

```

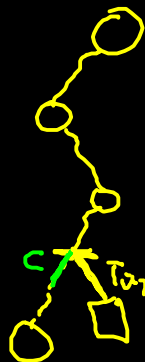
for  $i := 1$  to  $n - 1$  do begin
  // Phase  $i + 1$ 
  for  $j := 1$  to  $i + 1$  do begin
    Traverse  $IB_i$  along the path  $T[j..i]$ ;
    If necessary extend the tree at the
    position reached this way by  $T[i + 1]$ ;
    // Details follow
  end;
end;

```

**Rule 1:** If the path labeled  $T_{j,i}$  ends in a leaf,  $T_{i+1}$  is appended to the label of the edge leading to the leaf.

**Rule 2:** If the path does **not** end in a leaf and there is no possibility to continue it with  $T_{i+1}$ , a new edge to a new leaf is created and labeled with  $T_{i+1}$ . The leaf is labeled  $j$ . If  $T_{j,i}$  ends amidst an edge, additionally a new internal node has to be created at the respective position.

**Rule 3:** If the path can be continued with  $T_{i+1}$  nothing is done.



**Example:**  $T = cbacb$ .

**Caution:** Nested loops + traversal along  $T_{j,i} \Rightarrow$  running time cubic in length of text.

### Tricks:

1.) If for given  $j$  Rule 3 is applied for the first time:

$\Rightarrow$  The path labeled  $T_{j,i}$  can be continued with  $T_{i+1}$ , created when inserting word  $w$ .

$\Rightarrow$  Suffixes of  $w$  inserted in an earlier phase guarantee existence of a continuation of  $T_{j',i}$ ,  $j' > j$  with  $T_{i+1}$ .

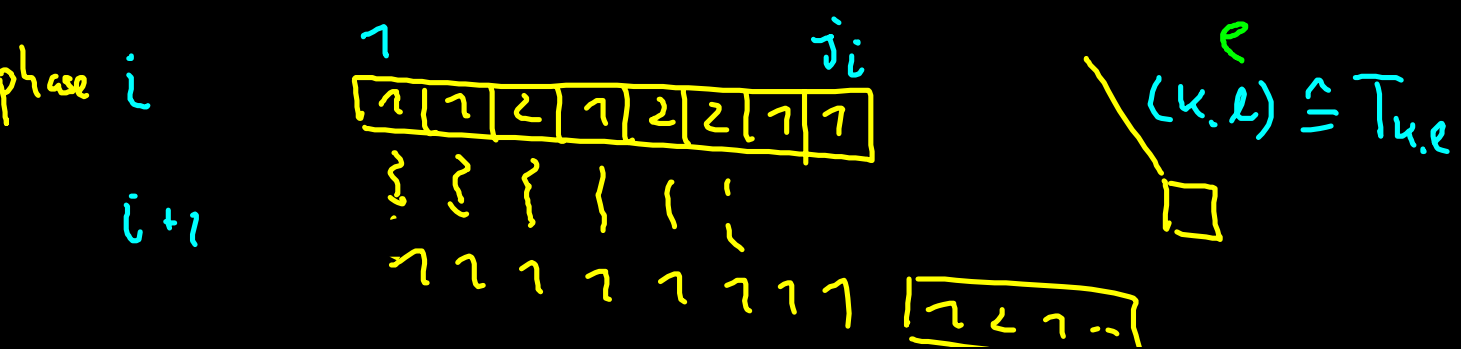
$\Rightarrow$  Rule 3 implies termination of the current phase.

2.) A leaf always remains a leaf. If labeled  $j$  it represents all suffixes of  $T$  starting at position  $j$ . Whenever we insert  $T_{j,i+1}$  in a later phase, we reach leaf  $j$  and apply rule 1. With 1.) it follows that:

- ▶ Each phase  $i + 1$  starts with a sequence of extensions (beginning with  $j = 1$ ), in which only rules 1 and 2 are applied.
- ▶ Let  $j_i$  denote the last extension in phase  $i$  in which rule 1 or rule 2 is applied. As each application of rule 2 generates a new leaf  $j_i \leq j_{i+1}$  holds. Thus the initial sequence of applications of rules 1 and 2 can't become shorter in later phases.

