

Burrows Wheeler Transformation – Algorithmen II

Timo Bingmann, Simon Gog

http://algo2.iti.kit.edu/AlgorithmenII_WS19.php

Institut für Theoretische Informatik - Algorithmik II

```
    result = current_weight;
    return true;
}

for( EdgeID eid = graph.edgeBegin( current ); eid != graph.edgeEnd( current ); ++eid ){
    const Edge & edge = graph.getEdge( eid );
    COUNTING( statistic_data.inc( DijkstraStatisticData::TOUCHED_EDGES ); )
    if( edge.forward ){
        COUNTING( statistic_data.inc( DijkstraStatisticData::RELAXED_EDGES ); )
        weight new_weight = edge.weight + current_weight;
        GUARANTEE( new_weight >= current_weight, std::runtime_error, "Weight overflow detected." );
        if( !priority_queue.isReached( edge.target ) ){
            COUNTING( statistic_data.inc( DijkstraStatisticData::SUCCESSFULLY_RELAXED_EDGES ); )
            COUNTING( statistic_data.inc( DijkstraStatisticData::REACHED_NODES ); )
            priority_queue.push( edge.target, new_weight );
        } else {
            if( priority_queue.getCurrentKey( edge.target ) > new_weight ){
                COUNTING( statistic_data.inc( DijkstraStatisticData::SUCCESSFULLY_RELAXED_NODES ); )
                priority_queue.decreaseKey( edge.target, new_weight );
            }
        }
    }
}
```

Burrows, Wheeler (1983, 1994)

- “nur” eine Umordnung des Textes,
→ gruppiert Zeichen mit ähnlichem Kontext, reversibel
- ursprünglich eingeführt zur **Textkompression**
→ Teilschritt von **bzip2**
- später auch für **Textindizierung** verwendet
→ z.B. Rückwärtssuche, Berechnung des LCP-Array

Burrows-Wheeler-Transformation

Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

1 2 3 4 5 6 7 8 9
 $T = \text{l a l a n g n g \$}$

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)

Burrows-Wheeler-Transformation

Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	
										\$
										\$
							n	g	g	\$
							n	g	g	\$
						n	g	g	n	\$
					a	n	g	n	g	\$
				l	a	n	g	n	g	\$
		a	l	a	n	g	n	g	g	\$
l	a	l	a	n	g	n	g	g	\$	

Burrows-Wheeler-Transformation

Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	
		a	l	a	n	g	n	g	\$	9
			a	n	g	n	g	\$	\$	2
				a	n	g	n	\$	\$	4
					n	g	n	\$	\$	8
						n	g	\$	\$	6
	l	a	l	a	n	g	n	\$	\$	1
		l	a	n	g	n	\$	\$	\$	3
					n	g	\$	\$	\$	7
						n	g	\$	\$	5

Burrows-Wheeler-Transformation

Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)

1	2	3	4	5	6	7	8	9		
$T =$	l	a	l	a	n	g	n	g	\$	SA
									\$	9
	a	l	a	n	g	n	g	\$	\$	2
			a	n	g	n	g	\$	\$	4
									\$	8
						g	n	g	\$	6
l	a	l	a	n	g	n	g	\$	\$	1
		l	a	n	g	n	g	\$	\$	3
						n	g	\$	\$	7
					n	g	n	g	\$	5

Burrows-Wheeler-Transformation

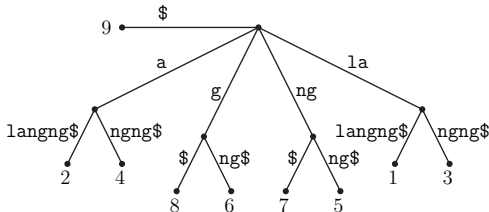
Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	SA
		a	l	a	n	g	n	g	\$	9
			a	n	g	n	g	\$	\$	2
				a	n	g	n	g	\$	4
					a	n	g	\$	\$	8
						a	n	\$	\$	6
							a	\$	\$	1
								\$	\$	3
									\$	7
										5



Burrows-Wheeler-Transformation

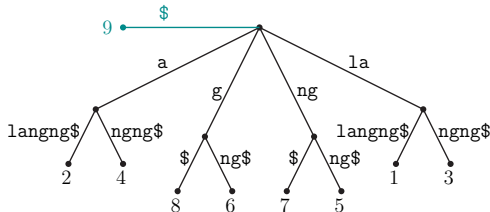
Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	SA
		a	l	a	n	g	n	g	\$	9
			a	n	g	n	g	\$	\$	2
				a	n	g	n	g	\$	4
					a	n	g	\$	\$	8
						a	n	\$	\$	6
	l	a	l	a	n	g	n	g	\$	1
		l	a	n	g	n	g	\$	\$	3
						n	g	\$	\$	7
							n	\$	\$	5



Burrows-Wheeler-Transformation

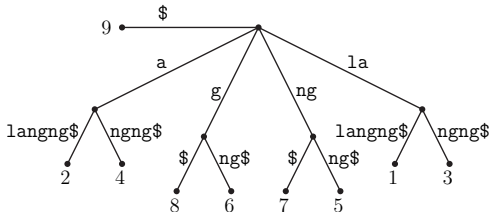
Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	SA
		a	l	a	n	g	n	g	\$	9
			a	n	g	n	g	\$	\$	2
				a	n	g	n	\$	\$	4
					a	n	\$	\$	\$	8
						g	n	\$	\$	6
	l	a					g	n	\$	1
		l	a					n	\$	3
							n	\$	\$	7
								n	\$	5

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)



Burrows-Wheeler-Transformation

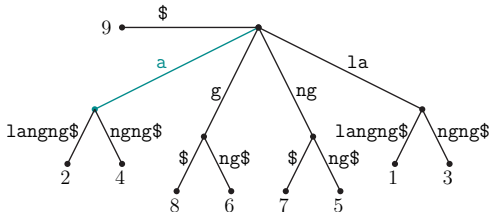
Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	SA
										9
										2
										4
										8
										6
										1
										3
										7
										5

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)



Burrows-Wheeler-Transformation

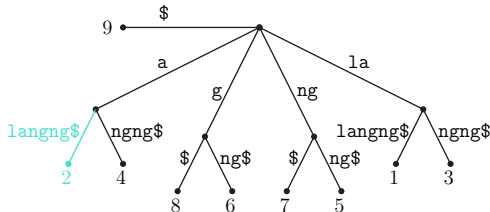
Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	
										SA
										9
										2
										4
										8
										6
										1
										3
										7
										5



Burrows-Wheeler-Transformation

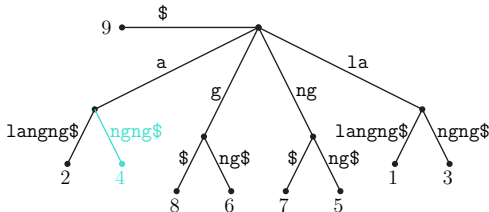
Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	SA
										9
										2
										4
										8
										6
										1
										3
										7
										5

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)



Burrows-Wheeler-Transformation

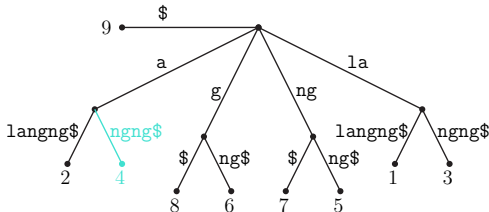
Wiederholung: Suffix-Arrays (und Suffixbäume)

- Textindizierung
→ schnelle Suche in Texten
- kann in $O(n)$ erzeugt werden
(DC3-Algorithmus)

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	SA
										9
										2
										4
										8
										6
										1
										3
										7
										5

Definition

- $SA[i] \hat{=}$ Index des i -ten sortierten Suffix
(einer Zeichenkette T)



USW.

Burrows-Wheeler-Transformation

Transformation

1 2 3 4 5 6 7 8 9
 $T = l a l a n g n g \$$

Burrows-Wheeler-Transformation

Transformation

	1	2	3	4	5	6	7	8	9
$T =$	l	a	l	a	n	g	n	g	\$
	a	l	a	n	g	n	g	\$	l

Burrows-Wheeler-Transformation

Transformation

	1	2	3	4	5	6	7	8	9
$T =$	l	a	l	a	n	g	n	g	\$
	a	l	a	n	g	n	g	\$	l
	l	a	n	g	n	g	\$	l	a

Burrows-Wheeler-Transformation

Transformation

	1	2	3	4	5	6	7	8	9
$T =$	l	a	l	a	n	g	n	g	\$
	a	l	a	n	g	n	g	\$	l
	l	a	n	g	n	g	\$	l	a
	a	n	g	n	g	\$	l	a	l

Burrows-Wheeler-Transformation

Transformation

	1	2	3	4	5	6	7	8	9
$T =$	l	a	l	a	n	g	n	g	\$
	a	l	a	n	g	n	g	\$	l
	l	a	n	g	n	g	\$	l	a
	a	n	g	n	g	\$	l	a	l
				⋮					

Burrows-Wheeler-Transformation

Transformation

	1	2	3	4	5	6	7	8	9
$T =$	l	a	l	a	n	g	n	g	\$
	a	l	a	n	g	n	g	\$	l
	l	a	n	g	n	g	\$	l	a
	a	n	g	n	g	\$	l	a	l
	n	g	n	g	\$	l	a	l	a
	g	n	g	\$	l	a	l	a	n
	n	g	\$	l	a	l	a	n	g
	g	\$	l	a	l	a	n	g	n
	\$	l	a	l	a	n	g	n	g

Burrows-Wheeler-Transformation

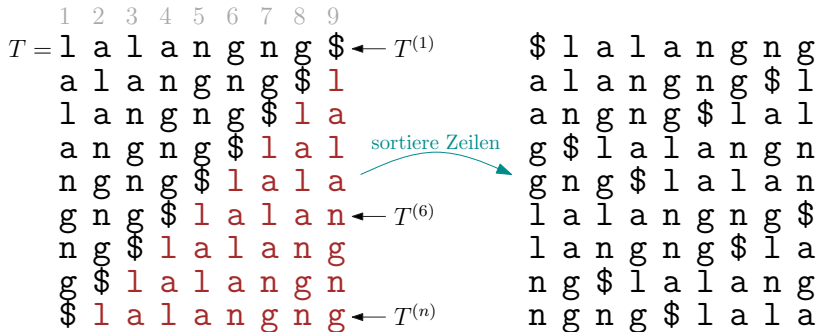
Transformation

	1	2	3	4	5	6	7	8	9	
$T =$	l	a	l	a	n	g	n	g	\$	$\leftarrow T^{(1)}$
	a	l	a	n	g	n	g	\$	l	
	l	a	n	g	n	g	\$	l	a	
	a	n	g	n	g	\$	l	a	l	
	n	g	n	g	\$	l	a	l	a	
	g	n	g	\$	l	a	l	a	n	$\leftarrow T^{(6)}$
	n	g	\$	l	a	l	a	n	g	
	\$	l	a	l	a	n	g	n	g	$\leftarrow T^{(n)}$

