



Contents:

1. Neural Language modelling
2. RNN's and LSTM's
3. Transformers
4. Decoding
5. Pre-training
6. Scaling Laws of LLMs
7. LM Assistants and RLHF policies
8. PBFIT Techniques (LORA, Adapters, Pruning)

Neural Language modeling RNNs :-

- language model is a system that assigns probability to a piece of text.
- let us have texts $x^{(1)}, \dots, x^{(T)}$, then probability of this text is

$$P(x^{(1)}, \dots, x^{(T)}) = P(x^{(1)}) \times P(x^{(2)} | x^{(1)}) \times \dots \times P(x^{(T)} | x^{(T-1)}, \dots, x^{(1)})$$

$$= \prod_{t=1}^T P(x^{(t)} | x^{(t-1)}, \dots, x^{(1)}).$$

N-gram language modeling :-

- $x^{(t+1)}$ depends only on the preceding $m-1$ words.
- $P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)}) = P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-m+2)}).$

$$= \frac{P(x^{(t+1)}, x^{(t)}, \dots, x^{(t-m+2)})}{P(x^{(t)}, \dots, x^{(t-m+2)})}$$

statistical
approximation

$$\approx \frac{\text{count}(x^{(t+1)}, \dots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \dots, x^{(t-n+2)})}$$

Ex :- students opened their w
condition on this

use a 4-gram language model

$$P(w / \text{students opened their}) = \frac{\text{Count students opened their } w)}{\text{Count students opened their}}$$

$V^n \rightarrow$ need to store count
for all n-grams you
saw in the corpus.

A fixed-window Neural language model :-

- It operates on a fixed-length sequences of words
- the model predicts the target word based on the probability distribution generated by it.

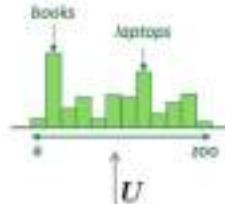
A fixed-window neural language model

$$\text{Number of parameters} = 4vh + vh = svh$$

output distribution

$$\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$$

V



hidden layer

$$h = f(We + b_1)$$

h

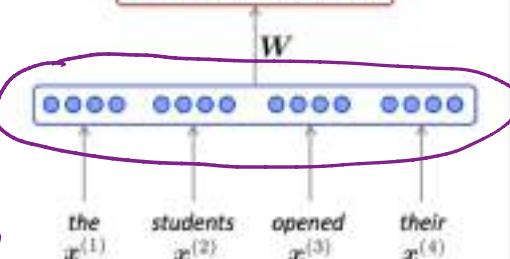
concatenated word embeddings
 $e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$

words / one-hot vectors
 $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$

4V
One-hot



W



Jasmin Ghosh (IIT Kharagpur)

Neural Language Model, RNNs

CSE6043

12/30

problem in one-hot encoding :- the vectors are orthogonal, and there is no notion of similarity between one-hot vectors.

Word2Vec representation :-

Any word w_i in the corpus is given a distributed representation by an embedding

$$w_i \in \mathbb{R}^d$$

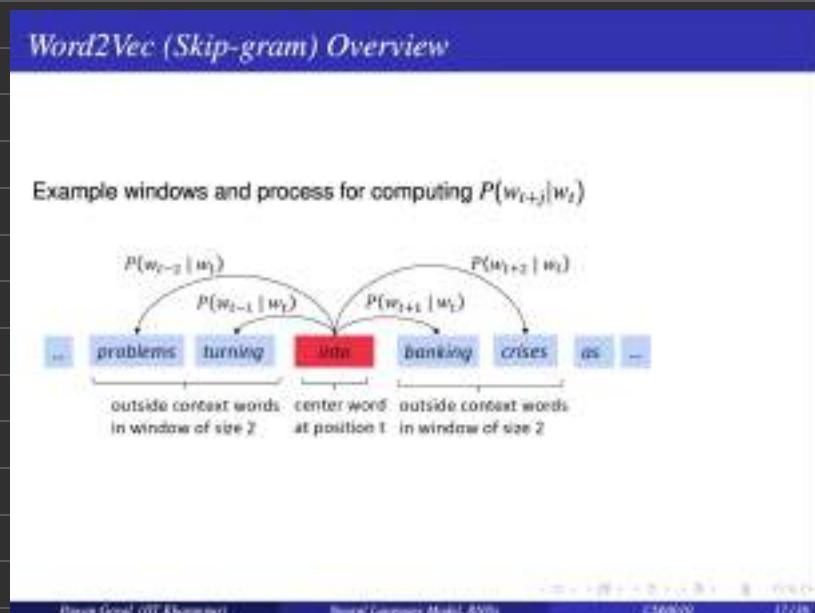
→ a d-dimensional vector

Basic idea! \rightarrow in a large corpus of text every word is represented by a vector.

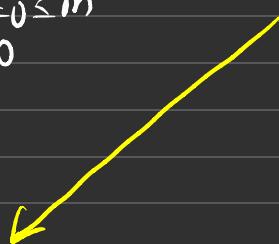
\rightarrow go through each position ' t ' in the text which has a center word ' c ' and context words ' o '.

\rightarrow use the similarity of the word vectors for c and o to calculate the probability of o given c .

\rightarrow keep adjusting the word vectors to maximize this probability.



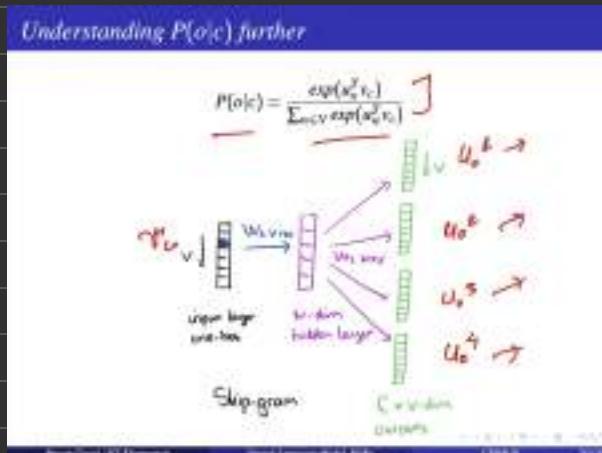
$$J(w) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log p(w_{t+j} | w_t; \theta)$$



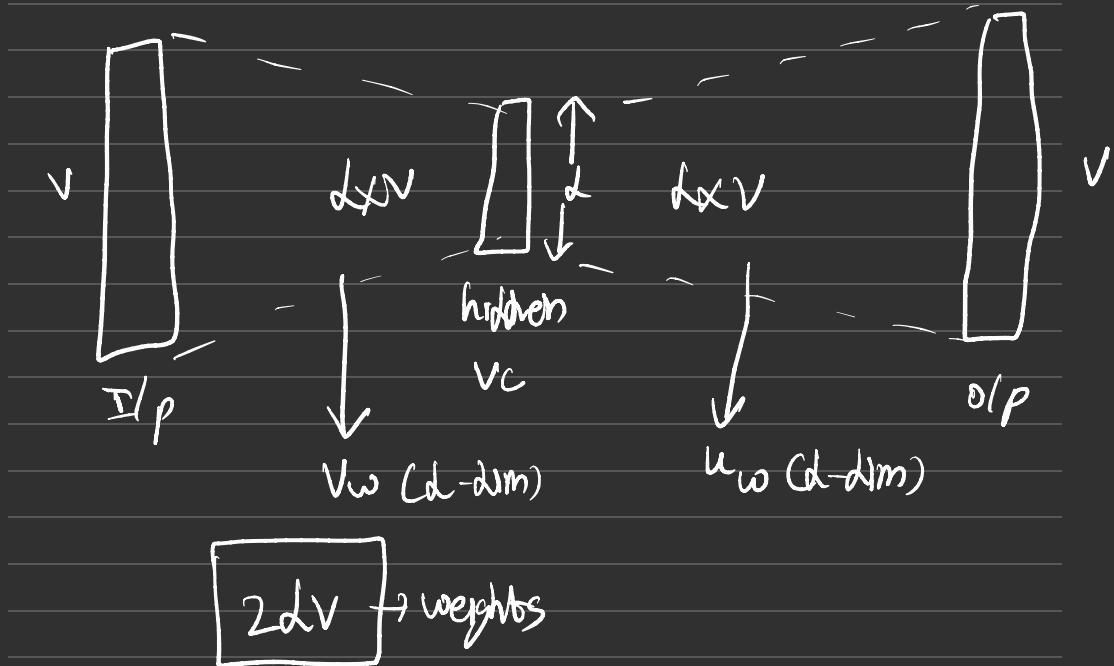
we will use two vectors per word below:-

- v_w when w is the center word
- u_w when w is the context word

$$p(c|w) = \frac{\exp(u_w^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$



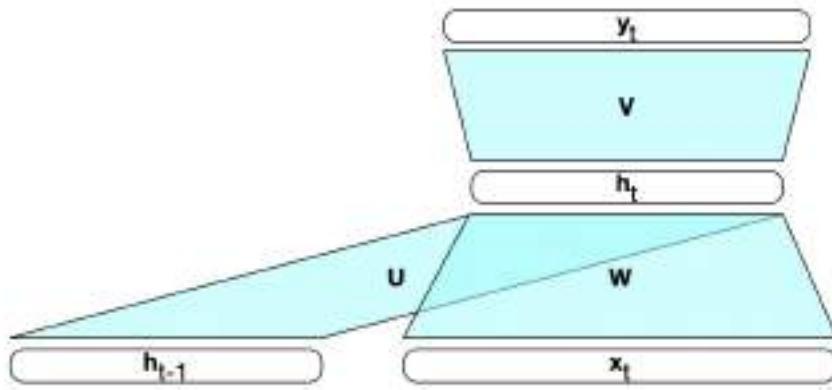
for a target word,



Recurrent Neural Networks :-

- Apply the same weights repeatedly.
- We can process a sequence of vectors \times by applying a recurrence formula at each step.

$$\overbrace{h_t}^{\text{new state}} = f_w \left(\underbrace{h_{t-1}}_{\text{old state}}, x_t \right)$$

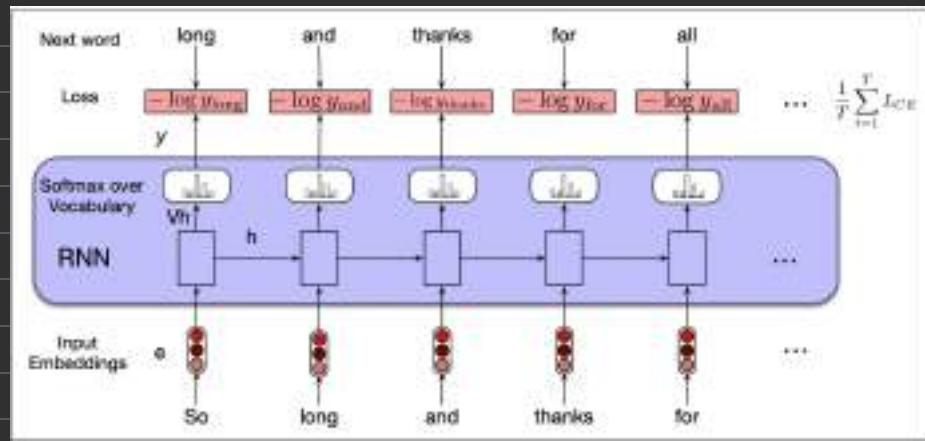


$$h_t = f(u h_{t-1} + w x_t) \quad w : d_h \times d_n$$

$$y_t = \text{softmax}(v h_t) \quad v : d_h \times d_n$$

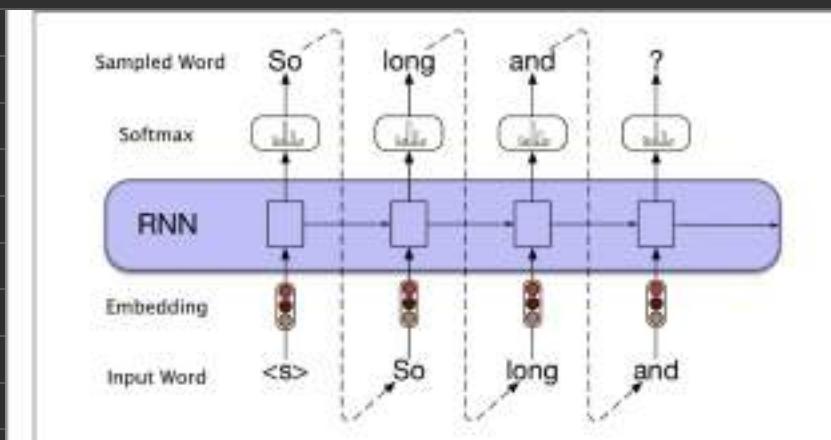
$$f : d_{out} \times d_h$$

- we simply train the model to minimize the error in predicting the true next word in the training sequence.



