

# BWT - Burrows Wheeler Transform

A block-sorting lossless  
data compression algorithm

- Introduction
- Ciphering phase
- Deciphering phase
- Deciphering phase – improvements
- Relevant compression techniques
- Effectiveness of BWT
- Comparison with other compression programmes
- References

Alberto Chiusole, s249223

Advanced Programming and Algorithmic Design @ uniTS – A.Y. 2017-2018

# Introduction to BWT

- Used to preprocess strings before actual compression
- Produces output with high Index of Coincidence (IC)

==

Identical characters are often in groups

- Compression is easy, i.e. with *run-length* encoding

# Ciphering phase (1)

- Encode the string ``mississippi``, length ``N``
- Add separator symbol (End Of Line) at the end; it must not be present in the string
- ``S = 'mississippi' + `$``

# Ciphering phase (2)

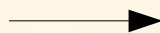
- Create all the cyclic rotations/circular shifts of the string `S`
- Single character permutations; move the first character after the last

```
m i s s i s s i p p i $  
i s s i s s i p p i $ m  
s s i s s i p p i $ m i  
s i s s i p p i $ m i s  
i s s i p p i $ m i s s  
s s i p p i $ m i s s i  
s i p p i $ m i s s i s  
i p p i $ m i s s i s s  
p p i $ m i s s i s s i  
p i $ m i s s i s s i p  
i $ m i s s i s s i p p  
$ m i s s i s s i p p i
```

# Ciphering phase (3)

- Sort all the permutations from top to bottom
- F and L are the First and Last columns of the matrix

```
m i s s i s s i p p i $  
i s s i s s i p p i $ m  
s s i s s i p p i $ m i  
s i s s i p p i $ m i s  
i s s i p p i $ m i s s  
s s i p p i $ m i s s i  
s i p p i $ m i s s i s  
i p p i $ m i s s i s s  
p p i $ m i s s i s s i  
p i $ m i s s i s s i p  
i $ m i s s i s s i p p  
$ m i s s i s s i p p i
```

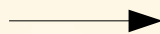


```
      F                               L  
$ m i s s i s s i p p i  
i $ m i s s i s s i p p  
i p p i $ m i s s i s s  
i s s i p p i $ m i s s  
i s s i s s i p p i $ m  
m i s s i s s i p p i $  
p i $ m i s s i s s i p  
p p i $ m i s s i s s i  
s i p p i $ m i s s i s  
s i s s i p p i $ m i s  
s s i p p i $ m i s s i  
s s i s s i p p i $ m i
```

# Ciphering phase (4)

- Take only the last column `L = ipssm\$piissii`
- Encoding done!

m i s s i s s i p p i \$  
i s s i s s i p p i \$ m  
s s i s s i p p i \$ m i  
s i s s i p p i \$ m i s  
i s s i p p i \$ m i s s  
s s i p p i \$ m i s s i  
s i p p i \$ m i s s i s  
i p p i \$ m i s s i s s  
p p i \$ m i s s i s s i  
p i \$ m i s s i s s i p  
i \$ m i s s i s s i p p  
\$ m i s s i s s i p p i



F L  
\$ m i s s i s s i p p i **i**  
i \$ m i s s i s s i p **p**  
i p p i \$ m i s s i s **s**  
i s s i p p i \$ m i s **s**  
i s s i s s i p p i \$ **m**  
m i s s i s s i p p i \$  
p i \$ m i s s i s s i **p**  
p p i \$ m i s s i s s **i**  
s i p p i \$ m i s s i **s**  
s i s s i p p i \$ m i **s**  
s s i p p i \$ m i s s **i**  
s s i s s i p p i \$ m **i**

# Deciphering phase (1)

- The decipher only sees the encoded string `L` but,
- `sort `L`, obtain `F = sort(L) = $iiiiimppssss``

F		L									
\$	m	i	s	s	i	s	s	i	p	p	i
i	\$	m	i	s	s	i	s	s	i	p	p
i	p	p	i	\$	m	i	s	s	i	s	s
i	s	s	i	p	p	i	\$	m	i	s	s
i	s	s	i	s	s	i	p	p	i	\$	m
m	i	s	s	i	s	s	i	p	p	i	\$
p	i	\$	m	i	s	s	i	s	s	i	p
p	p	i	\$	m	i	s	s	i	s	s	i
s	i	p	p	i	\$	m	i	s	s	i	s
s	i	s	s	i	p	p	i	\$	m	i	s
s	s	i	p	p	i	\$	m	i	s	s	i
s	s	i	s	s	i	p	p	i	\$	m	i

