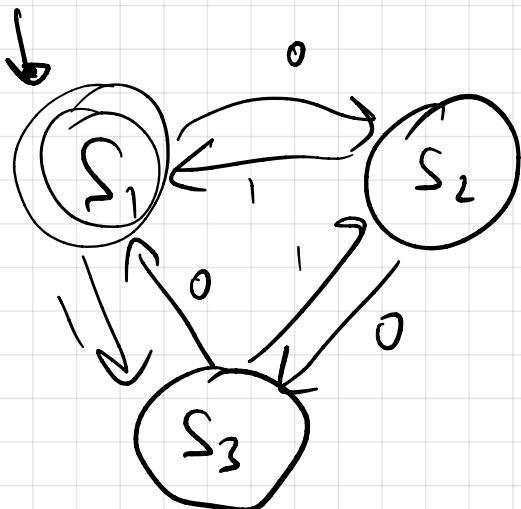


Suffix Automaton

Automata Theory



↓ - starting state

○ accepting st.

→ transitions

usual convention - ; if link is missing - it goes to "garbage" node

$\text{Aut} = (S, S_0 \in S, \delta: S \times \Sigma \rightarrow S, T \subseteq S)$

states

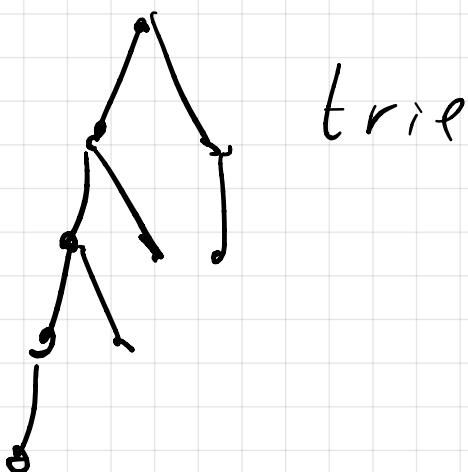
$L_{\text{Aut}} = \{ t \in \Sigma^*: \delta(S_0, t) \in T \}$

Suffix automaton ($S \in \Sigma^*$) - smallest

automata accepting suffixes of S
& only them

Some automata accepting only suffixes

if S



trie

Right Contexts:

$$R_{Aut}(t \in \Sigma^*) = \{ \alpha \in \Sigma^* : t\alpha \in L_{Aut} \}$$

$$R_{Aut}(v \in S) = \{ \alpha \in \Sigma^* : \delta(v, \alpha) \in T \}$$

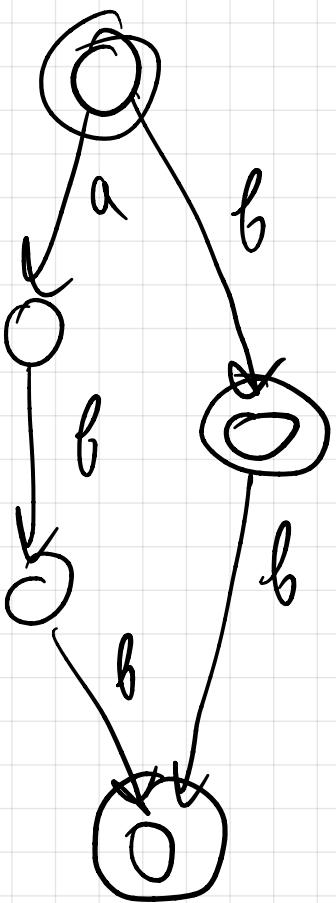
Claim: it's possible to compress
states of some $R(v)$

\Rightarrow so smallest automaton consists

of exactly $\{R(s) \mid s \in \Sigma^*\}$ states

Automaton edges: $R(s) \xrightarrow{\cdot} R(sc)$

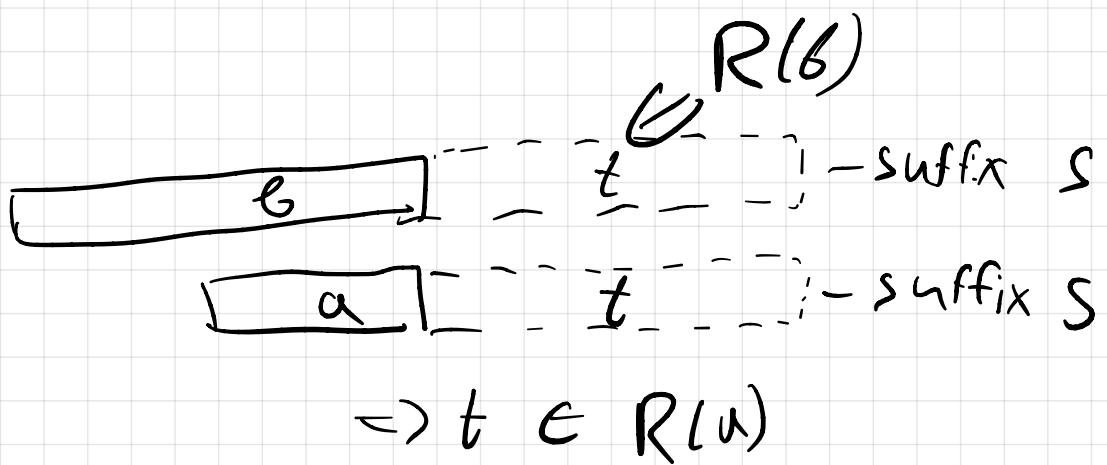
Suffix Automaton



$SA(S)$ - Suffix aut.
of S .