⇒ All extensions for  $j \in [1 : j_i]$  in phase  $i + 1$  will be by rule 1 because either there was already a leaf labeled  $j$  in phase  $i$  (in this case only its edge label is extended with  $T_{i+1}$ ) or in phase  $i$  the second rule has been applied for  $T_{j,i}$  (we consider  $j \in [1 : j_i]$ , so rule 3 is ruled out for extension  $j$ ) so now there exists a leaf labeled  $j$  for which the edge label is extended.

⇒ We don't have to do extensions  $1..j_i$  **explicitly** if we mark the leaves' edges with  $(p, e)$ ,  $e$  a global symbol meaning current end of text. ( $e$  is set to  $i + 1$  in phase  $i + 1$ .)



3.) For  $j \in [j_i + 1, i + 1]$  use rule 2 or 3.

- ▶ rule 3  $\Rightarrow j_{i+1} := j - 1$  (i.e.  $j = j_{i+1} + 1$ ); terminate.
- ▶ phase terminated by different rule  $\Rightarrow j_{i+1} := i + 1$ .

**Observation:** Two consecutive phases have at most (WC is that rule 3 ends phase) one  $j$  in common, for which both do **explicit** extensions:

**Phase 2:** compute explicit extensions for  $j = j_1 + 1 \dots \underline{j_2 + 1}$ ,

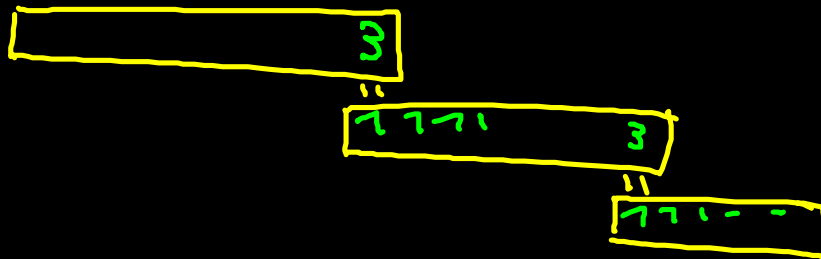
**Phase 3:** compute explicit extensions for  $j = \underline{j_2 + 1} \dots j_3 + 1$ ,

...

**Phase  $i-1$ :** compute explicit extensions for  $j = j_{i-2} + 1 \dots j_{i-1} + 1$ ,

**Phase  $i$ :** compute explicit extensions for  $j = j_{i-1} + 1 \dots j_i + 1$ .

$\Rightarrow |T|$  many phases and  $j_i \leq n$  imply at most  $2n$  explicit extensions.



```

Create tree  $IB_1$ ;
 $j[1] := 1$ ; // leaf 1 already exists
for  $i := 1$  to  $n - 1$  do begin
    // Phase  $i+1$ 
     $e := i + 1$ ; // all implicit extensions
     $j[i + 1] := i + 1$ ; // no application of rule 3
    for  $j := j[i] + 1$  to  $i + 1$  do begin
        Traverse  $IB_i$  along path  $T[j..i]$ ;
        If necessary, extend  $IB_i$  by  $T[i + 1]$ ;
        if (rule 3 was used) then begin
             $j[i + 1] := j - 1$ ;
            End phase  $i + 1$ ;
        end;
    end;
end;

```

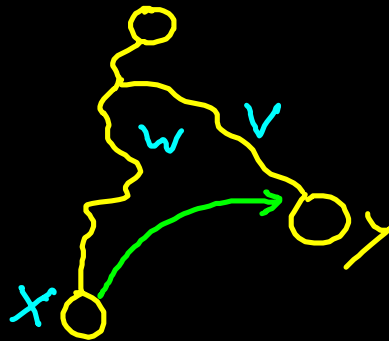
**Example:**  $T = acacag$ .

4.) Speed up traversal of edges labeled with more than one character by only evaluating the first character. The position at which to continue in the *traversal word* is determined from the indices saved for the edge. (We already noticed that in phase  $i + 1$  each word  $T_{j,i}$  is present in the tree.)

