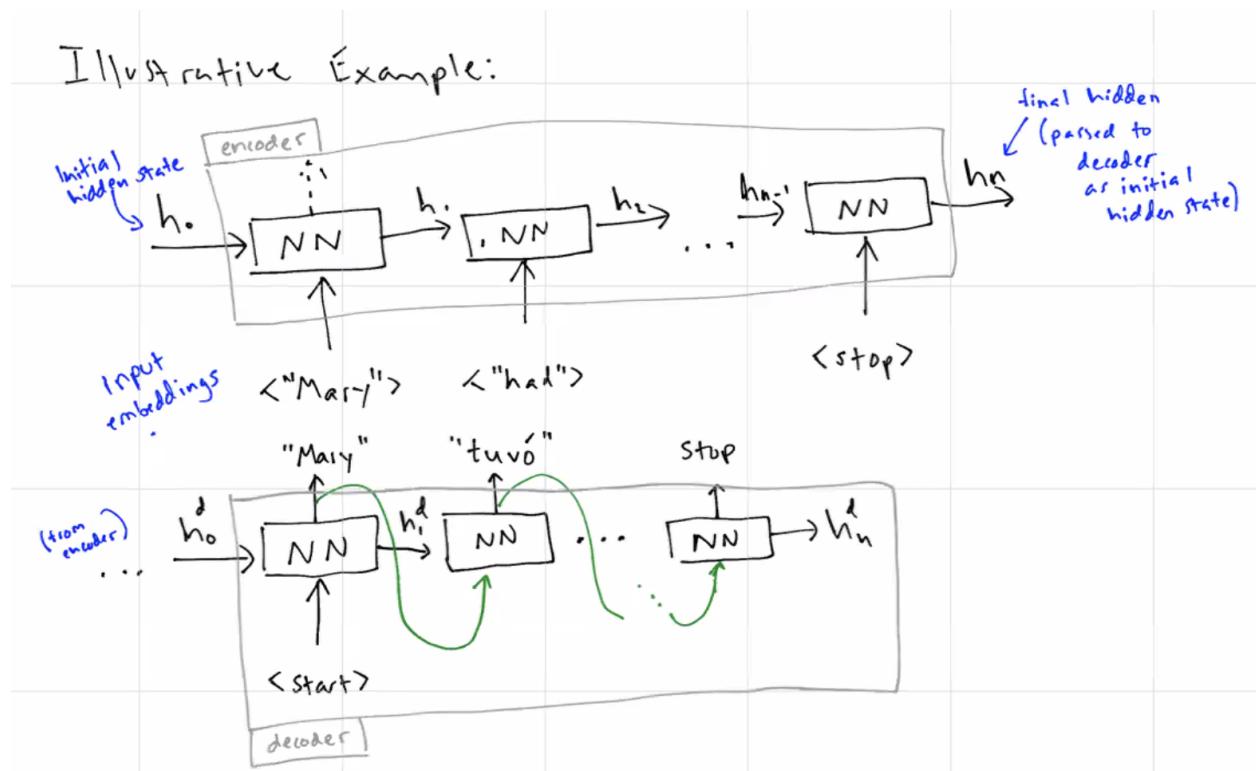


Day 12: Recurrent Neural Networks (continued)

Practical Challenges of RNNs

1. Non-parallel operations that scale with sequence length
2. Longer sequences require deeper networks
3. Fixed size hidden state to support long-range dependencies (relationships between words that are far apart)

"The cat was a wild one... He liked to..."



Above: recurrent NN encoder & decoder for translation - with optional teacher forcing while training decoder

One solution: Design specialized recurrent units that are better at remembering hidden state info.

