

Attention (P-9) (LLM notes)

- Multi headed attention. (all run in parallel)

• in the ex's from P1-P7, like the Ex of adjectives updating nouns is one single head of attention, but there's many more heads needed to add context like for ex updating words that are spelled the same but have different meaning, so we do what we did in single headed attention 10,000 times giving LLM's the level of context they really need.

Ex a upside-down fluffy blue creature

- is one head of attention, saying creature is fluffy and blue

- is another head of attention telling the position and orientation of the creature



• This together is multi headed attention, so this + 10,000 = full context

★ Each head has its own Key, Query and Value matrices. These matrices parameters dictate the context its adding like adjectives updating nouns, and is learned in training. There are 10,000's of these heads each learning its "attention pattern" in training. Each attention head proposes a change to the embedding of the word, these changes from each head are added up and added to the original embedding in that attention block.

