Lossless compression algorithms

Guillaume TOCHON

guillaume.tochon@lrde.epita.fr

LRDE, EPITA





Data before compression and after decompression are strictly identical.



Data before compression and after decompression are strictly identical.



Lossless compression algorithms exploit data statistical redundancy in two steps : 1) Derive a statistical model for data to compress :

- \rightarrow Static models : analyze whole data and build model according to it (Huffman, bzip2...).
- \rightarrow *Adaptive* models : start with a trivial model and improve it during compression (adapative Huffman, LZW...).

Data before compression and after decompression are strictly identical.



Lossless compression algorithms exploit data statistical redundancy in two steps : 1) Derive a statistical model for data to compress :

- \rightarrow Static models : analyze whole data and build model according to it (Huffman, bzip2...).
- \rightarrow *Adaptive* models : start with a trivial model and improve it during compression (adapative Huffman, LZW...).
- 2) Use statistical model to encode input data :

probable symbol \Leftrightarrow short encoding improbable symbol \Leftrightarrow longer encoding

Data before compression and after decompression are strictly identical.



Lossless compression algorithms exploit data statistical redundancy in two steps : 1) Derive a statistical model for data to compress :

- \rightarrow Static models : analyze whole data and build model according to it (Huffman, bzip2...).
- \rightarrow *Adaptive* models : start with a trivial model and improve it during compression (adapative Huffman, LZW...).
- 2) Use statistical model to encode input data :

probable symbol \Leftrightarrow short encoding

improbable symbol \Leftrightarrow longer encoding

 \Rightarrow Super effective on noise-free data (text documents, source codes, synthetic images).

Lossless compression algorithms

- RLE compression algorithm
 - General idea
 - Example: fax transmission
- 2 Huffman compression algorithm
 - General idea
 - An illustrative example
 - Properties and limits of Huffman encoding
- 3 bzip2 compression algorithm
 - General idea
 - Burrows-Wheeler transform
 - Move-to-front transform

4 LZW compression algorithm

- General idea
- Example of LZW encoding and decoding

What would be the most naive way to compress the following file F?

$$F = \left\{ AAAAAABBBBBCCCCCCC \right\}$$

What would be the most naive way to compress the following file F?

$$F = \left\{ \begin{array}{c} AAAAAABBBBBCCCCCCC \\ \Rightarrow F^{\#} = \left\{ \begin{array}{c} 6A4B7C \end{array} \right\} \end{array}$$

RLE compression algorithm

Run-length encoding

Encode *runs* of symbols (*i.e.* when a symbol is repeated several times in a row) as the number of occurrences in the run followed by the single symbol.

$$\dots s_j \underbrace{s_i \, s_i \dots s_j}_{N \text{ times}} s_k \dots \stackrel{\mathsf{RLE}}{\Longrightarrow} \dots s_j \, N \, s_i \, s_k \dots$$

Run-length encoding

Encode *runs* of symbols (*i.e.* when a symbol is repeated several times in a row) as the number of occurrences in the run followed by the single symbol.

$$\dots s_j \underbrace{s_i s_i \dots s_i}_{N \text{ times}} s_k \dots \overset{\mathsf{RLE}}{\Longrightarrow} \dots s_j N s_i s_k \dots$$

Not suited for text compression: HELLO becomes 1H1E2L10 or HE2LO

- $\rightarrow\,$ The runs must be sufficiently long.
- $\rightarrow\,$ Suitable for synthesis/noise-free images. Used in bitmap (.bmp) format.

Run-length encoding

Encode runs of symbols (i.e. when a symbol is repeated several times in a row) as the number of occurrences in the run followed by the single symbol.

$$\dots s_j \underbrace{s_i s_i \dots s_i}_{N \text{ times}} s_k \dots \overset{\mathsf{RLE}}{\Longrightarrow} \dots s_j N s_i s_k \dots$$

Not suited for text compression: HELLO becomes 1H1E2L10 or HE2LO

- \rightarrow The runs must be sufficiently long.
- \rightarrow Suitable for synthesis/noise-free images. Used in bitmap (.bmp) format.



 \rightarrow How to decide whether it is a number of occurrence or a symbol?

Run-length encoding

Encode *runs* of symbols (*i.e.* when a symbol is repeated several times in a row) as the number of occurrences in the run followed by the single symbol.

$$\dots s_j \underbrace{s_i \, s_j \, \dots \, s_j}_{N \text{ times}} s_k \dots \stackrel{\mathsf{RLE}}{\Longrightarrow} \dots s_j \# N \, s_i \, s_k \dots$$

Not suited for text compression: HELLO becomes 1H1E2L10 or HE2LO

- $\rightarrow\,$ The runs must be sufficiently long.
- $\rightarrow\,$ Suitable for synthesis/noise-free images. Used in bitmap (.bmp) format.

Problem when encoding numbers: 22231444457 $\stackrel{\mathsf{RLE}}{\Longrightarrow}$



- $\rightarrow\,$ How to decide whether it is a number of occurrence or a symbol?
- $\rightarrow\,$ Use a special character # to inform the decompression stage.

 $\underline{Fax:}$ prehistoric machine that used to send documents by wire transmission (dated circa 1500 BC).



 $\underline{Fax:}$ prehistoric machine that used to send documents by wire transmission (dated circa 1500 BC).

The document is scanned line by line by the fax. Each line is binarized and converted into an analog signal (electric impulses) that is sent through telephone lines.



 $\frac{Fax:}{mission}$ (dated circa 1500 BC).

The document is scanned line by line by the fax. Each line is binarized and converted into an analog signal (electric impulses) that is sent through telephone lines.



Rough idea: compress each line independently

 $\frac{Fax:}{mission}$ (dated circa 1500 BC).

The document is scanned line by line by the fax. Each line is binarized and converted into an analog signal (electric impulses) that is sent through telephone lines.



Rough idea: compress each line independently

Fax: prehistoric machine that used to send documents by wire transmission (dated circa 1500 BC).

The document is scanned line by line by the fax. Each line is binarized and converted into an analog signal (electric impulses) that is sent through telephone lines.



Rough idea:compress each line independently111111 \rightarrow no text, RLE is super effective. $11\cdots11001111011\cdots11$ \Rightarrow text, but RLE remains effective.

 $\underline{Fax:}$ prehistoric machine that used to send documents by wire transmission (dated circa 1500 BC).

The document is scanned line by line by the fax. Each line is binarized and converted into an analog signal (electric impulses) that is sent through telephone lines.



Smarter idea: compress difference between consecutive lines



 $\underline{Fax:}$ prehistoric machine that used to send documents by wire transmission (dated circa 1500 BC).

The document is scanned line by line by the fax. Each line is binarized and converted into an analog signal (electric impulses) that is sent through telephone lines.



Smarter idea: compress difference between consecutive lines



Huffman compression algorithm

- \rightarrow Proposed in 1952 by David Huffman (during its Ph.D).
- $\rightarrow\,$ Exploit the statistical distribution of the symbols to encode. \Rightarrow entropy encoding.
- $\rightarrow\,$ Frequent symbols are given shorter encoding support.
 - \Rightarrow variable-length code.



David Huffman

Huffman compression algorithm

- $\rightarrow\,$ Proposed in 1952 by David Huffman (during its Ph.D).
- $\rightarrow\,$ Exploit the statistical distribution of the symbols to encode. \Rightarrow entropy encoding.
- $\rightarrow\,$ Frequent symbols are given shorter encoding support.
 - \Rightarrow variable-length code.



David Huffman

Huffman encoding

The encoding is obtained by the following two-steps procedure:

- 1) Compute a binary tree whose leaves are the symbols, by iteratively merging the two symbols with lowest probabilities of occurence.
- 2) Label each left branch by 0 and each right branch by 1 (or the other way around).

The encoding of each symbol is given by the path from the root to the leaf corresponding to the symbol.

Take the following alphabet Σ with known probability distribution:

Si	Рi
*	0.3
0	0.06
0	0.1
Ð	0.1
•	0.4
\$	0.04

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.

Take the following alphabet Σ with known probability distribution:

Si	\mathbb{P}^i	Si	\mathbb{P}^i
*	0.3		0.4
•	0.06	*	0.3
0	0.1	0	0.1
	0.1	S	0.1
	0.4	•	0.06
♪	0.04	ſ	0.04

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.

Take the following alphabet $\boldsymbol{\Sigma}$ with known probability distribution:

Si	\mathbb{P}^i	Si	\mathbb{P}^i
*	0.3		0.4
	0.06	*	0.3
0	0.1	0	0.1
	0.1	S	0.1
	0.4	0	0.06
ل ا	0.04	ĥ	0.04

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Si	Рi	Si	₽i	
*	0.3		0.4	
	0.06	*	0.3	
0	0.1	0	0.1	
	0.1	Ð	0.1	
8	0.4	N_1	0.1	
♪	0.04			

Take the following alphabet Σ with known probability distribution:

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Take the following alphabet $\boldsymbol{\Sigma}$ with known probability distribution:

Si	\mathbb{P}^i	Si	\mathbb{P}^i
*	0.3		0.4
•	0.06	*	0.3
0	0.1	0	0.1
	0.1	•	0.1
	0.4	N_1	0.1
₽	0.04		



- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.

