

ECG 1786 Lecture #7

Last Day: Language Generation Using Transformers
≡ Project Ideation

Work-in-flight: Assignment 3 - tonight.
Approval + In-Principle due Thurs ^{but do sooner!}
Proposal Document } due Oct 31 @ 9pm
Proposal Slides }

Today: Understanding Transformers
Tokenization
Assignment 4
Proposal Doc / Presentations next week.
No lecture week of Nov 8 (reading week)

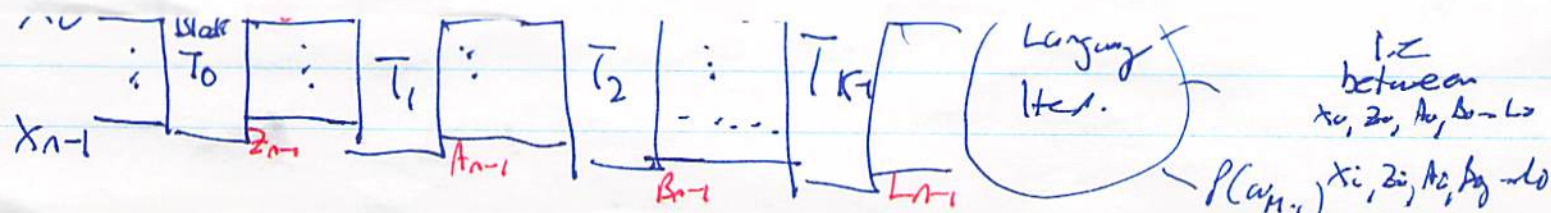
Recall, last lecture discussion ≡ demo of "zero-shot" prompting of GPT-3 - you tell it what you want it to do and it does what you asked! Remarkable!
Bad name → *few-shot* → *"priming"* → *psychology*

- a key goal for me in teaching this course was to dive in & try to understand what is happening in these models that give rise to these capabilities; and then to share it with you.
zero-shot *not all*

- for generation, I feel that a part of the explanation was given last week:
so it does what you tell it to do

- the next word generated must be consistent with *(or tries)* all the words input → the initial context + everything generated.
→ the "state" of this system is that full input context

→ I think that state "points" to higher level concepts represented in the model, and word-by-word generation brings those concepts out - e.g. fear of injury in last week's hockey example
with respect to word-by-word generation



- this understanding is just part of the story; these models are somehow very good at learning these concepts, today want to look and see what we might understand about that, ^{in part} based on lecture 5's discussion of attention.

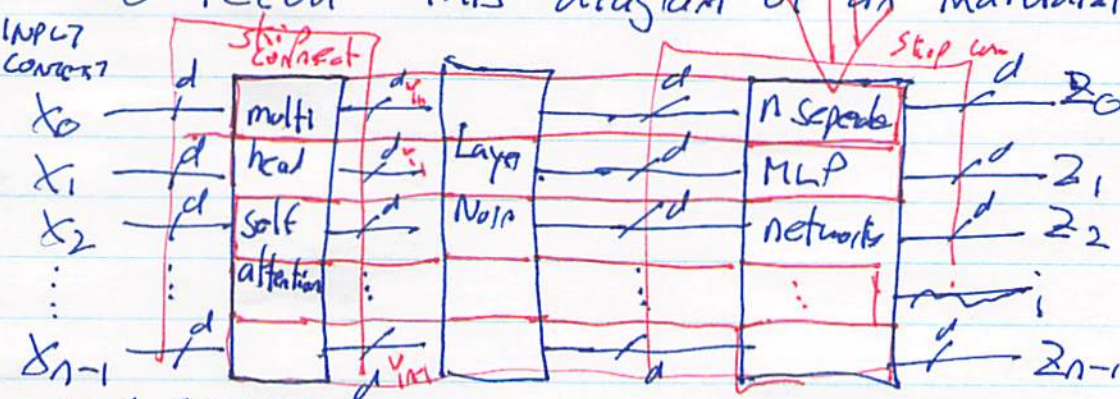
- but first, I need to correct a mistake from lecture 5 that is relevant to this discussion

TL;DR: I thought I presented the feed-forward fully connected MLP in the transformer block as one big MLP that fully connected all $n \times d$ inputs to all $n \times d$ outputs
 - that's wrong; there are actually n separate MLPs, each with d inputs and d outputs.

→ I thought the original way scrambled everything, but the ^{correct way does not} ~~scramble~~

an $n \times d$.

