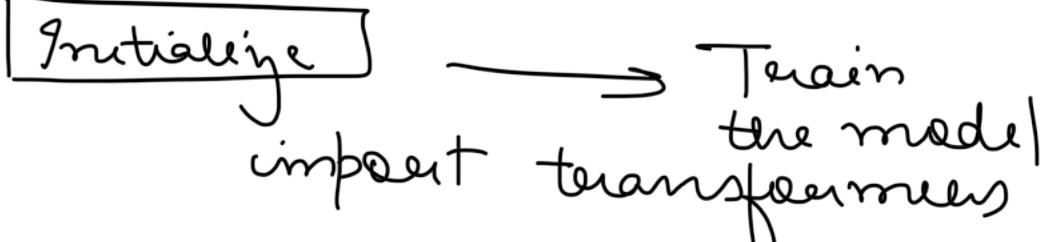


# DLP-Week-3 Notes

## L1 Intro to LLMs



- how to train transformer based models using the HF transformers module.
- Over the years, the dataset keeps growing, & also the model size.
- transformers based models can be of 3 types - Encoder only, decoder only, & encoder & decoder models.  
But these days decoder only models are emerging for the tasks of lang. modelling. (GPT based)
- all the LLMs go thru pretraining phase.
- we try to use unlabelled data to train the LLM, the entire language. After that, model can perform anything like sentiment, NLP, etc.

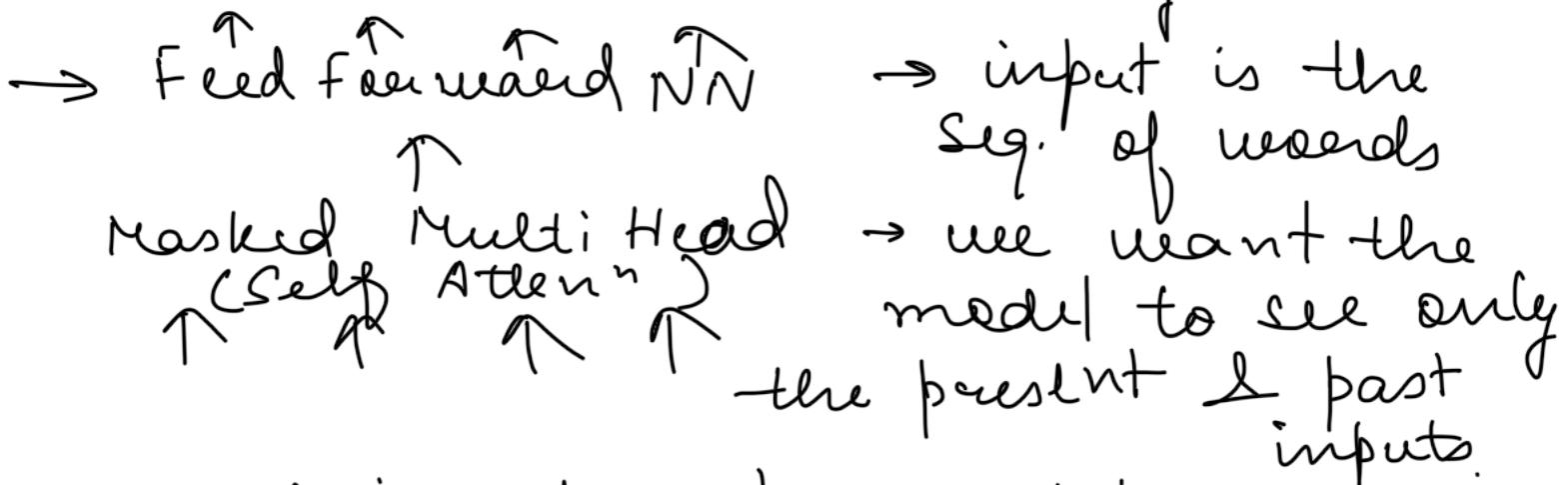
## L2 Causal Language Modelling

$$P(x_1, \dots, x_T) = \prod_{i=1}^T P(x_i | x_1, \dots, x_{i-1})$$

$$P(x_i | x_1, \dots, x_{i-1}) = f_\theta(x_i | x_1, \dots, x_{i-1})$$

can  $f_\theta$  be a transformer?

