

# Huffman: Lossless Compression

We saw:  $E L_H(X) < H(X) + 1 \text{ bits}$

↙ One bit gap (overhead) ↘

- Encode over blocks:
1. The overhead becomes negligible.
  2. Take advantage of correlation.

(the angry bird) e.g. Code:  $qa \rightarrow 00100$   
 $ab \rightarrow 010$   
 $\vdots$

Example without correlation:

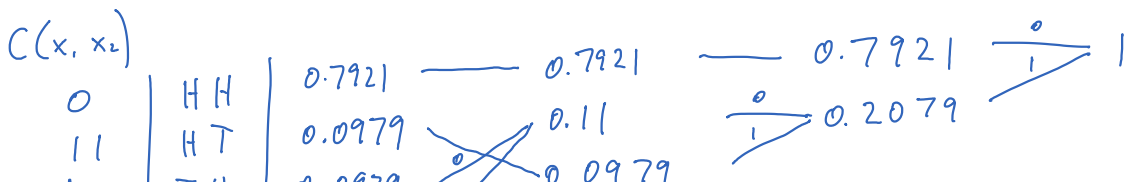
$\Omega = \{H, T\}, \text{ Prob}(X=H) = 0.89, \text{ Prob}(X=T) = 0.11$

$H(X) = \frac{1}{2} \text{ bit}$

Huffman code: Code:  $H \rightarrow 0, T \rightarrow 1$   $E L_u(X) = 1$

Encode over pairs:  $(X_1, X_2) \in \Omega, \Omega = \{HH, HT, TH, TT\}$

$p(HH) = (0.89)^2 = 0.7921$   
 $p(HT) = 0.89 \cdot 0.11 = 0.0979$  *by independence*  
 $p(TH) = 0.11 \cdot 0.89 = 0.0979$   
 $p(TT) = (0.11)^2 = 0.0121$



0	HH	0.11	$\begin{matrix} 0 \\ 1 \end{matrix} \rightarrow 0.2079$
11	HT	0.0979	
100	TH	0.0979	
101	TT	0.0121	

$$E L_H(X) = 0.7921 \cdot 1 + 0.0979 \cdot 2 + 0.11 \cdot 3 \approx 1.3 \text{ bits}$$

$$\text{bits/symbol} \approx \frac{1.3}{2} = 0.65 \text{ bits}$$

Bounds for longer blocks:

Let  $L_n^*$  be the normalized expected length of the Huffman code for blocklength  $n$ .

$$L_n^* = \frac{1}{n} E L_H(X_1, X_2, \dots, X_n)$$

$$n H(X) = H(X_1, X_2, \dots, X_n) \leq n L_n^* < H(X_1, X_2, \dots, X_n) + 1 \text{ bit} = n H(X) + 1 \text{ bit}$$

↑  
Assuming independent and identically distributed

$$\Rightarrow H(X) \leq L_n^* < H(X) + \frac{1}{n} \text{ bits.}$$

In general, "entropy rate" is the fundamental compression limit.

## JPEG:

Color is 3-dimensional as perceived by the human eye.

red, green, blue  $\leftarrow$  basis  
 red, yellow, blue  $\leftarrow$

$Y, C_B, C_R$        $Y$  is intensity

$Y, C_B, C_R$

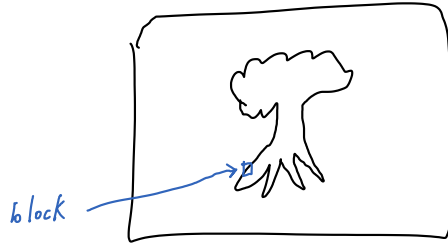
$Y$  is intensity

$C_B$ , color

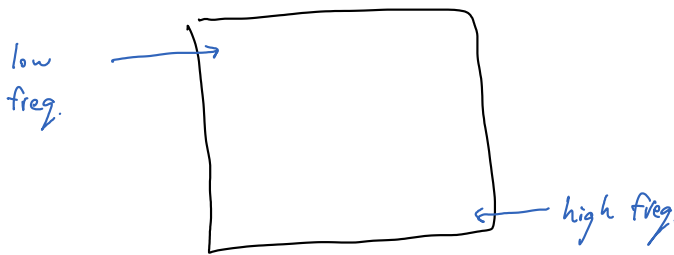
$C_R$ , color

} Downsample by 2

Split in  $8 \times 8$  pixel blocks:



DCT: (think of as DFT)



Uniform Quantization:

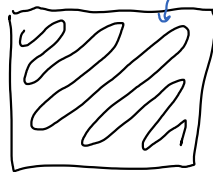
- Scale each entry (by different values)
- Round to nearest integer

16	18	
18		50
	50	100

← Matrix of scaling constants  
Fixed part of the protocol.

