

# Tries & Suffix Tries

# Tries

A trie (“try”) is a tree representing a collection of strings (keys):  
the smallest tree such that

- Each edge is labeled with a character  $c \in \Sigma$

- For given node, at most one child edge has label  $c$ , for any  $c \in \Sigma$

- Each key is “spelled out” along some path starting at root

Helpful for implementing a *set* or *map* when the keys are strings

# Tries

Keys: instant, internal, internet

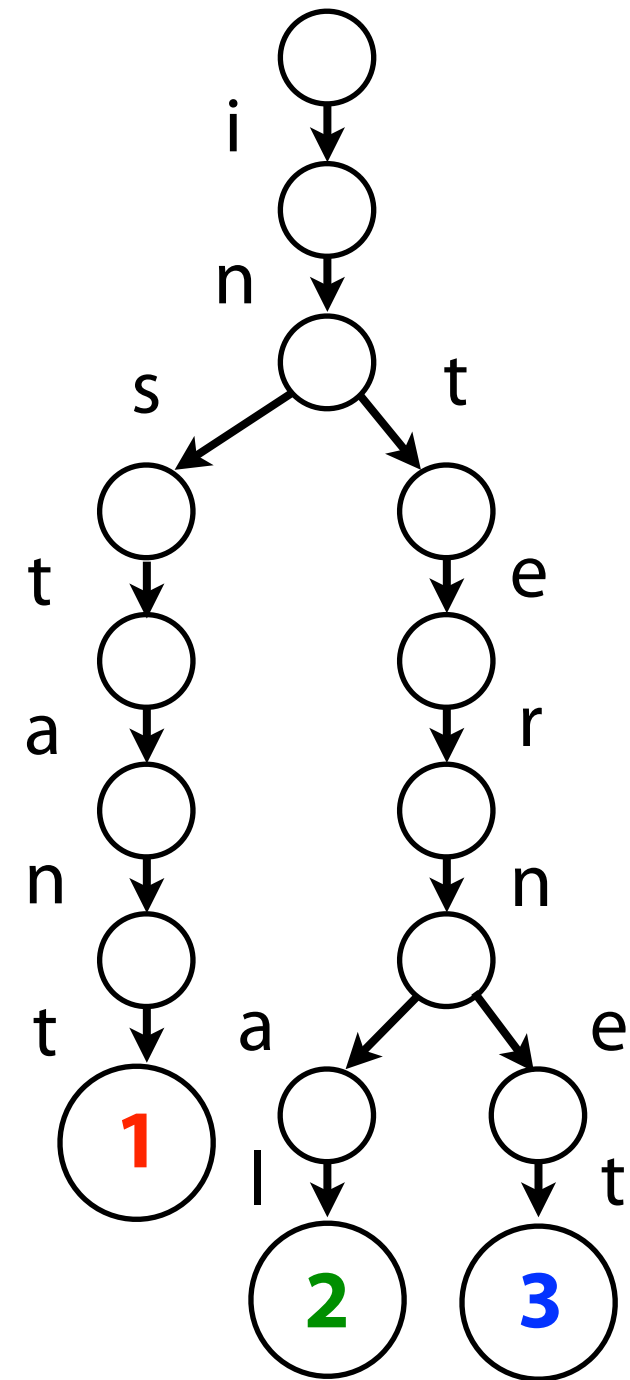
Key	Value
instant	<b>1</b>
internal	<b>2</b>
internet	<b>3</b>

Smallest tree such that:

Each edge is labeled with a character  $c \in \Sigma$

For given node, at most one child edge has label  $c$ , for any  $c \in \Sigma$

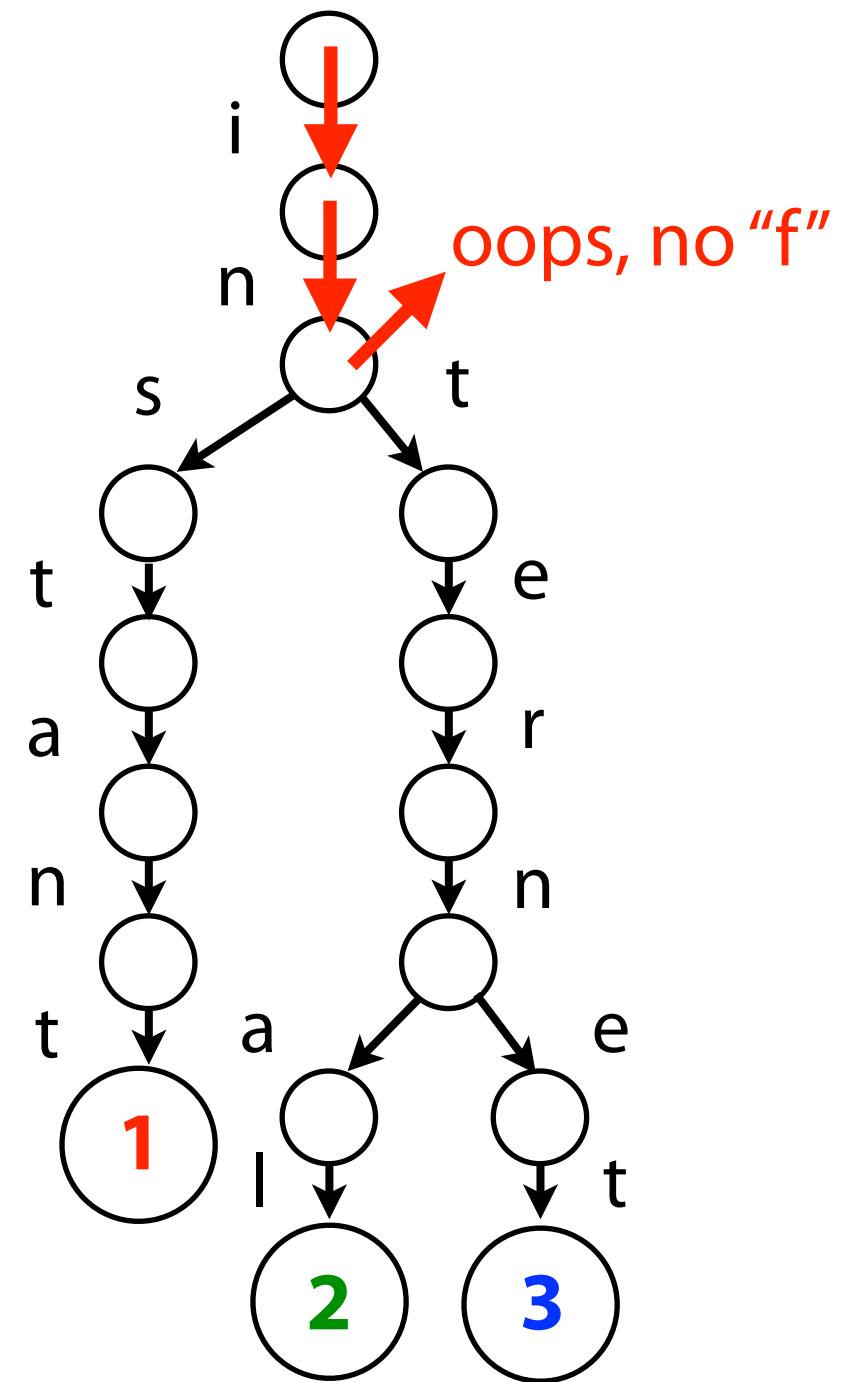
Each key is “spelled out” along some path starting at root



# Tries

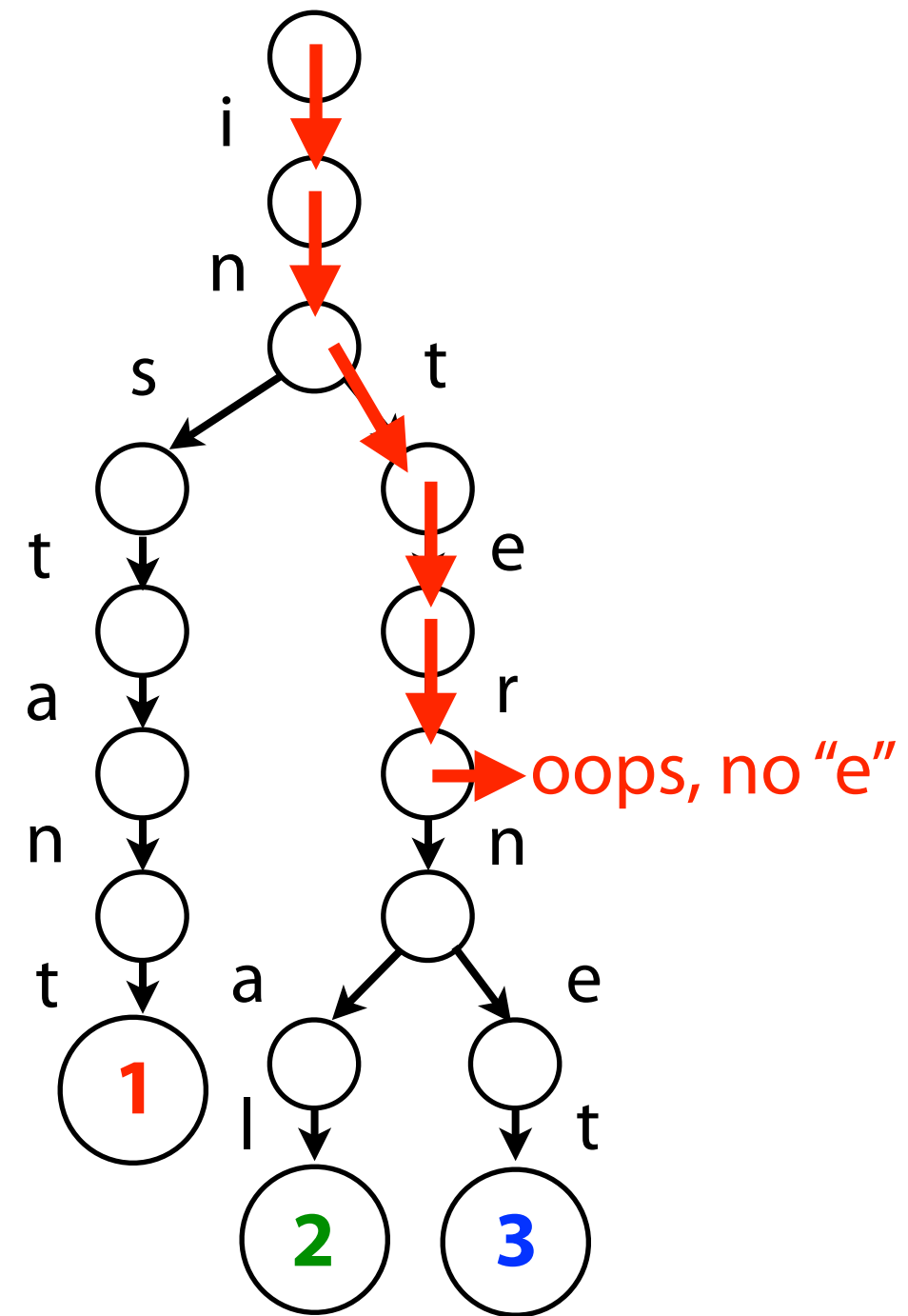
How do we check whether “infer” is in the trie?

