

On the Burrows-Wheeler Transform of string collections

Zsuzsanna Lipták

University of Verona (Italy)

Primavera dell'Informatica Teorica

11 Jan. 2024

The Burrows-Wheeler Transform (BWT)

Ex.: $T = \text{banana}$. The BWT is a permutation of T : **nbaaa**

The Burrows-Wheeler Transform (BWT)

Ex.: $T = \text{banana}$. The BWT is a permutation of T : nbbaaa

all rotations (conjugates)

banana
ananab
nanaba
anaban
nabana
abanan

→
lexicographic
order

all rotations, sorted

abanan
ananab
anaban
banana
nabana
nanaba

The Burrows-Wheeler Transform (BWT)

Ex.: $T = \text{banana}$. The BWT is a permutation of T : **nbbaaa**

all rotations (conjugates)

banana
ananab
nanaba
anaban
nabana
abanan

→
lexicographic
order

all rotations, sorted

abanan
ananab
anaban
banana
nabana
nanaba

A (non-efficient) algorithm: List all of rotations of string T , sort them lexicographically, concatenate last characters: $\text{bwt}(\text{banana}) = \text{nbbaaa}$



Michael Burrows



Paolo Ferragina



Giovanni Manzini

AWARDS & RECOGNITION

Inventors of BW-transform and the FM-index Receive Kanellakis Award

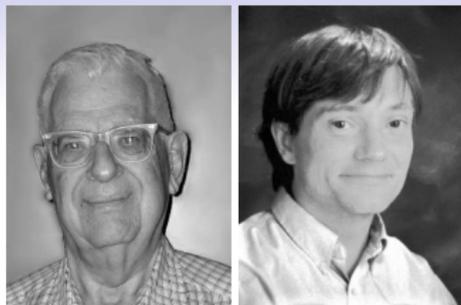
2022

Michael Burrows , Google; **Paolo Ferragina** , University of Pisa; and **Giovanni Manzini** , University of Pisa, receive the **ACM Paris Kanellakis Theory and Practice Award**  for inventing the BW-transform and the FM-index that opened and influenced the field of Compressed Data Structures with fundamental impact on Data Compression and Computational Biology. In 1994, Burrows and his late coauthor David Wheeler published their paper describing revolutionary data compression algorithm based on a reversible transformation of the input—the “Burrows-Wheeler Transform” (BWT). A few years later, Ferragina and Manzini showed that, by orchestrating the BWT with a new set of mathematical techniques and algorithmic tools, it became possible to build a “compressed index,” later called the FM-index. The introduction of the BW Transform and the development of the FM-index have had a profound impact on the theory of algorithms and data structures with fundamental advancements.

source: <https://awards.acm.org/kanellakis>

The BWT

- introduced by M. Burrows and D. Wheeler in 1994 as a **lossless text compression** algorithm
- P. Ferragina and G. Manzini showed later how to use it for **pattern matching**, leading to the **FM-index** [FOCS, 2000; JACM 2005]
- recent: **r-index** [Gagie et al, JACM 2020; Bannai et al. TCS 2020]



source: Adjeroh, Bell, Mukerjee (2008)

The BWT

- introduced by M. Burrows and D. Wheeler in 1994 as a **lossless text compression** algorithm
- P. Ferragina and G. Manzini showed later how to use it for **pattern matching**, leading to the **FM-index** [FOCS, 2000; JACM 2005]
- recent: **r-index** [Gagie et al, JACM 2020; Bannai et al. TCS 2020]



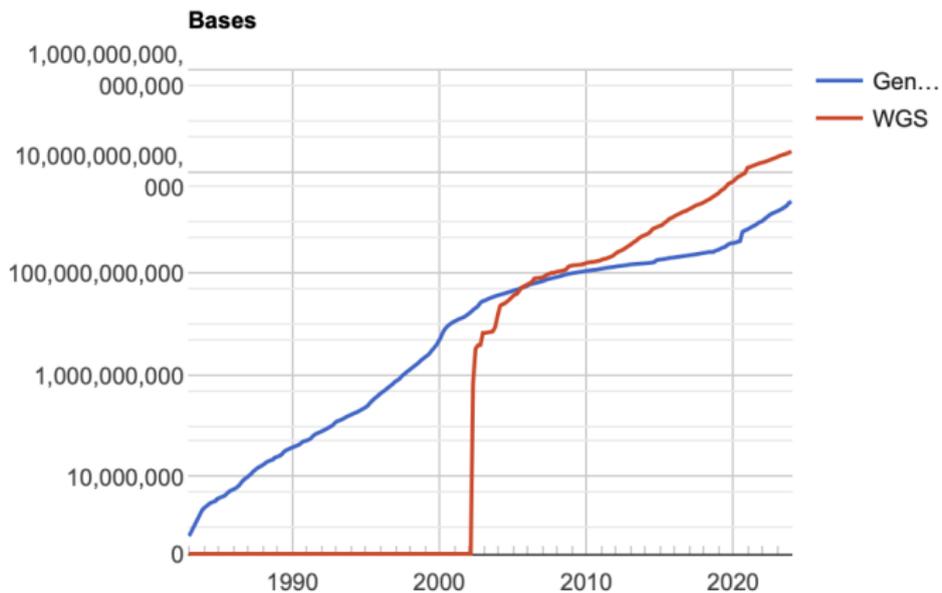
source: Adjeroh, Bell, Mukerjee (2008)

Some properties of the BWT:

- **computable** in linear time $\mathcal{O}(n)$
- **reversible** in linear time $\mathcal{O}(n)$
- uncompressed: **same space as text**
- if T **repetitive**, good for compression (**see later**)

$$n = |T|$$

GenBank and WGS Statistics



From strings to string collections

- Human Genome Project (first draft: 2000, completion: 2003)

From strings to string collections

- Human Genome Project (first draft: 2000, completion: 2003)
- Studying variation:
 - 1,000 Genomes Project (human): 2008-2015
 - 1001 Genomes (*Arabidopsis thaliana*)
 - 3,000 Rice Genomes Project
 - 100,000 Genomes Project (human, completed 2018)

From strings to string collections

- Human Genome Project (first draft: 2000, completion: 2003)
- Studying variation:
 - 1,000 Genomes Project (human): 2008-2015
 - 1001 Genomes (*Arabidopsis thaliana*)
 - 3,000 Rice Genomes Project
 - 100,000 Genomes Project (human, completed 2018)
- Population-wide:
 - Faroe Genome Project: sequence all 50,000 people
 - Sequencing Iceland (325,000 people): > 57,000 sequenced

From strings to string collections

- Human Genome Project (first draft: 2000, completion: 2003)
- Studying variation:
 - 1,000 Genomes Project (human): 2008-2015
 - 1001 Genomes (*Arabidopsis thaliana*)
 - 3,000 Rice Genomes Project
 - 100,000 Genomes Project (human, completed 2018)
- Population-wide:
 - Faroe Genome Project: sequence all 50,000 people
 - Sequencing Iceland (325,000 people): > 57,000 sequenced
- Human diversity:
 - Genes & Health in East London: 100,000 people of Bangladeshi and Pakistani origin
 - Sequencing African genomes (Nature 2020)
 - Sequencing indigenous Australian genomes (Nature 2023)

