

5. Übungsblatt zu Algorithmen II im WS 2023/2024

<https://algo2.iti.kit.edu/AlgorithmenII.WS23.php>
{sanders, moritz.laupichler, nikolai.maas}@kit.edu

Musterlösungen

Aufgabe 1 (Rechnen+Analyse: Suche in Strings)

- Zur Ausführung des KMP-Algorithmus muss zunächst ein sogenanntes *border-array* berechnet werden. Geben Sie das *border-array* für das Suchmuster $P = \text{abacababc}$ an.
- Führen Sie den KMP-Algorithmus auf dem Text $T = \text{abacababbabacababc}$ mit obigem Suchmuster durch.
- Wie oft muss der KMP-Algorithmus ein Muster P der Länge $|P|$ maximal an einen Text T der Länge $|T|$ anlegen, falls P nicht in T vorkommt? Wie oft minimal? Geben Sie das Ergebnis in Abhängigkeit von $|P|$ und $|T|$ an. Geben Sie außerdem jeweils ein Beispiel für P und T an.
- Zeigen oder widerlegen Sie:
Das *border-array* kann keine drei aufeinanderfolgenden Einträge enthalten, die jeweils um eins kleiner sind als ihr Vorgänger, falls der erste von diesen drei Einträgen auf ein Suffix der Größe $k \geq 3$ verweist, welches gleichzeitig echtes Präfix ist.
Beispiel: $\text{border}[10] = 3, \text{border}[11] = 2, \text{border}[12] = 1$
Kein Beispiel: $\text{border}[10] = 2, \text{border}[11] = 1, \text{border}[12] = 0$

Musterlösung:

a) Es ergibt sich folgendes *border-array*:

P		a	b	a	c	a	b	a	b	c
border[]		-1	0	0	1	0	1	2	3	2

b) Das Suchmuster wird insgesamt 4 mal angelegt:

T		a	b	a	c	a	b	a	b	b	a	b	a	c	a	b	a	b	c	
		a	b	a	c	a	b	a	b	c										Stelle 1
							a	b	a	c	a	b	a	b	c					Stelle 7
								a	b	a	c	a	b	a	b	c				Stelle 9
									a	b	a	c	a	b	a	b	c			Stelle 10

c) Das Muster muss maximal $|T| - |P| + 1$ mal angelegt werden. Dieser Fall tritt z.B. auf, falls das erste Zeichen von P nicht in T auftaucht. Das Muster muss minimal $\lceil |T| / (|P| - 1) \rceil$ mal angelegt werden. Dieser Fall tritt z.B. auf, falls P aus einer Folge paarweise unterschiedlicher Zeichen besteht und T aus Konkatenationen von P ohne das letzte Zeichen.

d) Zu zeigen oder widerlegen ist die Existenz von drei Einträgen:

$$\text{border}[i] = k, \text{border}[i + 1] = k - 1, \text{border}[i + 2] = k - 2.$$

Einträge dieser Art können nicht existieren, da in diesem Fall gelten würde:

$$\text{für } \text{border}[i + 0] = k - 0: P_{k-2} = P_{i-3}, P_{k-1} = P_{i-2}, P_k = P_{i-1},$$

$$\text{für } \text{border}[i + 1] = k - 1: P_{k-2} = P_{i-1}, P_{k-1} = P_i \text{ und}$$

$$\text{für } \text{border}[i + 2] = k - 2: P_{k-2} = P_{i+1}.$$

Damit würde sich ergeben:

$$P_{k-2} = P_{i-3} = P_{i-1} = P_{i+1} = P_k,$$

$$P_{k-1} = P_{i-2} = P_i.$$

In diesem Fall wäre aber $\text{border}[i + 2] = k$, da

$$P_{k-2} = P_{i-1} = P_{i-3}, P_{k-1} = P_i = P_{i-2} \text{ und } P_k = P_{i+1} = P_{i-1}.$$

Andernfalls wäre auch $\text{border}[i + 0] \neq k$.

Aufgabe 2 (Rechnen: Burrows-Wheeler-Transformation)

Die Entropie $H(T)$ einer Zeichenkette T gibt eine untere Schranke für die Anzahl Bits pro Zeichen an, die zur Encodierung von T nötig wären, falls T einer Zufallsquelle entspränge. Sie ist definiert als

$$H(T) = \sum_{c \in \Sigma} p(c) \log_2(1/p(c))$$

wobei $p(c) := |\{s_i : s_i = c\}|/|T|$ die relative Häufigkeit von c in T ist.

