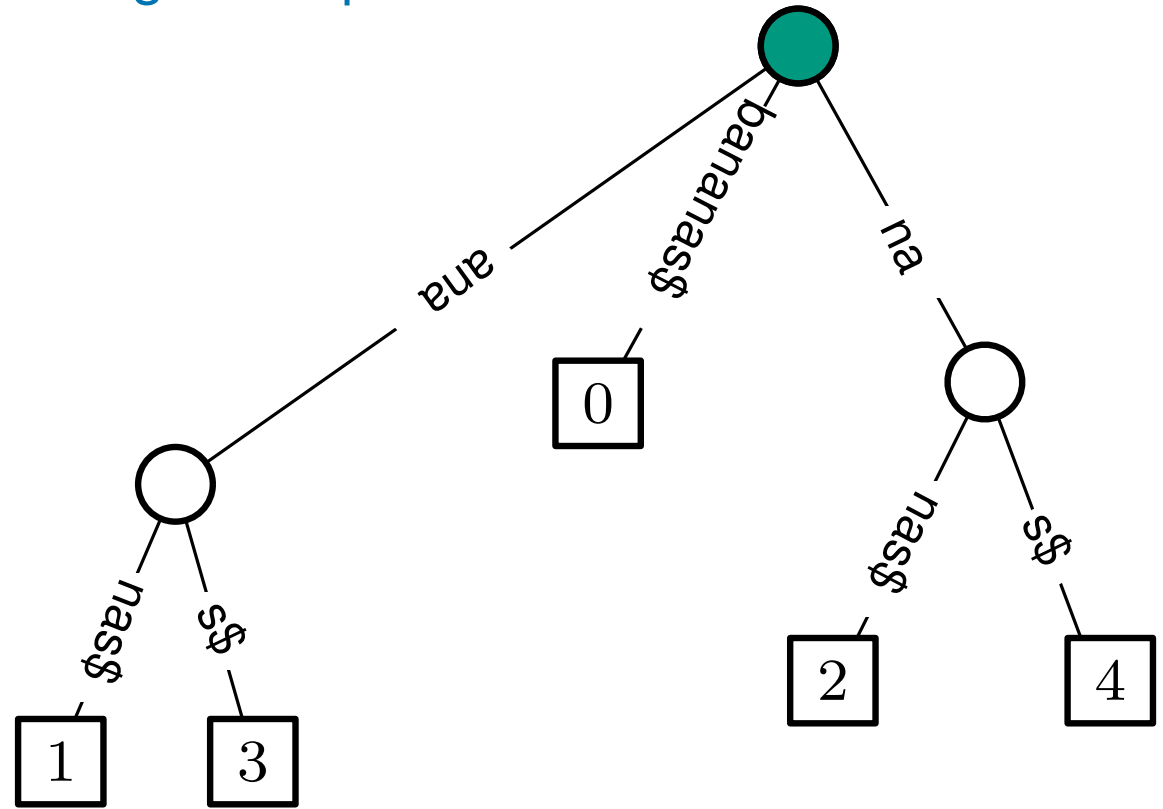
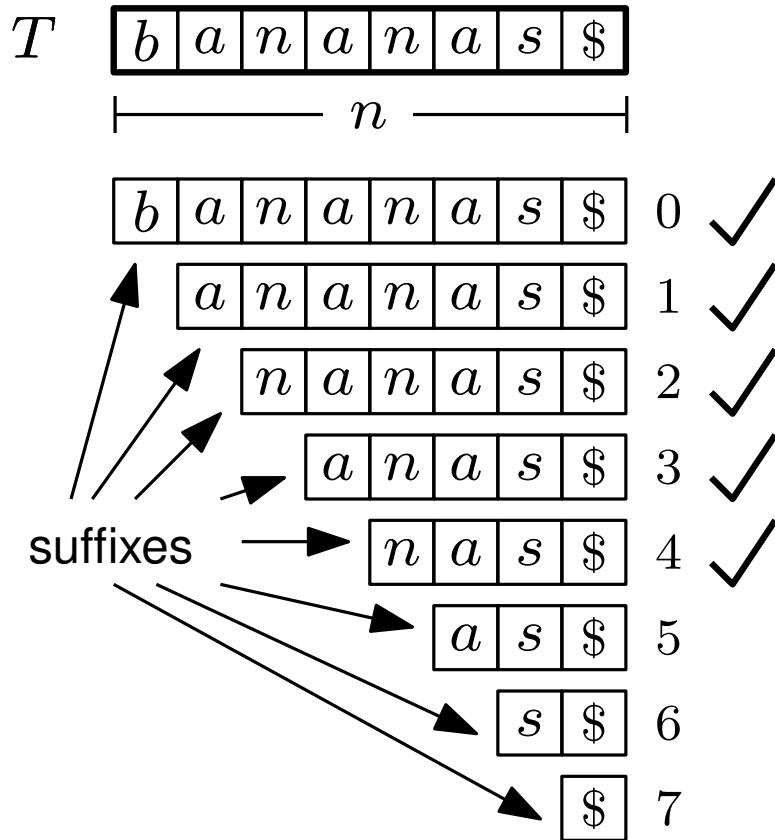


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

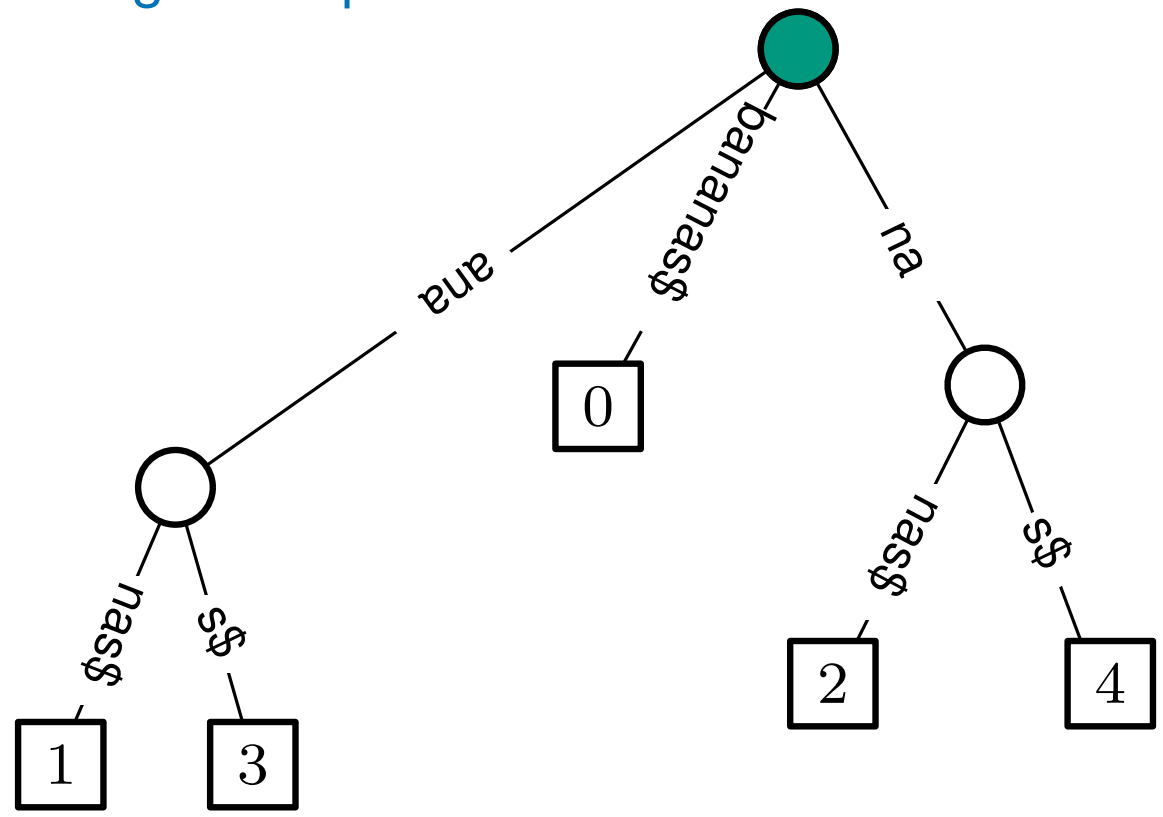
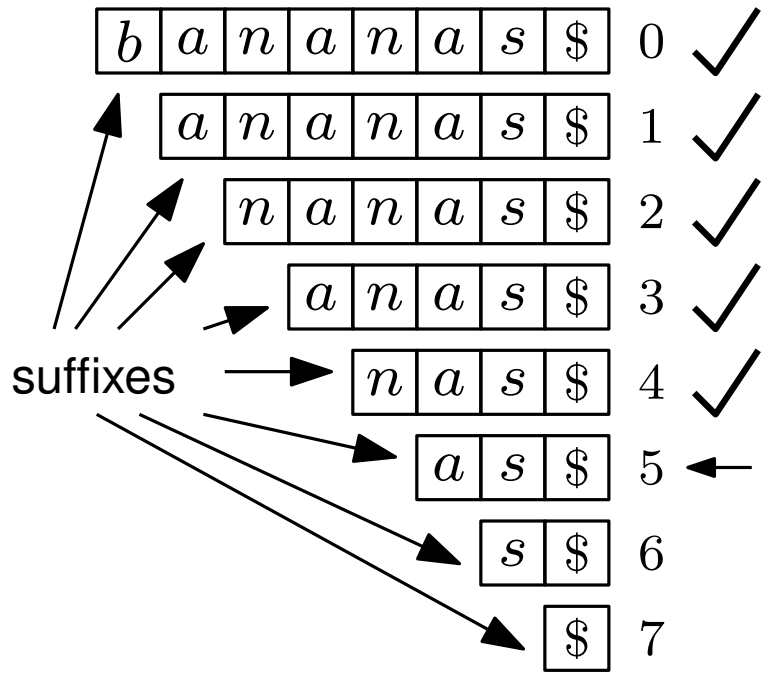
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

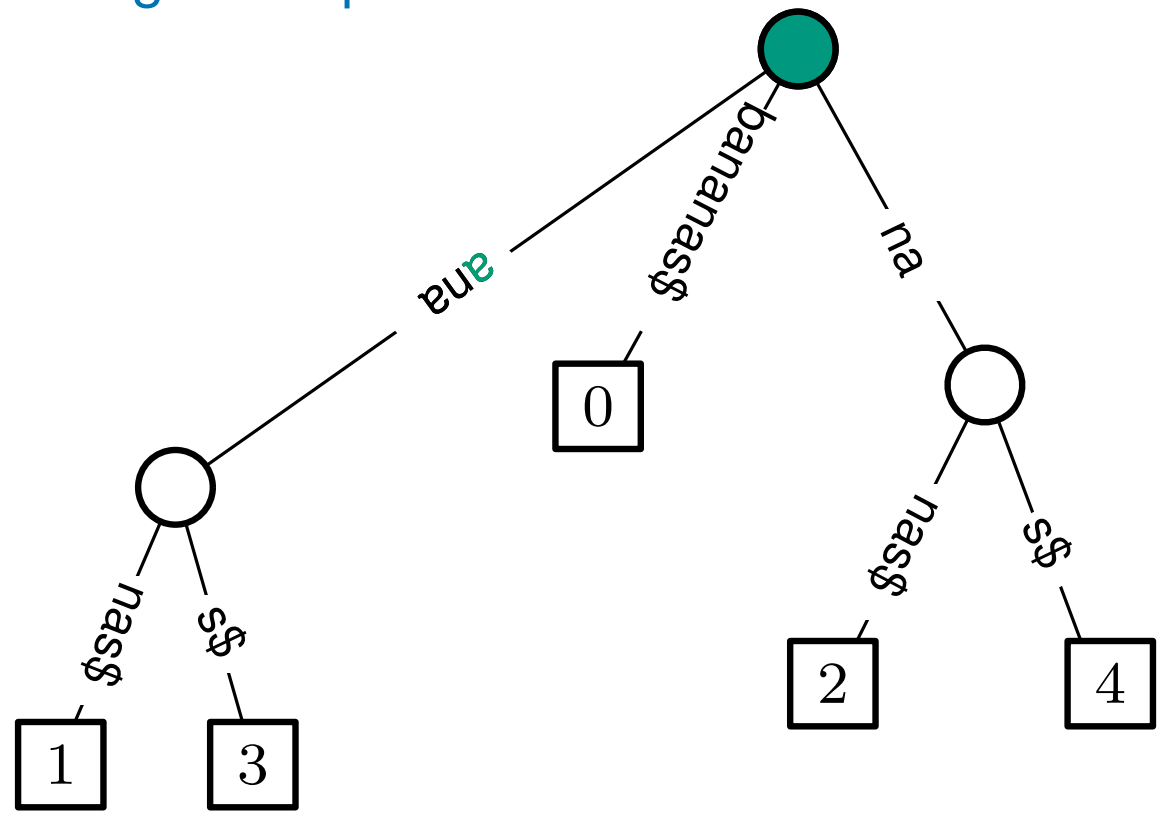
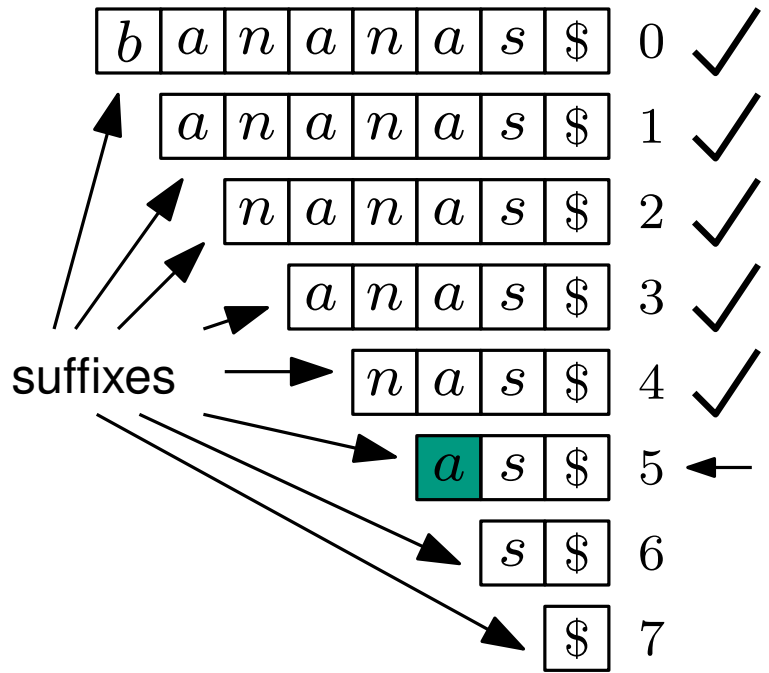
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|----- n -----|



Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
(as if you were matching a pattern)
- Add a new edge and leaf for the new suffix  
(this may require you to break an edge in two)

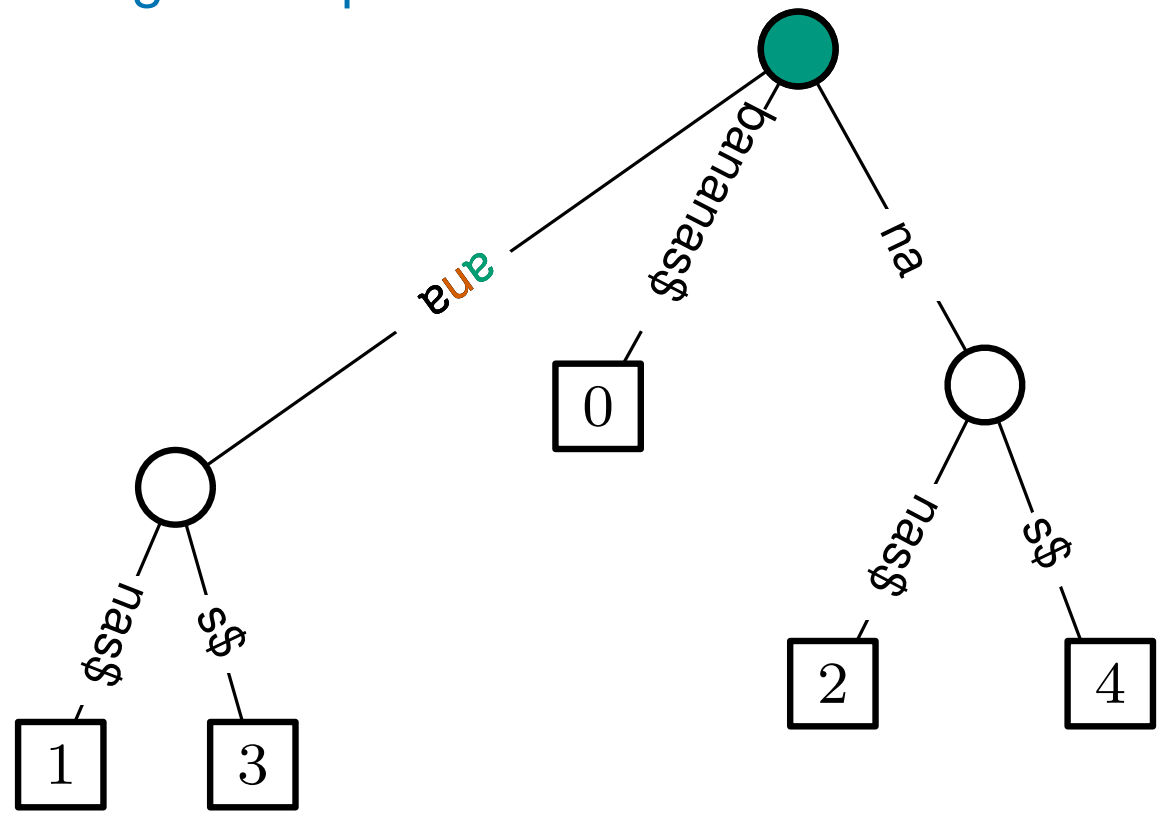
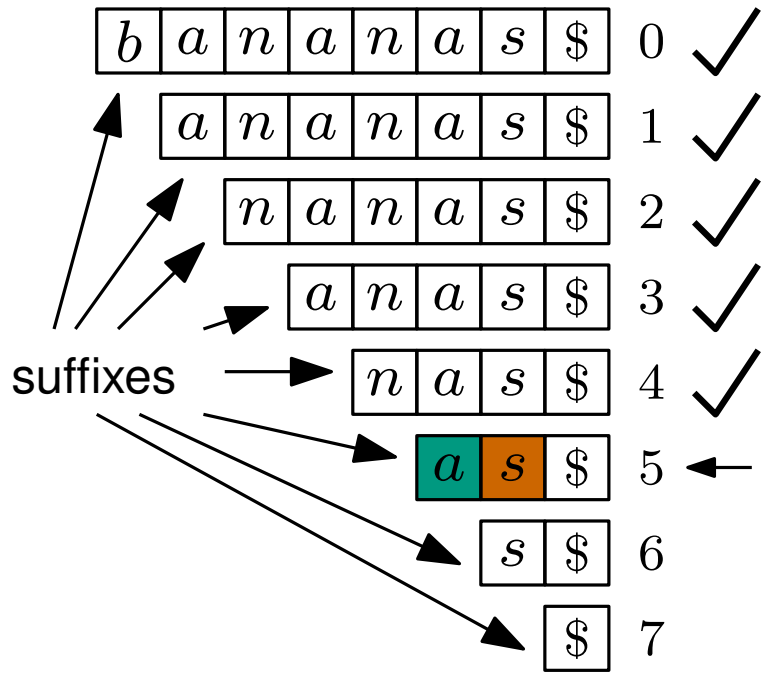
you should never actually do it like this

# Naively constructing a compacted suffix tree

$T$ 

b	a	n	a	n	a	s	\$
---	---	---	---	---	---	---	----

  
|-----  $n$  -----|



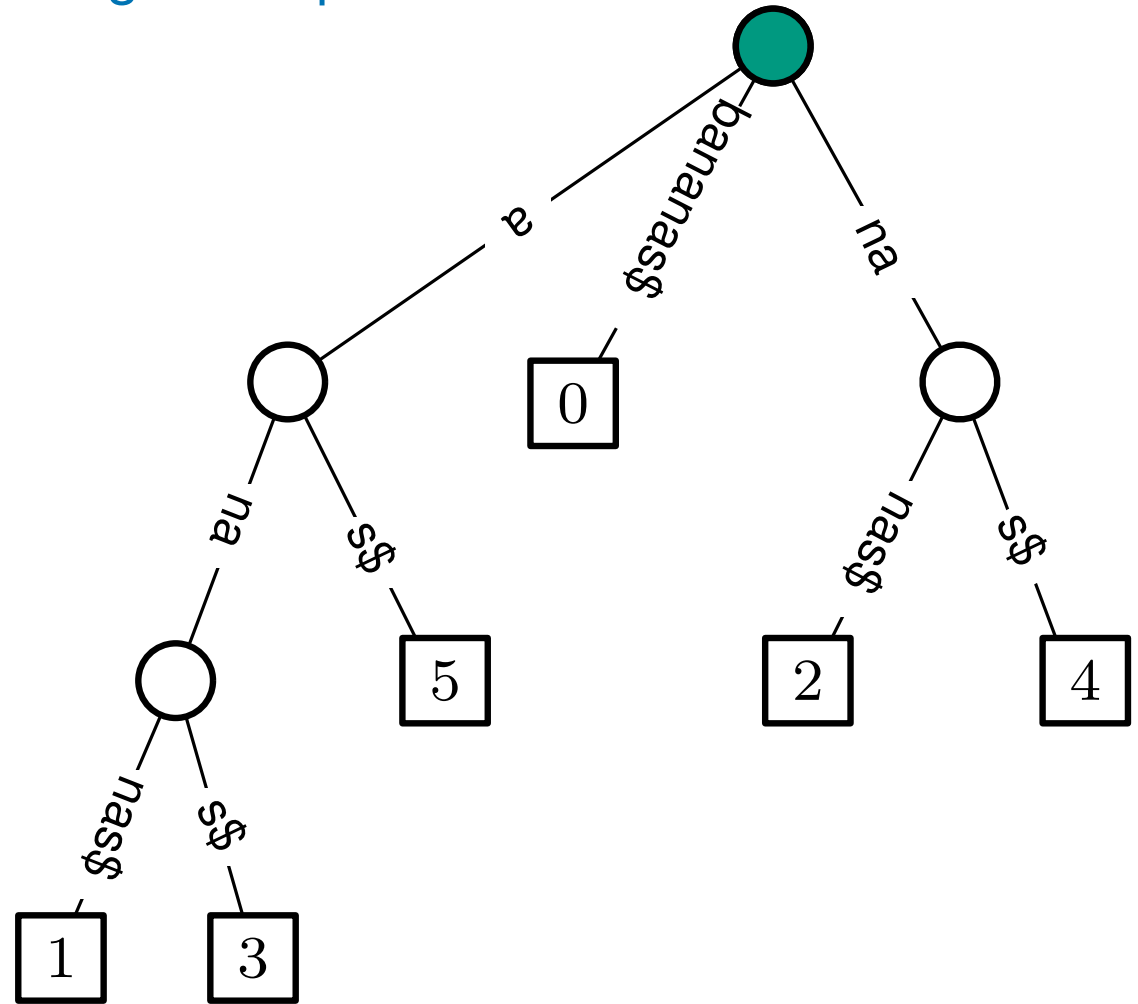
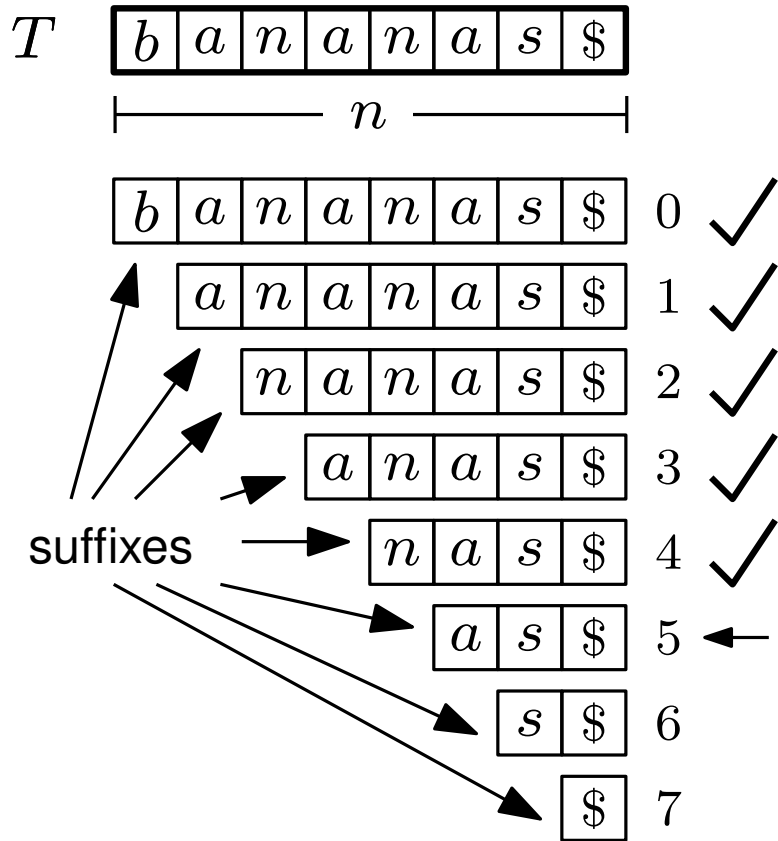
Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
(as if you were matching a pattern)
- Add a new edge and leaf for the new suffix  
(this may require you to break an edge in two)



you should never actually do it like this

# Naively constructing a compacted suffix tree



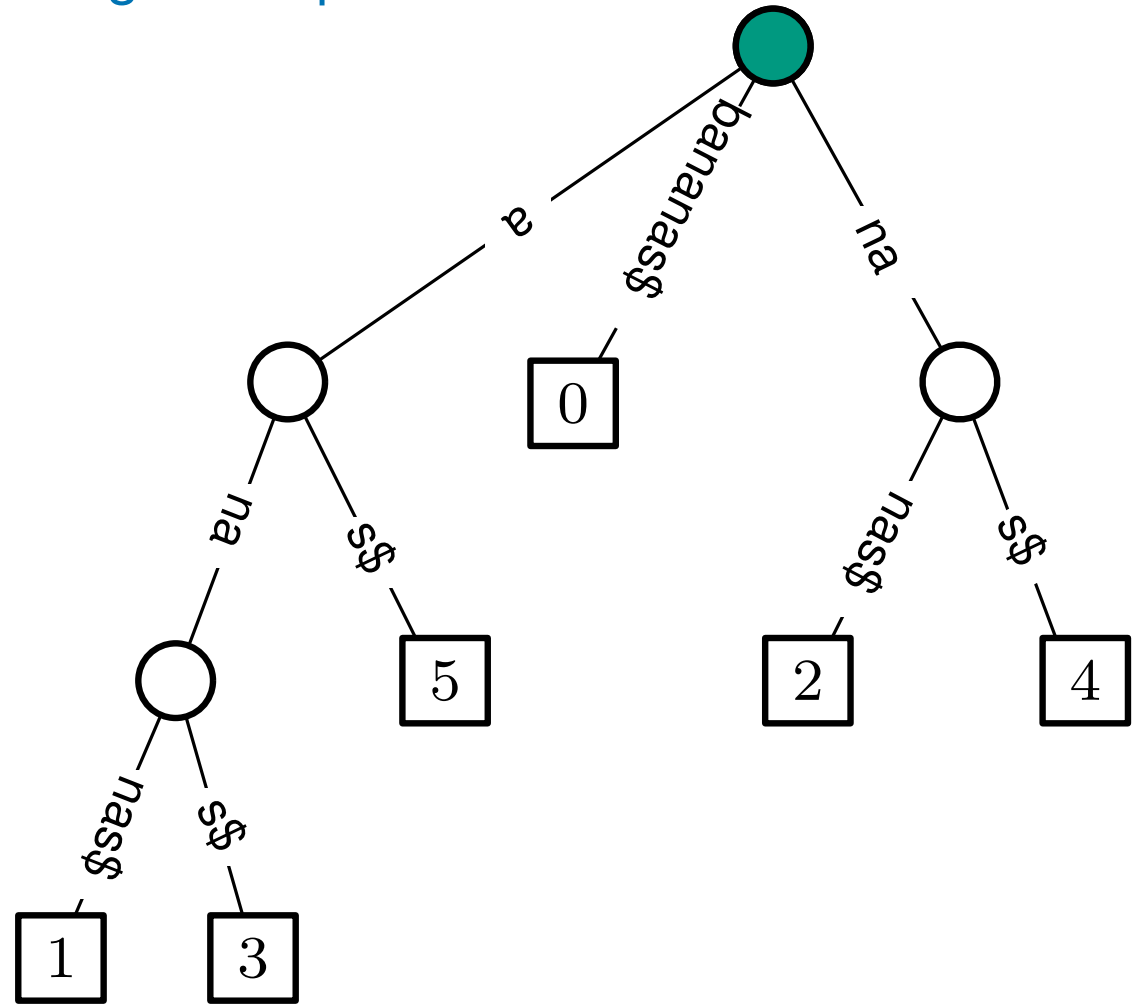
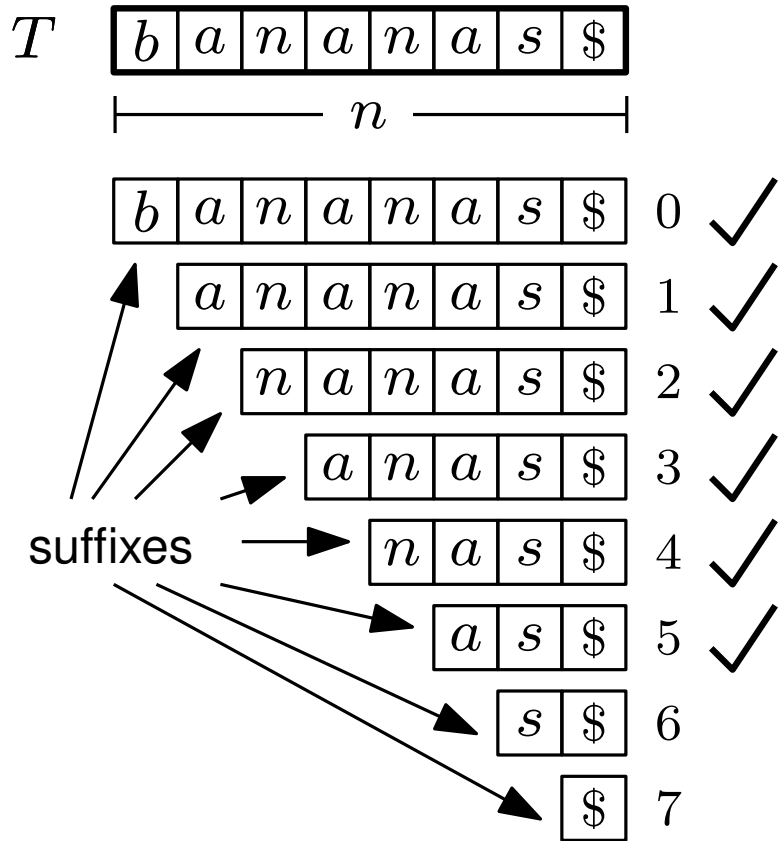
Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*



you should never actually do it like this

# Naively constructing a compacted suffix tree

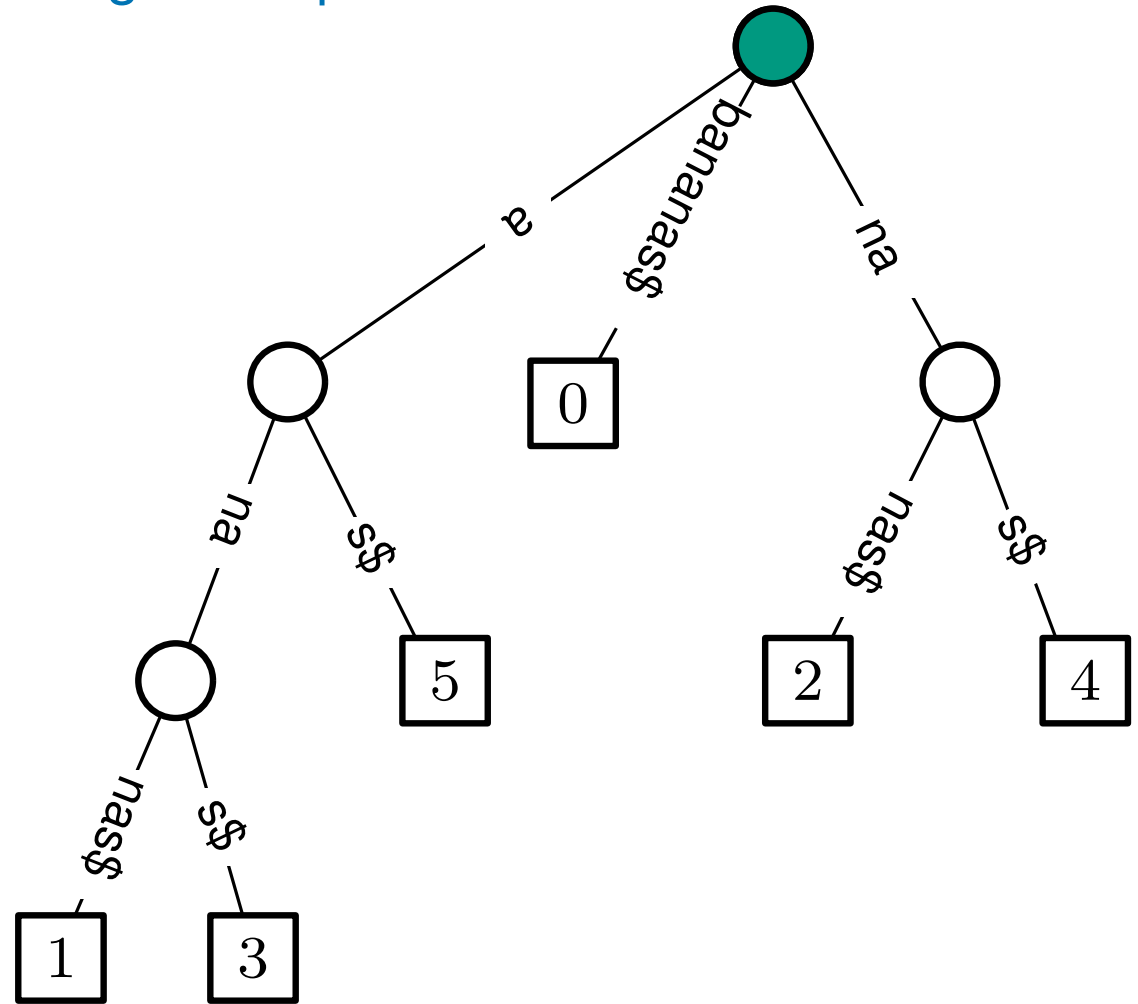
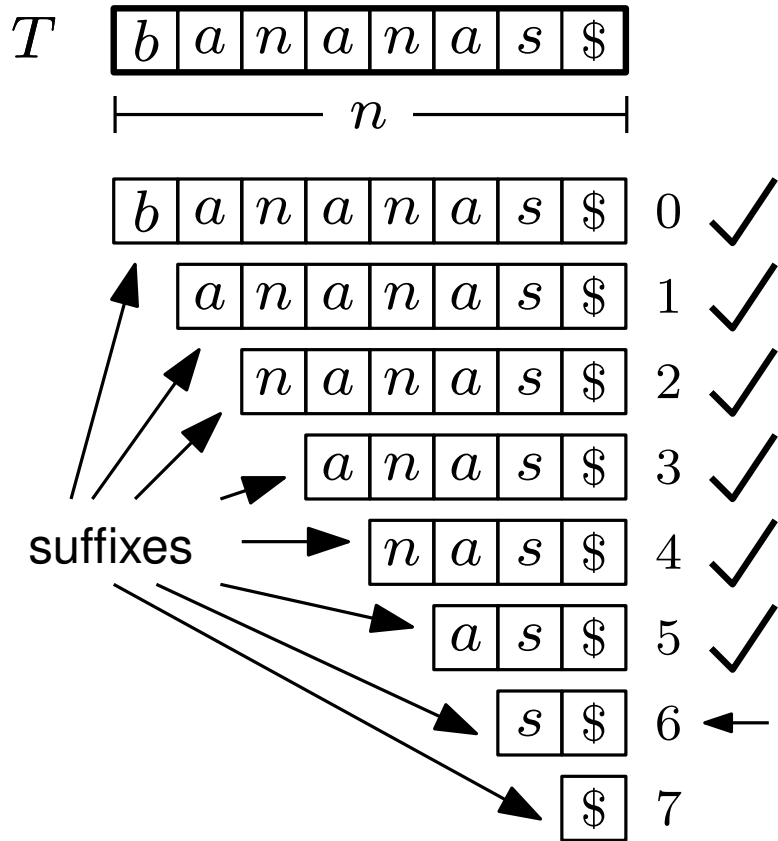