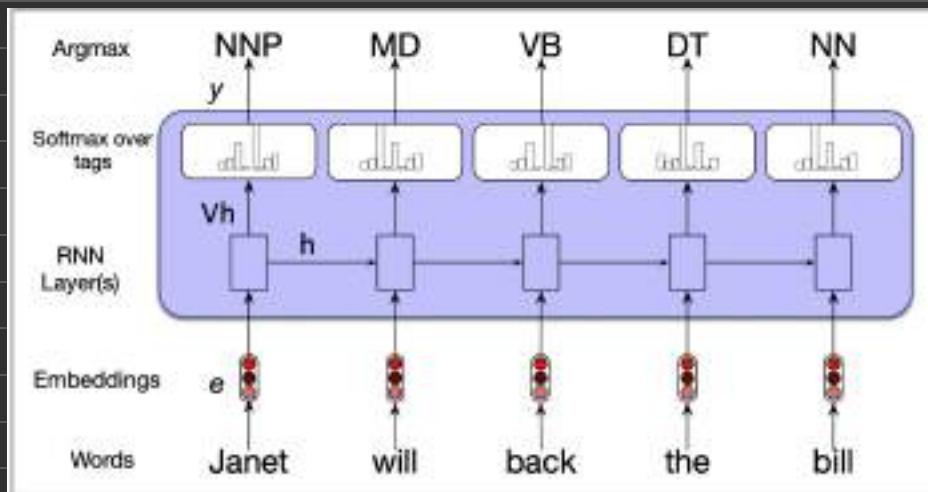
$$L_{CE}(\hat{y}_t, y_t) = -\log \hat{y}_t[w_{t+1}]$$

Auto regressive generation with RNNs :-



- All your parameters have already been trained.
- start with special token $\langle s \rangle$
- through forward propagation, obtain the probability distribution and sample a word.
- Feed the corresponding word vector as input for next time step
- continue generating until the end of sentence token is sampled or a fixed length of sentence has been reached.

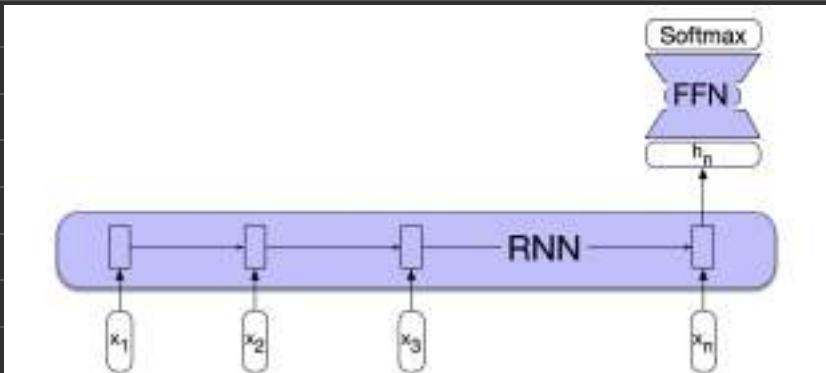
Sequence labeling and Sequence classification with RNN's :-



parts of speech tagging

The network's task is to assign a label chosen from a fixed set of labels to each element of a sequence.

Sequence classification :-



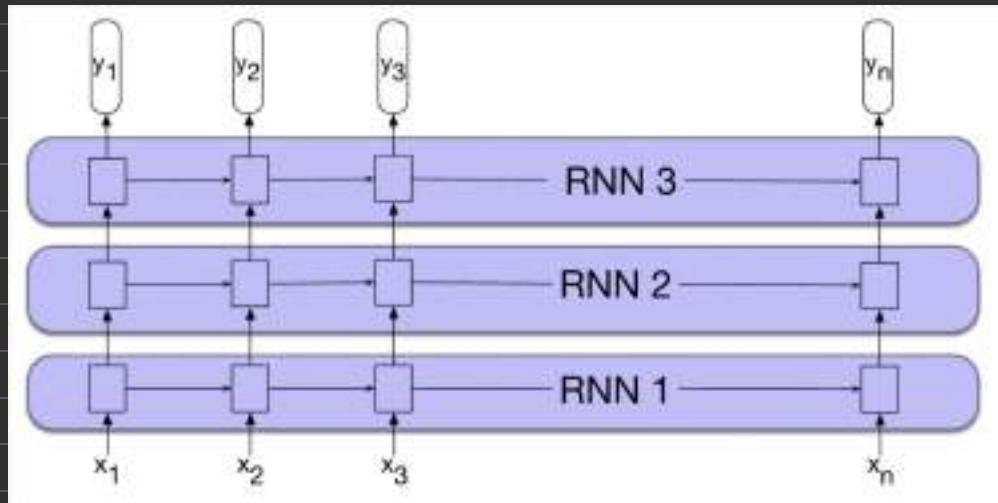
we pass the text to be classified through the RNN, a word at a time generating a new hidden layer at each time step.

- the hidden layer of the last token h_n , to constitute a compressed representation of entire sequence.

Alternate method :-

$$h_{\text{mean}} = \frac{1}{n} \sum_{i=1}^n h_i \rightarrow \text{feed it into FFN}$$

Stacked RNN's :-



the output of lower level serves as input to the higher levels with the output of the last network serving as final output.

Bi-Directional RNN :-

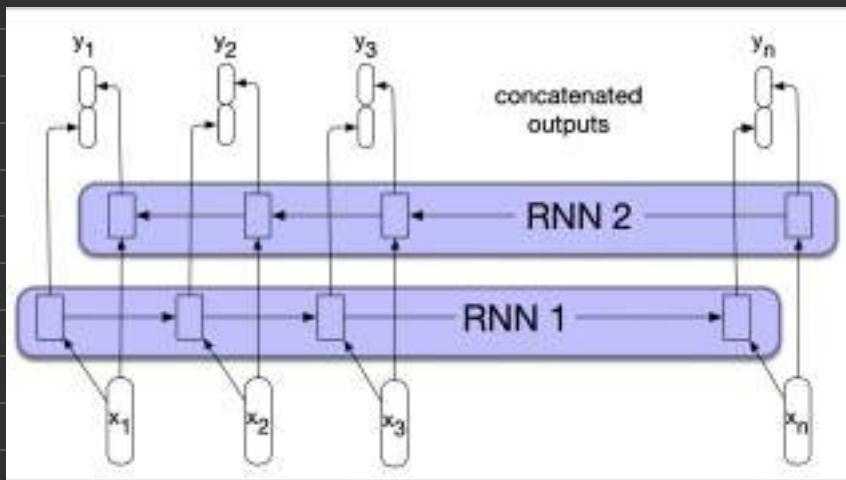
- in the left-to-right RNN

$$h_t^f = \text{RNN}_{\text{forward}}(x_1, \dots, x_t)$$

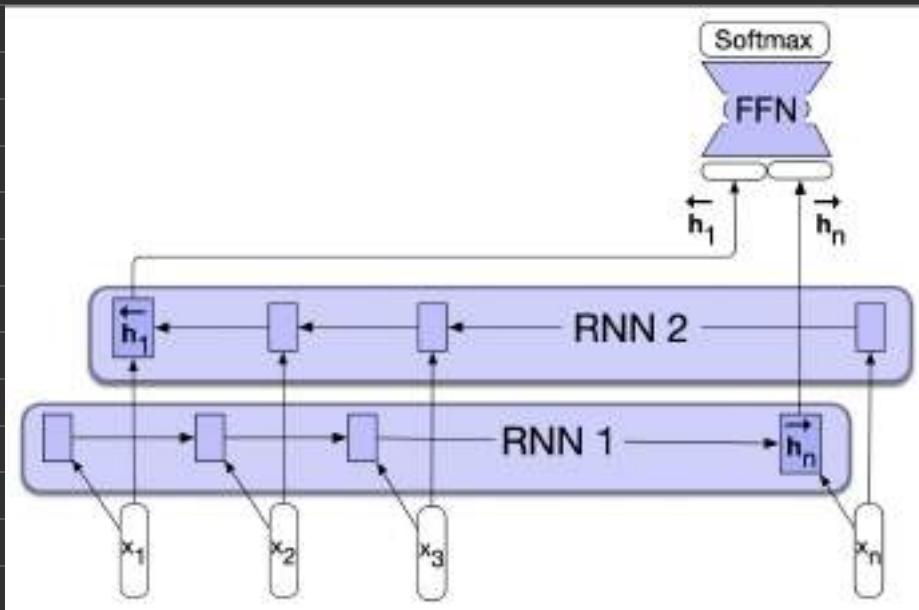
- in the right-to-left RNN

$$h_t^b = \text{RNN}_{\text{backward}}(x_t, \dots, x_n)$$

$$h_t = h_t^f \oplus h_t^b \Rightarrow y_t = \text{softmax}(Vh_t)$$



A bidirectional RNN for sequence classification:-



- A difficulty with the general approach is that the final state naturally reflects more information about the end of sentence than its beginning. So, bidirectional-RNN will solve this problem.

Problems in RNNs :-

- model weights are updated only with respect to near effects, not long-term effects.

- Gradient signal from faraway is lost because its much smaller than gradient signal from close by.

- Syntactic recency :- The writer of the books is ✓
- Sequential recency :- The writer of the books are ✗
- Due to vanishing gradient RNN-LM's are better at learning sequential recency than syntactic recency.

- In RNN's, the hidden state is constantly being rewritten

Long-Short term memory :- (LSTM)

- LSTMs divide the context management problem into two subproblems :-
 - ① removing information no longer needed from the context.
 - ② Adding information likely to be needed for later decision making.

- it uses gates to achieve the above two tasks. It uses additional weights.
- the gates in LSTM share a common design pattern; each consists of a feed forward layer followed by a sigmoid activation function, followed by a pointwise multiplication.

forget gate :-

$$f_t = \sigma(v_f h_{t-1} + w_f x_t)$$

$$k_t = c_{t-1} \odot f_t$$

- extract the actual information from previous hidden state.

$$g_t = \tanh(v_g h_{t-1} + w_g x_t)$$

add gate :-

$$i_t = \sigma(v_i h_{t-1} + w_i x_t)$$

$$j_t = g_t \odot i_t$$

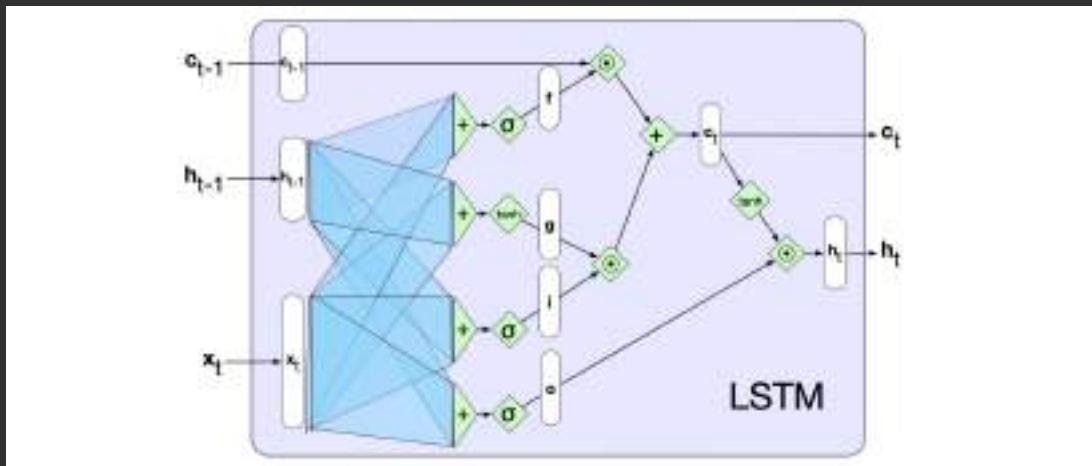
$$\Rightarrow C_t = j_t + k_t$$

output gate :-

$$o_t = \sigma(v_o h_{t-1} + w_o x_t)$$

$$h_t = o_t \odot \tanh(c_t)$$

- An LSTM as input the context layer, and hidden layer from the previous step.