Si	₽i	Si	\mathbb{P}^i
*	0.3		0.4
•	0.06	*	0.3
0	0.1	0	0.1
	0.1	N_2	0.2
	0.4		
_ ۲	0.04		

Take the following alphabet Σ with known probability distribution:

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Si	₽i	Si	\mathbb{P}^i
*	0.3		0.4
•	0.06	*	0.3
0	0.1	N_2	0.2
I	0.1	0	0.1
	0.4		
♪	0.04		

Take the following alphabet Σ with known probability distribution:

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Take the following alphabet Σ with known probability distribution:

Si	\mathbb{P}^i	Si	\mathbb{P}^i
*	0.3	*	0.4
•	0.06	*	0.3
0	0.1	N_2	0.2
	0.1	0	0.1
	0.4		
♪	0.04		

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Take the following alphabet Σ with known probability distribution:

Si	\mathbb{P}_i	Si	Рi
*	0.3	*	0.4
9	0.06	×	0.3
0	0.1	N_3	0.3
	0.1		
	0.4		
J	0.04		

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Take the following alphabet Σ with known probability distribution:

Si	Рi	Si	Рi
*	0.3		0.4
9	0.06	*	0.3
0	0.1	N_3	0.3
	0.1		
	0.4		
♪	0.04		

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Take the following alphabet $\boldsymbol{\Sigma}$ with known probability distribution:

Si	Рi	Si	Рi
*	0.3	•	0.4
•	0.06	N_4	0.6
0	0.1		
	0.1		
8	0.4		
♪	0.04		

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Take the following alphabet $\boldsymbol{\Sigma}$ with known probability distribution:

Si	Рi	Si	Рi
*	0.3	N ₄	0.6
	0.06		0.4
0	0.1		
	0.1		
8	0.4		
ا	0.04		

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Take the following alphabet $\boldsymbol{\Sigma}$ with known probability distribution:

			1			
Si	₽i	Si	₽i	Huffman encoding: 1 st step		
*	0.3	N_4	0.6	1) Sort table (if necessary)		
	0.06	*	0.4	2) Marga the two symbols with lowest proba		
0	0.1			bility of occurence.		
I	0.1			3) Update table.		
	0.4			4) Repeat if more than one symbol remaining.		
5	0.04					



Take the following alphabet $\boldsymbol{\Sigma}$ with known probability distribution:

c .	T D -	C 1	D .	
5,	Pi	51	Pi	
*	0.3	N _r	1	
	0.06			
0	0.1			
	0.1			
	0.4			
ل ا	0.04			

- 1) Sort table (if necessary).
- 2) Merge the two symbols with lowest probability of occurence.
- 3) Update table.
- 4) Repeat if more than one symbol remaining.



Si	\mathbb{P}^i	Encoding	Huffman encoding: 2 nd step
*	0.3		1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1		2) Always assign 1 to the left child and 0 to
	0.1		the right one (or the other way around).
•	0.4		3) Read encoding on the whole branch.
ک	0.04		


S _i	₽i	Encoding	Huffman encoding: 2 nd step
*	0.3		1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1		2) Always assign 1 to the left child and 0 to
	0.1		the right one (or the other way around).
•	0.4		3) Read encoding on the whole branch.
_	0.04		



Si	₽i	Encoding	Huffman encoding: 2 nd step
×	0.3		1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1		2) Always assign 1 to the left child and 0 to
	0.1		the right one (or the other way around).
•	0.4	0	3) Read encoding on the whole branch.
ئ	0.04		



Si	₽i	Encoding	Huffman encoding: 2 nd step
×	0.3		1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1		2) Always assign 1 to the left child and 0 to
	0.1		the right one (or the other way around).
•	0.4	0	3) Read encoding on the whole branch.
ک	0.04		



Si	Рi	Encoding	Huffman encoding: 2 nd step
×	0.3	11	1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1		2) Always assign 1 to the left child and 0 to
	0.1		the right one (or the other way around).
•	0.4	0	3) Read encoding on the whole branch.
_	0.04		



Si	Рi	Encoding	Huffman encoding: 2 nd step
×	0.3	11	1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1		2) Always assign 1 to the left child and 0 to
	0.1		the right one (or the other way around).
•	0.4	0	3) Read encoding on the whole branch.
_	0.04		



Si	Рi	Encoding	Huffman encoding: 2 nd step
×	0.3	11	1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1	100	2) Always assign 1 to the left child and 0 to
	0.1		the right one (or the other way around).
•	0.4	0	3) Read encoding on the whole branch.
ک	0.04		



Si	Рi	Encoding	Huffman encoding: 2 nd step
*	0.3	11	1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1	100	2) Always assign 1 to the left child and 0 to
Ð	0.1		the right one (or the other way around).
•	0.4	0	3) Read encoding on the whole branch.
ل	0.04		



Si	Рi	Encoding	Huffman encoding: 2 nd step
*	0.3	11	1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1	100	2) Always assign 1 to the left child and 0 to
	0.1	1011	the right one (or the other way around).
	0.4	0	3) Read encoding on the whole branch.
د	0.04		



Si	Рi	Encoding	Huffman encoding: 2 nd step
*	0.3	11	1) Traverse the tree from the root to the leaves
	0.06		(the symbols).
0	0.1	100	2) Always assign 1 to the left child and 0 to
	0.1	1011	the right one (or the other way around).
	0.4	0	3) Read encoding on the whole branch.
د	0.04		



Si	Рi	Encoding	Huffman encoding: 2 nd step
×	0.3	11	1) Traverse the tree from the root to the leaves
	0.06	10101	(the symbols).
0	0.1	100	2) Always assign 1 to the left child and 0 to
	0.1	1011	the right one (or the other way around).
•	0.4	0	3) Read encoding on the whole branch.
ل	0.04	10100	



The entropy of $\Sigma = \{ \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet} \}$ with previously given probability distribution is $H \simeq 2.144$ *Sh/symb*.

 \Rightarrow A theoretically optimal encoding needs an average length of 2.144 bits per symbol.

The entropy of $\Sigma = \{ \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet} \}$ with previously given probability distribution is $H \simeq 2.144$ *Sh/symb*. \Rightarrow A theoretically optimal encoding needs an average length of 2.144 bits per symbol.

The average length of Huffman encoding is

$$L = \mathbb{E}[\ell(s_i)] = \sum_{s_i \in \Sigma} p_i \ell(s_i)$$
 bits/symb.

where $\ell(s_i)$ is the encoding length of symbol s_i . $\Rightarrow L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2 \text{ bits/symb.}$

The entropy of $\Sigma = \{ \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet} \}$ with previously given probability distribution is $H \simeq 2.144$ *Sh/symb*. \Rightarrow A theoretically optimal encoding needs an average length of 2.144 bits per symbol.

The average length of Huffman encoding is

$$L = \mathbb{E}[\ell(s_i)] = \sum_{s_i \in \Sigma} p_i \ell(s_i)$$
 bits/symb.

where $\ell(s_i)$ is the encoding length of symbol s_i . $\Rightarrow L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2 \text{ bits/symb.}$

Take some file F with $N_F = 1000$ symbols drawn from Σ .

 \rightarrow 3000 bits are necessary to encode *F* without compression.

The entropy of $\Sigma = \{ \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet} \}$ with previously given probability distribution is $H \simeq 2.144$ *Sh/symb*. \Rightarrow A theoretically optimal encoding needs an average length of 2.144 bits per symbol.

The average length of Huffman encoding is

$$L = \mathbb{E}[\ell(s_i)] = \sum_{s_i \in \Sigma} p_i \ell(s_i)$$
 bits/symb.

where $\ell(s_i)$ is the encoding length of symbol s_i . $\Rightarrow L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2 \text{ bits/symb.}$

Take some file F with $N_F = 1000$ symbols drawn from Σ .

- \rightarrow 3000 bits are necessary to encode *F* without compression.
- \rightarrow The minimal compressed size of F is $H \times N_F = 2144$ bits.

The entropy of $\Sigma = \{ \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet}, \overset{\bullet}{\bullet} \}$ with previously given probability distribution is $H \simeq 2.144$ *Sh/symb*. \Rightarrow A theoretically optimal encoding needs an average length of 2.144 bits per symbol.

The average length of Huffman encoding is

$$L = \mathbb{E}[\ell(s_i)] = \sum_{s_i \in \Sigma} p_i \ell(s_i)$$
 bits/symb.

where $\ell(s_i)$ is the encoding length of symbol s_i . $\Rightarrow L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2 \text{ bits/symb.}$

Take some file F with $N_F = 1000$ symbols drawn from Σ .

