# Word Embeddings (transformers & LLM cont)

- In reality tokenization frequently divides words, not whole words every time but we will focus on it by thinking tokens are clean whole words

truth : [This] [pro]kess] [(known] [Fanci] fully]   Lie: [this] [precces] [(] [known] [fanci fully]

token ⟵  └→ and same for output token it can be any token not always a full word

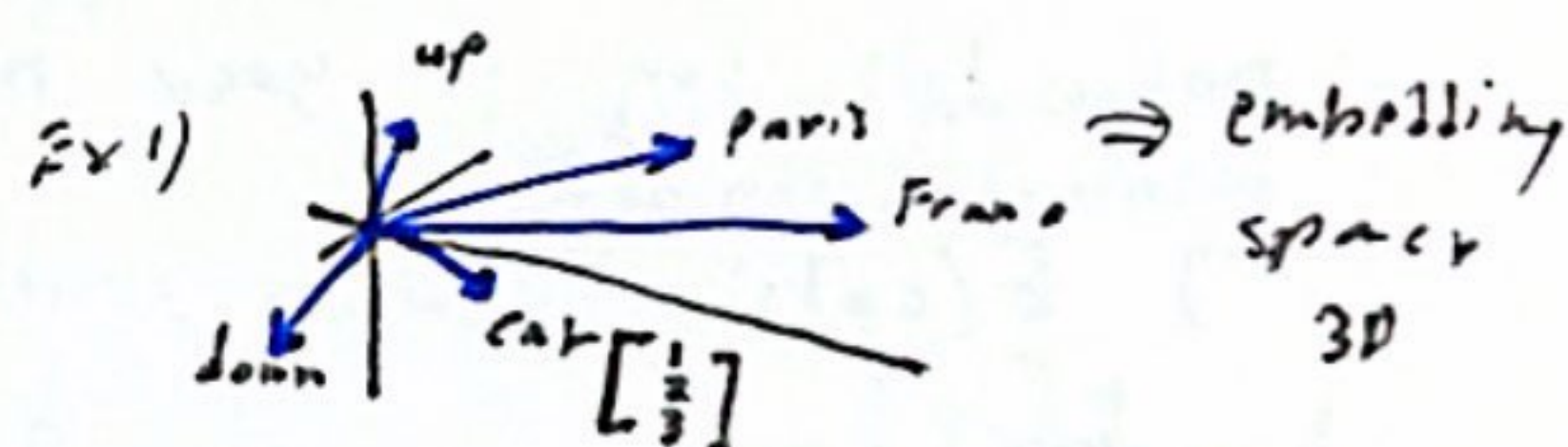- The model has a pre defined Vocabulary, some list of all possible tokens (words, symbols, spaces, number etc) → Not always whole and the first matrix we will see is the embedding matrix has a single column for each of these words these columns determine what vector each word get in that first step We label $W_E$ = Embedding Matrix.

All words ≈ 50k = 50k tokens

$$\begin{bmatrix} aback & blue & boy & \cdots & zing & ; \\ 2.1 & 2.1 & 5 & 7 \cdots & 2 & 5 \\ 3.2 & 2 & 6 & 8 \cdots & 7 & 7 \\ 5.1 & 3 & 7 & 9 \cdots & 6 & 8 \\ 6 & 9 & 8 & 1 & 8 & 9 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 9 & 9.11 & 2 & 5 \cdots 7 & & 9 \end{bmatrix}$$

Embedding Matrix $W_E$

the sky was blue

$$\begin{bmatrix} : \\ . \end{bmatrix} \begin{bmatrix} : \\ : \end{bmatrix} \begin{bmatrix} : \\ : \end{bmatrix} \begin{bmatrix} 5 \\ 6 \\ 7 \\ 8 \end{bmatrix}$$ ⟵ EM picks the vector for blue

✴ Values begin random and learned through data ✴

- We embed these words as tokens which means its a vector in some very high dimension space like 12,288 dimentions in GPT-3. So we visulize in 3D. as it learn these embedding in training it settles on embedding's whos direction have meaning for Ex paris and France's embedding are

Ex 1)


⇒ embedding space 3D

Ex 2)



$E(man) - E(woman) \approx E(King) - E(woman)$

close as there related but up and down opposite in meaning and space in the Ex 2 we see how Embedding are related and why using vectors is good say we did not know what a female leader is we could do King + (woman - man) and finding the vector (embedding) closest to that. $E(queen) \approx E(King) + E(woman) - E(man)$

- this is a perfect Ex in reality they might be further apart and are impossible to visulize. they are leaned in training.

- further more "Queen" in training data would not always be a female Leader Ex "Queen the band" hence the 50 embedding for Queen wont be the gives formula ... close still