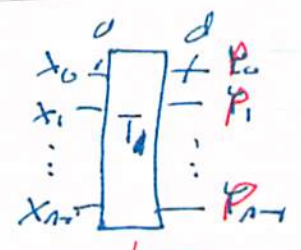


- Now, here is what is in each transformer block T_i :
- recall it is the same size in as out.

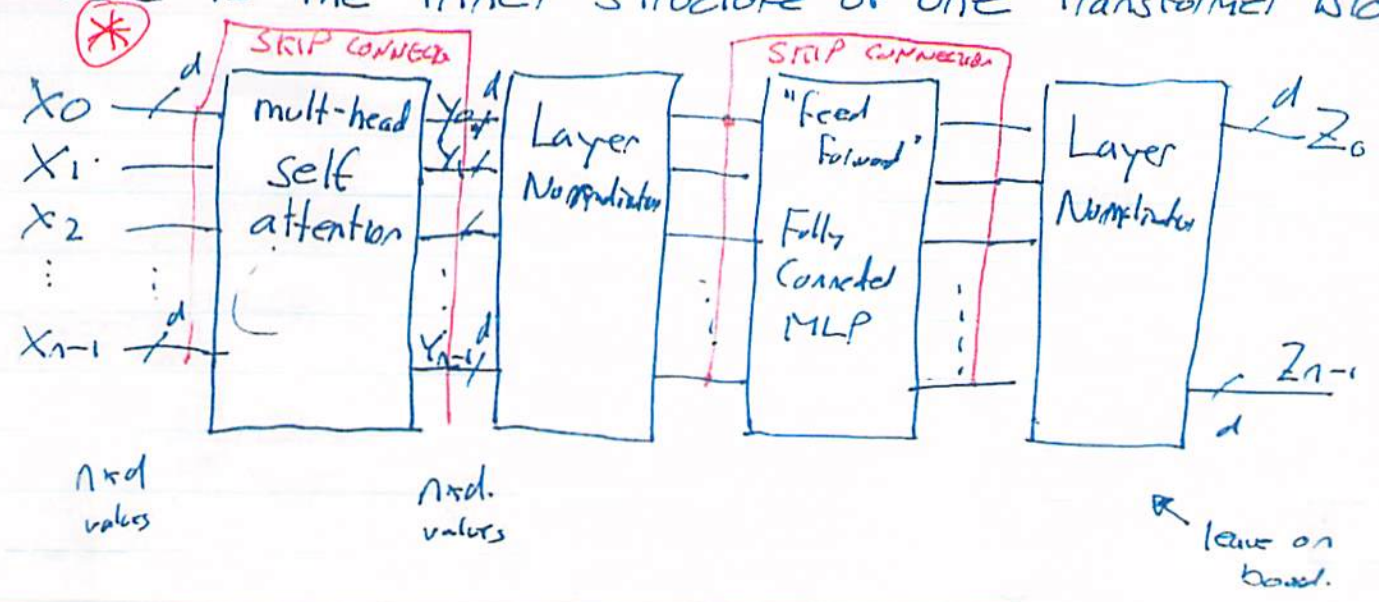


skeptical aside: we could ignore the inner structure of and just say "it is a big network with many parameters"

- it may be, to my mind, the underlying truth, $\hat{=}$ maybe a large MLP or CNN could do this same thing. (just bigger is better.)
- it may be that the particular structure is better.
- will present the structure $\hat{=}$ understanding per Jurafsky 9.7.
- however, the explanations only work for block T_0 , to me.

it does work.
Warning: This is complex!

- Here is the inner structure of one Transformer Block



Here is the core intuition given of what Transformer Block is ^{Attention} doing

- the input word embeddings are transformed from their initial, very general meanings (across all uses/contexts of the words) to something more specific to the context → i.e. the other words in sequence makes some sense.
 ⇒ they call these "Contextual embeddings"

- e.g. the embedding for ~~river~~ "bank" would become different in these ~~two~~ contexts

... ^{she sat on the} river bank ...
 ... ^{she emptied his} bank account ...
 ... they should not bank on the result ...

- from ~~*~~, consider the output y_i , which we say is a transformation of ^{word} x_i the input word.

e.g. ~~The~~ ~~river~~ x_3 bank ~~was~~ ~~steep~~. e.g. $i=3$
 He ^{emptied} his account

- self attention asks the question: how similar is ~~x_3 (bank)~~ ^{each word} to all preceding words and itself!
 [See Jurafsky Section 9.7]

e.g. how similar is x_3 (bank) to x_0 (~~The~~ ^{He} ~~emptied~~)?
 x_3 (bank) to x_1 (~~river~~ ^{his})?
 x_3 (bank) to x_2 (~~was~~)?
 x_3 (bank) to x_3 (bank)?

