

# Les données sont D<sup>1</sup>

Étude de méthodes de compression sans pertes

Thomas BAGREL

Lycée Henri POINCARÉ, Nancy

TIPE session 2018

# Aperçu

## Régularités et gains

Théorie

zip récursif

## Composantes de la compression

## Codage

Problème résolu

Inefficacité de HUFFMAN – pourquoi

# Régularités et gains

## Théorie

Compression des données sans pertes

# Régularités et gains

## Théorie

Compression des données sans pertes

- ▶ exploiter les régularités des données

# Régularités et gains

## Théorie

### Compression des données sans pertes

- ▶ exploiter les régularités des données
- ▶ données aléatoires : pas de gain

# Régularités et gains

## Théorie

### Compression des données sans pertes

- ▶ exploiter les régularités des données
- ▶ données aléatoires : pas de gain

Théorème. (Entropie de SHANNON)

$$H(S) = - \sum_{i=1}^n p_i \log_2(p_i)$$

$H(S)$  : *nb. de bits moyen par symbole de la source*

# Régularités et gains

## Théorie

### Compression des données sans pertes

- ▶ exploiter les régularités des données
- ▶ données aléatoires : pas de gain

### Théorème. (Entropie de SHANNON)

$$H(S) = - \sum_{i=1}^n p_i \log_2(p_i)$$

$H(S)$  : *nb. de bits moyen par symbole de la source*

- ▶ une fois les données compressées (dérivées) une fois, plus aucune régularité

# Régularités et gains

zip récursif

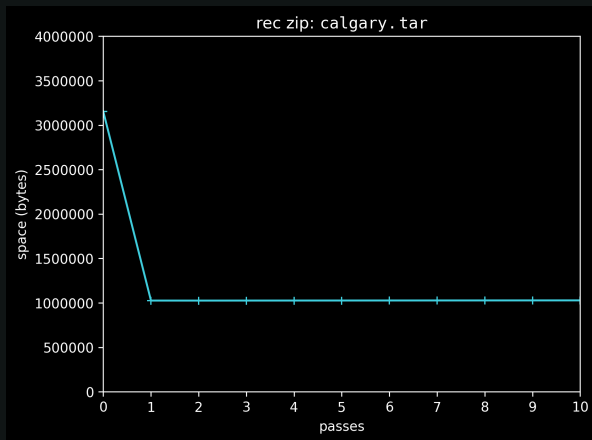
Expérience : compresser récursivement un fichier avec le même algorithme (zip)



# Régularités et gains

zip récursif

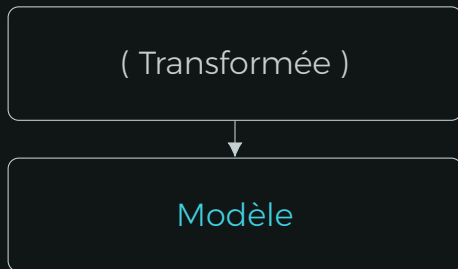
Expérience : compresser récursivement un fichier avec le même algorithme (zip)



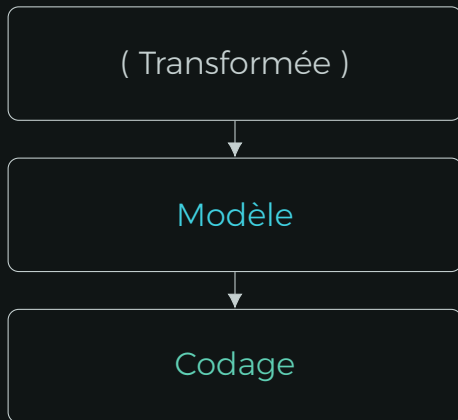
# Composantes de la compression

( Transformée )

# Composantes de la compression



# Composantes de la compression



# Codage

## Problème résolu

Le codage, contrairement aux apparences, est un problème **résolu**

- ▶ Si les  $p_i$  sont connus, la limite de compression théorique est donnée par SHANNON
- ▶ HUFFMAN permet d'approcher cette limite
- ▶ Le codage arithmétique l'atteint

# Codage

Problème résolu

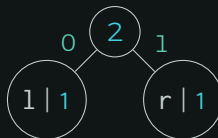
Le codage arithmétique sera détaillé plus loin.  
Codage d'HUFFMAN pour turlututu



# Codage

## Problème résolu

Le codage arithmétique sera détaillé plus loin.  
Codage d'HUFFMAN pour turlututu

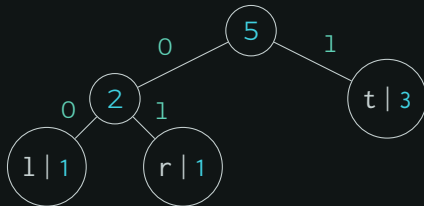


# Codage

## Problème résolu

Le codage arithmétique sera détaillé plus loin.

