# Information Theory

↳ Information Theory : movement & transformations are constrained by mathematical laws.
  ↳ Complexity of patterns is the length of the shortest program to generate this pattern
  ↳ Ultimate Data Compression = Entropy ($H$)
  ↳ Ultimate Rate of Reliable Communication = Channel Capacity ($C$)

## Foundations

Random Variables : take on values determined by their probability distributions

Probability has two meanings :
  ↳ ① RELATIVE FREQUENCY : (1) sample random variable multiple times
    (frequentist/operationalist view) (2) take fraction each outcome occurs
  ↳ ② DEGREE - OF - BELIEF : plausability of proposition - no experiment - likelihood that a particular state might occur & event outcome can only be determined once.

Less probable an event is, more information is gained by seeing occurence.

PRODUCT RULE $\begin{cases} p(A, B) = \text{joint prob of A and B} = p(A|B)\, p(B) = p(B|A)\, p(A) \\ \text{if } A \text{ and } B \text{ are independent, } P(A|B) = P(A) \text{ & } P(B|A) = p(B) \end{cases}$

SUM RULE $\left[ p(A) = \sum_B p(A, B) = \sum_B p(A|B)\, p(B) \right.$

BAYES THEOREM $\left[ P(B|A) = \dfrac{P(A|B)\, P(B)}{P(A)} \right.$

| Prior Prob | = prob of observing something before any data is collected
| Posterior Prob | = computed after observation of data .

## Entropy

both measured in bits

$\begin{cases} \text{Information Measure}(I) = \boxed{\log_2 p} \\ \text{Entropy } (H) = -I \end{cases}$ ← $p$ is probability of event occurring
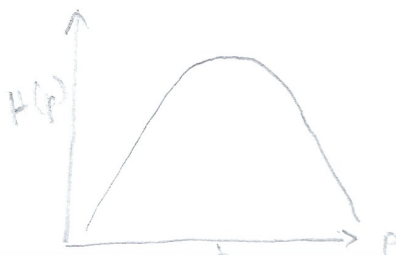
↳ uncertainty / disorder

↳ log as we want information to be additive. Joint probs (independent) are multiplicative, so logs makes them additive

↳ $= 1.443 \log_e p = 3.322 \log_{10} p$

## Asymptotic Equipartition Theorem

Given uniform probabilities, and question of choosing one of N items, entropy $= \log_2 N$

$$H = -I = - \sum_i p_i \log (p_i)$$

② Notation

$X, Y$ - random variables
$\quad \hookrightarrow x \in \{a_1, a_2, \ldots, a_J\} = \mathcal{A}$
$\quad\quad y \in \{b_1, b_2, \ldots, b_k\} = \mathcal{B}$
$\quad\quad\underbrace{\quad}$
$\quad\quad$ instances

Ensemble is a random variable. Joint Ensemble is an ensemble whose outcomes are ordered pairs $x, y$.
$\quad\quad\quad\quad\quad\quad \hookrightarrow XY$ defines $P(x, y)$ over all outcomes

MARGINAL PROB ABILITIES
$$p(x = a_i) = \sum_y p(x = a_i, y)$$
$$p(x) = \sum_y p(x, y)$$

CONDITIONAL PROBABILITIES
$$p(x = a_i \mid y = b_j) = \frac{p(x = a_i, y = b_j)}{p(y = b_j)}$$
$$p(x \mid y) = \frac{p(x, y)}{p(y)}$$

JOINT ENTROPY
$$\boxed{H(X, Y) = \sum_{x, y} p(x, y) \log \frac{1}{p(x, y)}}$$

If $X$ and $Y$ are independent random variables
$$H(X, Y) = H(X) + H(Y)$$

Conditional Entropy
$$H(X \mid y = b_j) = \sum_x p(x \mid y = b_j) \log \frac{1}{p(x \mid y = b_j)}$$

$$H(X \mid Y) = \sum_y p(y) \left[ \sum_x p(x \mid y) \log \frac{1}{p(x \mid y)} \right]$$

Chain Rule for Entropy

JOINT ENTROPY $(H(X, Y))$ = MARGINAL ENTROPY $(H(X))$ + CONDITIONAL ENTROPY $(H(Y \mid X))$

$$H(X, Y) = H(X) + H(Y \mid X) = H(Y) + H(X \mid Y)$$

$H(X)$ maximised when $p_i = 1/N$
$$H(X) = -\sum_i p_i \log_2 p_i = -\sum_1^N \frac{1}{N} \log_2 \frac{1}{N} = \log_2 N$$
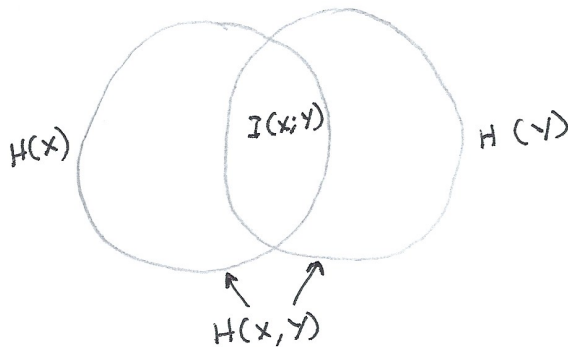
③

$$I(x;y) = \text{Mutual Information between } x \text{ and } Y = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} \geqslant 0$$

↳ if perfectly correlated, $I(x;Y) = H(x) = H(Y)$

$$I(x;Y) = H(x) - H(x|Y)$$
$$I(x;Y) = H(Y) - H(Y|x) = I(Y;X)$$
$$I(x;Y) = H(x) + H(Y) - H(x,Y)$$



H(X)    I(x;Y)    H(Y)

H(x,Y)

## Cross Entropy

↳ two different distributions ($p(x)$ and $q(x)$) over the same set of outcomes for random variable $X$.

$$H(p,q) = -\sum_x p(x) \log(q(x))$$

Used in coding theory to get cost of wrong representation. i.e. if coding scheme designed under assumption $q(x)$, $H(p,q)$ gets length of codewords on average given actual $p(x)$.

