

SA aus der BWT - Rücktransformation

0 1 2 3 4 5 6
 T = ANANAS\$
 BWT = S\$NNA\$A\$A\$
 LF = 6 0 4 5 1 2 3

Wir wissen wo das letzte Suffix in SA steht

(*) $LF(i) = j \Leftrightarrow A[i] = A[j-1]$

lesen \rightarrow $A[0] = 6$

Schreiben
↓

Mithilfe von (*) und dem Wissen dass $A[0] = n-1$ ist, können wir das SA aus der BWT und dem LF-Mapping konstruieren

$LF(0) = 6 \rightarrow A[6] = 5$
 $LF(6) = 3 \rightarrow A[3] = 4$
 $LF(3) = 5 \rightarrow A[5] = 3$
 $LF(5) = 2 \rightarrow A[2] = 2$
 $LF(2) = 4 \rightarrow A[4] = 1$
 $LF(4) = 1 \rightarrow A[1] = 0$

	0	1	2	3	4	5	6
A =	6	0	2	4	1	3	5
	\$	A	A	A	N	N	S
		N	N	S	A	A	\$
		A	A	\$	N	A	
		N	S		A	S	
		S	\$		S	\$	
		\$			\$		

Können wir LF-Mapping überschreiben?

Ja, können wir
 ↳ Wir greifen auf jeden Eintrag im LF-Mapping nur einmal zu
 ↳ Der Zugriff erfolgt genau nach dem Schreiben ins SA an die entsprechende Stelle (*)

Wie erstellen wir das LF-Array effizient

↳ Wir möchten nur auf der BWT arbeiten und uns das LF-Array erstellen, wenn wir es benötigen

↳ Die Zeichen behalten die gleiche Reihenfolge in L und F

↳ F können wir aus L generieren indem wir L sortieren

BWT = S \$ N N A A A L

↓ stabil(!) sortieren nach dem Zeichen

Die Text-Reihenfolge in F

1 4 5 6 2 3 0
\$ A A A N N S F

↳ ANANAS\$

0 1 2 3 4 5 6

↓ sortiere nach grünen Positionen

0 1 2 3 4 5 6

S \$ N N A A A

6 0 4 5 1 2 3

→ LF-Array

Konstruktion der BWT

- ① Konstruiere SA
- ② Lese SA und schreibe $\epsilon A[i]-1$

alternativ
↳ Konstruiere BWT direkt
z.B. mit DivSufSort

Rücktransformation

- ↳ Nutzen Trick von oben
- ↳ Also ohne LF-Array oder Generierung aller Rotationen.

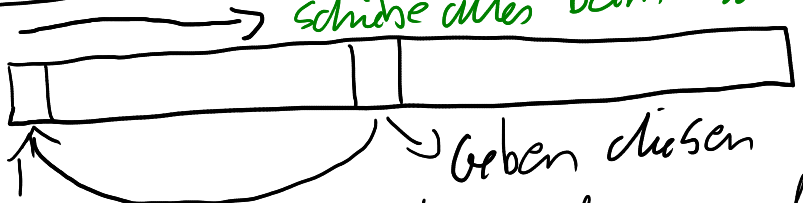
Wichtig: Das ganze funktioniert nur, so lange wir \$ am Ende haben

Alternativ BWST

↳ "Scott"

