

Suffix tree

- Build a tree from the text
- Used if the text is expected to be the same during several pattern queries
- Tree building is $O(m)$ where m is the size of the text. This is preprocessing.
- Given any pattern of length n , we can answer if it occurs in text in $O(n)$ time
- Suffix tree = “modified” keyword tree of all suffixes of text

Construct a suffix tree

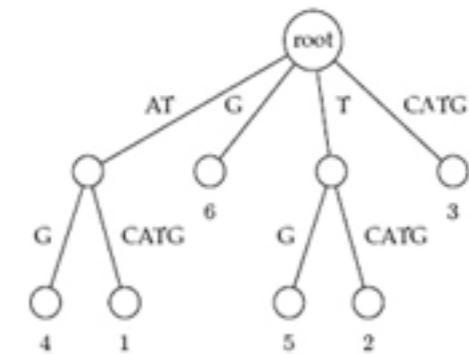
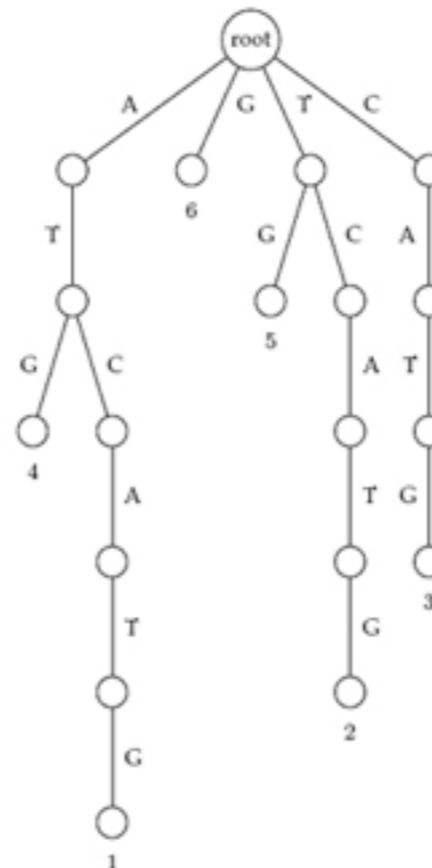
Text: ATCATG

suffixes

ATCATG
TCATG
CATG
ATG
TG
G

Keyword
Tree

Suffix
Tree

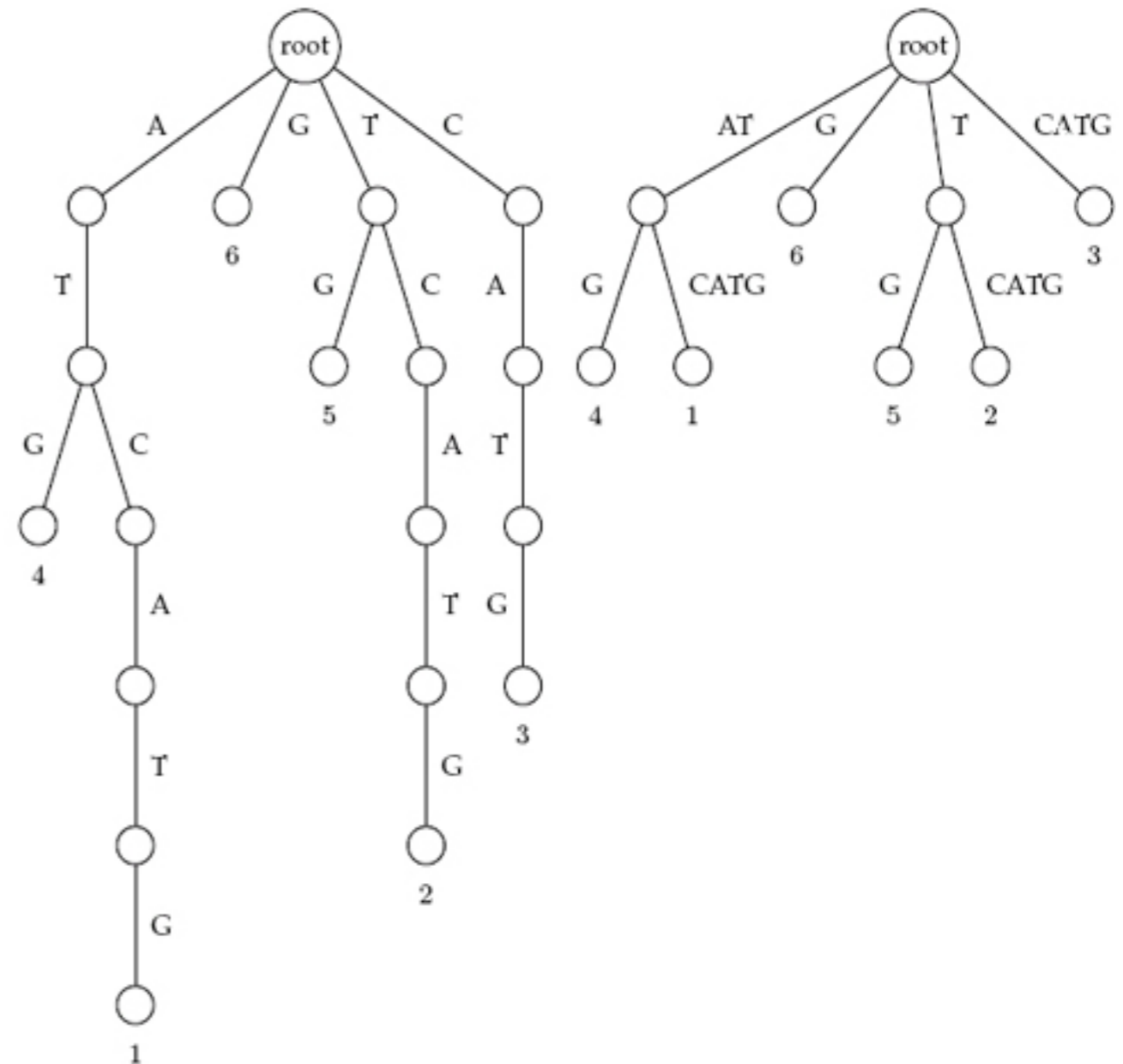


Suffix tree = Collapsed Keyword Tree on Suffixes

Similar to keyword trees, except edges that form paths are collapsed

- Each edge is labeled with a *substring* of a text for less space
- All internal edges have at least two outgoing edges
- Leaves labeled by the location of the suffix on the text.