- $T^{(i)} \hat{=} T$ zyklisch ab Position i , Länge $n = |T|$

Burrows-Wheeler-Transformation

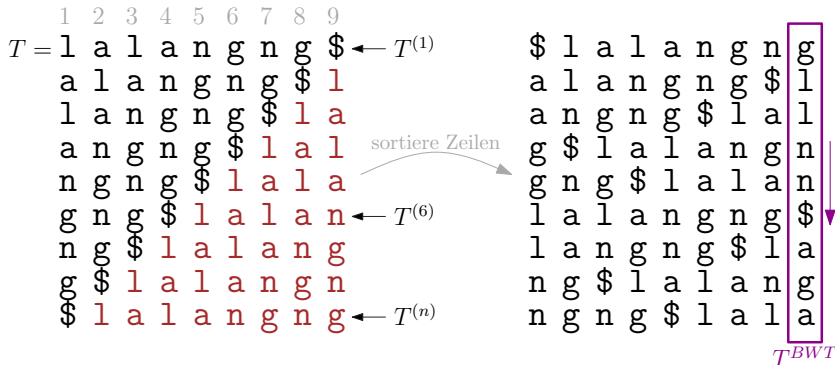
Transformation



- $T^{(i)} \hat{=} T$ zyklisch ab Position i , Länge $n = |T|$

Burrows-Wheeler-Transformation

Transformation

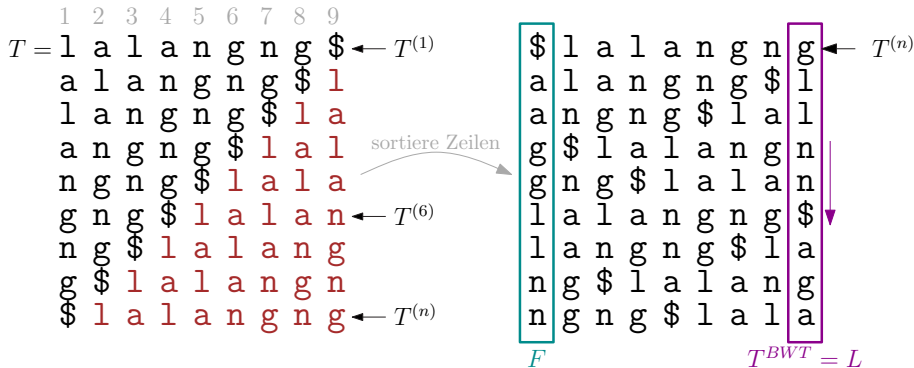


- $T^{(i)} \hat{=} T$ zyklisch ab Position i , Länge $n = |T|$
- $T = \text{lalangng\$} \rightarrow T^{BWT} = \text{gllnn\$aga}$

$O(n^2 + n \log n)$

Burrows-Wheeler-Transformation

Transformation



- $T^{(i)} \hat{=} T$ zyklisch ab Position i , Länge $n = |T|$
- $T = \text{lalangng\$} \rightarrow T^{BWT} = \text{gllnn\$aga}$

$$O(n^2 + n \log n)$$

Burrows-Wheeler-Transformation

Eigenschaften

1 2 3 4 5 6 7 8 9
 $T = l a l a n g n g \$$

- BWT in $\mathcal{O}(n)$ berechenbar
- Zeilen enthalten sortierte Suffixe
- Zeichen in letzter Spalte entspricht Zeichen vor Suffix in T

\$	l	a	l	a	n	g	n	g	1
a	l	a	n	g	n	g	\$	l	2
a	n	g	n	g	\$	l	a	l	3
g	\$	l	a	l	a	n	g	n	4
g	n	g	\$	l	a	l	a	n	5
l	a	l	a	n	g	n	g	\$	6
l	a	n	g	n	g	\$	l	a	7
n	g	\$	l	a	l	a	n	g	8
n	g	n	g	\$	l	a	l	a	9

Burrows-Wheeler-Transformation

Eigenschaften

1 2 3 4 5 6 7 8 9
 $T = l a l a n g n g \$$

- BWT in $\mathcal{O}(n)$ berechenbar

- Zeilen enthalten sortierte Suffixe

- Zeichen in letzter Spalte
entspricht Zeichen vor Suffix in T

\$	l	a	l	a	n	g	n	g	1
a	l	a	n	g	n	g	\$	l	2
a	n	g	n	g	\$	l	a	l	3
g	\$	l	a	l	a	n	g	n	4
g	n	g	\$	l	a	l	a	n	5
l	a	l	a	n	g	n	g	\$	6
l	a	n	g	n	g	\$	l	a	7
n	g	\$	l	a	l	a	n	g	8
n	g	n	g	\$	l	a	l	a	9

Burrows-Wheeler-Transformation

Eigenschaften

1 2 3 4 5 6 7 8 9
 $T = l a l a n g n g \$$

- BWT in $\mathcal{O}(n)$ berechenbar
- Zeilen enthalten sortierte Suffixe
- Zeichen in letzter Spalte entspricht Zeichen vor Suffix in T

\$	l	a	l	a	n	g	n	g	1
a	l	a	n	g	n	g	\$	l	2
a	n	g	n	g	\$	l	a	l	3
g	\$	l	a	l	a	n	g	n	4
g	n	g	\$	l	a	l	a	n	5
l	a	l	a	n	g	n	g	\$	6
l	a	n	g	n	g	\$	l	a	7
n	g	\$	l	a	l	a	n	g	8
n	g	n	g	\$	l	a	l	a	9

Burrows-Wheeler-Transformation

Eigenschaften

1 2 3 4 5 6 7 8 9
 $T = l a l a n g n g \$$

- BWT in $\mathcal{O}(n)$ berechenbar
- Zeilen enthalten sortierte Suffixe
- Zeichen in letzter Spalte entspricht Zeichen vor Suffix in T

\$	l	a	l	a	n	g	n	g	l	1
a	l	a	n	g	n	g	\$	l	l	2
a	n	g	n	g	\$	l	a	l	l	3
g	\$	l	a	l	a	n	g	n	g	4
g	n	g	\$	l	a	l	a	n	g	5
l	a	l	a	n	g	n	g	\$	l	6
l	a	n	g	n	g	\$	l	a	l	7
n	g	\$	l	a	l	a	n	g	l	8
n	g	n	g	\$	l	a	l	a	l	9

$T^{BWT} = L$

Burrows-Wheeler-Transformation

Eigenschaften

1 2 3 4 5 6 7 8 9
 $T = l a l a n g n g \$$

- BWT in $\mathcal{O}(n)$ berechenbar

- Zeilen enthalten sortierte Suffixe

- Zeichen in letzter Spalte entspricht Zeichen vor Suffix in T

- $T^{BWT}[i] = L[i] = T[SA[i] - 1] = T^{(SA[i])}[n]$
($T^{BWT}[i]$ ist das Zeichen vor dem i -ten Suffix in T)

\$	l	a	l	a	n	g	n	g	n	g	l	1	SA
a	l	a	n	g	n	g	\$	l	l	l	l	2	2
a	n	g	n	g	\$	l	a	l	l	l	l	3	4
g	\$	l	a	l	a	n	g	n	g	n	g	4	8
g	n	g	\$	l	a	l	a	n	g	n	g	5	6
l	a	l	a	n	g	n	g	\$	l	a	l	6	1
l	a	n	g	n	g	\$	l	a	l	l	l	7	3
n	g	\$	l	a	l	a	n	g	n	g	n	8	7
n	g	n	g	\$	l	a	l	a	l	l	l	9	5

$T^{BWT} = L$

Burrows-Wheeler-Transformation

Rücktransformation – Vorüberlegungen

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

■ betrachte Matrix aus Transformation

- erste Zeile enthält Lösung $T^{(n)}$
- $T^{BWT} = L$ gegeben
- F leicht zu bestimmen (sortiere L)
- L ist immer Spalte vor F (zyklisch)

\$	l	a	l	a	n	g	n	g	l
a	l	a	n	g	n	g	\$	l	l
a	n	g	n	g	\$	l	a	l	l
g	\$	l	a	l	a	n	g	n	l
g	n	g	\$	l	a	l	a	n	l
l	a	l	a	n	g	n	g	\$	l
l	a	n	g	n	g	\$	l	a	l
n	g	\$	l	a	l	a	n	g	l
n	g	n	g	\$	l	a	l	a	l

$$T^{BWT} = L$$

Burrows-Wheeler-Transformation

Rücktransformation – Vorüberlegungen

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

■ betrachte Matrix aus Transformation

- erste Zeile enthält Lösung $T^{(n)}$
- $T^{BWT} = L$ gegeben
- F leicht zu bestimmen (sortiere L)
- L ist immer Spalte vor F (zyklisch)

\$	l	a	l	a	n	g	n	g	← $T^{(n)}$
a	l	a	n	g	n	g	\$	l	
a	n	g	n	g	\$	l	a	l	
g	\$	l	a	l	a	n	g	n	
g	n	g	\$	l	a	l	a	n	
l	a	l	a	n	g	n	g	\$	
l	a	n	g	n	g	\$	l	a	
n	g	\$	l	a	l	a	n	g	
n	g	n	g	\$	l	a	l	a	

$T^{BWT} = L$

Burrows-Wheeler-Transformation

Rücktransformation – Vorüberlegungen

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

■ betrachte Matrix aus Transformation

- erste Zeile enthält Lösung $T^{(n)}$
- $T^{BWT} = L$ gegeben
- F leicht zu bestimmen (sortiere L)
- L ist immer Spalte vor F (zyklisch)

