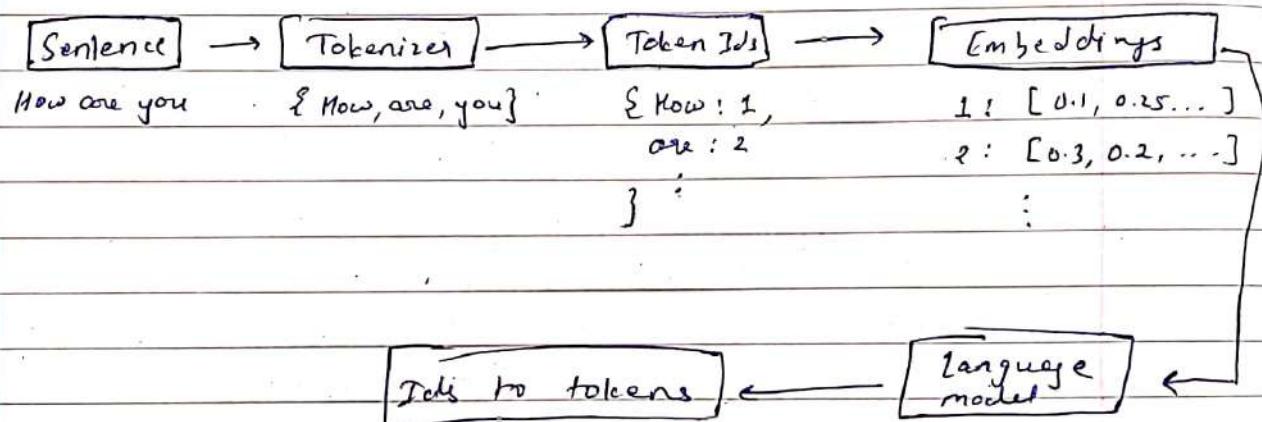


* Tokenization

Given a sentence, breaking it down into words → sub-words and then constructing an embedding matrix that represents all of those words/sub-words(tokens) is called tokenization.

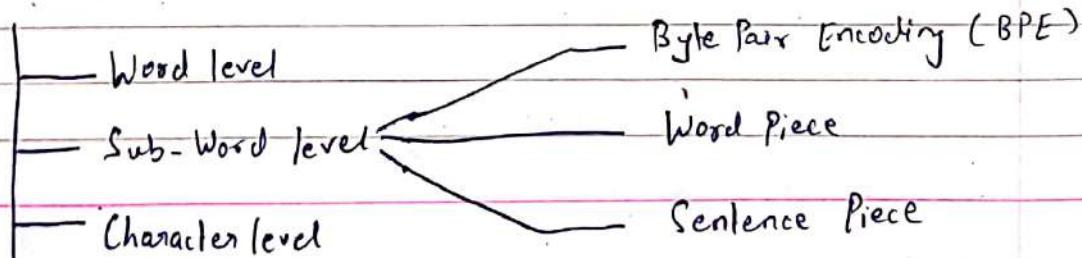
For LLMs following is the flow



→ Challenges with tokenization (we will try to deal with these in the approaches that we discuss)

- Picking optimal vocabulary size
- Out of Vocabulary
- Handling misspelled words
- Relationship between similar words (enjoy, enjoyed, enjoying)

→ Types of Tokenization



• Word Level Tokenization

Given a sentence we pick all the space separated unique words as our vocabulary.

(lets say instead of a sentence we have a huge text dataset that gives us 500k words)

The size of vocab/embedding matrix will be $500k \times d_{model}$ so with every word d_{model} no. of parameters get added which is huge

Other problem is the output matrix is sequence length $\times 500k$ and then we do softmax for each row, so softmax is calculated over 500k size which is computationally very expensive

• Character level Tokenization

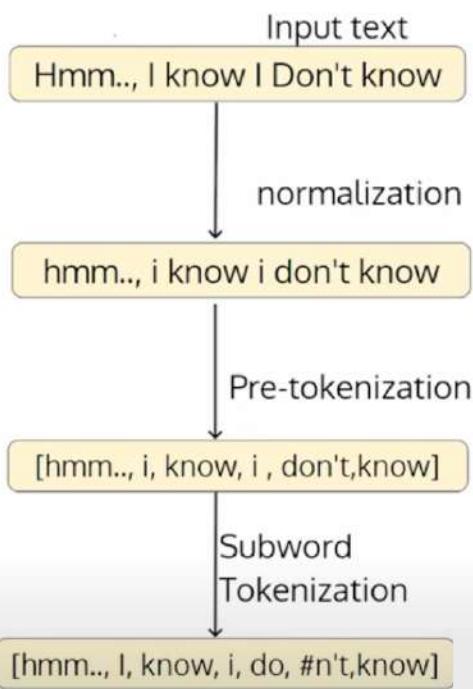
If we do this we would have < 100 total tokens in our vocabulary as no. of characters are limited.

Here the problem would be sequence length, for mere 5 words you would need to have approx sequence length of 25 & due to this all the weight matrices in the encoder-decoder block also increase & computations also increase.

==
All the challenges remain as it is, we now need something in middle like sub-word tokenization. Neither word nor character.

Ishan Modi

• Sub-Word Tokenization



What sub-word tokenization will give us is the following

text → I enjoyed playing cricket

Vocab = [I, enjoy, #ed, play, #ing, cricket]

means it is a split

this will be useful at the end, when model predicts words with # you merge it to previous word to get original sentence

How do you split data into tokens?

- vrc ⇒ BPE, Word Piece, Sentence Piece

Ishan Modi

• Byte Pair Encoding (BPE)

Let's say we have the following dummy vocabulary textual data.

low lower newest widest

In real world we might have a lot of words with meaningful sentences (e.g. a paragraph etc.)

Step 1

Break the sentence into words by splitting on space, remove special characters & make everything lowercase etc.

- low

- lower

- newest

- widest

Following are all the unique words in the

Step 2

Count the character frequencies (We start off with character & move towards word A stop somewhere in the middle)

l - 2

o - 2

w - 4

e - 5

r - 1

n - 1

s - 2

t - 2

i - 2

d - 1

This is also the initial vocabulary

that we have. We keep adding to this

Step 3

Now, starting from the first ~~word~~ character in the above table start building pairs

Ishan Modi

~~Merge¹~~

$(l, o) - 2$	$(e, \omega) - 1$
$(o, \omega) - 2$	$(e, s) - 2$
$(\omega, e) - 2$	$s - [s, t] - 2$
$(e, r) - 1$	$t - [w, i] - 1$
$(n, e) - 1$	$i - [i, d] - 1$

$\Sigma = 9$, $\lambda = 1$, $\mu = 1$

Highlighted pairs have the maximum frequency, so we pick the first pair that occurred in the corpus ie (l, o)

Now, we do the merge operation on the selected pair to add 'lo' to the vocabulary. Since we are adding 'lo' we need to subtract the frequency from individual 'l' & 'o'

Thus, Vocab looks like following

$l - (2-2) = 0$	$n - 1$	$\text{lo} = \text{M12}$
$o - (2-2) = 0$	$s - 2$	
$w - 4$	$t - f - 2$	$g - 3$
$e - 5$	$i - j - 1$	$l - 2$
$r - 1$	$d - k - 1$	$c - 8$

$\Sigma = 9$, $\lambda = 1$, $\mu = 1$

This is the completion of merge ①

~~Merge²~~

Pairs,

$(lo, \omega) - 2$	$r - (e, \omega) - 1$	$\lambda = 2$
$(\omega, e) - 2$	$s - [e, s] - 2$	$\lambda = 0$
$(e, r) - 1$	$t - [s, t] - 2$	$\lambda = 0$
$(n, e) - 1$	$(w, i) - 1$	

$\Sigma = 4$, $\lambda = 1$, $\mu = 1$

Ishan Modi

Note- we removed the pair (l, o) & (o, w) because we have a larger string that represents both. Also the frequency of 'l' & 'o' are zero, it means they do not exist independently if 'lo' is present

Merged,

$$l = 0$$

$$o = 0$$

$$w = (4-2) = 2$$

$$e = 5$$

$$r = 0 \text{ (1 - 0)}$$

$$n - 1$$

$$s - 2$$

$$t - 2$$

$$i - 1$$

$$d - 1 - n$$

$$lo = (2-2) = 0$$

$$low = 2$$

$$ew = 0$$

$$o = 0$$

$$o = 0$$

$$s = 0$$

Merge³

Pairs,

$$(low, e) = 1$$

$$(e, r) = 1$$

$$(w, e) = 1$$

$$[(e, s) = 2]$$

$(l, o) = 0$ (as $l = 0$) $\rightarrow [l, o] = 0$ (as $l = 0$) $\rightarrow [l, o] = 0$ (as $l = 0$)

$(e, w) = 1$ (as $e = 5$) $\rightarrow [e, w] = 1$ (as $e = 5$) $\rightarrow [e, w] = 1$ (as $e = 5$)

$$(i, d) = 1$$

Note- Their order has changed because we started building pairs using 'low lower newest oldest' using the latest vocabulary. \leftarrow low = recent \leftarrow high = old

Merged,

resulted after merging all sorted stuff

$$l = 0$$

$$o = 0$$

$$w = 2$$
 (as $w = 2$)

$$e = (5-2) = 3$$

$$r = 1$$

$$s = n - 1$$

$$t = 2$$

$$i = 1$$

$$d = 1$$

$$l = 0$$

$$lo = 0$$

$$low = 2$$

$$ew = 0$$

$$o = 0$$

$$s = 0$$

Ishan Modi

~~Merge⁴~~

and pairs, and change their parts merged to final

Final state part merged to final state

(low, e) - 1

(w, e) - 1

(c, r) - 1

[es, t] - 2

(n, e) - 1

(w, i) - 1

(e, w) - 1

(i, o) - 1

$S = (x, z) - 0$

$L = \emptyset$

$D = \emptyset$

Merged,

$L = \emptyset$

$S = (x, z) - 0$

$L = \emptyset$

$D = \emptyset$

$l = 0$

$n - 1 - 1$

$lo - 10 - 0$

$o - 0$

$s - 0$

$low - 2$

$w - 2$

$t - (2-2) = 0$

$es - (2-2) = 0$

$e - 3$

$i - 1$

$est - (2-2)$

$r - 1$

$d - 1$

$L = (s, e)$

$L = (s, w)$

$S = (z, g)$

$L = (s, z)$

Similarly you can continue with as many merges as you would like / till there are no more pairs to construct

In practice we have </w> character following each word.

