

# Algorithmen II

## Vorlesung am 19.12.2013

INSTITUT FÜR THEORETISCHE INFORMATIK · PROF. DR. DOROTHEA WAGNER

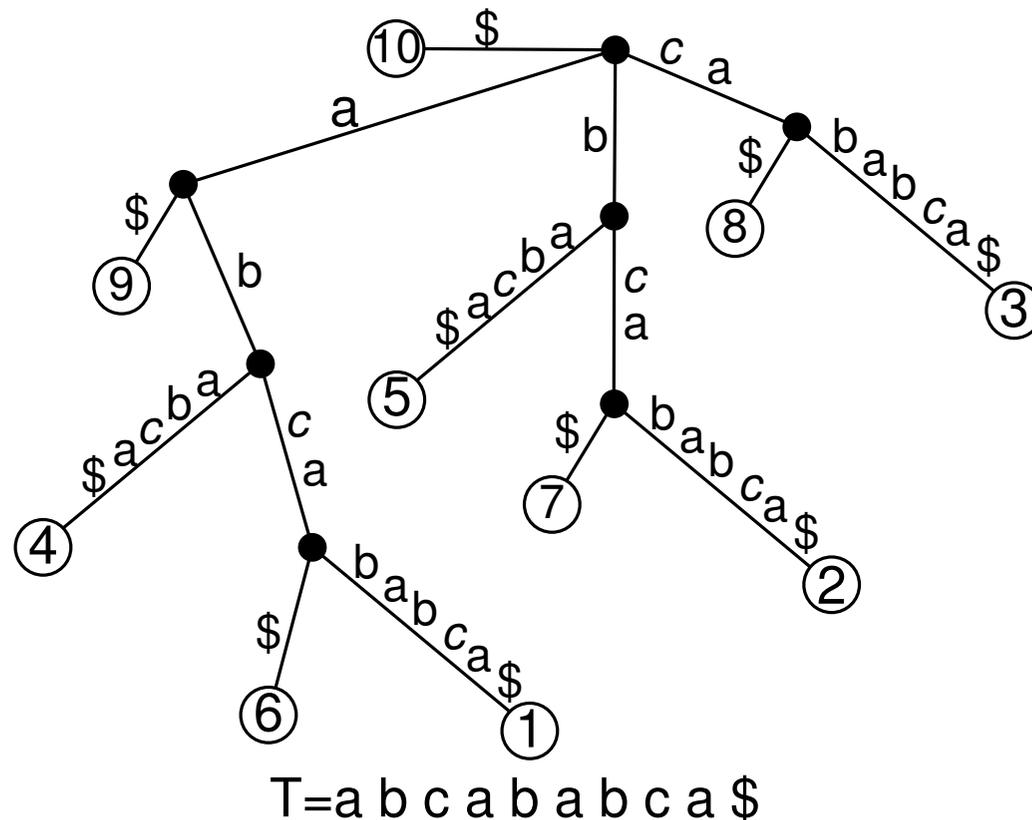


# Problemstellung

**Gegeben:** Text  $T$  der Länge  $n$ .

**Fragestellung:** Wie kann  $T$  vorverarbeitet werden, so dass Suchanfragen auf  $T$  schnell möglich sind?

**Idee:** Repräsentiere  $T$  als Suchbaum.

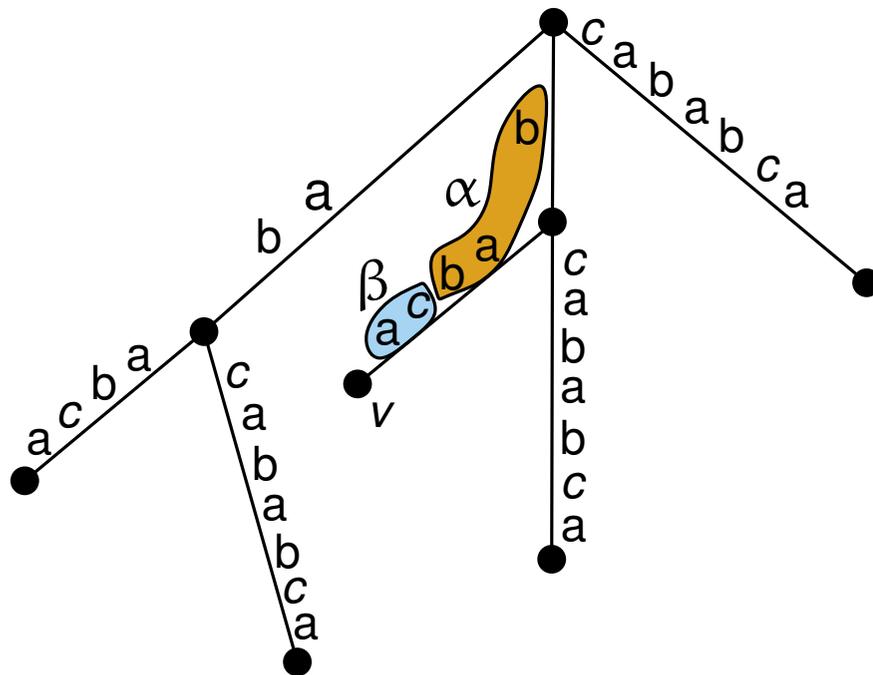




## Definition 8:

- $\bar{v} :=$  Konkatenation der Kantenbeschriftungen auf dem Pfad von  $r$  zu  $v$
- $d(v) := |\bar{v}|$  heißt String-Tiefe von  $v$ .
- $\mathcal{T}$  enthält  $\alpha \in \Sigma^*$ , falls es Knoten  $v \in V$  und Wort  $\beta \in \Sigma^*$  gibt, so dass  $\bar{v} = \alpha\beta$ .
- $\text{words}(\mathcal{T}) := \{\alpha \in \Sigma^* \mid \mathcal{T} \text{ enthält } \alpha\}$

$T = a b c a b a b c a$



$\bar{v} = babca$

$d(v) = |\bar{v}| = 5$

$\mathcal{T}$  enthält  $bab$ , da für  $\alpha = bab$  und  $\beta = ca$  es Knoten  $v$  gibt mit  $\bar{v} = \alpha\beta$ .

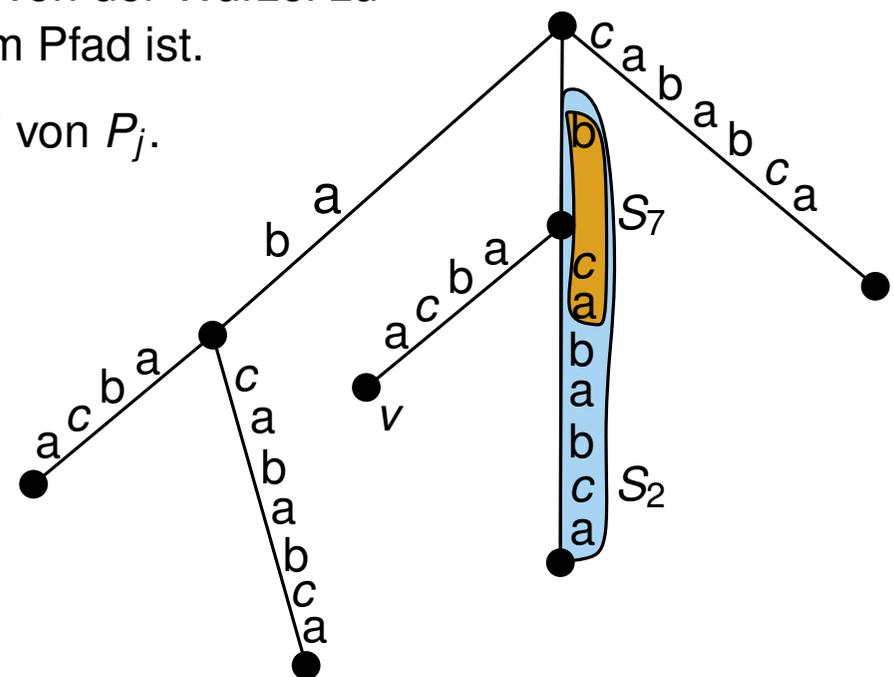
**Definition:** Ein *Suffixbaum* ist ein kompakter  $\Sigma^*$ -Baum  $S$ , sodass  $S$  die Infixe von  $T$  enthält:  

$$\text{words}(S) = \{ T[i, j] \mid 1 \leq i \leq j \leq n \}$$

## Beobachtung:

- Für jedes Suffix  $S_i$  gibt es einen Pfad  $P_i$ , der von der Wurzel zu einem Blatt führt, sodass  $S_i$  Präfix von diesem Pfad ist.
- Wenn  $S_i$  Präfix von  $S_j$  ist, dann ist  $P_i$  Teilpfad von  $P_j$ .

$T = a \text{ b c a b a b c a}$



**Später:** Es gibt für jeden Text  $T$  einen Suffixbaum.

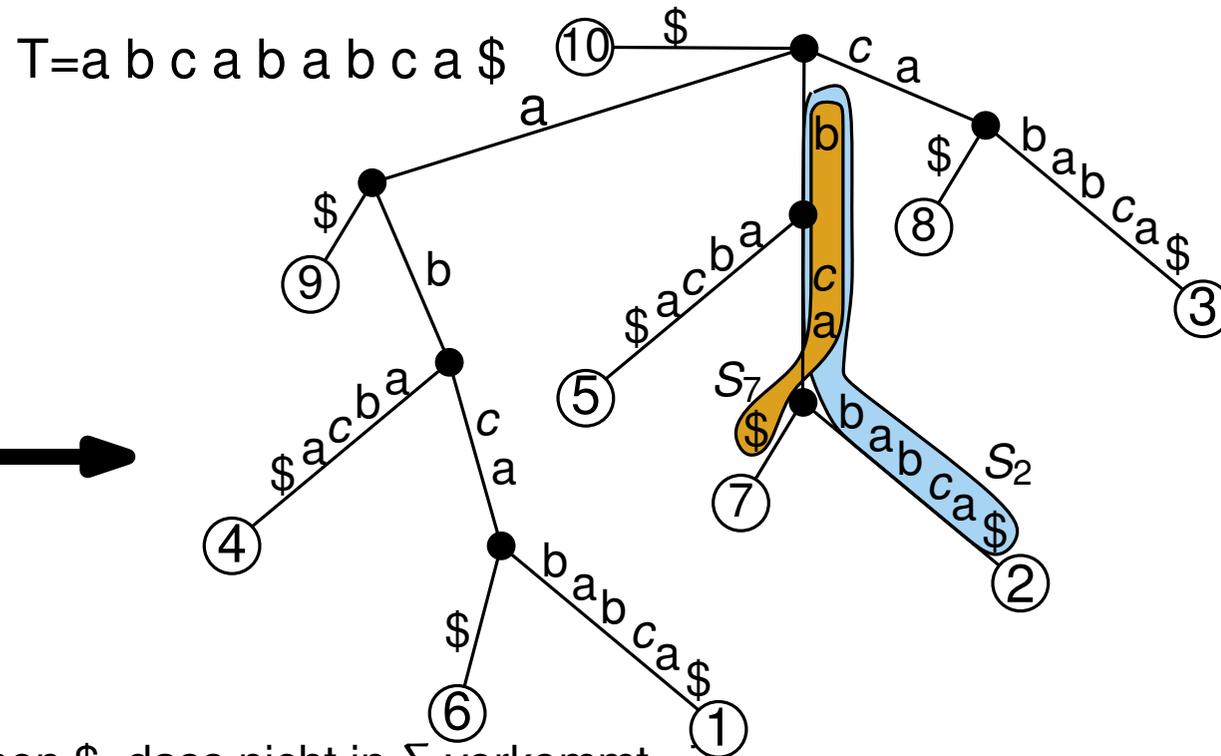
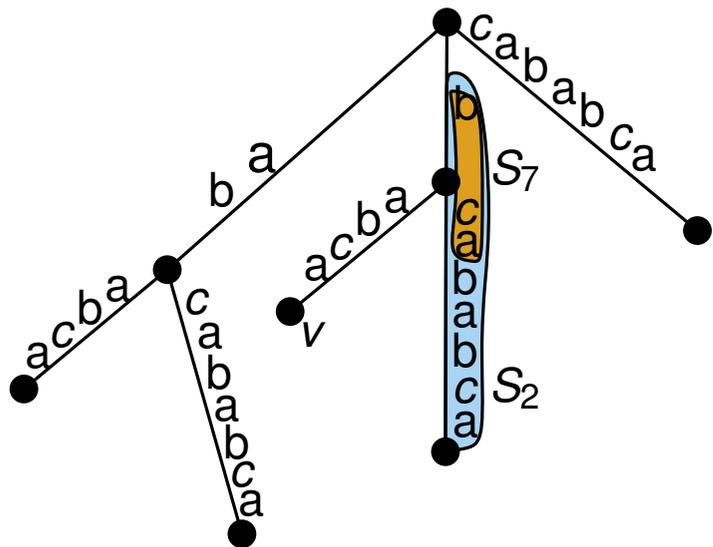
## Notation:

1.  $S_i$  bezeichnet das Suffix von  $T$  beginnend ab der  $i$ -ten Stelle.
2.  $T[i, j]$  bezeichnet den Teilstring von  $T$  von der  $i$ -ten bis zur  $j$ -ten Stelle.

# Suffixbaum

**Definition:** Ein *Suffixbaum* ist ein kompakter  $\Sigma^*$ -Baum  $S$ , sodass  $S$  die Infixe von  $T$  enthält:  
 $\text{words}(S) = \{ T[i, j] \mid 1 \leq i \leq j \leq n \}$

$T = a b c a b a b c a$



Nützlicher Trick: Hänge an  $T$  ein Zeichen  $\$$ , das nicht in  $\Sigma$  vorkommt.

1. Suffix  $S_i$  kann nicht Präfix eines anderen Suffixes  $S_j$  sein.
2. Jedes Suffix endet an einem Blatt.
3. Es gibt so viele Blätter wie Suffixe.

⇒ nummeriere Blätter so, dass  $i$ -te Blatt dem Suffix  $S_i$  entspricht.

Ab jetzt:  
 $S_i := i$ -te Suffix  
 $i :=$  Pfad von Wurzel zu Blatt  $i$ .

# Suchen im Suffixbaum

T = a b c a b a b c a

**Gegeben:** Suffixbaum  $S$  und Muster  $P$  mit  $|P| = m$

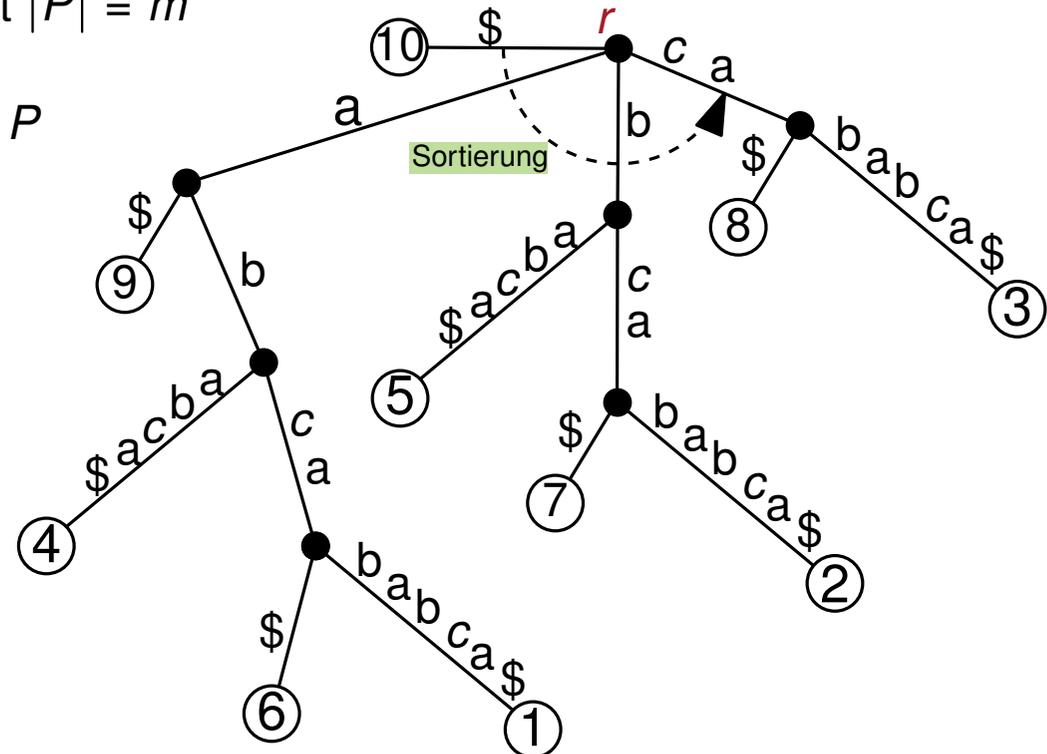
**Gesucht:** Anzahl Vorkommen von  $P$

**Idee:** Suche in  $S$  alle Pfade  $S_i$ , so dass  $P$  Präfix von  $S_i$  ist.

Ausgehend von Wurzel  $r$ :

