

# TRANSFORMERS

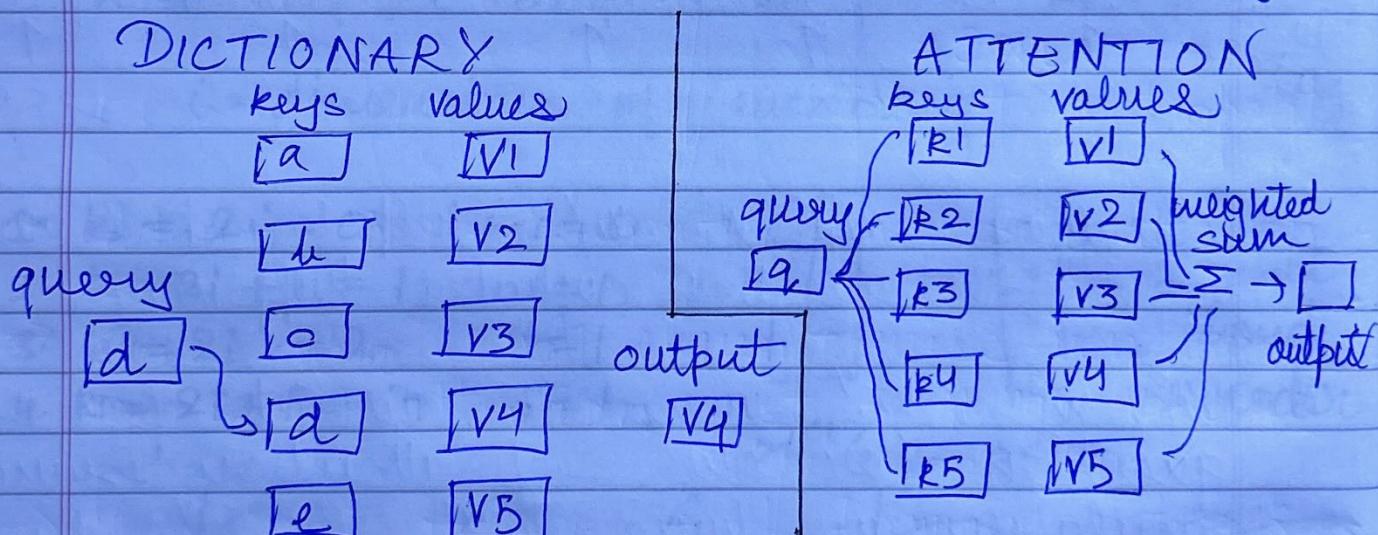
## \* Attention :

- Every word's representation as a query to access & incorporate info from a set of values.

~~transformer~~ is fast for long sequences. No. of unparallelizable operations do not increase with sequence length.

- Max. Interaction Dist :  $O(1)$ , as all words interact at every layer.

Analogy: Attention similar to look up table (dictionary).



query raised, matches a key & we get output as value.

query raised, matches all keys w/r to 1, weighted sum of all values calculated, software output displayed.

Self Attention: Keys, queries, values from same sequence.

For each  $w_i$ , let  $\pi_i = E w_i$  where  $E \in \mathbb{R}^{d \times |V|}$  is an embedding matrix.  $(d, 1) \rightarrow$  word, not entire sentence; for one word, not entire sentence.

1. Transform each word embedding with weight matrices  $Q, K, V$  each in  $\mathbb{R}^{d \times d}$ .

$$q_i = Q \pi_i \quad k_i = K \pi_i \quad v_i = V \pi_i$$

2. Compute pairwise similarities w/w keys & queries; normalise with softmax.

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

softmax  
basically

3. Compute output for each word as weighted sum of values

$$o_i = \sum_j \alpha_{ij} v_j$$

## \* Tailoring Self Attention for Sequence Models

→ Sequence Order of Word:

• Use vectors representing each sequence index.

•  $p_i \in \mathbb{R}^d$ , for  $i \in \{1, 2, \dots, n\}$  are position vectors

$$\tilde{\pi}_i = \pi_i + p_i$$

embedding vector  $\pi_i$

• We do it once at the beginning.  
The position vectors can be:

## → Sinusoidal position representations

$$\mathbf{P}_i = \begin{pmatrix} \sin(i/10000^{2+1/d}) \\ \cos(i/10000^{2+1/d}) \\ \vdots \\ \sin(i/10000^{2+\frac{d}{2}/d}) \\ \cos(i/10000^{2+\frac{d}{2}/d}) \end{pmatrix}$$

- Periodicity indicates absolute position isn't as important
- However, extrapolation isn't possible.

