

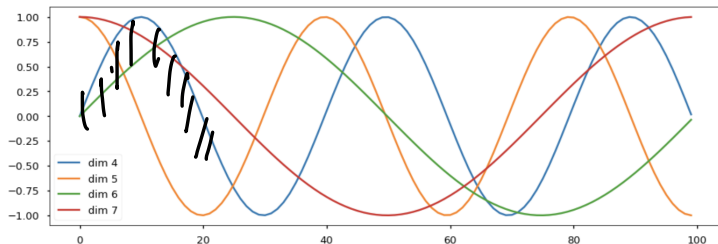
Positional Encoding: Use Sinusoids

Use Sinusoids of varying periods:

$$p_i = \begin{pmatrix} \sin(i/10000^{2\cdot 1/d}) \\ \cos(i/10000^{2\cdot 1/d}) \\ \vdots \\ \sin(i/10000^{2\cdot \frac{d}{2}/d}) \\ \cos(i/10000^{2\cdot \frac{d}{2}/d}) \end{pmatrix}$$

Handwritten notes: 1512 tokens and dim

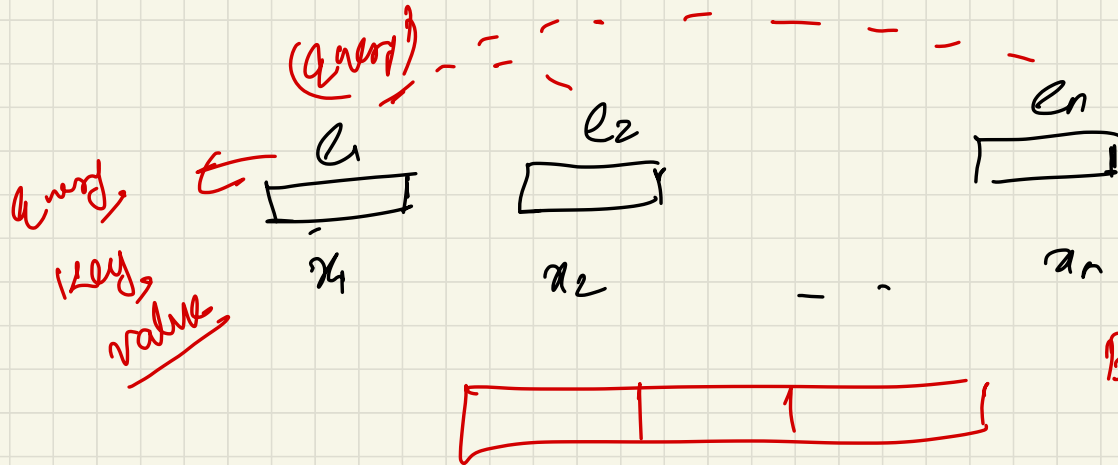
	0	1	0	1	0	1
1	0	1	0	1	0	1
2	0.1	0.9	0.05	0.95		



multi-head attention
single-head?

grammatical structure } $Q^{(1)}$ $K^{(1)}$ $V^{(1)}$ }
 morphologies } attⁿ head 1 } I

Semantic Hops



$$\sum w_{ij} e_j$$

$$w_{ij} = e_i \cdot e_j$$

$$E \quad V \times D$$

Position representations through learned embeddings

- We can posit a new parameter matrix $P \in R^{N \times d}$, where N is the maximum length of any sequence that the model can process
- These position representations are then added to the word embeddings

$$\hat{x}_i = x_i + P_i$$






[BERT]

A quick detail about the vocabulary

Let's take a look at the assumptions we've made about a language's vocabulary.

We assume a fixed vocab of tens of thousands of words, built from the training set.

All *novel* words seen at test time are mapped to a single UNK.

	word		vocab mapping	embedding
Common words	hat	→	pizza (index)	
	learn	→	tasty (index)	
Variations	taaaaasty	→	UNK (index)	
misspellings	laern	→	UNK (index)	
novel items	Transformerify	→	UNK (index)	

Handwritten notes:

- An orange arrow points from "learn" to "taaaaasty".
- An orange bracket underlines "Transformerify".
- An orange arrow points from "Transformerify" to "laern".
- An orange arrow points from "Transformerify" to "UNK (index)".
- The word "nationality" is written in cursive below the embeddings, with red underlines under "national" and "ity".