- a) Bestimmen Sie die Entropie von Zeichenkette $s = \text{ababababab}$.
- b) Führen Sie die *Burrows-Wheeler-Transformation* auf Zeichenkette s aus. Wie groß ist jetzt die Entropie der Zeichenkette?
- c) Führen Sie eine *Move-to-Front* Kodierung auf dem Ergebnis durch.
(das Abschlusszeichen $\$$ aus der BWT muss bei der Kompression nicht berücksichtigt werden)
- d) Vergleichen Sie das Ergebnis mit einer direkten *Move-to-Front* Kodierung von s . Wie groß ist jeweils die Entropie der beiden kodierten Zeichenketten?
- e) Führen Sie eine inverse *Burrows-Wheeler-Transformation* auf Zeichenkette $s^{BWT} = \text{bc}\$\text{aab}$ aus. Erzeugen und verwenden Sie hierfür das LF-Array der Zeichenkette (siehe Übung).

Musterlösung:

a)

$$\begin{aligned}
 H(s) &= p(a) \log_2(1/p(a)) + p(b) \log_2(1/p(b)) \\
 &= 1/2 \log_2(2) + 1/2 \log_2(2) \\
 &= 1
 \end{aligned}$$

b) Vorgehen:

- schreibe zyklische Permutationen von s als Spalten auf (Endzeichen \$ nicht vergessen!)
- sortiere die Permutationen/Spalten lexikographisch (\$ ist kleiner als alle anderen Zeichen!)
- lese das Ergebnis s^{BWT} aus der letzten Zeile ab

Ausführung:

ababababab\$	\$aaaaabbbbb		
babababab\$a	abbbbb\$aaaa		
abababab\$ab	b\$aaaaabbbb		
bababab\$aba	aabbbbb\$aaa		
ababab\$abab	bb\$aaaaabbb		
babab\$ababa	→ aaabbbbb\$aa	→ $s^{BWT} =$	bbbbbb\$aaaaa
abab\$ababab	bbb\$aaaaabb		
bab\$abababa	aaaabbbbb\$a		
ab\$abababab	bbbb\$aaaaab		
b\$ababababa	aaaaabbbbb\$		
\$ababababab	bbbbbb\$aaaaa		

Da die *Burrows-Wheeler-Transformation* lediglich die Zeichen der Zeichenkette permutiert, ändert sich die Entropie nicht.

c) Die kodierte Zeichenkette zu s^{BWT} lautet: $c^{BWT} = 2, 1, 1, 1, 1, 2, 1, 1, 1, 1$
 Sie baut sich wie folgt auf:

Index i	Y	$s^{BWT}[i]$	$c^{BWT}[i]$
1	ab	b	2
2	ba	b	1
3	ba	b	1
4	ba	b	1
5	ba	b	1
6	ba	a	2
7	ab	a	1
8	ab	a	1
9	ab	a	1
10	ab	a	1

Entropie von c^{BWT} :

$$\begin{aligned}
 H(s) &= p(1) \log_2(1/p(1)) + p(2) \log_2(1/p(2)) \\
 &= 8/10 \log_2(10/8) + 2/10 \log_2(10/2) \\
 &\approx 0.722
 \end{aligned}$$

Musterlösung:

- d) Die kodierte Zeichenkette zu s lautet: $c = 1, 2, 2, 2, 2, 2, 2, 2, 2, 2$
 Sie baut sich wie folgt auf:

Index i	Y	$s^{BWT}[i]$	$c^{BWT}[i]$
1	ab	a	1
2	ab	b	2
3	ba	a	2
4	ab	b	2
5	ba	a	2
6	ab	b	2
7	ba	a	2
8	ab	b	2
9	ba	a	2
10	ab	b	2

Entropie von c^{BWT} :

$$\begin{aligned}
 H(s) &= p(1) \log_2(1/p(1)) + p(2) \log_2(1/p(2)) \\
 &= 1/10 \log_2(10) + 9/10 \log_2(10/9) \\
 &\approx 0.4670
 \end{aligned}$$

Erstaunlicherweise ist die *Move-to-Front* Kodierung von die Zeichenkette s erfolgreicher als von Zeichenkette s^{BWT} nach *Burrows-Wheeler-Transformation*. Dies kann bei einigen Spezialfällen auftreten. Eine anschließende Komprimierung von c (z.B. mit *Huffman-Kodierung*) würde c stärker verkleinern können als c^{BWT} .

- e) Vorgehen:

- Erzeuge das LF-Array von s auf Basis von s^{BWT}
- Verwende LF und s^{BWT} zum Umkehren der BWT nach dem Algorithmus aus der Übung

Sei M die Zeichen-Matrix, welche beim Erzeugen der BWT von s entsteht, also die sortierte Konkatenation der zyklischen Permutationen von s als Spaltenvektoren. Das LF-Array von s bildet die letzte Zeile L von M auf die erste Zeile F ab.

Die letzte Zeile L ist gerade die BWT $L = s^{BWT} = \text{bc\$aab}$ von s .

Die erste Zeile F enthält alle Zeichen von s in sortierter Reihenfolge, also $F = \text{\$aabb}$.