- \rightarrow 3000 bits are necessary to encode *F* without compression.
- \rightarrow The minimal compressed size of F is $H \times N_F = 2144$ bits.
- \rightarrow Huffman encoding allows to reach $L \times N_F = 2200$ bits (in average).

The entropy of $\Sigma = \{ \overset{\bullet}{\clubsuit}, \overset{\bullet}{\textcircled{\baselineskiplimits}}, \overset{\bullet}{\clubsuit} \}$ with previously given probability distribution is $H \simeq 2.144$ Sh/symb. \Rightarrow A theoretically optimal encoding needs an average length of 2.144 bits per symbol.

The average length of Huffman encoding is

$$L = \mathbb{E}[\ell(s_i)] = \sum_{s_i \in \Sigma} p_i \ell(s_i)$$
 bits/symb.

where $\ell(s_i)$ is the encoding length of symbol s_i . $\Rightarrow L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2 \ bits/symb.$

Take some file F with $N_F = 1000$ symbols drawn from Σ .

- \rightarrow 3000 bits are necessary to encode F without compression.
- \rightarrow The minimal compressed size of F is $H \times N_F = 2144$ bits.
- \rightarrow Huffman encoding allows to reach $L \times N_F = 2200$ bits (in average).

Properties of Huffman encoding

Huffman encoding is a *prefix code* (no symbol encoding is a prefix of another symbol encoding). \Rightarrow it can be decoded on the fly.



Properties of Huffman encoding

Huffman encoding is a *prefix code* (no symbol encoding is a prefix of another symbol encoding). \Rightarrow it can be decoded on the fly.



Huffman encoding is *optimal* in terms of average encoding length with respect to any symbol-by-symbol encoding with prefix property.

That is, if \mathfrak{L} is another encoding strategy at the symbol level with prefix property, then $L_{Huffman} \leq L_{\mathfrak{L}}$.

Properties of Huffman encoding

Huffman encoding is a prefix code (no symbol encoding is a prefix of another symbol encoding). \Rightarrow it can be decoded on the fly.



Huffman encoding is *optimal* in terms of average encoding length with respect to any symbol-by-symbol encoding with prefix property.

That is, if \mathfrak{L} is another encoding strategy at the symbol level with prefix property, then $L_{Huffman} < L_{\mathfrak{L}}$.

 $H = L_{Huffman}$ if the probabilities of symbols are natural powers of 1/2, otherwise $H < L_{Huffman} < H + 1$.

Huffman encoding reaches the entropy H if $\forall s_i \in \Sigma$, $\mathbb{D}_i = (1/2)^k$, $k \in \mathbb{N}$. Otherwise, Huffman encoding is "close" to the entropy.

Huffman encoding strategy suffers from two mains drawbacks:

- \checkmark It is a symbol-by-symbol encoding \rightarrow does not handle *words* (as sequences of symbols).
 - Ex: In a piece of code, words such as for, while, end, are likely to frequently appear, and should directly be encoded at the word scale.

Huffman encoding strategy suffers from two mains drawbacks:

- **X** It is a symbol-by-symbol encoding \rightarrow does not handle *words* (as sequences of symbols).
 - Ex: In a piece of code, words such as for, while, end, are likely to frequently appear, and should directly be encoded at the word scale.
- X It relies (too) strongly on the probability distribution of the symbols to encode.
 - $\rightarrow~$ The file to compress must be scanned first to estimate the probability distribution.
 - $\rightarrow\,$ The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.

Huffman encoding strategy suffers from two mains drawbacks:

- \checkmark It is a symbol-by-symbol encoding \rightarrow does not handle *words* (as sequences of symbols).
 - Ex: In a piece of code, words such as for, while, end, are likely to frequently appear, and should directly be encoded at the word scale.
- X It relies (too) strongly on the probability distribution of the symbols to encode.
 - $\rightarrow~$ The file to compress must be scanned first to estimate the probability distribution.
 - \rightarrow The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.

 $\hookrightarrow {\sf Classical text statistics can} \\ {\sf be used (language-dependent)}. \\$



Huffman encoding strategy suffers from two mains drawbacks:

- \checkmark It is a symbol-by-symbol encoding \rightarrow does not handle *words* (as sequences of symbols).
 - Ex: In a piece of code, words such as for, while, end, are likely to frequently appear, and should directly be encoded at the word scale.
- X It relies (too) strongly on the probability distribution of the symbols to encode.
 - $\rightarrow~$ The file to compress must be scanned first to estimate the probability distribution.
 - \rightarrow The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.
 - $\hookrightarrow {\sf Classical \ text \ statistics \ can} \\ {\sf be \ used \ (language-dependent)}. }$



 $\rightarrow\,$ And what if the frequencies vary with time?

Huffman encoding strategy suffers from two mains drawbacks:

- \checkmark It is a symbol-by-symbol encoding \rightarrow does not handle *words* (as sequences of symbols).
 - Ex: In a piece of code, words such as for, while, end, are likely to frequently appear, and should directly be encoded at the word scale.
- X It relies (too) strongly on the probability distribution of the symbols to encode.
 - $\rightarrow~$ The file to compress must be scanned first to estimate the probability distribution.
 - \rightarrow The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.
 - $\hookrightarrow {\sf Classical \ text \ statistics \ can} \\ {\sf be \ used \ (language-dependent)}. }$



 $\rightarrow\,$ And what if the frequencies vary with time?

 $\hookrightarrow \mathsf{Adaptive} \ \mathsf{Huffman} \ \mathsf{algorithm}.$

Huffman encoding strategy suffers from two mains drawbacks:

- \checkmark It is a symbol-by-symbol encoding \rightarrow does not handle *words* (as sequences of symbols).
 - Ex: In a piece of code, words such as for, while, end, are likely to frequently appear, and should directly be encoded at the word scale.
- X It relies (too) strongly on the probability distribution of the symbols to encode.
 - $\rightarrow~$ The file to compress must be scanned first to estimate the probability distribution.
 - \rightarrow The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.
 - $\hookrightarrow {\sf Classical text statistics can} \\ {\sf be used (language-dependent)}. \\$



- $\rightarrow~\mbox{And}$ what if the frequencies vary with time?
 - $\hookrightarrow \mathsf{Adaptive} \ \mathsf{Huffman} \ \mathsf{algorithm}.$

Huffman encoding alone is nowadays hardly used for text compression, but serves as the last level of more advanced compression algorithms (bzip2, JPEG, etc).

bzip2 compression algorithm

- \rightarrow Created between 1996 and 2000 by Julian Seward.
- \rightarrow Initially proposed to tackle patenty issues with LZW compression algorithm.
- \rightarrow Free and open-source compression algorithm.

bzip2 encoding

Idea: strenghten the anisotropy of symbol probabilities.

- 1) Apply Burrows-Wheeler transform to convert frequently-recurring symbol sequences into runs of identical symbols.
- 2) Use Move-to-front (MTF) transform to replace each symbol by its index in a dynamic table.
- 3) Encode index sequence with Huffman algorithm.







The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

Ex: Consider the file banana $(\$ \equiv EOF)$

1st step: Construct a table containing all rotations of the string to transform.

banana \$

The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

Ex: Consider the file banana $(\$ \equiv EOF)$

1st step: Construct a table containing all rotations of the string to transform.

b	а	n	а	n	а	\$
\$	b	а	n	а	n	а

The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

Ex: Consider the file banana $(\$ \equiv EOF)$

1st step: Construct a table containing all rotations of the string to transform.

b	а	n	а	n	а	\$
\$	b	а	n	а	n	а
а	\$	b	а	n	а	n

The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

Ex: Consider the file banana[§] ($\$ \equiv EOF$)

1st step: Construct a table containing all rotations of the string to transform.



The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

Ex: Consider the file banana[§] ($\$ \equiv EOF$)

 1^{st} step: Construct a table containing all rotations of the string to transform. 2^{nd} step: Sort rows into lexicographic order.



The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

Ex: Consider the file banana[§] ($\$ \equiv EOF$)

 1^{st} step: Construct a table containing all rotations of the string to transform. 2^{nd} step: Sort rows into lexicographic order.

b	а	n	а	n	а	\$		\$	b	а	n	а	n	а
\$	b	а	n	а	n	а		а	\$	b	а	n	а	n
а	\$	b	а	n	а	n		а	n	а	\$	b	а	n
n	а	\$	b	а	n	а	lexicographic	а	n	а	n	а	\$	b
а	n	а	\$	b	а	n	sort	b	а	n	а	n	а	\$
n	а	n	а	\$	b	а		n	а	\$	b	а	n	а
а	n	а	n	а	\$	b		n	а	n	а	\$	b	а

The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

Ex: Consider the file banana[§] ($\$ \equiv EOF$)

1st step: Construct a table containing all rotations of the string to transform.
2nd step: Sort rows into lexicographic order.
BWT: row number of initial string + last column of sorted table.