# Deciphering phase (2)

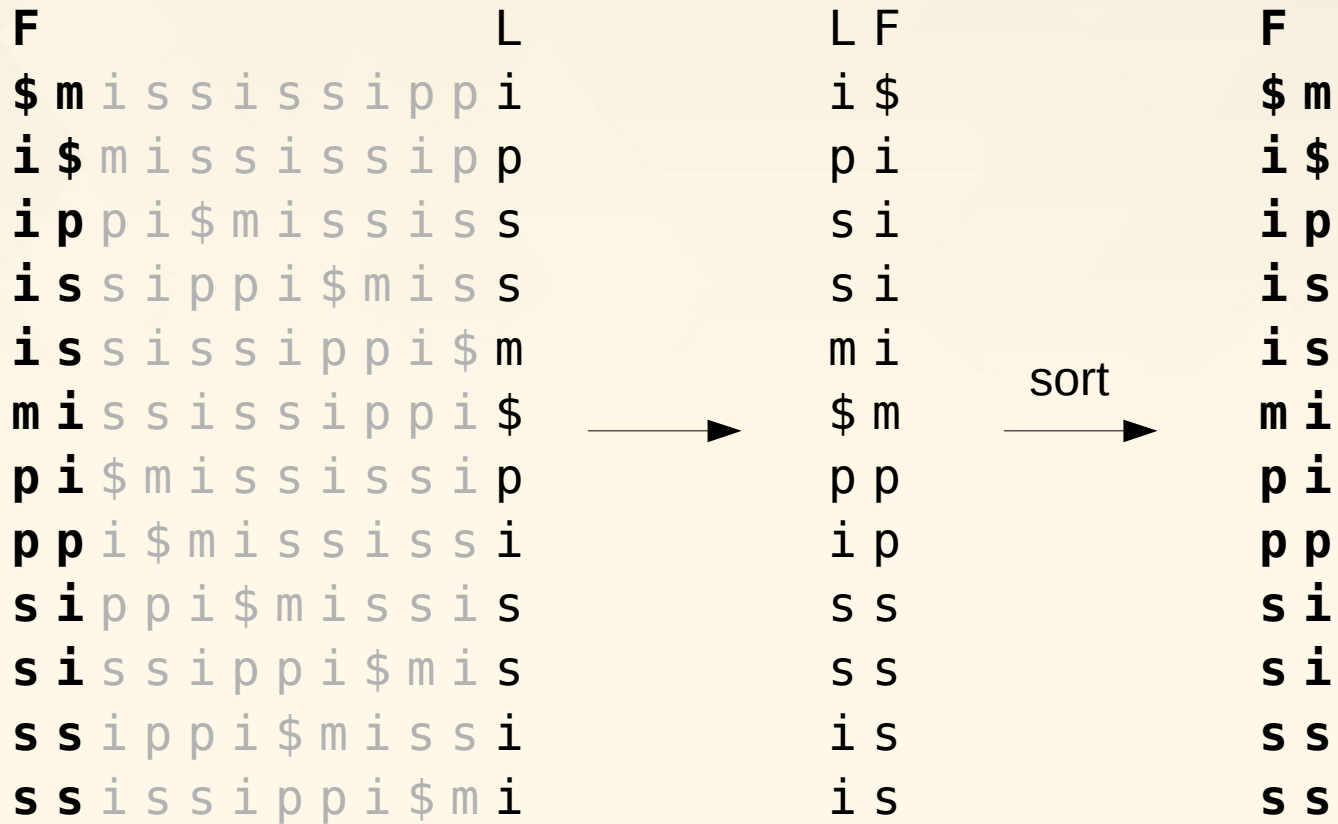
- To rebuild the initial matrix:  
stick together the columns `L - F`

F		L		L F										
\$	m	i	s	s	i	s	s	i	p	p	i	i	\$	
i	\$	m	i	s	s	i	s	s	i	p	p		p	i
i	p	p	i	\$	m	i	s	s	i	s			s	i
i	s	s	i	p	p	i	\$	m	i	s			s	i
i	s	s	i	s	s	i	p	p	i	\$	m		m	i
m	i	s	s	i	s	s	i	p	p	i	\$		\$	m
p	i	\$	m	i	s	s	i	s	s	i	p		p	p
p	p	i	\$	m	i	s	s	i	s				i	p
s	i	p	p	i	\$	m	i	s	s	i	s		s	s
s	i	s	s	i	p	p	i	\$	m	i	s		s	s
s	s	i	p	p	i	\$	m	i	s	s	i		i	s
s	s	i	s	s	i	p	p	i	\$	m	i		i	s