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should  
never actually  
do it like this

# Naively constructing a compacted suffix tree

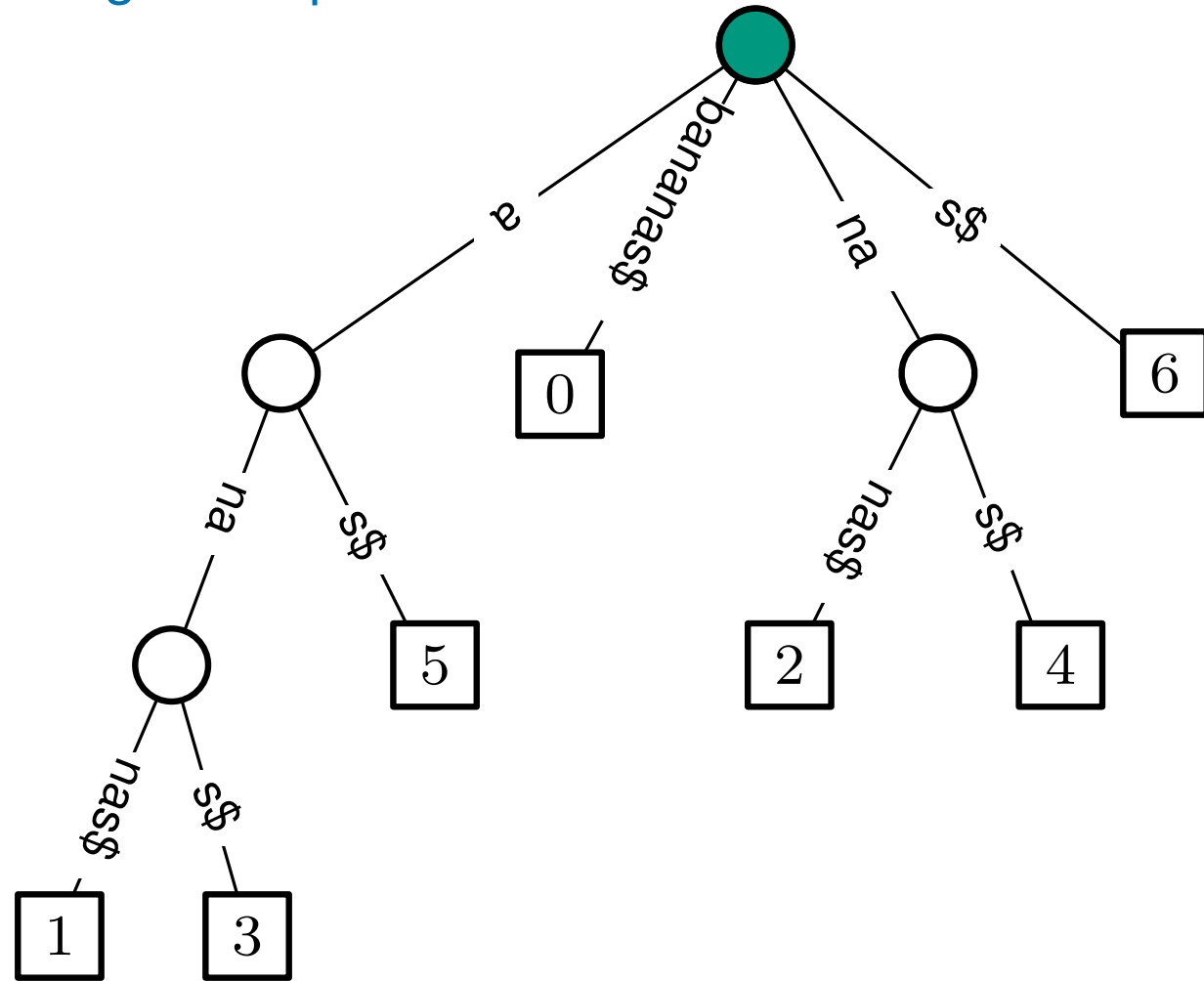
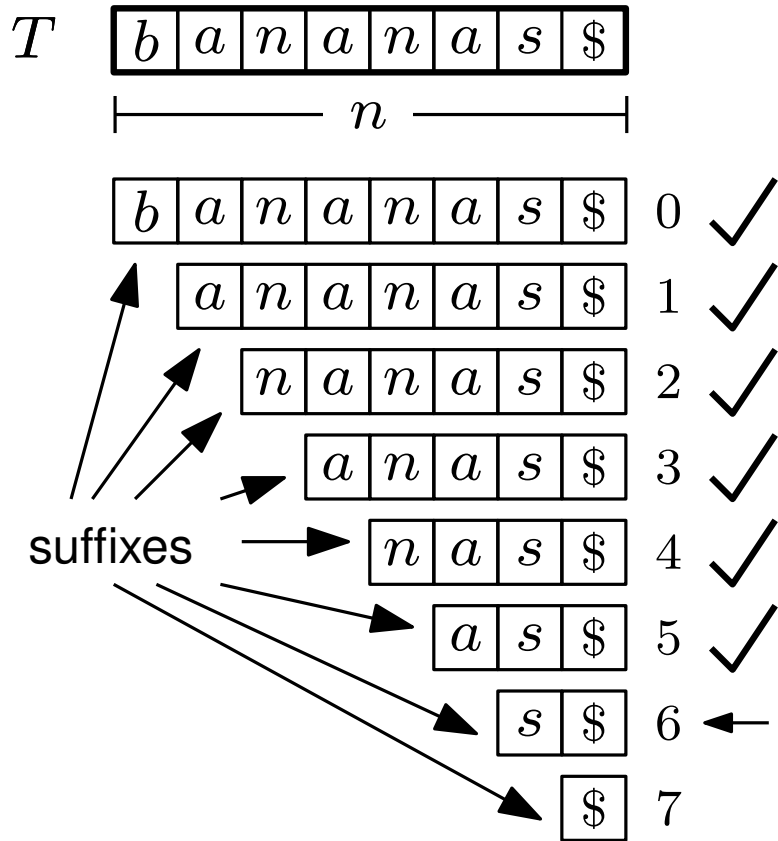


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

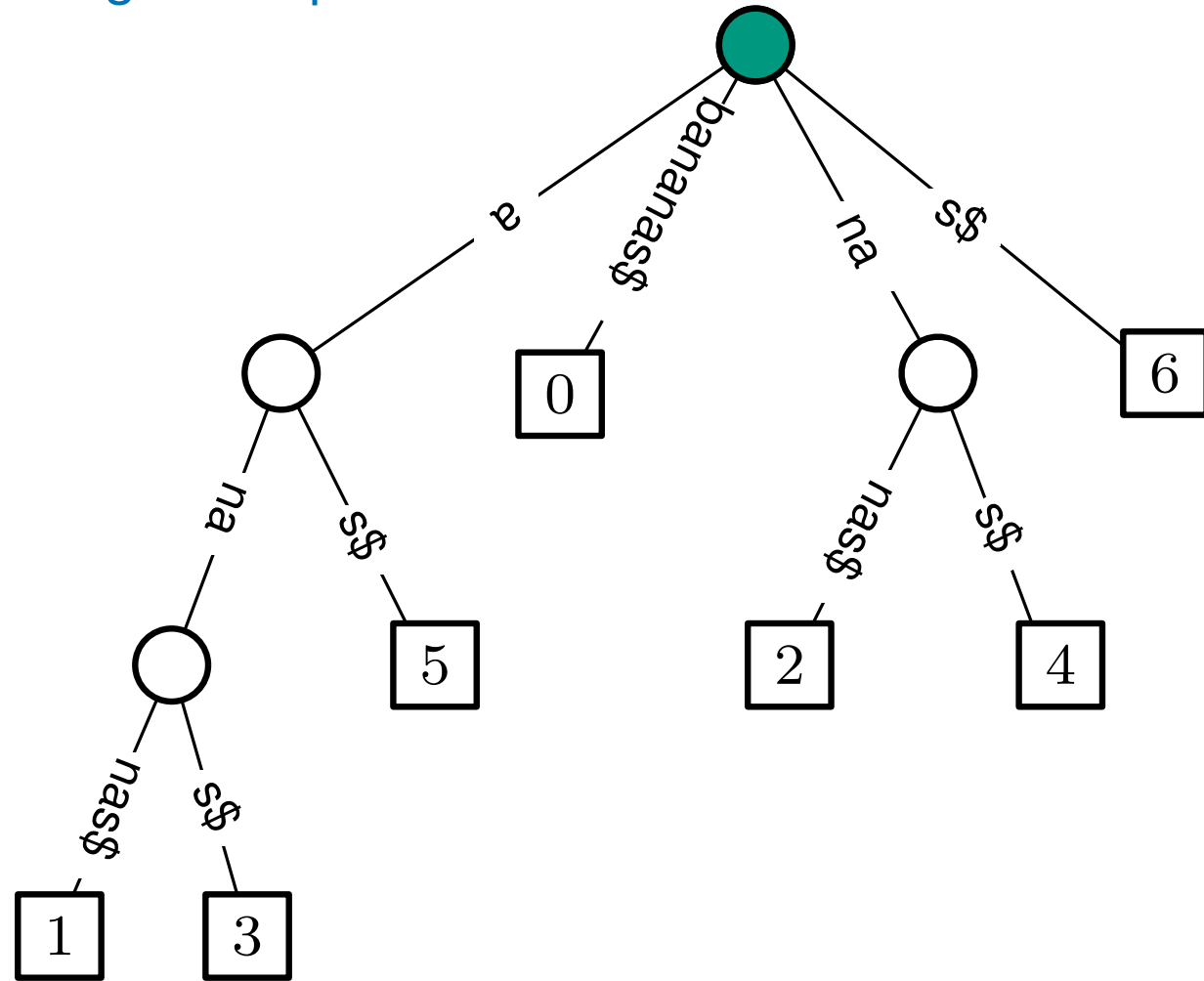
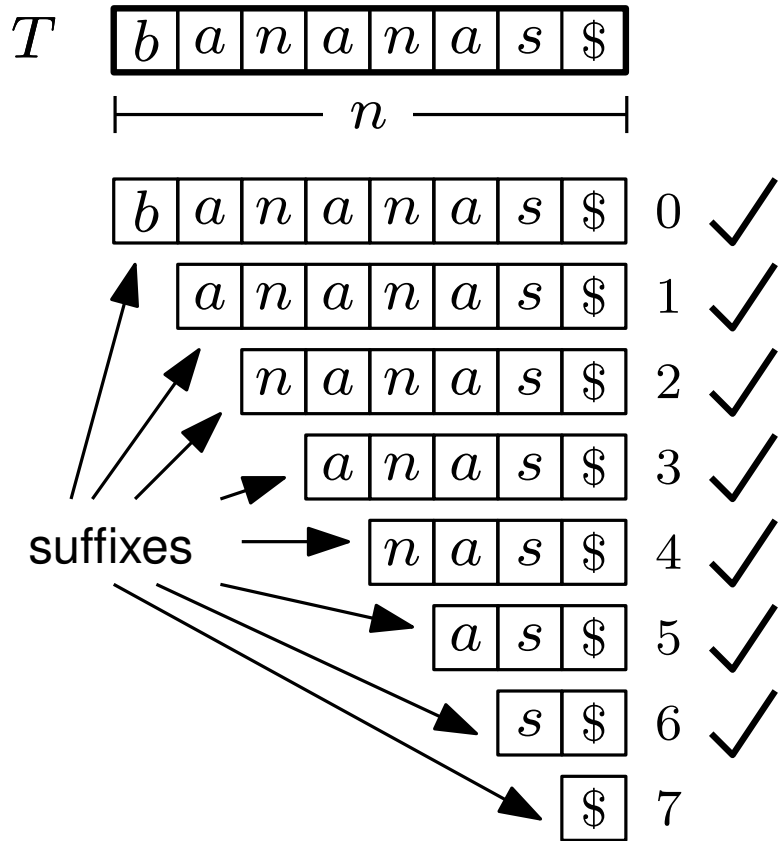


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

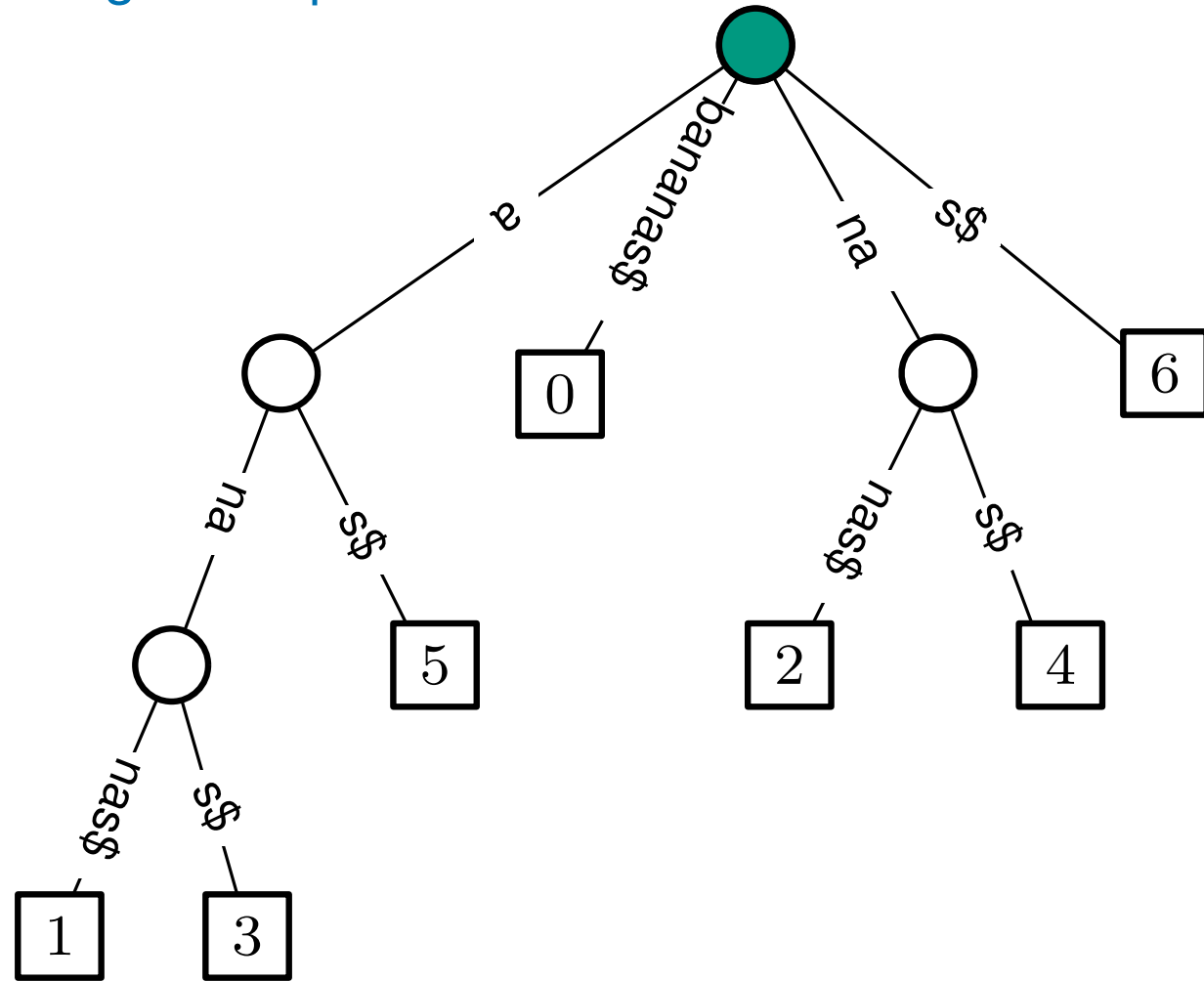
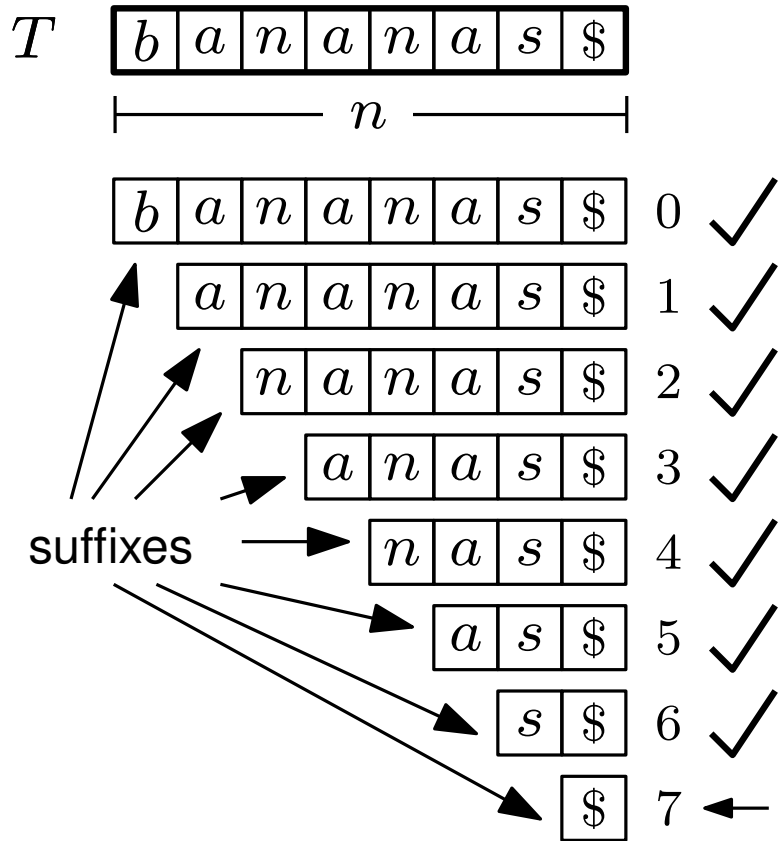


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

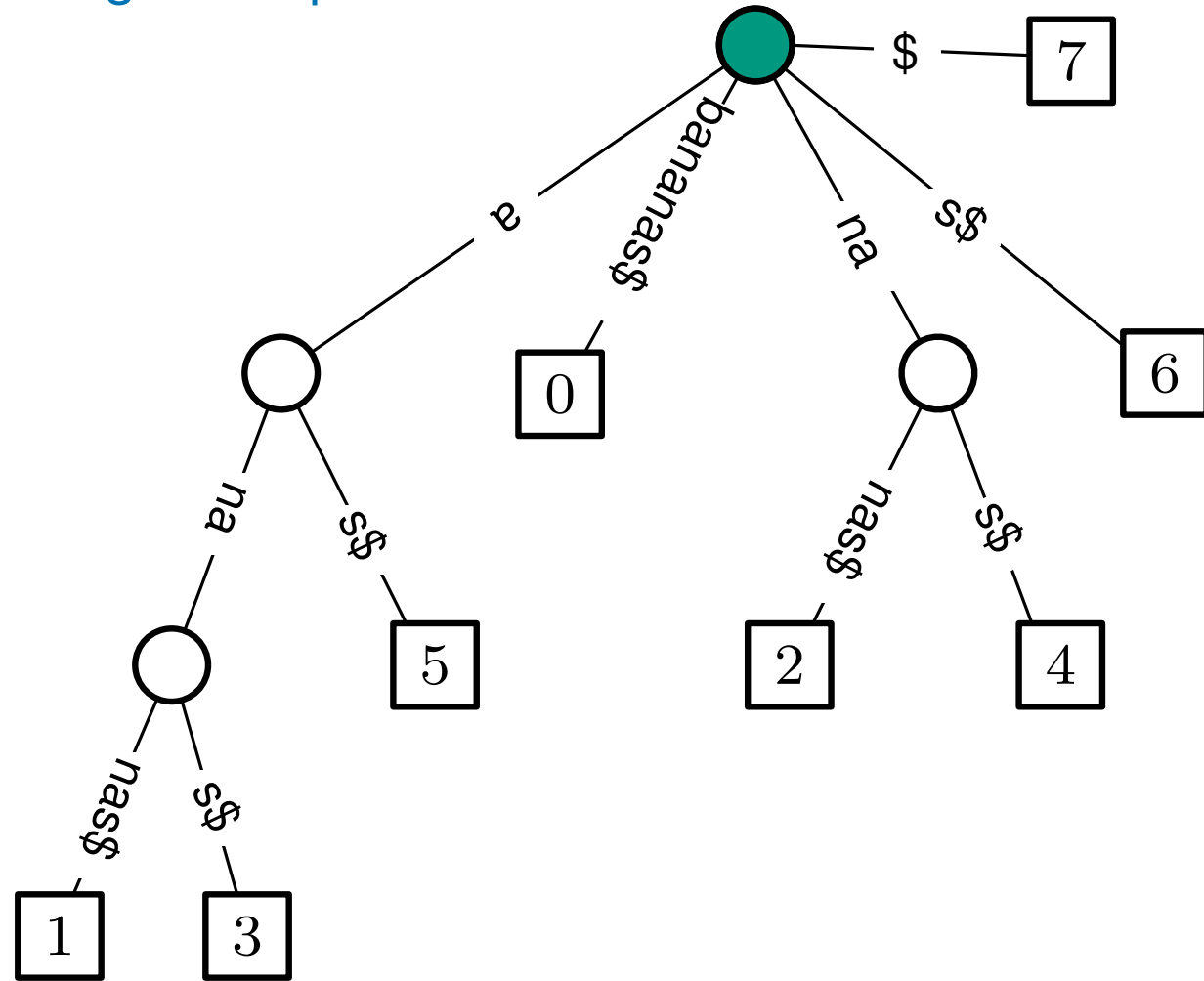
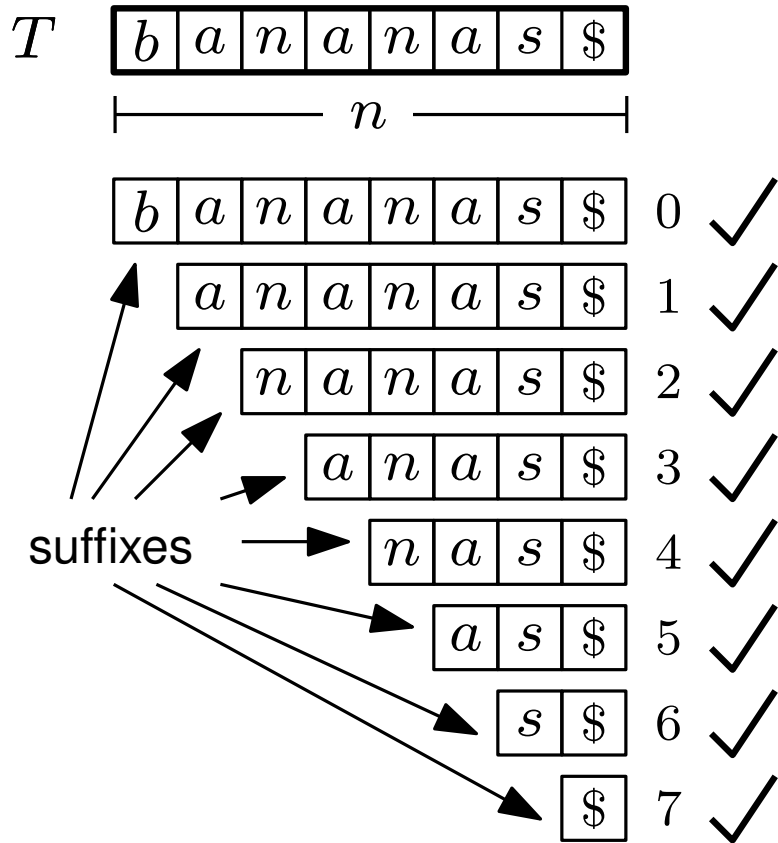


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

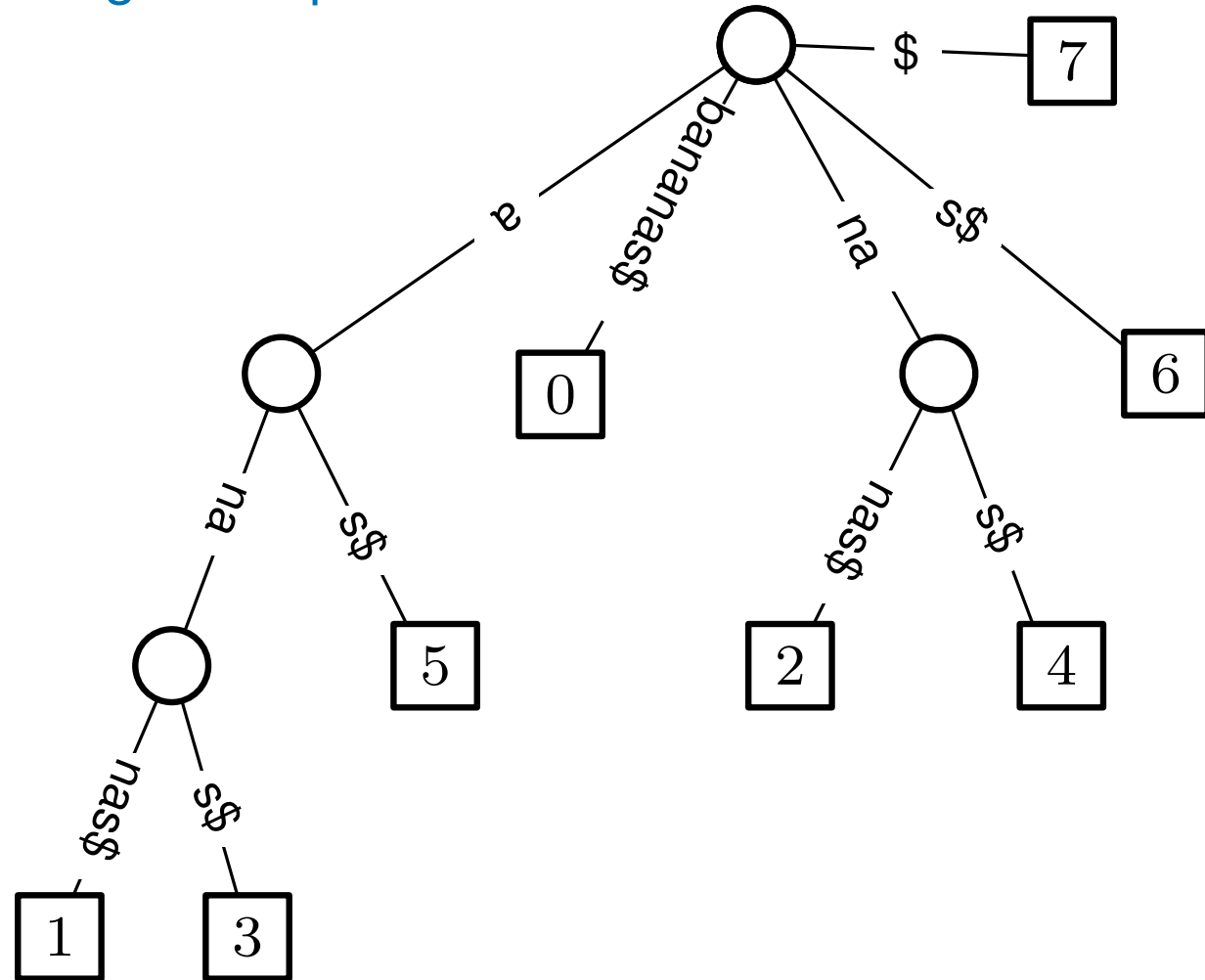
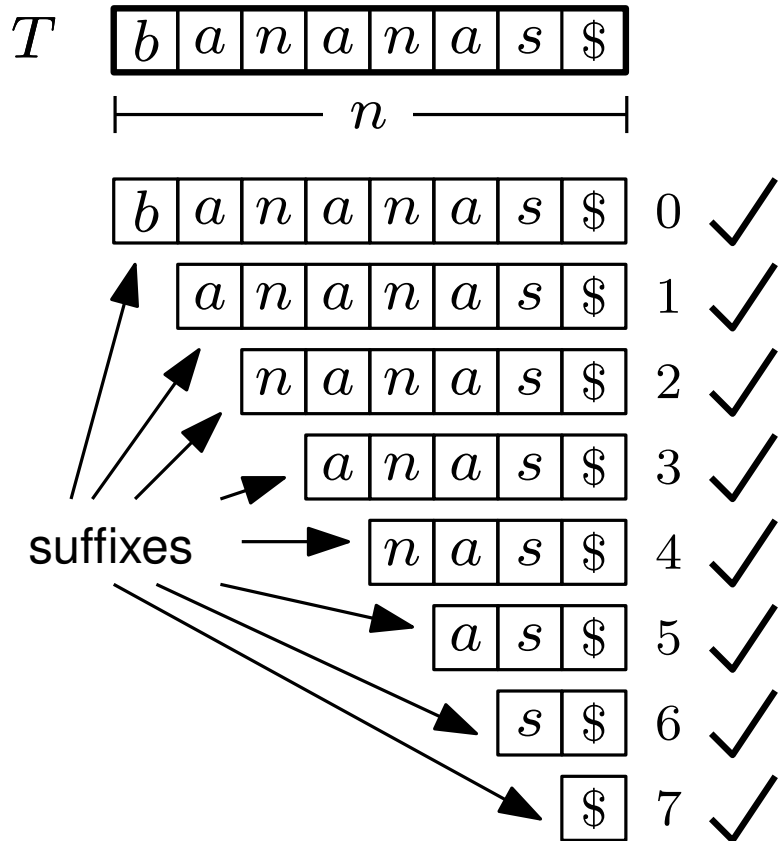


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree

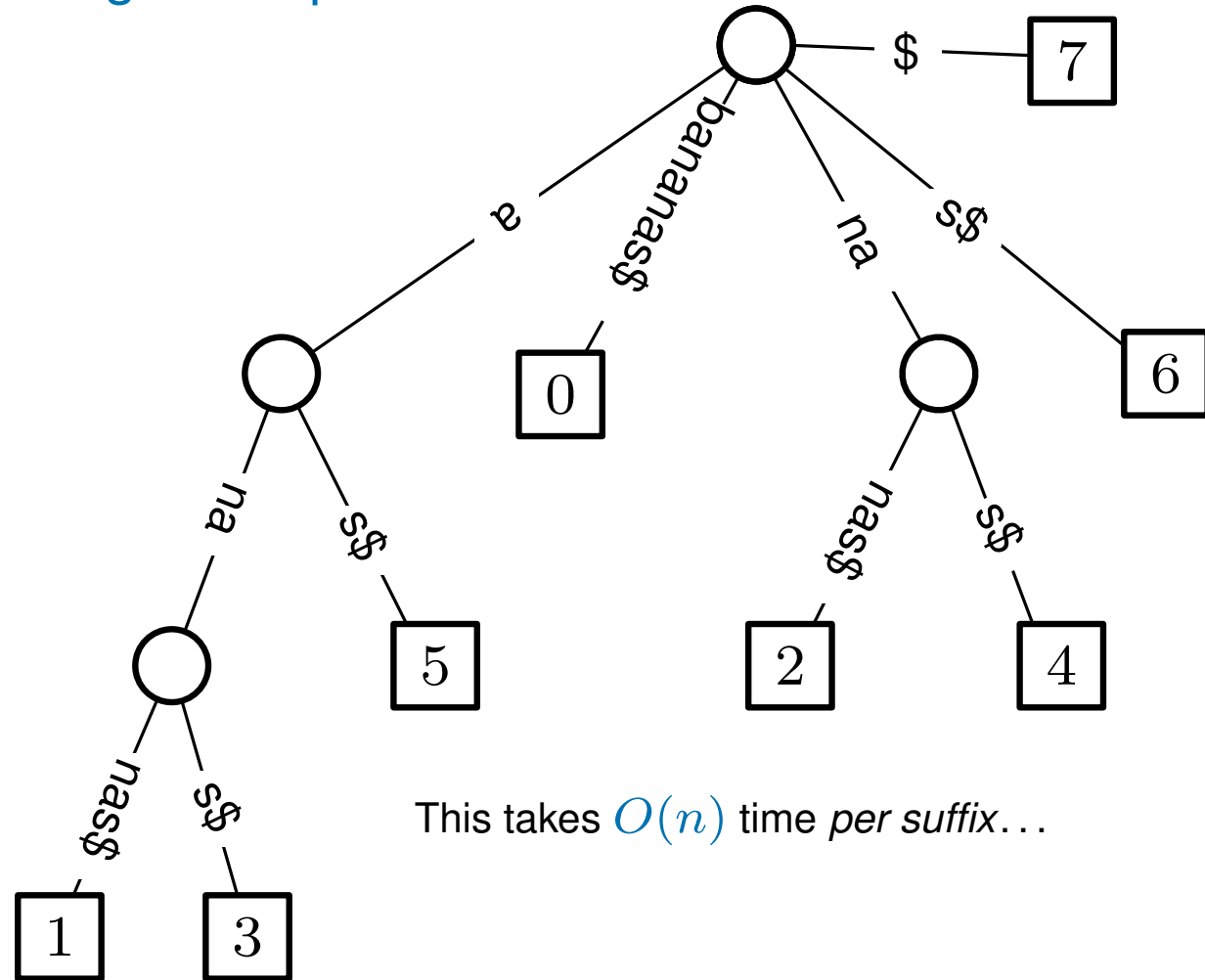
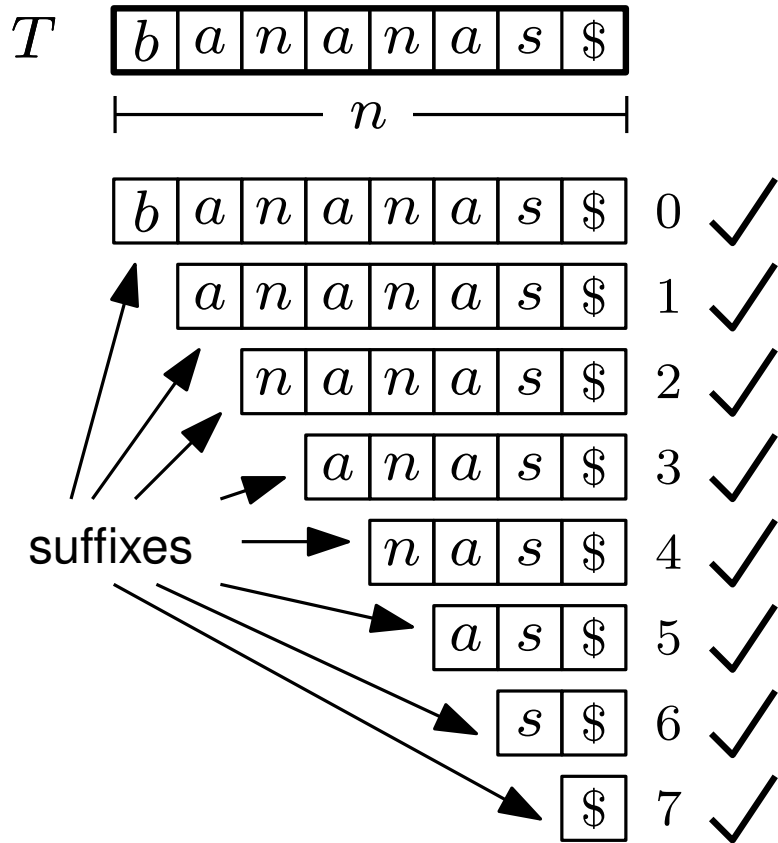


Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

you should never actually do it like this

# Naively constructing a compacted suffix tree



This takes  $O(n)$  time per suffix...

