

SUFFIX TREES

TEXT (T) : SHELLSSEASHELLSBYTHESEASHORE

PATTERN (P) : ASH

GOAL: see if P is in T

↳ count/find all matches

1) as fast as possible,
after pre-processing T

2) but also minimize
pre-processing time/space

(1) $\Omega(P) \rightarrow$ multiple searches : $\Omega(\Sigma(P_i))$

(2) $\Omega(T)$

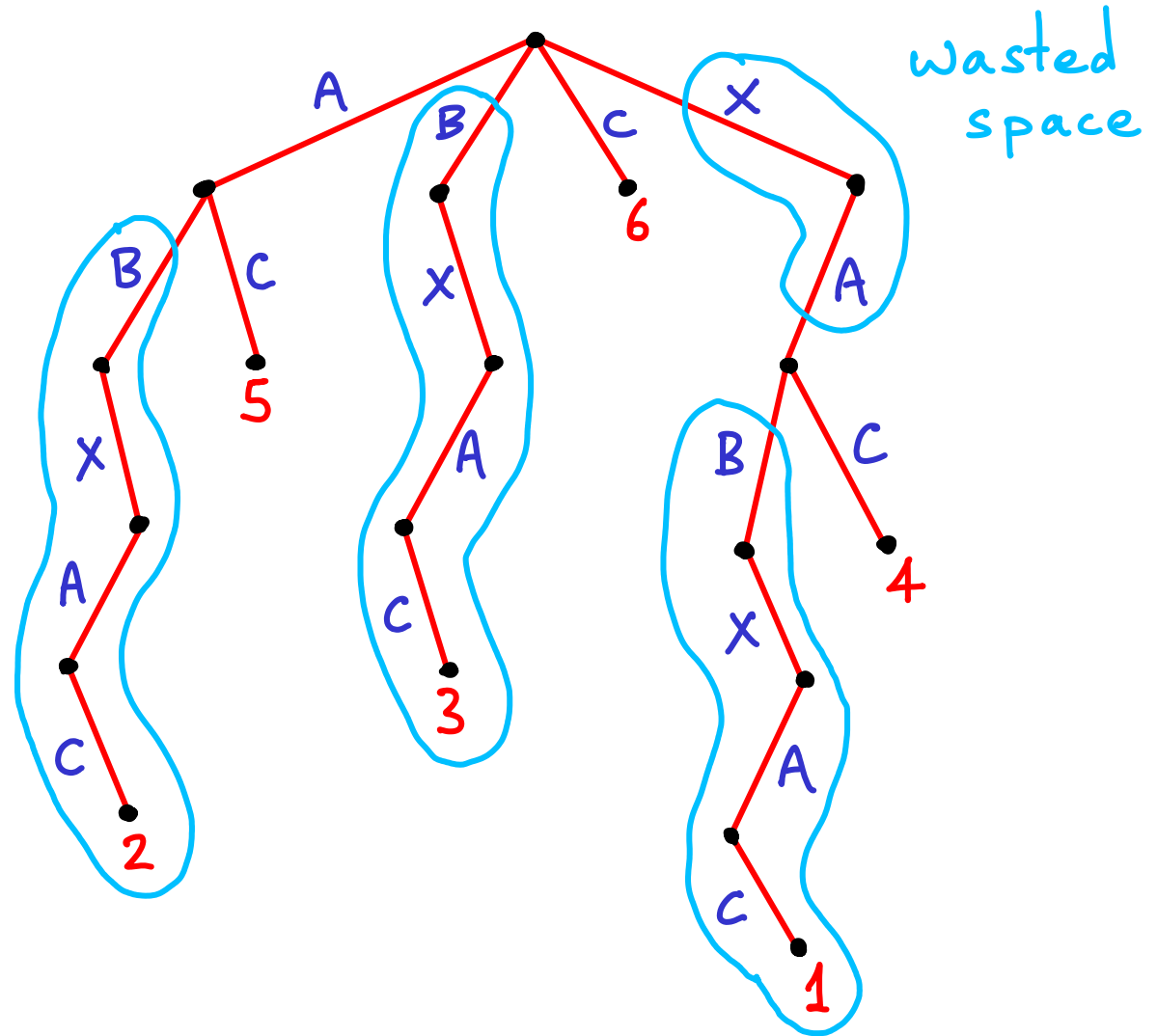
Text : 1 2 3 4 5 6
 X A B X A C

suffixes: 1 X A B X A C
 2 A B X A C
 3 B X A C
 4 X A C
 5 A C
 6 C

$|T| \sim "T": 6$

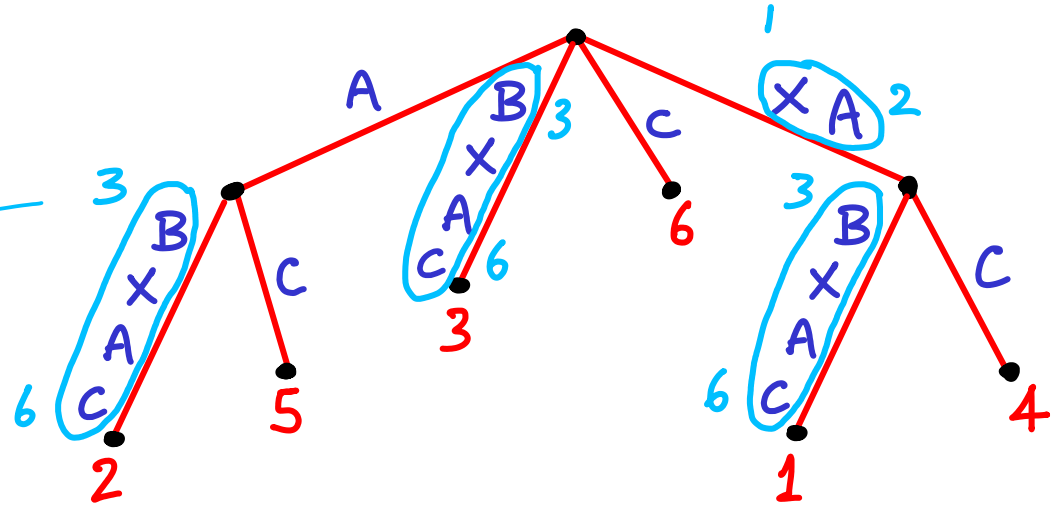
Search for pattern P:
it starts at some index i,
so is in suffix starting at i.

alphabet : A, B, C, X



Text : 1 2 3 4 5 6
 X A B X A C

alphabet : A, B, C, X



using start/end indices,

$$\text{size}(\text{tree}) = O(T)$$

assuming $|\text{alphabet}| = \text{const.}$

search time: $O(P)$

(to find if \exists match)

(augment: subtree sizes for # matches)

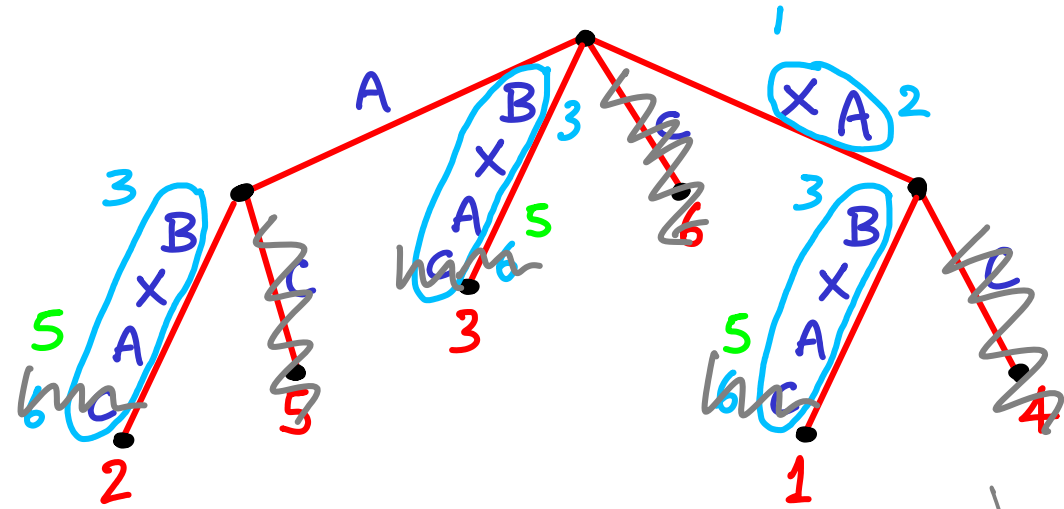
- every suffix is in tree
- # leaves = $T \Rightarrow \underline{\text{size}}(\text{tree}) = O(T)$