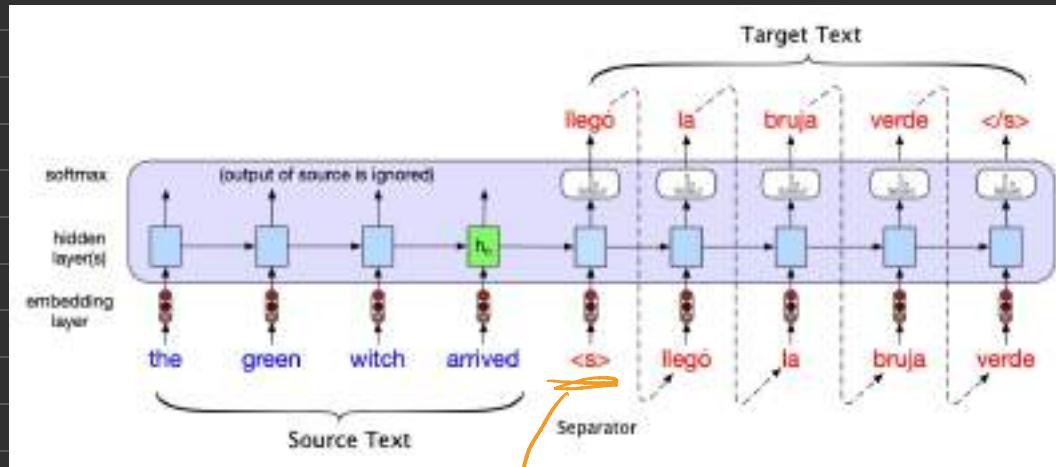
Encoder- Decoder Networks

- Encoder :- accepts an input sequence $x_{1:n}$ and generates an corresponding sequence of $h_{1:n}$.
- Context :- A Context Vector c which is a function of $h_{1:n}$, acts as input to

decoder.

- Decoder :- who accepts c as input and generates $h_{n:m}$ and corresponding output $y_{n:m}$ is generated

Machine Translation :-

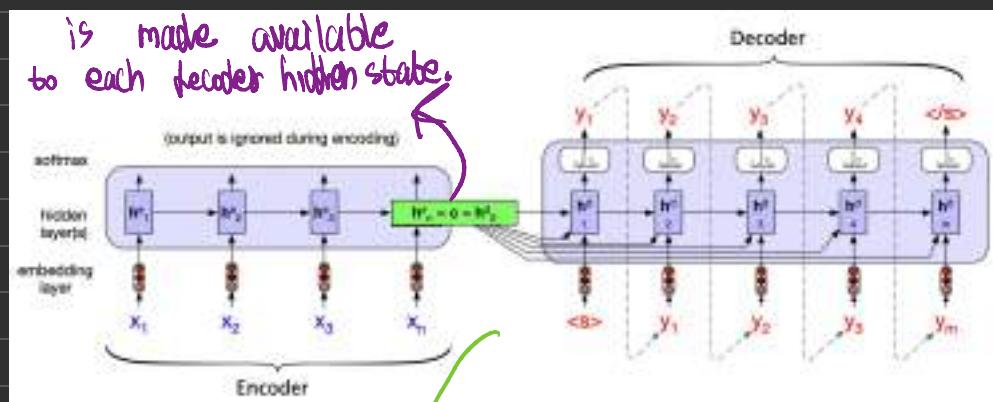
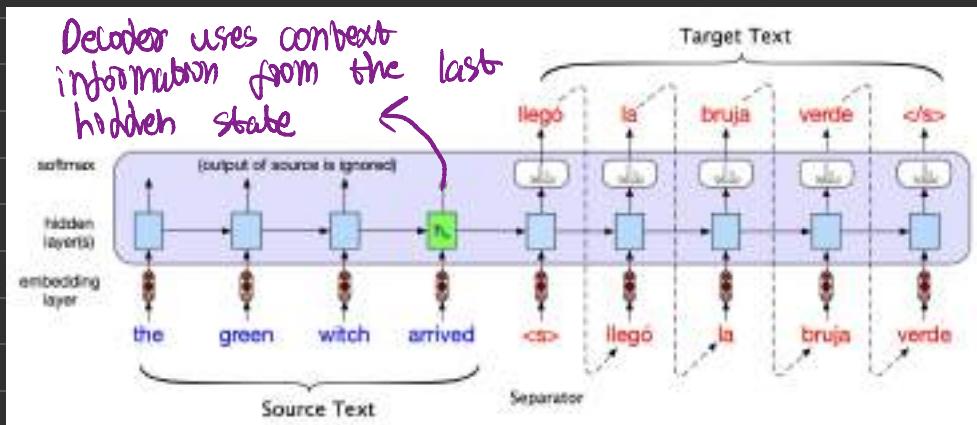


$x \rightarrow$ Source text
 $y \rightarrow$ target text

sentences separator token

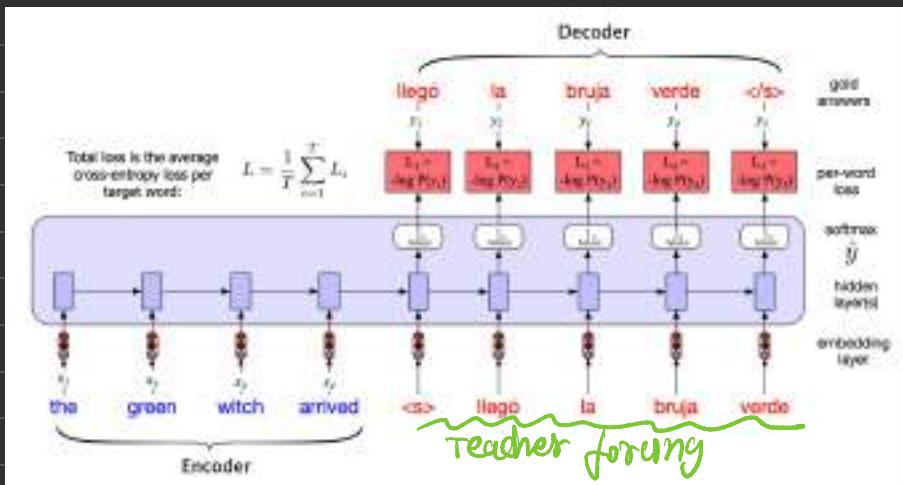
\Leftrightarrow sentences separator.

$$p(y|x) = p(y_1|x) \cdot p(y_2|y_1, x) \cdot p(y_3|y_1, y_2, x) \cdots p(y_m|y_1, \dots, y_{m-1}, x)$$



$$h_t^d = g(\hat{y}_{t-1}, h_{t-1}^d, c)$$

Training an Encoder-Decoder model :-



- it is trained autoregressively to predict the next word.
- Teacher forcing is used during training i.e. the system is forced to use the gold target tokens from training as the next inputs x_{t+1} rather than relying on the last decoder output.

Not during inference !

* The last hidden state acts as the bottleneck, as it has to represent absolutely everything about the source text.

Solution := Attention mechanism :-

→ context vector, c is $c = f(h_1^e, \dots, h_n^e)$, we can't directly use the entire set of encoder hidden states as they vary with the size of inputs

for each decoder step ;

$$h_i^d = g(\hat{y}_{i-1}, h_{i-1}^d, c_i)$$

Calculation of c_i :-

① Score $(h_{i-1}^d, h_j^e) = h_{i-1}^d \cdot h_j^e \rightarrow$ for each encoder hidden state j

$$\alpha_{ij} = \text{softmax}(\text{score}(h_{i-1}^d, h_j^e))$$

• finally, given the distribution

$$c_i = \sum_j \alpha_{ij} h_j^e$$

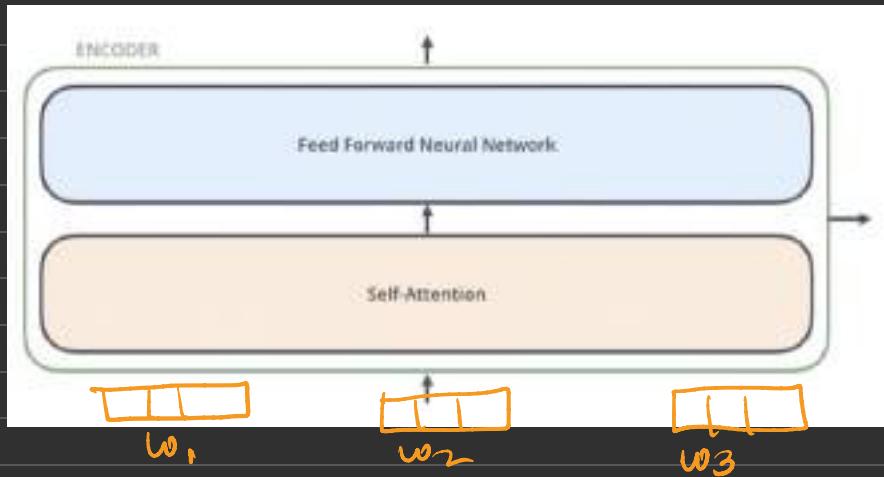
we can get a more powerful scoring function by parameterizing the slope with its own set of weights, w_s .

$$\text{Score } (h_i^d, h_j^e) = h_i^d \underbrace{w_s}_{\text{are trainable}} h_j^e$$

are trainable

issues with RNN's :-

- Hard to learn long-distance dependencies.
- Linear order of words ! Not right way to think
- Future RNN hidden states can't be computed in full before past RNN hidden states have been computed. lack of parallelizability !
- In Attention, all words can interact with each other and computation can be done in parallel.



- Each encoder receives a list of vectors each of size 512 .
- In the bottom encoder, this would be the word embeddings.

★★ The word in each position flows through its own path in the encoder!

There are dependencies between these paths in the self attention layer but no in FFN.

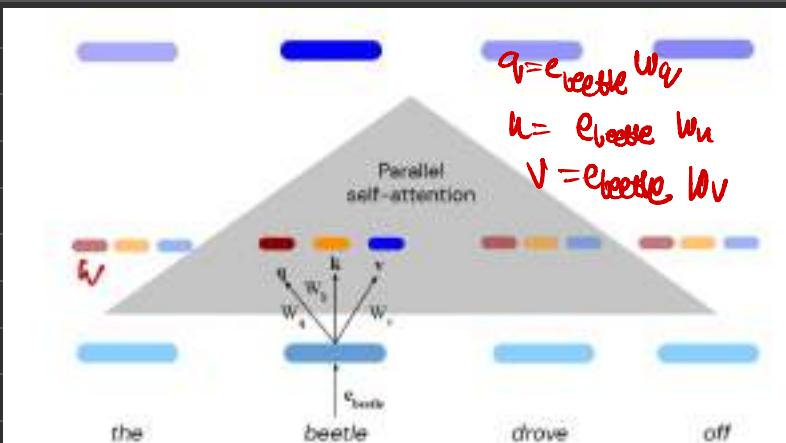
Attention :-

- in self-Attention the queries keys, and values are drawn from the same source.
- The dot product self-Attention is

$$e_{ij} = q_i^T k_j \quad \alpha_{ij} = \frac{\exp(e_{ij})}{\sum_j \exp(e_{ij})}$$

$$\text{output}_i = \sum_j \alpha_{ij} v_j$$

* The number of queries can differ from the number of keys in practice.



- Attention for all the words can be done in parallel!

Multi-Attention Heads :-

- The different words in a sentence can relate to each other in many different ways simultaneously.
- it would be difficult for a single self-Attention model to learn to capture all of the different kinds of parallel relations among its inputs

Concept of multi-head Attention :-

We use distinct set of parameters, each head can learn different aspects of the relationships among inputs at the same level of abstraction.