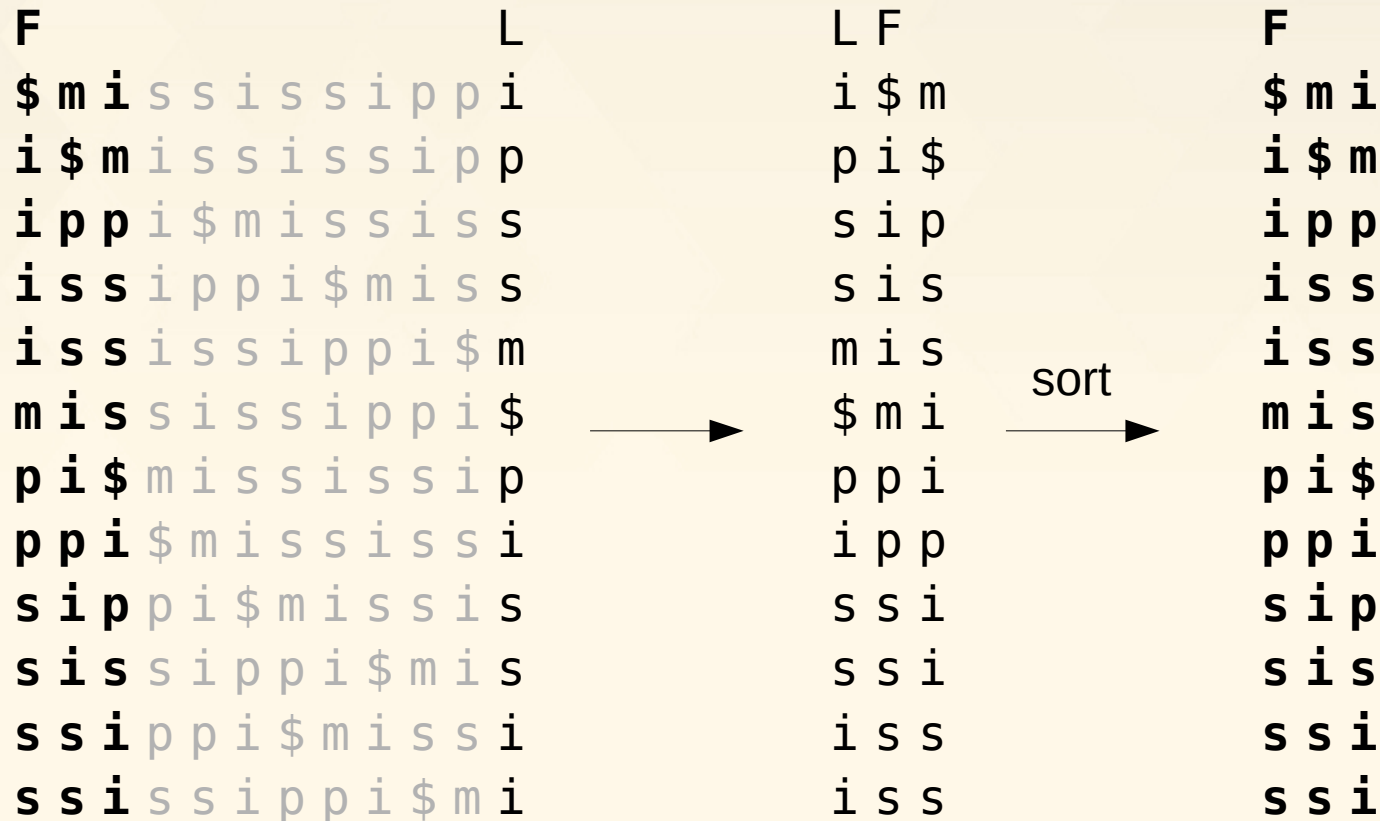
# Deciphering phase (3)

- The 2-mer `L - F` is made of substrings of `S`
- Sort lexicographically → first 2 columns original matrix!



# Deciphering phase (4)

- Take the 3-mer made of the first 2 columns + `L`
- Sort lexicographically → first 3 columns original matrix!