↓
nodes

(every internal node has ≥ 2 children)

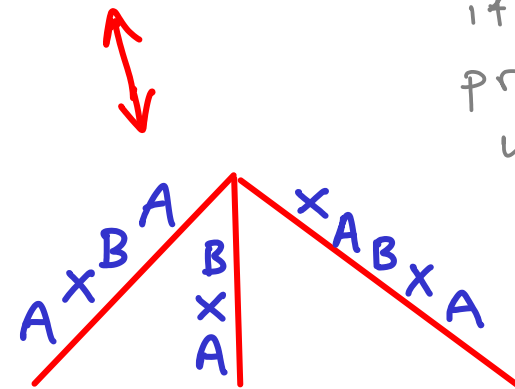
new
Text : 1 2 3 4 5 6
 X A B X A ~~C~~

alphabet : A, B, ~~C~~, X



The suffixes are still there,
but we can't enumerate matches.

Also we can't answer other queries
e.g. "find words starting with AB..."



if we follow
previous rules
we get

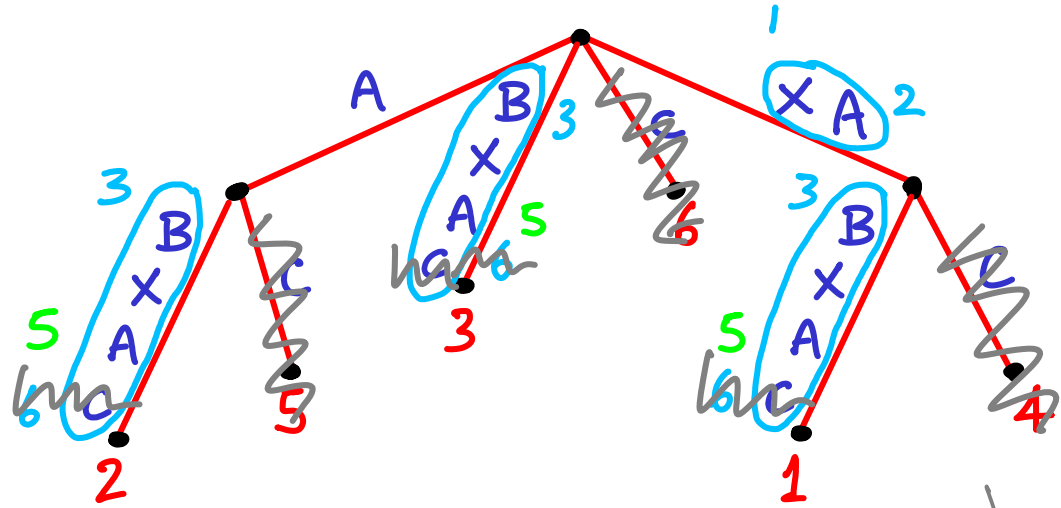
new
Text : 1 2 3 4 5 6
 X A B X A C

alphabet : A, B, ~~C~~, x

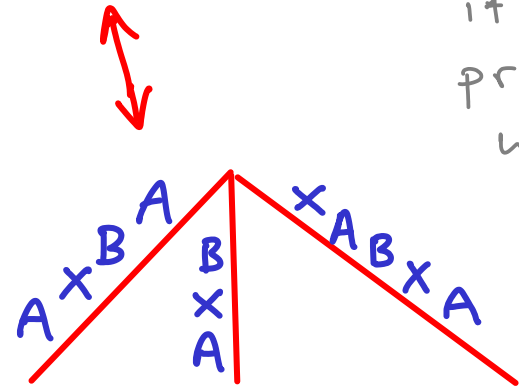
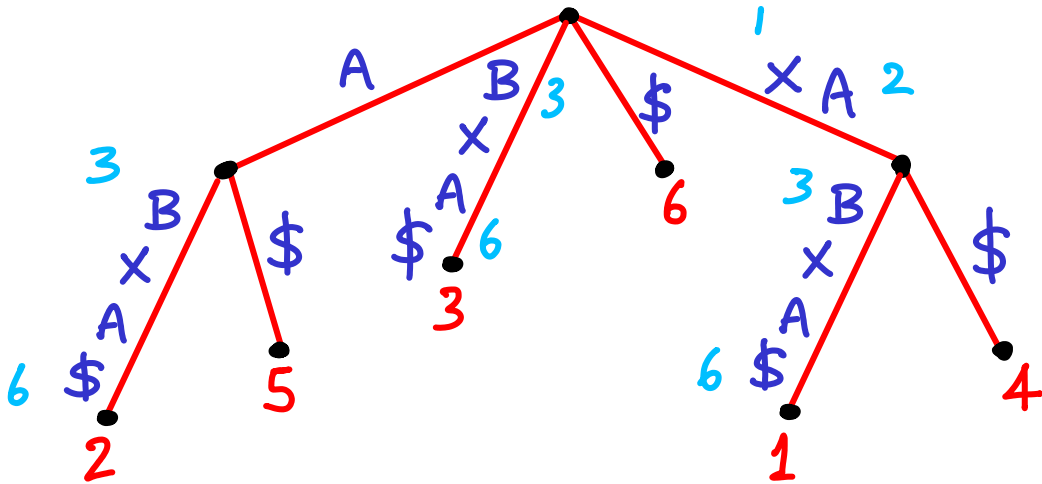
↙ solution

1 2 3 4 5 6
X A B X A \$

alphabet : A, B, x, \$



if we follow
previous rules
we get



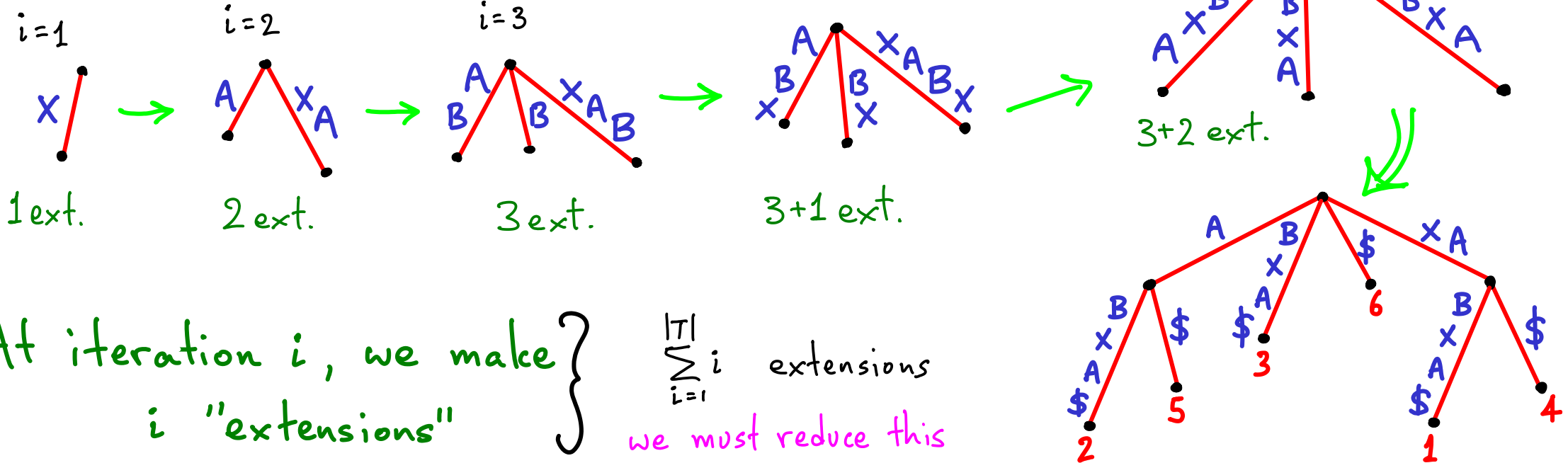
Ukkonen's Algorithm

to build a suffix tree in $O(T)$ time

- iteratively build suffix tree: iteration i builds tree on $T[1...i]$

↳ note: we actually build the "implicit" tree
(& make it proper at last iteration)

