

# An Extension of the Burrows Wheeler Transform and Applications to Sequence Comparison and Data Compression

Sabrina Mantaci

Antonio Restivo

Giovanna Rosone

Marinella Sciortino

Università di Palermo

# Outline

- A short description of the Burrows Wheeler Transform (BWT)
- Extension of the BWT to  $k$  words
- Sequences comparison by the extended BWT
- Simultaneous compression of  $k$  texts
- Experimental validation
  - of similarity on biological sequences
  - of compression on some files of the Calgary Corpus

# How does BWT work?

- **INPUT:**  $w = abra\color{red}{ca}$
- lexicographically sort all the cyclic shifts or **conjugates** of  $w$

	<b>F</b>		<b>L</b>	
	↓		↓	
0	<i>a a b r a c</i>		<i>c</i>	
1	<i>a b r a c a</i>	←	<i>a</i>	<b>I</b>
2	<i>a c a a b r</i>		<i>r</i>	
3	<i>b r a c a a</i>		<i>a</i>	
4	<i>c a a b r a</i>		<i>a</i>	
5	<i>r a c a a b</i>		<i>b</i>	

The following properties hold:

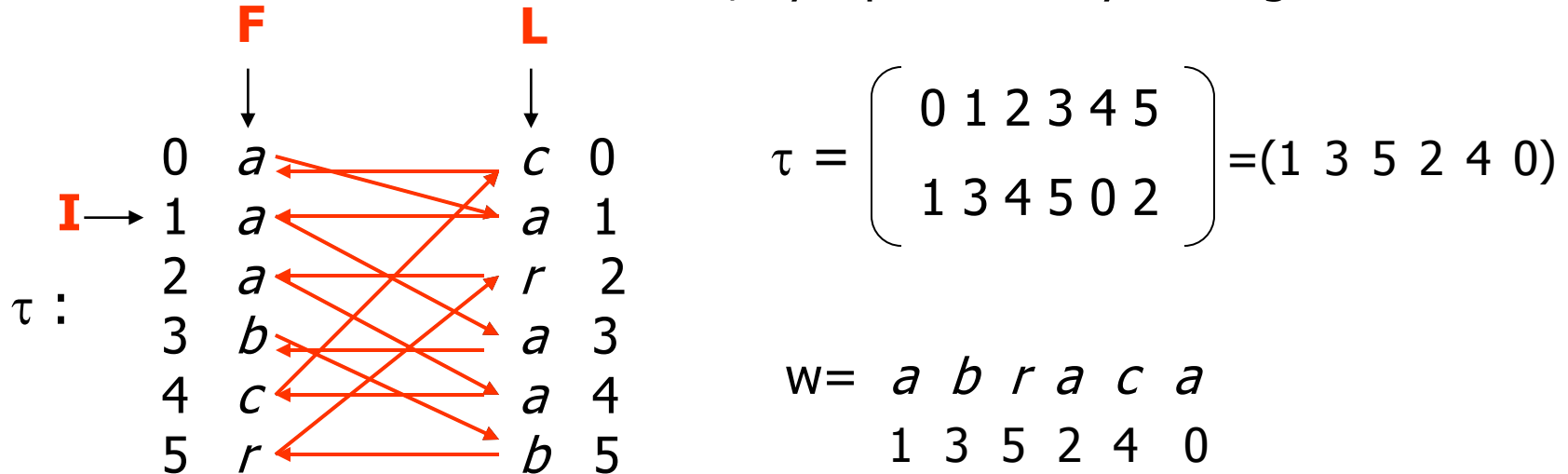
- For all  $i \neq I$ , the character  $L[i]$  is followed in  $w$  by  $F[i]$ ;
  - for any character  $ch$ , the  $i$ -th occurrence of  $ch$  in  $F$  corresponds to the  $i$ -th occurrence of  $ch$  in  $L$ .
- **OUTPUT:**  $BWT(w) = L = \color{red}{caraab}$  and the index  $I = \color{red}{1}$ , that denotes the position of the original word  $w$  after the lexicographical sorting of its conjugates.

