

Projekt Algorithm Engineering

Themenvorstellung

Uwe Baier

Institut für theoretische Informatik
Universität Ulm

Sommersemester 2020

Inhaltsübersicht

Einführung

Thema 1: Runminimierung in einer Multi-String BWT

Zusammenfassung und Ziele

Burrows Wheeler Transformation

Geschichte

- ▶ reversible Texttransformation [Burrows and Wheeler, 1994]
- ▶ ehemaliges Einsatzgebiet: Datenkompression

BWT heute

- ▶ zum Teil noch Einsatz in der Datenkompression, z.B. `bzip2` [Seward, 1996]
Problem: langsame Dekompressionsgeschwindigkeit
- ▶ Vielfältiger Einsatz in der Bioinformatik, z.B. `BWA` für read alignment [Li and Durbin, 2010]
- ▶ Einsatz in komprimierter Volltextindizierung, z.B. FM-Index [Ferragina and Manzini, 2005]

BWT in der Zukunft

- ▶ Fortschritte in der Volltextindizierung [Gagie et al., 2018]
- ▶ Fortschritte bei Kompressionsraten [Baier, 2018]

BWT

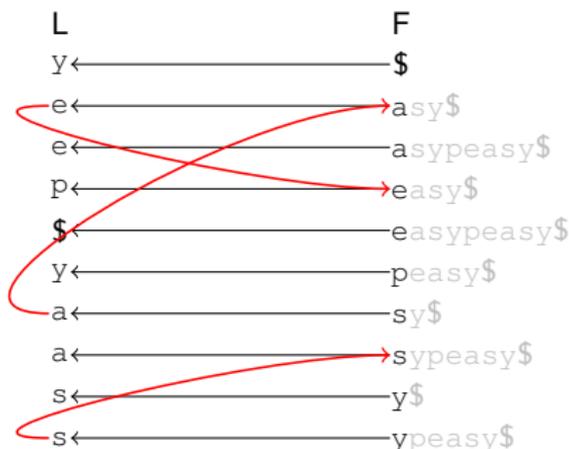
“Die BWT L ist ein String, der durch das konkatenieren der zyklisch vorhergehenden Zeichen aller lexikographisch sortierten Suffixe eines Strings S entsteht.”

BWT Generierung für $S = \text{easyeasy}\$$

	Suffixe		L	sortierte Suffixe
y	\$		y	\$
s	y\$		e	asy\$
a	sy\$		e	asypeasy\$
e	asy\$	sortieren	p	easy\$
p	easy\$	→	\$	easypeasy\$
y	peasy\$		y	peasy\$
s	ypeasy\$		a	sy\$
a	sypeasy\$		a	sypeasy\$
e	asypeasy\$		s	y\$
\$	easypeasy\$		s	ypeasy\$

BWT - Rückwärtsschritt

- ▶ F - Spalte: erste Spalte der sortierten Suffixe
Berechenbar durch die Sortierung der Zeichen in L
- ▶ k -tes Vorkommen des Zeichens c in L korrespondiert mit k -tem Vorkommen des Zeichens c in F (**LF-mapping**)



- ▶ Eine Tour, bei der dem LF-mapping gefolgt und die besuchten Buchstaben in L nacheinander aufgeschrieben werden ergibt den reversen Ausgangsstring

Multi-String BWT

- ▶ Mehrere Strings S_1, \dots, S_m werden mit Trennsymbol \$ zu neuem String $S := S_1\$S_2\$ \dots \$S_m\$$ konkateniert
- ▶ sortierte Suffixe werden nur bis zum Trennsymbol betrachtet
- ▶ restliche Definition und Rückwärtsschritt analog

Stringkonkatenation

$S_1 = GCA$ $S_2 = GGTGA$

$S_3 = GGTGC$ $S_4 = GG$

$S = GCA\$GGTGA\$GGTGC\$GG\$$

Multi-String BWT

L	sortierte Suffixe
A	\$
A	\$
C	\$
G	\$
C	A\$
G	A\$
G	C\$
G	CA\$
G	G\$
T	GA\$
T	GC\$
\$	GCA\$
\$	GG\$
\$	GGTGA\$
\$	GGTGC\$
G	GTGA\$
G	GTGC\$
G	TGA\$
G	TGC\$