Lm a - suffix b
 $\Rightarrow R(a) \supseteq R(b)$

Proof



Lm

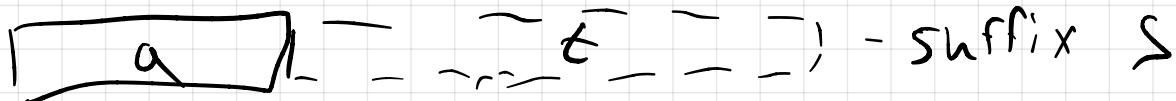
$$R(a) \cap R(b) \neq \emptyset$$

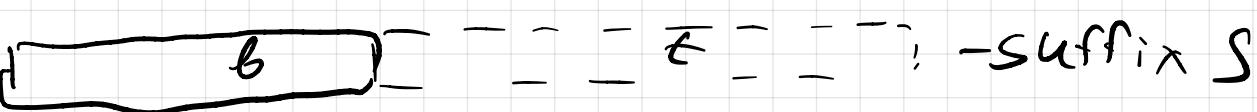
$$\Rightarrow R(a) \supseteq R(b) \vee R(b) \supseteq R(a)$$

(a - suffix b) (b - suffix a)

Proof

Suppose $|a| \leq |b|$, $t \in R(a) \cap R(b)$

 - suffix s

 - suffix s

$\Rightarrow a$ - suffix b

$\Rightarrow R(a) \supseteq R(b)$

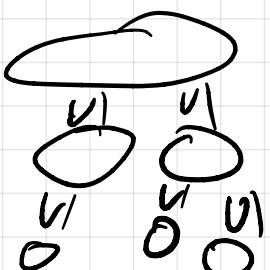
□

i.e. R sets are either nested or don't

intersect.

We can consider tree of R sets

(excl. empty set)