Start at root and try to match successive characters of “infer” to edges in trie



# Tries

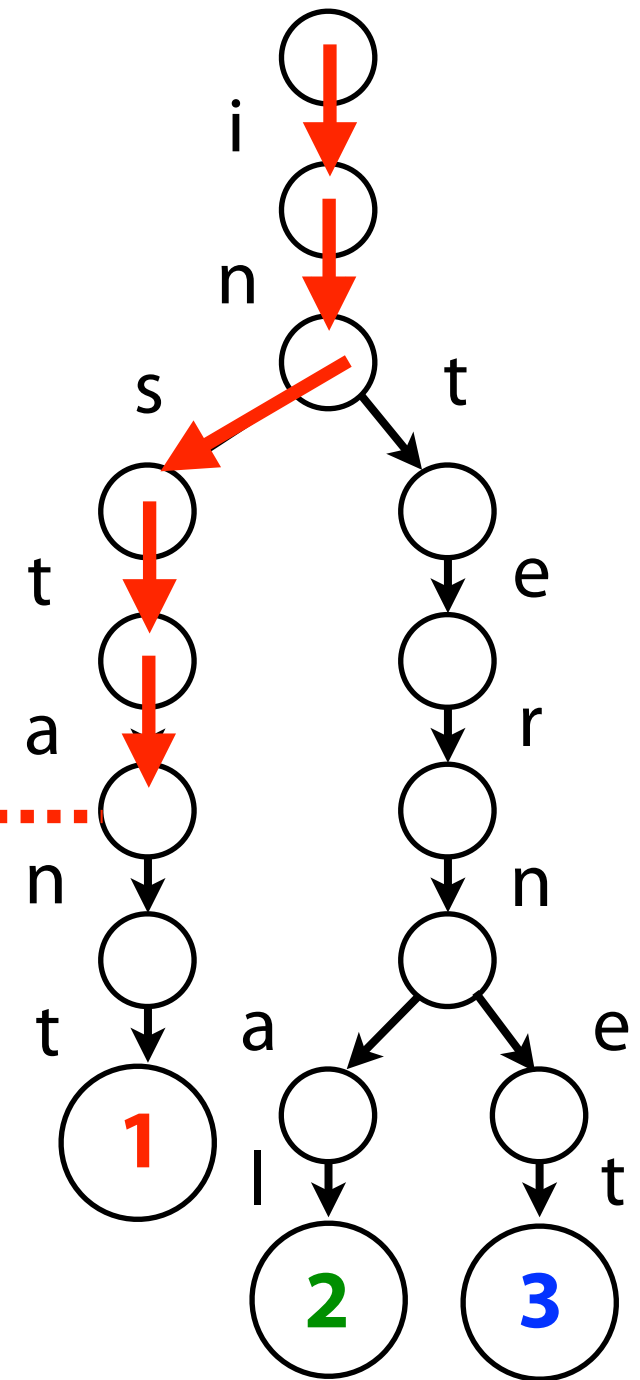
Matching “interesting”



# Tries

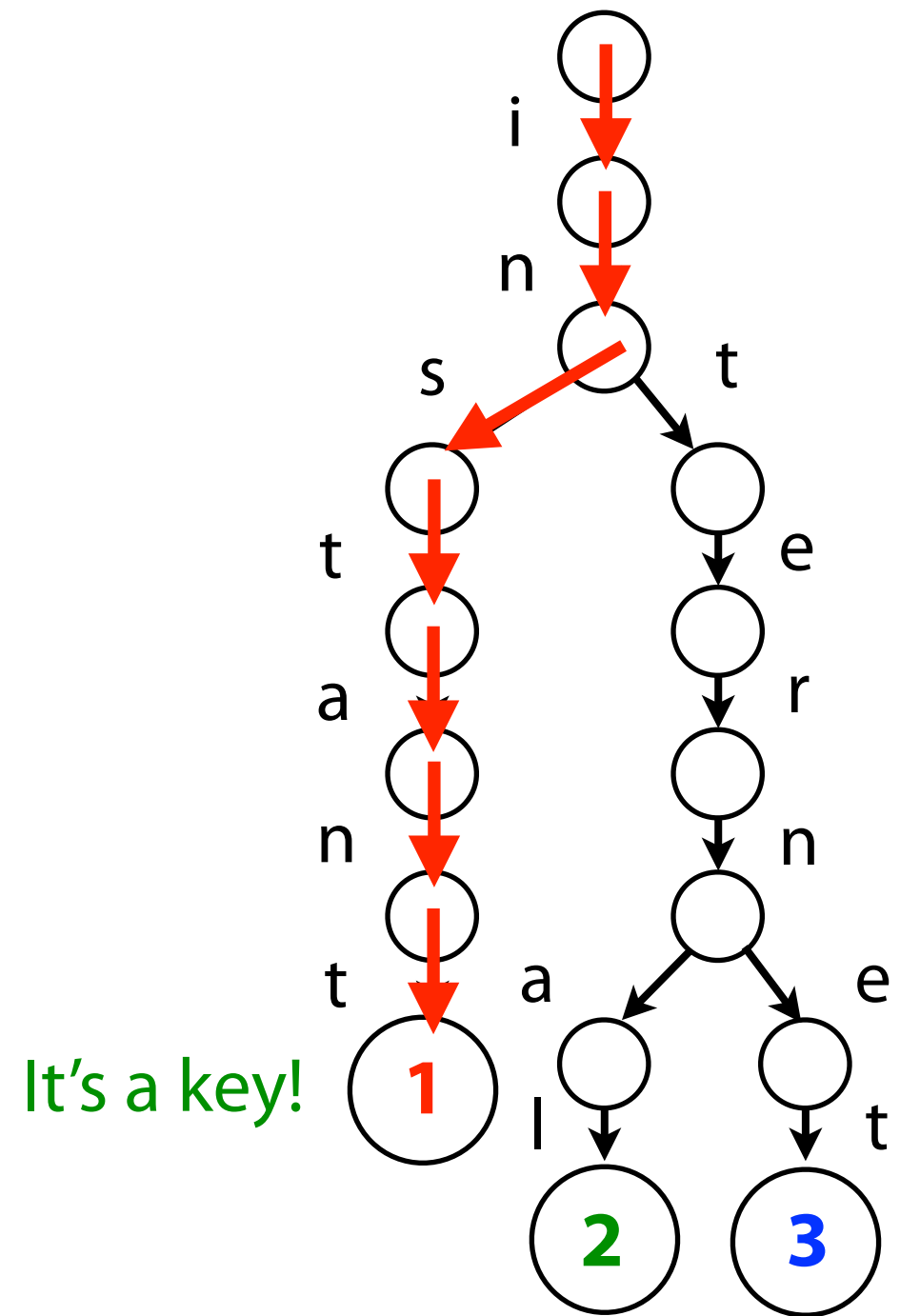
# Matching “insta”

No value associated  
with node, so “insta”  
wasn’t a key

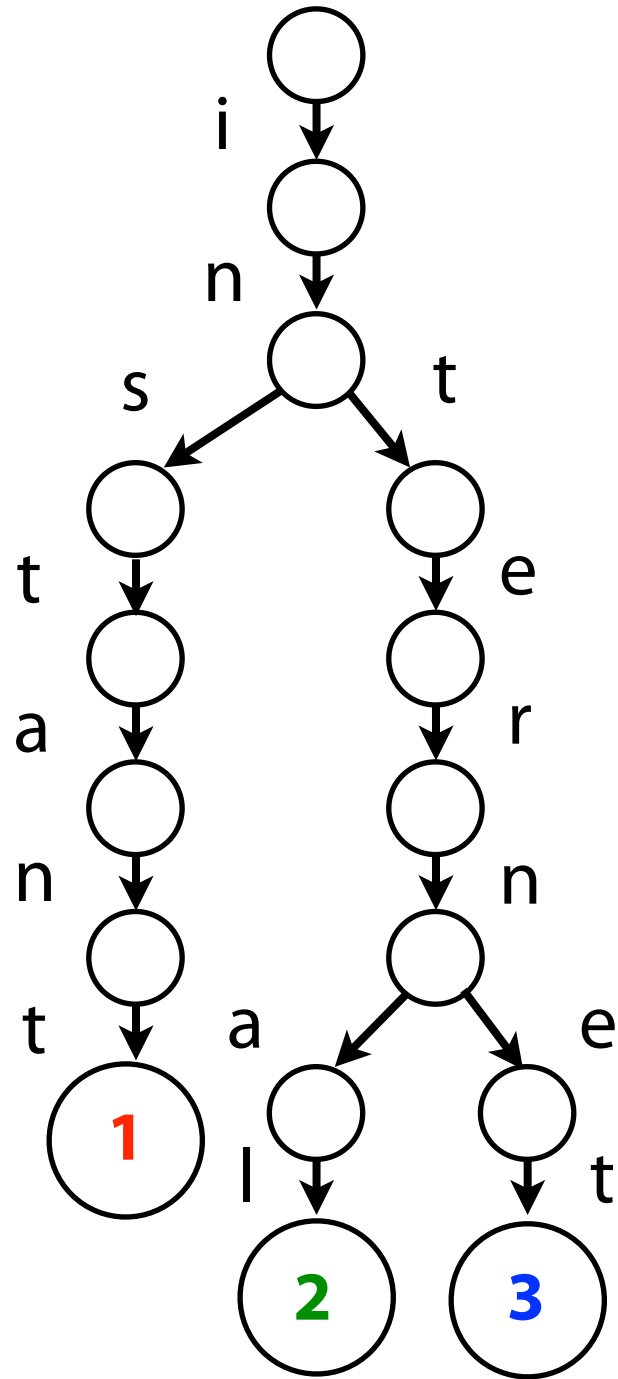


# Tries

Matching “instant”