Assymetric and minimised if $p(x) = q(x)$

## Distance (between two Random Variables)

$$D(X,Y) = H(x,Y) - I(x;Y)$$

↳ follows axioms → ① $P(X,Y) \geqslant 0$
② $P(X,X) = 0$
③ $D(X,Y) = D(Y,X)$
④ $D(X,Z) \leqslant D(X,Y) + D(Y,Z)$ — triangle inequality

→ measure of inefficiency

## Kullback-Leibler Distance (Relative Entropy) — information for discrimination

$\geqslant 0$, $p(x) = q(x) \Rightarrow D_{KL} = 0$

$$D_{KL}(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$
$$= H(p,q) - H(p)$$

Given optimal code for distribution $p(x)$ (needs $H(p)$ bits) then n(additional bits) needed if described $p(x)$ using optimal code as $q(x) = D_{KL}(p||q)$

## Fano's Inequality

Relates $P_e$ (prob of error) in guessing X from knowledge of Y to conditional entropy $H(X|Y)$ when $|A|$ is number of possible outcomes

$$P_e \geq \frac{H(X|Y) - 1}{\log |A|}$$

↳ length of symbol alphabet

## Data Processing Inequality

If X, Y, Z form Markov Chain (Z depends on Y and Y depends on X)
↳ X → Y → Z

$$\Rightarrow I(X; Y) \geq I(X; Z)$$

## SOURCE CODING THEOREM

<u>Markov Process</u>: source of symbols - letters emitted with known probabilities for each state.

<u>Entropy</u> of Markov Process expressed as bits per symbol — if emitted at known rate can get into bits per second. (baud rate)
↳ calculate entropy for each state then take weighted average given by occupancy probabilities $P_i$:

$$H = \sum_i P_i H_i = - \sum_i P_i \sum_j p_i(j) \log (p_i(j))$$

## Fixed Length Codes

Given set of N symbols, with entropy H, need fixed length block,
$R = \lceil \log_2 (N) \rceil + 1$. Code rate is R bits per symbol. $(H \leq R)$

Efficiency $\eta = \frac{H}{R}$

↳ aim is H = R

## Variable Length Codes

Can achieve more compressed using variable-length codes.
Features of codes:
  ↳ ① Uniquely Decodable: codes cannot be multiple combinations
  ↳ ② Prefix: No codeword is the prefix of a longer codeword.

↗ Shannon Source Coding Theorem: $R \geq H$

<span style="margin-left:0">noiseless coding theorem</span>

"For a discrete source with entropy H, for any $\varepsilon > 0$, possible to encode symbols into uniquely decodable code at average rate R s.t: $R = H + \varepsilon$ as an asymptotic limit"

# Huffman Codes

Algorithm to get ~~any given~~ an optimal prefix code for given probability distribution of symbols → illustration of Shannon Source - Coding Theorem

Ⓐ Assign bit in reverse sequence corresponding to increasing symbol prob

Ⓑ More probable symbols encoded with shorter codewords.

1. Find two symbols with lowest probs + assign bit to distinguish them
   ↳ defines branch in binary tree
2. Combine these two into virtual node with prob as sum
3. Start again and go until one symbol node

N.B. no unique Huffman code for a symbol alphabet ⟹ Huffman code is as efficient as possible

## Kraft-McMillan Inequality

Any instantaneous code (with prefix property) must have conditions :
→ necessary but not sufficient
   ↳ if N codewords have length $c_1 \leq c_2 \leq .. \leq c_n$ then :

$$\sum_{i=1}^{N} \frac{1}{2^{c_i}} \leq 1$$

## DISCRETE CHANNEL CAPACITY

Input alphabet = $X = \{x_1, ..., x_n\}$
Output alphabet = $Y = \{y_1, ..., y_m\}$

$\nearrow \sum_{k=1}^{K} p(y_k | x_j) = 1$

Channel represented as transition probabilities $p(y_k | x_j)$ — these can form the channel matrix

$$p(x_j, y_k) = p(y_k | x_j) \, p(x_j) \quad ] \text{joint prob. dist.}$$

$$p(y_k) = \sum_{j=1}^{J} p(x_j, y_k) = \sum_{j=1}^{J} p(y_k | x_j) p(x_j) \quad ] \begin{array}{l}\text{marginal prob} \\ \text{for output symbol} \\ y_k\end{array}$$

Av prob of symbol error $\left[ P_e = \sum_{j=1}^{J} \sum_{\substack{n=1, n \neq j}}^{K} p(y_k | x_j) p(x_j) \right.$
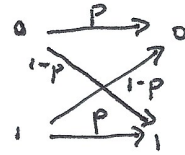   ↳ $1-P_e$ is prob of correction reception

## Binary Symmetric Channel

Two input and output symbols $\{0,1\}$ with channel matrix $\begin{pmatrix} 1-p & p \\ p & 1-p \end{pmatrix}$

$$H(X) = 1 \text{ bit}, \quad H(Y) = 1 \text{ bit}$$

$$H(X|Y) = -\sum_{x,y} p(x,y) \log p(x|y)$$

$$= -p \log(p) - (1-p) \log(1-p)$$

$$I(X;Y) = H(X) - H(X|Y)$$
$$= 1 + p \log p + (1-p) \log(1-p)$$

(binary)
assuming symmetric channel
↳ flips equally likely

$$\text{Channel Capacity (C)} = \max_{\{p(x_i)\}} I(X;Y) \longrightarrow \text{equiprobable symbols maximise } I(X;Y)$$

$I(X;Y)$ optimal with $p=0$ or $1 \Rightarrow$ channel capacity is <u>1 bit</u> per transmitted bit

Maximal for equiprobable case for probability distribution. (mutual information)

## Overcoming noise

① Repetition Codes

↳ if we transmit every symbol $N = 2m+1$ times + do majority voting.
↳ require if $m+1$ bits in error

$$P_e = \sum_{i=m+1}^{2m+1} \binom{2m+1}{i} p^i (1-p)^{2m+1-i}$$

## Shannon Second Theorem — Channel Coding Theorem

"For channel of capacity $C$ and symbol source of entropy $H$ provided $H \leq C$ ∃ coding scheme st source is reliably transmitted through channel with residual error rate $< \varepsilon$ (arbitrarily small)" — existence proof, no algorithm

$\left(b_1, b_2, \ldots\right.$
$= b_1, b_2,$
$= b_1, \widehat{b_2}, \ldots )$

**Ex** Given 7 bit blocks, no flip and one flip in each is equiprobable

$$C_S = \max_{p(x_j)} I(X;Y) = \max \left( H(Y) - \cancel{H(XY)} H(Y|X) \right)$$

$H(Y) = 7$ because $2^7$ equiprobable symbols

$$p(y_k|x_j) = \frac{1}{8}$$

~~$\frac{1}{8} \sum \sum 8 \times$~~

$$⑦ \frac{1}{7}\left( 7 - \sum_j \sum_k p(y_k|x_j) \log \frac{1}{p(y_k|x_j)} p(x_j)\right)$$

$$= \frac{1}{7}\left( 7 + \sum_j 8 \left(\frac{1}{8} \log \frac{1}{8}\right) \frac{1}{N}\right)$$

$$= \frac{4}{7} \text{ bits per bit}$$

## Syndromes

Given source entropy $H \le C = \frac{4}{7}$ bit per bit. Construct new 7-bit codewords with 4 bits of symbol encoding and 3 error correction bits.

