Suffix trees

Definition and applications

Recall tries

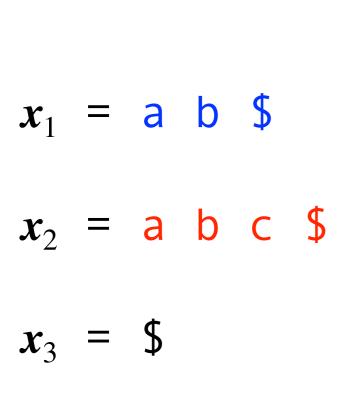
- We can search in a trie containing multiple strings.
 - As a function of the key-length what is the complexity?

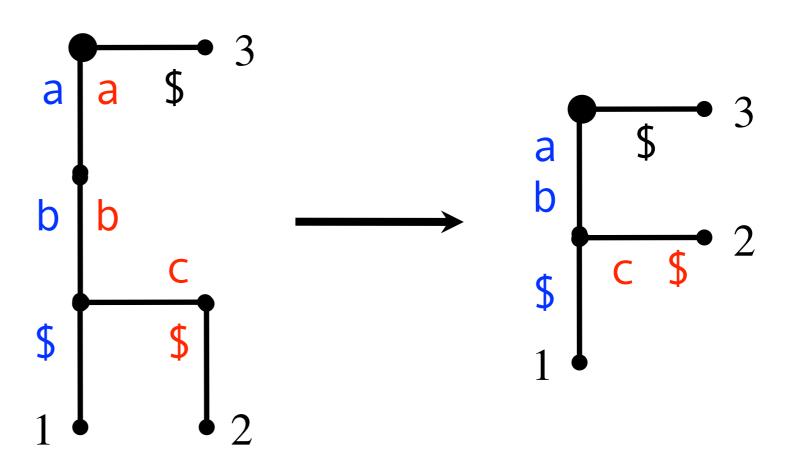
How would you use it to search for all occurrences of a key in a single string?

What is the space complexity?

Compacted tries

Saving space: Eliminate all internal nodes of degree 2 ...

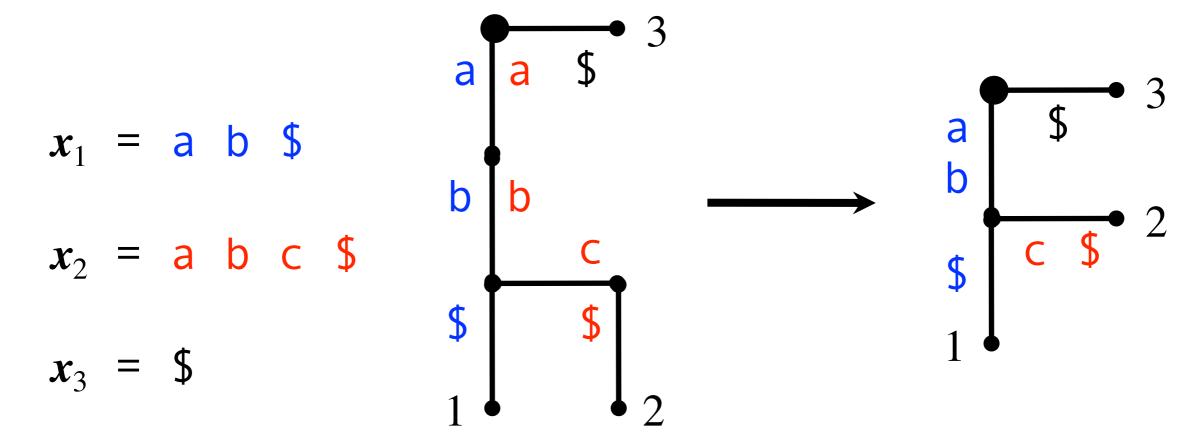




If we have n input-strings, then the trie has n+1 leaves and at most n internal nodes, i.e space O(n) for the tree. What about the labels?

Compacted tries

Saving space: Eliminate all internal nodes of degree 2 ...



If we have n input-strings, then the trie has n+1 leaves and at most n internal nodes, i.e space O(n) for the tree. What about the labels?

Labels can be represented in space O(1), i.e. "ab" \Rightarrow (1,0,2)

Suffix tree

The *suffix tree* T(x) of string x is the **compacted trie** of all suffixes x[i:n] for i = 0,...,n, i.e. including the empty suffix

Suffix tree

The *suffix tree* T(x) of string x is the **compacted trie** of all suffixes x[i:n] for i = 0,...,n, i.e. including the empty suffix

Example for x = tatat

```
tatat$
atat$
atat$
tat$

tat$

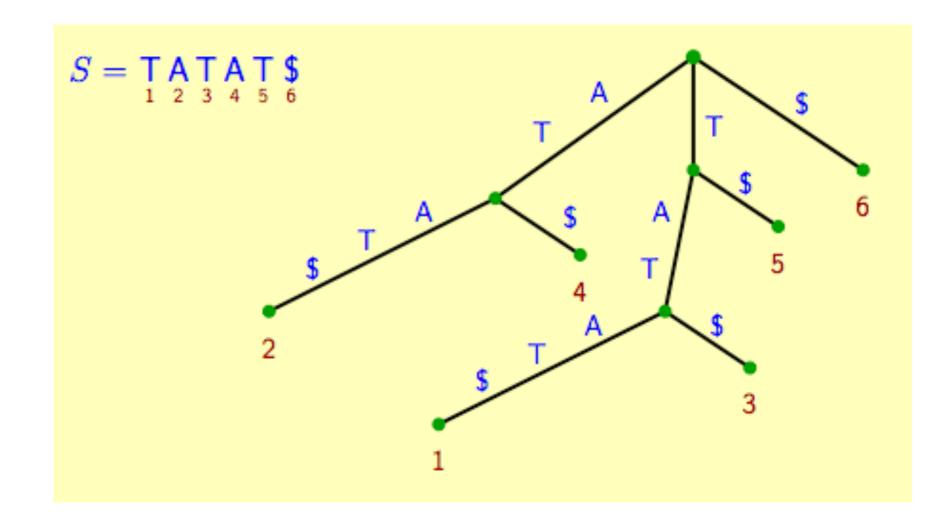
at$
```

Suffix tree

The *suffix tree* T(x) of string x is the **compacted trie** of all suffixes x[i:n] for i = 0,...,n, i.e. including the empty suffix