1. Wähle nächste ausgehende Kante mithilfe von binärer Suche.
2. Vergleiche  $P$  und Kantenbeschriftungen schrittweise.

**Beispiel:**  $P = a b c a$



# Suchen im Suffixbaum

T = a b c a b a b c a

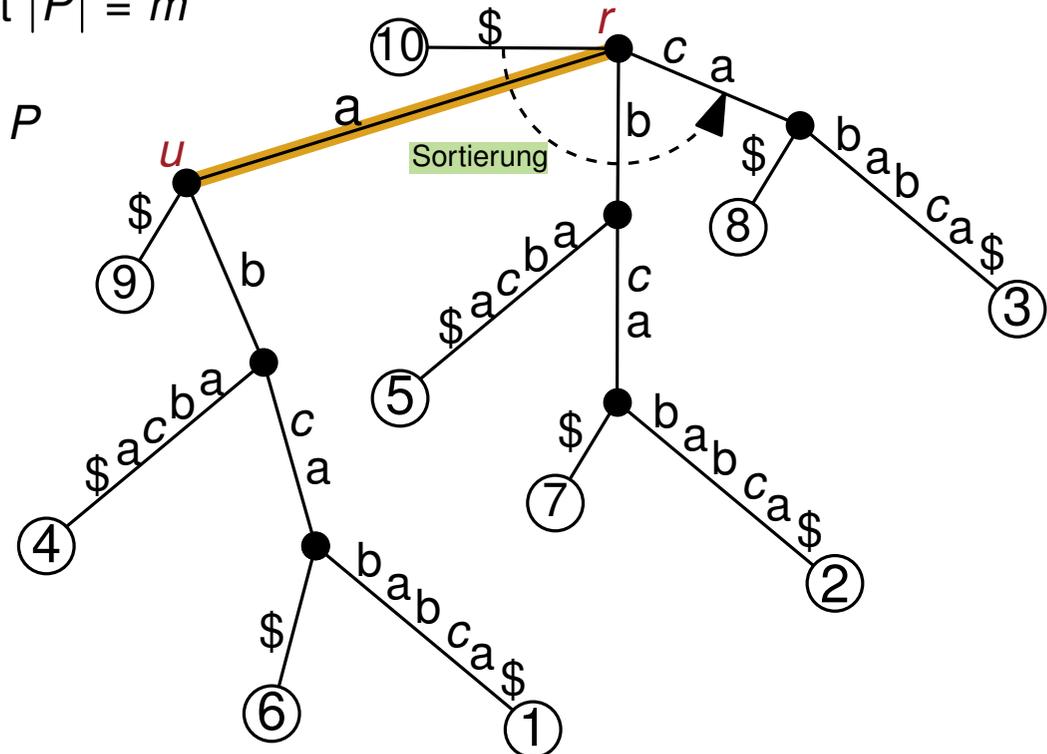
**Gegeben:** Suffixbaum  $S$  und Muster  $P$  mit  $|P| = m$

**Gesucht:** Anzahl Vorkommen von  $P$

**Idee:** Suche in  $S$  alle Pfade  $S_i$ , so dass  $P$  Präfix von  $S_i$  ist.

Ausgehend von Wurzel  $r$ :

1. Wähle nächste ausgehende Kante mithilfe von binärer Suche.
2. Vergleiche  $P$  und Kantenbeschriftungen schrittweise.



**Beispiel:**  $P = a b c a$

**1. Schritt:** Wähle Kante  $(r, u)$  ausgehend von  $r$ , deren Beschriftung  $B$  mit  $P[1] = a$  beginnt.

# Suchen im Suffixbaum

T = a b c a b a b c a

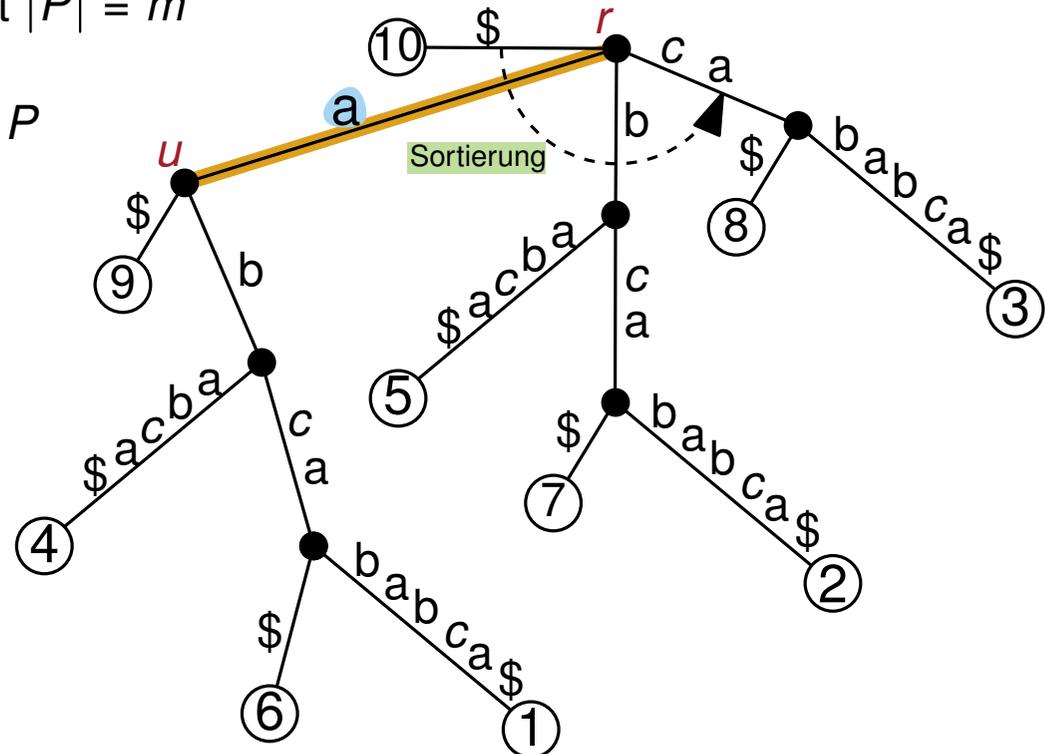
**Gegeben:** Suffixbaum  $S$  und Muster  $P$  mit  $|P| = m$

**Gesucht:** Anzahl Vorkommen von  $P$

**Idee:** Suche in  $S$  alle Pfade  $S_i$ , so dass  $P$  Präfix von  $S_i$  ist.

Ausgehend von Wurzel  $r$ :

1. Wähle nächste ausgehende Kante mithilfe von binärer Suche.
2. Vergleiche  $P$  und Kantenbeschriftungen schrittweise.



**Beispiel:**  $P = a b c a$

**1. Schritt:** Wähle Kante  $(r, u)$  ausgehend von  $r$ , deren Beschriftung  $B$  mit  $P[1] = a$  beginnt.

**2. Schritt:** Vergleiche  $P[1, 4]$  und  $B$  schrittweise.

- Wenn Vergleich negativ, dann gebe 0 zurück.
- Wenn  $P[1, 4]$  Präfix von  $B$ , gebe Anzahl Blätter im Unterbaum  $S_u$  aus.
- Wenn  $B = P[1, j]$  für ein  $j$  mit  $1 \leq j \leq m$ , betrachte Knoten  $u$  und  $P[j + 1, m]$ .
  - Im Beispiel: Gilt für  $j = 1$ .

# Suchen im Suffixbaum

T = a b c a b a b c a

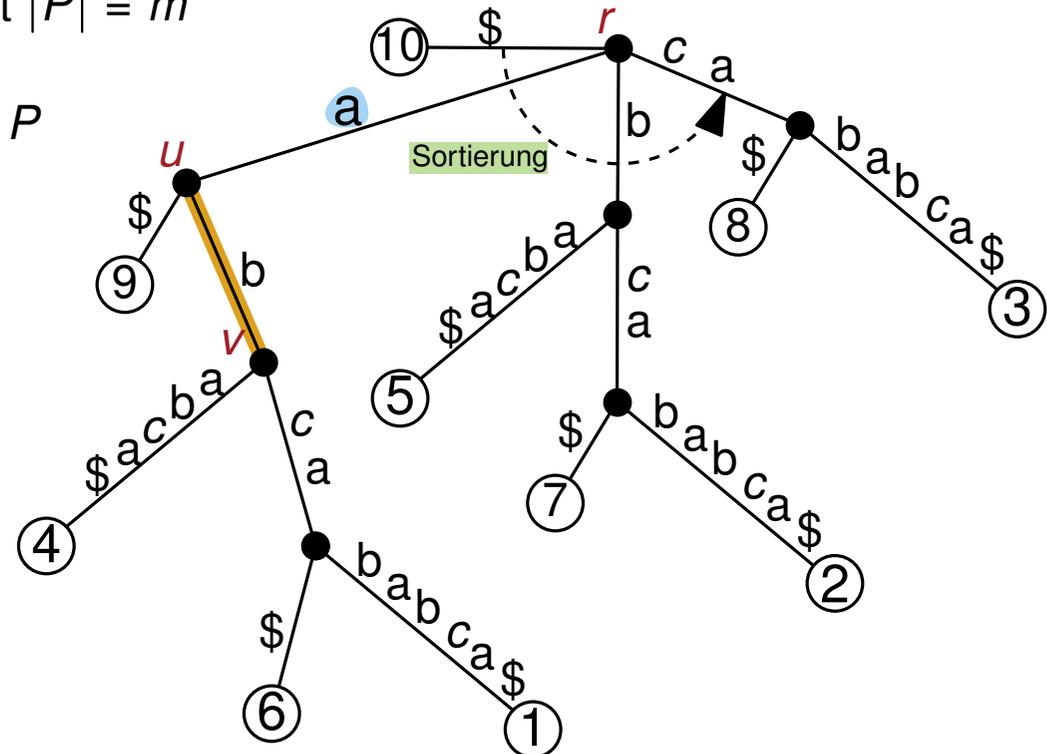
**Gegeben:** Suffixbaum  $S$  und Muster  $P$  mit  $|P| = m$

**Gesucht:** Anzahl Vorkommen von  $P$

**Idee:** Suche in  $S$  alle Pfade  $S_i$ , so dass  $P$  Präfix von  $S_i$  ist.

Ausgehend von Wurzel  $r$ :

1. Wähle nächste ausgehende Kante mithilfe von binärer Suche.
2. Vergleiche  $P$  und Kantenbeschriftungen schrittweise.



**Beispiel:**  $P = \mathbf{a} b c a$

**3. Schritt:** Wähle Kante  $(u, v)$  ausgehend von  $u$ , deren Beschriftung  $B$  mit  $P[2] = b$  beginnt.

# Suchen im Suffixbaum

T = a b c a b a b c a

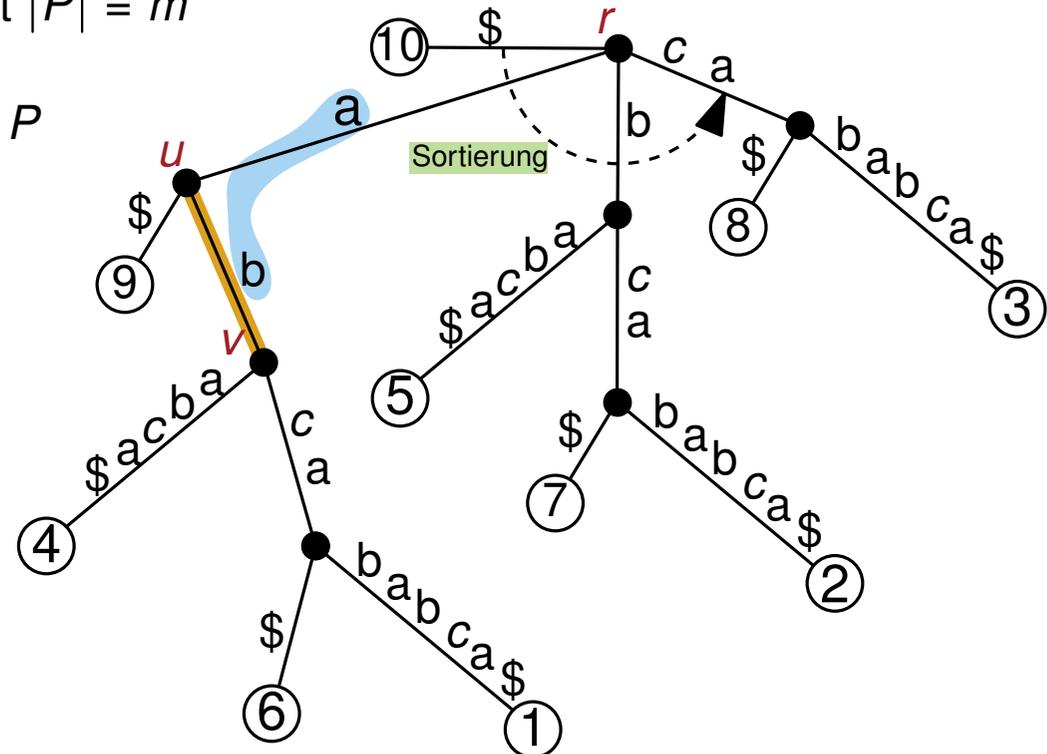
**Gegeben:** Suffixbaum  $S$  und Muster  $P$  mit  $|P| = m$

**Gesucht:** Anzahl Vorkommen von  $P$

**Idee:** Suche in  $S$  alle Pfade  $S_i$ , so dass  $P$  Präfix von  $S_i$  ist.

Ausgehend von Wurzel  $r$ :

1. Wähle nächste ausgehende Kante mithilfe von binärer Suche.
2. Vergleiche  $P$  und Kantenbeschriftungen schrittweise.



**Beispiel:**  $P = \mathbf{a b} c a$

**3. Schritt:** Wähle Kante  $(u, v)$  ausgehend von  $u$ , deren Beschriftung  $B$  mit  $P[2] = b$  beginnt.

**4. Schritt:** Vergleiche  $P[2, 4]$  und  $B$  schrittweise.

- Wenn Vergleich negativ, dann gebe 0 zurück.
- Wenn  $P[2, 4]$  Präfix von  $B$ , gebe Anzahl Blätter im Unterbaum  $S_v$  aus.
- Wenn  $B = P[2, j]$  für ein  $j$  mit  $1 \leq j \leq m$ , betrachte Knoten  $u$  und  $P[j + 1, m]$ .
  - Im Beispiel: Gilt für  $j = 2$ .

# Suchen im Suffixbaum

T = a b c a b a b c a

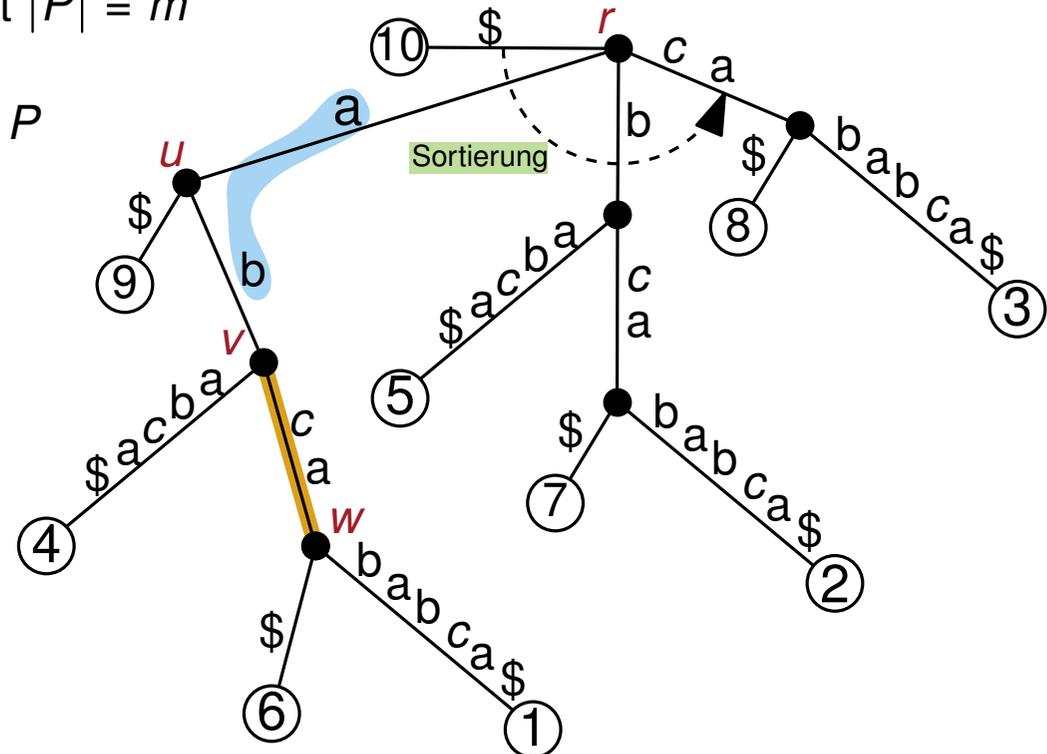
**Gegeben:** Suffixbaum  $S$  und Muster  $P$  mit  $|P| = m$

**Gesucht:** Anzahl Vorkommen von  $P$

**Idee:** Suche in  $S$  alle Pfade  $S_i$ , so dass  $P$  Präfix von  $S_i$  ist.

Ausgehend von Wurzel  $r$ :

1. Wähle nächste ausgehende Kante mithilfe von binärer Suche.
2. Vergleiche  $P$  und Kantenbeschriftungen schrittweise.



**Beispiel:**  $P = \text{a b c a}$

**5. Schritt:** Wähle Kante  $(v, w)$  ausgehend von  $v$ , deren Beschriftung  $B$  mit  $P[3] = c$  beginnt.

# Suchen im Suffixbaum

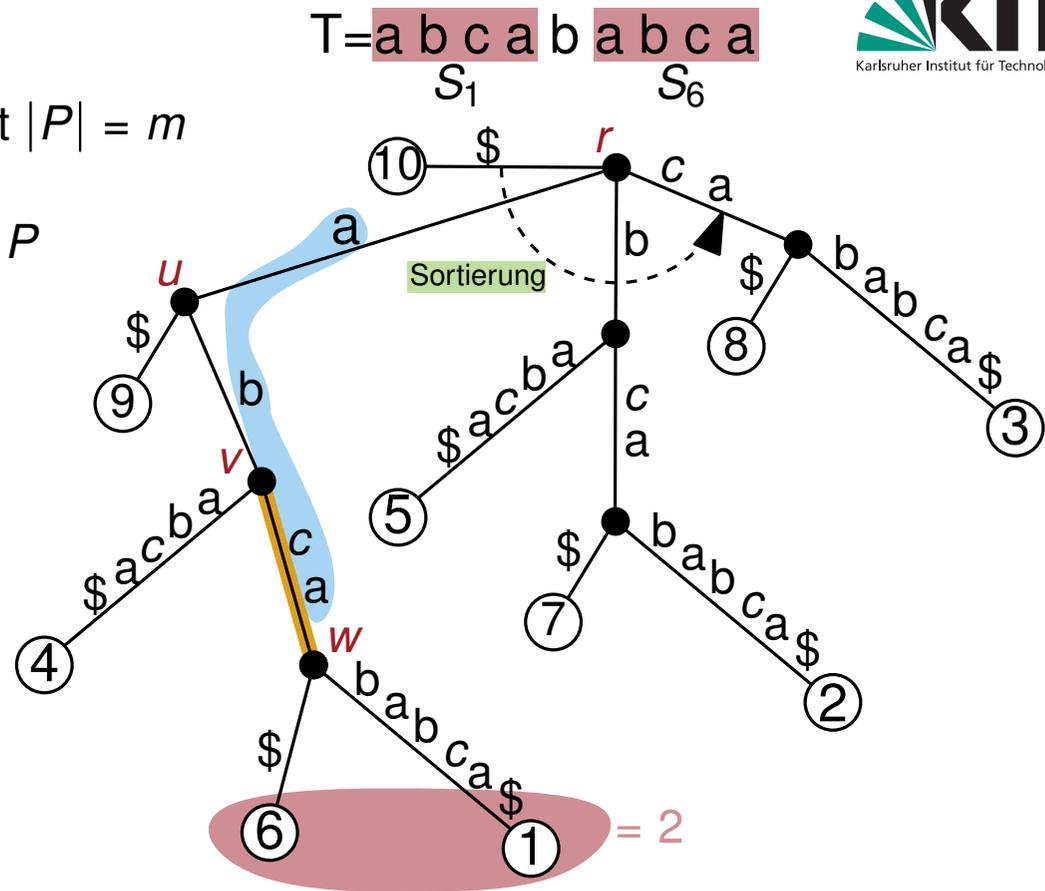
**Gegeben:** Suffixbaum  $S$  und Muster  $P$  mit  $|P| = m$

**Gesucht:** Anzahl Vorkommen von  $P$

**Idee:** Suche in  $S$  alle Pfade  $S_i$ , so dass  $P$  Präfix von  $S_i$  ist.

Ausgehend von Wurzel  $r$ :

1. Wähle nächste ausgehende Kante mithilfe von binärer Suche.
2. Vergleiche  $P$  und Kantenbeschriftungen schrittweise.



**Beispiel:**  $P =$ a b c a

**5. Schritt:** Wähle Kante  $(v, w)$  ausgehend von  $v$ , deren Beschriftung  $B$  mit  $P[3] = c$  beginnt.

**6. Schritt:** Vergleiche  $P[3, 4]$  und  $B$  schrittweise.

- Wenn Vergleich negativ, dann gebe 0 zurück.
- Wenn  $P[3, 4]$  Präfix von  $B$ , gebe Anzahl Blätter im Unterbaum  $S_w$  aus.
- Wenn  $B = P[3, j]$  für ein  $j$  mit  $1 \leq j \leq m$ , betrachte Knoten  $u$  und  $P[j + 1, m]$ .