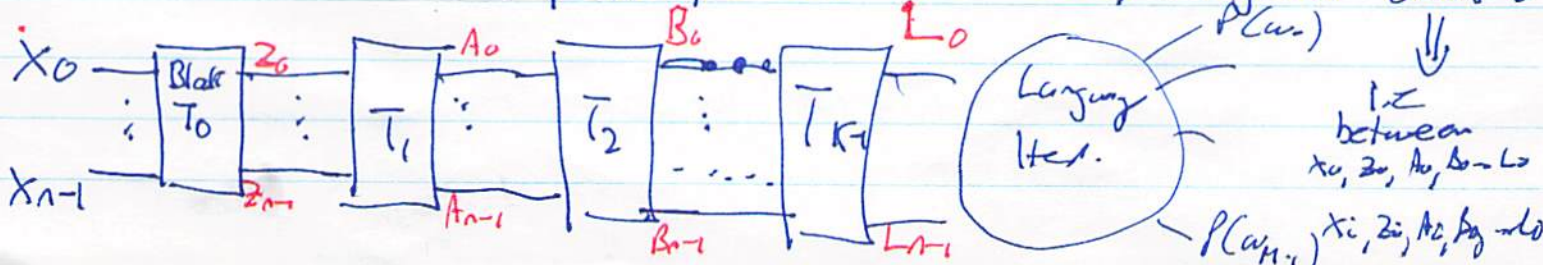
So recall this diagram of an individual Transformer Block:



- can see "mlp" in mingpt (nano) Block module in model.py

$n=2048$ for GPT-3; 1024 GPT-2

- this means that ① every red rectangle can be independently computed
 ② there is a relationship between z_i and x_i
 = perhaps all the way through all blocks



- In the literature e.g. Jurafsky text the vectors $D_i, A_i, B_i, \dots, L_i$ are what are called "contextual" vectors \rightarrow because the attention mechanism transforms the first set of vectors, x_i by looking at the surrounding x_i e.g. merging them into x_i . \Rightarrow using learned matrices W^q, W^k, W^v for each head.
 + multiple heads
- these trained weights somehow contain/are able to make those good predictions, along with the MLP
- bert viz is one way to look at them.

recall $x_0 \xrightarrow{q} v_0$ $y_i = \sum_{j < i} \alpha_{ij} v_j$

$x_1 \rightarrow y_1$

\vdots

$x_{n-1} \rightarrow y_{n-1}$

$\alpha_{ij} = \text{SOFTMAX}(g_i \cdot k_j)$

$g_i = W^q x_i$

$k_i = W^k x_i$

$v_i = W^v x_i$

- Vig's code visualizes an individual α_{ij} , showing, for a given head the link between x_i e.g. x_j after being ~~trans~~ multiplied by W^q, W^k → I think
- inputs are on the ~~right~~ left, outputs ~~left~~ right.
 in demo \rightarrow pull up colab demo.
- notice the original words are left ~~everywhere~~ x_i, y_i, z_i, k, \dots on later layers, but no guarantee that they survive intact.
- not exceedingly clear; too microscopic \Rightarrow perhaps other ways to analyze
- would like to find a way to uncover "concepts" that seem to be learned.

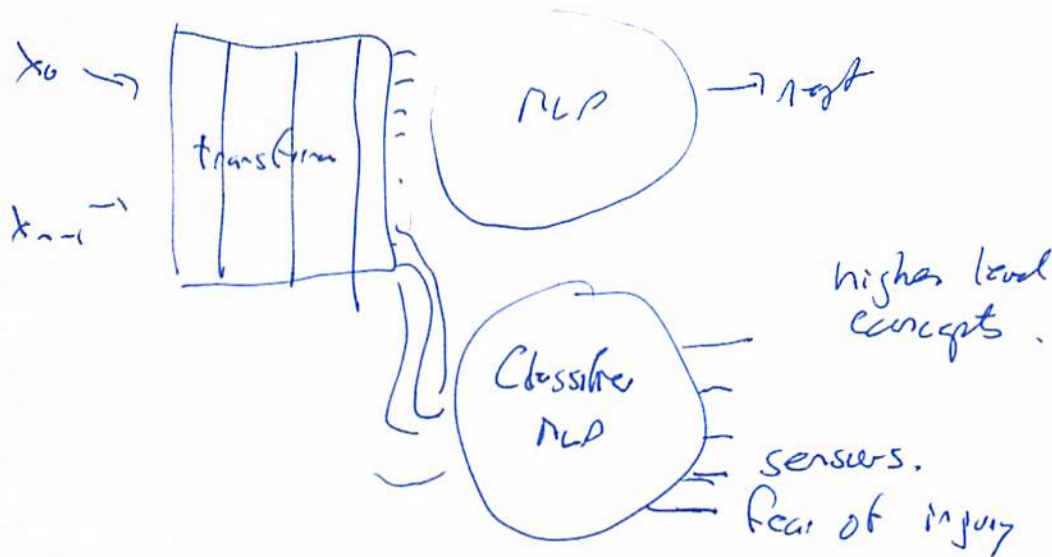
Demo from Vig's colab - post link on Quercus

Instead, let's think about higher-level concepts that are in play in GP7-3. 7-3a.

Write a creative ad for the following product on Google aimed at hockey players:

Product: A system that warns the hockey player that they are about to run into the boards.

Hockey players are always at risk of running into the boards and sustaining serious injuries. But with our new board-warning system, they can avoid these dangerous collisions! This system uses sensors to detect when a player is getting close to the boards and sends out a warning signal, giving them time to adjust their course and avoid a dangerous collision. Don't let your next game be your last! Protect yourself with our board-warning system.