\$	l	a	l	a	n	g	n	g	l	← $T^{(n)}$
a	l	a	n	g	n	g	\$	l	l	
a	n	g	n	g	\$	l	a	l	l	
g	\$	l	a	l	a	n	g	n	g	
g	n	g	\$	l	a	l	a	n	g	
l	a	l	a	n	g	n	g	\$	l	
l	a	n	g	n	g	\$	l	a	g	
n	g	\$	l	a	l	a	n	g	a	
n	g	n	g	\$	l	a	l	a	a	

$T^{BWT} = L$

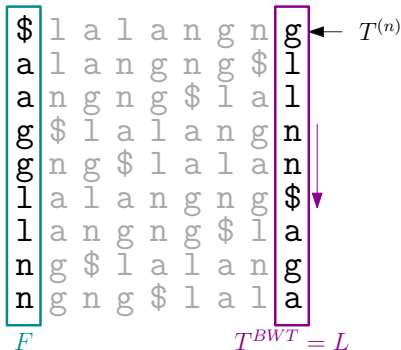
Burrows-Wheeler-Transformation

Rücktransformation – Vorüberlegungen

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

■ betrachte Matrix aus Transformation

- erste Zeile enthält Lösung $T^{(n)}$
- $T^{BWT} = L$ gegeben
- F leicht zu bestimmen (sortiere L)
- L ist immer Spalte vor F (zyklisch)



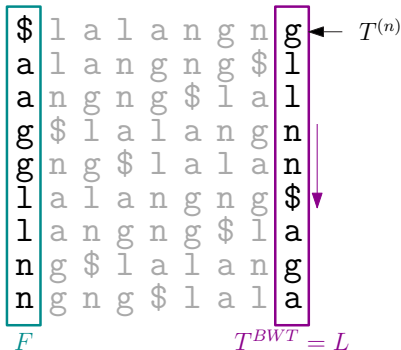
Burrows-Wheeler-Transformation

Rücktransformation – Vorüberlegungen

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

■ betrachte Matrix aus Transformation

- erste Zeile enthält Lösung $T^{(n)}$
- $T^{BWT} = L$ gegeben
- F leicht zu bestimmen (sortiere L)
- L ist immer Spalte vor F (zyklisch)



Burrows-Wheeler-Transformation

Rücktransformation – Vorüberlegungen

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

■ betrachte Matrix aus Transformation

- erste Zeile enthält Lösung $T^{(n)}$
- $T^{BWT} = L$ gegeben
- F leicht zu bestimmen (sortiere L)
- L ist immer Spalte vor F (zyklisch)

l	a	l	a	n	g	n	g	l	a	← $T^{(1)}$
l	a	n	g	n	g	\$	l	a	l	a
n	g	n	g	\$	l	a	l	a	n	g
\$	l	a	l	a	n	g	n	g	l	a
n	g	\$	l	a	l	a	n	g	l	a
a	l	a	n	g	n	g	\$	l	l	a
a	n	g	n	g	\$	l	a	l	g	n
g	\$	l	a	l	a	n	g	n	a	n
g	n	g	\$	l	a	l	a	n	g	n

L F

Burrows-Wheeler-Transformation

Rücktransformation – Vorüberlegungen

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

■ betrachte Matrix aus Transformation

- erste Zeile enthält Lösung $T^{(n)}$
- $T^{BWT} = L$ gegeben
- F leicht zu bestimmen (sortiere L)
- L ist immer Spalte vor F (zyklisch)

a	n	g	n	g	\$	l	a	l	← $T^{(4)}$
g	n	g	\$	l	a	l	a	n	
g	\$	l	a	l	a	n	g	n	
l	a	n	g	n	g	\$	l	a	
l	a	l	a	n	g	n	g	\$	
n	g	n	g	\$	l	a	l	a	
n	g	\$	l	a	l	a	n	g	
a	l	a	n	g	n	g	\$	l	
\$	l	a	l	a	n	g	n	g	

L F

Burrows-Wheeler-Transformation

Rücktransformation – Vorüberlegungen

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

■ betrachte Matrix aus Transformation

- erste Zeile enthält Lösung $T^{(n)}$
- $T^{BWT} = L$ gegeben
- F leicht zu bestimmen (sortiere L)
- L ist immer Spalte vor F (zyklisch)

n	g	\$	l	a	l	a	n	g	← $T^{(7)}$
\$	l	a	l	a	n	g	n	g	
a	l	a	n	g	n	g	\$	l	
g	n	g	\$	l	a	l	a	n	
a	n	g	n	g	\$	l	a	l	
g	\$	l	a	l	a	n	g	n	
l	a	l	a	n	g	n	g	\$	
n	g	n	g	\$	l	a	l	a	
l	a	n	g	n	g	\$	l	a	
	L	F							

Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} =$

	1	2	3	4	5	6	7	8	9
	g	l	l	n	n	\$	a	g	a

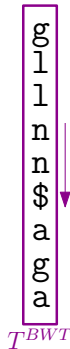
- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert

Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n \$ a g a

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert



Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n \$ a g a

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert

\$
a
a
a
a
l
l
n
n

sortiert nach letzter Spalte

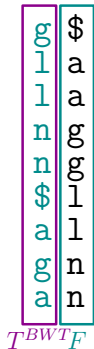
F

Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n \$ a g a

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert



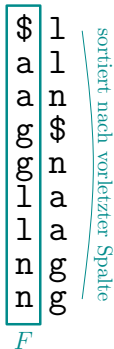
$T^{BWT}F$

Burrows-Wheeler-Transformation

Rücktransformation

$$T^{BWT} = \begin{array}{cccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \text{g} & \text{l} & \text{l} & \text{n} & \text{n} & \$ & \text{a} & \text{g} & \text{a} & \end{array}$$

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert



Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} =$

	1	2	3	4	5	6	7	8	9
	g	l	l	n	n	\$	a	g	a

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert

g	\$	l
l	a	l
l	a	n
n	g	\$
n	g	n
\$	l	a
a	l	a
g	n	g
a	n	g

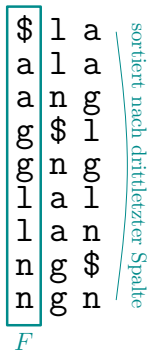
T^{BWT}_F

Burrows-Wheeler-Transformation

Rücktransformation

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert



Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n \$ a g a

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert

g	\$	l	a
l	a	l	a
l	a	n	g
n	g	\$	l
n	g	n	g
\$	l	a	l
a	l	a	n
g	n	g	\$
a	n	g	n

$T^{BWT}F$

Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n \$ a g a

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert

...

g	\$	l	a
l	a	l	a
l	a	n	g
n	g	\$	l
n	g	n	g
\$	l	a	l
a	l	a	n
g	n	g	\$
a	n	g	n

T^{BWT}_F

Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n \$ a g a

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert

\$	l	a	l	a	n	g	n	g
a	l	a	n	g	n	g	\$	l
a	n	g	n	g	\$	l	a	l
g	\$	l	a	l	a	n	g	n
g	n	g	\$	l	a	l	a	n
l	a	l	a	n	g	n	g	\$
l	a	n	g	n	g	\$	l	a
n	g	\$	l	a	l	a	n	g
n	g	n	g	\$	l	a	l	a

Burrows-Wheeler-Transformation

Rücktransformation

$T^{BWT} = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ g & l & l & n & n & \$ & a & g & a \end{matrix}$

- schreibe T^{BWT} in Spaltenform
- sortiere zeilenweise
- schreibe T^{BWT} in Spaltenform davor
- wiederhole bis $|T^{BWT}|$ mal sortiert

\$	l	a	l	a	n	g	n	g	← $T^{(n)}$
a	l	a	n	g	n	g	\$	l	
a	n	g	n	g	\$	l	a	l	
g	\$	l	a	l	a	n	g	n	
g	n	g	\$	l	a	l	a	n	
l	a	l	a	n	g	n	g	\$	
l	a	n	g	n	g	\$	l	a	
n	g	\$	l	a	l	a	n	g	
n	g	n	g	\$	l	a	l	a	

■ $T^{BWT} = gllnn$aga \rightarrow T = lalngng$$

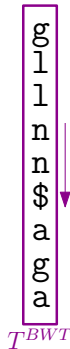
$O(n^2 \log n)$

Burrows-Wheeler-Transformation

Rücktransformation – weitere Überlegungen

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \\ T & = & \text{-----} & ? & \text{-----} & & & & & & \end{array}$$

- Wie lautet das Vorgängerzeichen?
→ *last-to-front mapping* $LF[\cdot]$
(Position in $L[\cdot]$ an der Vorgänger steht)

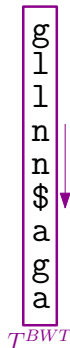


Burrows-Wheeler-Transformation

Rücktransformation – weitere Überlegungen

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \\ T = & & & & & & & & & \mathbf{\$} \end{array}$$

- Wie lautet das Vorgängerzeichen?
→ *last-to-front mapping* $LF[\cdot]$
(Position in $L[\cdot]$ an der Vorgänger steht)



Burrows-Wheeler-Transformation

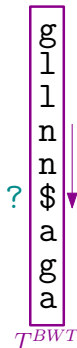
Rücktransformation – weitere Überlegungen

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \\ T = & & & & & & & ? & \$ & \end{array}$$

- Wie lautet das Vorgängerzeichen?

→ *last-to-front mapping* $LF[\cdot]$

(Position in $L[\cdot]$ an der Vorgänger steht)



Burrows-Wheeler-Transformation

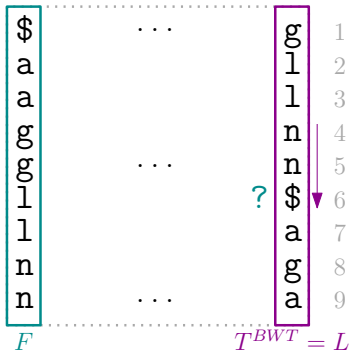
Rücktransformation – weitere Überlegungen

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \\ T = & & & & & & & ? & \$ & \end{array}$$

- Wie lautet das Vorgängerzeichen?

