

Move-to-Front, Distance Coding and Inversion Frequencies Revisited

Travis Gagie
and
Giovanni Manzini

CPM '07

<http://www.mfn.unipmn.it/~manzini/cpm07>

~~Move to Front,~~
Distance Coding
and
~~Inversion Frequencies~~
Revisited

Travis Gagie
and
Giovanni Manzini

CPM '07

<http://www.mfn.unipmn.it/~manzini/cpm07>

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

Computing the BWT

abracadabra

Computing the BWT

abracadabra\$

Computing the BWT

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

abracadabra\$

Computing the BWT

abracadabra\$		abracadabra\$
abracadabra\$		bracadabra\$a
abracadabra\$		racadabra\$ab
abracadabra\$		acadabra\$abr
abracadabra\$		cadabra\$abra
abracadabra\$	⇒	adabra\$abrac
abracadabra\$		dabra\$abraca
abracadabra\$		abra\$abracad
abracadabra\$		bra\$abracada
abracadabra\$		ra\$abracadab
abracadabra\$		a\$abracadabr
abracadabra\$		\$abracadabra

Computing the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a
abracadabra\$		bracadabra\$a		racadabra\$a b
abracadabra\$		racadabra\$ab		abra\$abraca d
abracadabra\$		acadabra\$abr		ra\$abracada b
abracadabra\$		cadabra\$abra		adabra\$abra c
abracadabra\$	⇒	adabra\$abrac	⇒	abracadabra \$
abracadabra\$		dabra\$abraca		acadabra\$ab r
abracadabra\$		abra\$abracad		a\$abracadab r
abracadabra\$		bra\$abracada		dabra\$abrac a
abracadabra\$		ra\$abracadab		bra\$abracad a
abracadabra\$		a\$abracadabr		cadabra\$abr a
abracadabra\$		\$abracadabra		\$abracadabr a

Computing the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	⇒	adabra\$abrac	⇒	abracadabra \$	\$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Computing the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	⇒	adabra\$abrac	⇒	abracadabra \$	⇒ \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	⇒	adabra\$abrac	⇒	abracadabra \$	\$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	⇒	adabra\$abrac	⇒	abracadabra \$	\$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$	a	a
abracadabra\$		bracadabra\$a		racadabra\$a	b	b
abracadabra\$		racadabra\$ab		abra\$abraca	d	d
abracadabra\$		acadabra\$abr		ra\$abracada	b	b
abracadabra\$		cadabra\$abra		adabra\$abra	c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow	\$
abracadabra\$		dabra\$abraca		acadabra\$ab	r	r
abracadabra\$		abra\$abracad		a\$abracadab	r	r
abracadabra\$		bra\$abracada		dabra\$abrac	a	a
abracadabra\$		ra\$abracadab		bra\$abracad	a	a
abracadabra\$		a\$abracadabr		cadabra\$abr	a	a
abracadabra\$		\$abracadabra		\$abracadabr	a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabr a	a

Inverting the BWT

abracadabra\$		abracadabra\$		bracadabra\$ a	a
abracadabra\$		bracadabra\$a		racadabra\$a b	b
abracadabra\$		racadabra\$ab		abra\$abraca d	d
abracadabra\$		acadabra\$abr		ra\$abracada b	b
abracadabra\$		cadabra\$abra		adabra\$abra c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow \$
abracadabra\$		dabra\$abraca		acadabra\$ab r	r
abracadabra\$		abra\$abracad		a\$abracadab r	r
abracadabra\$		bra\$abracada		dabra\$abrac a	a
abracadabra\$		ra\$abracadab		bra\$abracad a	a
abracadabra\$		a\$abracadabr		cadabra\$abr a	a
abracadabra\$		\$abracadabra		\$abracadabra	a

A useful property of the BWT

abracadabra\$		abracadabra\$		bracadabra\$	a	a
abracadabra\$		bracadabra\$a		racadabra\$a	b	b
abracadabra\$		racadabra\$ab		abra\$abraca	d	d
abracadabra\$		acadabra\$abr		ra\$abracada	b	b
abracadabra\$		cadabra\$abra		adabra\$abra	c	c
abracadabra\$	\Rightarrow	adabra\$abrac	\Rightarrow	abracadabra \$	\Leftrightarrow	\$
abracadabra\$		dabra\$abraca		acadabra\$ab	r	r
abracadabra\$		abra\$abracad		a\$abracadab	r	r
abracadabra\$		bra\$abracada		dabra\$abrac	a	a
abracadabra\$		ra\$abracadab		bra\$abracad	a	a
abracadabra\$		a\$abracadabr		cadabra\$abr	a	a
abracadabra\$		\$abracadabra		\$abracadabr	a	a

