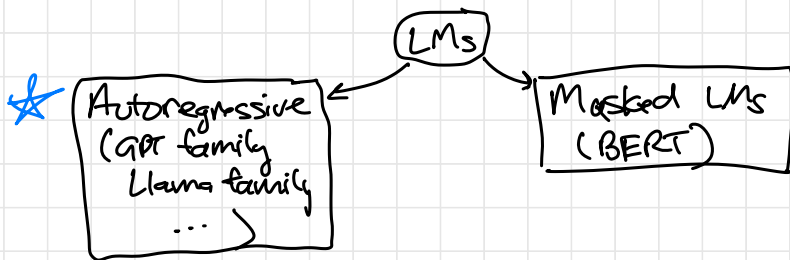


# 8/28/2024: Transformer LMs

1. Autoregressive LM
2. Using Transformers as LMs



Autoregressive LM: Takes input =  $w_1, \dots, w_T$  ( $T$  tokens)  
Output: Distribution over the text token  $w_{T+1}$   
*represent first  $T$  tokens of a document*

$$= P_{LM}(w_{T+1} \mid \underbrace{w_1, \dots, w_T}_{\text{prefix}}; \theta)$$

$\uparrow$  model parameters

**Big Assumption**: All tokens come from a fixed set  $V$  ("vocabulary")

Under this assumption, this distribution is over  $|V|$  possibilities  
 $\Rightarrow$  represent it as vector of dimension  $|V|$

Note: Implicitly, this defines probability distribution over longer sequences.  
For any string  $S = w_1, \dots, w_T$ :

$$P_{LM}(S; \theta) = \prod_{t=1}^T p(w_t \mid w_1, \dots, w_{t-1}; \theta)$$

**How to train?** To pretrain an autoregressive LM:

Need dataset  $D$  of many strings  $S = S_1, \dots, S_T$

Choose  $\theta$  to minimize

$$L(\theta) = \frac{1}{|D|} \sum_{S \in D} \frac{1}{T} \sum_{t=1}^T -\log P_{LM}(S_t \mid S_1, \dots, S_{t-1})$$

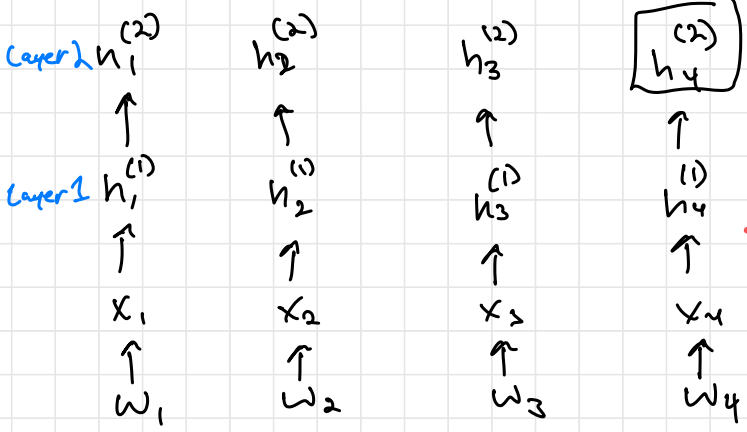
*every string* · *every token in string* · *log prob of the real next token*

By stochastic gradient descent

# Transformers

predicted  $p(w_5 | w_1, \dots, w_4)$

✓ (3) Generate final prediction = vector in  $\mathbb{R}^{|V|}$



✓ (2) Generate hidden states for layers 1, 2, ..., L  $\in \mathbb{R}^d$

✓ (1) Generate token embedding  $\in \mathbb{R}^d$

Input tokens  $\in V$

## Step 1: Token Embeddings

What is a token?

