

# Transformers (LLM Notes cont)

- Transformers can be used for img gen, text to speech translation of language etc but in our case ~~transformers~~ are next word predictors they take in a piece of text as input (along with image or sound for context) and outputs the next word in that piece of text. This predicts all possible next words as a probability distribution and the word with highest prob is chosen. ~~LLM~~ LLM are a transformer based Model trained on a large amount of text to understand and generate language like GPT. So while transformers can be used to generate many things LLM are used to generate text then the LLM repeatedly prompts the transformer with the additional new predicted word and does this until a stopping point is reached like a stop word "EoS" or "[End]" is reached or a token limit. This is autoregressive generation

## - Step by step inside a transformer

(Embedding  $\rightarrow$  Attention  $\rightarrow$  MLP  $\rightarrow$  unembedding)  $\rightarrow$  general transformer architecture

- Embedding: first the input is broken up into individual pieces:

Ex: To date, the smartest person is  $\rightarrow$  Tokenization  $\rightarrow$  [To] [date] [the] [smartest] [person] [is]

in LLM these pieces are the words that make up the prompt =  $\square$  these pieces are called tokens. in img's they can be pixels in voice, fractals

then each token is associated with a vector, a list of numbers which is meant to encode the meaning of that piece

Ex: To date, the ... ?  
 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$   
 $\begin{bmatrix} 5.2 \\ 7.1 \\ 1.1 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 3.5 \\ 6.1 \\ 7 \\ \vdots \\ 2 \end{bmatrix} \begin{bmatrix} 9 \\ 7 \\ 14 \\ \vdots \\ 41 \end{bmatrix} \begin{bmatrix} 12 \\ 13 \\ 16 \\ \vdots \\ 1 \end{bmatrix}$  also encode position

• these vectors are like coordinates in high dimension space and similar words tend to land on vectors that are close to each other in that 48 space and vice versa

