

Lossless compression algorithms

Guillaume TOCHON

guillaume.tochon@lrde.epita.fr

LRDE, EPITA



Lossless compression

Data before compression and after decompression are strictly identical.



Lossless compression

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 :

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

Lossless compression

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 :

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

2) Use statistical model to encode input data :

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

Lossless compression

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 :

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

2) Use statistical model to encode input data :

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

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

Lossless compression algorithms

- 1 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 ?

$$F = \{ \text{AAAAAABBBBCCCCCCC} \}$$

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

$$F = \{ \text{AAAAAABBBBCCCCCCC} \}$$
$$\Rightarrow F^\# = \{ 6A4B7C \}$$

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.

$$\dots s_j \underbrace{s_j s_j \dots s_j}_{N \text{ times}} s_k \dots \xRightarrow{\text{RLE}} \dots s_j N s_j s_k \dots$$

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.

$$\dots s_j \underbrace{s_j s_j \dots s_j}_{N \text{ times}} s_k \dots \xrightarrow{\text{RLE}} \dots s_j N s_j s_k \dots$$

Not suited for text compression: **HELLO** becomes **1H1E2L1O** or **HE2LO**.

→ The runs must be sufficiently long.

→ 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.

$$\dots s_j \underbrace{s_j s_j \dots s_j}_{N \text{ times}} s_k \dots \xrightarrow{\text{RLE}} \dots s_j N s_j s_k \dots$$

Not suited for text compression: **HELLO** becomes **1H1E2L1O** or **HE2LO**.

→ The runs must be sufficiently long.

→ Suitable for synthesis/noise-free images. Used in bitmap (.bmp) format.

Problem when encoding numbers: **22231444457** $\xrightarrow{\text{RLE}}$ **32314457**.

→ 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.

$$\dots s_j \underbrace{s_j s_j \dots s_j}_{N \text{ times}} s_k \dots \xrightarrow{\text{RLE}} \dots s_j \# N s_j s_k \dots$$

Not suited for text compression: **HELLO** becomes **1H1E2L1O** or **HE2LO**.

→ The runs must be sufficiently long.

→ Suitable for synthesis/noise-free images. Used in bitmap (.bmp) format.

Problem when encoding numbers: **22231444457** $\xrightarrow{\text{RLE}}$ **#3231#4457**.

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

→ Use a special character **#** to inform the decompression stage.

Application to fax transmission of binarized documents

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



Application to fax transmission of binarized documents

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.

DECLARATION DE PHOBIE ADMINISTRATIVE
dans le cadre d'un impôt

☐ Impôt sur le revenu ☐ TMI ☐ RIR ☐ URSSAF
☐ Taxe foncière ☐ TVA ☐ Autre précompte

Identité du déclarant
Nom, Prénom
Coordonnées du déclarant
Adresse
Profession du déclarant
Contrôle des raisons de la phobie

Préciser la date de la dernière échéance fiscale
CARRÉ RÉSERVÉ À L'ADMINISTRATION
Signature du déclarant

DECLARATION DE PHOBIE ADMINISTRATIVE
dans le cadre d'un impôt

☐ Impôt sur le revenu ☐ TMI ☐ RIR ☐ URSSAF
☐ Taxe foncière ☐ TVA ☐ Autre précompte

Identité du déclarant
Nom, Prénom
Coordonnées du déclarant
Adresse
Profession du déclarant
Contrôle des raisons de la phobie



Préciser la date de la dernière échéance fiscale
CARRÉ RÉSERVÉ À L'ADMINISTRATION
Signature du déclarant


Application to fax transmission of binarized documents

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.


 REPUBLIQUE FRANÇAISE	<h2 style="margin: 0;">DECLARATION FISCAL ADMINISTRATIVE</h2>	 REPUBLIQUE FRANÇAISE
<p>Donnez le cadre de cet impôt : <i>cochez les cases correspondantes à votre dernière situation (sans précision)</i></p> <p> <input type="checkbox"/> sans lien avec le revenu <input type="checkbox"/> foyer <input type="checkbox"/> PMA <input type="checkbox"/> isolé(a) <input type="checkbox"/> sans domicile <input type="checkbox"/> PVA <input type="checkbox"/> Autre (préciser) : _____ </p> <p style="text-align: right; font-size: small;">code de l'art. 161 A du code de l'impôt sur le revenu</p>		
<p>Identité du déclarant</p> <p>Nom : _____ Prénom : _____</p>		
<p>Coordonnées du déclarant</p> <p>Adresse : _____ CP : _____ Ville : _____</p>		
<p>Profession du déclarant</p> <p><input type="checkbox"/> Activité professionnelle <input type="checkbox"/> Autre (préciser quelle activité, sur quel territoire) : _____</p>		
<p>Conciter les raisons de la pénalité</p>		
<p><input type="checkbox"/> J'ai pu le donner d'avance à cet impôt</p> <p><input type="checkbox"/> J'ai pu le déclarer le 15 mai au moment de la déclaration</p> <p><input type="checkbox"/> J'ai pu le payer sans pénalité au service de cet impôt</p> <p><input type="checkbox"/> J'ai pu le verser tout de suite <input type="checkbox"/> Autre (PRÉCISER) : _____</p>		
<p>Préciser la date de la dernière déclaration l'année</p> <p> <input type="checkbox"/> japonaise <input type="checkbox"/> impériale à 3 ans <input type="checkbox"/> inférieure à 3 ans </p>		
<p>CADRE RÉSUMÉ À L'ENCAISSEMENT</p>		
<p>DECLARATION TRANSMISE</p> <p style="text-align: center; border: 1px solid black; padding: 2px;"> <i>(cochez)</i> <input type="checkbox"/> OUI <input type="checkbox"/> NON </p>	<p>COULEUR DE L'ENCAISSEMENT</p> <p style="text-align: center; border: 1px solid black; padding: 2px;"> <i>(cochez)</i> <input type="checkbox"/> ROUGE <input type="checkbox"/> VERT </p>	<p>NUMÉRO D'ENCAISSEMENT</p> <p style="text-align: center; border: 1px solid black; padding: 2px;"> <i>(cochez)</i> <input type="checkbox"/> OUI <input type="checkbox"/> NON </p>
<p>date</p>	<p>Signature du déclarant (préciser la date de la dernière déclaration et de l'approbation)</p>	<p>Signature du receveur (préciser la date de la réception et de l'approbation)</p>
<p>Document à déposer au bureau de l'impôt sur le revenu ou au Centre de gestion. Ne pas conserver.</p>		



REPUBLIQUE FRANÇAISE

DECLARATION D'IMPOT

FORMULAIRE ADMINISTRATIF



REPUBLIQUE FRANÇAISE

dans le cadre d'un Impôt

(Cocher la case correspondante à votre situation personnelle)

☐ Impôt sur le revenu

☐ Impôt

☐ RSI

☐ IRRS24

☐ Taxe foncière

☐ TIS

☐ Autre impôt : _____

Date de fin de validité de la déclaration : _____

Identité du déclarant

Nom (M/M) _____ Prénom _____

Coordonnées du déclarant

Adresse _____ CP _____ Ville _____

Profession du déclarant

☐ Activité professionnelle _____

☐ Autre (préciser votre métier, ou sans métier) : _____

Donner les raisons de la pénalité

☐ J'ai pué de rétroceder des impôts et des déductions

☐ J'ai pué de rétroceder les cotisations et déductions de la déduction

☐ J'ai pué de rétroceder mes versements de cotisations

☐ J'ai pué de rétroceder _____

☐ Autre (préciser) _____

Pénalité sur le revenu de la déduction d'impôt

☐ Annulée

☐ Supplémentaire à 2 ans

☐ Supplémentaire à 3 ans

CHACUN SEULEMENT A VALIDER/ANNULER

DROITS THÉORIQUES

CODE DE CONTRÔLE

NUMÉRO D'IDENTIFICATION

Signature du contribuable (préciser de la manière ou sans signature)

Signature du mandataire habilité (préciser de la manière ou sans signature)

Signature du contrôleur (préciser de la manière ou sans signature)

La déclaration est soumise à la vérification de l'administration fiscale.