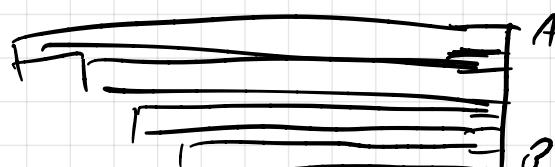
Lm

for any $R \neq \emptyset$

Consider

$$\text{Str}_R = \{\alpha \mid R(\alpha) = R\}$$

Then $\text{Str}_R = \text{continuous set of suffixes}$



i.e. if A is longest in STR_R
and β is shortest in STR_R

Then $STR_R = \{ t \mid t\text{-suffix } A,\text{ }\beta\text{-suffix } t \}$

and β -suffix A

Proof ① $R(A) = R(\beta)$

$\Rightarrow A$ suffix β

or β suffix A

Since $|\beta| \leq |A| \Rightarrow \beta$ -suffix A

② $R(t) = R(\beta) = R(A)$

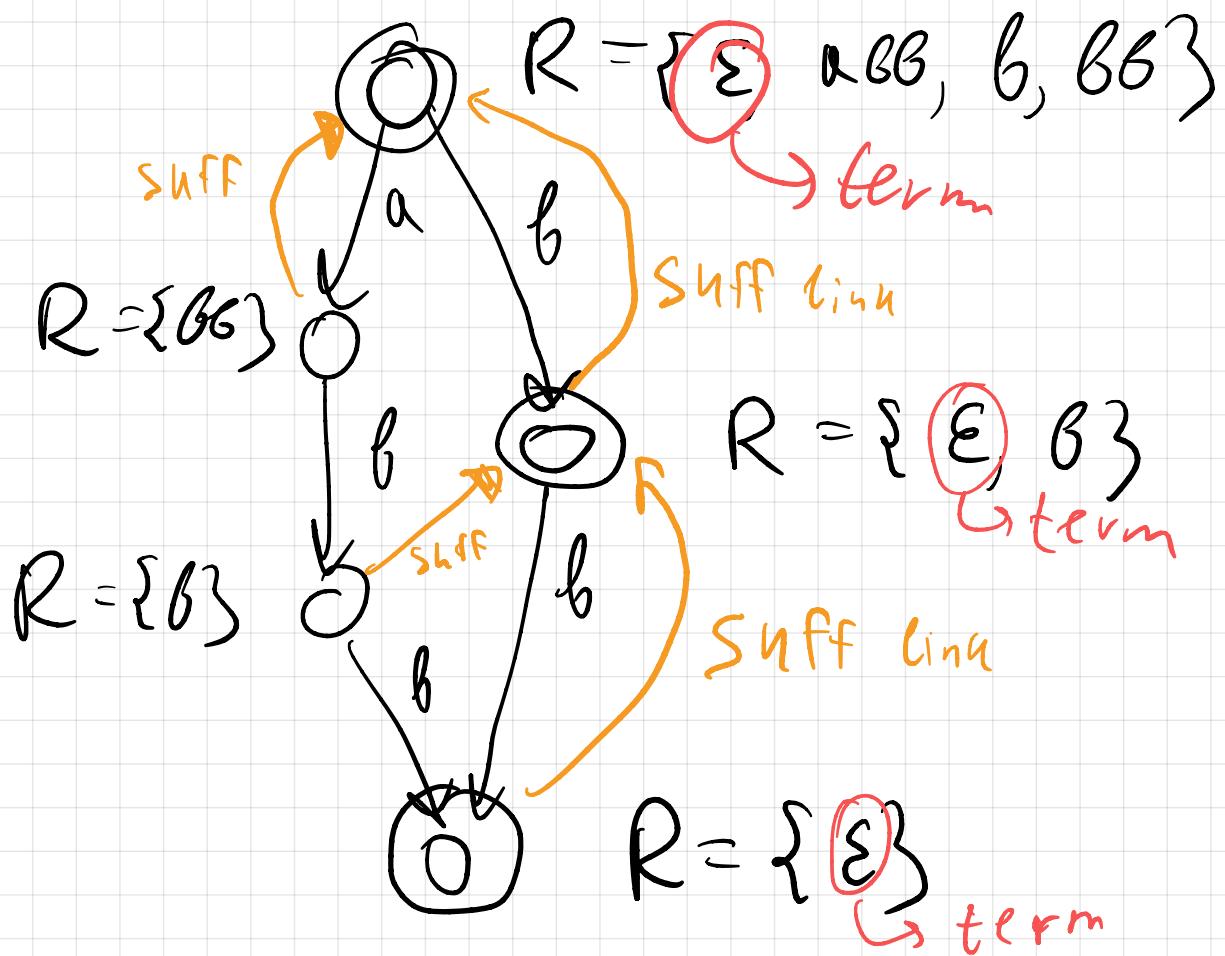
$\Rightarrow t$ suffix A ,

β -suffix t

③ if t -suffix A ,
 β -suffix t

$\Rightarrow R(A) \subseteq R(t) \subseteq R(\beta)$

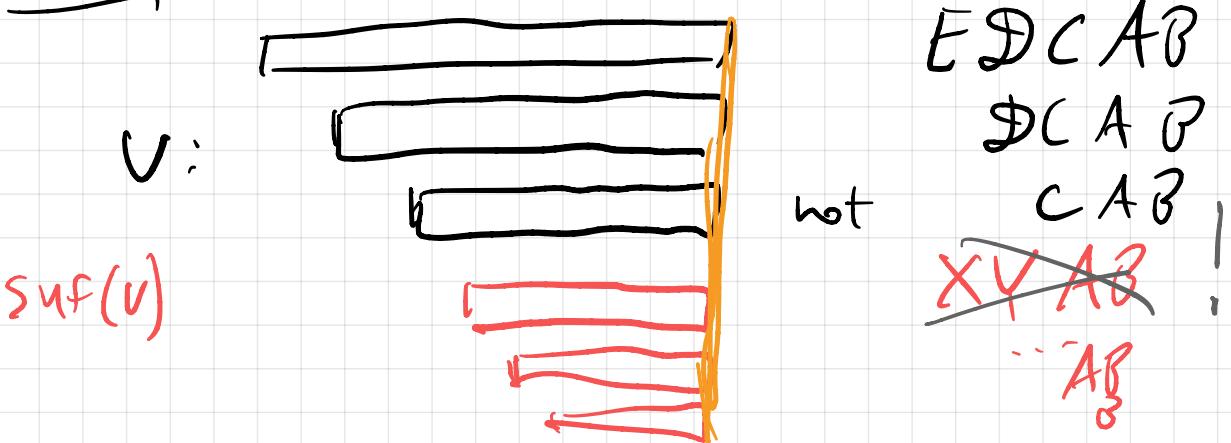
$\Rightarrow R(t) = R$



Def: Let $\text{len}(v)$ = length of longest t
s.t. $R(t)$ corresponds to v

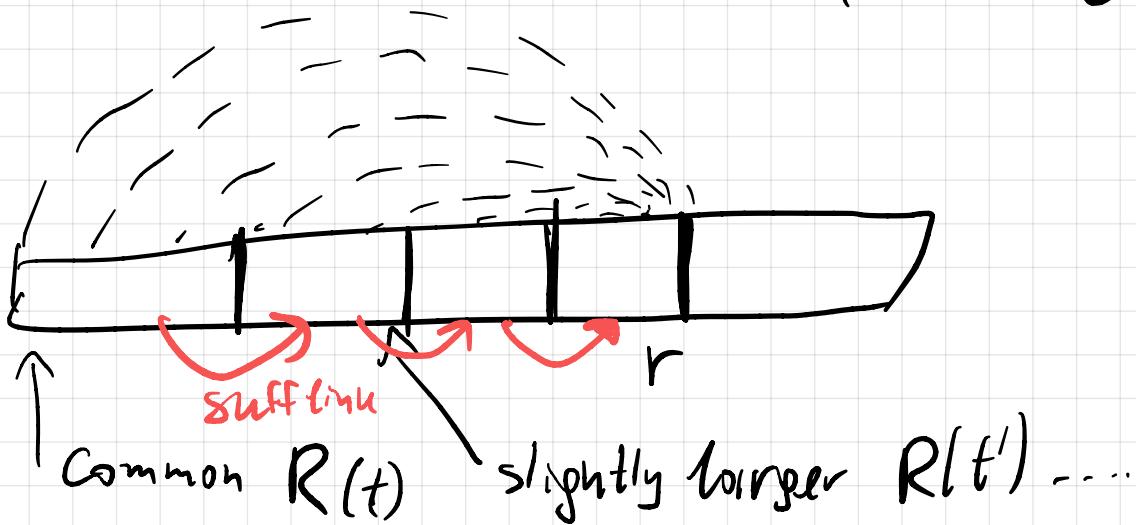
Let $\text{minlen}(v)$ = length of shortest t
s.t. $R(t)$ corresponds to v

Lm: $\text{minlen}(v) = \text{len}(\text{Suf}(v)) + 1$



The Journey of Shorter Strings

Consider strings with fixed r & different ℓ



$SA(s)$ vs $SA(s_c)$?

