

Stanford CS24N - I

* Issues with RNN's

- RNN's unroll $L \rightarrow R / R \rightarrow L$
- This induces linear locality
i.e. Nearby words affect each other's meaning
- RNN's are sequential thus, they take $O(\text{seq len})$ steps to let far distant words to interact.
- Forward & backward passes have $O(\text{seq len})$ of unparallelizable operations.
- This is bcz future RNN hidden states can't be computed without computing all the previous states.

* How does attention treat this?

Attention for a single sentence-

It treats each word step.
as a query to access.

→ No. of parallelizable operations do not increase with seq len

The nos. in boxes rep the no. of unparallelizable operations

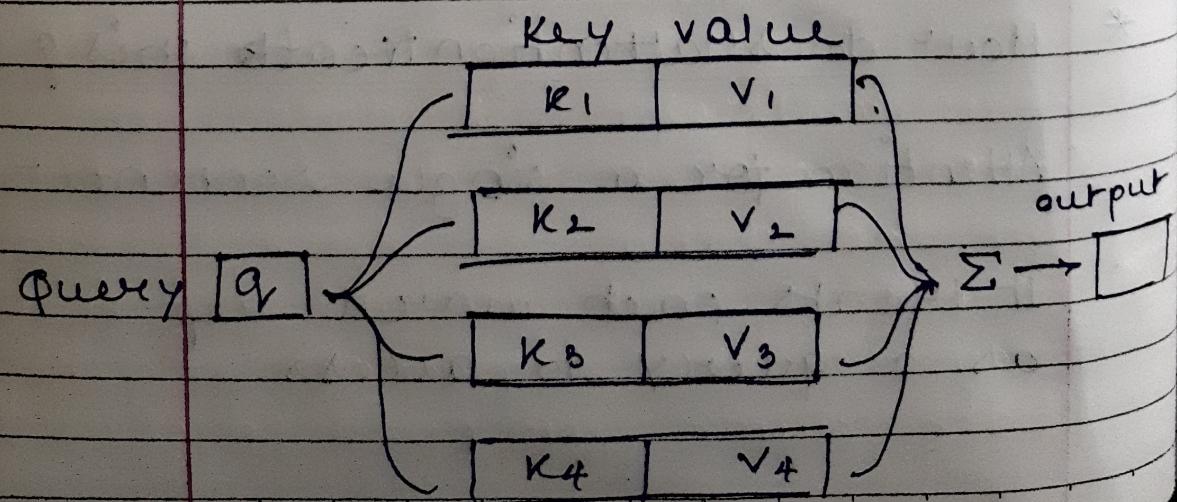
Embedding: $\boxed{0} \quad \boxed{0} \quad \boxed{0} \dots \boxed{0}$
 $h_1 \quad h_2 \quad h_3 \quad \dots \quad h_t$

Attention: $\boxed{1} \quad \boxed{1} \quad \boxed{1} \dots \boxed{1}$

Attention: $\boxed{2} \quad \boxed{2} \quad \boxed{2} \dots \boxed{2}$

* Query, Key & Value

In attention, the query matches all the keys softly to weigh between 0 & 1. The key's 'values' are multiplied by weights & summed.



* Self-Attention

Let $w_{1:N}$ be a seq. of words in vocab V , like zuko made his uncle tea.

For each w_i , $x_i = Ew_i$ where E is an embedding matrix

1. Transform each word emb with weight matrices K, Q, V

$$q_i^* = \Phi x_i \quad \text{query}$$

$$k_i^* = \Psi x_i \quad \text{key}$$

$$v_i^* = \Upsilon x_i \quad \text{value}$$

2. compute pairwise similarity btw keys and queries; normalize with softmax.

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

This basically means, how much i should look at j .

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

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

* Sequence Order

Self attention doesn't build in order information, so we encode the position of words in our keys, queries & values

We add p_i to our input of embedding matrix x_i . Thus, the 'positional' encoding becomes:

$$\tilde{x}_i = x_i + p_i$$

* Learned absolute position info:

We let all p_i to be learnable parameters. This way each position gets to be learned to fit the data, so most systems use this. But it can't handle beyond the sequence length. Each column map position of word

it has seen during training
this is why models have a
input token limit.



