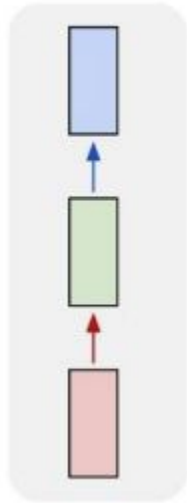


Sequence Modeling

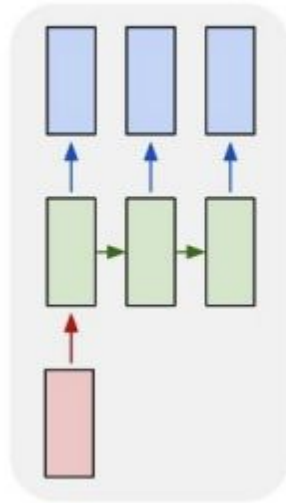
Attention, Transformers, and BERT

N. Rich Nguyen, PhD
SYS 6016

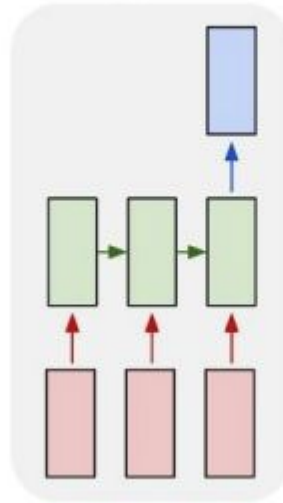
Neural network architectures are very flexible



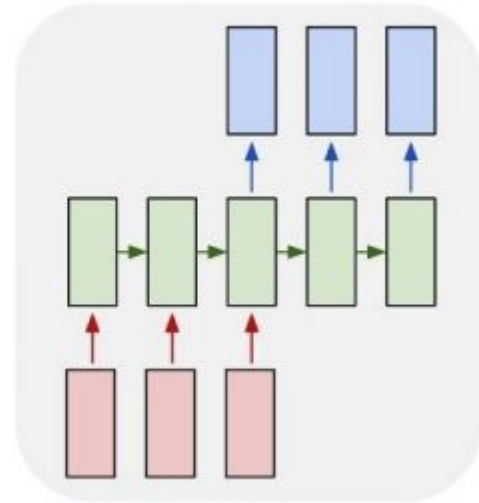
vector-to-vector



vector-to-sequence



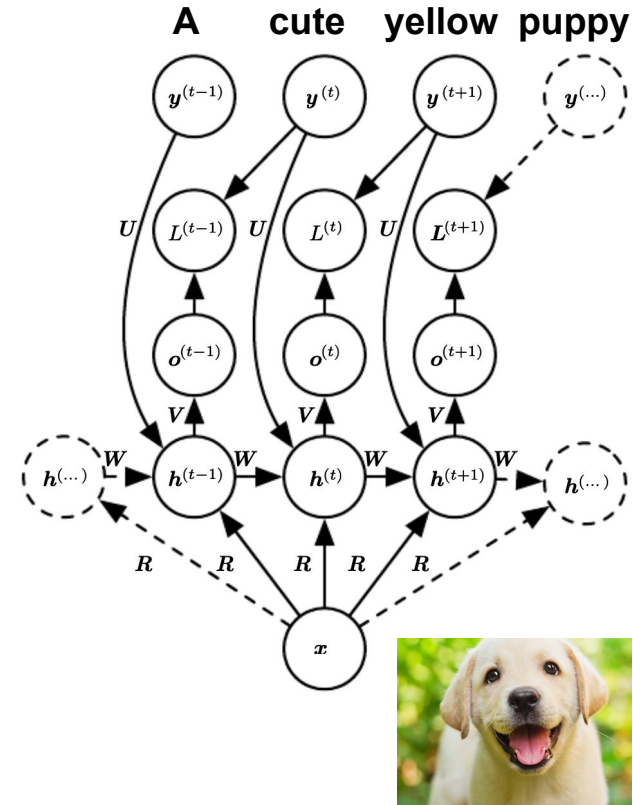
sequence-to-vector



sequence-to-sequence

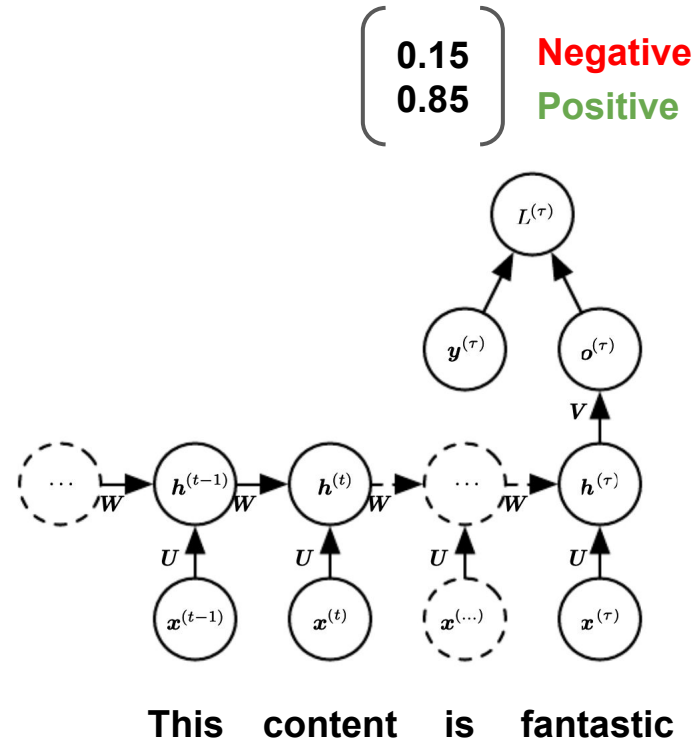
Vector-to-sequence Model

- Maps a fixed-length vector \mathbf{x} into a distribution over sequences $\mathbf{y} \rightarrow$ **Image Captioning**.
- Each element $y^{(t)}$ of the output sequence serves both as label (for the previous time step) and as input (for the current time step).