'low' </w> lower </w> newest </w> widest </w>

This helps us differentiate between

newest </w> & estate

c - est \rightarrow $g - (s - 5) \rightarrow$

$g - 1$

$g - 0$

Generally we do 30k - 50k merges this means vocab size = (30k to 50k) + initial 100 characters.

The number of merges = vocab size because each merge adds an element to vocab

Ishank Modi

Step 4

Get the splits in the data using the vocabulary
 and follow newest widest till now ($</w>$ is used)
 in practice

Let just pick one word 'newest' for simplicity. We
 iterate through the vocab in the exact order.

n | e | w | e | s | t

char based split

s - (e, e)

n | e | w | e | s | t

encountered 'e' so
 merged

n | e | w | e | s | t

encountered 'est' so merged

So the word 'newest' is broken/ tokenized to 4
 tokens 'n', 'e', 'w', 'est'

```

import re
from collections import defaultdict

def get_stats(vocab):
    pairs = defaultdict(int)
    for word, freq in vocab.items():
        symbols = word.split()
        for i in range(len(symbols)-1):
            pairs[symbols[i], symbols[i+1]] += freq
    return pairs

def merge_vocab(pair, v_in):
    v_out = {}
    bigram = re.escape(''.join(pair))
    p = re.compile(r'(?<!\S)' + bigram + r'(?!\S)')
    for word in v_in:
        w_out = p.sub(''.join(pair), word)
        v_out[w_out] = v_in[word]
    return v_out

def get_vocab(data):
    vocab = defaultdict(int)
    for line in data:
        for word in line.split():
            vocab[''.join(list(word)) + '</w>'] += 1
    return vocab

def character_to_word_bpe(data, n):
    vocab = get_vocab(data)
    for _ in range(n):
        pairs = get_stats(vocab)
        best = max(pairs, key=pairs.get)
        vocab = merge_vocab(best, vocab)
    return vocab

# Example usage
corpus = "low lower newest widest"
data = corpus.split('.')
n = 4 # Number of merge operations
bpe_vocab = character_to_word_bpe(data, n)
print(bpe_vocab)

```

BPE Algorithm →

Ishan Modi

Word Piece Tokenizer

Almost everything is similar to BPE, there is only one difference

If you recall in the 1st merge in BPE we had following pairs

file bread red.

$$(l, o) - 2$$

$$(o, u) - 2$$

because

higher freq of 'red'

file | bread | red | .

file | bread | .

(i, d) - 1

Here we see there was a tie in the pre frequency of many pairs but we picked the 1st occurrence.

For word Piece, instead of 1st occurrence we use the following formula

$$\text{score} = \frac{\text{count}(\alpha, \beta)}{\text{count}(\alpha) \text{ count}(\beta)}$$

we calculate this score for all the tied frequency (2) & this case I pick a value with maximum score

e.g -

$$\text{so score } (l, o) = \frac{\text{count}(l, o)}{\text{count}(l) \text{ count}(o)}$$

$$= \frac{2}{2 \times 2} = \frac{1}{2}$$

Ishan Modi

$$\text{Score}(c, w) = \frac{2}{2 \times 4} = \frac{1}{4}$$

1st onwards following fl
4 wordlengths total with

Similar for all other same frequency

so, (l, o) & (e, l) have same score $\frac{1}{2}$

and same highest frequency 2. So if score is also same again we pick 1st ie - (l, o)

Sentence Piece Tokenizer

Following is a small problem with BPE & wordpiece, that we want to address using sentence piece.

- lets say $V = \{h, e, l, l, o, he, el, lo, ll, hell\}$

where word 'hello' can be split as follows

h | e | l | l | o since 'he' occurred first so we merge

now we get he | l | l | o since 'el' is not there so we ignore

process next word 'll' has been merged & 'lo' occurred so we merge

he | l | lo } this is the tokenized value

- lets say $V = \{h, e, l, l, o, el, he, lo, ll, hell\}$

here word 'hello' is not merged

h | el | lo } this is the tokenized value

Ishan Modi

It greedily chooses the segments, and doesn't think about the best segmentation.

Till now, we were starting off with characters for building words.

The idea of Sentence piece is to start off with a larger vocabulary & reduce it probabilistically.

Step 1

Construct a reasonably large seed vocabulary using BPE or extended suffix array algorithm

test 92919, based on 2198 tokens, vocabulary size is 21, including

biggest prefix length of four words

Step 2

E-step

Estimate the probability for every token in the given vocabulary using frequency counts in training corpus

lets take as dummy example

Corpus - I love playing volleyball and football. I also love meeting people sharing ideas and learning from their experiences

Vocab from BPE

I	volley	also	their
love	ball	meet	experiences
play	volleyball	people	
ing	and	share	
vol	foot	ideas	
le	football	learn	
y	al	from	

Ishan Modi

Frequency in corpus

I	-	2	volley	-	1	also	-	1	their	-	1
love	-	2	ball	-	2	meet	-	1	experiences	-	1
play	-	1	volleyball	-	1	people	-	1			
ing	-	4	and	-	2	shar	-	1			
vol	-	1	foot	-	1	ideas	-	1			
le	-	3	football	-	1	learn	-	1			
y	-	2	al	-	3	from	-	1			

Total Frequencies $\rightarrow 35$

Probability

$$P(\text{token}) = \frac{\text{frequency of the token in the corpus}}{\sum \text{frequencies of all the tokens in the corpus}}$$

Calculate the above for all tokens in the vocab.
For the sake of sanity we will only do for one

$$P(\text{ing}) = \frac{4}{35}$$

Step 3

Use Viterbi Algorithm to segment the corpus & return optimal segment that maximizes the log likelihood

Pick all words from corpus one by one & start segmenting
we will do it for one word \rightarrow vol le y ball

Segment 1 \rightarrow vol le y ball

$$\begin{aligned} p(\text{volleyball}) &= p(\text{vol}) \cdot p(\text{le}) \cdot p(\text{y}) \cdot p(\text{ball}) \\ \log(p(\text{volleyball})) &= \log(p(\text{vol})) + \log(p(\text{le})) + \log(p(\text{y})) + \log(p(\text{ball})) \\ &= \log\left(\frac{1}{35}\right) + \log\left(\frac{3}{35}\right) + \log\left(\frac{2}{35}\right) + \log\left(\frac{2}{35}\right) \\ &= -5.09 \end{aligned}$$

Segment 2 \rightarrow Volley ball

$$\begin{aligned} p(\text{Volleyball}) &= p(\text{volley}) \cdot p(\text{ball}) \\ \log(p(\text{Volleyball})) &= \log(p(\text{volley})) + \log(p(\text{ball})) \\ &= \log\left(\frac{1}{35}\right) + \log\left(\frac{2}{35}\right) \\ &= -2.71 \end{aligned}$$

Segment 3 \rightarrow Volleyball

$$\begin{aligned} p(\text{Volleyball}) &= p(\text{volleyball}) \\ \log(p(\text{Volleyball})) &= \log(p(\text{volleyball})) \\ &= \log\left(\frac{1}{35}\right) \\ &= -1.54 \end{aligned}$$

We pick and store the third segment that has only one token 'volleyball' because it has the highest log likelihood

Ishan Modi

Similarly we pick up and store ~~and find~~ ~~the best segment~~
 for all words in the corpus.

At the end of these steps we would have optimal segments for all words.

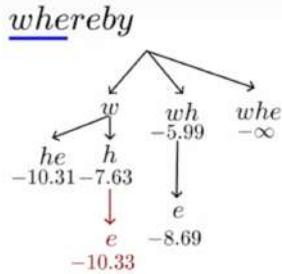
Now these segments become our vocab and we repeat steps 2 & 3. The (only) change in subsequent iterations here is that after step 2, we will remove $n\%$ of tokens that have smallest likelihoods.

We do this until the vocab is reduced to the required size.

There is still a question, how do we get segments for a given word in the corpus?

→ Brute force - calculating for all possible segmentations would be computationally very expensive

→ Use trie & prune as you traverse.



Token	$\log(p(x))$
b	-4.7
e	-2.7
h	-3.34
r	-3.36
w	-4.29
wh	-5.99
er	-5.34
where	-8.21
by	-7.34
he	-6.02
ere	-6.83
here	-7.84
her	-7.38

Note - When it comes to tokenization there is something known as Net fertility rate.

Net fertility rate = tokens in corpus / no. of words

Fertility rate → Measurement of at an average how many tokens each word

gets. Now if a word is broken down into clustering

of tokens then the tokens of each word after tokenization in words are same.

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

∴ Fertility rate = $\frac{\text{Total no. of tokens after tokenization of all words in the corpus}}{\text{Total no. of words in the corpus}}$

Ishan Modi

* Positional Encoding

- Till now, we have only looked at what is called absolute positional encoding (PE). Sinusoidal PE & learnable PE are two examples we have looked so far. What absolute PE does is that it looks at the entire sequence length and gives # tokens positions starting from 0 to end.
- There are two problems with PEs

- Let's say we have a five token sentence that starts at position 0 & ends at position 4.

Now, if we shift this sentence to position 10 to 14 the results were observed to change. The accuracy went down.

So there is a dependence of absolute position. This can be removed if we pick a word and give positions relative to that word eg - 0, 1, 2, -1, -2 with respect to a word.

- Second problem is the sequence length, if we have sequence length 512 & thus we learn 512 Positional embedding,

we would not be able to deal with 1024 embedding sequence length at inference. Predicting larger sequence won't be possible

Also if during the training all samples provided were less than 300 sequence length, model won't even generalize to sequence length above 300 at inference.