# Suchen im Suffixbaum

COUNT( $P, S$ )

**Ausgabe:** Anzahl der Vorkommen von  $P$  im Text repräsentiert von Suffixbaum  $S$ .

$u \leftarrow$  Wurzel von  $S$

$i \leftarrow 1$

**solange**  $u$  nicht Blatt ist **tue**

Suche ausgehende Kante  $e = (u, v)$ , deren Beschriftung  $B$  mit  $P[i]$  beginnt.

**wenn**  $e$  nicht existiert **dann return** 0

**Vergleiche**  $B$  schrittweise mit  $P[i, m]$  **und analysiere Ergebnis**

**wenn**  $P[i, m]$  ist Präfix von  $B$  **dann**

| **return** Anzahl Blätter im Teilbaum  $S_v$

**sonst wenn**  $P[i, j] = B$  für ein  $j < m$  **dann**

|  $i \leftarrow j + 1$

|  $u \leftarrow v$

**sonst**

| **return** 0

**return** 0

# Suchen im Suffixbaum

COUNT( $P, S$ )

**Ausgabe:** Anzahl der Vorkommen von  $P$  im Text repräsentiert von Suffixbaum  $S$ .

$u \leftarrow$  Wurzel von  $S$

$i \leftarrow 1$

**solange**  $u$  nicht Blatt ist **tue**

Suche ausgehende Kante  $e = (u, v)$ , deren Beschriftung  $B$  mit  $P[i]$  beginnt.  $O(\log |\Sigma|)$

**wenn**  $e$  nicht existiert **dann return** 0

**Vergleiche**  $B$  schrittweise mit  $P[i, m]$  **und analysiere Ergebnis**

**wenn**  $P[i, m]$  ist Präfix von  $B$  **dann**

| **return** Anzahl Blätter im Teilbaum  $S_v$

**sonst wenn**  $P[i, j] = B$  für ein  $j < m$  **dann**

|  $i \leftarrow j + 1$

|  $u \leftarrow v$

**sonst**

| **return** 0

**return** 0

Beschriftung jeder ausgehenden Kante eines Knoten beginnt unterschiedlich.

⇒ Sortierung eindeutig möglich.

⇒ Binäre Suche auf ausgehenden Kanten eines Knoten möglich.

# Suchen im Suffixbaum

COUNT( $P, S$ )

**Ausgabe:** Anzahl der Vorkommen von  $P$  im Text repräsentiert von Suffixbaum  $S$ .

$u \leftarrow$  Wurzel von  $S$

$i \leftarrow 1$

**solange**  $u$  nicht Blatt ist **tue**

Suche ausgehende Kante  $e = (u, v)$ , deren Beschriftung  $B$  mit  $P[i]$  beginnt.  $O(\log |\Sigma|)$

**wenn**  $e$  nicht existiert **dann return** 0

**Vergleiche**  $B$  schrittweise mit  $P[i, m]$  **und analysiere Ergebnis**

**wenn**  $P[i, m]$  ist Präfix von  $B$  **dann**

**return** Anzahl Blätter im Teilbaum  $S_v$   $O(1)$

**sonst wenn**  $P[i, j] = B$  für ein  $j < m$  **dann**

$i \leftarrow j + 1$

$u \leftarrow v$

**sonst**

**return** 0

**return** 0

Anzahl Blätter kann in  $O(1)$  Zeit bestimmt werden, wenn an jedem Knoten Anzahl Blätter seines Teilbaums gespeichert ist.

# Suchen im Suffixbaum

COUNT( $P, S$ )

**Ausgabe:** Anzahl der Vorkommen von  $P$  im Text repräsentiert von Suffixbaum  $S$ .

$u \leftarrow$  Wurzel von  $S$

$i \leftarrow 1$

**solange**  $u$  nicht Blatt ist **tue**  $O(m \cdot \log |\Sigma|)$

Suche ausgehende Kante  $e = (u, v)$ , deren Beschriftung  $B$  mit  $P[i]$  beginnt.  $O(\log |\Sigma|)$

**wenn**  $e$  nicht existiert **dann return** 0

**Vergleiche**  $B$  schrittweise mit  $P[i, m]$  **und analysiere Ergebnis** max.  $m - i$  Schritte

**wenn**  $P[i, m]$  ist Präfix von  $B$  **dann**

**return** Anzahl Blätter im Teilbaum  $S_v$   $O(1)$

**sonst wenn**  $P[i, j] = B$  für ein  $j < m$  **dann**

$i \leftarrow j + 1$

$u \leftarrow v$

**sonst**

**return** 0

**return** 0

Beim Abstieg im Baum wird jedes Zeichen von  $P$  maximal einmal betrachtet.

**Theorem 23:** Gegeben ein Suffixbaum  $S$  und ein Muster  $P$  über dem Zeichenvorot  $\Sigma$ . Die Anzahl der Vorkommen von  $P$  kann in  $O(m \cdot \log |\Sigma|)$  Zeit bestimmt werden, wenn  $|P| = m$ .

## Verallgemeinerung:

Wenn die Vorkommen direkt gefragt sind, dann gebe entsprechend die Blätter aus.

**Theorem 23:** Gegeben ein Suffixbaum  $S$  und ein Muster  $P$  über dem Zeichenvorot  $\Sigma$ . Die Vorkommen von  $P$  können in  $O(m \cdot \log |\Sigma| + k)$  Zeit bestimmt werden, wenn  $|P| = m$  und  $k$  die Anzahl der Vorkommen von  $P$  in  $S$  bezeichnet.

Hinweis: In dem Fall ist Laufzeit *ausgabesensitiv* beschrieben.



**Implementierungsdetail:** Jede Kantenbeschriftung  $B$  ist durch ein Paar  $(i, j)$  mit  $1 \leq i \leq j \leq n$  repräsentiert, sodass  $B = T[i, j]$ .  
 $\Rightarrow$  Da  $S$  genau  $n$  Blätter hat, ergibt sich damit  $O(n)$  Speicherverbrauch.

# Konstruktion

**Gegeben:** Text  $T$

**Gesucht:** Suffixbaum von  $T$

**Idee:** Konstruiere  $N_1, \dots, N_n$  Suffixbäume, sodass  $N_i$  die Suffixe  $S_1, \dots, S_i$  enthält.

**Initialisierung:**  $N_1$  besteht aus zwei Knoten und einer Kante beschriftet mit  $S_1$ .

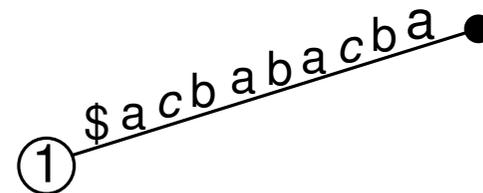
**Konstruktion von  $N_{i+1}$  aus  $N_i$ :**

Suche längsten Pfad in  $N_i$ , der mit Präfix von  $S_{i+1}$  übereinstimmt.

**1. Fall:** Wenn Übereinstimmung in der Mitte von  $e$  endet, dann spalte  $e$  danach in zwei Kanten auf, indem neuer Knoten  $w$  eingeführt wird.

**2. Fall:** Übereinstimmung endet an Knoten  $w$ : Führe neue ausgehende Kante  $e$  für  $w$  ein.

$T = a \text{ b c a b a b c a } \$$   
 $S_2$



**Nächster Schritt:**

Füge  $S_2 = b c a b a b c a \$$  ein:

- Matching hört bereits bei Wurzel auf.
- Folglich: zweiter Fall.



# Konstruktion

**Gegeben:** Text  $T$

**Gesucht:** Suffixbaum von  $T$

**Idee:** Konstruiere  $N_1, \dots, N_n$  Suffixbäume, sodass  $N_i$  die Suffixe  $S_1, \dots, S_i$  enthält.

**Initialisierung:**  $N_1$  besteht aus zwei Knoten und einer Kante beschriftet mit  $S_1$ .

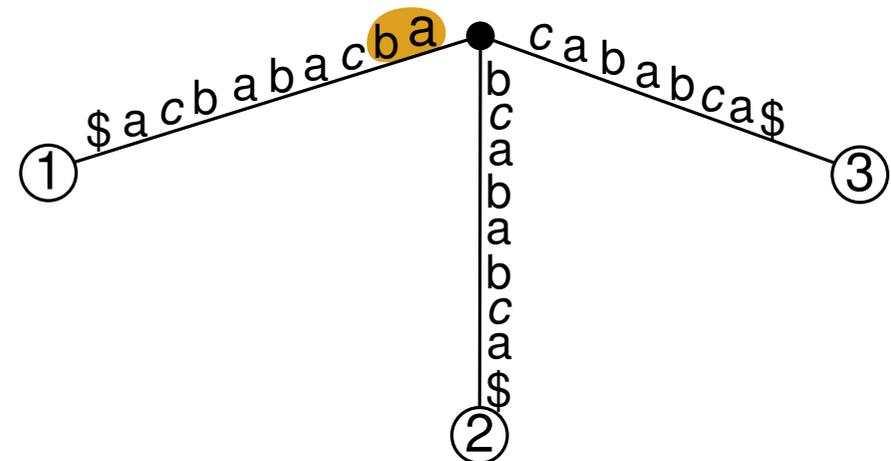
**Konstruktion von  $N_{i+1}$  aus  $N_i$ :**

Suche längsten Pfad in  $N_i$ , der mit Präfix von  $S_{i+1}$  übereinstimmt.

**1. Fall:** Wenn Übereinstimmung in der Mitte von  $e$  endet, dann spalte  $e$  danach in zwei Kanten auf, indem neuer Knoten  $w$  eingeführt wird.

**2. Fall:** Übereinstimmung endet an Knoten  $w$ : Führe neue ausgehende Kante  $e$  für  $w$  ein.

$T = a b c a b a b c a \$$   
 $S_4 = a b a b c a \$$



**Nächster Schritt:**

Füge  $S_4 = a b a b c a \$$  ein:

- Matching hört auf Pfad  $S_1$  nach 2 Zeichen auf.
- Folglich: erster Fall.

# Konstruktion

**Gegeben:** Text  $T$

**Gesucht:** Suffixbaum von  $T$

**Idee:** Konstruiere  $N_1, \dots, N_n$  Suffixbäume, sodass  $N_i$  die Suffixe  $S_1, \dots, S_i$  enthält.

**Initialisierung:**  $N_1$  besteht aus zwei Knoten und einer Kante beschriftet mit  $S_1$ .

**Konstruktion von  $N_{i+1}$  aus  $N_i$ :**

Suche längsten Pfad in  $N_i$ , der mit Präfix von  $S_{i+1}$  übereinstimmt.

**1. Fall:** Wenn Übereinstimmung in der Mitte von  $e$  endet, dann spalte  $e$  danach in zwei Kanten auf, indem neuer Knoten  $w$  eingeführt wird.

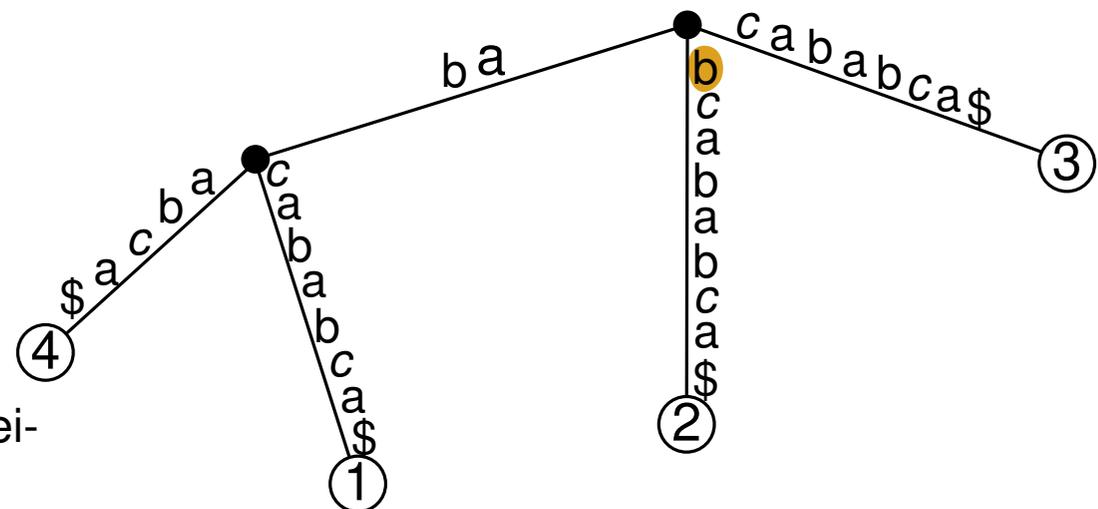
**2. Fall:** Übereinstimmung endet an Knoten  $w$ : Führe neue ausgehende Kante  $e$  für  $w$  ein.

$T = a b c a \mathbf{b a b c a} \$$   
 $S_5$

**Nächster Schritt:**

Füge  $S_5 = b a b c a \$$  ein:

- Matching hört auf Pfad  $S_2$  nach einem Zeichen auf.
- Folglich: erster Fall.



# Konstruktion

**Gegeben:** Text  $T$

**Gesucht:** Suffixbaum von  $T$

**Idee:** Konstruiere  $N_1, \dots, N_n$  Suffixbäume, sodass  $N_i$  die Suffixe  $S_1, \dots, S_i$  enthält.

**Initialisierung:**  $N_1$  besteht aus zwei Knoten und einer Kante beschriftet mit  $S_1$ .

**Konstruktion von  $N_{i+1}$  aus  $N_i$ :**

Suche längsten Pfad in  $N_i$ , der mit Präfix von  $S_{i+1}$  übereinstimmt.

**1. Fall:** Wenn Übereinstimmung in der Mitte von  $e$  endet, dann spalte  $e$  danach in zwei Kanten auf, indem neuer Knoten  $w$  eingeführt wird.

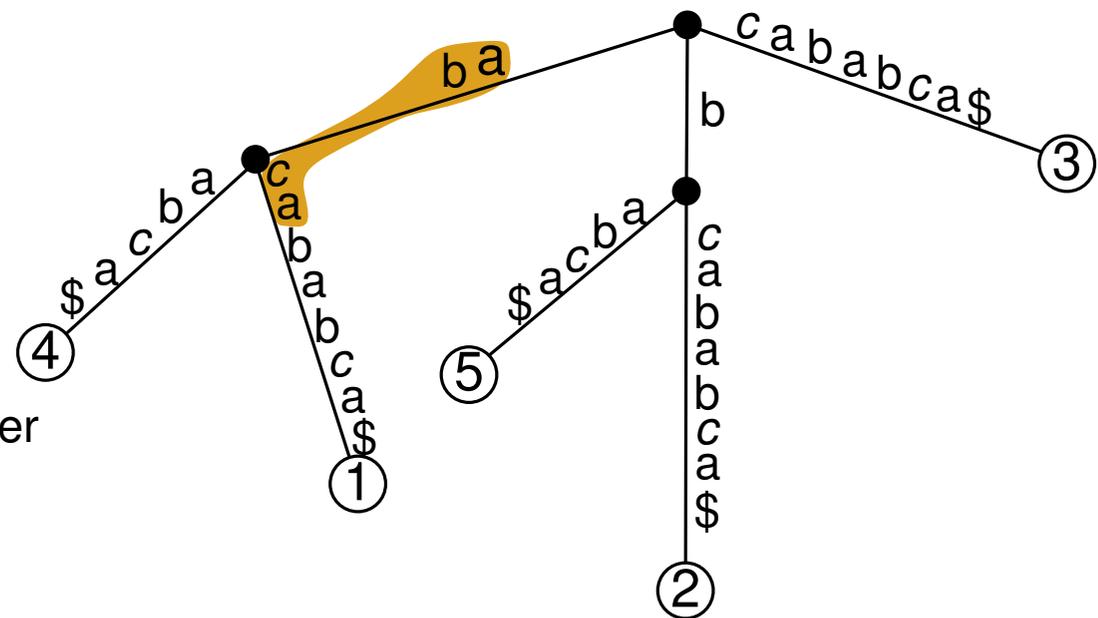
**2. Fall:** Übereinstimmung endet an Knoten  $w$ : Führe neue ausgehende Kante  $e$  für  $w$  ein.

