

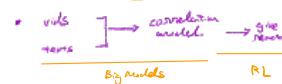
Generalized agents (Yuke Zhu)

△ Big data & RL agents

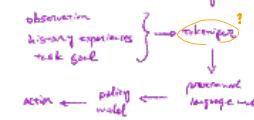
• web data → similar data

↙ real data

encoders?



△ LLM → task planning

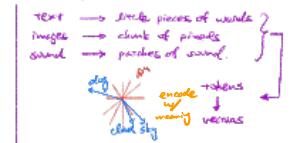


e.g.
input: get me a coffee
output: series of actions ...

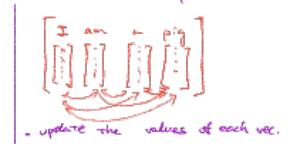
△ GPT

Generative
Pre-trained
Transformer

△ token



△ attention block operation



△ MLP like adding last digit of question

[] ↓ show answer

△ Attention & MLP repeat

↓ predict next objective

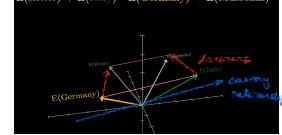
△ word embedding

△ Embedding matrix

$W_E = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$ @ 100, random when

GPT maps word to token
R 12,288 W 50257 tokens

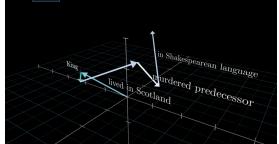
$$E(\text{Hitler}) + E(\text{Italy}) - E(\text{Germany}) \approx E(\text{Mussolini})$$



dot products = connection

△ training

The King doth wake tonight and takes his rouse ...



w/ data the token "King" could start from 2 place; then w/ different data drawing in, we can then pull the token around, then establish connection, like above, perfect visualization.

then, we can do prediction, i.e., generative AI

also, we need "attention", the relation between words.

△ unembedding matrix

$$W_U = [\quad] \in \mathbb{R}^{50257, 128, 6}$$

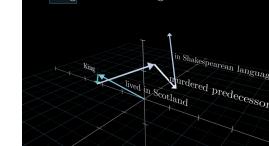
softmax

$$\text{values} \rightarrow \text{prob. distribution} = \frac{e^{x_i}}{\sum e^{x_j}}$$

temperature can add some randomness to the prediction

△ Transformer

The King doth wake tonight and takes his rouse ...



- get attention! between words
- allow words to evoke other memory from other words.

△ Query layer

e.g. A slightly blue creature around the red one.

$$E_1 \downarrow E_2 \downarrow \dots \downarrow E_T$$

multiple w/ W_q, query weights

- these represent the query map.

for instance, Qq's query might be "searching the adjectives".

& this is also "embed w/ a vector".

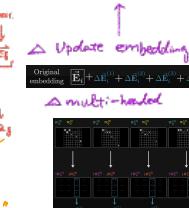
- GPT3 is using 128 query space

△ multi-layer

△ Update embedding

$$\text{Original embedding } E_1 + \Delta E_1 + \Delta E_2 + \Delta E_3 + \dots$$

△ multi-headed



△ Key layer

$$E_1 \downarrow E_2 \downarrow \dots \downarrow E_T$$

$$W_k$$

$$K_1 K_2 \dots K_T$$

- this "answers" the "query"

- Q1K1 → longer, align!

△ Construct a table of Q & K to predict

softmax

ATTEND TO !!!

$$\text{Info} \rightarrow E_1 \rightarrow K_1$$

$$\text{Info} \rightarrow E_2 \rightarrow K_2$$

$$\dots$$

softmax

[Qk] → normalized

Attention Pattern! (longer = stronger)

same diff.

what to add to reduce other words?

French, e.g.

"one head of attention"

Sur transform

$$E_1 E_2 \dots E_T \rightarrow P_1 P_2 \dots P_T$$