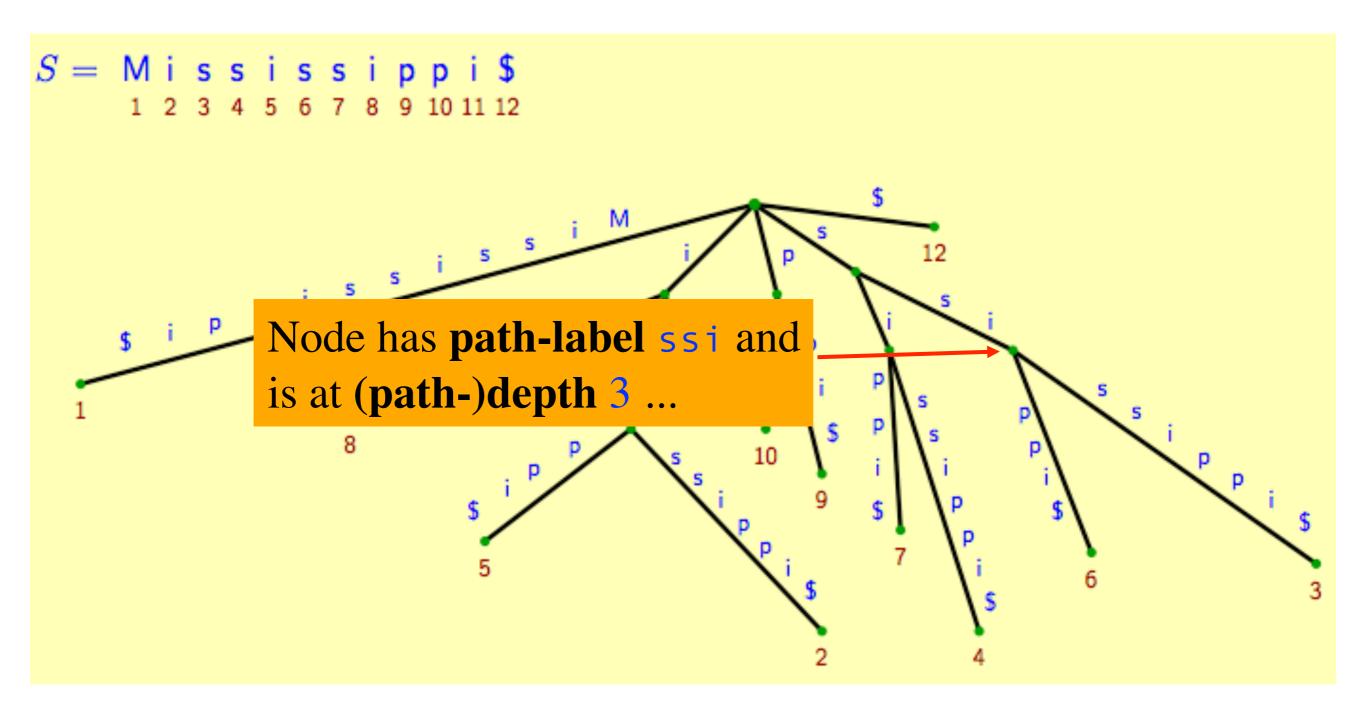
Example for x = tatat

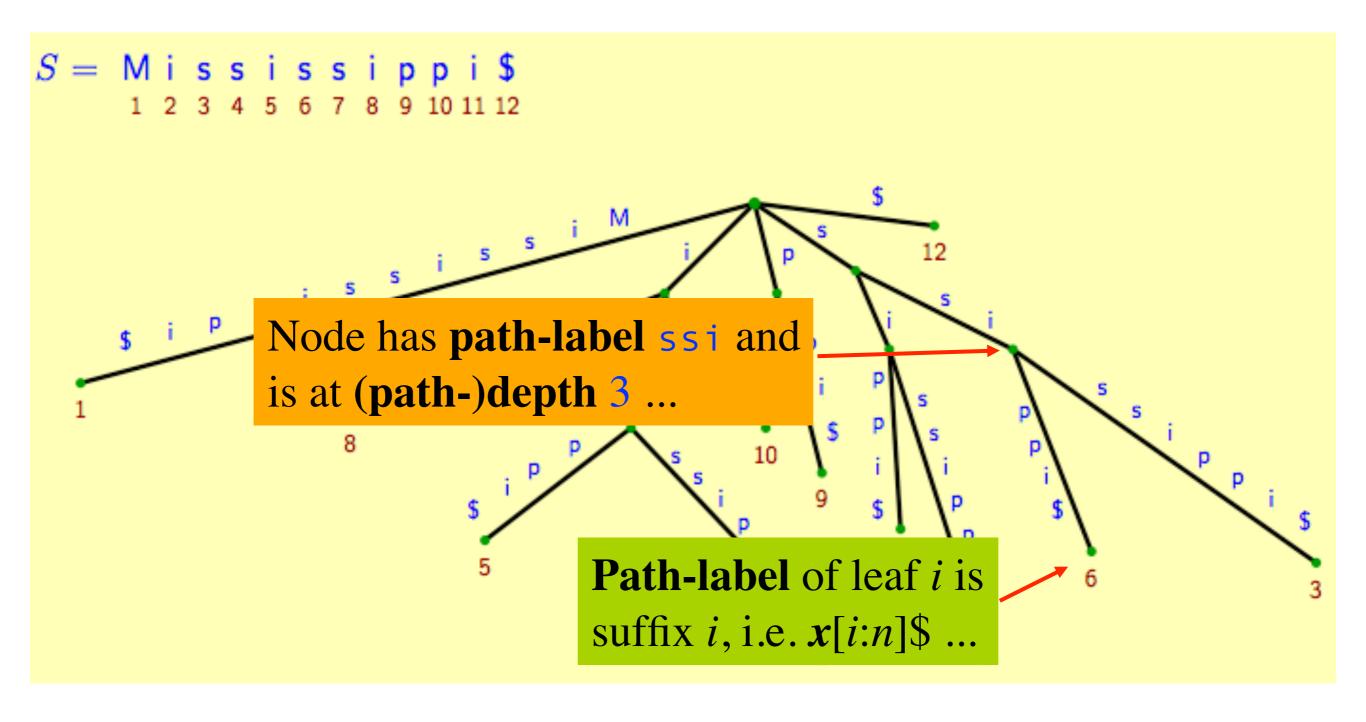
tatat\$
atat\$
atat\$
tat\$
tat\$
e\$

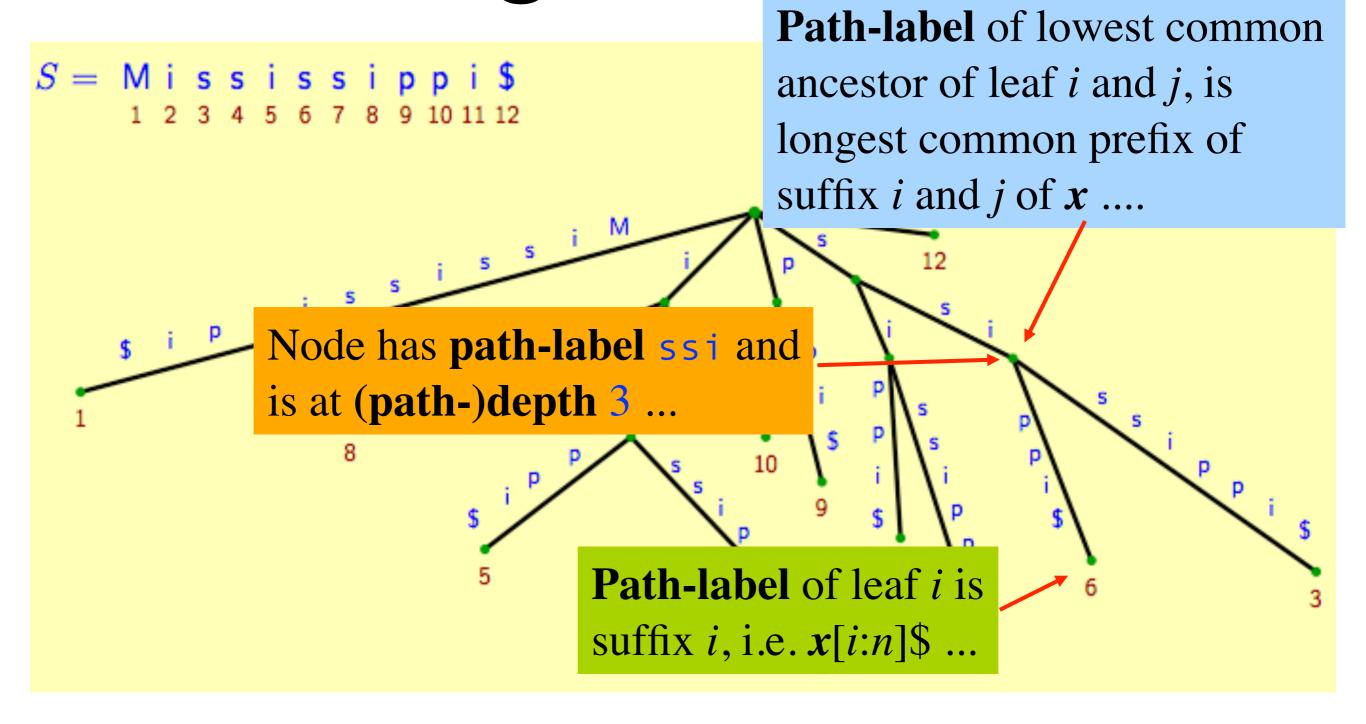


```
S = Mississippi
```

> gsa show suffixtree x | idot



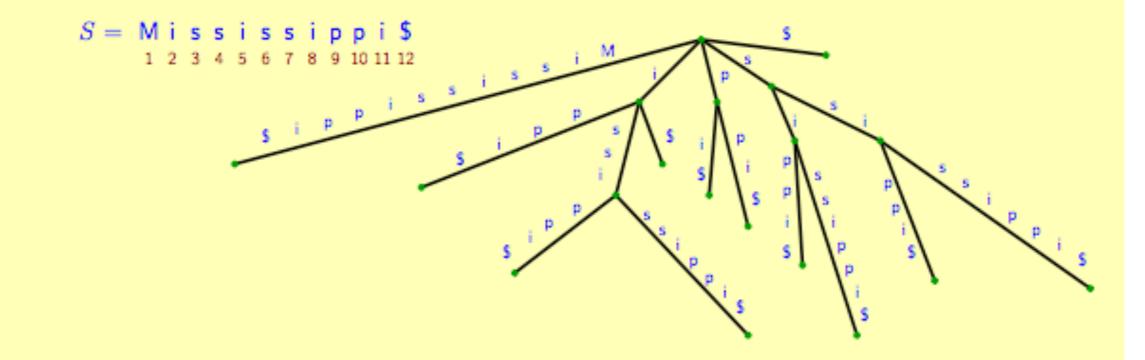




What is the space complexity?

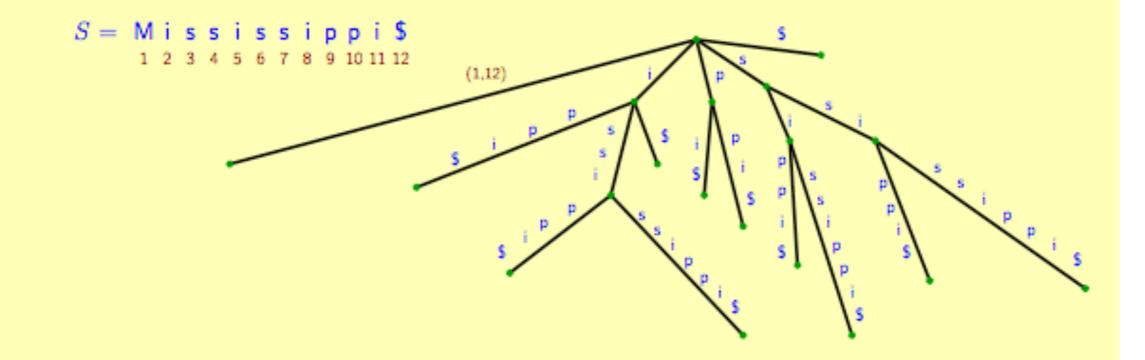
Observation: T(S) requires $\mathcal{O}(n)$ space.

- 1. T(S) has at most n leaves.
- 2. Each internal node is branching \Rightarrow at most n-1 internal nodes.
- 3. A tree with at most 2n-1 nodes has at most 2n-2 edges.
- 4. Each node requires constant space.
- 5. Each edge label is a substring of $S \Rightarrow$ pair of pointers (i, j) into S.



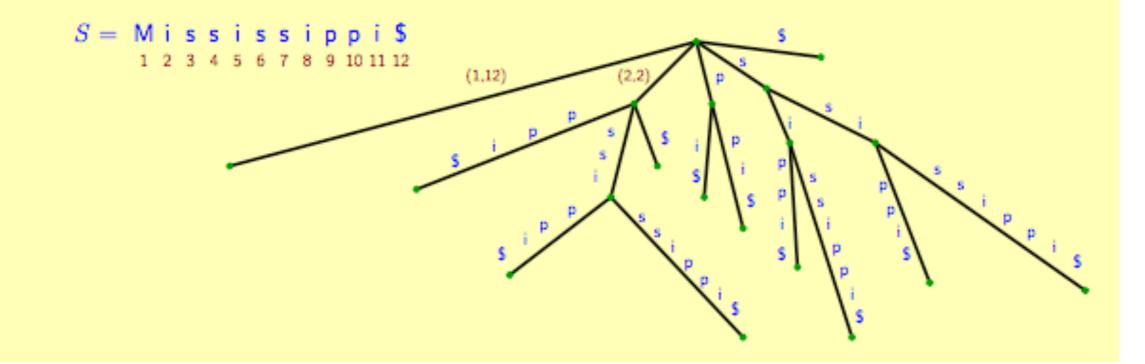
Observation: T(S) requires $\mathcal{O}(n)$ space.

- 1. T(S) has at most n leaves.
- 2. Each internal node is branching \Rightarrow at most n-1 internal nodes.
- 3. A tree with at most 2n-1 nodes has at most 2n-2 edges.
- 4. Each node requires constant space.
- 5. Each edge label is a substring of $S \Rightarrow$ pair of pointers (i, j) into S.



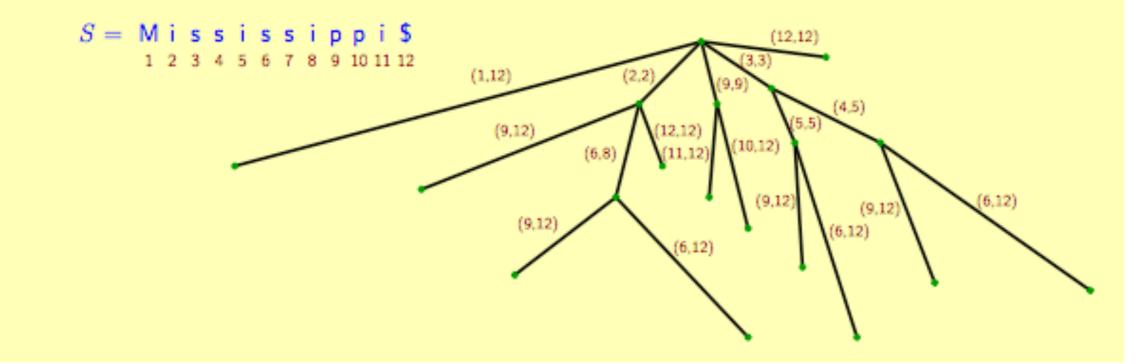
Observation: T(S) requires $\mathcal{O}(n)$ space.

- 1. T(S) has at most n leaves.
- 2. Each internal node is branching \Rightarrow at most n-1 internal nodes.
- 3. A tree with at most 2n-1 nodes has at most 2n-2 edges.
- 4. Each node requires constant space.
- 5. Each edge label is a substring of $S \Rightarrow \text{pair of pointers } (i, j) \text{ into } S$.



Observation: T(S) requires $\mathcal{O}(n)$ space.

- 1. T(S) has at most n leaves.
- 2. Each internal node is branching \Rightarrow at most n-1 internal nodes.
- 3. A tree with at most 2n-1 nodes has at most 2n-2 edges.
- 4. Each node requires constant space.
- 5. Each edge label is a substring of $S \Rightarrow$ pair of pointers (i, j) into S.



Constructing suffix trees

Constructing T(x) by inserting each suffix one by one takes time $O(n^2)$

Can we do better?

Constructing suffix trees

Constructing T(x) by inserting each suffix one by one takes time $O(n^2)$

Can we do better?

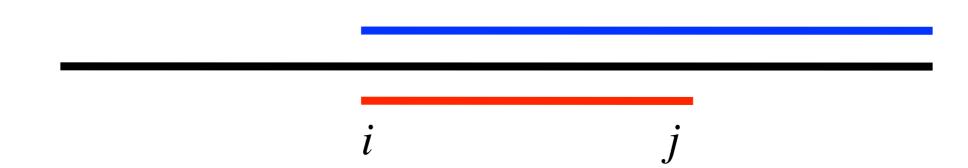
We will see a linear time construction algorithm next week.

What about applications?

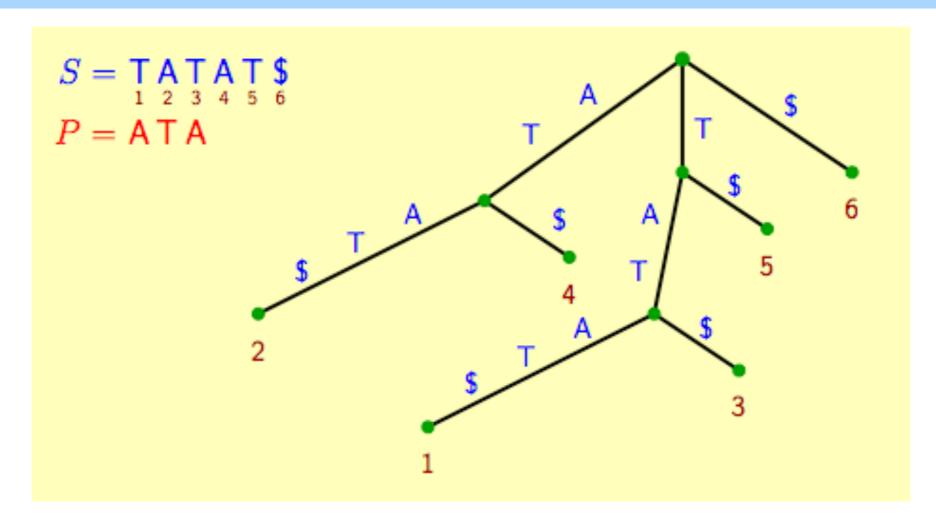
... exact matching, finding repeats, longest common substring ...