Rough idea: compress each line independently

Application to fax transmission of binarized documents

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 independently

1	1	1	1	1
---	---	---	---	---

1	1	1	1	1
---	---	---	---	---

 \Rightarrow no text, RLE is super effective.

Application to fax transmission of binarized documents

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 independently

1	1	1	1	1
---	---	---	---	---

1	1	1	1	1
---	---	---	---	---

 \Rightarrow no text, RLE is super effective.

1	1	...	1	1	0	0	1	1	1	1	0	0	1	1	0	1	1	...	1	1
---	---	-----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	-----	---	---

 \Rightarrow text, but RLE remains effective.

Application to fax transmission of binarized documents

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

1	1	...	1	1	0	0	1	1	1	1	0	0	1	1	0	1	1	...	1	1
1	1	...	1	1	0	0	1	1	1	0	0	1	1	1	0	0	1	...	1	1

\Rightarrow Very few differences between two lines.

Application to fax transmission of binarized documents

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

$$\begin{array}{c}
 l_i \\
 l_{i+1} \\
 l_i \oplus l_{i+1}
 \end{array}
 \begin{array}{cccccccccccccccccccc}
 1 & 1 & \cdots & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 1 & 1 & \cdots & 1 & 1 \\
 1 & 1 & \cdots & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 1 & \cdots & 1 & 1 \\
 0 & 0 & \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & \cdots & 0 & 0
 \end{array}$$

⇒ Compress the line difference instead, RLE is super effective.

Huffman compression algorithm

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



David Huffman

Huffman compression algorithm

- Proposed in 1952 by David Huffman (during its Ph.D).
- Exploit the statistical distribution of the symbols to encode.
 - ⇒ *entropy encoding*.
- Frequent symbols are given shorter encoding support.
 - ⇒ *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 occurrence.
- 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.

Example (1/3)

Computing a binary tree

Take the following alphabet Σ with known probability distribution:

s_i	P_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04







Huffman encoding: 1st step







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

Example (1/3)

Computing a binary tree

Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
	0.1
	0.1
	0.06
	0.04







Huffman encoding: 1st step







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

Example (1/3)

Computing a binary tree


Take the following alphabet Σ with known probability distribution:


s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
	0.1
	0.1
	0.06
	0.04

Huffman encoding: 1st step

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








0.06






0.04

Example (1/3)

Computing a binary tree

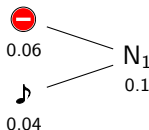
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
	0.1
	0.1
N_1	0.1

Huffman encoding: 1st step







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







Example (1/3)

Computing a binary tree

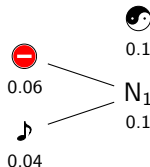
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
	0.1
	0.1
N_1	0.1

Huffman encoding: 1st step







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






Example (1/3)

Computing a binary tree

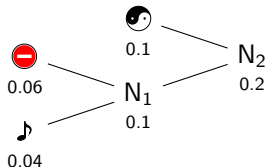
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
	0.1
N_2	0.2

Huffman encoding: 1st step







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






Example (1/3)

Computing a binary tree

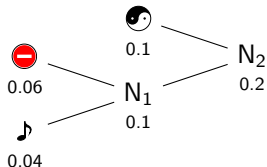
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
N_2	0.2
	0.1

Huffman encoding: 1st step







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






Example (1/3)

Computing a binary tree

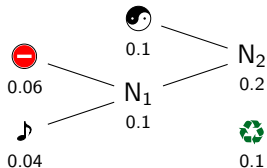
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
N_2	0.2
	0.1

Huffman encoding: 1st step







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





Example (1/3)

Computing a binary tree

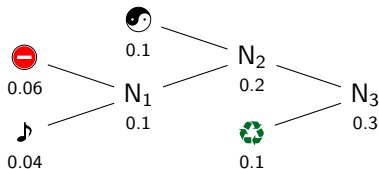
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
N_3	0.3

Huffman encoding: 1st step







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





Example (1/3)

Computing a binary tree

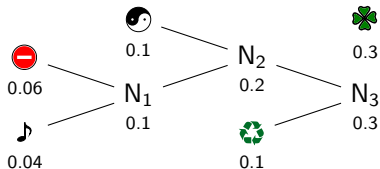
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
	0.3
N_3	0.3

Huffman encoding: 1st step







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




Example (1/3)

Computing a binary tree

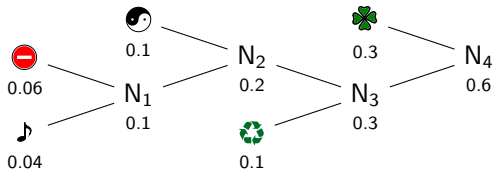
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
	0.4
N_4	0.6

Huffman encoding: 1st step







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




Example (1/3)

Computing a binary tree

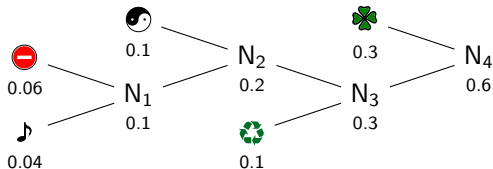
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
N_4	0.6
	0.4

Huffman encoding: 1st step







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




Example (1/3)

Computing a binary tree

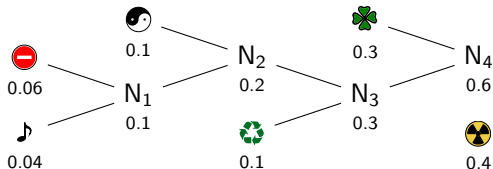
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
N_4	0.6
	0.4

Huffman encoding: 1st step







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



Example (1/3)

Computing a binary tree

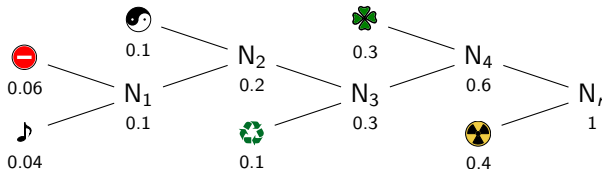
Take the following alphabet Σ with known probability distribution:

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

s_i	p_i
N_r	1







Huffman encoding: 1st step

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



Example (2/3)

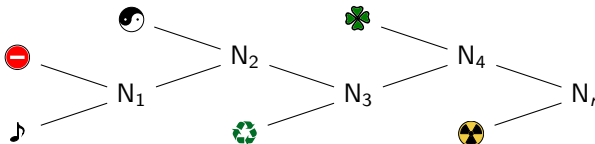
Retrieving the encoding for each symbol

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

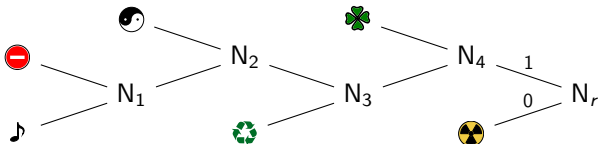
Retrieving the encoding for each symbol

s_i	P_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

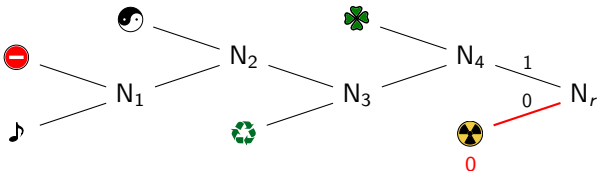
s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding

0







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

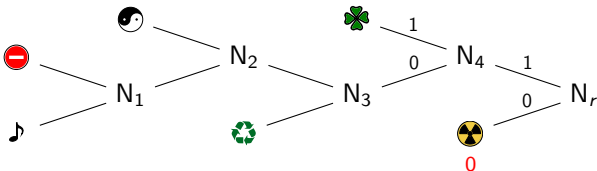
s_i	P_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding

0







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

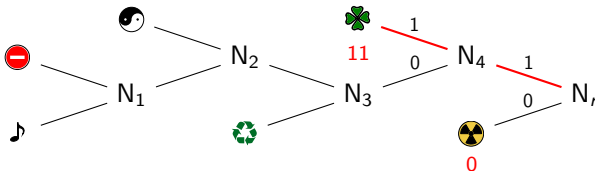
Encoding

11

0







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

s_i	P_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

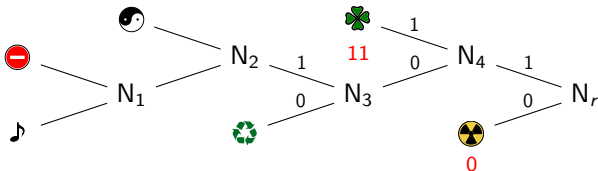
Encoding

11

0







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

s_i	P_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding

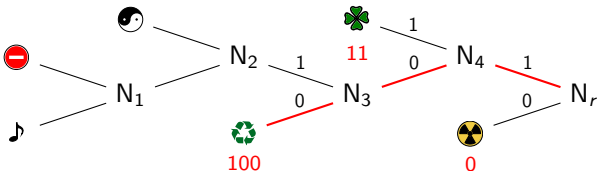
11

100

0







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding

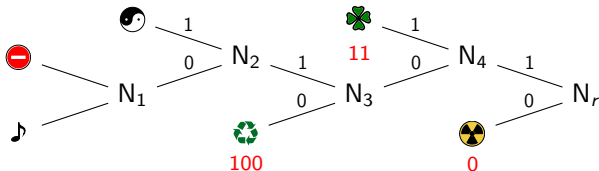
11

100

0







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

s_i	P_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding

11

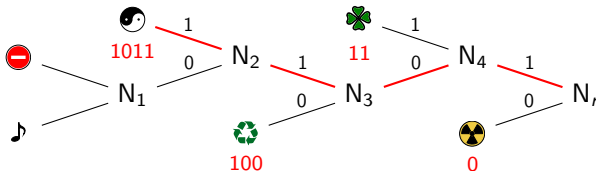
100

1011

0







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

s_i	P_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding

11

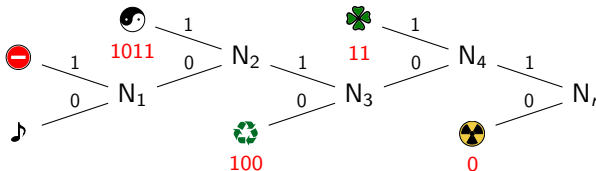
100

1011

0







Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (2/3)

Retrieving the encoding for each symbol

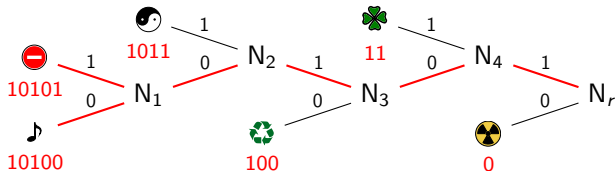
s_i	p_i
	0.3
	0.06
	0.1
	0.1
	0.4
	0.04

Encoding

11
10101
100
1011
0
10100

Huffman encoding: 2nd step

- 1) Traverse the tree from the root to the leaves (the symbols).
- 2) Always assign 1 to the left child and 0 to the right one (or the other way around).
- 3) Read encoding on the whole branch.



Example (3/3)

Performances of the encoding

The entropy of $\Sigma = \{\text{🍀}, \text{🚫}, \text{♻️}, \text{☯️}, \text{☢️}, \text{🎵}\}$ with previously given probability distribution is $H \simeq 2.144 \text{ Sh/symb.}$

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

Example (3/3)

Performances of the encoding

The entropy of $\Sigma = \{\text{♣}, \text{⊖}, \text{♻}, \text{☯}, \text{☢}, \text{♪}\}$ with previously given probability distribution is $H \simeq 2.144 \text{ Sh/symb.}$

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

The average length of Huffman encoding is

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

where $\ell(s_i)$ is the encoding length of symbol s_i .

⇒ $L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2 \text{ bits/symb.}$

Example (3/3)

Performances of the encoding

The entropy of $\Sigma = \{\text{♣}, \text{⊖}, \text{♻}, \text{☯}, \text{☢}, \text{♪}\}$ with previously given probability distribution is $H \simeq 2.144 \text{ Sh/symb.}$

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

The average length of Huffman encoding is

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

where $\ell(s_i)$ is the encoding length of symbol s_i .

⇒ $L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2 \text{ bits/symb.}$

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

→ 3000 bits are necessary to encode F without compression.

Example (3/3)

Performances of the encoding

The entropy of $\Sigma = \{\text{♣}, \text{⊖}, \text{♻}, \text{☯}, \text{☢}, \text{♪}\}$ with previously given probability distribution is $H \simeq 2.144 \text{ Sh/symb.}$

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

The average length of Huffman encoding is

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

where $\ell(s_i)$ is the encoding length of symbol s_i .

⇒ $L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2 \text{ bits/symb.}$

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

→ 3000 bits are necessary to encode F without compression.

→ The minimal compressed size of F is $H \times N_F = 2144$ bits.

Example (3/3)

Performances of the encoding

The entropy of $\Sigma = \{\text{♣}, \text{⊖}, \text{♻}, \text{☯}, \text{☢}, \text{♪}\}$ with previously given probability distribution is $H \simeq 2.144$ *Sh/symb.*

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

The average length of Huffman encoding is

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

where $\ell(s_i)$ is the encoding length of symbol s_i .

⇒ $L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2$ *bits/symb.*

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

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

Example (3/3)

Performances of the encoding

The entropy of $\Sigma = \{\text{♣}, \text{⊖}, \text{♻}, \text{☯}, \text{☢}, \text{♪}\}$ with previously given probability distribution is $H \simeq 2.144$ *Sh/symb.*

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

The average length of Huffman encoding is

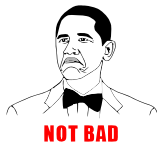
$$L = \mathbb{E}[\ell(s_i)] = \sum_{s_i \in \Sigma} p_i \ell(s_i) \text{ bits/symb.}$$

where $\ell(s_i)$ is the encoding length of symbol s_i .

⇒ $L = 0.3 \times 2 + 0.06 \times 5 + 0.1 \times 3 + 0.1 \times 4 + 0.4 \times 1 + 0.04 \times 5 = 2.2$ *bits/symb.*

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

- 3000 bits are necessary to encode F without compression.
- The minimal compressed size of F is $H \times N_F = 2144$ bits.
- 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).
⇒ it can be decoded on the fly.

Ex: **11001010100** ⇒ **|11|0|0|10101|100|** ⇒ 

Properties of Huffman encoding

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

Ex: **11001010100** ⇒ **|11|0|0|10101|100|** ⇒ 

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

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

Properties of Huffman encoding

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

Ex: **11001010100** ⇒ **11|0|0|10101|100** ⇒ 

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

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

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

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

Limits of Huffman encoding

Huffman encoding strategy suffers from two main drawbacks:

X It is a symbol-by-symbol encoding → 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.

