

# Types of Compression

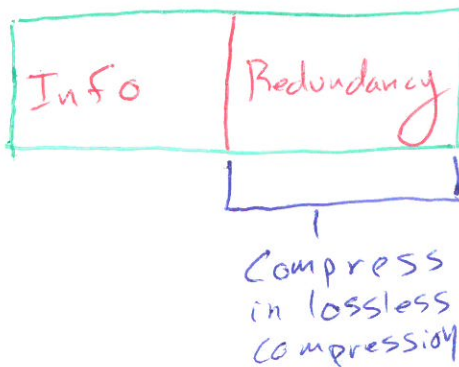
## Lossy v. Lossless

Lossy means we can lose information  
(media (images, video, sound))

Lossless means you get back exactly  
what you started with.

**Theorem #1** There is no "perfect"  
compression algorithm that can ~~and~~ take  
any data and make it smaller (always).

Information Theory:



Compression as a pipeline

Plain Data



Domain Specific  
Redundancy Removal



General Techniques  
(LZ & a.k.a zip)



Compressed data



General Decompress



Domain Specific  
Decompress



Plain Data

# Huffman I

A A-Z 5 bits  
 B  
 C  
 D  
 ...  
 z

An improvement.

EO  
 T 10  
 M 110  
 A 1110  
 R 11110  
 S 111110

ATMARS  
 11101011011011101110

BAD Way To assign few bytes to frequent letters & many bytes to infrequent letters.

ATMARS

EO  
 T 1  
 M 01  
 A 10  
 R 11  
 S 00

BAD ENCODING: 10101101100

7?  
 First half of A  
 First half of R?

# Huffman II

AAAAAABBBBCCCD 2 bits each

Histogram  
A 9 - 2 = 18  
B 5 - 2 = 10  
C 3 - 2 = 6  
D 2 - 2 = 4  
38

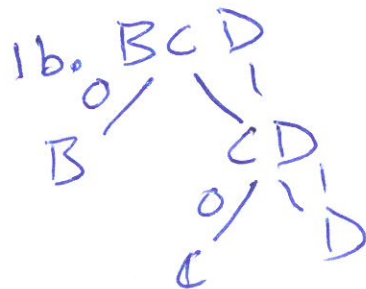
1. Find the two least frequent letters to assign them to 0 and 1.



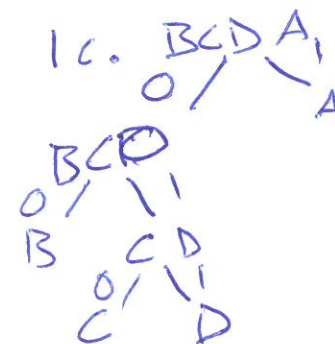
2. Reinsert combined letters into histogram

A 9  
B 5  
CD 5

3. Repeat.



2b. BCD 10  
A 9



2c. Done

A: 1 | 9 - 1 9  
B: 00 | 5 - 2 10  
C: 000 | 3 - 3 9  
D: 011 | 2 - 3 6  
34

1000000000000100100101111  
A B B B B B C C C D D

Compression Ratio

$$\frac{34}{38} = \frac{17}{19}$$

LZ77, LZ78, (zip), etc.

First binary called PKZIP  
Successor gzip

ABABACAB

Triples

- (- Look back
- Read Forward
- Next Byte)

(0, 0, A) A

(0, 0, B) AB

(2, B, C) AB, ABAC

(4, 2, -) ABABACAB