5.) Add utility links:

### Definition

Let  $w = u \cdot v$ ,  $u \in \Sigma$ ,  $v \in \Sigma^*$ , and  $IB$  an implicit suffix tree. If  $IB$  contains an internal node  $x$ , reached from the root by  $w$  and a node  $y$  with path-label  $v$ , the suffix link of  $x$  points to  $y$ .



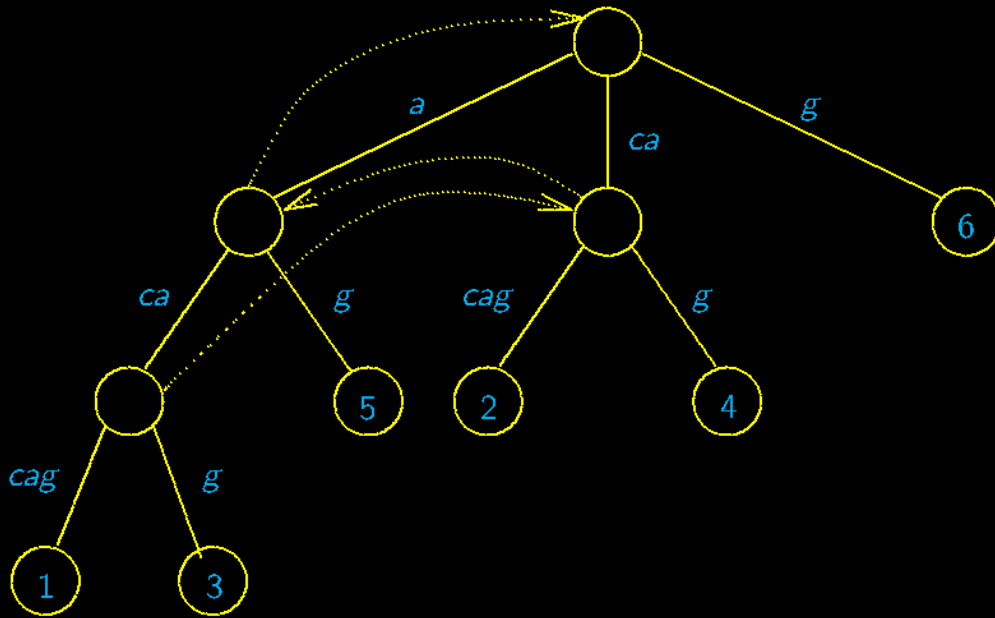
**Advantage:** The places of explicit extensions can be found without traversing the tree (starting at the root) each time. If in phase  $i + 1$  the extension for  $T_{j,i+1}$  has to be done, the tree has to be traversed along  $T_{j,i}$  and the existence of  $T_{i+1}$  at the position  $x$  reached this way has to be checked.

The next extension then considers  $T_{j+1,i+1}$ . The node reached by  $T_{j+1,i}$  is found via the suffix link of  $x$  (if there is an internal node). The same holds for the next extension (a suffix link leads possibly directly to the node reached by  $T_{j+2,i}$ ) and so on.

$x$  no internal node: Return to  $y$ , the last node on the way to  $x$  (we can always save this node) and follow this node's suffix link.  $\Rightarrow$  The node reached thusly is reached by a prefix of  $T_{j+1,i}$ . From there we have to continue using the symbols *between*  $y$  and  $x$ .

**Consequence:** Assuming each internal node has a suffix link, an **explicit** extension has constant amortized running time.

**Example:** Extend tree by  $T_{j,i+1} = acag$  and its suffixes:



### Lemma

If a new internal node  $x$  with path label  $u\alpha$ ,  $u \in \Sigma$ , is created by extension  $j$  of phase  $i + 1$  then either  $\alpha$  already ends at an internal node of the current tree or this node is created at the next extension (extension  $j + 1$  of phase  $i + 1$ ).

**Proof:** New internal node  $\Leftrightarrow$  rule 2  $\Leftrightarrow$  path  $T_{j,i}$  ends amidst an edge label with next character  $c$  not equal  $T_{i+1}$ .

$\Rightarrow$  Earlier phase inserted  $T_{j,i} \cdot c$ .

$\Rightarrow$  The same phase afterwards processed  $T_{j+1,i} \cdot c$ .