① LSTM (Long short-term memory [unit])

Hochreiter & Schmidhuber, 1997

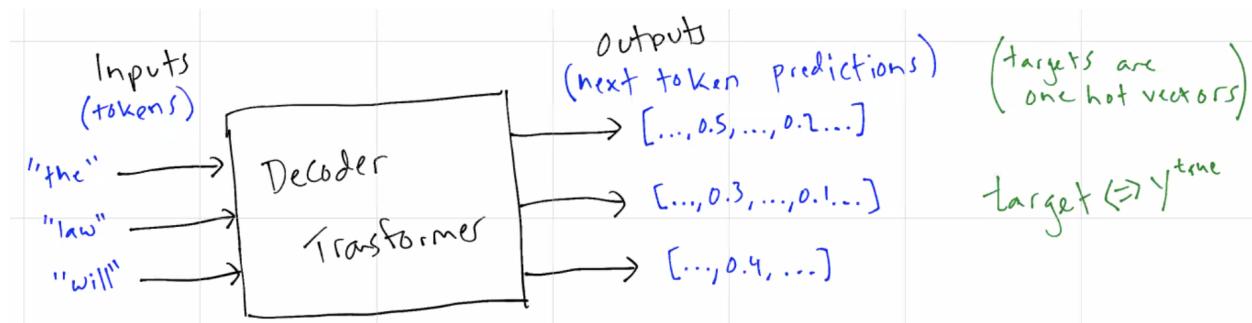
② GRU (Gated recurrent units)

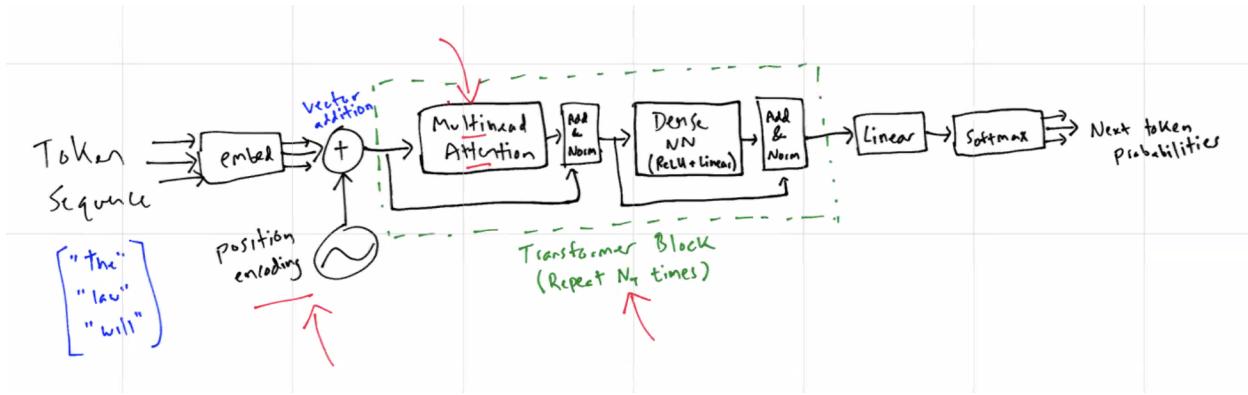
Kyunghyun Cho et al. 2014

Fundamentally, difficulties remain → can one create a non-recurrent architecture for sequential data?

Anatomy of a Decoder Transformer

- Inputs are individual tokens fed into decoder in parallel
- Outputs are pseudo-probabilities (next token predictions) over the next token set
- Attention/Transformer block “learns” what words/phrases are most important to be considered





- Address concerns of RNNs:

1. Fixed number non-parallelizable steps
 2. Fixed network depth
 3. No language dependencies
1. Decoder Transformer Objective: predict next token given part-tokens
 - a. More explicitly, let network model the following probability function:

$$F(\vec{w}; \vec{x} = [x_{t-1}, x_{t-2}, \dots, x_0]) \approx P(x_t | x_{t-1}, x_{t-2}, \dots, x_0)$$

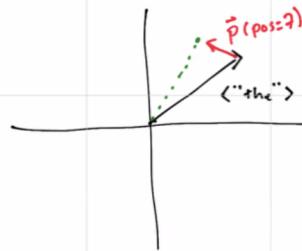
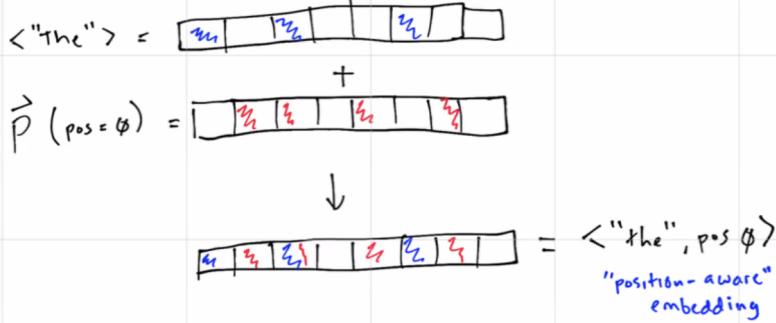
Annotations for the equation:

