

Algorithmen auf Sequenzen

Volltext-Indexdatenstrukturen: BWT, FM-Index

Sven Rahmann

Genominformatik
Universitätsklinikum Essen
Universität Duisburg-Essen
Universitätsallianz Ruhr

Die Burrows-Wheeler-Transformation (BWT)

- BWT wurde von Michael Burrows und David Wheeler 1994 zur verlustfreien Kompression von Daten entwickelt, basierend auf Vorarbeiten von 1983.
- Lange Zeit wurde die Bedeutung der BWT nicht vollständig erkannt.
- Heute: Basis für Kompression, Datenstrukturen, Algorithmen

Die Burrows-Wheeler-Transformation (BWT)

- BWT wurde von Michael Burrows und David Wheeler 1994 zur verlustfreien Kompression von Daten entwickelt, basierend auf Vorarbeiten von 1983.
- Lange Zeit wurde die Bedeutung der BWT nicht vollständig erkannt.
- Heute: Basis für Kompression, Datenstrukturen, Algorithmen
- Die BWT zu einem Text T oder $T\$$ ist stets eine Permutation der Zeichen des ursprünglichen Texts.
- informell: Folge der Zeichen **vor** lexikographisch sortierten Suffixen
- Verschiedene Varianten (z.B. zyklisch oder mit Wächter)
(Zyklisch: Nach Textende, betrachte Zeichen am Anfang des Texts)

Definition (BWT)

Gegeben T mit $|T| = n$,
bilde die $n \times n$ Matrix M , deren Zeilen zyklische Verschiebungen des Texts sind.
Sortiere die Zeilen lexikographisch. Die BWT entspricht der letzten Spalte.

Die zyklische Burrows-Wheeler Transformation

Definition (BWT)

Gegeben T mit $|T| = n$,
bilde die $n \times n$ Matrix M , deren Zeilen zyklische Verschiebungen des Texts sind.
Sortiere die Zeilen lexikographisch. Die BWT entspricht der letzten Spalte.

	F						L
	\$	b	a	n	a	n	a
	a	\$	b	a	n	a	n
	a	n	a	\$	b	a	n
	a	n	a	n	a	\$	b
	b	a	n	a	n	a	\$
	n	a	\$	b	a	n	a
	n	a	n	a	\$	b	a

(in der Regel angewendet auf Texte ohne Wächter)

Definition (BWT)

Sei Text $T\$$ mit $n := |T\$|$ gegeben und pos das Suffixarray.

Dann ist die BWT die Abbildung $r \mapsto bwt[r] := T[pos[r] - 1]$ mit $T[-1] := \$$.

Die BWT mit Wächter oder Stringende

Definition (BWT)

Sei Text $T\$$ mit $n := |T\$|$ gegeben und pos das Suffixarray.

Dann ist die BWT die Abbildung $r \mapsto bwt[r] := T[pos[r] - 1]$ mit $T[-1] := \$$.

r	$pos[r]$	$bwt[r]$	$T[pos[r] :]$
0	6	a	\$
1	5	n	a\$
2	3	n	ana\$
3	1	b	anana\$
4	0	\$	banana\$
5	4	a	na\$
6	2	a	nana\$

- Da das Suffixarray in Linearzeit konstruiert werden kann, beträgt die Laufzeit für die Konstruktion des BWT auch $\mathcal{O}(n)$.
- Bei natürlichsprachlichen Texten enthält die BWT häufig lange Läufe desselben Zeichens, was sich gut komprimieren lässt.
- Das k -te Auftreten eines Zeichens in der BWT entspricht dem k -ten Auftreten des Zeichens im Suffixarray.
- Die BWT lässt sich ohne zusätzliche Informationen wieder in ihren ursprünglichen Text rücktransformieren.

Eigenschaften der BWT

<i>bwt</i>	<i>Sortierter Text</i>
d	ie BWT lässt sich ohne
d	ie Indizes der sortierten
d	ie LMS-Substrings bilden
d	ie hinterste freie Stelle
d	ie selbe Anordnung
d	ie selbe Länge
d	ie sortierte Reihenfolge der
d	iese ursprüngliche Definition
d	ieser Ansatz hat eine

- Bei natürlichsprachlichen Texten kann ausgenutzt werden, dass bestimmte Buchstabenkombinationen, wie 'die', 'sch', 'tz' oder 'ie' häufig vorkommen.
- Dementsprechend beinhaltet eine BWT oft lange Blöcke von gleichen Zeichen, z.B. ... ddddddddddddddd...
- Hierauf lässt sich *run-length-encoding* (RLE) sinnvoll anwenden, um solche Buchstabenblöcke mit wenig Bits darzustellen.

Das Tool bzip2 von Julian Seward basiert auf der BWT. Als Eingabe bekommt bzip2 einen Text und eine Blockgröße und arbeitet wie folgt:

Für jeden Block separat:

- 1 Berechne die BWT.
- 2 Wende auf jede BWT die Move-to-front Transformation an. Die so entstehende Sequenz beinhaltet viele Nullen und kleine Werte.
- 3 Spezielles Run-length encoding (RLE) der Sequenz
- 4 Huffman-Codierung der Sequenz: Häufige Zeichen werden durch kurze Bitfolgen kodiert, seltene Zeichen durch lange Bitfolgen.

Alle Transformationsschritte sind invertierbar.

Lemma

Das k -te Auftreten des Zeichens $c \in \Sigma$ in der BWT entspricht dem k -ten Auftreten des Zeichens im Suffixarray.