Reminder:

$$R_s(\alpha) = \{ \beta \mid \alpha\beta - \text{suffix of } s \}$$

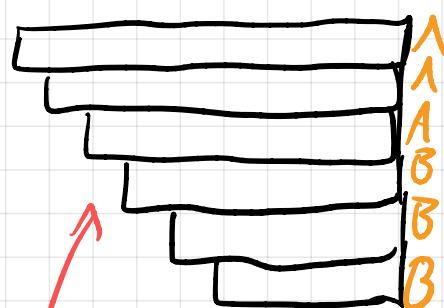
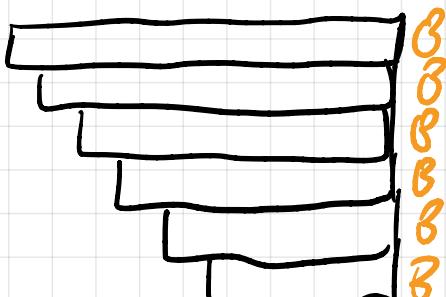
Then

$$R_{sc}(\alpha) = \begin{cases} \{\beta c \mid \beta \in R_s(\alpha)\} & \text{type A} \\ \{\beta c \mid \beta \in R_s(\alpha)\} \cup \{\epsilon\} & \text{type B} \end{cases}$$

↑ iff
 α -suffix sc

Classification

$$\{R_s(t) \mid t\text{-suffix sc}\}$$



'if some string $\{t\}$
is type-B, then t -
suffix sc
 \Rightarrow ALL SHORTER suffixes
ARE TYPE-B too.'

Remember edges:

$$R_s(t) \xrightarrow{b} R_s(tb)$$

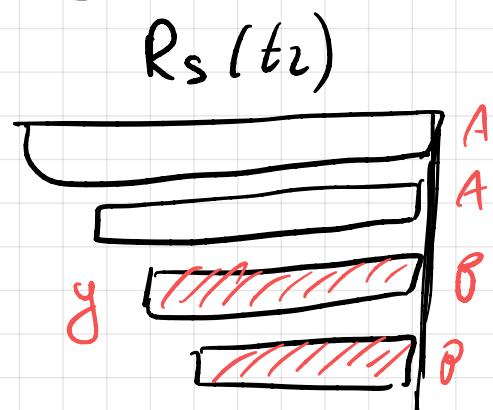
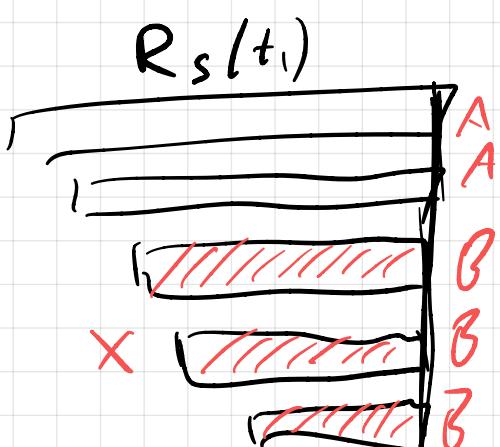
\Rightarrow if states $R_s(t)$ & $R_s(tb)$
change by rule A only
or rule B only,
no need to do
anything.

Otherwise we have to split vertex.

Lm: There is at most one vertex which splits

Proof

Suppose Contrary



Let x, y - any type-B string from

$R_s(t_1), R_s(t_2)$

$\leftarrow (SA(s) \rightarrow SA(sc))$

x - suffix

sc

y - suffix sc

wLOG assume $|x| < |y|$

$\Rightarrow x$ - suffix y

$\Rightarrow R_s(t_1) \supseteq R_s(t_2)$

$\Rightarrow R_s(t_2) \rightarrow R_s(t_1)$

by chain of suf links

\Rightarrow all strings in t_1 are suffixes of y

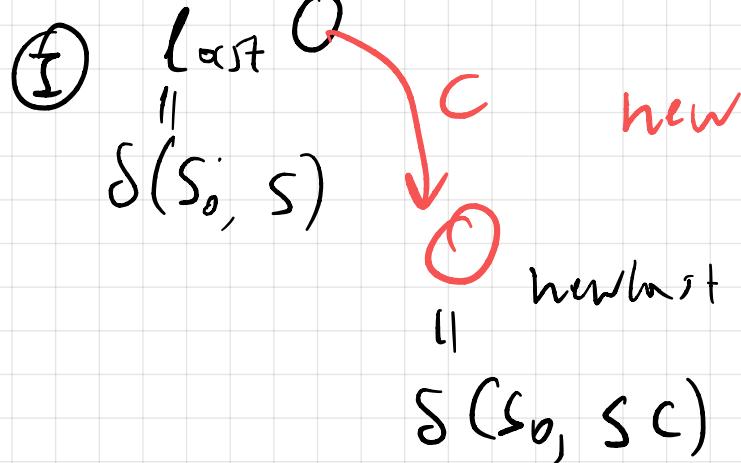
\Rightarrow type-B

Implementation Note

```
struct Node {  
    Node* suf;  
    (unordered_map<char, Node*>) go;  
    int len;  
};
```

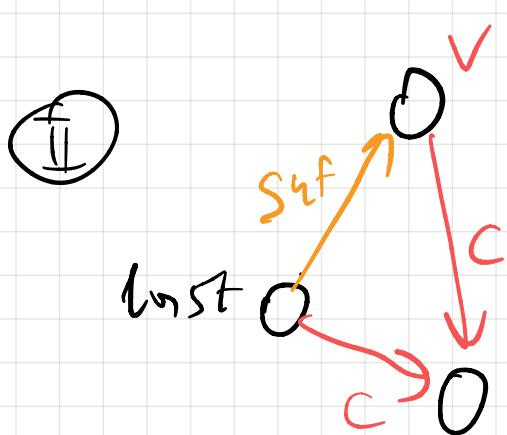
Algorithm

Extend (c)