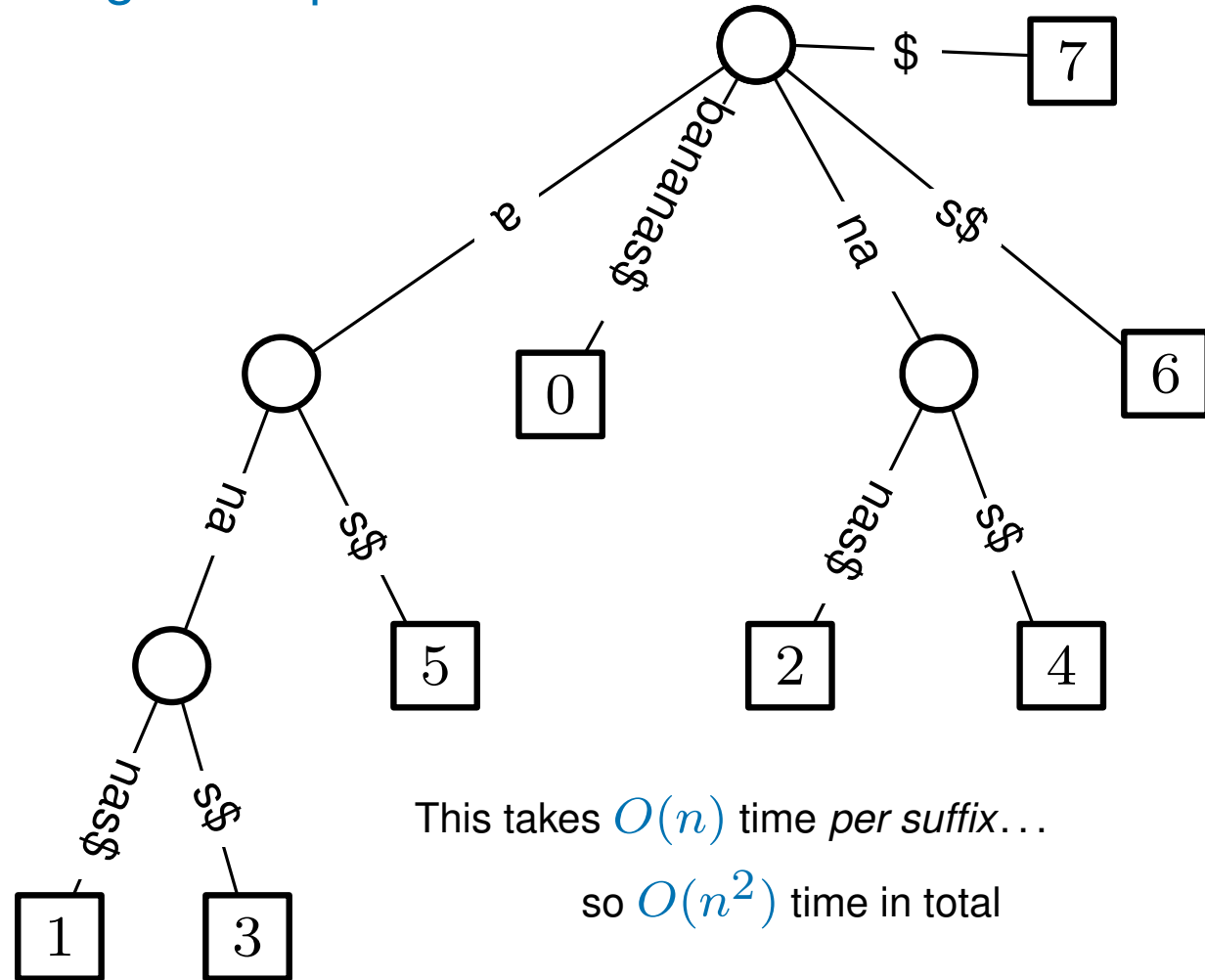
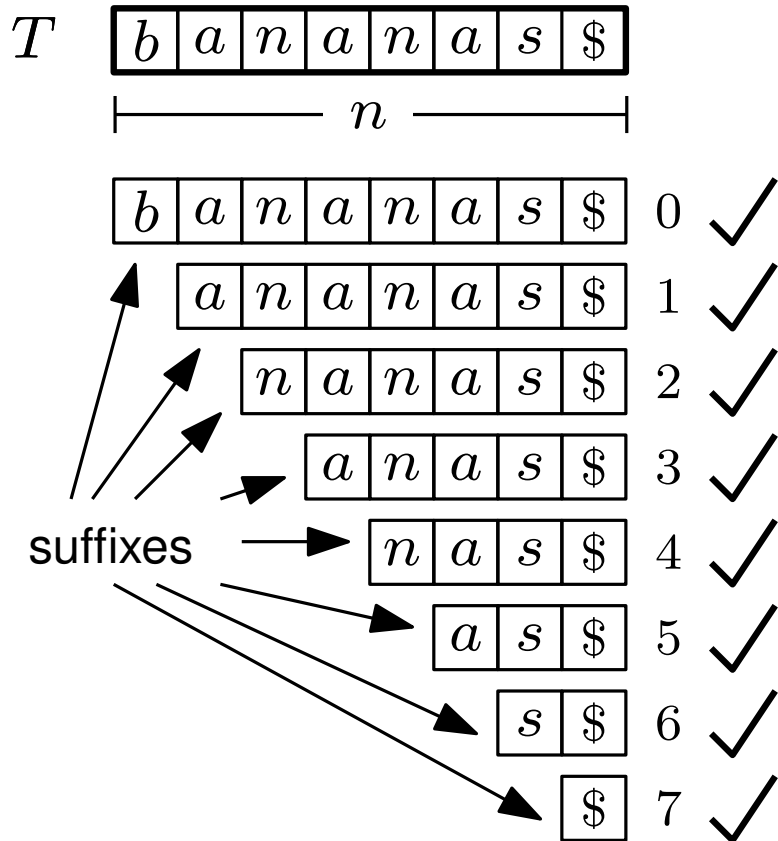
Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*



you should never actually do it like this

# Naively constructing a compacted suffix tree

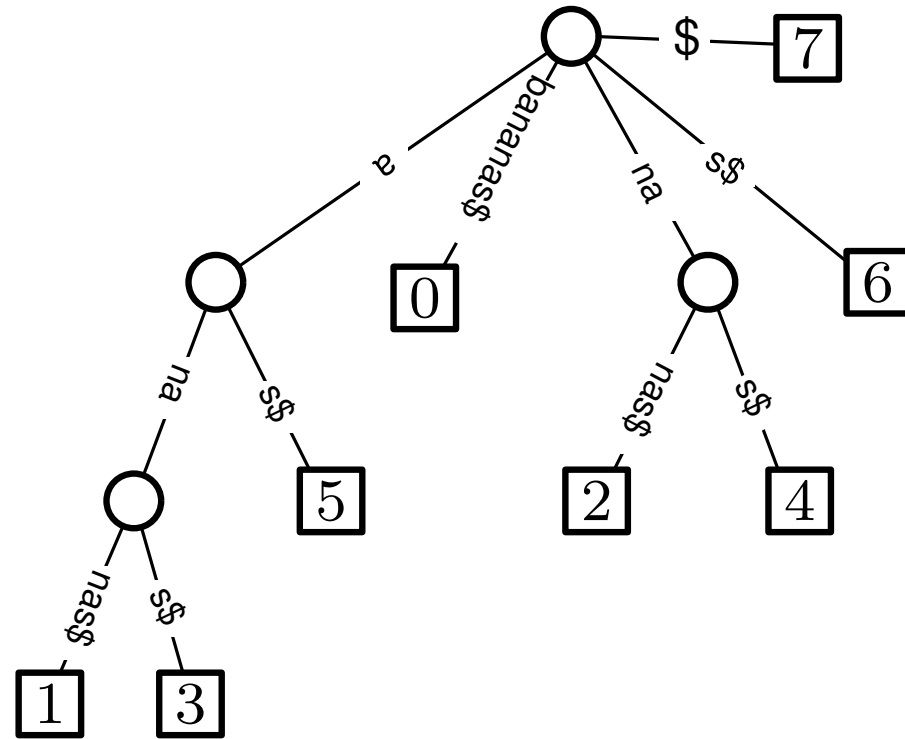
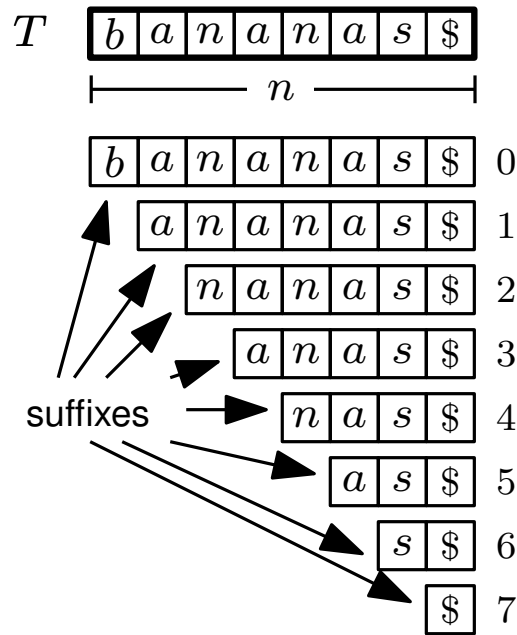


This takes  $O(n)$  time per suffix...  
so  $O(n^2)$  time in total

Insert the suffixes one at a time (longest first)

- Search for the new suffix in the partial suffix tree  
*(as if you were matching a pattern)*
- Add a new edge and leaf for the new suffix  
*(this may require you to break an edge in two)*

# Suffix tree summary



- The (compacted) suffix tree of a (length  $n$ ) text uses  $O(n)$  space
- Finding all matches of a pattern  $P$  of length  $m$  takes  $O(m + occ)$
- Suffix trees can be built in  $O(n)$  time

*where  $occ$  is the number of matches*

you should actually do it like this (or build a *suffix array* instead)

*but we have only seen the  $O(n^2)$  time method*

*we assumed that the alphabet contained a constant number of symbols*

# Multiple text indexing

$T_1$ 

b	a	n	a	n	a	s
---	---	---	---	---	---	---

  
|----- $n_1$ -----|

$T_2$ 

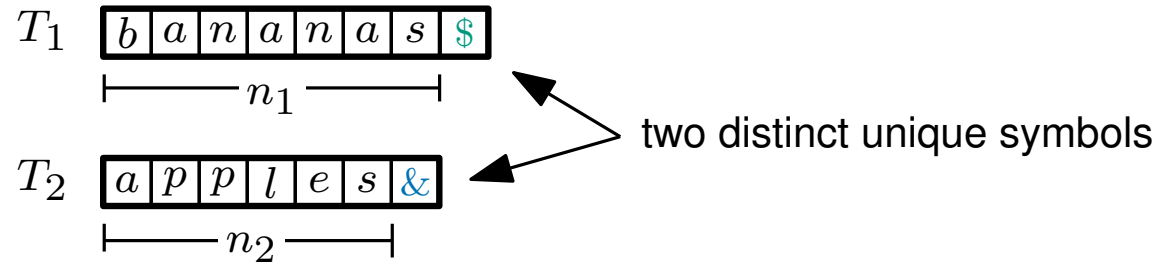
a	p	p	l	e	s
---	---	---	---	---	---

  
|----- $n_2$ -----|

---

How can we index multiple texts?

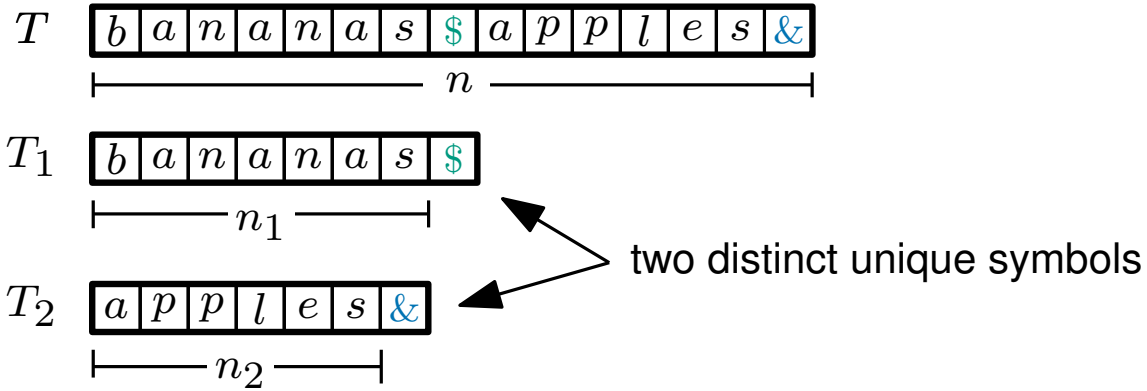
# Multiple text indexing




---

How can we index multiple texts?

# Multiple text indexing




---

How can we index multiple texts?

# Multiple text indexing

$T$ 

<i>b</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>s</i>	\$	<i>a</i>	<i>p</i>	<i>p</i>	<i>l</i>	<i>e</i>	<i>s</i>	&
----------	----------	----------	----------	----------	----------	----------	----	----------	----------	----------	----------	----------	----------	---

  
|-----  $n$  -----|

$T_1$ 

<i>b</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>s</i>	\$
----------	----------	----------	----------	----------	----------	----------	----

  
|-----  $n_1$  -----|

$T_2$ 

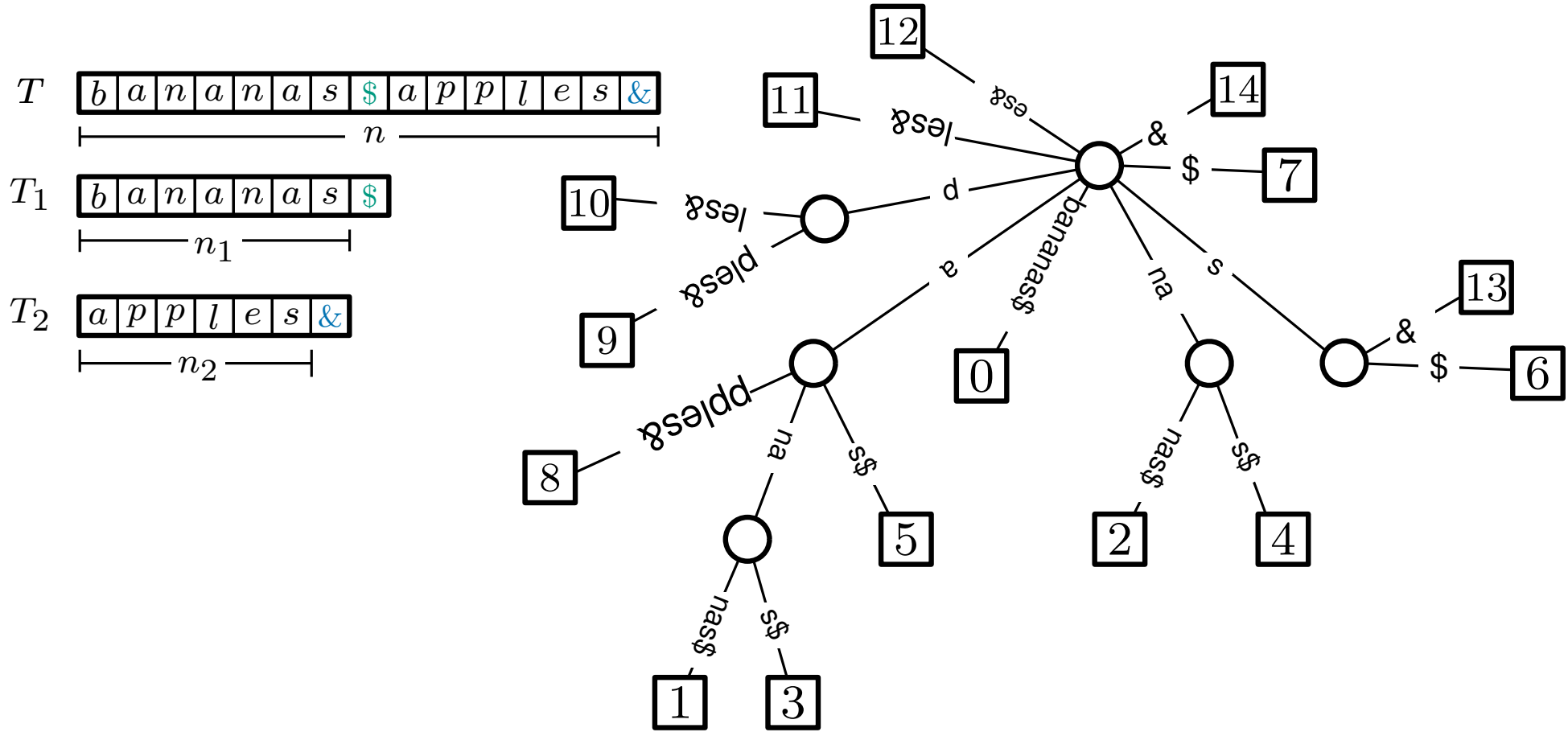
<i>a</i>	<i>p</i>	<i>p</i>	<i>l</i>	<i>e</i>	<i>s</i>	&
----------	----------	----------	----------	----------	----------	---

  
|-----  $n_2$  -----|

---

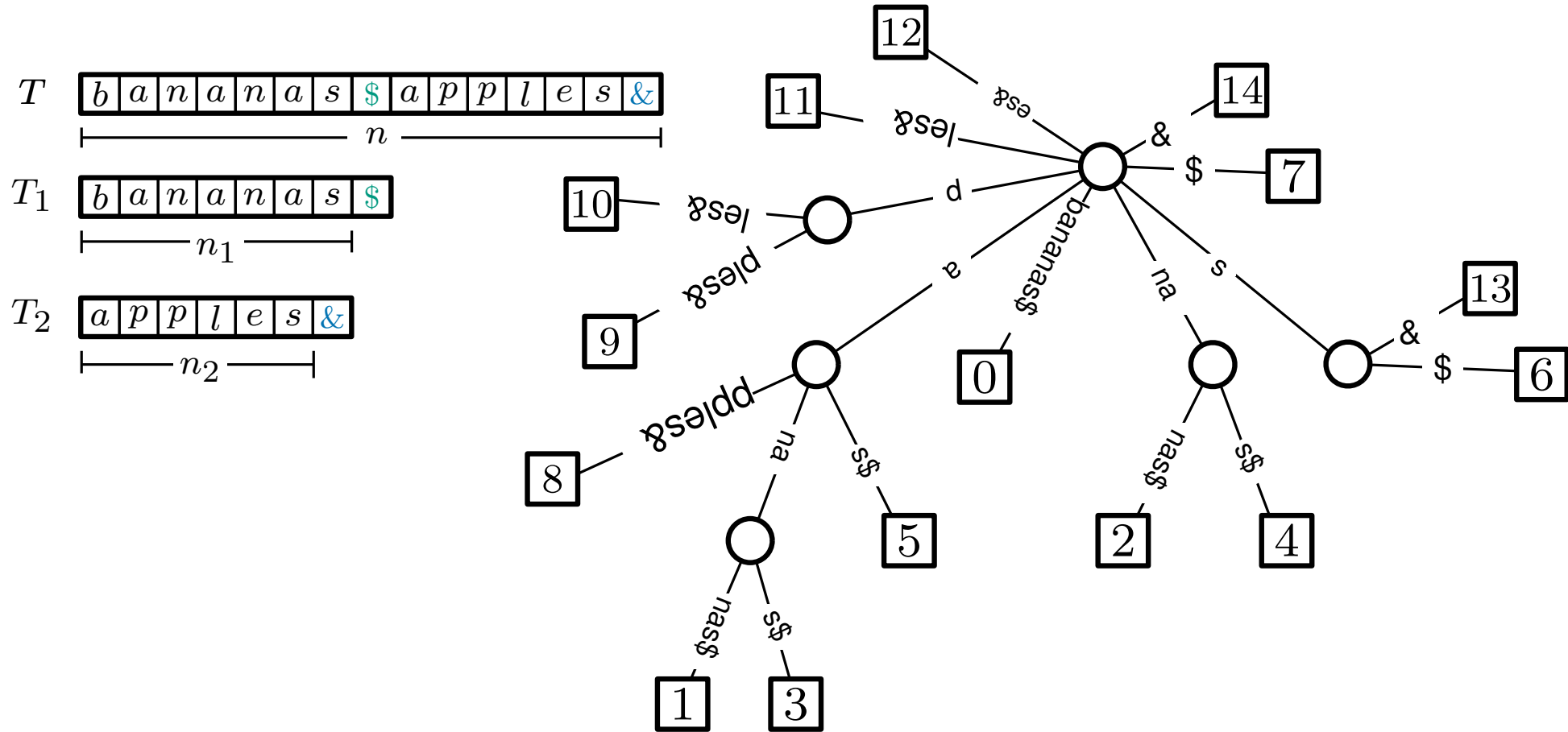
How can we index multiple texts?

# Multiple text indexing



How can we index multiple texts?

# Multiple text indexing



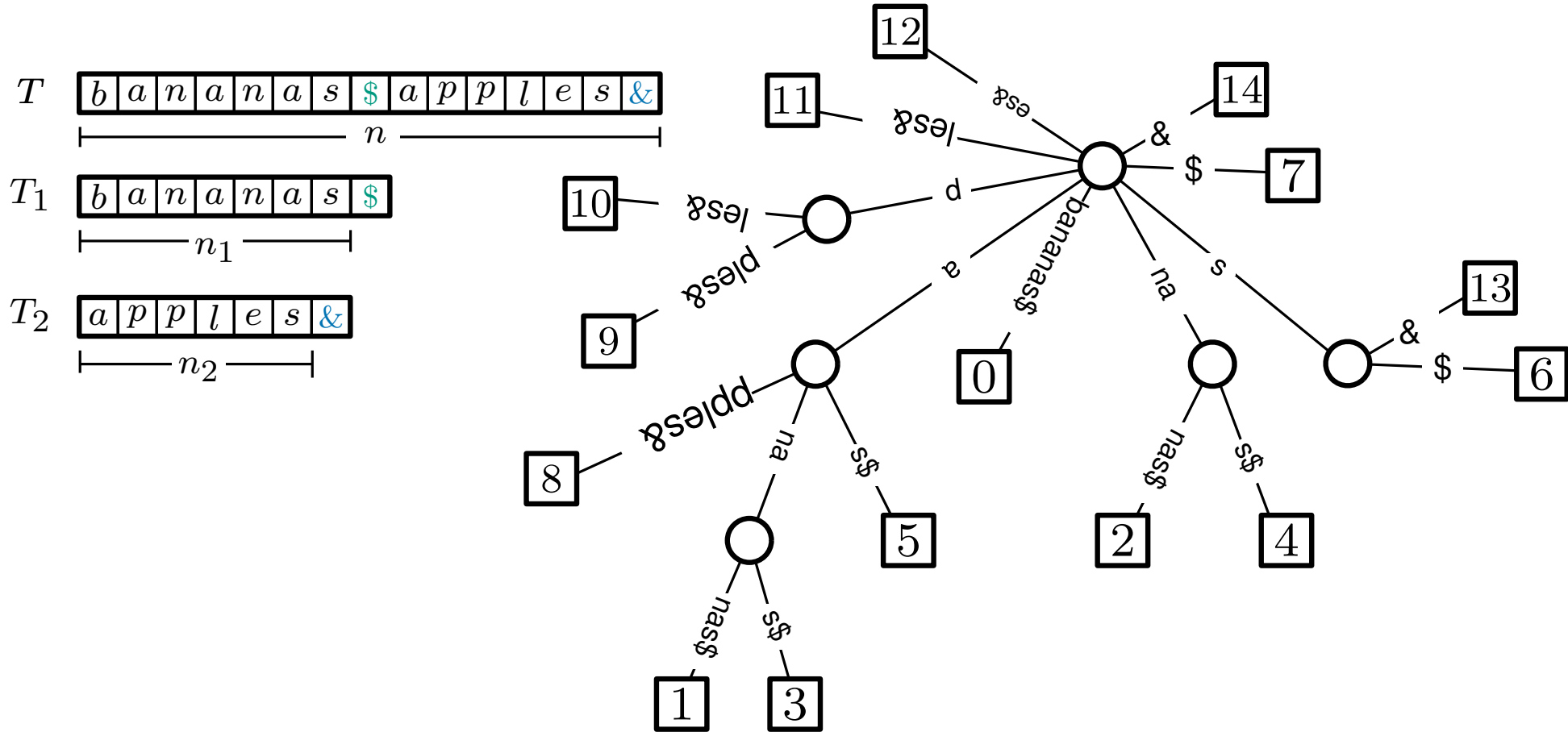
How can we index multiple texts?

- Build a generalised suffix tree in  $O(n_1 + n_2)$  space





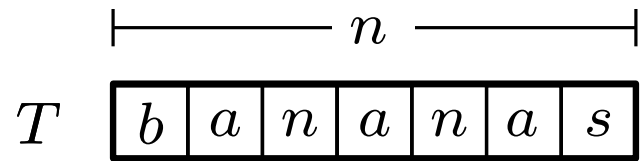
# Multiple text indexing



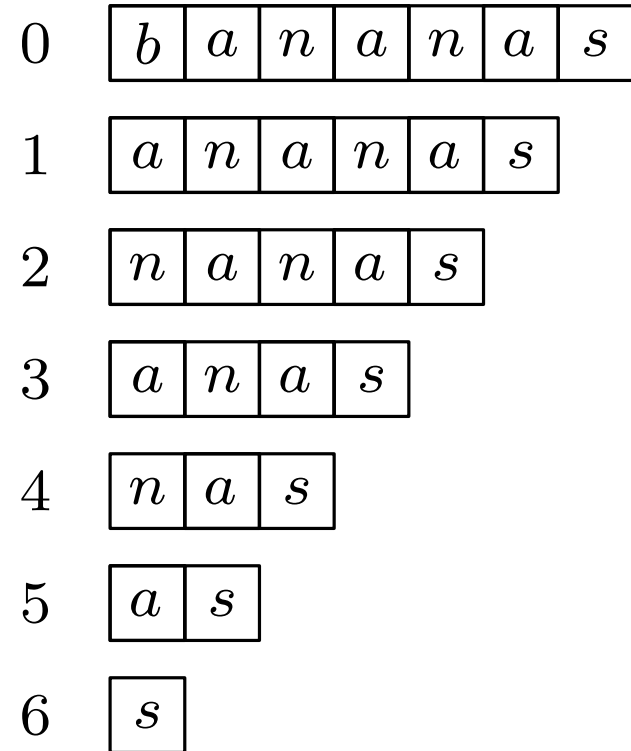
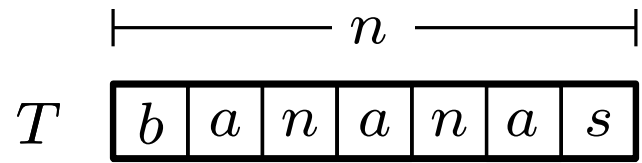
How can we index multiple texts?