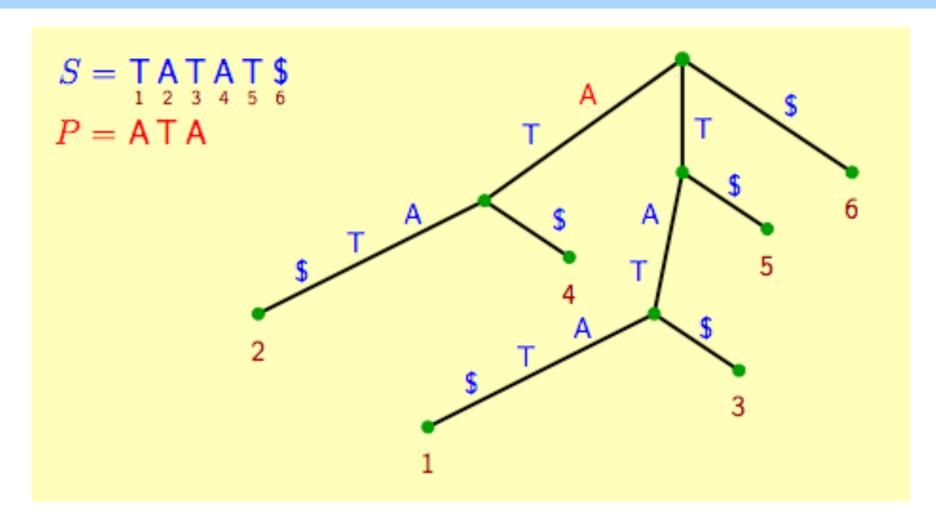
Given string x and pattern y, report where y occurs in x

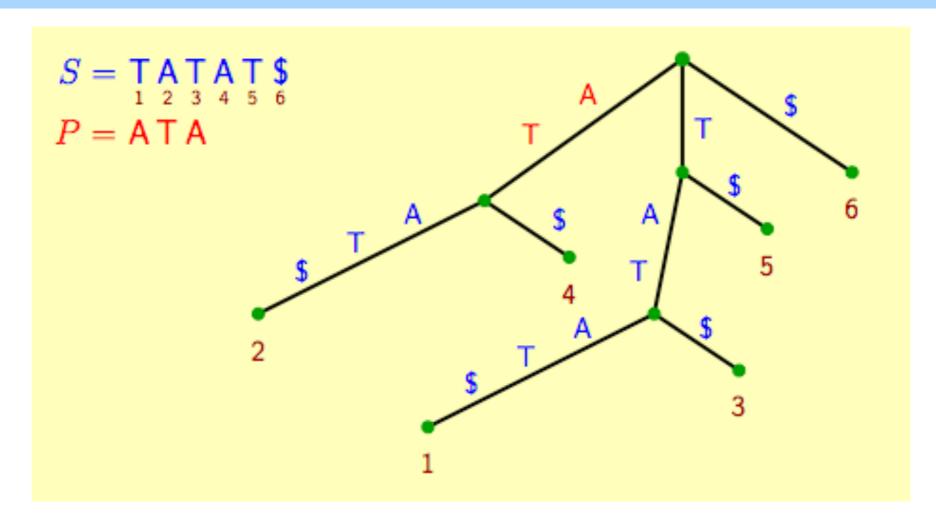
If y occurs in x at position i, then y is a prefix of suffix i of x

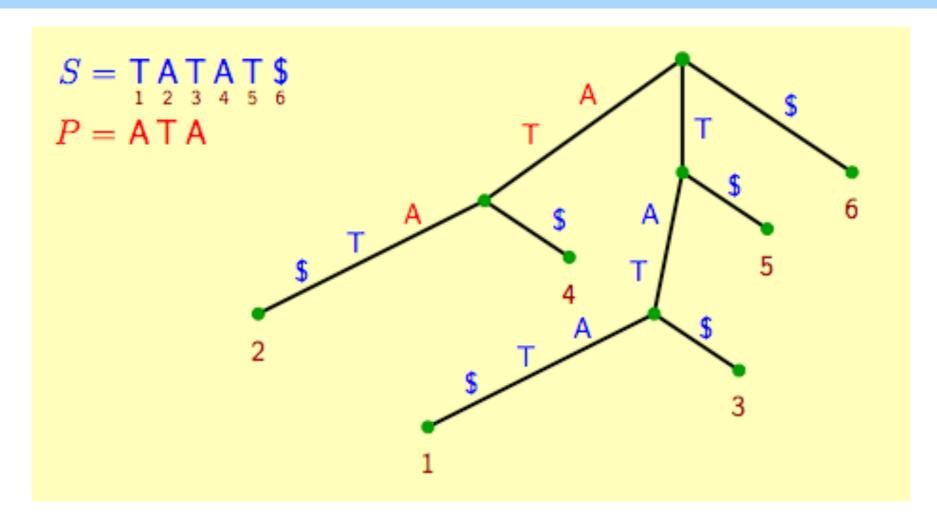


y is spelled by an initial part of the path from the root to leaf i in T(x)

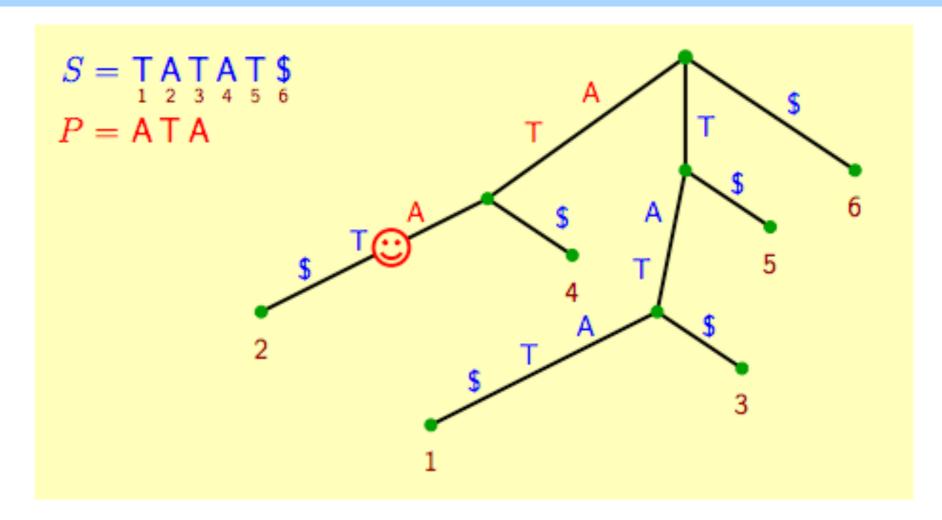






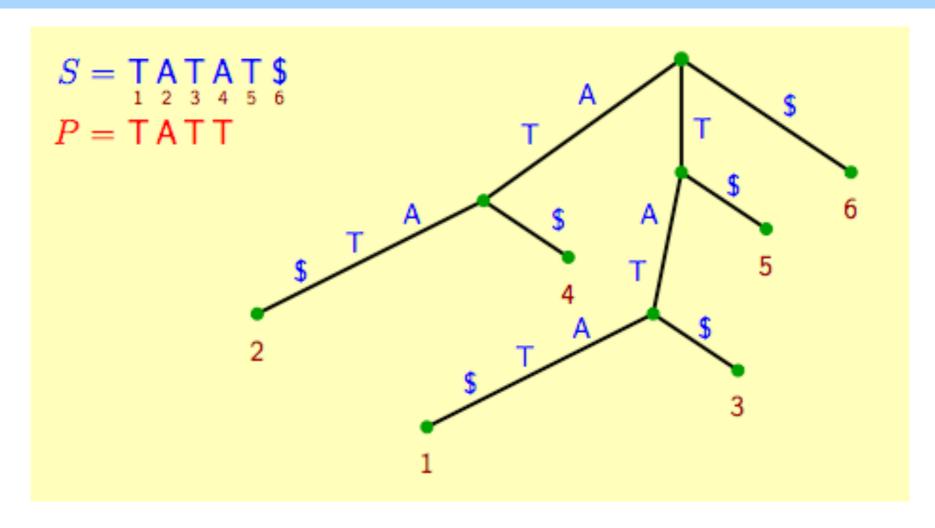


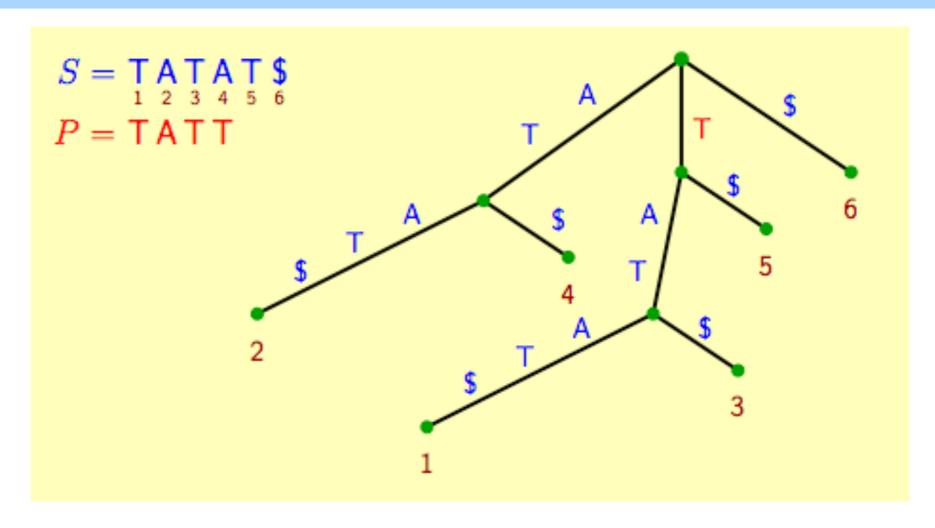
Given string x and pattern y, report where y occurs in x

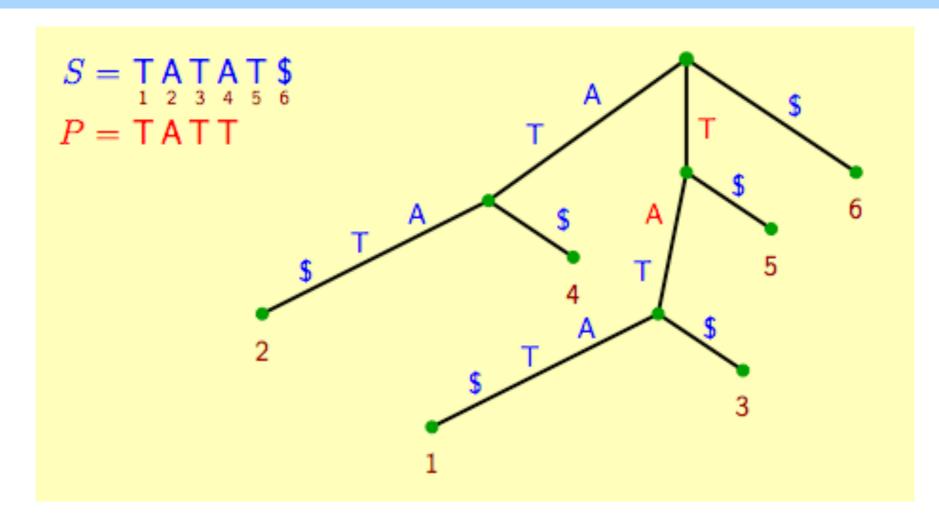


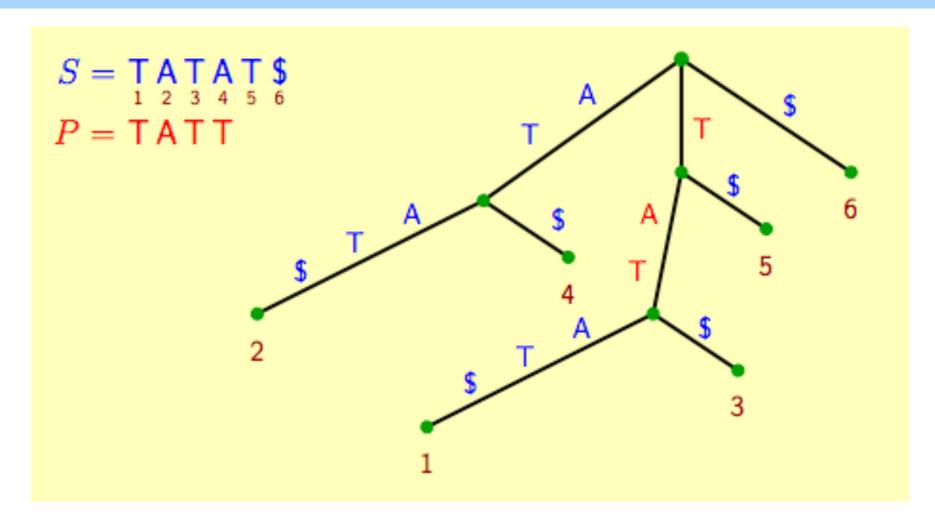
Pattern at a occurs at position 2 in tatat

Time: O(m) using the suffix tree T(S)

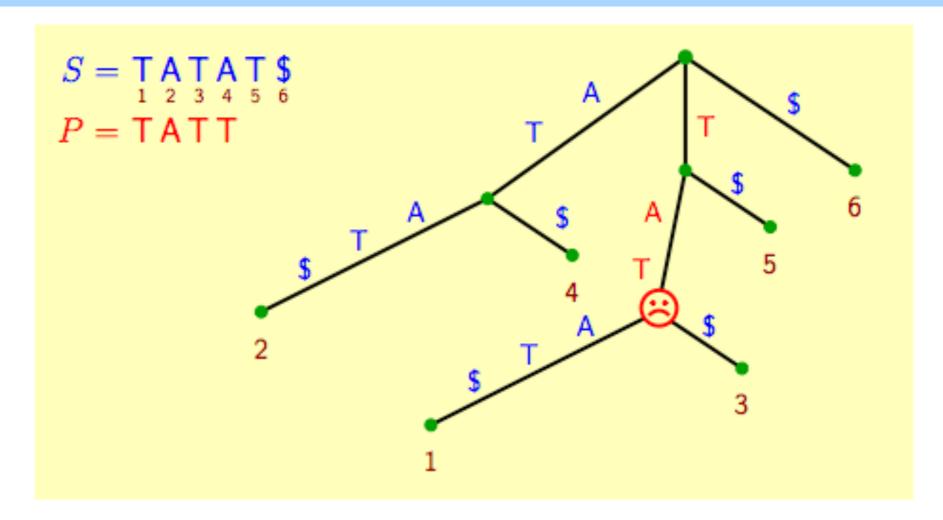








Given string x and pattern y, report where y occurs in x

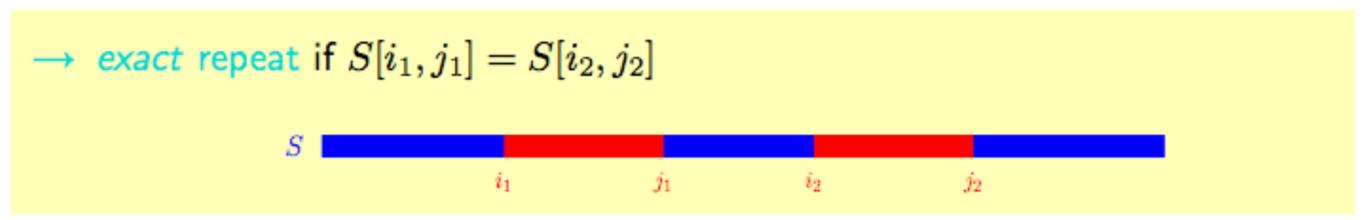


Pattern tatt does not occur in tatat

Time: O(m) using the suffix tree T(S)