Limits of Huffman encoding

Huffman encoding strategy suffers from two main drawbacks:

- ✗ It is a symbol-by-symbol encoding → 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.
- ✗ It relies (too) strongly on the probability distribution of the symbols to encode.
 - The file to compress must be scanned first to estimate the probability distribution.
 - The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.

Limits of Huffman encoding

Huffman encoding strategy suffers from two main drawbacks:

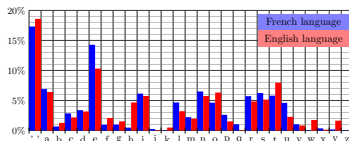
X It is a symbol-by-symbol encoding → does not handle *words* (as sequences of symbols).

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.

- The file to compress must be scanned first to estimate the probability distribution.
- The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.

↔ Classical text statistics can be used (language-dependent).



Limits of Huffman encoding

Huffman encoding strategy suffers from two main drawbacks:

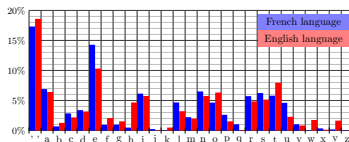
X It is a symbol-by-symbol encoding → does not handle *words* (as sequences of symbols).

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.

- The file to compress must be scanned first to estimate the probability distribution.
- The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.

↔ Classical text statistics can be used (language-dependent).



→ And what if the frequencies vary with time?

Limits of Huffman encoding

Huffman encoding strategy suffers from two main drawbacks:

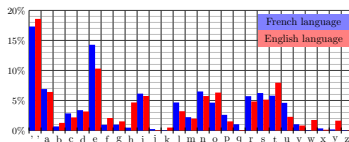
X It is a symbol-by-symbol encoding → does not handle *words* (as sequences of symbols).

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.

- The file to compress must be scanned first to estimate the probability distribution.
- The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.

↪ Classical text statistics can be used (language-dependent).



→ And what if the frequencies vary with time?

↪ Adaptive Huffman algorithm.

Limits of Huffman encoding

Huffman encoding strategy suffers from two main drawbacks:

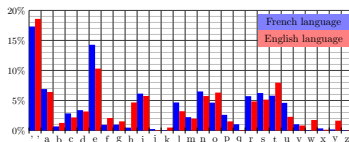
✗ It is a symbol-by-symbol encoding → does not handle *words* (as sequences of symbols).

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.

✗ It relies (too) strongly on the probability distribution of the symbols to encode.

- The file to compress must be scanned first to estimate the probability distribution.
- The frequency table (or the Huffman tree) must be stored with the text for the decoding stage.

↔ Classical text statistics can be used (language-dependent).



→ And what if the frequencies vary with time?

↔ Adaptive Huffman 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).

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



Julian Seward

bzip2 encoding

Idea: strengthen 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.

Burrows-Wheeler 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.

Burrows-Wheeler 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 \text{EOF}$)

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

b a n a n a \$

Burrows-Wheeler 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)

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

b	a	n	a	n	a	\$
\$	b	a	n	a	n	a

Burrows-Wheeler 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 \text{EOF}$)

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

b	a	n	a	n	a	\$
\$	b	a	n	a	n	a
a	\$	b	a	n	a	n

Burrows-Wheeler 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)

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

b	a	n	a	n	a	\$
\$	b	a	n	a	n	a
a	\$	b	a	n	a	n
n	a	\$	b	a	n	a
a	n	a	\$	b	a	n
n	a	n	a	\$	b	a
a	n	a	n	a	\$	b

Burrows-Wheeler 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)

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

b	a	n	a	n	a	\$
\$	b	a	n	a	n	a
a	\$	b	a	n	a	n
n	a	\$	b	a	n	a
a	n	a	\$	b	a	n
n	a	n	a	\$	b	a
a	n	a	n	a	\$	b

Burrows-Wheeler 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)

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

2nd step: Sort rows into lexicographic order.

b	a	n	a	n	a	\$		\$	b	a	n	a	n	a
\$	b	a	n	a	n	a		a	\$	b	a	n	a	n
a	\$	b	a	n	a	n		a	n	a	\$	b	a	n
n	a	\$	b	a	n	a	<div style="text-align: center;"> <div>lexicographic</div> <div>sort</div> <div>→</div> </div>	a	n	a	n	a	\$	b
a	n	a	\$	b	a	n		b	a	n	a	n	a	\$
n	a	n	a	\$	b	a		n	a	\$	b	a	n	a
a	n	a	n	a	\$	b		n	a	n	a	\$	b	a

Burrows-Wheeler 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)

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.

b	a	n	a	n	a	\$		1	\$	b	a	n	a	n	a
\$	b	a	n	a	n	a		2	a	\$	b	a	n	a	n
a	\$	b	a	n	a	n		3	a	n	a	\$	b	a	n
n	a	\$	b	a	n	a	<div>lexicographic sort</div>	4	a	n	a	n	a	\$	b
a	n	a	\$	b	a	n		5	b	a	n	a	n	a	\$
n	a	n	a	\$	b	a		6	n	a	\$	b	a	n	a
a	n	a	n	a	\$	b		7	n	a	n	a	\$	b	a

Burrows-Wheeler 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)

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.

b	a	n	a	n	a	\$		1	\$	b	a	n	a	n	a
\$	b	a	n	a	n	a		2	a	\$	b	a	n	a	n
a	\$	b	a	n	a	n		3	a	n	a	\$	b	a	n
n	a	\$	b	a	n	a		4	a	n	a	n	a	\$	b
a	n	a	\$	b	a	n		5	b	a	n	a	n	a	\$
n	a	n	a	\$	b	a		6	n	a	\$	b	a	n	a
a	n	a	n	a	\$	b		7	n	a	n	a	\$	b	a

$$\Rightarrow \text{BWT}(\text{banana\$}) = \text{5annb\$aa}$$

Burrows-Wheeler transform

Not convinced?

$$\text{BWT} \left(\begin{array}{l} \text{How much wood would} \\ \text{a woodchuck chuck if a} \\ \text{woodchuck could chuck} \\ \text{wood?} \end{array} \right) = \begin{array}{l} \text{dfkdkkwhkaad?d\$ udd} \\ \text{uuuu llooooiccccc ccc-} \\ \text{cuu oooowwwwwcHmh-} \\ \text{hhhoou} \end{array}$$

And that sure looks way better. How much better?