- \vec{w} : network trainable parameters
- $\vec{x} = [x_{t-1}, x_{t-2}, \dots, x_0]$: sequence of tokens
- x_0 : most recent token
- x_{t-1} : first token in sequence
- $P(x_t | x_{t-1}, x_{t-2}, \dots, x_0)$: approximate this distribution

$$F(\vec{w}; \vec{x} = [x_{t-1}, x_{t-2}, \dots, x_0]) \approx P(x_t | x_{t-1}, x_{t-2}, \dots, x_0)$$

- b. We rely on a cross-entropy loss function:

$$L(X = [\vec{x}^0, \vec{x}^1, \dots, \vec{x}^N]) = - \sum_i \log(F(\vec{W}_1 \vec{x}^i))$$



$$\approx - \sum_i \log(P(x_t^i | x_{t-1}^i, \dots, x_o))$$

$$= \sum_i \log(P_{c=c_{true}}^{\text{pred}})$$

Addendum: Cross-entropy term

Recall:

$$\mathcal{L}_{CE} = - \sum_{j=0}^{C-1} P_i^{\text{true}} \log(P_i^{\text{pred}}) \quad (\text{for 2 discrete prob dists})$$

) let only one P_i^{true} be 1

$$= - \log(P_{c=c_{true}}^{\text{pred}})$$

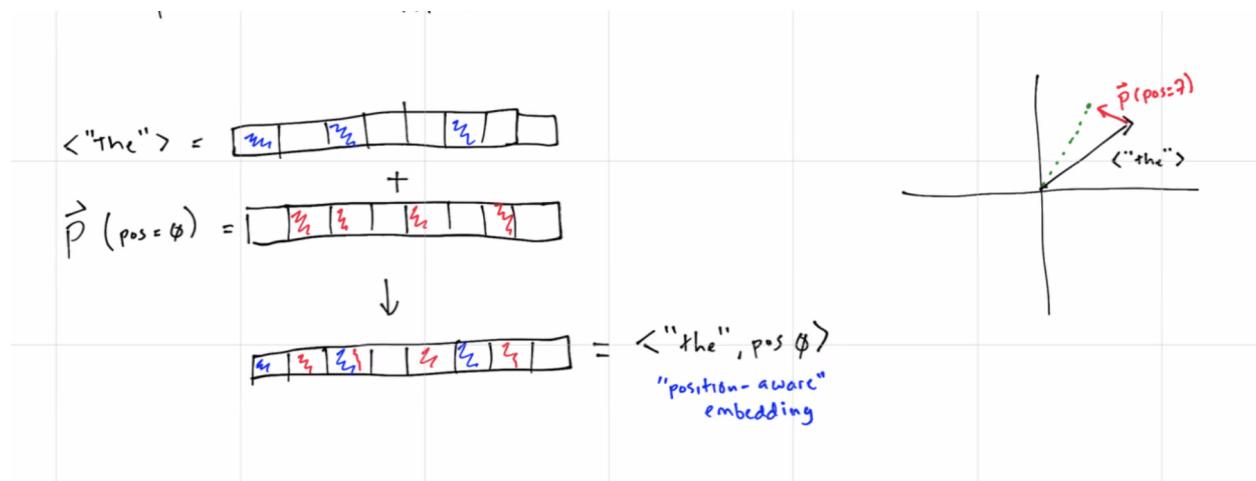
) sum over examples in some dataset

$$= - \sum_{i=0}^{\text{Nexamples}} \log(P_{c=c_{true}}^{\text{pred}})$$

- Outline
 1. Positional encoding
 2. Attention ("masked", "multiheaded", "sequential")
 3. How can we use decoder transformer architecture outside of text generation

Positional Encoding

Concept: Add position specific vector each input embedding to impart position information