Calculate 3 bits as XOR of 3 of the four bits.

$b_1, b_2, b_4$ are error correction bits:
$$b_4 = b_5 \oplus b_6 \oplus b_7$$
$$b_2 = b_3 \oplus b_6 \oplus b_7$$
$$b_1 = b_3 \oplus b_5 \oplus b_6$$

On receipt, calculate syndromes:
$$S_4 = b_4 \oplus b_5 \oplus b_6 \oplus b_7$$
$$S_2 = b_2 \oplus b_3 \oplus b_6 \oplus b_7$$
$$S_1 = b_1 \oplus b_3 \oplus b_5 \oplus b_7$$
$H = \Big\}$ if all 0, no error

Hence, source entropy goes to $\frac{4}{7}$ bit per bit of data ⇒ satisfies requirement that $H \le C$

Hamming codes perfect as use $m$ bits to correct $2^m - 1$ error patterns and transmit $2^m - 1 - m$ useful bits.

$$P_e = \sum_{i=2}^{7} \binom{7}{i} p^i (1-p)^{7-i}$$

## Information in Vector Spaces

**Example**  $\underline{u}$ is vector of data or samples and $\{e_1, \ldots, e_n\}$ are basis vectors of orthonormal system. Projection of $\underline{u}$ into that space gets coefficients $a_i$ as inner product of $\underline{u}$ with each basis vector.

$$u = \sum_{i=1}^{n} a_i e_i = \sum_{i=1}^{n} \langle u, e_i \rangle e_i$$

↑ inner product

**Example**  $f(x)$ represented as linear combination of functions $\psi_i(x)$ s.t.

$$f(x) = \sum_{i=1}^{n} a_i \psi_i(x)$$

where:
$$a_i = \langle f, \psi_i \rangle = \int_{-\infty}^{\infty} f(x)\, \psi_i(x)\, dx$$

$\hookrightarrow$ could be Fourier coefficient   $\hookrightarrow$ inner product

Frequently is the Fourier Transform

**Propagation Function**: inner product between vector of inputs and vector of learned weights.

## Definitions

$\boxed{V_1, \ldots, V_n}$

Vector Space $V$ is a linear combination of vectors in $V$ if $\exists$ coefficients to make that be true

vectors are a space basis for vectors

space → only if linearly independent, $n$ is dimension of the space.

$\boxed{\text{span}}\{v_1, \ldots v_n\} = \{u \in V : u \text{ is linear comb of } v_1, \ldots, v_n\}$

↳ everything that can be represented by linear combination of span vectors

Subset $W \subset V$ is linear subspace of $V \Rightarrow$ involves projecting onto that subspace of vectors - dimensionality reduction

For $V$, $v_1, \ldots, v_n \in V$ are $\boxed{\text{linearly independent}}$ if for scalars $a_1, a_2, \ldots$

$$a_1 v_1 + a_2 v_2 + \ldots + a_n v_n = 0 \text{ ; then } a_1 = a_2 = \ldots = a_n = 0$$

## Inner Product

$\langle u, v \rangle$ is scalar value satisfying

↳① $\forall v \in V, \langle v, v \rangle \in \mathbb{R} \geq 0$

↳② $\forall v \in V, \langle v, v \rangle = 0 \Leftrightarrow v = 0$

↳③ $\forall u, v, w \in V$, scalars $a, b$; $\langle au + bv, w \rangle = a \langle u, w \rangle + b \langle v, w \rangle$

↳④ $\forall u, v \in V \Rightarrow \langle u, v \rangle = \overline{\langle v, u \rangle}$

    ↳ if vectors are real,

    $\langle u, v \rangle = \langle v, u \rangle$

    $\overline{a}$ is complex conjugate

↳ Properties

    ↳① $\forall v \in V$ and scalar $a \Rightarrow \langle av, av \rangle = |a^2| \langle v, v \rangle$

    ↳② $\forall v \in V$ , $\langle 0, v \rangle = 0$

    ↳③ $\forall v \in V$ $u_1, u_2, \ldots, u_n \in V$ scalars $a_1, \ldots, a_n \Rightarrow$

$$\left\langle \sum_{j=1}^{n} a_j u_j , v \right\rangle = \sum_{j=1}^{n} a_j \langle u_j, v \rangle$$

$$\left\langle v, \sum_{j=1}^{n} a_j u_j \right\rangle = \sum_{j=1}^{n} \overline{a_j} \langle v, u_j \rangle$$

↳ Examples

    ↳ Euclidean Space $\mathbb{R}^n$ : $\langle x, y \rangle = \sum_{i=1}^{n} x_i y_i$

    ↳ Complex Vectors $\mathbb{C}^n$ : $\langle x, y \rangle = \sum_{i=1}^{n} x_i \overline{y_i}$

## Norm

Norm on Vector Space $V$: $V \to \mathbb{R}_+$ has properties:

    ↳① $\forall v \in V, \|v\| \geq 0$

    ↳② $\|v\| = 0 \Leftrightarrow v = 0$

    ↳③ $\forall v \in V$, scalar $a$, $\|av\| = |a| \|v\|$

    ↳④ $\forall u, v \in V \|u + v\| = \|u\| + \|v\|$

] generalisation of the notion of the distance between two vectors

# Orthogonal & Orthonormal Systems

$u, v \in V$ are orthogonal ($u \perp v$) if $\langle u, v \rangle = 0$. Sequence of vectors (finite or infinite) $\{v_i\}$ is orthogonal system if:

↳① $v_i \neq \bar{0} \quad \forall v_i$

↳② $v_i \perp v_j \quad \forall i \neq j$

Orthogonal system is orthonormal if $\|v_j\| = 1 \quad \forall i$

↳ call unit vectors ($e_i$)

