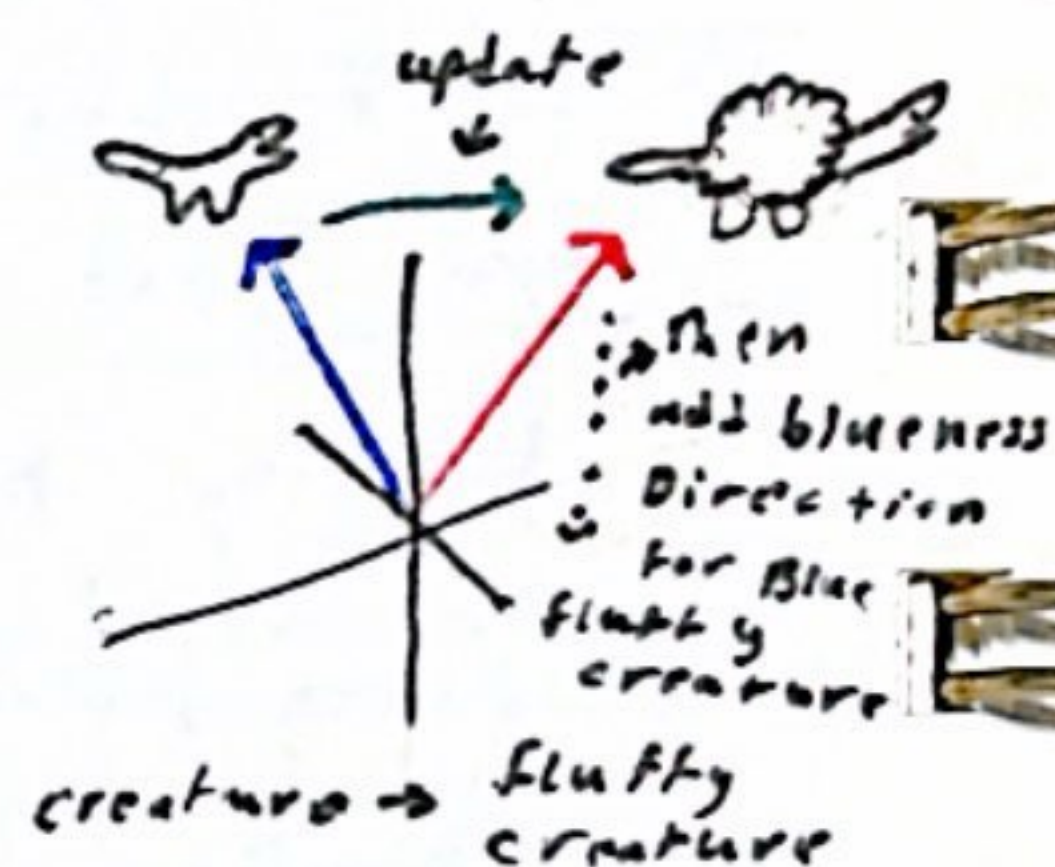
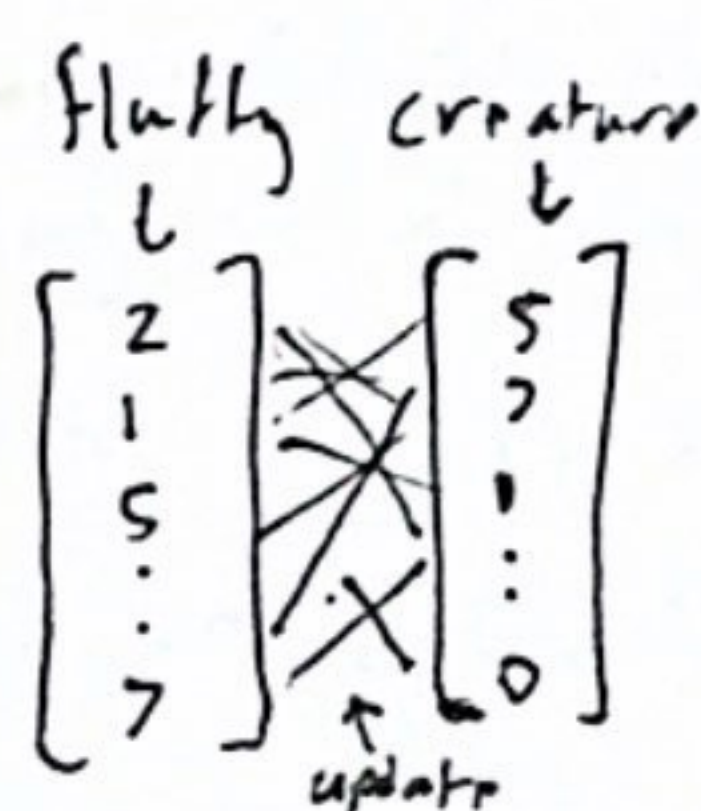


Attention (P-6) (LLM notes)

- in the attention pattern its size is the square of context size in our case $\text{Attention pattern} = 8^2 = 64$ and hence context size is a huge bottle neck. even though we can compute dot products in parallel with gpus scaling context size in LLM's is not easy. to achieve larger context sizes we use new techniques like Blockwise attention Ex.