→ how would we know what concepts to look for?

→ we can label examples with concepts - (create positive & negative concepts)
→ train (fine tune) classifier MLP.

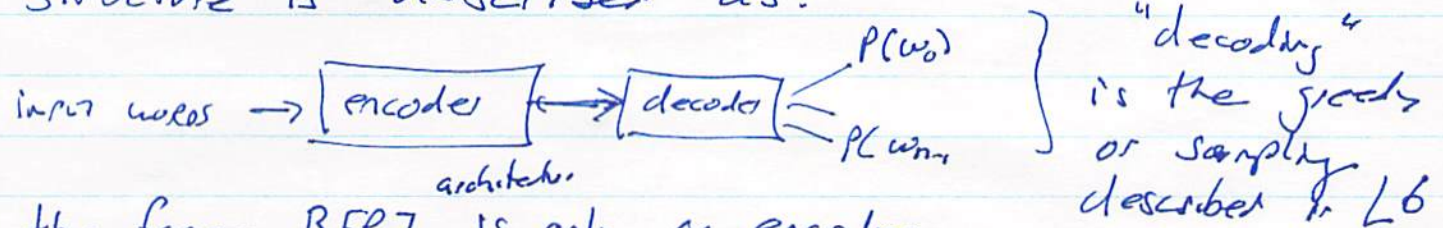
→ how do test the result?

→ with other inputs that do or do not contain the concepts

→ project?

Now, an important side point about Transformer architectures as they are commonly described.

- because of historical NLP architectures being based on RNNs, in the transformer architecture world ^{for seq-to-seq} the higher-level structure is described as:



- the famous BERT is only an encoder.
- Vashwani's original transformer was encoder-decoder.
- GPT-x is ~~said to be~~ ^{is} decoder only (works for draft)

- however, I believe one only needs 1 block (either encoder or decoder they are the same).

- you can for sure stick whatever head you want on it - classification or language gen. as shown in A3

- exactly as you did in Assignment 3.

- So why talk about both?

① The RNN history, that need encoder → decoder to do sequence → sequence

② The training is different? maybe, maybe not

↳ The essence of training BERT is to predict words missing (what are intentionally masked) in sentences

- that actually is little different than predicting the next words, as described for LTM & done in A3

- especially when context is $\gg 50$ words,
or like 2018 in GPT-3.

- what really matters, it seems is the model size
(in # parameters) & how much data it is
trained on

March 2022

\Rightarrow See Hoffman et. al. paper "Training Compute-
Optimal Large Language Models"

- bottom line GPT-3 is under-trained \rightarrow needs more
data $> 1T$.

- Chinchilla is 70B parameters vs 175B for GPT3
but was trained on 10^x data.

- is better across the board. - See paper linked
on Quercus

\Rightarrow I think once you've got a good LM, you can
use it with whichever "head" you want to
do generation or classification.

byte-pair encoding

Tokenization

\rightarrow look up BPE
in Jurafsky 2.4.3

- have not really discussed - what is it?
- OK to begin to think of words as tokens
- but some subtleties are connected to tokenization.

- however, Assign 3 + GPT-x use byte-pair
encoding which not only encodes full
words, but parts of words and even characters.
- same tokenization in A3 as GPT-1

↓
chop
words
in pieces
↓
tokens
↓
embeddings
↓
many

- the latter allows unknown words to be
tokenized. \rightarrow allows idk to work ~~something~~

→ character level tokenization may not bring meaning in its associated vectors
 - what could 'i' or 'd' or 'k' mean?

- but it does allow idk to be represented.
 - if idk shows up enough in training data then the sequence idk will have meaning encoded in the trained network maybe?

→ experimentally seems to be true:

- ① idk is not a token in GPT-3
- ② when asked GPT-3 says: the meaning of idk is
I don't know ← generated ↑ context

- by contrast, lol is a token, and it means laugh out loud, and is also "understood"

- other sub-words are produced in byte-pair encoding - e.g. hold, which could be used to tokenize "threshold"
 - it seems that sub-words are distinct from full words in tokenization.

- discuss Assignment 4 due in 3 weeks
- next week's ^{Project} Proposal → doc + presentation
- week after is reading week
→ no lecture.
- Approval - in - Principle due by Thursday
 - means I've said yes.
 - you should email me well before then in case I say no! (which I have already & means you need to email me again).
- Proposal doc + slide due Monday @ 9pm.
 - don't be late.
- Peer review: see next page.

Recall: Peer Review of Proposals

- You will be assigned another to another group to provide feedback on their written proposal and in-person presentation.
- The goal is to give you the opportunity to think critically about others' work
 - Typical and important in engineering workplaces

Peer Review Questions

1. Answer the following questions about their proposal:
 2. State the goal of the project in your own words.
 3. What was the best part of the proposal's technical contribution?
 4. What is one suggestion you have (that was not already made in class) that might help the project?
 5. Provide feedback on the quality of the oral presentation - what was good, and what could be improved.
-
- Submit PDF document of about 300 words
 - Due Friday November 4 at 6pm