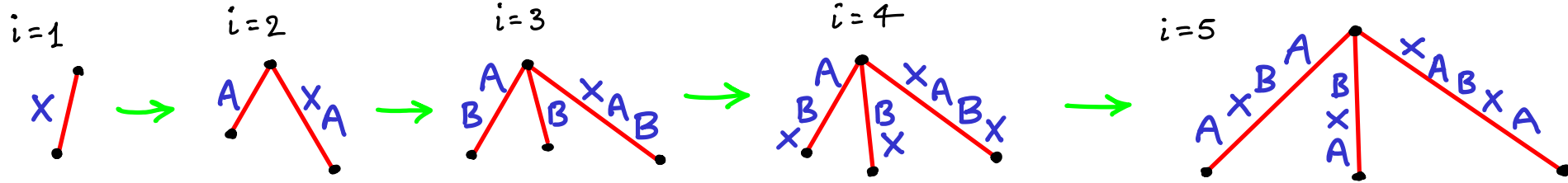
XABXA



At iteration i , we make } $\sum_{i=1}^{|T|} i$ extensions
 i "extensions" } we must reduce this

XABXA

At iteration i , we make i "extensions"

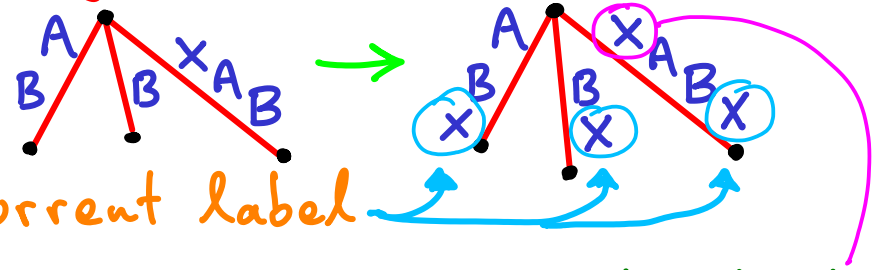


Each extension handles one of the existing suffixes (including empty suffix)

3 cases depending on how current suffix "ends"

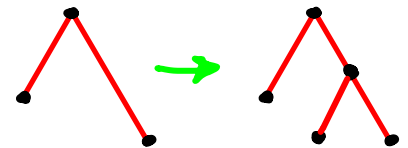
XABX \rightarrow XABXA
ABX \rightarrow ABXA
BX \rightarrow BXA
X \rightarrow XA
 \emptyset \rightarrow A

1) at a leaf



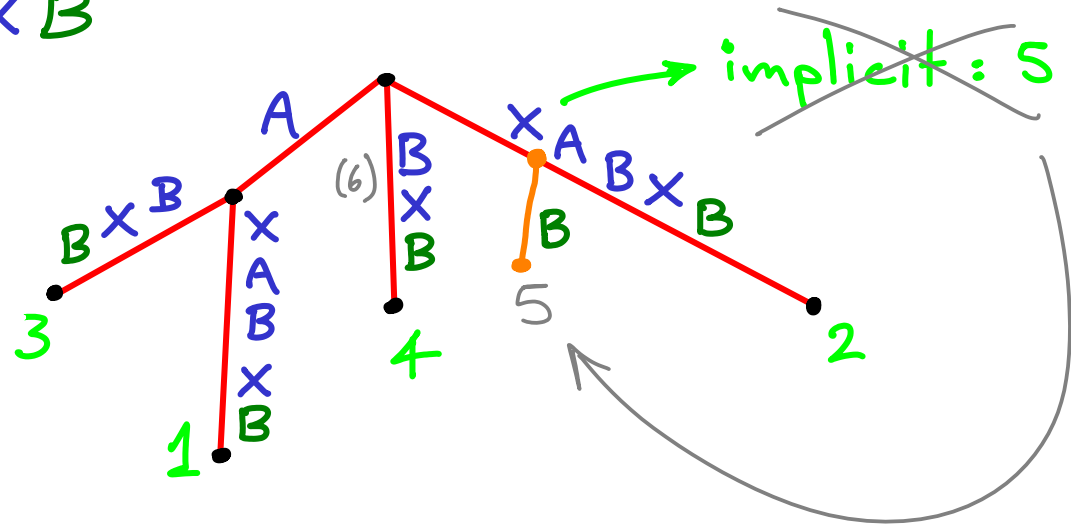
2) not at a leaf, but new character is already there
 \hookrightarrow do nothing

3) not at a leaf, new char. not there
 \hookrightarrow make new edge w/ new character



Another example : $A \times A B \times B$

suffix	1	$A \times A B \times B$
2	2	$\times A B \times B$
3	3	$A B \times B$
4	4	$B \times B$
5	5	$\times B$
6	6	B



Let $T[6] = B$
(iteration 6)

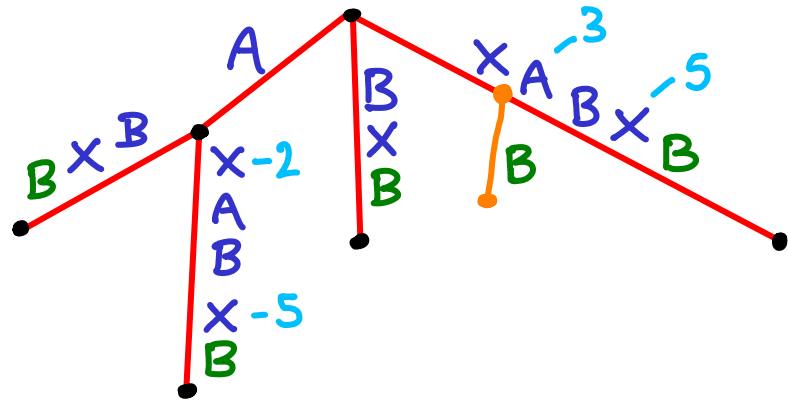
extension	1	2	3	4	5	6
case	1	1	1	1	3	2

- 1) at a leaf \rightarrow extend current label
- 2) not at a leaf, but new character is already there \rightarrow do nothing
- 3) not at a leaf, new char. not there \rightarrow make new edge w/ new character

At iteration i , we make i "extensions"

For each extension: 3 cases depending on how current suffix "ends"
each looks easy & quick

how do we determine this?