# Tries

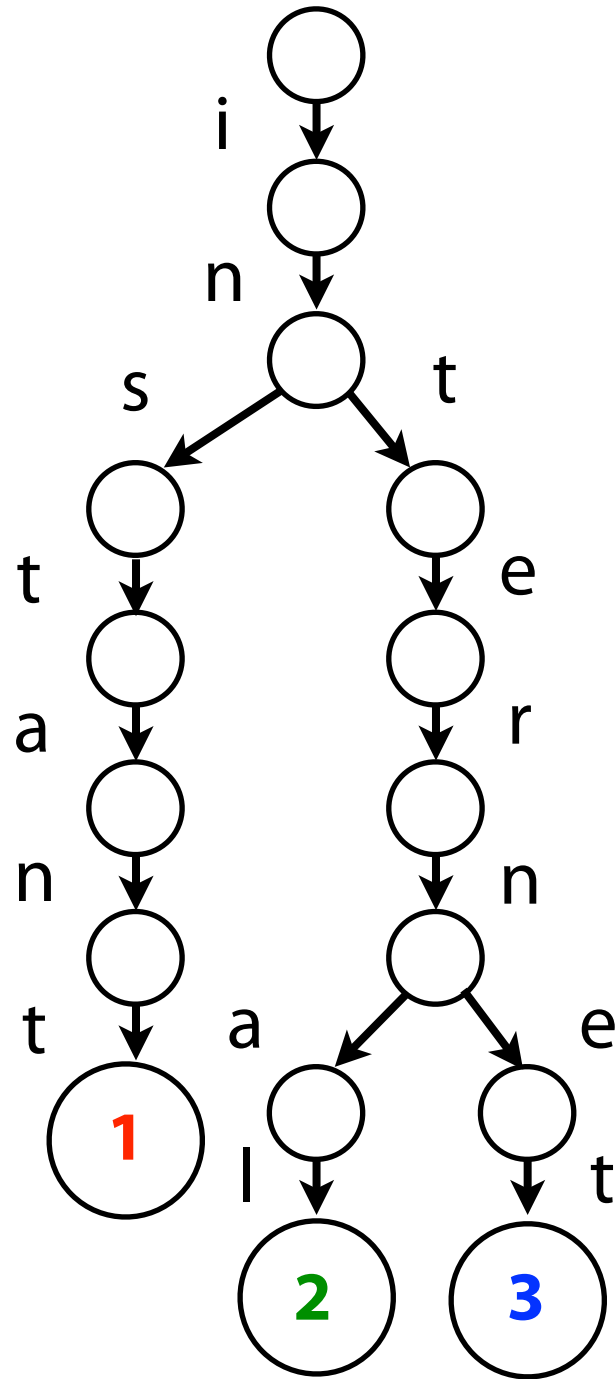


Checking for presence of key  $P$ ,  
where  $|P| = n$  traverses  $\leq n$  edges

If total length of all keys is  $N$ , trie  
has  $\leq N$  edges



# Tries



How to represent edges between a node and its children?

*Map* (from characters to child nodes)

Idea 1: Hash table

Idea 2: Sorted lists

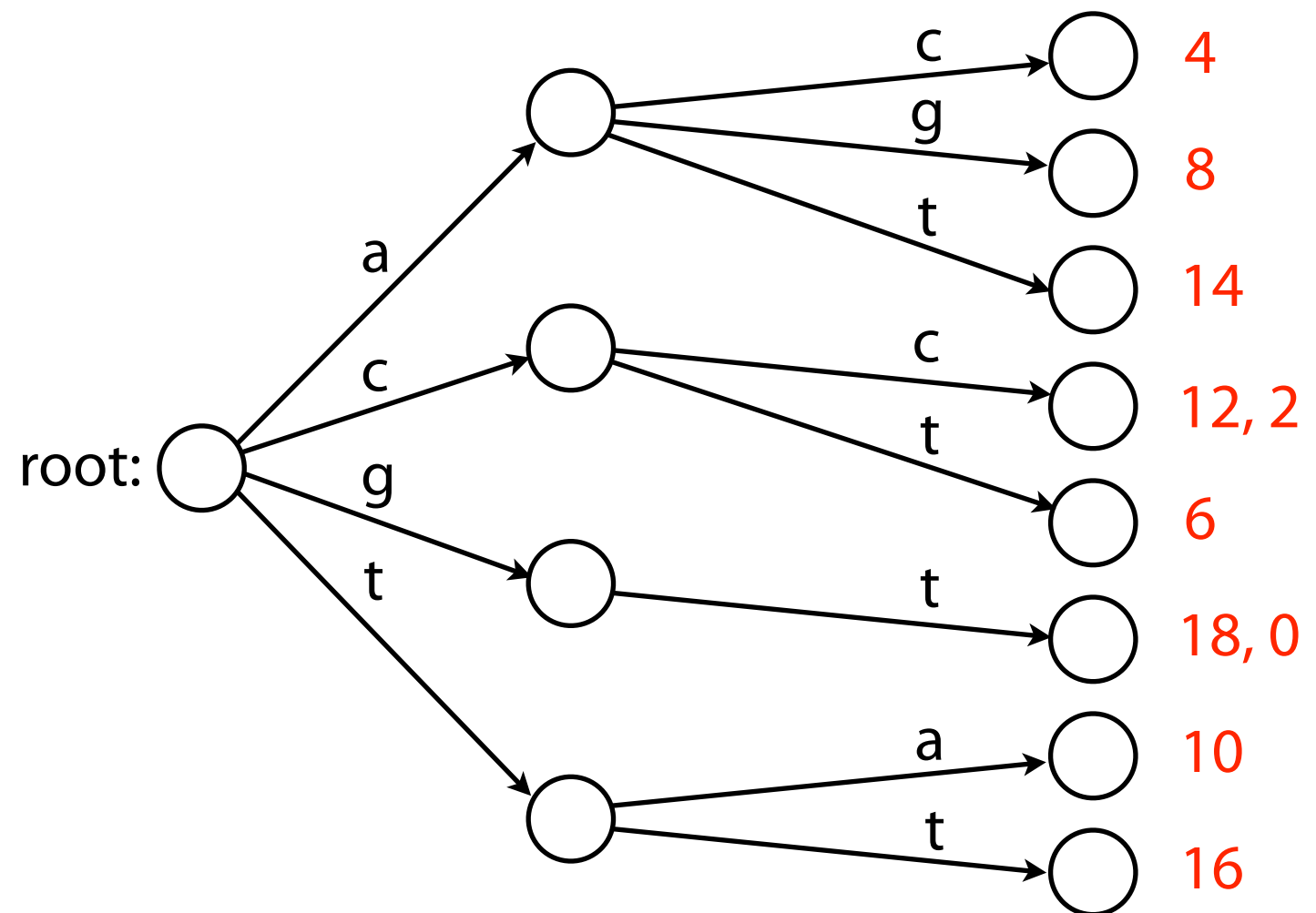
Assuming hash table, it's reasonable to say querying with  $P$ ,  $|P| = n$ , is  $O(n)$  time

# Tries

Could use trie to represent  $k$ -mer index.  
Map  $k$ -mers to offsets where they occur

ac	4
ag	8
at	14
cc	12
cc	2
ct	6
gt	18
gt	0
ta	10
tt	16

Index



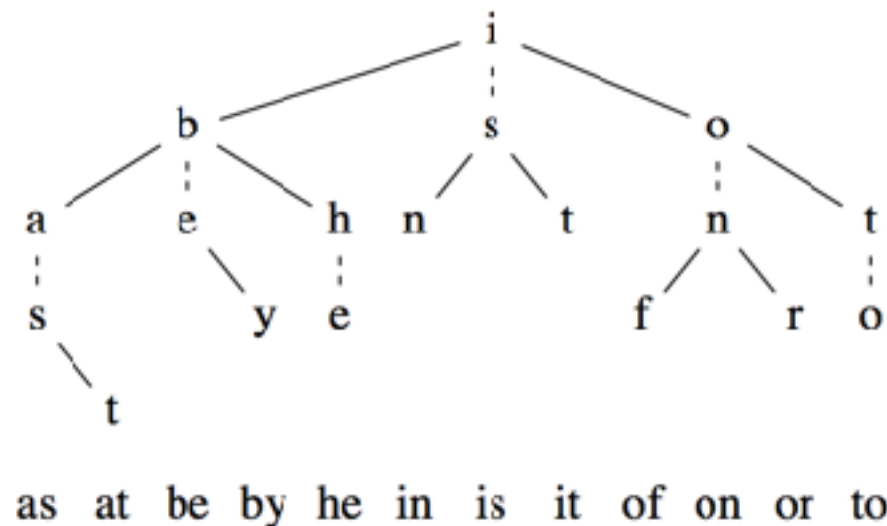
# Tries: implementation

[Refer to trie-map.ipynb](#)

# Tries: alternatives

Tries aren't the only way to encode sets or maps over strings using a tree.

