# Encoders and Decoders

- Encoders and decoders are two key components in many ML models, especially for sequence to sequence tasks like CLM

  • encoders take inputs (like a sentence, img or audio) and transform it into a compact, informative representation that captures its meaning or features

  • the decoders then uses that representation to generate or reconstruct the desired output, such as a translated sentence, summarized text or reconstructed img.

Ex (in transformers): in transformers used for translation the encoder reads and processes the input sentence (eg in english) to produce contextual embeddings that capture the meaning of each word in relation to the whole sentence. The decoder than takes these embeddings and generates the output sentence (ie in French) one word at a time using attention to focus on the most relevant parts on the encoded input

## Normalization and standardization

• Normalization is the process of scaling data so all features are within a specific range, typically $[0,1]$ or $[-1,1]$. This ensures that no feature dominates others simply because it has large numeric values

- for min-max normalization (most common) $X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$ → Scale $[0,1]$

- for feature scaling to $[-1,1]$ do: $X_{norm} = 2 \times \frac{X - X_{min}}{X_{max} - X_{min}} - 1$ → do for each data point

• Standardization (also called z-score norm) rescales data so that it has a mean $(M)$ of 0 and a standard deviation $(\sigma)$ of 1, this centers the data and ensures all features have comparable variance formula: $Z = \frac{X - M}{\sigma}$ → for each point

- Both together help model convergence/stable training by ensuring features contribute equally

## Rotary Positional Embeddings (RoPE)

• RoPE are a a way to encode token positions (positional embeddings see pg 112 for what those are) in transformers by rotating the query and key vectors in multi-headed attention, allowing the model to capture relitive positions efficiently without adding seperate position vectors. it works because the dot product of rotated K,Q vectors naturally encodes relitive position. <u>115</u> why: gen well to long er text sequences
  • is more efficient, cuts extra vectors