Lemma

Das k -te Auftreten des Zeichens $c \in \Sigma$ in der BWT entspricht dem k -ten Auftreten des Zeichens im Suffixarray.

r	$pos[r]$	$bwt[r]$	$T[pos[r] :]$
0	6	a	\$
1	5	n	a\$
2	3	n	ana\$
3	1	b	anana\$
4	0	\$	banana\$
5	4	a	na\$
6	2	a	nana\$

Lemma

Das k -te Auftreten des Zeichens $c \in \Sigma$ in der BWT entspricht dem k -ten Auftreten des Zeichens im Suffixarray.

r	$pos[r]$	$bwt[r]$	$T[pos[r] :]$
0	6	a	\$
1	5	n	a\$
2	3	n	ana\$
3	1	b	anana\$
4	0	\$	banana\$
5	4	a	na\$
6	2	a	nana\$

Betrachte zwei Ränge $r_1 < r_2$ mit $bwt[r_1] = bwt[r_2]$. Da die Suffixe $pos[r_1]$ und $pos[r_2]$ lexikographisch sortiert sind, die jeweils auf die BWT-Zeichen folgen, entspricht die Reihenfolge der BWT-Zeichen der im Suffixarray.

Stabile Sortierung (mit Speichern der r -Werte) der BWT liefert Nachfolger zu jedem Zeichen und erlaubt Rekonstruktion des Texts.

r	0	1	2	3	4	5	6
$bwt[r]$	a	n	n	b	\$	a	a
sortiert	\$	a	a	a	b	n	n
$orig_r[r]$	4	0	5	6	3	1	2

1) Der Wächter steht bei $r = 4$. Sein Nachfolger ist b . Der Originaltext beginnt mit diesem Zeichen.

$$T = b$$

Invertierbarkeit der BWT

Stabile Sortierung (mit Speichern der r -Werte) der BWT liefert Nachfolger zu jedem Zeichen und erlaubt Rekonstruktion des Texts.

r	0	1	2	3	4	5	6
$bwt[r]$	a	n	n	b	\$	a	a
sortiert	\$	a	a	a	b	n	n
$orig_r[r]$	4	0	5	6	3	1	2

2) Das b stand bei $r = 3$. Der Nachfolger von b ist ein a . Das Zeichen wird zum Text hinzugefügt.

$$T = ba$$

Invertierbarkeit der BWT

Stabile Sortierung (mit Speichern der r -Werte) der BWT liefert Nachfolger zu jedem Zeichen und erlaubt Rekonstruktion des Texts.

r	0	1	2	3	4	5	6
$bwt[r]$	a	n	n	b	\$	a	a
sortiert	\$	a	a	a	b	n	n
$orig_r[r]$	4	0	5	6	3	1	2

3 - n) Es wird insgesamt n mal die Position aktualisiert und der Nachfolger zum Text hinzugefügt.

$$T = \text{ban}$$

Invertierbarkeit der BWT

Stabile Sortierung (mit Speichern der r -Werte) der BWT liefert Nachfolger zu jedem Zeichen und erlaubt Rekonstruktion des Texts.

r	0	1	2	3	4	5	6
$bwt[r]$	a	n	n	b	\$	a	a
sortiert	\$	a	a	a	b	n	n
$orig_r[r]$	4	0	5	6	3	1	2

3 - n) Es wird insgesamt n mal die Position aktualisiert und der Nachfolger zum Text hinzugefügt.

$$T = \text{bana}$$

Invertierbarkeit der BWT

Stabile Sortierung (mit Speichern der r -Werte) der BWT liefert Nachfolger zu jedem Zeichen und erlaubt Rekonstruktion des Texts.

r	0	1	2	3	4	5	6
$bwt[r]$	a	n	n	b	\$	a	a
sortiert	\$	a	a	a	b	n	n
$orig_r[r]$	4	0	5	6	3	1	2

3 - n) Es wird insgesamt n mal die Position aktualisiert und der Nachfolger zum Text hinzugefügt.

$$T = \text{banan}$$

Invertierbarkeit der BWT

Stabile Sortierung (mit Speichern der r -Werte) der BWT liefert Nachfolger zu jedem Zeichen und erlaubt Rekonstruktion des Texts.

r	0	1	2	3	4	5	6
$bwt[r]$	a	n	n	b	\$	a	a
sortiert	\$	a	a	a	b	n	n
$orig_r[r]$	4	0	5	6	3	1	2

3 - n) Es wird insgesamt n mal die Position aktualisiert und der Nachfolger zum Text hinzugefügt.

$$T = \text{banana}$$

Invertierbarkeit der BWT

Stabile Sortierung (mit Speichern der r -Werte) der BWT liefert Nachfolger zu jedem Zeichen und erlaubt Rekonstruktion des Texts.

r	0	1	2	3	4	5	6
$bwt[r]$	a	n	n	b	\$	a	a
sortiert	\$	a	a	a	b	n	n
$orig_r[r]$	4	0	5	6	3	1	2

3 - n) Es wird insgesamt n mal die Position aktualisiert und der Nachfolger zum Text hinzugefügt.

$$T = \text{banana\$}$$

Sortierung in $\mathcal{O}(n)$ Zeit mit bucket sort, $|\Sigma| = \mathcal{O}(1)$

Exakte Mustersuche

- Bei der Mustersuche mit Hilfe der BWT sucht man die Intervallgrenzen $[L, R]$ für den Index r des Suffixarrays pos , so dass alle Vorkommen des Musters an $pos[r]$ mit $L \leq r \leq R$ beginnen.
- **Backward Search**: Muster wird zeichenweise rückwärts bearbeitet.
- Backward Search simuliert **nicht** die Suche im Suffixbaum.