E.g. ternary search tree:



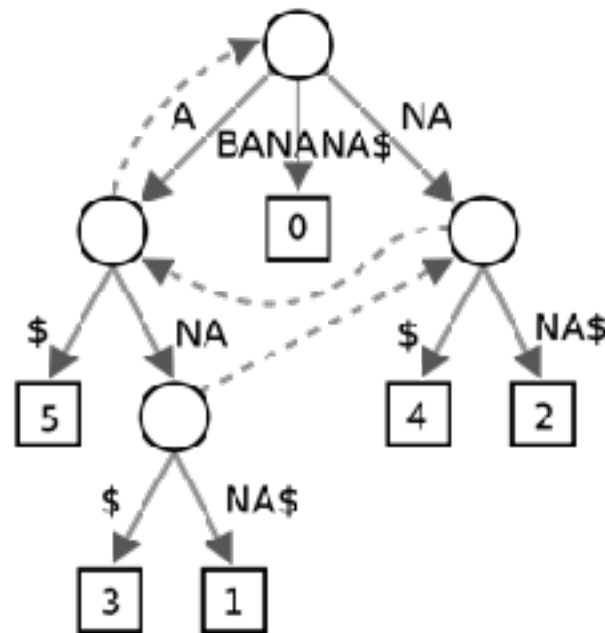
*Figure 2. A ternary search tree for 12 two-letter words*

Bentley, Jon L., and Robert Sedgwick. "Fast algorithms for sorting and searching strings." *Proceedings of the eighth annual ACM-SIAM symposium on Discrete algorithms*. Society for Industrial and Applied Mathematics, 1997

# Indexing with suffixes

We studied indexes built over substrings of  $T$

Different approach is to index *suffixes* of  $T$ . This yields surprisingly economical & practical data structures:



Suffix Tree

6	\$
5	A\$
3	ANA\$
1	ANANA\$
0	BANANA\$
4	NA\$
2	NANA\$

Suffix Array

\$	B	A	N	A	N	A
A	\$	B	A	N	A	N
A	N	A	\$	B	A	N
A	N	A	N	A	\$	B
B	A	N	A	N	A	\$
N	A	\$	B	A	N	A
N	A	N	A	\$	B	A

FM Index

# Suffix trie

Build a **trie** containing all **suffixes** of a text  $T$

$T$ : GTTAGCTGATCGCGGCGTAGCGG\$

GTTAGCTGATCGCGGCGTAGCGG\$

TTAGCTGATCGCGGCGTAGCGG\$

TATAGCTGATCGCGGCGTAGCGG\$

ATAGCTGATCGCGGCGTAGCGG\$

TAGCTGATCGCGGCGTAGCGG\$

AGCTGATCGCGGCGTAGCGG\$

GCTGATCGCGGCGTAGCGG\$

CTGATCGCGGCGTAGCGG\$

TGATCGCGGCGTAGCGG\$

GATCGCGGCGTAGCGG\$

ATCGCGGCGTAGCGG\$

TCGCGGCGTAGCGG\$

CGCGGCGTAGCGG\$

GCGGCGTAGCGG\$

CGGCGTAGCGG\$

GGCGTAGCGG\$

GCGTAGCGG\$

CGTAGCGG\$

GTAGCGG\$

TAGCGG\$

AGCGG\$

GCGG\$

CGG\$

GG\$

G\$

\$

$m(m+1)/2$   
chars

# Suffix trie

First add special *terminal character* **\$** to the end of  $T$

**\$** is a character that does not appear elsewhere in  $T$ , and we define it to be less than other characters (**\$** < **A** < **C** < **G** < **T**)

**\$** enforces a familiar rule: e.g. “as” comes before “ash” in the dictionary.

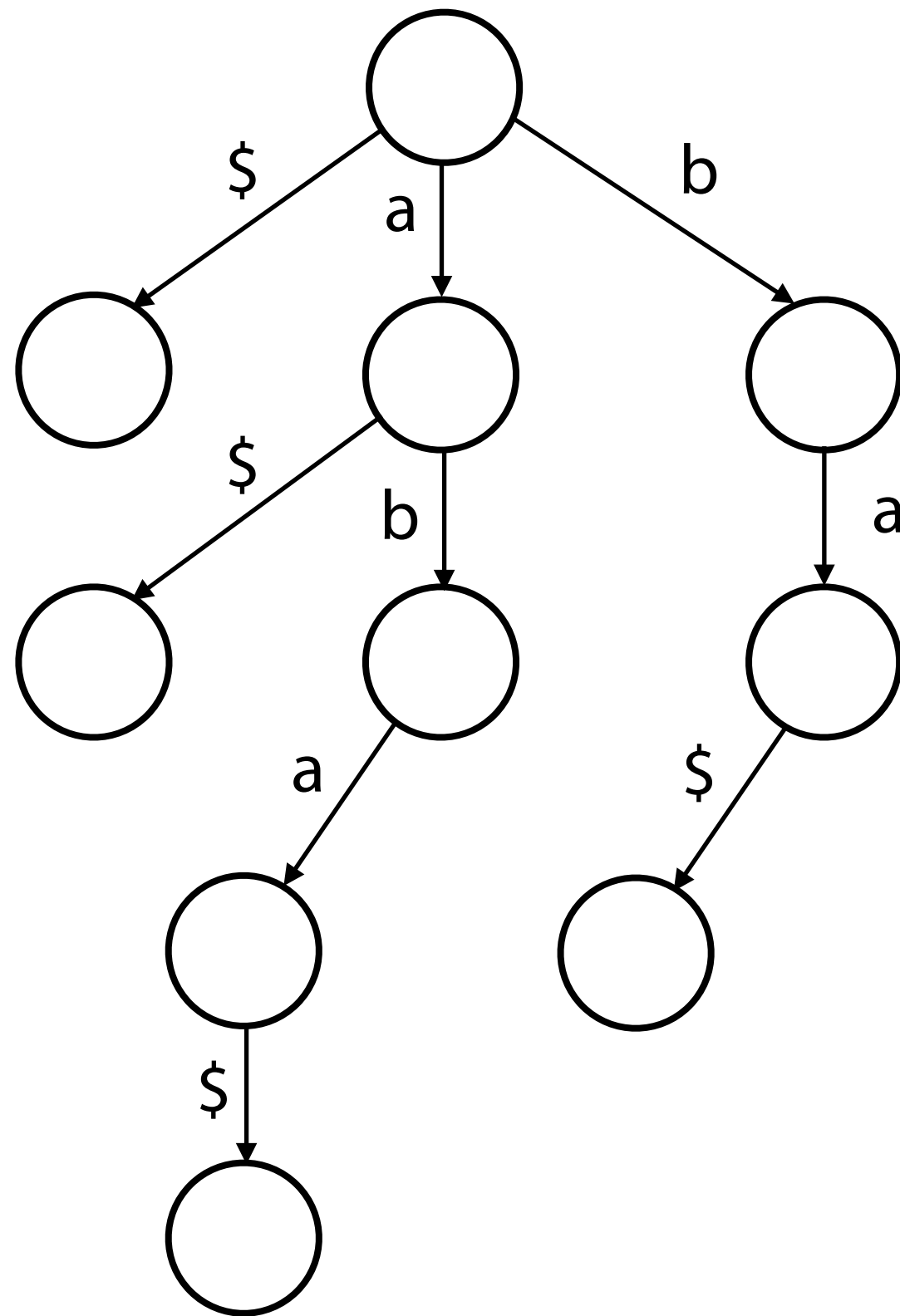
**\$** also guarantees no suffix is a prefix of any other suffix.

$T$ : G T T A T A G C T G A T C G C G G C G T A G C G G \$  
G T T A T A G C T G A T C G C G G C G T A G C G G \$  
T T A T A G C T G A T C G C G G C G T A G C G G \$  
T A T A G C T G A T C G C G G C G T A G C G G \$  
A T A G C T G A T C G C G G C G T A G C G G \$  
T A G C T G A T C G C G G C G T A G C G G \$  
A G C T G A T C G C G G C G T A G C G G \$  
G C T G A T C G C G G C G T A G C G G \$  
C T G A T C G C G G C G T A G C G G \$  
T G A T C G C G G C G T A G C G G \$  
G A T C G C G G C G T A G C G G \$  
A T C G C G G C G T A G C G G \$  
T C G C G G C G T A G C G G \$  
C G C G G C G T A G C G G \$  
G C G G C G T A G C G G \$  
C G G C G T A G C G G \$

# Suffix trie

*T*: aba\$

What's the suffix trie?





# Suffix trie

Excerpt from SuffixTrie class: [suffix-trie.ipynb](#)

```
def __init__(self, t):  
    """ Make suffix trie from t """  
    t += '$' # add terminator  
    self.root = {}  
    for i in range(len(t)): # for each suffix  
        cur = self.root  
        for c in t[i:]: # for each character in i'th suffix  
            if c not in cur:  
                cur[c] = {} # add outgoing edge if necessary  
            cur = cur[c] # follow the edge and continue
```