$T = a b c a b a b c a \$$   
 $S_6 = a b c a \$$

**Nächster Schritt:**

Füge  $S_6 = a b c a \$$  ein:

- Matching hört auf Pfad  $S_1$  nach vier Zeichen auf.
- Folglich: erster Fall.



# Konstruktion

**Gegeben:** Text  $T$

**Gesucht:** Suffixbaum von  $T$

**Idee:** Konstruiere  $N_1, \dots, N_n$  Suffixbäume, sodass  $N_i$  die Suffixe  $S_1, \dots, S_i$  enthält.

**Initialisierung:**  $N_1$  besteht aus zwei Knoten und einer Kante beschriftet mit  $S_1$ .

**Konstruktion von  $N_{i+1}$  aus  $N_i$ :**

Suche längsten Pfad in  $N_i$ , der mit Präfix von  $S_{i+1}$  übereinstimmt.

**1. Fall:** Wenn Übereinstimmung in der Mitte von  $e$  endet, dann spalte  $e$  danach in zwei Kanten auf, indem neuer Knoten  $w$  eingeführt wird.

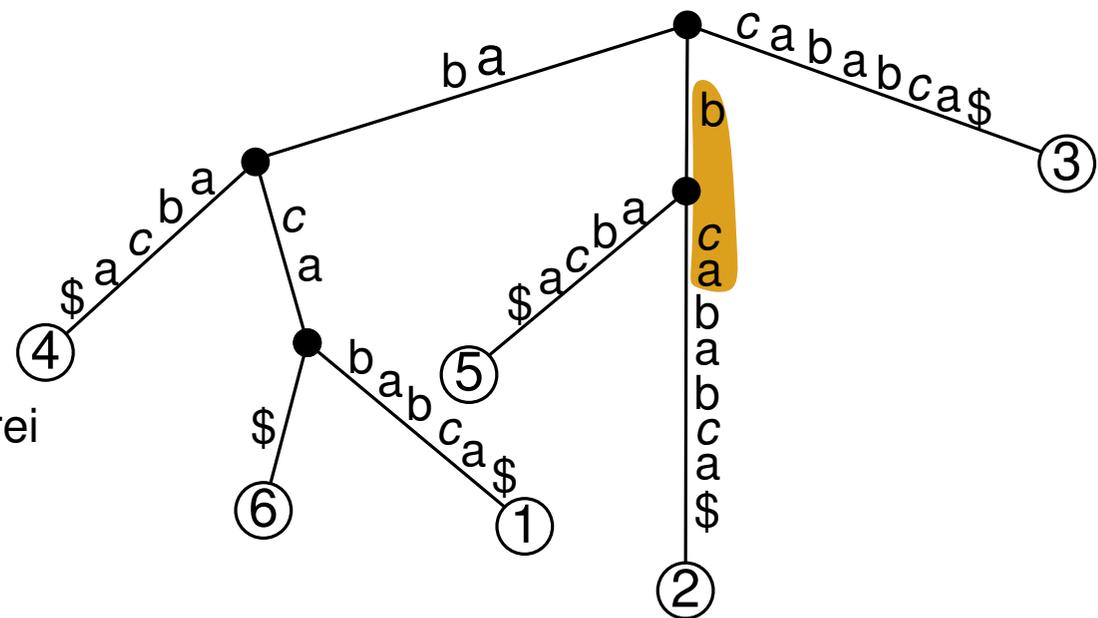
**2. Fall:** Übereinstimmung endet an Knoten  $w$ : Führe neue ausgehende Kante  $e$  für  $w$  ein.

$T = a b c a b a b c a \$$   
 $S_7 = b c a \$$

**Nächster Schritt:**

Füge  $S_7 = b c a \$$  ein:

- Matching hört auf Pfad  $S_2$  nach drei Zeichen auf.
- Folglich: erster Fall.



# Konstruktion

**Gegeben:** Text  $T$

**Gesucht:** Suffixbaum von  $T$

**Idee:** Konstruiere  $N_1, \dots, N_n$  Suffixbäume, sodass  $N_i$  die Suffixe  $S_1, \dots, S_i$  enthält.

**Initialisierung:**  $N_1$  besteht aus zwei Knoten und einer Kante beschriftet mit  $S_1$ .

**Konstruktion von  $N_{i+1}$  aus  $N_i$ :**

Suche längsten Pfad in  $N_i$ , der mit Präfix von  $S_{i+1}$  übereinstimmt.

**1. Fall:** Wenn Übereinstimmung in der Mitte von  $e$  endet, dann spalte  $e$  danach in zwei Kanten auf, indem neuer Knoten  $w$  eingeführt wird.

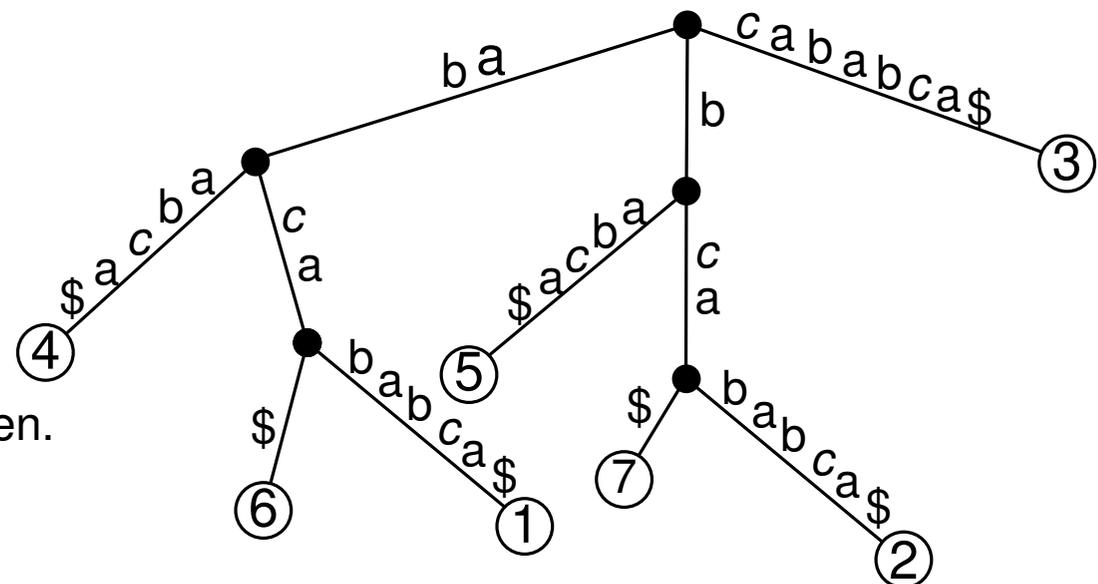
**2. Fall:** Übereinstimmung endet an Knoten  $w$ : Führe neue ausgehende Kante  $e$  für  $w$  ein.

$T = a b c a b a b c a \$$

$S_8, S_9$  und  $S_{10}$  analog.

Kann in  $O(n^2 \log |\Sigma|)$  implementiert werden.

Geht es auch besser?



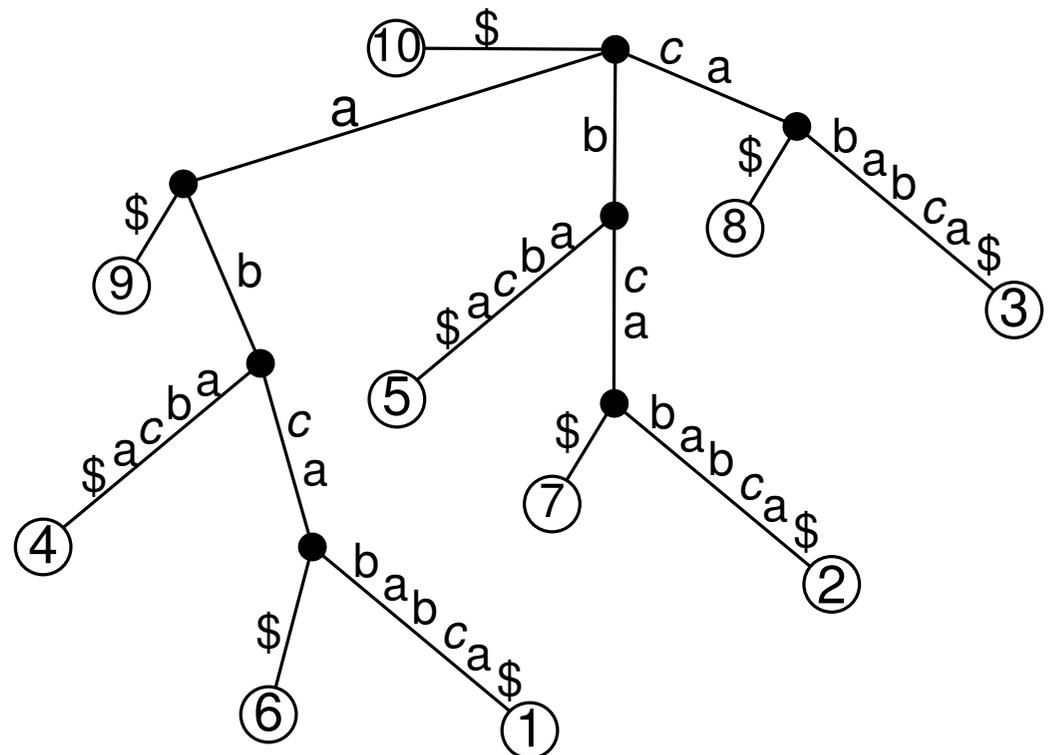
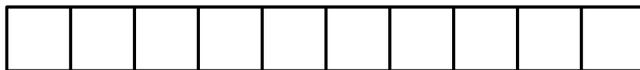
# Suffix-Arrays

**Definition:** Ein *Suffix-Array*  $A$  eines Textes  $T$  ist eine Permutation von  $\{1, \dots, n\}$ , sodass  $S_{A[i]}$  das  $i$ -kleinste Suffix in lexikographischer Ordnung ist:  $S_{A[i-1]} < S_{A[i]}$  für alle  $1 < i \leq n$ .



1. Suffix-Array gibt lexikographische Sortierung der Suffixe an.
2. Sei  $S$  Suffixbaum von  $T$ .  $A$  entspricht der Reihenfolge der Blätter, die sich ergibt, wenn Tiefensuche auf  $S$  mit lexikographischer Ordnung angewendet wird.

**Beispiel:**  $T = a b c a b a b c a \$$



**Definition:** Ein *Suffix-Array*  $A$  eines Textes  $T$  ist eine Permutation von  $\{1, \dots, n\}$ , sodass  $S_{A[i]}$  das  $i$ -kleinste Suffix in lexikographischer Ordnung ist:  $S_{A[i-1]} < S_{A[i]}$  für alle  $1 < i \leq n$ .



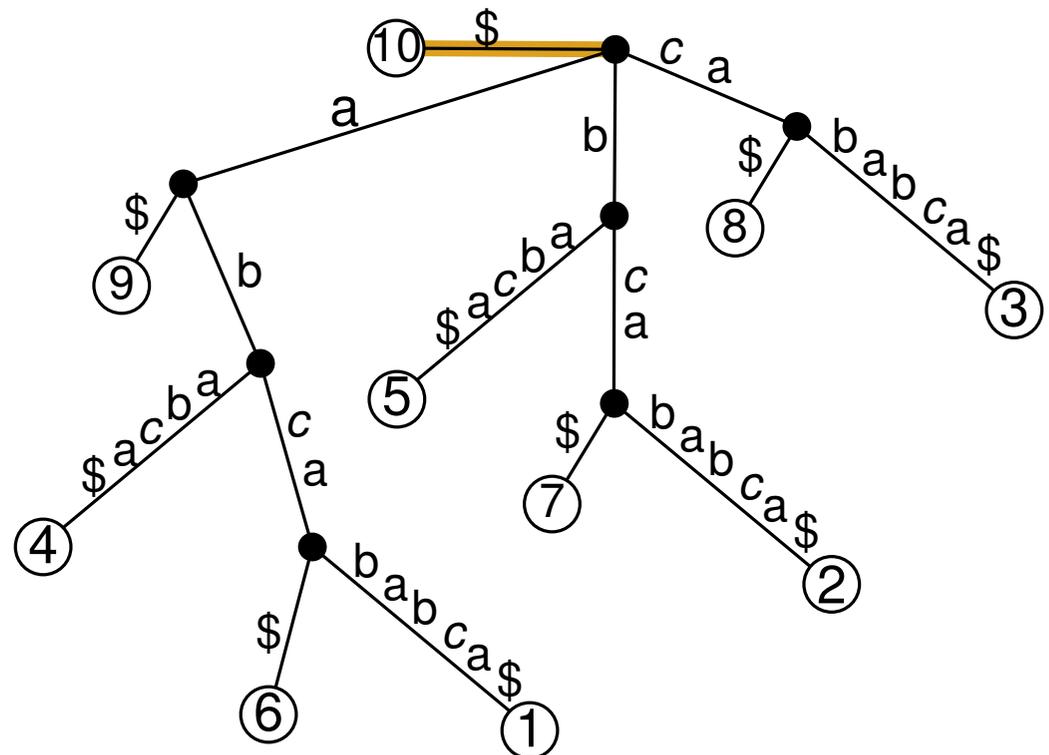
1. Suffix-Array gibt lexikographische Sortierung der Suffixe an.
2. Sei  $S$  Suffixbaum von  $T$ .  $A$  entspricht der Reihenfolge der Blätter, die sich ergibt, wenn Tiefensuche auf  $S$  mit lexikographischer Ordnung angewendet wird.

**Beispiel:**  $T = a b c a b a b c a \$$

$A =$ 

10									
----	--	--	--	--	--	--	--	--	--

  
\$





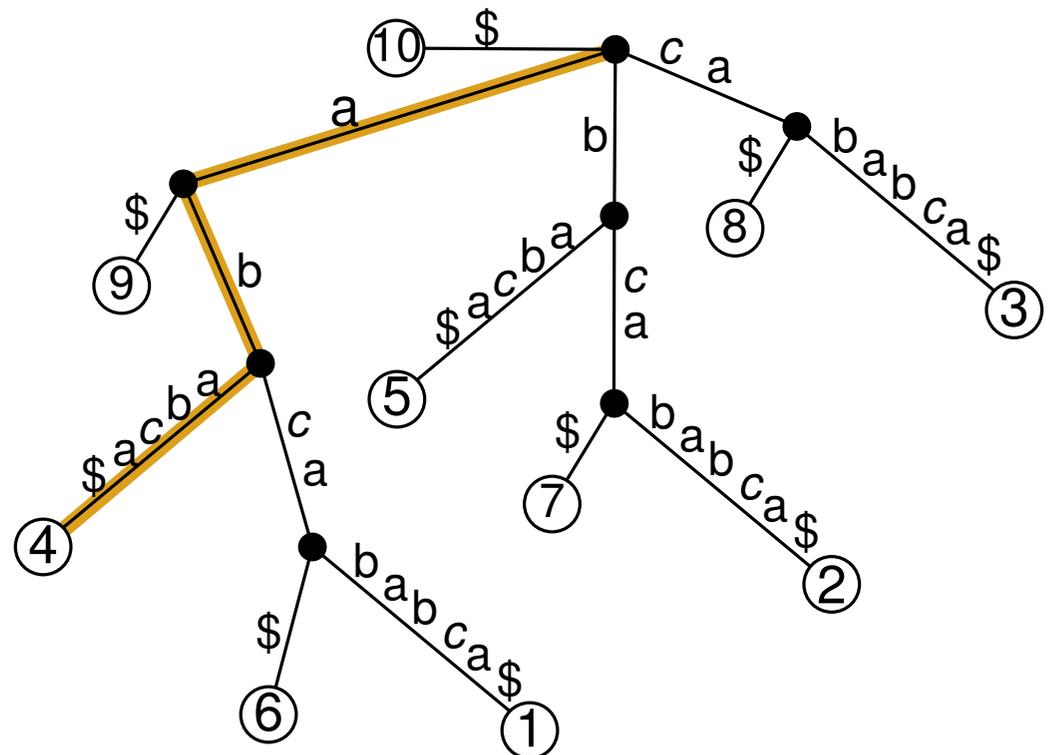
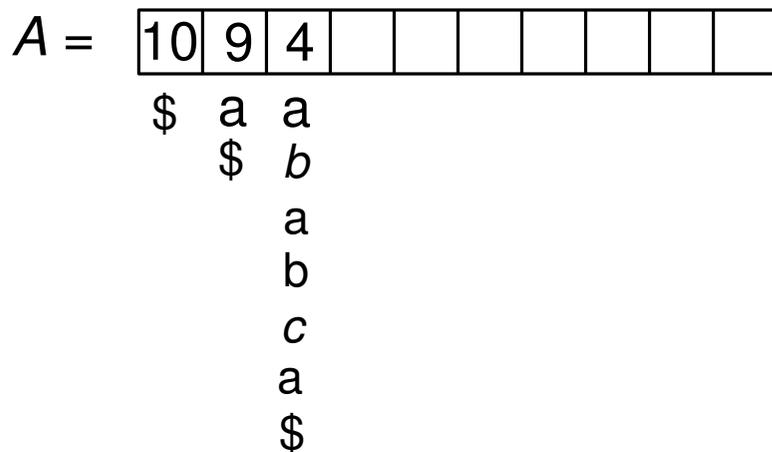
# Suffix-Arrays

**Definition:** Ein *Suffix-Array*  $A$  eines Textes  $T$  ist eine Permutation von  $\{1, \dots, n\}$ , sodass  $S_{A[i]}$  das  $i$ -kleinste Suffix in lexikographischer Ordnung ist:  $S_{A[i-1]} < S_{A[i]}$  für alle  $1 < i \leq n$ .



1. Suffix-Array gibt lexikographische Sortierung der Suffixe an.
2. Sei  $S$  Suffixbaum von  $T$ .  $A$  entspricht der Reihenfolge der Blätter, die sich ergibt, wenn Tiefensuche auf  $S$  mit lexikographischer Ordnung angewendet wird.