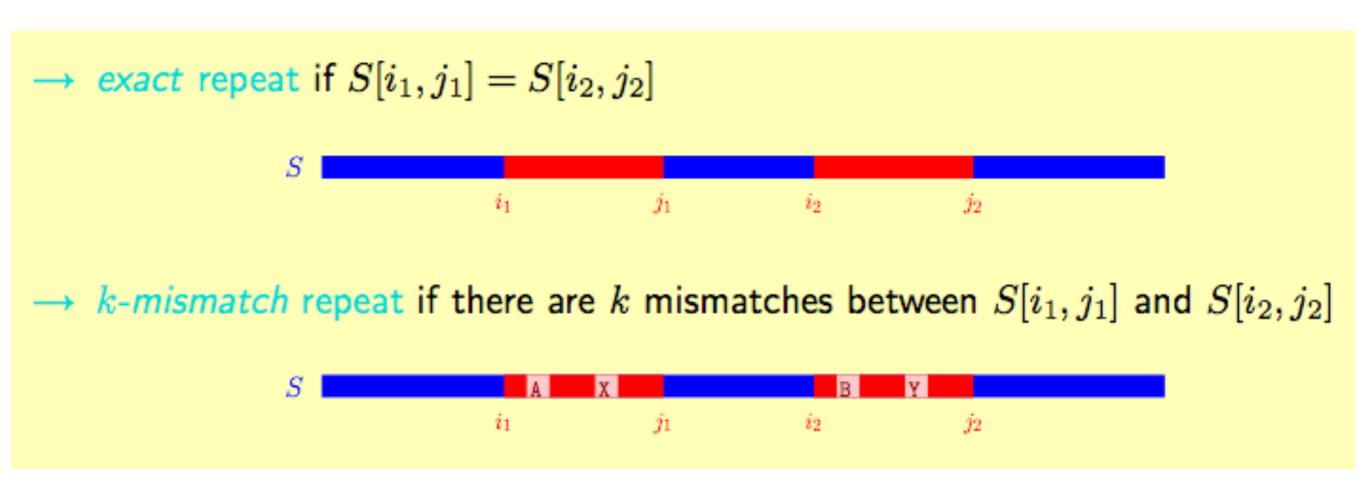
Repeats

A pair of substrings $R=(S[i_1..j_1], S[i_2..j_2])$ is a ...



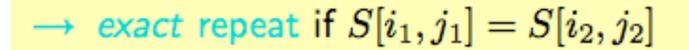
Repeats

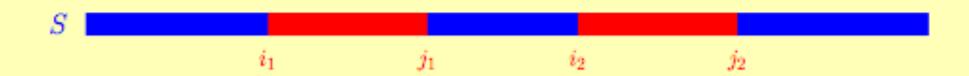
A pair of substrings $R = (S[i_1..j_1], S[i_2..j_2])$ is a ...



Repeats

A pair of substrings $R=(S[i_1..j_1], S[i_2..j_2])$ is a ...





ightarrow k-mismatch repeat if there are k mismatches between $S[i_1,j_1]$ and $S[i_2,j_2]$



 \rightarrow k-differences repeat if there are k differences (mismatches, insertions, deletions) between $S[i_1,j_1]$ and $S[i_2,j_2]$



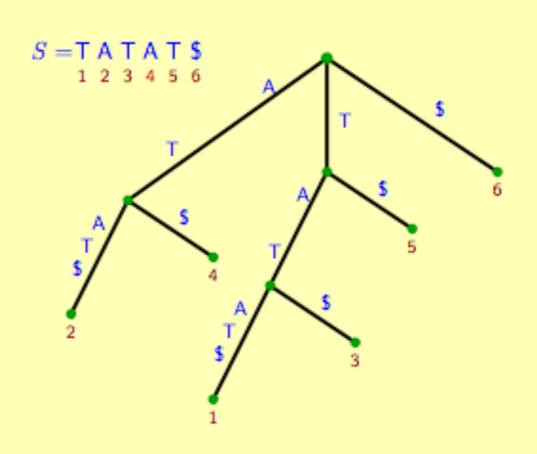
Finding exact repeats

Folklore: (see e.g. Gusfield, 1997)

It is possible to find all pairs of repeated substrings (repeats) in S in linear time.

Idea:

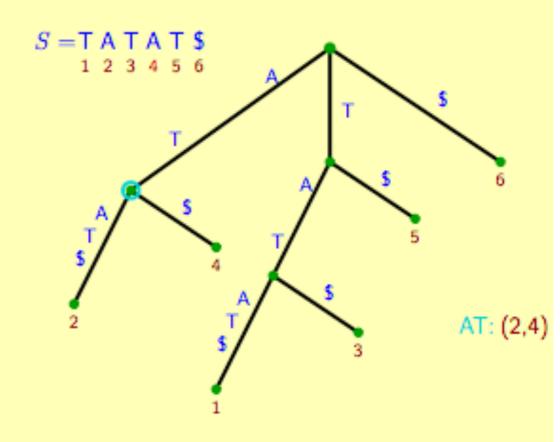
- consider string S and its suffix tree T(S).
- repeated substrings of S correspond to internal locations in T(S).
- leaf numbers tell us positions where substrings occur.



Folklore: (see e.g. Gusfield, 1997)

It is possible to find all pairs of repeated substrings (repeats) in S in linear time.

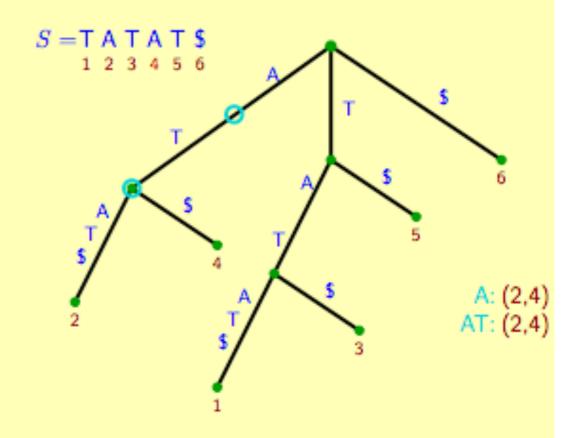
- consider string S and its suffix tree T(S).
- repeated substrings of S correspond to internal locations in T(S).
- leaf numbers tell us positions where substrings occur.



Folklore: (see e.g. Gusfield, 1997)

 \bullet It is possible to find all pairs of repeated substrings (repeats) in S in linear time.

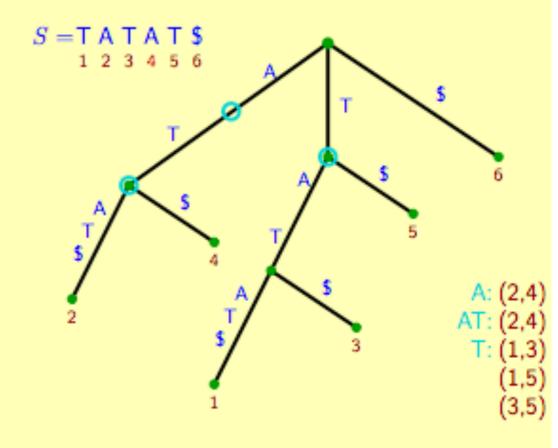
- consider string S and its suffix tree T(S).
- repeated substrings of S correspond to internal locations in T(S).
- leaf numbers tell us positions where substrings occur.



Folklore: (see e.g. Gusfield, 1997)

 \bullet It is possible to find all pairs of repeated substrings (repeats) in S in linear time.

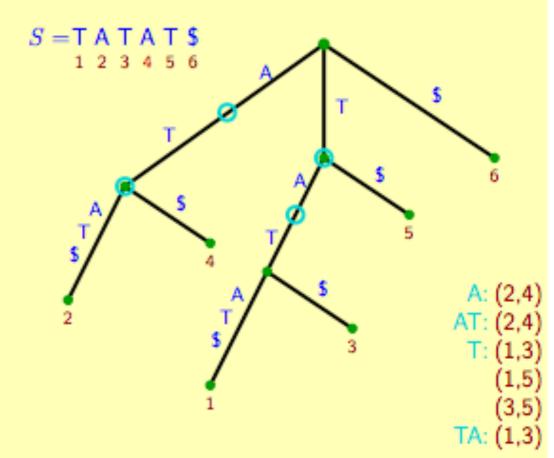
- consider string S and its suffix tree T(S).
- repeated substrings of S correspond to internal locations in T(S).
- leaf numbers tell us positions where substrings occur.



Folklore: (see e.g. Gusfield, 1997)

 \bullet It is possible to find all pairs of repeated substrings (repeats) in S in linear time.

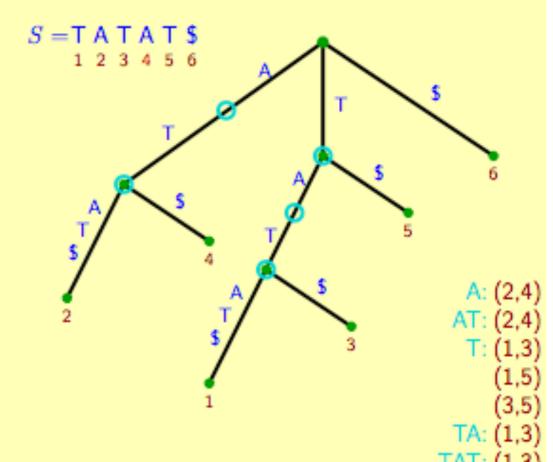
- consider string S and its suffix tree T(S).
- repeated substrings of S correspond to internal locations in T(S).
- leaf numbers tell us positions where substrings occur.



Folklore: (see e.g. Gusfield, 1997)

 \bullet It is possible to find all pairs of repeated substrings (repeats) in S in linear time.

- consider string S and its suffix tree T(S).
- repeated substrings of S correspond to internal locations in T(S).
- leaf numbers tell us positions where substrings occur.

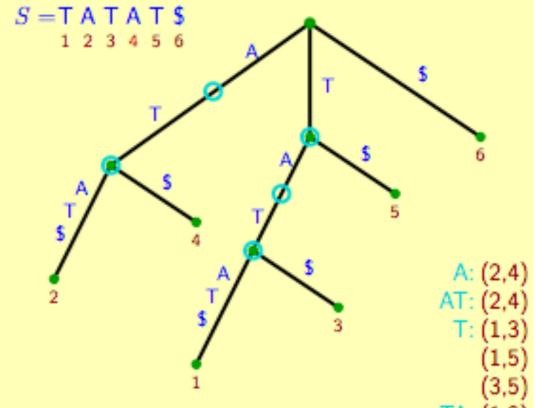


Folklore: (see e.g. Gusfield, 1997)

 \bullet It is possible to find all pairs of repeated substrings (repeats) in S in linear time.

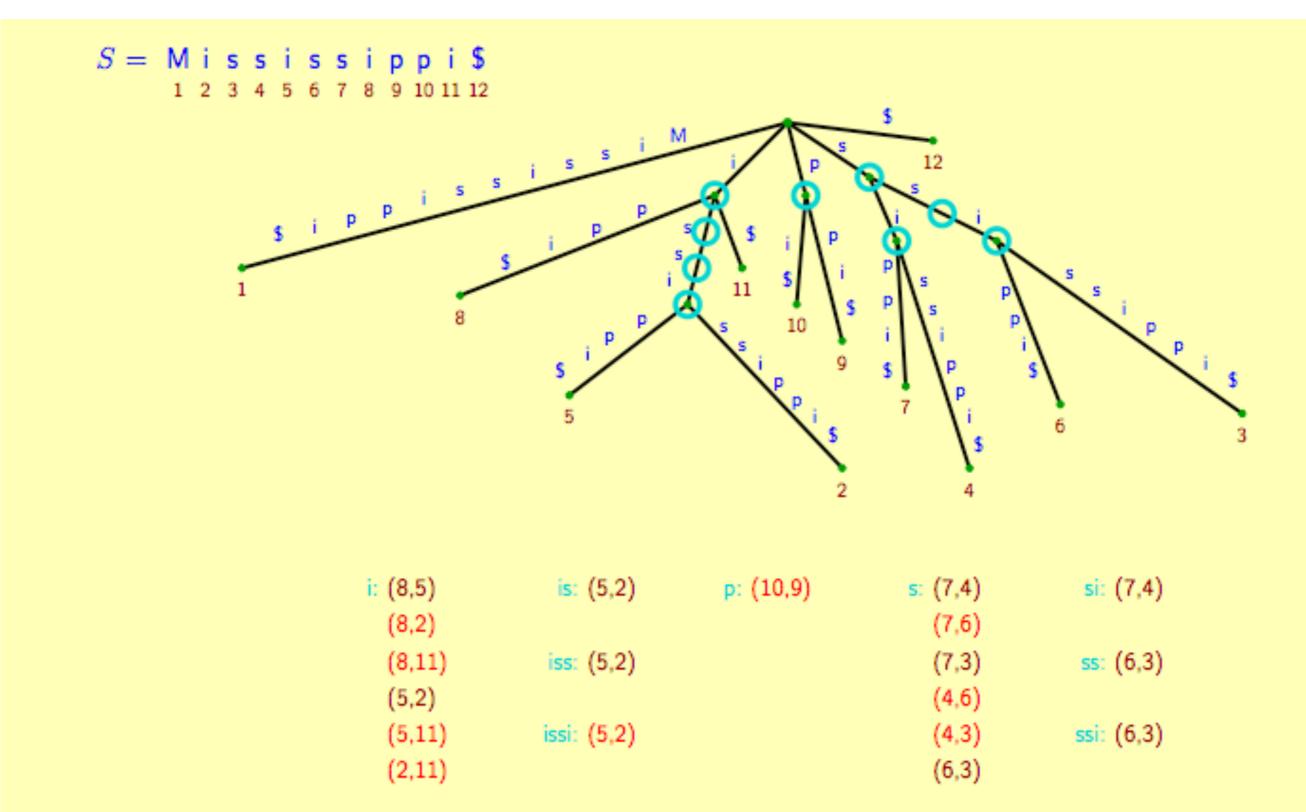
Idea:

- consider string S and its suffix tree T(S).
- repeated substrings of S correspond to internal locations in T(S).
- leaf numbers tell us positions where substrings occur.