These are the reasons for us to use Relative Positional Encodings

→ Relative Positional Encodings

As we discussed earlier we take each word as a center point give it position zero & give the next words positions from 1 and previous words positions from -1.....

	somewhere	something	incredible	is	waiting	to	be	known
somewhere	0	1	2	3	4	5	6	7
something	-1	0	1	2	3	4	5	6
incredible	-2	-1	0	1	2	3	4	5
is	-3	-2	-1	0	1	2	3	4
waiting	-4	-3	-2	-1	0	1	2	3
to	-5	-4	-3	-2	-1	0	1	2
be	-6	-5	-4	-3	-2	-1	0	1
known	-7	-6	-5	-4	-3	-2	-1	0

Now, earlier we just had 1 embedding for one token in absolute PE, which we added to token embedding

Here, you have T (sequence length) PE for one token that you add to token embedding.

Ishan Modi

• Speed: In Absolute PE we have to process all tokens. If we have news about tokens, we need to message box and not do

For Absolute PE we add 1 PE for 1 token embedding & we do it for T (sequence length) tokens. So speed is $O(T)$ because all are $O(T)$ processing.

For Relative PE up till now we add T PE to 1 token embeddings & we do it for T tokens. So speed is $O(T^2)$.

• Memory: This is also with locality of 1. For a sequence of size T , we have to store T PE in memory.

For Absolute PE we have T PE in memory.

For Relative PE as seen in above figure -7 to 7 for sequence length 8, ie we have ~~relative~~ PE $(2T-1)$

• Formula

For $T = 8$

$$h_3 = n_3 + (p_{-3} + p_{-2} + \dots + p_4)$$

$$h_i = n_i + \sum_{j=0}^{T-1} p_{j-i}$$

Now, in the next section we try to simplify the above

method by using accumulators. In last step, we can see that

matrix was not in (square matrix) so sum was not

partitioned into a block matrix, so

Ishan Modi

Let's look at how these PE look inside the transformer block computation.

attention matrix before softmax looks like following

$$e_{ij} = q_i k_j^T$$

recall that,

$$q_i = h_i w_\alpha = (n_i + p_i) w_\alpha$$

$$k_j = h_j w_k = (n_j + p_j) w_k$$

Thus, attention score can be written as

$$e_{ij} = (n_i + p_i) w_\alpha w_k^T (n_j + p_j)^T$$

$$\text{expanding } e_{ij} = (n_i w_\alpha + p_i w_\alpha) (w_k^T n_j^T + w_k^T p_j^T)$$

$$= [n_i w_\alpha w_k^T n_j^T] + [n_i w_\alpha w_k^T p_j^T] \\ + [p_i w_\alpha w_k^T n_j^T] + [p_i w_\alpha w_k^T p_j^T]$$

Interaction between key word embedding and query word embedding

comes from PE

correlation between word & position embedding

comes from PE

Research shows that the middle two terms don't contribute a lot to the overall attention score.

Ishan Modi

What if we use the below formula during attention computation rather than adding P_B at the beginning?

$$e_{ij} = \alpha_i w_k (z_j w_k + p_{ij})$$

p_{ij}^k is also known as positional bias.

Note - The d dimension that you see here is the dimension of head in such heads + concatenate to form context.

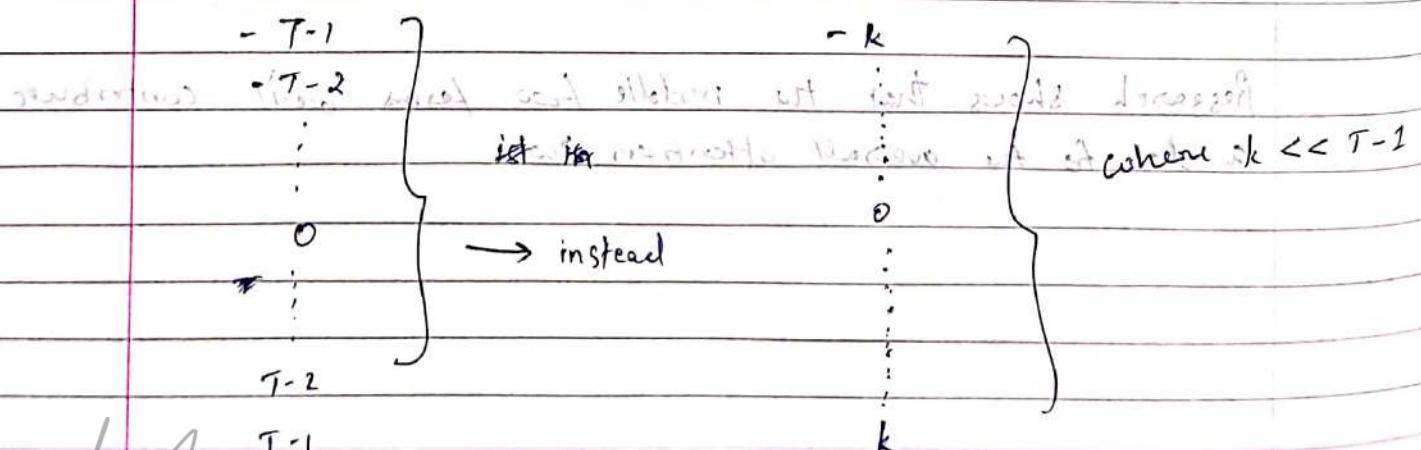
Here in the above we have removed the addition part

$$(e_{ij}^k + p_{ij}^k) = e_{ij}^k + w_k w_k^T (q_i + v_j)$$

($e_{ij}^k + p_{ij}^k$ is e_{ij}) (we are just picking one element from the previous summation formula)

But, this still doesn't generalize to arbitrary length, what if we have larger sequence length at test time?

We do a simple trick, currently we have $2T-1$ PE



Ishan Modi

Thus even if we have more sequence length we clip it to k , so we will only have $2k+1$ PE

following is formula for clipping

$$P_{ij}^k = w_{\text{clip}}^k (j-i, k)$$

$$\text{clip}(j-i, k) = \max(-k, \min(j-i, k))$$

So the formula for attention becomes

$$\begin{aligned} e_{ij} &= \pi_i w_\alpha (z_j w_k + P_{ij}^k)^T \\ &= \pi_i w_\alpha w_k^T z_j^T + \pi_i w_\alpha (P_{ij}^k)^T \end{aligned}$$

Then, $(\pi_i w_\alpha)$ is shared across heads but not layers

$$\alpha_{ij} = \text{softmax}(e_{ij})$$

Now, add relative PE for value matrix

$$z_i = \sum_j \alpha_{ij} (\pi_j w_v + P_{ij}^v)$$

Note - P_{ij}^k , P_{ij}^v are shared across heads but not layers

at every layer we add separate P_{ij}^k & P_{ij}^v

RPE Variations

recall the decomposition

$$e_{ij} = x_i W_Q W_K^T x_j^T + x_i W_Q W_K^T p_j^T + p_i W_Q W_K^T x_j^T + p_i W_Q W_K^T p_j^T$$

Transformer-XL (recursively extend context)

$$e_{ij} = x_i W_Q W_K^T x_j^T + x_i W_Q R_{j-i}^T + u W_K^T x_j^T + v W_K^T R_{j-i}^T$$

Where,

R_{i-j}^T is the relative distance between the position i and j computed from sinusoidal function (non-learnable)

u, v are learnable vectors.

T5-bias

All that is required is adding a constant, why not just do that?

$$e_{ij} = x_i W_Q W_K^T x_j^T + r_{j-i}$$

Where,

$r_{i-j} \in \mathbb{R}$ are learnable scalars, shared across layers. Maximum relative distance is clipped to

It also greatly reduces number of learnable parameters when the model size is scaled

→ Rotary Positional Embeddings (ROPE)

Following are the relative PE we have seen so far

Shaw $e_{ij} = x_i W_Q W_K^T x_j^T + x_i W_Q (p_{ij}^K)^T$

T-XL $e_{ij} = x_i W_Q W_K^T x_j^T + x_i W_Q R_{j-i}^T + u W_K^T x_j^T + v W_K^T R_{j-i}^T$

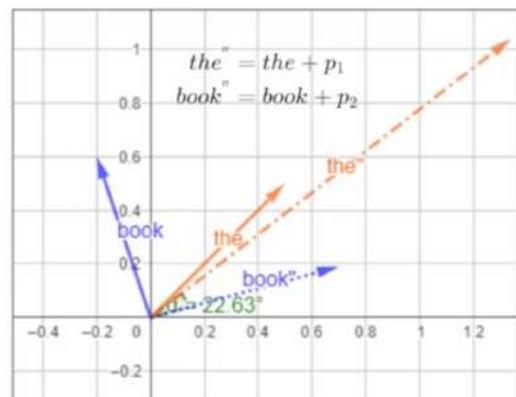
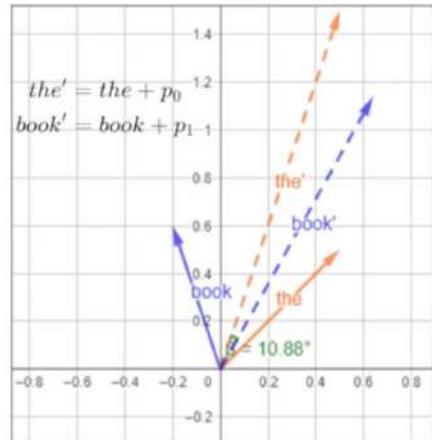
T5 $e_{ij} = x_i W_Q W_K^T x_j^T + r_{j-i}$

There is still a small problem with relative PE,
lets see it in detail

Ishan Modi

the book ...

read the book ...



As seen in the above images the relative distance between the words "the" & "book" is same in both sentences i.e. they are next to one another.