															\frown
b	а	n	а	n	а	\$		1	\$	b	а	n	а	n	а
\$	b	а	n	а	n	а		2	а	\$	b	а	n	а	n
а	\$	b	а	n	а	n		3	а	n	а	\$	b	а	n
n	а	\$	b	а	n	а	lexicographic	4	а	n	а	n	а	\$	b
а	n	а	\$	b	а	n	sort	5	b	а	n	а	n	а	\$
n	а	n	а	\$	b	а		6	n	а	\$	b	а	n	а
а	n	а	n	а	\$	b		7	n	а	n	а	\$	b	а

The encoding part

Invented by Michael Burrows and David Wheeler in 1994 to sort strings of characters into runs of similar characters to ease their compression, in a totally reversible way.

Ex: Consider the file banana^{\$} ($\$ \equiv EOF$)

1st step: Construct a table containing all rotations of the string to transform.
2nd step: Sort rows into lexicographic order.
BWT: row number of initial string + last column of sorted table.



Not convinced?

BWT How much wood would a woodchuck chuck if a woodchuck could chuck wood?\$

dfkdkkwhkaad?d\$ udd uuuu llooooiccccc ccccuu oooowwwwwcHmhhhhooo

=

And that sure looks way better. How much better?
Not convinced?

BWT (How much wood would a woodchuck chuck if a woodchuck could chuck wood?\$ = dfkdkkwhkaad?d\$ udd uuuu llooooiccccc ccc-cuu oooowwwwcHmh-

And that sure looks way better. How much better?

 $\mathsf{RLE}\left(\begin{array}{c}\mathsf{dfkdkkwhkaad?d\$ udd}\\\mathsf{uuuu}\ \mathsf{llooooiccccc\ ccc-}\\\mathsf{cuu\ oooowwwwcHmh-}\\\mathsf{hhbccc}\\\mathsf{wcHm}\#4h\#3o\end{array}\right) = \begin{array}{c}\mathsf{dfkd\#2kwhk\#2ad?d\clubsuit}\\\mathsf{u\#2d\ \#4u\ \#2l\#4oi}\\\mathsf{wcHm}\#4h\#3o\end{array}$

Not convinced?

BWT (How much wood would a woodchuck chuck if a woodchuck could chuck wood?§

dfkdkkwhkaad?d\$ udd = uuuu llooooiccccc ccc-cuu oooowwwwcHmhhhhooo

And that sure looks way better. How much better?

RLE $\begin{pmatrix} dfkdkkwhkaad?d\$ udd \\ uuuu llooooiccccc ccc-$ cuu oooowwwwcHmh- $hbbbcco \\ \end{pmatrix} = \begin{pmatrix} dfkd#2kwhk#2ad?d\$ \\ u#2d #4u #2l#4oi \\ #5c #4c#2u #4o#5 \\ wcHm#4h#3o \\ \end{pmatrix}$ hhhooo

wcHm#4h#3o

6 3 8

Transforming 65 characters into ... 61 characters!

Not convinced?

BWT How much wood would a woodchuck chuck if a woodchuck could chuck wood?

dfkdkkwhkaad?d\$ udd = uuuu llooooiccccc ccc-cuu oooowwwwcHmhhhhooo

And that sure looks way better. How much better?

 $\mathsf{RLE}\left(\begin{array}{c}\mathsf{dfkdkkwhkaad?d\$} \mathsf{udd}\\\mathsf{uuuu} \mathsf{llooooiccccc} \mathsf{ccc-}\\\mathsf{cuu} \mathsf{oooowwwwcHmh-}\\\mathsf{bbb} \mathsf{c} \mathsf{udd} \mathsf{udd}\\\mathsf{udd} \mathsf{udd} \mathsf{udd}$

wcHm#4h#3o

Transforming 65 characters into ... 61 characters! 👍 🛎 👌 🤘

BWT is applied in practice on much longer strings thanks to some efficient and optimized implementations \Rightarrow initial "pre-processing" for the MTF transform.

Guillaume TOCHON (LRDE)

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb\$aa before the last column of the table.

a n b \$ a a

The decoding part

It is easily possible to recover **bananas** from **5annb\$aa** by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb\$aa before the last column of the table. 2nd step: Sort the rows in lexicographic order.

> \$ a a b n n

The decoding part

It is easily possible to recover **bananas** from **5annb\$aa** by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb\$aa before the last column of the table. 2nd step: Sort the rows in lexicographic order.

а	\$
n	а
n	а
b	а
\$	b
а	n
Э	n

The decoding part

It is easily possible to recover banana^{\$} from 5annb^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb[§]aa before the last column of the table. 2nd step: Sort the rows in lexicographic order.

 \Rightarrow Repeat until the whole table is recreated.

\$b a\$ an an ba na na

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb\$aa before the last column of the table. 2nd step: Sort the rows in lexicographic order.

а	\$	b
n	а	\$
n	а	n
b	а	n
\$	b	а
а	n	а
а	n	а

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb\$aa before the last column of the table. 2nd step: Sort the rows in lexicographic order.

\$	b	а
а	\$	b
а	n	а
а	n	а
b	а	n
n	а	\$
n	а	n

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annbsaa before the last column of the table.
 2nd step: Sort the rows in lexicographic order.

а	\$	b	а
n	а	\$	b
n	а	n	а
b	а	n	а
\$	b	а	n
а	n	а	\$
а	n	а	n

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb a before the last column of the table.
 2nd step: Sort the rows in lexicographic order.

\$	b	а	n
а	\$	b	а
а	n	а	\$
а	n	а	n
b	а	n	а
n	а	\$	b
n	а	n	а

The decoding part

It is easily possible to recover banana^{\$} from 5annb^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb a before the last column of the table. 2nd step: Sort the rows in lexicographic order.

а	\$	b	а	n
n	а	\$	b	а
n	а	n	а	\$
b	а	n	а	n
\$	b	а	n	а
а	n	а	\$	b
а	n	а	n	а

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb a before the last column of the table. 2nd step: Sort the rows in lexicographic order.

\$	b	а	n	а
а	\$	b	а	n
а	n	а	\$	b
а	n	а	n	а
b	а	n	а	n
n	а	\$	b	а
n	а	n	а	\$

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb a before the last column of the table. 2nd step: Sort the rows in lexicographic order.

а	\$	b	а	n	а
n	а	\$	b	а	n
n	а	n	а	\$	b
b	а	n	а	n	а
\$	b	а	n	а	n
а	n	а	\$	b	а
а	n	а	n	а	\$

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annbsaa before the last column of the table. 2nd step: Sort the rows in lexicographic order.

\$	b	а	n	а	n
а	\$	b	а	n	а
а	n	а	\$	b	а
а	n	а	n	а	\$
b	а	n	а	n	а
n	а	\$	b	а	n
n	а	n	а	\$	b

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb a before the last column of the table. 2nd step: Sort the rows in lexicographic order.

а	\$	b	а	n	а	n
n	а	\$	b	а	n	а
n	а	n	а	\$	b	а
b	а	n	а	n	а	\$
\$	b	а	n	а	n	а
а	n	а	\$	b	а	n
а	n	а	n	а	\$	b

The decoding part

It is easily possible to recover banana^{\$} from 5annb^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb a before the last column of the table. 2nd step: Sort the rows in lexicographic order.

\$	b	а	n	а	n	а
а	\$	b	а	n	а	n
а	n	а	\$	b	а	n
а	n	а	n	а	\$	b
b	а	n	а	n	а	\$
n	а	\$	b	а	n	а
n	а	n	а	\$	b	а

The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb[®]aa before the last column of the table.
2nd step: Sort the rows in lexicographic order.
⇒ Repeat until the whole table is recreated.
BWT⁻¹: Read the corresponding row in the table.



The decoding part

It is easily possible to recover **banana**^{\$} from **5annb**^{\$}aa by reconstructing the sorted table that was constructed during the computation of the BWT.

1st step: Concatenate annb[®]aa before the last column of the table.
2nd step: Sort the rows in lexicographic order.
⇒ Repeat until the whole table is recreated.
BWT⁻¹: Read the corresponding row in the table.



Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.



Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

$$\underline{\mathsf{Ex:}} \ \mathsf{BWT}\Big(\ \underline{\mathsf{banana}} \Big) = \ \underline{\mathsf{5annb}} \\ \underline{\mathsf{5annb}} \\ \underline{\mathsf{aa}} \ \Rightarrow \ \mathsf{let's} \ \mathsf{encode} \ \underline{\mathsf{annb}} \\ \underline{\mathsf{annb}} \\ \underline{\mathsf{aa}} \ \mathcal{L}_0 = [\mathsf{a,b,n,\$}]$$

 $a \rightarrow 0$ first symbol of \mathcal{L}_0

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

$$\begin{array}{l} \underline{\mathsf{Ex:}} \ \mathsf{BWT}\Big(\begin{array}{c} \underline{\mathsf{banans}} \Big) = \begin{array}{c} \overline{\mathsf{5annb}} \\ \overline{\mathsf{saa}} \end{array} \Rightarrow \mathsf{let's encode annb} \\ \underline{\mathsf{aa}} & \mathcal{L}_0 = [\mathsf{a},\mathsf{b},\mathsf{n},\mathsf{s}] \end{array} \\ \\ \begin{array}{c} \underline{\mathsf{a}} \ \rightarrow 0 \\ \mathbf{n} \ \rightarrow 2 \end{array} & \begin{array}{c} \text{first symbol of } \mathcal{L}_0 \\ \mathcal{L}_0 \ \rightarrow \mathcal{L}_1 = [\mathsf{n},\mathsf{a},\mathsf{b},\mathsf{s}] \end{array} \end{array}$$

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

$$\begin{array}{l} \underline{\mathsf{Ex:}} \ \mathsf{BWT}\Big(\begin{array}{c} \underline{\mathsf{banans}} \Big) = \begin{array}{c} \overline{\mathsf{5annb\$aa}} \Rightarrow \mathsf{let's} \ \mathsf{encode} \ \begin{array}{c} \underline{\mathsf{annb\$aa}} \\ \mathcal{L}_0 = [\mathsf{a},\mathsf{b},\mathsf{n},\mathsf{s}] \end{array} \\ \\ \begin{array}{c} \underline{\mathsf{a}} \to 0 \\ \underline{\mathsf{n}} \to 2 \\ \overline{\mathsf{n}} \to 0 \end{array} & \begin{array}{c} \mathsf{first} \ \mathsf{symbol} \ \mathsf{of} \ \mathcal{L}_0 \\ \mathcal{L}_0 \to \mathcal{L}_1 = [\mathsf{n},\mathsf{a},\mathsf{b},\mathsf{s}] \\ \\ \underline{\mathsf{n}} \to 0 \end{array} & \begin{array}{c} \mathsf{first} \ \mathsf{symbol} \ \mathsf{of} \ \mathcal{L}_1 \end{array} \end{array}$$

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

 $\begin{array}{l} \underline{\mathsf{Ex:}} \ \mathsf{BWT}\left(\begin{array}{c} \underline{\mathsf{banana}}\\ \end{array}\right) = \begin{array}{c} \underline{\mathsf{5annb}}\underline{\mathsf{saa}}\\ \Rightarrow \ \mathsf{let's} \ \mathsf{encode} \ \begin{array}{c} \underline{\mathsf{annb}}\underline{\mathsf{saa}}\\ \end{array}\right) \mathcal{L}_0 = [\mathsf{a},\mathsf{b},\mathsf{n},\mathsf{s}] \\ \hline \mathsf{n} \to 0 & \text{first symbol of } \mathcal{L}_0 \\ \hline \mathsf{n} \to 0 & \text{first symbol of } \mathcal{L}_1 \\ \hline \mathsf{b} \to 2 & \text{third symbol of } \mathcal{L}_1 \\ \end{array} \mathcal{L}_1 \to \mathcal{L}_2 = [\mathsf{b},\mathsf{n},\mathsf{a},\mathsf{s}] \end{array}$

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

 $\begin{array}{l} \underline{\mathsf{Ex:}} \ \mathsf{BWT}\left(\begin{array}{c} \underline{\mathsf{banans}}\right) = \begin{array}{c} \underline{\mathsf{5annb\$aa}} \Rightarrow \mathsf{let's} \ \mathsf{encode} \ \begin{array}{c} \underline{\mathsf{annb\$aa}} \\ \mathcal{L}_0 = [\mathsf{a},\mathsf{b},\mathsf{n},\$] \\ \end{array} \\ \begin{array}{c} \underline{\mathsf{a}} \to 0 & \mathsf{first} \ \mathsf{symbol} \ \mathsf{of} \ \mathcal{L}_0 \\ \underline{\mathsf{n}} \to 2 & \mathsf{third} \ \mathsf{symbol} \ \mathsf{of} \ \mathcal{L}_0 \\ \end{array} \\ \begin{array}{c} \mathcal{L}_0 \to \mathcal{L}_1 = [\mathsf{n},\mathsf{a},\mathsf{b},\$] \\ \underline{\mathsf{n}} \to 0 & \mathsf{first} \ \mathsf{symbol} \ \mathsf{of} \ \mathcal{L}_1 \\ \end{array} \\ \begin{array}{c} \underline{\mathsf{b}} \to 2 & \mathsf{third} \ \mathsf{symbol} \ \mathsf{of} \ \mathcal{L}_1 \\ \end{array} \\ \begin{array}{c} \mathcal{L}_1 \to \mathcal{L}_2 = [\mathsf{b},\mathsf{n},\mathsf{a},\$] \\ \end{array} \\ \begin{array}{c} \underline{\mathsf{s}} \to 3 & \mathsf{fourth} \ \mathsf{symbol} \ \mathsf{of} \ \mathcal{L}_2 \\ \end{array} \end{array} \\ \begin{array}{c} \mathcal{L}_2 \to \mathcal{L}_3 = [\$,\mathsf{b},\mathsf{n},\mathsf{a}] \end{array} \end{array}$

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

Ex: BWT(banana^{\$}) = 5annb^{\$}aa \Rightarrow let's encode annb^{\$}aa $\mathcal{L}_0 = [a,b,n,^$]$

first symbol of \mathcal{L}_0 $a \rightarrow 0$

- $n \rightarrow 2$ third symbol of \mathcal{L}_0
 - ightarrow 0 first symbol of \mathcal{L}_1
- ${f b}
 ightarrow 2$ third symbol of ${\cal L}_1$
 - \rightarrow 3 fourth symbol of \mathcal{L}_2
 - \rightarrow 3 fourth symbol of \mathcal{L}_3

$$\mathcal{L}_0 \to \mathcal{L}_1 = [\mathsf{n},\mathsf{a},\mathsf{b},\$]$$

$$\mathcal{L}_1 \to \mathcal{L}_2 = [\mathsf{b},\mathsf{n},\mathsf{a},\textcolor{red}{\$}]$$

$$\mathcal{L}_2 \rightarrow \mathcal{L}_3 = [\$, b, n, a]$$

$$\mathcal{L}_3 \rightarrow \mathcal{L}_4 = [a, \$, b, n]$$

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

Ex: BWT (banana) = 5annbsaa \Rightarrow let's encode annbsaa



$$\mathcal{L}_0 = [\mathsf{a},\mathsf{b},\mathsf{n},\$]$$

- first symbol of \mathcal{L}_0 $a \rightarrow 0$
- $n \rightarrow 2$ third symbol of \mathcal{L}_0
- ightarrow 0 first symbol of \mathcal{L}_1
- ${f b}
 ightarrow 2$ third symbol of ${\cal L}_1$
 - \rightarrow 3 fourth symbol of \mathcal{L}_2
 - \rightarrow 3 fourth symbol of \mathcal{L}_3
 - first symbol of \mathcal{L}_4 $\rightarrow 0$

$$\mathcal{L}_0 \to \mathcal{L}_1 = [\mathsf{n},\mathsf{a},\mathsf{b},\$]$$

$$\mathcal{L}_1 \to \mathcal{L}_2 = [\mathsf{b},\mathsf{n},\mathsf{a},\textcolor{red}{\$}]$$

$$\mathcal{L}_2 \to \mathcal{L}_3 = [\$, b, n, a]$$

$$\mathcal{L}_3 \rightarrow \mathcal{L}_4 = [a, \mathbf{b}, \mathbf{b}, \mathbf{n}]$$

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

Ex: BWT(banana[§]) = 5annb[§]aa \Rightarrow let's encode annb[§]aa $\mathcal{L}_0 = [a,b,n,s]$ first symbol of \mathcal{L}_0 $a \rightarrow 0$ $n \rightarrow 2$ third symbol of $\mathcal{L}_0 \longrightarrow \mathcal{L}_1 = [n,a,b,\$]$ ${\sf n}
ightarrow 0 ~~$ first symbol of ${\cal L}_1$ **b** \rightarrow 2 third symbol of \mathcal{L}_1 $\mathcal{L}_1 \rightarrow \mathcal{L}_2 = [b,n,a,\$]$ $\begin{array}{ccc} & \bullet & \mathsf{3} & \text{ fourth symbol of } \mathcal{L}_2 & \mathcal{L}_2 \to \mathcal{L}_3 = [\$, \mathsf{b}, \mathsf{n}, \mathsf{a}] \\ & \mathsf{a} & \to \mathsf{3} & \text{ fourth symbol of } \mathcal{L}_3 & \mathcal{L}_3 \to \mathcal{L}_4 = [\mathsf{a}, \$, \mathsf{b}, \mathsf{n}] \end{array}$ $a \rightarrow 0$ first symbol of \mathcal{L}_4 $\mathsf{MTF}(\mathsf{annb}\) \Rightarrow 0202330$

0 becomes much more probable than other digits \rightarrow Huffman likes that!