- Each head has its own w_i^U , w_i^Q and w_i^V
- $d \rightarrow$ dimension of input and output

for each head i ,

$$w_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$$

$$w_i^K \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$$

$$w_i^V \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$$

These get multiplied by input x

outputs :-

$$\begin{matrix} Q \in \mathbb{R}^{N \times d_{\text{model}}} \\ K \in \mathbb{R}^{N \times d_{\text{model}}} \\ V \in \mathbb{R}^{N \times d_{\text{model}}} \end{matrix}$$

for

Each head, output shape? $N \times d_{\text{model}}$

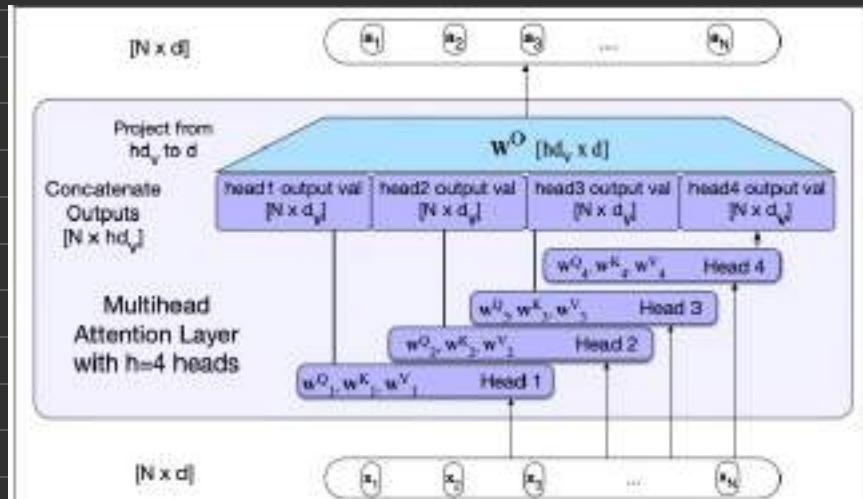
all the outputs are concatenated $\rightarrow N \times d_{\text{model}}$

This is multiplied by $w^O \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$ \rightarrow to reshape the output to original dimensions

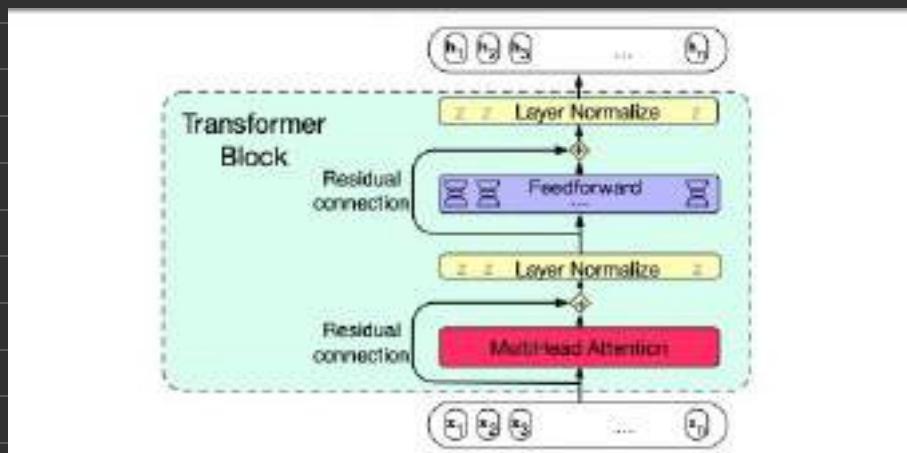
$$Q = x \cdot w_i^Q; \quad K = x \cdot w_i^K; \quad V = x \cdot w_i^V$$

head $_i$ = Self Attention (Q, K, V)

$$\text{MultiHeadAttention}(x) = (\text{head}_1 \oplus \text{head}_2 \dots \oplus \text{head}_n) \cdot w^O$$



A transformer block showing all layers ↴



$t_i^1 = \text{Multihead Attention} (x_i, x_1, \dots, x_N)$

$$t_i^2 = t_i^1 + x_i$$

$t_i^3 = \text{LayerNorm}(t_i^2)$

$t_i^4 = \text{PE}N(t_i^3)$

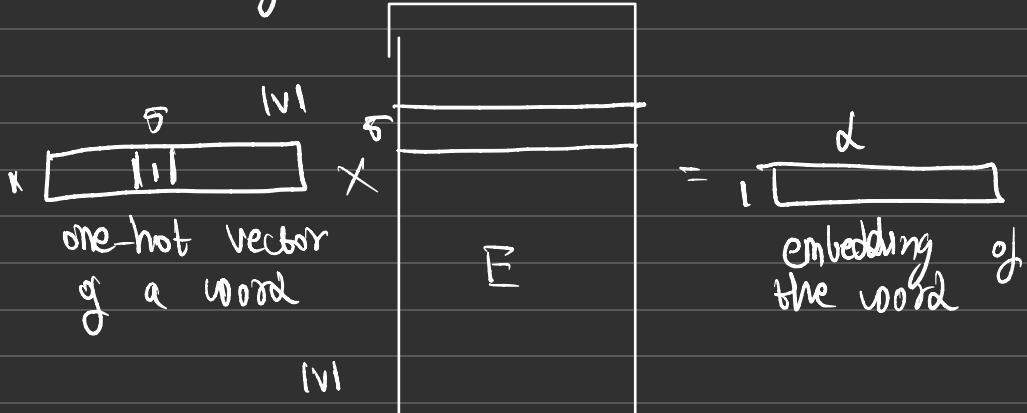
$$t_i^5 = t_i^4 + t_i^3$$

$h_i = \text{LayerNorm}(t_i^5)$

↳ output of transformer block

Embeddings :-

$\Sigma \rightarrow \text{Embedding Matrix}$



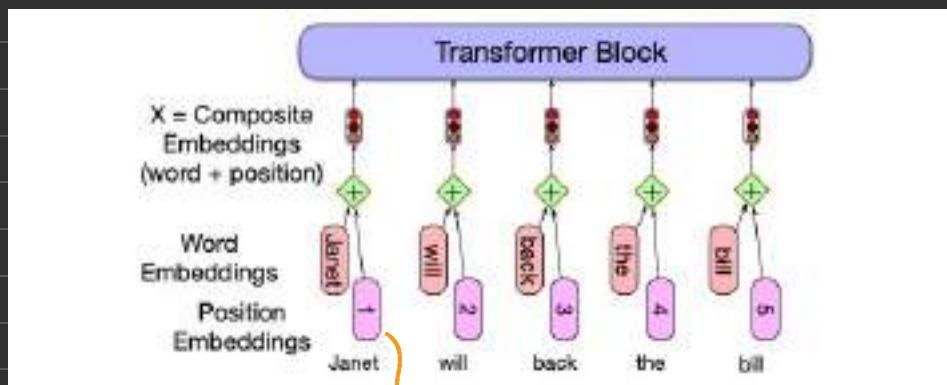
token-embeddings

These token embeddings are not position dependent. To represent the position of each token, we combine these token embeddings with positional embeddings.

Simple Approach :-

- Store the absolute position of the words.
- These positional embeddings are learned along with other parameters during training.

$$E_{pos} \rightarrow 1 \times N$$



both are of $[d \times d]$ size
problem with the above approach :-

- There will be plenty of training examples for the initial positions in our inputs and correspondingly fewer at the outer length limits.

- This leads to poor generalization for latter embeddings.

Solution :- choose a static function that maps integer inputs to real-valued vectors that captures the inherent relationship among the positions.

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

- These position embeddings are added to word embeddings.

Vocabulary :-

→ we assume a fixed Vocab of tens of thousands of words, built from the training set.

→ All novel words seen at best time are mapped to a single UNK.
↑
unknown type

Byte-pair encoding Algorithm :-

- start with a vocabulary containing only characters and an end of word symbol.
- using a corpus of text, find the most common adjacent characters "ab"; add "ab" as a subword.
- replace instances of the character pair with the new subword; repeat until desired vocab size.

Layer Normalization :-

- let $x = [x_1, \dots, x_3]$ be an individual word vector in the model.

$$u = \frac{1}{d} \sum_{i=1}^d x_i, \quad u \in \mathbb{R}$$

$$\sigma = \sqrt{\frac{1}{d} \sum_{i=1}^d (x_i - u)^2}, \quad \sigma \in \mathbb{R}$$

Let $\gamma \in \mathbb{R}^d$ and $\beta \in \mathbb{R}^d$ can be learnable parameters.

$$\text{LayerNorm} = \gamma \left(\frac{x - u}{\sigma + \epsilon} \right) + \beta$$

not important

Masked-Attention :-

The calculation in ~~out~~ results in a score for each query value to every key value including those that follow the query.

this is inappropriate in the setting of language modelling : guessing the next word is simple if you already knew it.

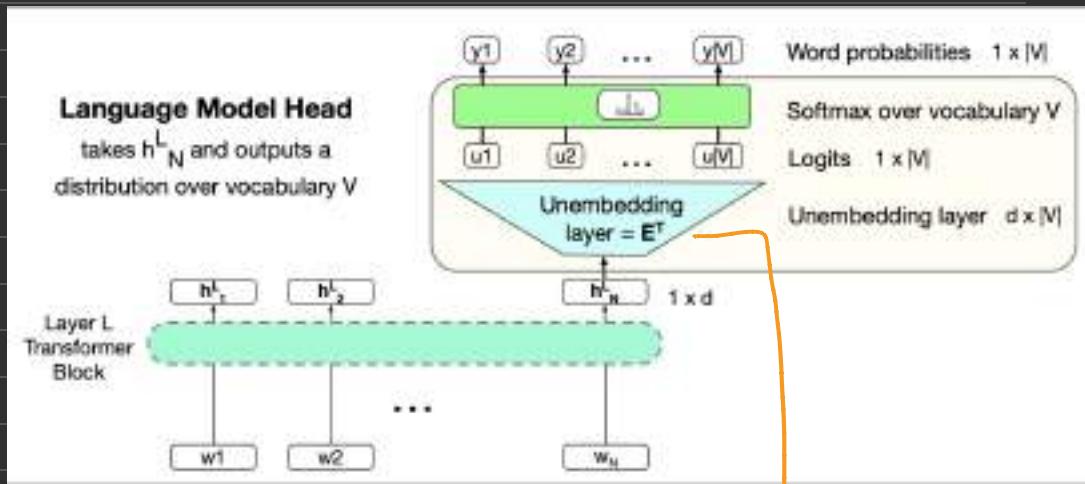
Solution

q1·k1	-∞	-∞	-∞	-∞
q2·k1	q2·k2	-∞	-∞	-∞
q3·k1	q3·k2	q3·k3	-∞	-∞
q4·k1	q4·k2	q4·k3	q4·k4	-∞
q5·k1	q5·k2	q5·k3	q5·k4	q5·k5

Encoder - Decoder Attention :-

- Let h_1, \dots, h_n be output vectors from the transformer encoder ; $h_i \in \mathbb{R}^d$
- Let z_1, \dots, z_n be input vectors from the transformer decoder ; $z_i \in \mathbb{R}^d$
- keys and values $\Rightarrow k_i = K h_i$, $v_i = V h_i$ \Rightarrow from encoder
- queries are drawn from decoder $\Rightarrow q_i = Q z_i$

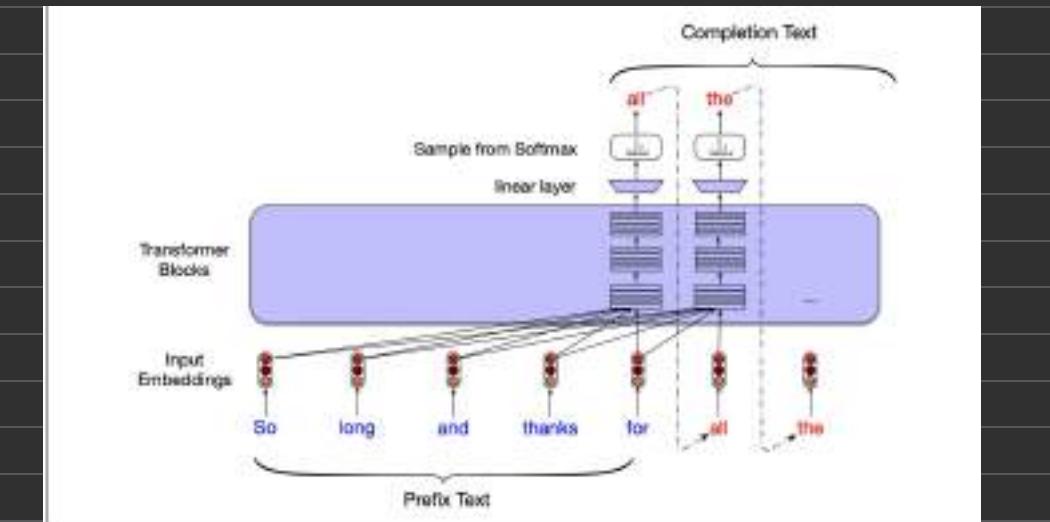
Transformers for language Modeling :-



→ Language modelling head

This linear layer can be learned but more commonly we tie this matrix to the transpose of the embedding matrix, (E^T)

- In the learning process E will be optimized to be good at doing both of these mappings.