have to add
this

$$R(S) \rightarrow R(S_c)$$



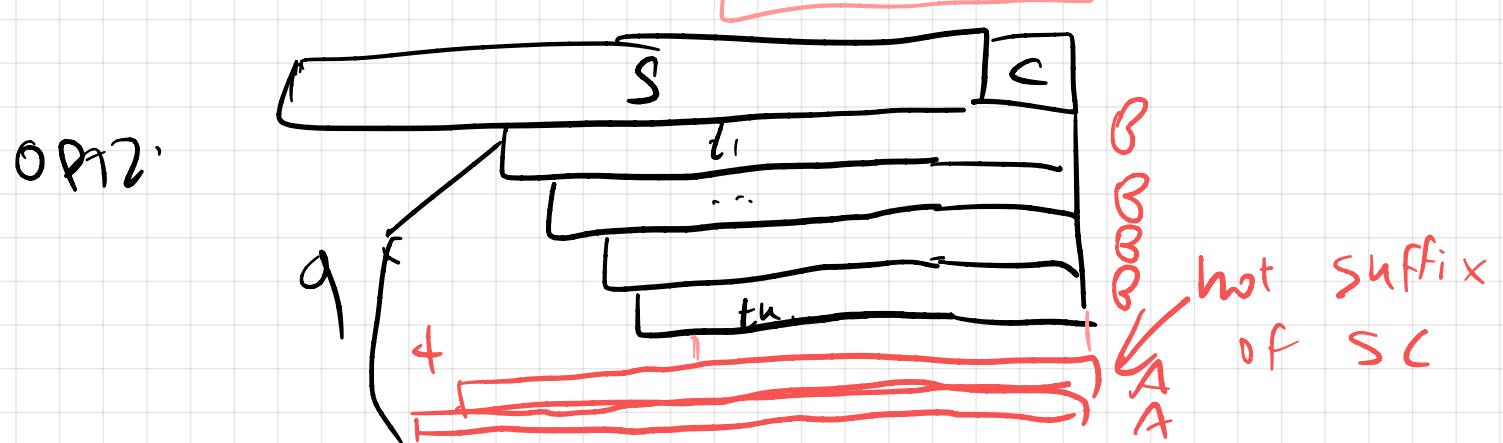
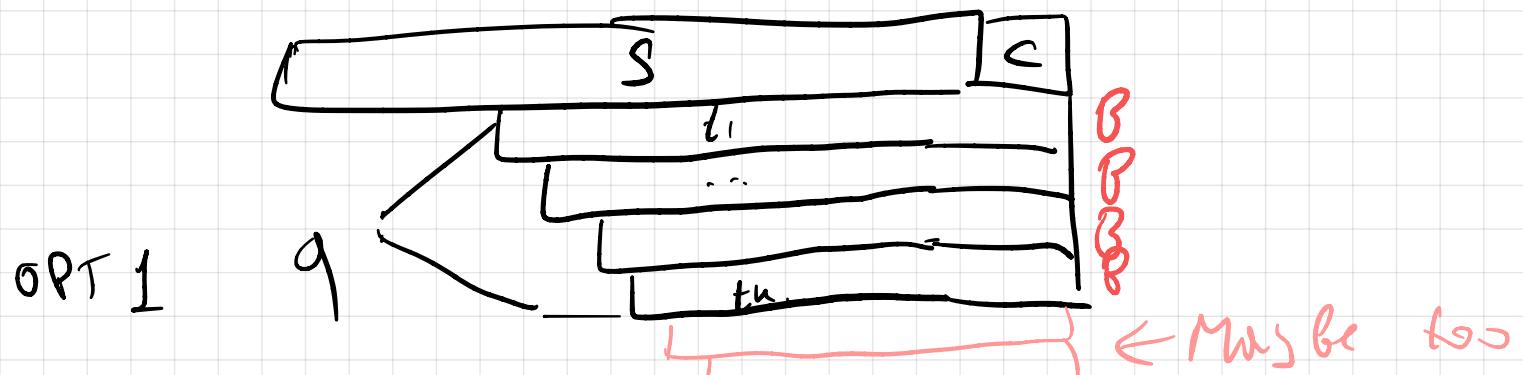
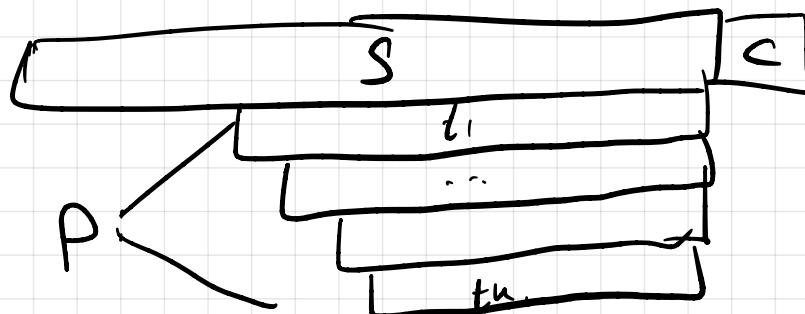
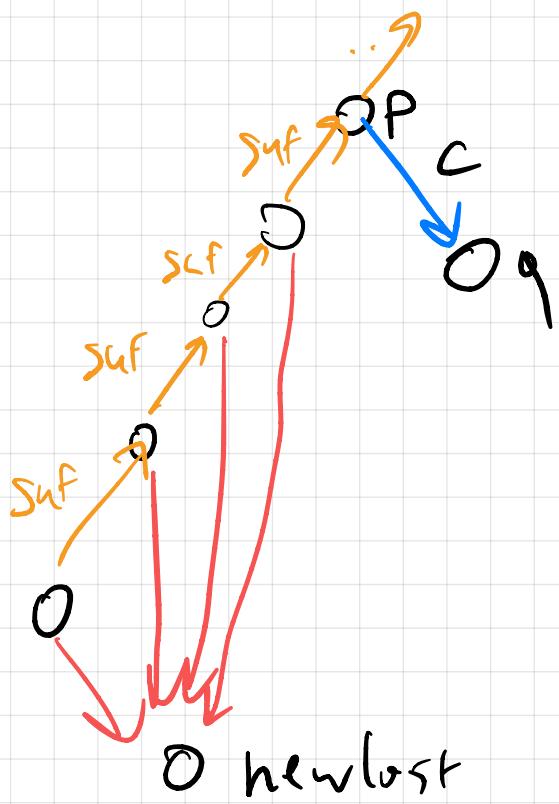
(While there is no
edge c, add it)