$\Rightarrow \exists$  path  $\mathcal{P}$  for  $T_{j+1,i}$  in the tree.

a)  $\mathcal{P}$  can only be continued with  $c$ :

$\Rightarrow$  Extension by  $T_{j+1,i+1}$  creates an internal node at the position considered.

b)  $\mathcal{P}$  can be continued with various symbols:

$\Rightarrow$  At the position considered an internal node must be present.  $\square$

## Corollary

For Ukkonen's algorithm we can assure that each internal node created has a suffix link after the following extension. This requires constant extra time.

**Proof:** Induction:  $IB_1$  has no internal nodes relevant for suffix links (at the root  $u = \epsilon$ ).

Assumption: Assumption holds after phase  $i$ .

Lemma 2  $\Rightarrow$  The target node of a node created in the  $j$ -th extension of phase  $i + 1$  will be present after the  $(j + 1)$ -th extension of the same phase. (This gives a simple method to create suffix links: We remember nodes created by rule 2 and add the suffix link after reaching or creating the target node during the next extension.)

Since it is impossible that the last extension of a phase (considering the single character  $T_{i+1}$ ) creates a new internal node all new internal nodes will have a suffix link after the  $(i + 1)$ -th phase. □

**Running time:** At most  $2 \cdot |T|$  explicit extensions with constant amortized cost lead to running time linear in  $|T|$ .

**Compact suffix tree:** Add a  $|T| + 1$ -th phase to the algorithm with  $T' = T \cdot \$$ . Afterward replace  $\epsilon$  with  $|T'|$  by a tree traversal (linear time)

$\Rightarrow$  compact suffix tree for  $T'$ .

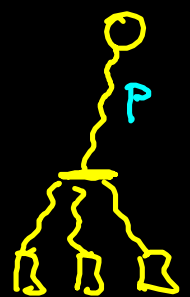
## Applications

### String-Matching:

Text  $T$  fixed, String  $P$  varies (suffix tree reasonable only in this case!).

$\Rightarrow$  Running time in  $\mathcal{O}(|P| + k)$  for  $k$  the number of occurrences of  $P$  in  $T$ .

(By assumption the subtree reached via  $P$  has  $k$  leaves and thus at most  $2k - 1$  nodes and can be traversed in time  $\mathcal{O}(k)$ .)



**Set-Matching:** Sequential Search of all  $P_i$  in the suffix tree has the same running time bound as the Aho-Corasick algorithm.

**But:** Aho-Corasick creates search term tree of size  $\mathcal{O}(m)$ ,  
 $m := \sum_{1 \leq i \leq N} |P_i|$ , in time  $\mathcal{O}(m)$ , and searches in time  $\mathcal{O}(n)$ ,  
 $n := |T|$ .

Suffix tree has size  $\mathcal{O}(n)$ , construction time  $\mathcal{O}(n)$  and search time  $\mathcal{O}(m)$ .

$\Rightarrow$  If all  $P_i$  together are longer than the text, the suffix tree solution needs less space but more time (preprocessing ignored). If the set of the strings is shorter than the text, the Aho-Corasick needs less space but more time.

$\Rightarrow$  We observe a place/time-trade-off, as no solution is superior in place and time consumption at the same time.

## Substring Problem for a Set of Texts:

Look for string  $P$  in a set of texts  $\mathcal{T} := \{T_i \mid 1 \leq i \leq N\}$ .

Example: Newly sequenced DNA fragment ( $P$ ) is to be searched amongst DNA sequences in the database ( $\mathcal{T}$ ) (Identification by mitochondrial DNA).

Generalized compact suffix tree: compact suffix tree for text  $T' = T_1\$_1 T_2\$_2 \cdots T_N\$_N$  (all  $\$_i$  different).

Leaves are labeled with pairs (text, position), symbols following an end mark are removed (as the  $\$_i$  are different, this only happens at leaf edges).

**So:** Create generalized suffix tree for  $T'$  in time and space  $\mathcal{O}(|T'|)$  and traverse along  $P$  to solve the substring problem for the given set of texts.

Dead end  $\Leftrightarrow P$  not contained.