## → $\mathbf{P}_i$ be learnable parameters

$\mathbf{P} \in \mathbb{R}^{d \times n}$ , each  $\mathbf{P}_i$  is a column of that matrix.

Pros

Con

Flexibility

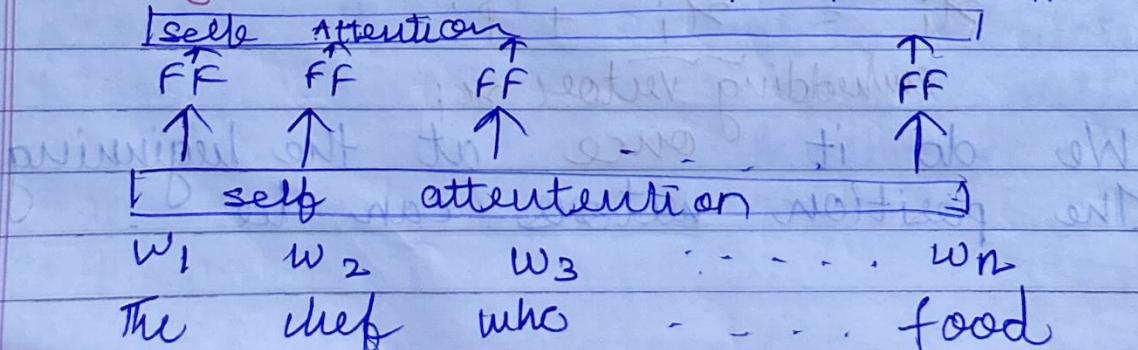
cannot extrapolate

- Adding Non-linearities in Self Attention:

→ There are no element wise non linearities in self-attention, just re-average value vector

→ Thus, we add a feed forward network to post process each output vector.

$$\star \star \quad \left\{ \begin{aligned} \mathbf{m}_i &= \text{MLP}(\text{output}_i) \\ &= \mathbf{w}_2 * \text{ReLU}(\mathbf{w}_1 \text{output}_i + \mathbf{b}_1) + \mathbf{b}_2 \end{aligned} \right.$$



The FF network processes the result of attention.

- Masking the future in Self-Attention
  - To use self attention in decoder, we need to ensure we can't peek in the future;
  - Thus, to enable parallelization, we mask out attention to future words by setting attention scores to  $-\infty$ .

$$e_{ij} = \begin{cases} q_i^T k_j & j \leq i \\ -\infty & j > i \end{cases}$$

	Start	The	Chef	Who
Start		$-\infty$	$-\infty$	$-\infty$
The			$-\infty$	$-\infty$
Chef				$-\infty$
Who				

$-\infty$  part  
word will  
not be  
visible to  
the corres-  
ponding  
row's word.

\* Necessities of Self Attention: Softmax  
Linear (sum)

• Self Attention

this block  
may be repeated  
multiple times  
for encoder

Feed Forward

- Position

Representations

three arrows  
query, key,  
value

Masked Self Attention

• Non Linearities

Add position embeddings

• Masking  
(decoder)

Embeddings  
Inputs

## \* Sequence Stacked form of Attention:

- Let  $x = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{n \times d}$  be the concatenation of input vectors.
- Note that  $XK, XQ, XV \in \mathbb{R}^{n \times d}$  ( $K, Q, V \in \mathbb{R}^{d \times d}$ )
- Output = softmax  $(XQ(XK)^T)XV \in \mathbb{R}^{n \times d}$

Visually:

$$\begin{matrix} XQ \\ n \times d \end{matrix} \quad \begin{matrix} KT XT \\ d \times n \end{matrix} = \begin{matrix} XQ KT XT \\ n \times n \end{matrix} \in \mathbb{R}^{n \times n}$$

$$\text{softmax} \left( \begin{matrix} XQ KT XT \\ n \times n \end{matrix} \right) \begin{matrix} XV \\ n \times d \end{matrix} = \begin{matrix} \text{output} \\ n \times d \end{matrix} \in \mathbb{R}^{n \times d}$$

## \* Multi Head Attention:

What if we want to look in multiple places in the sentence at once?