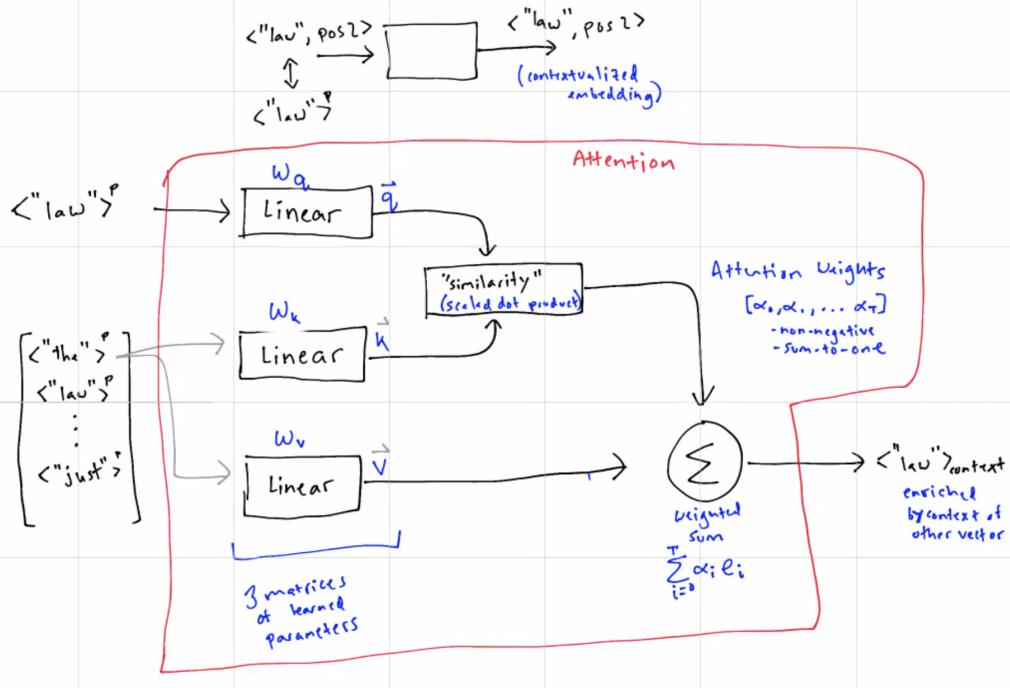
Are we ok to scatter info of index into vector space? Isn't it information inefficient?

- Approaches to calculate
 1. Learn positional embedding matrix
 - a. The "position-aware" embeddings are calculated as the sum off the word embeddings \vec{e}_i and position vectors \vec{p}_i (the latter of which is extracted from the learned function weights)
 2. Prescribe a function