Last steps: Lossless

Quantized DCT



Vectorize in this order

- Run-length code:

Example: 11111 0000000111111 0000

(5,1), (7,0), (6,1), (4,0)

5, 7, 6, 4

- Huffman Code

Expect many zeros in a row

Error Detection and Correction: Add redundancy in a smart way.

- Reed-Solomon:

1960

Many applications - in much of today's technology

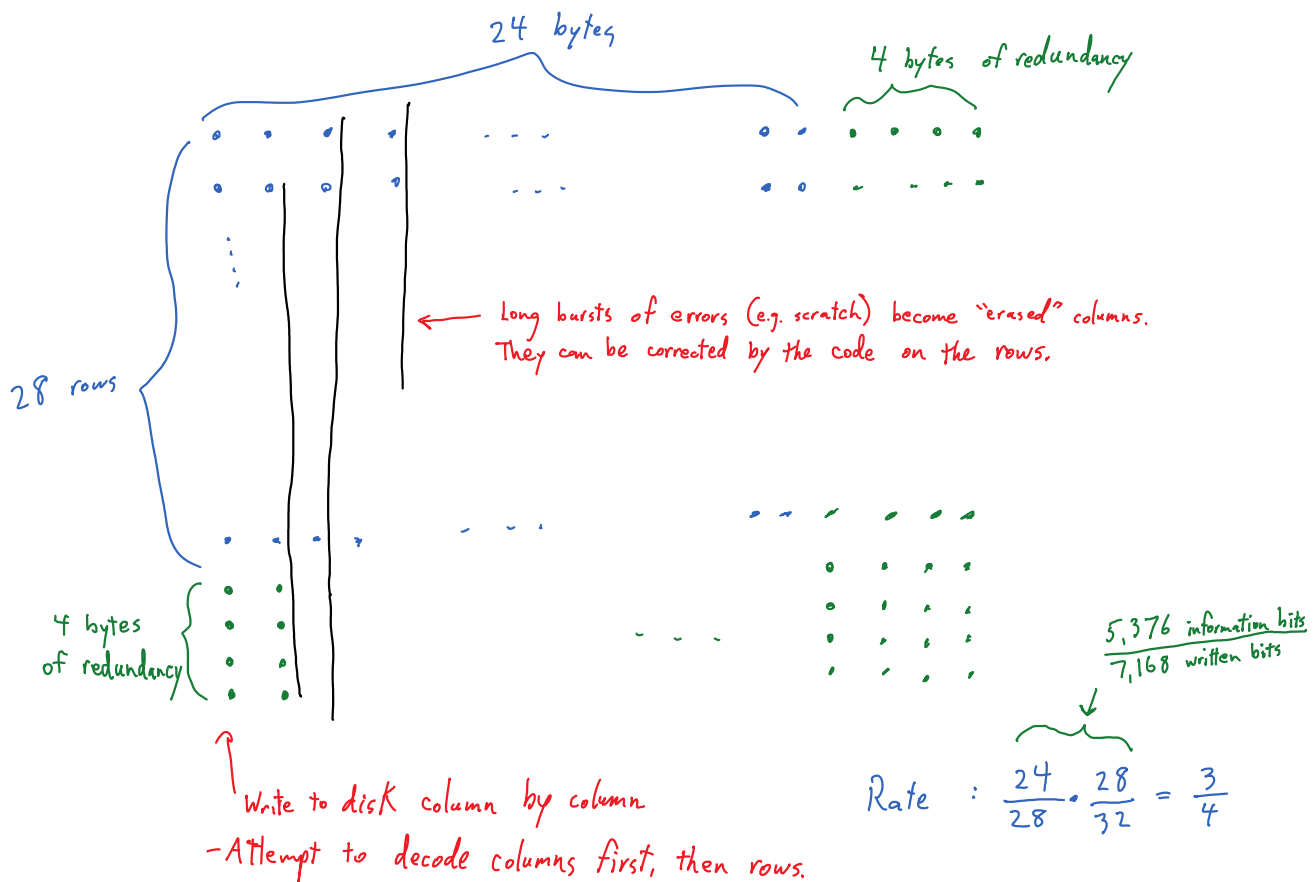
- Specifics:
- Take  $n$  symbols of information
  - Add  $t$  symbols of redundancy ← (used finite-field arithmetic)
  - Correct errors and erasures.

$$2 \cdot N_{\text{errors}} + N_{\text{erasures}} \leq t$$

- If it cannot correct, it usually declares failure.
- Decoder can be understood with the Fourier Transform

As Seen on CD!

- symbols are bytes (8 bits)
- 2 layers with interleaving.



Probability of error:

- noise model is important for analysis
- good for bursty errors
- not so good for scattered noise (15 well placed bit-errors breaks the error correction)