We define multiple attention "heads" through multiple  $Q, K, V$  matrices.

- Let  $Q_h, V_h, K_h \in \mathbb{R}^{d \times d}$  &  $h$  ranges from 1 to  $h$ .
- Each attention heads performs attention independently

Multihead:

$Q, K, V$  divided into  $Q_h, K_h, V_h$

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

- Then the outputs of all heads are combined
- $\text{output}_h = \text{softmax}(X Q_h K_h^T X^T) \cdot X V_h$  where  $\text{output}_h \in \mathbb{R}^{d/h}$
- $\text{output} = [\text{output}_1, \text{output}_2, \dots, \text{output}_n]^T$  where  $y \in \mathbb{R}^{d \times d}$

Multihed self attention is computationally efficient even though we have  $h$  heads.

for our understanding (&  $XV, XK$ )  
→ We compute  $XQ \in \mathbb{R}^{n \times d}$  & then reshape to  $\mathbb{R}^{n \times h \times \frac{d}{h}}$   
in two steps, → Then transpose to  $\mathbb{R}^{h \times n \times \frac{d}{h}}$ , now head axis is  
actually like batch axis  
→ Almost everything else is identical & matrices  
restored directly are the same size

$$\begin{matrix} n \times d \\ n \end{matrix} \xrightarrow{\text{XQ}} \boxed{XQ} \quad \boxed{K^T X} = \boxed{XQK^T X^T}$$

$3 = h$ , here, meaning results  
the same by reshaping, just that  
 $Q, K, V, X$  for each head are separated now.

$$\text{softmax} \left( \boxed{XQK^T X^T} \right) \boxed{XV} = \boxed{\text{P}}_{\text{mix}} = \boxed{\text{output}} \in \mathbb{R}^{n \times d}$$

- Generally (not always) it is good to have same no. of heads throughout & the no. of heads set to the dimensionality of  $Q, K, V$ .

- When dimensionality  $d$  becomes too large, dot product become large & so input softmax will be large, making gradients small.

Thus we divide attention scores by  $\sqrt{d/n}$  to stop the score from becoming large.

$$\text{output} = \text{softmax} \left( \frac{\mathbf{x}^T \mathbf{e}_k \mathbf{e}_k^T \mathbf{x}^T}{\sqrt{d/n}} \right) * \mathbf{V}_e$$

scaled dot product

- \* Residual Connections:
- Help models train better
- Instead of  $\mathbf{x}^{(i)} = \text{Layer}(\mathbf{x}^{(i-1)})$ , where  $i$  represents layer.

- We let  $\mathbf{x}^{(i)} = \mathbf{x}^{(i-1)} + \text{Layer}(\mathbf{x}^{(i-1)})$   
(so we only have to learn "the residual" from previous layer).

- Gradient through the residual connection is 1

no residual  
GD (Gradient Descent)

residuals  
(GP)

Residual Connection: Implemented in ~~other~~  
 Add & Norm Layer after ~~each attention~~  
 & FFT layer.

Advantages:

- In a deep NN, where outputs of deeper layers are close to unity, it helps the layer to learn faster.
- Avoids vanishing gradient as helps gradient flow as gradient is 1.
- Improves training.

## \* Layer Normalisation

Trick to help models train faster.

- Intuition: Cut down uninformative variation in hidden vector values by normalising to unit mean & standard deviation within each layer.

Let  $x \in \mathbb{R}^d$  be an individual (word) vector in the model.

Let  $\mu = \frac{1}{d} \sum_{j=1}^d x_j$ ; this is the mean,  $\mu \in \mathbb{R}$

$$\sigma = \sqrt{\frac{1}{d} \sum_{j=1}^d (x_j - \mu)^2} \rightarrow \text{Standard dev.}$$

$$\sigma \in \mathbb{R}$$

$\gamma \in \mathbb{R}^d$ ,  $\beta \in \mathbb{R}^d$  be learned gain & bias.