$$R_{S_c}(\text{newlast}) = \{\emptyset\},$$

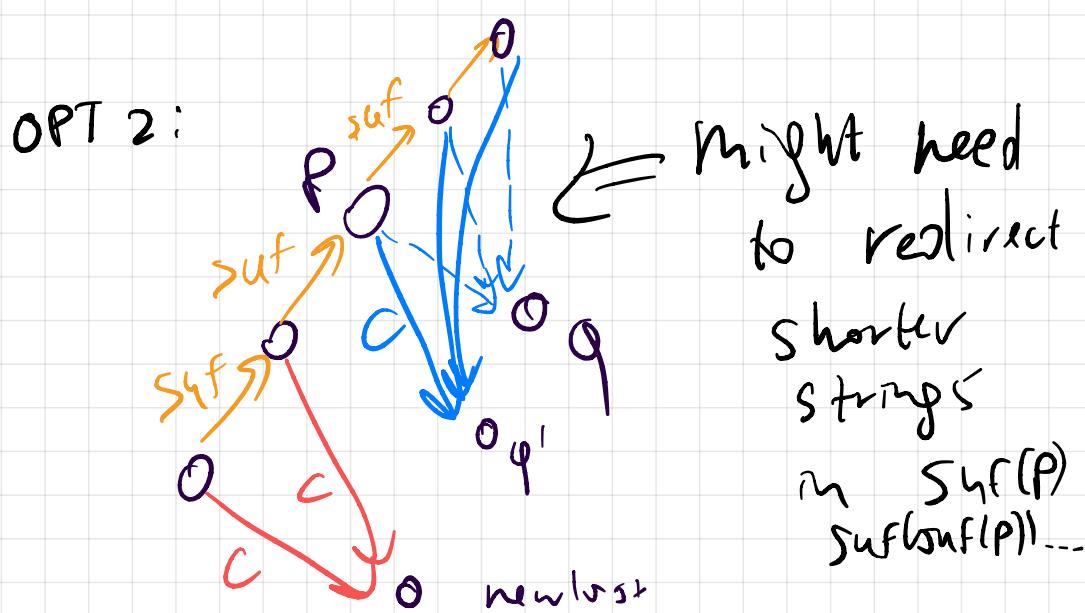
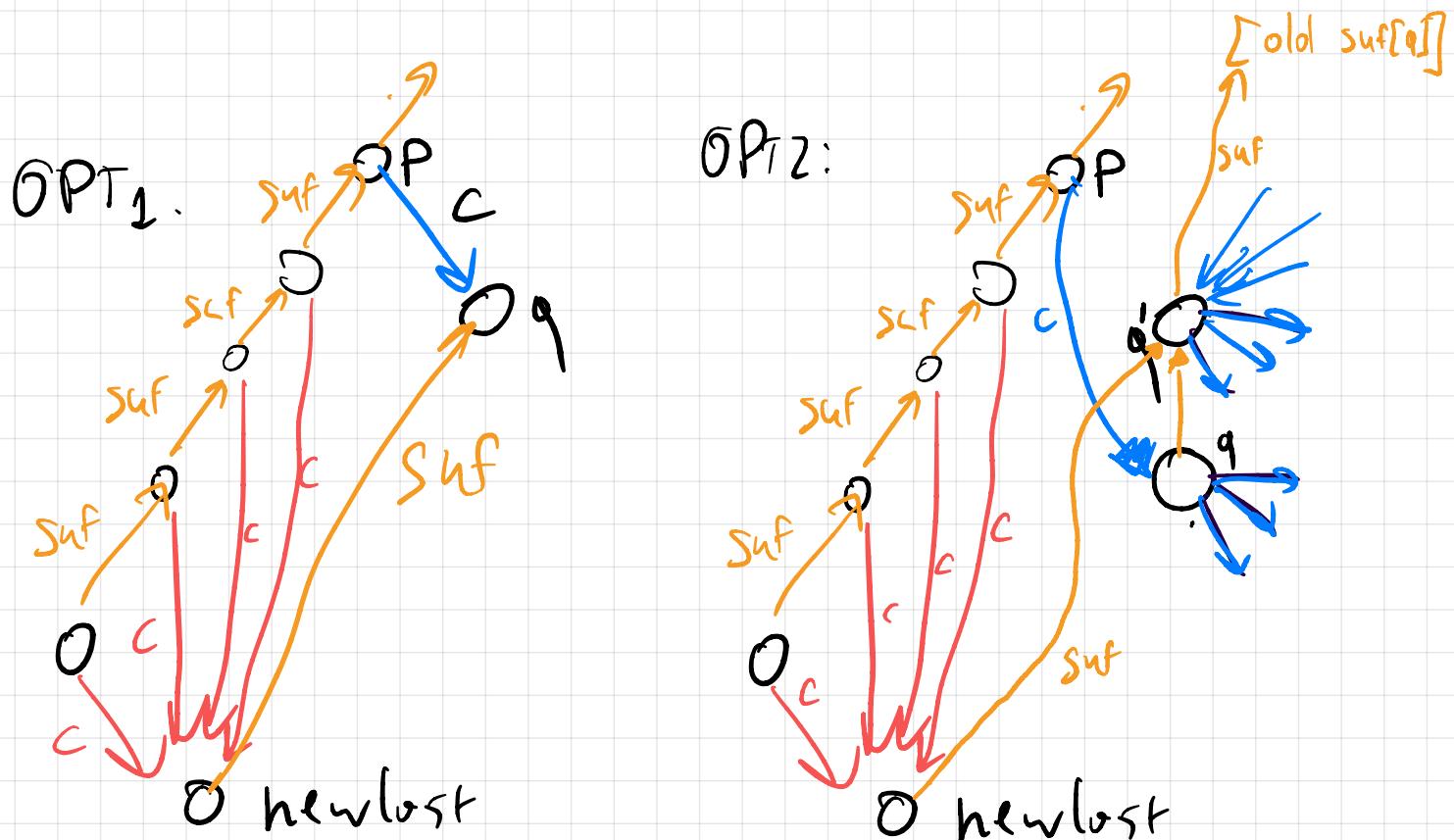
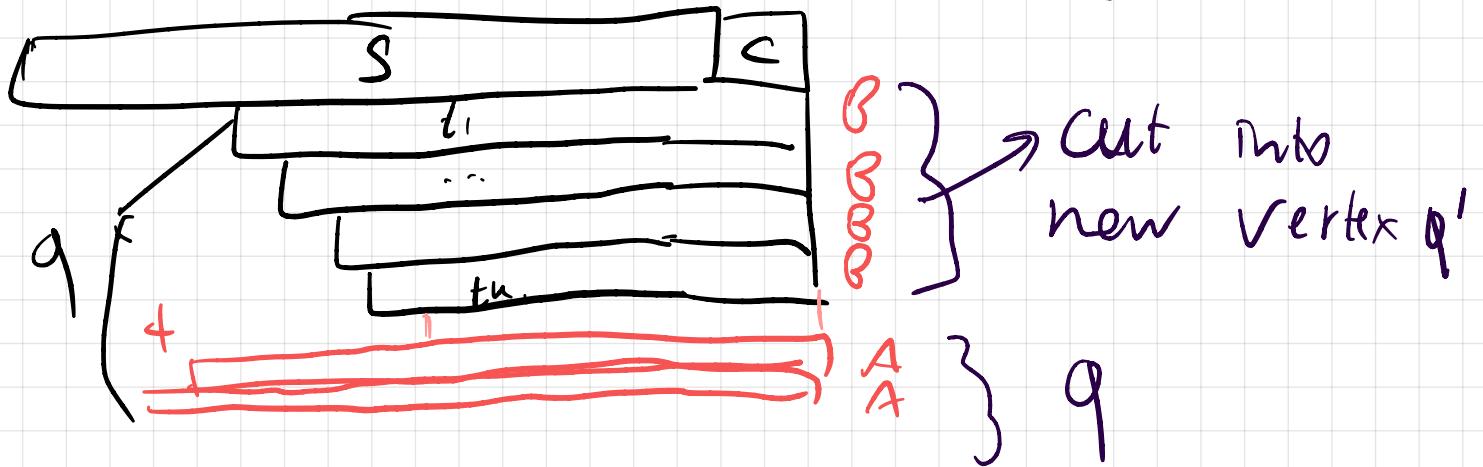
let $t : S(S_0, t) = \checkmark$