The byte-pair encoding algorithm

BPE

—

- Start with a vocabulary containing only characters and an “end of word” symbol.
- Using a corpus of text, find the most common adjacent characters ‘a,b’; add “ab” as a subword
- Replace [↑] instances of the character pair with the new subword ; repeat until desired vocab

Byte Pair Encoding

A **word segmentation** algorithm:

- Start with a vocabulary of **characters**
- Most frequent **ngram pairs** \mapsto a new **ngram**

Dictionary

5 low -
2 lower -
6 new est -
3 widest -

Vocabulary

l, o, w, e, r, n, ~~w~~, s, t, i, d -

Start with all characters
in vocab

ow 7
er 9
st 9

Byte Pair Encoding

A word segmentation algorithm:

- Start with a vocabulary of characters
- Most frequent ngram pairs \mapsto a new ngram

8

Dictionary

5 low
2 lower
6 new es t
3 wid es t

↑ ↑
es

es

Vocabulary

l, o, w, e, r, n, w, s, t, i, d, es

Add a pair (e, s) with freq 9

A word segmentation algorithm:

- Start with a vocabulary of characters
- Most frequent ngram pairs \mapsto a new ngram

Dictionary

5 low
2 lower
6 new est
3 wid est

Vocabulary

l, o, w, e, r, n, w, s, t, i, d, es, **est**

Add a pair (es, t) with freq 9

Byte Pair Encoding

- A **word segmentation** algorithm:
 - Start with a vocabulary of **characters**
 - Most frequent **ngram pairs** \mapsto a new **ngram**

Dictionary

5 **lo** w
2 **lo** w e r
6 n e w e s t
3 w i d e s t

Vocabulary

l, o, w, e, r, n, w, s, t, i, d, e s, e s t, **lo**

Add a pair (l, o) with freq 7

7 ab cd m
7 xyz

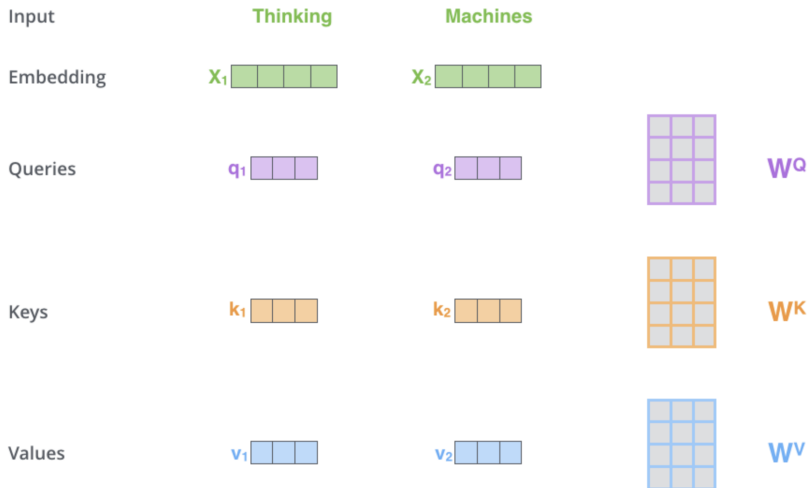
Byte Pair Encoding

32K BPE

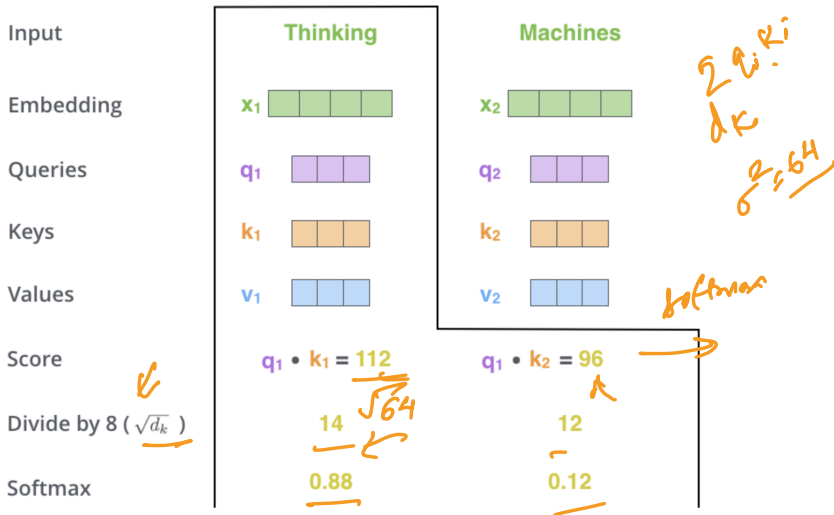
When do we stop

When the desired vocabulary size is met

More Details: Self-attention



Scaled Dot Product for Attention



Why is scaling required?

