# ECE 1786 Lecture #6

<u>Last Day</u>: Core Mechanisms of Transformers ≜ Assignment 3

<u>Work-in-flight</u>: - Assignment 3
- Project Approval-in-Principle.
→ Project Proposal Document ≜ Presentation

<u>Today</u>: ① Language Generation using Transformers
② Project ideation

| Bidir vs Unidir |
| BERT vs. GPT |

## <u>Language Generation Using Transformers</u>

Recall: A language model is trained to predict the <u>next word</u> after an input sequence of words.

So, if you can do that, then you can predict a whole sequence of words, one at a time, by taking each individual predicted word that is produced, appending it to the input sequence ≜ generating the next word after that — and so on.

→ this is called "Auto-regressive" generation
(is a general term used in many kinds of sequence predictors)
- see Jurafsky, section 9.4.3, in context of RNNs.

e.g. START with: "The clean river flowed.
Generate : into.
Then next input: The clean river flowed into.
Generate : the
Next input: ~~large~~ The clean river flowed into the
Generate : ...

1. From GPT-2 medium (mid-sized GPT-2 in terms of # parameters).

   Input: The clean river flowed

   Generation Parameters (Unknown! i.e. Temperature, Top-P, etc., max tokens generated)

   Output:       The clean river flowed in steady but strong currents.

   After a week of hard-boiled fish hunt, all was tranquil again. They went out in a group on the small river and ate their dinner along the banks with some of the fish

2. From GPT-2 XL (largest GPT-2).

   Input: The clean river flowed

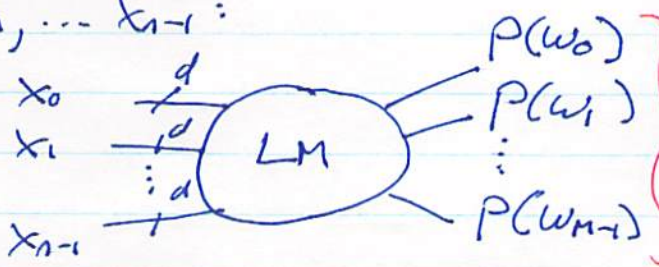   Generation Parameters (Unknown! i.e. Temperature, Top-P, etc., max tokens generated)

   Output: The clean river flowed. We walked on to the other side with the people we left behind. We found a small restaurant with a bench near the river — a small oasis at the end of the world, really.

   In the evening we

— see examples from GPT-2 on previous page [read]

— recall the specific output of the language model, given an input sequence of embeddings:
$x_0, x_1, \ldots x_{n-1}$:

$x_0 \xrightarrow{d}$
$x_1 \xrightarrow{d}$ **LM**
$\vdots \; d$
$x_{n-1} \xrightarrow{d}$

$\longrightarrow P(w_0)$
$\longrightarrow P(w_1)$
$\vdots$
$\longrightarrow P(w_{M-1})$

where M is the size of the vocabulary.

the probability that each word in vocabulary is the next word

→ DOESN'T CHOOSE THE WORD.

→ not related to decoder in model

→ unfortunate choice.

— So for a given sequence input, which word is selected as output? [This] is called **Decoding**

**N.B.** ⟹ WHAT IS THE BEST ANSWER? → THINK CAREFULLY — WHAT ARE WE LOOKING FOR?

① "Greedy" method: select the highest probability word
   — [you'll see this in Assignment 3 part 2.]
   — does not work well in general → obvious, but boring/un-interesting *sequences of* words are chosen; repetitive.
   → may choose the most likely <u>next</u> word, but does <u>not</u> result in <u>the most likely sequence</u> of generated words
        — gets stuck in a "highly" local optimum