Idea: encode a flow of symbols by their index in a list. The list is dynamically updated so that most recently seen symbol is moved to the front of the list.

Ex: BWT(banana[§]) = 5annb[§]aa \Rightarrow let's encode annb[§]aa $\mathcal{L}_0 = [a,b,n,s]$ first symbol of \mathcal{L}_0 $a \rightarrow 0$ $n \rightarrow 2$ third symbol of $\mathcal{L}_0 \qquad \mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n,a,b,\$]$ ${\sf n}
ightarrow 0$ first symbol of ${\cal L}_1$ **b** \rightarrow 2 third symbol of \mathcal{L}_1 $\mathcal{L}_1 \rightarrow \mathcal{L}_2 = [b,n,a,\$]$ $\begin{array}{ccc} \$ & \rightarrow 3 & \text{fourth symbol of } \mathcal{L}_2 & \mathcal{L}_2 \rightarrow \mathcal{L}_3 = [\$, \mathsf{b}, \mathsf{n}, \mathsf{a}] \\ \texttt{a} & \rightarrow 3 & \text{fourth symbol of } \mathcal{L}_3 & \mathcal{L}_3 \rightarrow \mathcal{L}_4 = [\mathsf{a}, \$, \mathsf{b}, \mathsf{n}] \end{array}$ $a \rightarrow 0$ first symbol of \mathcal{L}_4 $\mathsf{MTF}(\mathsf{annb}\) \Rightarrow 0202330$

0 becomes much more probable than other digits \rightarrow Huffman likes that!

 $\mathcal{L} \equiv$ the alphabet Σ from which are drawn the symbols (in general the ASCII table).

Guillaume TOCHON (LRDE)

 1^{st} step: ${\sf BWT} \to {\sf sort}$ data to create runs of similar symbols.

 2^{nd} step: MTF \rightarrow encode sorted data into digits with high anisotropy.

3rd step: Huffman \rightarrow 0 gets a short encoding \Rightarrow huge compression.

✓ Most of the time more efficient than other compression algorithms (.zip and .gz, based on DEFLATE algorithm).

-rw-r--r-- 1 gtochon lrde 20K mars 23 17:54 random_english.txt -rw-r--r-- 1 gtochon lrde 6,6K mars 23 17:54 random_english.txt.bz2 -rw-r--r-- 1 gtochon lrde 7,9K mars 23 17:54 random_english.txt.gz -rw-r--r-- 1 gtochon lrde 8,0K mars 23 18:38 random_english.txt.zip

X Notably slower for the compression stage (but not for decompression).

LZW compression algorithm

- $\rightarrow\,$ Improvement of the LZ78 algorithm proposed by Abraham Lempel and Jacob Ziv in 1978.
- $\rightarrow\,$ Published by Terry Welch in 1984.
- $\rightarrow\,$ Can be considered with LZ78 as the first unsupervised dictionary learning algorithms.





Abraham Lempel

Jacob Ziv

LZW encoding

Idea: Create a dictionary \mathcal{D} based on the sequences of symbols encountered in the data, and replace known sequences of symbols by their address in \mathcal{D} .

 $F = \{s_i\}_{i=1}^{N_F} \xrightarrow{s_i} \overbrace{s_{i-1}s_i \in \mathcal{D}?} \xrightarrow{\mathcal{D} = \mathcal{D} \cup \{s_{i-1}s_i\}} \underbrace{\mathcal{D} = \mathcal{D} \cup \{s_{i-2}s_{i-1}s_i\}}_{\mathcal{V}_{\mathcal{U}_{\mathcal{D}}}} \xrightarrow{\mathcal{D} = \mathcal{D} \cup \{s_{i-2}s_{i-1}s_i\}} \underbrace{\mathcal{D} = \mathcal{D} \cup \{s_{i-2}s_{i-1}s_i\}}_{s_{i-1}s_i \to \mathfrak{O}\mathcal{D}[s_{i-1}s_i]}$

Example

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output

Example

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			

What is happening?

- 1. T is the first read character.
- 2. Buffer is empty
- 3. T is already in \mathcal{D} .

Example

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output	W
	Т				
Т	0	то	256	©⊅ [T]	

What is happening?

- 1. T is put in the buffer.
- 2. O is read.
- 3. The sequence **TO** is created and tested to belong to \mathcal{D} .
- 4. TO is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 256.
- 5. **T** is encoded by its address in \mathcal{D} .
LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	@₽ [T]
0	-	0_	257	@D [O]

- 1. O is put in the buffer.
- 2. _ is read.
- 3. The sequence O_{-} is created and tested to belong to \mathcal{D} .
- 4. O. is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 257.
- 5. O is encoded by its address in \mathcal{D} .

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	@₽ [T]
0	-	0_	257	@D [0]
-	В	_B	258	@D [_]

- 1. _ is put in the buffer.
- 2. B is read.
- 3. The sequence $_B$ is created and tested to belong to \mathcal{D} .
- 4. **B** is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 258.
- 5. \Box is encoded by its address in \mathcal{D} .

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	$@\mathcal{D}[T]$
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
	•	•		

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	@₽ [T]
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	@₽ [T]
0	-	0_	257	@D [0]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	_	E_	260	@₽ [E]
-	0	_0	261	©⊅[_]

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	@₽ [T]
0	-	0_	257	@D [0]
-	В	_B	258	@D [_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	©⊅[_]
0	R	OR	262	@D [O]
1				

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

Note that some sequences should come back pretty often, such as TO, BE, OR...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $D \equiv$ ASCII table.

Duffer	Laura	D U()	0D[]	O tt
Buffer	Input	$\mathcal{D}\cup\{\cdot\}$	$\mathbb{Q}D[\cdot]$	Output
	Т			
Т	0	то	256	©⊅ [T]
0	-	0_	257	@D [0]
-	В	_B	258	©⊅[_]
В	E	BE	259	©⊅[B]
E	-	E_	260	@₽ [E]
-	0	_0	261	©⊅[_]
0	R	OR	262	@D [0]
R	-	R_	263	@D [R]

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

Note that some sequences should come back pretty often, such as TO, BE, OR...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	©⊅ [T]
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	@D [_]

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

Note that some sequences should come back pretty often, such as TO, BE, OR...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	©⊅ [T]
0	-	0_	257	@D [0]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	_	E.	260	@₽ [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [0]
R	_	R_	263	$@\mathcal{D}[R]$
-	N	_N	264	@D [_]
N	0	NO	265	@D [N]

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

Note that some sequences should come back pretty often, such as TO, BE, OR, NO...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	@D[·]	Output
	Т			
Т	0	то	256	$@\mathcal{D}[T]$
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	©⊅[_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

Note that some sequences should come back pretty often, such as TO, BE, OR, NO...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $D \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	©⊅ [T]
0	-	0_	257	@D [O]
-	В	_B	258	©₽ [_]
В	E	BE	259	@D [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [0]
R	-	R_	263	@D [R]
-	N	_N	264	@D [_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]
Т	-	T_	267	@D [T]

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

Note that some sequences should come back pretty often, such as TO, BE, OR, NO...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $D \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	@D [T]
0	-	0_	257	@D [0]
_	В	_B	258	@D [_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
_	0	_0	261	@D [_]
0	R	OR	262	@D [0]
R	_	R_	263	@D [R]
_	N	_N	264	@D [_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [0]
Т	-	T.	267	©⊅[T]
-	Т	_T	268	@D [_]

What is happening?

It keeps adding in $\ensuremath{\mathcal{D}}$ new sequences of two symbols that were never encountered before...

Note that some sequences should come back pretty often, such as TO, BE, OR, NO...

And surely enough...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	©⊅ [T]
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	_	E.	260	@₽ [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [O]
R	_	R_	263	@D [R]
-	N	_N	264	@D [_]
N	0	NO	265	@D [N]
0	Т	ОТ	266	@D [0]
Т	-	Τ.	267	© <i>D</i> [T]
_	Т	_T	268	@D [_]
Т	0	то		

What is happening?

- 1. T is put in the buffer.
- 2. O is read.

3. The sequence **TO** is created and tested to belong to \mathcal{D} .

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $D \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	©⊅ [T]
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	©⊅[B]
E	-	E.	260	@₽ [E]
-	0	_0	261	©⊅ [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	©⊅ [_]
N	0	NO	265	@D [N]
0	Т	ОТ	266	@D [O]
Т	-	T_	267	©⊅ [T]
-	Т	T_	268	©₽ [_]
Т	0	то		

- 1. T is put in the buffer.
- 2. O is read.
- 3. The sequence **TO** is created and tested to belong to \mathcal{D} .
- 4. TO is already in \mathcal{D} , at the address $256 \rightarrow$ it is not inserted in \mathcal{D} again.