But if we look at the vectors they are completely changed originally both had some vectors

Now, $q_k^T \Rightarrow q \cdot k$ is the operation happening here

$q \cdot k = \cos(\theta)$, can we keep somehow keep the same for these words is the idea behind rotary PE

Earlier, we had the following eqⁿ for relative PE

$$e_{ij} = \left(\frac{(z_i w_{0j})_{\text{rel}} - (z_i w_{0i})_{\text{rel}}}{\sqrt{(z_i w_{0j})_{\text{rel}}^2 + (z_i w_{0i})_{\text{rel}}^2}} \right)^T$$

$$= z_i w_{0j} w_{0i}^T z_i^T \quad (\text{term } z_i w_{0i} \text{ is ignored})$$

$$\left. \begin{aligned} & (z_m w_{0j})_{\text{rel}} - (z_n w_{0i})_{\text{rel}} \\ & (z_m w_{0j})_{\text{rel}} \cdot (z_n w_{0i})_{\text{rel}} \end{aligned} \right\} \text{we ignore this term}$$

$$z_m w_{0j} e^{jm\theta} \cdot (z_n w_{0i} e^{jn\theta})^T$$

(here we have replaced i, j with m, n because j also represents imaginary number)

Ishan Modi

$$\therefore e_{mn} = \underbrace{n_m w_\alpha e^{j m \theta}}_{\text{Hermitian}} \underbrace{(n_n w_k e^{j n \theta})^T}_{\text{Transpose}}$$

now instead of transpose we use Hermitian transpose

$$\mathbf{A} = \begin{bmatrix} 1 & -2-i & 5 \\ 1+i & i & 4-2i \end{bmatrix}$$

We first transpose the matrix:

$$\mathbf{A}^T = \begin{bmatrix} 1 & 1+i \\ -2-i & i \\ 5 & 4-2i \end{bmatrix}$$

Then we conjugate every entry of the matrix:

$$\mathbf{A}^H = \begin{bmatrix} 1 & 1-i \\ -2+i & -i \\ 5 & 4+2i \end{bmatrix}$$

associated binomial entries will remain same
 conjugate and multiply by $e^{-j\theta}$
 nothing else will change

Example: $e_{mn} = n_m w_\alpha e^{j m \theta} (n_n w_k e^{j n \theta})^H$

and conjugate $n_m w_\alpha (n_n w_k)^T e^{-j(m-n)\theta}$

This is all about formulation, how to solve
 in real life situations

$$e^{j n \theta} = \cos(n\theta) + j \sin(n\theta)$$

$$e^{j n \theta} = \begin{bmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{bmatrix}$$

replacing n by $(m-n)$

$$e^{j(m-n)\theta} \therefore \begin{bmatrix} \cos((m-n)\theta) & -\sin((m-n)\theta) \\ \sin((m-n)\theta) & \cos((m-n)\theta) \end{bmatrix}$$

in this is border guard is open

representing boundary condition is forward

Ishan Modi



Thus, unit impulse response in bivariate case is given by

$$e_{mn} = n_m w_\alpha + \begin{bmatrix} \cos(m\theta) & \sin(m\theta) \\ \sin(m\theta) & -\cos(m\theta) \end{bmatrix} \begin{bmatrix} \cos(n\theta) & \sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{bmatrix}$$

(unit impulse response) \times (bilateral $w_k^T z_n^T$)

$$= n_m w_\alpha \begin{bmatrix} \cos(m-n)\theta & -\sin(m-n)\theta \\ \sin(m-n)\theta & \cos(m-n)\theta \end{bmatrix} w_k^T z_n^T$$

(unit impulse response) \downarrow 2×2 \downarrow 2×2 \downarrow 2×2

Now if dimension increases beyond 2, it looks like following

$$e_{mn} = \begin{bmatrix} n_m w_\alpha R_m \end{bmatrix} \cdot \begin{bmatrix} n_n w_k R_n \end{bmatrix}^T$$

(unit impulse response) \downarrow $1 \times d$ \downarrow $d \times d$ \downarrow $d \times d$ \downarrow $d \times d$ \downarrow $d \times d$

$$((\cos(\theta_1), \dots, \cos(\theta_d)) \text{ and } (\sin(\theta_1), \dots, \sin(\theta_d)))$$

$$R_m = \begin{bmatrix} \cos(m\theta_1) & -\sin(m\theta_1) & 0 & 0 & \dots & 0 & 0 \\ \sin(m\theta_1) & \cos(m\theta_1) & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos(m\theta_2) & -\sin(m\theta_2) & \dots & 0 & 0 \\ 0 & 0 & \sin(m\theta_2) & \cos(m\theta_2) & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \cos(m\theta_{d/2}) & -\sin(m\theta_{d/2}) \\ 0 & 0 & 0 & 0 & \dots & \sin(m\theta_{d/2}) & \cos(m\theta_{d/2}) \end{bmatrix}$$

$$\left\{ \theta_i = 10000^{-\frac{(i-1)}{d}}, i \in \{1, 2, \dots, d\} \right\}$$

Ishan Modi

Note - RoPE is applied in all layers individually but unlike relative RE it doesn't calculate or add PE to the value matrix, it is just calculated for key/Query transformation.

\rightarrow Attention With Linear Bias (AlBi)

Following is the PE we saw with TS bias

$$c_{ij} = \text{softmax}(q_i^T w_k + w_q^T k_j + r_{j-i})$$

where r_{j-i} , G , R are learnable scalars

$$\text{AlBi makes these scalar fixed instead of learnable}$$

- Formula

$$\text{softmax}(q_i^T K + m \cdot (-1, -2, \dots, -1, 0))$$

- Understanding

$$\begin{matrix} q_1 \cdot k_1 & \\ q_2 \cdot k_1 & q_2 \cdot k_2 \\ q_3 \cdot k_1 & q_3 \cdot k_2 & q_3 \cdot k_3 \\ q_4 \cdot k_1 & q_4 \cdot k_2 & q_4 \cdot k_3 & q_4 \cdot k_4 \\ q_5 \cdot k_1 & q_5 \cdot k_2 & q_5 \cdot k_3 & q_5 \cdot k_4 & q_5 \cdot k_5 \end{matrix} + \begin{matrix} 0 \\ -1 & 0 \\ -2 & -1 & 0 \\ -3 & -2 & -1 & 0 \\ -4 & -3 & -2 & -1 & 0 \end{matrix} \cdot m$$

Here,

m is scalar that follows geometric progression

for $n=8$ heads, $m_i = \frac{1}{2^i}$, $i = 1, \dots, n$

and since this is relative PE, this is applied at all layers individually.

Ishan Modi

* Types of Attention Mechanisms

→ Time & Space complexity of Attention Mechanism

As seen in earlier parts, following is the formulation for attention.

$$\text{Attn} = \text{softmax}(\underline{QK^T}) \cdot V$$

Time Complexity

$$QK^T$$

$$Q =$$

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}_T$$

$$K^T =$$

$$\begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}_d$$

Here T is number of tokens & d is the dimension (size of embedding) each token

Let's pick a row of Q matrix & multiply by all columns of K^T

Here total $(T \cdot d)$ multiplications happen & $T \cdot (d-1)$ additions happen

∴ Total complexity for one row = $O(Td)$

So for T rows of Q matrix

Overall Complexity = $O(T^2d)$

Ishan Modi

- divide by $\sqrt{d_k}$ and then softmax for output

The entire operation of normalization has a complexity of $O(T^2)$

- Now, Attn. $\times V \rightarrow \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right)$

$$A = \begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix}^T \quad V = \begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix}^T$$

The overall complexity is similar to the first case
 $O(T^2d)$

~~Total complexity~~ $\therefore \text{Total} = O(T^2d) + O(T^2) + O(T^2d)$

parameters $\Rightarrow O(T^2d)$. when d is very large, it is slow

~~Space complexity~~ \rightarrow required working memory $((b, T))$ for b batch size
required storage $((b, d, T))$

lets (say), b = one row of pfivologram batch \therefore

- B = Batch Size
- V = Number of tokens in Vocab
- T = Sequence length

Ishan Modi

- d_{ff} = dimension of feed forward layer
- d_{model} = dimension of input embedding
- n_h = no. of heads
- N = No. of model parameters.

Now, the two parts of

Multihed Attention has parts

- $(mha) = (Q_{TxL}^T + K_{TxL}^T + V_{TxL}) n_h$
- $= 3T(C_{TxL}) n_h$

~~For multihed attention~~
~~for decoding all to~~ $= 3T \times n_h \times d$

~~For decoding all to~~ $= 3T \times d_{model}$ ~~the~~ $n_h \times d = d_{model}$

For batch size B

~~Handling batch size stuff~~

~~Handling batch size stuff~~ $\Rightarrow 3BTd_{model}$

~~Handled by an addition~~

- $(Attn)_{all} = \text{softmax}(QK^T)$
- $\text{Time} = Bn_h T^2$

~~Handled by an addition~~

~~softmax(QK^T)~~

~~Total = $3BTd_{model} + Bn_h T^2$~~

~~For each of batch~~

~~Thus it can be said that memory cost scales up linearly with batch size & quadratically with sequence length~~

~~(LSTMs)~~

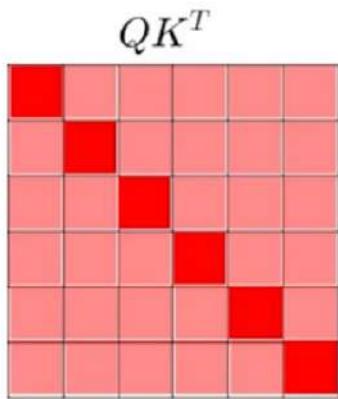
~~So generated attention head takes one multiplication with each~~

~~to process each of sequence element to predict next~~

Ishan Modi

→ Attention Mechanisms

(1) Full Attention



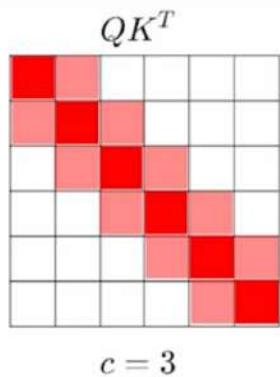
Here every word attends to itself and all the other words in the context.