Text: **ATCATG**



(a) Keyword tree

(b) Suffix tree

All suffixes of text T

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
G G C G T A G C G G
G C G T A G C G G
C G T A G C G G
G T A G C G G
T A G C G G
A G C G G
G C G G
C G G
G G
G

$m(m+1)/2$
chars

Example: suffix keyword tree

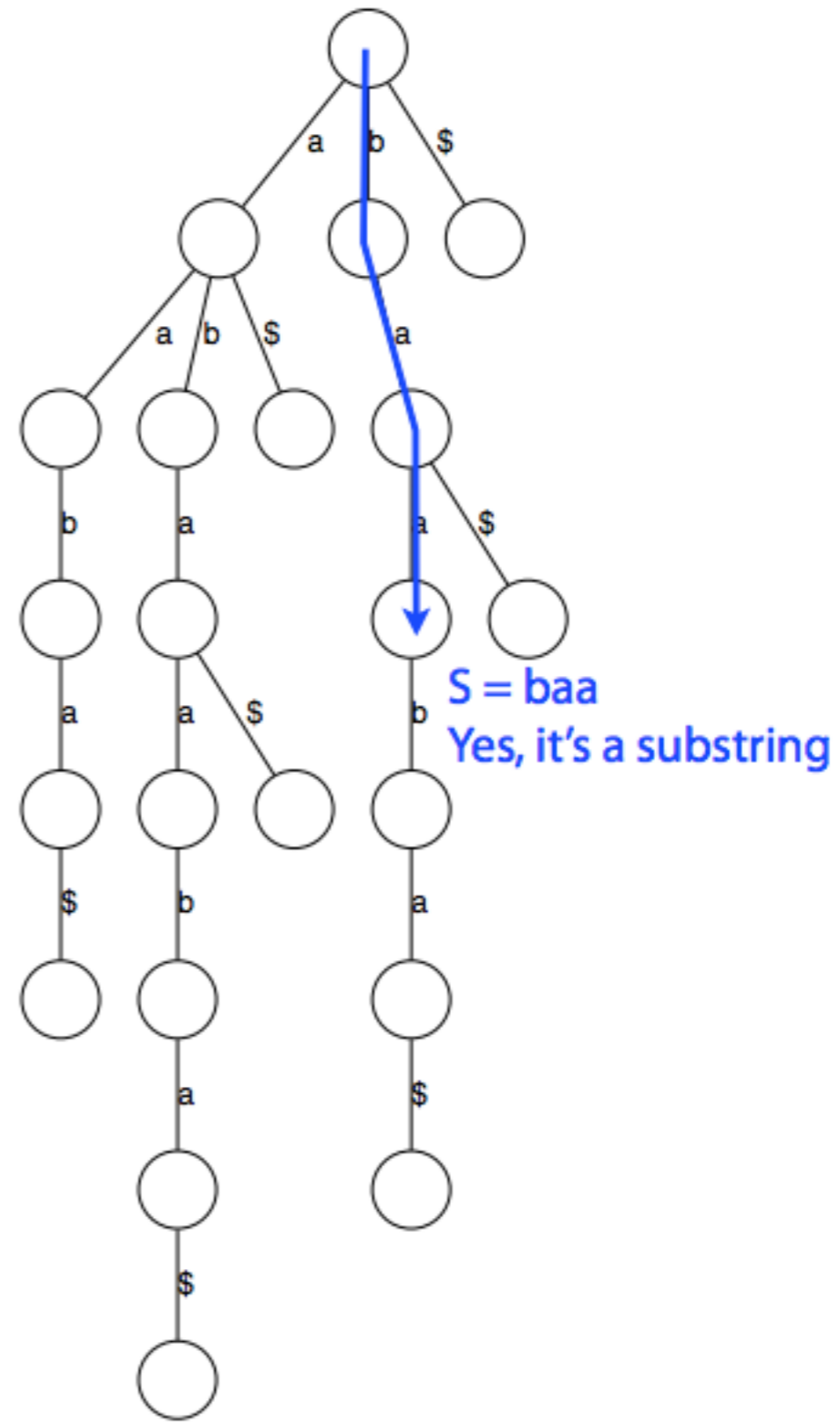
How do we check whether a string S is a substring of T ?