Exakte Mustersuche

- Bei der Mustersuche mit Hilfe der BWT sucht man die Intervallgrenzen $[L, R]$ für den Index r des Suffixarrays pos , so dass alle Vorkommen des Musters an $pos[r]$ mit $L \leq r \leq R$ beginnen.
- **Backward Search**: Muster wird zeichenweise rückwärts bearbeitet.
- Backward Search simuliert **nicht** die Suche im Suffixbaum.
- Das Intervall $[L, R]$ ist wie folgt definiert:
 (\leq_m : lexikographisch-kleiner-gleich auf den ersten $m = |P|$ Zeichen)

$$L := \min\{r \mid T[pos[r] :] \leq_m P\}$$

$$R := \max\{r \mid P \leq_m T[pos[r] :]\}$$

- Es gibt $R - L + 1$ Treffer des Musters im Text.
- Im Fall $L > R$ müsste man P zwischen R und L einfügen.

Exakte Mustersuche

- Ist $P = \epsilon$, dann gilt $L = 0$ und $R = n - 1$ (Initialisierung)
- Zur Suche sind zwei Hilfstabellen nötig:
 - $less[c]$: Anzahl Zeichen im Text, die kleiner sind als c
 - $Occ[c, i]$: zählt die Vorkommen von Zeichen c im Präfix $bwt[..i]$.
- Beide Hilfstabellen zusammen werden als **FM-Index** bezeichnet (“Ferragina-Manzini” oder “full text in minute space”)

Exakte Mustersuche

- Ist $P = \epsilon$, dann gilt $L = 0$ und $R = n - 1$ (Initialisierung)
- Zur Suche sind zwei Hilfstabellen nötig:
 - $less[c]$: Anzahl Zeichen im Text, die kleiner sind als c
 - $Occ[c, i]$: zählt die Vorkommen von Zeichen c im Präfix $bwt[..i]$.
- Beide Hilfstabellen zusammen werden als **FM-Index** bezeichnet (“Ferragina-Manzini” oder “full text in minute space”)

Beispiel: $T = \text{banana}\$, \text{bwt} = \text{annb}\aa

less	\$	a	b	n	Occ	0	1	2	3	4	5	6
	0	1	4	5	\$	0	0	0	0	1	1	1
					a	1	1	1	1	1	2	3
					b	0	0	0	1	1	1	1
					n	0	1	2	2	2	2	2

Lemma

Sei $P^+ := cP$; $[L, R]$ das Intervall zu P und $[L^+, R^+]$ das gesuchte Intervall zu P^+ .

$$L^+ = \text{less}[c] + \text{Occ}[c, L - 1]$$

$$R^+ = \text{less}[c] + \text{Occ}[c, R] - 1$$

Lemma

Sei $P^+ := cP$; $[L, R]$ das Intervall zu P und $[L^+, R^+]$ das gesuchte Intervall zu P^+ .

$$L^+ = \text{less}[c] + \text{Occ}[c, L - 1]$$

$$R^+ = \text{less}[c] + \text{Occ}[c, R] - 1$$

- Zu Beginn ist $P = \epsilon$, somit ist das Muster im gesamten Intervall vorhanden:
 $L = 0$ und $R = n - 1$.
- Da $P^+ = c$, ist das Subintervall $[L_c, R_c]$ gesucht, in dem die Suffixe mit c beginnen.
- Das Intervall ist durch $L_c = \text{less}[c]$ und $R_c = \text{less}[c] + \text{Occ}[c, n - 1] - 1$ gegeben.
Das Lemma stimmt für $|P^+| = 1$.

Lemma

Sei $P^+ := cP$; $[L, R]$ das Intervall zu P und $[L^+, R^+]$ das gesuchte Intervall zu P^+ .

$$L^+ = \text{less}[c] + \text{Occ}[c, L - 1]$$

$$R^+ = \text{less}[c] + \text{Occ}[c, R] - 1$$

- Allgemein: $[L^+, R^+]$ ist Subintervall von $[L_c, R_c]$
- Innerhalb $[L, R]$ zählen, wie viele Vorgänger ein c sind.
- Da das k -te c in der BWT auch dem k -ten c im Suffixarray entspricht, muss an Stelle L gezählt werden, wie oft c im Präfix $\text{bwt}[L]$ vorkam, also $\text{Occ}[c, L - 1]$. Diese Zahl muss nach $\text{less}[c]$ übersprungen werden, also $L^+ = L_c + \text{Occ}[c, L - 1]$.
- Neue rechte Grenze: Zählen, wie oft c bis zur Position R vorkommt: $\text{Occ}[c, R]$.

Backward Search

0	1	2	3
012345678901234567890123456789012			
<i>T</i> = einsameeselessennassenesseln g ern\$			
<i>P</i> = less			

bwt[*r*]= nsnm\$ssssgenlneeearneleienesseae
T[*pos*[*r*]]= \$aaeeeeeeeeeeegillmnnnnnrsssssss
L= 0
R= 32

Backward Search

0	1	2	3
012345678901234567890123456789012			
<i>T</i> = einsameeselessennassenesseln g ern\$			
<i>P</i> = les s			

$bwt[r] =$ nsnm\$ssssgenlneeearneleienesseae
 $T[pos[r]] =$ \$aaeeeeeeeeegillmnnnnnrsssssss
 $L = 0$ $less[s] + Occ[s, -1]$
 $R = 32$ $less[s] + Occ[s, 32] - 1$

Backward Search

0	1	2	3
012345678901234567890123456789012			
$T =$ einsameeselessennassenesseln g ern\$			
$P =$ less			

$bwt[r] =$	nsnm\$ssssgenlneeearneleienesseae
$T[pos[r]] =$	\$aaeeeeeeeeeeegillmnnnnnrsssssss
$L =$	0 25 + $Occ[s, -1]$
$R =$	32 25 + $Occ[s, 32] - 1$