Burrows-Wheeler transform

Not convinced?

$$\text{BWT} \left(\begin{array}{l} \text{How much wood would} \\ \text{a woodchuck chuck if a} \\ \text{woodchuck could chuck} \\ \text{wood?} \end{array} \right) = \begin{array}{l} \text{dfkdkkwhkaad?d\$ udd} \\ \text{uuuu llooooiccccc ccc-} \\ \text{cuu oooowwwwwcHmh-} \\ \text{hhhooo} \end{array}$$

And that sure looks way better. How much better?

$$\text{RLE} \left(\begin{array}{l} \text{dfkdkkwhkaad?d\$ udd} \\ \text{uuuu llooooiccccc ccc-} \\ \text{cuu oooowwwwwcHmh-} \\ \text{hhhooo} \end{array} \right) = \begin{array}{l} \text{dfkd\#2kwhk\#2ad?d\$} \\ \text{u\#2d \#4u \#2l\#4oi} \\ \text{\#5c \#4c\#2u \#4o\#5} \\ \text{wcHm\#4h\#3o} \end{array}$$

Burrows-Wheeler transform

Not convinced?

$$\text{BWT} \left(\begin{array}{l} \text{How much wood would} \\ \text{a woodchuck chuck if a} \\ \text{woodchuck could chuck} \\ \text{wood?} \end{array} \right) = \begin{array}{l} \text{dfkdkkwhkaad?d\$ udd} \\ \text{uuuu llooooiccccc ccc-} \\ \text{cuu oooowwwwwcHmh-} \\ \text{hhhoou} \end{array}$$

And that sure looks way better. How much better?

$$\text{RLE} \left(\begin{array}{l} \text{dfkdkkwhkaad?d\$ udd} \\ \text{uuuu llooooiccccc ccc-} \\ \text{cuu oooowwwwwcHmh-} \\ \text{hhhoou} \end{array} \right) = \begin{array}{l} \text{dfkd\#2kwhk\#2ad?d\$} \\ \text{u\#2d \#4u \#2l\#4oi} \\ \text{\#5c \#4c\#2u \#4o\#5} \\ \text{wcHm\#4h\#3o} \end{array}$$

Transforming 65 characters into ... 61 characters!



Burrows-Wheeler transform

Not convinced?

$$\text{BWT} \left(\begin{array}{l} \text{How much wood would} \\ \text{a woodchuck chuck if a} \\ \text{woodchuck could chuck} \\ \text{wood?} \end{array} \right) = \begin{array}{l} \text{dfkd kkw hkaad?d\$ udd} \\ \text{uuuu llooooiccccc ccc-} \\ \text{cuu oooowwwwwcHmh-} \\ \text{hhho o} \end{array}$$

And that sure looks way better. How much better?

$$\text{RLE} \left(\begin{array}{l} \text{dfkd kkw hkaad?d\$ udd} \\ \text{uuuu llooooiccccc ccc-} \\ \text{cuu oooowwwwwcHmh-} \\ \text{hhho o} \end{array} \right) = \begin{array}{l} \text{dfkd\#2kw h k\#2ad?d\$} \\ \text{u\#2d \#4u \#2l\#4oi} \\ \text{\#5c \#4c\#2u \#4o\#5} \\ \text{wcHm\#4h\#3o} \end{array}$$

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.

Burrows-Wheeler transform

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.

Burrows-Wheeler transform

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
n
b
\$
a
a

Burrows-Wheeler transform

The decoding part

It is easily possible to recover **banana\$** from **5annb\$a** 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.

\$
a
a
a
b
n
n

Burrows-Wheeler transform

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.

a	\$
n	a
n	a
b	a
\$	b
a	n
a	n

Burrows-Wheeler transform

The decoding part

It is easily possible to recover **banana\$** from **5annb\$a** 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.

⇒ Repeat until the whole table is recreated.

\$	b
a	\$
a	n
a	n
b	a
n	a
n	a

Burrows-Wheeler transform

The decoding part

It is easily possible to recover **banana\$** from **5annb\$a** 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.

⇒ Repeat until the whole table is recreated.

a	\$	b
n	a	\$
n	a	n
b	a	n
\$	b	a
a	n	a
a	n	a

Burrows-Wheeler transform

The decoding part

It is easily possible to recover **banana\$** from **5annb\$a** 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.

⇒ Repeat until the whole table is recreated.

\$	b	a
a	\$	b
a	n	a
a	n	a
b	a	n
n	a	\$
n	a	n

Burrows-Wheeler transform

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.

a	\$	b	a
n	a	\$	b
n	a	n	a
b	a	n	a
\$	b	a	n
a	n	a	\$
a	n	a	n

Burrows-Wheeler transform

The decoding part

It is easily possible to recover **banana\$** from **5annb\$a** 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.

⇒ Repeat until the whole table is recreated.

\$	b	a	n
a	\$	b	a
a	n	a	\$
a	n	a	n
b	a	n	a
n	a	\$	b
n	a	n	a

Burrows-Wheeler transform

The decoding part

It is easily possible to recover **banana\$** from **5annb\$a** 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.

⇒ Repeat until the whole table is recreated.

a	\$	b	a	n
n	a	\$	b	a
n	a	n	a	\$
b	a	n	a	n
\$	b	a	n	a
a	n	a	\$	b
a	n	a	n	a

Burrows-Wheeler transform

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.

\$	b	a	n	a
a	\$	b	a	n
a	n	a	\$	b
a	n	a	n	a
b	a	n	a	n
n	a	\$	b	a
n	a	n	a	\$

Burrows-Wheeler transform

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.

a	\$	b	a	n	a
n	a	\$	b	a	n
n	a	n	a	\$	b
b	a	n	a	n	a
\$	b	a	n	a	n
a	n	a	\$	b	a
a	n	a	n	a	\$

Burrows-Wheeler transform

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.

\$	b	a	n	a	n
a	\$	b	a	n	a
a	n	a	\$	b	a
a	n	a	n	a	\$
b	a	n	a	n	a
n	a	\$	b	a	n
n	a	n	a	\$	b

Burrows-Wheeler transform

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.

a	\$	b	a	n	a	n
n	a	\$	b	a	n	a
n	a	n	a	\$	b	a
b	a	n	a	n	a	\$
\$	b	a	n	a	n	a
a	n	a	\$	b	a	n
a	n	a	n	a	\$	b

Burrows-Wheeler transform

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.

\$	b	a	n	a	n	a
a	\$	b	a	n	a	n
a	n	a	\$	b	a	n
a	n	a	n	a	\$	b
b	a	n	a	n	a	\$
n	a	\$	b	a	n	a
n	a	n	a	\$	b	a

Burrows-Wheeler transform

The decoding part

It is easily possible to recover **banana\$** from **5annb\$a** 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.

⇒ Repeat until the whole table is recreated.

BWT⁻¹: Read the corresponding row in the table.

1	\$	b	a	n	a	n	a
2	a	\$	b	a	n	a	n
3	a	n	a	\$	b	a	n
4	a	n	a	n	a	\$	b
5	b	a	n	a	n	a	\$
6	n	a	\$	b	a	n	a
7	n	a	n	a	\$	b	a

Burrows-Wheeler transform

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.

1	\$	b	a	n	a	n	a
2	a	\$	b	a	n	a	n
3	a	n	a	\$	b	a	n
4	a	n	a	n	a	\$	b
5	b	a	n	a	n	a	\$
6	n	a	\$	b	a	n	a
7	n	a	n	a	\$	b	a

$$\Rightarrow \text{BWT}^{-1} \left(\text{5annb\$aa} \right) = \text{banana\$}$$

Move-to-front transform

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.