# Reversibility of the BWT

The Burrows-Wheeler Transform is reversible, in the sense that, given  $BWT(w)$  and an index  $I$ , it is possible to recover the original word  $w$ .

Given  $L=BWT(w)=\textit{caraab}$  and  $I=1$ :

- We construct the first column  $F$ , by alphabetically sorting the letters in  $L$ .



- A permutation  $\tau$  is defined on the set  $\{0,1,\dots,n-1\}$ , establishing a correspondence between the positions of the same letter in  $F$  and in  $L$ ;
- Starting from position  $I$ , we can recover  $w = F(\tau^0(I))F(\tau^1(I))\dots F(\tau^{n-1}(I))$  where  $\tau^0(x)=x$ ,  $\tau^{i+1}(x) = \tau(\tau^i(x))$

# Remarks

If we do not care about the index, the BWT defines a correspondence between the set of conjugacy classes of words over  $A$  and the set of words over  $A$ . In this sense we have that:

- 1:** The BWT **is injective**. Every conjugacy class have a different transformation.
- 2:** The Burrows-Wheeler transform **is not surjective**, in the sense that there exist words in  $A^*$  that are not image by the BWT of any conjugacy class. Take for instance *bccaaab*.
- 3:** The permutation defining the BWT **is always a cycle with a descent number less than or equal to the size of the alphabet minus one**.

# The BWT and Combinatorics on words

- Relationship with Standard words.

[S. Mantaci, A. Restivo, M. Sciortino, "Burrows-Wheeler Transform and Sturmian Words", Information Processing Letters, 2003].

- Relationship with combinatorics on permutations.

[M. Crochemore, J. Désarménien, D. Perrin, "A note on the Burrows-Wheeler Transform", Theoretical Computer Science, 2005];

[M. Gessel, C. Reutenauer, "Counting permutations with given cycle structure and descent set", J. Comb. Theory, 1993].

# A new order relation

Let  $A$  be a finite alphabet and let  $A^*$  denote the set of the words over  $A$ .

A word  $u$  in  $A^*$  is **primitive** if the condition  $u=w^n$  implies  $u=w$  and  $n=1$ .

If  $u \in A^*$ , we denote by  $u^\omega = uuuu\dots$

For every word  $v$ , there exists a unique primitive word  $w$  and an integer  $k$  such that  $v=w^k$ . By notation,  $w = \text{root}(v)$  and  $k = \text{exp}(v)$

**Definition:** Let  $u$  and  $v$  be two primitive words.

$$u \leq_{\omega} v \iff \begin{cases} \text{exp}(u) < \text{exp}(v) \text{ if } \text{root}(u) = \text{root}(v) \\ u^\omega <_{\text{lex}} v^\omega \text{ otherwise} \end{cases}$$

Where  $u^\omega <_{\text{lex}} v^\omega$  denotes the usual lexicographic order between infinite words.

# A new order relation

Notice that the  $<_{\omega}$  order is different than the lexicographic one. For instance  $ab <_{\text{lex}} aba$ , but  $aba \leq_{\omega} ab$ , since  $abaabaaba... <_{\text{lex}} ababab...$

**Proposition:** Given two primitive words  $u$  and  $v$ ,  
 $u \leq_{\omega} v \iff \text{pref}_k(u^{\omega}) <_{\text{lex}} \text{pref}_k(v^{\omega})$   
where  $k = |u| + |v| - \text{gcd}(|u|, |v|)$ .

The bound is tight: indeed  $abaababa \leq_{\omega} abaab$  because

$\overbrace{abaababaaba} \overbrace{bab...}$   
 $\overbrace{abaababaaba} \overbrace{a...}$

differ for the character in position  $12 = 8 + 5 - 1$ .



# The extended transform E

INPUT:  $S = \{abac, cbab, bca, cba\}$ .