**Beispiel:**  $T = a b c a b a b c a \$$







# Suffix-Arrays

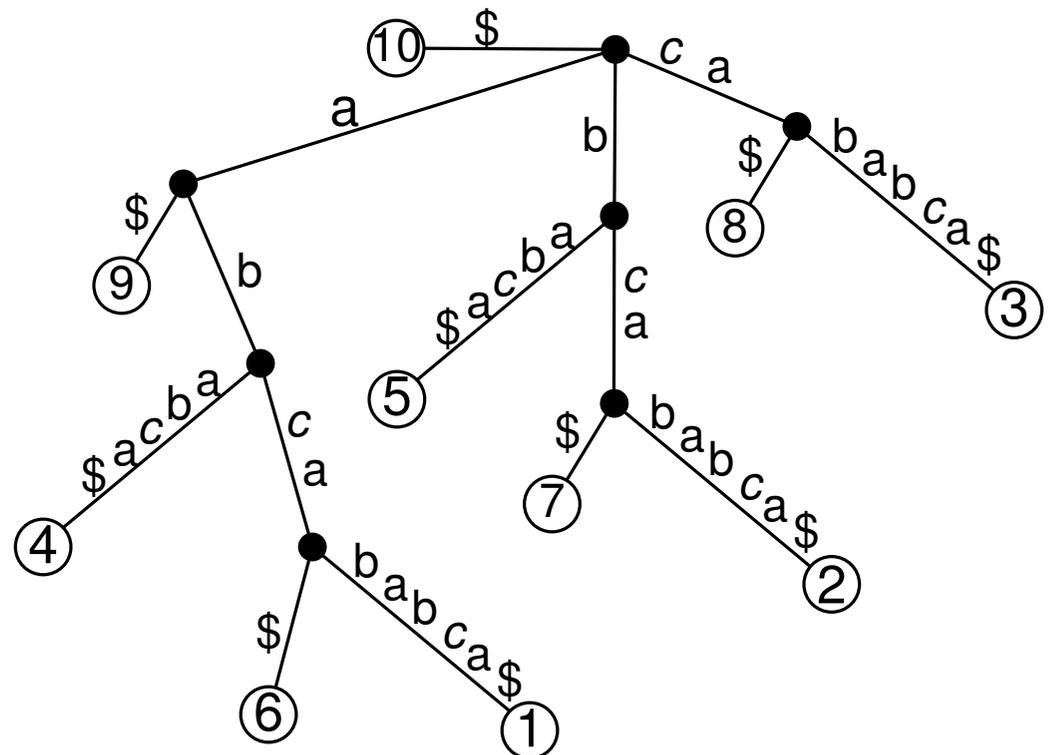
**Definition:** Ein *Suffix-Array*  $A$  eines Textes  $T$  ist eine Permutation von  $\{1, \dots, n\}$ , sodass  $S_{A[i]}$  das  $i$ -kleinste Suffix in lexikographischer Ordnung ist:  $S_{A[i-1]} < S_{A[i]}$  für alle  $1 < i \leq n$ .



1. Suffix-Array gibt lexikographische Sortierung der Suffixe an.
2. Sei  $S$  Suffixbaum von  $T$ .  $A$  entspricht der Reihenfolge der Blätter, die sich ergibt, wenn Tiefensuche auf  $S$  mit lexikographischer Ordnung angewendet wird.

**Beispiel:**  $T = a b c a b a b c a \$$

$A =$	10	9	4	6	1	5	7	2	8	3
	\$	a	a	a	a	b	b	b	c	c
		\$	b	b	b	a	c	c	a	a
			a	c	c	b	a	a	\$	b
			b	a	a	c	\$	b	a	b
			c	\$	b	a	b	c	a	\$
			a		a	\$				
			\$		b	c	a			
					\$					



# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b			c		a
				c			a		\$
				a			\$		

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				c	a		\$		\$
				a					

Benutze binäre Suche um linke Intervallgrenze zu finden:

```

FINDELINKEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme unteren Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch größer als  $S_{A[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[i + 1, size(A')]$  (rechtes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[1, i]$  (linkes Teilarray)
 $\ell \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $\ell$ 
    
```

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				c	a		\$		\$

Benutze binäre Suche um linke Intervalgrenze zu finden:

```

FINDELINKEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme unteren Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch größer als  $S_{A[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[i + 1, size(A')]$  (rechtes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[1, i]$  (linkes Teilarray)
 $\ell \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $\ell$ 
    
```

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				c	a		\$		\$

Benutze binäre Suche um linke Intervalgrenze zu finden:

```

FINDELINKEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme unteren Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch größer als  $S_{A[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[i + 1, size(A')]$  (rechtes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[1, i]$  (linkes Teilarray)
 $\ell \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $\ell$ 
    
```

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		a		a
				c	a		\$		\$

Benutze binäre Suche um linke Intervalgrenze zu finden:

```

FINDELINKEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme unteren Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch größer als  $S_{A[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[i + 1, size(A')]$  (rechtes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[1, i]$  (linkes Teilarray)
 $\ell \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $\ell$ 
    
```

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				a			\$		\$

Benutze binäre Suche um linke Intervalgrenze zu finden:

```

FINDELINKEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    Bestimme unteren Median  $i$  von  $A'$ .
    wenn  $P$  lexikographisch größer als  $S_{A[i]}[1, m]$  dann
        |  $A' \leftarrow A'[i + 1, size(A')]$  (rechtes Teilarray)
    sonst
        |  $A' \leftarrow A'[1, i]$  (linkes Teilarray)
 $\ell \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $\ell$ 
    
```

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				c	a		\$		\$

Benutze binäre Suche um linke Intervalgrenze zu finden:

```

FINDELINKEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme unteren Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch größer als  $S_{A[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[i + 1, size(A')]$  (rechtes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[1, i]$  (linkes Teilarray)
 $\ell \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $\ell$ 
    
```

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				c	a		\$		\$

Benutze binäre Suche um linke Intervallgrenze zu finden:

```

FINDELINKEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme unteren Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch größer als  $S_{A[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[i + 1, size(A')]$  (rechtes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[1, i]$  (linkes Teilarray)
 $\ell \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $\ell$ 
    
```

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				a			\$		\$

Benutze binäre Suche um linke Intervalgrenze zu finden:

```

FINDELINKEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme unteren Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch größer als  $S_{A[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[i + 1, size(A')]$  (rechtes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[1, i]$  (linkes Teilarray)
 $\ell \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $\ell$ 
    
```

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

Analog binäre Suche für rechte Intervalgrenze:

```

FINDERECHEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme oberen Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch kleiner als  $S_{A'[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[1, i - 1]$  (linkes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[i, size(A')]$  (rechtes Teilarray)
 $r \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $r$ 
    
```

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				c	a		\$		\$

Lexikographischer Vergleich kann in  $O(m)$  Zeit durchgeführt werden  $\Rightarrow O(m \cdot \log n)$  Zeit.  
**Vergleiche:** Suche im Suffixbaum in  $O(m \log |\Sigma|)$  möglich. Normalerweise gilt  $|\Sigma| \leq m \leq n$ .

# Suchen in Suffix-Arrays

**Gegeben:** Suffix-Array  $A$  für Text  $T$

**Gesucht:** Anzahl Vorkommen von Muster  $P$  in  $T$

**Beobachtung:** Gesuchte Stellen bilden Intervall im Suffix-Array.

Analog binäre Suche für rechte Intervalgrenze:

```

FINDERECHTEGRENZE(Suffix-Array  $A$ )
 $A' \leftarrow A$ 
solange  $size(A') > 1$  tue
    | Bestimme oberen Median  $i$  von  $A'$ .
    | wenn  $P$  lexikographisch kleiner als  $S_{A'[i]}[1, m]$  dann
    | |  $A' \leftarrow A'[1, i - 1]$  (linkes Teilarray)
    | sonst
    | |  $A' \leftarrow A'[i, size(A')]$  (rechtes Teilarray)
 $r \leftarrow$  Index von  $A'[1]$  in  $A$ .
return  $r$ 
    
```

**Beispiel:**  $T = a b c a b a b c a \$$ ,  $P = a b$

$A =$

10	9	4	6	1	5	7	2	8	3
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		c		a
				c	a		\$		\$

Lexikographischer Vergleich kann in  $O(m)$  Zeit durchgeführt werden  $\Rightarrow O(m \cdot \log n)$  Zeit.  
**Vergleiche:** Suche im Suffixbaum in  $O(m \log |\Sigma|)$  möglich. Normalerweise gilt  $|\Sigma| \leq m \leq n$ .

# Konstruktion von Suffix-Arrays

**Eingabe:** Text  $T := t_0 t_1 \dots t_{n-1}$

	0	1	2	3	4	5	6	7	8	9	10	11
T=	y	a	b	b	a	d	a	b	b	a	d	o

**Gesucht:** Suffix-Array  $A$  von  $T$ , d.h. lexikographische Sortierung von Suffixen von  $T$ .

# Konstruktion von Suffix-Arrays

**Eingabe:** Text  $T := t_0 t_1 \dots t_{n-1}$

	0	1	2	3	4	5	6	7	8	9	10	11
T=	y	a	b	b	a	d	a	b	b	a	d	o

**Gesucht:** Suffix-Array  $A$  von  $T$ , d.h. lexikographische Sortierung von Suffixen von  $T$ .

Kürze  $x \bmod y = z$  mit  $x \equiv z(y)$  ab.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

$\mathcal{S} =$  Suffixe von  $T$

$\mathcal{S}_0 =$  Suffixe mit Index  $i \equiv 0(3)$

$\mathcal{S}_1 =$  Suffixe mit Index  $i \equiv 1(3)$

$\mathcal{S}_2 =$  Suffixe mit Index  $i \equiv 2(3)$

# Konstruktion von Suffix-Arrays

**Eingabe:** Text  $T := t_0 t_1 \dots t_{n-1}$

	0	1	2	3	4	5	6	7	8	9	10	11
T=	y	a	b	b	a	d	a	b	b	a	d	o

**Gesucht:** Suffix-Array  $A$  von  $T$ , d.h. lexikographische Sortierung von Suffixen von  $T$ .

Kürze  $x \bmod y = z$  mit  $x \equiv z(y)$  ab.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

$\mathcal{S}$  = Suffixe von  $T$

$\mathcal{S}_0$  = Suffixe mit Index  $i \equiv 0(3)$

$\mathcal{S}_1$  = Suffixe mit Index  $i \equiv 1(3)$

$\mathcal{S}_2$  = Suffixe mit Index  $i \equiv 2(3)$

SUFFIXARRAY(Text  $T = t_0 t_1 \dots t_{n-1}$ )

**wenn**  $n = O(1)$  **dann**

| Konstruiere  $A$  in  $O(1)$  Zeit.

**sonst**

Berechne Suffix-Array  $A_{12}$  für  $\mathcal{S}_1 \cup \mathcal{S}_2$ .

Berechne Suffix-Array  $A_0$  für  $\mathcal{S}_0$  basierend auf  $A_{12}$ .

Vermenge  $A_{12}$  mit  $A_0$ .

# 1. Schritt: Erstelle Suffix-Array $A_{12}$

**Eingabe:** Text  $T := t_0 t_1 \dots t_{n-1}$

	0	1	2	3	4	5	6	7	8	9	10	11
T=	y	a	b	b	a	d	a	b	b	a	d	o

Fasse Tripel  $[t_i t_{i+1} t_{i+2}]$  mit  $i \not\equiv 0(3)$  als eigenes Zeichen auf. Fülle gegebenenfalls mit \$ auf.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

Konstruiere für  $k \in \{1, 2\}$  den String

$$R_k = [t_k t_{k+1} t_{k+2}][t_{k+3} t_{k+4} t_{k+5}][t_{k+6} t_{k+7} t_{k+8}] \dots$$

$$R_1 = [abb][ada][bba][do\$]$$

$$R_2 = [bba][dab][bad][o\$\$]$$

$S_0$  = Suffixe mit Index  $i \equiv 0(3)$

$S_1$  = Suffixe mit Index  $i \equiv 1(3)$

$S_2$  = Suffixe mit Index  $i \equiv 2(3)$

# 1. Schritt: Erstelle Suffix-Array $A_{12}$

**Eingabe:** Text  $T := t_0 t_1 \dots t_{n-1}$

	0	1	2	3	4	5	6	7	8	9	10	11
T=	y	a	b	b	a	d	a	b	b	a	d	o

Fasse Tripel  $[t_i t_{i+1} t_{i+2}]$  mit  $i \not\equiv 0(3)$  als eigenes Zeichen auf. Fülle gegebenenfalls mit \$ auf.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

Konstruiere für  $k \in \{1, 2\}$  den String

$$R_k = [t_k t_{k+1} t_{k+2}][t_{k+3} t_{k+4} t_{k+5}][t_{k+6} t_{k+7} t_{k+8}] \dots$$

$$R_1 = [abb][ada][bba][do\$]$$

$$R_2 = [bba][dab][bad][o\$\$]$$

Konkatenation von  $R_1$  und  $R_2$ .

$$R = R_1 \cdot R_2$$

$$R = [abb][ada][bba][do\$][bba][dab][bad][o\$\$]$$

$S_0$  = Suffixe mit Index  $i \equiv 0(3)$

$S_1$  = Suffixe mit Index  $i \equiv 1(3)$

$S_2$  = Suffixe mit Index  $i \equiv 2(3)$

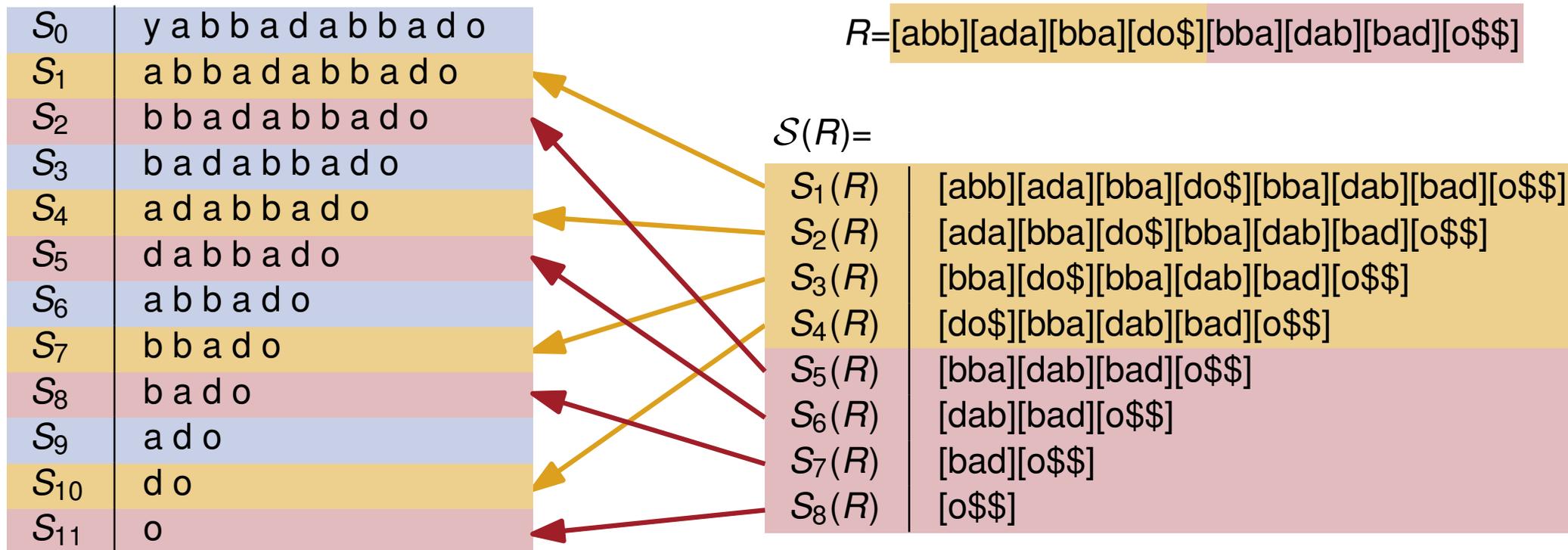
# 1. Schritt: Erstelle Suffix-Array $A_{12}$

**Beobachtung** Jedes Suffix von  $R$  hat eindeutige Entsprechung in  $S_1 \cup S_2$ .

$S_i$  entspricht Suffix  $S'_i = [t_i t_{i+1} t_{i+2}][t_{i+3} t_{i+4} t_{i+5} \dots]$  von  $R'$

Sortierung von Suffixen  $\mathcal{S}(R)$  induziert gesuchte Sortierung von  $S_1 \cup S_2$ .

Wie  $\mathcal{S}(R)$  effizient sortieren?



**Erinnerung:**  $t_i$  erstes Zeichen von  $S_i$ .

# 1. Schritt: Erstelle Suffix-Array $A_{12}$

Wie  $S(R)$  effizient sortieren?

$R = [abb][ada][bba][do\$][bba][dab][bad][o\$\$]$

$S(R) =$

$S_1(R)$	[abb][ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_2(R)$	[ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_3(R)$	[bba][do\$][bba][dab][bad][o\$]\$]
$S_4(R)$	[do\$][bba][dab][bad][o\$]\$]
$S_5(R)$	[bba][dab][bad][o\$]\$]
$S_6(R)$	[dab][bad][o\$]\$]
$S_7(R)$	[bad][o\$]\$]
$S_8(R)$	[o\$]\$]

# 1. Schritt: Erstelle Suffix-Array $A_{12}$

Wie  $S(R)$  effizient sortieren?

$R = [abb][ada][bba][do\$][bba][dab][bad][o\$\$]$

Wende Radixsort auf Zeichen von  $R$  an.  
(siehe Übung)

$S(R) =$

$S_1(R)$	[abb][ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_2(R)$	[ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_3(R)$	[bba][do\$][bba][dab][bad][o\$]\$]
$S_4(R)$	[do\$][bba][dab][bad][o\$]\$]
$S_5(R)$	[bba][dab][bad][o\$]\$]
$S_6(R)$	[dab][bad][o\$]\$]
$S_7(R)$	[bad][o\$]\$]
$S_8(R)$	[o\$]\$]

Rang	Zeichen
1	[abb]
2	[ada]
3	[bad]
4	[bba]
5	[dab]
6	[do\$]
7	[o\$]\$]

# 1. Schritt: Erstelle Suffix-Array $A_{12}$

Wie  $S(R)$  effizient sortieren?