* Scanning whole suffix $\rightarrow O(i^2)$ per iteration
 $\left(\sum_{j=1}^i i-j\right)$

* Using indexing doesn't help (yet)

A X A B X B

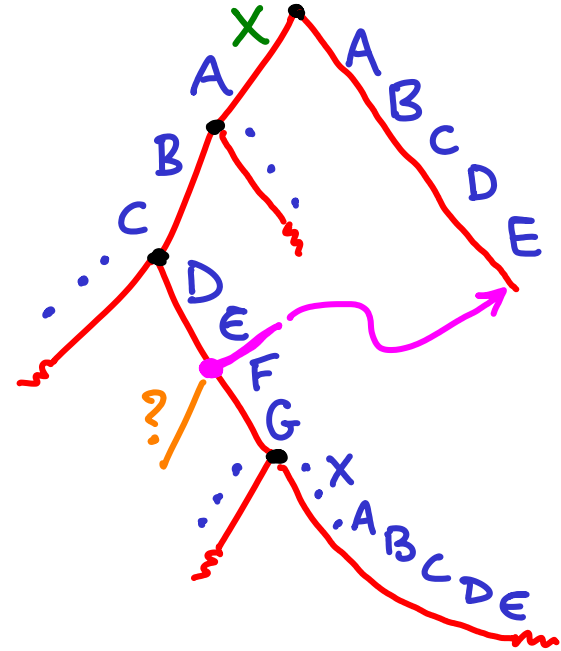
- 1) at a leaf \rightarrow extend current label
- 2) not at a leaf, but new character is already there \rightarrow do nothing
- 3) not at a leaf, new char. not there \rightarrow make new edge w/ new character

How do we determine where a current suffix "ends" ? (quickly)

Suffix Links

Say we just found the end of a suffix, **XABCDE**

(maybe while extending it w/ \mathbb{F} or something else)



Then **ABCDE** is also a suffix,
so it is also in the tree
and it is the next extension

so it would help to have a link

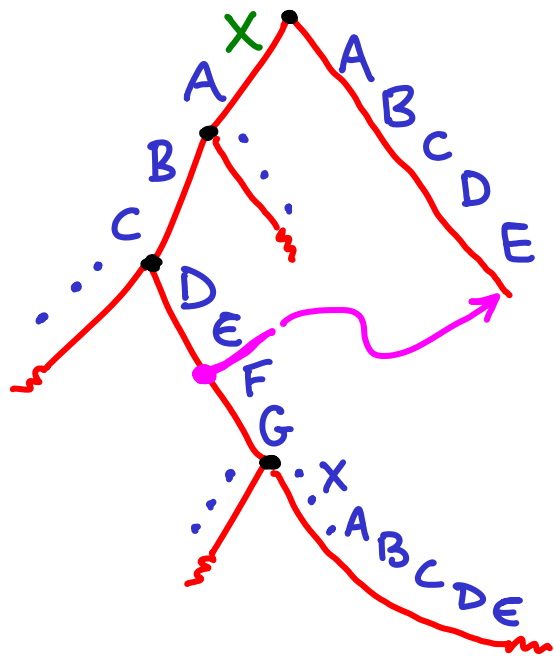
from every x_α to α ... except maybe there was no node at x_α

$T[1 \dots i-1] = \text{XABCDEFGXABCDE} \boxed{?}$
 $T[i]$

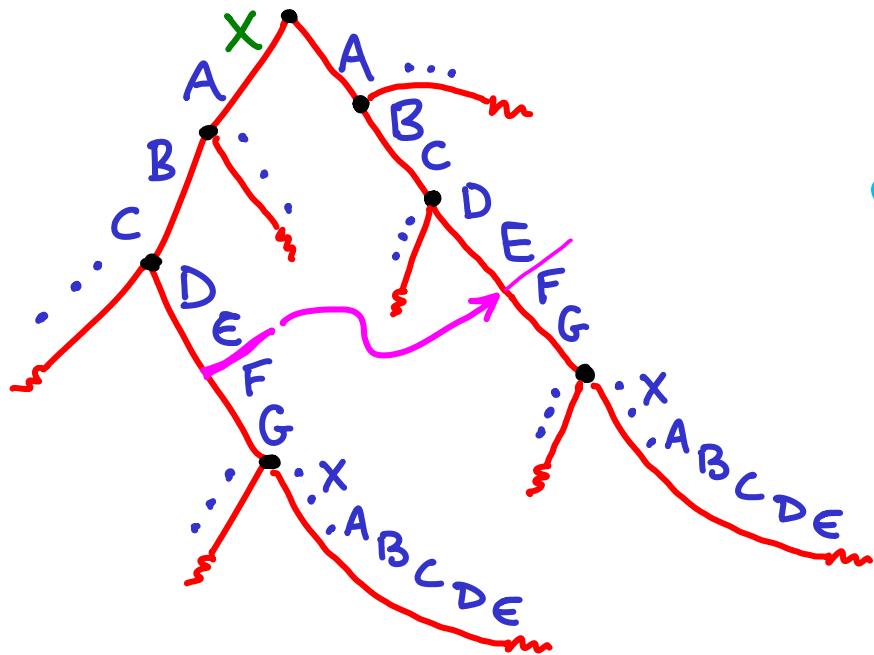
$$T[1 \dots i-1] = \text{XABCDEFGXABCDE}$$

actual
suffix tree

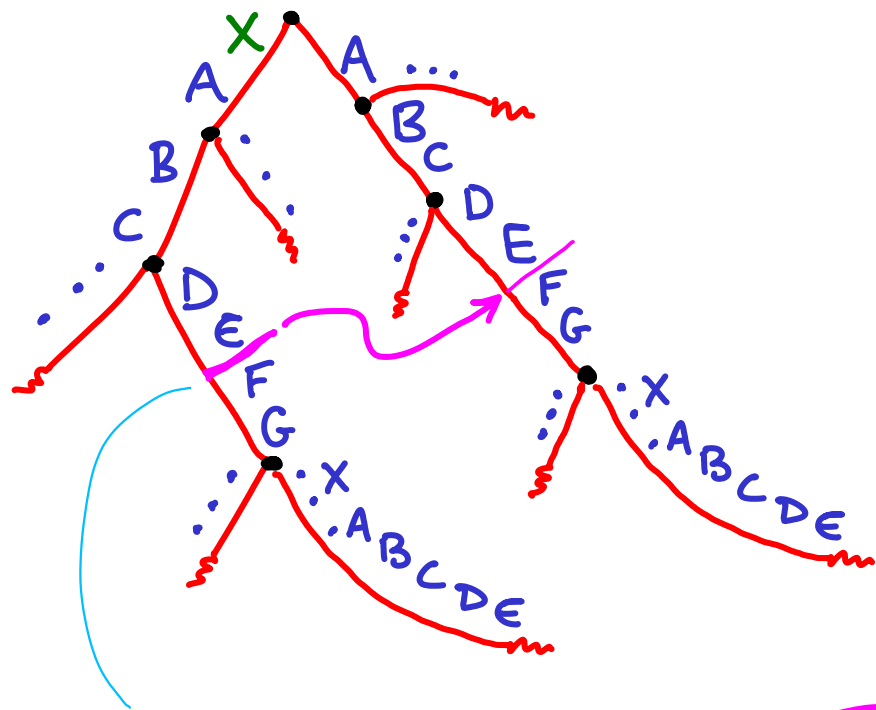
(now enjoy this unused slide)