Thus: output =  $\frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} * \gamma + \beta$

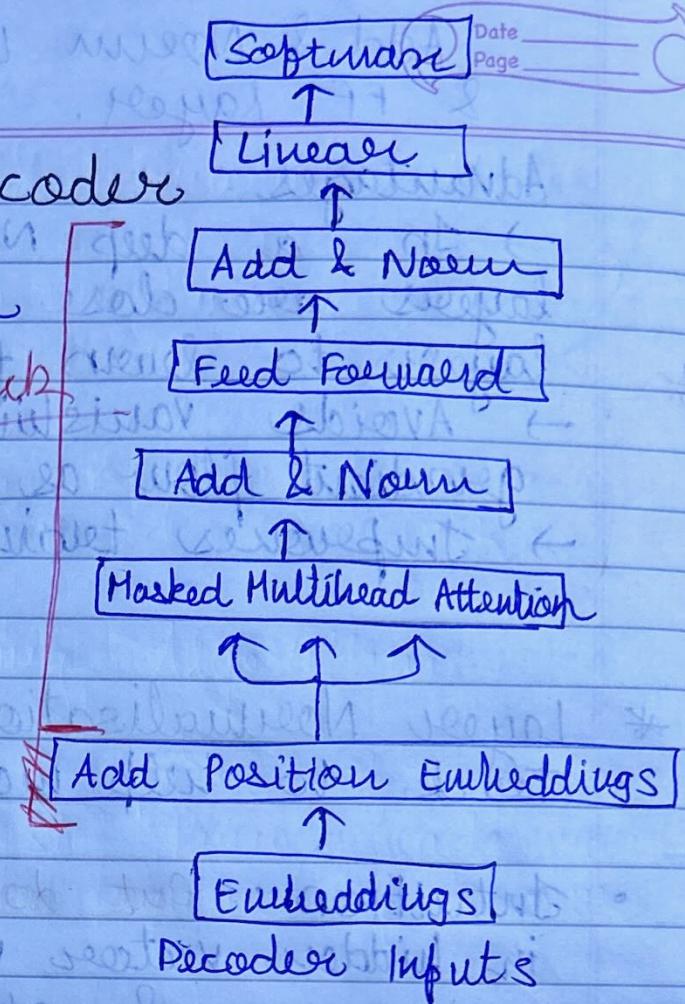
$\mu$  is broadcasted (const) to dimensionality

where  $x_j = x^{(j-1)} + \text{layer } (\gamma^{(j-1)})$  {Residual connection.}

## \* Transformer Decoder

- Decoder is a stack of blocks.
- Each block consists of:
  - Self Attention
  - Add & Norm
  - Feed Forward
  - Add & Norm