If $\{e_1, e_2, \ldots, e_n\}$ is an orthonormal system in $V$.

$$u = \sum_{i=1}^{\hat{n}} a_i e_i \Rightarrow a_i = \langle u, e_i \rangle$$

/ in orthonormal system, expansion coefficients are same as projection coefficients

$$u = \sum_{j=1}^{\hat{n}} a_i e_i = \sum_{i=1}^{\hat{n}} \langle u, e_i \rangle e_i$$

## Infinite Orthonormal Systems : $V$ with $\dim(V) = \infty$

Let $\{v_1, v_2, \ldots\}$ be infinite sequence of vectors in normed vector space $V$.
Let $\{a_1, a_2, \ldots\}$ be sequence of scalars.

$$\sum_{n=1}^{\tilde{\infty}} a_n v_n \quad \text{converges in norm to } w \in V \text{ if}$$

$$\lim_{m \to \infty} \left\| w - \sum_{n=1}^{m} a_n v_n \right\| = 0$$

$w$ would be exactly represented by a linear combination of vectors $\{v_i\}$ in space $V$ in limit that we could use all of them. This property of infinite orthonormal system in inner product space is a closure

## Karhunen-Loeve Transform : numerical method for constructing orthonormal system such that any set of vectors can be represented with best possible accuracy using any specified finite number of terms

also Hotelling Transform, Dimensionality Reduction or Principal Components Analysis

## FOURIER REPRESENTATIONS

Decomposition of functions into superpositions of elementary sinusoidal functions — consider information in piecewise continuous functions $f$ for convenience over interval $[-\pi, \pi]$ by projecting to vector space of:

$$\left\{ \frac{1}{\sqrt{2}}, \sin(x), \cos(x), \sin(2x), \cos(2x), \ldots \right\}$$

| average value.

↳ Fourier components

$$F(f) = \frac{a_0}{2} + \sum_{n=1}^{\infty} \left[ a_n \cos(nx) + b_n \sin(nx) \right]$$

$$a_n = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos(nx) \, dx$$

$$b_n = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin(nx) \, dx$$

if even symmetry, only cosine terms, if odd symmetry, then only sine terms

for $2\pi$-periodic function, ie $g(x+2\pi) = g(x)$

If f even, $a_n = \frac{2}{\pi} \int_0^{\pi} f(x) \cos(nx) \, dx$

$b_n = 0$

If f odd, $a_n = 0$

$b_n = \frac{2}{\pi} \int_0^{\pi} f(x) \sin(nx) \, dx$

even: $f(-x) = f(x)$
odd: $f(-x) = -f(x)$

① if f,g are even then fg even
② if f,g are odd then fg even
③ if f even and g odd then fg odd
④ If g odd, for any $h > 0$, $\int_{-h}^{h} g(x) \, dx = 0$

N.B. If derivative is infinite, Fourier coefficient dies very quickly and convergence is faster

⑤ If g even, for any $h > 0$, $\int_{-h}^{h} g(x) \, dx = 2 \int_0^{h} g(x) \, dx$

### Complex Fourier Series

$$\{1, e^{ix}, e^{-ix}, \dots\}$$

$$\sum_{n=-\infty}^{\infty} c_n e^{inx}$$

where $c_n = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-inx} \, dx$

N.B. $\langle f, g \rangle = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) \overline{g(x)} \, dx$

$$c_n = \frac{a_n - ib_n}{2}$$

$$c_{-n} = \frac{a_n + ib_n}{2}$$

$$\left( c_0 = \frac{a_0}{2} \right)$$

### SPECTRAL PROPERTIES OF CONTINUOUS-TIME CHANNELS

Information bands assigned spectral band - carrier signal modulated inside in a parameter (frequency of sine wave, amplitude, phase). Results in a complex $f(t)$
①         ③         ⑤

$$f(t) = \sum_n c_n e^{iwnt}$$     (by Fourier Analysis)

Channels are linearly time-invariant systems whose eigenfunctions are complex exponentials

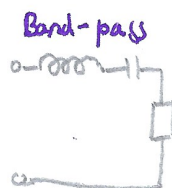Linear time-invariant systems obey superposition and can be described by linear $h(t)$

Complex exp is never changed, only multiplied by complex constant $\alpha$

Can use amplitude and phase changing $\alpha_n$ for $\omega_n$ can incorporate into $f(t)$ to understand how $f(t)$ is affected by transmission through channel, hence:
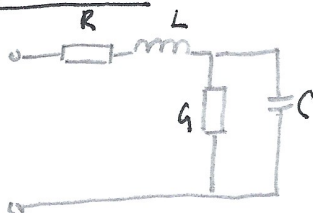
$$f(t) = \sum_n \alpha_n c_n e^{i\omega_n t}$$

## Filters

↳ Resistors - just have constant impedance : $Z = R$

↳ Capacitors - low impedance at high frequencies and high impedance at low frequencies. $Z(\omega) = \dfrac{1}{i\omega C}$

↳ Inductors : $Z(\omega) = i\omega(L)$ — in Henrys

Low-pass filter          High-pass          Band-pass          Band-reject



## Coax Cable



↳ non-zero series resistance $R$
↳ non-zero series inductance $L$
↳ non-zero shunt conductance $G$
↳ non-zero capacitance $C$

Effectively like a low pass filter - restricted to finite bandwidth $W$. Frequency components $> W \approx 1/RC$ attenuated by signal pathway

But add wideband noise to signal during transmission ⇒ normally shot noise
    ↳ if uniform, called white noise

## CONTINUOUS INFORMATION

Given value $X$ has probability density $p(x)$ with $\int_{-\infty}^{\infty} p(x)\,dx = 1$, differential entropy:

$$h(X) = \int_{-\infty}^{\infty} p(x) \log_2\left(\frac{1}{p(x)}\right) dx$$

$$h(X,Y) = \int_{-\infty}^{\infty}\int_{-\infty}^{\infty} p(x,y) \log_2\left(\frac{1}{p(x,y)}\right) dx\,dy$$

$$h(X|Y) = \int_{-\infty}^{\infty}\int_{-\infty}^{\infty} p(x,y) \log_2\left(\frac{p(y)}{p(x,y)}\right) dx\,dy$$

$$h(X,Y) \leq h(X) + h(Y)$$

$$i(X;Y) = \int_{-\infty}^{\infty}\int_{-\infty}^{\infty} p(x,y) \log_2\left(\frac{p(x|y)}{p(x)}\right) dx\,dy = h(Y) - h(Y|X)$$

$$C = \max_{p(x)} i(X;Y)$$

Entropy maximised with equiprobable symbols — in continuous case if single value $x$ is limited to range $v$, prob density $p(x) = 1/v$