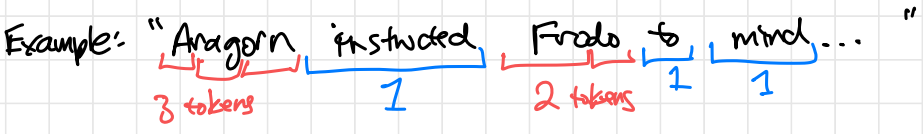
Problem with word-level: A LOT of possible words

- New names
- Typos

Massive Vocabulary  
"unknown" word problem

Either a word or part of a word!  
Uncommon words will be split into multiple tokens!

subword tokenization



Common strategy for subword tokenization  $\rightarrow$  Byte Pair Encoding (BPE)  
Sentence Piece tokenization

How to map tokens  $w_1, \dots, w_T$  to vectors  $x_1, \dots, x_T$

One common strategy: Absolute positional embeddings

Goal: each  $x_t$  must encode 2 things:

- Identity of  $t$ -th token  $w_t$
- Order of all the tokens

Solution: add together 2 vectors,  
1 for each piece of info

① Encode identity of the token:

Learn a parameter matrix  $W^{\text{embed}} \in \mathbb{R}^{|V| \times d}$

corresponding vector for  $w_t$  is  $W^{\text{embed}} [w_t]$

= " $w_t$ "-th row of matrix

if we number all tokens from  $1, 2, \dots, |V|$

② Encode position of the token:

Learn a vector  $p_t$  for each  $t=1, \dots, T$

Final embedding  $x_t = W^{\text{embed}} [w_t] + p_t$

Note of caution: Some modern models don't use absolute position embeds.

Instead have tricks to encode "relative position"

### Step 3: Generating Predictions

Input: Final hidden state

$h_T^{(L)} \in \mathbb{R}^d$   
L = total # of layers  
Last timestep

Output: Probability distribution over  $V$

$$\text{Softmax} \left( \underbrace{W^{\text{embed}}}_{\in \mathbb{R}^{|V| \times d}} \times \underbrace{\text{LN} \left( h_T^{(L)} \right)}_{\in \mathbb{R}^d} \right)$$

Layer Normalization

$$\text{Softmax}(v) = \left[ \frac{e^{v_1}}{\sum_{i=1}^n e^{v_i}}, \dots, \frac{e^{v_n}}{\sum_{i=1}^n e^{v_i}} \right]$$

exp

Key Properties: • Every entry is positive

• Sums to 1

⇒ Valid probability distribution!

What is Layer Norm?

Trick to ensure that vector is at a "reasonable" scale

Input: Vector  $x \in \mathbb{R}^n$

**Step 1:** Normalize  $x$ 's entries to have 0 mean & variance 1

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i, \quad \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

$$\text{Normalized } x = \tilde{x} = \frac{x - \mu}{\sigma}$$

**Step 2:** Rescale & shift by learned parameters  $a \in \mathbb{R}^n, b \in \mathbb{R}^n$

$$\text{LN}(x) = a \odot \tilde{x} + b$$

↑ elementwise re-scaling    ↑ normalized  $x$     ↑ elementwise shift

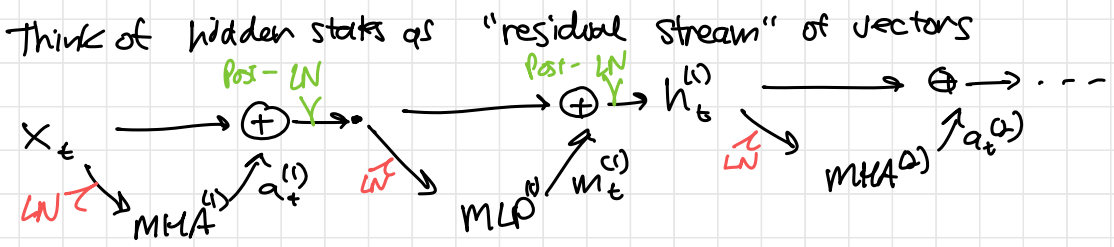
Why? Want some flexibility in scale of vector entries