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



Example: suffix keyword tree

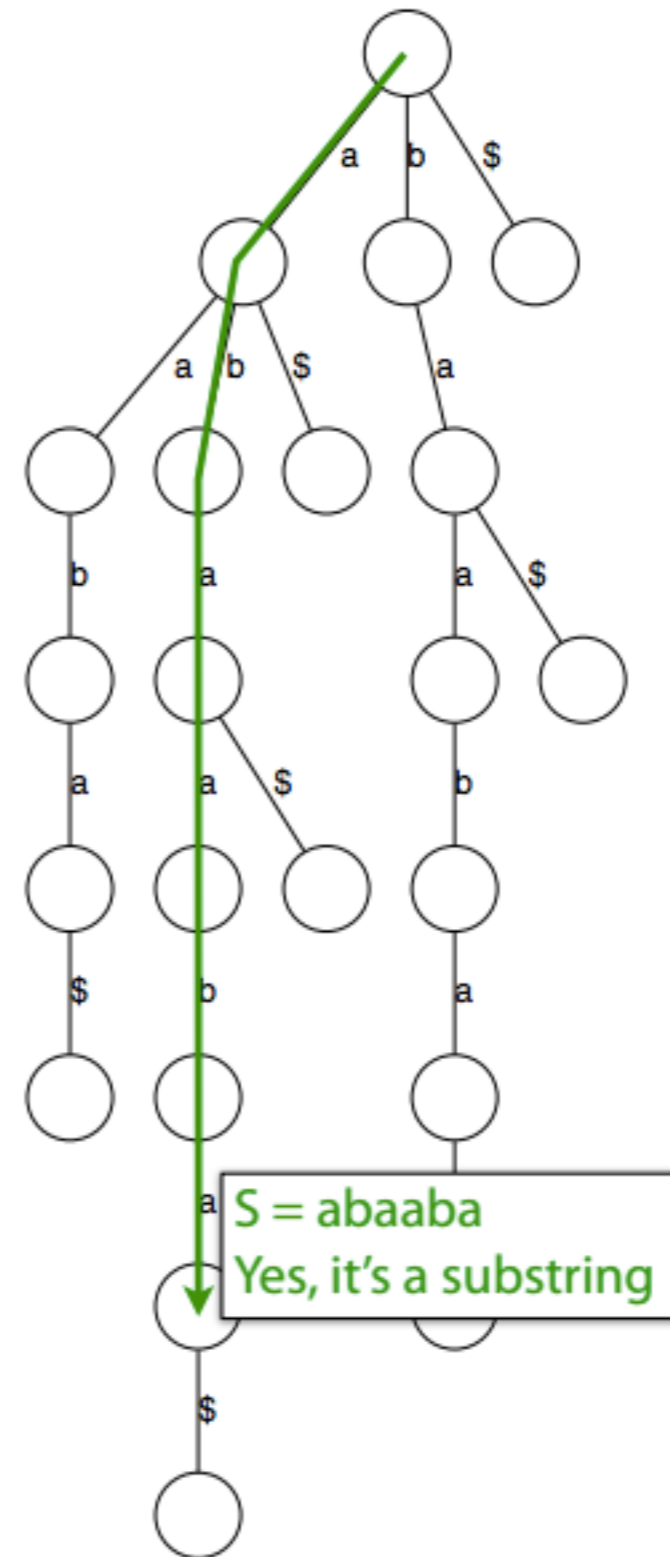
How do we check whether a string S is a substring of T ?

Note: Each of T 's substrings is spelled out along a path from the root. I.e., every *substring* is a *prefix* of some *suffix* of T .

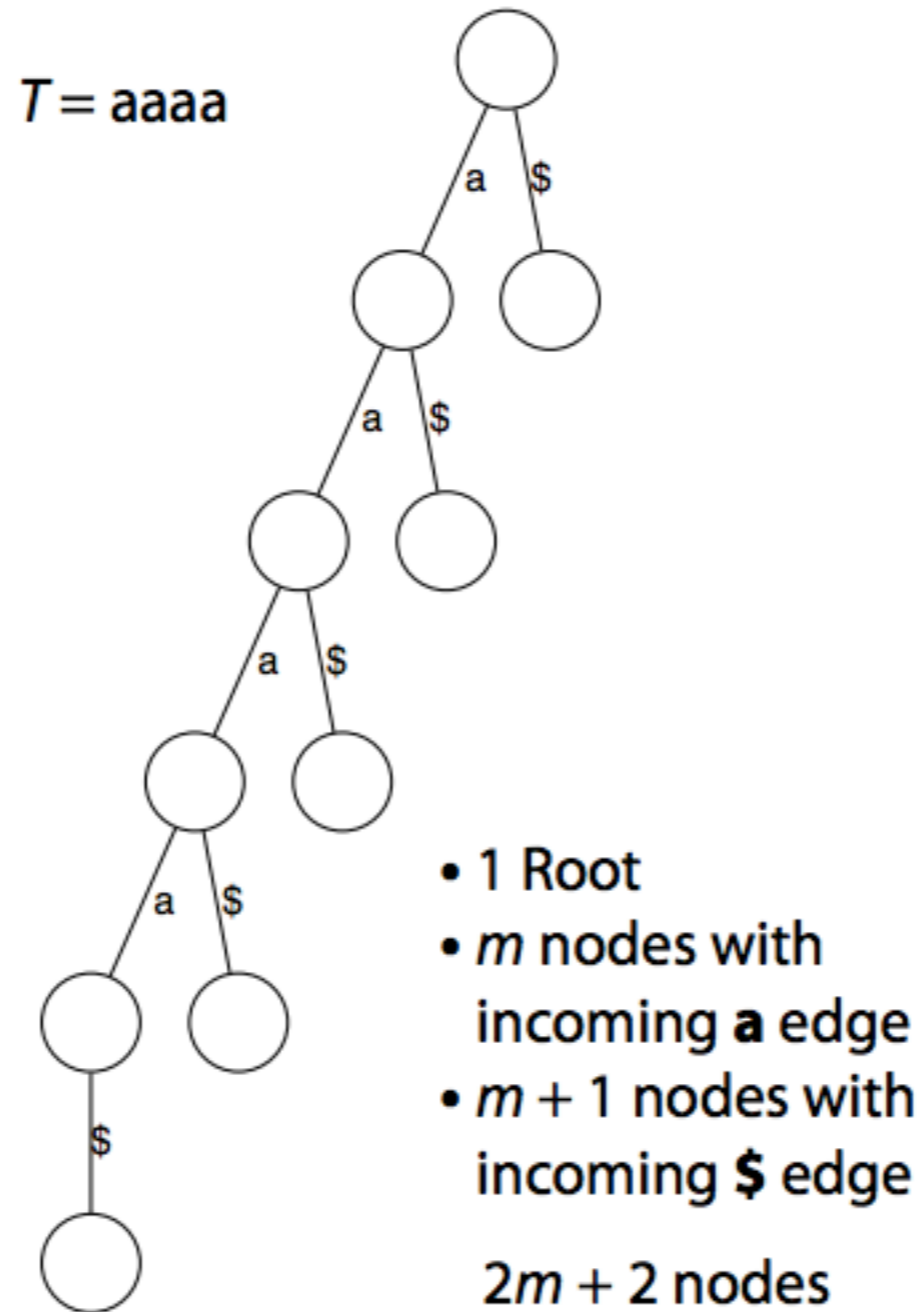
Start at the root and follow the edges labeled with the characters of S

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



How many nodes does a suffix keyword tree have?