- Updating embeddings: Now that we have the attention pattern, which tells us which words are relevant to which other words. Now we need to update those words by updating its embedding i.e. allowing words to pass info

to other words that are relevant too for ex you work embedding of fluffy to some how cause a change to creature that moves it in this 12K embedding space

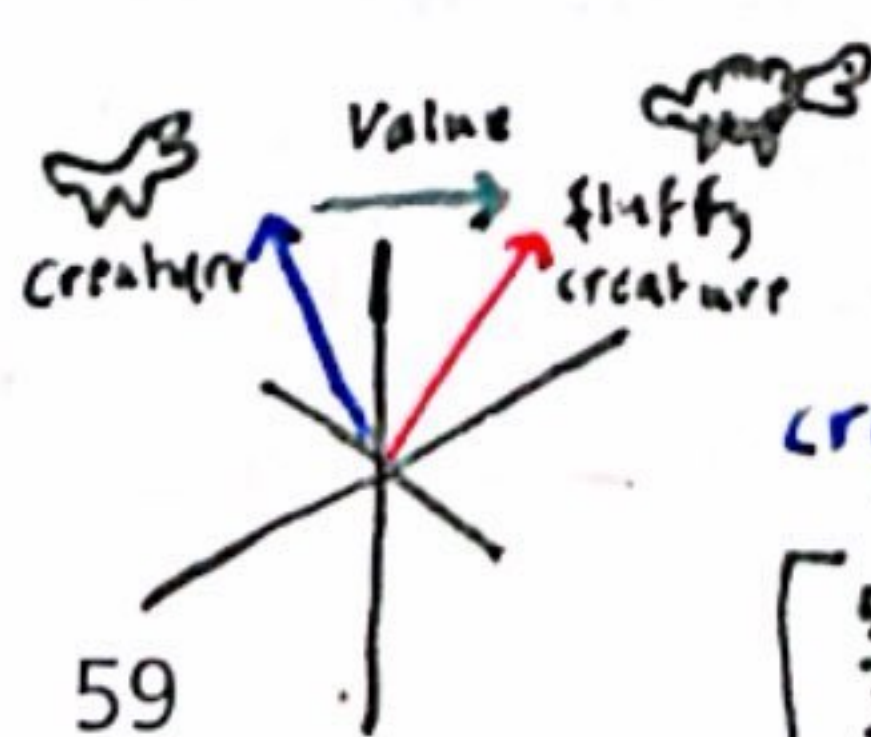


that more specifically encodes a fluffy creature.

- for single headed attention: this is the straight forward way and is to use another matrix, Value matrix (W_v). which you multiply by the first word "fluffy" the result of this is a value vector and you add this to the second word "creature"

Think of it as saying: if this word is relevant for adjusting the meaning of another word what exactly should be added to the embedding of that other word to reflect this.

$$\begin{matrix}
 12,289 & 112,219 \\
 \begin{bmatrix} 1 & 2 & 5 & 7 & \dots & 9 \\ 2 & \dots & \dots & \dots & \dots & \dots \\ 5 & \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & \dots & \dots & \dots \end{bmatrix} & W_v & \begin{bmatrix} 2 \\ 1 \\ 5 \\ \vdots \\ 7 \end{bmatrix} = \begin{bmatrix} -21 \\ 7 \\ 6.1 \\ \vdots \\ 0 \end{bmatrix} \\
 \text{learned in training} & & \text{same Dim}
 \end{matrix}$$



$$\begin{aligned}
 \text{fluffy} &= \text{new embedding of creature} \\
 &= W_v \cdot \text{fluffy} + \text{creature embedding old} \\
 &= \begin{bmatrix} 5 \\ 7 \\ 0 \end{bmatrix} \text{ update } \begin{bmatrix} 2.1 \\ 8.2 \\ 2.1 \end{bmatrix} \\
 &\text{was its before creature, the word to update}
 \end{aligned}$$