Analysis: $\mathcal{O}(n+z)$ time with $z=|\mathsf{output}|$, $\mathcal{O}(n)$ space

A larger example

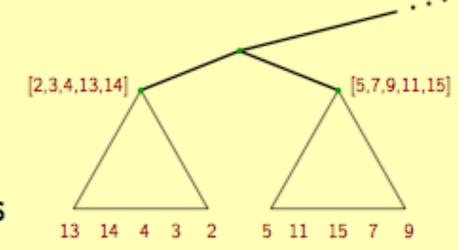








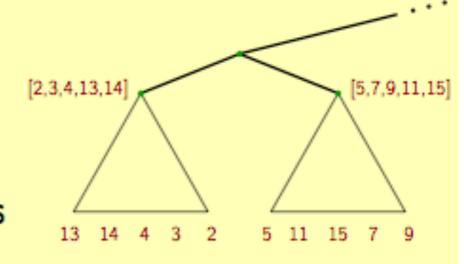
- For right-maximality $(X \neq Y)$
 - consider only internal nodes of T(S)
 - report only pairs of leaves from different subtrees (or from different leaf-lists)



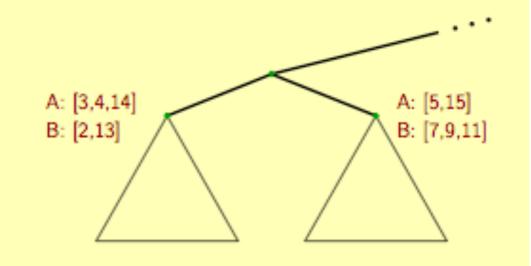


Idea:

- For right-maximality $(X \neq Y)$
 - consider only internal nodes of T(S)
 - report only pairs of leaves from different subtrees (or from different leaf-lists)



- For left-maximality $(A \neq B)$
 - keep lists for the different left-characters
 - report only pairs from different lists



Analysis: $\mathcal{O}(n+z)$ time with z = |output|, $\mathcal{O}(n)$ space

Other repeats

Maximal repeats with bounded gap in time $O(n \log n + z)$



Tandem repeats in time $O(n \log n + z)$



Palindromic repeats in O(n + z)



... all using suffix trees ...

The *longest common substring* of x and y is the longest string z which occurs in both x and y ...

Can this be found efficiently using a suffix tree?

The *longest common substring* of x and y is the longest string z which occurs in both x and y ...

Can this be found efficiently using a suffix tree?

z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]

z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]

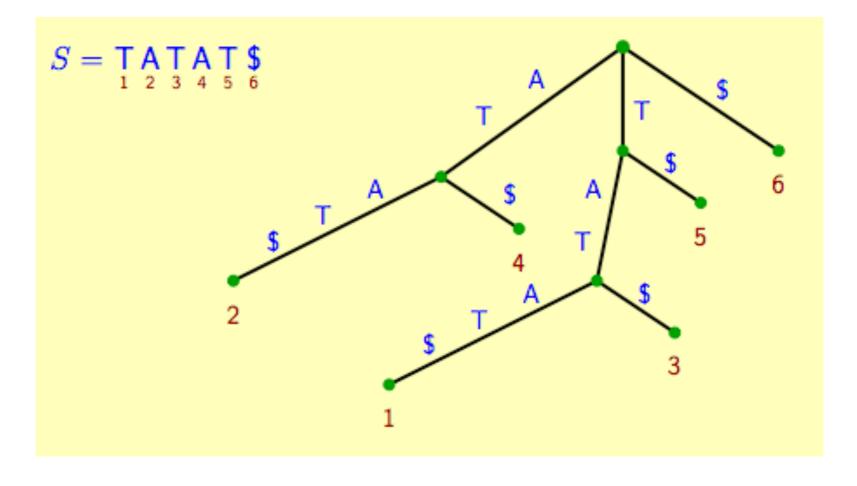
Idea: Build a compacted trie of all suffixes of x and y, such that each suffix of x and y corresponds to unique root-to-leaf paths ...

tatat\$
atat\$
atat\$
tat\$

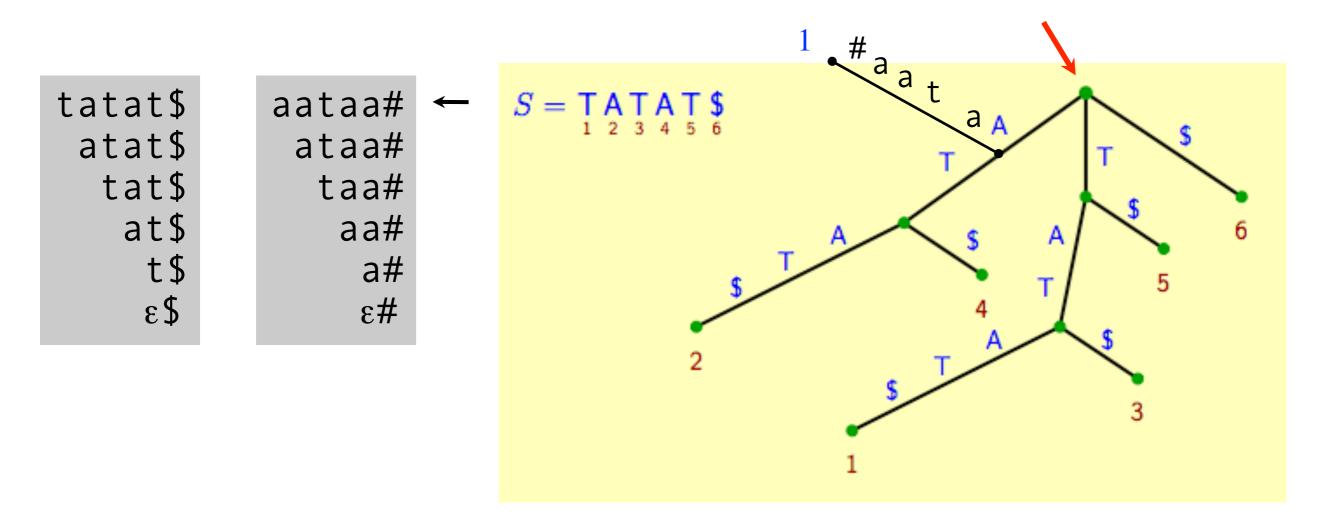
tat\$

at\$

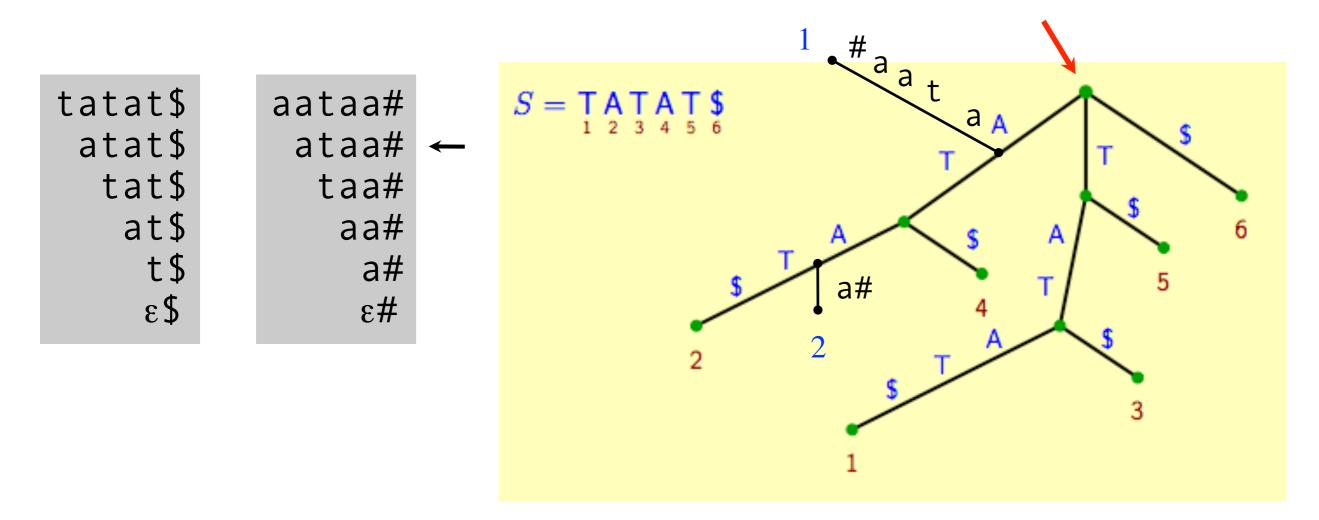
aataa# ataa# taa# aa# a# ε#



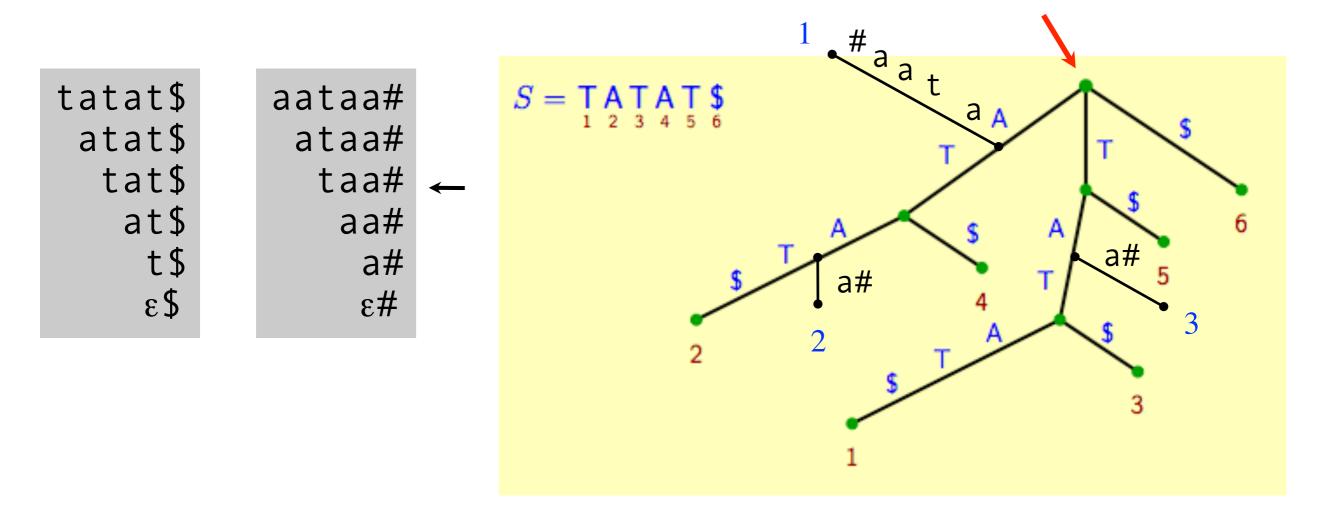
z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]



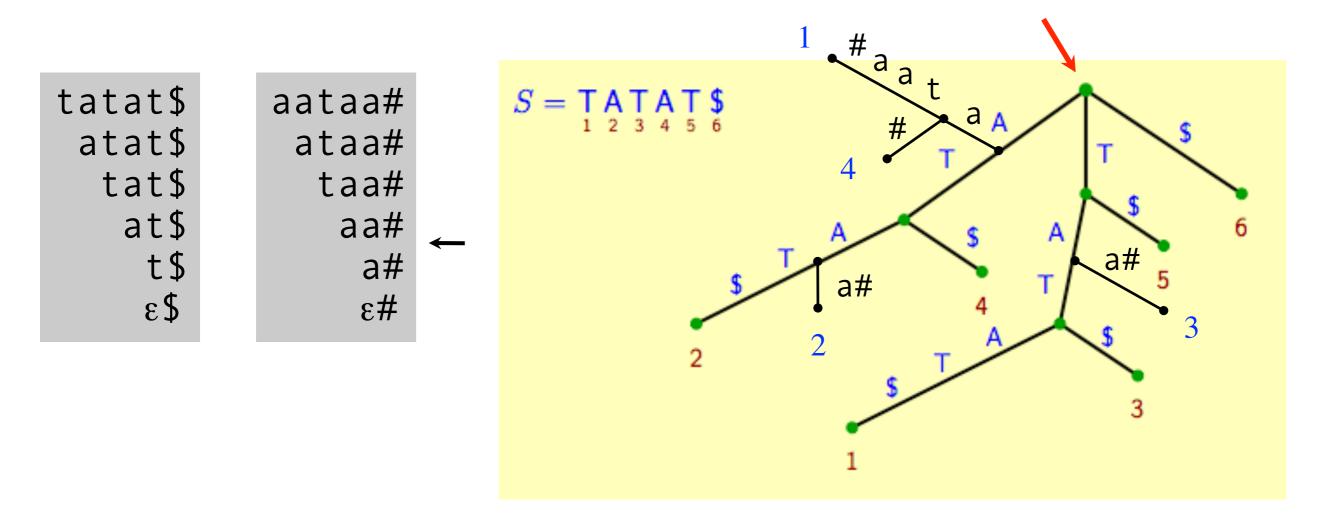
z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]