Backward Search

0	1	2	3
012345678901234567890123456789012			
$T =$ einsameeselessennassenesseln g ern $\$$			
$P =$ less			

$bwt[r] =$	n	s	n	m	§	s	s	s	s	g	e	n	l	n	e	e	e	a	r	n	e	l	e	i	e	n	e	s	s	e	a	e			
$T[pos[r]] =$	§	a	a	e	e	e	e	e	e	e	e	e	e	e	e	e	e	e	g	i	l	m	n	n	n	n	n	n	r	s	s	s	s	s	s
$L =$	0																		25	+	0														
$R =$	32																		25	+	$Occ[s, 32]$	-	1												

Backward Search

0 1 2 3
012345678901234567890123456789012
 $T =$ einsameeselessennassenesseln g ern $\$$
 $P =$ less

$bwt[r] =$ nsnm $\$$ ssssgenlneeearneleienesseae
 $T[pos[r]] =$ $\$$ aaeeeeeeeeeeegillmnnnnnrsssssss
 $L =$ 25
 $R =$ 32

Backward Search

0	1	2	3
012345678901234567890123456789012			
<i>T</i> = einsameeselessennassenesseln g ern\$			
<i>P</i> = le s s			

$bwt[r] =$ nsnm\$ssssgenlneeearneleienesseae
 $T[pos[r]] =$ \$aaeeeeeeeeegillmnnnnnr~~sssssss~~
 $L = 25$ ~~less[s]~~ + Occ[s, 24]
 $R = 32$ ~~less[s]~~ + Occ[s, 32] - 1

Backward Search

0	1	2	3
012345678901234567890123456789012			
$T =$ einsameeselessennassenesseln g ern\$			
$P =$ less			

$bwt[r] =$	n	s	n	m	\$	s	s	s	s	g	e	n	l	n	e	e	e	a	r	n	e	l	e	i	e	n	e	s	s	e	a	e				
$T[pos[r]] =$	\$	a	a	e	e	e	e	e	e	e	e	e	e	e	e	e	e	e	g	i	l	m	n	n	n	n	n	n	r	s	s	s	s	s	s	s
$L =$	25																										$25 + Occ[s, 24]$									
$R =$	32																										$25 + Occ[s, 32] - 1$									

Backward Search

0	1	2	3
012345678901234567890123456789012			
<i>T</i> = einsameeselessennassenesseln g ern\$			
<i>P</i> = less			

$bwt[r] =$ nsnm\$~~ssss~~genlneeearneleiene~~ss~~seae
 $T[pos[r]] =$ \$aaeeeeeeeeeeegillmnnnnnr~~sssssss~~
 $L =$ 25 25 + 5
 $R =$ 32 25 + ~~Occ[s, 32]~~ - 1

Backward Search

0	1	2	3																				
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	
<i>T</i> =																							
einsameeselessennassenesseln g ern\$																							
<i>P</i> =																							
le s s																							

bwt[*r*]= nsnm\$ssssgenlneeearneleienesseae

T[*pos*[*r*]]= \$aaeeeeeeeeeeegillmnnnnnnr~~sssssss~~

L= 30

R= 32

Backward Search

0	1	2	3
012345678901234567890123456789012			
$T =$ einsameeselessennassenesseln g ern\$			
$P =$ less			

$bwt[r] =$ nsnm\$ssssgenlneeearneleienesseae
 $T[pos[r]] =$ \$aaeeeeeeeeegillmnnnnnrsssss sss
 $L = 30$ $less[e] + Occ[e, 29]$
 $R = 32$ $less[e] + Occ[e, 32] - 1$

Backward Search

0	1	2	3																				
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	
<i>T</i> = einsameeselessennassenesseln g ern\$																							
<i>P</i> = le s s																							

<i>bwt</i> [<i>r</i>]=	nsnm\$ssssg e nlnee e arnele e ieness e ae
<i>T</i> [<i>pos</i> [<i>r</i>]]=	\$aaeeeeeeeeegillmnnnnnrsssss sss
<i>L</i> = 30	3 + <i>Occ</i> [e, 29]
<i>R</i> = 32	3 + <i>Occ</i> [e, 32] - 1

Backward Search

0	1	2	3																				
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	
<i>T</i> =																							
einsameeselessennassenesseln g ern\$																							
<i>P</i> =																							
less																							

<i>bwt</i> [<i>r</i>]=	nsnm\$ssssg	enlneee	arnele	ieness	seae	
<i>T</i> [<i>pos</i> [<i>r</i>]]=	\$aa	eeeeeeee	egill	mnnnnnn	nrsssss	sss
<i>L</i> =	30		3 + 8			
<i>R</i> =	32		3 + <i>Occ</i> [e, 32] - 1			

Backward Search

0	1	2	3																				
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	
<i>T</i> =																							
einsameeselessennassenessselngern\$																							
<i>P</i> =																							
less																							

bwt[*r*]= nsnm\$ssssgenlneeearneleienesseae

T[*pos*[*r*]]= \$aaeeeeeeeeeeegillmnnnnnnrsssss**sss**

L= 11

R= 12

Backward Search

0 1 2 3
012345678901234567890123456789012
 $T =$ einsameeselessennassenesseln g ern $\$$
 $P =$ less

$bwt[r] =$ nsnm\$ssssgenlneeearneleienesseae
 $T[pos[r]] =$ \$aaeeeeeeeeeeegillmnnnnnrsssssss
 $L = 11$ $less[1] + Occ[1, 10]$
 $R = 12$ $less[1] + Occ[1, 12] - 1$

Backward Search

0	1	2	3
012345678901234567890123456789012			
$T =$ einsameeselessennassenesseln g ern\$			
$P =$ less			