- Build a generalised suffix tree in  $O(n_1 + n_2)$  space
- Using the linear time method (which we omitted), this takes  $O(n_1 + n_2)$  time
- Finding all matches of a pattern  $P$  of length  $m$  still takes  $O(m + occ)$  time  
where  $occ$  is the number of matches

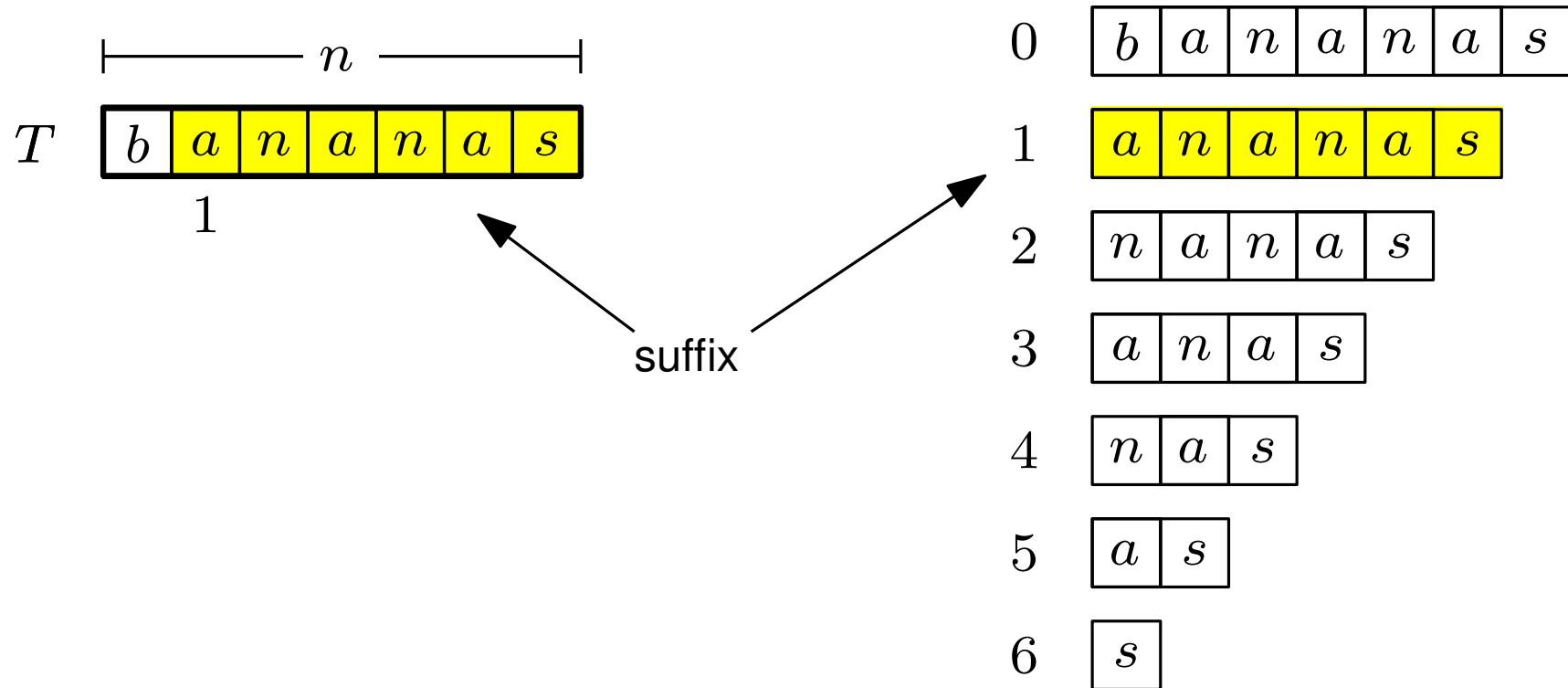
# The suffix array - a sneak preview



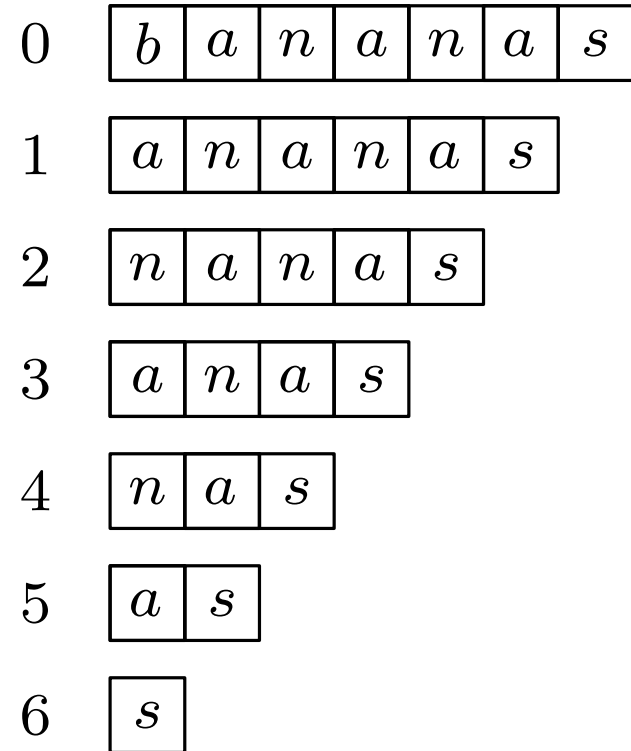
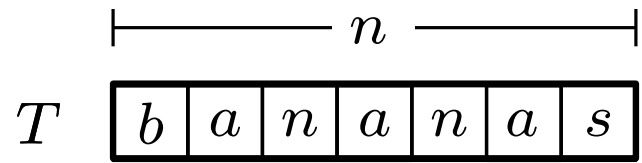
# The suffix array - a sneak preview



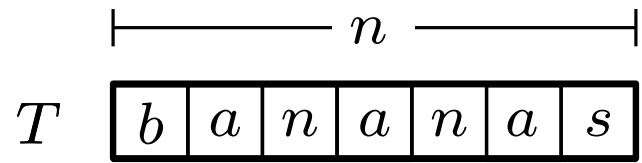
# The suffix array - a sneak preview



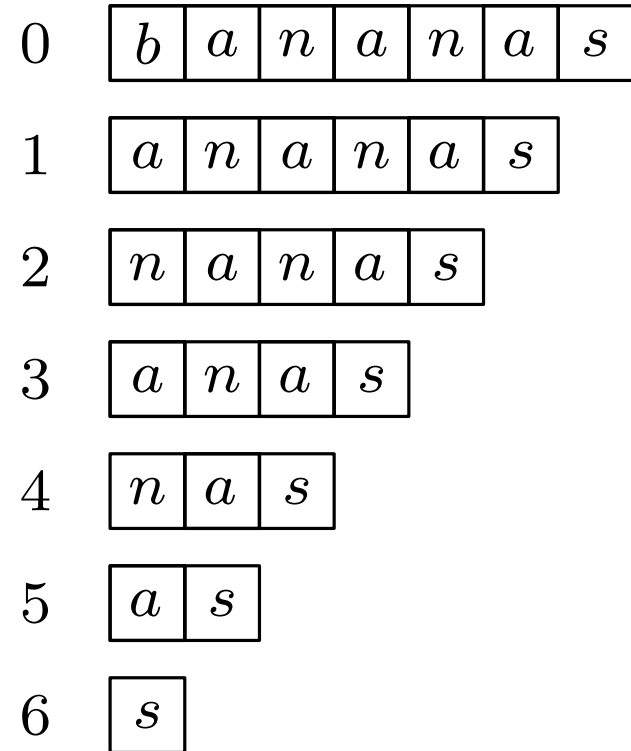
# The suffix array - a sneak preview



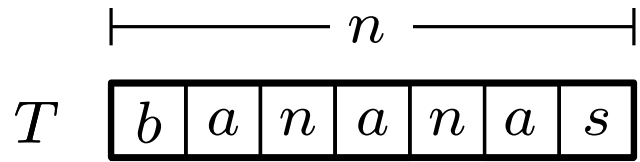
# The suffix array - a sneak preview



*Sort the suffixes  
lexicographically*

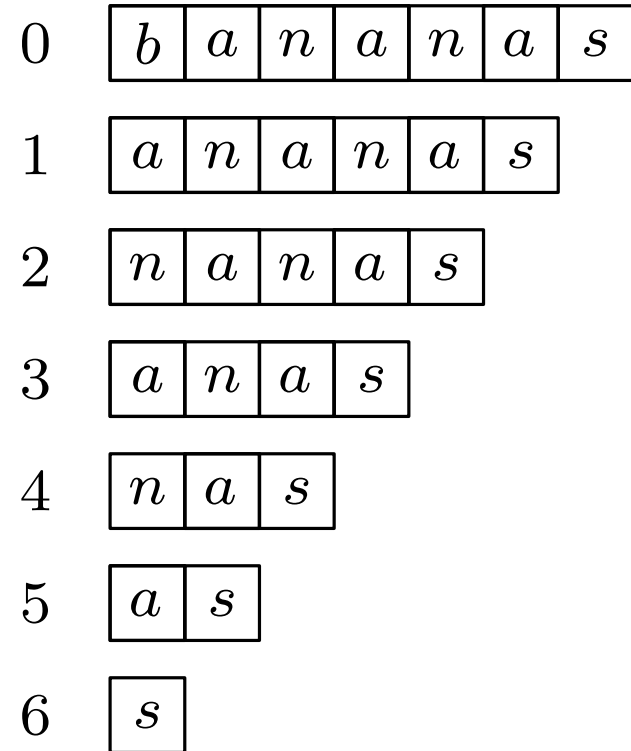


# The suffix array - a sneak preview



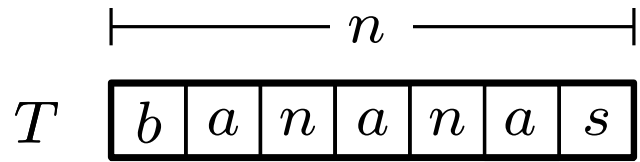
*Sort the suffixes  
lexicographically*

- The symbols themselves must have an order  
*throughout we will use alphabetical order*



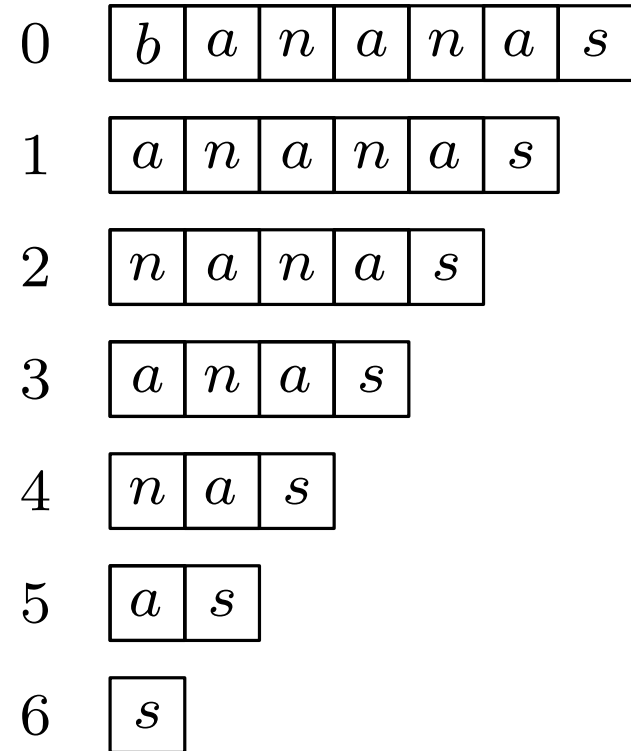


# The suffix array - a sneak preview



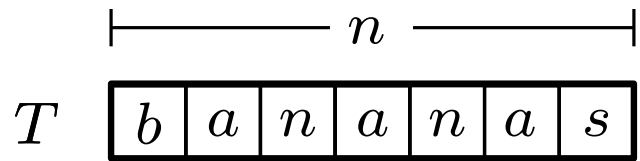
*Sort the suffixes  
lexicographically*

- The symbols themselves must have an order  
*throughout we will use alphabetical order*



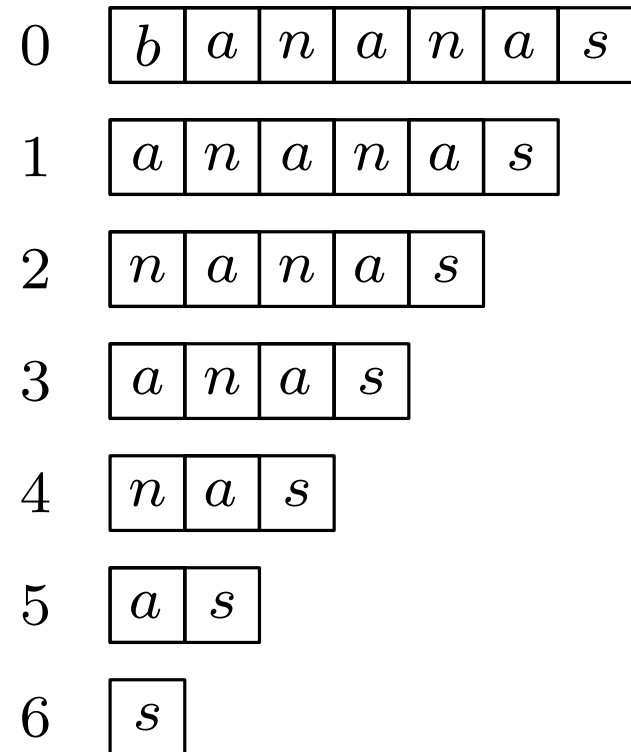
In lexicographical ordering we sort strings based on the first symbol that differs:

# The suffix array - a sneak preview

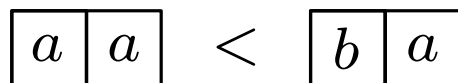


*Sort the suffixes  
lexicographically*

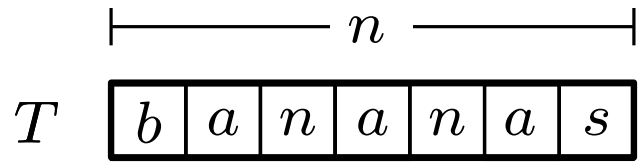
- The symbols themselves must have an order  
*throughout we will use alphabetical order*



In lexicographical ordering we sort strings based on the first symbol that differs:

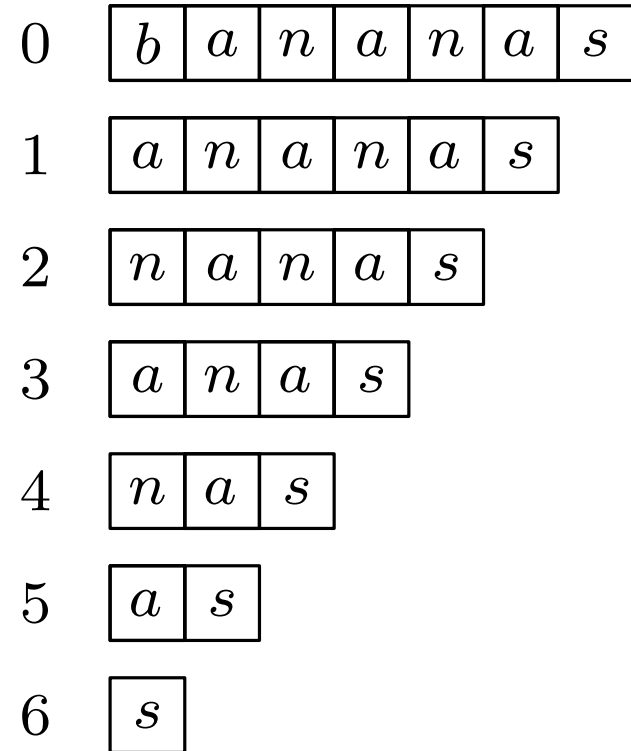


# The suffix array - a sneak preview

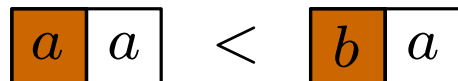


*Sort the suffixes  
lexicographically*

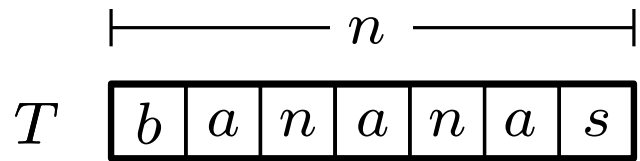
- The symbols themselves must have an order  
*throughout we will use alphabetical order*



In lexicographical ordering we sort strings based on the first symbol that differs:

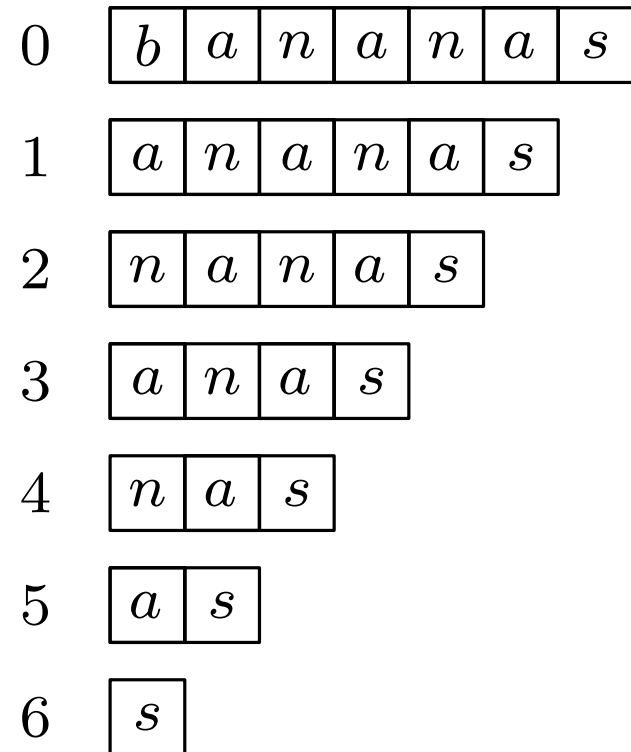


# The suffix array - a sneak preview

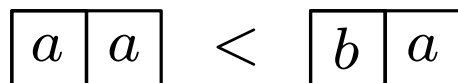


*Sort the suffixes  
lexicographically*

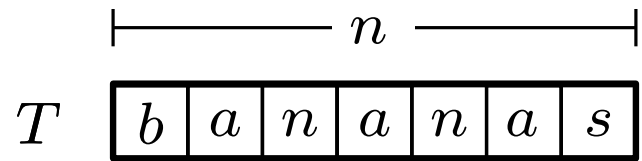
- The symbols themselves must have an order  
*throughout we will use alphabetical order*



In lexicographical ordering we sort strings based on the first symbol that differs:

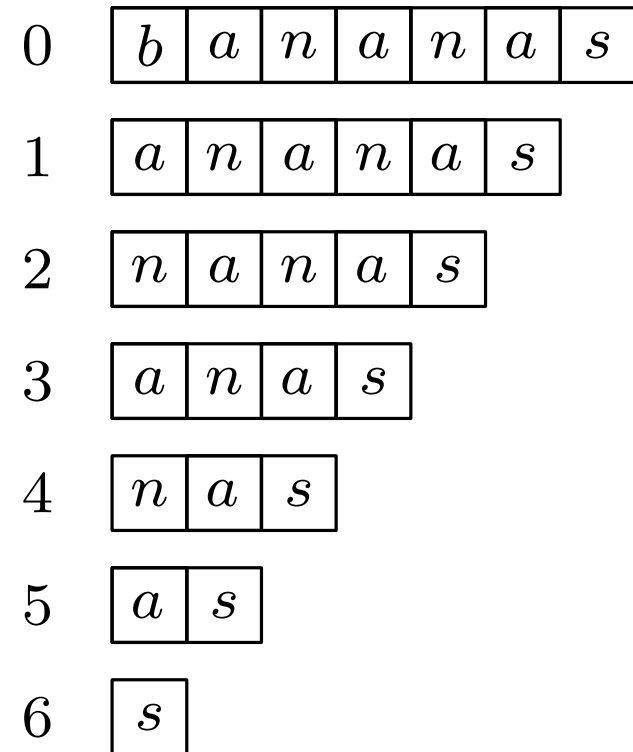


# The suffix array - a sneak preview

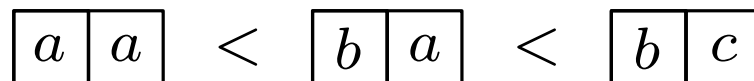


*Sort the suffixes  
lexicographically*

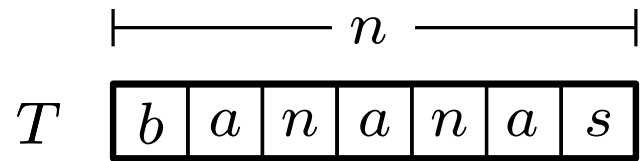
- The symbols themselves must have an order  
*throughout we will use alphabetical order*



In lexicographical ordering we sort strings based on the first symbol that differs:

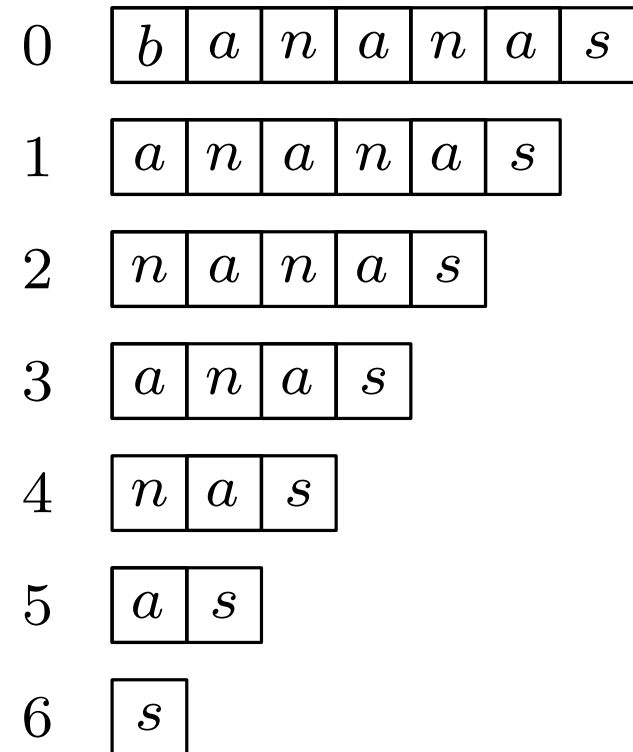


# The suffix array - a sneak preview

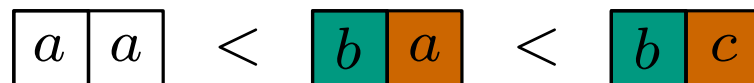


*Sort the suffixes  
lexicographically*

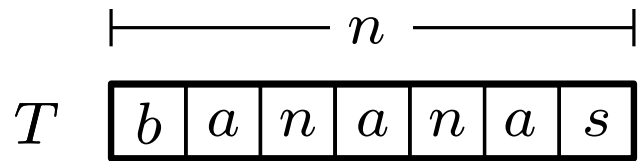
- The symbols themselves must have an order  
*throughout we will use alphabetical order*



In lexicographical ordering we sort strings based on the first symbol that differs:

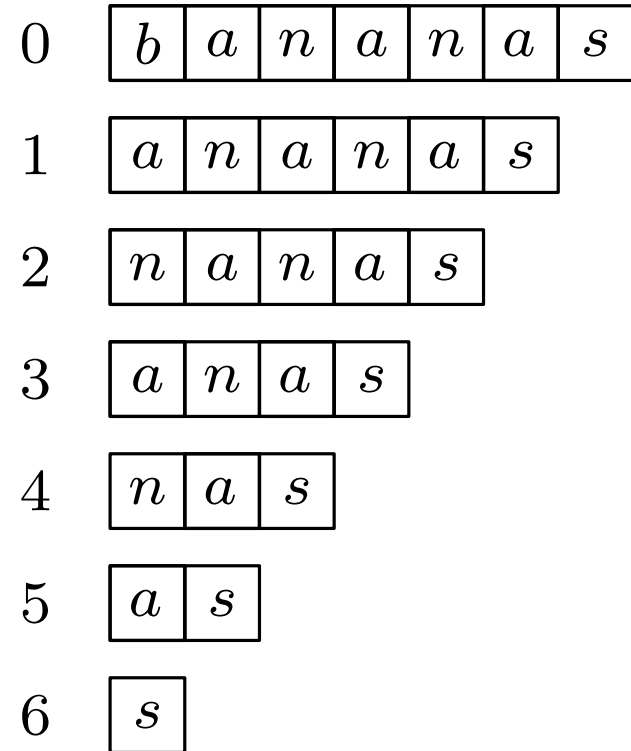


# The suffix array - a sneak preview

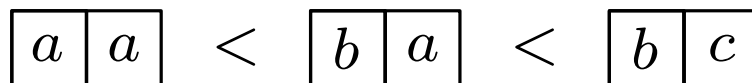


*Sort the suffixes  
lexicographically*

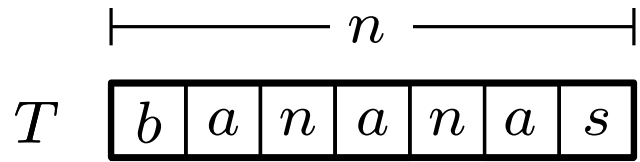
- The symbols themselves must have an order  
*throughout we will use alphabetical order*



In lexicographical ordering we sort strings based on the first symbol that differs:

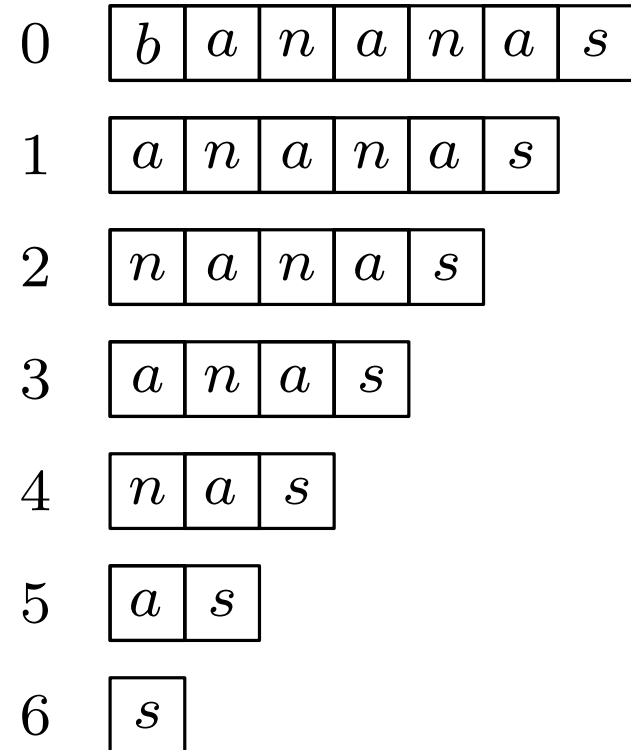


# The suffix array - a sneak preview

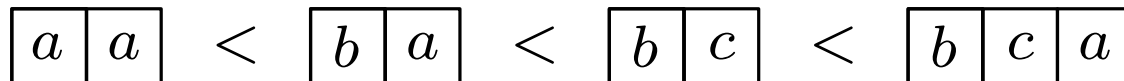


*Sort the suffixes  
lexicographically*

- The symbols themselves must have an order  
*throughout we will use alphabetical order*

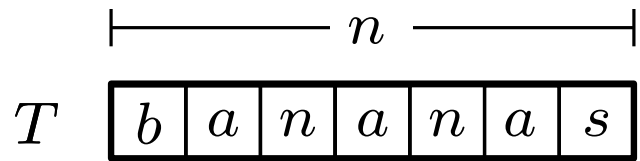


In lexicographical ordering we sort strings based on the first symbol that differs:



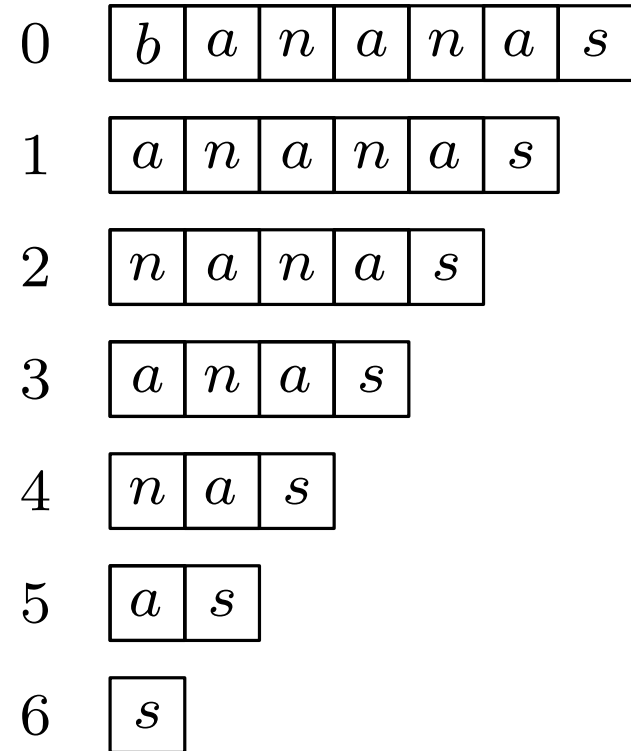


# The suffix array - a sneak preview

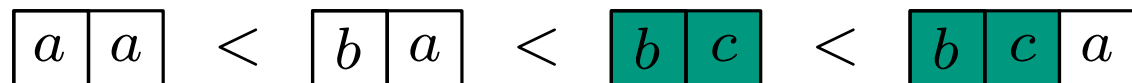


*Sort the suffixes  
lexicographically*

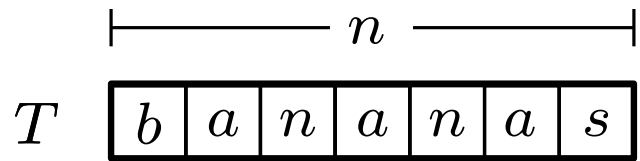
- The symbols themselves must have an order  
*throughout we will use alphabetical order*



In lexicographical ordering we sort strings based on the first symbol that differs:

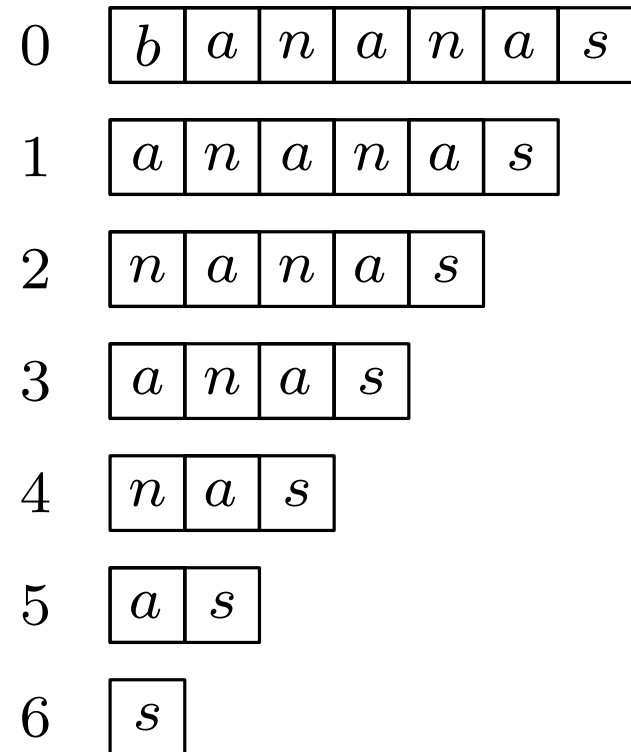


# The suffix array - a sneak preview

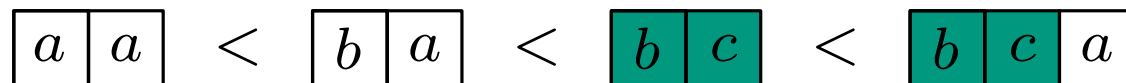


*Sort the suffixes  
lexicographically*

- The symbols themselves must have an order  
*throughout we will use alphabetical order*

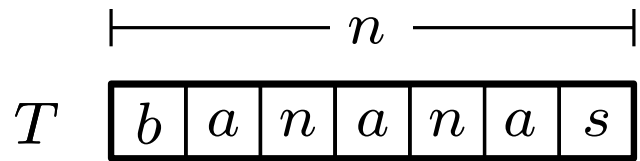


In lexicographical ordering we sort strings based on the first symbol that differs:



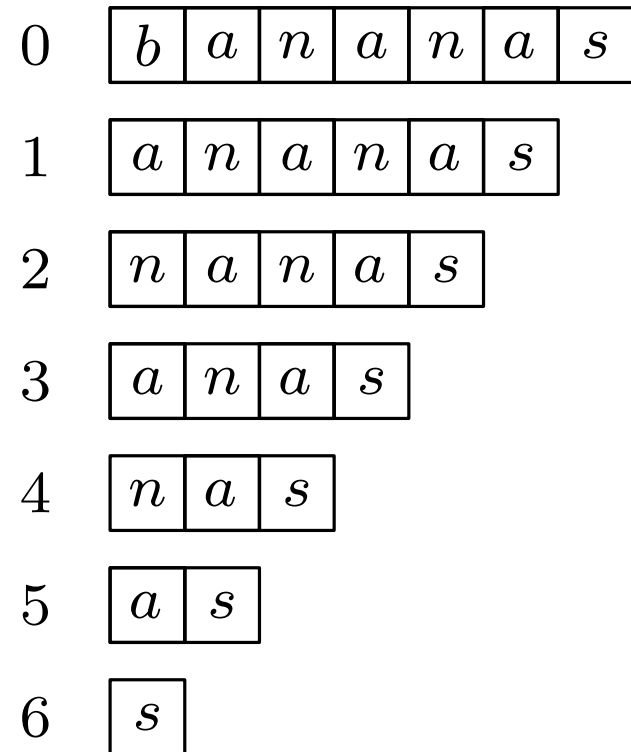
*(in a 'tie', the shorter string is smaller)*

# The suffix array - a sneak preview

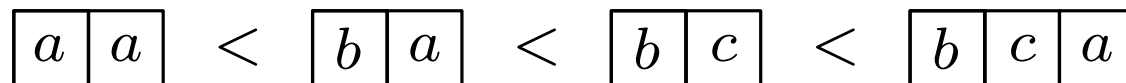


*Sort the suffixes  
lexicographically*

- The symbols themselves must have an order  
*throughout we will use alphabetical order*

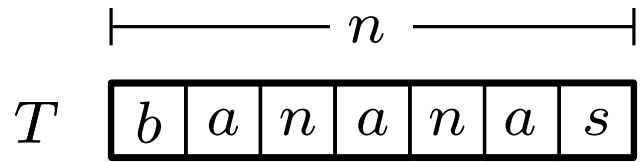


In lexicographical ordering we sort strings based on the first symbol that differs:



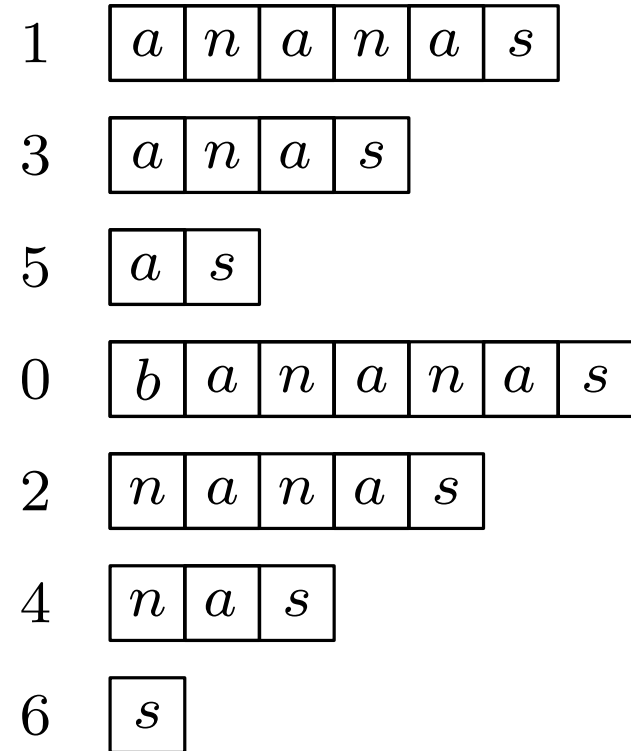
*(in a 'tie', the shorter string is smaller)*

# The suffix array - a sneak preview

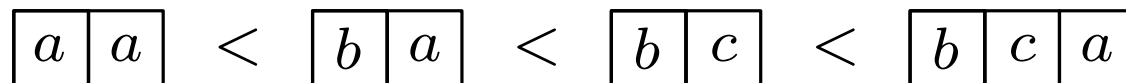


Sort the suffixes  
lexicographically

- The symbols themselves must have an order  
*throughout we will use alphabetical order*

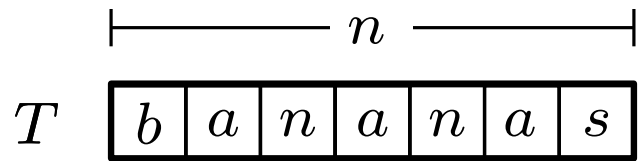


In lexicographical ordering we sort strings based on the first symbol that differs:



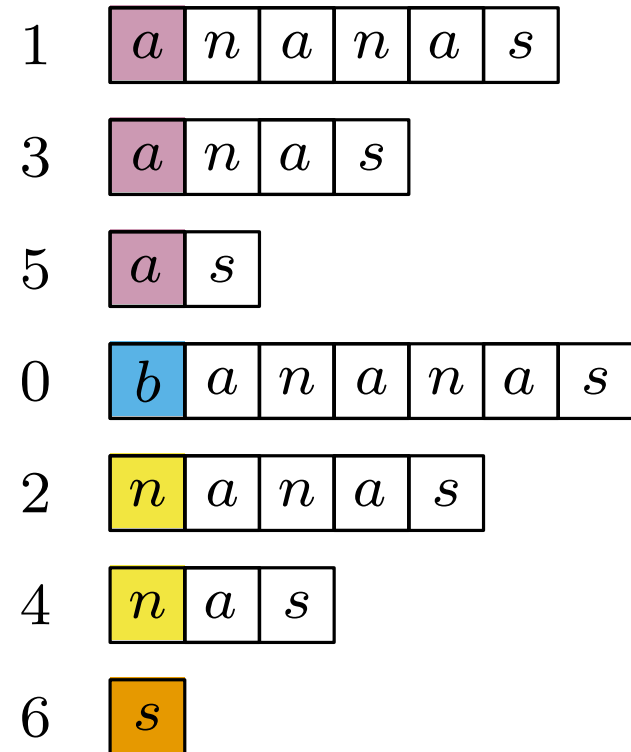
*(in a 'tie', the shorter string is smaller)*