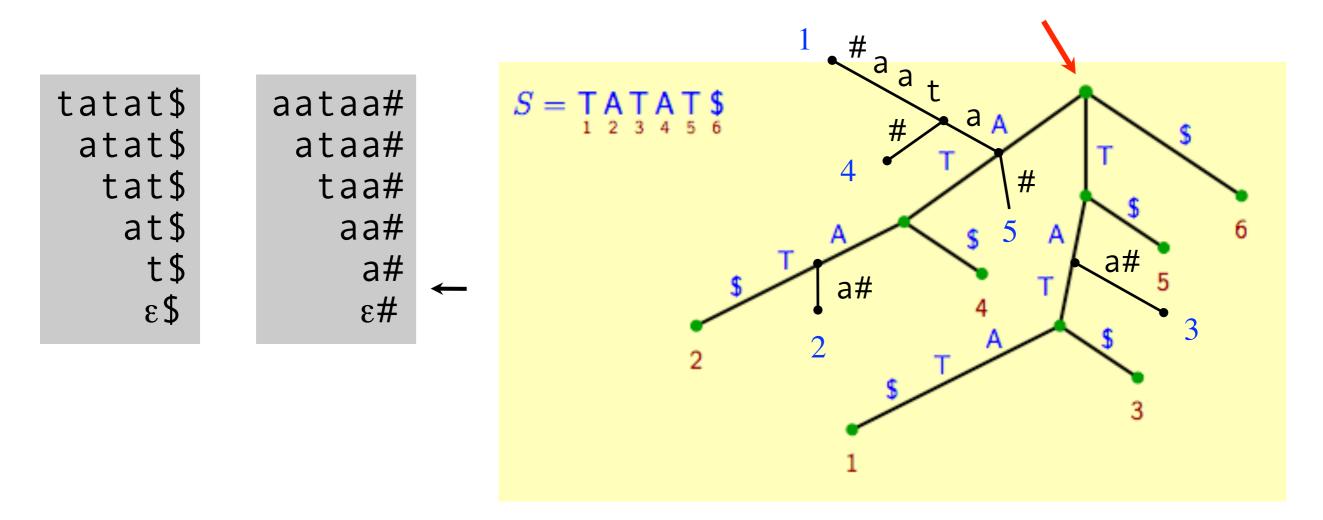
z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]



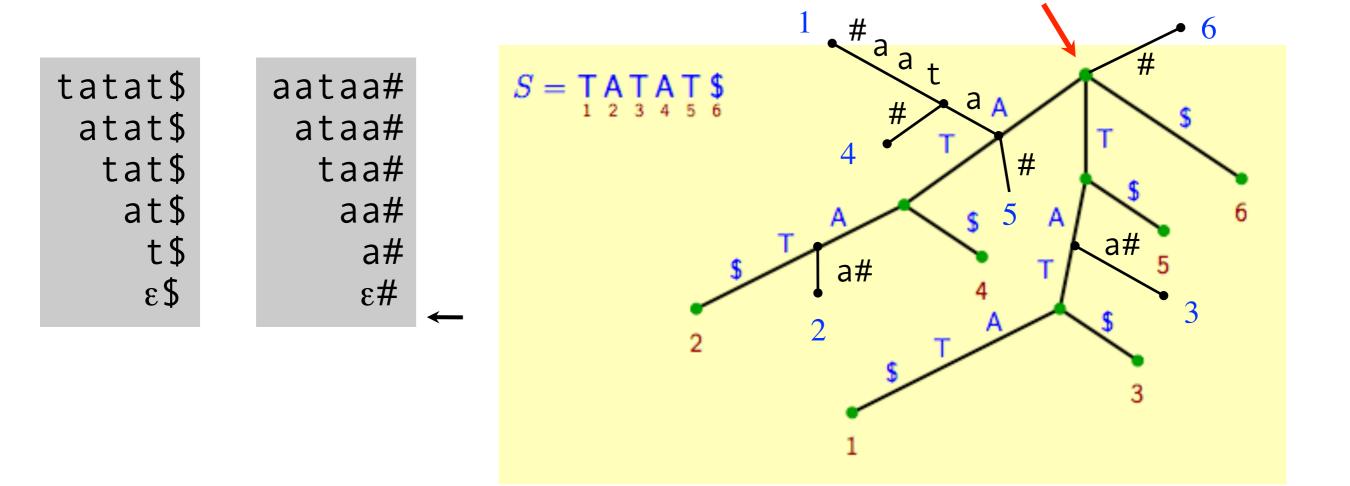
z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]



z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]



z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]



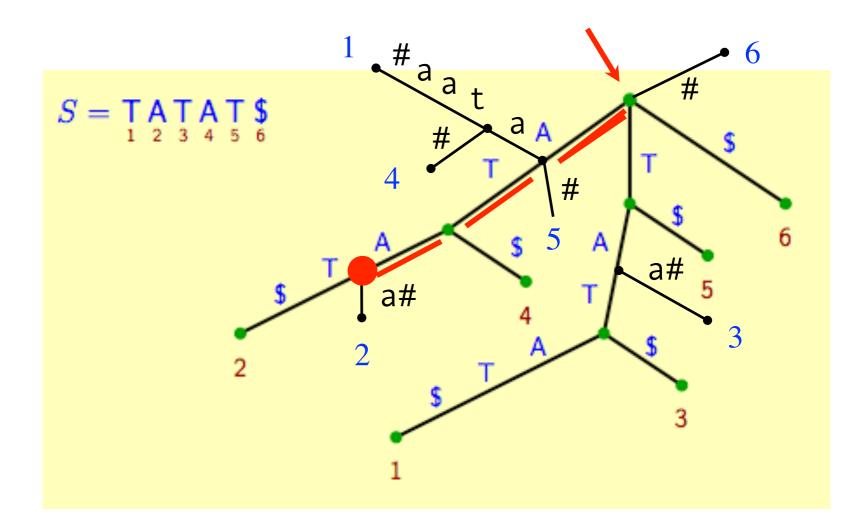
z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]

tatat\$
atat\$
tat\$
tat\$

tat\$

\$
\$
\$
\$
\$

aataa#
ataa#
taa#
aa#
aa#
a#
e#



Observe: z is the path-label of the deepest node with suffixes from both x and y as leaves in its sub-tree ...

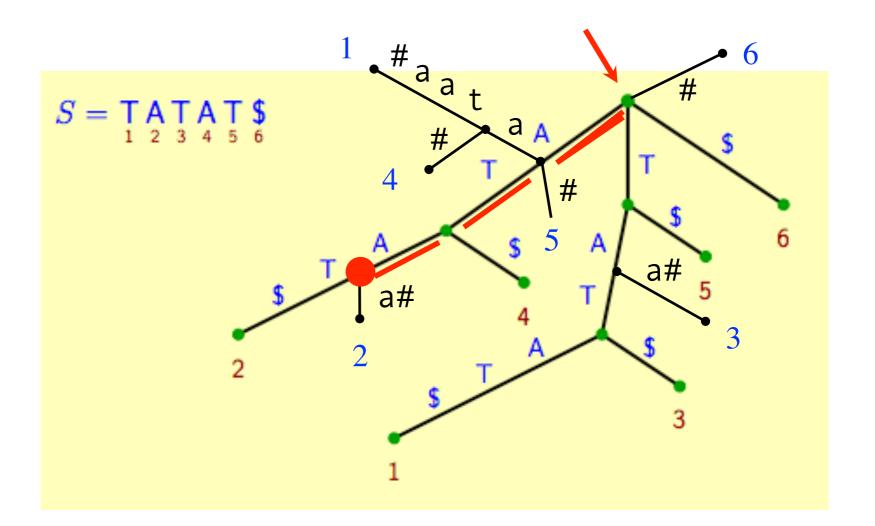
z is the longest common prefix of any pair of suffixes x[i:n] and y[j:m]

tatat\$
atat\$
tat\$
tat\$

tat\$

ε\$

aataa#
ataa#
taa#
aa#
aa#
a#
e#



Observe: z is the path-label of the deepest node with suffixes from both x and y as leaves in its sub-tree ... **Time**: O(n+m)

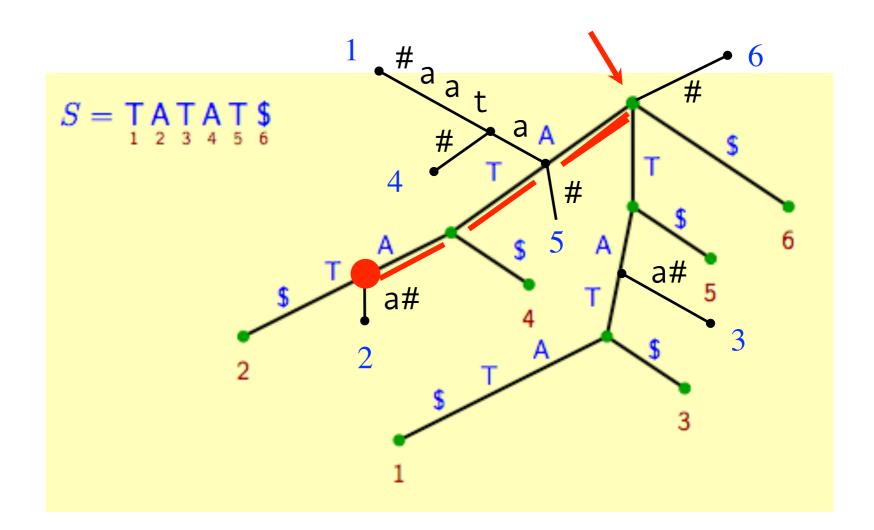
This is the generalized suffix tree of tatat and aataa