Variance of continuous random variable $x =$ power of sound signal $=$ differential entropy $h(X)$. Can be proven that $p(x)$ of excursions around mean $\mu$:

$$p(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-(x-\mu)^2 / 2\sigma^2\right)$$

$$h(X) = \frac{1}{2} \log_2 (2\pi e \sigma^2) \quad \text{(where } p(x) \text{ is maximised)}$$

$\hookrightarrow$ white noise since power spectrum is flat

## Channel with injected Gaussian noise

Limited Spectral Bandwidth $W \Rightarrow$ signal + noise are lowpass — no frequency components higher than $W$

Power Spectral Density $N_0$ is white noise

$\hookrightarrow$ noise power $= N_0 W$
$\hookrightarrow$ noise variance $= N_0 W$

Given noise $N$ independent of input signal $X$, $h(Y|X) = h(N)$

$$h(Y|X) = h(N) = \frac{1}{2} \log_2 (2\pi e \sigma^2) = \frac{1}{2} \log_2 (2\pi e N_0 W)$$

$\hookrightarrow$ channel output $Y = X + N$ has variance $P + N_0 W$

$$i(X;Y) = \frac{1}{2} \log_2 \left(1 + \frac{P}{N_0 W}\right)$$

$$C = \frac{1}{2} \log_2 \left(1 + \frac{P}{N_0 W}\right) \to \text{in bits per channel symbol}$$

Shannon's Third Theorem

$$\boxed{C = W \log_2 \left(1 + \frac{P}{N_0 W}\right) \text{ bits/second}}$$

$\downarrow$

Shannon - Hartley Theorem

$\downarrow$

Noisy Channel Coding Theorem

(+ Signal to Noise ratio (SNR) $\Rightarrow$ decibel
$\hookrightarrow 10 \times \log_{10}$ (SNR) — if ratio of power
$\hookrightarrow 20 \times \log_{10}$ (SNR) — if ratio of amplitudes

Increasing $W$ yields monotonic but asymptotically limited improvement in capacity: as $W \to \infty$, $C \to \frac{P}{N_0} \log_2 e$

straight line only in loglog plot



dB

violet noise

blue

white noise

pink noise (flicker noise $-\frac{1}{f}$) has power spectral density $\propto \frac{1}{f}$

red noise

$10^3 \qquad 10^4 \qquad 10^5$ $f$

$$C = \int_{w_1}^{w_2} \log_2 (1 + SNR(w)) \, dw \quad \text{bib/sec}$$

## ENCODINGS BASED ON FOURIER TRANSFORM PROPERTIES

Fourier transform extended to aperiodic functions by making range $(b-a) \to \infty$. Interval between frequency components becomes infinitesimial $\to$ acquires density of reals not integers. Fourier Transform not Fourier Series

$\left[\begin{array}{l} f(x) \text{ must be piecewise continuous} \\ \text{and absolutely integrable} \end{array}\right]$

$F(w)$ of $f(x)$
  ↳ ① $F(w)$ defined $\forall w \in \mathbb{R}$
  ↳ ② $F(w)$ is continuous
  ↳ ③ $\lim_{w \to \pm \infty} F(w) = 0$

For $f: \mathbb{R} \to \mathbb{C}$, $F: \mathbb{R} \to \mathbb{C}$ given by: $\left]\begin{array}{l} \text{local-to-global} \\ \text{and} \\ \text{global-to-local} \end{array}\right.$
  ↳ $F(w) = \mathcal{F}_{[f]}(w) = \frac{1}{2\pi} \int_{-\infty}^{\infty} f(x) e^{-iwx} \, dx$

### Properties

① $\mathcal{F}_{[af + bg]}(w) = a\mathcal{F}_{[f]}(w) + b\mathcal{F}_{[g]}(w)$

② Hermitian Symmetry: $F(-w) = \overline{F(w)}$

③ If $f(x)$ is even function $\Rightarrow F(w)$ is even and purely even valued

④ $f(x)$ odd $\Rightarrow F(w)$ is odd and purely imaginary

frequency shifting $\left[\begin{array}{l} \end{array}\right.$ ⑤ $\mathcal{F}_{[exp(icx) f(x)]}(w) = \mathcal{F}_{[f]}(w - c)$
  ↳ shift in $f(x)$ by $b$ causes $\mathcal{F}_{[F]}$ to be multiplied by $e^{iwb}$
  multiply $f(x)$ by $e^{icx}$ causes shift by $c$ in $\mathcal{F}_{[f]}(w)$

⑥ Modulation Property

$$\mathcal{F}_{[f(x) \cos(cx)]}(w) = \frac{1}{2}(\mathcal{F}_{[f]}(w-c) + \mathcal{F}_{[f]}(w+c))$$

$$\mathcal{F}_{[f(x) \sin(cx)]}(w) = \frac{1}{2i}(\mathcal{F}_{[f]}(w-c) - \mathcal{F}_{[f]}(w+c))$$

### Derivatives

Fourier transforms can be used to solve differential equations

$f$ has derivative $f' \Rightarrow \mathcal{F}_{[f']}(w) = iw \mathcal{F}_{[f]}(w)$

$\left.\begin{array}{l} \end{array}\right]$ can be used in the opposite direction for integration

$$\mathcal{F}_{[f^{(n)}]}(w) = (iw)^n \mathcal{F}_{[f]}(w)$$

Hence is effectively a filtering operation $\Rightarrow$ high-pass filtering
  ↳ can also generalise derivatives to non instagrams

## Inverse Fourier Transform

$$f(x) = \int_{-\infty}^{\infty} \mathcal{F}_{[f]}(\omega) e^{i\omega x} \, dx \qquad \left[\begin{array}{l}\text{again local-to-global and}\\ \text{global-to-local property}\end{array}\right]$$

### Convolution  (commutative)

$f, g : \mathbb{R} \to \mathbb{C}$

$$(f * g)(x) = \int_{-\infty}^{\infty} f(x-y) g(y) \, dy$$

$$\mathcal{F}_{[f*g]}(\omega) = 2\pi \mathcal{F}_{[f]}(\omega) \cdot \mathcal{F}_{[g]}(\omega)$$

N.B. convolution of two Gaussian
is a Gaussian.