Move-to-front transform

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

Move-to-front transform

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

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

Move-to-front transform

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

a	$\rightarrow 0$	first symbol of \mathcal{L}_0	
n	$\rightarrow 2$	third symbol of \mathcal{L}_0	$\mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n, a, b, \$]$

Move-to-front transform

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

a	$\rightarrow 0$	first symbol of \mathcal{L}_0	$\mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n, a, b, \$]$
n	$\rightarrow 2$	third symbol of \mathcal{L}_0	
n	$\rightarrow 0$	first symbol of \mathcal{L}_1	

Move-to-front transform

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

a	$\rightarrow 0$	first symbol of \mathcal{L}_0	
n	$\rightarrow 2$	third symbol of \mathcal{L}_0	$\mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n, a, b, \$]$
n	$\rightarrow 0$	first symbol of \mathcal{L}_1	
b	$\rightarrow 2$	third symbol of \mathcal{L}_1	$\mathcal{L}_1 \rightarrow \mathcal{L}_2 = [b, n, a, \$]$

Move-to-front transform

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

a	$\rightarrow 0$	first symbol of \mathcal{L}_0	
n	$\rightarrow 2$	third symbol of \mathcal{L}_0	$\mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n, a, b, \$]$
n	$\rightarrow 0$	first symbol of \mathcal{L}_1	
b	$\rightarrow 2$	third symbol of \mathcal{L}_1	$\mathcal{L}_1 \rightarrow \mathcal{L}_2 = [b, n, a, \$]$
\$	$\rightarrow 3$	fourth symbol of \mathcal{L}_2	$\mathcal{L}_2 \rightarrow \mathcal{L}_3 = [\$, b, n, a]$

Move-to-front transform

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

a	$\rightarrow 0$	first symbol of \mathcal{L}_0	
n	$\rightarrow 2$	third symbol of \mathcal{L}_0	$\mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n, a, b, \$]$
n	$\rightarrow 0$	first symbol of \mathcal{L}_1	
b	$\rightarrow 2$	third symbol of \mathcal{L}_1	$\mathcal{L}_1 \rightarrow \mathcal{L}_2 = [b, n, a, \$]$
\$	$\rightarrow 3$	fourth symbol of \mathcal{L}_2	$\mathcal{L}_2 \rightarrow \mathcal{L}_3 = [\$, b, n, a]$
a	$\rightarrow 3$	fourth symbol of \mathcal{L}_3	$\mathcal{L}_3 \rightarrow \mathcal{L}_4 = [a, \$, b, n]$

Move-to-front transform

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

a	$\rightarrow 0$	first symbol of \mathcal{L}_0	
n	$\rightarrow 2$	third symbol of \mathcal{L}_0	$\mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n, a, b, \$]$
n	$\rightarrow 0$	first symbol of \mathcal{L}_1	
b	$\rightarrow 2$	third symbol of \mathcal{L}_1	$\mathcal{L}_1 \rightarrow \mathcal{L}_2 = [b, n, a, \$]$
\$	$\rightarrow 3$	fourth symbol of \mathcal{L}_2	$\mathcal{L}_2 \rightarrow \mathcal{L}_3 = [\$, b, n, a]$
a	$\rightarrow 3$	fourth symbol of \mathcal{L}_3	$\mathcal{L}_3 \rightarrow \mathcal{L}_4 = [a, \$, b, n]$
a	$\rightarrow 0$	first symbol of \mathcal{L}_4	

Move-to-front transform

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

a	$\rightarrow 0$	first symbol of \mathcal{L}_0	
n	$\rightarrow 2$	third symbol of \mathcal{L}_0	$\mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n, a, b, \$]$
n	$\rightarrow 0$	first symbol of \mathcal{L}_1	
b	$\rightarrow 2$	third symbol of \mathcal{L}_1	$\mathcal{L}_1 \rightarrow \mathcal{L}_2 = [b, n, a, \$]$
\$	$\rightarrow 3$	fourth symbol of \mathcal{L}_2	$\mathcal{L}_2 \rightarrow \mathcal{L}_3 = [\$, b, n, a]$
a	$\rightarrow 3$	fourth symbol of \mathcal{L}_3	$\mathcal{L}_3 \rightarrow \mathcal{L}_4 = [a, \$, b, n]$
a	$\rightarrow 0$	first symbol of \mathcal{L}_4	

MTF(**annb\$aa**) \Rightarrow 0202330

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

Move-to-front transform

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

a	$\rightarrow 0$	first symbol of \mathcal{L}_0	
n	$\rightarrow 2$	third symbol of \mathcal{L}_0	$\mathcal{L}_0 \rightarrow \mathcal{L}_1 = [n, a, b, \$]$
n	$\rightarrow 0$	first symbol of \mathcal{L}_1	
b	$\rightarrow 2$	third symbol of \mathcal{L}_1	$\mathcal{L}_1 \rightarrow \mathcal{L}_2 = [b, n, a, \$]$
\$	$\rightarrow 3$	fourth symbol of \mathcal{L}_2	$\mathcal{L}_2 \rightarrow \mathcal{L}_3 = [\$, b, n, a]$
a	$\rightarrow 3$	fourth symbol of \mathcal{L}_3	$\mathcal{L}_3 \rightarrow \mathcal{L}_4 = [a, \$, b, n]$
a	$\rightarrow 0$	first symbol of \mathcal{L}_4	

MTF(**annb\$aa**) \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).

Final overview of bzip2 compression algorithm

1st step: BWT → sort data to create runs of similar symbols.

2nd step: MTF → encode sorted data into digits with high anisotropy.

