Introduction to Nature Language Processing 2nd lecture

**Boston University** 

CS 640, AI

Slides by Margrit Betke, Yiwen Gu



#### Outline (Recap)

- What is NLP?
- Key NLP tasks
  - Text Classification, NER, MT, Sentiment Analysis, QA, Summarization
- NLP Techniques and Approaches
  - Traditional techs,
  - RNNs
  - Attention Mechanism & Transformers



#### **RNNs:** Limitations

- Difficulty in Capturing Long-Term Dependencies: Even with improvements like LSTMs and GRUs, RNNs can struggle with very long sequences.
- Sequential Computation: RNNs process one time step at a time, which makes training slower, especially for long sequences.

→ Transformers: RNNs have been largely supplanted by transformers, which use self-attention mechanisms to process entire sequences in parallel, making them more efficient and better at capturing long-range dependencies.



#### Transformers





#### Transformers

























Attention(Q,K,V) = softmax(Q  $K^T$ ) V

Q = query vector = current English (or French) word

K key and V value = memory of words seen before

Goal: Find key(s) most similar to query and retrieve value(s) that correspond to this/these key(s)

Softmax =  $\sum_{i} e^{qki} / \sum_{j} e^{qkj} v_{j}$  produces probability distribution over keys with peaks for keys similar to query



Attention(Q,K,V) = softmax(Q K<sup>T</sup>) V

Acts as a weight mask over V

Q = query vector = current English (or French) word

K key and V value = memory of words seen before

Goal: Find key(s) most similar to query and retrieve value(s) that correspond to this/these key(s)

Softmax =  $\sum_{i} e^{qki} / \sum_{j} e^{qkj} v_{j}$  produces probability distribution over keys with peaks for keys similar to query



Very fast: 2 matrix multiplications & 1 softmax operation

Attention(Q,K,V) = softmax(Q K<sup>T</sup>) V

Acts as a weight mask over V

Q = query vector = current English (or French) word

K key and V value = memory of words seen before

Goal: Find key(s) most similar to query and retrieve value(s) that correspond to this/these key(s)

Softmax =  $\sum_{i} e^{qki} / \sum_{j} e^{qki} = v_{i}$  produces probability distribution over keys with peaks for keys similar to query



#### Very fast: 2 matrix multiplications & 1 softmax operation

Attention(Q,K,V) =  $softmax(Q K^T / sqrt(d_k)) V$ 

Acts as a weight mask over V

Technical detail: sqrt(d<sub>k</sub>) normalization needed for training

Q = query vector = current English (or French) word

K key and V value = memory of words seen before

Goal: Find key(s) most similar to query and retrieve value(s) that correspond to this/these key(s)

Softmax =  $\sum_{i} e^{qki} / \sum_{j} e^{qkj} v_{j}$  produces probability distribution over keys with peaks for keys similar to query



 $\operatorname{Attention}(Q,K,V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$ 

How to get Q,K,V?





# $\operatorname{Attention}(Q,K,V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$

How to get Q,K,V? In the Matrix Form





#### Multi-Head Self-Attention

Attention
$$(Q, K, V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

 $\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, ..., \text{head}_h) W^O \\ \text{where head}_i &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned}$ 

## **Question:**

Known: Number of heads: h=8, Model dim = Embedding dim: d=512, Weight matrices for K, Q, V are in the shape of 512 x 64, What is the shape of  $W^o$ ? d<sub>k</sub>

64\*8 x 512 See explanation on the next slide



#### Multi-Head Self-Attention

1) This is our2) We embedinput sentence\*each word\*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting Q/K/V matrices 5) Concatenate the resulting Z matrices, then multiply with weight matrix  $W^{0}$  to produce the output of the layer



\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



X





In this example, n = 2 (not relevant to  $W^o$ ), h = 8 d = 4  $d_k = 3$ The shape of the  $W^o$  is drawn accordingly, i.e. (3x8, 4)



#### Why Multi-Head Attention?

- Multiple attention layers (heads) in paraellel
- Each head uses different linear transformation
- Different heads can learn different relationships



#### Visualization of Attentions





Figure 3: An example of the attention mechanism following **long-distance dependencies** in the encoder selfattention in layer 5 of 6. Many of the attention heads attend to a distant dependency of the verb 'making', completing the phrase 'making...more difficult'. Attentions here shown only for the word 'making'. Different colors represent different heads. Best viewed in color.

cial Intelligence, 2024



#### • <u>Feed Forward</u>:

 Position-wise Feed-Forward Networks, process every word

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

• <u>Where facts in LLM live</u> (3Blue1Brown)

- <u>Add & Norm</u>:
  - Residual Connection and LayerNorm
- <u>N x</u>: Stacking
  - Stacking these attention models also allows for higher level reasoning
  - Lower layers focus more on word relationships and syntax, whereas higher layers more on contextual relationships and semantics





- Encoder-Decoder Attention (Cross-Attention)
  - This cross-attention allows the decoder to attend to the output of the encoder, combining information from both the input sequence (processed by the encoder) and the current decoder states.
  - Used in Visual Transformers (ViTs)  $\rightarrow$  Later
- Masked Self-Attention:
  - Each token can **only** attend to previous tokens in the sequence.
  - This prevents the decoder from "cheating" by looking at tokens it hasn't generated yet, ensuring autoregressive generation.