Codage d'HUFFMAN pour turlututu



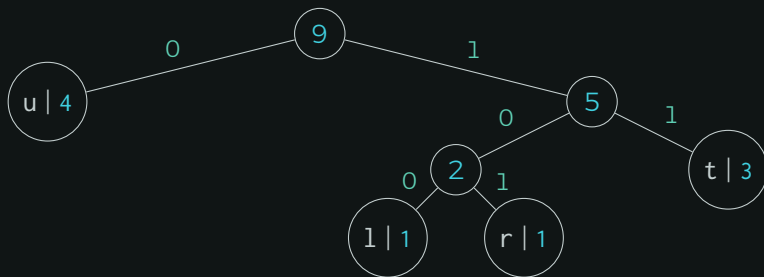


# Codage

Problème résolu

Le codage arithmétique sera détaillé plus loin.

Codage d'HUFFMAN pour turlututu

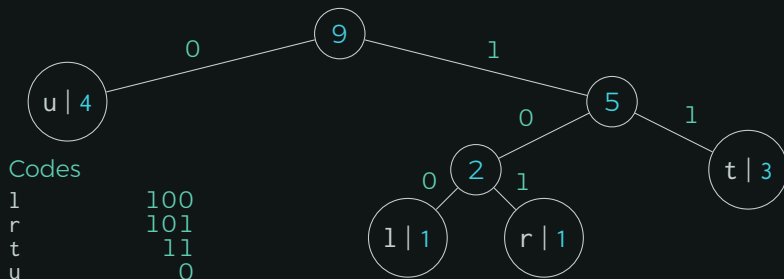


# Codage

Problème résolu

Le codage arithmétique sera détaillé plus loin.

Codage d'HUFFMAN pour turlututu

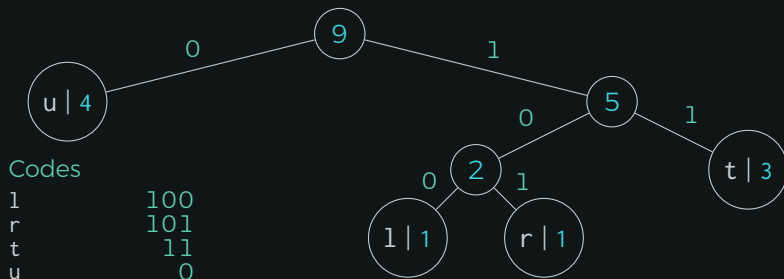


# Codage

Problème résolu

Le codage arithmétique sera détaillé plus loin.

Codage d'HUFFMAN pour turlututu



# Codage

## Inefficacité de HUFFMAN – pourquoi

En pratique, efficacité de 30 %

# Codage

## Inefficacité de HUFFMAN – pourquoi

En pratique, efficacité de 30 %

En appliquant directement SHANNON aux fréquences d'apparition, on commet des erreurs

# Codage

## Inefficacité de HUFFMAN – pourquoi

En pratique, efficacité de 30 %

En appliquant directement SHANNON aux fréquences d'apparition, on commet des erreurs

- ▶ un symbole n'est pas indépendant des précédents

# Codage

## Inefficacité de HUFFMAN – pourquoi

En pratique, efficacité de 30 %

En appliquant directement SHANNON aux fréquences d'apparition, on commet des erreurs

- ▶ un symbole n'est pas indépendant des précédents
- ▶ par exemple, en Français,  $q \rightarrow u$  est plus fréquent que  $q \rightarrow z$

# Codage

## Inefficacité de HUFFMAN – pourquoi

En pratique, efficacité de 30 %

En appliquant directement SHANNON aux fréquences d'apparition, on commet des erreurs

- ▶ un symbole n'est pas indépendant des précédents
- ▶ par exemple, en Français,  $q \rightarrow u$  est plus fréquent que  $q \rightarrow z$
- ▶ en quelque sorte, on oublie le caractère lipschitzien de nos données

Il faut donc un modèle