From strings to string collections

- Human Genome Project (first draft: 2000, completion: 2003)
- Studying variation:
 - 1,000 Genomes Project (human): 2008-2015
 - 1001 Genomes (*Arabidopsis thaliana*)
 - 3,000 Rice Genomes Project
 - 100,000 Genomes Project (human, completed 2018)
- Population-wide:
 - Faroe Genome Project: sequence all 50,000 people
 - Sequencing Iceland (325,000 people): > 57,000 sequenced
- Human diversity:
 - Genes & Health in East London: 100,000 people of Bangladeshi and Pakistani origin
 - Sequencing African genomes (Nature 2020)
 - Sequencing indigenous Australian genomes (Nature 2023)
- SARS-CoV-2 viral sequences

From strings to string collections

Our data is

- growing rapidly, and
- changing: **from individual strings to string collections**
- many of these consist of **many similar copies** of the same string

Outline of talk

- The Burrows-Wheeler Transform (BWT)
- The extended BWT (eBWT)
- Other variants of the BWT for string collections
- Why does it matter?
- Conclusions

The Burrows-Wheeler Transform

The Burrows-Wheeler Transform (BWT)

Recall: $T = \text{banana}$. The BWT is a permutation of T : nbaaaa

all rotations (conjugates)

banana
 ananab
 nanaba
 anaban
 nabana
 abanan

→
lexicographic
order

all rotations, sorted

abanan
 ananab
 anabab
 banana
 nabana
 nanaba

Why is the BWT useful in text compression?

BWT-matrix (F = first column, L = last column)

	F	L
1	abanan	a
2	anaban	a
3	ananab	b
4	banana	a
5	nabana	a
6	nanaba	a

Why is the BWT useful in text compression?

BWT-matrix (F = first column, L = last column)

	F	L
1	abanan	n
2	anaban	n
3	ananab	b
4	banana	a
5	nabana	a
6	nanaba	a

- **Obs. 1:** F = all characters of T in lex-order:
aaabnn

Why is the BWT useful in text compression?

BWT-matrix (F = first column, L = last column)

	F	L
1	abanan	n
2	anaban	n
3	ananab	b
4	banana	a
5	nabana	a
6	nanaba	a

- **Obs. 1:** F = all characters of T in lex-order:
aaabnn
- **Obs. 2:** for all i : L_i precedes F_i in T (cyclically):

$T =$ banana
1 2 3 4 5 6

Why is the BWT useful in text compression?

BWT-matrix (F = first column, L = last column)

	F	L
1	abanan	n
2	anaban	n
3	ananab	b
4	banana	a
5	nabana	a
6	nanaba	a

- **Obs. 1:** F = all characters of T in lex-order:
aaabnn
- **Obs. 2:** for all i : L_i precedes F_i in T (cyclically):
 $T = \text{banana}$
1 2 3 4 5 6
- **Obs. 3:** all occurrences of a substring appear in consecutive rows as prefix

Why is the BWT useful in text compression?

- **Obs. 1:** F = characters of T in lexicographic order
- **Obs. 2:** L_i precedes F_i in T
- **Obs. 3:** all occurrences of a substring appear in **consecutive** rows as prefix

Why is the BWT useful in text compression?

- **Obs. 1:** F = characters of T in lexicographic order
- **Obs. 2:** L_i precedes F_i in T
- **Obs. 3:** all occurrences of a substring appear in **consecutive** rows as prefix

Ex.: $T = \text{banana}$ has 2 occurrences of the substring ana

2 occ's of ana

abanan
anaban
ananab
banana
nabana
nanaba

Why is the BWT useful in text compression?

- **Obs. 1:** F = characters of T in lexicographic order
- **Obs. 2:** L_i precedes F_i in T
- **Obs. 3:** all occurrences of a substring appear in **consecutive** rows as prefix

Ex.: $T = \text{banana}$ has 2 occurrences of the substring ana

2 occ's of ana

abanan
anaban
ananab
banana
nabana
nanaba

2 occ's of na

preceded by a

abanan
anaban
ananab
banana
nabana
nanaba

Why is the BWT useful in text compression?

- **Obs. 1:** F = characters of T in lexicographic order
- **Obs. 2:** L_i precedes F_i in T
- **Obs. 3:** all occurrences of a substring appear in **consecutive** rows as prefix

Ex.: $T = \text{banana}$ has 2 occurrences of the substring **ana**

2 occ's of ana

abanan
anaban
ananab
banana
nabana
nanaba

*2 occ's of na
preceded by a*

abanan
anaban
ananab
banana
nabana
nanaba

*2 occ's of a
preceded by n*

abanan
anaban
ananab
banana
nabana
nanaba

Why is the BWT useful in text compression?

- **Obs. 1:** F = characters of T in lexicographic order
- **Obs. 2:** L_i precedes F_i in T
- **Obs. 3:** all occurrences of a substring appear in **consecutive** rows as prefix

Ex.: $T = \text{banana}$ has 2 occurrences of the substring **ana**

*2 occ's of **ana***

abanan
anaban
ananab
banana
nabana
nanaba

*2 occ's of **na**
preceded by **a***

abanan
anaban
ananab
banana
nabana
nanaba

*2 occ's of **a**
preceded by **n***

abanan
anaban
ananab
banana
nabana
nanaba

So: we get **a run** of **a**'s of length 2, and **a run** of **n**'s of length 2

Why is the BWT useful in text compression?

- **Obs. 1:** $F =$ characters of T in lexicographic order
- **Obs. 2:** L_i precedes F_i in T
- **Obs. 3:** all occurrences of a substring appear in **consecutive** rows as prefix

Ex.: $T = \text{banana}$ has 2 occurrences of the substring **ana**

*2 occ's of **ana***

abanan
anaban
ananab
banana
nabana
nanaba

*2 occ's of **na**
preceded by **a***

abanan
anaban
ananab
banana
nabana
nanaba

*2 occ's of **a**
preceded by **n***

abanan
anaban
ananab
banana
nabana
nanaba

So: we get **a run** of **a**'s of length 2, and **a run** of **n**'s of length 2 (2 = no. occ's).