Leaf  $\Leftrightarrow P$  contained exactly once (label identifies text and position).

Internal node  $\Leftrightarrow P$  contained multiple times (visit all leaves in time proportional to number of occurrences).

# Longest Common Substring:

Search a longest substring, common to all words in  $\mathcal{T} := \{T_i \mid T_i \in \Sigma^*, 1 \leq i \leq N\}$ .

Example: Identify important and thus in related organisms *mutationless* regions of DNA.

Solution: construct generalized suffix tree for  $\mathcal{T}$ .

Inner node  $x$  corresponds to

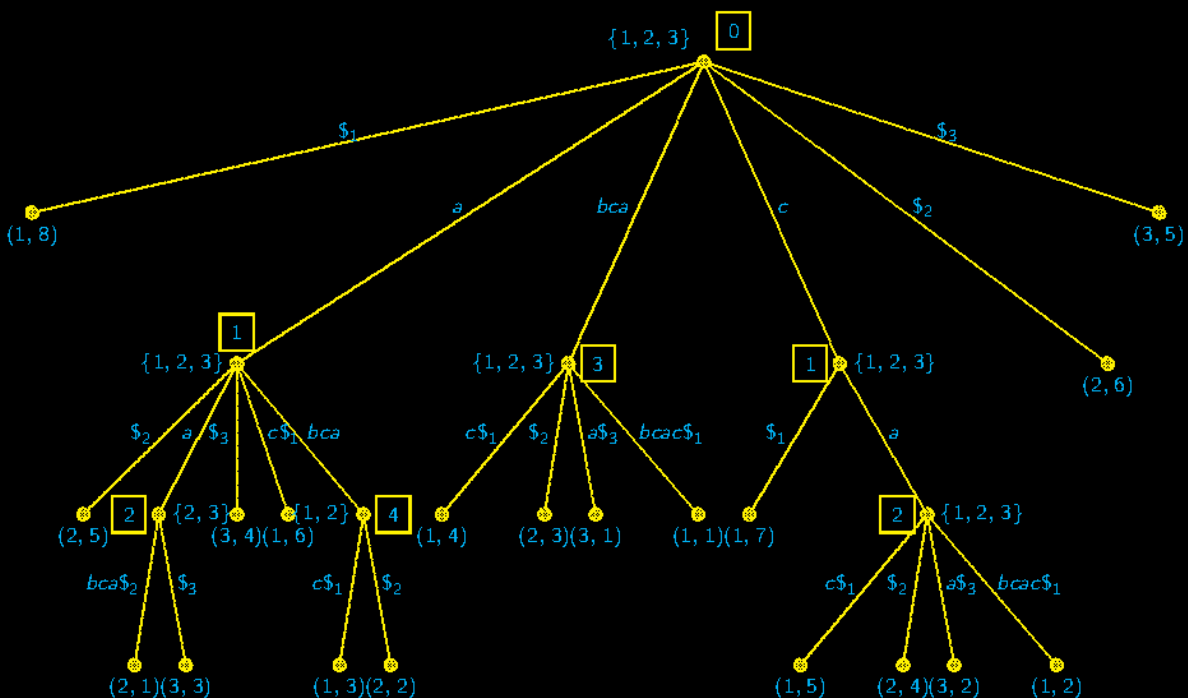
- ▶ prefix of a suffix of only one of the  $T_i$  (all leaves in the subtree with root  $x$  belong to the same  $T_i$ ), or
- ▶ prefix of a suffix of multiple texts (leaves in the subtree with root  $x$  belong to different  $T_i$ ).

Thus we label each internal node  $x$  with the set of text indices appearing in the subtree with root  $x$  (This can e.g. be achieved by traversing the tree in postorder and labeling each node with the union of the labels of its children).

Now only such nodes are solution candidates that are labeled with the complete set  $\{1, 2, \dots, N\}$  (the solution has to appear in each of the  $T_i$ ).

From these nodes we chose one with maximum string depth, i.e. a node  $y$  which has a path-label  $\alpha$  of maximum length. This can be done with a tree traversal.

Now the searched substring is given by  $\alpha$ .