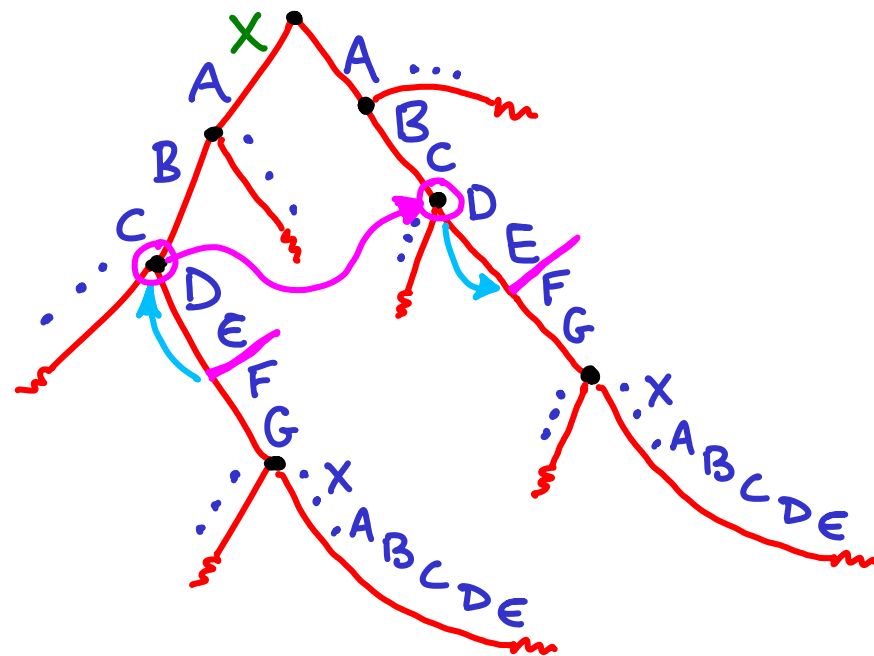
for illustration \rightarrow



lots of repeated structure



no node here



Instead move up until a node is found,
use a suffix link from node to node,
then "move down"

deal w/ this
later

$O(1)$

next issue:
how to make suffix links



When do we create a new node?

↳ case where we process $x \alpha$ to extend w
and $\exists x \alpha z$ ($z \neq w$)

(end of x_α is in mid-edge)

When do we create a suffix link from $x\alpha$ to α ?

↳ First we need to know that a node exists at α .



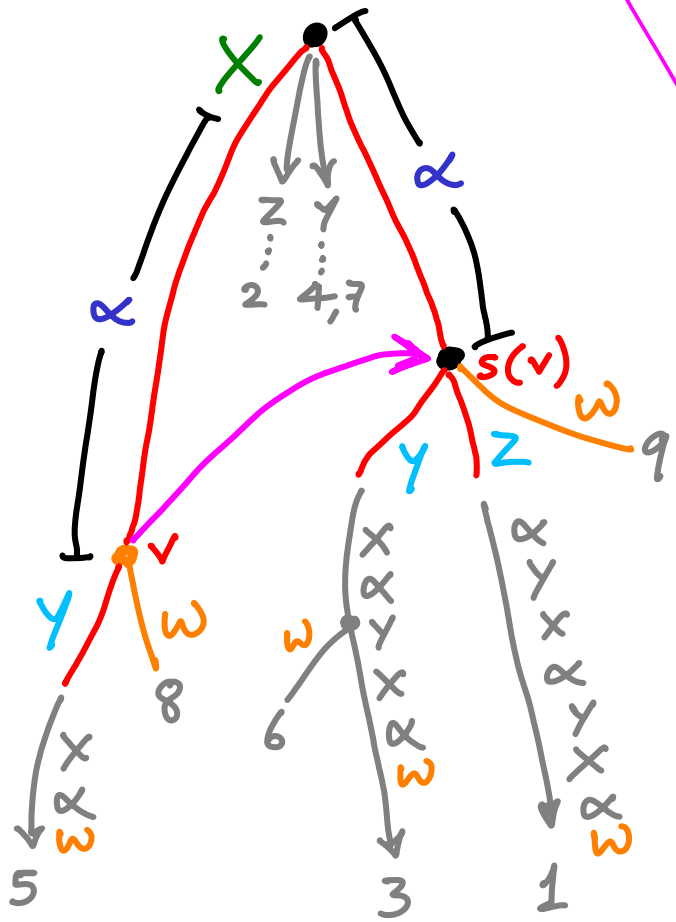
(end of x_α is in mid-edge)

↳ First we need to know that a node exists at α .

it will be created at next extension (we will process ∞)

So we can assume that all "old" nodes have outgoing suffix links

Example of suffix link target node existing before source node



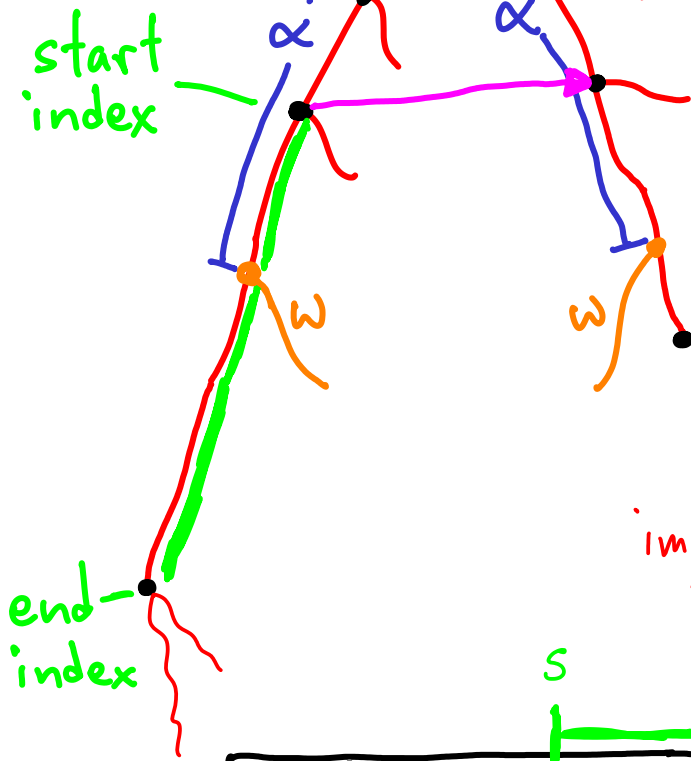
- to be created at extension 9

1 2 3 4 5 6 7 8 9

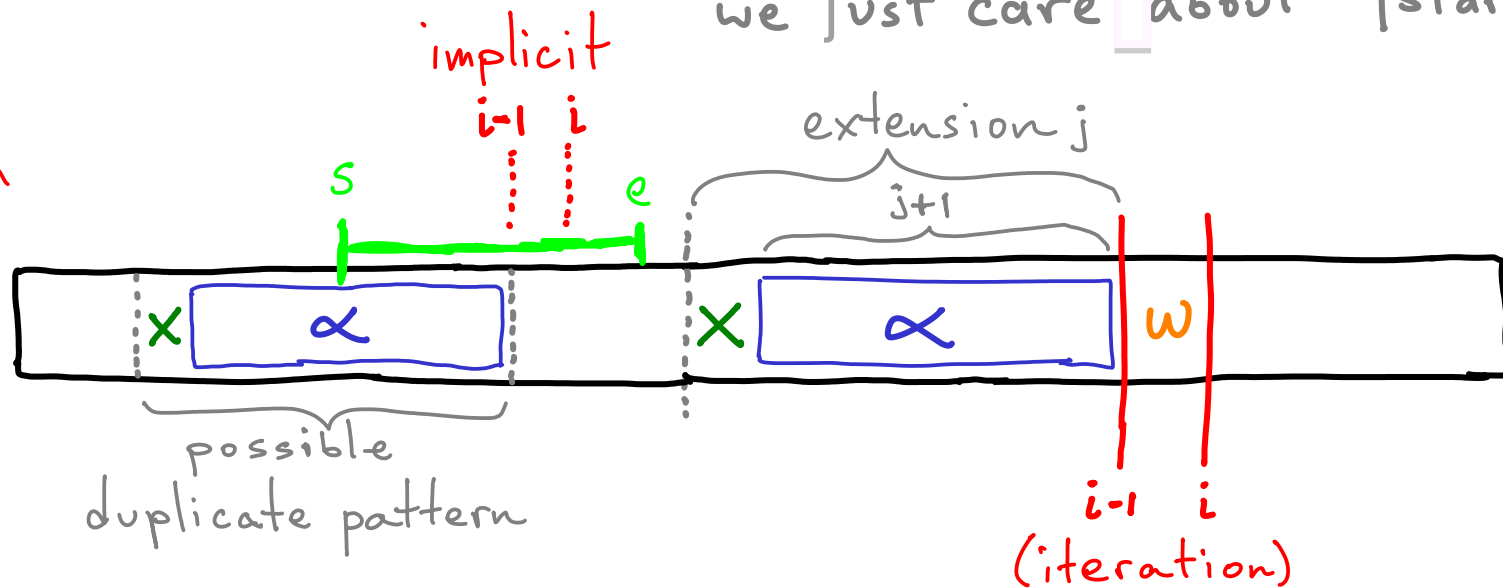
α Z α Y X α Y X α W

currently processing extension 8

Technical details



For each edge we store the start/end indices. While processing a suffix that is the prefix of a larger suffix we would be using indices implicitly (to find actual chars from T). Besides that, we just care about $|start - end| = \text{length}$