# The suffix array - a sneak preview

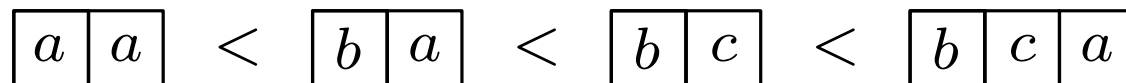


Sort the suffixes  
lexicographically

- The symbols themselves must have an order  
*throughout we will use alphabetical order*



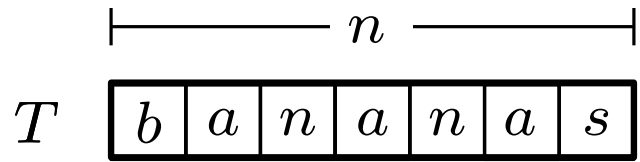
In lexicographical ordering we sort strings based on the first symbol that differs:



*(in a 'tie', the shorter string is smaller)*

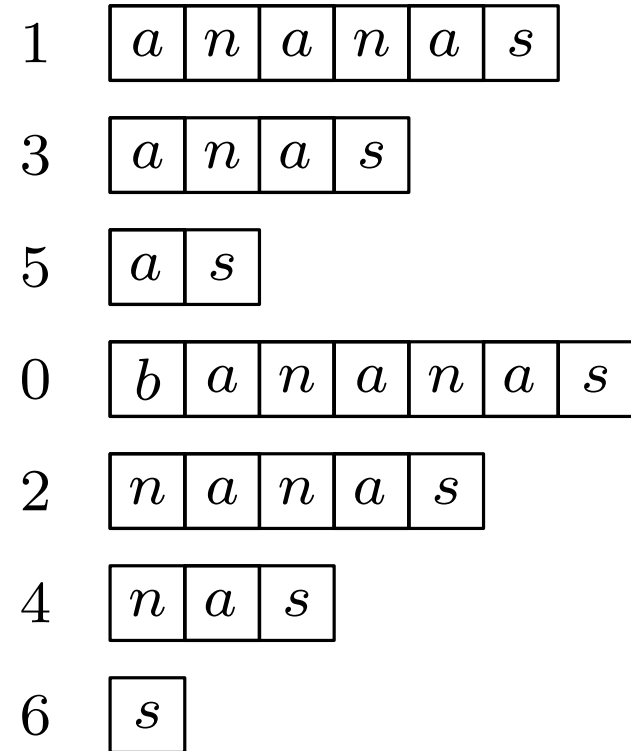


# The suffix array - a sneak preview

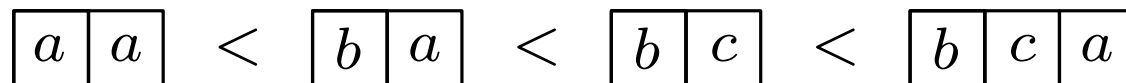


*Sort the suffixes  
lexicographically*

- The symbols themselves must have an order  
*throughout we will use alphabetical order*

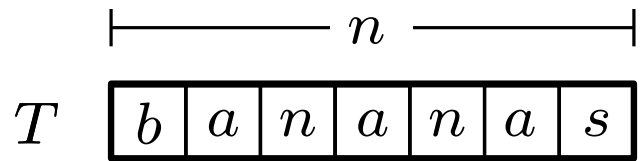


In lexicographical ordering we sort strings based on the first symbol that differs:



*(in a 'tie', the shorter string is smaller)*

# The suffix array - a sneak preview

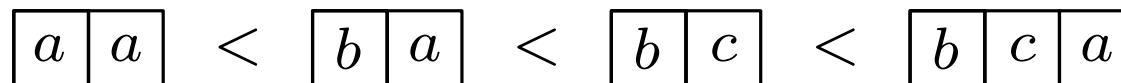


Sort the suffixes  
lexicographically

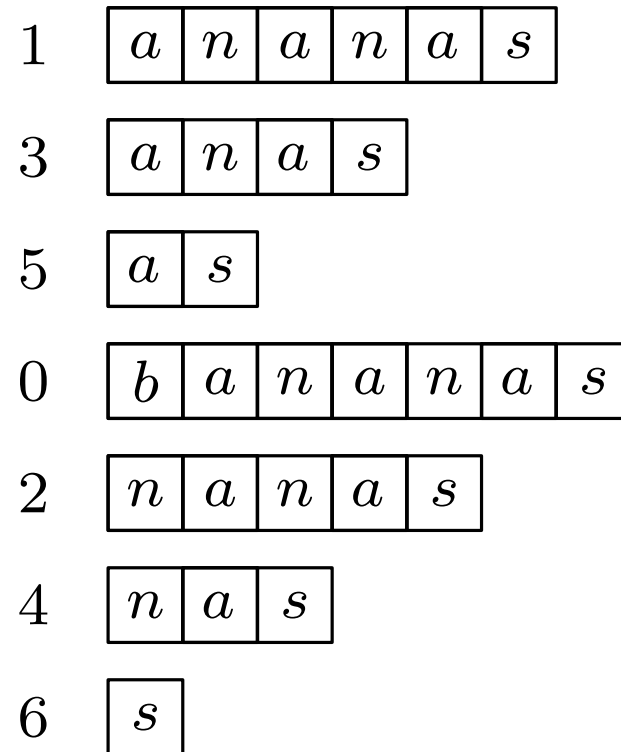
- The symbols themselves must have an order  
*throughout we will use alphabetical order*

just a fancy name for the order the strings would appear in a dictionary

In **lexicographical** ordering we sort strings based on the first symbol that differs:

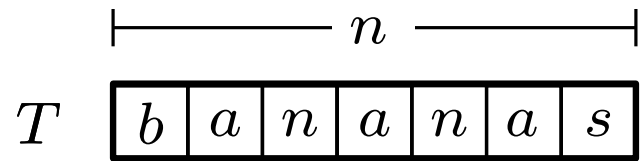


*(in a 'tie', the shorter string is smaller)*





# The suffix array - a sneak preview

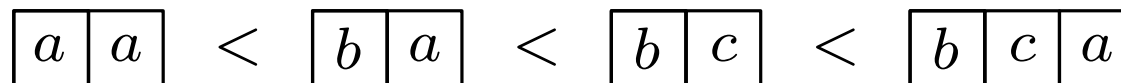


Sort the suffixes  
lexicographically

- The symbols themselves must have an order  
*throughout we will use alphabetical order*

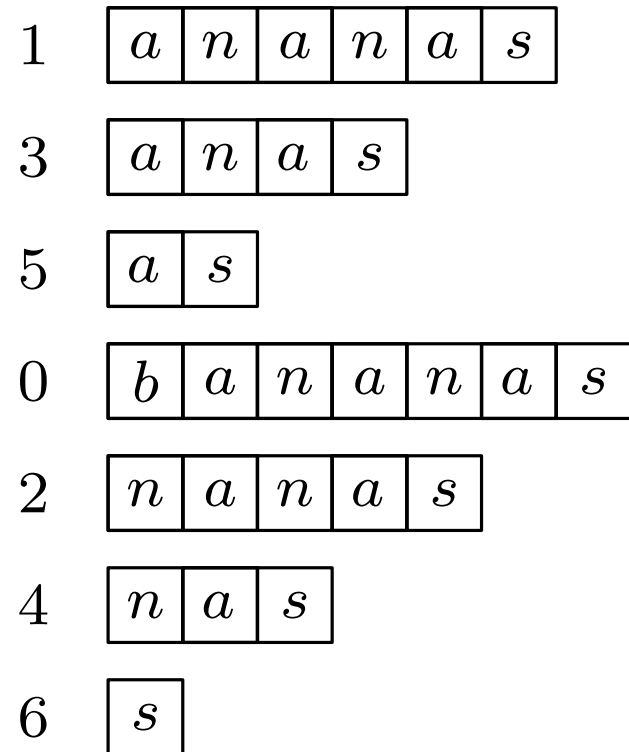
just a fancy name for the order the strings would appear in a dictionary

In **lexicographical** ordering we sort strings based on the first symbol that differs:

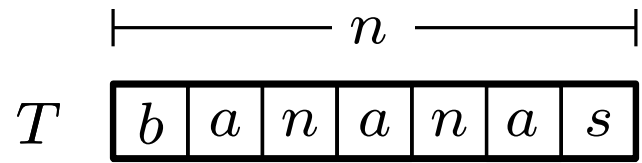


*(in a 'tie', the shorter string is smaller)*

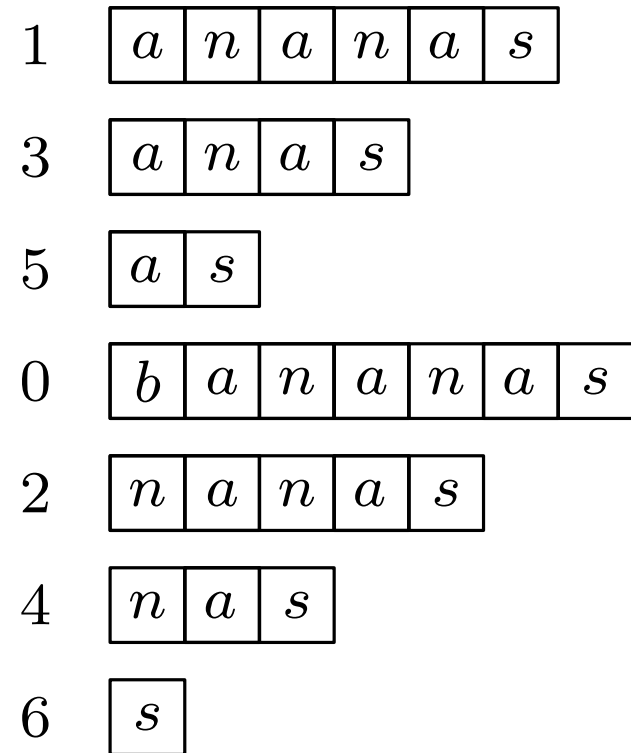
If the symbols don't have a natural order, we use their binary representation in memory



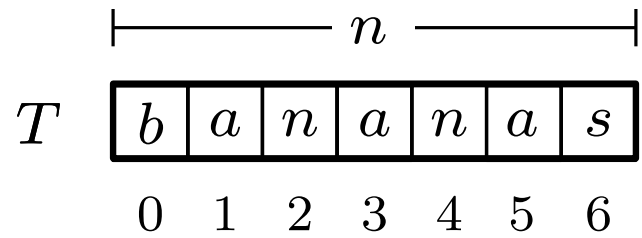
# The suffix array - a sneak preview



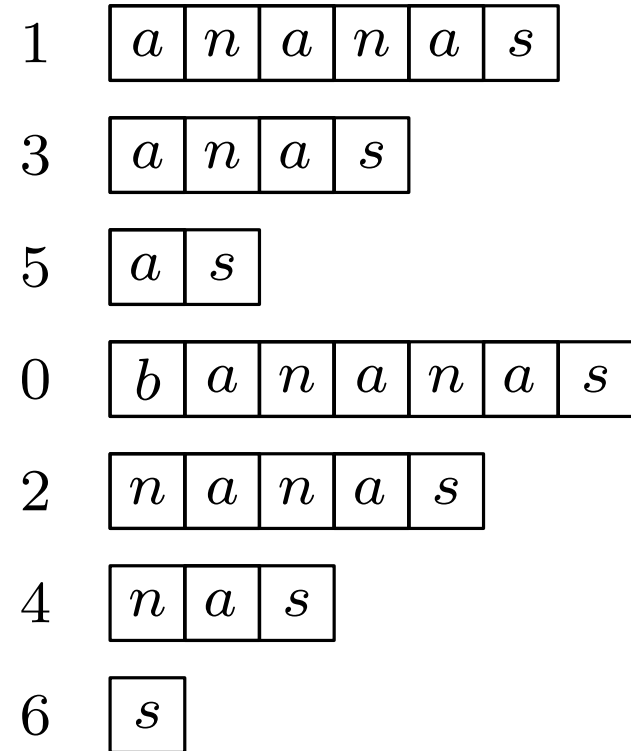
*Sort the suffixes  
lexicographically*



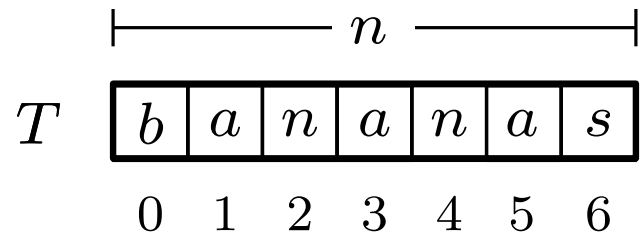
# The suffix array - a sneak preview



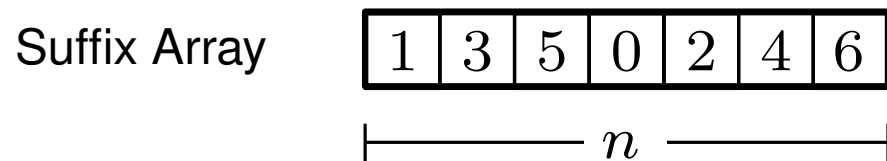
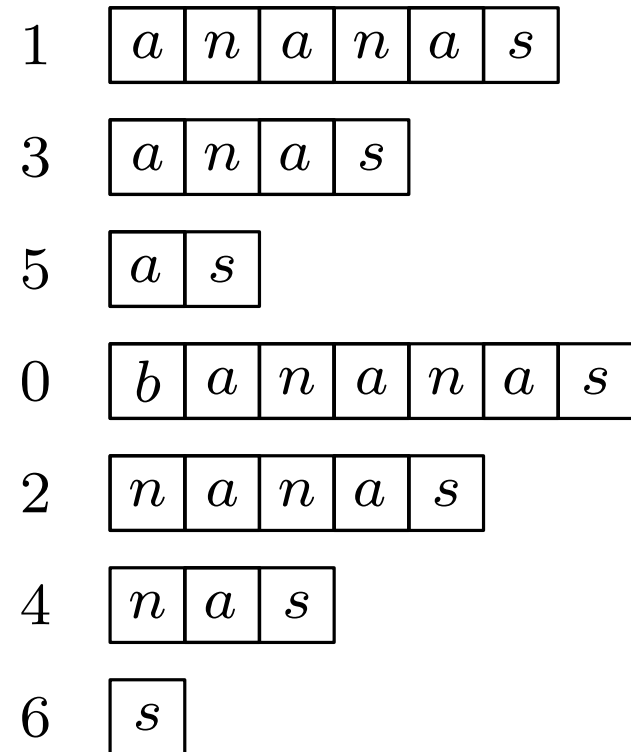
*Sort the suffixes  
lexicographically*



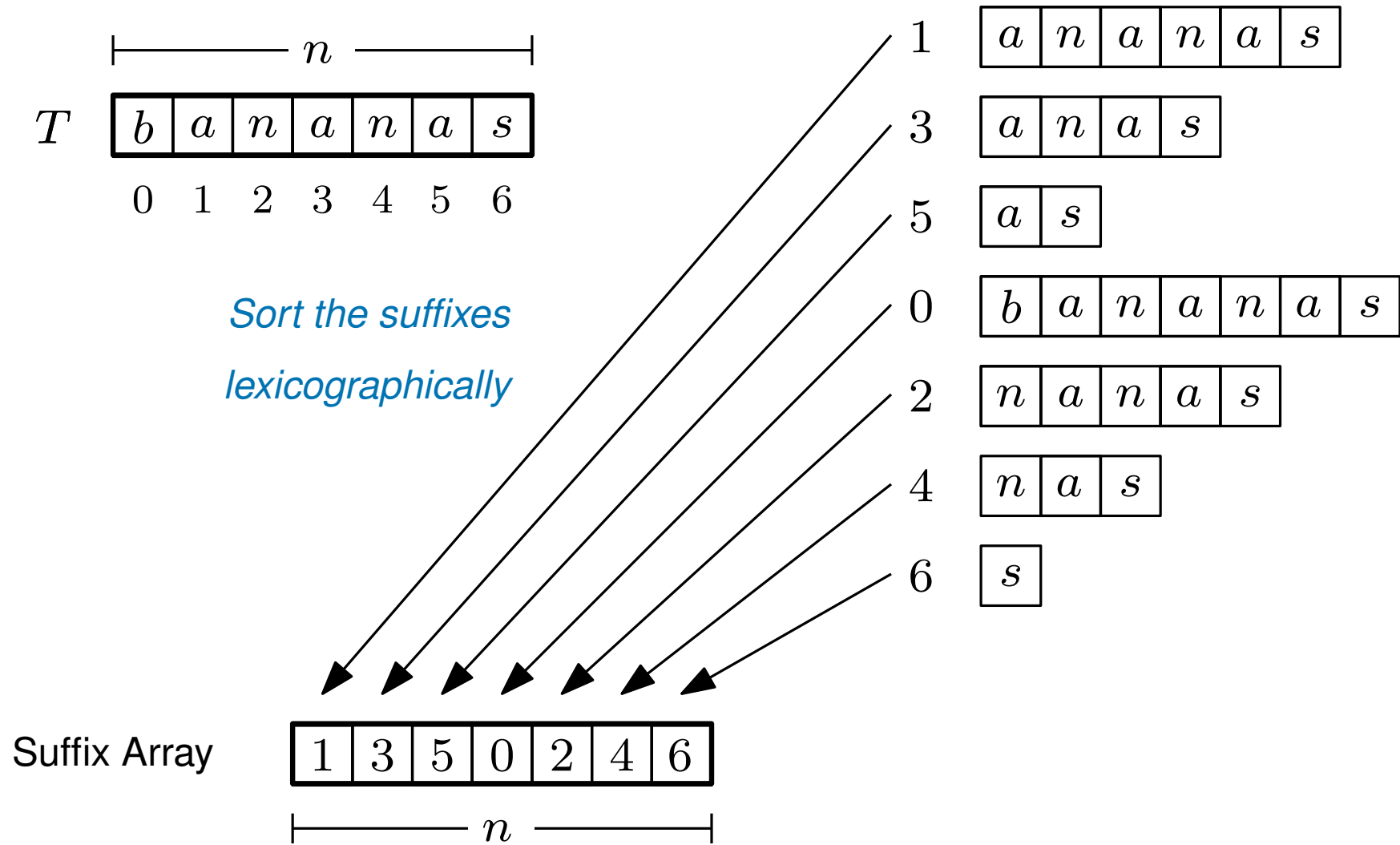
# The suffix array - a sneak preview



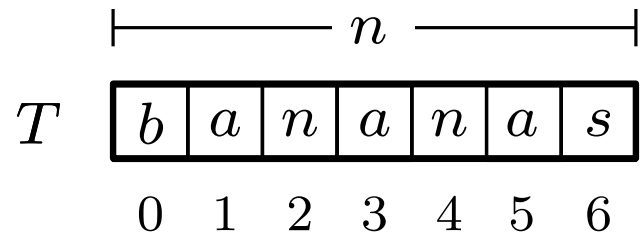
*Sort the suffixes  
lexicographically*



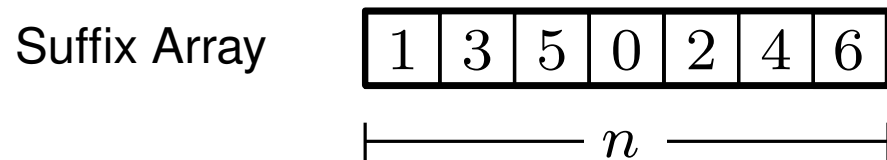
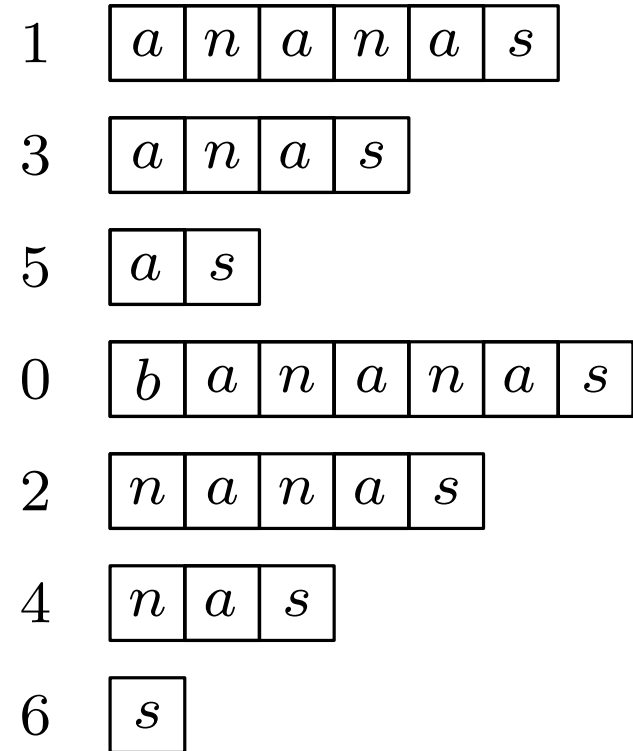
# The suffix array - a sneak preview



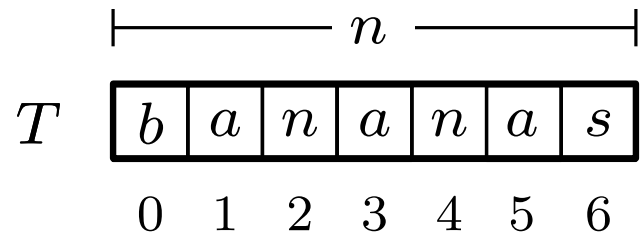
# The suffix array - a sneak preview



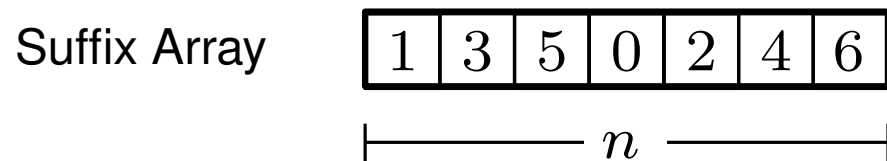
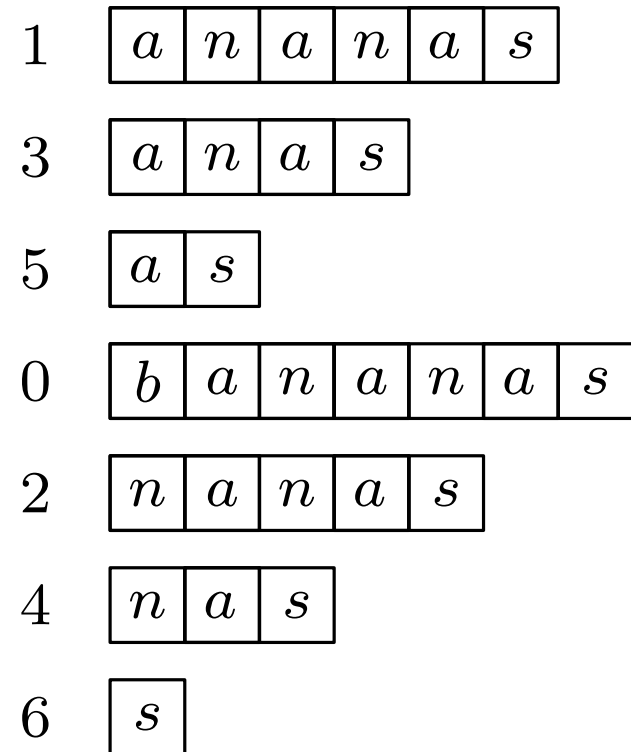
*Sort the suffixes  
lexicographically*



# The suffix array - a sneak preview

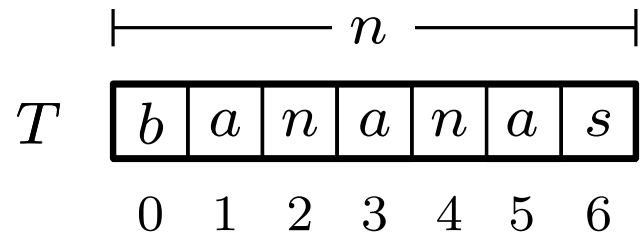


*Sort the suffixes  
lexicographically*

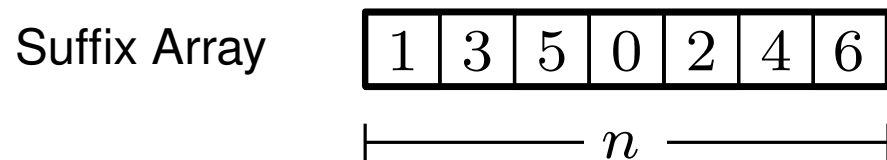
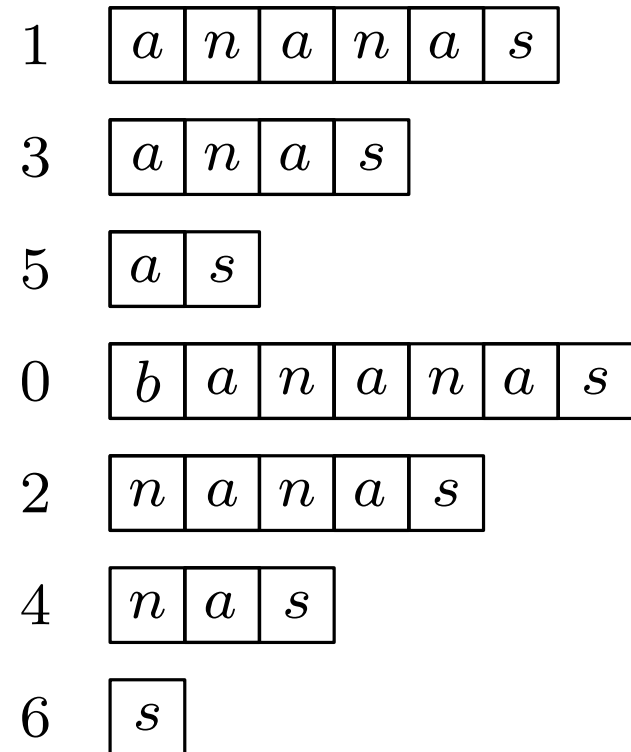


*The suffix array is much smaller than the suffix tree (in terms of constants)*

# The suffix array - a sneak preview



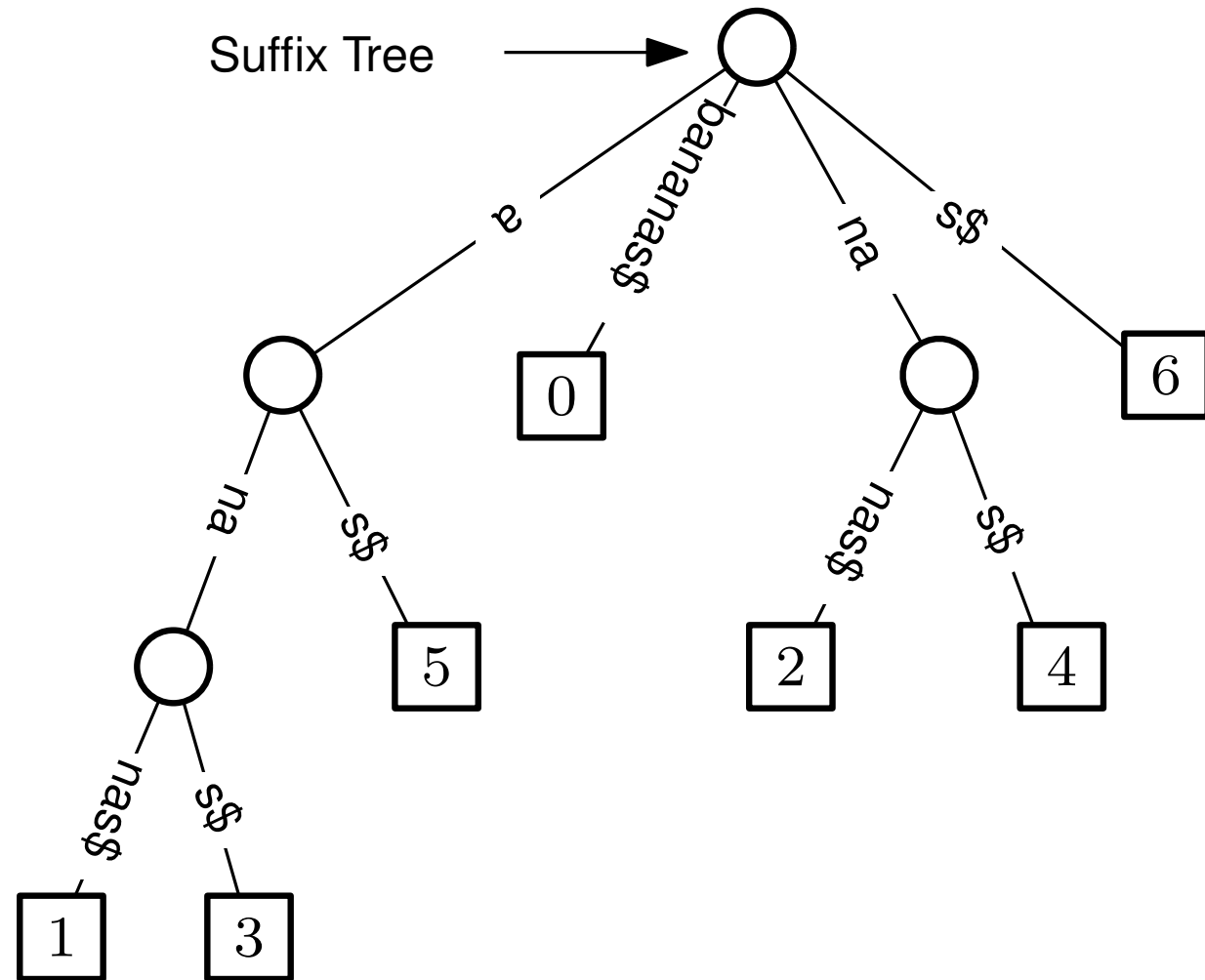
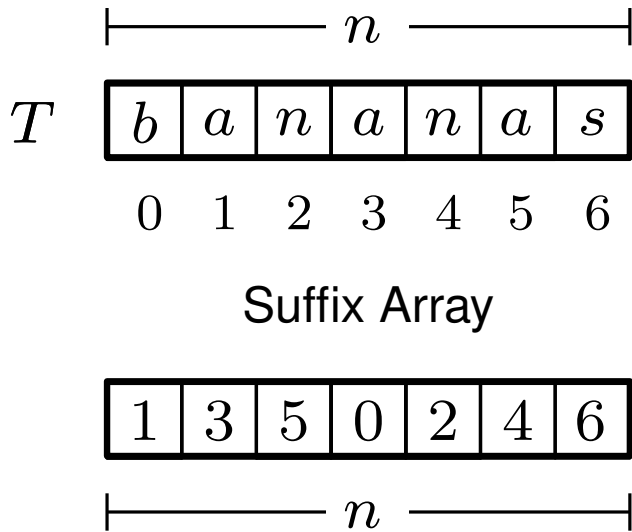
*Sort the suffixes  
lexicographically*