→ *last-to-front mapping* $LF[\cdot]$

(Position in $L[\cdot]$ an der Vorgänger steht)



- $LF[6] = ?$

Burrows-Wheeler-Transformation

Rücktransformation – weitere Überlegungen

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & g & l & l & n & n & \$ & a & g & a \\ T & = & & & & & & & ? & \$ & \end{array}$$

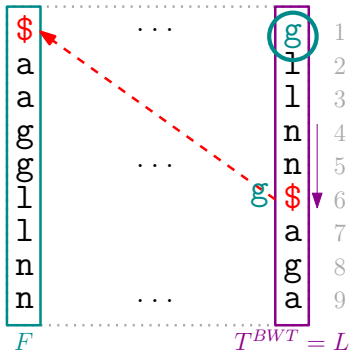
- Wie lautet das Vorgängerzeichen?

→ *last-to-front mapping* $LF[\cdot]$

(Position in $L[\cdot]$ an der Vorgänger steht)

- $LF[6] = 1$

($LF[6]$ ist die Position in $F[\cdot]$ an der $L[6]$ steht)



Burrows-Wheeler-Transformation

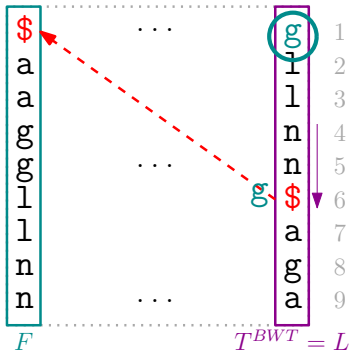
Rücktransformation – weitere Überlegungen

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & g & l & l & n & n & \$ & a & g & a \\ T & = & & & & & & & ? & \$ & \end{array}$$

- Wie lautet das Vorgängerzeichen?

→ *last-to-front mapping* $LF[\cdot]$

(Position in $L[\cdot]$ an der Vorgänger steht)



- $LF[i] = j \Leftrightarrow T^{(SA[j])} = (T^{(SA[i])})^{(n)} \Leftrightarrow SA[i] = SA[j] - 1 \pmod{n}$
($LF[i]$ ist die Position in $F[\cdot]$ an der $L[i]$ steht)

Burrows-Wheeler-Transformation

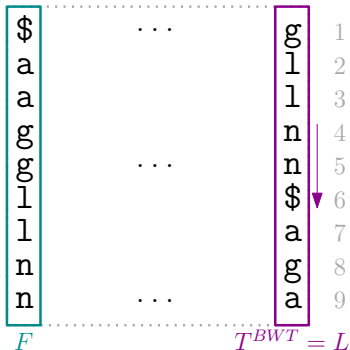
Rücktransformation – weitere Überlegungen

$$T^{BWT} = \begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{array}$$

- T^{BWT} hat Struktur, so dass gilt

→ gleiche Zeichen besitzen
gleiche Reihenfolge in $F[\cdot]$ und $L[\cdot]$

→ falls $L[i] = L[j]$ mit $i < j$,
dann $LF[i] < LF[j]$



- Grund: $\alpha < \beta$ lexikographisch (nach Konstruktion)

Burrows-Wheeler-Transformation

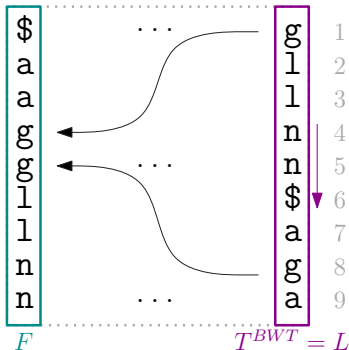
Rücktransformation – weitere Überlegungen

$$T^{BWT} = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \end{matrix}$$

- T^{BWT} hat Struktur, so dass gilt

→ gleiche Zeichen besitzen
gleiche Reihenfolge in $F[\cdot]$ und $L[\cdot]$

→ falls $L[i] = L[j]$ mit $i < j$,
dann $LF[i] < LF[j]$



- Grund: $\alpha < \beta$ lexikographisch (nach Konstruktion)

Burrows-Wheeler-Transformation

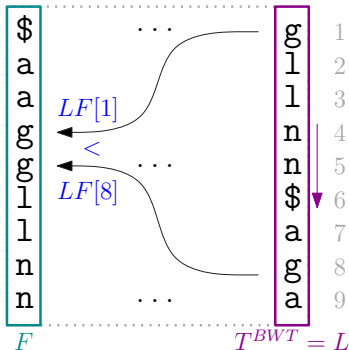
Rücktransformation – weitere Überlegungen

$$T^{BWT} = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \end{matrix}$$

- T^{BWT} hat Struktur, so dass gilt

→ gleiche Zeichen besitzen
gleiche Reihenfolge in $F[\cdot]$ und $L[\cdot]$

→ falls $L[i] = L[j]$ mit $i < j$,
dann $LF[i] < LF[j]$



- Grund: $\alpha < \beta$ lexikographisch (nach Konstruktion)

Burrows-Wheeler-Transformation

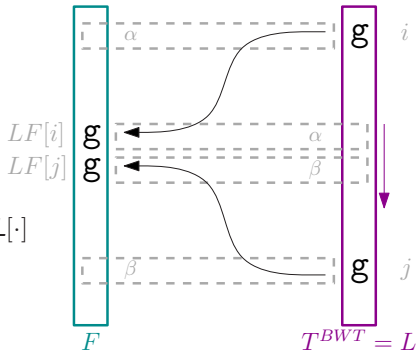
Rücktransformation – weitere Überlegungen

$$T^{BWT} = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \end{matrix}$$

- T^{BWT} hat Struktur, so dass gilt

→ gleiche Zeichen besitzen
gleiche Reihenfolge in $F[\cdot]$ und $L[\cdot]$

→ falls $L[i] = L[j]$ mit $i < j$,
dann $LF[i] < LF[j]$



- Grund: $\alpha < \beta$ lexikographisch (nach Konstruktion)

Burrows-Wheeler-Transformation

Berechnung von $LF[\cdot]$

$T^{BWT} =$

	1	2	3	4	5	6	7	8	9
	g	l	l	n	n	\$	a	g	a

- $LF[\cdot]$ nur aus $T^{BWT}[\cdot]$ berechenbar

- $LF[i] = C(L[i]) + occ[i]$
 - $C(a) = \#$ Zeichen kleiner als a
 - $occ[i] = \#$ Zeichen gleich $L[i]$ in $L[1..i]$

($LF[i]$ ist Position in $F[\cdot]$, an der $L[i]$ steht)



Burrows-Wheeler-Transformation

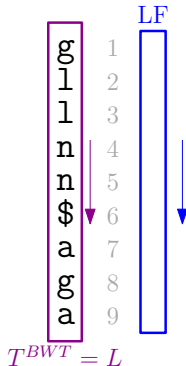
Berechnung von $LF[\cdot]$

$$T^{BWT} = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & g & l & l & n & n & \$ & a & g & a \end{matrix}$$

- $LF[\cdot]$ nur aus $T^{BWT}[\cdot]$ berechenbar

- $LF[i] = C(L[i]) + occ[i]$
 - $C(a) = \#$ Zeichen kleiner als a
 - $occ[i] = \#$ Zeichen gleich $L[i]$ in $L[1..i]$

($LF[i]$ ist Position in $F[\cdot]$, an der $L[i]$ steht)



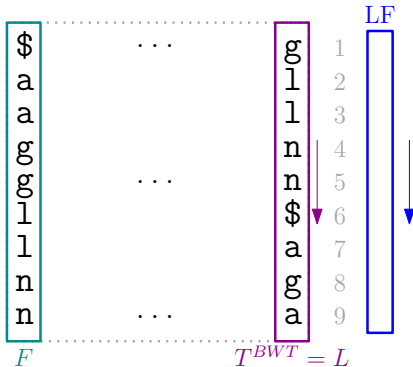
Burrows-Wheeler-Transformation

Berechnung von $LF[\cdot]$

$$T^{BWT} = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ g & l & l & n & n & \$ & a & g & a \end{matrix}$$

- $LF[\cdot]$ nur aus $T^{BWT}[\cdot]$ berechenbar

- $LF[i] = C(L[i]) + occ[i]$
 - $C(a) = \# \text{Zeichen kleiner als } a$
 - $occ[i] = \# \text{Zeichen gleich } L[i] \text{ in } L[1..i]$($LF[i]$ ist Position in $F[\cdot]$, an der $L[i]$ steht)



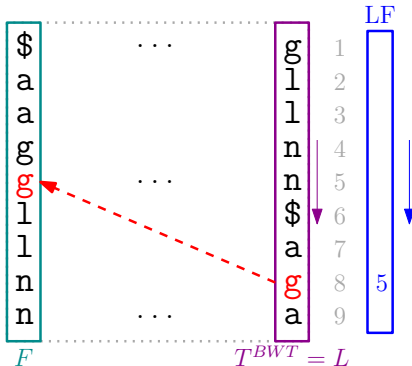
Burrows-Wheeler-Transformation

Berechnung von $LF[\cdot]$

$$T^{BWT} = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ g & l & l & n & n & \$ & a & g & a \end{matrix}$$

- $LF[\cdot]$ nur aus $T^{BWT}[\cdot]$ berechenbar

- $LF[i] = C(L[i]) + occ[i]$
 - $C(a) = \#$ Zeichen kleiner als a
 - $occ[i] = \#$ Zeichen gleich $L[i]$ in $L[1..i]$($LF[i]$ ist Position in $F[\cdot]$, an der $L[i]$ steht)



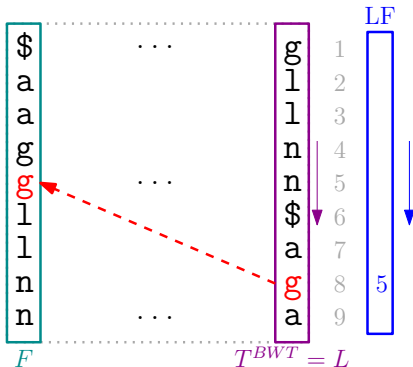
Burrows-Wheeler-Transformation

Berechnung von $LF[\cdot]$

$$T^{BWT} = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ g & l & l & n & n & \$ & a & g & a \end{matrix}$$

- $LF[\cdot]$ nur aus $T^{BWT}[\cdot]$ berechenbar

- $LF[i] = C(L[i]) + occ[i]$
 - $C(a) = \#$ Zeichen kleiner als a
 - $occ[i] = \#$ Zeichen gleich $L[i]$ in $L[1..i]$($LF[i]$ ist Position in $F[\cdot]$, an der $L[i]$ steht)