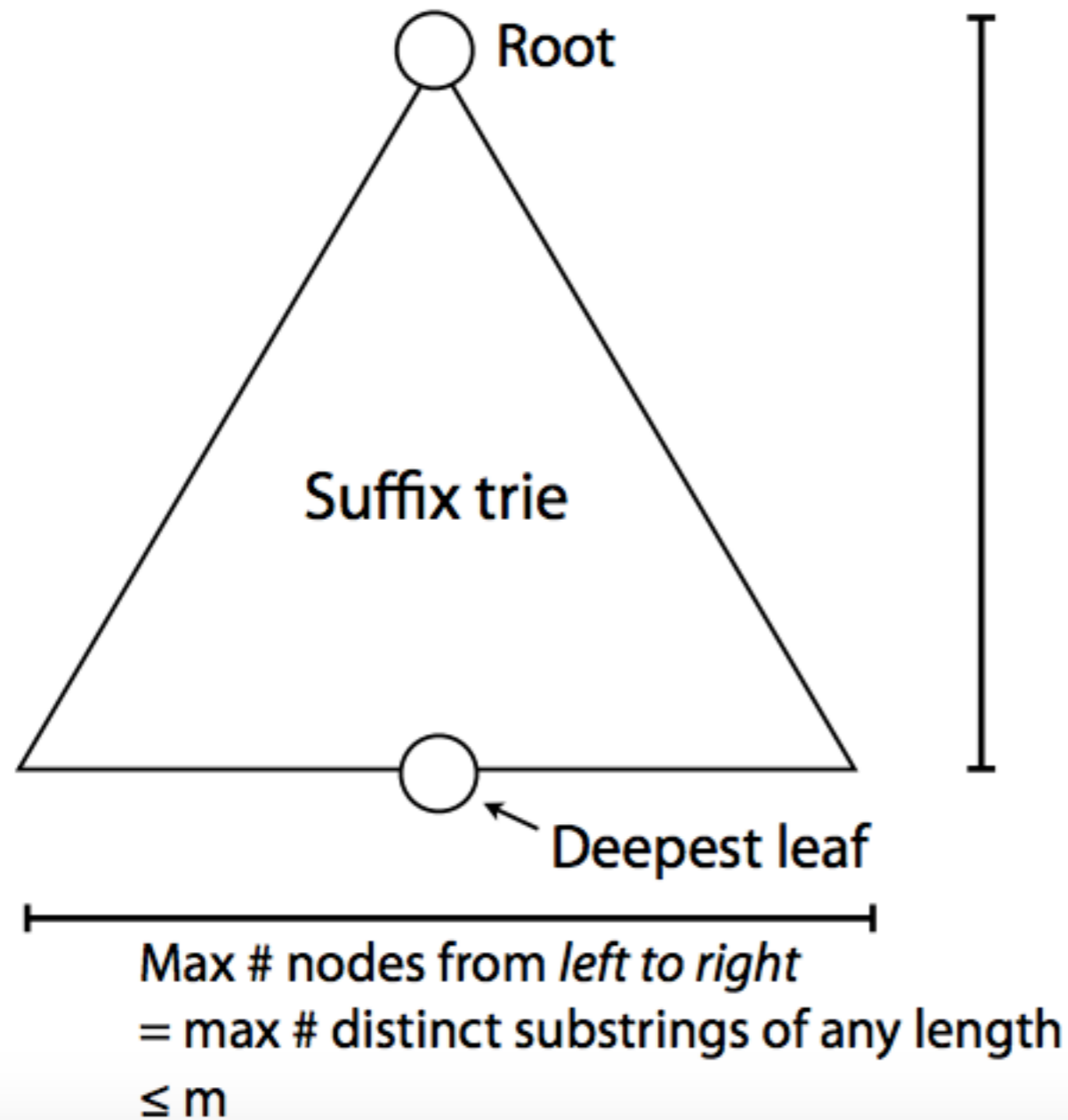
How many nodes does a suffix keyword tree have?

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
G G C G T A G C G G
G C G T A G C G G
C G T A G C G G
G T A G C G G
T A G C G G
A G C G G
G C G G
C G G
G G
G

$m(m+1)/2$
chars

How many nodes does a suffix keyword tree have?

Could worst-case # nodes be worse than $O(m^2)$?