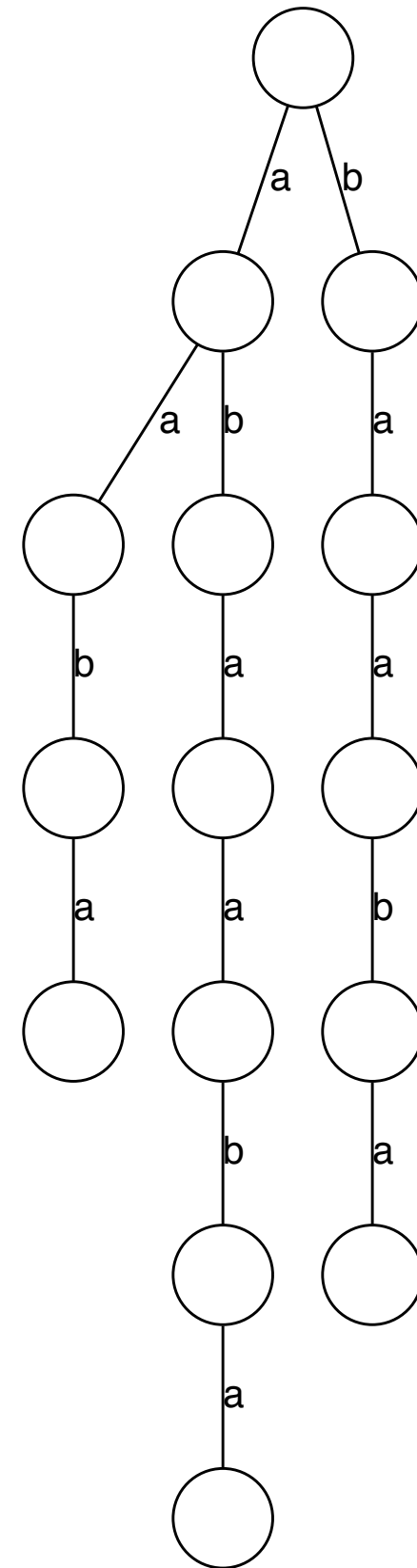


# Suffix trie

***T*: abaaba\$**

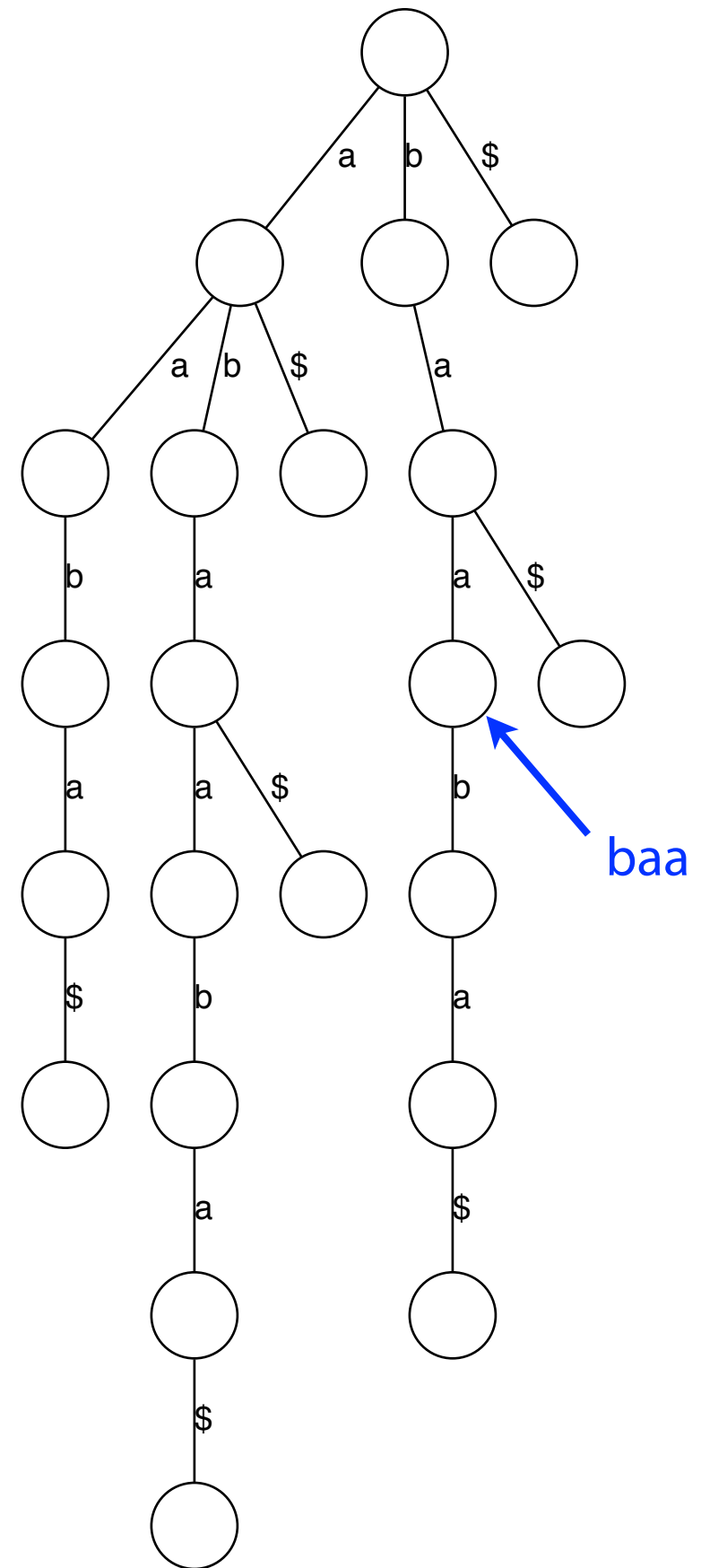
Each path from root to leaf represents a suffix; each suffix is represented by some path from root to leaf

Would this still be the case if we hadn't added \$? **No**



# Suffix trie

Think of each node as having a **label**, spelling out characters on path from root to node



# Suffix trie

How do we check whether a string  $S$  is a substring of  $T$ ?

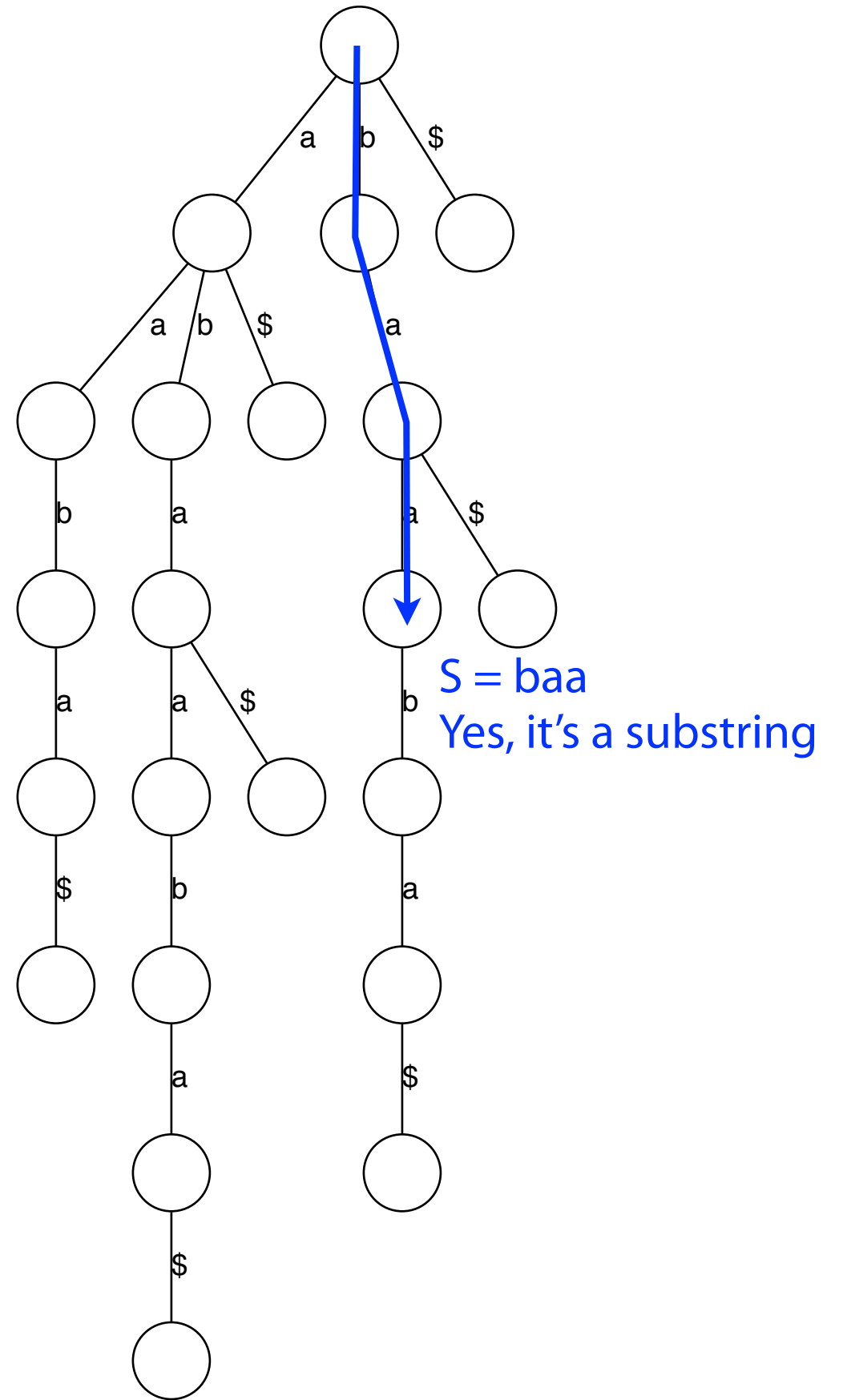
Note: Each of  $T$ 's substrings is spelled out along a path from the root.

**Every *substring* is a *prefix* of some *suffix* of T.**

Start at the root and follow the edges labeled with the characters of  $S$

If we “fall off” the trie -- i.e. there is no outgoing edge for next character of  $S$ , then  $S$  is not a substring of  $T$

If we exhaust  $S$  without falling off,  $S$  is a substring of  $T$



# Suffix trie

How do we check whether a string  $S$  is a substring of  $T$ ?