- $C(\cdot)$, $occ[\cdot]$ in $\mathcal{O}(n)$ berechenbar $\rightarrow LF[\cdot]$ auch in $\mathcal{O}(n)$ berechenbar

Burrows-Wheeler-Transformation

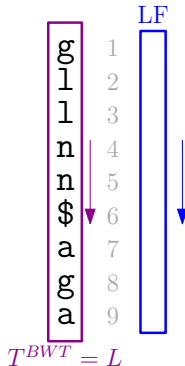
Ablauf der Berechnung von $LF[\cdot]$

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n \$ a g a

$occ =$

\$ a g l n
 $h = 0 0 0 0 0$
(zählt Zeichen)

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



Burrows-Wheeler-Transformation

Ablauf der Berechnung von $LF[\cdot]$

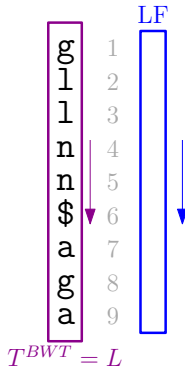
$$T^{BWT} = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \end{matrix}$$

$occ =$

$$h = \begin{matrix} \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{l} & \mathbf{n} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{matrix}$$

(zählt Zeichen)

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



Burrows-Wheeler-Transformation

Ablauf der Berechnung von $LF[\cdot]$

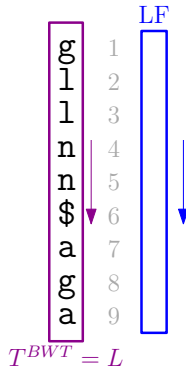
$$T^{BWT} = \begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ & \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \end{array}$$

$occ =$

$$h = \begin{array}{cccccc} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{l} & \mathbf{n} \\ & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} \end{array}$$

(zählt Zeichen)

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



Burrows-Wheeler-Transformation

Ablauf der Berechnung von $LF[\cdot]$

$$T^{BWT} = \begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ & g & l & l & n & n & \$ & a & g & a \\ occ = & 1 & & & & & & & & \end{array}$$

$$h = \begin{array}{cccccc} \$ & a & g & l & n \\ 0 & 0 & 1 & 0 & 0 \end{array}$$

(zählt Zeichen)

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



Burrows-Wheeler-Transformation

Ablauf der Berechnung von $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \\ occ = & \mathbf{1} & \mathbf{1} & & & & & & & \end{array}$$

$$\begin{array}{cccccc} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{l} & \mathbf{n} \\ h = & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{0} \\ \text{(zählt Zeichen)} & & & & & \end{array}$$

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



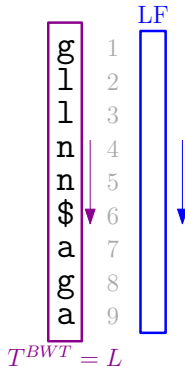
Burrows-Wheeler-Transformation

Ablauf der Berechnung von $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \\ occ = & \mathbf{1} & \mathbf{1} & \mathbf{2} & & & & & & \end{array}$$

$$\begin{array}{cccccc} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{l} & \mathbf{n} \\ h = & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{2} & \mathbf{0} \\ \text{(zählt Zeichen)} & & & & & \end{array}$$

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



Burrows-Wheeler-Transformation

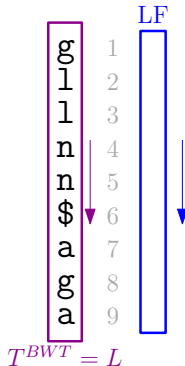
Ablauf der Berechnung von $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \\ occ = & \mathbf{1} & \mathbf{1} & \mathbf{2} & \dots & & & & & \end{array}$$

$$\begin{array}{cccccc} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{l} & \mathbf{n} \\ h = & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{2} & \mathbf{0} \end{array}$$

(zählt Zeichen)

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



Burrows-Wheeler-Transformation

Ablauf der Berechnung von $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \\ occ = & 1 & 1 & 2 & 1 & 2 & 1 & 1 & 2 & 2 \end{array}$$

$$\begin{array}{cccccc} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{l} & \mathbf{n} \\ h = & 1 & 2 & 2 & 2 & 2 \\ \text{(zählt Zeichen)} \end{array}$$

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



Burrows-Wheeler-Transformation

Ablauf der Berechnung von $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} = & \mathbf{g} & \mathbf{l} & \mathbf{l} & \mathbf{n} & \mathbf{n} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{a} \\ occ = & 1 & 1 & 2 & 1 & 2 & 1 & 1 & 2 & 2 \end{array}$$

$$\begin{array}{cccccc} & \mathbf{\$} & \mathbf{a} & \mathbf{g} & \mathbf{l} & \mathbf{n} \\ h = & 1 & 2 & 2 & 2 & 2 \\ C = & 0 & 1 & 3 & 5 & 7 \end{array}$$

(Präfixsumme von h)

- initialisiere occ und h
- laufe durch $T^{BWT} = L$ ($i = 1..n$)
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



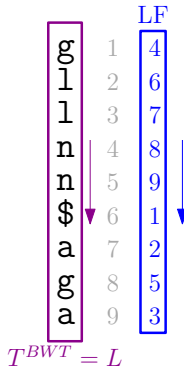
Burrows-Wheeler-Transformation

Ablauf der Berechnung von $LF[\cdot]$

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	\$	a	g	a
$occ =$	1	1	2	1	2	1	1	2	2

	\$	a	g	l	n
$h =$	1	2	2	2	2
$C =$	0	1	3	5	7

- initialisiere occ und h
- laufe durch $T^{BWT} = L (i = 1..n)$
 - $h(L[i]) = h(L[i]) + 1$
 - $occ(L[i]) = h(L[i])$
- bilde exkl. Präfixsumme von $h \rightarrow C$
- $LF[i] = C(L[i]) + occ[i]$



Burrows-Wheeler-Transformation

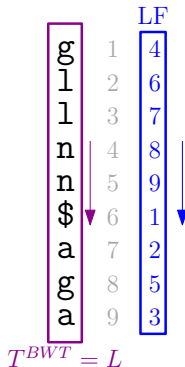
Rücktransformation mit $LF[\cdot]$

$$T^{BWT} = \begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ g & l & l & n & n & \$ & a & g & a & \end{array}$$

$T =$

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...



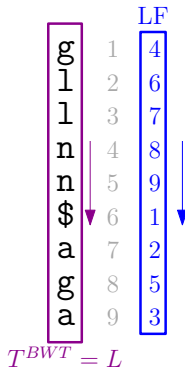
Burrows-Wheeler-Transformation

Rücktransformation mit $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & g & l & l & n & n & \$ & a & g & a \\ T & = & & & & & & & & & \$ \end{array}$$

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...



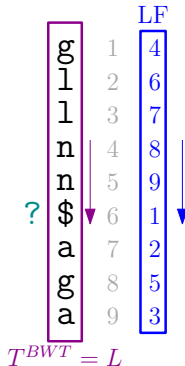
Burrows-Wheeler-Transformation

Rücktransformation mit $LF[\cdot]$

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	\$	a	g	a
$T =$?	\$	

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...



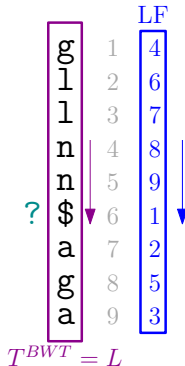
Burrows-Wheeler-Transformation

Rücktransformation mit $LF[\cdot]$

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	\$	a	g	a
$T =$?	\$	

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...



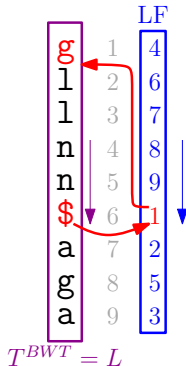
Burrows-Wheeler-Transformation

Rücktransformation mit $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & g & l & l & n & n & \$ & a & g & a \\ T & = & & & & & & & g & & \$ \end{array}$$

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...



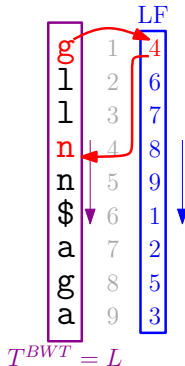
Burrows-Wheeler-Transformation

Rücktransformation mit $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & g & l & l & n & n & \$ & a & g & a \\ T & = & & & & & n & g & & & \$ \end{array}$$

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...



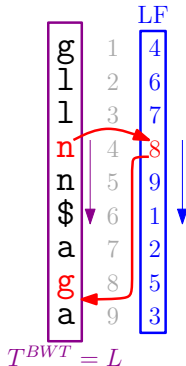
Burrows-Wheeler-Transformation

Rücktransformation mit $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & g & l & l & n & n & \$ & a & g & a \\ T & = & & & & & \dots & g & n & g & \$ \end{array}$$

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...



Burrows-Wheeler-Transformation

Rücktransformation mit $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & g & l & l & n & n & \$ & a & g & a \\ T & = & l & a & l & a & n & g & n & g & \$ \end{array}$$

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...

■ Allgemein: $T[n-i] = L[LF(LF(\dots(LF(1))\dots))] (i-1) LF$ Anwendungen

		LF
g	1	4
l	2	6
l	3	7
n	4	8
n	5	9
\$	6	1
a	7	2
g	8	5
a	9	3

$$T^{BWT} = L$$

Burrows-Wheeler-Transformation

Rücktransformation mit $LF[\cdot]$

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ T^{BWT} & = & g & l & l & n & n & \$ & a & g & a \\ T & = & l & a & l & a & n & g & n & g & \$ \end{array}$$

■ Berechnung von T von rechts nach links

- Initialisierung: $T[n] = \$ \Rightarrow LF(?) = 1$ (immer!)
- $L[1] = g \Rightarrow T[n-1] = g \Rightarrow LF(1) = 4$
- $L[4] = n \Rightarrow T[n-2] = n \Rightarrow LF(4) = 8$
- ...

■ Allgemein: $T[n-i] = L[LF(LF(\dots(LF(1))\dots))] (i-1)$ LF Anwendungen

■ Rücktransformation in $\mathcal{O}(n)$ (ohne Zusatzinformationen)

		LF
g	1	4
l	2	6
l	3	7
n	4	8
n	5	9
\$	6	1
a	7	2
g	8	5
a	9	3

$$T^{BWT} = L$$

Burrows-Wheeler-Transformation

Was bringt die BWT?