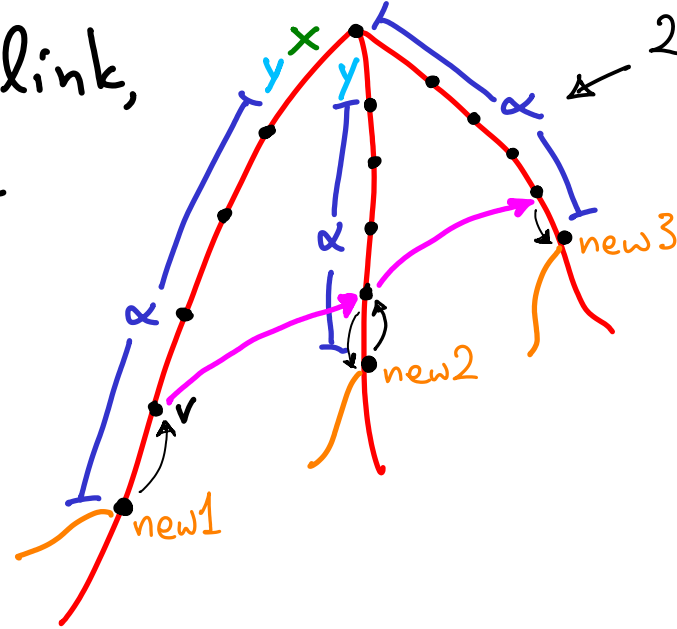
WHAT MATTERS:

using indices
we can move
between nodes in
 $O(i)$ time

When following a suffix link,
node depth (from root) can

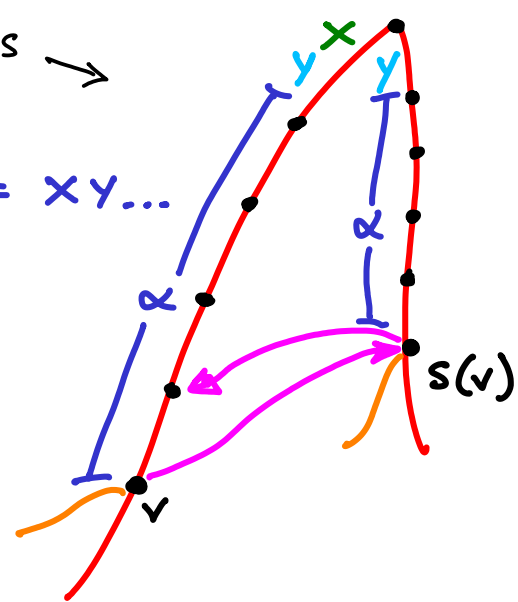
- stay same
- increase (down to ~leaf)
- decrease
(i.e. move up)

example
to follow

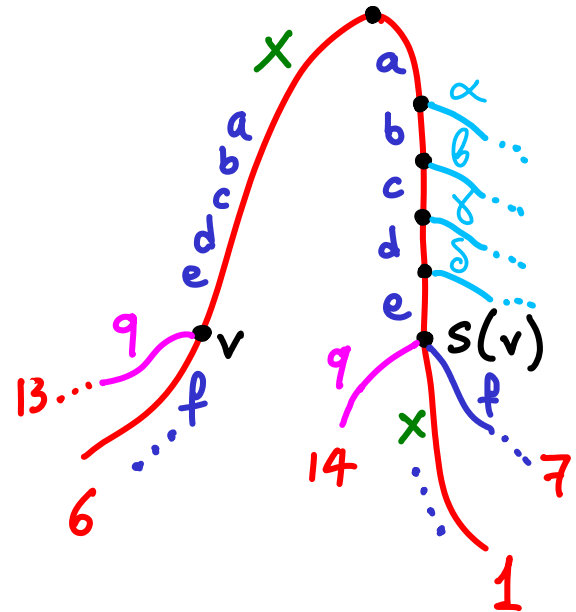


2 examples \rightarrow

$\alpha = xy \dots$



example of v with suffix link to $s(v)$
 where $\text{depth}(v) \ll \text{depth}(s(v))$



1
 abcde X abcdef X abcde q a x ab babcyabcd 5
 6
 X abcdef X abcde q a x ab babcyabcd 5
 7
 abcdef X abcde q a x ab babcyabcd 5
 13
 X abcde q a x ab babcyabcd 5
 14
 abcde q a x ab babcyabcd 5

When following a suffix link,
node depth (from root) can

- stay same
- increase (down to ~leaf)
- decrease (only by 1)
(i.e. move up)

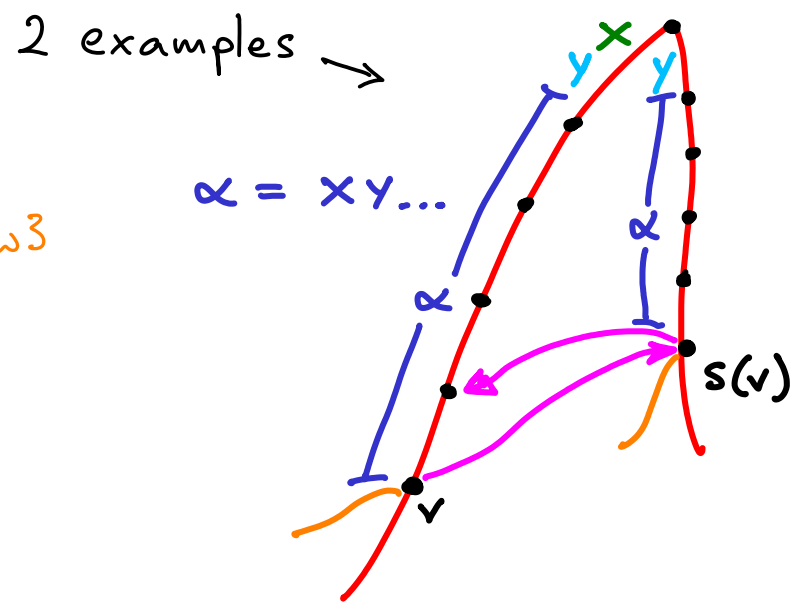
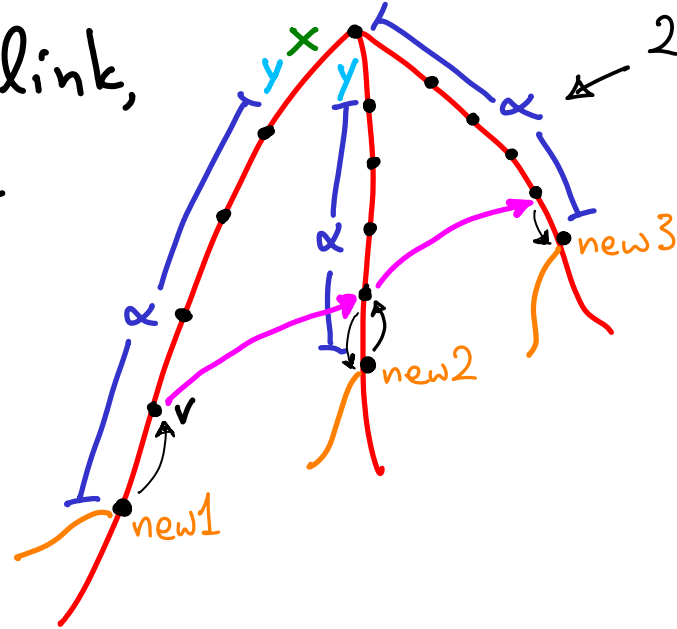
→ we've seen: if xy leads to v