This time complexity is

$$O(N^2 \cdot D)$$

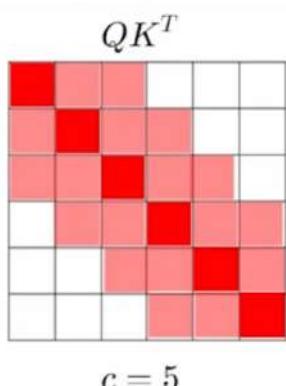
✓ → hidden dimension
no. of tokens D of the embedding in the sequence

(2) Strided Local Attention



Here each word attends to itself & 2 adjacent neighbours, one before & one after, thus $c=3$.

The complexity becomes linear $O(3Nd)$



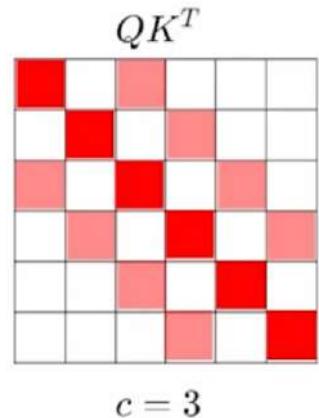
Similar to above but this is $c=5$ in this case

The complexity is still linear $O(5Nd)$

Note - These attentions are called local attentions because we are looking at / paying attention to only few words in the vicinity.

Shank Modis

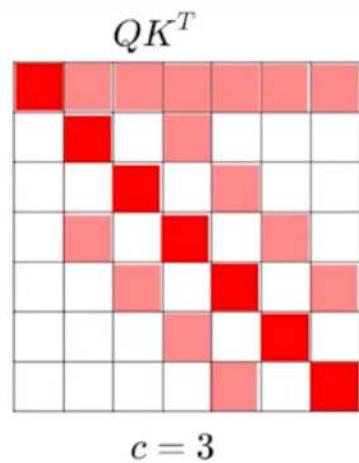
(3) Dilated Attention (hybrid local attention with dilations to capture broader context)



In this case we still look at adjacent neighbours, but pay attention to alternate words. Here $c=3$

The complexity is linear $O(3T_d)$

(4) Global + Local Attention (Is $O((T+T_d) \cdot T_d)$)



Here we can use any of the above local attention / dilated attention idea

But for some tokens we would also pay attention to all the other words

The complexity in this case will be linear

$$O((3T_f + T_d)(d_f + T_d))$$

It follows basic global and local attention

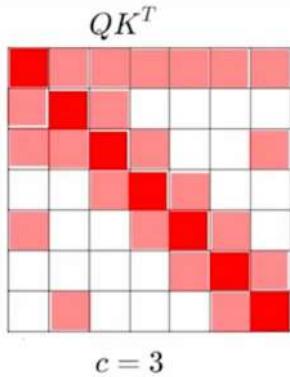
This can also be written as

$$O(3(T_f - 1)d_f + T_d)$$

Local attention component

Global attention component

(5) Global + Random + Local Attention

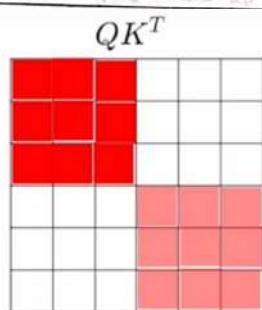


Similar to the previous, but for each word you pay attention to some random words, you may also not select any word, ie random words pick are between (0-T)

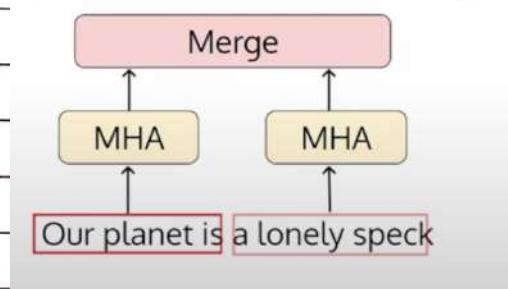
Complexity for the above is

$$O(C(3T + T + T)d)$$

(6) Local Block Attention



Computational complexity $O(\frac{T^2}{2})$



Here as we can see, we have divided the context length into 2 blocks

Thus, $n=2$ in this case.

In the real world if $n=2, 3, 4$

(4 context length) does 512

the block sizes would be

256, 2128, 64 respectively

$$T \rightarrow (2, 4, 8) \rightarrow 0$$

Let's understand in a bit more detail. Let's say we have the following matrices

$$\text{Q} = \begin{bmatrix} \text{Q}_0 \\ \text{Q}_1 \\ \text{Q}_2 \end{bmatrix} \quad \begin{array}{l} \text{columns are rows} \\ \text{of input matrix } \text{E} \in \mathbb{R}^{T \times d_k} \\ \text{rows are columns} \\ \text{of hidden layer} \end{array}$$

$$K^T = \begin{bmatrix} K_0^T & K_1^T & K_2^T \end{bmatrix} \quad \begin{array}{l} \text{columns are rows} \\ \text{of input matrix } \text{E} \in \mathbb{R}^{d_k \times T} \\ \text{rows are columns} \\ \text{of hidden layer} \end{array}$$

As seen earlier we divided both matrix Q & K^T into n blocks (where $n=3$), thus the attention matrix that will be formed will have n^2 , ie 9 blocks

$\text{Q}_0 K_0^T$	$\text{Q}_0 K_1^T$	$\text{Q}_0 K_2^T$
$\text{Q}_1 K_0^T$	$\text{Q}_1 K_1^T$	$\text{Q}_1 K_2^T$
$\text{Q}_2 K_0^T$	$\text{Q}_2 K_1^T$	$\text{Q}_2 K_2^T$

Note - We will use

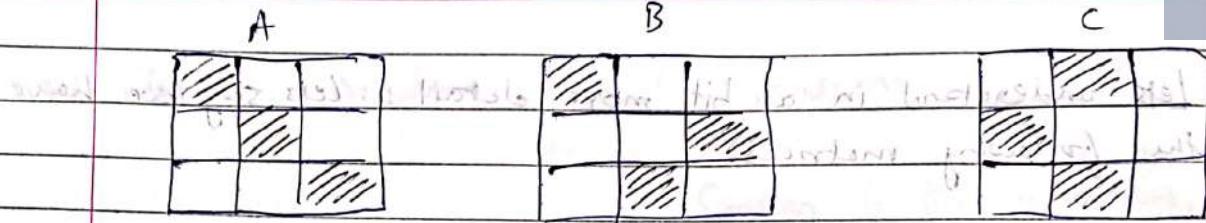
similar blocks for V matrix as K

$\text{Q}_i K_0^T$ will be multiplied by V_0 & same for others as well

Now, for each Q_i we can pick any K_j and this results into different combinations of the above metric

Now, for each Q_i we can pick any K_j and this results into different combinations of the above metric

Ishan Modi



Above are some examples, but there can be $n!$ such permutations.

Let's formulate

$$\pi = \text{perm}(0, 1, \dots, n-1)$$

Let k^{th} element of π be $\pi(k)$ & masking matrix be M of size $T \times T$

$$M_{ij} = \begin{cases} 1 & \text{if } \pi(\lfloor \frac{i-1}{T} \rfloor) = \lfloor \frac{j-1}{T} \rfloor \\ 0 & \text{otherwise} \end{cases}$$

Now we have $M = \text{softmax}(\epsilon Q K^T + \Omega, M)$

hadamard product

Now let's say

A has permutation $\pi = (0, 1, 2)$

B has permutation $\pi = (0, 2, 1)$

C has permutation $\pi = (1, 0, 2)$

We will be creating mask matrix for all the above cases using the formula. The mask matrix obtained will have 1 in place of highlighted regions & 0 otherwise.

Ishan Modi

Note - In the figure of A, B, C the highlighted regions are actually blocks and in practice each block has multiple rows & columns. The masked matrix is not formed at a block level but at the row/column level.

Note - In practice ~~these~~ different such permutations are used in different attention heads to obtain maximum coverage.

Let's discuss the complexity for block local attention

$$= O\left(\frac{T}{n} \times \frac{T}{n} \times n \times d\right)$$

~~= $O\left(\frac{T^2 d}{n}\right)$ get last words first, decay w.r.t. position number need a sort of ordering~~

(7) Linear Attention (Low Rank Approximation)

As we know, following is the formula for full attention

$$A = \text{softmax} \left(\frac{\alpha K^T}{\sqrt{d_k}} \right)$$

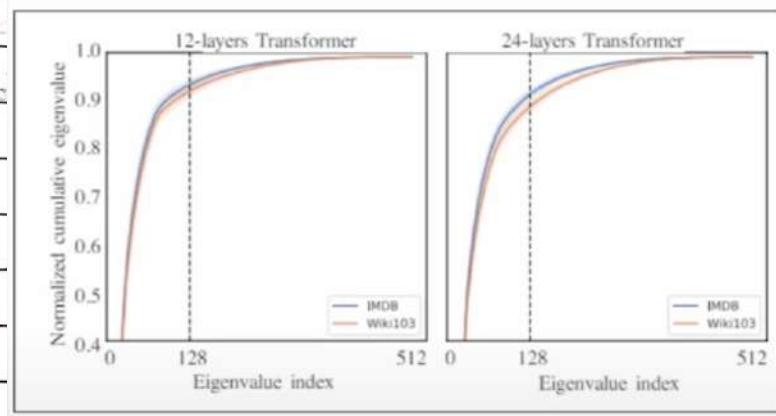
$$Z = A \cdot V$$

Let's ignore Z for now. $A \in \mathbb{R}^{T \times T}$ where T is the sequence length

Let's say $T = 512$, thus we have A which is 512×512 dimensional matrix.

Ishan Modi

We do SVD on this 512×512 matrix, this means we get 512 eigenvalues in diagonal and some null values b/c zero eigenvalue won't add up. So we normalize these values in $0-1$ scale such that all values add up to 1 & plotting cumulative eigen values.