$R = [abb][ada][bba][do\$][bba][dab][bad][o\$\$]$

Wende Radixsort auf Zeichen von  $R$  an.  
(siehe Übung)

**Transformation:** Ersetze Zeichen in  $R$  durch ihren Rang  $\rightarrow R'$

$R' = 1\ 2\ 4\ 6\ 4\ 5\ 3\ 7$

$S(R) =$

$S_1(R)$	[abb][ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_2(R)$	[ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_3(R)$	[bba][do\$][bba][dab][bad][o\$]\$]
$S_4(R)$	[do\$][bba][dab][bad][o\$]\$]
$S_5(R)$	[bba][dab][bad][o\$]\$]
$S_6(R)$	[dab][bad][o\$]\$]
$S_7(R)$	[bad][o\$]\$]
$S_8(R)$	[o\$]\$]

Rang	Zeichen
1	[abb]
2	[ada]
3	[bad]
4	[bba]
5	[dab]
6	[do\$]
7	[o\$]\$]

# 1. Schritt: Erstelle Suffix-Array $A_{12}$

Wie  $S(R)$  effizient sortieren?

$R = [abb][ada][bba][do\$][bba][dab][bad][o\$\$]$

Wende Radixsort auf Zeichen von  $R$  an.  
(siehe Übung)

**Transformation:** Ersetze Zeichen in  $R$  durch ihren Rang  $\rightarrow R'$

$R' = 1\ 2\ 4\ 6\ 4\ 5\ 3\ 7$

Suffix  $S_i(R)$  entspricht  $S_i(R')$

→ Sortierung von  $S(R')$  liefert gewünschte Sortierung von  $S(R)$ .

$S(R) =$

$S_1(R)$	[abb][ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_2(R)$	[ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_3(R)$	[bba][do\$][bba][dab][bad][o\$]\$]
$S_4(R)$	[do\$][bba][dab][bad][o\$]\$]
$S_5(R)$	[bba][dab][bad][o\$]\$]
$S_6(R)$	[dab][bad][o\$]\$]
$S_7(R)$	[bad][o\$]\$]
$S_8(R)$	[o\$]\$]

Rang	Zeichen
1	[abb]
2	[ada]
3	[bad]
4	[bba]
5	[dab]
6	[do\$]
7	[o\$]\$]

Suffixe von  $R'$

$S_1(R')$	1 2 4 6 4 5 3 7
$S_2(R')$	2 4 6 4 5 3 7
$S_3(R')$	4 6 4 5 3 7
$S_4(R')$	6 4 5 3 7
$S_5(R')$	4 5 3 7
$S_6(R')$	5 3 7
$S_7(R')$	3 7
$S_8(R')$	7

# 1. Schritt: Erstelle Suffix-Array $A_{12}$

Wie  $S(R)$  effizient sortieren?

$R = [abb][ada][bba][do\$][bba][dab][bad][o\$\$]$

Wende Radixsort auf Zeichen von  $R$  an.  
(siehe Übung)

**Transformation:** Ersetze Zeichen in  $R$  durch ihren Rang  $\rightarrow R'$

$R' = 1\ 2\ 4\ 6\ 4\ 5\ 3\ 7$

Suffix  $S_i(R)$  entspricht  $S_i(R')$

→ Sortierung von  $S(R')$  liefert gewünschte Sortierung von  $S(R)$ .

$S(R) =$

$S_1(R)$	[abb][ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_2(R)$	[ada][bba][do\$][bba][dab][bad][o\$]\$]
$S_3(R)$	[bba][do\$][bba][dab][bad][o\$]\$]
$S_4(R)$	[do\$][bba][dab][bad][o\$]\$]
$S_5(R)$	[bba][dab][bad][o\$]\$]
$S_6(R)$	[dab][bad][o\$]\$]
$S_7(R)$	[bad][o\$]\$]
$S_8(R)$	[o\$]\$]

Rang	Zeichen
1	[abb]
2	[ada]
3	[bad]
4	[bba]
5	[dab]
6	[do\$]
7	[o\$]\$]

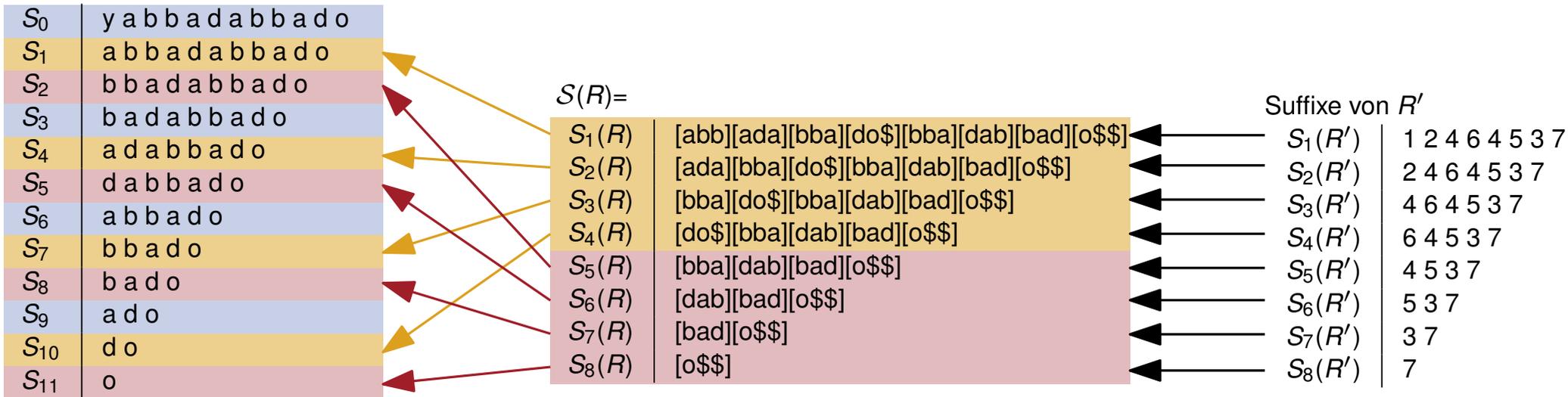
Suffixe von  $R'$

$S_1(R')$	1 2 4 6 4 5 3 7
$S_2(R')$	2 4 6 4 5 3 7
$S_3(R')$	4 6 4 5 3 7
$S_4(R')$	6 4 5 3 7
$S_5(R')$	4 5 3 7
$S_6(R')$	5 3 7
$S_7(R')$	3 7
$S_8(R')$	7

Erstelle Suffix-Array  $A'$  von  $R'$ :

Rekursiver Aufruf  $A' \leftarrow \text{SUFFIXARRAY}(R')$

# 1. Schritt: Erstelle Suffix-Array $A_{12}$

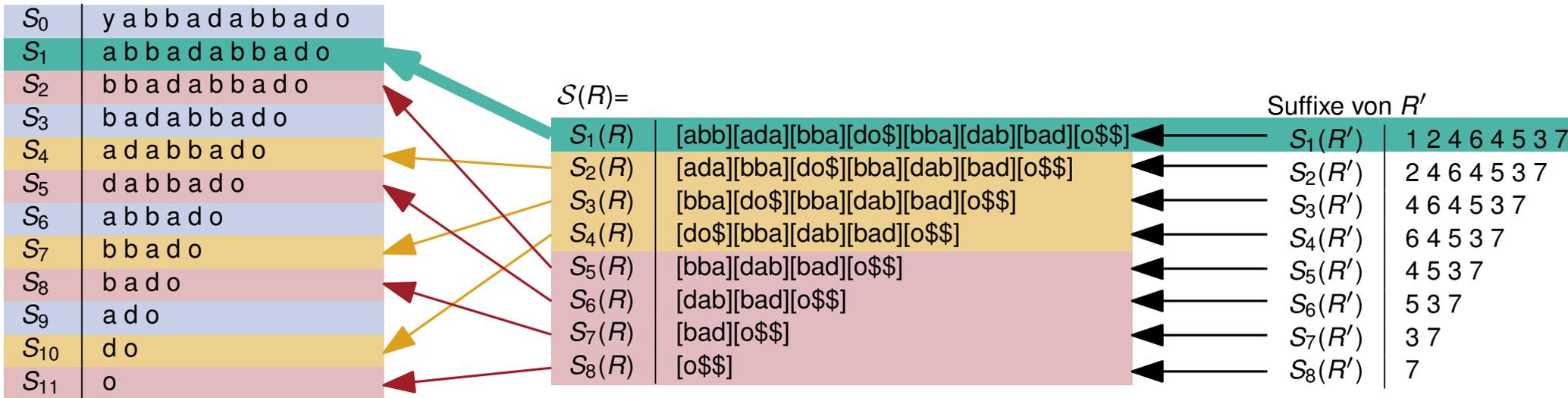


Sortierung von  $S(R')$   $\xrightarrow{\text{liefert}}$  Sortierung von  $S(R)$   $\xrightarrow{\text{liefert}}$  Sortierung von  $S_1 \cup S_2$

↓  
Suffix-Array  $A_{12}$

$A_{12} =$			
1	$S_1$	a b b a d a b b a d o	$S_1(R')$
2	$S_4$	a d a b b a d o	$S_2(R')$
3	$S_8$	b a d o	$S_7(R')$
4	$S_2$	b b a d a b b a d o	$S_5(R')$
5	$S_7$	b b a d o	$S_3(R')$
6	$S_5$	d a b b a d o	$S_6(R')$
7	$S_{10}$	d o	$S_4(R')$
8	$S_{11}$	o	$S_8(R')$

# 1. Schritt: Erstelle Suffix-Array $A_{12}$



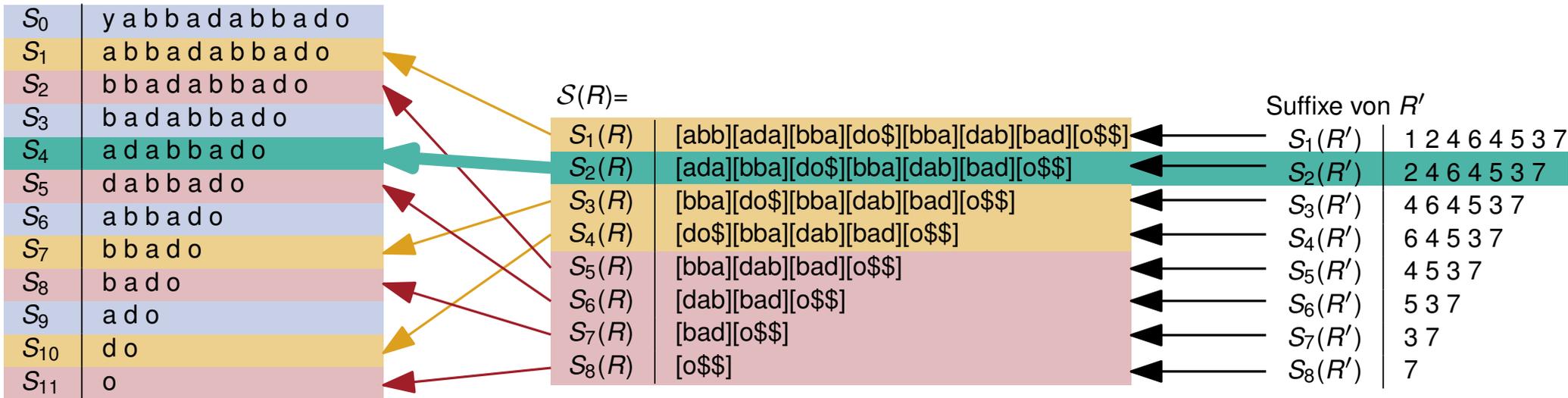
Sortierung von  $S(R')$   $\xrightarrow{\text{liefert}}$  Sortierung von  $S(R)$   $\xrightarrow{\text{liefert}}$  Sortierung von  $S_1 \cup S_2$

↓  
Suffix-Array  $A_{12}$

$A_{12} =$

1	$S_1$	a b b a d a b b a d o	$S_1(R')$
2	$S_4$	a d a b b a d o	$S_2(R')$
3	$S_8$	b a d o	$S_7(R')$
4	$S_2$	b b a d a b b a d o	$S_5(R')$
5	$S_7$	b b a d o	$S_3(R')$
6	$S_5$	d a b b a d o	$S_6(R')$
7	$S_{10}$	d o	$S_4(R')$
8	$S_{11}$	o	$S_8(R')$

# 1. Schritt: Erstelle Suffix-Array $A_{12}$



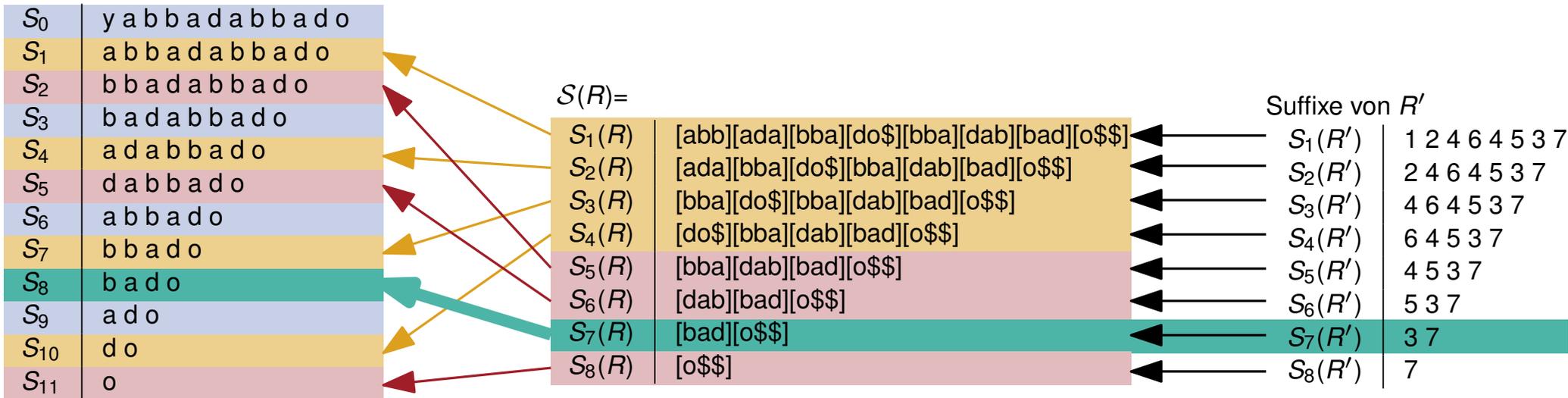
Sortierung von  $\mathcal{S}(R')$   $\xrightarrow{\text{liefert}}$  Sortierung von  $\mathcal{S}(R)$   $\xrightarrow{\text{liefert}}$  Sortierung von  $\mathcal{S}_1 \cup \mathcal{S}_2$

↓  
Suffix-Array  $A_{12}$

$A_{12} =$

1	$S_1$	a b b a d a b b a d o	$S_1(R')$
2	$S_4$	a d a b b a d o	$S_2(R')$
3	$S_8$	b a d o	$S_7(R')$
4	$S_2$	b b a d a b b a d o	$S_5(R')$
5	$S_7$	b b a d o	$S_3(R')$
6	$S_5$	d a b b a d o	$S_6(R')$
7	$S_{10}$	d o	$S_4(R')$
8	$S_{11}$	o	$S_8(R')$

# 1. Schritt: Erstelle Suffix-Array $A_{12}$



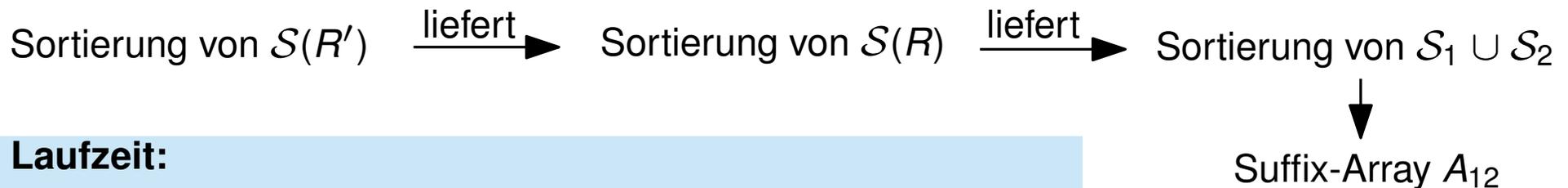
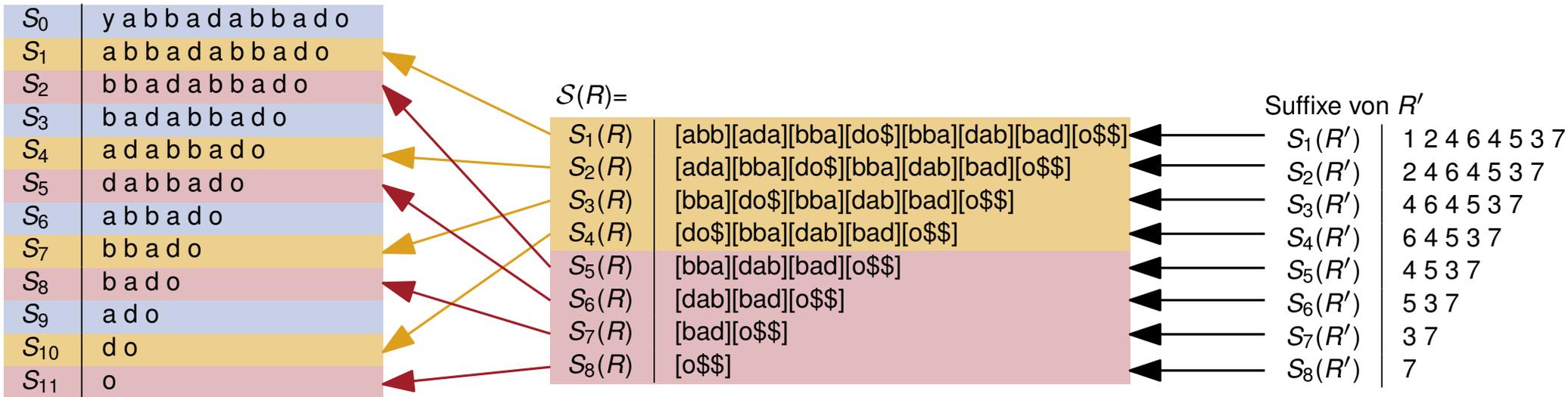
Sortierung von  $S(R')$   $\xrightarrow{\text{liefert}}$  Sortierung von  $S(R)$   $\xrightarrow{\text{liefert}}$  Sortierung von  $S_1 \cup S_2$

↓  
Suffix-Array  $A_{12}$

$A_{12} =$

1	$S_1$	a b b a d a b b a d o	$S_1(R')$
2	$S_4$	a d a b b a d o	$S_2(R')$
3	$S_8$	b a d o	$S_7(R')$
4	$S_2$	b b a d a b b a d o	$S_5(R')$
5	$S_7$	b b a d o	$S_3(R')$
6	$S_5$	d a b b a d o	$S_6(R')$
7	$S_{10}$	d o	$S_4(R')$
8	$S_{11}$	o	$S_8(R')$

# 1. Schritt: Erstelle Suffix-Array $A_{12}$



## Laufzeit:

$O(n)$  Zeit für Konstruktion von  $R$ .