Non-linearity in self attention

When we add multiple self attention layers, it just re-avg the value vector and ends up becoming one big self-attention layer.

To fix this, we add a FFN to post-process each value vector. Basically, the FFN processes the result of attention



Masking future self-attention

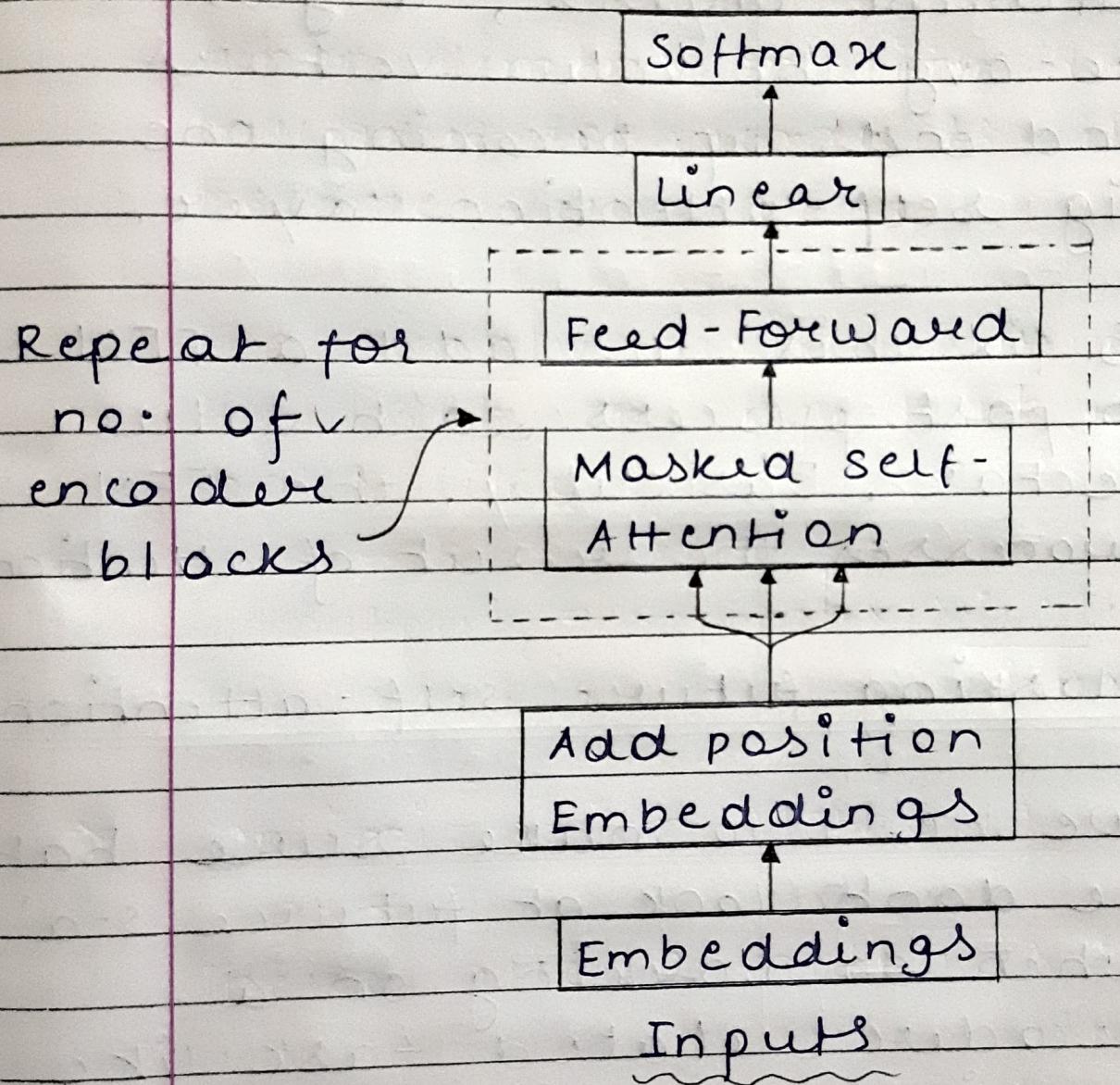
We need to make sure that we don't 'look' at future S-A when ~~is~~ working on a decoder related task like machine translation and language modelling.