*The suffix array is much smaller than the suffix tree (in terms of constants)*



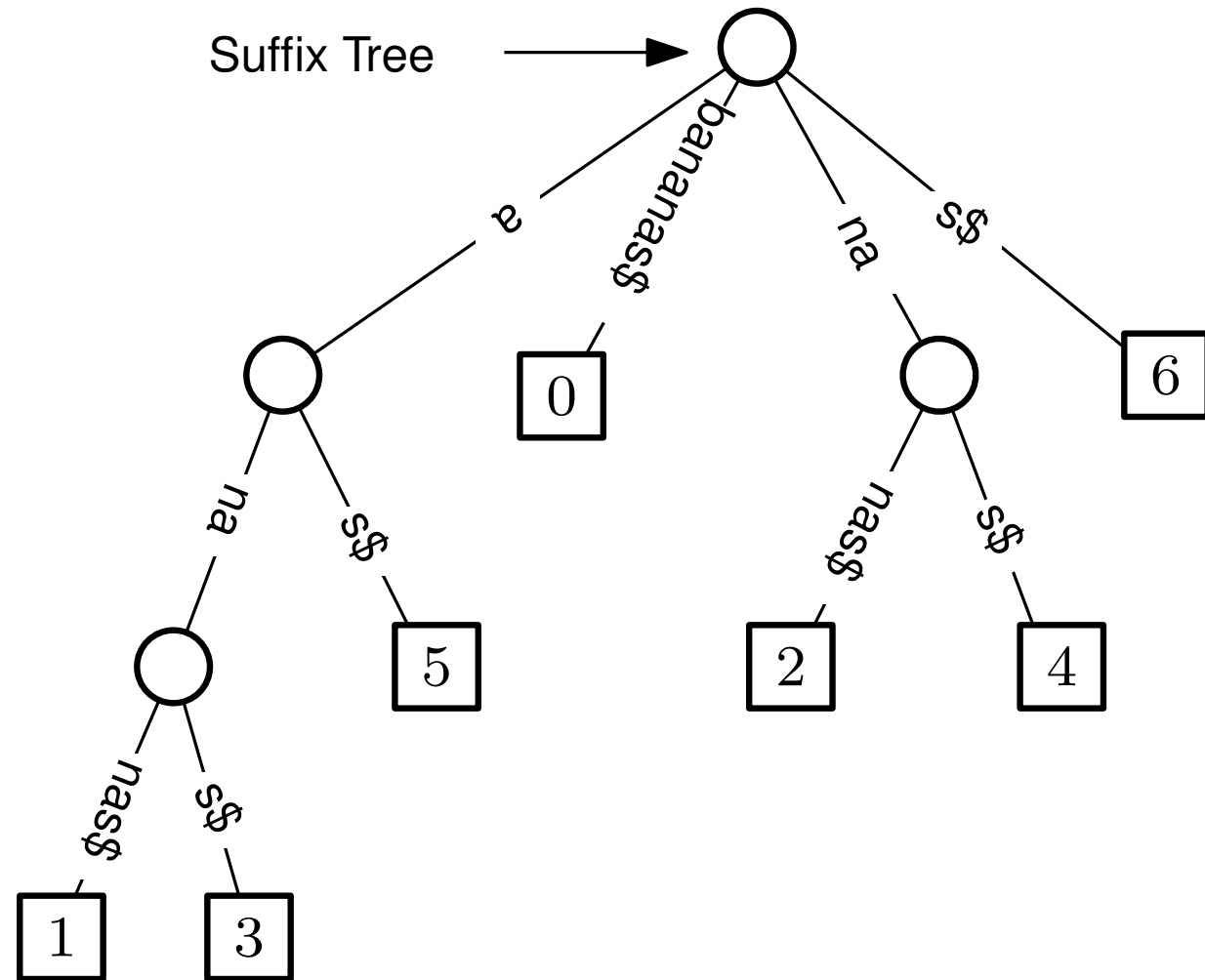
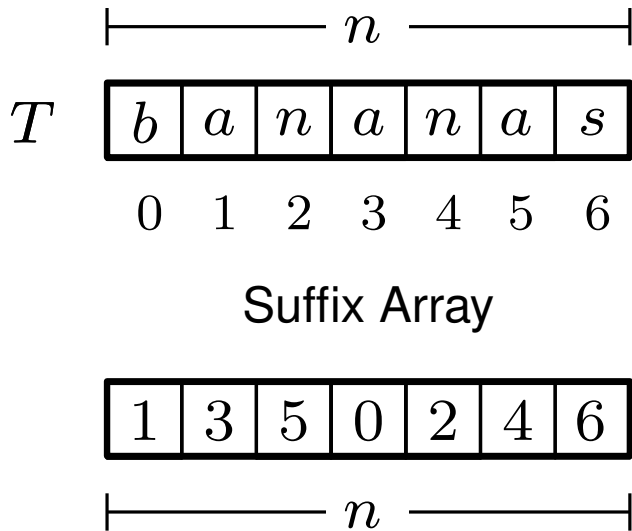
# Constructing the Suffix Array from the Suffix Tree



recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

# Constructing the Suffix Array from the Suffix Tree

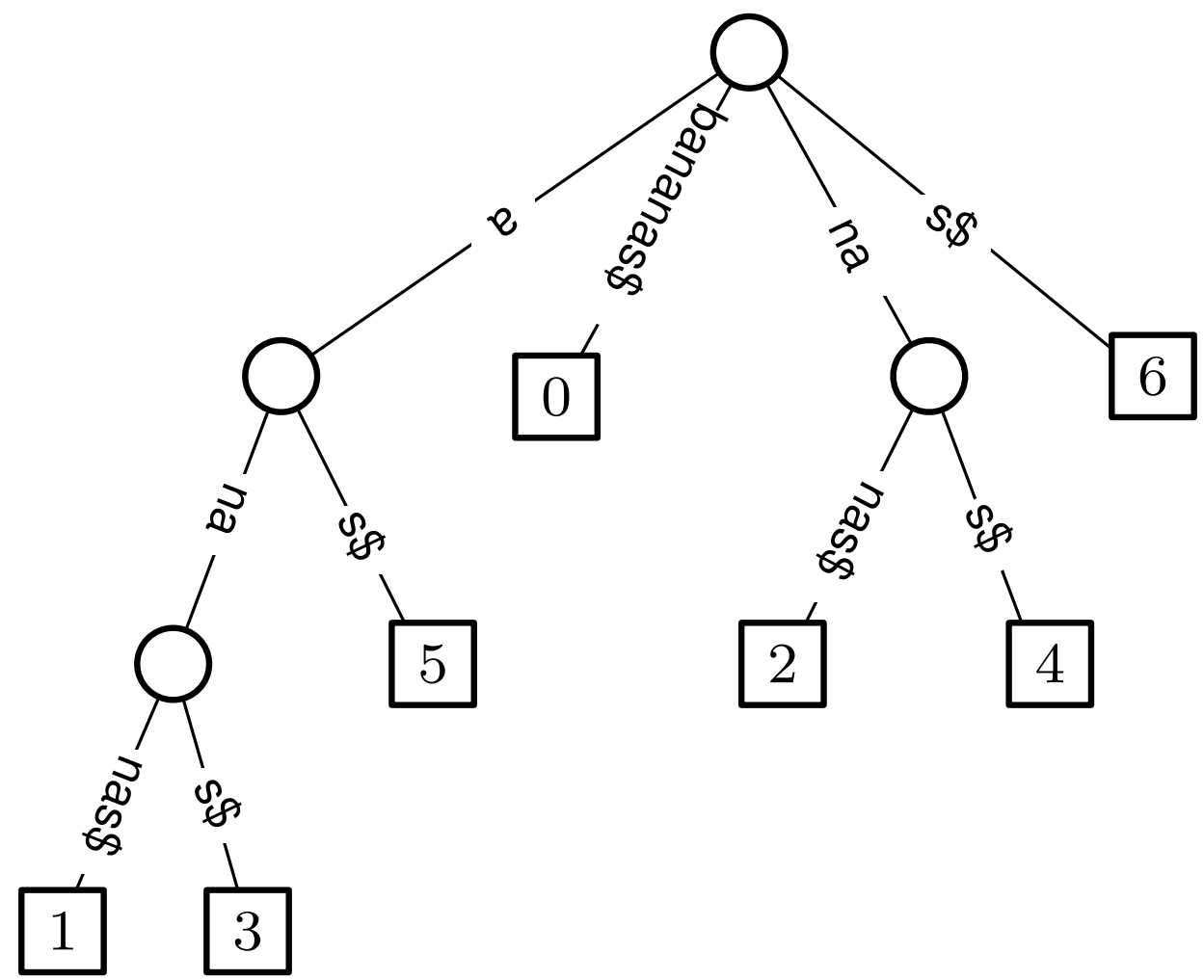
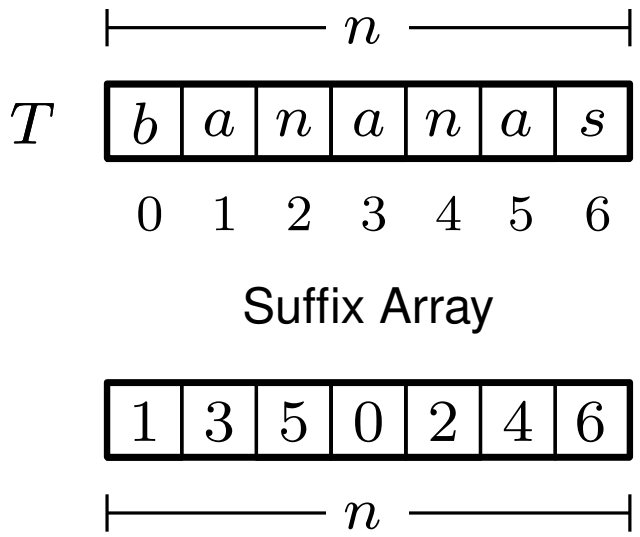


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

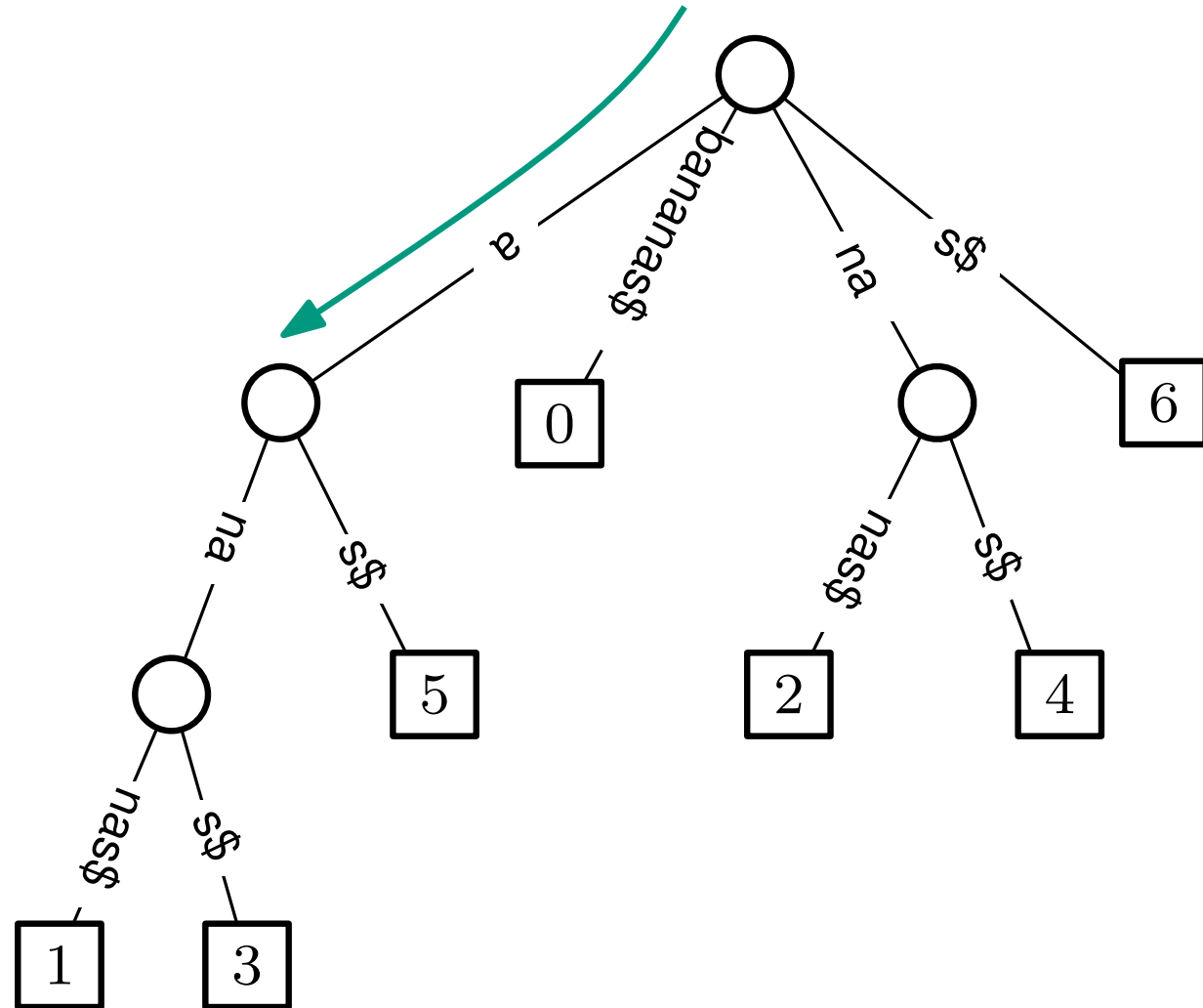
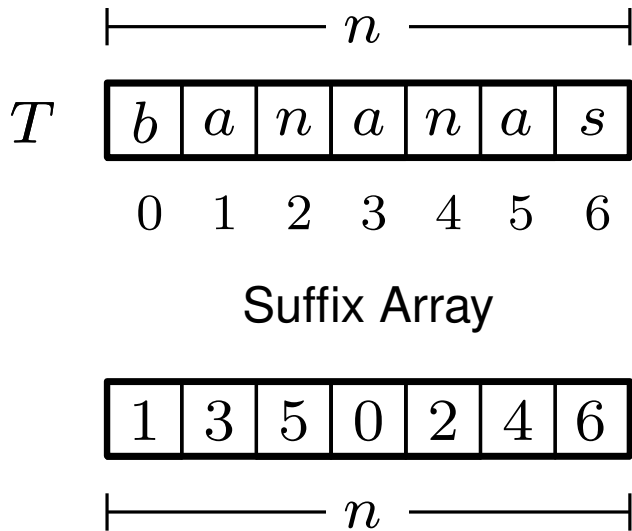


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

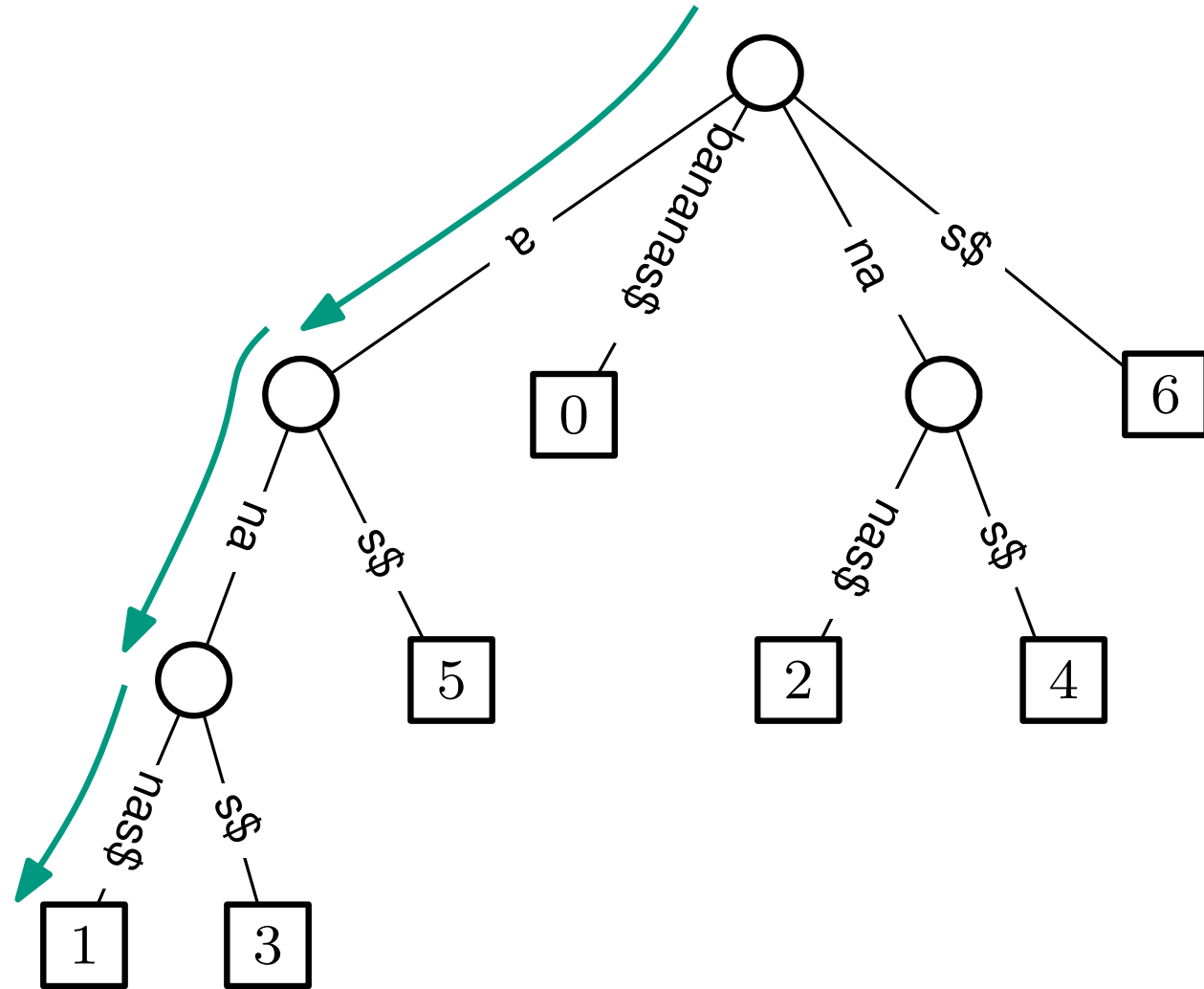
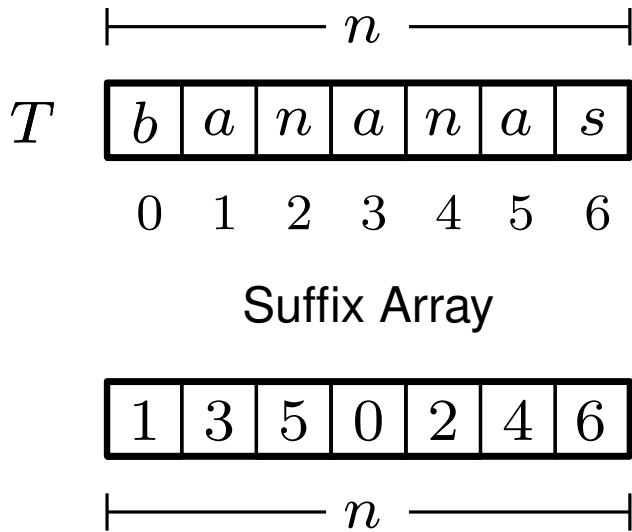


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree



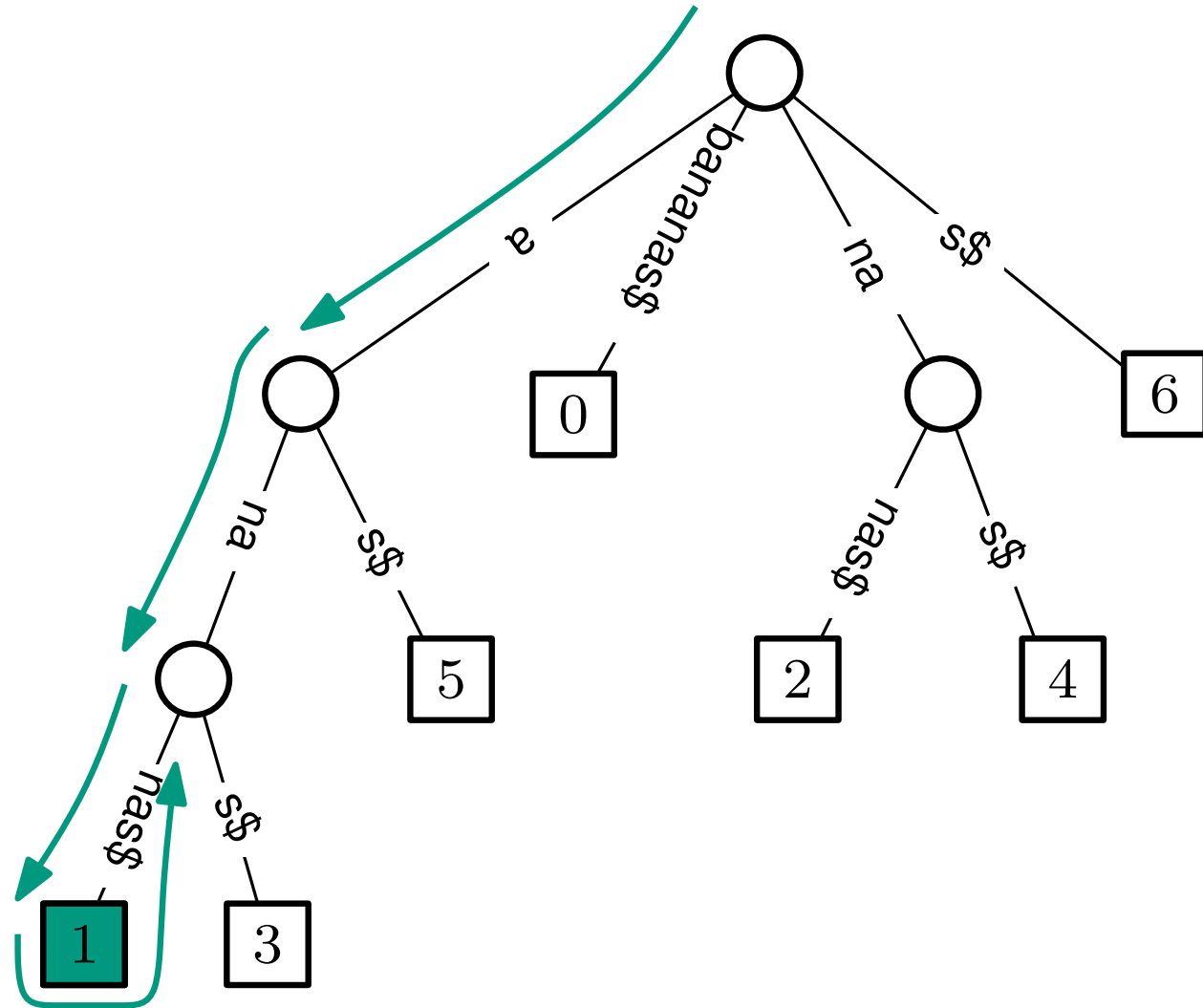
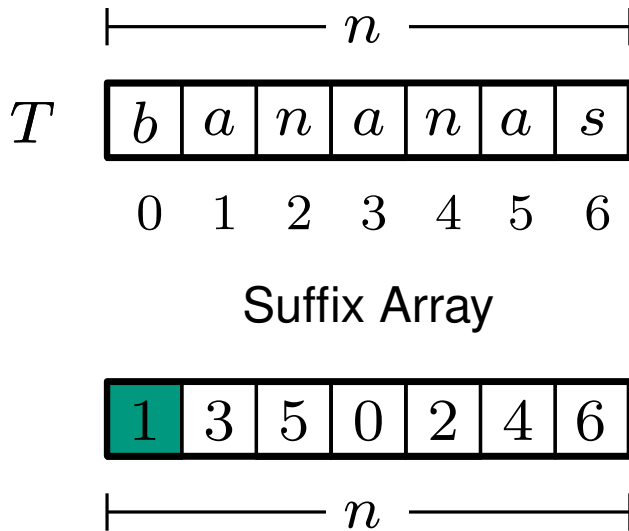
recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)



# Constructing the Suffix Array from the Suffix Tree



recall that we added a unique symbol \$ to make sure the tree exists

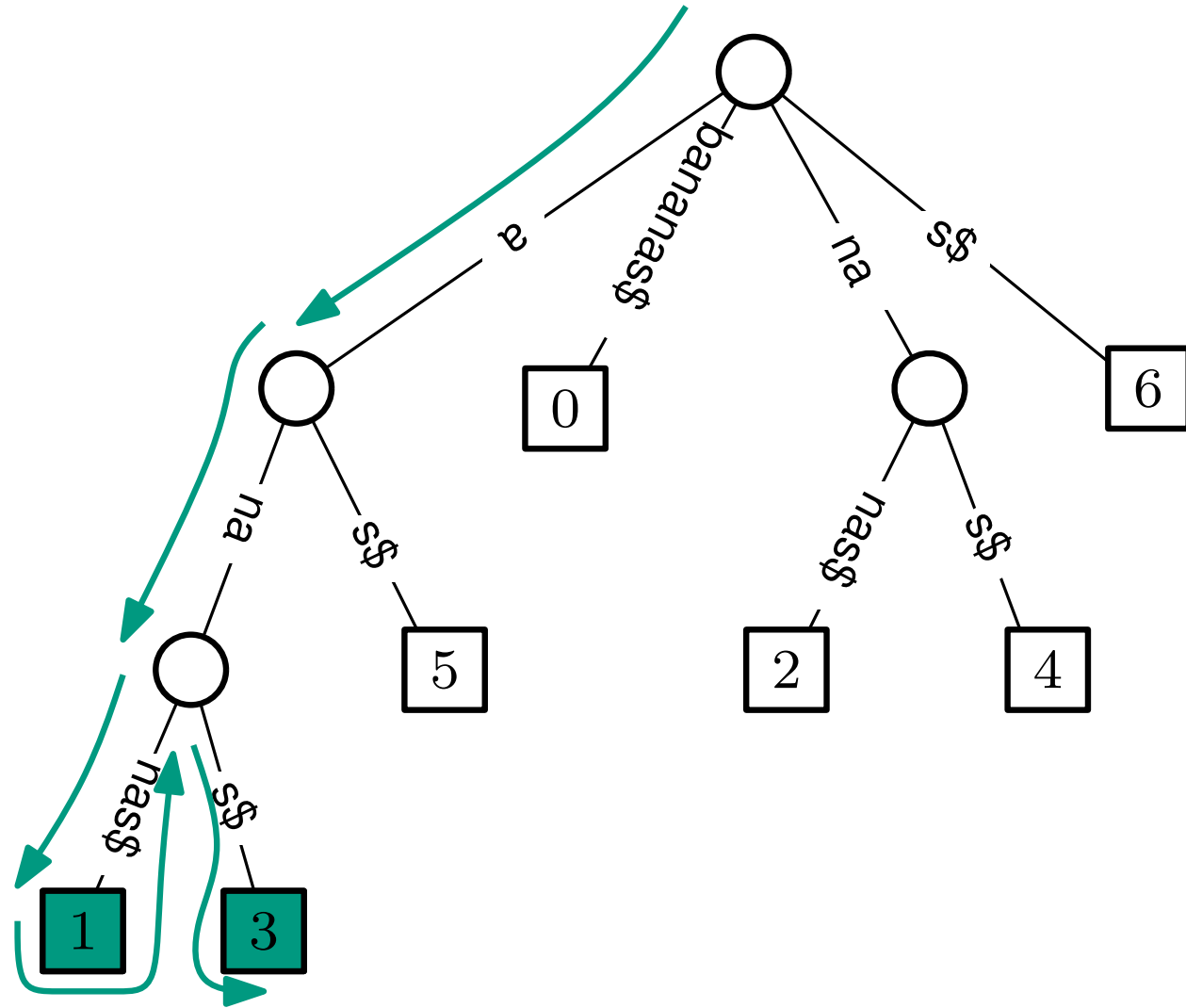
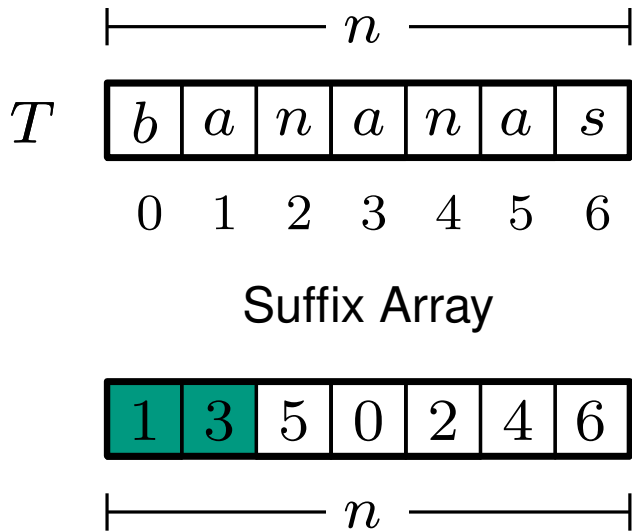
- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)





# Constructing the Suffix Array from the Suffix Tree



recall that we added a unique symbol \$ to make sure the tree exists

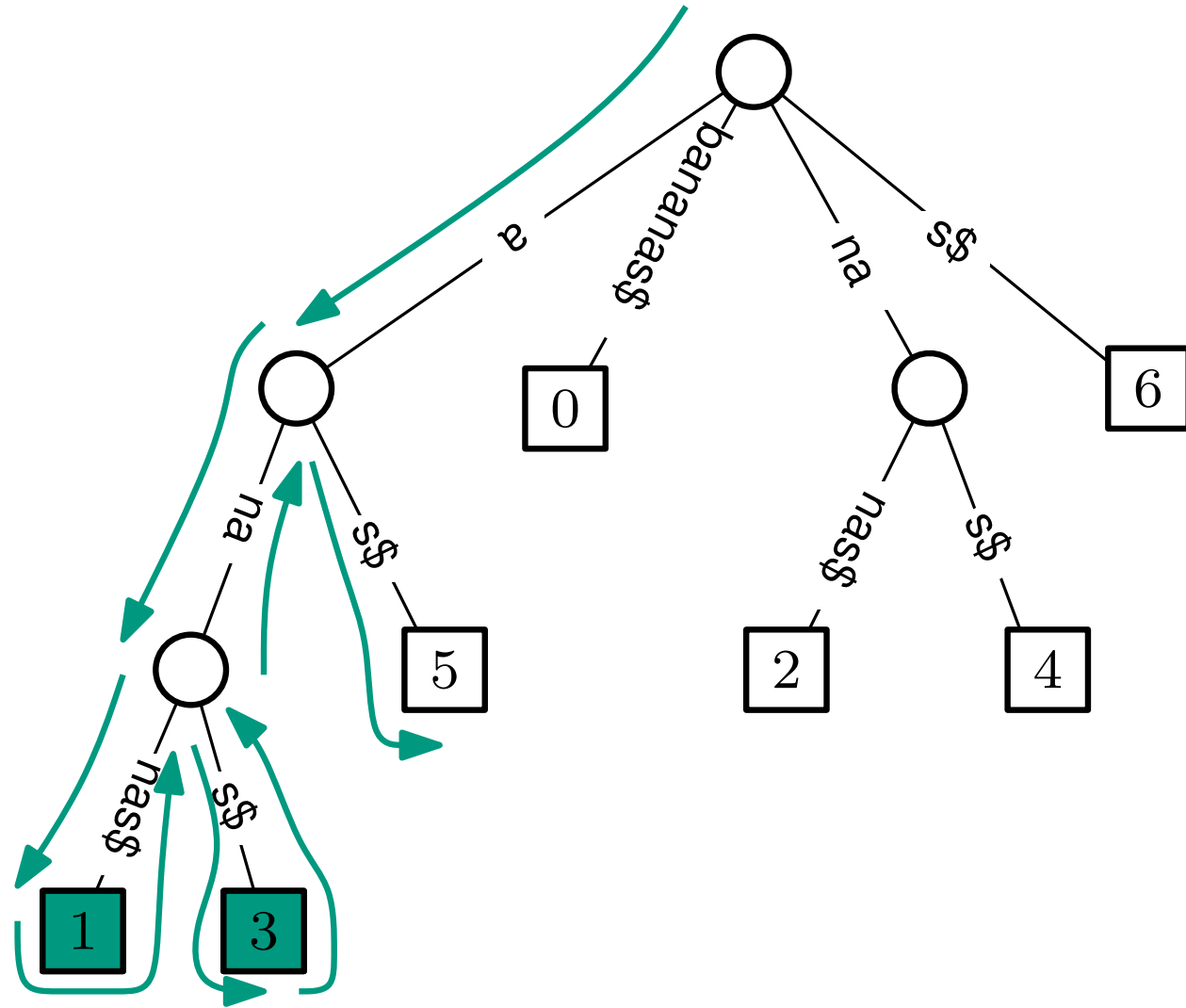
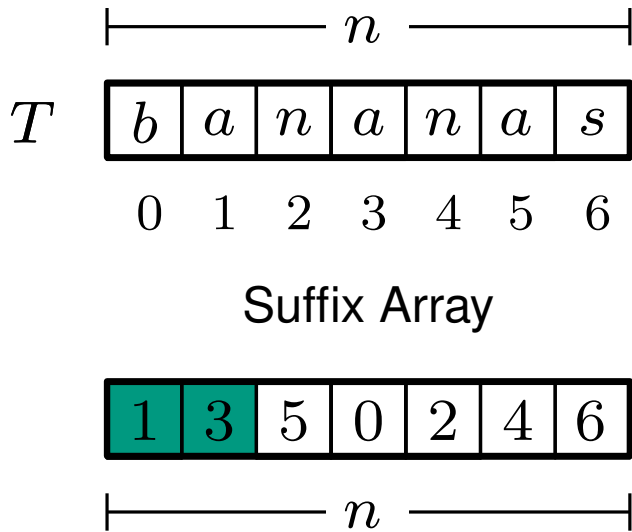
- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)





# Constructing the Suffix Array from the Suffix Tree



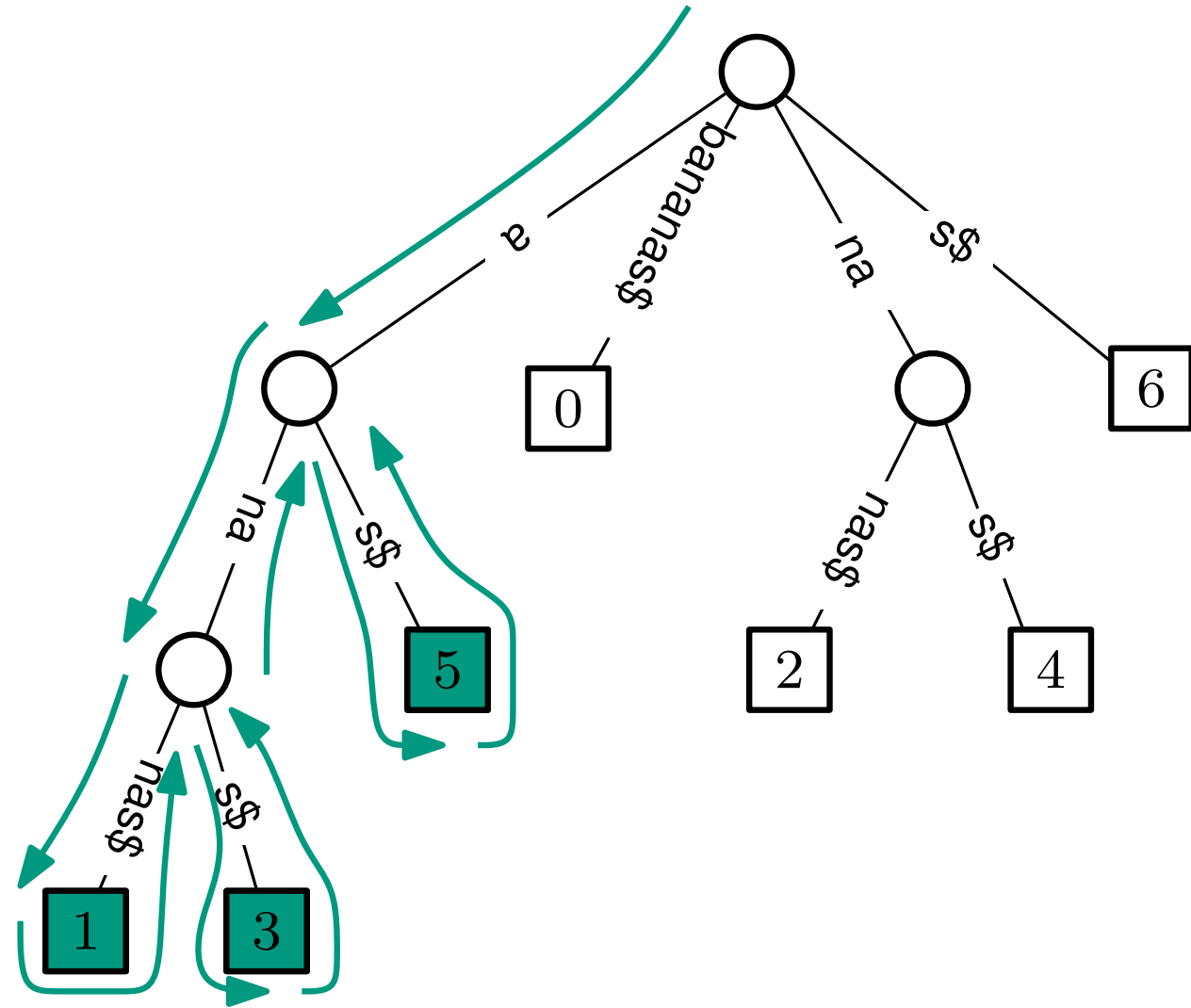
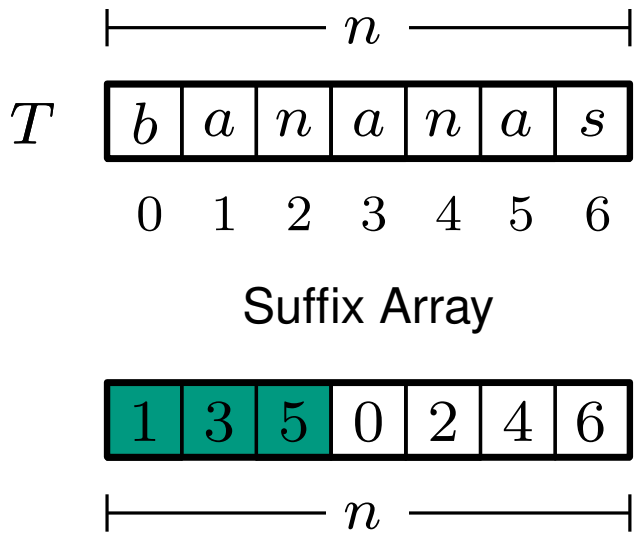
recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)