i.e. given input sequence $x_0 \ldots x_{n-1}$ want the generated sequence $y_0 \ldots y_{k+1}$ of $k$ words to be most likely.
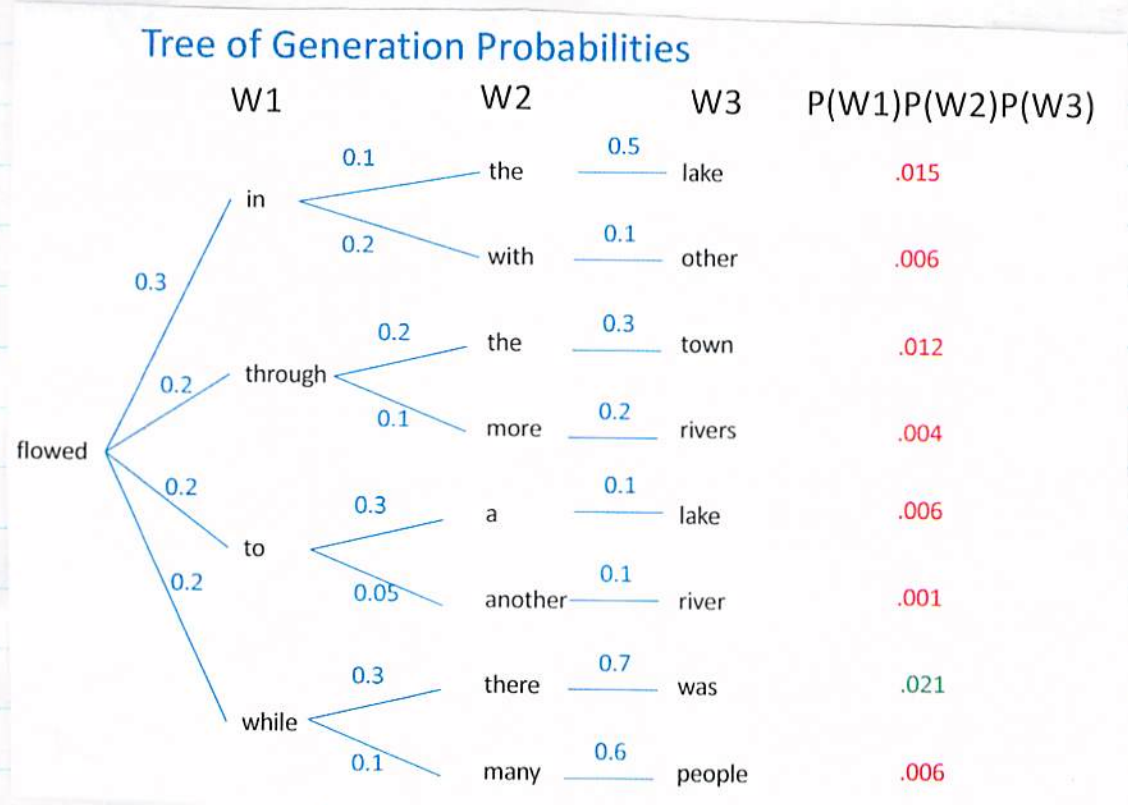   i.e. want $P(y_0) \times P(y_1) \times \cdots \times P(y_{k-1})$, to be maximized.
but we don't know $P(y_1)$ when selecting $y_0$ (or $y_j$ when selecting $y_i$)   UGLY SEARCH PROBLEM.
   — this is a tough problem because there are
   $M^k$ possible solutions *sequences* OR $K$ *words* (e.g $M = \cancel{5000} 50000$, $k = 20$, $50000^{20}$
   is a big number of large NN *model* inferences to make!)
        → i.e. way too many

- per Jurafsky Section 10.5, think of the $\overset{\text{AUTO}}{\text{REGRESSIO}}$ selection of the output sequence as a tree with probabilities at each layer

Input: The clean river flowed....
Outputs, produced from <u>many</u> invocations of inference $\overset{\text{using}}{\text{trained}}$ model:

give probabilities are made up.

### Tree of Generation Probabilities

| W1 | | W2 | | W3 | P(W1)P(W2)P(W3) |
|----|----|----|----|----|----|
| | 0.1 | the | 0.5 | lake | .015 |
| in | 0.2 | with | 0.1 | other | .006 |
| through 0.2 | 0.2 | the | 0.3 | town | .012 |
| | 0.1 | more | 0.2 | rivers | .004 |
| to 0.2 | 0.3 | a | 0.1 | lake | .006 |
| | 0.05 | another | 0.1 | river | .001 |
| while | 0.3 | there | 0.7 | was | .021 |
| | 0.1 | many | 0.6 | people | .006 |

the clean river flowed   (0.3 in, 0.2 through, 0.2 to, 0.2 while)

- so the most probable W1 — "in" doesn't lead to the most probable sequence of W1 W2 W3.
  $\overset{.015}{\underset{.01 .006}{}}$
- to get the optimal prob sequence is very hard. vs. .021

Method

(2) Beam Search is a heuristic that prunes the full search tree a lot.
- general description: walk ~~through~~ down the tree keeping the K-most probable sequences.
- at each level, look at V possible new words for each of the K sequences

- that produces K $*$ V new sequences
  - compute the probability for all by multiply probabilities.
- keep the K highest.