tatat\$
atat\$
tat\$
tat\$

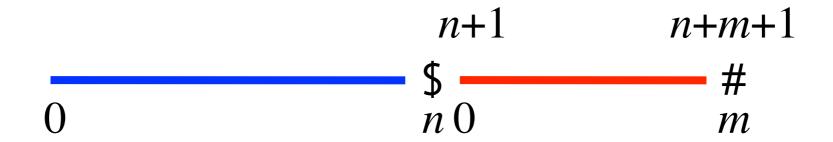
tat\$

tat\$

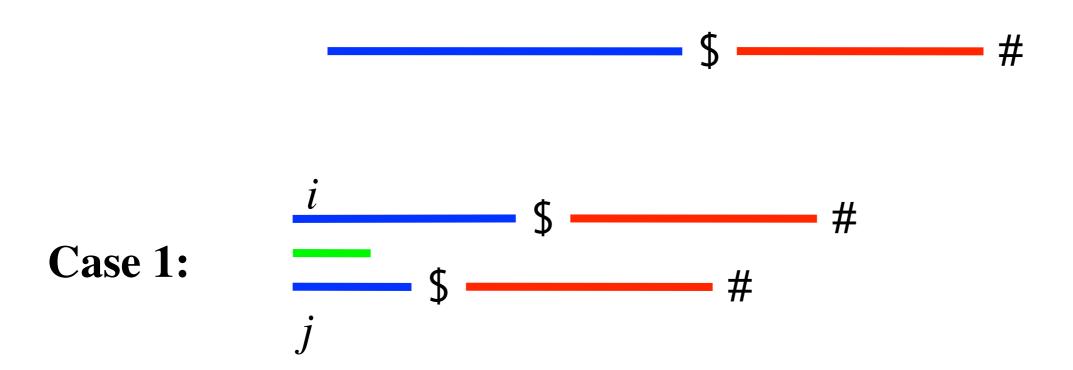
aataa#
ataa#
taa#
aa#
aa#
a#
e#

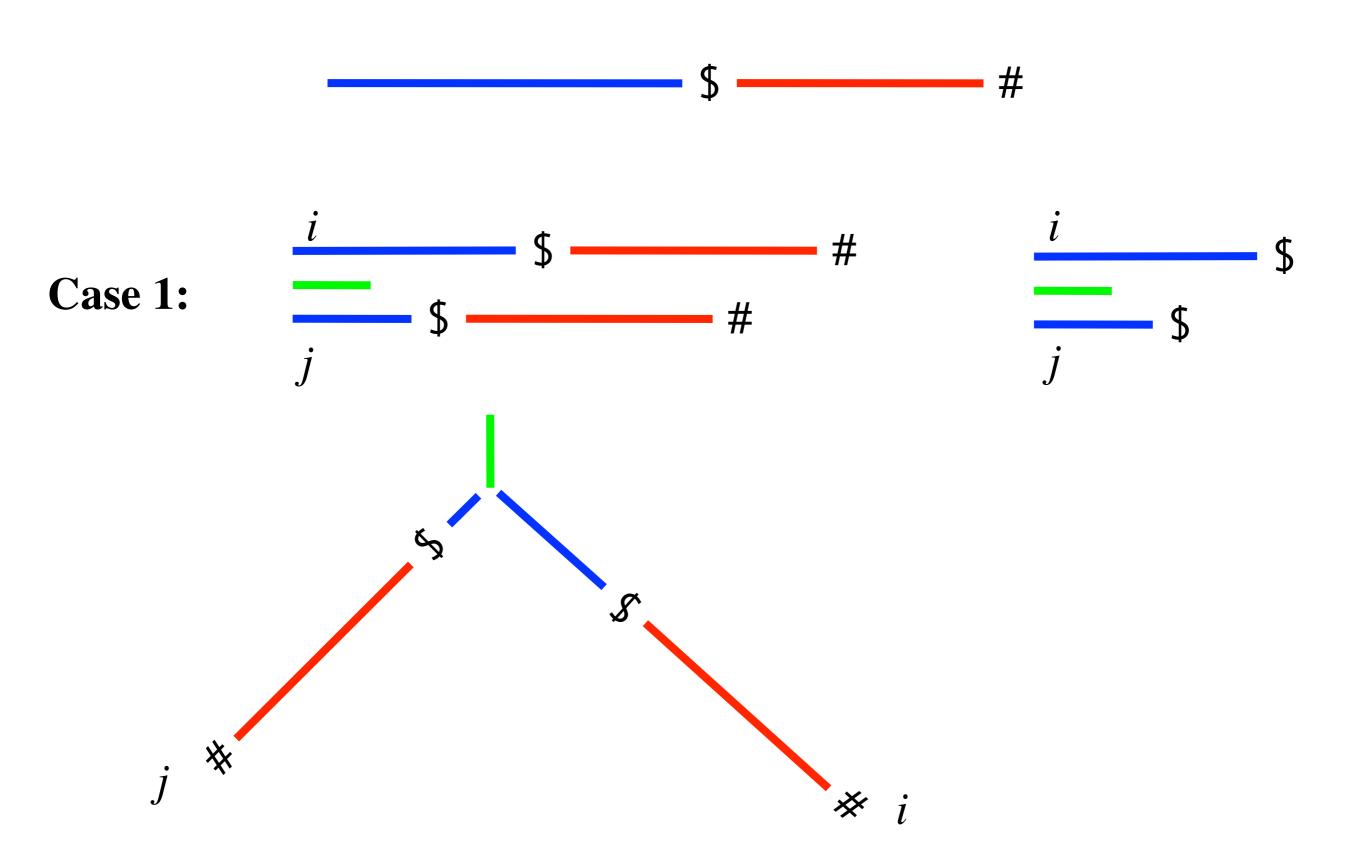


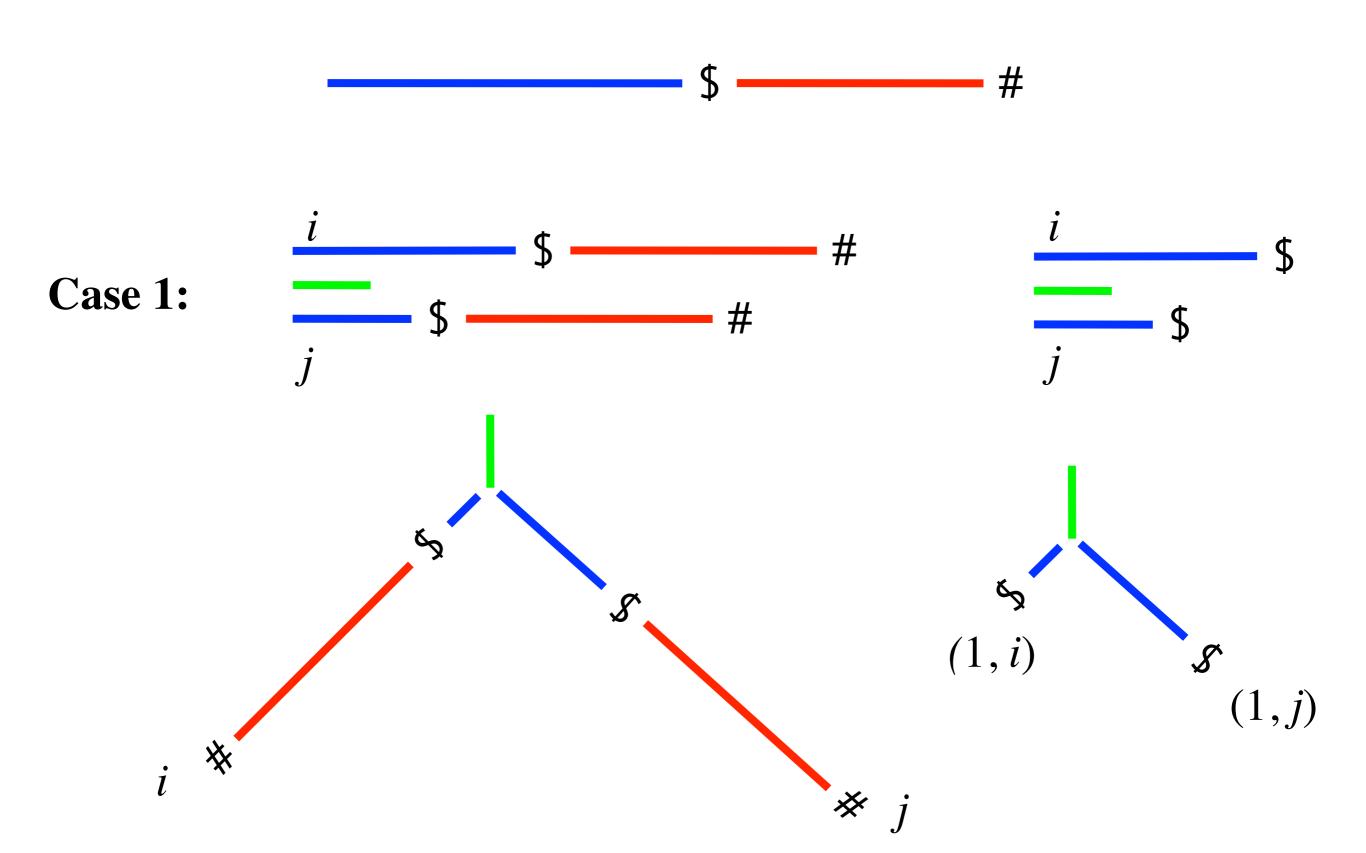
Can be constructed by constructing the suffix tree of ... tatat\$aataa#

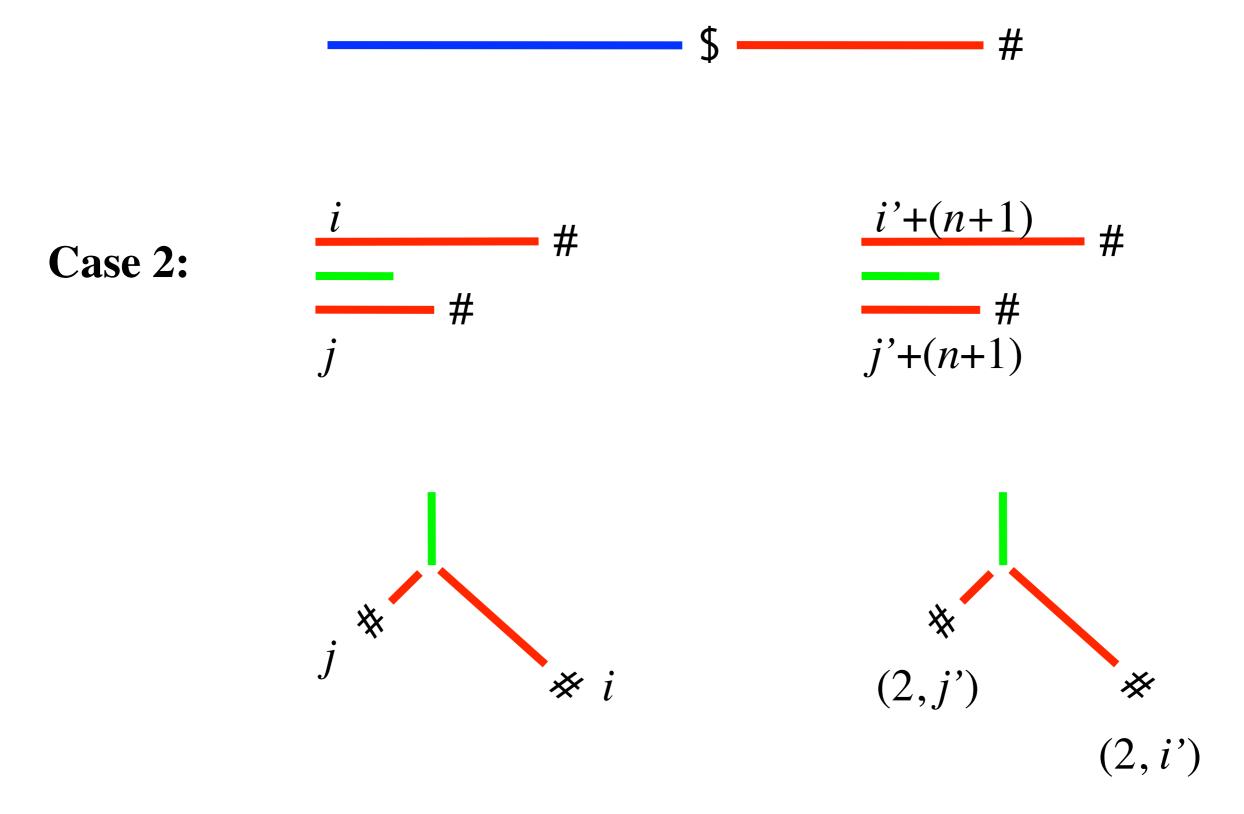


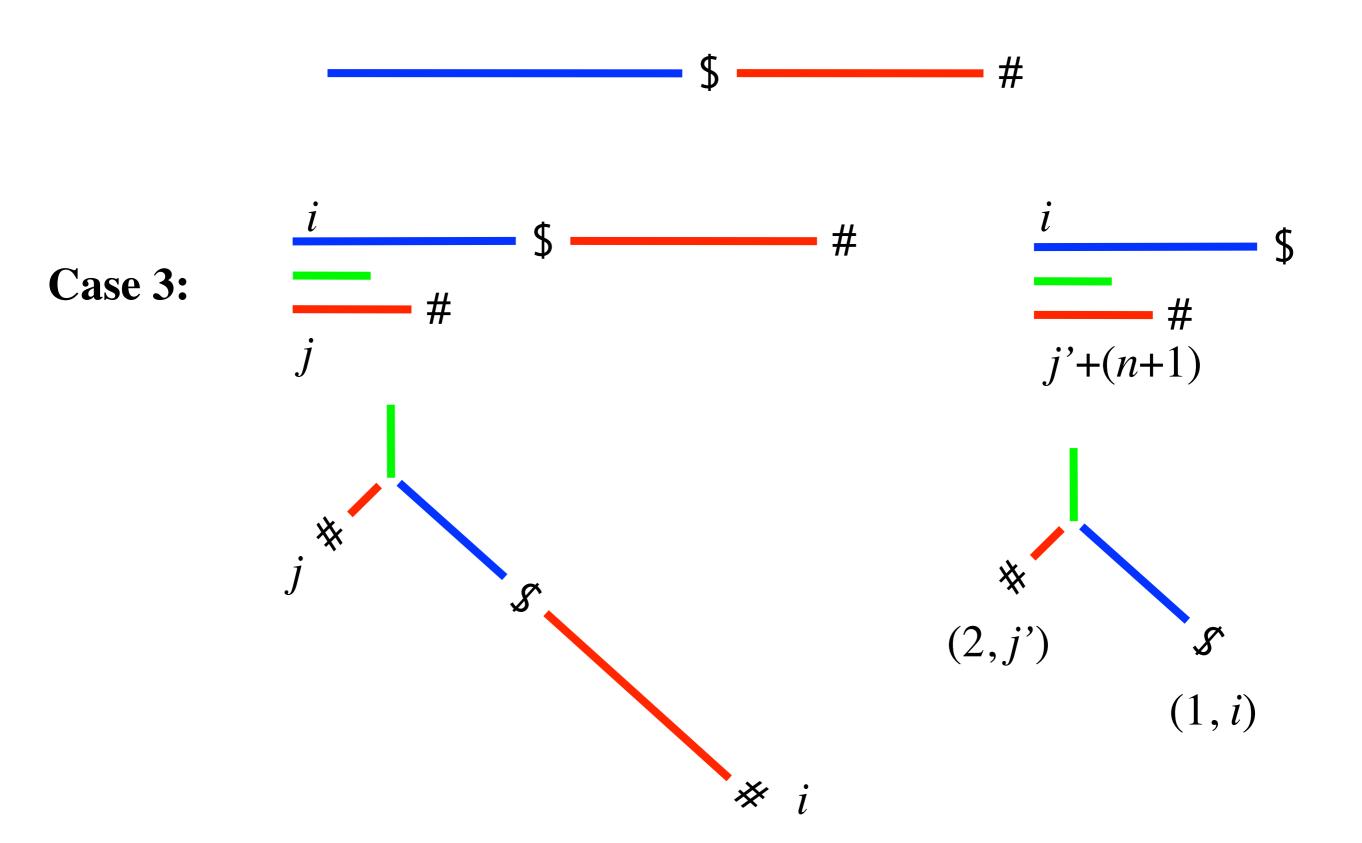
... we must argue that we get the same branching structure ...