Passband (modulated carrier) -- demodulated ---> baseband (encoded
information)

### Amplitude Modulation

Single sideband

$f(t)$ is baseband message
  $\hookrightarrow F(\omega) = \mathcal{FT}\{f(t)\} = \frac{1}{2\pi} \int_{-\infty}^{\infty} f(t) e^{-i\omega t} \, dt$
instead firstly modulate by $e^{ict}$:

$$\mathcal{FT}\{e^{ict} f(t)\} = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{ict} f(t) e^{-i\omega t} \, dt$$

$$= F(\omega - c)$$

Demodulation:

$$f(t) = \left[ e^{ict} f(t) \right] e^{-ict}$$

### Double Sideband

Just multiply $f(t)$ by one cosine wave

$$\mathcal{FT}\{\cos(ct) f(t)\} = \frac{1}{2}(F(\omega-c) + F(\omega+c))$$

NB FDMA and phase modulation
can be used for multiple
applications

## QUANTISED DEGREES OF FREEDOM IN CONTINUOUS SIGNAL

Bandlimiting continuous function means it has a finite, countable number of
degrees-of-freedom $\Rightarrow$ quantised

Nyquist's Sampling Theorem : if signal strictly bandlimited to some highest
frequency $\omega$, then sampling at rate $2\omega$ specifies it completely everywhere. ⇒

Logan's Theorem : if bandlimited to one octave, simply listing zero-crossing
determines it

## Nyquist's Theorem

Sampling Function $comb(t) = \delta_x(t)$ ~ endless sequence of regularly spaced ~~times~~ times separated by some sampling interval $X$.

Each ~~time~~ tine is a Dirac $\delta$-function. limit of a Gaussian who's width shrinks to 0. Multiplying signal with $\delta(t)$ samples value at $t=0$. Portrayed sequence of tines spaced by $X$ is sum of shifted $\delta$-functions, making sampling comb.

$$\delta_x(t) = \sum_n \delta(t - nX)$$

Sampling function $\delta_x(t)$ is self-Fourier $\longrightarrow \frac{1}{X}$ in frequency

$$FT(\delta_x(t)) = \Delta_x(w) = \frac{1}{X} \sum_m (wX - 2\pi m)$$

$$\delta(t) = \begin{cases} \infty & \text{if } t=0 \\ 0 & \text{if } t \neq 0 \end{cases}, \qquad \int_{-\infty}^{\infty} \delta(t) = 1$$

Multiplying $f(t)$ by $\delta(t-c)$ picks out value of $f(t)$ at $t=c$

$$\int_{-\infty}^{\infty} f(t)\, \delta(t-c)\, dt = f(c)$$
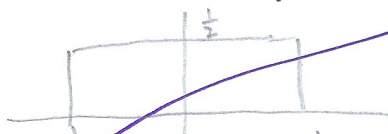
Given signals which have been bandlimited to $w$:

$$F(w) = 0 \quad \text{for } |w| > w$$

If not already true for $f(t)$, first run through lowpass filter - means Fourier Transform $F(w)$ becomes truncated.

If sampling $f(t)$ multiplying by $\delta_x(t)$ + use sampling rate $2w$ so sampling interval $X \leq 1/(2w)$ we know fourier transform of resulting sequence of samples = convolution of $F(w)$ with $\Delta_x(w)$ - Fourier transform of $\delta_x(t)$

Means $F(w)$ becomes reproduced at every tine of $\Delta_x(w)$. But given they did not overlap or superimpose get original by strict lowpass filtering.

$\downarrow$ since sampling rate is double $w$

ideal lowpass filter described as zero-centred pulse function $F(w)$ in frequency domain $w$.



if not, get aliasing, not separatable by lowpass filters

get weird effects - eg backwards moving wheels.

invese Fourier transform of $F(w)$ is $f(x)$

$$f(x) = \int_{-\infty}^{\infty} F(w)\, e^{iwx}\, dx = \frac{1}{2} \int_{-1}^{1} e^{iwx}\, dx$$

$$= \left[ \frac{1}{2} \frac{e^{iwx}}{ix} \right]_{w=-1}^{w=1} = \frac{\sin(x)}{x} = sinc(x)$$

For lowpass filter at $w$: $\dfrac{\sin(2\pi wt)}{2\pi wt}$

Effectively, sinc fills in space between sample points giving us back the same continuous function $f(t)$

<u>Logan's Theorem</u>: if signal is bandpass limited to one octave: $N_H \leq 2N_L$, then listing zero-crossings suffice for exact reconstruction (up to scale factor in amplitude)

↳ explains cartoons

Does not work with Amplitude Modulated signals — de Broglie principle

<u>Gabor's Information Diagram</u>: <u>Uncertainty principle</u> limits localisibility of a signal in time and frequency

Shorter duration = broader bandwidth

Diagrams constructed with axes: time & frequency - no function with area $> \frac{1}{4\pi}$ ⟹ (blobs) can be any shape

↳ logons

## <u>GABOR - HEISENBERG - WEYL PRINCIPLE</u>

Effective width $(\Delta x)$ of complex-valued functions $f(x)$ in terms of normalised variance:

— normalized first-moment of $\|f(x)\|$

$$\mu = \frac{\int_{-\infty}^{\infty} x f(x) f^*(x) \, dx}{\int_{-\infty}^{\infty} f(x) f^*(x) \, dx}$$

$$(\Delta x)^2 = \frac{\int_{-\infty}^{\infty} f(x) f^*(x) (x-\mu)^2 \, dx}{\int_{-\infty}^{\infty} f(x) f^*(x) \, dx}$$

Effective width of $(\Delta \omega)$ & $(\Delta \omega)^2 = \frac{\int_{-\infty}^{\infty} F(\omega) F^*(\omega) (\omega - v)^2 \, d\omega}{\int_{-\infty}^{\infty} F(\omega) F^*(\omega) \, d\omega}$   where $v$ is the mean value of $\|F(\omega)\|$

$F(\omega) = \mathcal{F}\{f(x)\}$

$$v = \frac{\int_{-\infty}^{\infty} \omega F(\omega) F^*(\omega) \, d\omega}{\int_{-\infty}^{\infty} F(\omega) F^*(\omega) \, d\omega}$$

using Schwartz Inequality, can show: $(\Delta x)(\Delta \omega) \geq \frac{1}{4\pi}$
↳ property of all functions and their Fourier Transforms