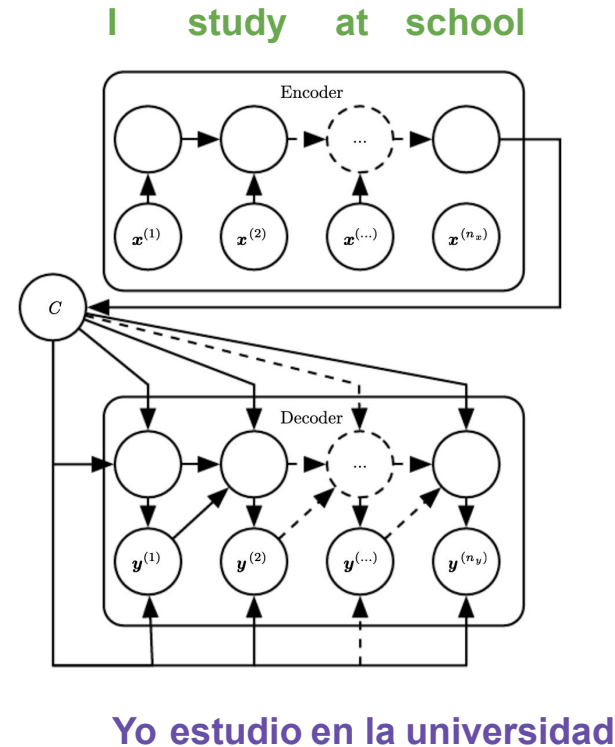
Sequence-to-vector Model

- Maps a sequence of inputs to a single output at the end of the sequence
- Such a network can be used to summarize a sequence and produce a label at the end → **Sentiment Analysis**



Sequence-to-sequence Model

- Aka encoder-decoder model: generates an output sequence given an input sequence → **Machine Translation**
- The final hidden state of the encoder is used to compute a fixed-size context vector c , which represents a ***semantic summary*** of the input sequence and is given as input to the decoder

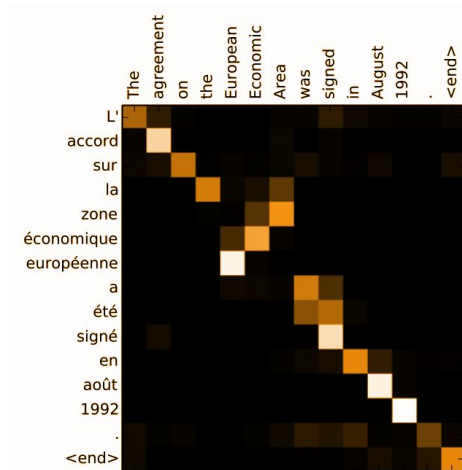


RNNs with Attention Mechanism



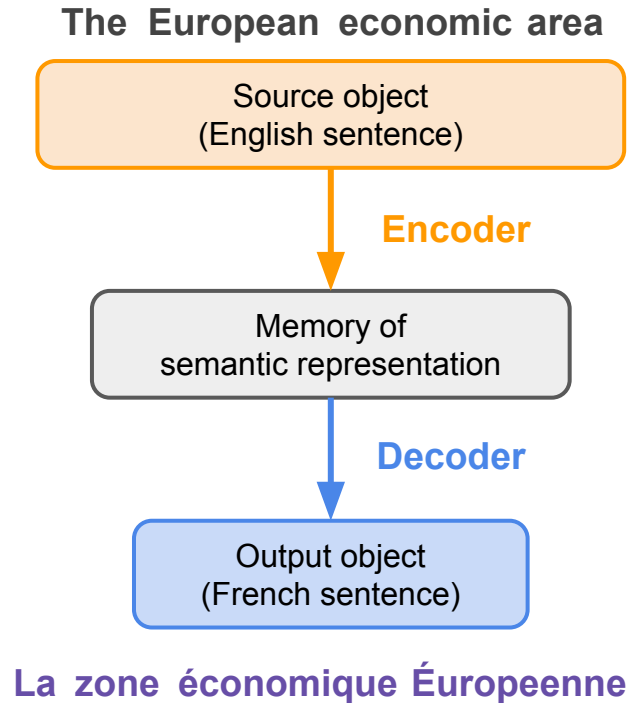
RNNs with Attention Mechanism

- A core idea of a 2014 paper by **Dzimitry Bahdanau** and **Yoshua Bengio**
 - Read the whole sentence or paragraph (to get the context),
 - Produce the translated words one at a time,
 - *Each time focusing on a different part of the input sentence to gather semantic details.*
- Attention mechanism **revolutionized** machine translation and NLP in general



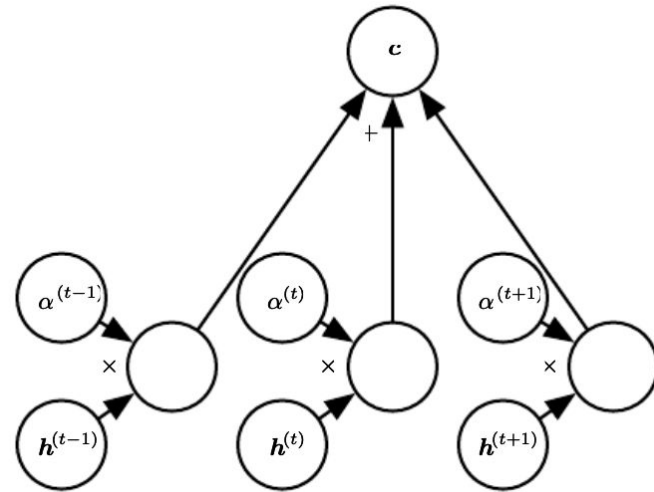
Components of an Attention model

1. **Encoder**: reads words in a sentence and converts them into a feature vector associated with **each word position in the sequence**.
2. **Memory**: A list of feature vectors storing output of the encoder → **memory** containing a sequence of semantic representation.
3. **Decoder**: exploits the content of the memory, at each time step put attention on the content of *one memory element* → translated word.