$bwt[r] =$	nsnm\$ssssgenlneeearneleienesseae
$T[pos[r]] =$	\$aaeeeeeeeeeeegillmnnnnnnrsssssss
$L = 11$	$15 + Occ[1, 10]$
$R = 12$	$15 + Occ[1, 12] - 1$

Backward Search

0	1	2	3
012345678901234567890123456789012			
$T =$ einsameeselessennassenesseln g ern\$			
$P =$ less			

$bwt[r] =$	nsnm\$ssssgenlneeearneleienesseae
$T[pos[r]] =$	\$aaeeeeeeeeeeegillmnnnnnnrsssssss
$L =$	11 15 + 0
$R =$	12 15 + $Occ[1, 12] - 1$

Backward Search

	0	1	2	3
	012345678901234567890123456789012			
$T =$	einsame	eselessennassenesseln	gern	\$
$P =$	less			

$bwt[r] =$ nsnm\$ssssgenlneeearneleienesseae
 $T[pos[r]] =$ \$aaeeeeeeeeeeegi1lmnnnnnrsssssss
 $L = 15$
 $R = 15$

$pos = [32, 4, 17, 6, 0, 9, 25, 20, 14, 29, 7, 22, 11, 28, 1, 10, 26,$
 $5, 31, 16, 21, 27, 15, 2, 30, 3, 8, 24, 19, 13, 23, 18, 12]$

- Pro Zeichen des Musters und pro Grenze zwei Lookups, also $\mathcal{O}(1)$.
- Insgesamt erfolgt die Mustersuche in $\mathcal{O}(m)$ Zeit.
- zusätzlicher Speicherplatz:
 - *less*: $\mathcal{O}(|\Sigma| \log(n))$ Bits
 - *Occ*: $\mathcal{O}(|\Sigma| n \log(n))$ Bits

- Pro Zeichen des Musters und pro Grenze zwei Lookups, also $\mathcal{O}(1)$.
- Insgesamt erfolgt die Mustersuche in $\mathcal{O}(m)$ Zeit.
- zusätzlicher Speicherplatz:
 - *less*: $\mathcal{O}(|\Sigma| \log(n))$ Bits
 - *Occ*: $\mathcal{O}(|\Sigma| n \log(n))$ Bits
- *Occ* ist nicht sehr komplex, ist der Speicher nötig?

Optimierung des Speicherplatzverbrauchs

BWT als Bitvektoren (z.B. einen für jeden Buchstaben) darstellen,
Rank-Datenstruktur für Bitsequenzen verwenden

Optimierung des Speicherplatzverbrauchs

BWT als Bitvektoren (z.B. einen für jeden Buchstaben) darstellen,
Rank-Datenstruktur für Bitsequenzen verwenden

Beispiel: $T = \text{banana}\$, bwt = \text{annb}\aa

Bits	0	1	2	3	4	5	6
\$	0	0	0	0	1	0	0
a	1	0	0	0	0	1	1
b	0	0	0	1	0	0	0
n	0	1	1	0	0	0	0

Optimierung des Speicherplatzverbrauchs

BWT als Bitvektoren (z.B. einen für jeden Buchstaben) darstellen,
Rank-Datenstruktur für Bitsequenzen verwenden

Beispiel: $T = \text{banana}\$, bwt = \text{annb}\aa

Bits	0	1	2	3	4	5	6
\$	0	0	0	0	1	0	0
a	1	0	0	0	0	1	1
b	0	0	0	1	0	0	0
n	0	1	1	0	0	0	0

Für jeden Buchstaben eigenes Bitfeld und Rank-Datenstruktur: $\mathcal{O}(n|\Sigma|) + o(n|\Sigma|)$ Bits.

Optimierung des Speicherplatzverbrauchs

BWT als Bitvektoren (z.B. einen für jeden Buchstaben) darstellen,
 Rank-Datenstruktur für Bitsequenzen verwenden

Beispiel: $T = \text{banana}\$, bwt = \text{annb}\aa

Bits	0	1	2	3	4	5	6
\$	0	0	0	0	1	0	0
a	1	0	0	0	0	1	1
b	0	0	0	1	0	0	0
n	0	1	1	0	0	0	0

Für jeden Buchstaben eigenes Bitfeld und Rank-Datenstruktur: $\mathcal{O}(n|\Sigma|) + o(n|\Sigma|)$ Bits.
 Text/bwt selbst benötigt eigentlich nur $\mathcal{O}(n \log |\Sigma|)$ Bits!

Definition (Waveletbaum)

Sei T ein Text mit $|T| = n$ und Alphabet $\Sigma = \{0, \dots, s - 1\}$ mit $|\Sigma| = s$.

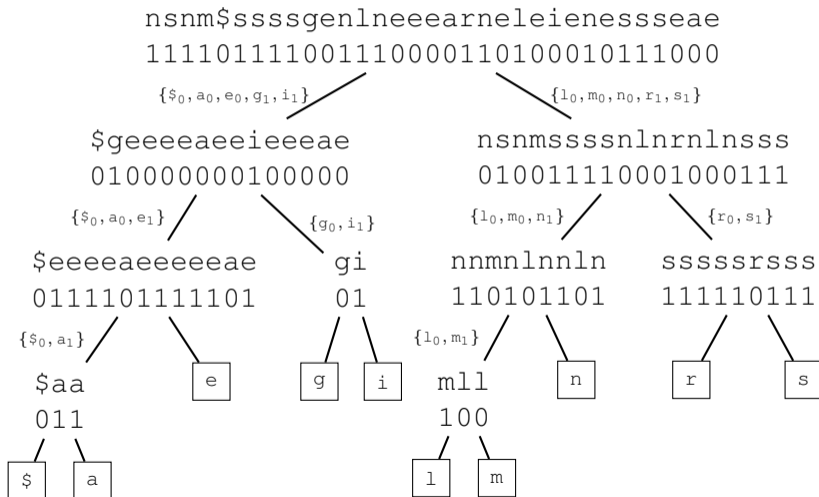
Der Waveletbaum zu T ist ein balancierter Binärbaum mit $|\Sigma|$ Blattknoten.