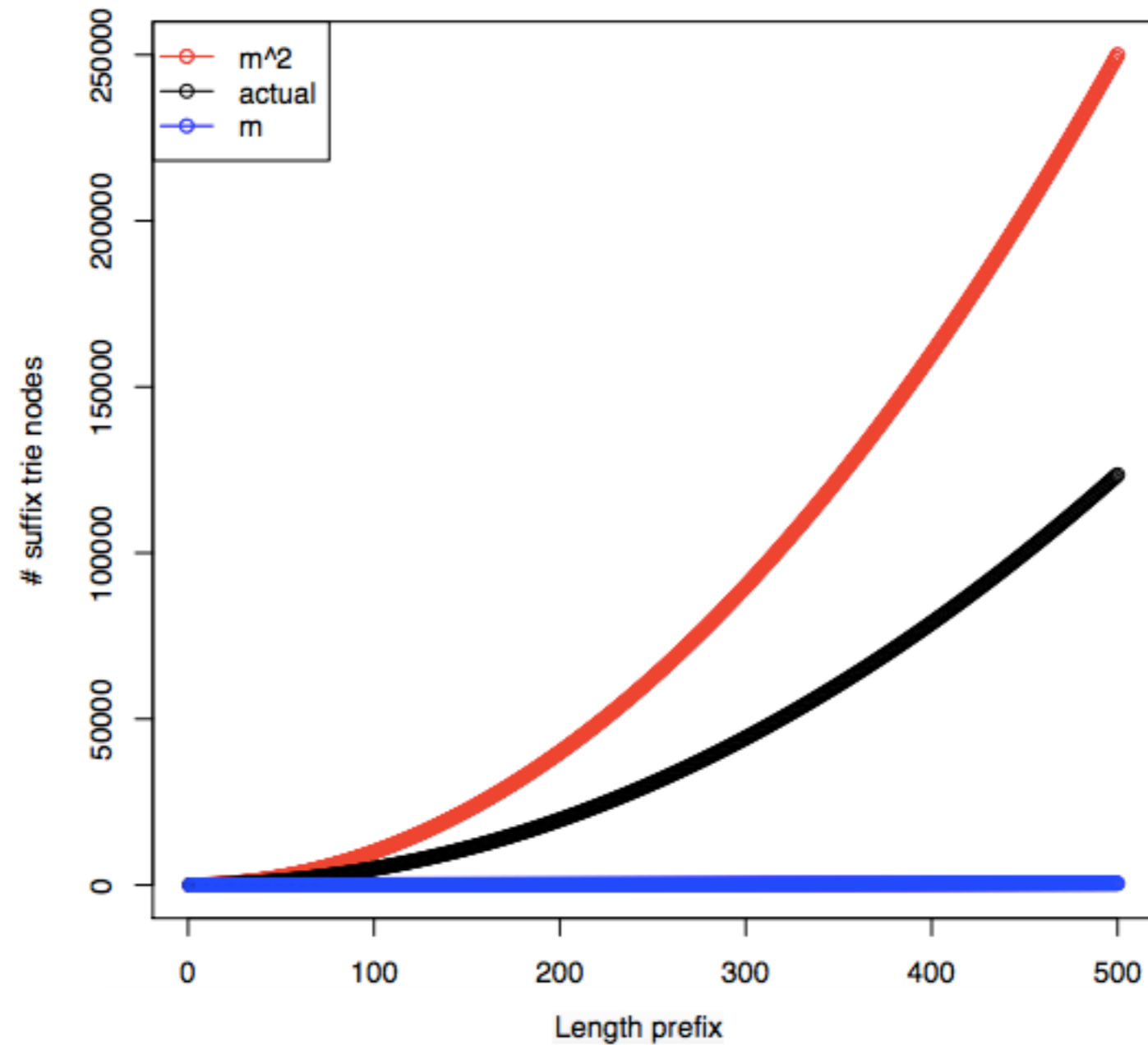


Max # nodes from *top to bottom*
= length of longest suffix + 1
= $m + 1$

$O(m^2)$ is worst case

Actual growth: an example

Trees built using the first 500 prefixes of the lambda phage virus genome



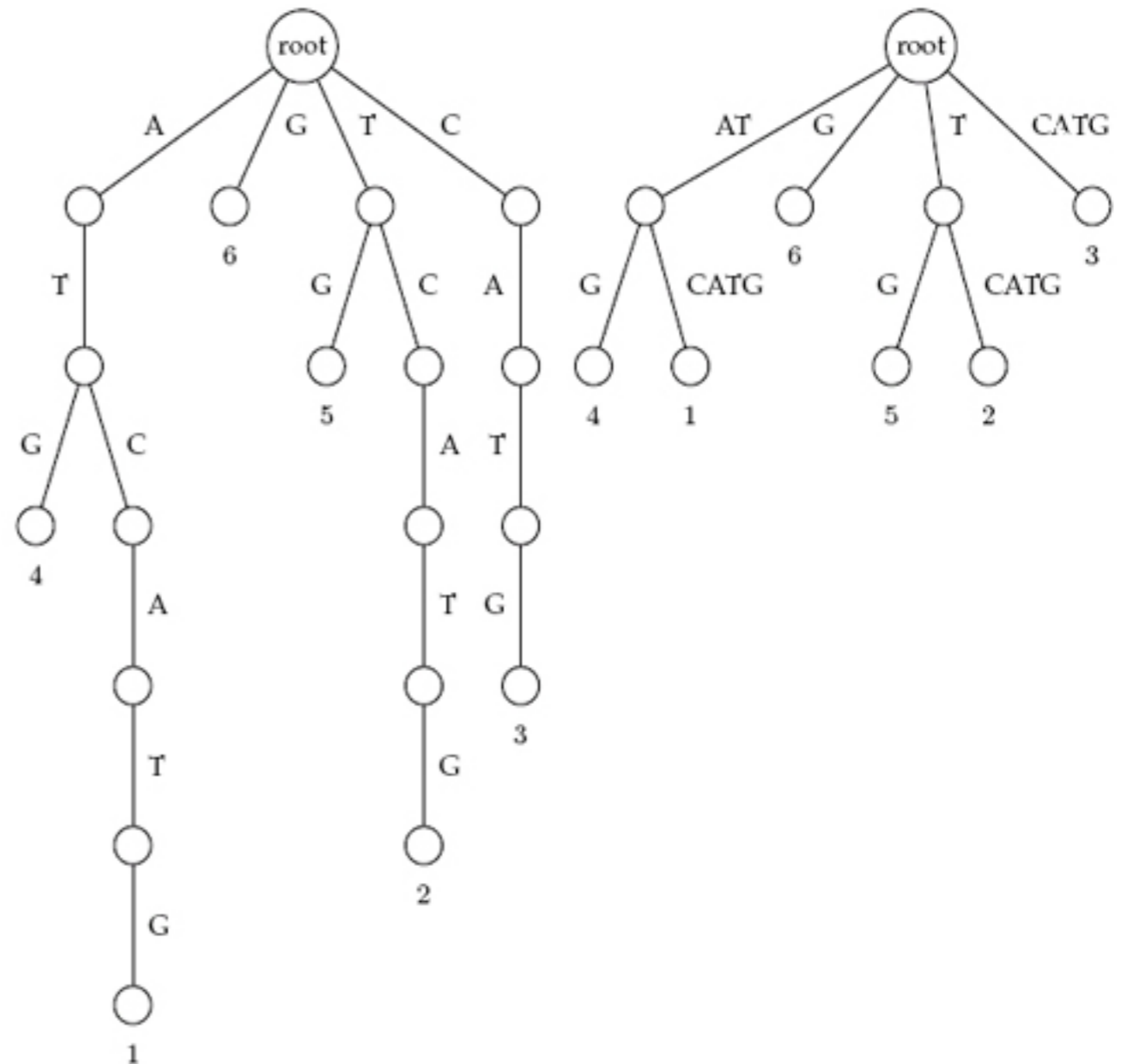
How to compress these trees?