$R_{S_c}(tc) = \{\emptyset\}$, otherwise
 $R_S(t)$ should have edge c.

III



SC automation:



Extend(c)

newlast = new node()

newlast.len = last.len + 1

P = last

while P & !(a ∈ P.go):

P.go[c] = newlast (*)

P = P → suf.

if P == null:

newlast.suf = root

return

q = P.go[c]

if P.len + 1 == q.len

newlast.suf = q

else:

q' = new node(len = P.len + 1)

go = P.go (*)

.suf = q.suf)

newlast.suf = q.suf = q'

while P & P.go[c] = q: (***)

P.go[c] = q'; P = P.suf

Linearity:

Goals:

$$\begin{array}{lll} \# \text{Vertices:} & \leq 2n - 1 & (n \geq 2) \\ \# \text{Edges} & \leq 3n - 4 & (n \geq 3) \end{array}$$

$$\text{Complexity: } O(n)$$

Vertices:

$$V(n) = \max \text{ Vert. for } |S|=n.$$

$$V(0) = 1$$

$$V(1) = 2$$

$$V(2) = 3 \leq 2 \cdot n - 1$$

$$V(n: n \geq 2) \leq 2 + V(n-1) \leq 2 \cdot n - 1$$

Edges

def: $P \rightarrow q$ continuous edge if
 $P.\text{len} + 1 = q.\text{len}$

$$\begin{aligned} \# \text{Edges} &= \# \text{Cont Edges} + \# \text{Non Cont Edges} \\ &= V(n) - 1 \end{aligned}$$

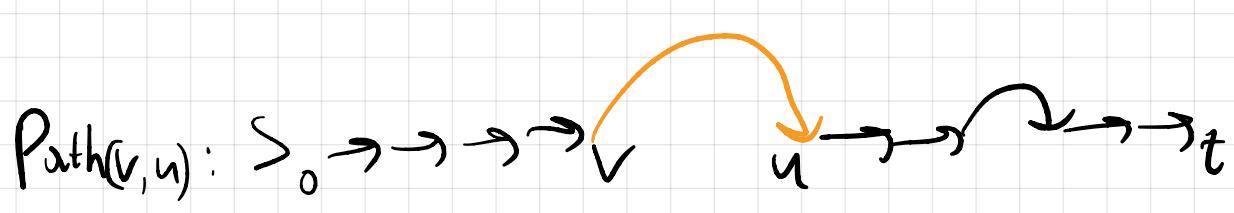
(Each vertex
has exactly 1
incoming edge)

Non Cont Edges:

Let $v \rightarrow_h$ non Cont.