Ein innerer Knoten repräsentiert die zu einem Teilalphabet $\{a, \dots, b\}$ gehörende Teilsequenz von T (Wurzel entspricht Σ).

Ein innerer Knoten teilt die repräsentierte Sequenz in die Teilsequenzen zu zwei (etwa gleich großen Teilalphabeten (das niedere $\{a, \dots, \lfloor (a + b)/2 \rfloor\}$ und das obere $\{\lfloor (a + b)/2 \rfloor + 1, \dots, b\}$), indem ein Bitvektor die entsprechenden Zugehörigkeiten anzeigt.

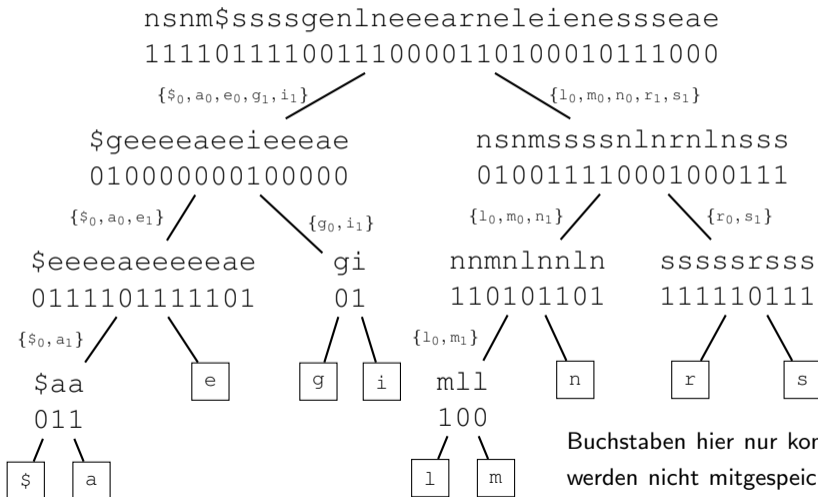
Waveletbaum: Beispiel

$$\Sigma = \{\$, a_0, e_0, g_0, i_0, l_1, m_1, n_1, r_1, s_1\}$$



Waveletbaum: Beispiel

$$\Sigma = \{\$, a_0, e_0, g_0, i_0, l_1, m_1, n_1, r_1, s_1\}$$



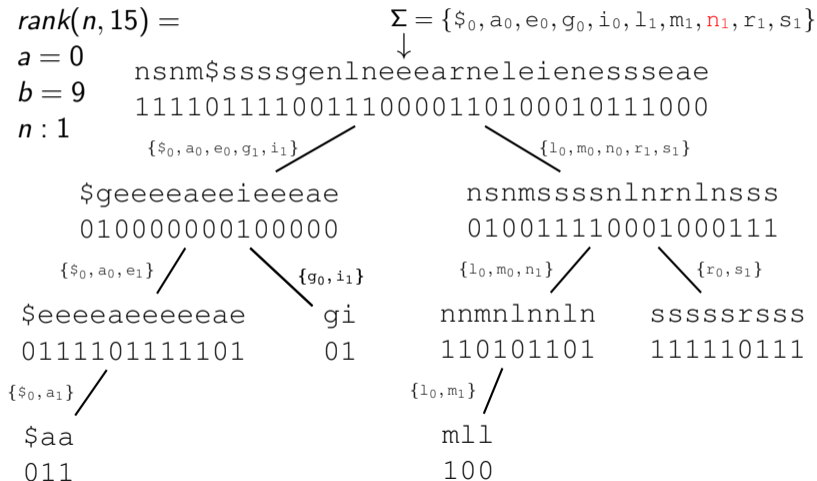
- Es gibt $\log |\Sigma|$ Ebenen in einem Waveletbaum.
- Jede Ebene speichert insgesamt (verteilt auf alle Knoten) maximal n Bits.
- Ein Waveletbaum benötigt so viel Speicher wie der ursprüngliche Text T , nämlich $n \cdot \lceil \log(|\Sigma|) \rceil$ Bits.
- $\mathcal{O}(|\Sigma|)$ Bits zusätzlich für die Baumstruktur
- Der Text selbst muss nicht mehr gespeichert werden.

- Es gibt $\log |\Sigma|$ Ebenen in einem Waveletbaum.
- Jede Ebene speichert insgesamt (verteilt auf alle Knoten) maximal n Bits.
- Ein Waveletbaum benötigt so viel Speicher wie der ursprüngliche Text T , nämlich $n \cdot \lceil \log(|\Sigma|) \rceil$ Bits.
- $\mathcal{O}(|\Sigma|)$ Bits zusätzlich für die Baumstruktur
- Der Text selbst muss nicht mehr gespeichert werden.
- Rank-Datenstrukturen für jeden Knoten, insgesamt $o(n \log |\Sigma|)$ Bits
- Zeichen- und Rank-Anfragen mit Wavelet-Baum in $\mathcal{O}(\log |\Sigma|)$ Zeit