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	$@\mathcal{D}[T]$
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@D [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	@D [_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]
Т	-	T_	267	©₽ [T]
-	Т	_T	268	@D [_]
Т	0			×
то				
		• / -		

- . T is put in the buffer.
- 2. O is read.
- 3. The sequence **TO** is created and tested to belong to \mathcal{D} .
- 4. **TO** is already in \mathcal{D} , at the address $256 \rightarrow$ it is not inserted in \mathcal{D} again.
- 5. **T** is not encoded by its address (nothing is output), but **TO** is put in the buffer instead.

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	$@\mathcal{D}[T]$
0	-	0_	257	@D [O]
-	В	_B	258	@D [_]
В	E	BE	259	@₽ [B]
Е	-	E_	260	@₽ [E]
-	0	_0	261	©⊅ [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	@D [_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]
Т	-	T_	267	©₽ [T]
-	Т	_T	268	©⊅[_]
Т	0			
то	-	TO_	269	©⊅ [TO]
		• ·		

- 1. TO is in the buffer.
- 2. _ is read.
- 3. The sequence **TO** is created and tested to belong to \mathcal{D} .
- 4. **TO.** is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 269.
- 5. **TO** is encoded by its address in \mathcal{D} .

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	@D[·]	Output
	Т			
Т	0	то	256	©₽ [T]
0	-	0_	257	@D [0]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	@D [_]
N	0	NO	265	@D [N]
0	Т	ОТ	266	@D [0]
Т	-	Τ.	267	©⊅ [T]
-	Т	_T	268	@D [_]
Т	0			
то	-	TO_	269	256
	•	•		

- 1. TO is in the buffer.
- 2. _ is read.
- 3. The sequence **TO** is created and tested to belong to \mathcal{D} .
- 4. **TO.** is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 269.
- 5. TO is encoded by its address in D, *i.e.*, 256.

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	@₽ [T]
0	-	0_	257	@D [O]
-	В	_B	258	@D [_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	@D [_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]
Т	-	T_	267	©₽ [T]
-	Т	_T	268	@D [_]
Т	0			@D [<mark>%</mark>]
то	-	TO_	269	256

What is happening?

- 1. TO is in the buffer.
- 2. _ is read.
- 3. The sequence **TO**₋ is created and tested to belong to \mathcal{D} .
- TO_ is not in D yet → it is inserted at the next available address, <u>i.e.</u> 269.
- 5. TO is encoded by its address in D, *i.e.*, 256.
- But wait, 256 requires 9 bits to be encoded (instead of 8 bits so far). → needs to emit a special character

to prevent the decompressor that it shall read further addresses on 9 bits and no longer on 8 bits.

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	$@\mathcal{D}[T]$
0	-	0_	257	@D [O]
-	В	_B	258	@D [_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	@D [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	©⊅[_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]
Т	-	T_	267	$@\mathcal{D}[T]$
-	Т	_T	268	@D [_]
Т	0			@D [<mark>%</mark>]
то	-	TO_	269	256
-	В			

- 1. _ is put in the buffer.
- 2. B is read.
- 3. The sequence $_B$ is created and tested to belong to \mathcal{D} .
- 4. **.** is already in \mathcal{D} , at the address $258 \rightarrow$ it is not inserted in \mathcal{D} again.
- 5. Nothing is output, and _B is put in the buffer.

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	$@\mathcal{D}[T]$
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	©⊅ [_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	©⊅[_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]
Т	-	T_	267	$@\mathcal{D}[T]$
-	Т	_T	268	©⊅[_]
Т	0			@D [<mark>%</mark>]
то	-	TO_	269	256
-	В			
_B	E	_BE	270	258
		•		

- 1. _B is in the buffer.
- 2. E is read.
- 3. The sequence $_BE$ is created and tested to belong to \mathcal{D} .
- 4. **LBE** is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 270.
- 5. **B** is encoded by its address in \mathcal{D} , *i.e.*, 258.

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	@₽ [T]
0	-	0_	257	@D [O]
-	В	_B	258	©⊅[_]
В	E	BE	259	@₽ [B]
E	-	E_	260	$@\mathcal{D}[E]$
-	0	_0	261	©⊅[_]
0	R	OR	262	@D [0]
R	-	R_	263	$@\mathcal{D}[R]$
-	Ν	_N	264	©⊅[_]
Ν	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]
Т	-	T_	267	@₽ [T]
-	Т	_T	268	©⊅[_]
Т	0			@D [<mark>%</mark>]
то	-	TO_	269	256
-	В			
_B	E	₋BE	270	258
E	\$			©⊅[E]

- 1. E is put in the buffer.
- 2. The EOF character \$ is read.
- The dictionary update stops, and
 is encoded by its address in D.

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $D \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	©₽ [T]
0	-	0_	257	@D [O]
-	В	_B	258	@D [_]
В	E	BE	259	@₽ [B]
E	-	E_	260	@₽ [E]
-	0	_0	261	©⊅[_]
0	R	OR	262	@D [O]
R	-	R_	263	@D [R]
-	N	_N	264	@D [_]
N	0	NO	265	$@\mathcal{D}[N]$
0	Т	ОТ	266	@D [O]
Т	-	T_	267	©₽ [T]
-	Т	_T	268	@D [_]
Т	0			@D [<mark>%</mark>]
то	-	TO_	269	256
-	В			
_B	E	_BE	270	258
E	\$			@D [E]

Compression performance:

- Without compression \rightarrow 18 \times 8 = 144 bits.
- With compression \rightarrow 14×8+3×9 = 139 bits.

OK, it is not so impressive for this example...

LZW encoding

Let's encode F = TO BE OR NOT TO BE with $\mathcal{D} \equiv$ ASCII table.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	Т			
Т	0	то	256	©⊅ [T]
0	-	0_	257	@D [O]
-	В	_B	258	@D [_]
В	E	BE	259	@D [B]
E	-	E_	260	@D [E]
-	0	_0	261	©⊅ [_]
0	R	OR	262	@D [O]
R	-	R_	263	$@\mathcal{D}[R]$
-	N	_N	264	©⊅ [_]
N	0	NO	265	@D [N]
0	Т	ОТ	266	@D [O]
Т	-	T_	267	@D [T]
-	Т	_T	268	©⊅ [_]
Т	0			@D [<mark>%</mark>]
то	-	TO_	269	256
-	В			
_B	E	_BE	270	258
E	\$			©⊅ [E]

Compression performance:

But on longer strings, when the dictionary \mathcal{D} has sufficiently grown, the gain becomes really significant.

	How much wood would a wood-
Ex:	chuck chuck if a woodchuck could
	chuck wood?

- Without compression \rightarrow 70 \times 8 = 560 bits.

- With compression $\rightarrow 14\times 8+9\times 31=391$ bits $_{(check it yourself).}$

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	©₽[T]			

- 1. $\mathcal{OD}[\mathsf{T}]$ is received by the decoder.
- 2. T is recognized by looking up in \mathcal{D} .

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output	W
	@D [T]				
Т	@D [O]	то	256	Т	

- 1. \mathbf{T} is put in the buffer.
- 2. **O** is recognized from @D[O].
- 3. The sequence **TO** is created and tested to belong to \mathcal{D} .
- 4. TO is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 256.
- 5. T is output.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	@D [T]			
Т	@D [O]	ТО	256	Т
0	@D [_]	0_	257	0

- 1. O is put in the buffer.
- 2. _ is recognized from $@D[_].$
- 3. The sequence O_{-} is created and tested to belong to \mathcal{D} .
- 4. O_ is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 257.
- 5. O is output.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	@D [T]			
Т	@D [0]	то	256	Т
0	©⊅[_]	0_	257	0
-	©⊅[B]	_B	258	-
В	©⊅[E]	BE	259	В
E	@D [_]	E_	260	E
-	@D [0]	_0	261	-
0	@D [R]	OR	262	0
R	@D [_]	R_	263	R
-	$@\mathcal{D}[N]$	_N	264	-
N	@D [0]	NO	265	N
0	©⊅[T]	ОТ	266	0
Т	@D [_]	T_	267	Т

What is happening?

And so on, as long as no special character is encountered.