- Sort all the conjugates of the words in  $S$  by the  $\leq_\omega$  order relation;
- Consider the sequence of the sorted words and take the word  $S'$  obtained by concatenating the last letter of each word;
- Take the set  $I$  containing the positions of the words corresponding to the ones in  $S$ ;
- The output of the E transformation is the pair  $(S', I)$ .

$a b a c a b \dots$	$\rightarrow$	1	$a b a c$
$a b c a b c \dots$		2	$a b c$
$a b c b a b \dots$		3	$a b c b$
$a c a b a c \dots$		4	$a c a b$
$a c b a c b \dots$		5	$a c b$
$b a b c b a \dots$		6	$b a b c$
$b a c a b a \dots$		7	$b a c a$
$b a c b a c \dots$		8	$b a c$
$b c a b c a \dots$	$\rightarrow$	9	$b c a$
$b c b a b c \dots$		10	$b c b a$
$c a b a c a \dots$		11	$c a b a$
$c a b c a b \dots$		12	$c a b$
$c b a b c b \dots$	$\rightarrow$	13	$c b a b$
$c b a c b a \dots$	$\rightarrow$	14	$c b a$

OUTPUT:  $(c b b b c a c a a b b a, \{1, 9, 13, 14\})$

# Properties of the E transform

INPUT:  $S = \{abac, cbab, bca, cba\}$ .

→ 1	<i>a b a c</i>
2	<i>a b c</i>
3	<i>a b c b</i>
4	<i>a c a b</i>
5	<i>a c b</i>
6	<i>b a b c</i>
7	<i>b a c a</i>
8	<i>b a c</i>
→ 9	<i>b c a</i>
10	<i>b c b a</i>
11	<i>c a b a</i>
12	<i>c a b</i>
→ 13	<i>c b a b</i>
→ 14	<i>c b a</i>

- In any row  $i \notin I$ , the first symbol follows the last one, in a word in  $S$ .
- For each character, the  $i$ -th occurrence in the first column corresponds to the  $i$ -th occurrence in the red column.

# The inverse transformation

Given  $E(S)=(\text{ccbbbcacaaabba}, \{1,9,13,14\})$ , the following permutation is defined:

$$\begin{array}{cccccccccccccccc}
 F = & a & a & a & a & a & b & b & b & b & b & c & c & c & c \\
 & \textcircled{1} & 2 & 3 & 4 & 5 & 6 & 7 & 8 & \textcircled{9} & 10 & 11 & 12 & \textcircled{13} & \textcircled{14} \\
 & \left[ \begin{array}{cccccccccccc}
 7 & 9 & 10 & 11 & 14 & 3 & 4 & 5 & 12 & 13 & 1 & 2 & 6 & 8
 \end{array} \right. \\
 L = & c & c & b & b & b & c & a & c & a & a & a & b & b & a
 \end{array}$$

Consider the cyclic decomposition of the permutation and read the corresponding letters in F:

$$\begin{array}{cccc}
 \textcircled{1} & 7 & 4 & 11) & \textcircled{2} & \textcircled{9} & 12) & (3 & 10 & \textcircled{13} & 6) & (5 & \textcircled{14} & 8) \\
 a & b & a & c & a & b & c & a & b & c & b & a & c & b
 \end{array}$$

# Bijectivity

Let  $M$  be the family of multisets of conjugacy classes of primitive words of  $A^*$ . Then, if we don't care about the indices

$$E: M \longrightarrow A^*$$

**Theorem:** The transformation  $E$  is injective.

**Theorem:** For each word  $u \in A^*$ , there exists a multiset  $S \in M$  such that  $E(S) = u$ .

For instance,  $E(ab, abcac) = (bccaaab)$ .

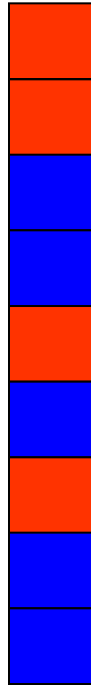
**Theorem [Gessel-Reutenauer]:** There exists a bijection between  $A^*$  and the family of multisets of conjugacy classes of primitive words in  $A^*$ .