**Step 2: Computing hidden states**

$$h_t^{(l)} = h_t^{(l-1)} + \underbrace{a_t^l}_{\text{output of } l\text{-th multi-headed attention block (MHA)}} + \underbrace{m_t^l}_{\text{output of } l\text{-th feedforward network/MLP}}$$

↑ layer    ↑ token index

$$\text{Base case: } h_t^{(0)} = x_t$$



"Stream": Constantly changing sequence of vectors,  
Read & write from this stream

"residual": Adding result of block to previous state is  
called "residual connection"

Post-LN: old way of applying LN - after adding to stream  
Pre-LN: new way - after read from stream

MHA: Retrieve relevant info from other tokens  
MLP: Does some processing for current token

### Multi-headed attention

First: what is single-headed attention?

- New hyperparam  $d_{\text{attn}}$  (Size of vector that each head uses)
- Input: sequence of vectors  $[u_1, \dots, u_T]$
- Output: vector  $\in \mathbb{R}^{d_{\text{attn}}}$

How? 3 parameter matrices  $W^Q, W^K, W^V \in \mathbb{R}^{d_{\text{attn}} \times d}$   
query key value

① Compute query vector  $q_T = W^Q \cdot u_T$  "what info are we looking for"

② Compute key vectors  $k_t = W^K \cdot u_t$  for all  $t = 1, \dots, T$   
"what info is available at each token"

③ Compute value vectors  $v_t = W^V \cdot u_t$  for all  $t = 1, \dots, T$   
"Actual info desired"

④ Try to match query with all keys:

$$\text{Compute } S_t = \underbrace{q_T^T K_t}_{\sqrt{d_{\text{data}}}} \text{ for each } t=1, \dots, T$$

⑤ Convert to probabilities by Softmax:

$$[P_1, \dots, P_T] = \text{Softmax}([S_1, \dots, S_T])$$

$$\text{⑥ output} = \sum_{t=1}^T P_t \cdot v_t$$

From single head to multi-head:

MHA layer has  $n_{\text{heads}}$  attention heads

(from this,  $d_{\text{data}} = d / n_{\text{heads}}$ )  
- Each head has separate  $W^Q, W^K, W^V$  parameters  
yields outputs  $O_1, \dots, O_{n_{\text{heads}}}$

Final output of MHA is

$$\text{MHA}(u_1, \dots, u_T) = \underbrace{W^O}_{\substack{\text{parameter} \\ \text{matrix of} \\ \text{size } d \times d}} \times \underbrace{\text{Concat}([O_1, \dots, O_{n_{\text{heads}}}] )}_{\in \mathbb{R}^d}$$

In Transformer:

There's  $\text{MHA}^l$  at each layer  $l=1, \dots, L$

$$a_t^l = \text{MHA}^l \left( \text{LN}(h_1^{(l-1)}), \dots, \text{LN}(h_t^{(l-1)}) \right)$$

Feedforward Network / Multi-Layer Perceptron (FFN) (MLP)

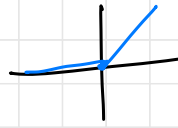
Most basic: 2-layer MLP:

Input vector  $x \in \mathbb{R}^d$ , output  $y \in \mathbb{R}^d$

$$y = \underbrace{W^{(2)}}_{d \times d \text{ hidden}} \cdot \underbrace{f}_{\substack{\text{elementwise} \\ \text{non-linearity}}} \left( \underbrace{W^{(1)}}_{d \text{ hidden} \times d} x + \underbrace{b^{(1)}}_{d \text{ hidden}} \right) + \underbrace{b^{(2)}}_d$$

What is  $f$ ?

- ReLU



- Swish



How used in TF?

$$m_t^{(l)} = \text{FFN}^{(l)} \left( \text{LN} \left( h_t^{(l-1)} + a_t^{(l)} \right) \right)$$

↑  
d-th layer's  
FFN