# Deciphering phase (5)

- Repeat until the original matrix is complete
- The string `S` is the one ending with `\$`
- Simple approach, but poor performances:
  - $O(N^2)$  *space* consumption for storing the  $N^2$  matrix
  - $O(N^2 \log N)$  computations:  $N$   $k$ -mers sorting operations with a  $O(\log N)$  sorting algorithm

# Deciphering phase - improvements (1)

- Curious *first-last* property:
- The *relative* order in the group of same letters is preserved: observe, e.g., the relative order of the `i`'s

F		L
\$	m i s s i s s i p p	<b>i</b> <sub>1</sub>
<b>i</b> <sub>1</sub>	\$ m i s s i s s i p p	p <sub>1</sub>
<b>i</b> <sub>2</sub>	p p i \$ m i s s i s s	s <sub>1</sub>
<b>i</b> <sub>3</sub>	s s i p p i \$ m i s s	s <sub>2</sub>
<b>i</b> <sub>4</sub>	s s i s s i p p i \$ m	m <sub>1</sub>
m <sub>1</sub>	i s s i s s i p p i \$	
p <sub>1</sub>	i \$ m i s s i s s i p	p <sub>2</sub>
p <sub>2</sub>	p i \$ m i s s i s s	<b>i</b> <sub>2</sub>
s <sub>1</sub>	i p p i \$ m i s s i s	s <sub>3</sub>
s <sub>2</sub>	i s s i p p i \$ m i s	s <sub>4</sub>
s <sub>3</sub>	s i p p i \$ m i s s	<b>i</b> <sub>3</sub>
s <sub>4</sub>	s i s s i p p i \$ m	<b>i</b> <sub>4</sub>

Proof:

- consider rows starting with a specific letter
- delete the first letter, the columns are still sorted
- apply those first letter after `L`, columns are still sorted
- these new rows are available inside the matrix, and they are still sorted relatively to one another

# Deciphering phase - improvements (2)

- Start from the EOL symbol `\$` in `F`, go backwards into `L`: find `i<sub>1</sub>`
- Search for `i<sub>1</sub>` in `F`, go backwards in `L` and find `p<sub>1</sub>`

F		L
\$	m i s s i s s i p p	<b>i<sub>1</sub></b>
<b>i<sub>1</sub></b>	\$ m i s s i s s i p	<b>p<sub>1</sub></b>
i <sub>2</sub>	p p i \$ m i s s i s	s <sub>1</sub>
i <sub>3</sub>	s s i p p i \$ m i s	s <sub>2</sub>
i <sub>4</sub>	s s i s s i p p i \$	m <sub>1</sub>
m <sub>1</sub>	i s s i s s i p p i	\$
p <sub>1</sub>	i \$ m i s s i s s i	p <sub>2</sub>
p <sub>2</sub>	p i \$ m i s s i s s	i <sub>2</sub>
s <sub>1</sub>	i p p i \$ m i s s i	s <sub>3</sub>
s <sub>2</sub>	i s s i p p i \$ m i	s <sub>4</sub>
s <sub>3</sub>	s i p p i \$ m i s s	i <sub>3</sub>
s <sub>4</sub>	s i s s i p p i \$ m	i <sub>4</sub>

S = .....**pi**

# Deciphering phase - improvements (3)

- Search for  $p_1$  in  $F$ , go backwards in  $L$  and find  $p_2$
- Search for  $p_2$  in  $F$ , go backwards in  $L$  and find  $i_2$

F		L
\$	m i s s i s s i p p i	$i_1$
$i_1$	\$ m i s s i s s i p p	$p_1$
$i_2$	p p i \$ m i s s i s s	$s_1$
$i_3$	s s i p p i \$ m i s s	$s_2$
$i_4$	s s i s s i p p i \$ m	$m_1$
$m_1$	i s s i s s i p p i \$	
<b><math>p_1</math></b>	<b>i</b> \$ m i s s i s s i	<b><math>p_2</math></b>
<b><math>p_2</math></b>	<b>p</b> i \$ m i s s i s s	<b><math>i_2</math></b>
$s_1$	i p p i \$ m i s s i s	$s_3$
$s_2$	i s s i p p i \$ m i s	$s_4$
$s_3$	s i p p i \$ m i s s i	
$s_4$	s i s s i p p i \$ m i	

$S = \dots\dots\dots\mathbf{ippi}$

# Deciphering phase - improvements (4)

- No need to rebuild the original matrix
- Linear memory consumption:  $2N \rightarrow O(N)$
- Linear number of computations to build the 2 first-last arrays