block

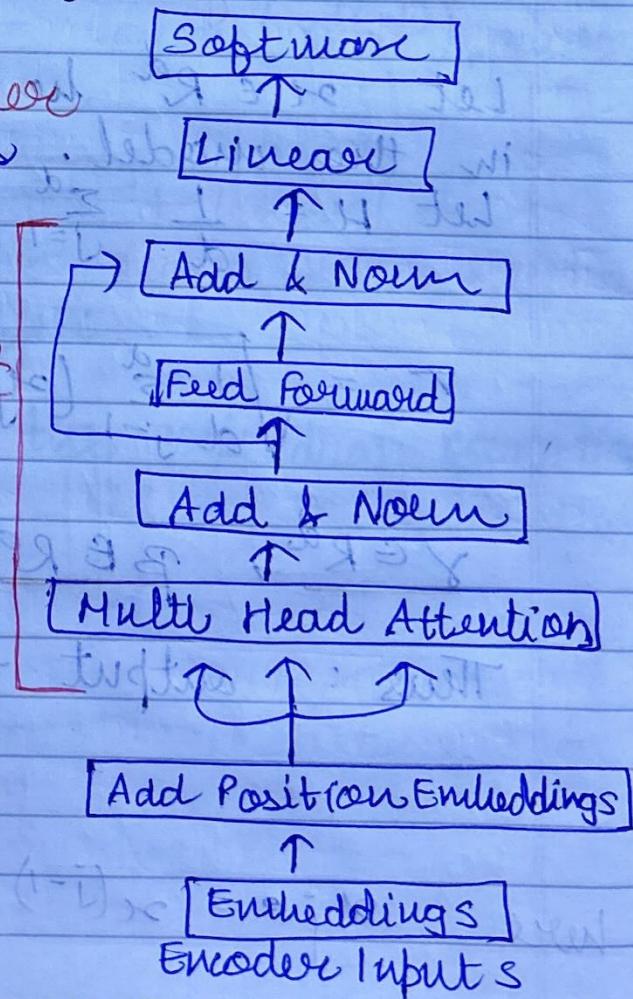


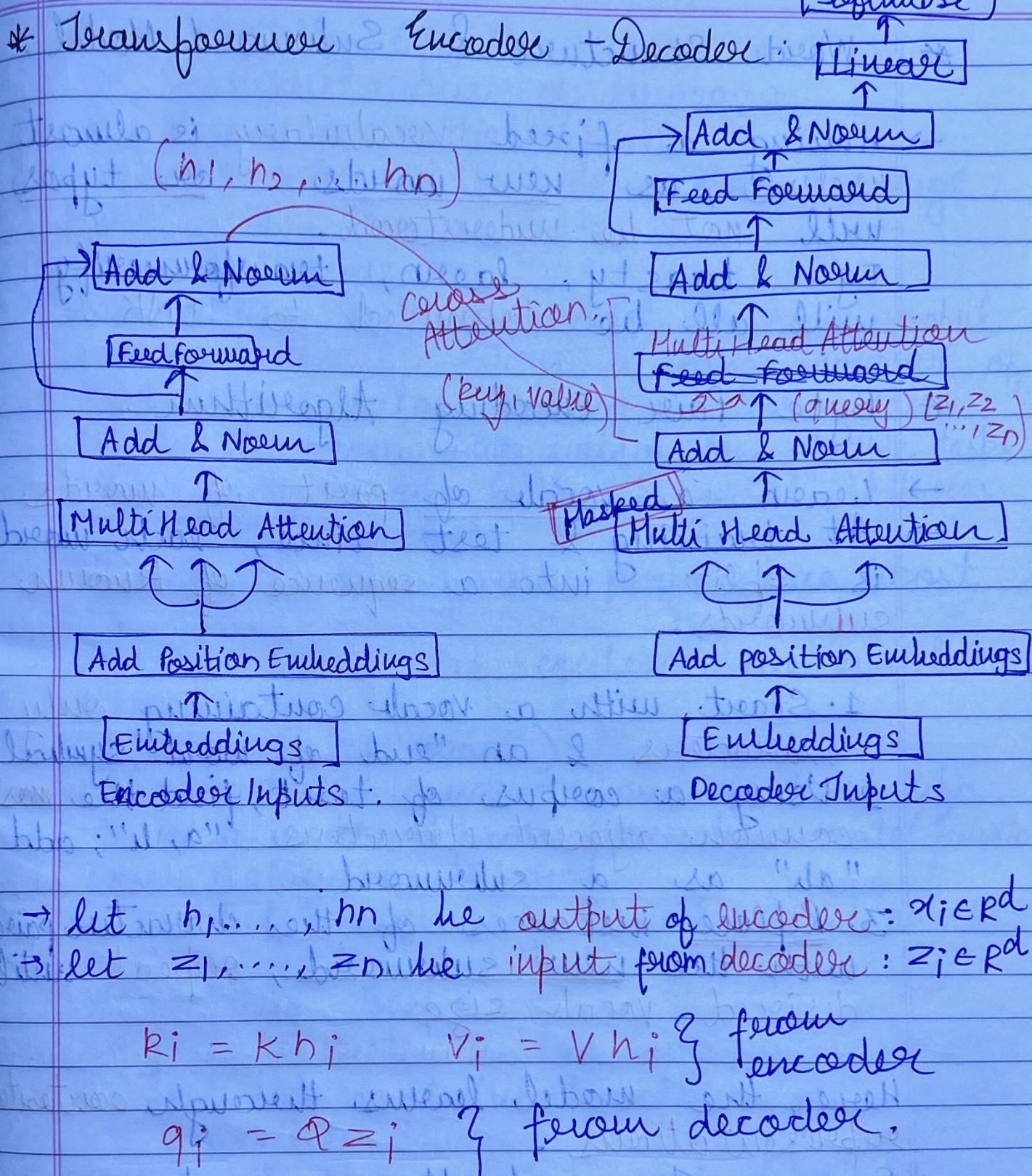
## \* Transformer Encoder

Now, transformer decoder is unidirectional as for language models.

- However, if we want **bidirectional context**, we do:
  - remove masking from encoder

block



Software

- \* Drawbacks of Transformer
- Quadratic Computation of  $\text{tf}_i$
  - Position representations

Solutions 2)

## \* Word Structure & Subword Model:

→ Having a fixed vocabulary is almost useless, as new words or types will not be understood.