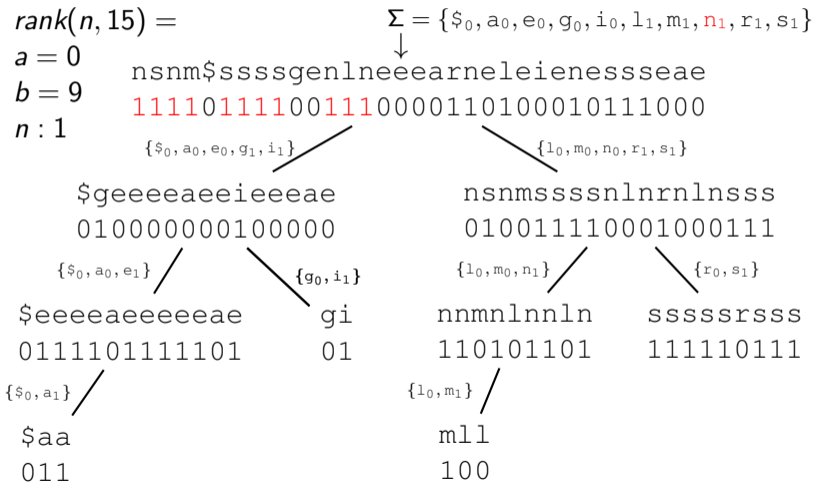
- In der Wurzelebene wird eine Rank-Anfrage $rank(\sigma, i)$ für $\sigma \in \Sigma, i < n$ aufgerufen.
- Es wird überprüft, durch welches Bit (0, 1) in dieser Ebene σ kodiert wird.
 - 1-Bit: Berechne $k = rank(i)$.
 - 0-Bit: Berechne $k = i - rank(i) + 1$.

- In der Wurzelebene wird eine Rank-Anfrage $rank(\sigma, i)$ für $\sigma \in \Sigma, i < n$ aufgerufen.
- Es wird überprüft, durch welches Bit (0, 1) in dieser Ebene σ kodiert wird.
 - 1-Bit: Berechne $k = rank(i)$.
 - 0-Bit: Berechne $k = i - rank(i) + 1$.
- Entsprechend der Kodierung ins linke (0) oder rechte (1) Kind gehen und die Grenzen Alphabetgrenzen a, b anpassen.
- In dieser Ebene $rank(\sigma, k)$ anfragen
- In der Ebene, wo zu σ kein Kind mehr existiert, hat man das Ergebnis.