$O(n)$  Zeit für Konstruktion von  $R'$ .

**Laufzeit  $T_1(n)$  für 1. Schritt**  $T_1(n) = O(n) + T(\frac{2}{3}n)$

wobei  $T(n)$  = Laufzeit für Konstruktion eines Suffix-Arrays

# Konstruktion von Suffix-Arrays

**Eingabe:** Text  $T := t_0 t_1 \dots t_{n-1}$

	0	1	2	3	4	5	6	7	8	9	10	11
T=	y	a	b	b	a	d	a	b	b	a	d	o

**Gesucht:** Suffix-Array  $A$  von  $T$ , d.h. lexikographische Sortierung von Suffixen von  $T$ .

Kürze  $x \bmod y = z$  mit  $x \equiv z(y)$  ab.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

$\mathcal{S}$  = Suffixe von  $T$

$\mathcal{S}_0$  = Suffixe mit Index  $i \equiv 0(3)$

$\mathcal{S}_1$  = Suffixe mit Index  $i \equiv 1(3)$

$\mathcal{S}_2$  = Suffixe mit Index  $i \equiv 2(3)$

**SUFFIXARRAY**(Text  $T = t_0 t_1 \dots t_{n-1}$ )

**wenn**  $n = O(1)$  **dann**

    Konstruiere  $A$  in  $O(1)$  Zeit.

**sonst**

    Berechne Suffix-Array  $A_{12}$  für  $\mathcal{S}_1 \cup \mathcal{S}_2$ .

    Berechne Suffix-Array  $A_0$  für  $\mathcal{S}_0$  basierend auf  $A_{12}$ .

    Vermenge  $A_{12}$  mit  $A_0$ .

## 2. Schritt: Konstruktion von $A_0$

**Beobachtung:** Für zwei  $S_i, S_j \in \mathcal{S}_0$  gilt

$$S_i < S_j \text{ genau dann wenn } t_i < t_j, \text{ oder } t_i = t_j \text{ und } S_{i+1} < S_{j+1}.$$

Beobachtung gibt Definition der lexikographischen Ordnung wieder.

↳ Beschreibt wie Suffixe  $\mathcal{S}_0$  sortiert werden müssen. Verwende Radixsort.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Erinnerung:**

$t_i =$  erstes Zeichen von  $S_i$

## 2. Schritt: Konstruktion von $A_0$

**Beobachtung:** Für zwei  $S_i, S_j \in \mathcal{S}_0$  gilt

$$S_i < S_j \text{ genau dann wenn } t_i < t_j, \text{ oder } t_i = t_j \text{ und } S_{i+1} < S_{j+1}.$$

Beobachtung gibt Definition der lexikographischen Ordnung wider.

↳ Beschreibt wie Suffixe  $\mathcal{S}_0$  sortiert werden müssen. Verwende Radixsort.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 1 folgt  $A_{12}$ :**

1	$S_1$	a b b a d a b b a d o
2	$S_4$	a d a b b a d o
3	$S_8$	b a d o
4	$S_2$	b b a d a b b a d o
5	$S_7$	b b a d o
6	$S_5$	d a b b a d o
7	$S_{10}$	d o
8	$S_{11}$	o

**Erinnerung:**

$t_i =$  erstes Zeichen von  $S_i$

## 2. Schritt: Konstruktion von $A_0$

**Beobachtung:** Für zwei  $S_i, S_j \in \mathcal{S}_0$  gilt

$$S_i < S_j \text{ genau dann wenn } t_i < t_j, \text{ oder } t_i = t_j \text{ und } S_{i+1} < S_{j+1}.$$

Beobachtung gibt Definition der lexikographischen Ordnung wider.

↳ Beschreibt wie Suffixe  $\mathcal{S}_0$  sortiert werden müssen. Verwende Radixsort.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 1 folgt  $A_{12}$ :**

1	$S_1$	a b b a d a b b a d o
2	$S_4$	a d a b b a d o
3	$S_8$	b a d o
4	$S_2$	b b a d a b b a d o
5	$S_7$	b b a d o
6	$S_5$	d a b b a d o
7	$S_{10}$	d o
8	$S_{11}$	o

$S_6 < S_9$  weil  $S_7 < S_{10}$   
 $S_9 < S_3$  weil  $a < b$   
 $S_3 < S_0$  weil  $b < y$

$\Rightarrow S_6 < S_9 < S_3 < S_0$   


 beschreibt  $A_0$

**Erinnerung:**

$t_i =$  erstes Zeichen von  $S_i$

# Konstruktion von Suffix-Arrays

**Eingabe:** Text  $T := t_0 t_1 \dots t_{n-1}$

	0	1	2	3	4	5	6	7	8	9	10	11
T=	y	a	b	b	a	d	a	b	b	a	d	o

**Gesucht:** Suffix-Array  $A$  von  $T$ , d.h. lexikographische Sortierung von Suffixen von  $T$ .

Kürze  $x \bmod y = z$  mit  $x \equiv z(y)$  ab.

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

$\mathcal{S}$  = Suffixe von  $T$

$\mathcal{S}_0$  = Suffixe mit Index  $i \equiv 0(3)$

$\mathcal{S}_1$  = Suffixe mit Index  $i \equiv 1(3)$

$\mathcal{S}_2$  = Suffixe mit Index  $i \equiv 2(3)$

SUFFIXARRAY(Text  $T = t_0 t_1 \dots t_{n-1}$ )

**wenn**  $n = O(1)$  **dann**

| Konstruiere  $A$  in  $O(1)$  Zeit.

**sonst**

| Berechne Suffix-Array  $A_{12}$  für  $\mathcal{S}_1 \cup \mathcal{S}_2$ .

| Berechne Suffix-Array  $A_0$  für  $\mathcal{S}_0$  basierend auf  $A_{12}$ .

| Vermenge  $A_{12}$  mit  $A_0$ .

### 3. Schritt: Vermengen von $A_0$ und $A_{12}$

**Beobachtung:** Sei  $S_i \in \mathcal{S}_0$

1. Für  $S_j \in \mathcal{S}_1$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $S_{i+1} < S_{j+1}$ .
2. Für  $S_j \in \mathcal{S}_2$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $t_{i+1} < t_{j+1}$ , oder  $t_i t_{i+1} = t_j t_{j+1}$  und  $S_{i+2} < S_{j+2}$ .

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Erinnerung:**

$t_i =$  erstes Zeichen von  $S_i$

### 3. Schritt: Vermengen von $A_0$ und $A_{12}$

**Beobachtung:** Sei  $S_i \in \mathcal{S}_0$

1. Für  $S_j \in \mathcal{S}_1$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $S_{i+1} < S_{j+1}$ .
2. Für  $S_j \in \mathcal{S}_2$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $t_{i+1} < t_{j+1}$ , oder  $t_i t_{i+1} = t_j t_{j+1}$  und  $S_{i+2} < S_{j+2}$ .

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 2 folgt:**

$$S_1 < S_4 < S_8 < S_2 < S_7 < S_5 < S_{10} < S_{11}$$

**Aus Schritt 2 folgt:**

$$S_6 < S_9 < S_3 < S_0$$

**Vermengen wie bei Merge-Sort:**

**Erinnerung:**

$t_i =$  erstes Zeichen von  $S_i$

### 3. Schritt: Vermengen von $A_0$ und $A_{12}$

**Beobachtung:** Sei  $S_i \in \mathcal{S}_0$

1. Für  $S_j \in \mathcal{S}_1$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $S_{i+1} < S_{j+1}$ .
2. Für  $S_j \in \mathcal{S}_2$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $t_{i+1} < t_{j+1}$ , oder  $t_i t_{i+1} = t_j t_{j+1}$  und  $S_{i+2} < S_{j+2}$ .

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 2 folgt:**

$S_1 < S_4 < S_8 < S_2 < S_7 < S_5 < S_{10} < S_{11}$

**Aus Schritt 2 folgt:**

$S_6 < S_9 < S_3 < S_0$

**Vermengen wie bei Merge-Sort:**

$S_1$

**Erinnerung:**

$t_i =$  erstes Zeichen von  $S_i$

### 3. Schritt: Vermengen von $A_0$ und $A_{12}$

**Beobachtung:** Sei  $S_i \in \mathcal{S}_0$

1. Für  $S_j \in \mathcal{S}_1$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $S_{i+1} < S_{j+1}$ .
2. Für  $S_j \in \mathcal{S}_2$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $t_{i+1} < t_{j+1}$ , oder  $t_i t_{i+1} = t_j t_{j+1}$  und  $S_{i+2} < S_{j+2}$ .

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 2 folgt:**

$$S_1 < S_4 < S_8 < S_2 < S_7 < S_5 < S_{10} < S_{11}$$

**Aus Schritt 2 folgt:**

$$S_6 < S_9 < S_3 < S_0$$

**Vermengen wie bei Merge-Sort:**

$$S_1 < S_6$$

**Erinnerung:**

$t_i$  = erstes Zeichen von  $S_i$

### 3. Schritt: Vermengen von $A_0$ und $A_{12}$

**Beobachtung:** Sei  $S_i \in \mathcal{S}_0$

1. Für  $S_j \in \mathcal{S}_1$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $S_{i+1} < S_{j+1}$ .
2. Für  $S_j \in \mathcal{S}_2$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $t_{i+1} < t_{j+1}$ , oder  $t_i t_{i+1} = t_j t_{j+1}$  und  $S_{i+2} < S_{j+2}$ .

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 2 folgt:**

$$S_1 < S_4 < S_8 < S_2 < S_7 < S_5 < S_{10} < S_{11}$$

**Aus Schritt 2 folgt:**

$$S_6 < S_9 < S_3 < S_0$$

**Vermengen wie bei Merge-Sort:**

$$S_1 < S_6 < S_4$$

**Erinnerung:**

$t_i$  = erstes Zeichen von  $S_i$

### 3. Schritt: Vermengen von $A_0$ und $A_{12}$

**Beobachtung:** Sei  $S_i \in \mathcal{S}_0$

1. Für  $S_j \in \mathcal{S}_1$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $S_{i+1} < S_{j+1}$ .
2. Für  $S_j \in \mathcal{S}_2$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $t_{i+1} < t_{j+1}$ , oder  $t_i t_{i+1} = t_j t_{j+1}$  und  $S_{i+2} < S_{j+2}$ .

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 2 folgt:**

$$S_1 < S_4 < S_8 < S_2 < S_7 < S_5 < S_{10} < S_{11}$$

**Aus Schritt 2 folgt:**

$$S_6 < S_9 < S_3 < S_0$$

**Vermengen wie bei Merge-Sort:**

$$S_1 < S_6 < S_4 < S_9$$

**Erinnerung:**

$t_i =$  erstes Zeichen von  $S_i$

### 3. Schritt: Vermengen von $A_0$ und $A_{12}$

**Beobachtung:** Sei  $S_i \in \mathcal{S}_0$

1. Für  $S_j \in \mathcal{S}_1$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $S_{i+1} < S_{j+1}$ .
2. Für  $S_j \in \mathcal{S}_2$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $t_{i+1} < t_{j+1}$ , oder  $t_i t_{i+1} = t_j t_{j+1}$  und  $S_{i+2} < S_{j+2}$ .

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 2 folgt:**

$$S_1 < S_4 < S_8 < S_2 < S_7 < S_5 < S_{10} < S_{11}$$

**Aus Schritt 2 folgt:**

$$S_6 < S_9 < S_3 < S_0$$

**Vermengen wie bei Merge-Sort:**

$$S_1 < S_6 < S_4 < S_9 < S_3$$

**Erinnerung:**

$t_i$  = erstes Zeichen von  $S_i$

### 3. Schritt: Vermengen von $A_0$ und $A_{12}$

**Beobachtung:** Sei  $S_i \in \mathcal{S}_0$

1. Für  $S_j \in \mathcal{S}_1$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $S_{i+1} < S_{j+1}$ .
2. Für  $S_j \in \mathcal{S}_2$  gilt:  $S_i < S_j$  genau dann, wenn  $t_i < t_j$ , oder  $t_i = t_j$  und  $t_{i+1} < t_{j+1}$ , oder  $t_i t_{i+1} = t_j t_{j+1}$  und  $S_{i+2} < S_{j+2}$ .

$S_0$	y a b b a d a b b a d o
$S_1$	a b b a d a b b a d o
$S_2$	b b a d a b b a d o
$S_3$	b a d a b b a d o
$S_4$	a d a b b a d o
$S_5$	d a b b a d o
$S_6$	a b b a d o
$S_7$	b b a d o
$S_8$	b a d o
$S_9$	a d o
$S_{10}$	d o
$S_{11}$	o

**Aus Schritt 2 folgt:**

$$S_1 < S_4 < S_8 < S_2 < S_7 < S_5 < S_{10} < S_{11}$$

**Aus Schritt 2 folgt:**

$$S_6 < S_9 < S_3 < S_0$$

**Vermengen wie bei Merge-Sort:**

$$S_1 < S_6 < S_4 < S_9 < S_3 <$$

$$S_8 < S_2 < S_7 < S_5 < S_{10} < S_{11} < S_0$$

**Erinnerung:**

$t_i =$  erstes Zeichen von  $S_i$

**Theorem:** Ein Suffix-Array kann für einen Text  $T$  der Länge  $n$  in  $O(n)$  Zeit konstruiert werden.

SUFFIXARRAY(Text  $T = t_0 t_1 \dots t_{n-1}$ )

**wenn**  $n = O(1)$  **dann**

| Konstruiere  $A$  in  $O(1)$  Zeit.

**sonst**

| Berechne Suffix-Array  $A_{12}$  für  $\mathcal{S}_1 \cup \mathcal{S}_2$ .

| Berechne Suffix-Array  $A_0$  für  $\mathcal{S}_0$  basierend auf  $A_{12}$ .

| Vermenge  $A_{12}$  mit  $A_0$ .

**Laufzeit:**

$$T(n) = \begin{cases} O(1), & \text{falls } n=O(1) \\ O(n) + T\left(\frac{2}{3}n\right), & \text{sonst} \end{cases}$$

Kann mithilfe des Master-Theorems aufgelöst werden.

# Konstruktion von Suffixbäumen

**Beobachtung:** Bisher sind Suffix-Arrays im Aufbau schneller, aber in der Abfrage langsamer als Suffix-Bäume.

**Idee:** Konstruiere Suffixbäume aus Suffix-Arrays in  $O(n)$ -Zeit.

**Definition 11:** Gegeben ein Text  $T$  der Länge  $n$  und dessen Suffix-Array  $A$ . Ein *LCP-Array*  $L$  ist ein Array der Größe  $n$ , sodass

- $L[1] = 0$
- $L[i] =$  Länge des längsten gemeinsamen Präfixes von  $A[i]$  und  $A[i - 1]$  für  $i > 1$

LCP = Longest Common Prefix

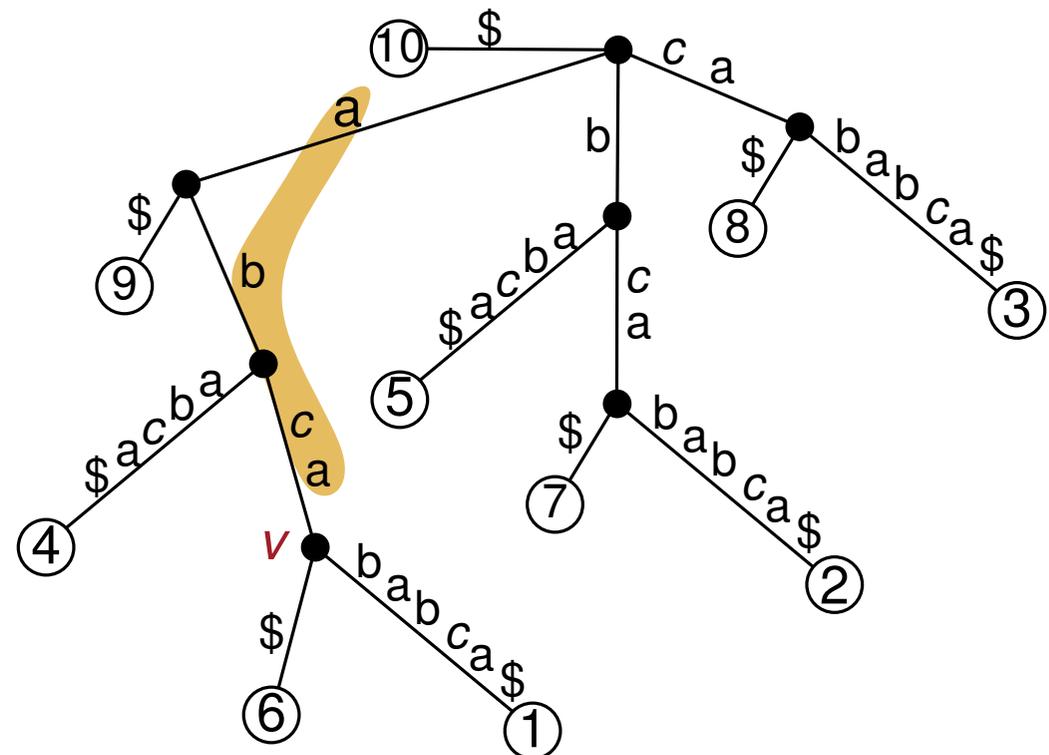
# Beispiel LCP-Array

**Definition 11:** Gegeben ein Text  $T$  der Länge  $n$  und dessen Suffix-Array  $A$ . Ein  $LCP$ -Array  $L$  ist ein Array der Größe  $n$ , sodass

- $L[1] = 0$
- $L[i] =$  Länge des längsten gemeinsamen Präfixes von  $A[i]$  und  $A[i - 1]$  für  $i > 1$

**Beispiel:**  $T = a b c a b a b c a \$$

$A =$	10	9	4	6	1	5	7	2	8	3
$L =$	0	0	1	2	4	0	1	3	0	2
	\$	a	a	a	a	b	b	b	c	c
		\$	b	b	b	a	c	c	a	a
			a	c	c	b	a	a	\$	b
			b	a	a	c	\$	b		b
			c	\$	b	a		a	b	c
			a		a	\$		b	c	a
			\$		b	c	a		\$	
					c	a				\$
					a					



**Bsp.:**  $v$  ist niedrigster gemeinsamer Vorgänger von Knoten 6 und 1.  $d(v)$  entspricht der Länge des kleinsten gemeinsamen Präfixes von  $S_1$  und  $S_6$ .