- benötigt gleichen Platz, verwendet gleiche Zeichen wie T
→ **scheinbar keine Vorteile !?**
- Permutation einfach **umkehrbar**
(benötigt keine Zusatzinformationen, $\mathcal{O}(n)$)
- Zeichen mit ähnlichem Kontext gruppiert
→ **vereinfacht Komprimierung**
- besonders gut auf Texten mit vielen gleichen Substrings
→ Beispiel: englischer Text
(u.a. viele "the", ...)

(Außerdem: Vorgänger von Suffixen einfach bestimmbar → Indizierung / Suche)

Burrows-Wheeler-Transformation

Was bringt die BWT?

- benötigt gleichen Platz, verwendet gleiche Zeichen wie T
→ **scheinbar keine Vorteile !?**
- Permutation einfach **umkehrbar**
(benötigt keine Zusatzinformationen, $\mathcal{O}(n)$)
- Zeichen mit ähnlichem Kontext gruppiert
→ **vereinfacht Komprimierung**
- besonders gut auf Texten mit vielen gleichen Substrings
→ Beispiel: englischer Text
(u.a. viele "the", ...)

(Außerdem: Vorgänger von Suffixen einfach bestimmbar → Indizierung / Suche)

Burrows-Wheeler-Transformation

Was bringt die BWT?

- benötigt gleichen Platz, verwendet gleiche Zeichen wie T
→ **scheinbar keine Vorteile !?**
- Permutation einfach **umkehrbar**
(benötigt keine Zusatzinformationen, $\mathcal{O}(n)$)
- Zeichen mit ähnlichem Kontext gruppiert
→ **vereinfacht Komprimierung**
- besonders gut auf Texten mit vielen gleichen Substrings
→ **Beispiel: englischer Text**
(u.a. viele "the", ...)

(Außerdem: Vorgänger von Suffixen einfach bestimmbar → Indizierung / Suche)

Burrows-Wheeler-Transformation

Was bringt die BWT?

- benötigt gleichen Platz, verwendet gleiche Zeichen wie T
→ **scheinbar keine Vorteile !?!**
- Permutation einfach **umkehrbar**
(benötigt keine Zusatzinformationen, $\mathcal{O}(n)$)
- Zeichen mit ähnlichem Kontext gruppiert
→ **vereinfacht Komprimierung**
- besonders gut auf Texten mit vielen gleichen Substrings
→ **Beispiel: englischer Text**
(u.a. viele "the", ...)

⋮				⋮	
t	h	e	–
t	h	e	–
t	h	e	–
⋮					⋮

(Außerdem: Vorgänger von Suffixen einfach bestimmbar → Indizierung / Suche)

Burrows-Wheeler-Transformation

Was bringt die BWT?

- benötigt gleichen Platz, verwendet gleiche Zeichen wie T
→ **scheinbar keine Vorteile !?!**
- Permutation einfach **umkehrbar**
(benötigt keine Zusatzinformationen, $\mathcal{O}(n)$)
- Zeichen mit ähnlichem Kontext gruppiert
→ **vereinfacht Komprimierung**
- besonders gut auf Texten mit vielen gleichen Substrings
→ **Beispiel: englischer Text**
(u.a. viele "the", ...)

⋮		⋮
h e	t
h e	t
h e	t
⋮		⋮
t h e	-
t h e	-
t h e	-
⋮		⋮

(Außerdem: Vorgänger von Suffixen einfach bestimmbar → Indizierung / Suche)

Burrows-Wheeler-Transformation

Was bringt die BWT?

- benötigt gleichen Platz, verwendet gleiche Zeichen wie T
→ **scheinbar keine Vorteile !?**
- Permutation einfach **umkehrbar**
(benötigt keine Zusatzinformationen, $\mathcal{O}(n)$)
- Zeichen mit ähnlichem Kontext gruppiert
→ **vereinfacht Komprimierung**
- besonders gut auf Texten mit vielen gleichen Substrings
→ **Beispiel: englischer Text**
(u.a. viele "the", ...)

T^{BWT}

⋮			⋮
h e	t	
h e	t	
h e	t	
⋮		⋮	
t h e	-	
t h e	-	
t h e	-	
⋮		⋮	

(Außerdem: Vorgänger von Suffixen einfach bestimmbar → Indizierung / Suche)

Burrows-Wheeler-Transformation

Was bringt die BWT?

- benötigt gleichen Platz, verwendet gleiche Zeichen wie T
→ **scheinbar keine Vorteile !?**
- Permutation einfach **umkehrbar**
(benötigt keine Zusatzinformationen, $\mathcal{O}(n)$)
- Zeichen mit ähnlichem Kontext gruppiert
→ **vereinfacht Komprimierung**
- besonders gut auf Texten mit vielen gleichen Substrings
→ **Beispiel: englischer Text**
(u.a. viele "the", ...)

T^{BWT}

⋮			⋮
h e	t	
h e	t	
h e	t	
⋮		⋮	
t h e	-	
t h e	-	
t h e	-	
⋮		⋮	

(Außerdem: Vorgänger von Suffixen einfach bestimmbar → Indizierung / Suche)

gegeben: Text T

gesucht : komprimierter Text C

bzip2 (1996)

- (Huffmann Kodierung)
- erzeuge *Burrows-Wheeler-Transformation*
- *Move-To-Front* (MTF) Kodierung
- Huffmann Kodierung
(eigentliche Kompression)

bzip2

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt **lokale Redundanz**
- erzeugt **kleine Zahlen** für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n a g a \$

$R =$

 1 2 3 4 5
 $Y =$ \$ a g l n

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9
g l l n n a g a \$

$R =$ 3

 1 2 3 4 5
 $Y =$ \$ a g l n

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	a	g	a	\$
$R =$	3								

	1	2	3	4	5
	\$	a	g	l	n
$Y =$	g	\$	a	l	n

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	a	g	a	\$
$R =$	3	4							

	1	2	3	4	5
	\$	a	g	l	n
$Y =$	g	\$	a	l	n

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	a	g	a	\$
$R =$	3	4							

	1	2	3	4	5
	\$	a	g	l	n
	g	\$	a	l	n
$Y =$	l	g	\$	a	n

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	a	g	a	\$
$R =$	3	4	1						

	1	2	3	4	5
$Y =$	\$	a	g	l	n
	g	\$	a	l	n
	l	g	\$	a	n

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	a	g	a	\$
$R =$	3	4	1						

	1	2	3	4	5
	\$	a	g	l	n
	g	\$	a	l	n
	l	g	\$	a	n
$Y =$	l	g	\$	a	n

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

$T^{BWT} =$

	1	2	3	4	5	6	7	8	9
	g	l	l	n	n	a	g	a	\$

$R =$

3	4	1	...
---	---	---	-----

$Y =$

					⋮
	1	2	3	4	5
	\$	a	g	l	n
	g	\$	a	l	n
	l	g	\$	a	n
	l	g	\$	a	n

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

	1	2	3	4	5	6	7	8	9
$T^{BWT} =$	g	l	l	n	n	a	g	a	\$
$R =$	3	4	1	5	1	5	4	2	5

	1	2	3	4	5
$Y =$	\$	a	g	l	n
	g	\$	a	l	n
	l	g	\$	a	n
	l	g	\$	a	n
	n	l	g	\$	a
	n	l	g	\$	a
	a	n	l	g	\$
	g	a	n	l	\$
	a	g	n	l	\$
	\$	a	g	n	l

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

$T^{BWT} =$ 1 2 3 4 5 6 7 8 9 10 11 12 13
a a a a b b b b c c c c d

$R =$

 1 2 3 4
 $Y =$ a b c d

Burrows-Wheeler-Transformation

Kompression: *Move-To-Front* (MTF) Kodierung

- nutzt lokale Redundanz
- erzeugt kleine Zahlen für gleiche Zeichen, die nahe beieinander sind

Ablauf

- initialisiere Y mit Alphabet von T^{BWT}
- durchlaufe T^{BWT} ($i = 1..n$), generiere $R[1..n]$
 - $R[i] =$ Position von $T^{BWT}[i]$ in Y
 - Schiebe $T^{BWT}[i]$ an den Anfang von Y

	1	2	3	4	5	6	7	8	9	10	11	12	13
$T^{BWT} =$	a	a	a	a	b	b	b	b	c	c	c	c	d
$R =$	1	1	1	1	2	1	1	1	3	1	1	1	4

	1	2	3	4
	a	b	c	d
	⋮			
	b	a	c	d
	⋮			
	c	b	a	d
	⋮			
$Y =$	d	c	b	a

Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

- präfixfreie Codes **variabler Länge**
- können **greedy** konstruiert werden

Ablauf

- erzeuge binären Baum *bottom-up*
 - nimm **seltenste** 2 Zeichen(-gruppen)
 - erzeuge neuen Knoten, der beide Zeichen(-gruppen) repräsentiert, neue Häufigkeit $\hat{=}$ Summe beider Häufigkeiten
- Beschriftungen der Baumkanten (links:0, rechts:1) ergeben Zeichenkodierungen

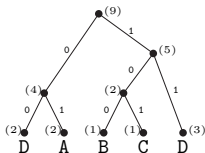
Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

- präfixfreie Codes **variabler Länge**
- können **greedy** konstruiert werden

Ablauf

- erzeuge binären Baum *bottom-up*
 - nimm **seltenste** 2 Zeichen(-gruppen)
 - erzeuge neuen Knoten, der beide Zeichen(-gruppen) repräsentiert, neue Häufigkeit $\hat{=}$ Summe beider Häufigkeiten
- Beschriftungen der Baumkanten (links:0, rechts:1) ergeben Zeichenkodierungen



Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$$\begin{array}{l} T = 1 \ a \ 1 \ a \ n \ g \ n \ g \ \$ \\ R = 3 \ 4 \ 1 \ 5 \ 1 \ 5 \ 4 \ 2 \ 5 \end{array}$$

Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$T = 1 a 1 a n g n g \$$
 $R = 3 4 1 5 1 5 4 2 5$

Symbol	Häufigkeit
1	2
2	1
3	1
4	2
5	3

Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$T = 1 a 1 a n g n g \$$
 $R = 3 4 1 5 1 5 4 2 5$