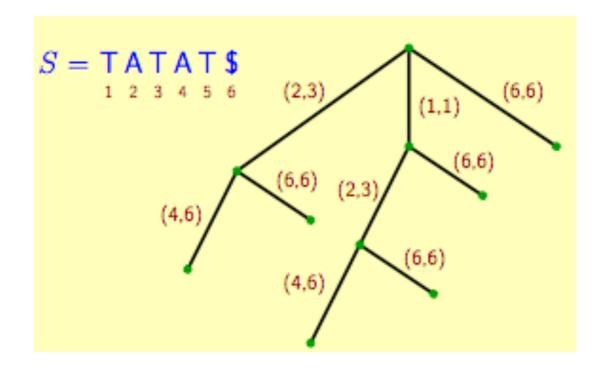


Is everything great?

What are the caveats of suffix trees?

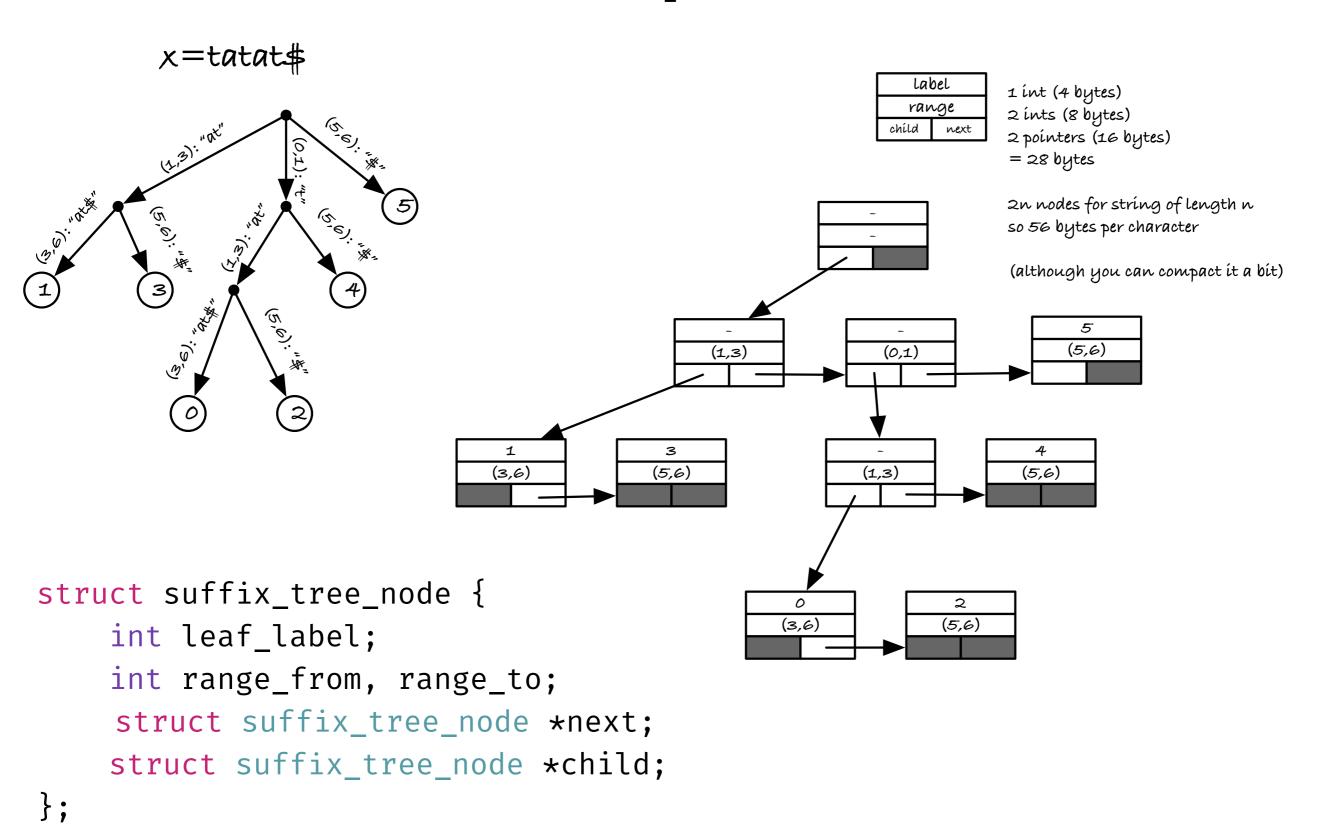
Space consumption

Fact: T(x) requires O(n) space, where n=|x|

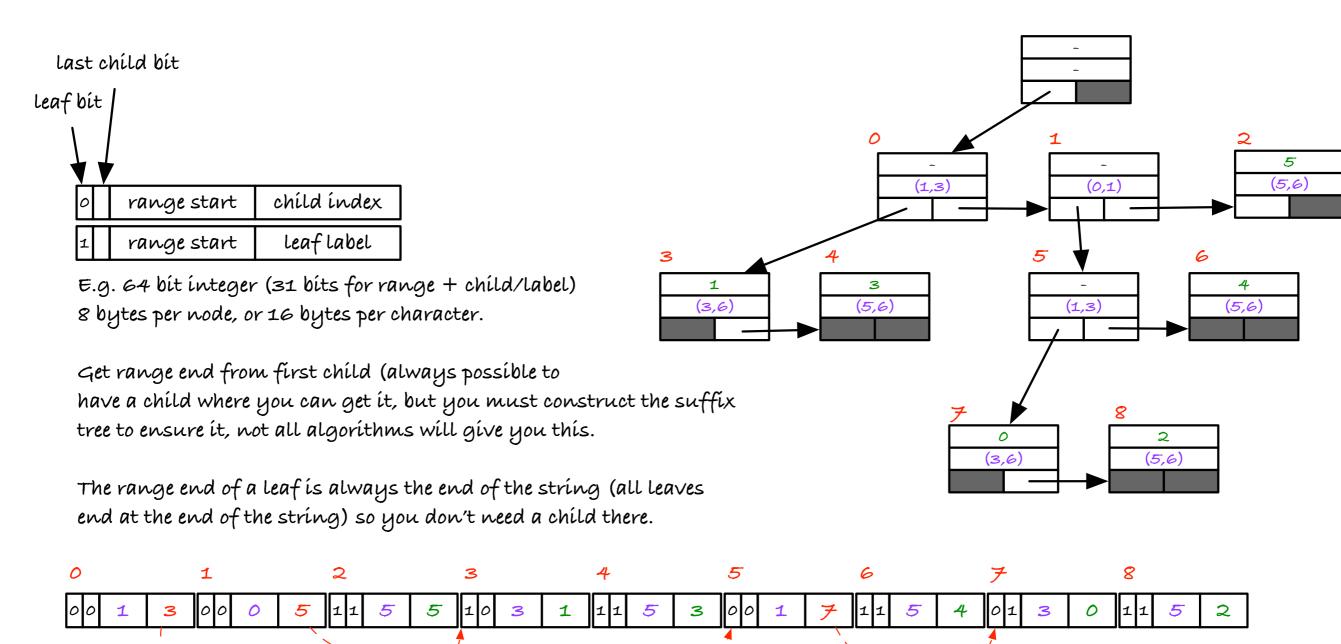


... but how much space does it consume in "practice"?

Linked list representation



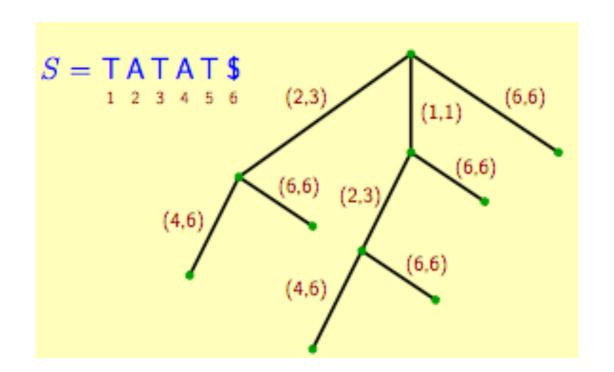
Compaction is possible



(show code)

Space consumption

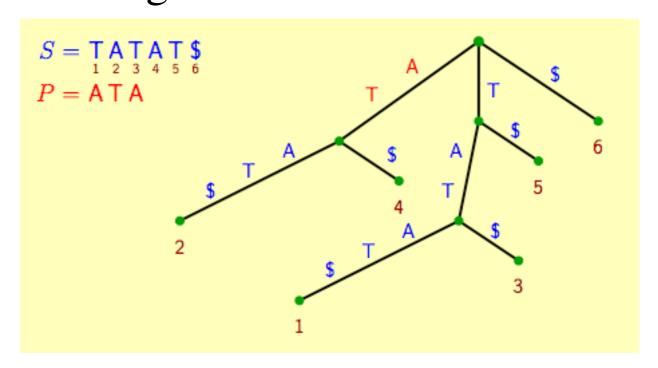
Fact: T(x) requires O(n) space, where n=|x|, but



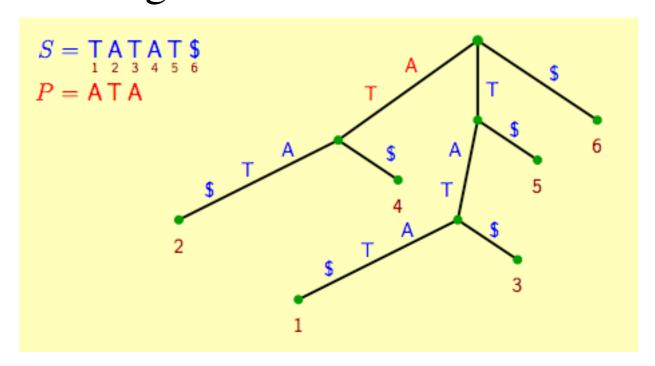
... in practice somewhere between 10 and 60 bytes per letter in x ... (depending on representation)

Is this a problem? Depends on n, if $\approx 500,000,000$ then yes... (and then we have other ways to represent a string)

How much time does it take to find the proper edge out from a node when searching in a suffix tree?



How much time does it take to find the proper edge out from a node when searching in a suffix tree?

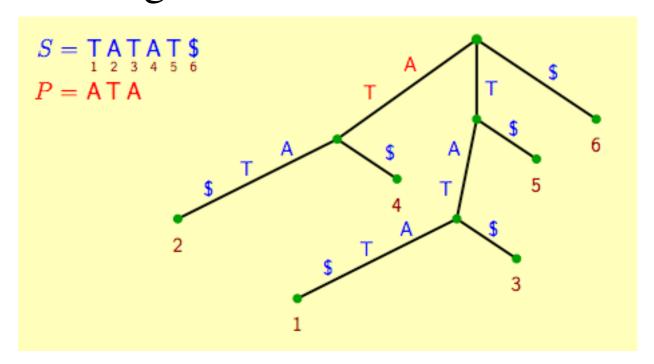


Time proportional to the out-degree of the node $\leq \sigma$...

... search time in "practice" could be $O(\sigma \cdot m)$...

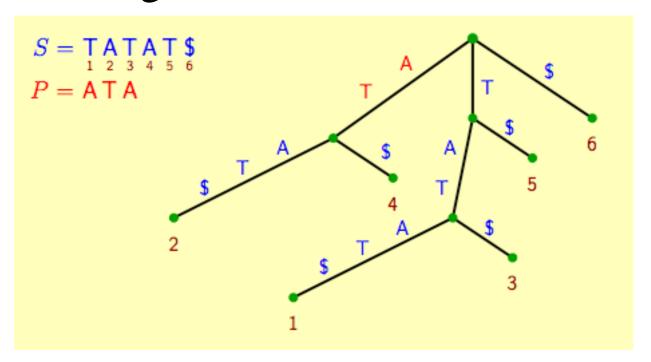
If σ is large, e.g. 256, this matters!!

How much time does it take to find the proper edge out from a node when searching in a suffix tree?



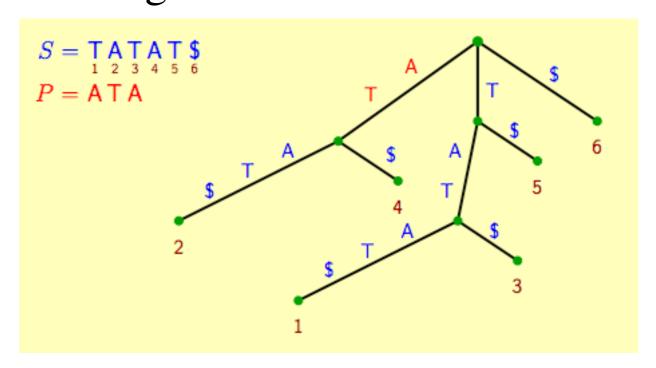
Idea 1: Organising children in a search-tree, reduces search time from σ to $O(\log \sigma)$... (requires an ordered alphabet)

How much time does it take to find the proper edge out from a node when searching in a suffix tree?



Idea 2: Organising children in an array of size σ indexed by letters, reduces search time from σ to O(1) ... (requires a finite alphabet)

How much time does it take to find the proper edge out from a node when searching in a suffix tree?



Idea 3: Use some other dictionary for mapping letters to children ...

... the alphabet size matters in practice ... (but for our readmapper, $\sigma=4$ (or 5 with \$) so we won't worry much)

Enough about applications...

 Next week, we learn how to construct suffix trees in linear time! That's all Folks/