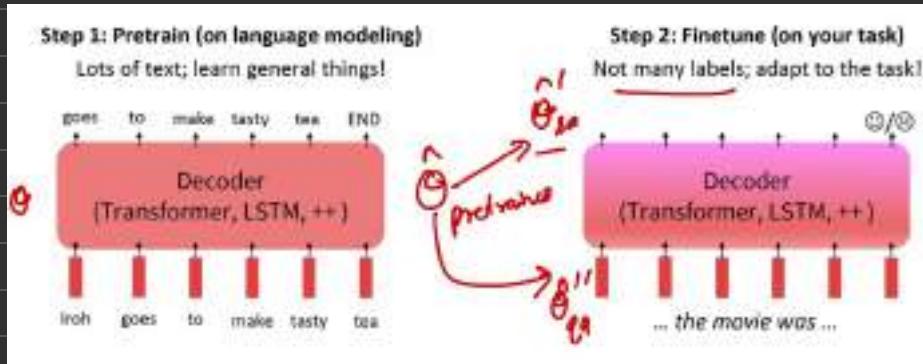


- As the generation proceeds, model has access to the context as well as its previously generated tokens.

Pre-Training :-

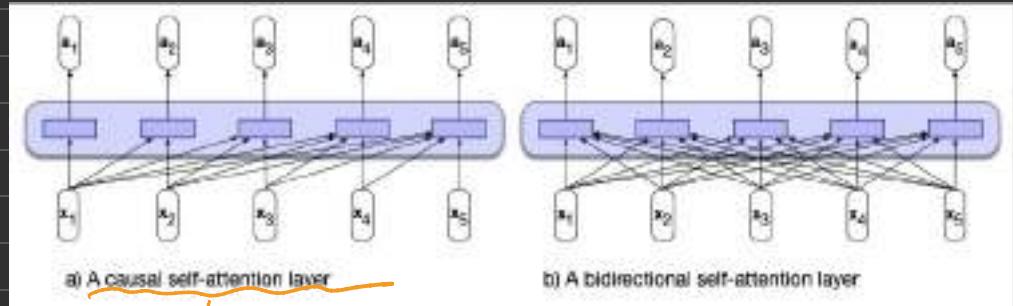
- Pre-training methods take parts of the input from the model and train the model to reconstruct these parts
- This worked exceptionally well for :-
 → representations of language
 → parameter initializations for strong NLP models

→ probability distributions over language that we can sample from.



- pretraining loss \Rightarrow profiles approximating parameters $\hat{\theta}$ by $\min_{\theta} L_{\text{pretrain}}(\theta)$
- Now, finetuning loss \Rightarrow approximate starting at $\min_{\theta} L_{\text{finetune}}(\theta)$,
- the pretraining may matter because stochastic gradient descent sticks relatively close to $\hat{\theta}$ during finetuning.
- maybe the finetuning local minima near $\hat{\theta}$ tend to generalize well!
- maybe the gradients of finetuning loss near $\hat{\theta}$ propagate nicely.

Bi-Directional Transformer Encoders :-



These models will not be able to solve problems like self-Attention and labelling which involves processing text from right to left also.

- Bi-directional encoders use self-Attention to map sequences of input embeddings (x_1, \dots, x_n) to sequences of output embeddings.

These output embeddings are contextualized representations of each input token that are generally useful across a range of downstream tasks.

The architecture for Bi-directional models :-

Bidirectional models use the same self-attention mechanism as causal models

$$q = XW^Q, \quad k = XW^K; \quad v = XW^V$$

$$X \in R^{N \times d_h}$$

q1·k1	q1·k2	q1·k3	q1·k4	q1·k5
q2·k1	q2·k2	q2·k3	q2·k4	q2·k5
q3·k1	q3·k2	q3·k3	q3·k4	q3·k5
q4·k1	q4·k2	q4·k3	q4·k4	q4·k5
q5·k1	q5·k2	q5·k3	q5·k4	q5·k5

↳ QK^T matrix

$$\text{Self-Attention } (Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_h}} \right) V$$

The key architecture difference is in bidirectional

models, we don't mesh the future. we simply skip the mesh each token using information from the future.

- Hidden layers of size 768
- 12 layers of transformers blocks, with 12 multihead attention layers each.
- The resulting model has about 100M parameters.
- Bert and all its descendants are based on subwords tokens rather than words.

Masking words :-

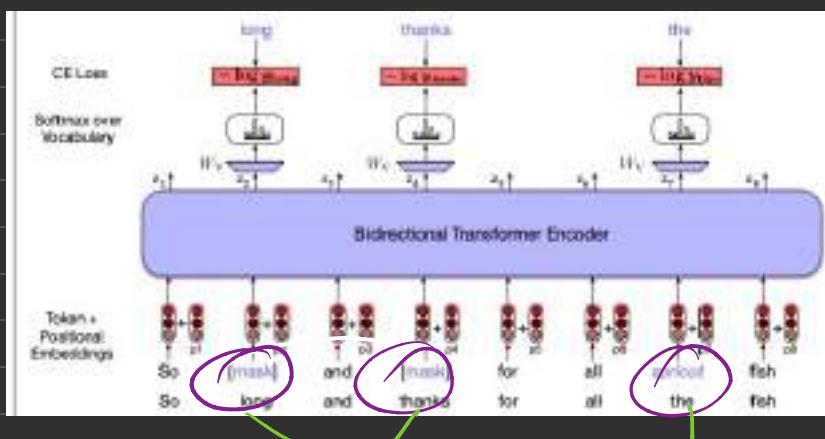
- The original approach to training bidirectional encoders is called masked language modelling
- The model is presented with a series of sentences from the training corpus where a random sample of tokens from each training sequence is selected for use in the learning task.

Once chosen, a token is used as

- ① it is replaced with unique vocabulary token [mask]
- ② it is replaced with another token from the vocabulary, randomly sampled on the token unigram probabilities.
- ③ it is left unchanged.

In BERT, 15% of the input tokens are sampled for learning

of these 80% replaced with [mask], 10% are selected tokens, and the remaining 10% are left unchanged.



replaced with mask

replaced with
other word

- All the inputs tokens play a role in the self Attention process, but only the sampled tokens are used for learning

$$y_i = \text{softmax}(W_v z_i)$$

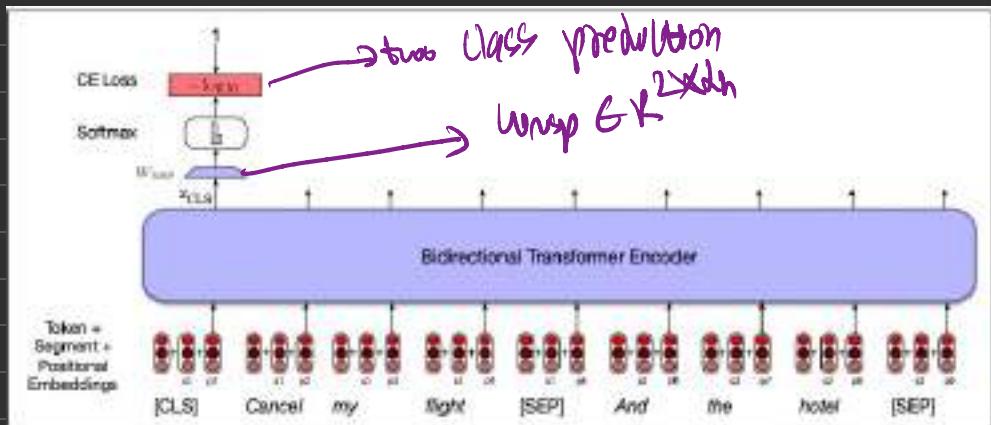
$$L_{MLM} = -\frac{1}{|M|} \sum_{i \in M} \log p(x_i | z_i)$$

Next sentence prediction :-

In this task, the model is presented with pairs of sentences and is asked to predict whether each pair consists of an actual pair of adjacent sentences from the training corpus or a pair of unrelated sentences.

- To facilitate Nsp training, BERT introduces two new tokens to input representation.
 - The token [CLS] is prepended to the input sequence
 - The token [sep] is placed between the sentences and after the final token of the second

Sentence.



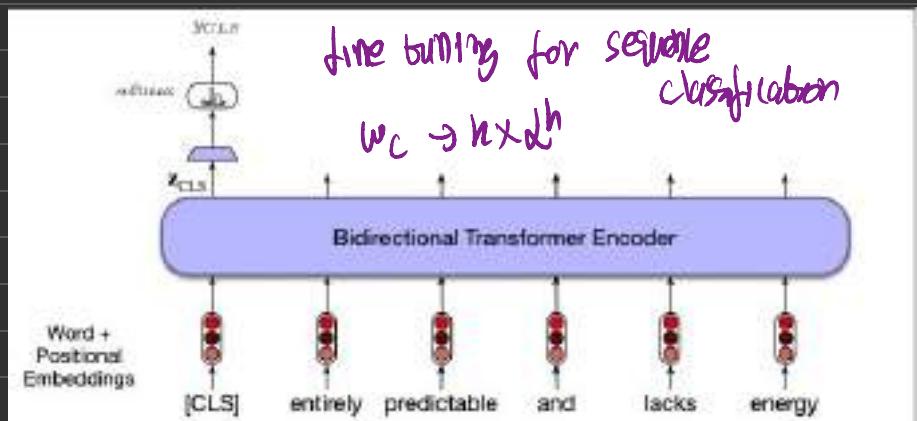
Cross entropy is used to train Nsp loss

why Nsp?

- Detecting if two sentences have similar meaning
- Detecting if the meanings of two sentences contradict each other
- Deciding if two sentences form a coherent disclosure.

Fine - tuning :-

- Fine - tuning facilitates the creation of Applications on top of pretrained models through the possible addition of small-set of application specific parameters.
- The fine-tuning process consists of using labeled data from the application to train those additional application-specific parameters.
- Typically this training will either freeze or make only minimal adjustments to the pretrained language model parameters.



Sequence labeling :-

Parts of speech tagging

(2)

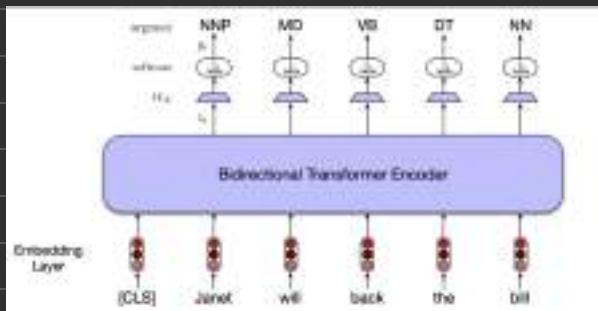
BIO-based named entity recognition

B → indicates the first token of a multi-token entity

I → indicates a token that is inside a multi-token entity

O → indicates a token that is not part of any entity.

Allows us to understand the types of entities in the text



$$y_i = \text{Softmax}(w_i z_i)$$

$$t_i = \text{argmax}_k (y_i)$$

Alternative :-

the distribution over labels provided by the softmax for each input token can be passed to a conditional random field (CRF) layer which can take global tag-level transitions into account.

→ They capture dependences between input features and output labels in sequence data.

[loc Mt. Sanitas] is in [loc sunshine canyon]

The set of per-word BIO tags.

Mt. Sanitas is in sunshine Canyon.
B-loc I-loc O O B-loc I-loc → O

The word piece tokenization for this sentence yields the following sequence of tokens

"Mt", " ", "San", "#nitas", "is", "in", "sunshine", "canyon", "

To deal with misalignment, we need a way to assign BIO tags to subword tokens

during training and a corresponding way to recover word-level tags from subwords during decoding.

- for training we can just assign the gold-standard tag associated with each word to all derived from it.
- of the subword tokens
- for decoding, the simplest approach is to use $\text{argmax}_{\text{subword tokens}}$ BIO tag associated with the first word.