Symbol	Häufigkeit
1	2
2	1
3	1
4	2
5	3

(2)•
4

(2)•
1

(1)•
2

(1)•
3

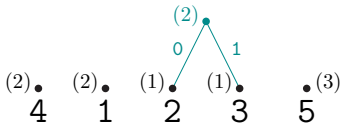
•(3)
5

Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$T = 1 \ a \ 1 \ a \ n \ g \ n \ g \ \$$
 $R = 3 \ 4 \ 1 \ 5 \ 1 \ 5 \ 4 \ 2 \ 5$

Symbol	Häufigkeit
1	2
2	1
3	1
4	2
5	3

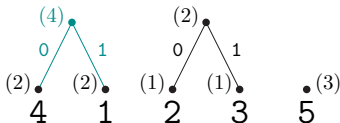


Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$T = 1 \ a \ 1 \ a \ n \ g \ n \ g \ \$$
 $R = 3 \ 4 \ 1 \ 5 \ 1 \ 5 \ 4 \ 2 \ 5$

Symbol	Häufigkeit
1	2
2	1
3	1
4	2
5	3

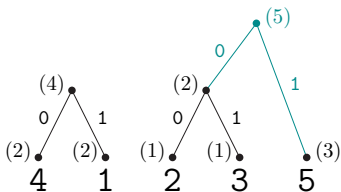


Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$T = 1 \ a \ 1 \ a \ n \ g \ n \ g \ \$$
 $R = 3 \ 4 \ 1 \ 5 \ 1 \ 5 \ 4 \ 2 \ 5$

Symbol	Häufigkeit
1	2
2	1
3	1
4	2
5	3

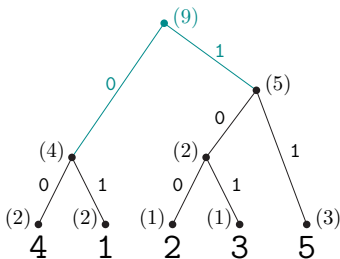


Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$T = 1 \ a \ 1 \ a \ n \ g \ n \ g \ \$$
 $R = 3 \ 4 \ 1 \ 5 \ 1 \ 5 \ 4 \ 2 \ 5$

Symbol	Häufigkeit
1	2
2	1
3	1
4	2
5	3

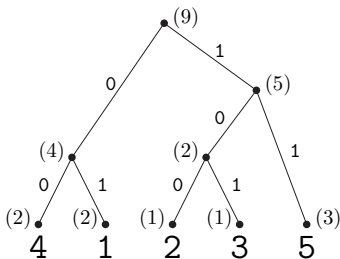


Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$T = 1 \ a \ 1 \ a \ n \ g \ n \ g \ \$$
 $R = 3 \ 4 \ 1 \ 5 \ 1 \ 5 \ 4 \ 2 \ 5$

Symbol	Häufigkeit	Code
1	2	01
2	1	100
3	1	101
4	2	00
5	3	11



Burrows-Wheeler-Transformation

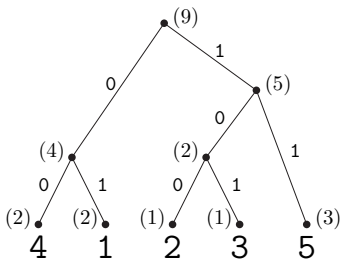
Kompression: Huffman Kodierung

$T = 1 a 1 a n g n g \$$

$R = 3 4 1 5 1 5 4 2 5$

101 00 01 11 01 11 00 100 11 20 Bits

Symbol	Häufigkeit	Code
1	2	01
2	1	100
3	1	101
4	2	00
5	3	11

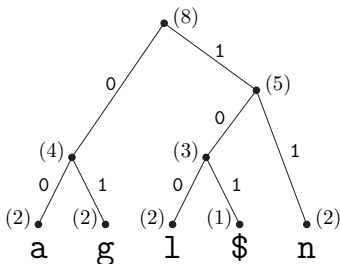


Burrows-Wheeler-Transformation

Kompression: Huffman Kodierung

$T = l a l a n g n g \$$
100 00 100 00 11 01 11 01 101 21 Bits

Symbol	Häufigkeit	Code
\$	1	101
a	2	00
g	2	01
l	2	100
n	2	11



Burrows-Wheeler-Transformation

Zusammenfassung

- erzeugt (sinnvolle) **Permutation** der Eingabe
(gruppiert Zeichen mit ähnlichem Kontext nahe beieinander)
- **keine Zusatzinformation** für Rücktransformation nötig
(alle Informationen in Struktur der Permutation)
- Hin- und Rücktransformation in $\mathcal{O}(n)$
(einfache Papier-und-Bleistift-Methode existiert auch)
- **Vorverarbeitung** (statischer) Texte
(Komprimierung, Indizierung, Suche)

Suche in der Burrows-Wheeler-Transformation

Ferragina &Manzini (2000)

- Index basierend auf der Burrows-Wheeler-Transformation (BWT)
- Vergleich des Musters von rechts nach links
- Zeitkomplexität: $\mathcal{O}(m \log \sigma)$

BWT

- $BWT[i] = \mathcal{T}[SA[i] - 1 \bmod n]$
- Unkomprimierte Größe: $n \log \sigma$ Bits
- Komprimierte Größe: $nH_k(\mathcal{T})$ Bits (+Kontextinformation)

Backward Search

i	$SA[i]$	BWT	$\mathcal{T}[SA[i]..n-1]$
0	18	a	\$
1	17	r	a\$
2	10	r	abarbara\$
3	7	d	abrabarbara\$
4	0	\$	abracadabrabarbara\$
5	3	r	acadabrabarbara\$
6	5	c	adabrabarbara\$
7	15	b	ara\$
8	12	b	arbara\$
9	14	r	bara\$
10	11	a	barbara\$
11	8	a	brabarbara\$
12	1	a	bracadabrabarbara\$
13	4	a	cadabrabarbara\$
14	6	a	dabrabarbara\$
15	16	a	ra\$
16	9	b	rabarbara\$
17	2	b	racadabrabarbara\$
18	13	a	rbara\$

- $BWT[i] = \mathcal{T}[SA[i] - 1]$, for $SA[i] > 0$
- $BWT[i] = \mathcal{T}[n - 1]$, for $SA[i] = 0$
- I.e. $BWT[i]$ is the character preceding suffix $SA[i]$

Backward Search

i	BWT	$\mathcal{T}[SA[i..n-1]]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abababababara\$
5	r	abababababara\$
6	c	ababababara\$
7	b	abara\$
8	b	ababara\$
9	r	ababara\$
10	a	abababara\$
11	a	ababababara\$
12	a	abababababara\$
13	a	abababababara\$
14	a	abababababara\$
15	a	abara\$
16	b	abababara\$
17	b	abababababara\$
18	a	ababara\$

Array C contains for each $c \in \Sigma$ the position of the first suffix in SA which starts with c :

\$	a	b	c	d	r	r+1
0	1	9	13	14	15	19

- Operation $rank(i, X, BWT)$ returns how often character $X \in \Sigma$ occurs in the prefix $BWT[0..i-1]$.
- Example: search for $\mathcal{P} = bar$.

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	ababababababara\$
5	r	abababababara\$
6	c	abababababara\$
7	b	ababara\$
8	b	ababara\$
9	r	ababara\$
10	a	ababara\$
11	a	abababababara\$
12	a	ababababababara\$
13	a	ababababababara\$
14	a	abababababara\$
15	a	ababara\$
16	b	abababababara\$
17	b	ababababababara\$
18	a	ababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for bar .
- Initial interval: $[sp_0, ep_0] = [0..n-1]$
- Determine interval for r :
 $sp_1 = C[r] + rank(sp_0, r, BWT)$
 $ep_1 = C[r] + rank(ep_0 + 1, r, BWT) - 1$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	abarbara\$
3	d	abrabarbara\$
4	\$	abracadabrabarbara\$
5	r	acadabrabarbara\$
6	c	adabrabarbara\$
7	b	ara\$
8	b	arbara\$
9	r	bara\$
10	a	barbara\$
11	a	brabarbara\$
12	a	bracadabrabarbara\$
13	a	cadabrabarbara\$
14	a	dabrabarbara\$
15	a	ra\$
16	b	rabarbara\$
17	b	racadabrabarbara\$
18	a	rbara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Initial interval: $[sp_0, ep_0] = [0..n-1]$
- Determine interval for *r*:

$$sp_1 = 15 + \text{rank}(0, r, \text{BWT})$$

$$ep_1 = 15 + \text{rank}(19, r, \text{BWT})$$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	ababababababara\$
5	r	abababababara\$
6	c	abababababara\$
7	b	ababara\$
8	b	ababara\$
9	r	ababara\$
10	a	ababara\$
11	a	abababababara\$
12	a	ababababababara\$
13	a	ababababababara\$
14	a	abababababara\$
15	a	ababara\$
16	b	abababababara\$
17	b	ababababababara\$
18	a	ababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Initial interval: $[sp_0, ep_0] = [0..n-1]$
- Determine interval for *r*:
 $sp_1 = 15+0$
 $ep_1 = 15 + \text{rank}(19, r, \text{BWT}) - 1$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	ababababababara\$
5	r	ababababababara\$
6	c	ababababababara\$
7	b	abababababara\$
8	b	abababababara\$
9	r	abababababara\$
10	a	abababababara\$
11	a	abababababara\$
12	a	abababababara\$
13	a	abababababara\$
14	a	abababababara\$
15	a	abababababara\$
16	b	abababababara\$
17	b	abababababara\$
18	a	abababababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Initial interval: $[sp_0, ep_0] = [0..n-1]$
- Determine interval for *r*:
 $sp_1 = 15 + 0 = 15$
 $ep_1 = 15 + 4 - 1 = 18$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abababababara\$
5	r	abababababara\$
6	c	abababababara\$
7	b	ababara\$
8	b	ababara\$
9	r	ababara\$
10	a	ababara\$
11	a	ababababara\$
12	a	abababababara\$
13	a	abababababara\$
14	a	ababababara\$
15	a	abara\$
16	b	abababara\$
17	b	abababababara\$
18	a	ababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_1, ep_1] = [15..18]$
- Determine interval for *ar*:
 $sp_2 = C[a] + rank(sp_1, a, BWT)$
 $ep_2 = C[a] + rank(ep_1 + 1, a, BWT) - 1$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abababababara\$
5	r	abababababara\$
6	c	abababababara\$
7	b	ababara\$
8	b	ababara\$
9	r	ababara\$
10	a	ababara\$
11	a	ababababara\$
12	a	abababababara\$
13	a	abababababara\$
14	a	ababababara\$
15	a	ababara\$
16	b	abababara\$
17	b	abababababara\$
18	a	ababara\$