# A distance measure between words

$u = bcaaa$

$v = ccbab$

$aabc$   
 $abca$   
 $abccb$   
 $babcc$   
 $bcaa$   
 $bccba$   
 $caab$   
 $cbabc$   
 $ccbab$



u  
u  
v  
v  
u  
v  
u  
v  
v

↓  
 $\gamma(u, v) = u^2 v^2 uvuv^2$

$\delta(u, v) = 3$

In general:

$$\gamma(u, v) = u^{n_1} v^{n_2} u^{n_3} \dots v^{n_k}$$

Definition

$$\delta(u, v) = \sum_{\substack{i=1, \\ n_i \neq 0}}^k (n_i - 1)$$

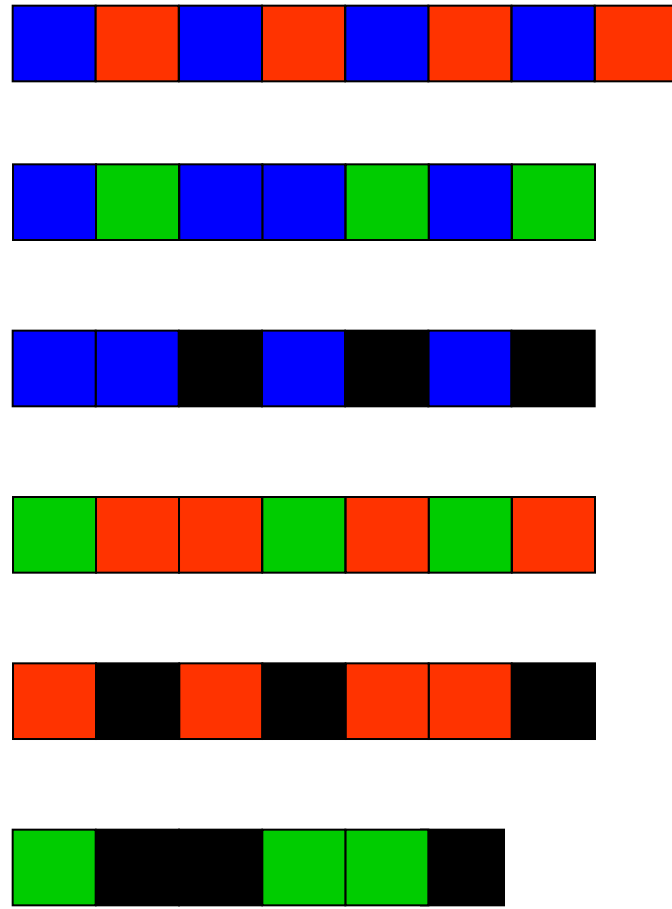
# Properties of the distance measure $\delta$

- $\delta(u,v)=\delta(v,u)$ , that is the distance measure  $\delta$  is symmetric.
- If  $u$  and  $v$  are conjugates, then  $\delta(u,v)=0$ . In fact in this case  $\gamma(u,v)=(uv)^{|u|}$ , then  $n_i=1$  for all  $i=1,\dots,2|u|$ .
- If  $u'$  is a conjugate of  $u$  and  $v'$  is a conjugate of  $v$ , then  $\delta(u,v)=\delta(u',v')$ . Then,  $\delta$  is a similarity measure between conjugacy classes.
- $\delta(u,v)=0$  does not imply that  $u$  and  $v$  are conjugate. For instance, if  $u=abc$  and  $v=abcb$ ,  $\gamma(u,v)=(uv)^4$ . Then  $\delta(u,v)=0$ .
- $\delta$  does not satisfy the triangle inequality. For instance if  $u=abaab$ ,  $v=babab$ ,  $z=abbba$ , we have  $\delta(u,v)=6$ ,  $\delta(v,z)=3$  and  $\delta(u,z)=2$ .