Fine-tuning for span based Applications :-

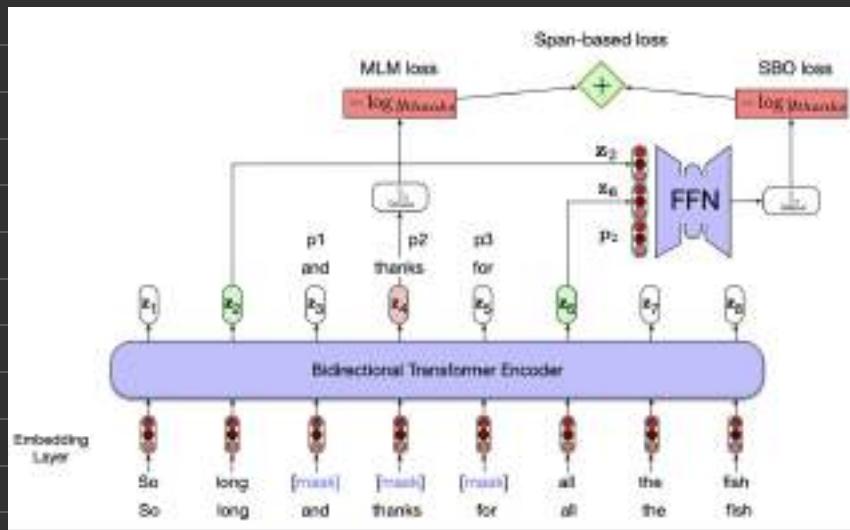
- Once a span is chosen for masking, all the tokens within the span are substituted according to the same regime used in BERT:-
 - 80% are replaced with [MASK]
 - 10% are replaced by randomly sampled tokens
 - 10% are kept as it is
- The spanBERT learning objective augments the MLM objective with a boundary oriented component called the span boundary objective (SBO).

- $x_1, x_2, \dots, x_n \rightarrow \text{input}$
- $z_1, z_2, \dots, z_n \rightarrow \text{output}$

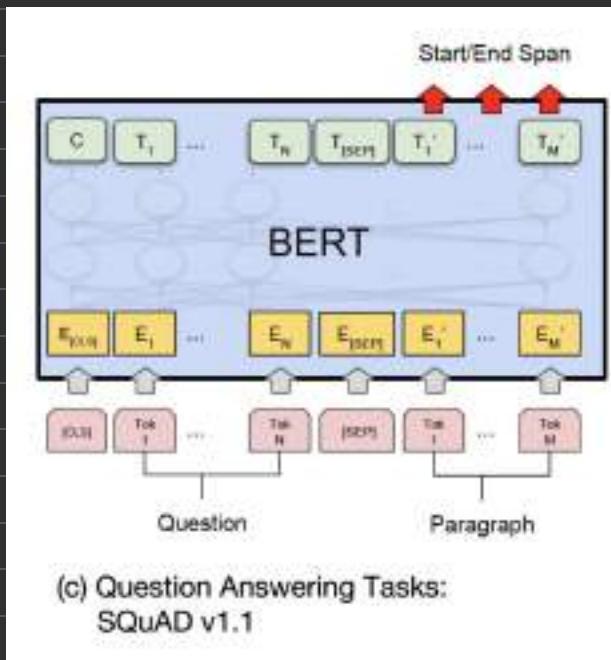
→ A token x_i in a masked span of tokens

(x_s, \dots, x_e) •

$$L_{SBO}(x_i) = -\log p(x_i | x_{s-1}, x_{e+1}, p_{i-s+1})$$



Fine-tuning for SQuAD Dataset :-



(c) Question Answering Tasks:
SQuAD v1.1

- We only introduce a start vector $S \in \mathbb{R}^n$ and end vector $E \in \mathbb{R}^n$ during fine tuning
- The probability of word i being the start of answer span is computed as dot product between T_i and S followed by softmax over all of the words in the paragraph.

$$P_i = \frac{e^{S \cdot T_i}}{\sum_j e^{S \cdot T_j}}$$

The score of a candidate span from position i to position j is defined as $S.T_i + E.T_j$.
 And the maximum scoring span where $j \geq i$ is used as prediction.

* for practical purposes, span-based models often impose an application-specific length limit L , so the legal spans are limited to those where $j-i < L$.

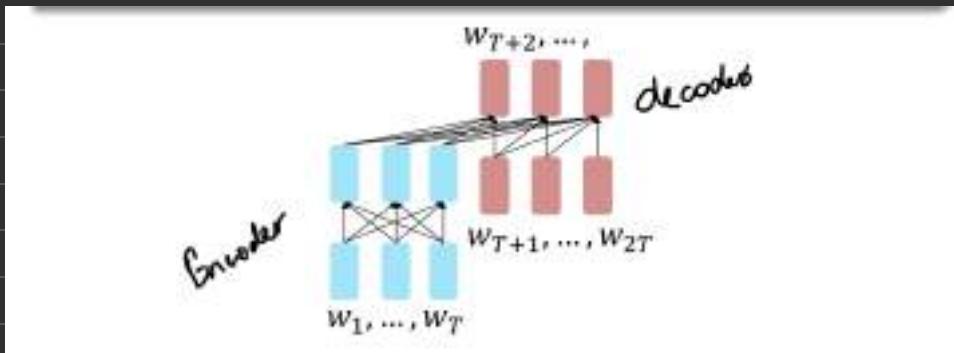
GLUE Benchmark :-

GLUE						
General Language Understanding Evaluation						
Corpus	Train	Test	Task	Metrics	Domain	
Single-Sentence Tasks						
CoLA	8.5k	1k	acceptability	Matthews corr.	misc.	
SST-2	67k	1.8k	sentiment	acc.	movie reviews	
Similarity and Paraphrase Tasks						
MRPC	3.1k	1.7k	paraphrase	acc./F1	news	
STS-B	7k	1.4k	sentence similarity	Pearson/Spearman corr.	misc.	
QQP	36.4k	391k	paraphrase	acc./F1	social QA questions	
Inference Tasks						
MNLI	39k	20k	NLI	mached acc./mismatched acc.	misc.	
QNLI	102k	5.4k	QNLP	acc.	Wikipedia	
RTE	2.5k	3k	NLI	acc.	news, Wikipedia	
WNLI	634	145	coference/NLI	acc.	fiction books	

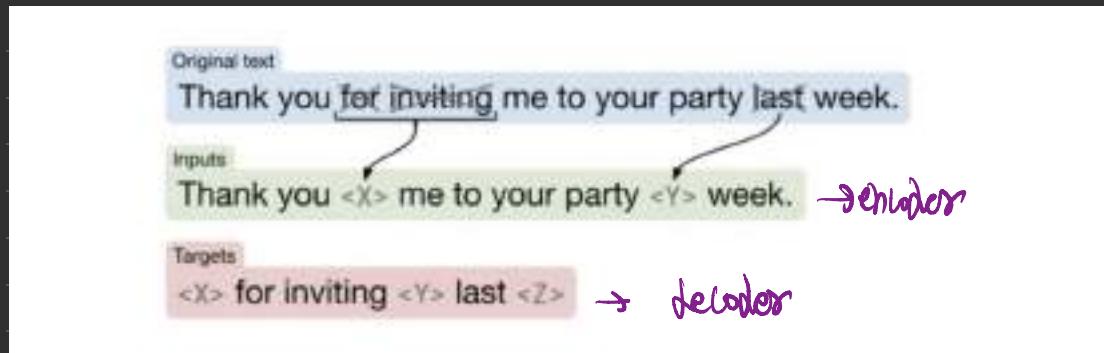
A key takeaway from Roberta paper :-

more compute, more data can improve pretrain even when not changing the underlying Transformer encoder.

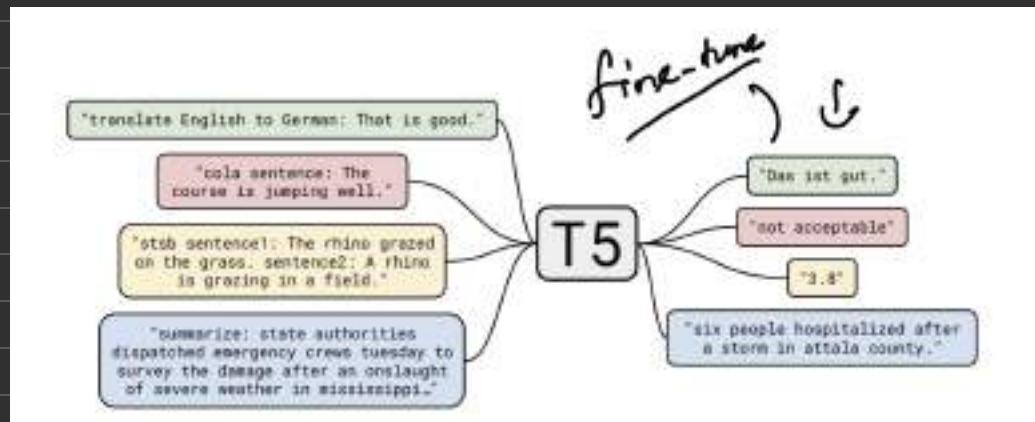
Pretraining encoder-Decoder models :-



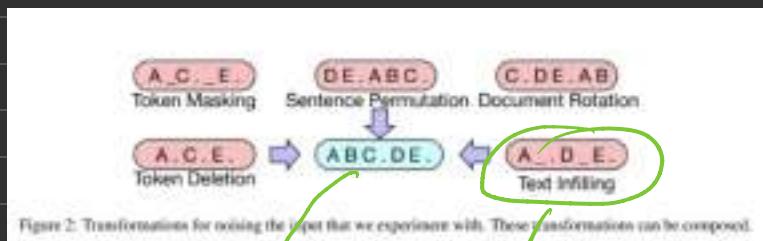
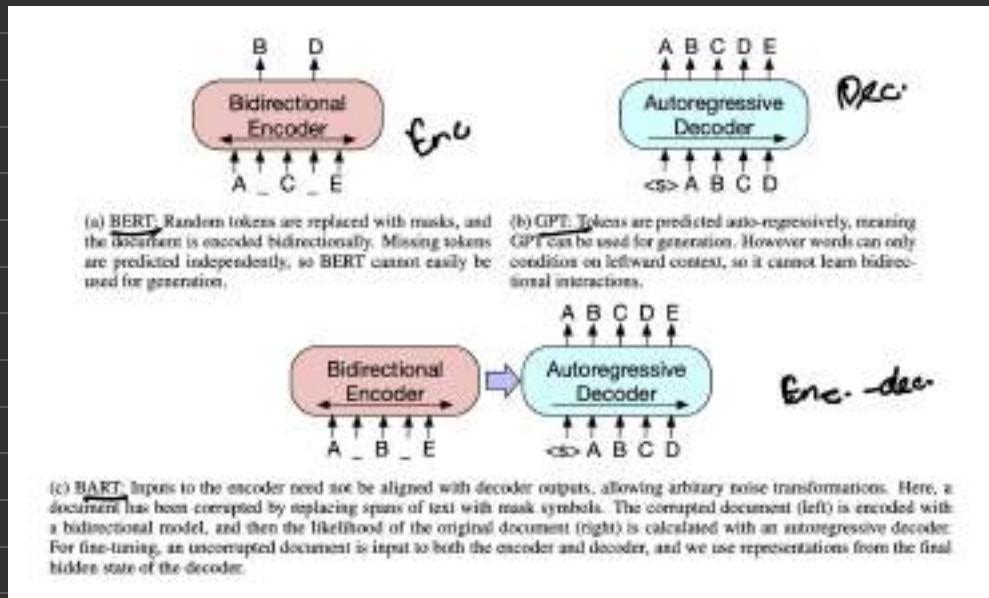
The encoder portion benefits from bi-directional context. The decoder portion is used to train the whole model through language modeling.



The same model loss function, hyperparameters, etc. can be used across all tasks



BART



How to encode an image as a sequence :-

- original image shape :- $H \times W \times C$
- N patches of dimensions $P \times P$: $N = H \times W / P \times P$
- Each patch is mapped to D dimensions through a trainable layer $(P^2 \times C) \times D$
- Standard learnable 1-D positional embeddings are used $B_{pos} \in R^{(N+1) \times D}$
- A learnable class token is prepended, and its representation at last layer is used as image representation.

Decoding :-

- The task of choosing a word to generate based on model probability is called decoding.
- Methods that give bit more weight to the middle probability words tend to be more creative and diverse but likely to be incoherent and less factual.