- so we had 3 types of transformers?
  - ↪ we look into decoder only model.



→ we achieve this by masking. So we need the multi-head attention layer is required. The output represents each term in the chain rule -

$$= P(x_1) P(x_2 | x_1) P(x_3 | x_2, x_1)$$

- the probab. are determined by the parameters of the model  $\rightarrow f_\theta$ .  
 We want to maximize the likelihood of  $L(\theta)$

## L3 Generating Pre-Trained Transformers

- stack of  $n$  modified decoders layers.

$X$  is input seq.

$$h_0 = x^T E R^{-1} \hat{x} + \text{model}$$

$h_t = \text{transformer-block}(h_{t-1})$ ,

$$P(x_i) = \text{softmax}(h_n[i] W_v)$$

$$L = \sum_{i=1}^T \log(P(x_i | x_1, \dots, x_n))$$

- Book corpus data  $\rightarrow$  7000 unique books, 74 million sentences, 1 billion words 16 genres, uses long-range contiguous text. BPE Tokenizer. Vocab size = 40478 Embedding dim  $\rightarrow$  768
  - Model - 12 decoder layers (transformer blocks)  
Context size = 512  
Attention heads = 12  
FFN Hidden Layer Size =  $768 \times 4 = 3072$   
GELU Activation (Gaussian)  $\rightarrow$  FFNN  
Multi head attention mask  
concatenate

Moduli

Softmax

all week parallelly

$$\frac{1}{\sqrt{d_e}}$$

- is such attempt mask heads

$$Q^T K + M$$

- after that a residual connection (here a dropout is used).
  - FFNN
- 768
- 3072

— 768

- no. of params.

$$\begin{aligned}\text{token embedding} &= V \times \text{embedding} \\ &= 40478 \times 768 = 31 \text{ M}\end{aligned}$$

$$\begin{aligned}\text{Position embedding} &= \text{content len} \times \text{em.} \\ &= 512 \times 768 = 0.3 \text{ M}\end{aligned}$$

attention params per block

$$\begin{aligned}W_Q = W_K = W_V &= 768 \times 64 \times 3 \times 12 \\ &= 1.7 \text{ M}\end{aligned}$$

$$\text{Linear layers} = 768 \times 768 \approx 0.6$$

$$2 \cdot 3 \text{ M} \times 12 = 27.6 \text{ M} \quad \text{for attention block.}$$

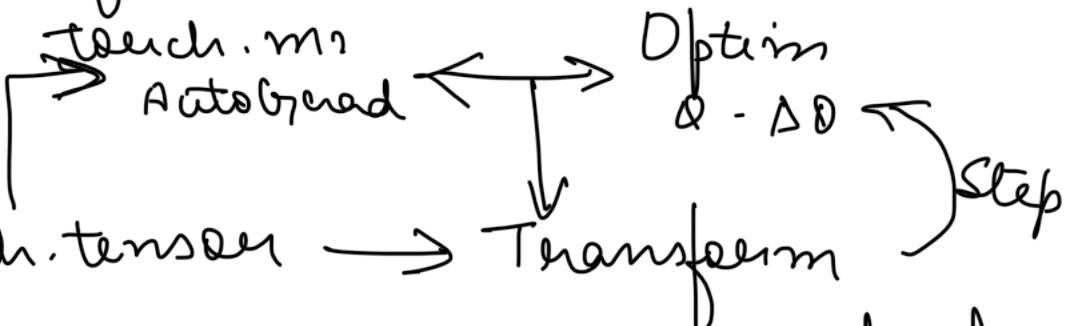
FFN per block

$$\begin{aligned}2 \times (768 \times 3072) + 3072 + 768 &= 4.7 \text{ M} \\ &\times 12 \\ &= 56.4 \text{ M}\end{aligned}$$

Total = 117 million params

## By HF Transformers

import transformers



- the module abstracts the seq. of opera<sup>n</sup>
- Use the data in batches.
- Please beware of shuffling of data.  
Because of Python GIL.  
↓ to avoid this  
use built-in data loader