# Relevant compression techniques (1)

- *Move-to-front transform* encoding can be applied after BWT and can improve Index Coincidence
- using the alphabet a . . z, take a string (e.g.,  
`S = ipssm\$piissi` ), note the index of each element  
but also modify the alphabet by moving to the front the last char of the alphabet used
- There will be lots of low index values

$$\text{MTF}(S, \text{"a..z+$"}) = [8, 15, 18, 0, 14, 26, 3, 4, 4, 0, 1, 0]$$



## Relevant compression techniques (2)

- An encoded `L` string can be compressed with, e.g., *run-length encoding*.
- For long strings, it reduces the memory needed, especially with short alphabets (e.g., DNA strings)

`RLE('ipssm$piissii') = i1 p1 s2 m1 $1 p1 s2 i2`

# Effectiveness of BWT

- Example: Normal text contains many occurrences of the word `the`
- create the BWT matrix -> lots of rows starting with `he` and ending with `t`

mlpeeHehhemmigsdmgrgfmeeeeeeeetmierpvttpsssHHchcmthlllrrcvivfnteleDeMMMMMMMMWlphmwhhmMhpw  
chtwmmempwWrWwrrWrtlcehhsthhhhhhhhhhhhhhbwwwwtlrrbfhhHfrrrrrsmeHeerllwdmldwtwwdrmmibrmumia  
nnnnonirrnaieutiuiiuaaiiuttrsiiaieseieaienaa'eeeoleen'nennnlnninna'n'annineno'a'nnannoiinonnnnnne"n'nnnnn  
o'oreneanseiiidnnoelnrrsvthhhhchhrltrlnhuhcnhmhhhlmggchmskhhhhmvslhsgshrtkhhhhhhhhhrhbctspmdvgce  
hhnktthhhhhhHsheebrhmmhrrhrhkkkkdrktggukbhshchhmc nmtldrfrercsyiltgkmPdhdurpppppsfpp lwahr  
rjrtidigtvmtldlchhhSwtbhriBbrfsebwWtWfferld vvvzmh'vehmtufDst tmvrNhhpbhvhrhhdBBBBBBhihh hhhtuttg  
svtmmoodvtyrhdrlli uujlilepggmmremrii d oooooooooeoOooooooooOoofol i r fooin ee n nnnnnnnnnnn  
aAanaarnanooiuiiiuiiiinoiiccctgcttctctctccttcgttcttctccctccSWTttTtWwWtttTttttttttttTttttttttttttTttttttsWw  
ttttswtTntwwWwwwctttstTttttttttwWwcttWwwGTttttgggggggTwWWtrtlgswhhhhhhstrtDaDDDsvrrfhftdollthnm  
nmmnresfflllllwtoahltnaatmmKhmdKkkKKbphhaaoTTTTnttlhaspfkhhhhhHThhhhhhhhTh h hDruhlox s  
erwwwWwwwWbltdebarsrcalnroaaonaaacraeiaiiiaoislocllllallallouousurabbcttbraspwtmmeaerccuaaeieiaee  
uaaobiufoitroeioioonoiooaeaoaioaoroyeoiooaaelsruoeiiiooelooeaeaeoeewaeaoagioeaiioiellioelilowirwgioo  
oweueeiiiieaoaAaaaaAaaaaaaAaaaAaaaaaaAaaiaueoaaaooiuiiiuoiiiiiaaiauiiieankkkkoaioeioaeaaiaeaU  
ueueetStShtsSnntttttTtstttthsiitsioooooGbDgtTlmtjmoopoPhrsrtrRcccssRcCssccconiiipiwgggrcimfGttlLLdcFfF  
HHHHHHHHrHHHHfsmffffwwffspFFFFNNNthhhhhrrmmMmmphlhnrrrrrrrrmmDd mhhyyhYrvrhwc m  
hsbithrcnnHnnNnt nrdrnumrme ossssSppsmommsU aa m r miuemisn naiaeueoeuoeuuueeuouuaauuuueea

*BWT of ~700 words of Shakespeare's Hamlet*

# Comparison with other compression progr.

Program name	CPU Time (s)		Average bit/char
	Compression	Decompression	
compress	9.6	5.2	3.63
gzip	42.6	4.9	2.71
<i>BWT + Huffman coder</i>	<i>51.1</i>	<i>9.4</i>	<i>2.43</i>
comp-2	603.2	614.1	2.47

BWT is widely used, e.g., in the bzip2 compressor, for its efficiency and great advantages.

# References

- Burrows, Michael; Wheeler, David J. (1994),  
[A block sorting lossless data compression algorithm](#),  
Technical Report 124, Digital Equipment Corporation