The graph clearly shows that top 128 eigen values contribute 90%, thus a lower dimensional / ranked matrix can represent this 512×512 matrix.

$$Q \in \mathbb{R}^{T \times d} \quad K \in \mathbb{R}^{T \times d} \quad V \in \mathbb{R}^{T \times d}$$

Introduce two learnable linear projection matrices

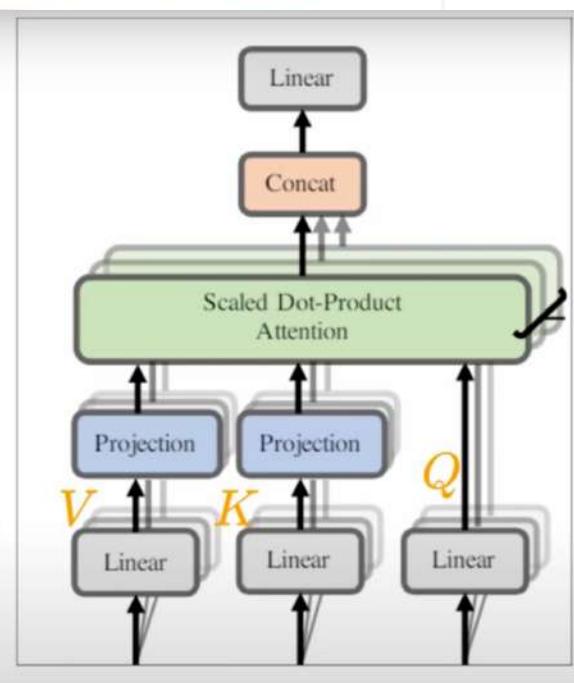
$$E, F \in \mathbb{R}^{k \times T}$$

$$A = \text{softmax}\left(\frac{Q(EK)^T}{\sqrt{d}}\right)$$

$$A \in \mathbb{R}^{T \times k}$$

$$O = AFV$$

The computational complexity is $O(kTd)$, where k is the rank and can be set according to the error ϵ



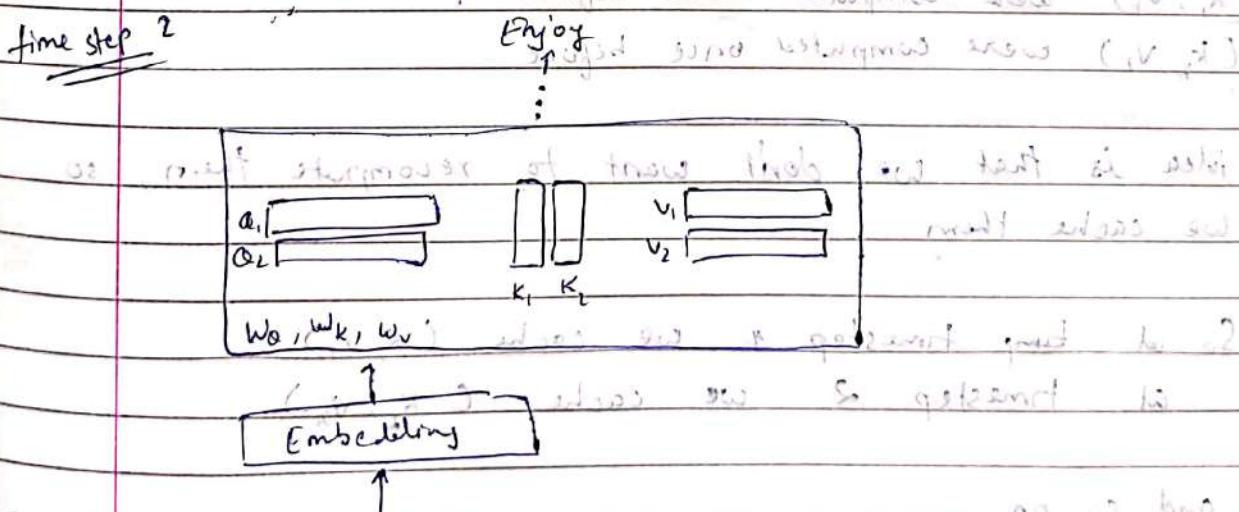
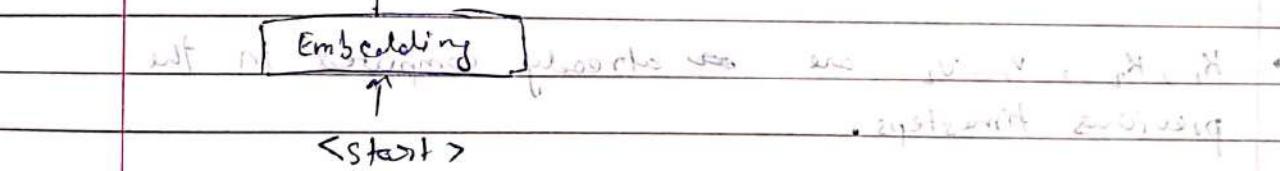
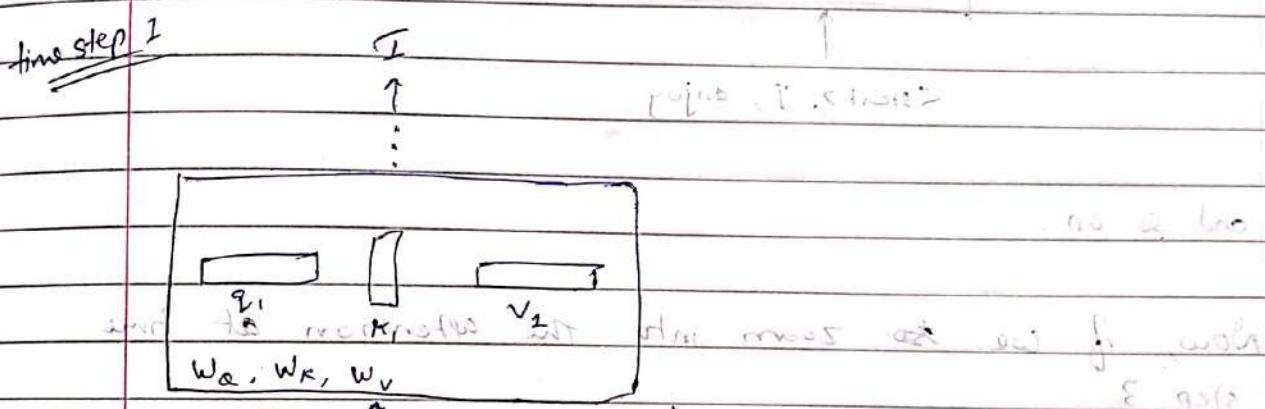
Ishan Modi

* Inference Mechanisms

→ KV Cache

The concept of KV cache is mainly used in transformer decoder.

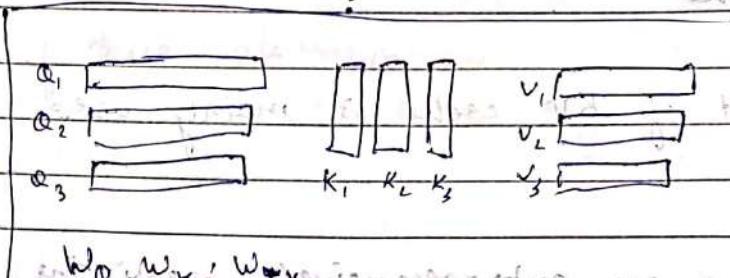
Let's say we are auto-regressively predicting tokens at inference



timestep 3

dancing

to end of



Embedding

<start>, I, enjoy

and so on.

Now, if we zoom into the attention at time step 3

- K_1, K_2, v_1, v_2 are already computed in the previous timesteps.

(K_1, v_1) were computed twice before

(k_2, v_2) were computed once before

idea is that we don't want to recompute them, so we cache them

So at time timestep 1 we cache (K_1, v_1)

at timestep 2 we cache (K_2, v_2)

and so on ...

This is called KV Cache

Ishan Modi

- Can we also do something for q_1, q_2, q_3 because they also keep consuming memory in addition to KV cache at time step B^1 we have q_1 , at time step 2^1 we have q_1, q_2 & so on at time step 3^1 we have q_1, q_2, q_3 & so on

The answer is only use q_i at i^{th} time step other q 's are not needed

$$\begin{matrix} a_1 \\ a_2 \\ a_3 \end{matrix} \left[\begin{matrix} \quad \\ \quad \\ \quad \end{matrix} \right] \left[\begin{matrix} K_1 & K_2 & K_3 \end{matrix} \right] \begin{matrix} v_1 \\ v_2 \\ v_3 \end{matrix} \left[\begin{matrix} \quad \\ \quad \\ \quad \end{matrix} \right] = \begin{matrix} a_1 \\ a_2 \\ a_3 \end{matrix} \left[\begin{matrix} \quad \\ \quad \\ \quad \end{matrix} \right]$$

$\underbrace{\alpha}_{\text{softmax}}$ K^T v \underbrace{A}

If we observe closely a_3 is the output of the timestep 3 it can also be called q_4 (dancing) in this case

So for computing a_3 , we only need $K_1, K_2, K_3, v_1, v_2, v_3$

$$q_3 \left[\begin{matrix} \quad \end{matrix} \right] \left[\begin{matrix} K_1 & K_2 & K_3 \end{matrix} \right] \begin{matrix} v_1 \\ v_2 \\ v_3 \end{matrix} \left[\begin{matrix} \quad \\ \quad \\ \quad \end{matrix} \right] = a_3 \left[\begin{matrix} \quad \\ \quad \\ \quad \end{matrix} \right]$$

$\underbrace{\alpha}_{\text{softmax}}$ K v \underbrace{A}

Ishan Modi

So KV cache reduces computation complexity from $O(n^2)$ to $O(n)$ but it increases memory requirement

(Let's compute memory requirement for KV cache)

KV cache per token in bytes =

~~In each layer $2 \times \text{num_layers} \times (\text{num_heads} \times \text{dim_head})$~~
 ~~\times precision in bytes~~

GPT-3

num_layers = 96

num_heads = 96

dim_head = 128

precision_in_bytes = 4 = 32 bit precision

$$\begin{aligned}\text{KV cache per token} &= 2 \times 96 \times (96 \times 128) \times 4 \\ &= 9.4 \text{ MB}\end{aligned}$$

~~if context length = 2000~~

$$\text{KV cache} = 2000 \times 9.4 \text{ MB} = 19.2 \text{ GB}$$

This is huge, so we need some modifications to deal with this

padding, per token closest address in memory and swap between memory

No padding, swap between memory and swap between memory

swap between memory and swap between memory

swap between memory and swap between memory

Ishan Modi

→ Multi-Query Attention (MQA)

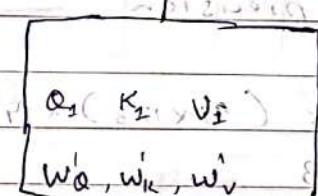
KV cache per token in Bytes =

$$\text{bytes per token} \times \text{num-layers} \times (\text{num-heads} \times \text{dim-head}) \\ \times \text{precision in bytes}$$

The idea is to share KV across all heads and thus reduce num-heads to 1

Let's see with example

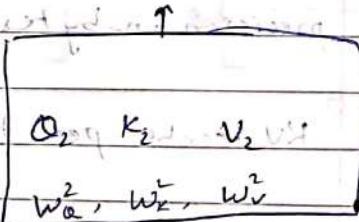
head 1



$38 = \text{embedding size}$

head 2

$38 = \text{embedding size}$



So here for KV weights, w_k^1, w_v^1 & w_k^2, w_v^2 are different.

are different in different heads. Thus they produce different K_1, V_1 & K_2, V_2 .