Scaled dot-product

- Assume that the components of q and k are independent random variables with mean 0 and variance 1.
- Then their dot product $q \cdot k = \sum_{i=1}^{d_k} q_i k_i$ has mean 0 and variance d_k .
- Hence the scaling ensures that the resultant dot product has mean 0 and variance 1.

$\sqrt{d_k}$

Matrix Calculation of Self Attention

X W^Q Q

\times $=$


X W^K K

\times $=$

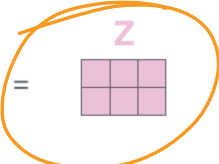
X W^V V

\times $=$

Matrix Calculation of Self Attention



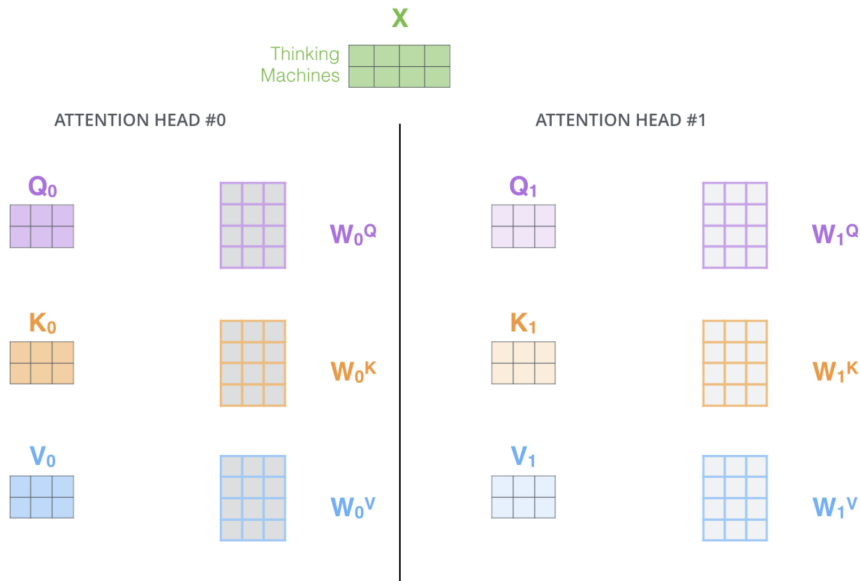
$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$



$$= \begin{matrix} \text{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$



Multi-headed Attention



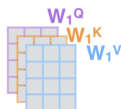
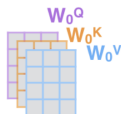
Multi-headed Self Attention Block

- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

Thinking
Machines



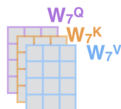
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



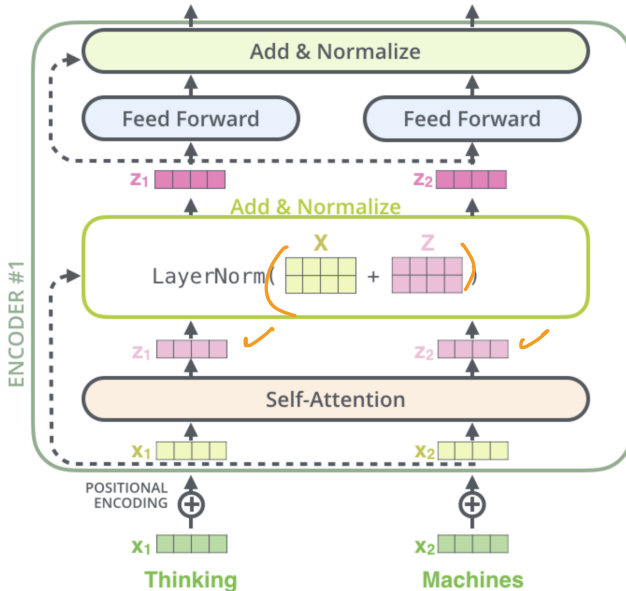
...

...

...



Encoder Block



What is Layer Normalization?

Basic Idea

Cut down on uninformative variation in hidden layer vectors by normalization

How does the normalization work

Let $x = \{x_1, \dots, x_d\}$ be an individual word vector in the model.