Intuitions of Attention models

- A attention mechanism is a **weighted average**: a context c is formed by taking weighted average of encoder's hidden state $h^{(t)}$ with weights $\alpha^{(t)}$.
- Weights $\alpha^{(t)}$ are produced by the model itself by applying a `softmax()` over scores measuring *how well each output is aligned with the decoder's previous hidden state*.
- The attention of weighted average is a smooth, differentiable approximation that can be trained with backprop.



Limitations of RNNs and LSTM

1. **Very hardware intensive**: difficult and slow to train
2. **Not feasible for longer sequence**: LSTM on a 100 word document has gradients like a 100-layer network
3. **No parallel processing**: Input data needs to be passed in sequentially
4. **No available pretrained model**: needs specific labelled data for every task

If only we have a **pretrained** network that can utilize **parallelization** for **long sequential** data and use **attention** mechanism?

⇒ **Transformers!**

Transformers



No, not these Transformers!

Transformers - The next big thing, for now!

- In a groundbreaking 2017 paper “**Attention is All You Need**”, a team of Google researchers take the **attention mechanism** idea to create an architecture called the “**Transformer**”.
- Transformer architecture is much *faster to train* and *easier to parallelize*.
- Transformer significantly improve the state of the art **without** using any recurrent or convolutional layer, just the attention mechanism.
- **Sequence Modeling:** Bag-of-words → RNNs / LSTM → Transformers



Demo of “Writing With Transformer”

As we are learning about machine learning, machine learning systems will continue to evolve at an accelerating pace and we must prepare for it. We must create a well defined and cohesive education and training system that is accessible to a wide variety of students , including those in the developing world, by teaching the skills they need in order to take advantage of new technologies.

Written by Transformer · transformer.huggingface.co 

Transformer Architecture

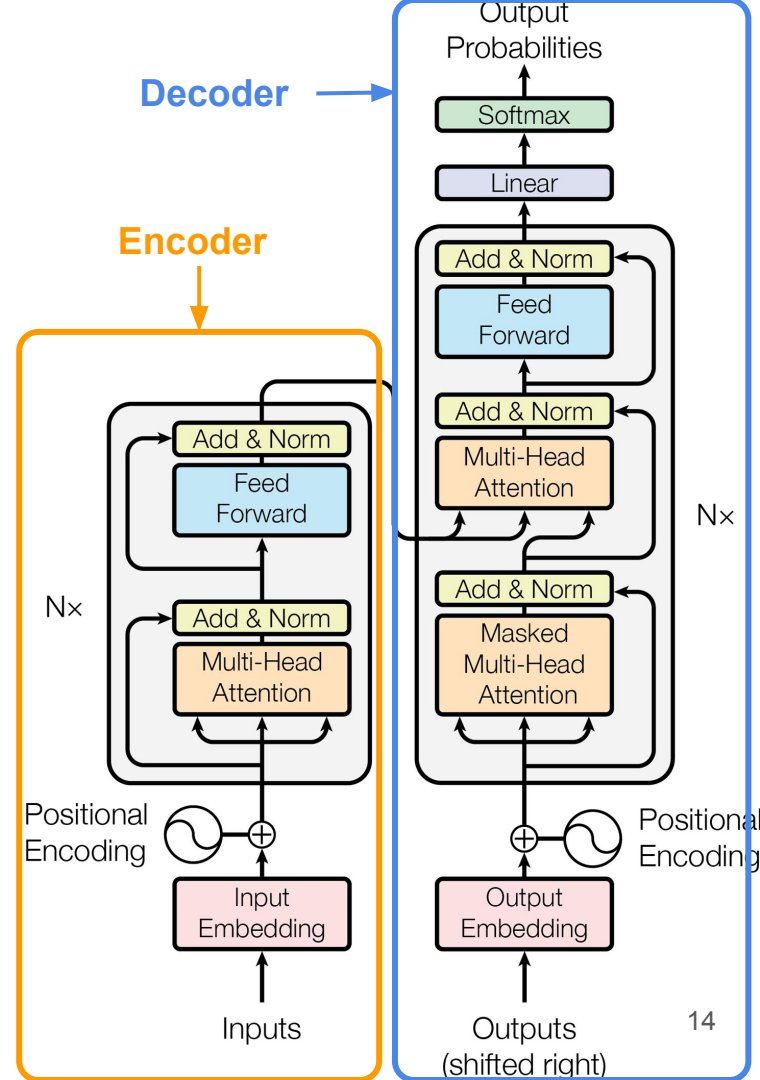
Encoder:

- Takes as input a batch of sequence of words
- Encodes each word into a 512-dimension embedding (from GloVe or Word2Vec).
- Is stacked N times.

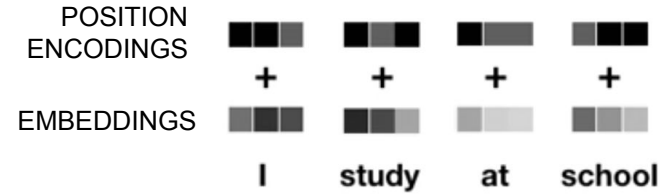
I study at school

Decoder:

- Takes the target sentence as input, shifted one time step to the right b/c it should be from a previous time step.
- Receives the output of encoder.
- Is stacked N times.
- Outputs a probability for each possible word, at each time step.