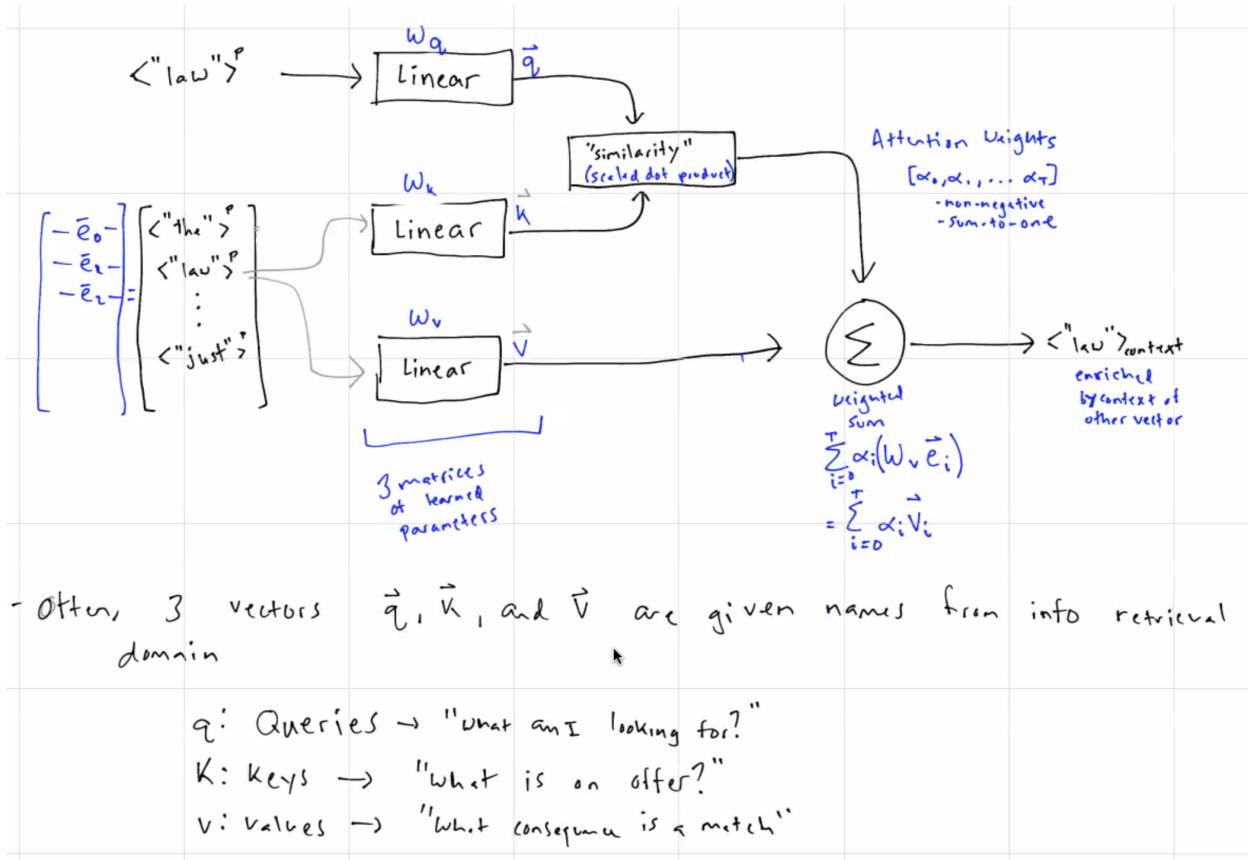
Attention Mechanisms

- Concept: Let network learn "contextualized" embedding for each input as a weighted sum of non-contextualized embeddings
- 3 vectors, \vec{q} , \vec{k} , \vec{v} are given name from info retrieval domain
 - \vec{q} : Queries → "What am I looking for?"
 - \vec{k} : Keys → "What if "

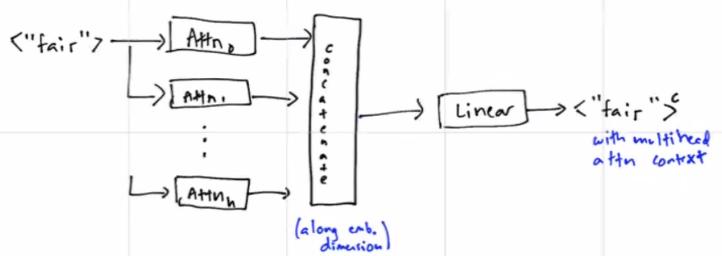
Concept: Let network learn a weighted sum "contextualized" embedding for each input as of non-contextualized embeddings



Circled in red above: the mechanisms within a "head" or attention block



Multihead Attention (basically multiple attention blocks ["heads"] running in parallel)



- Each attention block can specialize (analogous to CNN filters)

- Sequential Attention

- Single attention block contextualizes inputs with pair associations

- 2 sequential attention blocks contextualize inputs with pair-of-pair associations