## Repeats in Words:

### Definition

Let  $T \in \Sigma^n$  and  $P \in \Sigma^m$ ,  $0 < m \leq n$  and let  $T_0 := \$ =: T_{n+1}$ ,  $\$ \notin \Sigma$ .  $P$  is called an exact repeat of  $T$ , if and only if

$$(\exists i, j \in [0 : n - 1])(i \neq j \wedge P = T_{i+1, i+m} = T_{j+1, j+m).$$

If additionally  $T_i \neq T_j$  and  $T_{i+m+1} \neq T_{j+m+1}$ ,  $P$  is called a maximum repeat in  $T$ .

This definition does not rule out the possibility of several different maximum repeats starting at the same position of  $T$ .

**Example:** In  $T = aabcbabacabcc$ ,  $ab$  (because of  $aabcbabacabcc$ ) and  $abc$  (because of  $aabcbabacabcc$ ) are maximum repeats, one of which starts at  $T_2$ .

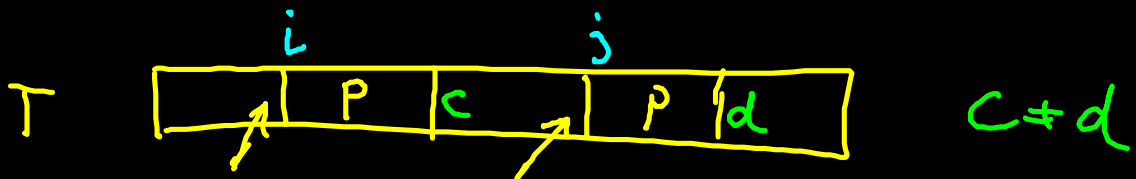
### Lemma

Let  $T$  be a word given as compact suffix tree and  $P$  a maximum repeat in  $T$ . Then there exists an internal node  $x$  with path label  $P$  in the tree.

Note that this lemma implies a maximum number of  $n - 1$  maximum repeats in  $T \in \Sigma^n$ .

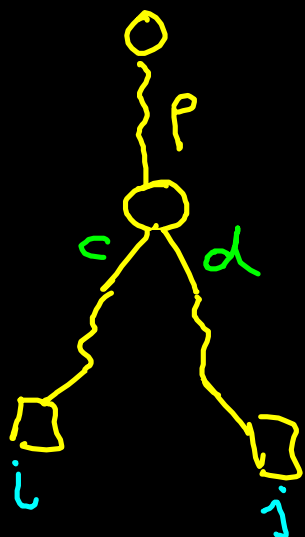
Proof:  $P$  max. repeat  $\rightarrow \exists 2$  positions (different)

for  $T$  where  $P$  occurs as substring



$\rightarrow \exists 2$  different suffixes with prefix  $P$

→ suffix tree has 2 leaves with path labels  $l_1$  and  $l_2$  both starting with  $P$



### Definition

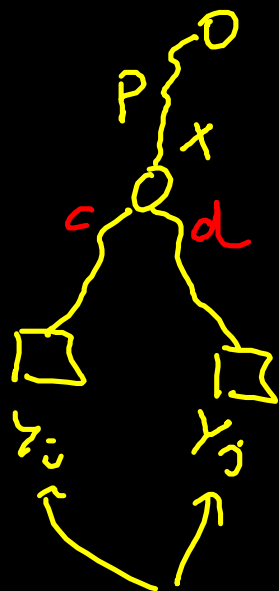
Given a compact suffix tree  $t$  for text  $T$  an internal node  $x$  of  $t$  is called left divers, if the subtree with root  $x$  contains two leaves with labels  $i$  and  $j$  for which  $T_{i-1} \neq T_{j-1}$  holds (Remember that the label of the leaves denotes the position in the text where the respective suffix starts).

### Theorem

Let  $T \in \Sigma^n$  be given as compact suffix tree. Then  $P \in \Sigma^+$  is a maximum repeat in  $T$ , if and only if traversal along  $P$  leads to a left divers internal node.

Proof:

1. case



$T_{i-1} P c$  resp.

$T_{j-1} P d$  are prefixes

of  $T_{0..n}$  resp.  $T_{j-1..n}$   
with  $c \neq d$

→ max. report.