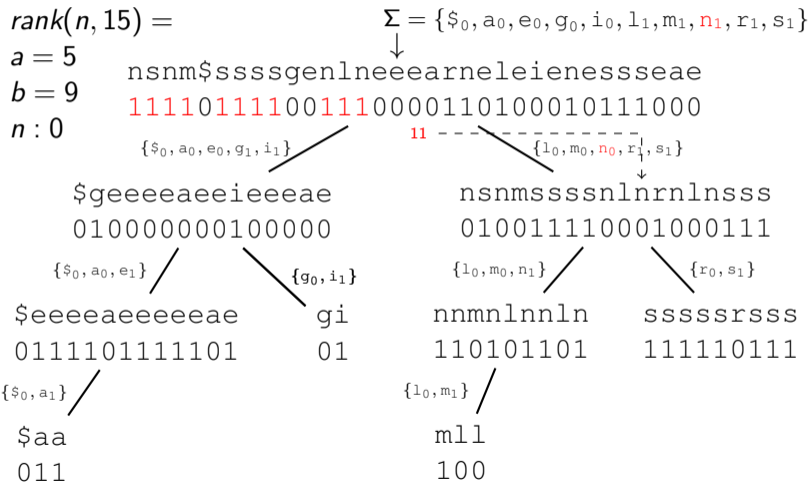
Waveletbaum Rank-Anfrage



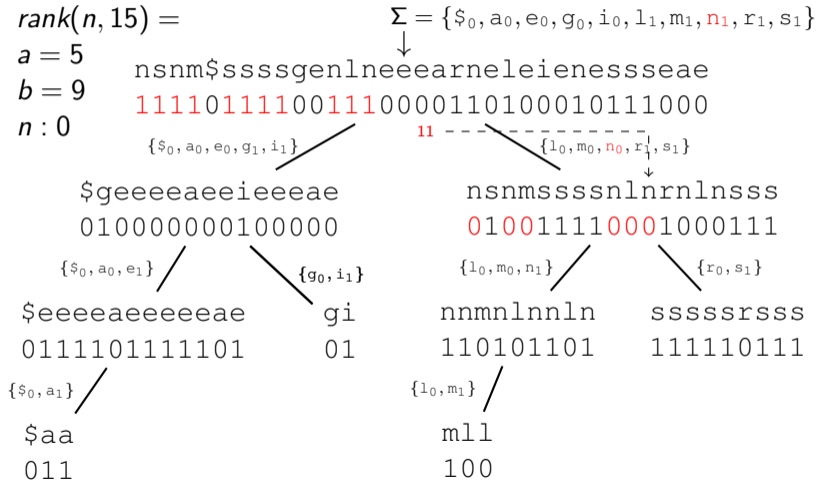
Waveletbaum Rank-Anfrage



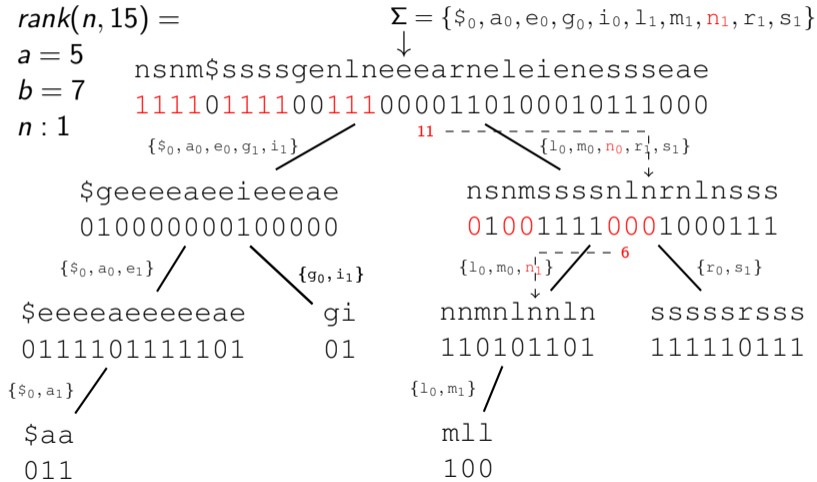
Waveletbaum Rank-Anfrage



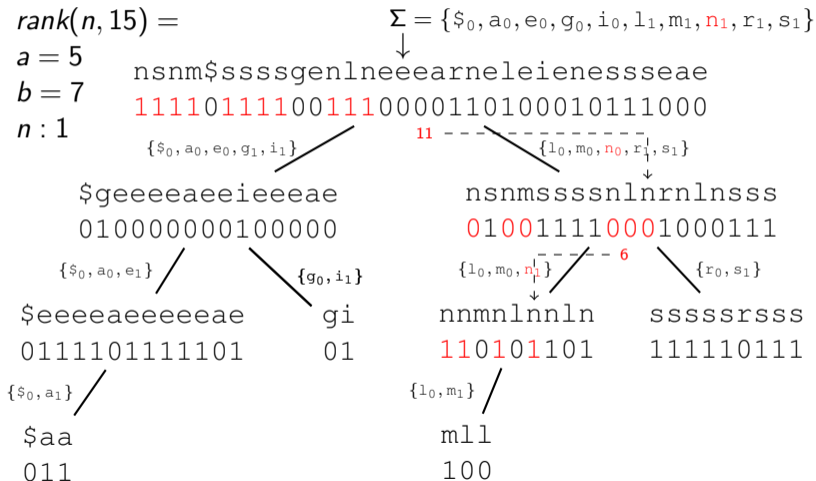
Waveletbaum Rank-Anfrage



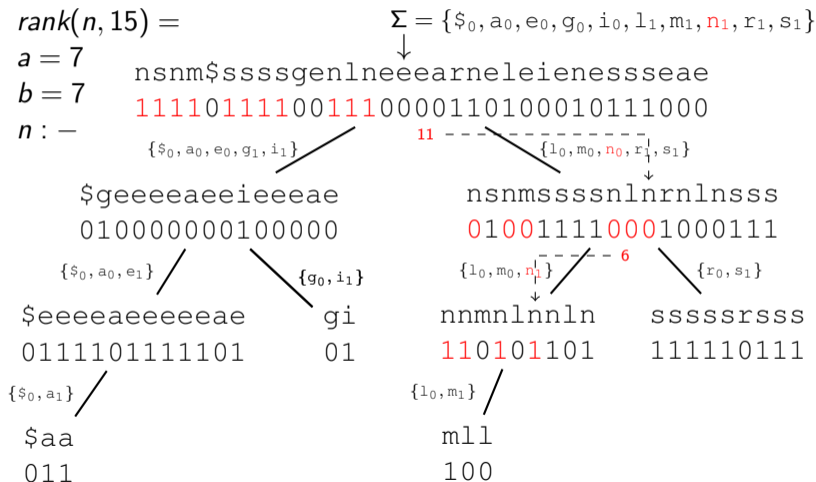
Waveletbaum Rank-Anfrage



Waveletbaum Rank-Anfrage



Waveletbaum Rank-Anfrage



Waveletbaum Rank-Anfrage

$rank(n, 15) = 4$

$a = 7$

$b = 7$

$n : -$

$\Sigma = \{\$0, a_0, e_0, g_0, i_0, l_1, m_1, n_1, r_1, s_1\}$

n s n m \$ s s s s g e n l n e e e a r n e l e i e n e s s e a e

1 1 1 1 0 1 1 1 1 0 0 1 1 1 0 0 0 0 1 1 1 0 0 0 1 0 1 1 1 1 0 0 0

$\{\$0, a_0, e_0, g_1, i_1\}$

11

$\{l_0, m_0, n_0, r_1, s_1\}$

\$ g e e e e a e e i e e e a e

0 1 0 0 0 0 0 0 0 1 0 0 0 0 0

$\{\$0, a_0, e_1\}$

$\{g_0, i_1\}$

n s n m s s s s n l n r n l n s s s

0 1 0 0 1 1 1 1 1 0 0 0 1 0 0 0 1 1 1

$\{l_0, m_0, n_1\}$

6

$\{r_0, s_1\}$

\$ e e e e a e e e e e a e

0 1 1 1 1 0 1 1 1 1 1 1 0 1

$\{\$0, a_1\}$

g i

0 1

n n m l n l n l n

1 1 0 1 0 1 1 0 1

$\{l_0, m_1\}$

s s s s r s s s

1 1 1 1 1 0 1 1 1

\$ a a

0 1 1

m l l

1 0 0

- Die BWT ist eine invertierbare Permutation des Ausgangstextes.
- Sie “verwandelt” wiederholte Teilstrings in Läufe desselben Buchstaben.
- Die BWT ist invertierbar.
- Das Erstellen oder Invertieren der BWT ist in $\mathcal{O}(n)$ Zeit möglich.
- Mit Hilfe des FM-Indexes (Hilfstabellen *less*, *Occ* zur BWT) lässt sich die Mustersuche in $\mathcal{O}(m)$ Zeit durchführen (Backward search).
- Der Waveletbaum der BWT reduziert den Speicherverbrauch der *Occ* Tabelle; insgesamt werden $\mathcal{O}(n \log |\Sigma|) + o(n \log |\Sigma|)$ Bits benötigt.