Thema 1: Runminimierung in einer Multi-String BWT

- ▶ Run: längenmaximaler Substring mit Wiederholung desselben Zeichens

... CC AAAAAAAAAAAA TT ...
Run mit 11 Zeichen

- ▶ Run durch Lauflängenkodierung komprimierbar

AAAAAA → A10
 $6 = 110_2$ $110_2 = 6$ A's

- ▶ Komprimierbarkeit der BWT hängt von der Anzahl an Runs ab

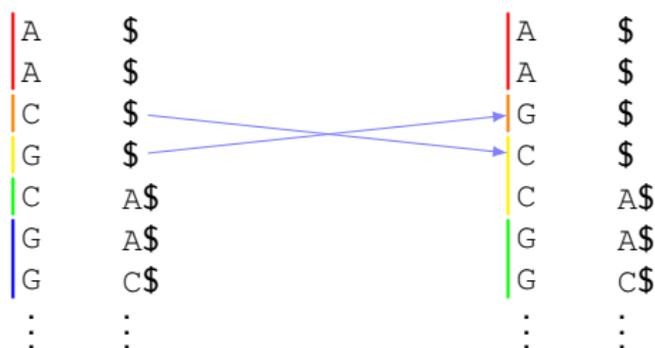
[Cazaux and Rivals, 2018]

In Multi-String BWT: Zeichen mit gleichem sortiertem Suffix können vertauscht werden

L	sortierte Suffixe
A	\$
A	\$
C	\$
G	\$
C	A\$
G	A\$
G	C\$
G	CA\$
G	G\$
T	GA\$
T	GC\$
\$	GCA\$
\$	GG\$
\$	GGTGA\$
\$	GGTGC\$
G	GTGA\$
G	GTGC\$
G	TGA\$
G	TGC\$

Thema 1: Runminimierung in einer Multi-String BWT

- ▶ Ziel: Zeichen mit gleichem Suffix derart vertauschen, dass die Anzahl der Runs minimiert wird



- ▶ Eingabe: BWT in Form eines Wavelet Tree [Foschini et al., 2006] sowie Bitvektor mit Markierungen bei Suffixtransitionen
- ▶ Algorithmik mittels einfacher Fallunterscheidung [Cazaux and Rivals, 2018]

Rahmenbedingungen

Bereitgestellte Materialien

- ▶ Projektkonzept
 - ▶ Grundlagen
 - ▶ Beschreibung sowie Pseudocode der zu implementierenden Algorithmen
- ▶ Quellcode als Programmiergrundlage (C++)
 - ▶ Enthält alle benötigten Basisoperationen
 - ▶ Beispiele, wie die Operationen benutzt werden

Aufgaben

- ▶ Implementierung (C++)
- ▶ Durchführen von Experimenten
- ▶ Anfertigen einer Ausarbeitung
 - ▶ Problemstellung, Lösung etc. in eigenen Worten
 - ▶ Darstellung der Experimente und derer Ergebnisse
- ▶ Projektvortrag

Zusammenfassung

- ▶ Implementierung + Experimente + Ausarbeitung + Vortrag
- ▶ Eigenständiges Arbeiten
- ▶ Hilfestellung bei Bedarf
- ▶ Auch mehrfache Themenvergabe möglich

Interesse? Einfach melden:
`uwe.baier@uni-ulm.de`

References I



Baier, U. (2018).

On Undetected Redundancy in the Burrows-Wheeler Transform.

In [Annual Symposium on Combinatorial Pattern Matching \(CPM 2018\)](#), CPM '18, pages 3:1–3:15.



Burrows, M. and Wheeler, D. J. (1994).

A block-sorting lossless data compression algorithm.

Technical Report 124, Digital Equipment Corporation.



Cazaux, B. and Rivals, E. (2018).

Strong link between BWT and XBW via Aho-Corasick automaton and applications to Run-Length Encoding.

[CoRR](#), abs/1805.10070.



Ferragina, P. and Manzini, G. (2005).

Indexing Compressed Text.

[Journal of the ACM](#), 52(4):552–581.

References II



Foschini, L., Grossi, R., Gupta, A., and Vitter, J. S. (2006).

When Indexing Equals Compression: Experiments with Compressing Suffix Arrays and Applications.

[ACM Transactions on Algorithms](#), 2(4):611–639.



Gagie, T., Navarro, G., and Prezza, N. (2018).

Optimal-Time Text Indexing in BWT-runs Bounded Space.

In [Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms](#), SODA '18, pages 1459–1477.



Li, H. and Durbin, R. (2010).

Fast and accurate long-read alignment with Burrows–Wheeler transform.

[Bioinformatics](#), 26(5):589–595.



Seward, J. (1996).

bzip2 File Compressor.

<http://bzip.org/>.

last visited January 2018.