						C
\$	a	b	c	d	r	
0	1	9	13	14	15	

- Search backwards for *bar*.
- Interval: $[sp_1, ep_1] = [15..18]$
- Determine interval for *ar*:
 $sp_2 = 1 + rank(15, a, BWT)$
 $ep_2 = 1 + rank(ep_1, a, BWT)$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abababababara\$
5	r	abababababara\$
6	c	abababababara\$
7	b	ababara\$
8	b	ababara\$
9	r	abara\$
10	a	ababara\$
11	a	abababara\$
12	a	abababababara\$
13	a	abababababara\$
14	a	abababababara\$
15	a	abara\$
16	b	abababara\$
17	b	abababababara\$
18	a	ababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_1, ep_1] = [15..18]$
- Determine interval for *ar*:
 $sp_2 = 1 + \text{rank}(15, a, \text{BWT})$
 $ep_2 = 1 + \text{rank}(ep_1, a, \text{BWT})$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abababababara\$
5	r	abababababara\$
6	c	abababababara\$
7	b	ababababara\$
8	b	abababara\$
9	r	ababara\$
10	a	ababara\$
11	a	abababara\$
12	a	abababababara\$
13	a	abababababara\$
14	a	abababababara\$
15	a	ababababara\$
16	b	ababababara\$
17	b	abababababara\$
18	a	ababababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_1, ep_1] = [15..18]$
- Determine interval for *ar*:
 $sp_2 = 1 + 6$
 $ep_2 = 1 + \text{rank}(19, a, \text{BWT}) - 1$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	ababababababara\$
5	r	ababababababara\$
6	c	ababababababara\$
7	b	abababababara\$
8	b	abababababara\$
9	r	abababababara\$
10	a	abababababara\$
11	a	abababababara\$
12	a	abababababara\$
13	a	abababababara\$
14	a	abababababara\$
15	a	abababababara\$
16	b	abababababara\$
17	b	abababababara\$
18	a	abababababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_1, ep_1] = [15..18]$
- Determine interval for *ar*:
 $sp_2 = 1 + 6 = 7$
 $ep_2 = 1 + 8 - 1 = 8$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abababababara\$
5	r	abababababara\$
6	c	abababababara\$
7	b	ababababara\$
8	b	ababababara\$
9	r	ababababara\$
10	a	ababababara\$
11	a	ababababara\$
12	a	ababababara\$
13	a	ababababara\$
14	a	ababababara\$
15	a	ababababara\$
16	b	ababababara\$
17	b	ababababara\$
18	a	ababababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_2, ep_2] = [7..8]$
- Determine interval for *bar*:
 $sp_3 = C[b] + rank(sp_2, b, BWT)$
 $ep_3 = C[b] + rank(ep_2 + 1, b, BWT) - 1$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abracadabababara\$
5	r	acadabababara\$
6	c	adabababara\$
7	b	ababara\$
8	b	ababara\$
9	r	ababara\$
10	a	ababara\$
11	a	abababara\$
12	a	abracadabababara\$
13	a	acadabababara\$
14	a	adabababara\$
15	a	ababara\$
16	b	abababara\$
17	b	abracadabababara\$
18	a	ababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_2, ep_2] = [7..8]$
- Determine interval for *bar*:
 $sp_3 = 9 + rank(7, b, BWT)$
 $ep_3 = 9 + rank(ep_1, b, BWT)$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	ababababababara\$
5	r	abababababara\$
6	c	abababababara\$
7	b	ababara\$
8	b	ababara\$
9	r	ababara\$
10	a	ababara\$
11	a	abababara\$
12	a	ababababababara\$
13	a	abababababara\$
14	a	abababababara\$
15	a	abara\$
16	b	abababara\$
17	b	ababababababara\$
18	a	ababara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_2, ep_2] = [7..8]$
- Determine interval for *bar*:
 $sp_3 = 9 + \text{rank}(7, b, \text{BWT})$
 $ep_3 = 9 + \text{rank}(ep_1, b, \text{BWT})$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abracadabababara\$
5	r	acadabababara\$
6	c	adabababara\$
7	b	abara\$
8	b	ababara\$
9	r	abara\$
10	a	ababara\$
11	a	abababara\$
12	a	abracadabababara\$
13	a	acadabababara\$
14	a	adabababara\$
15	a	abara\$
16	b	abababara\$
17	b	abracadabababara\$
18	a	abara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_2, ep_2] = [7..8]$
- Determine interval for *bar*:
 $sp_3 = 9+0$
 $ep_3 = 9 + \text{rank}(9, b, \text{BWT}) - 1$

Backward Search

i	BWT	$\mathcal{T}[SA[i]..n-1]$
0	a	\$
1	r	a\$
2	r	ababara\$
3	d	abababara\$
4	\$	abracadabababara\$
5	r	acadabababara\$
6	c	adabababara\$
7	b	ara\$
8	b	arabara\$
9	r	bara\$
10	a	barabara\$
11	a	brababara\$
12	a	bracadabababara\$
13	a	cadabababara\$
14	a	dabababara\$
15	a	ra\$
16	b	rababara\$
17	b	racadabababara\$
18	a	rbara\$

C					
\$	a	b	c	d	r
0	1	9	13	14	15

- Search backwards for *bar*.
- Interval: $[sp_2, ep_2] = [7..8]$
- Determine interval for *bar*:
 $sp_3 = 9 + 0 = 9$
 $ep_3 = 9 + 2 - 1 = 10$

- Only C and a data structure R supporting the *rank* operation on BWT are required for existence and count queries.
- Space: $\sigma \log n$ bits for C + space for R
- Time: $\mathcal{O}(m \cdot t_{rank})$, where t_{rank} is time for one rank operation. Independent from n ?
- Next: How to implement *rank*?

Rank operation

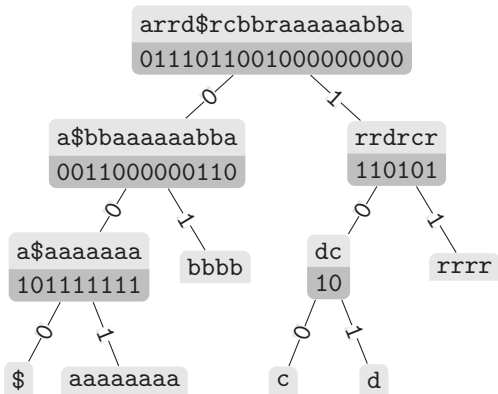
- Constant time and $o(n)$ extra space solution on bitvectors (Jacobson 1989)
- Solution on general sequences: Wavelet Tree (Grossi & Vitter 2003)

- Only C and a data structure R supporting the *rank* operation on BWT are required for existence and count queries.
- Space: $\sigma \log n$ bits for C + space for R
- Time: $\mathcal{O}(m \cdot t_{rank})$, where t_{rank} is time for one rank operation. Independent from n ? If t_{rank} is independent from n
- Next: How to implement *rank*?

Rank operation

- Constant time and $o(n)$ extra space solution on bitvectors (Jacobson 1989)
- Solution on general sequences: Wavelet Tree (Grossi & Vitter 2003)

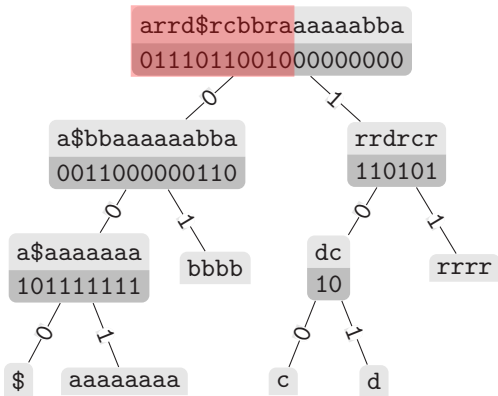
Wavelet Tree Example: Calculate Rank



$a = 001$

$$\text{rank}(11, a, WT) = \text{rank}(\text{rank}(\text{rank}(11, 0, b_e) = 5, 0, b_0) = 3, 1, b_{00}) = 2$$

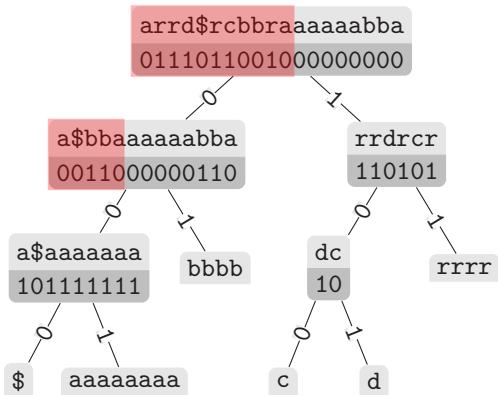
Wavelet Tree Example: Calculate Rank



$a = 001$

$$\text{rank}(11, a, WT) = \text{rank}(\text{rank}(\text{rank}(11, 0, b_e) = 5, 0, b_0) = 3, 1, b_{00}) = 2$$

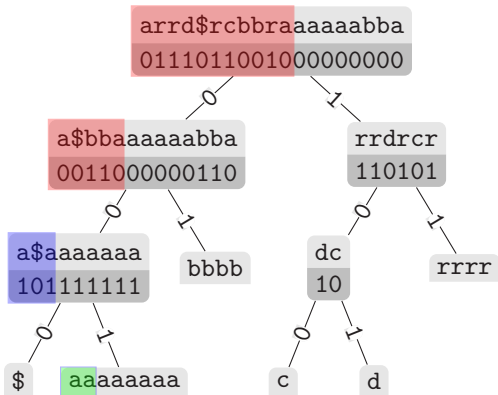
Wavelet Tree Example: Calculate Rank



$a = 001$

$$\text{rank}(11, a, WT) = \text{rank}(\text{rank}(\text{rank}(11, 0, b_e) = 5, 0, b_0) = 3, 1, b_{00}) = 2$$

Wavelet Tree Example: Calculate Rank



$a = 001$

$$\text{rank}(11, a, WT) = \text{rank}(\text{rank}(\text{rank}(11, 0, b_e) = 5, 0, b_0) = 3, 1, b_{00}) = 2$$