- ↳ Immer bijektive Transformation.
- Auch ohne \$ am Ende

Move-to-Front



Wir wissen, dass hier das gesuchte Zeichen rein geschrieben wird

Beim Dehoekera funktioniert es analog

Huffman - Kodierung

T = ANANAS\$

BWT = S\$NNAAA

MTF = 3130300

0 1 2 3

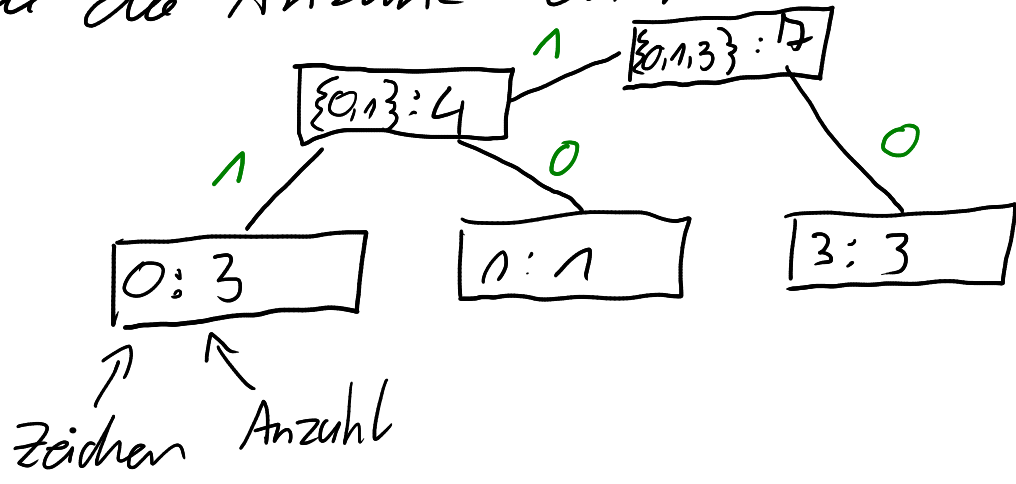
\$ANS → S\$AN → \$SAN → ...

(7 · 2 Bit = 14 Bit)

1. Erzeuge Histogramm → Zählen die Vorkommen der einzelnen Zeichen

Hist = $\begin{matrix} 0 & 1 & 2 & 3 \\ 3 & 1 & 0 & 3 \end{matrix}$

2. Erzeuge Knoten für jedes Zeichen, die die Anzahl enthalten



3. Verschmelze die beiden Knoten mit den kleinsten Anzahlen

Um die Huffman-Codes für einzelne Zeichen zu bekommen lesen wir die Zeichen auf dem Pfad von der Wurzel zum Knoten:

0 → 11
1 → 10
3 → 0

MTF → 01001101111 11 Bit

→ Nachteil: Zum Dehocoderen wird der Baum
oder ein Wörterbuch benötigt
↳ alternative zum Baum

Canonical - Huffman - Codes

1. Berechne die Länge der Huffman-Codes für die Zeichen
2. Starten wir mit dem Code 0 der Länge 1
 - ↳ Immer wenn wir ein neues Zeichen mit der aktuellen Huffman-Code-Länge kodieren wollen addieren wir 1
 - ↳ Wenn wir eine längere Länge benötigen, dann addieren wir auch 1 und fügen passend viele 0en an

Länge 1: 3
Länge 2: 1, 0

- Start mit 0 → Codewort für 3
- addiere 1 und füge eine 0 an → Codewort für 10
- addiere 1 → Codewort für 11

Was ist an diesen Codes besser?

↳ Codes innerhalb einer Länge sind immer aufsteigend

↳ Anstelle eines Baumes müssen wir nur die Längen und die Zeichen übergeben

↳ Für unser Beispiel:

$\{1, 2\}$

der Codes pro Länge

$\{3, 1, 0\}$

Zeichen in der Reihenfolge Länge und Zeitpunkt der Kodierung

Huffman-Codes implementieren

↳ Benutze Priority-Queue

↳ Speichere Knoten, die wir auch für den Baum erzeugen

↳ Verschmelzen wie beim Baum

ABER: Wir speichern die aktuelle Tiefe des, denn die brauchen wir für die Canonical-Huffman-Codes

→ Beim Erstellen speichern wir die Code-Wörter in der Reihenfolge der zu kodierenden Zeichen
↳ code-words [... abcd... ...]

→ Beim Erstellen der Canonical-HC brauchen wir die Reihenfolge der Längen → code-length-order

→ Zum Dekodieren übergen wir einfach die Tabelle denn die ist für ASCII-Texte klein und so müssen wir nicht neu berechnen.