Of course, things are a bit more complicated in general:

Of course, things are a bit more complicated in general:

rotation

BWT

...
he caverns measureless to man, And sank in tumult to a ... t
he caves. It was a miracle of rare device, A sunny pleasure-... t
he dome of pleasure Floated midway on the waves; Where was ... t
he fountain and the caves. It was a miracle of rare device,... t
he green hill athwart a cedarn cover! A savage place! as ... t
he hills, Enfolding sunny spots of greenery. But oh! that ... t
he milk of Paradise. t
he mingled measure From the fountain and the caves. It was a ... t
he on honey-dew hath fed, And drunk the milk of Paradise. ... t
he played, Singing of Mount Abora. Could I revive within me ... s
he sacred river ran, Then reached the caverns measureless ... t
he sacred river, ran Through caverns measureless to man ... t
he sacred river. Five miles meandering with a mazy motion ... t
he shadow of the dome of pleasure Floated midway on the waves ... T
he thresher's flail: And mid these dancing rocks at once and ... t
he waves; Where was heard the mingled measure From the ... t

*Kubla Kahn by Samuel Coleridge
(1998 characters)*

Of course, things are a bit more complicated in general:

rotation

BWT

...	
he caverns measureless to man, And sank in tumult to a ...	t
he caves. It was a miracle of rare device, A sunny pleasure-...	t
he dome of pleasure Floated midway on the waves; Where was ...	t
he fountain and the caves. It was a miracle of rare device,...	t
he green hill athwart a cedarn cover! A savage place! as ...	t
he hills, Enfolding sunny spots of greenery. But oh! that ...	t
he milk of Paradise.	t
he mingled measure From the fountain and the caves. It was a ...	t
he on honey-dew hath fed, And drunk the milk of Paradise.
he played, Singing of Mount Abora. Could I revive within me ...	s
he sacred river ran, Then reached the caverns measureless ...	t
he sacred river, ran Through caverns measureless to man ...	t
he sacred river. Five miles meandering with a mazy motion ...	t
he shadow of the dome of pleasure Floated midway on the waves ...	T
he thresher's flail: And mid these dancing rocks at once and ...	t
he waves; Where was heard the mingled measure From the ...	t

many **the's**, some **he**, **she**, **The**

Kubla Kahn by Samuel Coleridge
(1998 characters)

Compression with the BWT

- takes advantage of this 'clustering effect'

Compression with the BWT

- takes advantage of this 'clustering effect'
- **Def.:** $r(T)$ = number of runs of $\text{bwt}(T)$
(run: maximal equal-letter run)

Ex.: $r(\text{banana}) = 3$
 $\text{bwt}(\text{banana}) = \text{nbbaaa}$

Compression with the BWT

- takes advantage of this 'clustering effect'
- **Def.:** $r(T)$ = number of runs of $\text{bwt}(T)$
(run: maximal equal-letter run)
- **compression with BWT:**
uses runlength-encoding (RLE)

Ex.: $r(\text{banana}) = 3$
 $\text{bwt}(\text{banana}) = \text{nbaaa}$

Compression with the BWT

- takes advantage of this 'clustering effect'
- **Def.:** $r(T)$ = number of runs of $\text{bwt}(T)$
(run: maximal equal-letter run)
- **compression with BWT:**
uses runlength-encoding (RLE)

replace each run by (char,int)-pair

$\text{RLE}(\text{bbbbbbbbc1aaaaaaaaabb}) = \text{b8c1a11b2}$

Ex.: $r(\text{banana}) = 3$
 $\text{bwt}(\text{banana}) = \text{nbaaa}$

Compression with the BWT

- takes advantage of this 'clustering effect'
- **Def.:** $r(T)$ = number of runs of $\text{bwt}(T)$
(run: maximal equal-letter run)
- **compression with BWT:**
uses runlength-encoding (RLE)

replace each run by (char,int)-pair

$\text{RLE}(\text{bbbbbbbbc1a11b2}) = \text{b8c1a11b2}$

Compression: $T \mapsto \underbrace{\text{RLE}(\text{bwt}(T))}_{\text{storage space: } O(r)}$

Ex.: $r(\text{banana}) = 3$
 $\text{bwt}(\text{banana}) = \text{nb1aa}$

Ex.: $\text{banana} \mapsto \text{n2b1a3}$

Compression with the BWT

- takes advantage of this 'clustering effect'
- **Def.:** $r(T)$ = number of runs of $\text{bwt}(T)$
(run: maximal equal-letter run)

Ex.: $r(\text{banana}) = 3$
 $\text{bwt}(\text{banana}) = \text{nbaaa}$

- **compression with BWT:**
uses runlength-encoding (RLE)

replace each run by (char,int)-pair

$\text{RLE}(\text{bbbbbbbbcaaaaaaaaabb}) = \text{b8c1a11b2}$

Compression: $T \mapsto \underbrace{\text{RLE}(\text{bwt}(T))}_{\text{storage space: } O(r)}$

Ex.: $\text{banana} \mapsto \text{n2b1a3}$

- good if r is much smaller than $n = |T|$
(i.e. if few runs)

Compression with the BWT

- takes advantage of this 'clustering effect'
- **Def.:** $r(T)$ = number of runs of $\text{bwt}(T)$
(run: maximal equal-letter run)

Ex.: $r(\text{banana}) = 3$
 $\text{bwt}(\text{banana}) = \text{nbaaa}$

- **compression with BWT:**
uses runlength-encoding (RLE)

replace each run by (char,int)-pair

$\text{RLE}(\text{bbbbbbbbcaaaaaaaaabb}) = \text{b8c1a11b2}$

Compression: $T \mapsto \underbrace{\text{RLE}(\text{bwt}(T))}_{\text{storage space: } O(r)}$

Ex.: $\text{banana} \mapsto \text{n2b1a3}$

- good if r is much smaller than $n = |T|$
(i.e. if few runs)
- for **repetitive strings**, r is small
(repetitive: many repeated substrings)

Reversing the BWT (lossless compression)

input: `mnbaaa`, 4

output: (wanted) `banana`.

$\text{bwt}(T)$, i : where $1 \leq i \leq n$

T : i 'th rotation lex.ly

Reversing the BWT (lossless compression)

input: `nbbaaa`, 4

output: (wanted) `banana`.

$\text{bwt}(T)$, i : where $1 \leq i \leq n$

T : i 'th rotation lex.ly

Thm. (LF-property): The j 'th occurrence of character x in L is the j 'th occurrence of character x in F .

	F	L
1	abanan	
2	anaban	
3	ananab	
4	banana	
5	nabana	
6	nanaba	

$T = \text{banana}$
1 2 3 4 5 6

Reversing the BWT (lossless compression)

input: `mnbaaa`, 4

output: (wanted) `banana`.

$\text{bwt}(T)$, i : where $1 \leq i \leq n$

T : i 'th rotation lex.ly

Thm. (LF-property): The j 'th occurrence of character x in L is the j 'th occurrence of character x in F .

	F	L
1	abanan	
2	anaban	
3	ananab	
4	banana	
5	nabana	
6	nanaba	

$T = \text{banana}$
1 2 3 4 5 6

Recall:

Obs. 1: F = all characters of T in lex-order:

Obs. 2: for all i : L_i precedes F_i in T .