Also:

	1	2	3	4	5	6
L	b	c	\\$	a	a	b
LF	4	6	1	2	3	5

Dann gilt $s[n] = \$$ und $s[n - i] = L[\underbrace{LF(LF(\dots(LF(1))))}_{i-1 \text{ mal}}]$. Also:

- $s[6] = \$$ und $k = 1$
- $s[5] = L[1] = b$ und $k = LF(1) = 4$
- $s[4] = L[4] = a$ und $k = LF(4) = 2$
- $s[3] = L[2] = c$ und $k = LF(2) = 6$
- $s[2] = L[6] = b$ und $k = LF(6) = 5$
- $s[1] = L[5] = a$ und $k = LF(5) = 3$

Die gesuchte Zeichenfolge lautet $s = \text{abcab}$.

Aufgabe 3 (Rechnen: Suffixarrays und DC3-Algorithmus)

Gegeben sei die Zeichenkette $s = \text{aberakadabera}\$$.

- Geben Sie den Suffixbaum für s an.
- Geben Sie das Suffixarray für s an.

In der Vorlesung haben Sie einen Linearzeitalgorithmus zur Konstruktion von Suffixarrays kennengelernt. Dieser ist unter dem Namen *DC3-Algorithmus* bekannt. Im Folgenden soll der Algorithmus Schritt für Schritt per Hand ausgeführt werden.

Die Suffixe von s werden zunächst in drei Sequenzen $B_i = \langle s_k \mid (k \bmod 3) = i \rangle$ für $i \in \{0, 1, 2\}$ aufgeteilt. Danach müssen die Sequenzen $C = B_0 \cup B_1$ und B_2 lexikographisch sortiert werden.

Sortierung von C :

- Geben Sie die Tripelsequenzen $R_k = \langle s[i..i+2] \mid (i \bmod 3) = k \rangle$ für $k \in \{0, 1\}$ an. Für $i \geq |s|$ gelte $s[i] = \$$ (Auffüllen mit zusätzlichen Abschlusszeichen).
- Bestimmen Sie den Rang der Tripel von $R = R_0 \circ R_1$. Sortieren Sie dazu die Tripel und entfernen mehrfache Vorkommnisse. Die Position eines Tripels in dieser Sortierung gibt seinen Rang an.
- Die berechneten Ränge definieren eine eindeutige Bezeichnung für jedes Tripel in R . Drücken Sie R mit Hilfe dieser Ränge aus. Diese Darstellung ergibt die Zeichenkette s^{01} . Muss der DC3-Algorithmus eine Rekursion ausführen?
- Geben Sie das Suffixarray SA^{01} für s^{01} von Hand an (unabhängig, ob der DC3-Algorithmus eine Rekursion durchführt). Vergewissern Sie sich, dass SA^{01} eine Sortierung von C beschreibt.

Sortierung von C_2 :

- Erstellen Sie eine Zuordnung rank , die jedem i mit $s_i \in C$ den Index von s_i in der sortierten Sequenz C zuweist. Für alle anderen i sei $\text{rank}(i) = 0$.

Formale Berechnung von rank mit Hilfe des Suffixarray SA^{01} nach:

$$\begin{aligned} \text{rank}[3 \cdot (\text{SA}^{01}[i])] &= i \quad \text{falls } \text{SA}^{01}[i] < |B_0| \\ \text{rank}[3 \cdot (\text{SA}^{01}[i] - |B_0|) + 1] &= i \quad \text{falls } \text{SA}^{01}[i] \geq |B_0| \end{aligned}$$

Alle anderen Werte von $\text{rank}[i]$ können gleich 0 gesetzt werden.

- Erstellen Sie Tupel $(s[i], \text{rank}[i+1])$ f.a. $s_i \in B_2$ und sortieren diese lexikographisch. Vergewissern Sie sich, dass diese Sortierung einer Sortierung von B_2 entspricht.

Nachdem C und B_2 sortiert worden sind, kann das Suffixarray von s bestimmt werden:

- Führen Sie eine Mischen-Operation auf C und B_2 aus. Die resultierende Sequenz wird mit D bezeichnet. Sei $s_i \in C$ und $s_j \in B_2$. Es gelten folgende Sortierkriterien:

$$s_i \leq s_j \iff \begin{cases} (s[i], \text{rank}[i+1]) \leq (s[j], \text{rank}[j+1]) & \text{falls } s_i \in B_0 \\ (s[i], s[i+1], \text{rank}[i+2]) \leq (s[j], s[j+1], \text{rank}[j+2]) & \text{falls } s_i \in B_1 \end{cases}$$

Vergewissern Sie sich, dass D lexikographisch sortiert ist und damit das Suffixarray induziert.

Notation:

- Alle Indizes fangen bei 0 an – analog zu Kapitel 9.3.6, auf dem die Aufgabe basiert.