then y is in tree, leading to $s(v)$

So every node on path from x to v

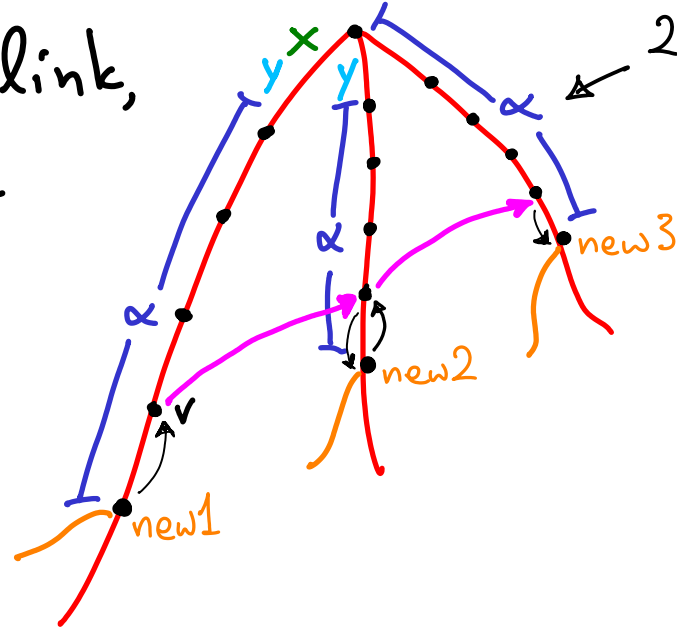
has a "mirror" node on path from root to $s(v)$

(reverse not necessarily true)



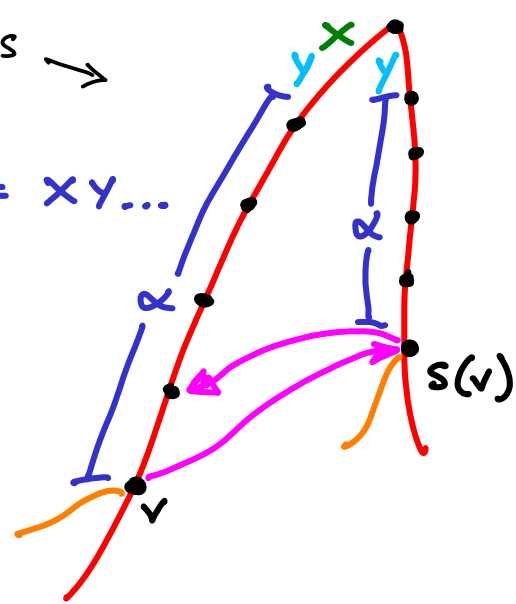
When following a suffix link, node depth (from root) can

- stay same
- increase (down to ~ leaf)
- decrease (only by 1)
(i.e. move up)



2 examples →

$$\alpha = xy \dots$$



amortize : total node hops per iteration = $O(i)$

$$[\text{total upward} = O(i) \quad \& \quad \text{max depth} = O(i)]$$

CONCLUSION

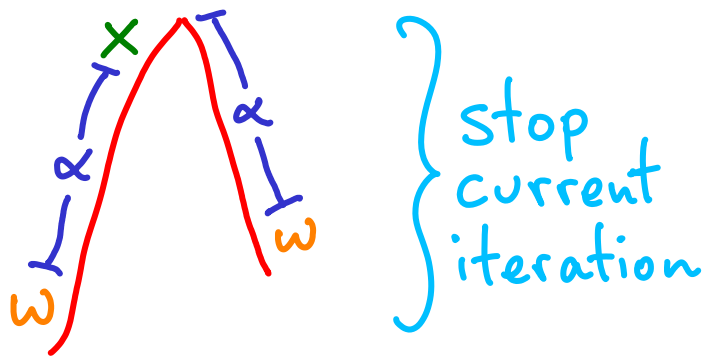
During iteration i , we visit $O(i)$ nodes

...spending $O(i)$ time to find each, excluding the first one: total work $O(i)$

which is always
the longest suffix,
i.e. same node

Story so far: find "end" of next extension using suffix link from end of current extension (plus a bit of up/down)

What if current extension makes no new node?



extend $x\alpha + w$
but $x\alpha w$ already in tree
(αw also in tree)

		extension					T
		1	2	3	...	i	
iteration + node (i)	1	2					
	2	3	1				
	3	1	1	1			
	...	3	2	-	-		
	...	1	3	1	3	3	
	...	3	3	3	1	1	3
	(i) T	1	3	1	2	-	-

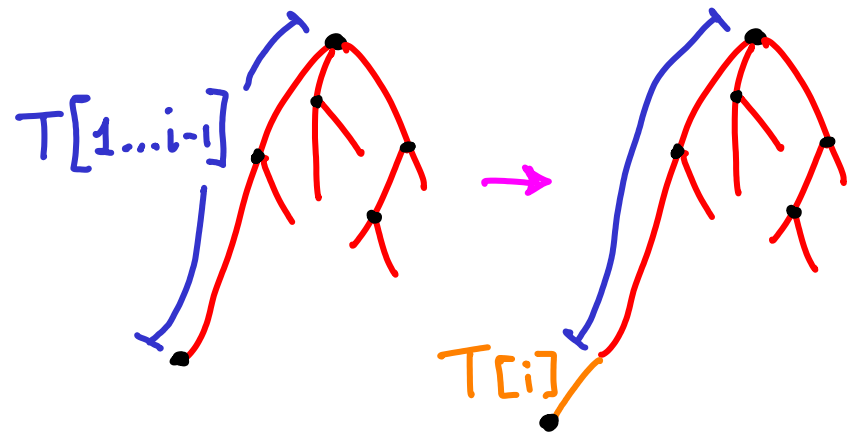
Table of cases applied in each extension of each iteration
↳ nothing right of 2

All future extensions in current iteration will also follow same case

- 1) at a leaf → extend current label
- 2) not at a leaf, but new character is already there → do nothing
- 3) not at a leaf, new char. not there → make new edge w/ new character

Story so far: find "end" of next extension using suffix link from end of current extension (plus a bit of up/down)

What if current extension makes no new node?