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_j precedes F_j in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_i precedes F_i in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: `nbbaaa`, 4

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_j precedes F_j in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: `mbaaa`, 4

	L
1	n
2	n
3	b
4	a
5	a
6	a

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_j precedes F_j in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: `nbbaaa`, 4

	F	L
1	a	n
2	a	n
3	a	b
4	b	a
5	n	a
6	n	a

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_i precedes F_i in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: **nb**aaaa, 4

	F	L	
1	a	n	
2	a	n	
3	a	b	a
4	b	a	
5	n	a	
6	n	a	

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_i precedes F_i in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: `nbbaaa`, 4

	F	L	
1	a	n	
2	a	n	
3	a	b	n a
4	b	a	
5	n	a	
6	n	a	

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_i precedes F_i in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: `nbaaa`, 4

	F	L	
1	a	n	
2	a	n	
3	a	b	a n a
4	b	a	
5	n	a	
6	n	a	

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_i precedes F_i in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: `nbbaaa`, 4

	F	L	
1	a	n	
2	a	n	
3	a	b	n a n a
4	b	a	
5	n	a	
6	n	a	

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_i precedes F_i in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: `nbbaaa`, 4

	F	L	
1	a	n	
2	a	n	
3	a	b	a n a n a
4	b	a	
5	n	a	
6	n	a	

Reversing the BWT

- **Obs. 1:** F = all characters of T in lex-order
- **Obs. 2:** L_i precedes F_i in T
- **LF-property:** The j 'th x in L is the j 'th x in F .

input: `nbbaaa`, 4

	F	L	
1	a	n	
2	a	n	
3	a	b	b a n a n a
4	b	a	
5	n	a	
6	n	a	

The BWT of string collections

- The BWT is good on repetitive strings.
- Our string collections are highly repetitive:
many similar copies of the same string
- But: **how** do we compute the BWT of a multiset?

The BWT of string collections

- The BWT is good on repetitive strings.
- Our string collections are highly repetitive:
many similar copies of the same string
- But: **how** do we compute the BWT of a multiset?

Generalization of the BWT to multisets:
the **extended BWT (eBWT)** (**next**)

The extended BWT

The extended BWT

[Mantaci, Restivo, Rosone, Sciortino, TCS, 2007]

Ex. $\mathcal{M} = \{\text{bana}, \text{an}\}$. The eBWT is a permutation of the characters of \mathcal{M} : $\text{eBWT}(\mathcal{M}) = \text{nbnaaa}$.

all rotations (conjugates)

bana
anab
naba
aban
an
na

→

omega order

all rotations, sorted

aban n
anab b
an n
bana a
naba a
na a

N.B. $\text{anab} <_{\omega} \text{an}$, since $\text{anab} \cdot \text{anab} \cdots <_{\text{lex}} \text{an} \cdot \text{an} \cdot \text{an} \cdot \text{an} \cdots$

The extended BWT

Def.(omega-order): $T <_{\omega} S$ if (a) $T^{\omega} <_{\text{lex}} S^{\omega}$, or
(b) $T^{\omega} = S^{\omega}$, $T = U^k$, $S = U^m$ and $k < m$

$\mathcal{M} = \{\text{bana}, \text{an}\}$

omega-order

lex-order

aban n

aban n

anab b

an n

an n

anab b

bana a

bana a

naba a

na a

na a

naba a

N.B. With the lex-order, the LF-property would not hold!

The extended BWT

- **omega-order** instead of lex-order
- same as lex-order if neither string is prefix of the other
- omega-order necessary for the **LF-property**
- the eBWT inherits **BWT properties**: clustering effect, reversibility, useful for lossless text compression, efficient pattern matching, . . .
- However, until recently no linear-time algorithm known.

The extended BWT

- **omega-order** instead of lex-order
- same as lex-order if neither string is prefix of the other
- omega-order necessary for the **LF-property**
- the eBWT inherits **BWT properties**: clustering effect, reversibility, useful for lossless text compression, efficient pattern matching, . . .
- However, until recently no linear-time algorithm known.

2021:

- linear-time algorithm [Bannai, Kärkkäinen, Köppl, Piatkowski, CPM 2021]
- We simplified this algorithm, and
- gave first efficient implementations of the eBWT: tools [pfpebwt](#), [cais](#) [Boucher, Cenzato, L., Rossi, Sciortino, SPIRE 2021]

Other BWT variants for string collections

The BWT of string collections

[Cenzato and L., CPM 2022, Arxiv 2023]

Question: How do dedicated tools compute the BWT of a string collection? (string collection: multiset of strings)

- We studied 18 publicly available tools.
- Only ours compute the eBWT (`pfpebwt`, `cais`).
- We identified 4 more non-equivalent approaches: the resulting BWTs are all different.
- Often the method is not explicitly stated.
- **Underlying assumption: they are all the same.**
- But they differ a lot (Hamming distance, number of runs).
- **N.B.:** all BWT variants are correct (LF-property, ...)

The other BWT variants for string collections

The different approaches are:

1. extended BWT of strings with terminator symbol \$ (dollarEBWT)
2. concatenate strings, separating them with different dollars (multidoBWT)
3. first sort colexicographically, then do 2. (colexBWT)
4. concatenate strings, separating them with same dollar (concatBWT)

All use terminator / separator symbols ('dollars'). So we call them **separator-based** BWT variants.

The BWT variants for string collections

Ex. $\mathcal{M} = \{ATATG, TGA, ACG, ATCA, GGA\}$

variant (our terminology)	result on example	tools
eBWT	CGGGATGTACGTTAAAAA	pfpebwt, cais
dollarEBWT	GGAAACGG\$\$\$\$TTACTGT\$AAA\$	G2BWT, msbwt
multidolBWT	GAGAAGCG\$\$\$\$TTATCTG\$AAA\$	gsufsort, ropebwt2, eGSA, Merge-BWT, eGAP, nvSetBWT, BCR-LCP-GSA, grlBWT, BEETL, bwt-lcp-parallel
colexBWT	AAAGGCGG\$\$\$\$TTACTGT\$AAA\$	ropebwt2, BCR-LCP-GSA
concatBWT	\$AAGAGGGC#\$TTACTGT\$AAA\$	BigBWT, r-pfbwt, CMS-BWT
		tools for single strings

The dollar-eBWT

1. $\text{dollarEBWT}(\mathcal{M}) = \text{eBWT}(\{T_i\$: T_i \in \mathcal{M}\})$, $\$ < c$ for all char's c

Now no string is prefix of another \implies omega-order same as lex-order.

$\mathcal{M} = \{\text{bana}\$, \text{an}\$\}$

dollarEBWT

\$an n

\$bana a

a\$bana n

an\$ \$

ana\$b b

bana\$ \$

n\$a a

na\$b a

nan\$b\$a

The dollar-eBWT

1. $\text{dollarEBWT}(\mathcal{M}) = \text{eBWT}(\{T_i\$: T_i \in \mathcal{M}\})$, $\$ < c$ for all char's c

Now no string is prefix of another \implies omega-order same as lex-order.