→  $L$  gibt für zwei Suffixe an, an welcher Stelle sich Baum aufgespaltet.

# Konstruktion von Suffixbäumen

**Gegeben:** Text  $T$ , Suffix-Array  $A$  und LCP-Array  $L$  von  $T$

**Gesucht:** Suffixbaum von  $T$

**Idee:** Konstruiere  $N_0, N_1, \dots, N_n$  Suffixbäume, sodass  $N_i$  Suffixbaum von  $S_{A[1]}, \dots, S_{A[i]}$  ist, d.h. Suffix-Array gibt Reihenfolge an, in der Suffixe eingefügt werden.

**Initialisierung:**  $N_0$  enthält nur die Wurzel und repräsentiert damit das leere Wort.

**Konstruktion von  $N_{i+1}$  aus  $N_i$ :**

1. Betrachte den Pfad  $P$  des Suffixes  $S_{A[i]}$  in  $N_i$ , also den rechtesten Pfad in  $N_i$ .
2. Wähle den tiefsten Knoten  $v$  auf  $P$  mit  $d(v) \leq L[i + 1]$ .

**1. Fall:**  $d(v) = L[i + 1]$ .

Hänge an  $v$  ein neues Blatt  $x$  und beschrifte  $(v, x)$  mit  $T[A[i + 1] + L[i + 1], n]$ . Blatt  $x$  wird mit  $A[i + 1]$  beschriftet.

# Konstruktion von Suffixbäumen

**Gegeben:** Text  $T$ , Suffix-Array  $A$  und LCP-Array  $L$  von  $T$

**Gesucht:** Suffixbaum von  $T$

**Idee:** Konstruiere  $N_0, N_1, \dots, N_n$  Suffixbäume, sodass  $N_i$  Suffixbaum von  $S_{A[1]}, \dots, S_{A[i]}$  ist, d.h. Suffix-Array gibt Reihenfolge an, in der Suffixe eingefügt werden.

**Initialisierung:**  $N_0$  enthält nur die Wurzel und repräsentiert damit das leere Wort.

**Konstruktion von  $N_{i+1}$  aus  $N_i$ :**

1. Betrachte den Pfad  $P$  des Suffixes  $S_{A[i]}$  in  $N_i$ , also den rechtesten Pfad in  $N_i$ .
2. Wähle den tiefsten Knoten  $v$  auf  $P$  mit  $d(v) \leq L[i + 1]$ .

**2. Fall  $d(v) < L[i + 1]$ .** Sei  $w$  das Kind  $v$  auf dem rechtesten Pfad von  $N_i$ .

1. Entferne  $(v, w)$ .
2. Füge neuen Knoten  $y$  mit Kante  $(v, y)$  ein. Beschriftung von  $(v, y)$ :  
 $T[A[i] + d(v), A[i] + L[i + 1] - 1]$
3. Füge  $(y, w)$  ein mit Beschriftung  
 $T[A[i] + L[i + 1], A[i] + d(w) - 1]$ .
4. Füge Blatt  $x$  mit Beschriftung  $A[i + 1]$  ein und Kante  $(y, x)$  mit Beschriftung  
 $T[A[i + 1] + L[i + 1], n]$ .

# Beispiel



T = a b c a b a b c a \$

A =	10	9	4	6	1	5	7	2	8	3
L =	0	0	1	2	4	0	1	3	0	2

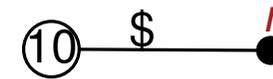
\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b			c		a
				c			a		\$
				a			\$		
				\$					

# Beispiel

$T = a b c a b a b c a \$$

$A =$	10	9	4	6	1	5	7	2	8	3
$L =$	0	0	1	2	4	0	1	3	0	2

\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b			c		a
				c			a		\$
				a					
				\$					



1. Betrachte den Pfad  $P$  des Suffixes  $S_{A[i]}$  in  $N_i$ , also den rechtesten Pfad in  $N_i$ .
2. Wähle den tiefsten Knoten  $v$  auf  $P$  mit  $d(v) \leq L[i + 1]$ .

**1. Fall:**  $d(v) = L[i + 1]$ .

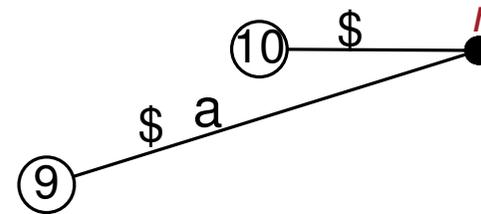
Hänge an  $v$  ein neues Blatt  $x$  und beschrifte  $(v, x)$  mit  $T[A[i + 1] + L[i + 1], n]$ . Blatt  $x$  wird mit  $A[i + 1]$  beschriftet.

# Beispiel

T = a b c a b a b c a \$

A =	10	9	4	6	1	5	7	2	8	3
L =	0	0	1	2	4	0	1	3	0	2

\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b			c		a
				c			a		\$
				a			\$		
				\$					



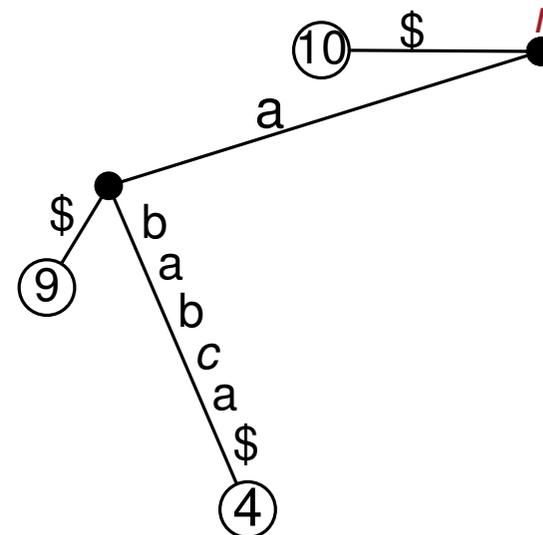
1. Betrachte den Pfad  $P$  des Suffixes  $S_{A[i]}$  in  $N_i$ , also den rechten Pfad in  $N_i$ .
2. Wähle den tiefsten Knoten  $v$  auf  $P$  mit  $d(v) \leq L[i + 1]$ .

**1. Fall:**  $d(v) = L[i + 1]$ .  
 Hänge an  $v$  ein neues Blatt  $x$  und beschrifte  $(v, x)$  mit  $T[A[i + 1] + L[i + 1], n]$ . Blatt  $x$  wird mit  $A[i + 1]$  beschriftet.

# Beispiel

$T = a b c a b a b c a \$$

$A =$	10	9	4	6	1	5	7	2	8	3
$L =$	0	0	1	2	4	0	1	3	0	2
	\$	a	a	a	a	b	b	b	c	c
		\$	b	b	b	a	c	c	a	a
			a	c	c	b	a	a	\$	b
			b	a	a	c	\$	b		a
			c	\$	b	a		a		b
			a		a	\$		b		c
			\$		b	c		a		\$
					c	a				
					a					
					\$					



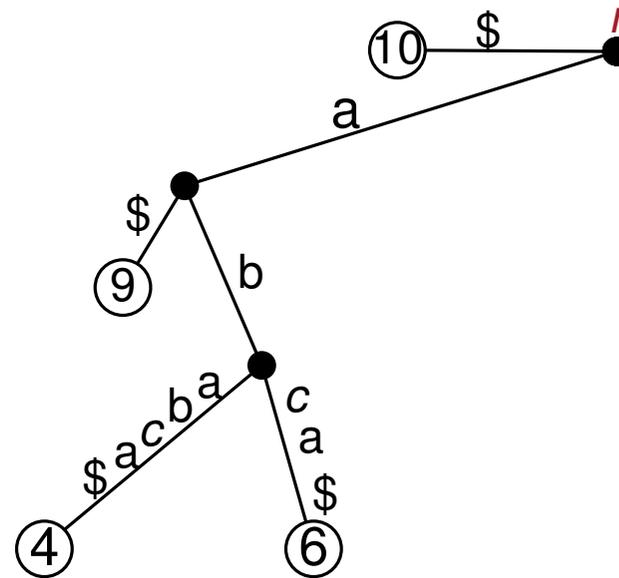
1. Betrachte den Pfad  $P$  des Suffixes  $S_{A[i]}$  in  $N_i$ , also den rechten Pfad in  $N_i$ .
2. Wähle den tiefsten Knoten  $v$  auf  $P$  mit  $d(v) \leq L[i + 1]$ .

- 2. Fall**  $d(v) < L[i + 1]$ . Sei  $w$  das Kind  $v$  auf dem rechten Pfad von  $N_i$ .
1. Entferne  $(v, w)$ .
  2. Füge neuen Knoten  $y$  mit Kante  $(v, y)$  ein. Beschriftung von  $(v, y)$ :  $T[A[i] + d(v), A[i] + L[i + 1] - 1]$
  3. Füge  $(y, w)$  ein mit Beschriftung  $T[A[i] + L[i + 1], A[i] + d(w) - 1]$ .
  4. Füge Blatt  $x$  mit Beschriftung  $A[i + 1]$  ein und Kante  $(y, x)$  mit Beschriftung  $T[A[i + 1] + L[i + 1], n]$ .

# Beispiel

$T = a b c a b a b c a \$$

$A =$	10	9	4	6	1	5	7	2	8	3
$L =$	0	0	1	2	4	0	1	3	0	2
	\$	a	a	a	a	b	b	b	c	c
		\$	b	b	b	a	c	c	a	a
			a	c	c	b	a	a	\$	b
			b	a	a	c	\$	b		a
			c	\$	b	a		a		b
			a		a	\$		b		c
			\$		b	c		a		a
					c	a		\$		
					a					
					\$					



1. Betrachte den Pfad  $P$  des Suffixes  $S_{A[i]}$  in  $N_i$ , also den rechten Pfad in  $N_i$ .
2. Wähle den tiefsten Knoten  $v$  auf  $P$  mit  $d(v) \leq L[i + 1]$ .

**2. Fall**  $d(v) < L[i + 1]$ . Sei  $w$  das Kind  $v$  auf dem rechten Pfad von  $N_i$ .

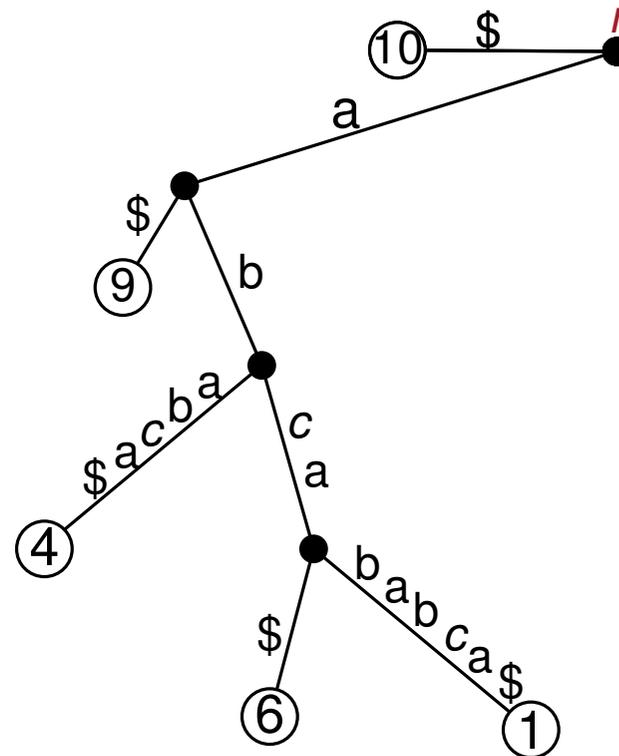
1. Entferne  $(v, w)$ .
2. Füge neuen Knoten  $y$  mit Kante  $(v, y)$  ein. Beschriftung von  $(v, y)$ :  $T[A[i] + d(v), A[i] + L[i + 1] - 1]$
3. Füge  $(y, w)$  ein mit Beschriftung  $T[A[i] + L[i + 1], A[i] + d(w) - 1]$ .
4. Füge Blatt  $x$  mit Beschriftung  $A[i + 1]$  ein und Kante  $(y, x)$  mit Beschriftung  $T[A[i + 1] + L[i + 1], n]$ .

# Beispiel

T = a b c a b a b c a \$

A =	10	9	4	6	1	5	7	2	8	3
L =	0	0	1	2	4	0	1	3	0	2

\$	a	a	a	a	b	b	b	c	c
	\$	b	b	b	a	c	c	a	a
		a	c	c	b	a	a	\$	b
		b	a	a	c	\$	b		a
		c	\$	b	a		a		b
		a		a	\$		b		c
		\$		b	c		a		\$
				c	a		\$		
				\$					



1. Betrachte den Pfad  $P$  des Suffixes  $S_{A[i]}$  in  $N_i$ , also den rechten Pfad in  $N_i$ .
2. Wähle den tiefsten Knoten  $v$  auf  $P$  mit  $d(v) \leq L[i + 1]$ .

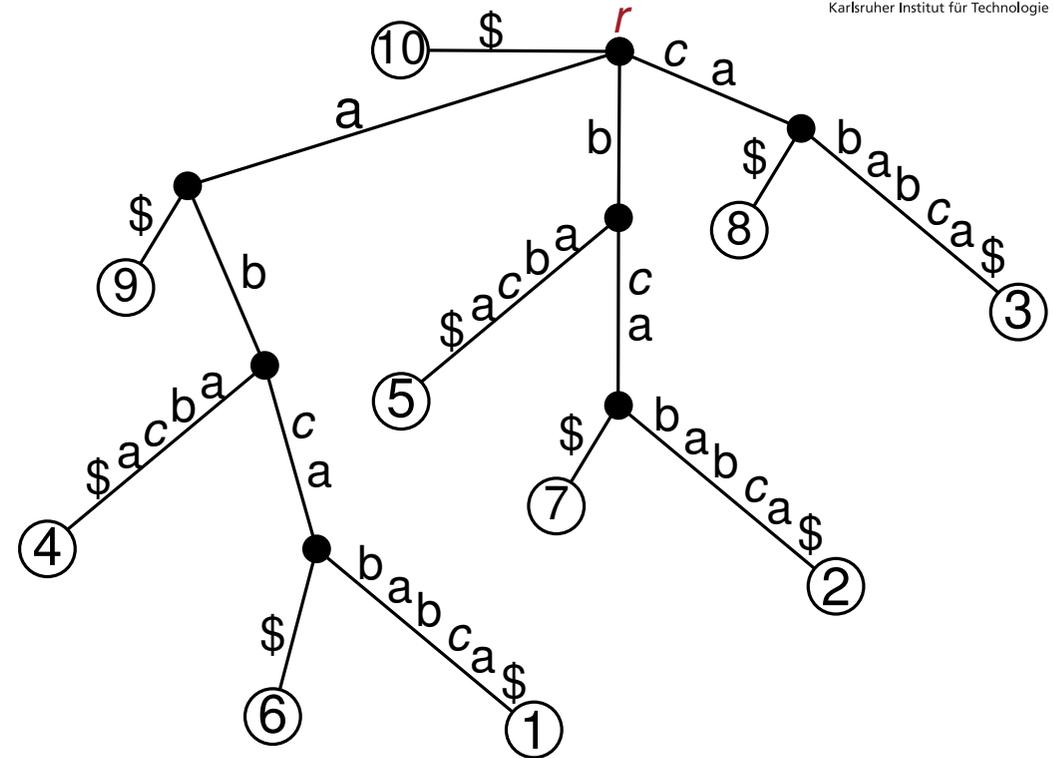
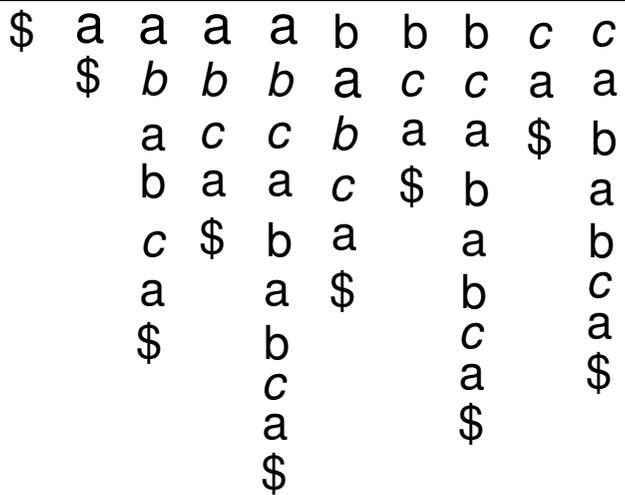
**2. Fall**  $d(v) < L[i + 1]$ . Sei  $w$  das Kind  $v$  auf dem rechten Pfad von  $N_i$ .

1. Entferne  $(v, w)$ .
2. Füge neuen Knoten  $y$  mit Kante  $(v, y)$  ein. Beschriftung von  $(v, y)$ :  $T[A[i] + d(v), A[i] + L[i + 1] - 1]$
3. Füge  $(y, w)$  ein mit Beschriftung  $T[A[i] + L[i + 1], A[i] + d(w) - 1]$ .
4. Füge Blatt  $x$  mit Beschriftung  $A[i + 1]$  ein und Kante  $(y, x)$  mit Beschriftung  $T[A[i + 1] + L[i + 1], n]$ .

# Beispiel

$T = a b c a b a b c a \$$

$A =$	10	9	4	6	1	5	7	2	8	3
$L =$	0	0	1	2	4	0	1	3	0	2



**Theorem 28:** Gegeben ein Suffix-Array  $A$  und ein LCP-Array  $L$  eines Texts  $T$ , dann kann der Suffixbaum von  $T$  in  $O(n)$  Zeit konstruiert werden.

**Ohne Beweis:** LCP-Array kann in  $O(n)$  Zeit konstruiert werden.