<u>Gabor Wavelets</u>: family of functions achieving lower bound in uncertainty - complex exponentials multiplied by Gaussians

Gaussians space parameter / phasors

Logons

$f(x) = \exp(-(x-x_0)^2 (\alpha^2)) [\exp(i\omega_0 (x-x_0))]$

⑰ Wavelets are self-Fourier: $F(w) = \exp(-(w-w_0)^2 \alpha^2) \exp(-ix_0(w-w_0))$

          ⤷ this is still an Gaussian

Helical functions of $x$, localised at epoch $x_0$, modulated with frequency $w_0$ and size or spread constant $\alpha$

For wavelet with epoch at $x_0 = 0$, Fourier Transform is Gaussian at modulation frequency $w_0$ and size $1/\alpha$

Gabor proposed them as expansion basis – but since non-orthogonal, difficult to compute expansion coefficients



as $w_0$ increases

$F(w)$

self similar wavelets with $a^{-b} \sim w_0$

$F(w)$

2D Gabor Wavelets are used in Computer Vision – form complete basis image structure with vocabulary of:

    ⤷① Location
    ⤷② Scale
    ⤷③ Spatial Frequency
    ⤷④ Orientation
    ⤷⑤ Phase

Particular example is iris recognition systems which is particularly used for phase structure
    ⤷ works because of high entropy, hence collision avoidance

## DATA COMPRESSION CODES

Redundancy = potential for compression

① Run Length Encoding : summarise repetition

② Predictive Coding : deviations from predictions encoded rather than just the information     Lempel and Ziv

③ Dictionary Based Methods : exploit fact that strings of symbols have probabilities that vary much more than probabilities of symbols individually – sparseness can be exploited
    ⤷ first construct dictionary then scan – most common words have shortest indices to record.
    ⤷ This can be done adaptively.

## Vector Quantization

Also uses dictionary lookup - exploits sparsely populated combinations of sample symbols - can be generalized to image structures for example by creating codebook of pixel combinations. However, ~~use~~ of codebook memory for ~~memory size~~ is an issue.

## DISCRETE AND FAST FOURIER TRANSFORMS

Describe functions of discrete values - $f[n] : f[0], f[1], \ldots$ => effectively a vector of data points. Discrete Transforms of $f[n]$ is matrix multiplication $(N \times N)$ matrix

$$(F[1], \ldots, F[N]) = (f[1], \ldots f[N]) \begin{pmatrix} e_1[1] & \cdots & e_n[1] \\ \vdots & & \vdots \\ e_1[N] & \cdots & e_n[N] \end{pmatrix}$$

Fast Fourier Transform takes this from $O(N^2)$ to $O(N \log_2 N)$

$N^{th}$ roots of unity: $e_k = \exp(-2\pi i n k/N)$

can represent any data sequence $f = (f[0], \ldots f[N-1]) \in \mathbb{C}^{\wedge}$ by vector sum:

$$f = \frac{1}{N} \sum_{k=0}^{N-1} \langle f, e_k \rangle e_k$$

Discrete Fourier Transform

$$\boxed{F[k]} = \langle f, e_k \rangle = \sum_{n=0}^{N-1} f[n] \exp(-2\pi i n k /N) \quad \text{and have inverse:}$$

$$f[n] = \frac{1}{N} \sum_{k=0}^{N-1} F[k] \exp(2\pi i n k /N)$$

Complete DFT requires as many Fourier coefficients as number of values in sequence of $f[n]$ that we're doing DFT for.

### Properties

↳① Cyclical Convolution: CC of $f[n]$ and $g[n]$ of period $N$ :

$$(f*g)[n] = \sum_{m=0}^{N-1} f[m] g[n-m]$$

↳ if negative, taken mod N

Hence, DFT of $f*g$ is product $F[k] G[k]$ where $F$ and $G$ are the DFTs of $f$ and $g$ respectively

### Fast Fourier Transform

Exploits ~~inefficienc~~ inefficiencies in computing DFT

$$F[k] = \sum_{n=0}^{N-1} f[n] \exp(-2\pi i n k /N)$$
$$= f[0] + f[1] \exp(-2\pi i k/N) + \ldots + f[N-1] \exp(-2\pi i k(N-1)/N)$$

Hence, need to do N multiplications and N additions. To do this for
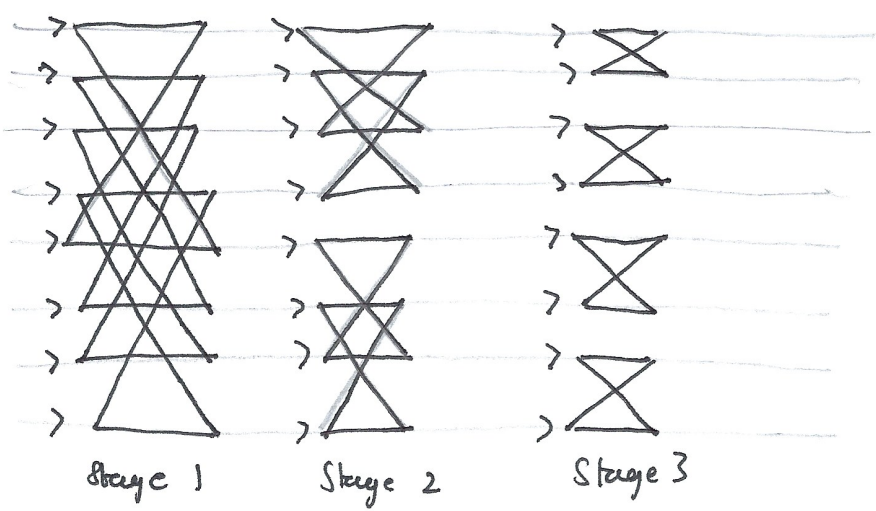$\lambda = 0, 1, ..., (N-1)$, requires $2N^2$ operations.

$$W = \exp(2\pi i / N)$$

Hence, $\exp(-2\pi i n k / N) = W^{-nk}$

Can rewrite Fourier coefficients in terms of powers of $W$. ⇒ can split into two
halves – each requiring only a quarter of the multiplication
    ↳ This (Danielson-Lanczos Lemma) can be done recursively
    ↳ Hence, $O(N \log_2 N)$