$\mathcal{M} = \{\text{bana}\$, \text{an}\$\}$

dollarEBWT

\$an n

\$bana a

a\$bana n

an\$ \$

ana\$b b

bana\$ \$

n\$a a

na\$b a

nan\$b\$a

eBWT of {bana, an}

aban n

anab b

an n

bana a

naba a

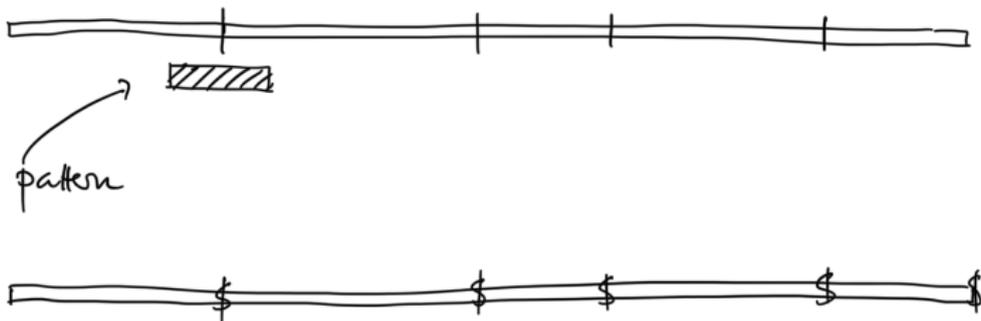
na a

nbnaaa

The different BWT variants

The other 3 methods concatenate the input strings, and then apply the classical BWT.

The main issue here is to avoid spurious substrings:



The multidollar BWT

2. $\text{multidollBWT}(\mathcal{M}) = \text{bwt}(T_1\$_1 T_2\$_2 \cdots T_k\$_k)$, where dollars are smaller than characters from Σ , and $\$_1 < \$_2 < \dots < \$_k$

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\} \rightsquigarrow$

$\text{bwt}(\text{ATATG}\$_1 \text{TGA}\$_2 \text{ACG}\$_3 \text{ATCA}\$_4 \text{GGA}\$_5) = \text{GAGAAGCG}\$\$\$\$\text{TTATCTG}\$\text{AAAA}\$$

The multidollar BWT

2. $\text{multidolBWT}(\mathcal{M}) = \text{bwt}(T_1\$_1 T_2\$_2 \cdots T_k\$_k)$, where dollars are smaller than characters from Σ , and $\$_1 < \$_2 < \dots < \$_k$

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\} \rightsquigarrow$

$\text{bwt}(\text{ATATG}\$_1\text{TGA}\$_2\text{ACG}\$_3\text{ATCA}\$_4\text{GGA}\$_5) = \text{GAGAAGCG}\$\$\$\$\text{TTATCTG}\$\text{AAA}\$\$

- most commonly used method
- analogous to Generalized Suffix Tree and Generalized Suffix Array
- dollars are different only conceptually (break ties by index)
- equivalent: concatenate without separators, use bitstring marking string beginnings

The colex BWT

3. **colexBWT**(\mathcal{M}): multidolBWT of the strings in colexicographic order

colex order = lexicographic order of the reverse strings

The colex BWT

3. $\text{colexBWT}(\mathcal{M})$: multidolBWT of the strings in colexicographic order

colex order = lexicographic order of the reverse strings

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\}$

colex order: $\text{ATCA}, \text{GGA}, \text{TGA}, \text{ACG}, \text{ATATG} \rightsquigarrow$

$\text{bwt}(\text{ATCA}\$1\text{GGA}\$2\text{TGA}\$3\text{ACG}\$4\text{ATATG}\$5) = \text{AAAGGCGG}\$3\$\$4\text{TTACTGT}\$5\text{AAA}\$$

The colex BWT

3. `colexBWT`(\mathcal{M}): multidolBWT of the strings in colexicographic order

colex order = lexicographic order of the reverse strings

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\}$

colex order: `ATCA`, `GGA`, `TGA`, `ACG`, `ATATG` \rightsquigarrow

`bwt(ATCA$1GGA$2TGA$3ACG$4ATATG$5) = AAAGGCGG$$$TTACTGTAAA`

- reduces number of runs (see later)
- implemented as an option in `ropebwt2`, `BCR-LCP-GSA`

The concat BWT

4. $\text{concatBWT}(\mathcal{M}) = \text{bwt}(T_1\$T_2\$ \cdots T_k\$\#)$, where $\# < \$$

The concat BWT

4. $\text{concatBWT}(\mathcal{M}) = \text{bwt}(T_1\$T_2\$ \dots T_k\$\#)$, where $\# < \$$

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\} \rightsquigarrow$

$\text{bwt}(\text{ATATG}\$\text{TGA}\$\text{ACG}\$\text{ATCA}\$\text{GGA}\$\#) = \$\text{AAGAGGGC}\$\#\text{TTACTGT}\$\text{AAAA}\$$

The concat BWT

4. $\text{concatBWT}(\mathcal{M}) = \text{bwt}(T_1\$T_2\$ \dots T_k\$\#)$, where $\# < \$$

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\} \rightsquigarrow$

$\text{bwt}(\text{ATATG}\$\text{TGA}\$\text{ACG}\$\text{ATCA}\$\text{GGA}\$\#) = \$\text{AAGAGGGC}\$\#\$\text{TTACTGT}\$\text{AAAA}\$$

(for easier comparison, we simplify to $\text{AAGAGGGC}\$\#\$\text{TTACTGT}\$\text{AAAA}\$$)

The concat BWT

4. $\text{concatBWT}(\mathcal{M}) = \text{bwt}(T_1\$T_2\$ \dots T_k\$\#)$, where $\# < \$$

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\} \rightsquigarrow$

$\text{bwt}(\text{ATATG}\$\text{TGA}\$\text{ACG}\$\text{ATCA}\$\text{GGA}\$\#) = \$\text{AAGAGGGC}\$\#\$\text{TTACTGT}\$\text{AAAA}\$$

(for easier comparison, we simplify to $\text{AAGAGGGC}\$\$\$\text{TTACTGT}\$\text{AAAA}\$$)

- very easy to implement
- used e.g. in BigBWT, CMS-BWT.

Interesting intervals

Q. Where exactly do these BWT variants differ? **A.** in interesting intervals

Interesting intervals

Q. Where exactly do these BWT variants differ? **A.** in interesting intervals

Ex. $\mathcal{M} = \{ATATG, TGA, ACG, ATCA, GGA\}$

BWT variant	example
<i>non-sep. based</i> eBWT(\mathcal{M})	CGGGATGTACGTTAAAAA
<i>separator-based</i> dollarEBWT(\mathcal{M})	GGAAACGG\$\$\$\$TTACTGT\$AAA\$
multidolBWT(\mathcal{M})	GAGAAGCG\$\$\$\$TTATCTG\$AAA\$
colexBWT(\mathcal{M})	AAAGGCGG\$\$\$\$TTACTGT\$AAA\$
concatBWT(\mathcal{M})	AAGAGGCG\$\$\$\$TTACTGT\$AAA\$

in color: **interesting intervals**

Interesting intervals

Lemma: If two separator-based BWTs differ in position i then $i \in [b, e]$ for some interesting interval $[b, e]$.

Interesting intervals

Lemma: If two separator-based BWTs differ in position i then $i \in [b, e]$ for some interesting interval $[b, e]$.

Def. U is called a **left-maximal shared suffix** if there exist two strings $S_1, S_2 \in \mathcal{M}$ such that U is a suffix of $S_1\$$ and $S_2\$$ and is preceded by different characters in S_1 and S_2 . An interval $[b, e]$ is **interesting** if it corresponds to all occurrences of some left-maximal shared suffix U (i.e., its SA-interval).

