

Why?

- Saves space.
- Reduces I/O.
- Reduces communication load. }
- Cryptography.

Saves Time

1. © Klein S. T. and Wiseman Y.

Data Compression

```

    graph LR
      A[DATA] --> B[Compression Algorithm]
      B --> C[DATA]
      D[DATA] --> E[Decompression Algorithm]
      E --> F[DATA]
  
```

1. © Klein S. T. and Wiseman Y.

Example: FAX

- Some fax machines use a version of Run-Length coding as a part of the compression procedure.
- Each row is encoded by a list of integers telling the number of white and black pixels in each group. e.g. 1,8,17,3,...
- Each row starts with a dummy white pixel.
- All rows have the same length.

1. © Klein S. T. and Wiseman Y.

How?

- Data Redundancy
 - ▶ Example: In an English text file, ftp sends just 7 bits since 8th bit is always 0.
 - ▶ Yet another example: A string with many 0s and very few 1s - 1 is an innovation. 0 doesn't tell much.

1. © Klein S. T. and Wiseman Y.

Classification I

- Statistic codes
 - ▶ Code length depends on probability of item.
- Dictionary codes
 - ▶ The compressed file consists of pointers to a collection of strings.

1. © Klein S. T. and Wiseman Y.

What?

- TRESOR DE LA LANGUE FRANCAISE (TLF) - 120,000,000 words.
- RESPONSA RETRIEVAL PROJECT (RRP) - 60,000,000 words.

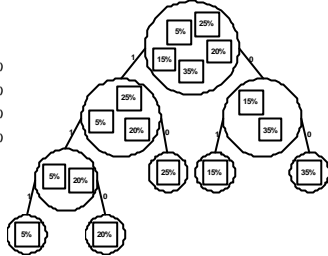
Category	original (MBs)	compressed (MBs)
Text	~700	~250
Dictionary	~10	~10
Concordance	~400	~250
Bit Maps	~700	~100

1. © Klein S. T. and Wiseman Y.

Example

A - 5%, B - 15%, C - 20%, D - 25%, E - 35%.

111 A (5%)
 01 B (15%)
 110 C (20%)
 10 D (25%)
 00 E (35%)



1. © Klein S. T. and Wiseman Y.

Shannon-Fano coding

- Divide the set of symbols into two equal or almost equal subsets based on the probability of occurrence of characters in each subset.
- One set is assigned 0.
- The other set is assigned 1.
- Repeat the procedure until all subsets have a single element.

1. © Klein S. T. and Wiseman Y.

Shannon-Fano is not optimal

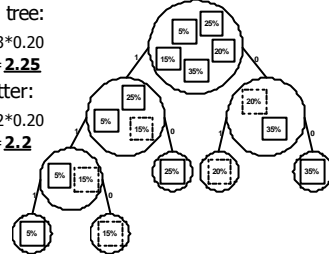
● This tree will give better compression:

● Shannon-Fano tree:

$$3 * 0.05 + 2 * 0.15 + 3 * 0.20 + 2 * 0.25 + 2 * 0.35 = \mathbf{2.25}$$

● This tree is better:

$$3 * 0.05 + 3 * 0.15 + 2 * 0.20 + 2 * 0.25 + 2 * 0.35 = \mathbf{2.2}$$



1. © Klein S. T. and Wiseman Y.