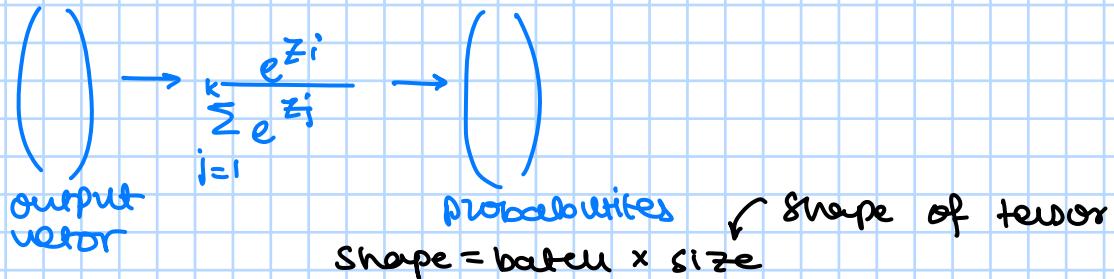


`referee_gpt2.to_tokens("something here")`
versus tensor [...]

logits: we want a probability distribution over next tokens.
we use softmax to convert vector into prob. distribution
so model outputs tensor of logits (vector of size d-vocab for each input)



Text \rightarrow tokens \rightarrow logits

`logits, cache = referee_gpt2.run_with_cache(tokens)`

shape = batch_size x size x d-vocab

↑ shape of logit

`log_probs = logits.log_softmax(dim=-1)`

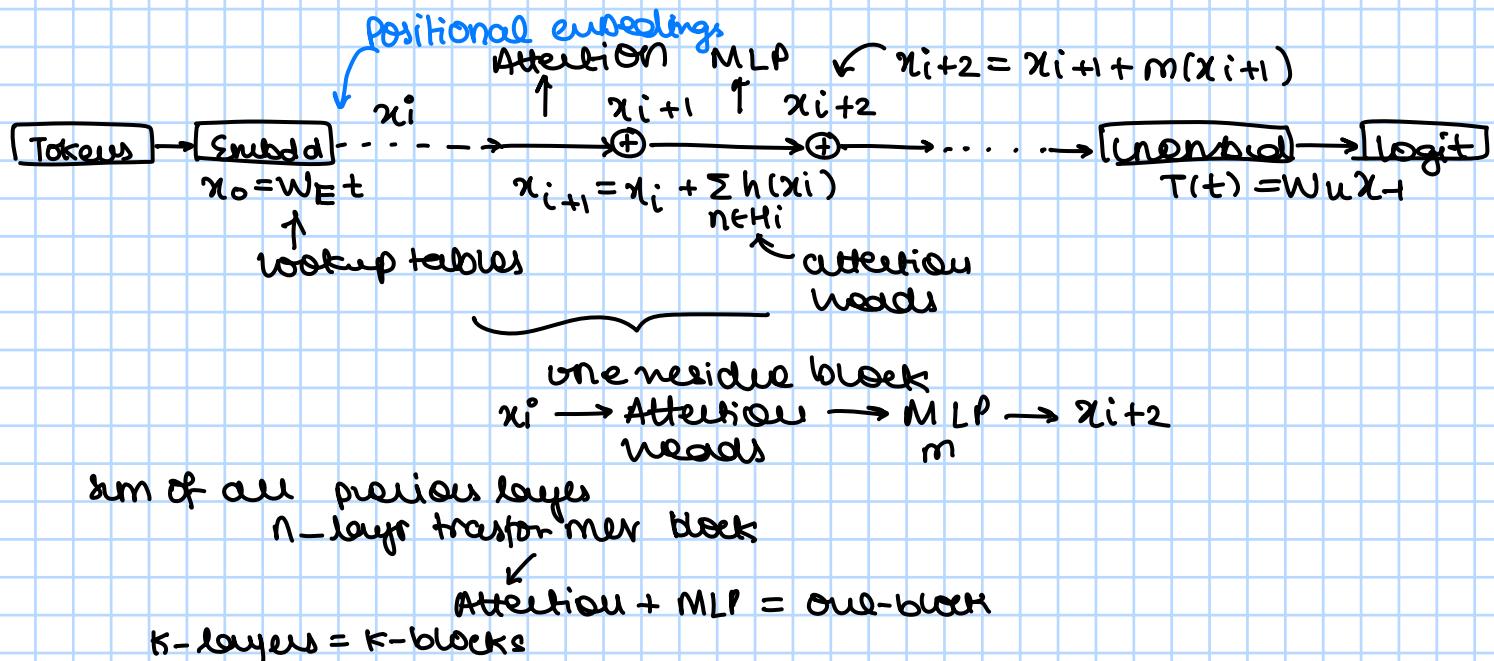
now we get the next token:

\downarrow
`next_token = logits[0, -1].argmax(dim=-1)`

add this in input and re-run

Architecture:

Input tokens, integers \rightarrow Embedding is a lookup table mapping tokens to vectors
vectors (lives in residual stream).



sum of all previous layers

n -layer transformer block

K -layers = K -blocks

Attention:

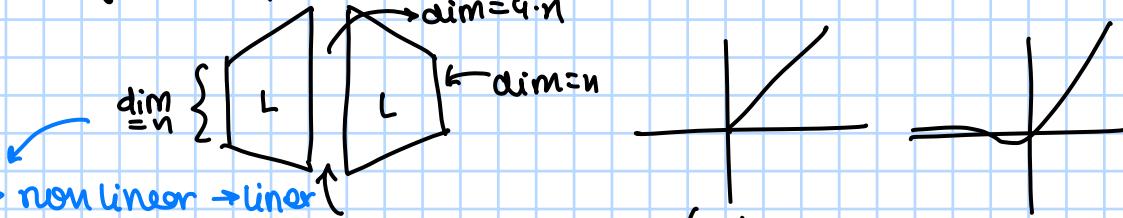
Attention moves information from prior position in seq to current token, we do this for every token, we only look backwade and is called attention pattern \rightarrow probability dist, which some token to get input from

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q K^T}{\sqrt{d_K}}\right) V$$

↑ key ↓ query value

$$Q_i \dots \dots \quad \begin{matrix} K_i & \cdot & \cdot & * \\ \vdots & \cdot & \cdot & \cdot \\ \vdots & 0 & \cdot & \cdot \end{matrix} \quad \left\{ \frac{Q K^T}{\sqrt{d_K}} \xrightarrow{\text{softmax}} V \right.$$

MLP: multilayer preseptron, Standard NN



activation gelu / relu

Note: one attention has moved relevant info to single position, MLP can do computation, reasoning, lookup info, etc.

embedding: final step of going vector to logits, and logit to text

layer norm: applied at each layer start, normalisation function, before Activation, MLP and embedding

all input layer - mean 0, variance 1

- then effective scaling and translation

$$y_0 \in \mathbb{R}^n$$

$$y_1 = y_0 - \frac{\sum y_0}{n}$$

$$\Rightarrow y_1 = y_0 - \langle \langle y_0, (1, \dots, 1) \rangle \frac{1}{\sqrt{n}}, (1, \dots, 1) \rangle \frac{1}{\sqrt{n}}$$

✓ equivalent to setting
first word to 0
(good 'aim')

$$\Rightarrow y_2 = \frac{y_1}{\|y_1\|} \text{ and this makes } \|y_2\| = \sqrt{n}$$

or var = 1

$$\Rightarrow y_3 = y_2 \odot w$$

↑ representation of residue stream

effective

multiplication

$$= y_2 \cdot w$$

diagonal with w as diagonal

$$\Rightarrow y_4 = y_3 + b$$

where $b \in \mathbb{R}^n$

Note: As layernorm is just scale and translate, we can apply layernorm before another linear mlp (linear \times linear = linear) so we can just fold this into one single effective linear layer

positional information: Attention operates over all pairs of positions, so the attention weight from token 5 to token 1 and 5 to 2 will be calculated same but

Nearest tokens are more relevant

so a new lookup table mapping index of the position of each token to a residual stream vector, add it to embed.