# Constructing the Suffix Array from the Suffix Tree



recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

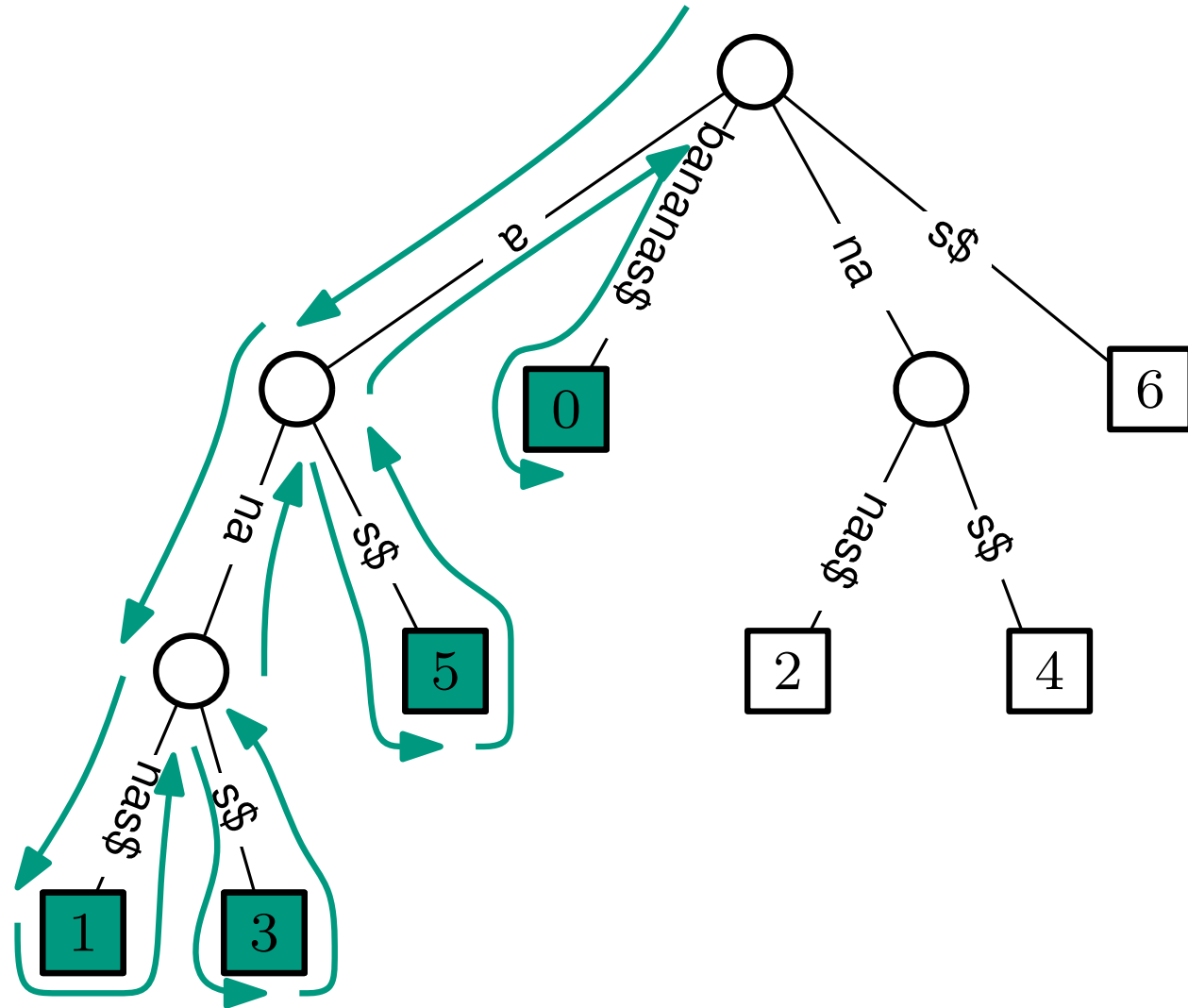
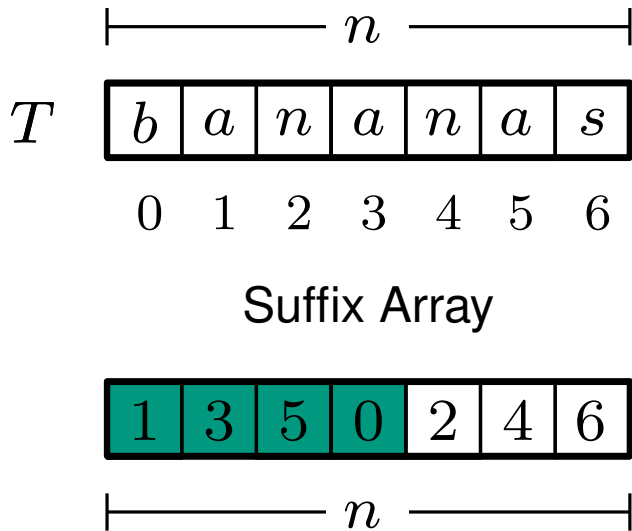
To get the Suffix array perform a depth-first search (in lexicographical order)







# Constructing the Suffix Array from the Suffix Tree

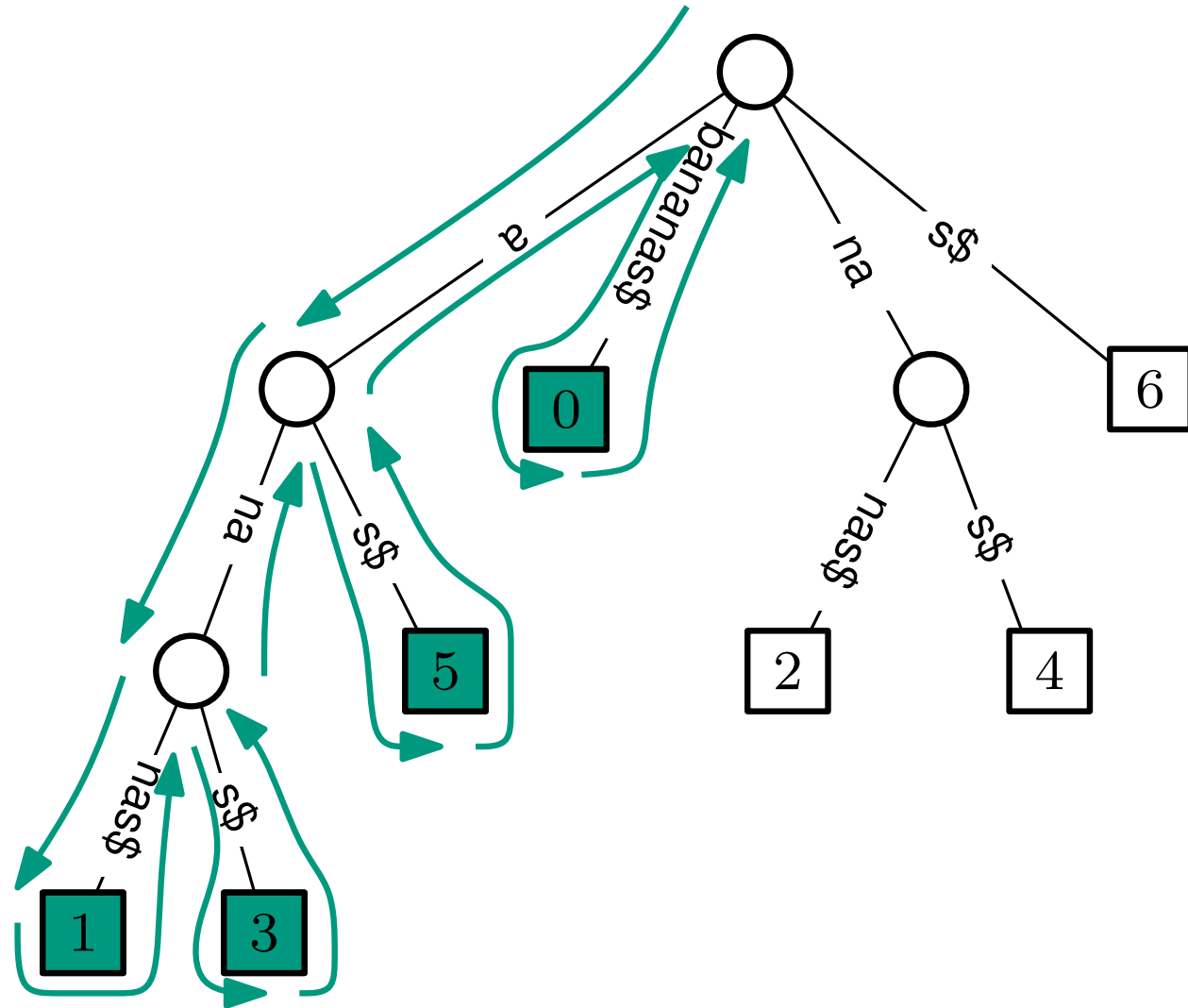
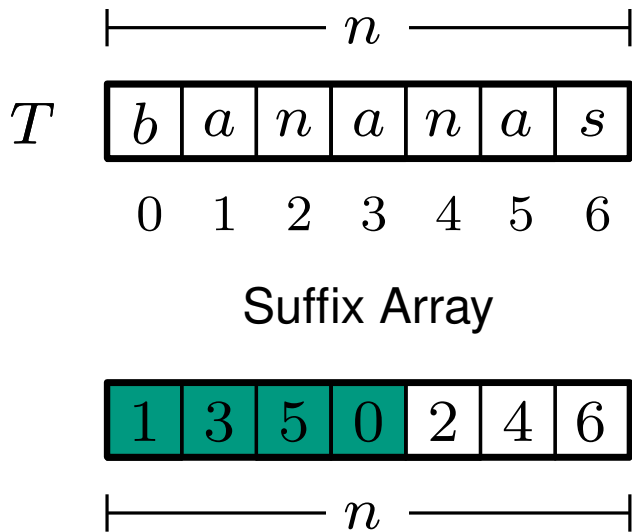


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

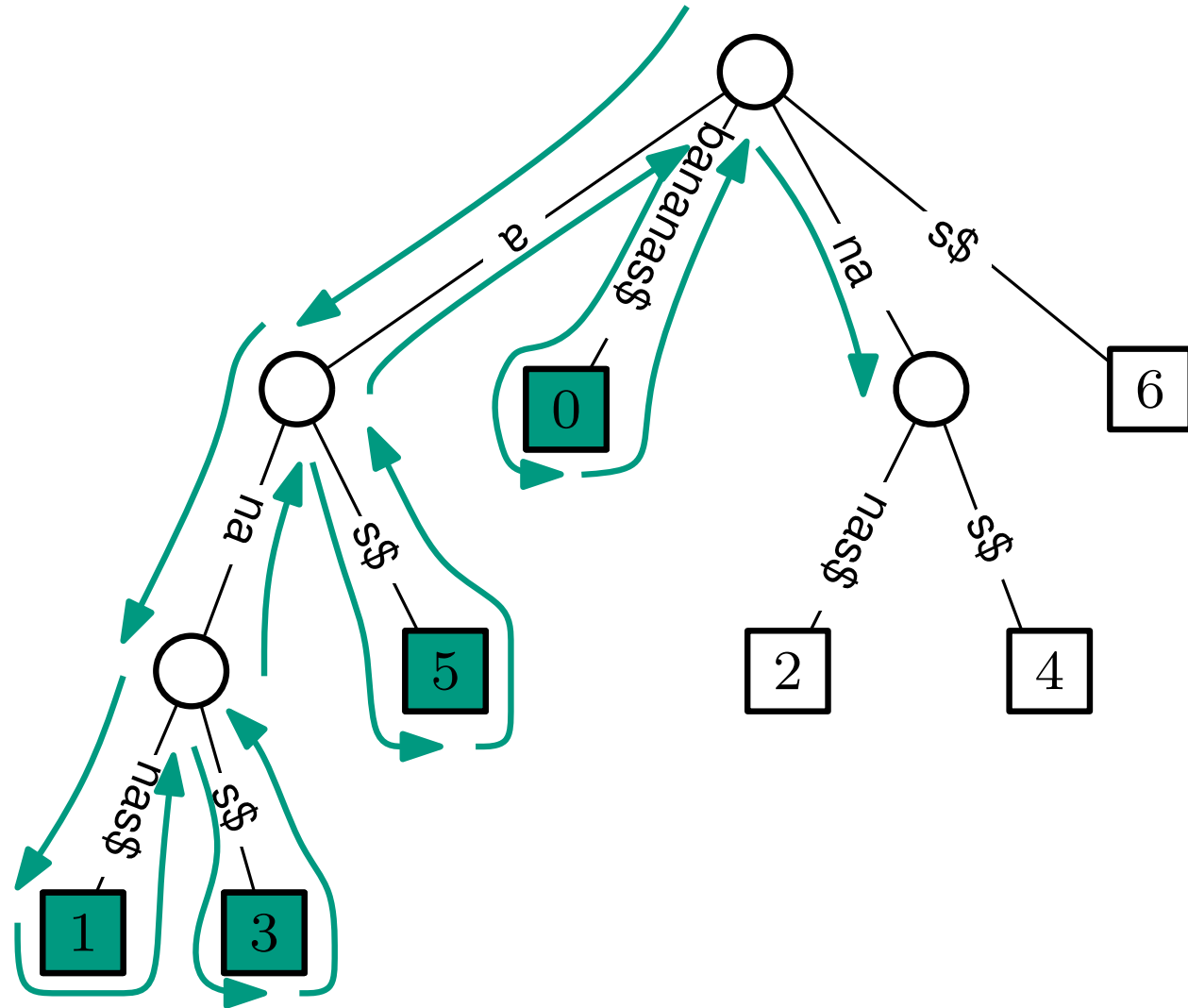
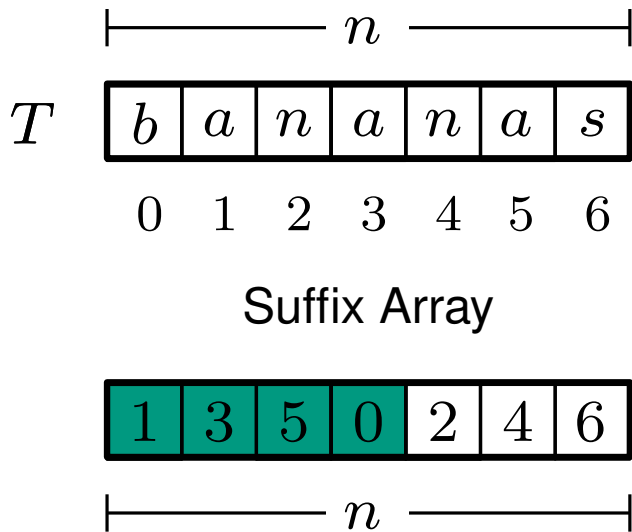


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

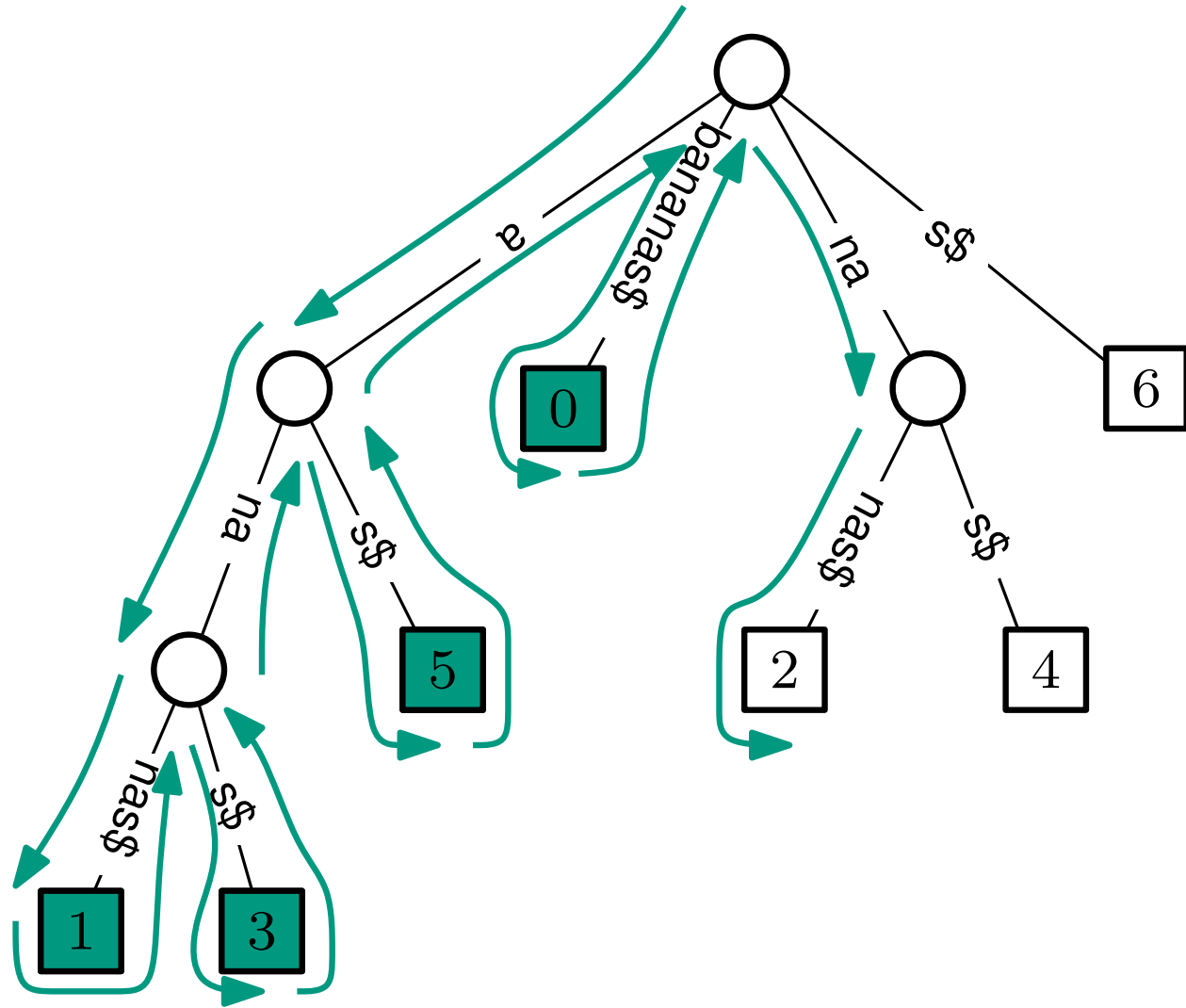
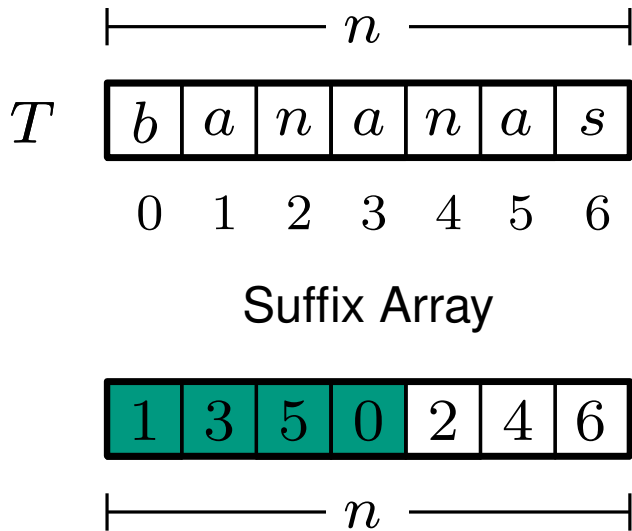


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

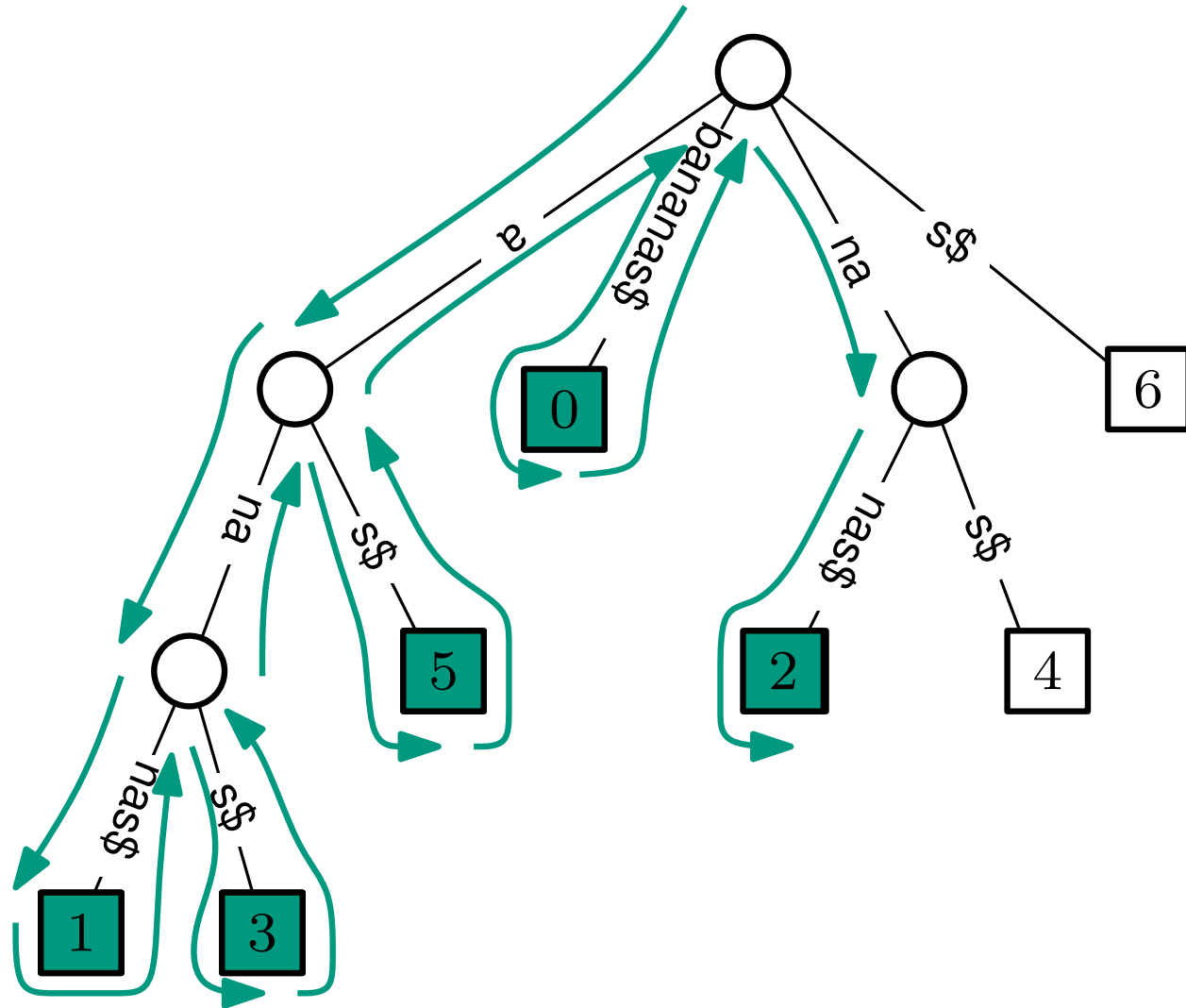
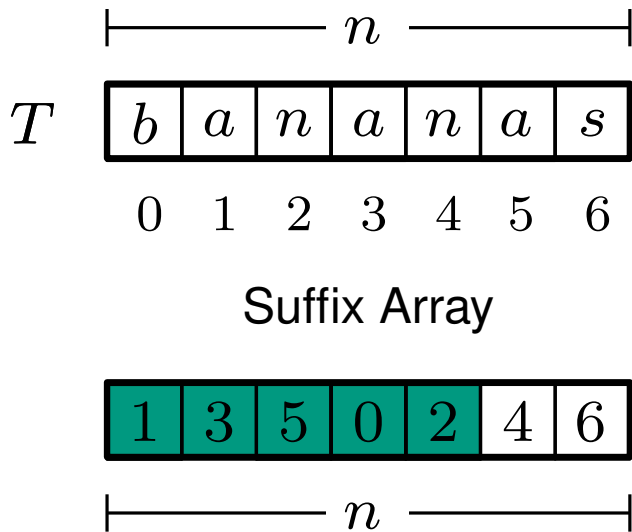


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

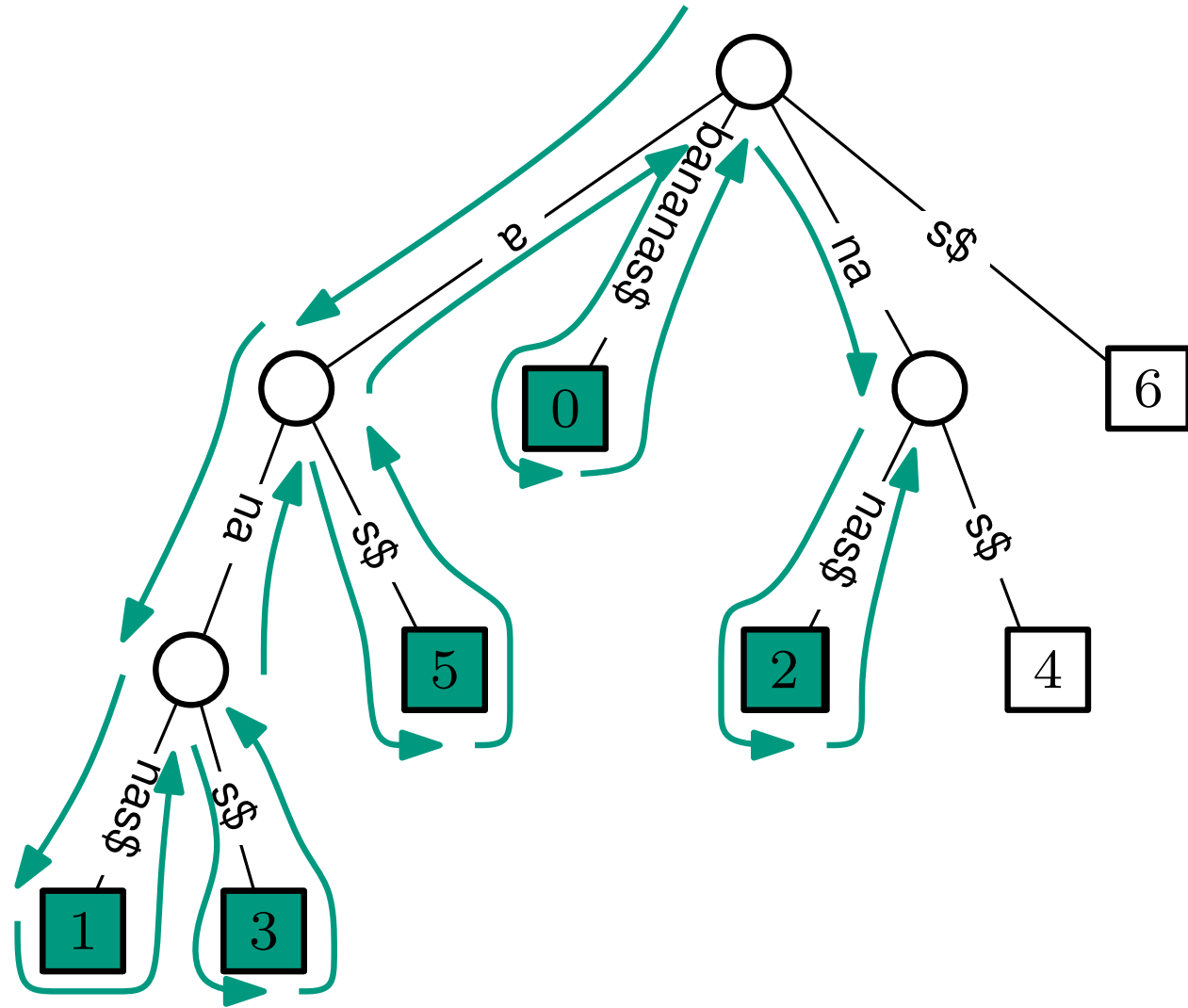
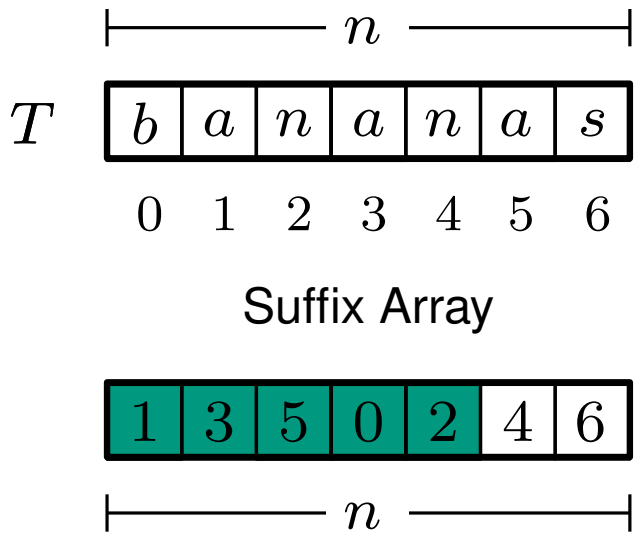


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

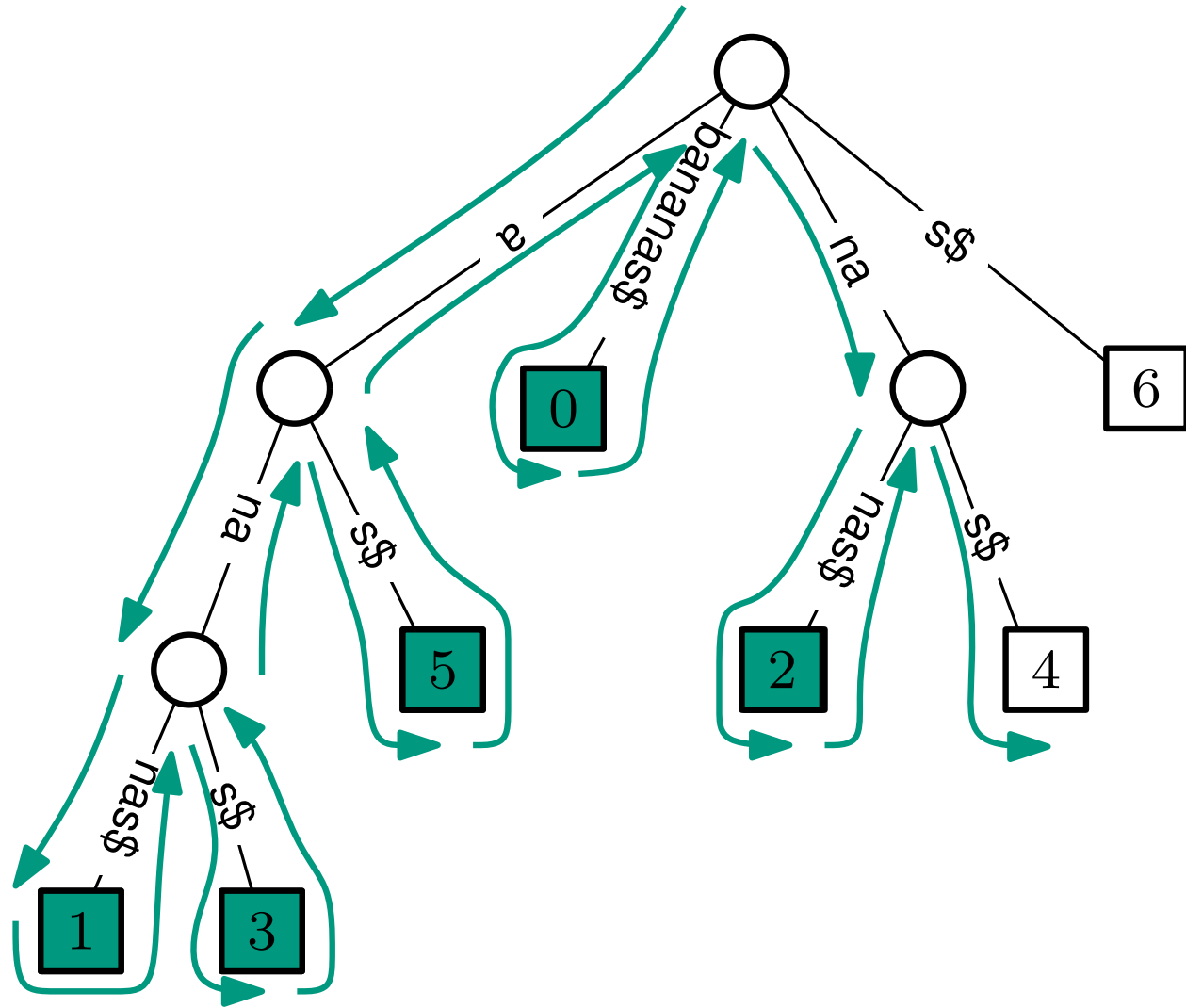
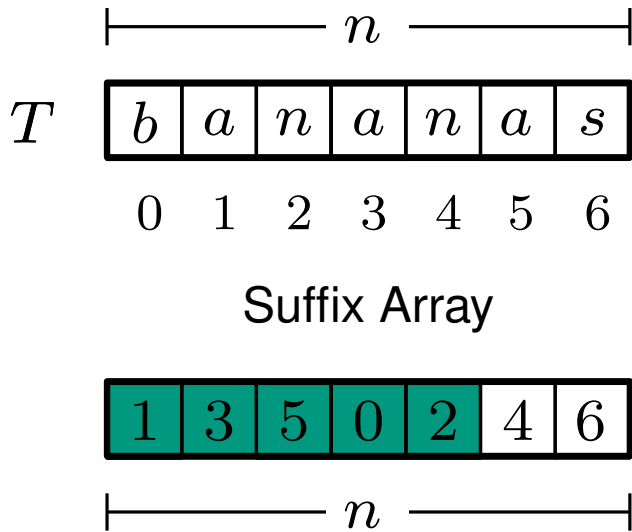