- how have we computed a ^{single} number that says how similar (related) two words are? ~~before~~
- ⇒ use dot product of word embeddings

Compute $x_0 \cdot x_3$ (can Define

$$\text{score}(x_i, x_j) = x_i \cdot x_j \quad \text{for all } j \leq i$$

$x_1 \cdot x_3$ \rightarrow but need to normalize it, so \rightarrow

$x_2 \cdot x_3$

$x_3 \cdot x_3$

- do this for every word

- when computing the relative similarity of x_i to everything that comes before, we calculate

$$d_{ij} = \text{softmax}(\text{score}(x_i, x_j)) \quad \text{for all } j \leq i$$

$$= \frac{\exp(\text{score}(x_i, x_j))}{\sum_{k=0}^i \exp(\text{score}(x_i, x_k))} \quad \text{for all } j \leq i$$

only previous & current word

- this score, α_{ij} , gives the relative importance of x_j to x_i , and we use it to compute ^{the} a new embedding, y_i that combines different proportions of the x_j like so:

$$y_i = \sum_{j \leq i} \alpha_{ij} x_j \quad (**)$$

original x_j \rightarrow including the α_{ij} to

\rightarrow we are adding a fraction of the meaning of those other words into the original embedding.

\rightarrow this is how "bank" gets more "river" in it

\rightarrow "contextual embeddings" - makes sense on 1st layer

~~Go to next slide~~

\rightarrow nothing is being learned; ML trick is to insert parameters to "enrich" this transformation.

Now, notice that there are no learned parameters so far

- also notice that x_i gets used in 3 ways:

- 1 as the focus x_i in $\text{score}(x_i, x_j)$ "query"
- 2 as the "searched" x_j in $\text{score}(x_i, x_j)$ "key" looks backwards
- 3 to compute y_i in $(**)$ above "value"

- for all 3 cases we will insert learned parameters to modify ^{propag} these three ~~cases~~ uses as follows, e.g. for 1

let $q_i = W^Q x_i$ where $W^Q = \begin{bmatrix} w_{0,0} & \dots & w_{0,d-1} \\ \vdots & & \vdots \\ w_{i,j} & \dots & w_{i,d-1} \\ \vdots & & \vdots \\ w_{d-1,0} & \dots & w_{d-1,d-1} \end{bmatrix} \begin{bmatrix} x_i^0 \\ \vdots \\ x_i^{d-1} \end{bmatrix} = d$

\rightarrow

- this is a ML trick
- think of W^Q like a CNN kernel \rightarrow something learned across the time axis
- w.o. (lets have multiple)

- if you multiply this out you'll see that

q_i is the same vector size as $x_i = d$

- but it has been: projected/transformed by W^Q .

- the W^Q are learned parameters through gradient descent

- Similarly, there are two other learned 2d ~~mat~~ matrices W^K & W^V to get

$$K_i = W^K x_i$$

$$V_i = W^V x_i$$

together q_i , k_i & v_i "look" for patterns in the input & express the output based on them - again like 1 bit kernel in a CNN.

- So the overall computation becomes

For each input embedding x_i , compute:

$$\alpha_{ij} = \text{Score}(x_i, x_j) = \frac{q_i \cdot k_j}{\sqrt{d}} \quad \text{for all } j \leq i$$

$$y_i = \sum_{j \leq i} \alpha_{ij} v_j$$

→ scaled to keep ~~stable~~ sizes under control

↑
the output embedding correspond to x_i

- the learned matrices W^Q , W^K , W^V learn which values of the x_i are important in these roles and presumably emphasize them.

Intuition:

- When the inputs look like this \circ - W^Q & W^K then focus the output on that: W^V - maybe.

Jurabsky: "Transformers allow us to create a more sophisticated way (using the learned weights) of representing how individual words can contribute to the representation of longer inputs"

- one set of weights W^K, W^Q, W^V can be thought of looking for one set of relationships
 → a single self-attention "head"

→ could be many such relationships. that will help NN succeed. like conv kernels

→ so rather than just one set of W^K, W^Q, W^V we use several

⇒ called "Multi-head Self-Attention"

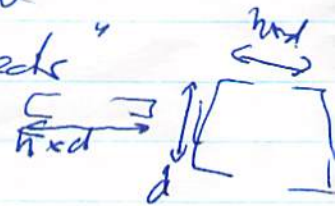
⇒ get W_i^K, W_i^Q, W_i^V $1 \leq i \leq h$ $h \leftarrow \# \text{heads}$

constant in case

in Assignment 3, use mingpt "nano" * Transformer blocks = 3 = 1-layer
* heads $h = 3 = 1\text{-head}$
 $d = 48 = n\text{-embed}$

- Every head, though, increases the size of the output → don't want larger output than input

- So also train another weight matrix W^O that is shape $hd \times d$ which "projects" the hd output back down to d . *learned parameter*



- the other parts of the transformer block are more common → see \otimes on 5-3.

(1) Skip connections → "insurance policy against failed optimization" $x \rightarrow [\text{block}] \rightarrow x + [\text{block}]$

(2) Normalization, Dropout, weight decay

Assignment A3

→ show Karpathy's
Battleship v. Speedboat.

- the code for a Transformer is too complex to write
- instead we use Karpathy's `mingpt`

- still quite complex! Configurable to large Transformers
- nicely written software, but will require careful reading, thinking & looking up PyTorch stuff.

- we have turned it into a language model + given you the code.

⇒ asked you questions about that code, to help you (prod you) to understand it.

- train it with:

① `SmallSampleCorpus.txt` Run A1.

- easier to see how Language Model works
- look at predictions.

② `LargeCorpus.txt`
→ to create pre-trained model

⇒ fine tune a classifier to do sentiment analysis

③ Use Huggingface → where the ~~red~~ models live, to also do sentiment analysis

- partially released tonight @ 9pm
- Fully released later this week.