→ GO DOWN TO NEXT LEVEL.

- each of the K $*$ V a each level <u>costs</u>
  one inference run through model → <u>expensive</u>!
  <span style="font-size:small">in time.</span>

- was thought to be best method, but isn't used it seems
- instead:

③ ~~Nucleus~~ Sampling.
 - given the set of output probabilities    M = vocab size
   of the next word — $P(w_0), P(w_1) \dots P(w_{M-1})$

 - select the next word through a <u>random</u> process
   in which the probability of selecting
   word $w_i$ <u>is</u> $P(w_i)$ — how?

 - how is that done? Show by example:
   suppose there are just 3 words with probabilities
   for the next word:

   UP     $P(up) = 0.5$           then: generate a <u>uniformly</u>
   down   $P(down) = 0.3$           distributed random
   left   $P(left) = 0.2$          <u>number</u> between
                                    0 and 1.0
   
   | 0 | | .5 | .8 | 1.0 |
   | --- | --- | --- | --- | --- |
   | | .5 | .3 | .2 | |
   | | UP | down | left | |

                                    - call it R

 - if $R \leq 0.5$ choose up          so  UP will be
   $0.5 < R \leq 0.8$ choose down       chosen with probability
   $0.8 < R \leq 1.0$ choose left       0.5; down with prob .3
                                        & left with prob 0.2

– Simplification: only select from the top K most
         probable words (top-k sampling)

– most widely used: top-p sampling: only
select from the 'top' words that all together
have the sum of probabilities $\geq p$     ($0 < p \leq 1$)
($p=1$ means use all; often $p > 0.8$ vastly
reduces the number considered)

– also a really unusual adjustment to the
sampling method involves an adjustment to
the softmax probability calculation @ end of each inference
   – recall $P(w_i)$ = softmax ( logits ( = network output).)
let $l_i$ be the logits output of word $i$.

$$P(w_i) = \frac{\exp\left(\frac{l_i}{t}\right)}{\sum_{all\ i=0,\ n-1} \exp\left(\frac{l_i}{t}\right)}$$

where $t$ is
a parameter
call __Temperature__
$0 < t$   ↑NOT

– $t=1$ is normal sampling; $t$ close to zero is great
– $t < 1$ makes the higher probability words __more__ likely
   (because exponential works against smaller probs).
→ $t$ closer (or greater than) 1 makes result __more__ diverse → the
   less likely words become more likely.
                             _greatly_

→ often use a combination of top-p + temperature.
→ +1 more repetition penalty - divide e.g. by [1.3] if word already used.

[ – Demonstration using GPT-3 playground.
  – show code in mingpt generate function
       – temperature
       – torch.multinomial | torch.topk.

- demo of GPT-3: basic continuation of words
  - show effect of T, P, $ probabilities.

Now,
- GPT3 is a 175B parameter model trained on
  ≈ 1 Trillion words
  $\rightarrow$ IT DOES WHAT IT IS TOLD TO DO.
- It is capable of "zero-shot" learning (a bit of a misnomer)
  meaning that it can do tasks that the model parameters are
  not trained on. How can that be?

- e.g The ~~class is listening intently~~ deaf river flowed.
                                                    — tokens are
                                                       interesting!
- hockey. ad.                                       → click on
                                                       them.
- write a poem / classify.        — note different settings.

- I find zero-shot truly remarkable. How could it work?
  Is the model writing a computer program to do what is asked?
   - don't think so; struggled to explain, but do have this
- . it is all based on the ability to predict the next word.
- to do so well, one has to know alot.
  e.g how to end: "The car was going faster than a___"
   - need to know alot

- with very powerful predictor, the language at the
  front becomes an instruction → the predicted
  words must be consistent with the front instruction
  & so the model does as asked if it ca..
- each time inference is run, those words are
  considered: the whole context is the state