# Multiple sequence comparison

INPUT:  $S = \{u = abac, v = cbab, w = bca, z = cba\}$ .

*a b a c*  
*a b c*  
*a b c b*  
*a c a b*  
*a c b*  
*b a b c*  
*b a c a*  
*b a c*  
*b c a*  
*b c b a*  
*c a b a*  
*c a b*  
*c b a b*  
*c b a*



$\delta(u, v) = 0$   
 $\delta(u, w) = 1$   
 $\delta(u, z) = 1$   
 $\delta(v, w) = 1$   
 $\delta(u, z) = 1$   
 $\delta(w, z) = 1$

# Multiple sequence comparison

- We can compute the distance  $\delta$  of all pairs taken out from a set  $S$  of  $k$  sequences of length  $n$  by simultaneously applying the transformation  $E$  to the entire set  $S$ .
- In order to obtain the  $k \times k$ -matrix of distances we can perform a single sorting of  $kn$  sequences of length  $n$  instead of  $k^2$  sortings of  $2n$  sequences of length  $n$ .
- We can define the notion of **distance between sets**  
If  $S$  and  $T$  are two sets, we encode each conjugate of an element of  $S$  and  $T$  by  $U$  and  $V$ , respectively. If  $\gamma(S, T) = U^{n_1} V^{n_2} U^{n_3} \dots V^{n_k}$  is the sequence so obtained, we define

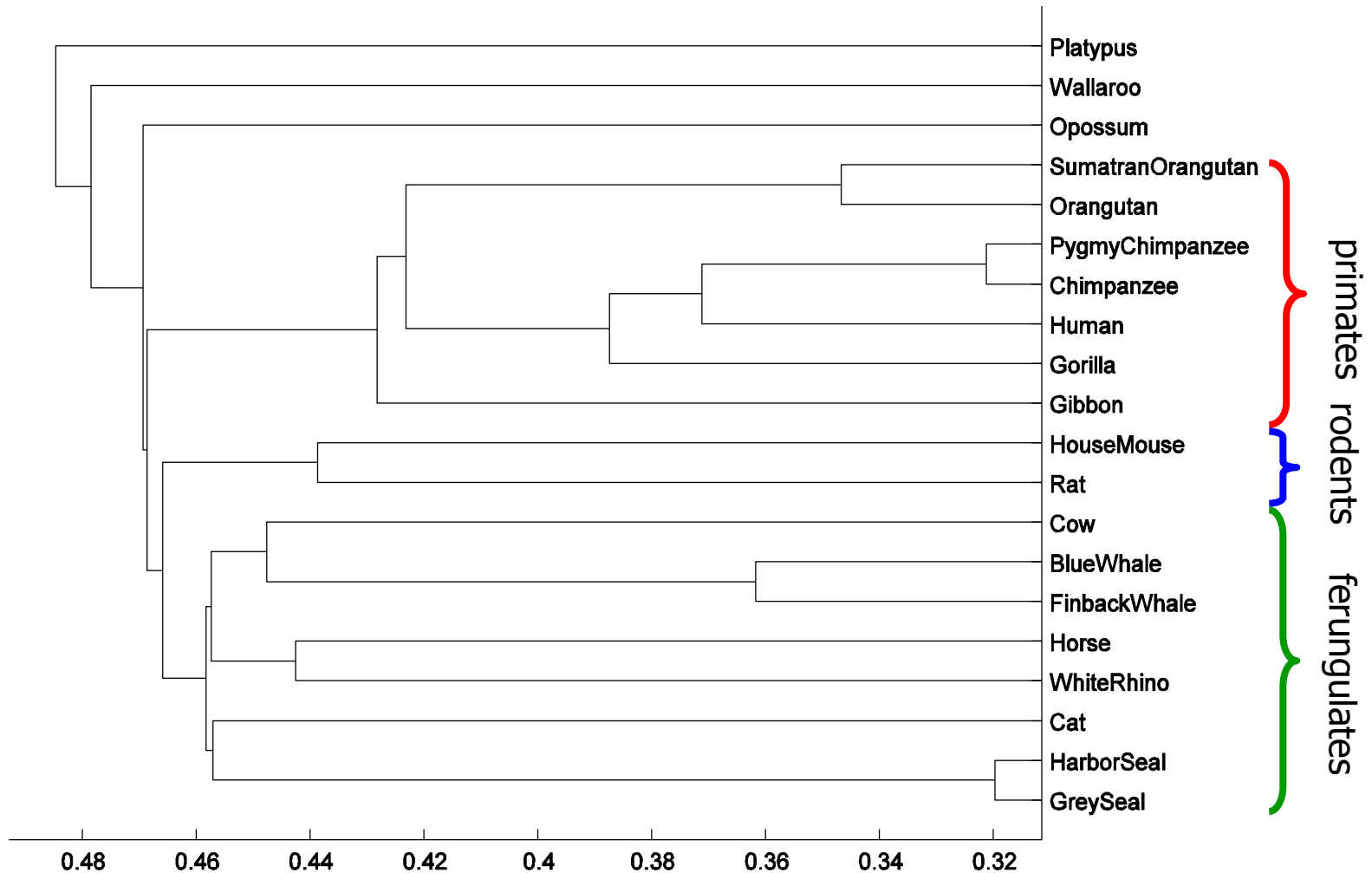
$$\delta(S, T) = \sum_{\substack{i=1, \\ n_i \neq 0}}^k (n_i - 1)$$