3rd step: Huffman → 0 gets a short encoding ⇒ 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
```

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

LZW compression algorithm

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



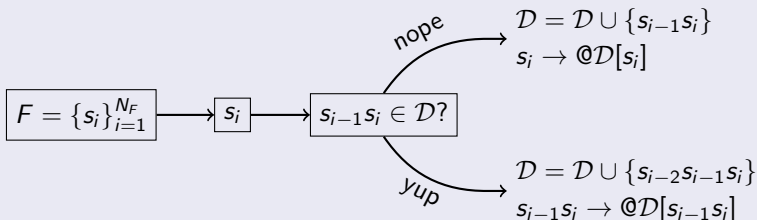
Abraham Lempel



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} .



Example

LZW encoding

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

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

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
	T			

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 $F = \text{TO BE OR NOT TO BE}$ with $\mathcal{D} \equiv \text{ASCII table}$.

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
T	O	TO	256	$@\mathcal{D}[\text{T}]$

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} .

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$

What is happening?

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} .

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$

What is happening?

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. **_** is encoded by its address in \mathcal{D} .

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[\text{-}]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$

What is happening?

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

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[T]$
O	-	O_	257	$@\mathcal{D}[O]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[B]$
E	-	E_	260	$@\mathcal{D}[E]$

What is happening?

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

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$

What is happening?

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

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[T]$
O	-	O_	257	$@\mathcal{D}[O]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[B]$
E	-	E_	260	$@\mathcal{D}[E]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[O]$

What is happening?

It keeps adding in \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...

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$

What is happening?

It keeps adding in \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...

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$

What is happening?

It keeps adding in \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...

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$

What is happening?

It keeps adding in \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...

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$

What is happening?

It keeps adding in \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...

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$

What is happening?

It keeps adding in \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...

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$

What is happening?

It keeps adding in \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...

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O	TO		

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} .

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O	TO		

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 already in \mathcal{D} , at the address 256 \rightarrow it is not inserted in \mathcal{D} again.

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O			X
TO				

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 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.

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O			
TO	-	TO_	269	$@\mathcal{D}[\text{TO}]$

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} .
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} .

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O			
TO	-	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} .
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} , *i.e.*, 256.

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O			$@\mathcal{D}[\%]$
TO	-	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} .
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} , *i.e.*, 256.
6. But wait, 256 requires 9 bits to be encoded (instead of 8 bits so far). \rightarrow 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.

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	B	258	$@\mathcal{D}[\text{B}]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[\text{O}]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[\text{N}]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[\text{B}]$
T	O			$@\mathcal{D}[\text{B}]$
TO	-	TO_	269	256
-	B			

What is happening?

1. **B** 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 already in \mathcal{D} , at the address **258** → it is not inserted in \mathcal{D} again.
5. Nothing is output, and **_B** is put in the buffer.

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O			$@\mathcal{D}[\%]$
TO	-	TO_	269	256
-	B			
_B	E	_BE	270	258

What is happening?

1. **_B** is in the buffer.
2. **E** is read.
3. The sequence **_BE** is created and tested to belong to \mathcal{D} .
4. **_BE** 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.

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[\text{_}]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[\text{_}]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[\text{_}]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[\text{_}]$
T	O			$@\mathcal{D}[\text{\%}]$
TO	-	TO_	269	256
-	B			
_B	E	_BE	270	258
E	\$			$@\mathcal{D}[\text{E}]$

What is happening?

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

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O			$@\mathcal{D}[\%]$
TO	-	TO_	269	256
-	B			
_B	E	_BE	270	258
E	\$			$@\mathcal{D}[\text{E}]$

Compression performance:

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

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

Example

LZW encoding

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

Buffer	Input	$\mathcal{D} \cup \{.\}$	$@\mathcal{D}[\cdot]$	Output
	T			
T	O	TO	256	$@\mathcal{D}[\text{T}]$
O	-	O_	257	$@\mathcal{D}[\text{O}]$
-	B	_B	258	$@\mathcal{D}[_]$
B	E	BE	259	$@\mathcal{D}[\text{B}]$
E	-	E_	260	$@\mathcal{D}[\text{E}]$
-	O	_O	261	$@\mathcal{D}[_]$
O	R	OR	262	$@\mathcal{D}[\text{O}]$
R	-	R_	263	$@\mathcal{D}[\text{R}]$
-	N	_N	264	$@\mathcal{D}[_]$
N	O	NO	265	$@\mathcal{D}[\text{N}]$
O	T	OT	266	$@\mathcal{D}[\text{O}]$
T	-	T_	267	$@\mathcal{D}[\text{T}]$
-	T	_T	268	$@\mathcal{D}[_]$
T	O			$@\mathcal{D}[\%]$
TO	-	TO_	269	256
-	B			
_B	E	_BE	270	258
E	\$			$@\mathcal{D}[\text{E}]$

Compression performance:

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

Ex:

How much wood would a woodchuck 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).

Example

LZW decoding

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

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

LZW decoding

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

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

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
T	$@\mathcal{D}[\text{T}]$ $@\mathcal{D}[\text{O}]$	TO	256	T

What is happening?

1. **T** is put in the buffer.
2. **O** is recognized from $@\mathcal{D}[\text{O}]$.
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.

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[T]$			
T	$@\mathcal{D}[O]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O

What is happening?

1. **O** is put in the buffer.
2. **_** is recognized from $@\mathcal{D}[_]$.
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.

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D} [T]$			
T	$@\mathcal{D} [O]$	TO	256	T
O	$@\mathcal{D} [-]$	O_	257	O
-	$@\mathcal{D} [B]$	_B	258	-
B	$@\mathcal{D} [E]$	BE	259	B
E	$@\mathcal{D} [-]$	E_	260	E
-	$@\mathcal{D} [O]$	_O	261	-
O	$@\mathcal{D} [R]$	OR	262	O
R	$@\mathcal{D} [-]$	R_	263	R
-	$@\mathcal{D} [N]$	_N	264	-
N	$@\mathcal{D} [O]$	NO	265	N
O	$@\mathcal{D} [T]$	OT	266	O
T	$@\mathcal{D} [-]$	T_	267	T

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.


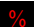
Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[\text{T}]$			
T	$@\mathcal{D}[\text{O}]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O
-	$@\mathcal{D}[\text{B}]$	_B	258	-