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

Computing DC

abdbcrraaaa

Computing DC

abdbccraaaa



first character: 1

Computing DC

ab**b**crraaaa



first character: 1

other characters: 2, 3, 5, 6

Computing DC

~~abdb~~raaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7

Computing DC

ab**b**ccrra~~aaa~~



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2

Computing DC

abdbcrraaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

Computing DC

abdbccraaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4

Computing DC

abdbccraaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4



1, 2, 3, 5, 6, 7, 2, 1, 1, 1, 1, 1, 4

Inverting DC

abdbccraaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4



1, 2, 3, 5, 6, 7, 2, 1, 1, 1, 1, 1, 4

Inverting DC

abdbccraaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4



1, 2, 3, 5, 6, 7, 2, 1, 1, 1, 1, 1, 4

Inverting DC

ab**b**cr^aaaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4



1, 2, 3, 5, 6, 7, 2, 1, 1, 1, 1, 1, 4

Inverting DC

~~abd~~ cr aaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4



1, 2, 3, 5, 6, 7, 2, 1, 1, 1, 1, 1, 4

Inverting DC

ab**ed**cr~~ra~~aaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4



1, 2, 3, 5, 6, 7, 2, 1, 1, 1, 1, 1, 4

Inverting DC

abdbcr aaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4



1, 2, 3, 5, 6, 7, 2, 1, 1, 1, 1, 1, 4

Inverting DC

abdbcrraaaa



first character: 1

other characters: 2, 3, 5, 6

distances: 7, 2, 1, 1, 1, 1, 1

last run: 4



1, 2, 3, 5, 6, 7, 2, 1, 1, 1, 1, 1, 4

Two useful properties of DC

1. Suppose P_{fx} is a prefix-free code for the positive integers such that $|P_{fx}(i)| \leq a \log i + b$ for some constants a and b ; then

$$|P_{fx}(DC(s))| \leq a (H_0(s)n + \log n) + b|DC(s)| + g$$

for some g depending only on the alphabet size.

2. $|DC(s)|$ is the number of maximal runs in s plus a term depending only on the alphabet size.

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

Manzini, 2001

Theorem (Man01)

For any string s of length n and $k \geq 0$,

$$|\text{BW0}(s)| \leq 8H_k(s)n + (0.08 + \mu)n + (1 + o(1)) \log n + g_k$$

for some g_k depending only on the alphabet size and k .

Manzini, 2001

Theorem (Man01)

For any string s of length n and $k \geq 0$,

$$|\text{BW}_{0\text{RL}}(s)| \leq (5 + 3\mu)H_k^*(s)n + (1 + o(1)) \log n + g_k$$

for some g_k depending only on the alphabet size and k .

Kaplan, Landau and Verbin, 2006

Theorem (KLV06)

For any string s of length n , $k \geq 0$ and $\lambda > 1$,

$$|\text{BW0}(s)| \leq \lambda H_k(s)n + (\log \zeta(\lambda) + \mu) n + (1 + o(1)) \log n + \lambda g_k$$

for some g_k depending only on the alphabet size and k .

Kaplan, Landau and Verbin, 2006

Theorem (KLV06)

For any string s of length n and $k \geq 0$,

$$|\text{BW0}(s)| \leq 4.45H_k(s)n + (0.08 + \mu)n + (1 + o(1)) \log n + g_k$$

for some g_k depending only on the alphabet size and k .

Kaplan, Landau and Verbin, 2006

Theorem (KLV06)

For any string s of length n , $k \geq 0$ and $\lambda > 1$,

$$\begin{aligned} |\text{Order0}(\text{DC}(\text{BWT}(s)))| \\ \leq \lambda H_k(s)n + (\log(\zeta(\lambda)) - 1) n + (1 + o(1)) \log n + \lambda g_k \end{aligned}$$

for some g_k depending only on the alphabet size and k .

Kaplan, Landau and Verbin, 2006

Theorem (KLV06)

For any string s of length n and $k \geq 0$,

$$|\text{Order}_0(\text{DC}(\text{BWT}(s)))| \leq 1.73H_k(s)n + \mu n + (1 + o(1)) \log n + g_k$$

for some g_k depending only on the alphabet size and k .