Take $S_0 \rightarrow v$ path with Cont. edges

$u \rightarrow +$ any Path



There $\leq h$ paths $S_0 \rightarrow t$

Note $\{\text{path}(v, u)\}$ are all diff. since
(v, u) - first non Cont edge on
path(v, u)

$$\begin{aligned} \# \text{Edges} &= \# \text{Cont Edges} + \# \text{Non Cont Edges} \\ &= V(h) - 1 \end{aligned}$$

$$\leq 2h - 1$$

$$\leq 3h - 2$$

To improve to $\leq 3h - 3$, note that

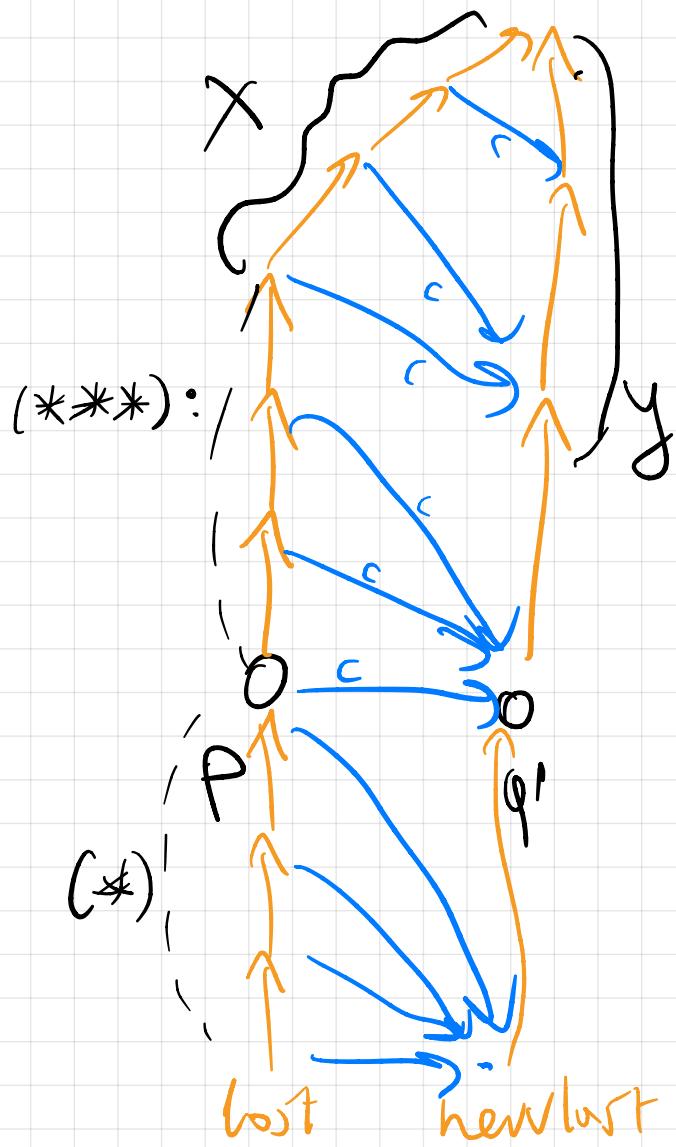
$S_0 \rightarrow t$ path $s \rightarrow \rightarrow \rightarrow \rightarrow t$ doesn't correspond to any long edge. ($3h - 4$ skipped proof)

Time Complexity

There are 3 places in algorithm's

loop not working in $O(1)$

(*) , (***) \leftarrow upper bounded by #Edges



φ : len of suff links
path lost \rightarrow root

$$\varphi_{\text{old}} = x + (****) + (*) + d$$

$$\varphi_{\text{new}} = y + O(1)$$

$$\begin{aligned} (****) + (*) + \varphi_{\text{new}} - \varphi_{\text{old}} &= \\ &= y - x + O(1) \end{aligned}$$

$$y \leq x + 1$$

due to surjectivity
of mapping

Applications

[1]

Check if t -substr s

Compute $S(S_{\text{state}_0}, t)$.

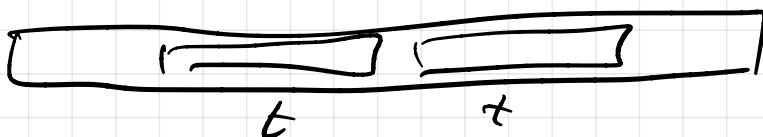
[2]

Number of substrings

$$\begin{aligned} \# \text{substrings} &+= \# \text{new substrings} = \\ &= \text{lost} \rightarrow l(n - \text{lost}) \rightarrow l(n) \end{aligned}$$

[2]

Find t : t has 2 non-overlapping occurrs.



Compute for each v :

$\text{first}(v), \text{lost}(v)$
i.e., if $R(t)$ corresponds to v

$\text{first}(v) = \text{smallest end of occurr.}$
of t

$\text{lost}(v) = \text{largest end of occurr.}$
 $\text{first}(v) \quad \text{of } t -$



$\text{ans} \leftarrow \min(\text{lost}(v) - \text{first}(v), \text{len}(v))$