Greedy decoding :-

at each time step t in generation, the output y_t is chosen by computing the probability for each word in vocabulary.

$$\hat{w}_t = \operatorname{argmax}_{w \in V} P(w / w < t)$$

Beam Search :-

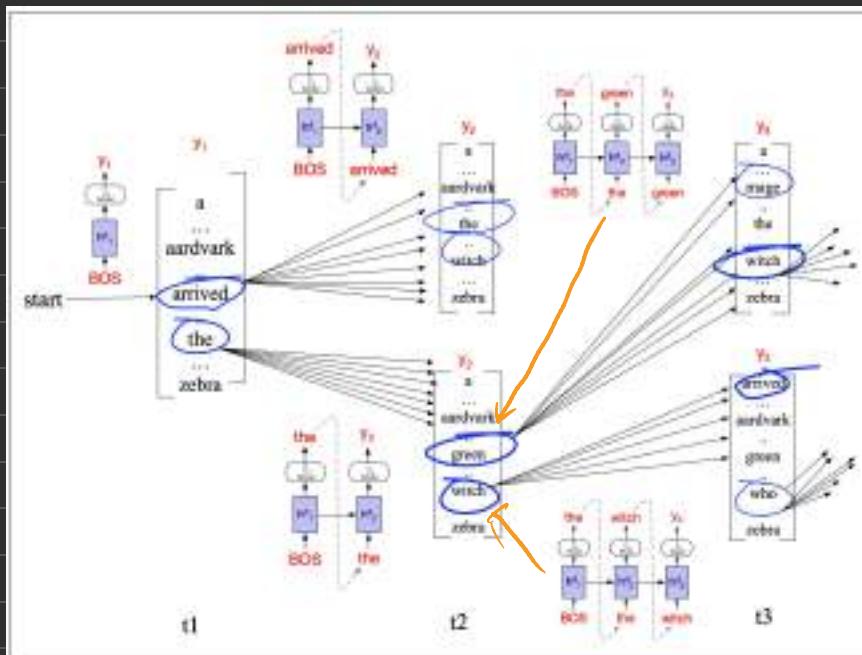
in beam search instead of choosing the best token to generate at each timestep, we keep κ possible tokens at each step.



beam width

- we select the k -best options from this softmax output.
- these k words are called hypotheses.

Schematic View of beam Search :-



- At each time step, $N \times V$ hypotheses are scored by $p(y_i | x, y_{<i})$.
- Select the best n hypotheses again. never more than n decoders or n hypotheses at frontier of search
- Once BOS is generated, we have found the outputs.

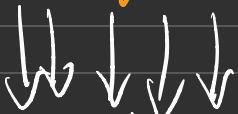
$$\text{Score}(y) = \log p(y|x)$$

$$= \sum_{i=1}^t \log (p(y_1|x) p(y_2|y_1, x) \dots)$$
$$= \sum_{i=1}^t \log p(y_i | y_1, \dots, y_{i-1}, x)$$

Top-n sampling :-

- choose in advance a number of words n
- for each word in the vocabulary V , use the language model to compute the likelihood of this word given the context $p(w_t | \text{rest})$
- sort the words by their likelihood and throw any word that is not one of the top n most probable words.
- normalise the score of n words to be a legitimate probability distribution
- Randomly sample a word from within these remaining n most-probable words according to its probability.

- keeping n-fixed is a big disadvantage, as the probability distribution may change with different contexts.



Top-p Sampling :-

Given a distribution $p(w_0/w_{\text{ct}})$ the top-p vocabulary $V^{(p)}$ is the smallest set of words such that

$$\sum_{w \in V^{(p)}} p(w_0/w_{\text{ct}}) \geq p$$

Temperature sampling :-

- we don't truncate the distribution, but reshape it

(Thermodynamics)

Intuition :- when a system at high temperature at a high temperature, it is flexible and explore many possible states

when a system at low temperature it is likely to explore a subset of lower energy states

modified softmax computation

$$y = \text{softmax}(u/\tau)$$

$\tau \rightarrow \text{temperature}$

$\tau \rightarrow 1$ doesn't make any change

$\tau \leq 1 \Rightarrow \text{inputs become larger} \rightarrow$

probability distribution become more greedy

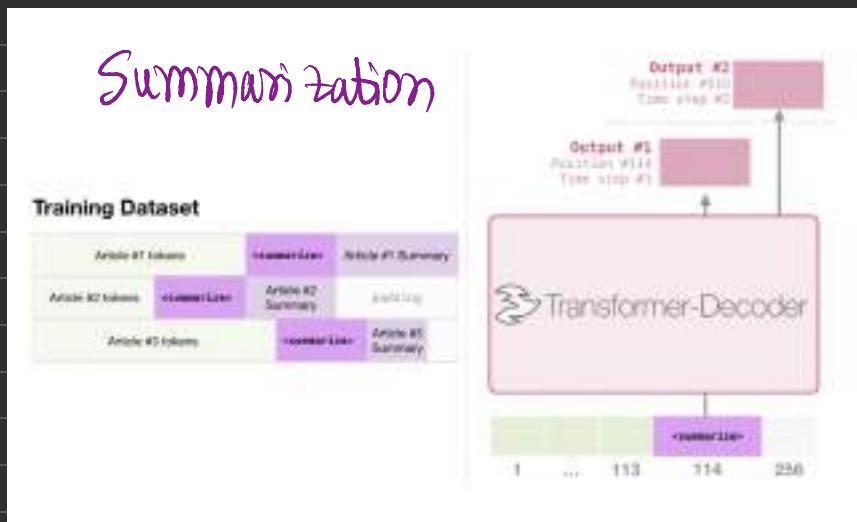
softmax tends to push

high values to 1 and low values to 0

Scaling Laws :-

- the performance of a large language model is hugely affected by
- model size • Training data size
- Amount of compute used for training

formatting the inputs for training on various tasks:



Rough estimation of parameters of a model :-

$$N \approx 2d n_{\text{layer}} (2d_{\text{atten}} + d_{\text{ff}})$$

$$N \approx 12n_{\text{layer}} d^2$$

$$\text{datttn} = \frac{d_{\text{ff}}}{4} = d$$

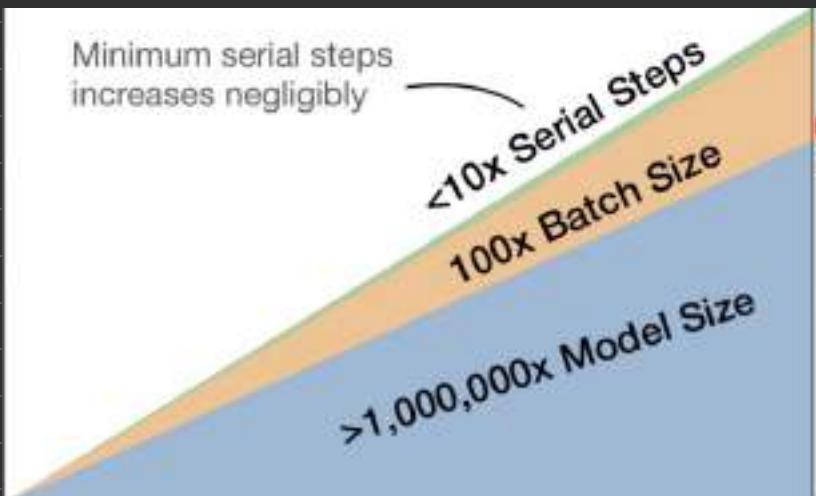
In-Context learning

- very large language models seem to perform some kind of learning without gradient steps simply from examples you provide within their contexts.

Prompting over finetuning

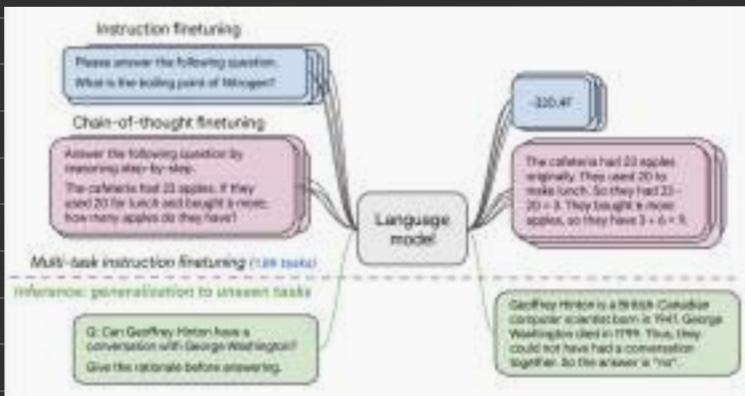
$$L(N) = \left(\frac{n_c}{N} \right)^{\alpha_N} \quad L(D) = \left(\frac{D_c}{D} \right)^{\alpha_D} \quad L(C) = \left(\frac{C_c}{C} \right)^{\alpha_C}$$

Scaling Laws !



as compute increase

Instruction fine-tuning



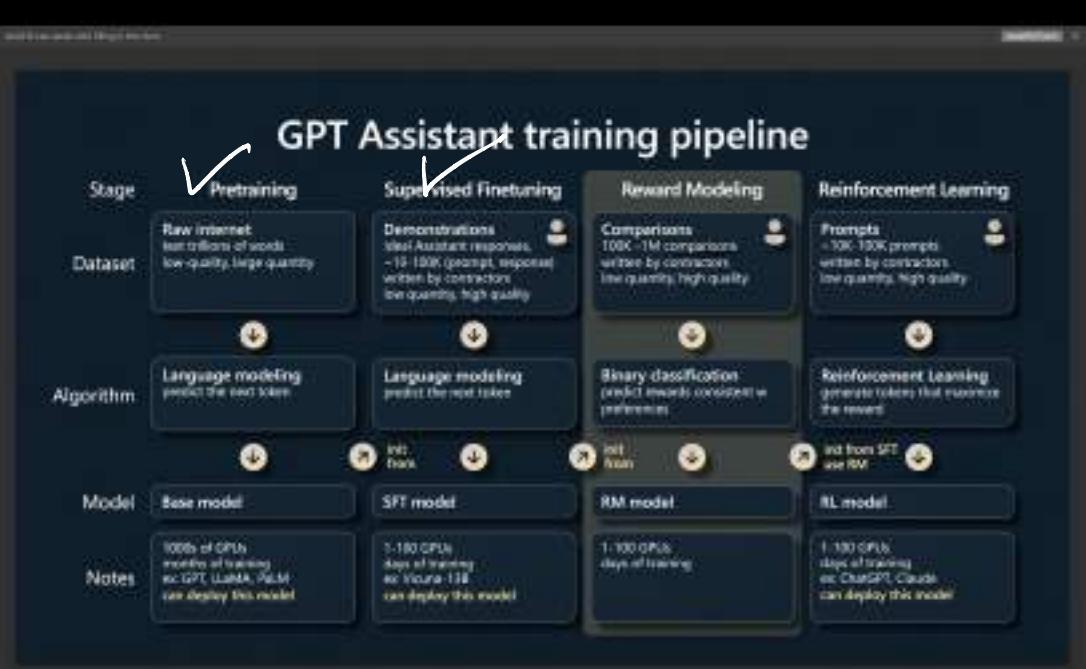
collect pairs of (instruction, output) across many tasks and finetune an LLM.

Major Limitations :-

- expensive to collect data
- No standard answer for open-ended creative generation
- Language modelling penalizes all token-level mistakes equally, but some errors are worse than others.

LM Assistants :-

The base models cannot be used as LM assistants.



Reward Modelling Training :-

RM datasets

prompt:	—	—	completion 1	—	—	<[reward]>	—	—	—	<[reward]>
prompt:	—	—	completion 2	—	—	—	—	—	—	<[reward]>
prompt:	—	—	completion 3	—	<[reward]>	—	—	—	—	—

same prompt tokens

only this output is used

Goal :- predict rewards consistent with preferences

preferences

RL Training :- (maximise the reward)

- The model is now trained on only yellow cells
- it predicts only the <reward>
- The sample tokens becomes labels (labels by the training objective is weighted by the advantage - normalized rewards)
 - ↓ optimise thus using human preferences

for each LM generated sample s , imagine we have a way to obtain $R(s)$ reward of that summary.