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\}$, $U = \text{A\$}$.

Interesting intervals

Lemma: If two separator-based BWTs differ in position i then $i \in [b, e]$ for some interesting interval $[b, e]$.

Def. U is called a **left-maximal shared suffix** if there exist two strings $S_1, S_2 \in \mathcal{M}$ such that U is a suffix of $S_1\$$ and $S_2\$$ and is preceded by different characters in S_1 and S_2 . An interval $[b, e]$ is **interesting** if it corresponds to all occurrences of some left-maximal shared suffix U (i.e., its SA-interval).

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\}$, $U = \text{A\$}$.

A\$ATC C

A\$GG G

A\$TG G

dollarEBWT

Interesting intervals

Lemma: If two separator-based BWTs differ in position i then $i \in [b, e]$ for some interesting interval $[b, e]$.

Def. U is called a **left-maximal shared suffix** if there exist two strings $S_1, S_2 \in \mathcal{M}$ such that U is a suffix of $S_1\$$ and $S_2\$$ and is preceded by different characters in S_1 and S_2 . An interval $[b, e]$ is **interesting** if it corresponds to all occurrences of some left-maximal shared suffix U (i.e., its SA-interval).

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\}$, $U = \text{A\$}$.

A\$ATC	C	A\$ ₂ ...	G
A\$GG	G	A\$ ₄ ...	C
A\$TG	G	A\$ ₅ ...	G

dollarEBWT

multidolBWT

Interesting intervals

Lemma: If two separator-based BWTs differ in position i then $i \in [b, e]$ for some interesting interval $[b, e]$.

Def. U is called a **left-maximal shared suffix** if there exist two strings $S_1, S_2 \in \mathcal{M}$ such that U is a suffix of $S_1\$$ and $S_2\$$ and is preceded by different characters in S_1 and S_2 . An interval $[b, e]$ is **interesting** if it corresponds to all occurrences of some left-maximal shared suffix U (i.e., its SA-interval).

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\}$, $U = \text{A\$}$.

A\$ATC	C	A\$ ₂ ...	G	A\$ ₁ ...	C
A\$GG	G	A\$ ₄ ...	C	A\$ ₂ ...	G
A\$TG	G	A\$ ₅ ...	G	A\$ ₃ ...	G

dollarEBWT

multidoBWT

colexBWT

Interesting intervals

Lemma: If two separator-based BWTs differ in position i then $i \in [b, e]$ for some interesting interval $[b, e]$.

Def. U is called a **left-maximal shared suffix** if there exist two strings $S_1, S_2 \in \mathcal{M}$ such that U is a suffix of $S_1\$$ and $S_2\$$ and is preceded by different characters in S_1 and S_2 . An interval $[b, e]$ is **interesting** if it corresponds to all occurrences of some left-maximal shared suffix U (i.e., its SA-interval).

Ex. $\mathcal{M} = \{\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}\}$, $U = \text{A\$}$.

A\$ATC	C	A\$ ₂ ...	G	A\$ ₁ ...	C	A\$#	G
A\$GG	G	A\$ ₄ ...	C	A\$ ₂ ...	G	A\$A...	G
A\$TG	G	A\$ ₅ ...	G	A\$ ₃ ...	G	A\$G...	C

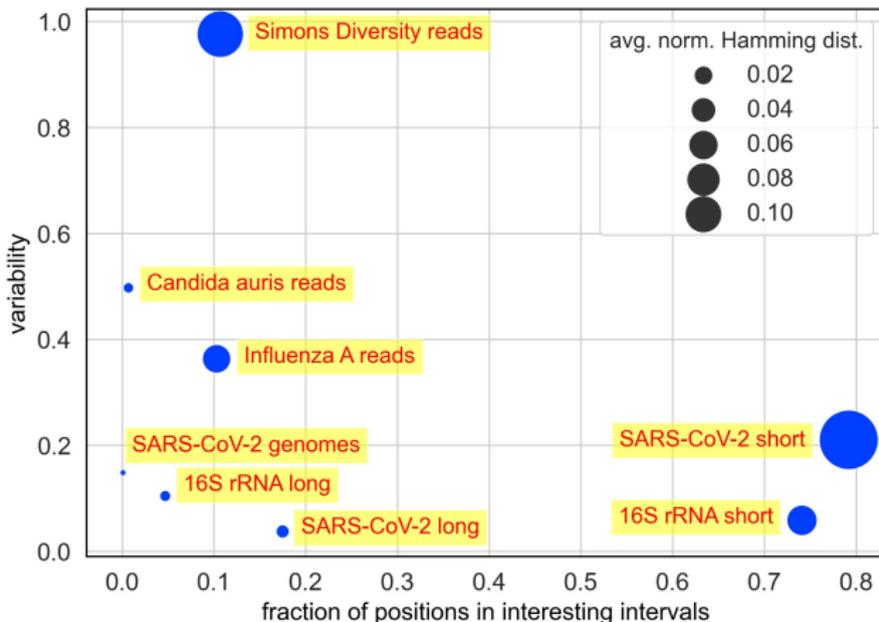
dollarEBWT

multidolBWT

colexBWT

concatBWT

Hamming distance between separator-based BWTs



Variability

$$\text{var}(\mathcal{M}) = \frac{\sum_{[b,e] \text{ interesting int.}} \text{var}([b, e])}{\sum_{[b,e] \text{ interesting int.}} (e - b + 1)}, \text{ where } \text{var}([b, e]) = \max \text{ no. runs in } [b, e]$$

(depends on Parikh vector)

Why does it matter?

Theoretician: You are all using different methods to compute the BWT of string collections, and the results are pretty different!

Theoretician: You are all using different methods to compute the BWT of string collections, and the results are pretty different!

Programmer: It doesn't matter, all I care about is that it's efficient.

Theoretician: You are all using different methods to compute the BWT of string collections, and the results are pretty different!

Programmer: It doesn't matter, all I care about is that it's efficient.

Theoretician: ...and correct?

Theoretician: You are all using different methods to compute the BWT of string collections, and the results are pretty different!

Programmer: It doesn't matter, all I care about is that it's efficient.

Theoretician: ... and correct?

Programmer: Ok, but you said yourself that it was all correct!

Theoretician: You are all using different methods to compute the BWT of string collections, and the results are pretty different!

Programmer: It doesn't matter, all I care about is that it's efficient.

Theoretician: ... and correct?

Programmer: Ok, but you said yourself that it was all correct!

Theoretician: But it's not nice that your tool computes a different thing from your competitor's.

Theoretician: You are all using different methods to compute the BWT of string collections, and the results are pretty different!

Programmer: It doesn't matter, all I care about is that it's efficient.

Theoretician: ... and correct?

Programmer: Ok, but you said yourself that it was all correct!

Theoretician: But it's not nice that your tool computes a different thing from your competitor's.

Programmer: I am never going to use her tool anyway!

Why you should care

1. number of runs
2. the parameter r is **not well-defined**
3. input order dependence

1. Number of runs

r = number of runs of the BWT.