Stage 1    Stage 2    Stage 3

Also called
the Butterfly

↓

⊕ only $O(N)$
in space
terms

Fourier Methods can be extended to higher dimensions of functions. Fourier
components of images (of N×M) are 2D complex exponentials $f[n,m] = \exp(2\pi i (kn/N + jm/M))$
    ↳ spatial frequency $= \sqrt{(\lambda^2 + j^2)}$
    ↳ phase (orientation) $= \arctan(j/\lambda)$

## WAVELET REPRESENTATIONS OF INFORMATION

Wavelets are size-specific local undulations – acting as bases for representing information.
Hence, used as basis for JPEG-2000. Image compression works since:
    ↳ neighbouring pixels are highly correlated
    ↳ Projecting this onto Fourier basis leads to highly decorrelated
        coefficients ⇒ many are 0 or small so don't need to be encoded
    ↳ done by Direct Cosine Transform on 8×8 tiles – coefficients are
    quantised ⇒ RLE is used
        ↳ more coarsely for higher frequencies and fewer bits used for this
        (defined by a quantisation table)

### Dyadic Wavelets
Dyadic transformations for generating wavelet $\psi(x)$ spawn orthonormal wavelet basis
$\psi_{jk}(x)$ for expansions of $f(x)$

$$f(x) = \sum_{j=-\infty}^{\infty} \sum_{k=-\infty}^{\infty} c_{jk} \boxed{\psi_{jk}(x)}$$ → generating by shifting and scaling operations
applied to another wavelet
also include θ for notation if → $\psi_{jk}(x) = 2^{j/2} \psi(2^j x - k)$

The compression works well but can suffer from block quantisation artifacts. 2000: using encoders like Daubechies 9/7 wavelet. Across multiple scales and lattice of positions, wavelet inner products with the image yield coefficients that constitute the Discrete Wavelet Transform (DWT). Implemented by simply recursively filtering and downsampling image vertically and horizontally in scale pyramid

## KOLMOGOROV COMPLEXITY

Any set of data can be created by a program - algorithmic complexity is the minimum length of this program => this Kolmogorov Complexity is the data string's minimal description length.

  ↳ this is approx equal to entropy of distribution from which data string drawn.
  ↳ Most sequences length n have Kolmogorov Complexity K close to n.
   ↳ if >n, sequence is algorithmically random
  ↳ Not computable - cannot know you've found shortest possible description

Infinite string is K-incompressible iff:

$$\lim_{n \to \infty} \frac{K(x_1, x_2 \ldots x_n | n)}{n} = 1$$

Strong Law of Large Numbers for Incompressible Sequences - proportions of 0s and 1s in any incompressible string must nearly be equal. Must also pass all statistical test for randomness.

## SCIENTIFIC APPLICATIONS

① <u>Astrophysics</u>

Pulsars are collapsed neutron stars - spinning causes emission of an electromagnetic radio beam. Signals faint and hidden in noise. Since coherent, can use auto-correlation integral $p_f(t)$ extracts this component from $f(t)$ as noise cancels itself out.

$$p_f(t) = \int_{-\infty}^{\infty} f(s) f(s+t) ds$$

Inverse Fourier transform of Power modulus.

If $f(t)$ has $\mathcal{FT}\{f(t)\} = F(\omega)$, then $FT\{p_f(t)\}$ is power spectral density of $f(t)$

$$FT\{p_f(t)\} = \boxed{F(\omega) F^*(\omega)} - \text{power function}$$

$$p_f(t) = FT^{-1}\{F(\omega) F^*(\omega)\}$$

② <u>Genomics</u>: translating genetic code into Amino Acids regarded as error-prone information channel → Entropy of gene inheritance is 1 bit
  ↳ $H = N$ bits per gene where $N$ is number of generations gone back
  ↳ Every person appears multiple times in family tree → Genetic Isopoint is time in history when everyone is ancestor of everyone else: $N \approx 1.77 \log_2(m)$ generations ago.

Information transmission rate by blowfly photoreceptors & SNR(w) analysis at synapses using $C = \int \log_2(1 + SNR(w)) dw$ is upto $1.65$ kB/s

## ③ BIOMETRIC PATTERN RECOGNITION
  ↳ Iris Kids
  ↳ Face Recognition
  ↳ Forensics

Discriminating power of a biometric reflects its entropy; surviving large database searches without false matches requires high entropy
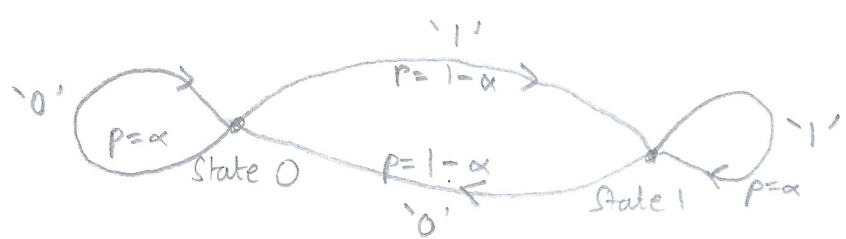
### Why phase?
↳ Achieves structural information, independent of amplitude. Hence achieves some invariances
↳ Higher entropy than amplitude — hence coarsely quantied
↳ Phase classification ≡ clustering algorithm

Can use Gabor wavelets (which encode phase naturally, but in frequency-specific way)
or in total way : → ① Analytic Functions: $f(x)$ – Hilbert Transforms i $f_{HT}$ cousin
                  ↳ $f(x) - i f_{HT}(x)$

### Iris Code

Regard Iris Code as channel Codes computed from natural and from white noise iris patterns are well-modelled by bits emitted from two-state Markov process, with differing values of $\alpha$:



$H(\alpha) = -\alpha \log_2(\alpha) - (1-\alpha) \log_2(1-\alpha)$ bits per bit emitted
Statistical distributions well computed with $\alpha = 0.867 \Rightarrow C = 0.566$ bits per bit.
  ↳ Correlations lead to this reduction in entropy

Lots of mutual information between adjacent rings of Iris Codes
    (0.311 bits per bit pair max)
Mash used to remove eyelashes, specular reflections from glasses, etc.
    ↳ eg. A and B are data words of two Iris Codes
    ↳ C, D are respective mash words
    result = $(A \oplus B)$ & C & D
    ↳ each executes in one clock tick