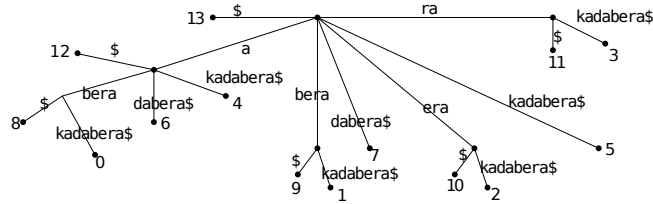
- $s[i..j]$: Zeichen an Stelle i (bis j) in s (z.B. $s[1..3] = \text{ber}$)
- s_i : Suffix von s ab Stelle i (z.B. $s_2 = \text{erakadabera}$)

Musterlösung:

Als Referenz zunächst Zeichenkette s mit ihren Indizes (beachten Sie das Abschlusszeichen \$!):

Index i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
$s[i]$	a	b	e	r	a	k	a	d	a	b	e	r	a	\$

a) Der Suffixbaum für s als kompakterer Trie:



Jeder Weg von der Wurzel zu einem Blatt gibt einen Suffix von s an. Der Wert am Blatt gibt den Index des Suffix an (d.h. die Stelle in der Zeichenkette an der der Suffix beginnt).

b) Das Suffixarray SA für s lautet:
 $SA = \langle 13, 12, 8, 0, 6, 4, 9, 1, 7, 10, 2, 5, 11, 3 \rangle$.

Index i	$SA[i]$	$s_{SA[i]}$
0	13	\$
1	12	a\$
2	8	abera\$
3	0	aberakadabera\$
4	6	adabera\$
5	4	akadabera\$
6	9	bera\$
7	1	berakadabera\$
8	7	dabera\$
9	10	era\$
10	2	erakadabera\$
11	5	kadabera\$
12	11	ra\$
13	3	rakadabera\$

In der rechten Spalte der Suffixtabelle sind die durch das Suffixarray indizierten Suffixe aufgetragen. Man sieht, dass sie durch das Suffixarray in alphabetischer Reihenfolge geordnet vorliegen.

Musterlösung:

c) Die Tripelsequenzen lauten:

$$R_0 = \langle \text{abe, rak, ada, ber, a}\$\$ \rangle$$

$$R_1 = \langle \text{ber, aka, dab, era, } \$ \$ \$ \rangle$$

$$R = \langle \text{abe, rak, ada, ber, a}\$\$, \text{ber, aka, dab, era, } \$ \$ \$ \rangle$$

Beachten Sie, dass das jeweils letzte Tripel mit weiteren Abschlusszeichen \$ ergänzt wurde.

d) Die sortierten Tripel von R ohne Duplikate ergeben folgende Ränge:

Index i	9	4	0	2	6	3, 5	7	8	1
$R[i]$	\$\$\$	a\$\$	abe	ada	aka	ber	dab	era	rak
Rang	0	1	2	3	4	5	6	7	8

e) Die Ränge definieren eine eindeutige Bezeichnung der Tripel von R . Damit kann R dargestellt werden als $s^{01} = \langle 2, 8, 3, 5, 1, 5, 4, 6, 7, 0 \rangle$.

Die Sequenz s^{01} enthält das selbe Zeichen (5) mehrfach. Ansonsten wäre über die Ränge bereits eine vollständige Sortierung von s^{01} –und damit auch von C – bestimmt. Um diese Sortierung zu erhalten, bestimmt man rekursiv mit dem DC3-Algorithmus das Suffixarray für s^{01} .

f) Das Suffixarray für R ist –nach Konstruktion– gleich dem Suffixarray für s^{01} . Es ergibt sich zu $SA^{01} = \langle 10, 9, 4, 0, 2, 6, 3, 5, 7, 8, 1 \rangle$.