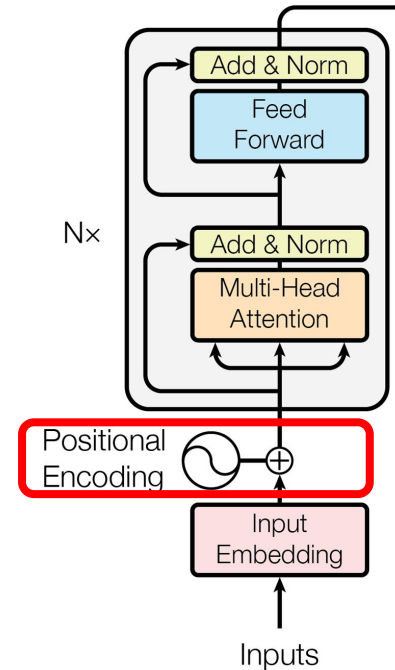
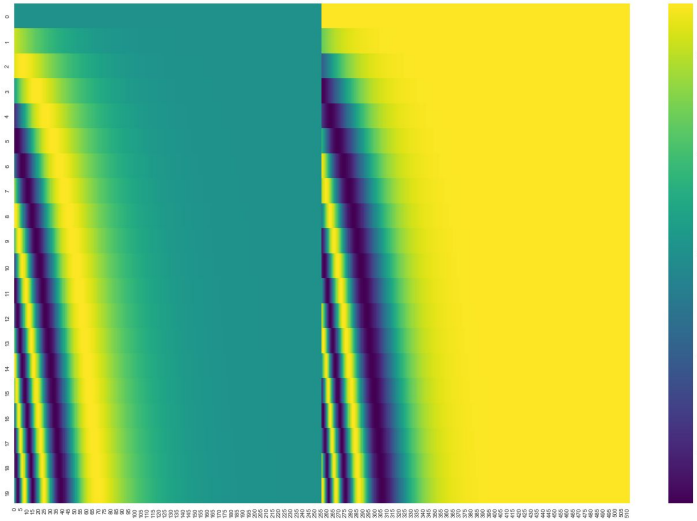
Positional Encodings



- Are simply dense vectors (same size d as word embeddings) that represent **the position of a word in the sentence**.
- Using sine and cosine functions, compute a matrix $\mathbf{PE}(p, i)$ where p is the word position in a sentence and i is the embedding dimension

$$\mathbf{PE}_{p,2i} = \sin\left(\frac{p}{10000^{2i/d}}\right)$$

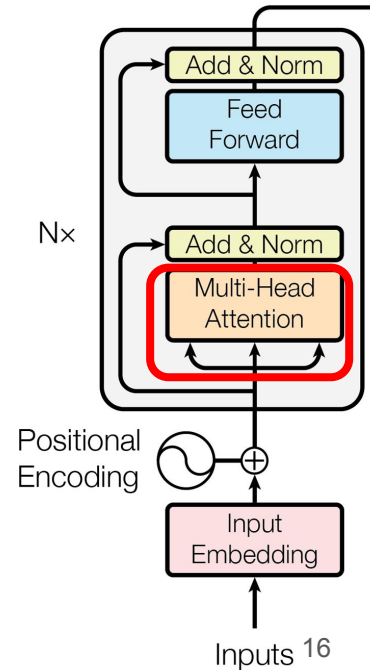
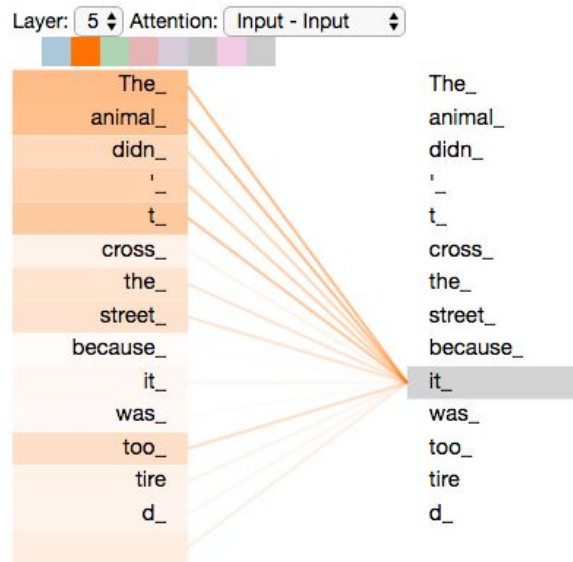
$$\mathbf{PE}_{p,2i+1} = \cos\left(\frac{p}{10000^{2i/d}}\right)$$



Self-attention

















First let's discuss self-attention! Consider this sentence:

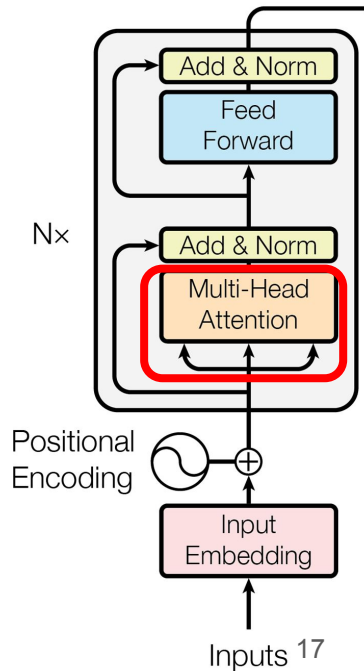
"The animal didn't cross the street because it was too tired"



