

# A Mathematical Theory of Communication

*By C. E. Shannon*

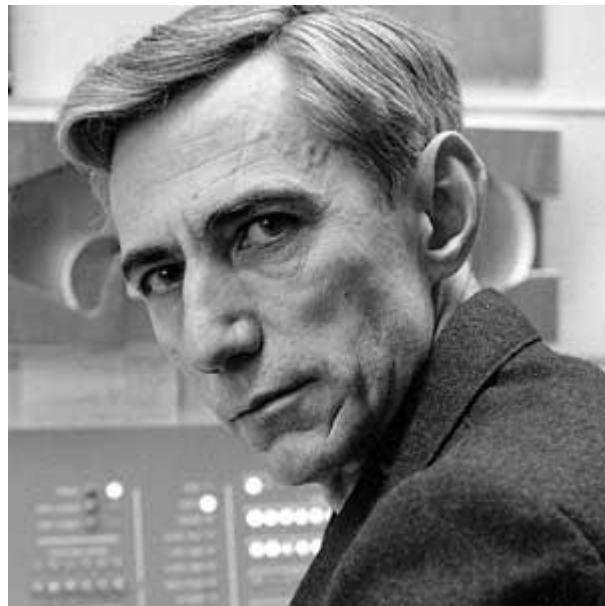
# Agenda

- Introduction
  - History, Who is Shannon?
- Basic Information Theory
  - What is information?
  - How to measure it?, Entropy
- Grams
  - Stochastic Process for generating symbols / information
  - n-gram

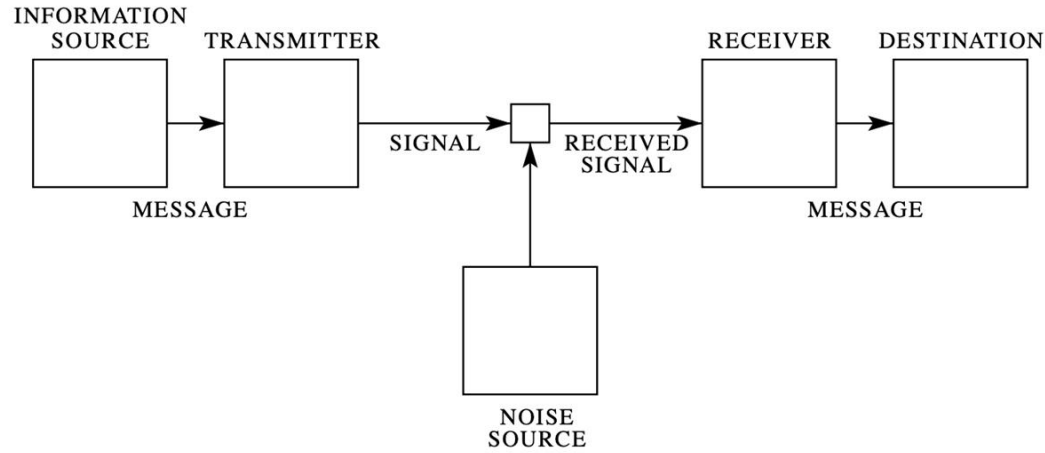
# History

Claude Elwood Shannon, often referred to as the "father of modern digital communication and information theory," was a pioneering American mathematician, electrical engineer, and cryptographer. He made significant contributions to various fields, including information theory, cryptography, and computer science.

- Shannon's most famous work, "A Mathematical Theory of Communication" (1948), laid the foundation for information theory.