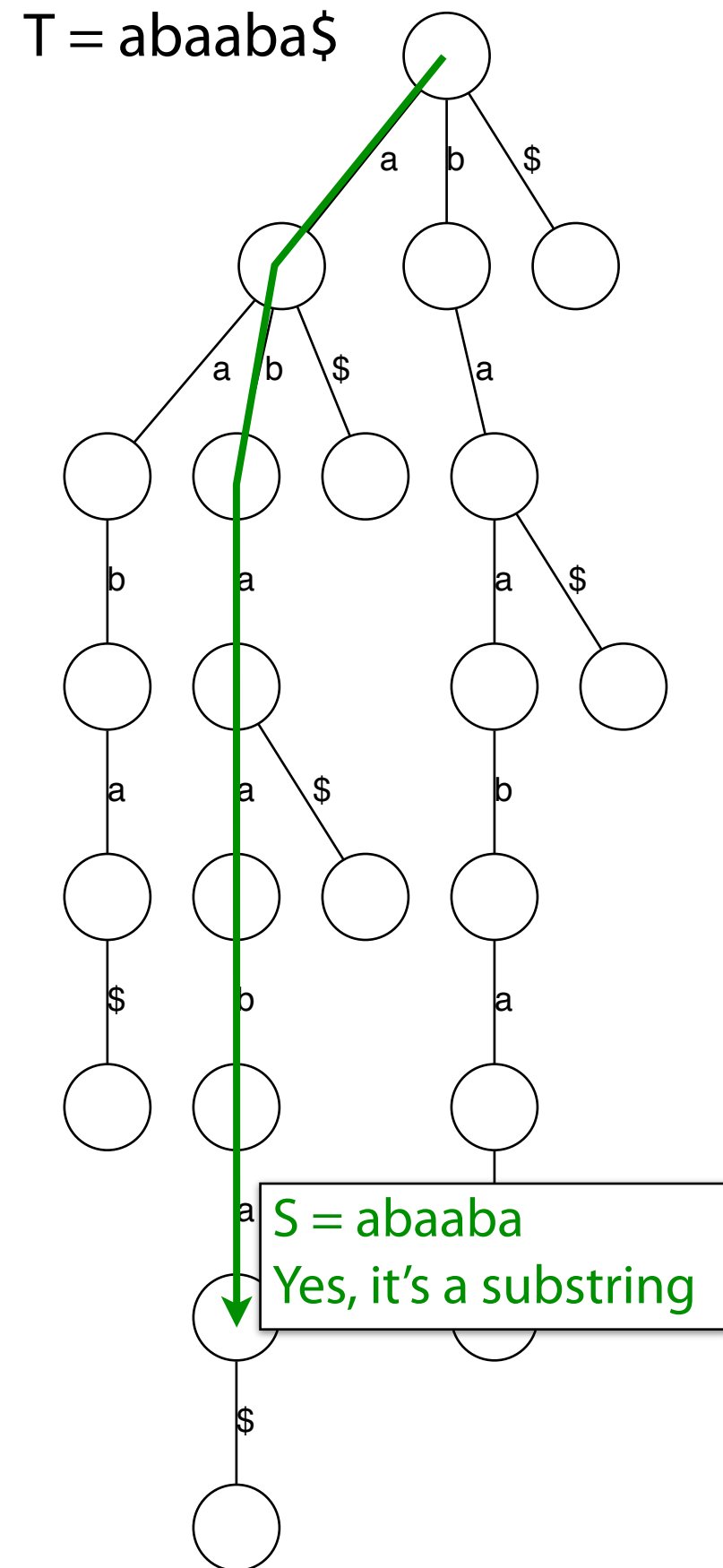
Note: Each of  $T$ 's substrings is spelled out along a path from the root.

**Every substring is a prefix of some suffix** of  $T$ .

Start at the root and follow the edges labeled with the characters of  $S$

If we “fall off” the trie -- i.e. there is no outgoing edge for next character of  $S$ , then  $S$  is not a substring of  $T$

If we exhaust  $S$  without falling off,  $S$  is a substring of  $T$



# Suffix trie

How do we check whether a string  $S$  is a substring of  $T$ ?

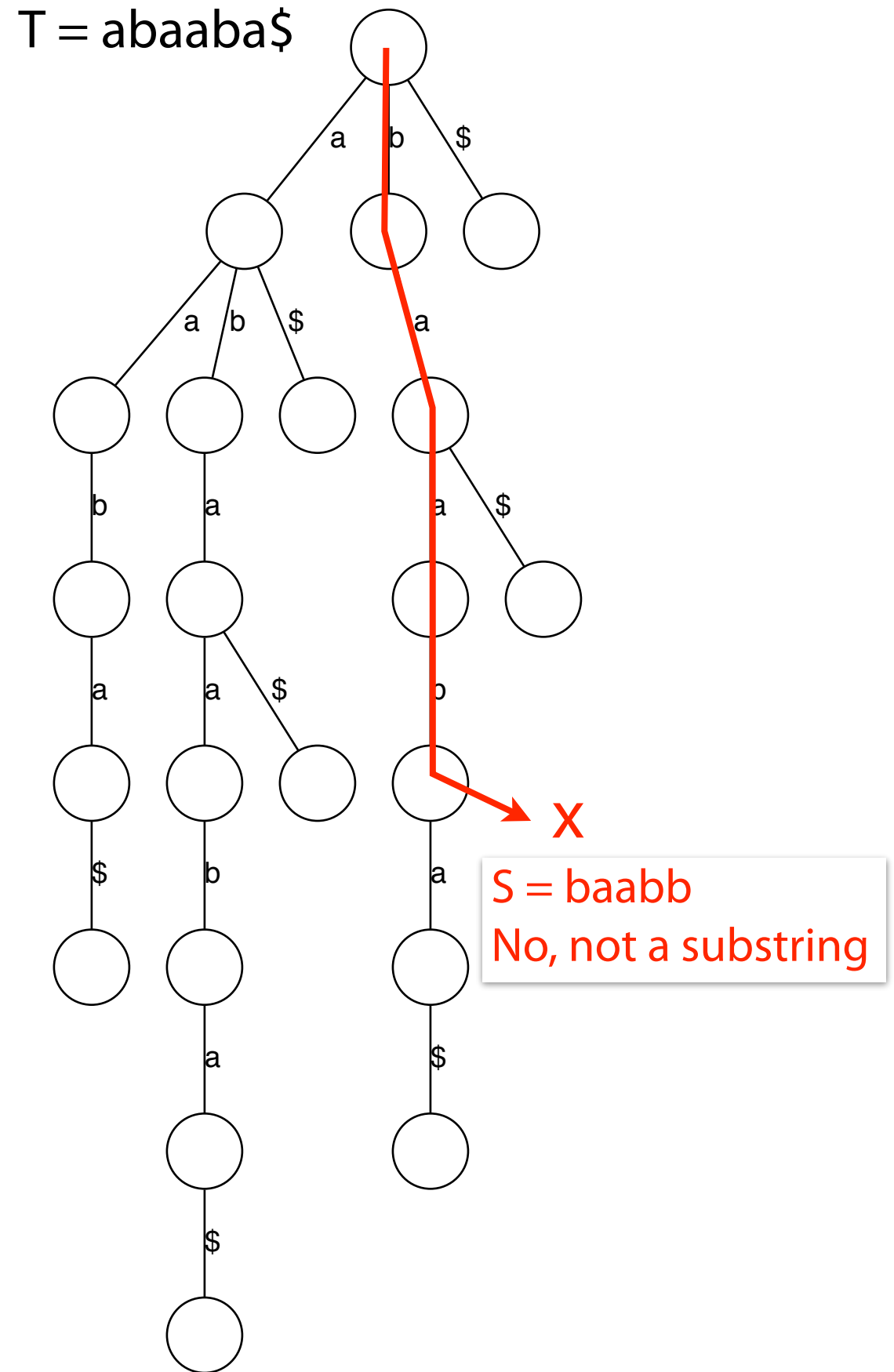
Note: Each of  $T$ 's substrings is spelled out along a path from the root.

**Every *substring* is a *prefix* of some *suffix* of T.**

Start at the root and follow the edges labeled with the characters of  $S$

If we “fall off” the trie -- i.e. there is no outgoing edge for next character of  $S$ , then  $S$  is not a substring of  $T$

If we exhaust  $S$  without falling off,  $S$  is a substring of  $T$



# Suffix trie

Excerpt from SuffixTrie class: [suffix-trie.ipynb](#)

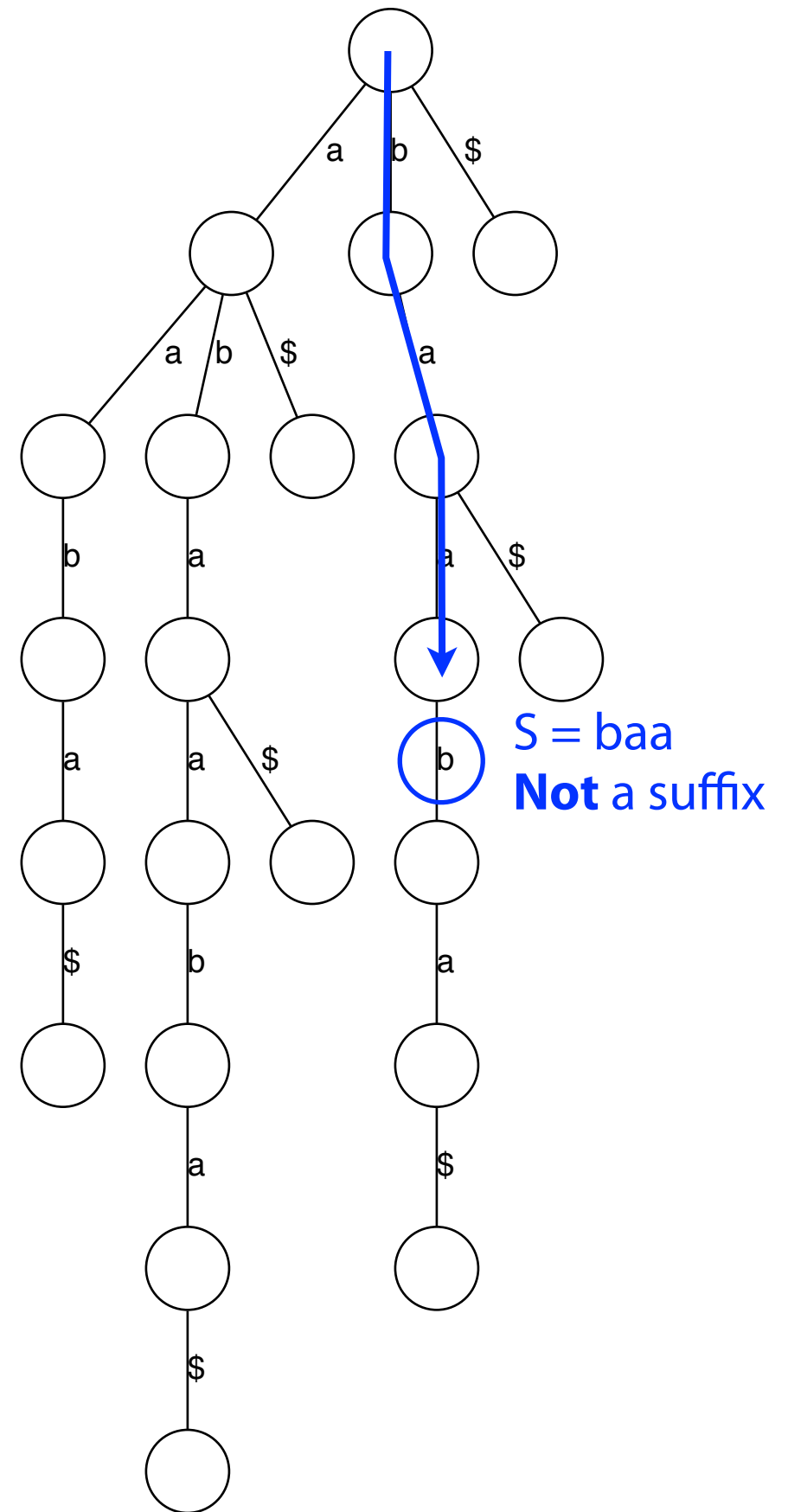
```
def follow_path(self, s):  
    """ Follow path given by characters of s. Return node at  
        end of path, or None if we fall off. """  
    cur = self.root  
    for c in s:  
        if c not in cur:  
            return None # no outgoing edge on next character  
        cur = cur[c] # descend one level  
    return cur  
  
def has_substring(self, s):  
    """ Return true if s appears as a substring of t """  
    return self.follow_path(s) is not None
```