Compute Self-attention with Query, Key, Value

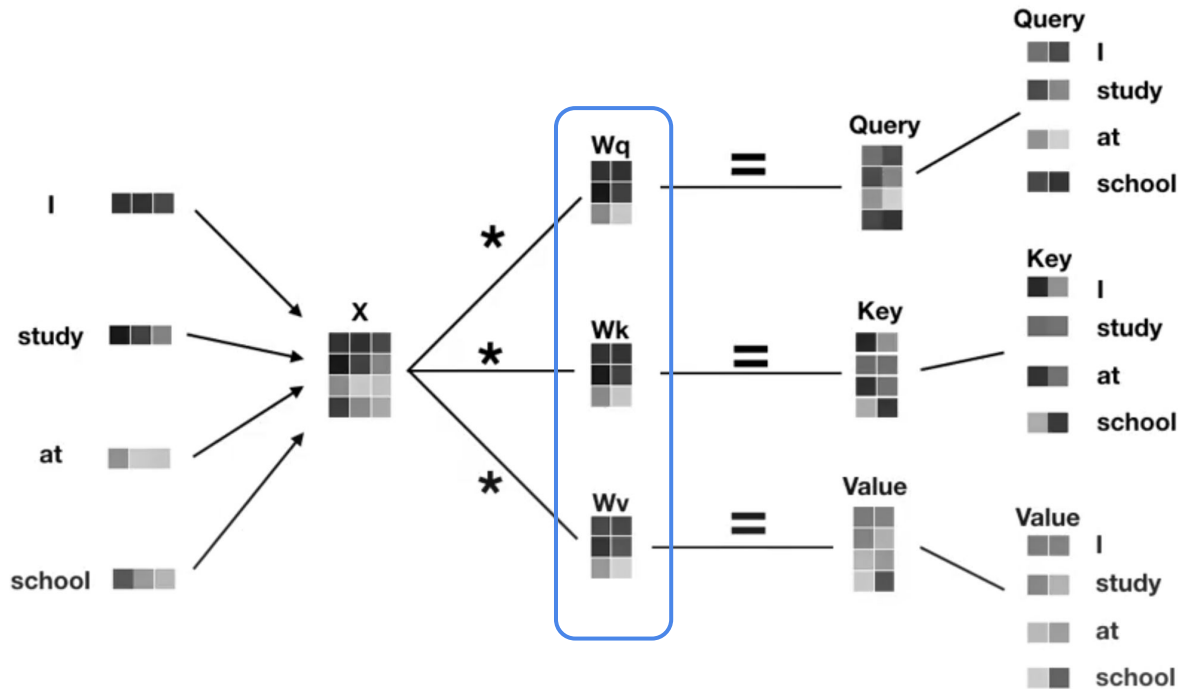
- Inspired from the information retrieval system: Query \rightarrow (Key, Value)
- For every encoded word, generate a (**Key**, **Value**) pair as well as a **Query**
- The **self-attention score** is the product of **Query** and **Key**