- j Sequential " " " " " " z^j word associations

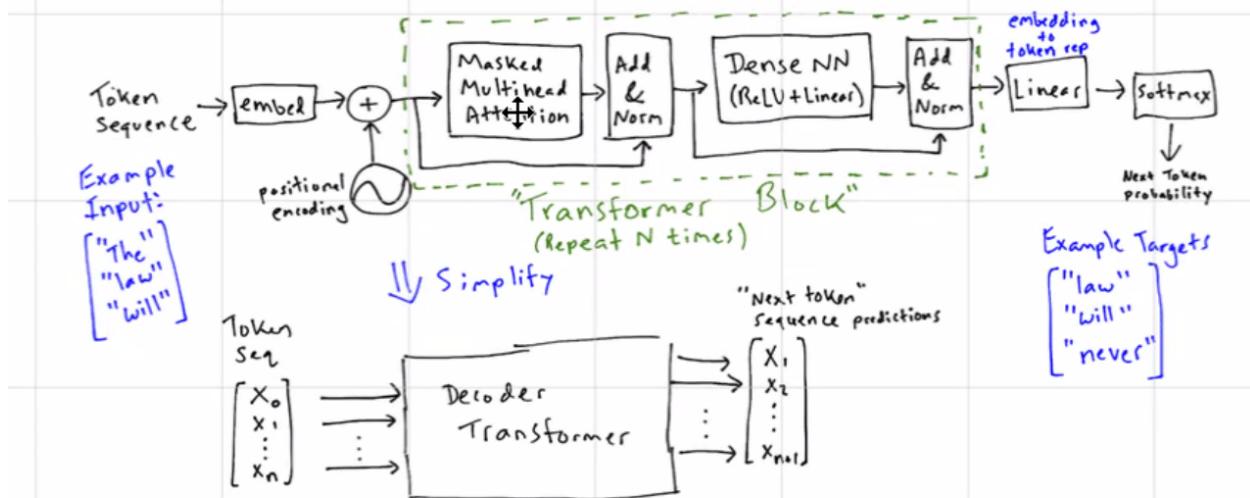
Hope that the different "heads" learn different contexts that may be helpful when used together

Example: "Jeff is a musician. Jeff went to the park, where he _____"

One head might make "he" pay attention to the earlier "Jeff" and another head might make "he" pay attention to "musician," where combined, they might help determine that the next word is "played".

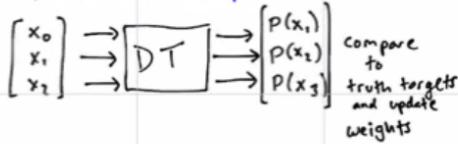
"Predictions for position i can only depend on positions less than i ."

Training and Generation Procedures



Training Procedure

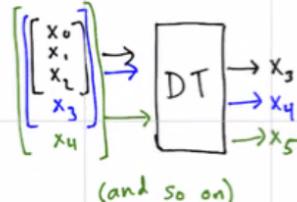
- Compare full output sequence to targets and update weights to reduce cross-entropy loss (over full text corpus)



compare to truth targets and update weights

Generation Procedure

(weights static)



(and so on)

How can we use the model for tasks beyond text generation?

A: We can prompt model for what we want using language

① News article generation

"Title: flying Saucers visit White House

Subtitle: Onlookers shocked by alien life forms

... [let model complete] "

② Conversational Question Answering

"Mary had a little lamb, her fleece was filthy

Question: What color was Mary's lamb?

Answer: ... [let model complete]

③ Grammar Correction

"poor English grammar: I eated purple berries

good English grammar: I ate purple berries

poor English grammar: The patient was died

good English grammar: ... [model completes]

Examples of other use cases

Network Sizes and Training Costs:

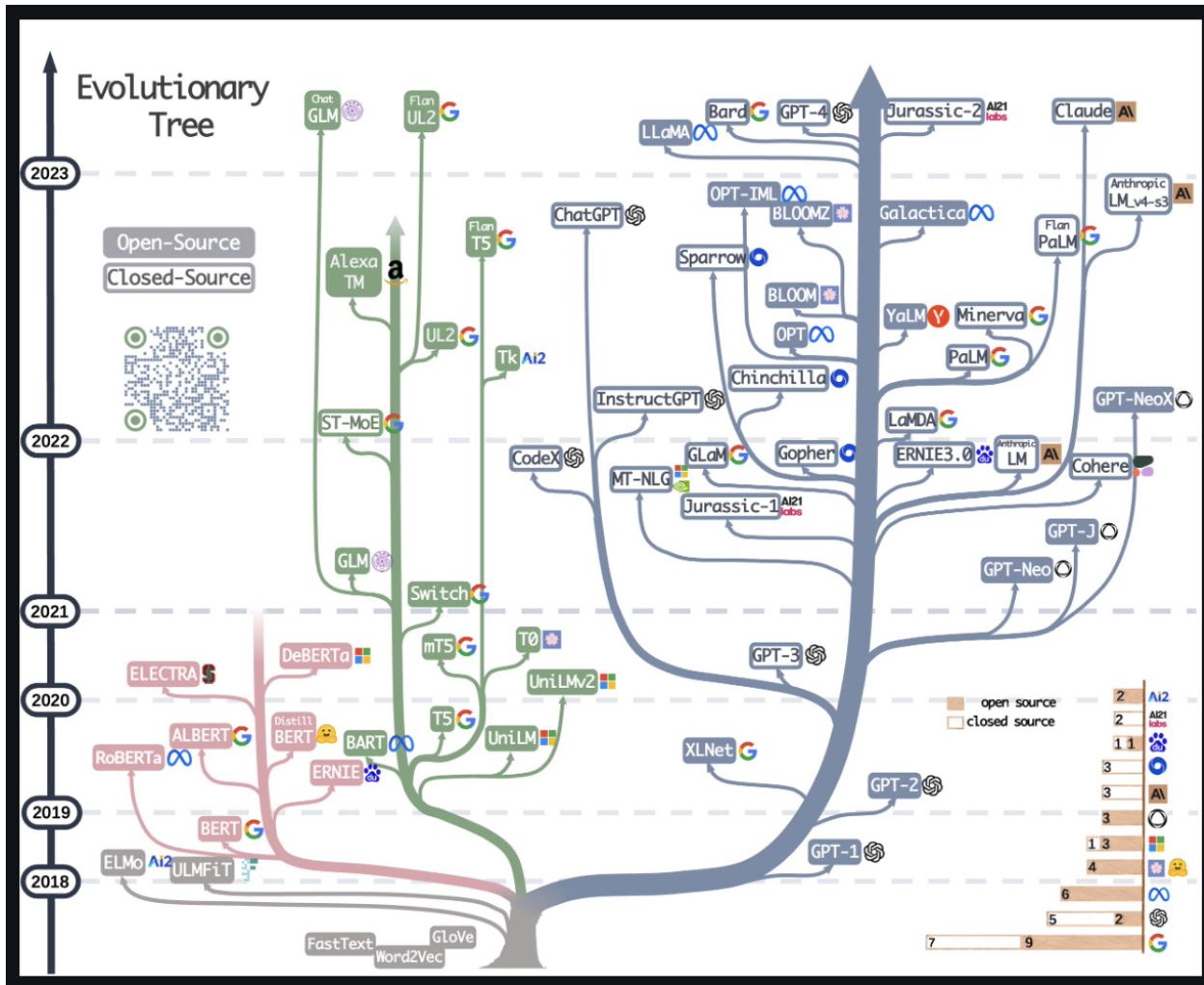
| (* Data from "Generative Pre-trained Transformer" wikipedia page) | | | | FLOPs | Effective* train time |
|---|----------------|-----------------------|--|----------------------|------------------------|
| <u>date</u> | <u># param</u> | <u>dataset (size)</u> | | | |
| GPT-1 | June, 2018 | ~100 million | BooksCorpus (4.5 GB) | 1.7×10^{19} | 1 month |
| GPT-2 | Feb, 2019 | 1.3 billion | WebText (40 GB text) | 1.5×10^{21} | ~100 months |
| GPT-3 | May, 2020 | 175 billion | CommonCrawl WebText English Wikipedia Books1 & Books2 (100s of GB) | 3.1×10^{23} | ~10000 months (~1000y) |
| GPT-3.5 | March, 2022 | 175 billion | Undisclosed | Undisclosed | Undisclosed |
| GPT-4 | March, 2023 | undisclosed | undisclosed | undisclosed | undisclosed |

* assume computer FLOPs used on GPT-1, i.e. $\frac{1.7 \times 10^{19} \text{ flop}}{1 \text{ month}} = \sim 6 \text{ peta-flops}$

Use of Large Language Models

Wednesday, July 26, 2023 3:52 PM

1. baseline: what are the current strengths and weaknesses of LLMs?
2. list of ≥ 3 areas that are likely to be disrupted by LLMs
3. how should people in your area use/not use LLMs? are there general principles you can extract from your area to apply more broadly?



Exam results (ordered by GPT-3.5 performance)

Estimated percentile lower bound (among test