# Comparison of biological sequences

- $\delta$  is an example of alignment free distance measure.
- $\delta$  measures how dissimilar two conjugacy classes of sequences are.
- In order to test our method we applied the normalized version of our distance to the whole mitochondrial genome phylogeny.
- The results we have obtained are very close to the ones derived, with other approaches in most of the papers in which the considered species are the same.

# Evolutionary tree built from mtDNA sequences of 20 species



# Simultaneous compression of k texts

- The BWT was introduced as a tool in order to get a word easier to compress.
- We use the transformation E as a preprocessing for the simultaneous compression of a set of k texts.
- If  $\{x_1, x_2, \dots, x_k\}$  is a multiset of words, we denote by  $C(x_1, x_2, \dots, x_k)$  the word obtained by applying a compressor C to the output of  $E(\{x_1, x_2, \dots, x_k\})$ .
- If  $\pi$  is a permutation on  $\{1, \dots, k\}$ , then

$$C(x_1, \dots, x_k) = C(x_{\pi(1)}, x_{\pi(2)}, \dots, x_{\pi(k)})$$

# Simultaneous compression of k texts

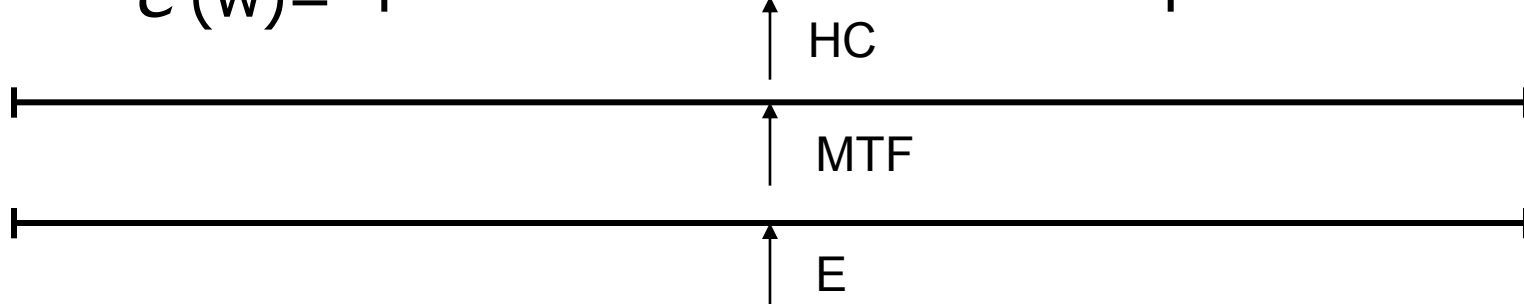
- In most of practical cases, if X and Y are two multisets of words, then