Goals :-

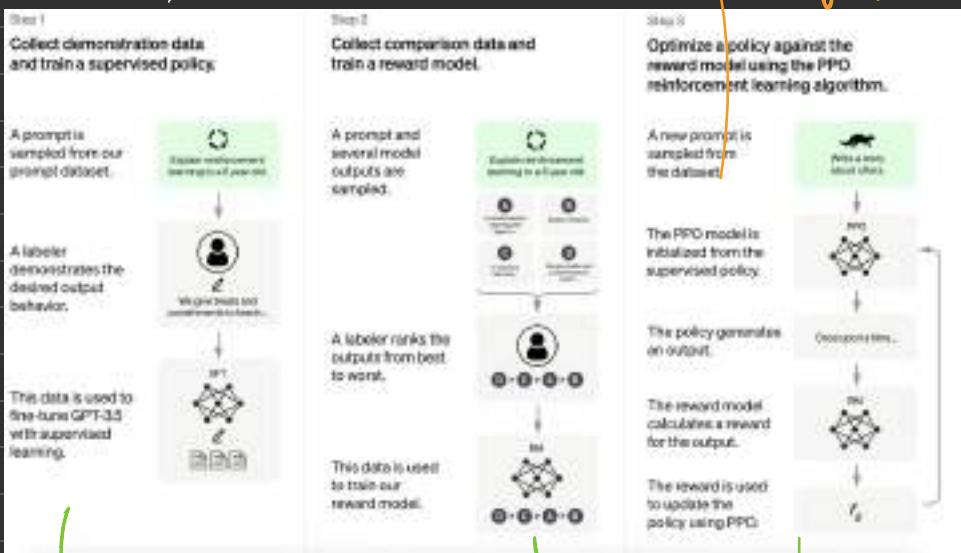
Maximize

$$E_{\hat{s} \sim p_{\theta}(s)} [R(\hat{s})]$$

But how do we do this?

RLHF pipeline :- [Fine-tuning the model on human preferences]

Step 1 :-



↳ Instruction tuning

Reward modelling

RL policy for maximizing the reward.

Reinforce Algorithm :- Cppo

Goal :- $E_{S \sim p_\theta(S)} [R(\hat{S})] \rightarrow \text{maximize}$ this
This may be a non-convex problem

so, better to use gradient ascent

$$\Theta_{t+1} = \Theta_t + \alpha \nabla_{\Theta_t} E_{S \sim p_\theta(S)} [R(S)]$$

↓
Reward function
may be non-differentiable

$$\rightarrow \nabla_{\Theta} E_{S \sim p_\theta(S)} [R(S)] = \sum_S R(S) \frac{\nabla_{\Theta} p_\theta(S)}{\downarrow}$$

$$p_\theta(S) \nabla_{\Theta} \log p_\theta(S)$$

$$\sum_S p_\theta(S) R(S) \nabla_{\Theta} \log p_\theta(S)$$

$$= E_{S \sim p_\theta(S)} [R(S) \nabla_{\Theta} \log p_\theta(S)]$$

↓

This can be averaged over all the prompts

$$\frac{1}{m} \sum_{i=1}^m R(s_i) \nabla_{\theta} \log p_{\theta}(s_i)$$

↓

if reward is +ve,
Take gradient steps to
maximize $p_{\theta}(s_i)$

if reward is -ve
Take gradient steps to
minimize $p_{\theta}(s_i)$

Major issue :- HUMAN in the process

Model the Human preferences as a separate NLP problem

$L^M \rightarrow R^M p(s) \rightarrow$ for optimising the reward

Problem with Human annotations :- Humans

annotations are highly noisy.

- provide annotations as so, that a relative pairwise comparisons judgement can be made.

Training objective

$$J_{RM}(\phi) = -\mathbb{E}_{(s^w, s^d) \sim D} \left[\log \sigma(RM_{\phi}(s^w) - RM_{\phi}(s^d)) \right]$$

winning sample
 ↓
 (should score higher)

losing sample

R_{LHP}

$\rho^T(s) \rightarrow$ Retrained LM

$RM_{\phi}(s) \rightarrow$ reward model (trained)

$\rho_s^L(s) \rightarrow$ copy of $\rho^T(s)$ for PPO.

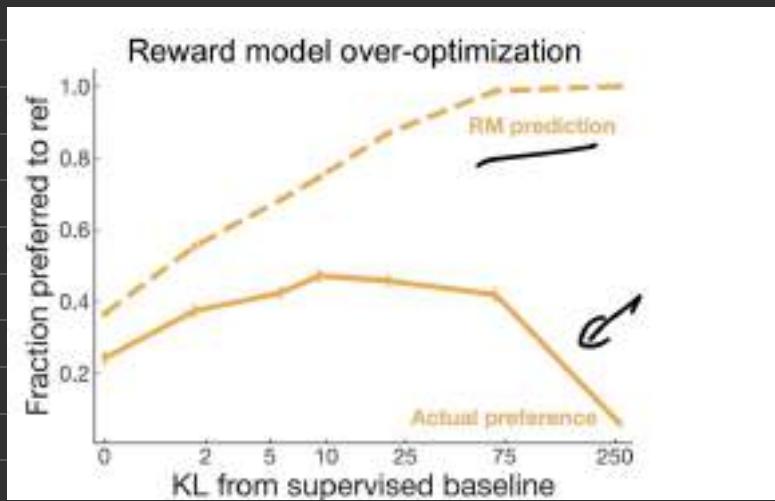
Use of KL Divergence here !

Making sure that the model doesn't deviate by a large amount from its original weights.

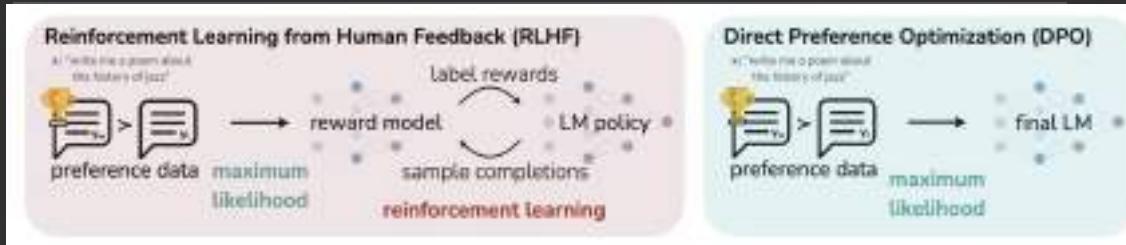
To implement this, a reward penalty is used

$$R(s) = RM_p(s) - \beta \log \left(\frac{p_\theta^{KL}(s)}{p_\theta^{PT}(s)} \right)$$

$p_\theta^{KL}(s) > p_\theta^{PT}(s)$ → Penalty



Direct preference optimization :-



Dpo directly optimises for the policy
vects satisfying the preferences with
a classification objective

$$L_{DPO} = -E_{(\pi_0; \pi_{ref})} [y_w, y_b] \sim D \left[\log \alpha (\beta \log \frac{\pi_0(y_w|x)}{\pi_{ref}(y_w|x)} - \beta \log \frac{\pi_0(y_b|x)}{\pi_{ref}(y_b|x)}) \right]$$

Dpo model

Parameter efficient fine-tuning techniques :-

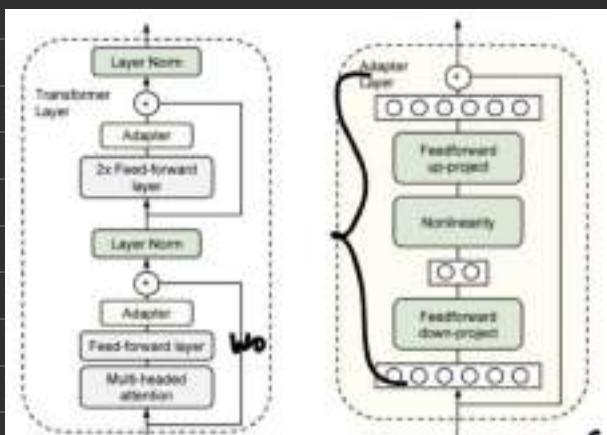
Adapters :-

Fine tuning over a specific task can be composed as

$$f_i(x) = f_{\theta_i}(x) \odot f_{p_i}(x)$$

The original parameters of the model are not changed.

- insert new function of pre-trained model downstream task. f_ϕ between layers to adapt to a



$$\begin{aligned} w_0 &\in \mathbb{R}^{K \times d} \\ w_1 &\in \mathbb{R}^{d \times K} \end{aligned}$$

$$f_\phi(x) = w^T (\sigma(w^0 x))$$

for one chapter :-

Feed-forward down projection : $d_{model} \times d_{model}$

Feed-forward Upward projection : $d_{model} \times d_{model}$

→ for L layers in the decoder, there would be
 $2L$ adapters

→ number of parameters → $2L \times (2 \times d_{model} \times r + d_{model} \times r)$

Pruning for pre-trained models :-

- A fraction of the lowest weights are removed.
- The non-pruned weights are re-trained.

Magnitude pruning : keep weights farthest from 0

Movement pruning : keep weights move the most away from 0.

Major objective of fine tuning :-

Maximize the conditional language modeling
objective :-

$$\max_{\phi} \sum_{(x,y)} \sum_{t=1}^{|y|} \log(P_{\phi}(y_{+|x}, y_{<t}))$$

LORA :-

$$\Delta\phi = \underbrace{\Delta\phi(\theta)}$$

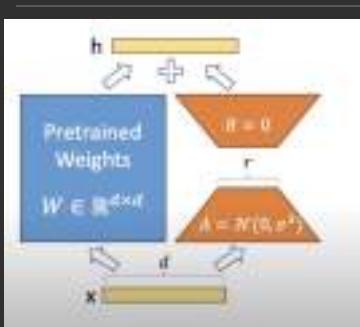
↓

$\theta \rightarrow$ Task specific parameters

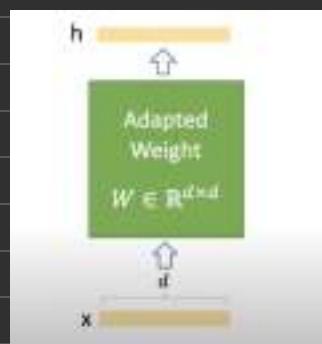
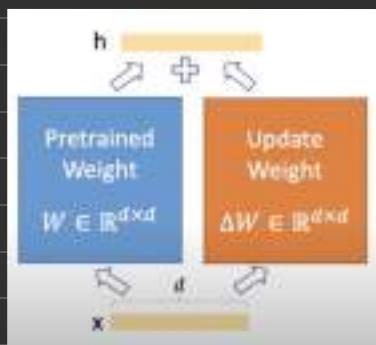
$$|\theta| \ll |\phi_0|$$

The task of finding $|\phi|$ becomes as

$$\max_{\theta} \sum_{(x,y)} \sum_{t=1}^{|y|} \log(P_{\phi_0 + \Delta\phi(\theta)}(y_{+|x}, y_{<t}))$$



Fine Tuning



inference time

$w_0 \in \mathbb{R}^{d \times h}$

$B \in \mathbb{R}^{d \times w}, A \in \mathbb{R}^{w \times h}$

only A and B contain trainable parameters

weight decomposition :- $w = \underbrace{w_0}_{\text{frozen}} + \underbrace{\Delta w}_{BA}$

Forward pass :- $w_0x + \Delta w x$

$$= w_0x + BAx$$

During inference time, $w = w_0 + BA$, no additional latency.

- If LoRA underperforms, increase the value of rank.

Quantization of a pretrained model :-

Ex:- Quantizing a FP32 number to Int8 integer $[-127, 127]$

$$x^{\text{int8}} = \text{round}\left(\frac{x^{\text{FP32}}}{\underbrace{\text{abs max}(x^{\text{FP32}})}_c}\right)$$

Example :- (2-bit)

Maintain absolute bits of values rather than minimize information loss.

map : {Index : 0, 1, 2, 3 → Values : -1, 0.3, 0.5, 103}

input tensor :- [10, -3, 5, 4]

Normalise with absmax : [10, -3, 5, 4]



[1, -0.3, 0.5, 0.4]

Newest values :- [1, 0.3, 0.5, 1.5]

Associated index :- [3, 1, 2, 2] → Store

↓
Dequantisation

[10, 3, 5, 5] → error

→ problem :- if an outlier appears in the input, then some of the bins would not utilize well with few or no numbers in those bins

Block-wise K-bit Quantization

- chunk into contiguous blocks that are independently quantized each with their quantization constant / c.
- NFL Data type ensures every bin has equal number of values assigned from the input tensor.