If we share w_k^1 & w_k^2 , w_v^1 & w_v^2 for all

heads we can compute K, V only once and use it with Q_i 's in all heads.

Ishan Modi

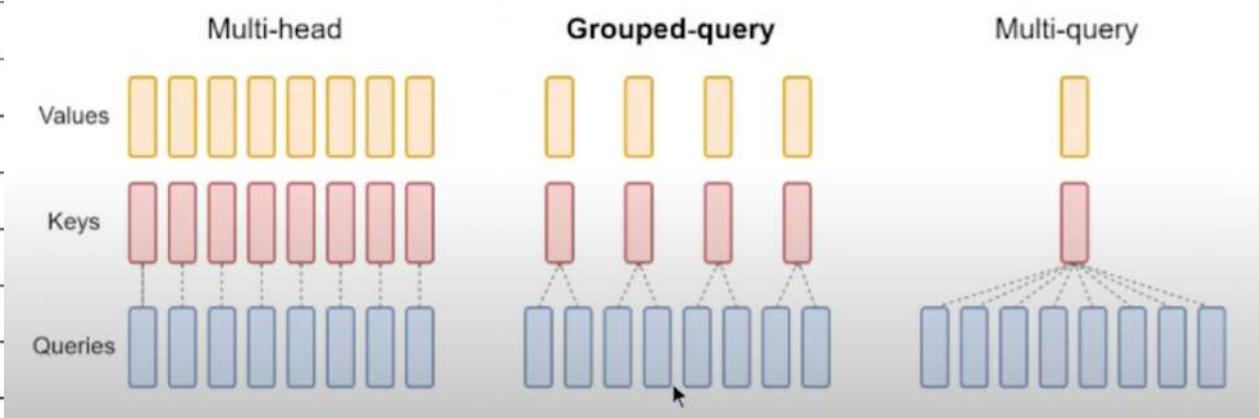
We are reducing parameters thus accuracy will degrade & we will also have to retain more data from training.

→ Grouped Query Attention (GQA)

This is a middle ground between multi-head attention & multi-query attention.

The idea here is to share KV weights & thus KV values for "n" number of queries.

Thus we could have groups of queries working with different KV values.



* Decoding Strategies (these strategies are used during inference)

Deterministic

- Exhaustive Search
- Greedy Search
- Beam Search

Stochastic

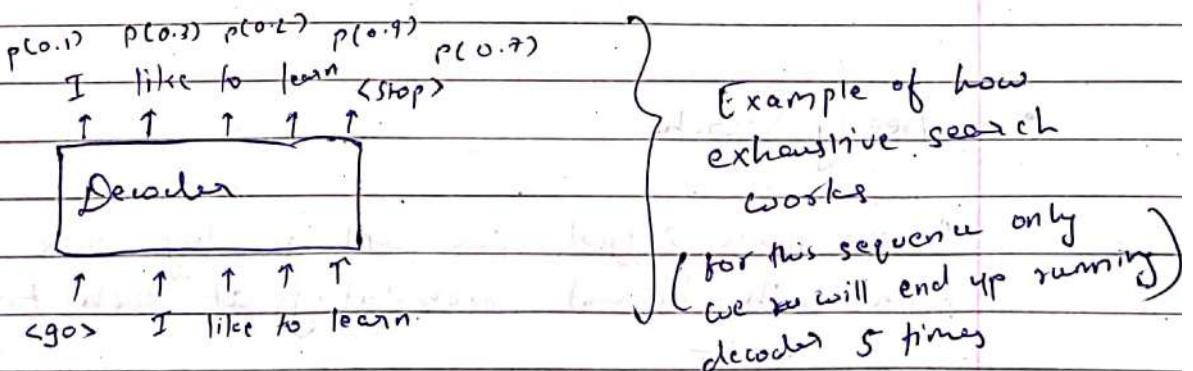
→ Deterministic Decoding (same output everytime)

• Exhaustive Search

Let's say we want a output from the decoder that contains 5 words. Our vocabulary size is 40k words. Then the total number of unique sequences that can be formed

$$|V|^5 = (\text{vocabulary size})^5 \approx 40,000^5$$

The number of permutations is huge & thus exhaustive search is practically not possible.



Now, we calculate the probability of the above sequence

$$P = p(I) p(\text{like}|I) p(\text{to}|\text{like}, I) \dots .$$

$$= 0.1 \times 0.3 \times 0.2 \dots \times 0.9 \times 0.7$$

Similarly we calculate the probabilities of all possible sequences and pick the sequence with highest probability

If we talk about time complexity you will end up running the transformer decoder $|V|^5$ times for this case

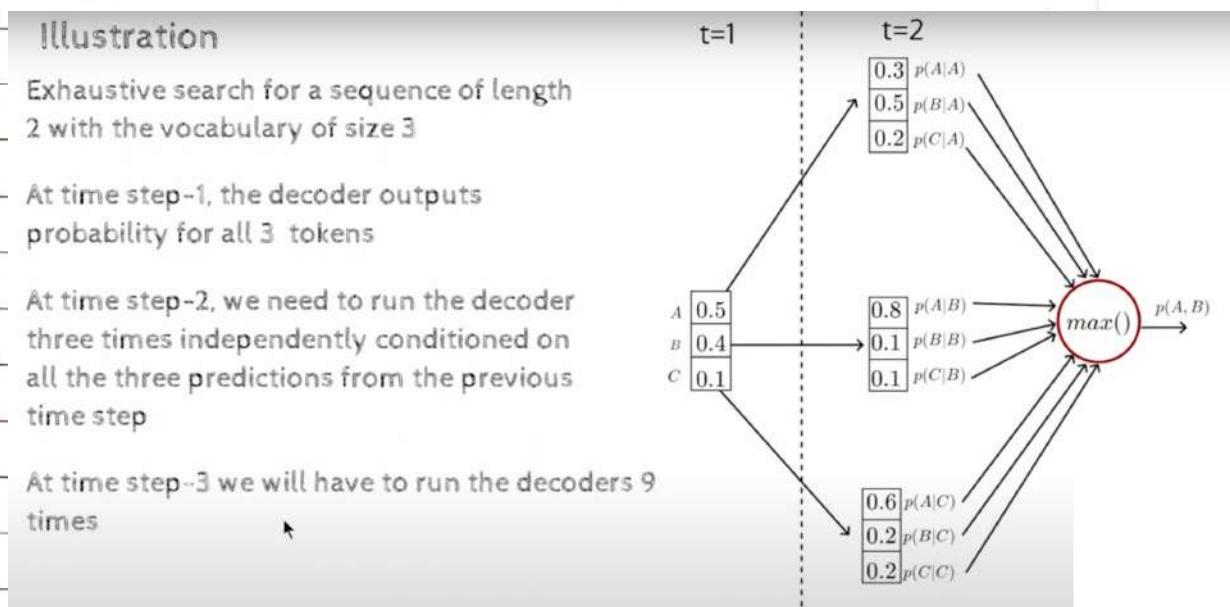
Illustration

Exhaustive search for a sequence of length 2 with the vocabulary of size 3

At time step-1, the decoder outputs probability for all 3 tokens

At time step-2, we need to run the decoder three times independently conditioned on all the three predictions from the previous time step

At time step-3 we will have to run the decoders 9 times



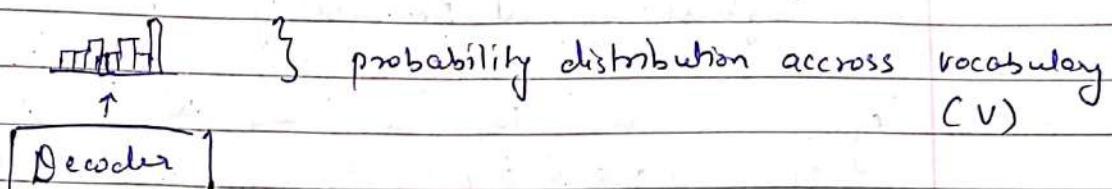
• Greedy Search

As the name suggests, we will greedily pick the word with highest probability at each time step

The time step here means every time we run the decoder to get the next word

Taking the same example of predicting a sequence of 5 words

time step 1



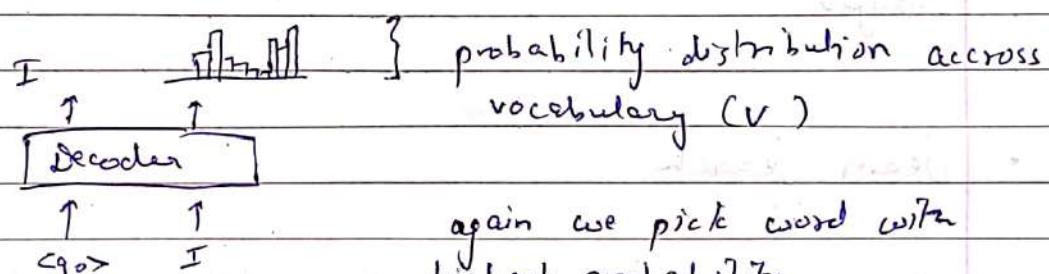
$\langle \text{go} \rangle$

since it is greedy search out of all probabilities in the distribution we pick word with highest probability

lets say at timestamp 1 "I" had highest probability

$$p(I) = 0.9$$

time step 2



$\langle \text{go} \rangle \quad I$

again we pick word with highest probability.