$$|C(XUY)| < |C(X)| + |C(Y)|$$

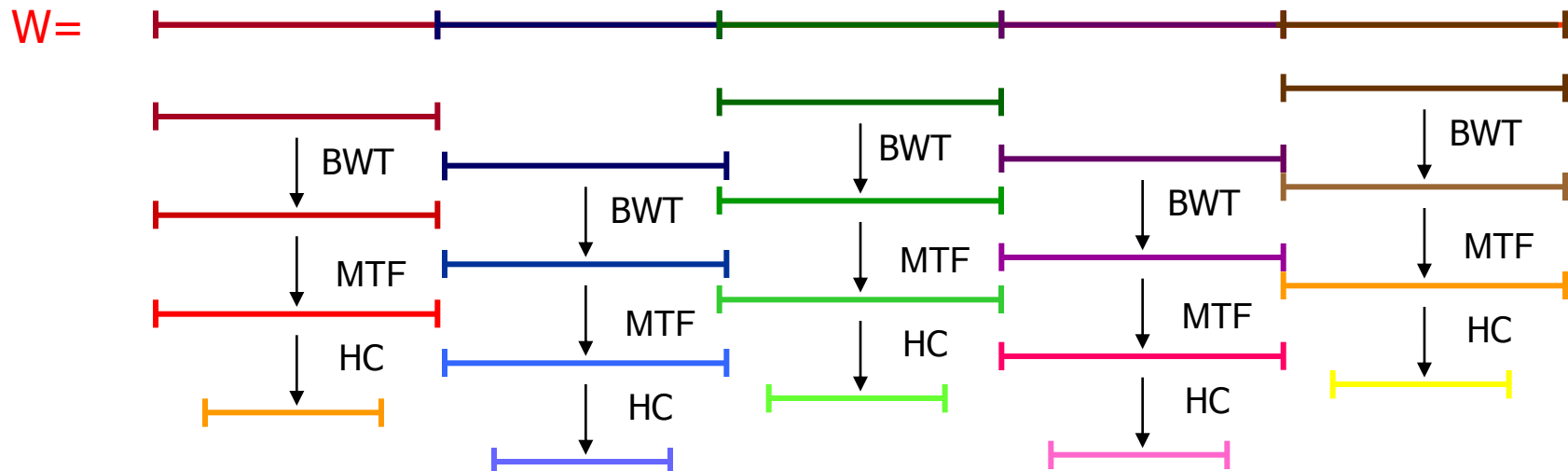
This means that the simultaneous compression of  $\{x_1, x_2, \dots, x_k\}$  by using E as preprocessing is better than compressing each word  $x_i$  separately and concatenating the outputs of the compressed words.

# An E-based compressor

$C(w) =$



{ }



$A(w) =$

# Simultaneous compression of k texts

*A* is a BWT-based compressor (BWT+MTF+HC).

*C* is the E-based compressor.

*M* and *N* are the size of the blocks and the whole text, respectively.

We compare the compression ratios (expressed as output bits per input character) obtained with some files of the Calgary Corpus.

File	Size(in bytes)	Alg	M=16K	M=64K	M=N
bib	111261	<i>C</i>	2.547	2.461	2.425
		<i>A</i>	3.204	2.634	
obj1	21504	<i>C</i>	4.743		4.740
		<i>A</i>	5.076		
paper2	82199	<i>C</i>	2.805	2.786	2.779
		<i>A</i>	3.330	2.917	
progl	71646	<i>C</i>	2.145	2.138	2.131
		<i>A</i>	2.440	2.200	
trans	93695	<i>C</i>	2.064	1.978	1.950
		<i>A</i>	2.667	2.123	