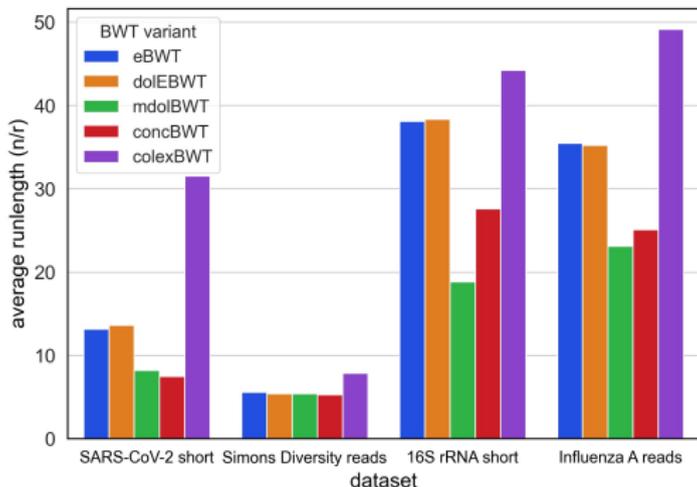
1. Number of runs

r = number of runs of the BWT.

Ex. $\mathcal{M} = \{ATATG, TGA, ACG, ATCA, GGA\}$

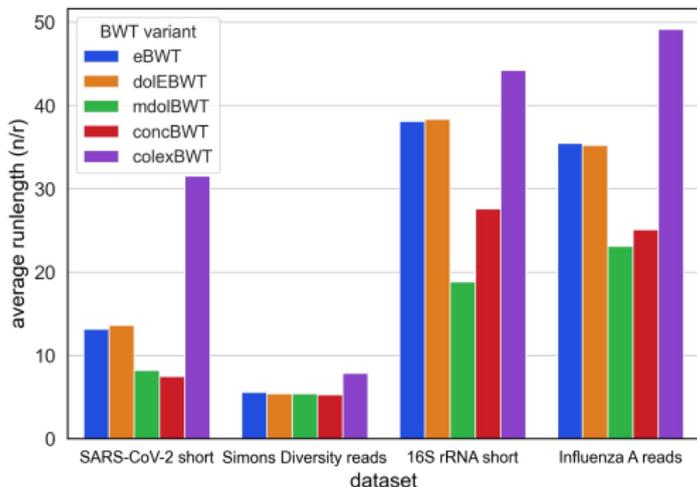
BWT variant	example	r	r w/o '\$'s
<i>non-sep. based</i>			
eBWT(\mathcal{M})	CGGGATGTACGTTAAAA	11	11
<i>separator-based</i>			
dollarEBWT(\mathcal{M})	GGAAACGG\$\$\$\$TTACTGT\$AAA\$	14	11
multidolBWT(\mathcal{M})	GAGAAGCG\$\$\$\$TTATCTG\$AAA\$	17	14
colexBWT(\mathcal{M})	AAAGGCGG\$\$\$\$TTACTGT\$AAA\$	14	11
concatBWT(\mathcal{M})	AAGAGGGC\$\$\$\$TTACTGT\$AAA\$	15	12

1. Number of runs



Average runlength (n/r) on four short sequence datasets, of all BWT variants.
(500,000 sequences each, of length between 50 and 301.)

1. Number of runs

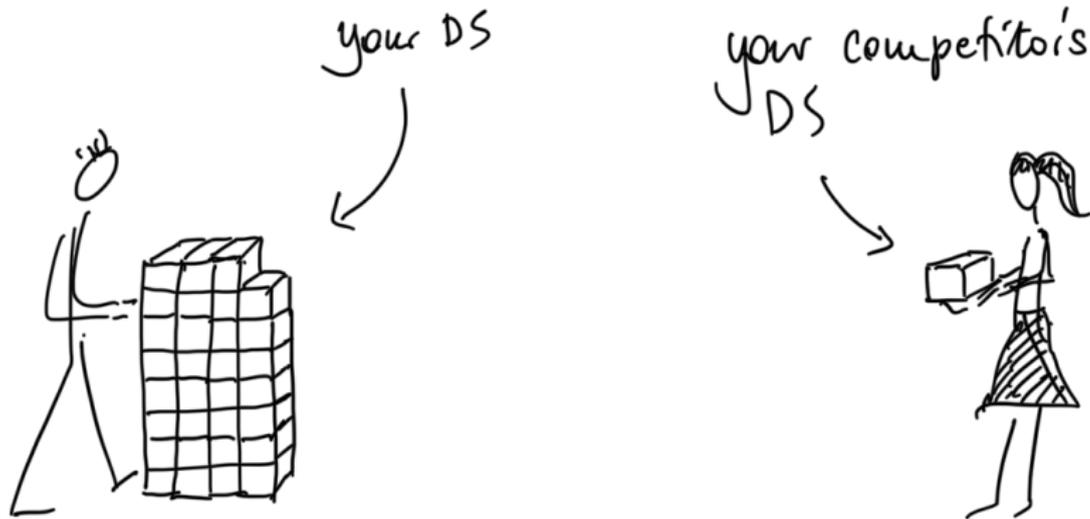


Average runlength (n/r) on four short sequence datasets, of all BWT variants.
(500,000 sequences each, of length between 50 and 301.)

- On these datasets, difference of a factor of **up to 4.2**.
- In a separate work, difference of a factor of **up to 31**.

[Cenzato, Guerrini, L., Rosone, DCC 2023]

size of data structures $O(r)$



So maybe you should care...

2. The parameter r

- size of data structures $\mathcal{O}(r)$ (r -index)

Gagie et al. [JACM 2020],
Bannai et al. [TCS 2020]

2. The parameter r

- size of data structures $\mathcal{O}(r)$ (r -index) Gagie et al. [JACM 2020],
Bannai et al. [TCS 2020]
- algorithms' running time ideally a function of r (not of $n = |T|$)

2. The parameter r

- size of data structures $\mathcal{O}(r)$ (r -index) Gagie et al. [JACM 2020],
Bannai et al. [TCS 2020]
- algorithms' running time ideally a function of r (not of $n = |T|$)
- increasingly used as a repetitiveness measure of T , similar to z (number of Lempel-Ziv phrases)
 - as a property of the dataset Bannai et al. [TCS 2020],
Boucher et al. [ALENEX 2021],
Cobas et al. [CPM 2021]

2. The parameter r

- size of data structures $\mathcal{O}(r)$ (r -index) Gagie et al. [JACM 2020],
Bannai et al. [TCS 2020]
- algorithms' running time ideally a function of r (not of $n = |T|$)
- increasingly used as a repetitiveness measure of T , similar to z (number of Lempel-Ziv phrases)
 - as a property of the dataset Bannai et al. [TCS 2020],
Boucher et al. [ALENEX 2021],
Cobas et al. [CPM 2021]
 - in theoretical work on repetitiveness measures Kempa and Kociumaka [FOCS 2020],
Navarro [ACM Comp. Surv., 2021],
Akagi et al. [Inf. Comp. 2023]

3. Input order dependence

3. Input order dependence

N.B. multidolBWT and concatBWT depend on the input order!