Eg: tasty, laun, transformerify will all be tokenised to UNK.

## \* Byte Pair Encoding Algorithm:

→ Learn a vocab of pair of words

→ At training & test time, each word is split into a sequence of known subwords.

1. Start with a vocab containing only characters & an "end of word" symbol.

2. Using a corpus of text, find the most common adjacent characters "a, u"; add "au" as a subword.

3. Replace instances of the character pair with the new subword, repeat until desired vocab size

Here, the model learns through contextual representations.

For eg: I record the record.

both the "records" will have their diff. meaning here, while word2vec would just fail.

Pretraining is

Exceptionally Effective in:

- representations of language
- parameter initialisation
- probability distribution over language.

labelled  
data is  
expensive.

This is because:

- Pretraining allows model to learn useful representations from a large amount of unlabeled data, making it more data efficient when fine tuning on specific tasks.
- Helps the model learn general features of the language.
- Faster convergence, performance boost.

By predicting the context the model can learn:

- patterns (fibonacci)
- understand scenario based off test info.
- using context can remember male/female or singular/plural.
- predict geography & understand world.

Pretraining Steps:

Step 1: Pretrain (on language modelling)  
lots of text (unstructured data), learn general things.

Step 2: Finetune (on your task)

Not many labels, adapt to the task

goes to make END

↑ ↑ ↑ ... ↑  
[ ]  
↑ ↑ ↑ ... ↑  
From goes to tea

Finetune

↑ - - - ↑ (sentiment analysis)  
[ ]  
↑ ↑ ↑ ... ↑  
... the movie was ...

Pretrain

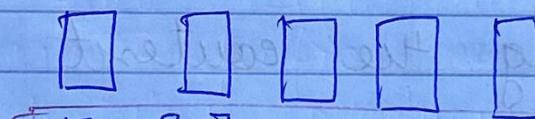
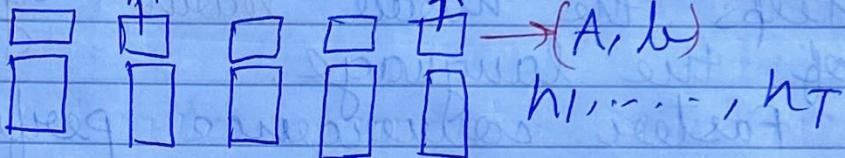
## \* Pretraining Encoders:

→ Encoders get bidirectional context, so language modelling is now possible.

Idea: replace some fraction of words in the input with a special [MASK] token; predict these words.

$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

$$y_i \sim A_{\text{multi}} h_i + b_{\text{store}}$$



I [M] to the [M]  $\approx$

Only add loss terms from words that are "masked out". If  $\tilde{x}$  is the masked version of  $x$  we're learning  $P(\tilde{x}|x)$  "masked LM"

\* BERT: Bidirectional Encoder Representations from Transformers

Masked LM for Bert!

- Predict a brandom 15% of (sub)word:
    - replace input with [MASK] 80%.
    - replace input with random token 10%.
    - leave input word unchanged 10%.

- In the pretraining input to BERT was two separate contiguous chunks of text.
- Although this was not necessary to predict the "text mask" it gave the model a hard time.

### Intuition:

We're coming up with hard problems for the network to solve such that by solving them it has to learn a lot about language. We're defining those problems by making simple transfer or removing info, or both randomly.

- The best usecase of BERT is NOT next-word prediction or generating a large sequence. It works well masked inputs.

- \* Fine Tuning & Light Weight Fine Tuning:
  - works great, → expensive
  - memory intensive

- \* Prefix Tuning, Prompt Tuning
  - initialising some parameters in the beginning & only fine tuning them.
  - We can also split input matrix into low rank matrices as  $A \in \mathbb{R}^{d \times k}$ ,  $B \in \mathbb{R}^{k \times d}$

- \* Span Corruption:

Masks a seq of words in the encoder-decoder structure & gives an output of masked seq. of words.

Majority of the language models use decoder only, the reason is speculated to be that's its easier (less complex) & does not require data to split divide in encoder / decoder.

GPT → BERT → GPT2 → GPT3

much larger  
more focussed on the generation part of sent than GPT2, up to 100 million more.