Index i	$SA^{01}[i]$	$s_{SA^{01}[i]}^{01}$	$R_{SA^{01}[i]}$
0	10	$\langle . \rangle$	$\langle \rangle$
1	9	$\langle 0. \rangle$	$\langle \$ \$ \$ \rangle$
2	4	$\langle 1, 5, 4, 6, 7, 0, . \rangle$	$\langle \text{a}\$\$, \text{ber, aka, dab, era, } \$ \$ \$, \rangle$
3	0	$\langle 2, 8, 3, 5, 1, 5, 4, 6, 7, 0, . \rangle$	$\langle \text{abe, rak, ada, ber, a}\$\$, \text{ber, aka, dab, era, } \$ \$ \$, \rangle$
4	2	$\langle 3, 5, 1, 5, 4, 6, 7, 0, . \rangle$	$\langle \text{ada, ber, a}\$\$, \text{ber, aka, dab, era, } \$ \$ \$, \rangle$
5	6	$\langle 4, 6, 7, 0, . \rangle$	$\langle \text{aka, dab, era, } \$ \$ \$, \rangle$
6	3	$\langle 5, 1, 5, 4, 6, 7, 0, . \rangle$	$\langle \text{ber, a}\$\$, \text{ber, aka, dab, era, } \$ \$ \$, \rangle$
7	5	$\langle 5, 4, 6, 7, 0, . \rangle$	$\langle \text{ber, aka, dab, era, } \$ \$ \$, \rangle$
8	7	$\langle 6, 7, 0, . \rangle$	$\langle \text{dab, era, } \$ \$ \$, \rangle$
9	8	$\langle 7, 0, . \rangle$	$\langle \text{era, } \$ \$ \$, \rangle$
10	1	$\langle 8, 3, 5, 1, 5, 4, 6, 7, 0, . \rangle$	$\langle \text{rak, ada, ber, a}\$\$, \text{ber, aka, dab, era, } \$ \$ \$, \rangle$

Beachten Sie das zusätzliche Abschlusszeichen . für s^{01} . Es ist funktional identisch zum \$ für s . Zur leichteren Unterscheidung wurde aber ein anderes Zeichen gewählt. Das Abschlusszeichen benötigt keine Entsprechung in R , daher ist hier nur ein leeres Tripel eingetragen. Es wird für die Bestimmung von SA^{01} benötigt, kann aber im weiteren Verlauf ignoriert werden.

Das Suffixarray gibt eine Sortierung der Elemente von R bzw. s^{01} und damit auch von C an (für C ist dies spätestens dann ersichtlich, wenn man in der rechten Spalte alle Tripel nach a \$\$ entfernt und die restlichen in einer Zeile konkateniert).

g) Aus dem Suffixarray lässt sich folgende Rang-Funktion ableiten:

Index i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	...
$s[i]$	a	b	e	r	a	k	a	d	a	b	e	r	a	\$	\$...
$rank[i]$	3	7	⊥	10	5	⊥	4	8	⊥	6	9	⊥	2	1	⊥	...

Beachten Sie, dass $rank[i] = \perp$ anstatt 0 gesetzt wurde f.a. $(i \bmod 3) = 0$. Dies ist zulässig, da diese Einträge von $rank$ nie verwendet werden.

h) Es ergeben sich die folgenden Tupel (für $s_i \in B_2$ gilt $(i \bmod 3) = 2, i < |s|$):

Index i	2	5	8	11
$(s[i], rank[i + 1])$	(e, 10)	(k, 4)	(a, 6)	(r, 2)

Sortiert ergibt sich: $\langle (\mathbf{a}, 6), (\mathbf{e}, 10), (\mathbf{k}, 4), (\mathbf{r}, 2) \rangle$. Diese Sortierung entspricht einer Sortierung von B_2 .

Musterlösung:

- i) Für das Mischen werden C und B_2 in sortierter Reihenfolge benötigt. Diese Reihenfolge wurde in den vorherigen Teilaufgaben berechnet. Für den Vergleich beim Mischen wird zusätzlich für jedes Suffix eine bestimmte Tupeldarstellung benötigt.

Für B_2 ergeben sich die sortierte Reihenfolge und die benötigten Tupel wie folgt:

Index i	s_i	$(s[i], rank[i + 1])$	$(s[i], s[i + 1], rank[i + 2])$
8	abera\$ $\in B_2$	(a, 6)	(a, b, 9)
2	erakadabera\$ $\in B_2$	(e, 10)	(e, r, 5)
5	kadabera\$ $\in B_2$	(k, 4)	(k, a, 8)
11	ra\$ $\in B_2$	(r, 2)	(r, a, 1)

Für C sind die sortierte Reihenfolge (z.B. aus $rank$ abzulesen) und die benötigten Tupel:

Index i	s_i	$(s[i], rank[i + 1])$	$(s[i], s[i + 1], rank[i + 2])$
13	\$ $\in B_1$		(\$, \$, 0)
12	a\$ $\in B_0$	(a, 1)	
0	abera\$ $\in B_0$	(a, 7)	
6	adabera\$ $\in B_0$	(a, 8)	
4	akadabera\$ $\in B_1$		(a, k, 4)
9	bera\$ $\in B_0$	(b, 9)	
1	berakadabera\$ $\in B_1$		(b, e, 10)
7	dabera\$ $\in B_1$		(d, a, 6)
10	era\$ $\in B_1$		(e, r, 2)
3	rakadabera\$ $\in B_0$	(r, 5)	

Zusammengemischt ergibt sich D zu:

Index i	s_i	$(s[i], rank[i + 1])$	$(s[i], s[i + 1], rank[i + 2])$
13	\$ $\in B_1$		(\$, \$, 0)
12	a\$ $\in B_0$	(a, 1)	
8	abera\$ $\in B_2$	(a, 6)	(a, b, 9)
0	abera\$ $\in B_0$	(a, 7)	
6	adabera\$ $\in B_0$	(a, 8)	
4	akadabera\$ $\in B_1$		(a, k, 4)
9	bera\$ $\in B_0$	(b, 9)	
1	berakadabera\$ $\in B_1$		(b, e, 10)
7	dabera\$ $\in B_1$		(d, a, 6)
10	era\$ $\in B_1$		(e, r, 2)
2	erakadabera\$ $\in B_2$	(e, 10)	(e, r, 5)
5	kadabera\$ $\in B_2$	(k, 4)	(k, a, 8)
11	ra\$ $\in B_2$	(r, 2)	(r, a, 1)
3	rakadabera\$ $\in B_0$	(r, 5)	

Die Suffixe in D sind offensichtlich lexikographisch sortiert (zweite Spalte). Damit erhält man das Suffixarray für s als $SA = \langle 13, 12, 8, 0, 6, 4, 9, 1, 7, 10, 2, 5, 11, 3 \rangle$.

Aufgabe 4 (*Rechnen+Analyse: LCP-Array*)

Gegeben sei die Zeichenkette $s = \text{salsadipp\$}$.

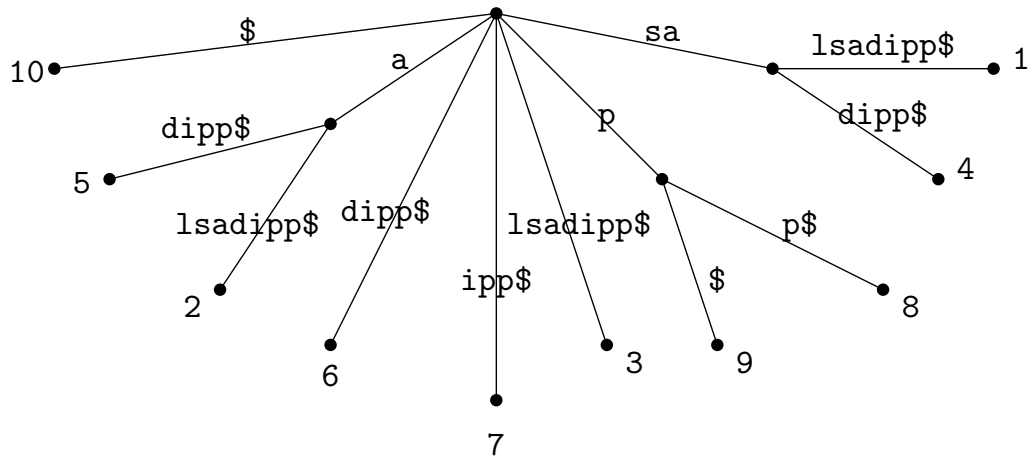
- a) Geben Sie den Suffixbaum für s an.
- b) Geben Sie das Suffixarray für s an.
- c) Geben Sie das LCP-Array für s an.

Im Folgenden sei ein String T sowie dessen Suffixarray $\text{SA}[\cdot]$ und dessen LCP-Array $\text{LCP}[\cdot]$ gegeben.

- d) Wie kann der längste sich wiederholende Substring in T effizient bestimmt werden?
(der Substring darf sich dabei selbst überlappen)
- e) Wie viele paarweise unterschiedliche Substrings kann ein String der Länge n maximal besitzen?
Wie kann die tatsächliche Anzahl für einen konkreten String T bestimmt werden?
- f) Ein String lässt sich schlecht komprimieren, wenn er wenig Redundanz besitzt. Ein Maß dafür ist die Anzahl paarweise unterschiedlicher Substrings normiert auf die mögliche Gesamtanzahl unterschiedlicher Substrings. Geben Sie an, wie dieses Maß für T berechnet werden kann.

Musterlösung:

a) Der Suffixbaum für s als kompaktierter Trie:



Jeder Weg von der Wurzel zu einem Blatt gibt einen Suffix von s an. Der Wert am Blatt gibt den Index des Suffix an (d.h. die Stelle in der Zeichenkette an der der Suffix beginnt).

b) Das Suffixarray SA für s lautet:
 $SA = \langle 10, 5, 2, 6, 7, 3, 9, 8, 4, 1 \rangle$.

Index i	$SA[i]$	$s_{SA[i]}$
1	10	\$
2	5	adipp\$
3	2	alsadipp\$
4	6	dipp\$
5	7	ipp\$
6	3	lsadipp\$
7	9	p\$
8	8	pp\$
9	4	sadipp\$
10	1	salsadipp\$

In der rechten Spalte der Suffixtabelle sind die durch das Suffixarray indizierten Suffixe aufgetragen. Man sieht, dass sie durch das Suffixarray in alphabetischer Reihenfolge geordnet vorliegen.

c) Das LCP-Array LCP für s lautet:
 $LCP = \langle 0, 0, 1, 0, 0, 0, 0, 1, 0, 2 \rangle$.