	X	Query	Key	Value
I				
study				
at				
school				

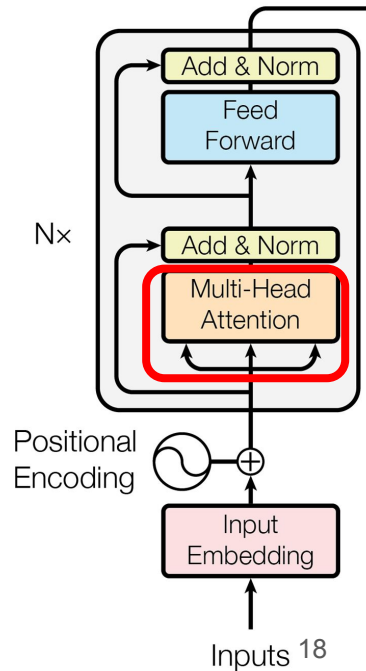


Generating Query, Key, Value

Learning the weight matrices to generate Query, Key, and Value

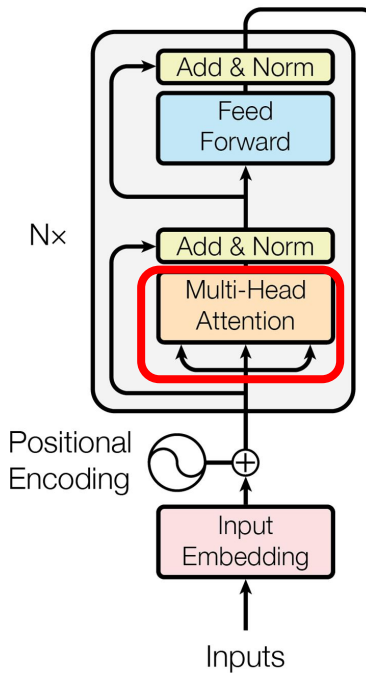


Learn these weights by backpropagation

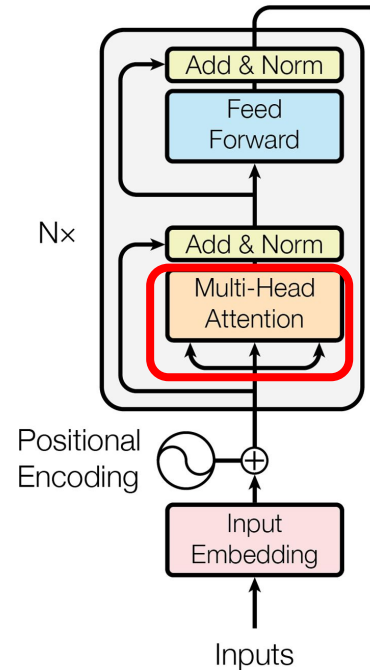
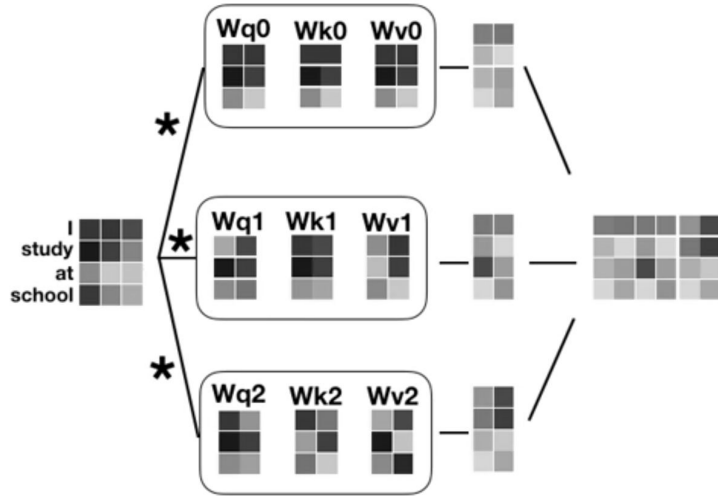


Computing one-to-one comparisons in parallel

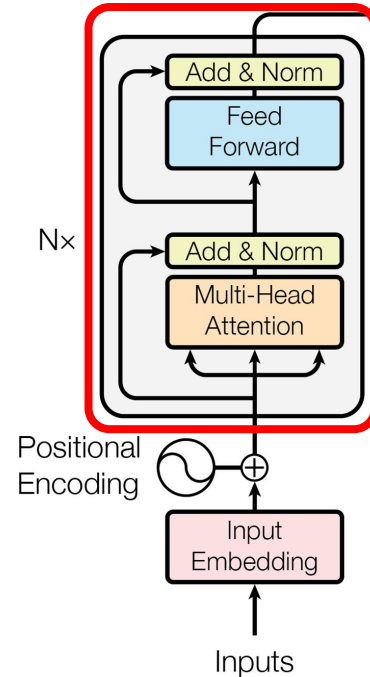
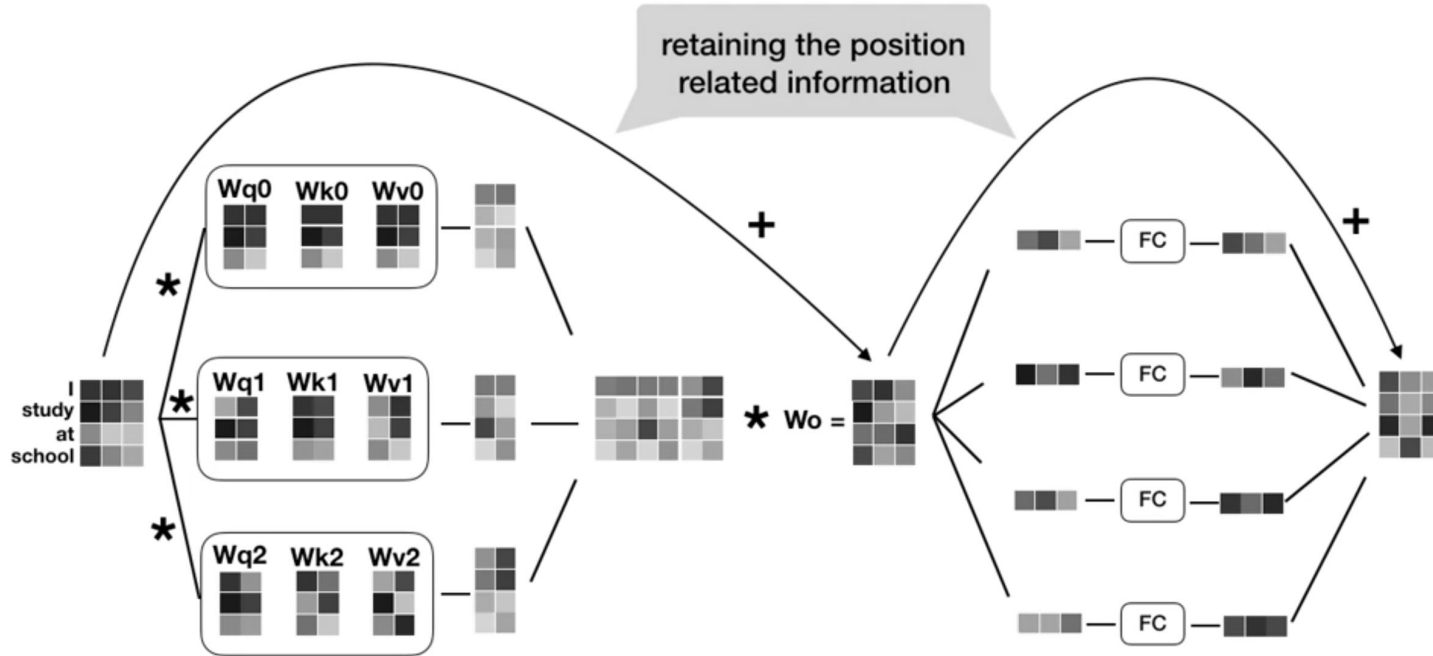
		Query * Key ^T	Score	Softmax	Value	Softmax * Value	Σ Softmax * Value (Context Vector)
I	I * I	= 130	0.92		I		}
	I * study	= 50	0.05		study		
	I * at	= 20	0.02		at		
	I * school	= 10	0.01		school		
study	study * I	= 30	0.02				}
	study * study	= 110	0.70				
	study * at	= 20	0.03				
	study * school	= 70	0.25				
at	at * I	= 30	0.03				}
	at * study	= 50	0.10				
	at * at	= 90	0.80				
	at * school	= 40	0.07				
school	school * I	= 30	0.01				}
	school * study	= 80	0.27				
	school * at	= 23	0.02				
	school * school	= 160	0.70				