First extension of every iteration:
case 1

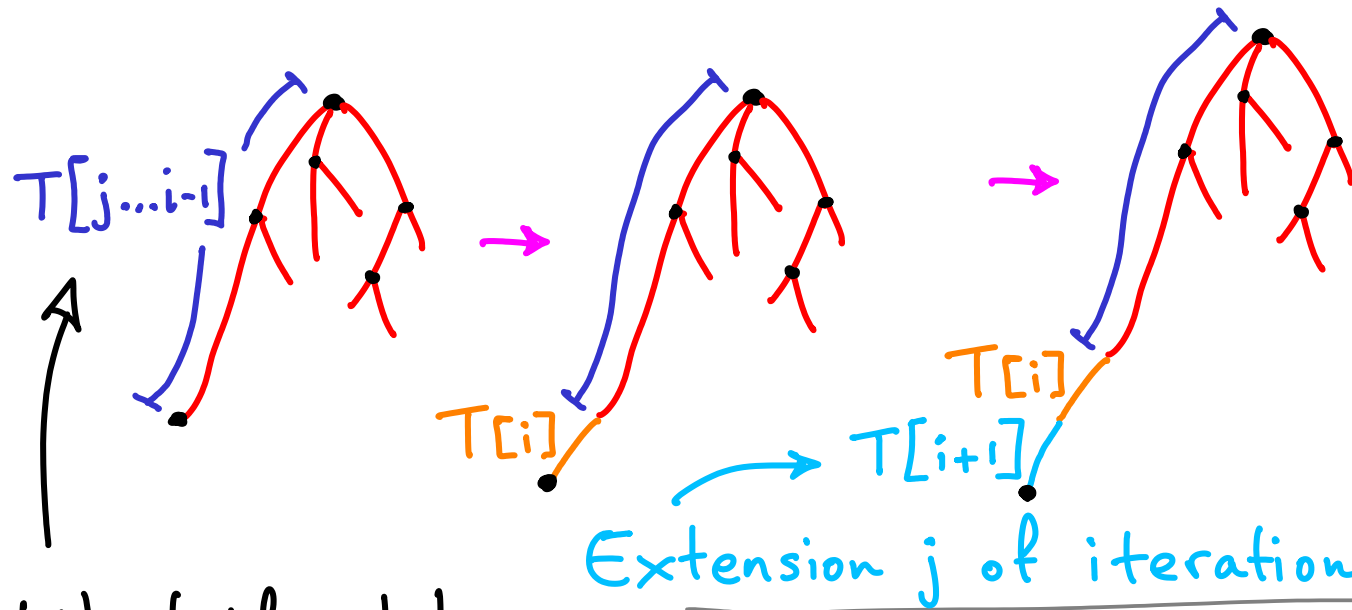
(longest suffix ends at leaf)

		extension					
		1	2	3...	i...		T
iteration non	1	1					
	2	1	1				
	3	1	1	1			
		1	2	-	-		
		1	3	1	3	3	
		1	3	3	1	1	3
		1	3	1	2	-	-
(i)	T	1	3	1	2	-	-

- 1) at a leaf → extend current label
- 2) not at a leaf, but new character is already there → do nothing
- 3) not at a leaf, new char. not there → make new edge w/ new character

Story so far: find "end" of next extension using suffix link from end of current extension (plus a bit of up/down)

What if current extension makes no new node?



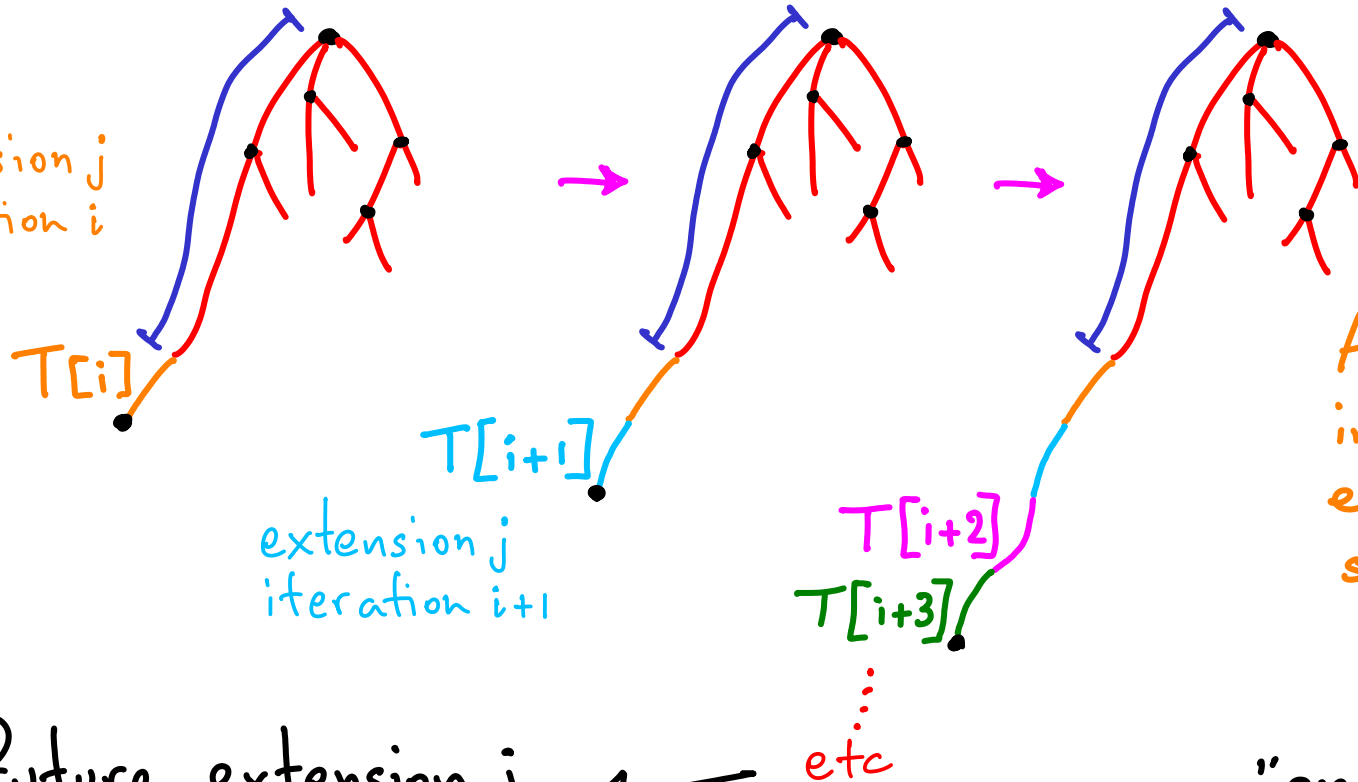
All suffixes considered in the future will divert elsewhere in the tree, or stop short on this path, or extend this path

"once a leaf, always a leaf"

What if arbitrary extension j of iteration i uses case 1?

- 1) at a leaf \rightarrow extend current label
- 2) not at a leaf, but new character is already there \rightarrow do nothing
- 3) not at a leaf, new char. not there \rightarrow make new edge w/ new character

extension j
iteration i



[for future iterations
extension $k < j$ diverts;
extension $k > j$ too short]

All suffixes considered
in the future will divert
elsewhere in the tree, or
stop short on this path,
or extend this path

"once a leaf, always a leaf"

Any future extension j
has same prefix, so
will follow same path
& use same rule

- 1) at a leaf \rightarrow extend current label
- 2) not at a leaf, but new character is already there \rightarrow do nothing
- 3) not at a leaf, new char. not there \rightarrow make new edge w/ new character

what if extension j uses case 3? → creates a leaf

... same logic,
must have 1's
below

We've determined that
below any 1 we can
have only 1's

		extension						
		1	2	3...	i...		T	
iteration (i)	1	1						
	2	1	2					
	3	1	3	2				
	...	1	1	2	-			
	...	1	1	2	-	-		
	...	1	1	3	1	3	3	~
	...	1	1	1	1	1	1	2
	T	1	1	1	1	1	1	2

1) at a leaf → extend current label

2) not at a leaf, but new character is already there → do nothing

3) not at a leaf, new char. not there → make new edge w/ new character

		extension					
		1	2	3...	i...		$ T $
iteration non- (i)	1	1					
	2	1	2				
	3	1	3	2			
	⋮	1	1	2	-		
	⋮	1	1	2	-	-	
	⋮	1	1	1	3	3	2
	⋮	1	1	1	1	1	2
	$ T $	1	1	1	1	1	-

if 1 or 3 \rightarrow all 1's below
 if 2 \rightarrow empty to right

$\leq 2|T|$ "interesting" entries
 (still $O(T^2)$ work!)

How to handle the 1's: use a global variable = iteration #

- recall, the 1's involve leaf (edge) extensions
- edges are represented by indices in T .
- the "end" index just gets incremented: $i-1 \rightarrow i$

Last step: add \$: makes all implicit suffixes proper

conclusion: suffix tree construction $\rightarrow O(T)$ time & space