B	$@\mathcal{D}[\text{E}]$	BE	259	B
E	$@\mathcal{D}[_]$	E_	260	E
-	$@\mathcal{D}[\text{O}]$	_O	261	-
O	$@\mathcal{D}[\text{R}]$	OR	262	O
R	$@\mathcal{D}[_]$	R_	263	R
-	$@\mathcal{D}[\text{N}]$	_N	264	-
N	$@\mathcal{D}[\text{O}]$	NO	265	N
O	$@\mathcal{D}[\text{T}]$	OT	266	O
T	$@\mathcal{D}[_]$	T_	267	T
-	$@\mathcal{D}[\%]$			

What is happening?

1.  is put in the buffer.
2. The special character  is recognized from $@\mathcal{D}[\%]$ \rightarrow all following addresses will have to be decoded on 9 bits instead of 8.


Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[\text{T}]$			
T	$@\mathcal{D}[\text{O}]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O
-	$@\mathcal{D}[\text{B}]$	_B	258	-
B	$@\mathcal{D}[\text{E}]$	BE	259	B
E	$@\mathcal{D}[_]$	E_	260	E
-	$@\mathcal{D}[\text{O}]$	_O	261	-
O	$@\mathcal{D}[\text{R}]$	OR	262	O
R	$@\mathcal{D}[_]$	R_	263	R
-	$@\mathcal{D}[\text{N}]$	_N	264	-
N	$@\mathcal{D}[\text{O}]$	NO	265	N
O	$@\mathcal{D}[\text{T}]$	OT	266	O
T	$@\mathcal{D}[_]$	T_	267	T
-	$@\mathcal{D}[\text{\%}]$			
-				

What is happening?

1.  is in the buffer.


Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D} [T]$			
T	$@\mathcal{D} [O]$	TO	256	T
O	$@\mathcal{D} [_]$	O_	257	O
-	$@\mathcal{D} [B]$	_B	258	-
B	$@\mathcal{D} [E]$	BE	259	B
E	$@\mathcal{D} [_]$	E_	260	E
-	$@\mathcal{D} [O]$	_O	261	-
O	$@\mathcal{D} [R]$	OR	262	O
R	$@\mathcal{D} [_]$	R_	263	R
-	$@\mathcal{D} [N]$	_N	264	-
N	$@\mathcal{D} [O]$	NO	265	N
O	$@\mathcal{D} [T]$	OT	266	O
T	$@\mathcal{D} [_]$	T_	267	T
-	$@\mathcal{D} [\%]$			
-	256			

What is happening?

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


Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[\text{T}]$			
T	$@\mathcal{D}[\text{O}]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O
-	$@\mathcal{D}[\text{B}]$	_B	258	-
B	$@\mathcal{D}[\text{E}]$	BE	259	B
E	$@\mathcal{D}[_]$	E_	260	E
-	$@\mathcal{D}[\text{O}]$	_O	261	-
O	$@\mathcal{D}[\text{R}]$	OR	262	O
R	$@\mathcal{D}[_]$	R_	263	R
-	$@\mathcal{D}[\text{N}]$	_N	264	-
N	$@\mathcal{D}[\text{O}]$	NO	265	N
O	$@\mathcal{D}[\text{T}]$	OT	266	O
T	$@\mathcal{D}[_]$	T_	267	T
-	$@\mathcal{D}[\%]$			
-	256			

What is happening?

1.  is in the buffer.
2. 256 is received, that is, $@\mathcal{D}[\text{TO}]$.

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[T]$			
T	$@\mathcal{D}[O]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O
-	$@\mathcal{D}[B]$	_B	258	-
B	$@\mathcal{D}[E]$	BE	259	B
E	$@\mathcal{D}[_]$	E_	260	E
-	$@\mathcal{D}[O]$	_O	261	-
O	$@\mathcal{D}[R]$	OR	262	O
R	$@\mathcal{D}[_]$	R_	263	R
-	$@\mathcal{D}[N]$	_N	264	-
N	$@\mathcal{D}[O]$	NO	265	N
O	$@\mathcal{D}[T]$	OT	266	O
T	$@\mathcal{D}[_]$	T_	267	T
-	$@\mathcal{D}[\%]$			
-	$@\mathcal{D}[TO]$	_T		

What is happening?

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

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[\text{T}]$			
T	$@\mathcal{D}[\text{O}]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O
-	$@\mathcal{D}[\text{B}]$	_B	258	-
B	$@\mathcal{D}[\text{E}]$	BE	259	B
E	$@\mathcal{D}[_]$	E_	260	E
-	$@\mathcal{D}[\text{O}]$	_O	261	-
O	$@\mathcal{D}[\text{R}]$	OR	262	O
R	$@\mathcal{D}[_]$	R_	263	R
-	$@\mathcal{D}[\text{N}]$	_N	264	-
N	$@\mathcal{D}[\text{O}]$	NO	265	N
O	$@\mathcal{D}[\text{T}]$	OT	266	O
T	$@\mathcal{D}[_]$	T_	267	T
-	$@\mathcal{D}[\text{\%}]$			
-	$@\mathcal{D}[\text{TO}]$	_T	268	

What is happening?

1. **_** is in the buffer.
2. 256 is received, that is, $@\mathcal{D}[\text{TO}]$.
3. The sequence **_T** is created using **_** in the buffer and the first letter of the received sequence **TO**.
4. **_T** is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 268.

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[\text{T}]$			
T	$@\mathcal{D}[\text{O}]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O
-	$@\mathcal{D}[\text{B}]$	_B	258	-
B	$@\mathcal{D}[\text{E}]$	BE	259	B
E	$@\mathcal{D}[_]$	E_	260	E
-	$@\mathcal{D}[\text{O}]$	_O	261	-
O	$@\mathcal{D}[\text{R}]$	OR	262	O
R	$@\mathcal{D}[_]$	R_	263	R
-	$@\mathcal{D}[\text{N}]$	_N	264	-
N	$@\mathcal{D}[\text{O}]$	NO	265	N
O	$@\mathcal{D}[\text{T}]$	OT	266	O
T	$@\mathcal{D}[_]$	T_	267	T
-	$@\mathcal{D}[\%]$			
-	$@\mathcal{D}[\text{TO}]$	_T	268	-

What is happening?

1. **_** is in the buffer.
2. 256 is received, that is, $@\mathcal{D}[\text{TO}]$.
3. The sequence **_T** is created using **_** in the buffer and the first letter of the received sequence **TO**.
4. **_T** is not in \mathcal{D} yet \rightarrow it is inserted at the next available address, *i.e.* 268.
5. **_** is output.

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D} [T]$			
T	$@\mathcal{D} [O]$	TO	256	T
O	$@\mathcal{D} [-]$	O_	257	O
-	$@\mathcal{D} [B]$	_B	258	-
B	$@\mathcal{D} [E]$	BE	259	B
E	$@\mathcal{D} [-]$	E_	260	E
-	$@\mathcal{D} [O]$	_O	261	-
O	$@\mathcal{D} [R]$	OR	262	O
R	$@\mathcal{D} [-]$	R_	263	R
-	$@\mathcal{D} [N]$	_N	264	-
N	$@\mathcal{D} [O]$	NO	265	N
O	$@\mathcal{D} [T]$	OT	266	O
T	$@\mathcal{D} [-]$	T_	267	T
-	$@\mathcal{D} [%]$			
-	$@\mathcal{D} [TO]$	_T	268	-
TO	258			

What is happening?

1. **TO** is put in the buffer.
2. 258 is received.

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D} [T]$			
T	$@\mathcal{D} [O]$	TO	256	T
O	$@\mathcal{D} [-]$	O_	257	O
-	$@\mathcal{D} [B]$	_B	258	-
B	$@\mathcal{D} [E]$	BE	259	B
E	$@\mathcal{D} [-]$	E_	260	E
-	$@\mathcal{D} [O]$	_O	261	-
O	$@\mathcal{D} [R]$	OR	262	O
R	$@\mathcal{D} [-]$	R_	263	R
-	$@\mathcal{D} [N]$	_N	264	-
N	$@\mathcal{D} [O]$	NO	265	N
O	$@\mathcal{D} [T]$	OT	266	O
T	$@\mathcal{D} [-]$	T_	267	T
-	$@\mathcal{D} [\%]$			
-	$@\mathcal{D} [TO]$	_T	268	-
TO	$@\mathcal{D} [-B]$			

What is happening?

1. **TO** is put in the buffer.
2. 258 is received, that is, **_B**.

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{\cdot\}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[\text{T}]$			
T	$@\mathcal{D}[\text{O}]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O
-	$@\mathcal{D}[\text{B}]$	_B	258	-
B	$@\mathcal{D}[\text{E}]$	BE	259	B
E	$@\mathcal{D}[_]$	E_	260	E
-	$@\mathcal{D}[\text{O}]$	_O	261	-
O	$@\mathcal{D}[\text{R}]$	OR	262	O
R	$@\mathcal{D}[_]$	R_	263	R
-	$@\mathcal{D}[\text{N}]$	_N	264	-
N	$@\mathcal{D}[\text{O}]$	NO	265	N
O	$@\mathcal{D}[\text{T}]$	OT	266	O
T	$@\mathcal{D}[_]$	T_	267	T
-	$@\mathcal{D}[\%]$			
-	$@\mathcal{D}[\text{TO}]$	_T	268	-
TO	$@\mathcal{D}[_ \text{B}]$	TO_	269	TO

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 \mathcal{D} yet \rightarrow it is inserted at the next available address, i.e. 269.
5. **TO** is output.

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D}[\text{T}]$			
T	$@\mathcal{D}[\text{O}]$	TO	256	T
O	$@\mathcal{D}[_]$	O_	257	O
-	$@\mathcal{D}[\text{B}]$	_B	258	-
B	$@\mathcal{D}[\text{E}]$	BE	259	B
E	$@\mathcal{D}[_]$	E_	260	E
-	$@\mathcal{D}[\text{O}]$	_O	261	-
O	$@\mathcal{D}[\text{R}]$	OR	262	O
R	$@\mathcal{D}[_]$	R_	263	R
-	$@\mathcal{D}[\text{N}]$	_N	264	-
N	$@\mathcal{D}[\text{O}]$	NO	265	N
O	$@\mathcal{D}[\text{T}]$	OT	266	O
T	$@\mathcal{D}[_]$	T_	267	T
-	$@\mathcal{D}[\%]$			
-	$@\mathcal{D}[\text{TO}]$	_T	268	-
TO	$@\mathcal{D}[_B]$		269	TO
_B	$@\mathcal{D}[\text{E}]$	_BE	270	_B

What is happening?

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

Example

LZW decoding

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

Buffer	Input	$\mathcal{D} \cup \{ \cdot \}$	$@\mathcal{D}[\cdot]$	Output
	$@\mathcal{D} [T]$			
T	$@\mathcal{D} [O]$	TO	256	T
O	$@\mathcal{D} [_]$	O_	257	O
-	$@\mathcal{D} [B]$	_B	258	-
B	$@\mathcal{D} [E]$	BE	259	B
E	$@\mathcal{D} [_]$	E_	260	E
-	$@\mathcal{D} [O]$	_O	261	-
O	$@\mathcal{D} [R]$	OR	262	O
R	$@\mathcal{D} [_]$	R_	263	R
-	$@\mathcal{D} [N]$	_N	264	-
N	$@\mathcal{D} [O]$	NO	265	N
O	$@\mathcal{D} [T]$	OT	266	O
T	$@\mathcal{D} [_]$	T_	267	T
-	$@\mathcal{D} [%]$			
-	$@\mathcal{D} [TO]$	_T	268	-
TO	$@\mathcal{D} [_B]$		269	TO
_B	$@\mathcal{D} [E]$	_BE	270	_B
E	\$			E

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 $@\mathcal{D} [E]$ and is output.

Concluding remarks on LZW

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

The dictionary \mathcal{D} cannot grow indefinitely.

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

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

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