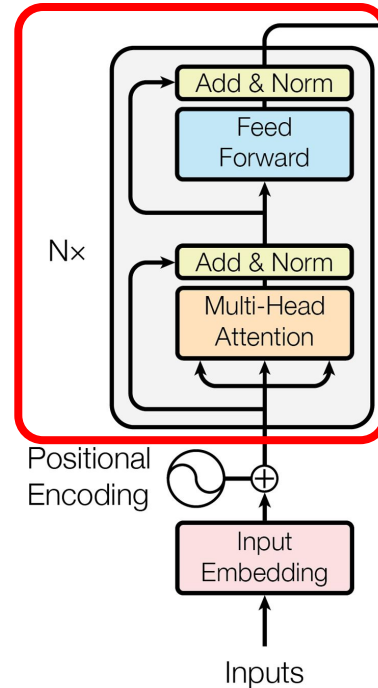
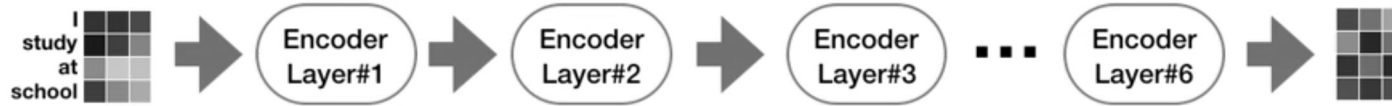
Multi-headed Attention - more parallelism