The sequences that are recreated by this process are strictly identical to those that were generated during the encoding stage.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	@D [T]			
Т	@D [0]	то	256	Т
0	@D [_]	0_	257	0
-	©⊅[B]	_B	258	-
В	@₽ [E]	BE	259	В
Е	@D [_]	E_	260	Е
-	@D [0]	_0	261	-
0	@₽ [R]	OR	262	0
R	@D [_]	R_	263	R
-	@D [N]	_N	264	-
Ν	@D [0]	NO	265	N
0	©⊅[T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			

- 1. _ is put in the buffer.
- The special character % is recognized from @D [%] → all following addresses will have to be decoded on 9 bits instead of 8.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	©₽[T]			
Т	@D [0]	то	256	Т
0	©⊅ [_]	0_	257	0
-	©⊅[B]	_B	258	-
В	@₽ [E]	BE	259	В
E	@D [_]	E_	260	E
-	@D [0]	_0	261	-
0	@D [R]	OR	262	0
R	@D [_]	R_	263	R
-	$@\mathcal{D}[N]$	_N	264	-
Ν	@D [0]	NO	265	N
0	©⊅[T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			
-				



1. _ is in the buffer.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	©₽[T]			
Т	@D [0]	то	256	Т
0	©⊅[_]	0_	257	0
-	@₽ [B]	_B	258	_
В	@₽ [E]	BE	259	В
Е	@D [_]	E_	260	E
-	@D [0]	_0	261	-
0	@D [R]	OR	262	0
R	@D [_]	R_	263	R
-	@D [N]	_N	264	-
Ν	@D [0]	NO	265	N
0	©⊅[T]	ОТ	266	0
Т	@D [_]	T.	267	Т
_	@D [<mark>%</mark>]			
-	256			



- 1. _ is in the buffer.
- 2. 256 is received.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	©⊅[T]			
Т	@D [0]	то	256	Т
0	@D [_]	0_	257	0
-	©⊅[B]	_B	258	-
В	@₽ [E]	BE	259	В
E	@D [_]	E_	260	Е
-	@D [0]	_0	261	-
0	@D [R]	OR	262	0
R	@D [_]	R_	263	R
-	©⊅ [N]	_N	264	-
Ν	@D [0]	NO	265	Ν
0	©⊅[T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	256			

What is happening?

1. _ is in the buffer.

2. $\overline{256}$ is received, that is, $@\mathcal{D}[\mathsf{TO}]$.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	@D [T]			
Т	@D [O]	то	256	Т
0	@D [_]	0_	257	0
-	©⊅[B]	_B	258	-
В	@₽ [E]	BE	259	В
Е	@D [_]	E_	260	E
-	@D [O]	_0	261	-
0	@D [R]	OR	262	0
R	@D [_]	R_	263	R
-	$@\mathcal{D}[N]$	_N	264	-
Ν	@D [O]	NO	265	N
0	©₽ [T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	©⊅[TO]	T_		
		•		

What is happening?
 is in the buffer.
 256 is received, that is, @D [TO].
 The sequence T is created using

 in the buffer and the first letter of the received sequence TO.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	@D [T]			
Т	@D [0]	то	256	Т
0	@D [_]	0_	257	0
-	©⊅[B]	_B	258	-
В	@₽ [E]	BE	259	В
E	@D [_]	E_	260	E
-	@D [0]	_0	261	-
0	@D [R]	OR	262	0
R	@D [_]	R_	263	R
-	@D [N]	_N	264	-
N	@D [0]	NO	265	N
0	©⊅[T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	@ <i>D</i> [TO]	T_	268	
		•	•	

What is happening?

- _ is in the buffer. 1.
- 2. $\overline{256}$ is received, that is, @D[TO].
- 3. The sequence _T is created using _ in the buffer and the first letter

of the received sequence TO

- \square is not in \mathcal{D} yet \rightarrow it is inserted 4. at the next available address. *i.e.* 268

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	@D [T]			
Т	@D [0]	то	256	Т
0	©⊅[_]	0_	257	0
-	©⊅[B]	_B	258	-
В	@₽ [E]	BE	259	В
E	@D [_]	E_	260	E
-	@D [0]	_0	261	-
0	@⊅ [R]	OR	262	0
R	@D [_]	R_	263	R
-	@D [N]	_N	264	-
N	@D [0]	NO	265	N
0	©⊅[T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	@D [TO]	_T	268	-
			•	

What is happening?
1. is in the buffer.
2. 256 is received, that is, @D [TO].
3. The sequence T is created using

in the buffer and the first letter of the received sequence TO.

4. T is not in D yet → it is inserted at the next available address, *i.e.* 268.
5. is output.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	©₽[T]			
Т	@D [0]	то	256	Т
0	©⊅[_]	0_	257	0
-	©⊅[B]	_B	258	-
В	©⊅[E]	BE	259	В
E	©⊅ [_]	E_	260	E
-	@D [O]	_0	261	-
0	@D [R]	OR	262	0
R	@D [_]	R_	263	R
-	$@\mathcal{D}[N]$	_N	264	-
Ν	@D [0]	NO	265	Ν
0	©⊅[T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	@ <i>D</i> [TO]	_T	268	-
то	258			
	-			

What is happening?

1. **TO** is put in the buffer.

2. 258 is received.
LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[T]$			
Т	@₽ [0]	то	256	Т
0	©⊅[_]	0_	257	0
-	@₽ [B]	_B	258	-
В	©⊅[E]	BE	259	В
E	©⊅[_]	E_	260	Е
-	@₽ [0]	_0	261	-
0	$@\mathcal{D}[R]$	OR	262	0
R	@D [_]	R_	263	R
-	$@\mathcal{D}[N]$	_N	264	-
Ν	@₽ [0]	NO	265	Ν
0	$@\mathcal{D}[T]$	ОТ	266	0
Т	©⊅[_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	©⊅[T0]	_T	268	-
то	@D [₋B]			

What is happening?

1. TO is put in the buffer.

2. 258 is received, that is, _B .

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	@D [T]			
Т	@D [0]	то	256	Т
0	@D [_]	0_	257	0
-	@D [B]	_B	258	-
В	@₽ [E]	BE	259	В
E	@D [_]	E_	260	E
-	@D [O]	_0	261	-
0	@D [R]	OR	262	0
R	@D [_]	R_	263	R
-	@D [N]	_N	264	-
N	@D [0]	NO	265	N
0	@₽ [T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	@ <i>D</i> [TO]	T_	268	-
то	@D [_B]	TO_	269	то
			•	

What is happening?
1. TO is put in the buffer.
2. 258 is received, that is, _B.
3. The sequence TO. is created using TO in the buffer and the first letter . of the received sequence _B.
4. TO. is not in D yet → it is inserted at the next available address, *i.e.* 269.
5. TO is output.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[T]$			
Т	@₽ [0]	то	256	Т
0	©⊅[_]	0_	257	0
-	@₽ [B]	_B	258	-
В	©⊅[E]	BE	259	В
E	@D [_]	E_	260	Е
-	@₽ [0]	_0	261	-
0	$@\mathcal{D}[R]$	OR	262	0
R	@D [_]	R_	263	R
-	$@\mathcal{D}[N]$	_N	264	-
Ν	@₽ [0]	NO	265	Ν
0	$@\mathcal{D}[T]$	ОТ	266	0
Т	©⊅[_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	©⊅[T0]	_T	268	-
то	@D [₋B]		269	то
_B	@D [E]	_BE	270	_B

What is happening?

- 1. _B is put in the buffer.
- 2. **E** is recognized from @D[E].
- 3. The sequence $_BE$ is created and tested to belong to \mathcal{D} .
- 4. **LBE** is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 270.
- 5. _B is output.

LZW decoding

The LZW decoding scheme reconstructs the dictionary ${\cal D}$ from the ASCII table in a very symmetric fashion.

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	@D [T]			
Т	@D [O]	то	256	Т
0	@D [_]	0_	257	0
-	©⊅[B]	_B	258	-
В	@₽ [E]	BE	259	В
E	@D [_]	E_	260	Е
-	@D [O]	_0	261	-
0	@₽ [R]	OR	262	0
R	@D [_]	R_	263	R
-	$@\mathcal{D}[N]$	_N	264	-
Ν	@D [O]	NO	265	N
0	©₽ [T]	ОТ	266	0
Т	@D [_]	T_	267	Т
-	@D [<mark>%</mark>]			
-	@ <i>D</i> [TO]	_T	268	-
то	@D [₋B]		269	то
_B	@D [E]	_BE	270	_B
Е	\$			Е

What is happening?

- 1. E is put in the buffer.
- 2. The EOF character **\$** is read.
- 3. The dictionary update stops, **E** is recognized from $\mathcal{OD}[E]$ and is output.

The dictionaries generated by the encoding and the decoding stages are strictly identical.

The dictionary \mathcal{D} cannot grow indefinitely.

- \rightarrow Most of the time limited to 12 bits, *i.e.*, 4096 entries.
- $\rightarrow\,$ Can be reset with another special character if needed.

LZW is used in the GIF (Graphics Interchange Format) encoding format for images.

- $\rightarrow\,$ Pixel values are between 0 and 255, hence all 8 bits combinaisons are required.
- $\rightarrow\,$ Directly mapped to 9 bits encoding (also to accomodate for special characters).
- ightarrow Dictionary size is precisely limited to 12 bits.