we add rather than concatenate

math logic:

LayerNorm:

$$X \in \mathbb{R}^{B \times P \times D}$$

B = Batch size

P = No of positions (tokens)

D = hidden dimension

$$\mu_{b,p,i} = \frac{1}{D} \sum_{i=1}^D x_{b,p,i}$$

$$x'_{b,p,i} = x_{b,p,i} - \mu_{b,p,i}$$

$$\sigma_{b,p}^2 = \frac{1}{D} \sum_{i=1}^D (x_{b,p,i} - \mu_{b,p,i})^2$$

$$\text{scale}_{b,p,i} = \sqrt{\sigma_{b,p}^2 + \epsilon}$$

$$\hat{x}_{b,p,i} = \frac{x_{b,p,i} - \mu_{b,p,i}}{\sqrt{\sigma_{b,p,i}^2 + \epsilon}} \quad \text{to prevent 0 division}$$

if $[2, 4, 768] \rightarrow [2, \underbrace{4}_{X_{b,p,d}}, 768]$ (every column is $\hat{x}_{b,p,i}$)

Embedding:

for each $X \in \mathbb{R}^{B \times P}$

↳ input token

$x_{b,p} \in [0, V-1]$

↳ no of words

$X = [x_{b,p}]$

↑
matrix X

↳ now, $\underbrace{\text{embed}_{b,p,*}}_{\text{vector (column)}} = W_E [x_{b,p}]$

and so $\text{embed}_{b,p,d} = \text{columns of } \underbrace{\text{embed}_{b,p,*}}_{\text{matrix}}$

$\Rightarrow \text{embed} \in \mathbb{R}^{B \times P \times d}$

so $X \in \mathbb{R}^{B \times P}$

$\Rightarrow W_E(X) = \text{embed}_{B \times P \times d}$

↳ matrix

d is model dimension

$\underbrace{\text{pos_embed}_{p,D}}_{\substack{\text{small matrix} \\ \text{for given batch}}} = W_{\text{pos}} [x_{b,*}]$

↑↑ all positions that batch
fixed batch

$\underbrace{\text{pos_emb}}_{\substack{\text{full matrix} \in \mathbb{R}^{B \times P \times D} \\ \text{all values are like } W_{\text{pos}}[x_{b,*}]}}$

fixed column

Attention:

$\hat{x} \in \mathbb{R}^{B \times P \times D}$ (normalised input)

$$Q[b, i, h, :] = \sum_{j=1}^D x[b, i, j] \cdot W_Q[h, j, :] + b_Q[h, :]$$

$$K[b, i, h, :] = \sum_{j=1}^D x[b, i, j] \cdot W_K[h, j, :] + b_K[h, :]$$

$$V[b, i, h, :] = \sum_{j=1}^D x[b, i, j] \cdot W_V[h, j, :] + b_V[h, :]$$

$Q, K, V \in \mathbb{R}^{B \times P \times H \times d}$

↑ ↑ dim of head

$$\text{Attention-score} = \underbrace{Q[b, i, h, :] \cdot K[b, j, h, :]^T}_{n \text{ word}} / \sqrt{d}$$

att-score $\in \mathbb{R}^{B \times H \times P \times P}$
for each pos i , need h

attn-score $[i, j] = -\infty$ if $j > i$

pattern $[b, h, i, :] = \text{softmax}(\text{att-score}(b, h, i, :))$
 $\text{pattern} \in \mathbb{R}^{B \times H \times P \times P}$

$$Z[b, i, h, :] = \sum_{j=1}^P \text{pattern}[b, h, i, j] \cdot v[b, j, h, :]$$

$$\text{attn_out}[b, i, :] = \sum_{h=1}^H \sum_{k=1}^d Z[b, i, h, k] \cdot w[h, k, :, :] + b_0$$

$\text{attn_out} \in \mathbb{R}^{B \times P \times d}$

MLP:

$$\hat{x} \xrightarrow{L} \text{Activation} \xrightarrow{L} x' \rightarrow \hat{x}'$$

Unlabelled: same as previous