Suffix tree = Collapsed Keyword Tree on Suffixes

Similar to keyword trees, except edges that form paths are collapsed

- Each edge is labeled with a *substring* of a text for less space
- All internal edges have at least two outgoing edges
- Leaves labeled by the location of the suffix on the text.

Text: **ATCATG**



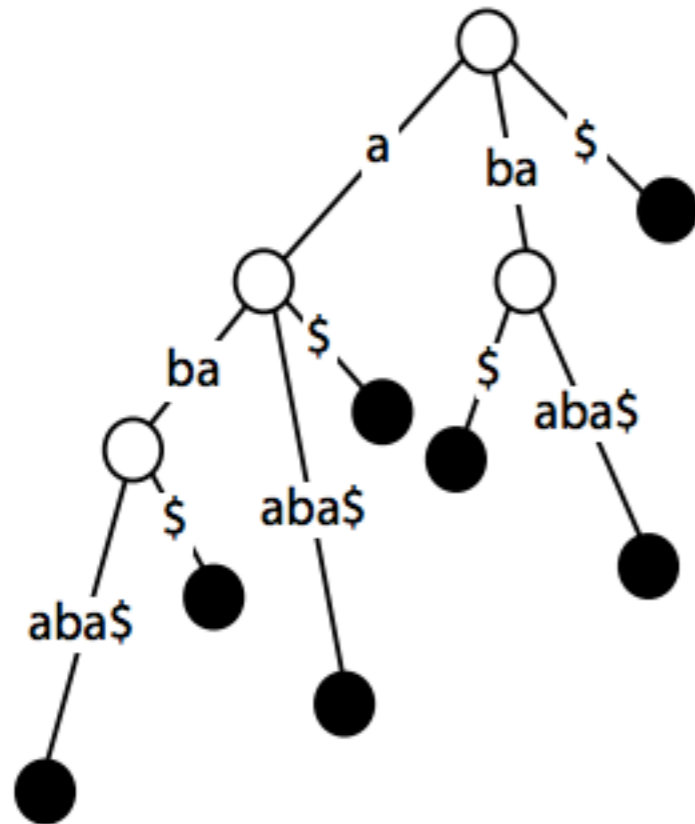
(a) Keyword tree

(b) Suffix tree

How many nodes does a suffix tree have?