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

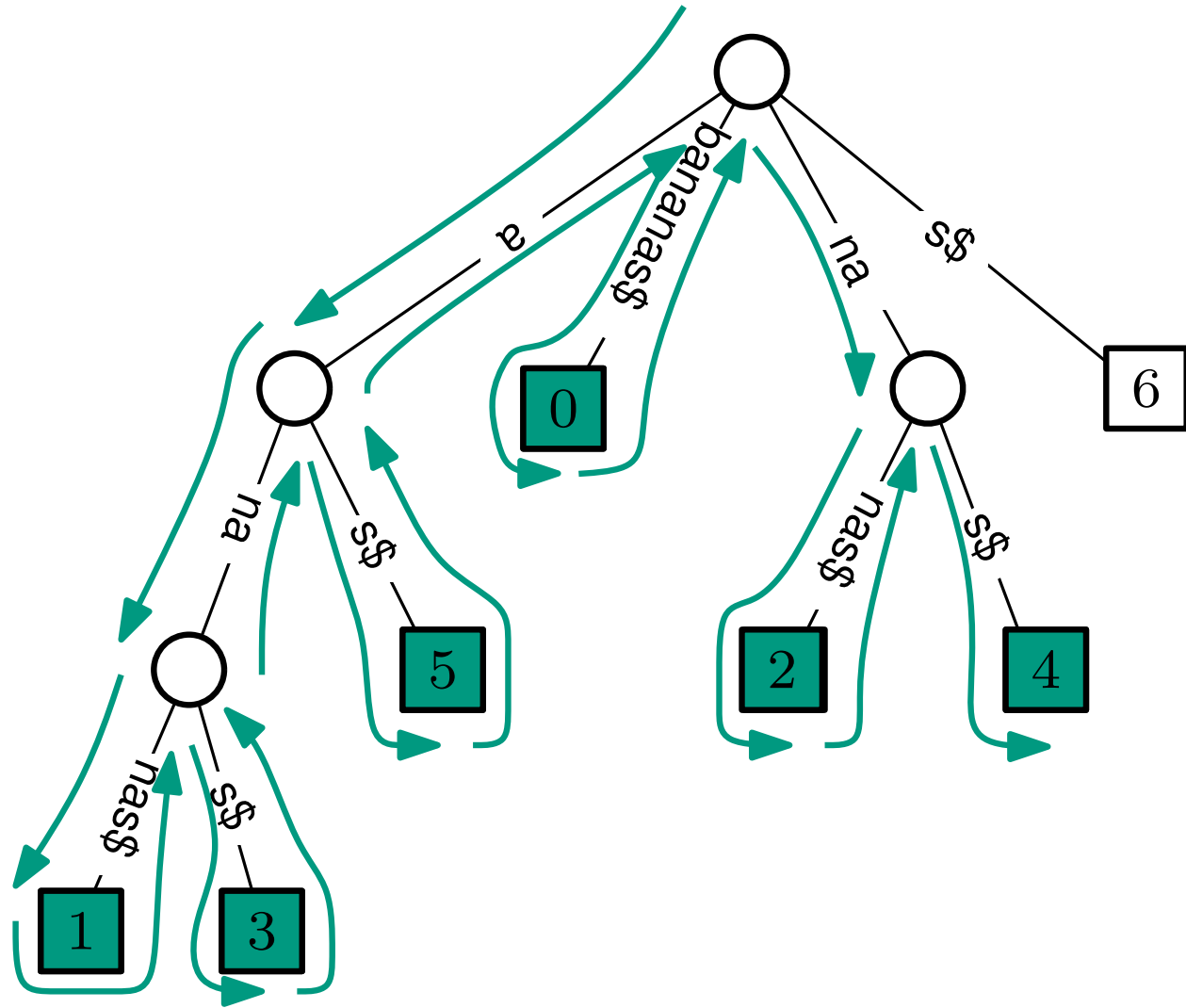
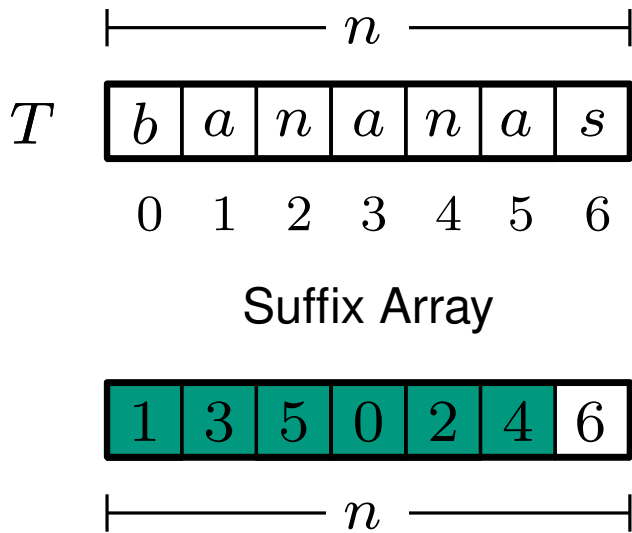


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree



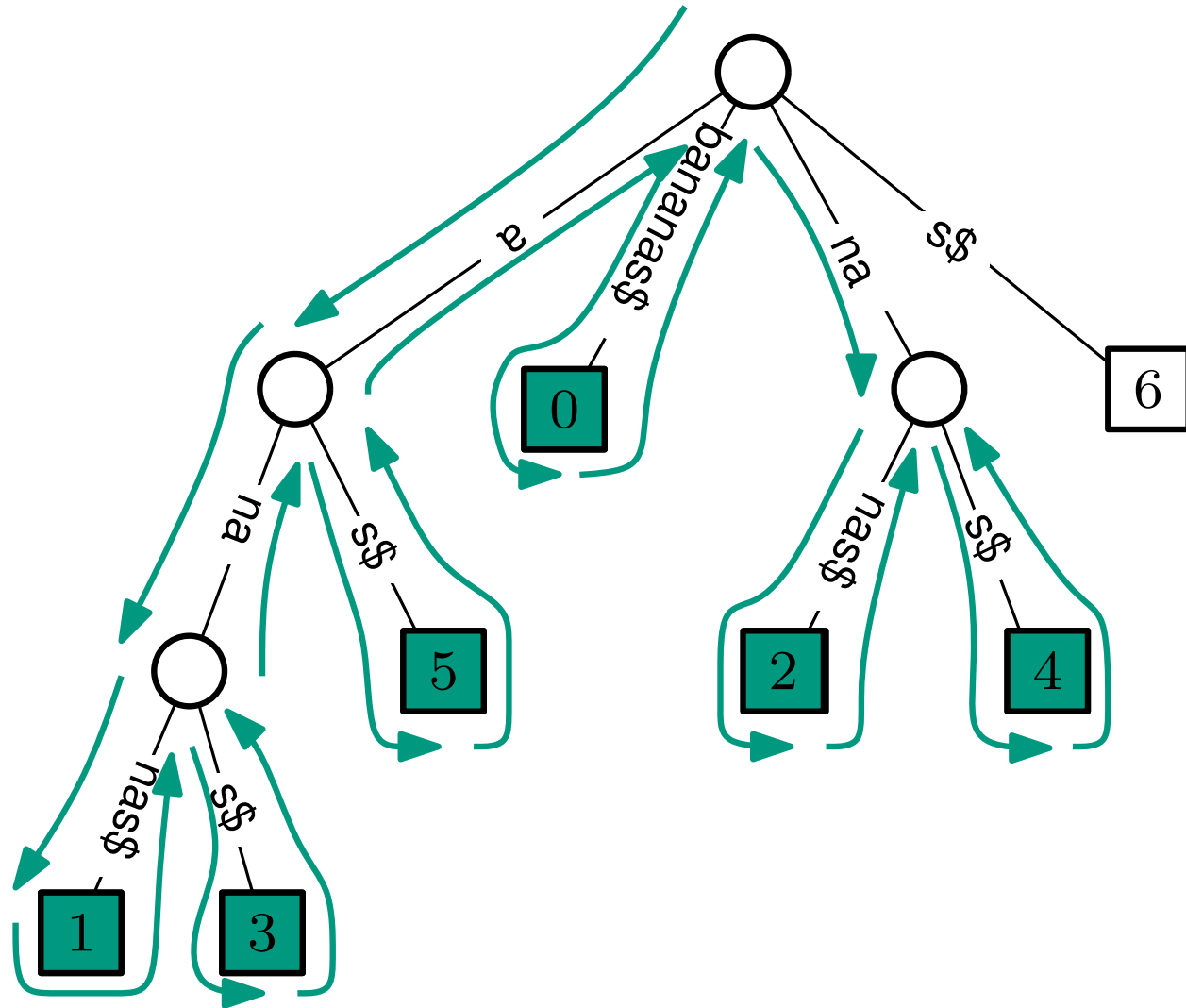
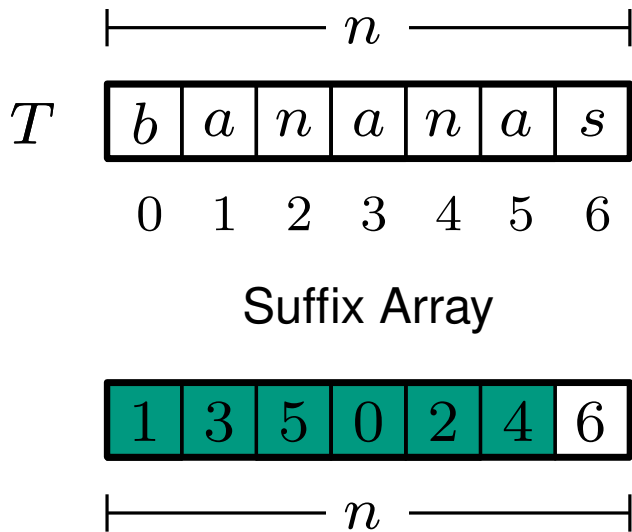
recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)



# Constructing the Suffix Array from the Suffix Tree

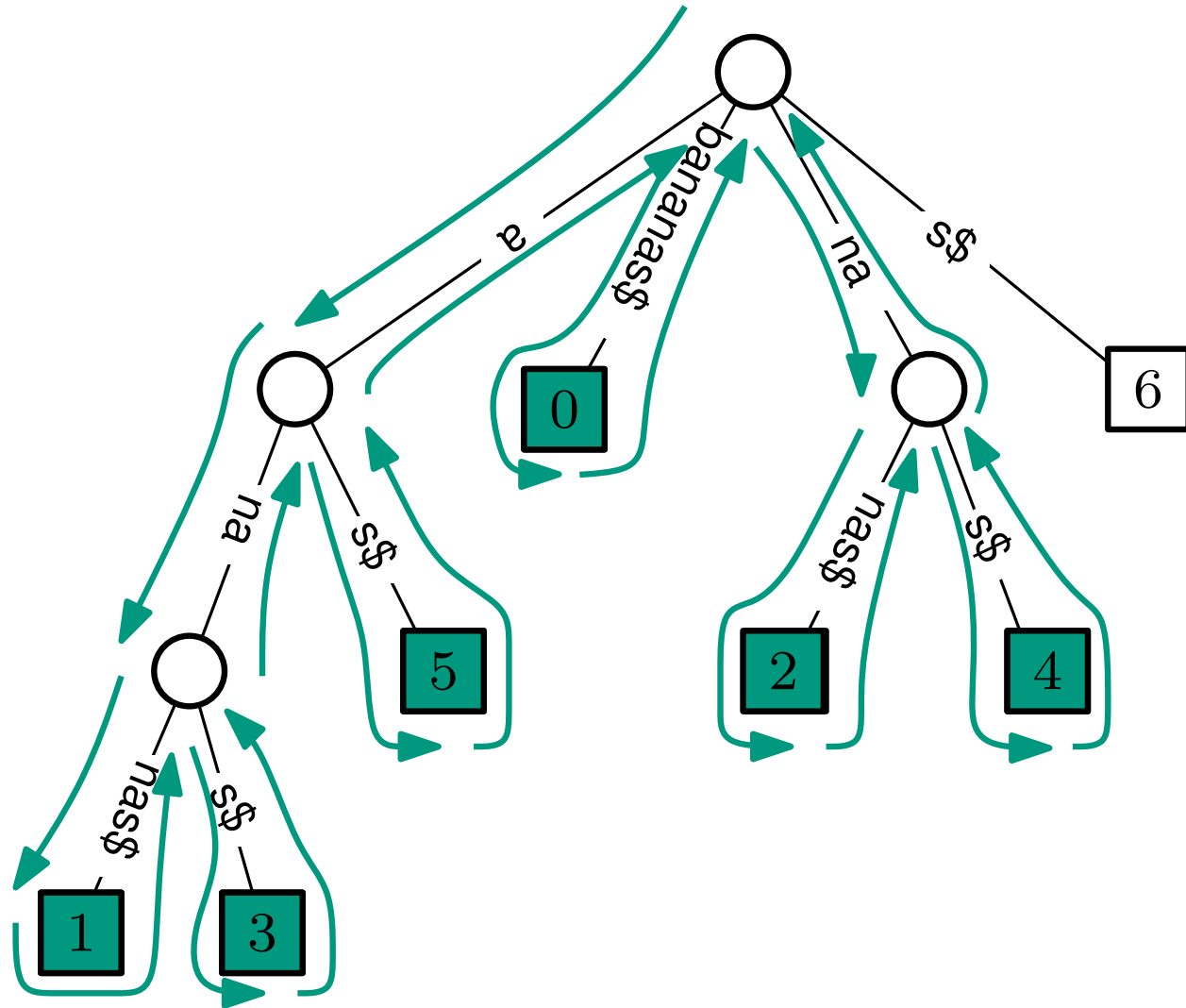
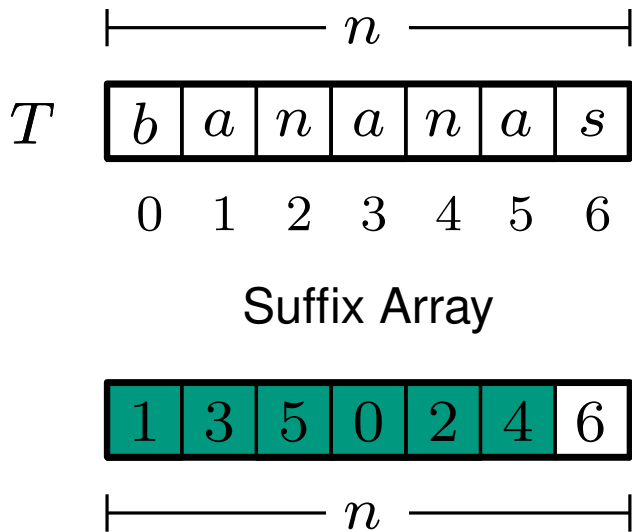


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

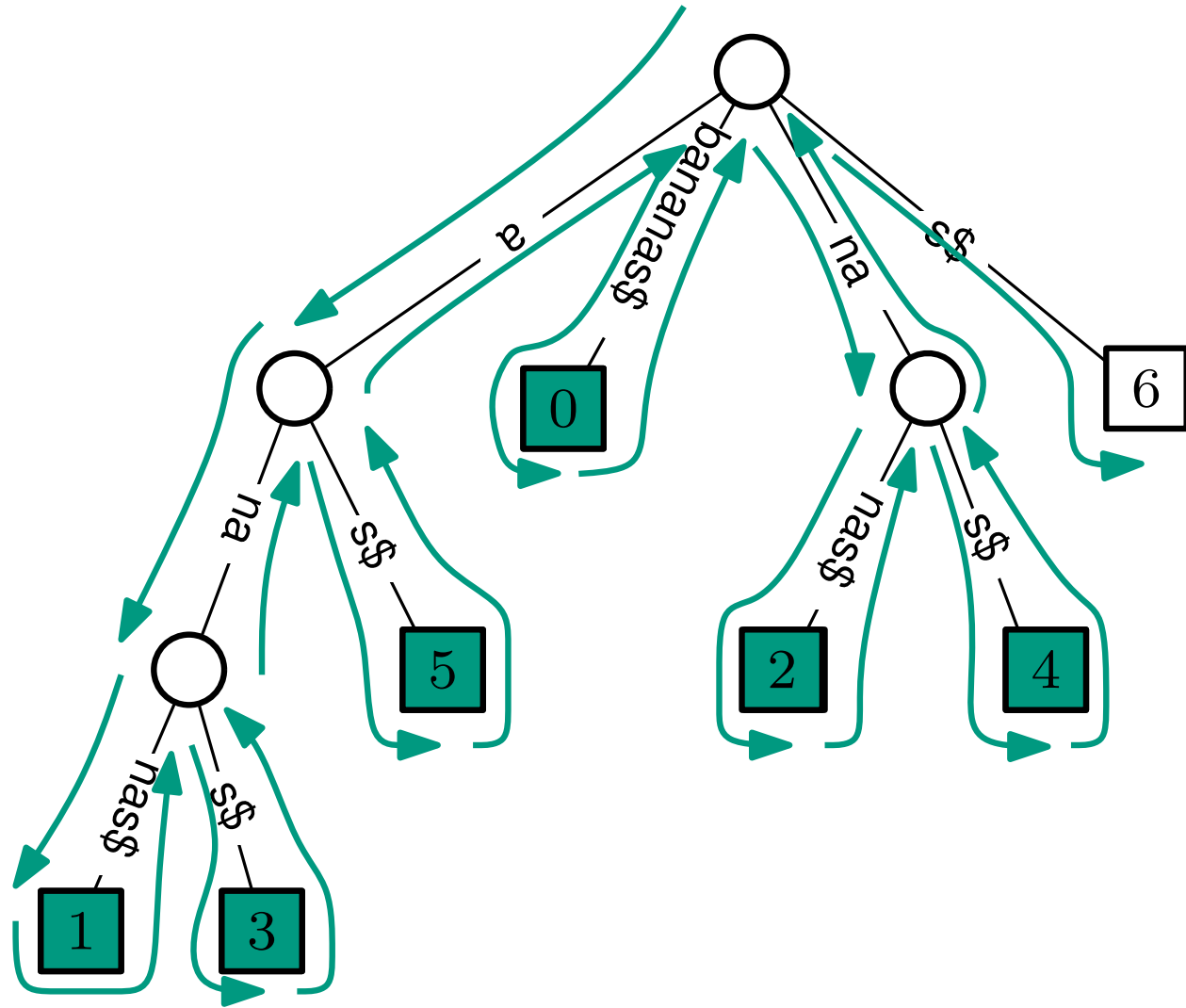
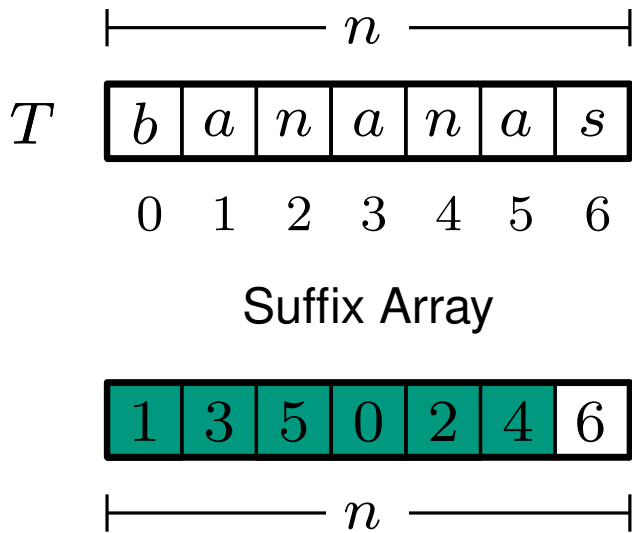


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

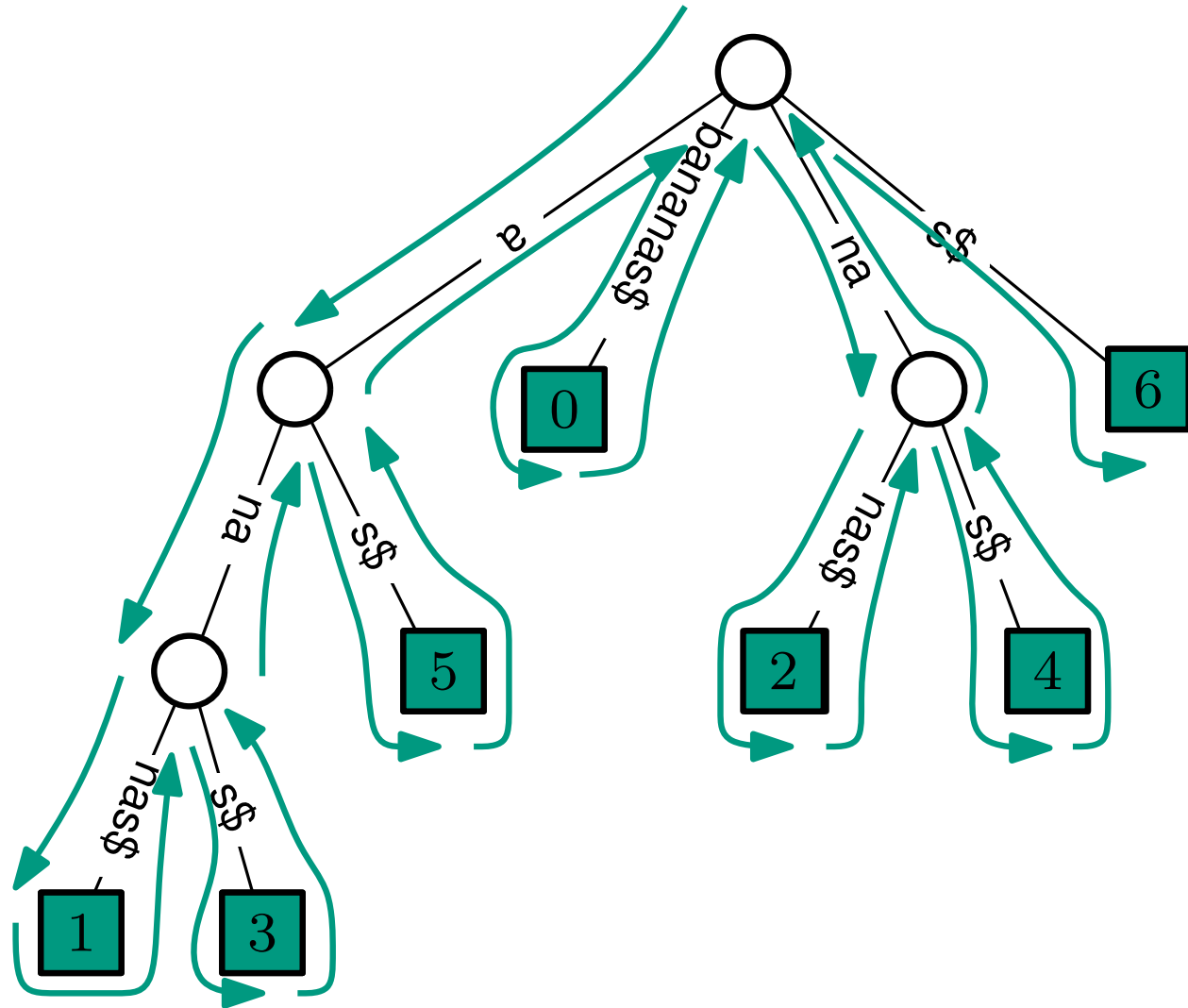
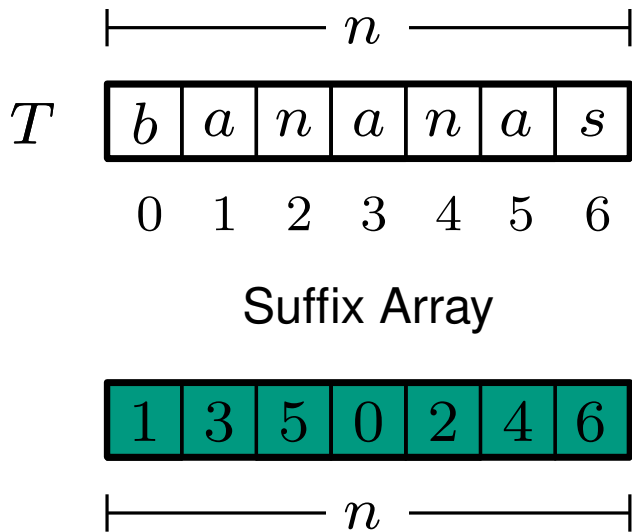


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree

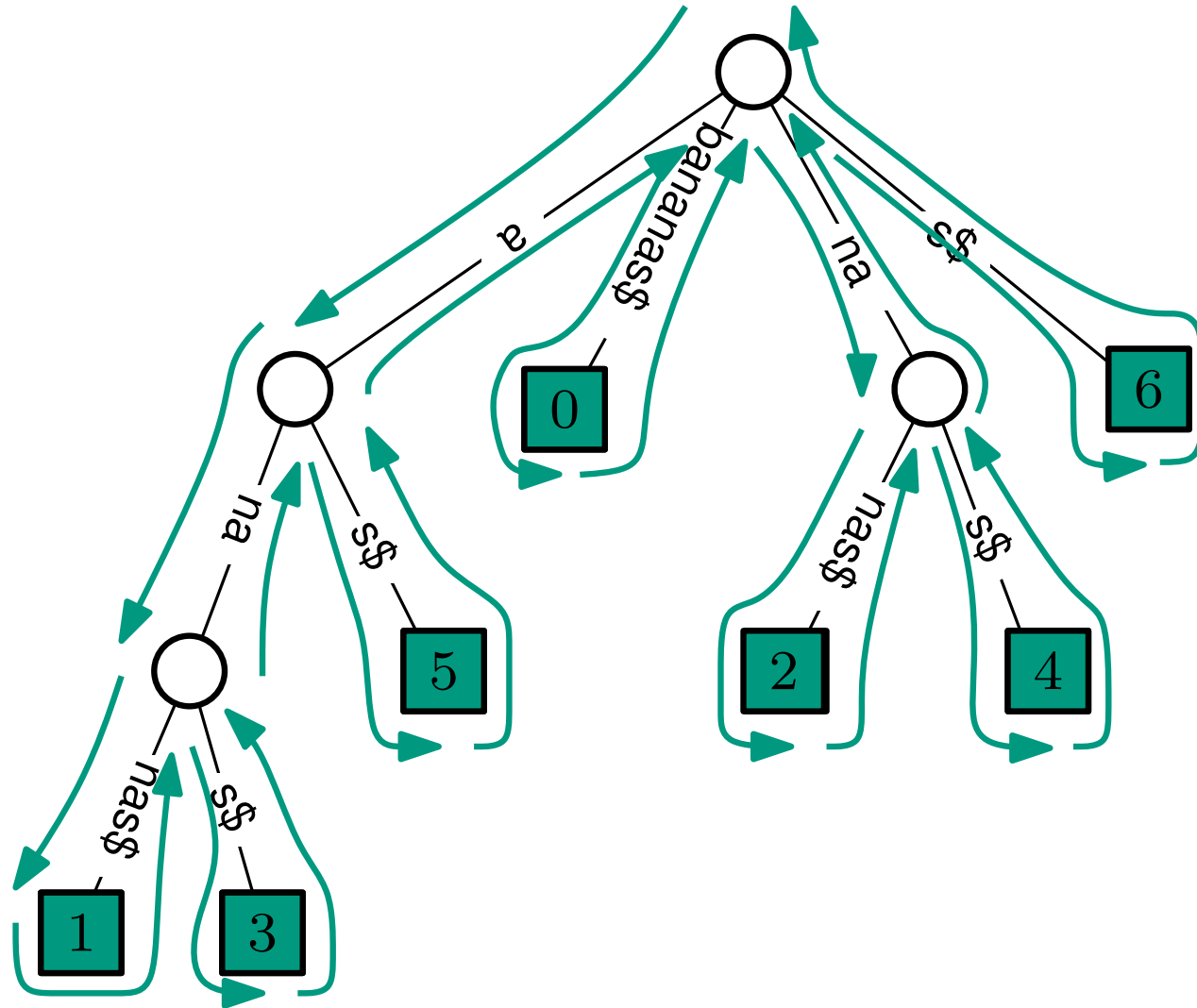
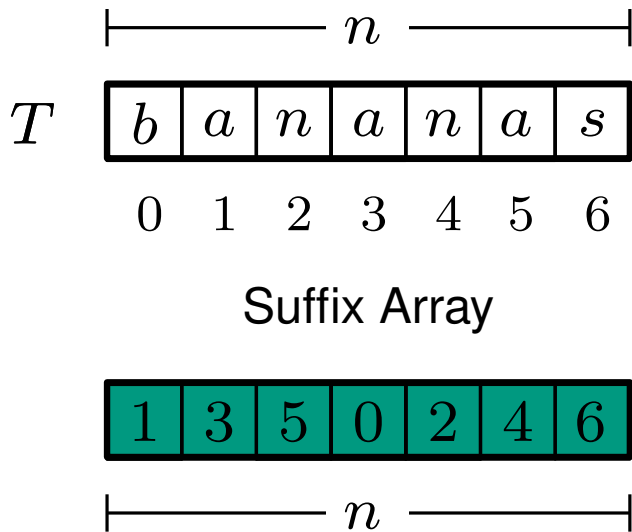


recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)

# Constructing the Suffix Array from the Suffix Tree



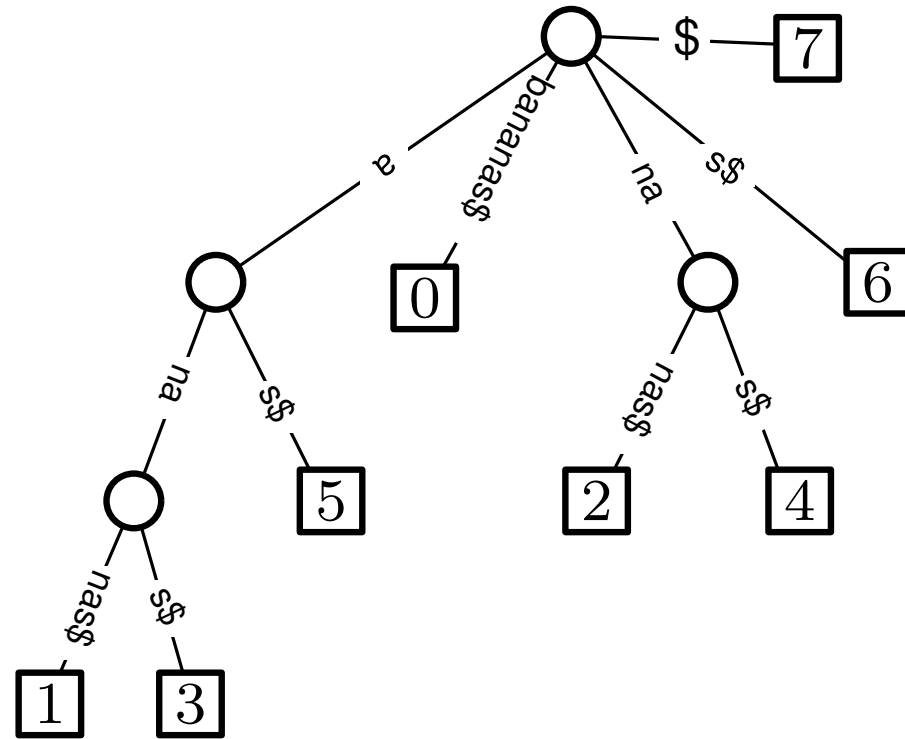
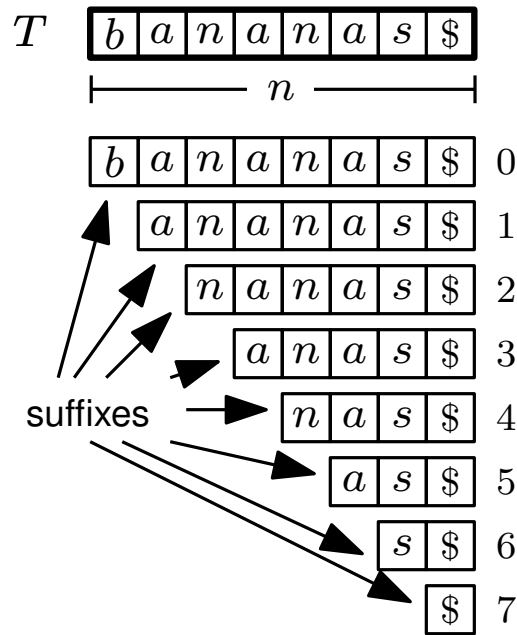
recall that we added a unique symbol \$ to make sure the tree exists

- the \$ is the smallest symbol in the alphabet

To get the Suffix array perform a depth-first search (in lexicographical order)



# Suffix tree summary



- The (compacted) suffix tree of a (length  $n$ ) text uses  $O(n)$  space
- Finding all matches of a pattern  $P$  of length  $m$  takes  $O(m + occ)$   
where  $occ$  is the number of matches
- Suffix trees can be built in  $O(n)$  time  
but we have only seen the  $O(n^2)$  time method

we assumed that the alphabet contains a constant number of symbols