Let the mean be $\mu = \frac{1}{d} \sum_{i=1}^d x_i$, $\mu \in R$

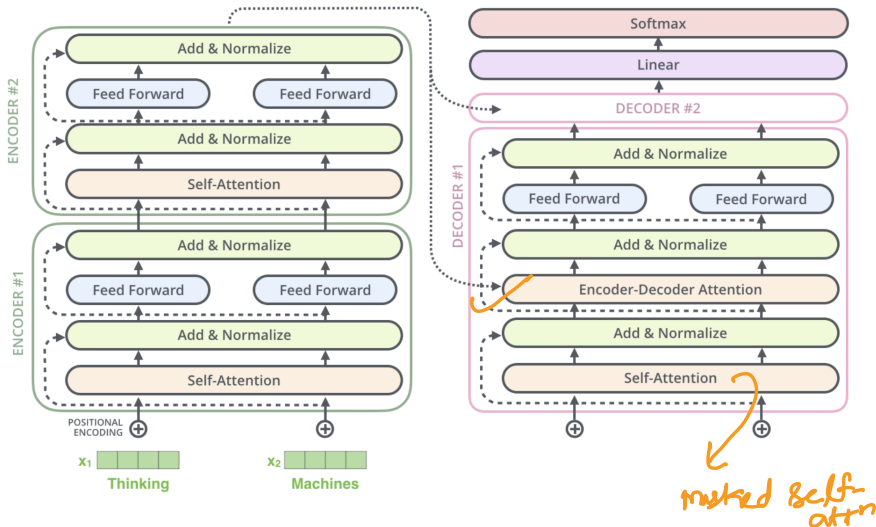
Let the variance be $\sigma = \sqrt{\frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2}$, $\sigma \in R$

Let $\gamma \in R^d$ and $\beta \in R^d$ be learnable parameters (can omit)

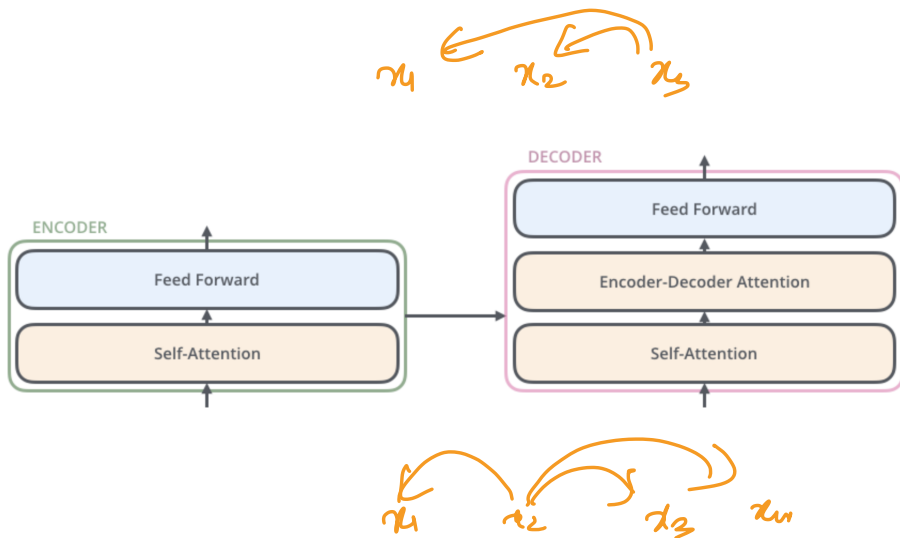
$$\text{LayerNorm} = \gamma \left(\frac{x - \mu}{\sigma + \epsilon} \right) + \beta$$

Original paper also used γ and β as learnable parameters, similar to BatchNorm, but later papers showed that those are not important (and are even harmful).

Transformer with 2 stacked encoders and decoders



How is the decoder different?

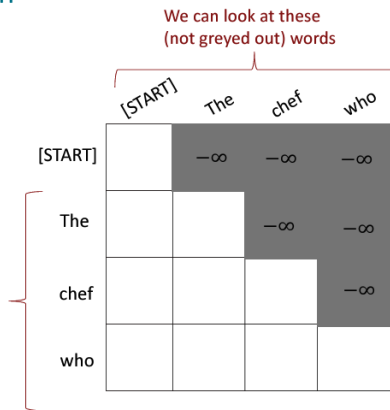


Masking the future in self-attention

- To use self-attention in **decoders**, we need to ensure we can't peek at the future.
- At every timestep, we could change the set of **keys and queries** to include only past words. (Inefficient!)
- To enable parallelization, we **mask out attention** to future words by setting attention scores to $-\infty$.

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

For encoding these words



I went to the market

h₁ went — —

P went to the market