L leaves, I internal nodes, E edges

$T = \text{abaaba}\$$



$$E = L + I - 1$$

$$E \geq 2I \text{ (each internal node branches)}$$

$$L + I - 1 \geq 2I \Rightarrow I \leq L - 1$$

but

$$L \leq m \text{ (at most } m \text{ suffixes)}$$

$$I \leq m - 1$$

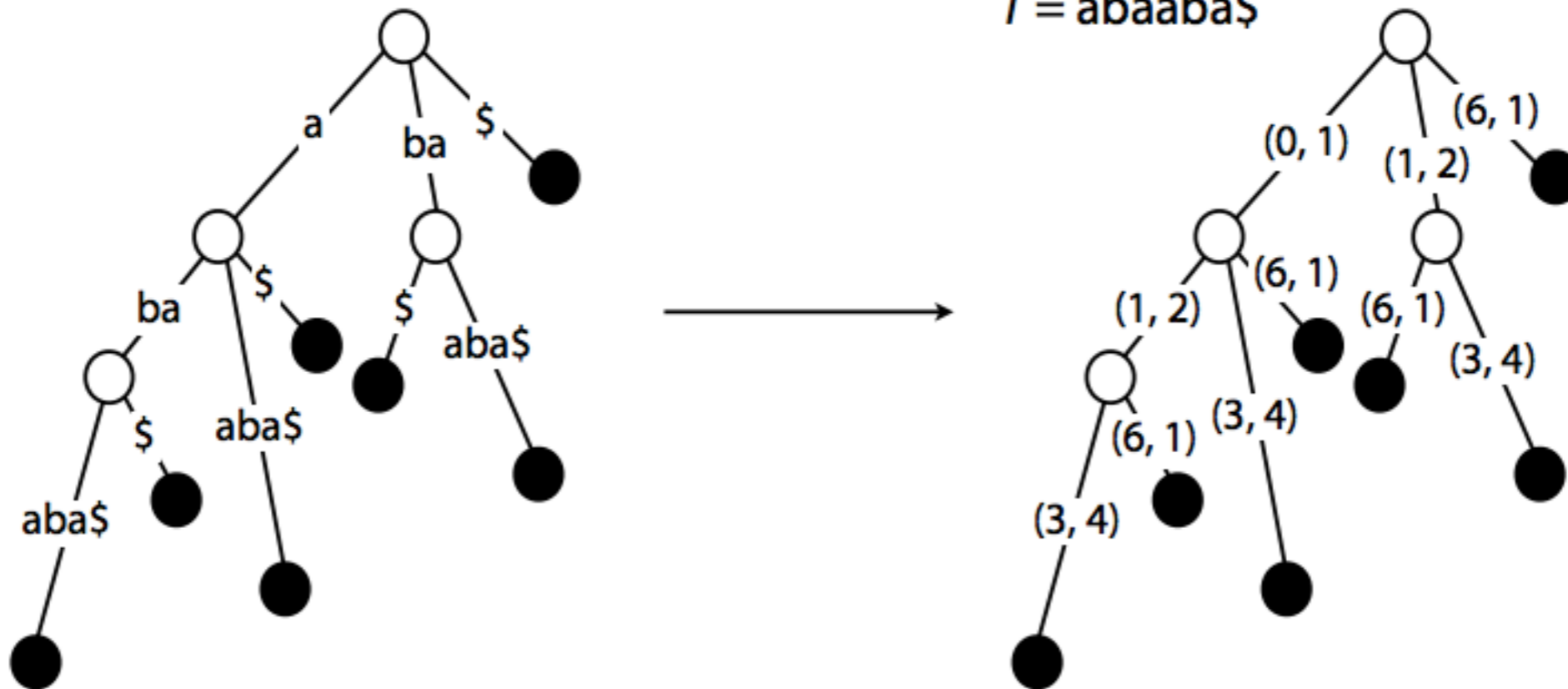
$$E = L + I - 1 \leq 2m - 2$$

$$E + L + I \leq 4m - 3 \in O(m)$$

Space complexity

$T = \text{abaaba}\$$

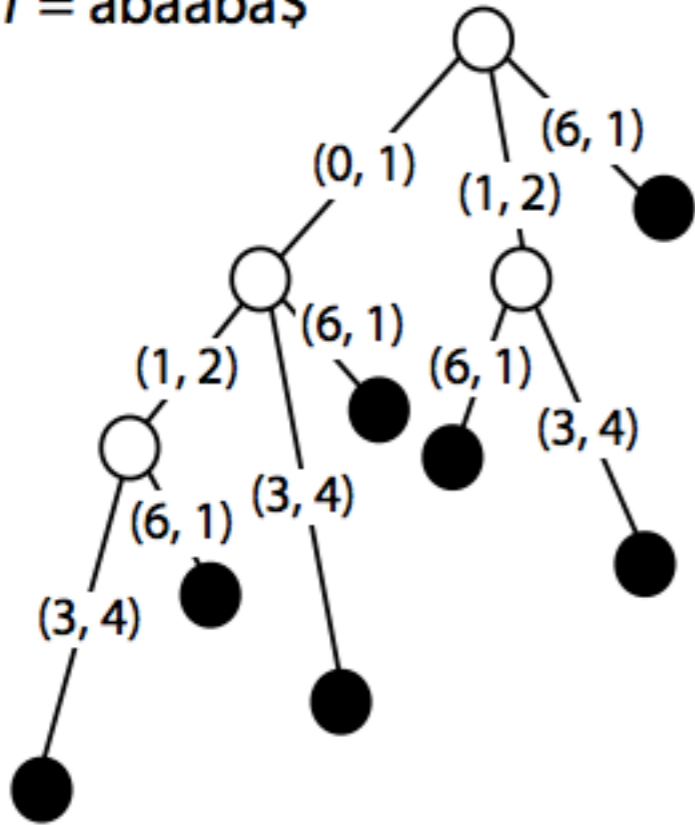
Idea 2: Store T itself in addition to the tree. Convert tree's edge labels to (offset, length) pairs with respect to T .



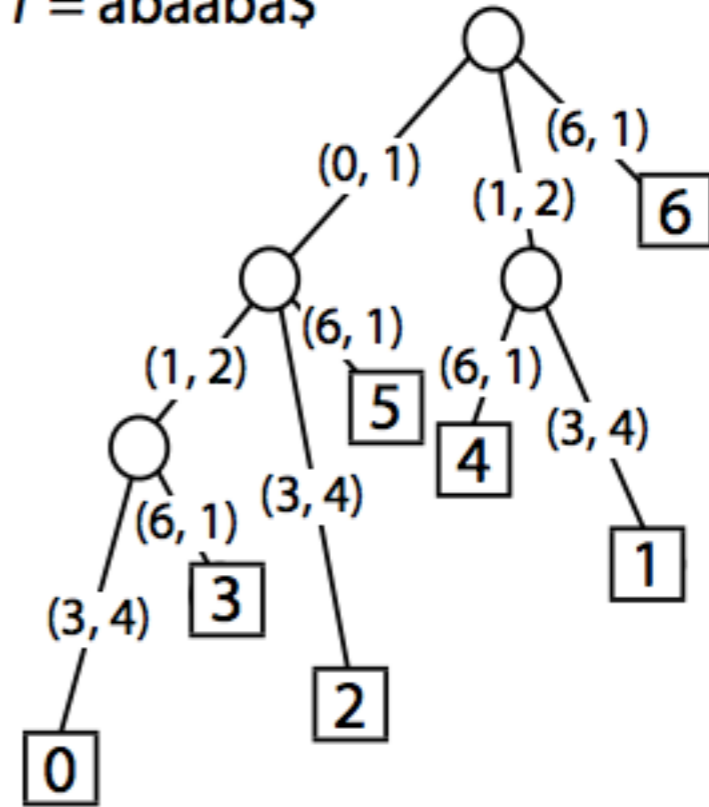
Space required for suffix tree is now $O(m)$