$\mathcal{M}_1 = [\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}]$

$\text{mdolBWT}(\mathcal{M}_1) = \text{GAGAAGCG} \$ \$ \$ \text{TTATCTG} \$ \text{AAA} \$$

$\mathcal{M}_2 = [\text{ACG}, \text{ATATG}, \text{GGA}, \text{TGA}, \text{ATCA}]$

$\text{mdolBWT}(\mathcal{M}_2) = \text{GGAAAGGC} \$ \$ \$ \text{TTACTGT} \$ \text{AAA} \$$

$\mathcal{M}_1 = [\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}]$

$\text{concBWT}(\mathcal{M}_1) = \text{AAGAGGC} \$ \$ \$ \text{TTACTGT} \$ \text{AAA} \$$

$\mathcal{M}_2 = [\text{ACG}, \text{ATATG}, \text{GGA}, \text{TGA}, \text{ATCA}]$

$\text{concBWT}(\mathcal{M}_2) = \text{AGAGCGG} \$ \$ \$ \text{TTACTGT} \$ \text{AAA} \$$

3. Input order dependence

N.B. multidolBWT and concatBWT depend on the input order!

$\mathcal{M}_1 = [\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}]$	$\text{mdolBWT}(\mathcal{M}_1) =$	
$\mathcal{M}_2 = [\text{ACG}, \text{ATATG}, \text{GGA}, \text{TGA}, \text{ATCA}]$	$\text{mdolBWT}(\mathcal{M}_2) =$	

$\mathcal{M}_1 = [\text{ATATG}, \text{TGA}, \text{ACG}, \text{ATCA}, \text{GGA}]$	$\text{concBWT}(\mathcal{M}_1) =$	
$\mathcal{M}_2 = [\text{ACG}, \text{ATATG}, \text{GGA}, \text{TGA}, \text{ATCA}]$	$\text{concBWT}(\mathcal{M}_2) =$	

Thus, giving the same dataset to the same tool but in different order can produce very different results! (incl. the number of runs)

The multidollar BWT can simulate all others

Prop. Let \mathcal{M} be given, and L some separator-based BWT on \mathcal{M} . Then there exists an input permutation π such that $\text{multidol}(\pi(\mathcal{M})) = L$.

The multidollar BWT can simulate all others

Prop. Let \mathcal{M} be given, and L some separator-based BWT on \mathcal{M} . Then there exists an input permutation π such that $\text{multidol}(\pi(\mathcal{M})) = L$.

Proof sketch: colexBWT: colex order, dollarEBWT: lex order, concatBWT: lex order of subseq strings

The multidollar BWT can simulate all others

Prop. Let \mathcal{M} be given, and L some separator-based BWT on \mathcal{M} . Then there exists an input permutation π such that $\text{multidol}(\pi(\mathcal{M})) = L$.

Proof sketch: colexBWT: colex order, dollarEBWT: lex order, concatBWT: lex order of subseq strings

- Prop. \implies any separator-based BWT variant can be computed using the multidollar method

The multidollar BWT can simulate all others

Prop. Let \mathcal{M} be given, and L some separator-based BWT on \mathcal{M} . Then there exists an input permutation π such that $\text{multidol}(\pi(\mathcal{M})) = L$.

Proof sketch: colexBWT: colex order, dollarEBWT: lex order, concatBWT: lex order of subseq strings

- Prop. \implies any separator-based BWT variant can be computed using the multidollar method
- Bentley, Gibney, and Thankachan [ESA 2020] gave a linear-time algorithm for the input order of multidollar BWT with **minimum r**

The multidollar BWT can simulate all others

Prop. Let \mathcal{M} be given, and L some separator-based BWT on \mathcal{M} . Then there exists an input permutation π such that $\text{multidol}(\pi(\mathcal{M})) = L$.

Proof sketch: colexBWT: colex order, dollarEBWT: lex order, concatBWT: lex order of subseq strings

- Prop. \implies any separator-based BWT variant can be computed using the multidollar method
- Bentley, Gibney, and Thankachan [ESA 2020] gave a linear-time algorithm for the input order of multidollar BWT with **minimum r**
- We implemented this algorithm in our tool **optimalBWT**
[Cenzato, Guerrini, L., Rosone, DCC 2023]

Conclusions

Conclusions

- there are different ways of computing the BWT of a string collection
- these are **non-equivalent**
- the most commonly used ones are **input-order dependent**
- the number of runs r **varies significantly**

Conclusions

- there are different ways of computing the BWT of a string collection
- these are **non-equivalent**
- the most commonly used ones are **input-order dependent**
- the number of runs r **varies significantly**
- \implies different tools on the same dataset can produce different size data structures

Conclusions

- there are different ways of computing the BWT of a string collection
- these are **non-equivalent**
- the most commonly used ones are **input-order dependent**
- the number of runs r **varies significantly**
- \implies different tools on the same dataset can produce different size data structures
- \implies the same tool on the same dataset can produce different size data structures

Conclusions

- there are different ways of computing the BWT of a string collection
- these are **non-equivalent**
- the most commonly used ones are **input-order dependent**
- the number of runs r **varies significantly**
- \implies different tools on the same dataset can produce different size data structures
- \implies the same tool on the same dataset can produce different size data structures
- **optBWT minimizes r** , and has been implemented

Conclusions

- there are different ways of computing the BWT of a string collection
- these are **non-equivalent**
- the most commonly used ones are **input-order dependent**
- the number of runs r **varies significantly**
- \implies different tools on the same dataset can produce different size data structures
- \implies the same tool on the same dataset can produce different size data structures
- **optBWT minimizes r** , and has been implemented
- definition of r **should be standardized** (**optBWT** or **colexBWT**)

Open Problems

- upper bound on differences between separator-based BWT variants
- characterize string collections for which differences highest
- analyze differences between eBWT and separator-based BWTs

Open Problems

- upper bound on differences between separator-based BWT variants
- characterize string collections for which differences highest
- analyze differences between eBWT and separator-based BWTs

My personal conclusion:

Definitions matter!

Acknowledgements

- Davide Cenzato and Zsuzsanna Lipták: *A survey of BWT variants for string collections*, arXiv:2202.13235 (conf. version: CPM 2022)
github.com/davidecenzato/BWT-variants-for-string-collections
- Davide Cenzato, Veronica Guerrini, Zsuzsanna Lipták, and Giovanna Rosone: *Computing the optimal BWT for very large string collections*, DCC 2023.
github.com/davidecenzato/optimalBWT

Acknowledgements

- Davide Cenzato and Zsuzsanna Lipták: *A survey of BWT variants for string collections*, arXiv:2202.13235 (conf. version: CPM 2022)
github.com/davidecenzato/BWT-variants-for-string-collections
- Davide Cenzato, Veronica Guerrini, Zsuzsanna Lipták, and Giovanna Rosone: *Computing the optimal BWT for very large string collections*, DCC 2023.
github.com/davidecenzato/optimalBWT

Thank you for your attention!