lets say at timestep 2 "Like" had highest probability

$$p(\text{Like} | I) = 0.8$$

Now, we carry on same thing across 5 time steps and pick highest probability words at each step given previous context. Thus you obtain the sequence

So, if we talk about time complexity, we will end up running the decoder 5 times, which is equal to the words in the sequence

number of

Ishan Modi

Note - Unlike exhaustive search, greedy search doesn't guarantee that the sequence that is selected has the highest overall probability

$$\text{eg - } \begin{matrix} w_1 & w_2 & w_3 & w_4 & w_5 \end{matrix} \\ P(S_1) \Rightarrow 0.8 \times 0.4 \times 0.25 \times 0.3 \times 0.22 \quad (\text{sequence 1 probs}) \\ w_1' \quad w_2' \quad w_3' \quad w_4' \quad w_5' \\ P(S_2) \Rightarrow 0.19 \times 0.65 \times 0.82 \times 0.65 \times 0.76 \quad (\text{sequence 2 probs})$$

$$\text{Here, } P(S_2) > P(S_1)$$

What we have actually done, at the first time step we picked the second wood with second highest probability. Thus all the subsequent distributions changed

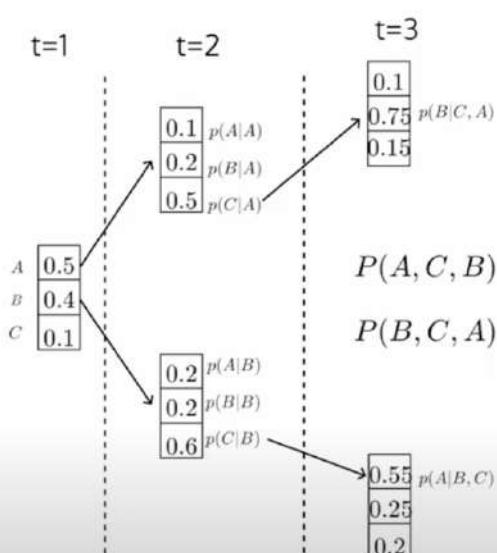
- Beam Search

Beam search is somewhere in between exhaustive & greedy search

- Exhaustive search looks at all probabilities at all time steps
- Greedy search looks at highest probabilities at all time steps autoregressively
- Beam search looks at "b" top probabilities at all time steps where $b \Rightarrow$ beam size.

Let's understand it with the following example.

here vocab size $V = 3$ & beam size $b = 2$



$$P(A, C, B) = 0.18$$

$$P(B, C, A) = 0.13$$

As seen in the figure, at first timestep decoder is called only once & for all other timesteps it is called "b" (beam size) ~~a~~ number of times.

time step 1 \rightarrow we pick the top 2 words with highest probabilities

time step 2

$$p(A) p(A|A) = 0.5 \times 0.1 = 0.05$$

$$p(A) p(A|B/A) = 0.5 \times 0.2 = 0.01$$

$$p(A) p(C|A) = 0.5 \times 0.5 = 0.25$$

$$p(B) p(A|B) = 0.4 \times 0.2 = 0.08$$

$$p(B) p(B|B) = 0.4 \times 0.2 = 0.08$$

$$p(B) p(C|B) = 0.4 \times 0.6 = 0.24$$

we pick the top 2 words with highest probabilities

We carry on same till desired words in the sequence & then out of b sequences we pick sequence with highest probability

Ishan Modi

Note - Beam search with $b=1$ is equivalent to greedy search

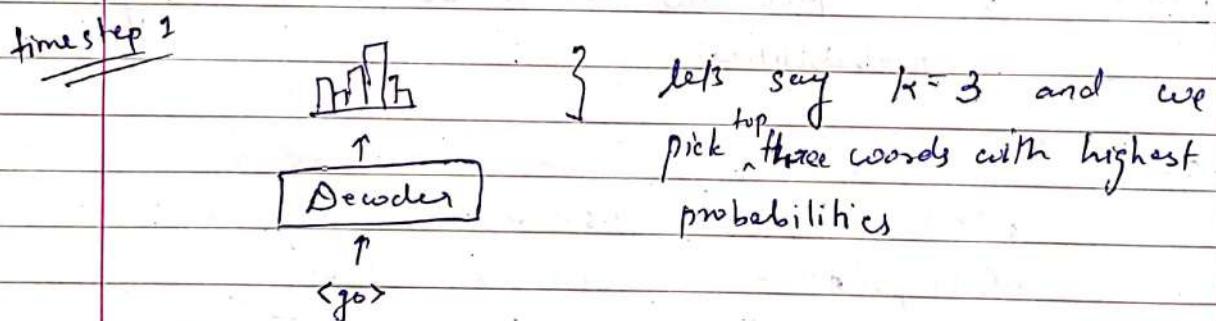
Beam search with $b=(V)$ is equivalent to exhaustive search.

→ Sampling Based Non Deterministic Decoding
(also called Stochastic Decoding) (output is not same everytime)

- Top-K Sampling

This approach is similar to beam search, but we add some random sampling at each time step.

$K \Rightarrow$ the number of words to pick based on highest probability



$$w_1 = 0.31, w_2 = 0.25, w_3 = 0.15$$

Now, we normalize all of the selected words

$$w_1 = \frac{w_1}{w_1 + w_2 + w_3} = \frac{0.31}{0.31 + 0.25 + 0.15} = 0.436$$

$$\omega_2 = \frac{0.25}{0.71} = 0.352$$

$$\omega_3 = \frac{0.15}{0.71} = 0.211$$

Now, we choose a random number between 0 & 1

- interval $\omega_1 = 0 \leftrightarrow 0.436$

- interval $\omega_2 = 0.436 \leftrightarrow (0.436 + 0.352)$
 $0.436 \leftrightarrow 0.788$

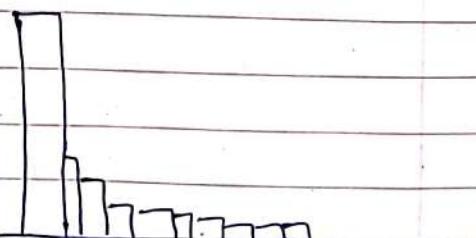
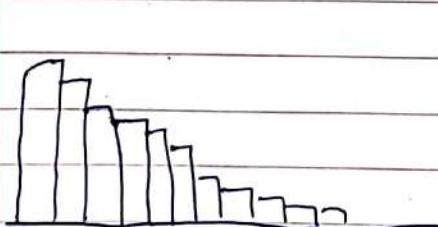
- interval $\omega_3 = 0.788 \leftrightarrow 1$

We pick the word within whose interval the random number lies

Follow the same thing for all the other time steps
 & you will get random sequences

- Temperature Sampling

Let's say we have the following two distributions



Ishan Modi

A

B

Each bar in the distributions represents the probability of a word

The distribution A has probabilities evenly distributed among words, whereas B is peaky distribution

If we pick top-k samples for distribution B it can lead to unwanted results because the probabilities in the k words other than the peak is very less. But the same thing would work fine with distribution A.

To deal with this something known as temperature is introduced

Instead of applying softmax to logits to get probability we can apply softmax with temperature

$$\text{softmax with temperature} = \frac{\exp(\frac{x}{T})}{\sum \exp(\frac{x}{T})}$$

The default value of $T=1$

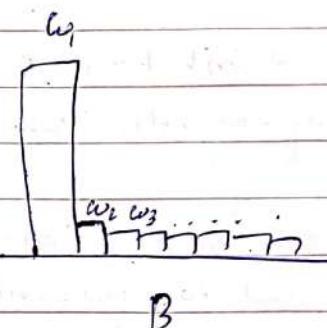
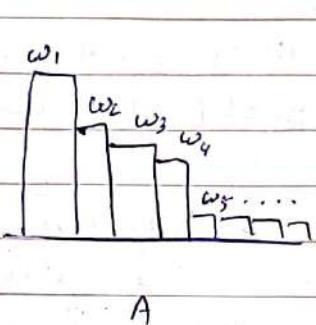
Now, the property of the above function is

- smaller the value of T , it makes larger softmax even larger & smaller softmax smaller. Thus it makes distribution peakier. A less no. of words means less creative

- higher the value of T , it makes smaller softmax larger & larger softmax small. Thus flattens the distribution & leads to more words & more creativity

Ishan Modi

- Top- P (Nucleus) Sampling



Let's say A & B are two distributions at 1st time step.

Now as per top-k sampling we pick the fixed number of words k.

The idea of top-p (nucleus) sampling is to pick dynamic value of k at each time step

Now, we calculate cumulative distribution (CDF) starting from w_1 till CDF becomes less than a threshold p, lets say $p = 0.6$

Example

- For distribution A

$$\begin{aligned}
 w_1 &= 0.3 \\
 w_1 + w_2 &= 0.45 \\
 w_1 + w_2 + w_3 &= 0.55 \\
 w_1 + w_2 + w_3 + w_4 &= 0.59 \\
 w_1 + w_2 + w_3 + w_4 + w_5 &= 0.62
 \end{aligned}$$

} so we include all words below threshold 0.6

$\therefore k = 4$ in this case
and w_1, w_2, w_3, w_4
move to next time step

Ishan Modi

- For distribution B:

$$w_1 = 0.8 \quad \left. \right\}$$

Since CDF > threshold (0.6) & there is no word

$\therefore k=1$ & w_1 goes to next time step

We do similar for all time steps & get randomised sequences