# Suffix trie

How do we check whether a string  $S$  is a **suffix** of  $T$ ?

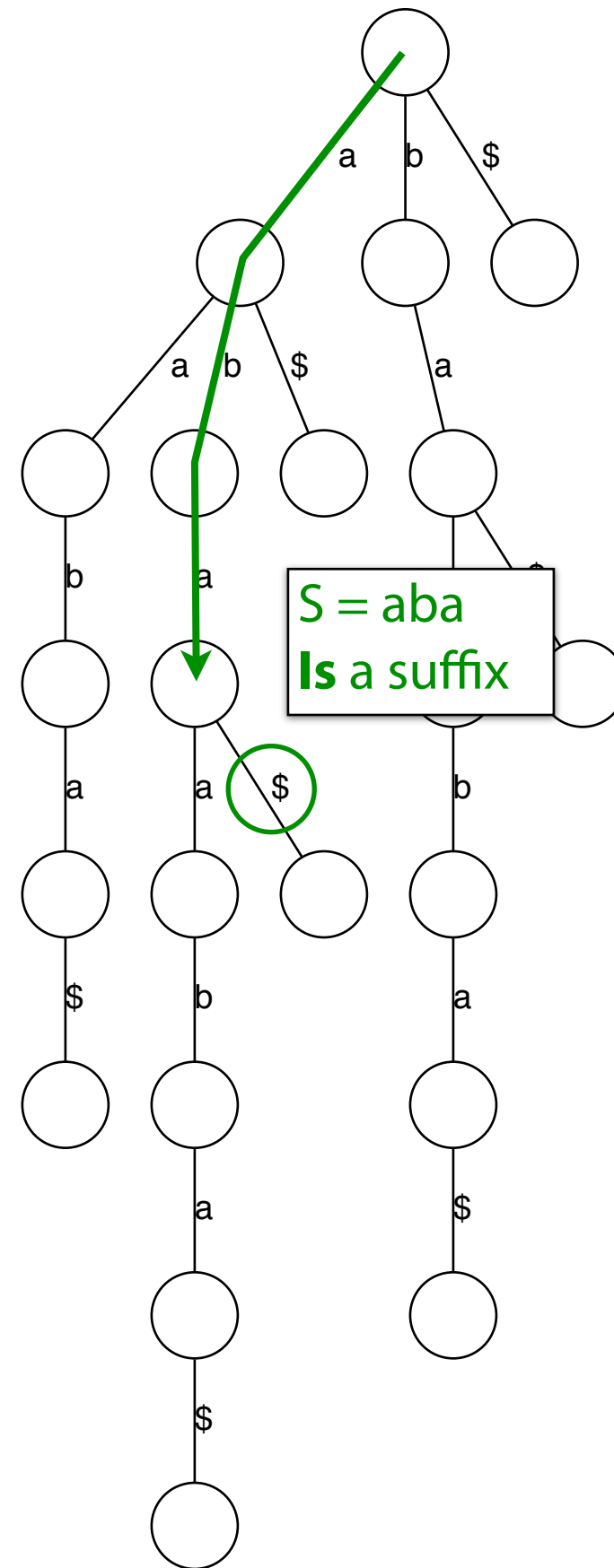
Same procedure as for substring, but additionally check terminal node for \$ child



# Suffix trie

How do we check whether a string  $S$  is a **suffix** of  $T$ ?

Same procedure as for substring, but additionally check terminal node for **\$** child



# Suffix trie

Excerpt from SuffixTrie class: [suffix-trie.ipynb](#)

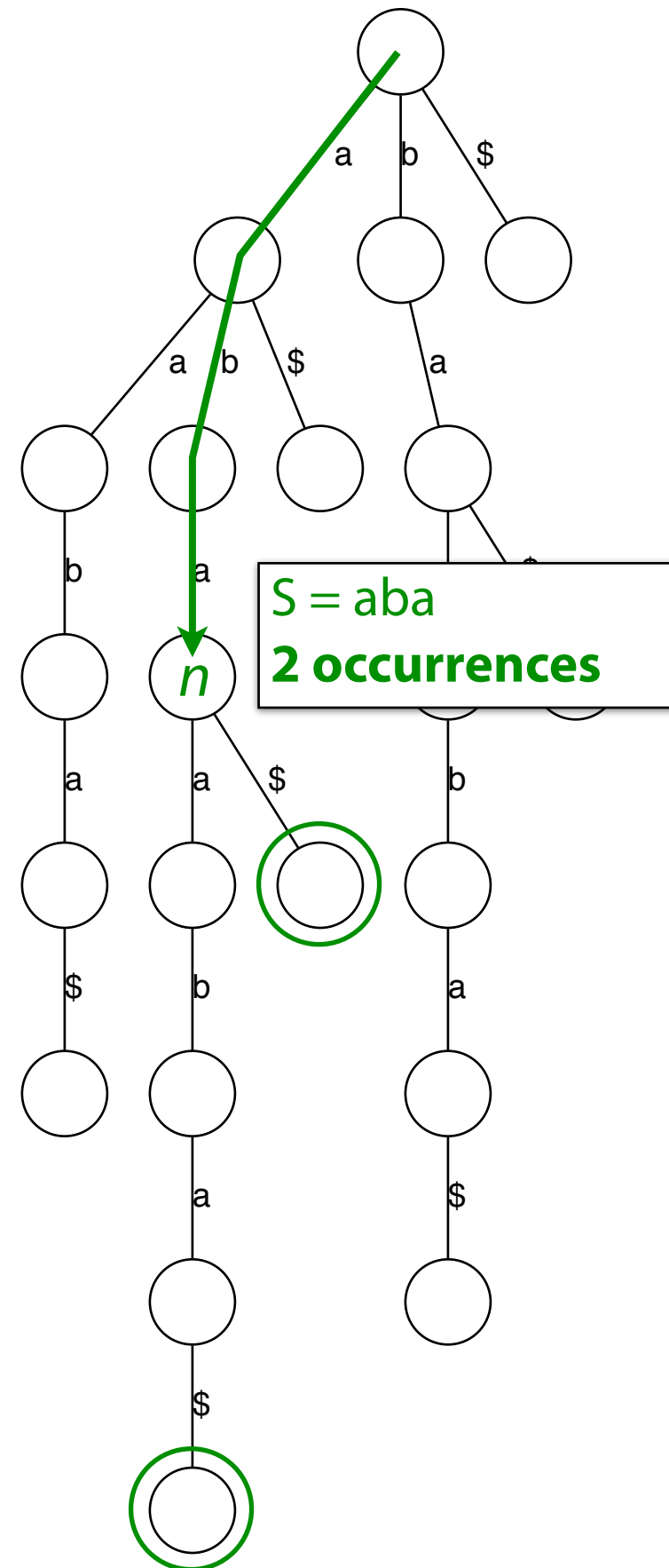
```
def has_suffix(self, s):  
    """ Return true if s is a suffix of t """  
    node = self.follow_path(s)  
    return node is not None and '$' in node
```

# Suffix trie

How do we count the **number of times**  
a string  $S$  occurs as a substring of  $T$ ?

Follow path labeled with  $S$ . If we fall off, answer is 0. If we end up at node  $n$ , answer equals # of leaves in subtree rooted at  $n$ .