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

*H vs. H**

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

H, the empirical entropy

Definition (see Man01 for details)

The *k*th-order empirical entropy of a string *s* of length *n* is

$$H_k(s) = \begin{cases} (1/n) \sum_i n_i \log(n/n_i) & \text{if } k = 0, \\ (1/n) \sum_{|w|=k} H_0(w_s) & \text{if } k \geq 1. \end{cases}$$

H , the empirical entropy

Example

$$\begin{aligned} H_0(\text{abracadabra}) \\ = \frac{1}{11} \left(5 \log \frac{11}{5} + 2 \log \frac{11}{2} + \log 11 + \log 11 + 2 \log \frac{11}{2} \right) \end{aligned}$$

H , the empirical entropy

Example

$$H_0(\text{abracadabra}) \approx 2.04,$$

$$H_1(\text{abracadabra})$$

$$= \frac{1}{11} (4H_0(\text{bcdb}) + 2H_0(\text{rr}) + H_0(\text{a}) + H_0(\text{a}) + 2H_0(\text{aa}))$$

H , the empirical entropy

Example

$$H_0(\text{abracadabra}) \approx 2.04,$$

$$H_1(\text{abracadabra}) \approx 0.55,$$

$$H_2(\text{abracadabra}) = 0.$$

H , the empirical entropy

Example

$$H_0(\text{abracadabra}) \approx 2.04,$$

$$H_1(\text{abracadabra}) \approx 0.55,$$

$$H_2(\text{abracadabra}) = 0.$$

Example

$$H_0(a^n) = 0.$$

H^* , the modified empirical entropy

Definition (see Man01 for details)

The 0th-order modified empirical entropy of a string s of length n is

$$H_0^*(s) = \begin{cases} 0 & \text{if } n = 0, \\ (\lfloor \log n \rfloor + 1)/n & \text{if } n \geq 1 \text{ and } H_0(s) = 0, \\ H_0(s) & \text{otherwise.} \end{cases}$$

Low-entropy strings

Usually, μn is a minor concern and $O(\log n)$ is negligible — but not always:

Low-entropy strings

Usually, μn is a minor concern and $O(\log n)$ is negligible — but not always:

- ▶ μn can be exponentially larger than $H_k^*(s)n$,

Low-entropy strings

Usually, μn is a minor concern and $O(\log n)$ is negligible — but not always:

- ▶ μn can be exponentially larger than $H_k^*(s)n$,
- ▶ the constant hidden in the $O(\log n)$ term is often about h^{k+1}

Low-entropy strings

Usually, μn is a minor concern and $O(\log n)$ is negligible — but not always:

- ▶ μn can be exponentially larger than $H_k^*(s)n$,
- ▶ the constant hidden in the $O(\log n)$ term is often about h^{k+1} (e.g., for 4th-order compression of ASCII, about 2^{40}).

Low-entropy strings

Usually, μn is a minor concern and $O(\log n)$ is negligible — but not always:

- ▶ μn can be exponentially larger than $H_k^*(s)n$,
- ▶ the constant hidden in the $O(\log n)$ term is often about h^{k+1} (e.g., for 4th-order compression of ASCII, about 2^{40}).

We want a *small* ratio.

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

Distance Coding's output

Lemma (following KLV06)

For any string s of length n and $\lambda > 1$,

$$\begin{aligned} H_0(\text{DC}(s))|\text{DC}(s)| \\ \leq (\lambda + 0.01) (H_0(s)n + \log n) + \log(\zeta(\lambda) - 1)|\text{DC}(s)| + g \end{aligned}$$

for some g depending only on the alphabet size.

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

The maximal runs in a string

Lemma (MN05)

The number of maximal runs in a string s of length n is at most $H_0(s)n + 1$.

Corollaries for DC

Corollary

For any string s of length n and $\lambda > 1$,

$$\begin{aligned} H_0(\text{DC}(s))|\text{DC}(s)| \\ \leq (\lambda + \log(\zeta(\lambda) - 1) + 0.01) (H_0(s)n + \log n) + g \end{aligned}$$

for some g depending only on the alphabet size.

Corollaries for DC

Corollary

For any string s of length n ,

$$H_0(\text{DC}(s))|\text{DC}(s)| \leq 1.74 (H_0(s)n + \log n) + g$$

for some g depending only on the alphabet size.

Corollaries for DC

Corollary

For any string s of length n and $k \geq 0$,

$$|\text{Order}_0(\text{DC}(\text{BWT}(s)))| \leq (1.74 + \mu)H_k(s)n + O(\log n)$$

with the hidden coefficient depending only on the alphabet size and k .

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

Local optimality

Definition (Man01)

A compression algorithm A is *locally λ -optimal* if, for any string s and partition s_1, \dots, s_t of s ,

$$|A(s)| \leq \lambda \sum_{i=1}^t H_0^*(s_i) |s_i| + gt$$

for some g depending only on the alphabet size.

Local optimality

Lemma (Man01)

If A is locally λ -optimal then, for any string s of length n and $k \geq 0$,

$$|A(\text{BWT}(s))| \leq \lambda H_k^*(s)n + (1 + o(1)) \log n + g_k$$

for some g_k depending only on the alphabet size and k .