Add starting location/offset at each leaf node

$T = \text{abaaba}\$$

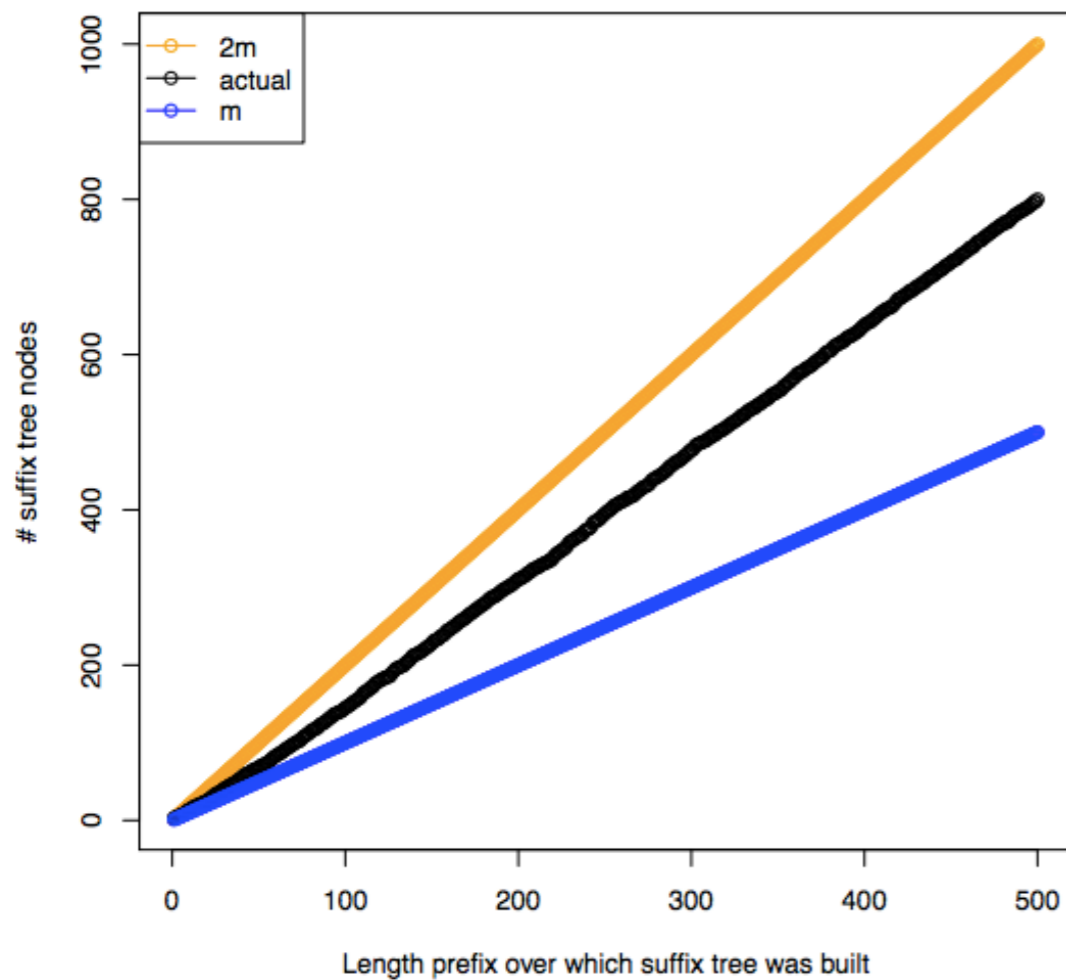


$T = \text{abaaba}\$$

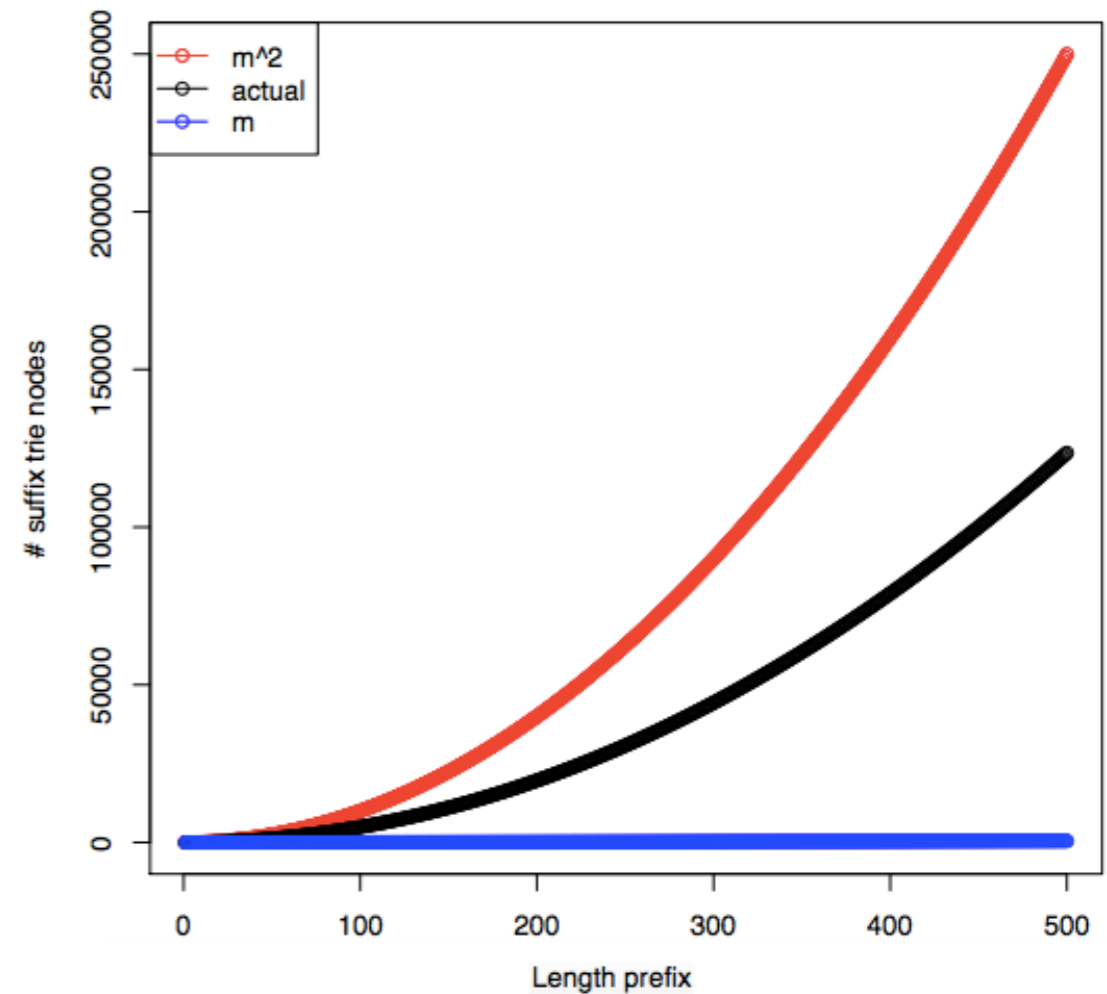


Actual growth: comparison

Trees built using the first 500 prefixes of the lambda phage virus genome



suffix tree



keyword tree

Summary

- Keyword and suffix trees are used to find patterns in a text
- Keyword trees:
 - Build keyword tree of patterns, and thread text through it
 - Usage: checking a set of patterns within various texts
- Suffix trees:
 - Build suffix tree of text, and thread patterns through it
 - Usage: checking various patterns in the same text