Claude Elwood Shannon



| LETTERS<br>FIGURES |   | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z | CARRIAGE<br>RETUN | LINE<br>FEED | LETTERS | FIGURES | SPACE | ALL-SPACE<br>NOT IN USE |
|--------------------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|-------------------|--------------|---------|---------|-------|-------------------------|
| CODE<br>ELEMENTS   | 1 | ● | ● |   | ● | ● | ● |   |   |   | ● | ● |   |   |   |   |   | ● |   | ● |   | ● |   | ● | ● | ● | ● |                   |              | ●       | ●       |       |                         |
|                    | 2 | ● |   | ● |   |   |   | ● |   | ● | ● | ● | ● |   |   |   | ● | ● | ● |   |   | ● | ● | ● | ● | ● |   |                   | ●            | ●       | ●       |       |                         |
|                    | 3 | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○                 | ○            | ○       | ○       | ○     | ○                       |
|                    | 4 |   | ● | ● | ● |   | ● | ● |   |   | ● | ● |   | ● | ● | ● | ● |   | ● |   |   |   | ● |   | ● | ● | ● | ●                 | ●            |         | ●       | ●     |                         |
|                    | 5 |   | ● |   |   |   | ● | ● | ● |   |   |   | ● |   |   | ● | ● | ● |   |   | ● |   | ● | ● | ● | ● | ● | ●                 |              |         | ●       | ●     |                         |

**The International Telegraph Alphabet**

- INDICATES A MARK ELEMENT (A HOLE PUNCHED IN THE TAPE)
- INDICATES POSITION OF A SPROCKET HOLE IN THE TAPE

# Information Theory

What is information ?

- Semantic meaning of a message is irrelevant to its transmission. A message should be conceived as a sequence with statistical properties. It is the message's statistics that could be captured and its coding minimized to allow for effective transmission.
- Information is measured in *bits*, and one bit of information allows you to choose between two equally alternatives.

Basic laws of information

1. Upper limit for channel capacity
2. Noise
3. Encoding of data

# Measure of information

$$H = - \sum_i^n p_i \log(p_i)$$

Consider **a**, **b**, **c**,

$$p(a) = 2/4 = 0.5$$

$$p(b) = 1/4 = 0.25$$

$$p(c) = 1/4 = 0.25$$

When,  $p(c) = 1/4$ , we need the following number of bits to represent them all.

$$\log_2(4) = \log_2(1/p) \text{ bits} = 2 \text{ bits}$$

Now, to determine the overall storage required

- a occurs half and need 1 bit, so  $0.5 \times 1$
- b occurs quarter of the time and needs 2 bits, so  $0.25 \times 2$
- c occurs quarter of the time and needs 2 bits, so  $0.25 \times 2$

In total  $(0.5 \times 1) + (0.25 \times 2) + (0.25 \times 2) = 1.5$  bits



# Stochastic Process for generating symbols / information

- Consider 5 letters A B C D E, let the probabilities be 0.4, 0.1, 0.2, 0.2, 0.1.
- A typical message constructed from this source is "AACDCBDCEAADADACEDAEADCBEDADDCECAAAAAD."
- Transition Probabilities  $p_i(j)$ 
  - The probability of  $i$  is followed by letter  $j$

N = 1 : This is a sentence *unigrams:* this, is, a, sentence

N = 2 : This is a sentence *bigrams:* this is, is a, a sentence

N = 3 : This is a sentence *trigrams:* this is a, is a sentence

Image Source : Stackoverflow

# $n$ - gram

- $n$ - grams is a sequence of  $n$  adjacent words or letters from the text corpus
- $\langle s \rangle$  I am Sam  $\langle /s \rangle$
- $\langle s \rangle$  Sam I am  $\langle /s \rangle$
- $\langle s \rangle$  I do not like green eggs and hams  $\langle /s \rangle$

Here are the calculations for some of the bigram probabilities from this corpus

|  |  |                                 |
|--|--|---------------------------------|
| $P(I   \langle s \rangle) = 2/3 = 0.67$          | $P(\text{Sam}   \langle s \rangle) = 1/3 = 0.33$ | $P(\text{am}   I) = 2/3 = 0.67$ |
| $P(\langle /s \rangle   \text{Sam}) = 1/2 = 0.5$ | $P(\text{Sam}   \text{am}) = 1/2 = 0.5$          | $P(\text{do}   I) = 1/3 = 0.33$ |

- For the general case of MLE  $n$ -gram parameter estimation:

$$P(w_n | w_{n-N+1:n-1}) = \frac{C(w_{n-N+1:n-1} w_n)}{C(w_{n-N+1:n-1})}$$

# $n$ - gram Applications

- Sentiment analysis
- Text classification
- Text generation

Thank You