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

DC + Pfx is *not* locally optimal

Example

Consider $s = s_1 s_2 s_3$ with $s_1 = s_3 = a_1 \dots a_{h-1}$ and $s_2 = a_h^n$:

$$\sum_{i=1}^3 H_0^*(s_i) |s_i| = \log n + O(h \log h)$$

but $|\text{Pfx}(\text{DC}(s))| > (h-1) \log n$.

Avoiding long jumps

To use escapes:

1. whenever DC would write 1, we write $\langle 1, 1 \rangle$;
2. whenever DC would write a distance, we compare the cost of writing the distance to the cost of writing an escape and a re-entry, and do whichever is cheaper;
3. an escape is $\langle 1, 1 \rangle$;
4. a re-entry is $\langle 1, 2, \ell, a_i \rangle$.

Computing DC_{ESC}

ab**b**crraaaa



first character: $\langle 1, 1 \rangle$

other characters: 2, 3, 5, 6

Computing DC_{ESC}

~~abdb~~raaaa



first character: $\langle 1, 1 \rangle$

other characters: 2, 3, 5, 6

distances: $\langle 1, 1 \rangle$

Computing DC_{ESC}

abdbcr~~r~~aaaa



first character: $\langle 1, 1 \rangle$

other characters: 2, 3, 5, 6

distances: $\langle 1, 1 \rangle, 2, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 1 \rangle$

Computing DC_{ESC}

abdbcrraaaa



first character: $\langle 1, 1 \rangle$

other characters: 2, 3, 5, 6

distances: $\langle 1, 1 \rangle, 2, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 2, 2, 1 \rangle$

Computing DC_{ESC}

abdbccraaaa



first character: $\langle 1, 1 \rangle$

other characters: 2, 3, 5, 6

distances: $\langle 1, 1 \rangle, 2, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 2, 2, 1 \rangle$

last run: 4

Computing DC_{ESC}

abdbccraaaa



first character: $\langle 1, 1 \rangle$

other characters: 2, 3, 5, 6

distances: $\langle 1, 1 \rangle, 2, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 1 \rangle, \langle 1, 2, 2, 1 \rangle$

last run: 4



1, 1, 2, 3, 5, 6, 1, 1, 2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 1, 4

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

A stepping-stone

Consider the algorithm DC_{ESC}^* that takes a partition s_1, \dots, s_t of s as an argument — whereas DC_{ESC} takes *only* s — and uses escapes and re-entries only to replace distances that cross between partition elements.

$$\begin{aligned} & |\text{Order0}(DC_{ESC}(s))| \\ & \leq |\text{Order0}(DC_{ESC}^*(s_1, \dots, s_t))| \end{aligned}$$

A stepping-stone

Consider the algorithm DC_{ESC}^* that takes a partition s_1, \dots, s_t of s as an argument — whereas DC_{ESC} takes *only* s — and uses escapes and re-entries only to replace distances that cross between partition elements.

$$\begin{aligned}
 & |\text{Order0}(\text{DC}_{\text{ESC}}(s))| \\
 & \leq |\text{Order0}(\text{DC}_{\text{ESC}}^*(s_1, \dots, s_t))| \\
 & \leq \sum_{i=1}^t |\text{Order0}(\text{DC}(s_i))| + gt \\
 & \leq \sum_{i=1}^t \left((\lambda + 0.01) (H_0(s_i)|s_i| + \log n) + \right. \\
 & \quad \left. (\log(\zeta(\lambda) - 1) + \mu) |\text{DC}(s_i)| \right) + gt.
 \end{aligned}$$

Another lemma

Lemma

For any string s of length n and positive constants a , b and ϵ ,

$$a(H_0(s)n + \log n) + b|DC(s)| + O(1) \leq \max(2a, a + b + \epsilon)H_0^*(s)n + O(1).$$

DC_{ESC} + Order0 is locally optimal

Lemma

For any string s and partition s_1, \dots, s_t of s ,

$$|\text{Order0}(\text{DC}_{\text{ESC}}(s))| \leq (2.7 + \mu) \sum_{i=1}^t H_0^*(s_i) |s_i| + gt$$

for some g depending only on the alphabet size.

Outline

Introduction

The Burrows-Wheeler Transform

Distance Coding

Previous work

H vs. H^*

Three useful lemmas

Distance Coding's output

The maximal runs in a string

Local optimality

Distance Coding with escapes

Algorithm

Analysis

Conclusion

An entropy-only bound

Theorem

For any string s of length n and $k \geq 0$,

$$|\text{Order0}(\text{DC}_{\text{ESC}}(\text{BWT}(s)))| \leq (2.7 + \mu)H_k^*(s)n + (1 + o(1)) \log n + g_k$$

for some g_k depending only on the alphabet size and k .