Der Eintrag $LCP[i]$ gibt die Länge des gemeinsamen Präfixes von $s_{SA[i-1]}$ und $s_{SA[i]}$ an. Damit ist $LCP[1]$ nicht definiert. Nach Konvention setzt man normalerweise $LCP[1] = 0$.

Musterlösung:

d) Das Suffixarray $SA[\cdot]$ enthält alle Suffixe von T lexikographisch sortiert. Damit sind die Einträge mit dem längstem gemeinsamen Präfix –und damit mit dem längsten Substring– benachbart. Da jeder Eintrag $LCP[i]$ im LCP-Array die Länge des gemeinsamen Präfixes von benachbarten Suffixen $s_{SA[i-1]}$ und $s_{SA[i]}$ angibt, genügt es, das Maximum von $LCP[\cdot]$ zu bestimmen, um die Position des längsten sich wiederholenden Substring zu erhalten.

e) Die Menge aller Substrings eines Strings ist durch die Menge aller Präfixe seiner Suffixe gegeben. In einem String mit maximal vielen unterschiedlichen Substrings besitzen keine zwei Suffixe ein gemeinsames Präfix. Damit steuert jedes Suffix s genau $|s|$ zur Menge der paarweise unterschiedlichen Substrings bei. Somit ist die gesuchte Anzahl $\sum_{i=1}^n |s_{SA[i]}| = \sum_{i=1}^n i = \frac{n \cdot (n+1)}{2}$.

Die Anzahl paarweiser unterschiedlicher Substrings für einen konkreten String T ist gegeben durch $\#_{pds} = \sum_{i=1}^n |s_{SA[i]}| - LCP[i]$. Dies folgt aus folgender Überlegung:

Sei $LCP[j] = k$. Dies bedeutet, dass $s_{SA[i-1]}$ und $s_{SA[i]}$ ein gemeinsames Präfix der Länge k besitzen. Das wiederum heißt, dass die k Suffixe dieses gemeinsamen Präfix nicht gezählt werden dürfen, da sie gemeinsam und nicht paarweise verschieden sind.

Die Länge eines Suffix lässt sich mit Hilfe des Suffix-Arrays bestimmen zu $|s_{SA[i]}| = n - SA[i]$.

f) Dies lässt sich direkt aus der vorherigen Teilaufgabe ableiten: $\frac{2}{n \cdot (n+1)} \sum_{i=1}^n P[i]$.

Aufgabe 5 (*RMQ in Wavelet Trees*)

Gegeben sei ein Universum von Zahlen \mathcal{U} und ein Feld A von Zahlen aus \mathcal{U} . Geben sie einen Algorithmus an, mit dem sich unter Benutzung eines Wavelet Trees in $\log |\mathcal{U}|$ Zeit $\arg \min_i \{A[i] \mid i \in [a, b]\}$ für Parameter a, b berechnen lässt.

Musterlösung:

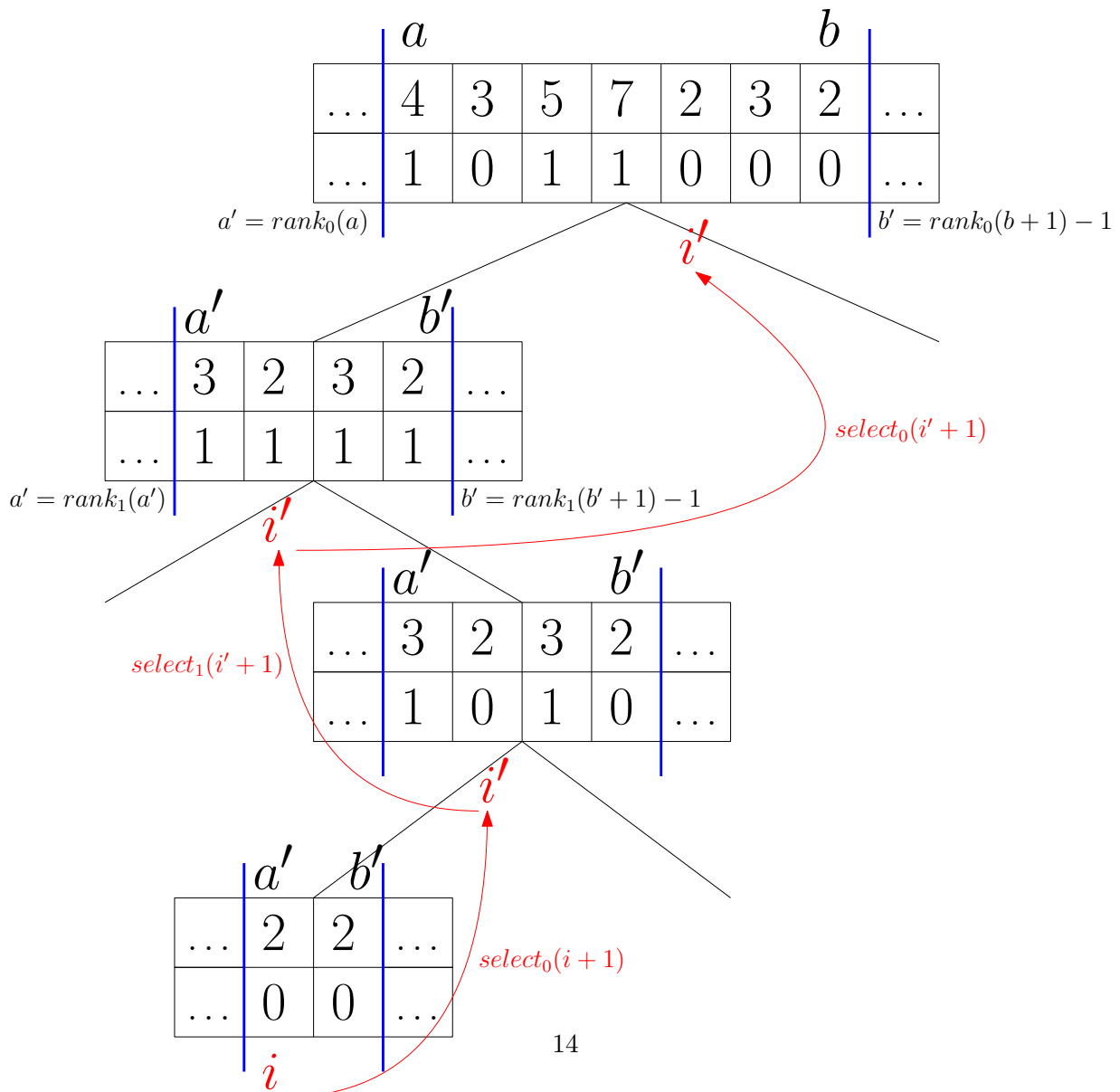
Sei WT der Wavelet Tree für unser Feld A .

```

1: function RMQ( $WT, a, b$ )
2:   if  $WT$  ist Blatt then
3:     return  $a$ 
4:   end if
5:    $a' \leftarrow rank_0(a)$ 
6:    $b' \leftarrow rank_0(b + 1) - 1$ 
7:   if  $b' - a' \geq 0$  then                                     ▷ Minimum ist in linkem Teilbaum
8:      $i' \leftarrow RMQ(WT \rightarrow child_{left}, a', b')$ 
9:      $i \leftarrow select_0(i' + 1)$ 
10:  else                                                         ▷ Minimum ist in rechtem Teilbaum
11:     $i' \leftarrow RMQ(WT \rightarrow child_{right}, a - a', b - b' - 1)$ 
12:     $i \leftarrow select_1(i' + 1)$ 
13:  end if
14:  return  $i$ 
15: end function

```

Dieser Algorithmus ist leicht modifizierbar um das linkeste, rechteste oder mittlere Minimum auszugeben.



Aufgabe 6 (Rechnen: Kompression)

- a) Bestimmen Sie die LZ77-Faktoren $f_i = (l_i, p_i)$ des Textes $T_1 = \text{abbababbbb}\$$.
- b) Gegeben seien folgende LZ77-Faktoren. Bestimmen Sie den Text T_2 , aus dem sie entstanden sind.

$(0, \mathbf{a}), (5, 1), (0, \mathbf{b}), (6, 1), (7, 7), (0, \$)$

- c) Angenommen, jedes Zeichen kann in 1 byte, sowie jeder LZ77-Faktor in 2 byte repräsentiert werden. Wie viel Prozent des Speicherplatzes von T_2 konnte somit eingespart werden?

Musterlösung:

- a) Die Lempel-Ziv-Faktoren lauten:

f_1	a	(0, a)
f_2	b	(0, b)
f_3	bb	(2, 2)
f_4	ab	(2, 1)
f_5	abbb	(4, 1)
f_6	b	(1, 2)
f_7	\$	(0, \$)

- b) $T_2 = \text{aaaaaabaaaaabaaaaaa}\$$

f_1	a	(0, a)
f_2	aaaaa	(5, 1)
f_3	b	(0, b)
f_4	aaaaaa	(6, 1)
f_5	baaaaaa	(7, 7)
f_6	\$	(0, \$)

- c) $|T_1| = 21$ benötigt 21 byte. $|f_i| = 6$ benötigt 12 byte. Dadurch ergibt sich eine Einsparung von $1 - 12/21 \approx 42\%$.