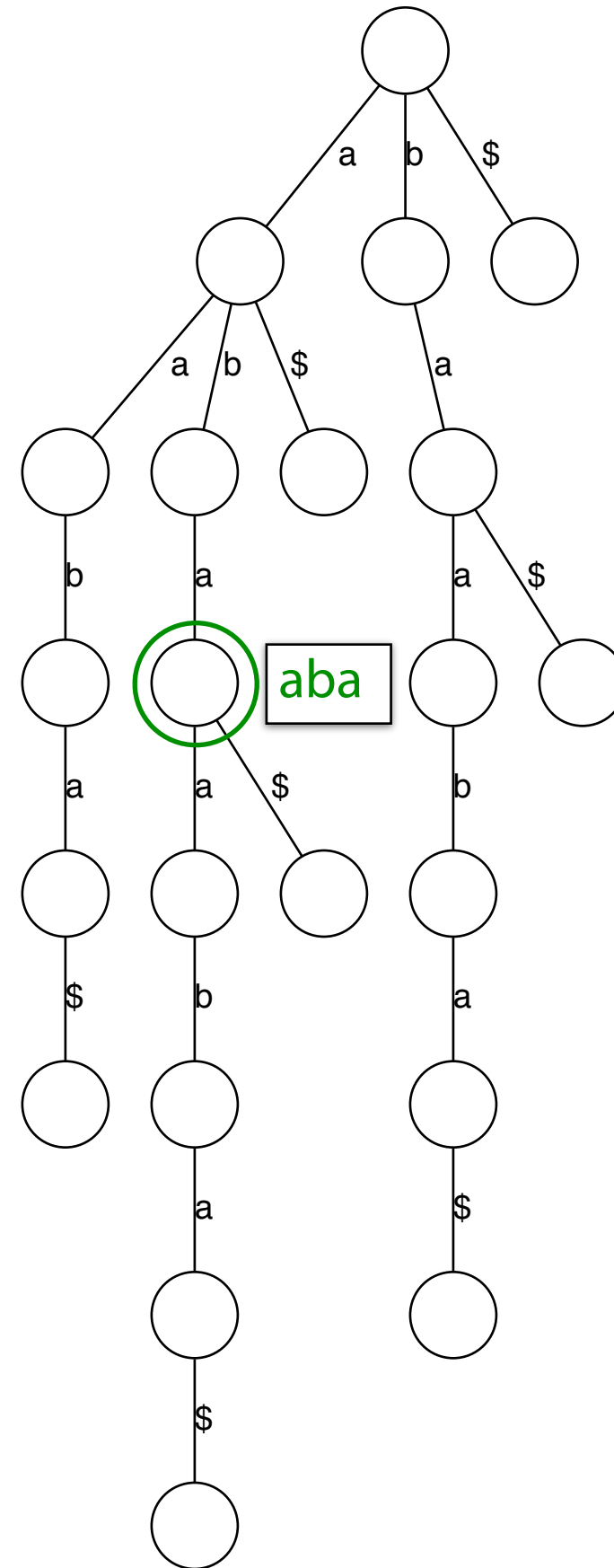
Leaves can be counted with depth-first traversal.



# Suffix trie

How do we find the **longest repeated substring** of  $T$ ?

## Find the deepest node with more than one child



# Suffix trie

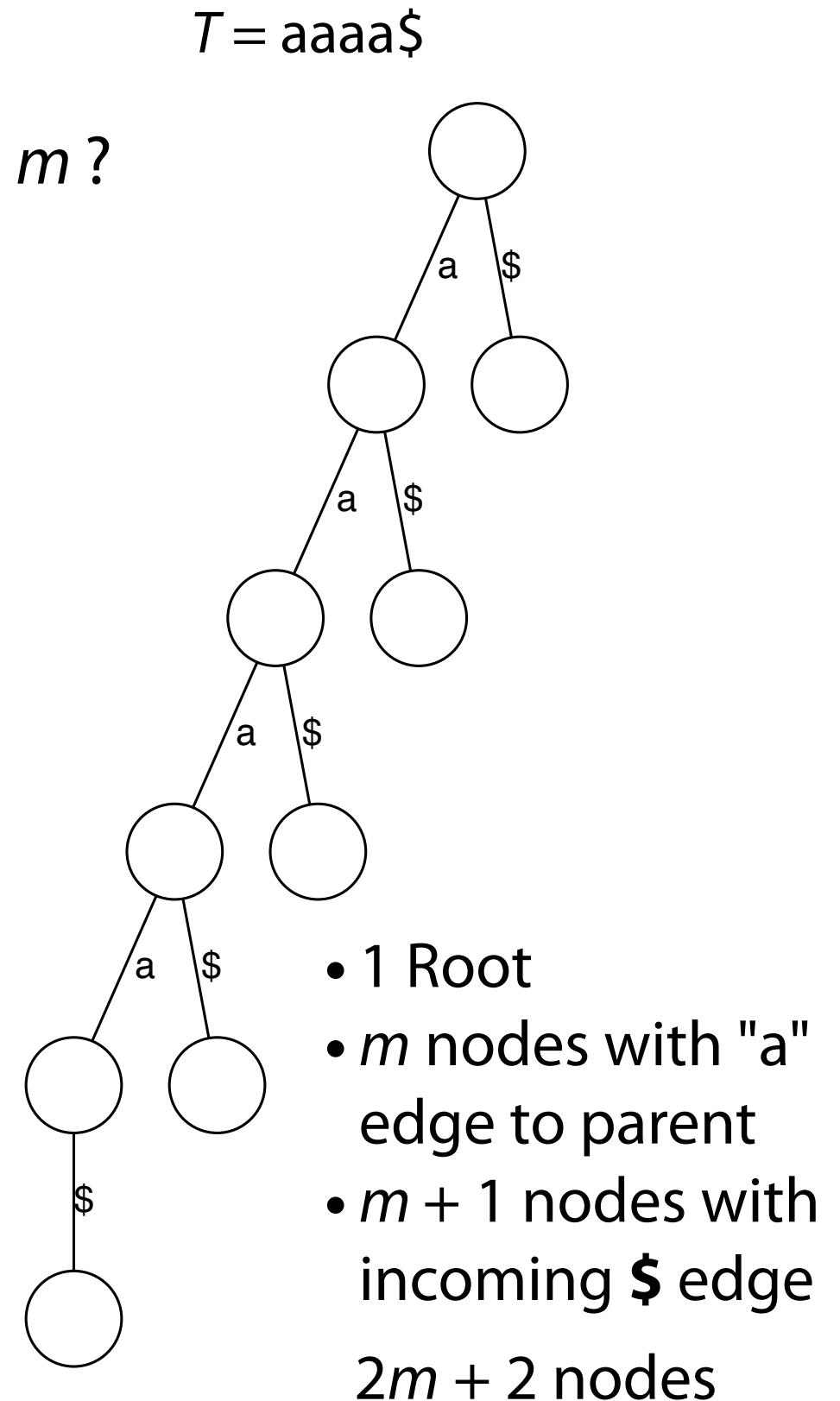
How does the suffix trie grow with  $|T| = m$  ?

# Suffix trie

How does the suffix trie grow with  $|T| = m$ ?

Is there a class of string where the number of suffix trie nodes grows linearly with  $m$ ?

Yes: a string of  $m$  a's in a row ( $a^m$ )



# Suffix trie

$T = \text{aaabbbb}\$$

How does the suffix trie grow with  $|T| = m$ ?

Is there a class of string where the number of suffix trie nodes grows with  $m^2$ ?

Yes:  $a^n b^n$  where  $2n = m$

- 1 root
- $n$  nodes along "b chain," right
- $n$  nodes along "a chain," middle
- $n$  chains of  $n$  "b" nodes hanging off "a chain" ( $n^2$  total)
- $2n + 1$  \$ leaves (not shown)

$n^2 + 4n + 2$  nodes, where  $m = 2n$

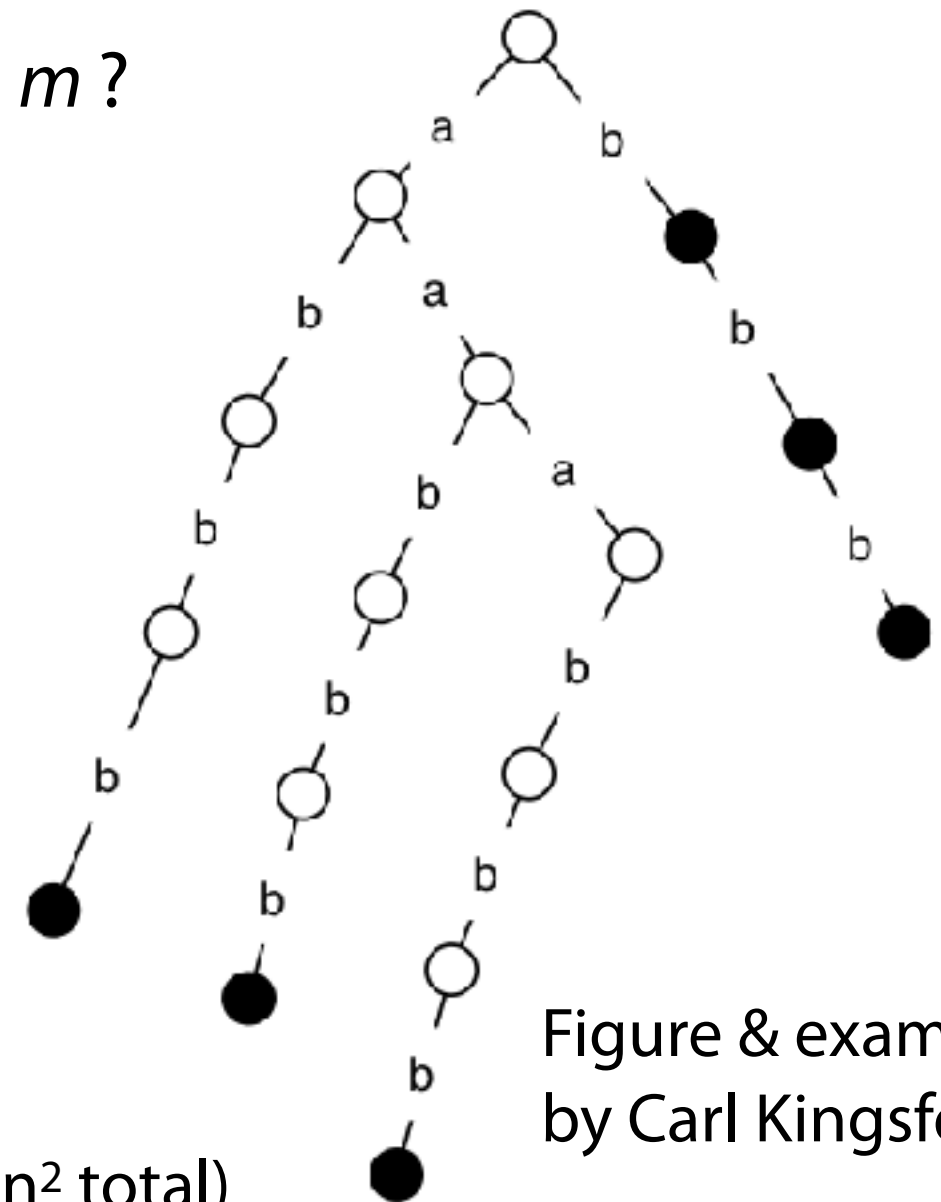


Figure & example  
by Carl Kingsford