To efficiently enable parallelization, we mask out the

attention to future words
in this manner :

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

* Self-Attention building block



* Multi-Head attention

This contains multiple attention heads tasked with diff types of attention. For ex, 1st head might attend to meaning of the words, 2nd head might check the relevance and so on.

$x = [x_1; \dots; x_n]$ be concatenation of input vectors.

$$\text{Output} = \text{softmax}(xQ(xK)^T) xV$$

I. We take query - Key dot prod in one matrix multiplication $xQ(xK)^T$

II. Softmax above step and matrix multiply with xV to compute weighted average

Let the no. of attention heads be ' h ', and ' i ' ranges 1 to h . We get $Q_i, K_i, V_i \in \mathbb{R}^{d \times d_h}$

We compute output, for each head and concatenate it in one final output.

$$\text{output} = [\text{output}_1; \dots; \text{output}_n]$$

* scaled dot product

When dimensionality d becomes large, dot products also become very large

We simply divide our ~~output~~^{scaled}, by $\sqrt{d/h}$ to stop that. It looks like :

$$\text{output}_i = \text{softmax}\left(\frac{x_i Q_i (x_k)^T}{\sqrt{d/h}}\right) x_k$$

* optimization Methods

There are 2 methods :

1. Residual connections

In deep neural networks, we usually face the problem of vanishing & exploding gradients.

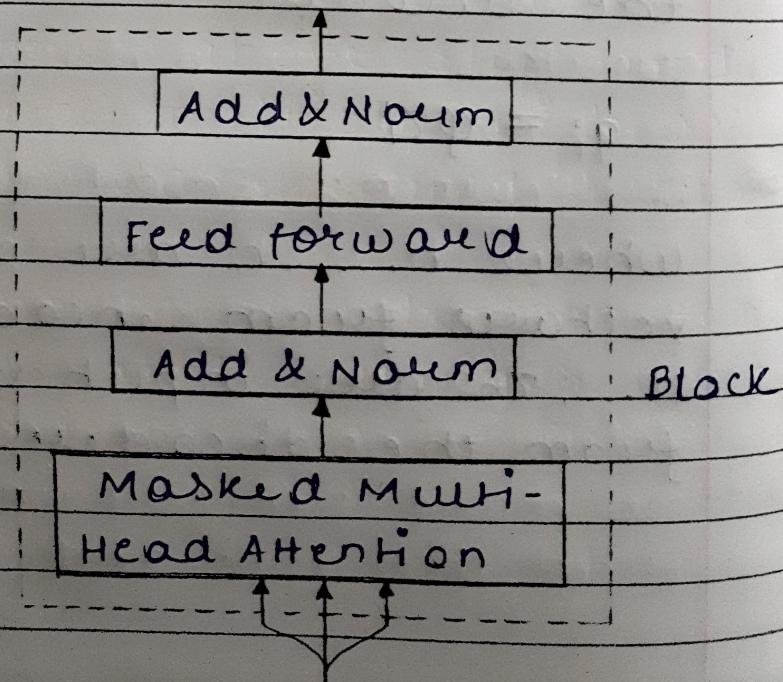
so instead of $x^{(i)} = \text{layer}_k(x^{(i)})$
we use $x^{(i)} = x^{(i)} + \text{layer}_k(x^{(i)})$

This way the model only has to learn how much to change the input to get correct output. This change is 'residue'.

2. Layer Normalization

The basic idea is to normalize gradients in each layer. This is achieved by normalizing the unit mean and standard deviation within each layer.

* The transformer decoder



* Transformer Encoder-Decoder

The encoder architecture is mostly similar to decoder. Main difference is that encoder uses multi-head and not masked multi-head attention because it is better to look at all attentions in input sequence for an encoder.

The keys and values are drawn from the encoder (like memory)

$$k_i^* = K h_i, \quad v_i^* = V h_i$$

The queries are drawn from the decoder

$$q_i^* = \varphi z_i^*$$

where h_1, \dots, h_n are output vectors from encoder and z_1, \dots, z_n are input vectors from the decoder.