#### Last Step linear + softmax, (sampling)

#### output token probabilities (logits)





#### Putting Together

Decoding time step: 1 (2) 3 4 5 6

OUTPUT





#### Putting Together

Decoding time step: 1 (2) 3 4 5 6

OUTPUT







Figure 1: The Transformer - model architecture.

#### Note:

The training dataset (WMT) consists of parallel corpora from multiple languages. This means **that the source and the target are paired in the sentence level during the training**. In this way, the source sentence is passed through the encoder to get some representation; the target sentence is used during training to guide the decoder in generating the correct translation, with the model learning to minimize the difference between its output and the actual target sentence. During inference (when the model is actually translating sth.), only the source sentence is provided. And the model generates the target sentence word by word, based on its learned patterns from the paired training data.



# Transformer Architecture Complexity (per layer)

- n= number of words in sequence
- d= network dimension

Number of operations:  $n^2 d$ Number of activations:  $n^2 + n d$ 

Much better than RNNs with number of operations n d<sup>2</sup>



# Transformer Architecture Complexity (per layer)

- n= number of words in sequence (<70 words per sentence)
- d= network dimension (easily go beyond 1000 in a transformer)

Every word attends to every word To get the dot product of Q, K (both of size **d/h**), then multiply by h for multi-head

Number of operations: n<sup>2</sup> d

e.g., 70x70x1000=4.9 mill

Number of activations:  $n^2 + n d$  From FFN, ReLU

From self-attention, softmax

Much better than RNNs with number of operations n d<sup>2</sup>

e.g., 70x1000x1000=70 mill



#### Training a Transformer

- ADAM optimizer
- Dropout during training at every layer
- Label smoothing
- Auto-regressive decoding with beam-search
- Checkpoint-averaging
- Library available: <a href="https://github.com/tensorflow/tensor2tensor">https://github.com/tensorflow/tensor2tensor</a>



# How much data to train?



CS 640: Artificial Intelligence, 2024

## All of it...



CS 640: Artificial Intelligence, 2024

## All text on the internet?

Is that legal?

AI & Ethics!





S<sup>®</sup> World ∨ Business ∨ Markets ∨ Sustainability ∨ Legal ∨ Breakingviews ∨ Technology ∨ Investic

## All text on the internet?

Is that legal?

AI & Ethics!

Litigation | Copyright | Litigation | Technology | Intellectual Property

#### John Grisham, other top US authors sue OpenAl over copyrights

By Blake Brittain

September 21, 2023 6:34 AM EDT · Updated 7 months ago







## All text on the internet?

December 27, 2023

Is that legal?

AI & Ethics!

### The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work

Millions of articles from The New York Times were used to train chatbots that now compete with it, the lawsuit said.



## All text on the internet?

Is that legal?

AI & Ethics!





96 layers (decoder blocks) 2048 tokens

## Training the 175 billion parameters of GPT-3 on "all text on the internet" on a single GPU or computer would take 355 years and \$4,600,000

Lambdalabs.com



CS 640: Artificial Intelligence, 2024

## How long did it take OpenAl to train GPT-3?

#### a month



CS 640: Artificial Intelligence, 2024

### What did OpenAl train on?

GP1-3 training data <sup>rtys</sup>		
Dataset	# tokens	Proportion within training
Common Crawl	410 billion	60%
WebText2	19 billion	22%
Books1	12 billion	8%
Books2	55 billion	8%
Wikipedia	3 billion	3%

CDT 2 training data [1].9

Source: Wikipedia



#### What about GPT-4?

- 1 trillion parameters
- Sam Altman stated that the cost of training GPT-4 was more than \$100 million.

Source: Wikipedia



CS 640: Artificial Intelligence, 2024

#### What about GPT-4?

- 1 trillion parameters
- Sam Altman stated that the cost of training GPT-4 was more than \$100 million.

#### Why are the lawsuits so costly?

- We don't know how to "untrain" neural networks.
- "Unlearning" is an exciting research area!

Source: Wikipedia



#### Learning Outcomes

- Understand how transformers work
  - Encoder-Decoder (e.g. THE Transformer, BERT)
  - Decoder Only (e.g. GPTs)
  - Attention Mechanism
    - self-attention,
    - Masked self-attention,
    - cross-attention