Add Residual & Layer Normalization



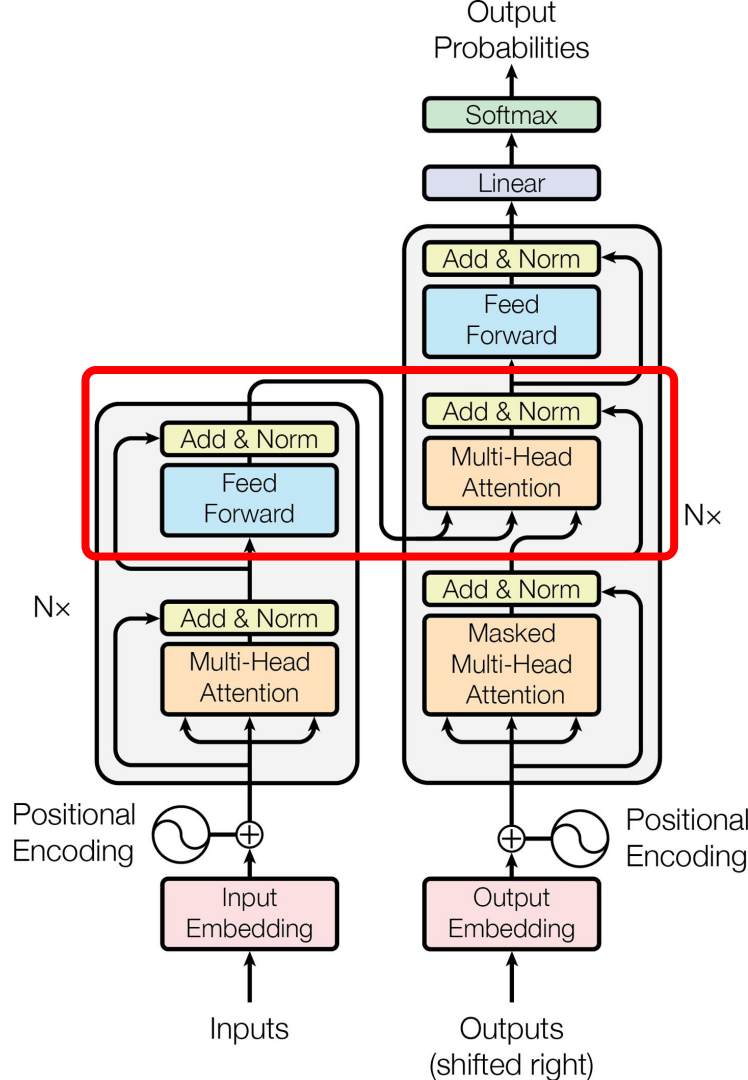
Stacking N encoder layers



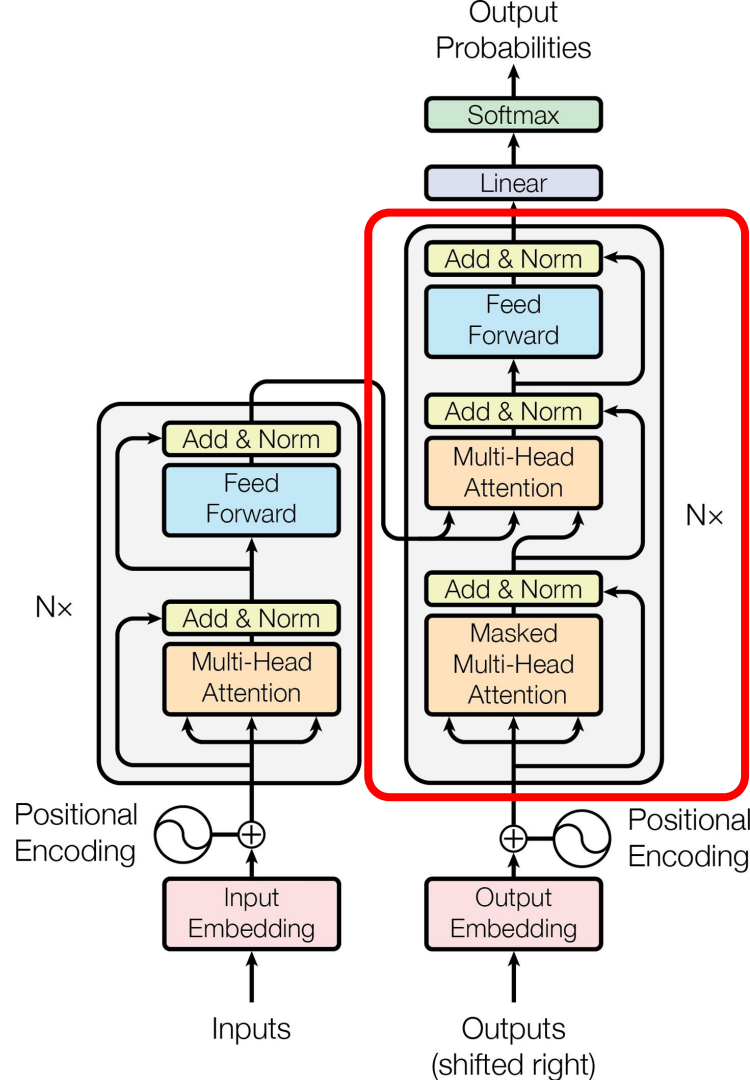
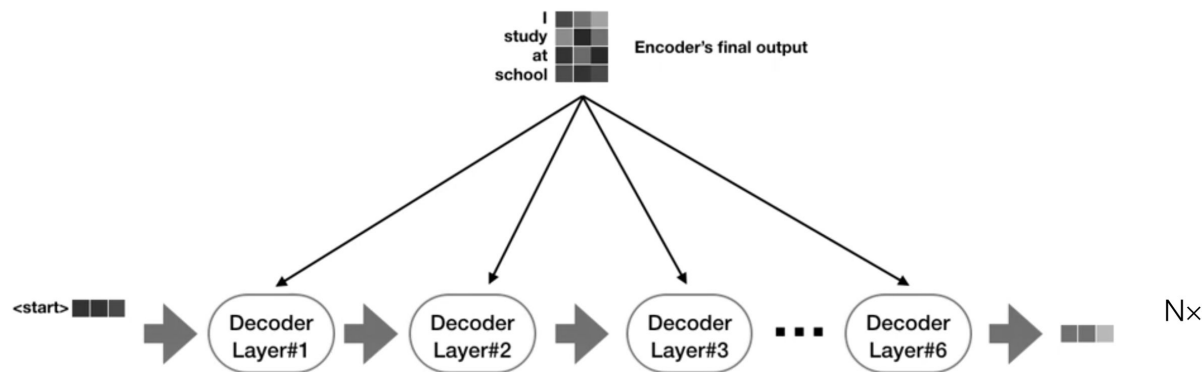
Encoder-Decoder Attention

Encoder-decoder attention is slightly different from self-attention:

- create its **Query** matrix from the layer below it, which is decoder self-attention,
- take the **Key** and **Value** matrices from the output of the encoder stack.
- help the decoder focus on appropriate places in the input sequence



Stacking decoder layers



Decoder side: Look-Ahead Mask

- To prevent future labels from being seen by current time step
- Top triangle of the matrix become $-\infty$ → after softmax it becomes 0 → no attention on those future cells

	<start>	I	am	fine
<start>	0.7	0.1	0.1	0.1
I	0.1	0.6	0.2	0.1
am	0.1	0.3	0.6	0.1
fine	0.1	0.3	0.3	0.3

Scaled Scores

0.7	0.1	0.1	0.1
0.1	0.6	0.2	0.1
0.1	0.3	0.6	0.1
0.1	0.3	0.3	0.3

Look-Ahead Mask

0	$-\infty$	$-\infty$	$-\infty$
0	0	$-\infty$	$-\infty$
0	0	0	$-\infty$
0	0	0	0

+

=

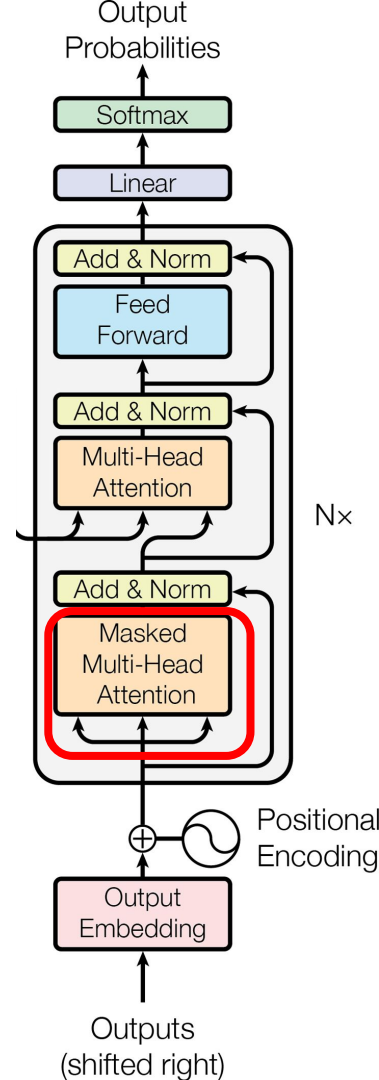
Masked Scores

0.7	$-\infty$	$-\infty$	$-\infty$
0.1	0.6	$-\infty$	$-\infty$
0.1	0.3	0.6	$-\infty$
0.1	0.3	0.3	0.3

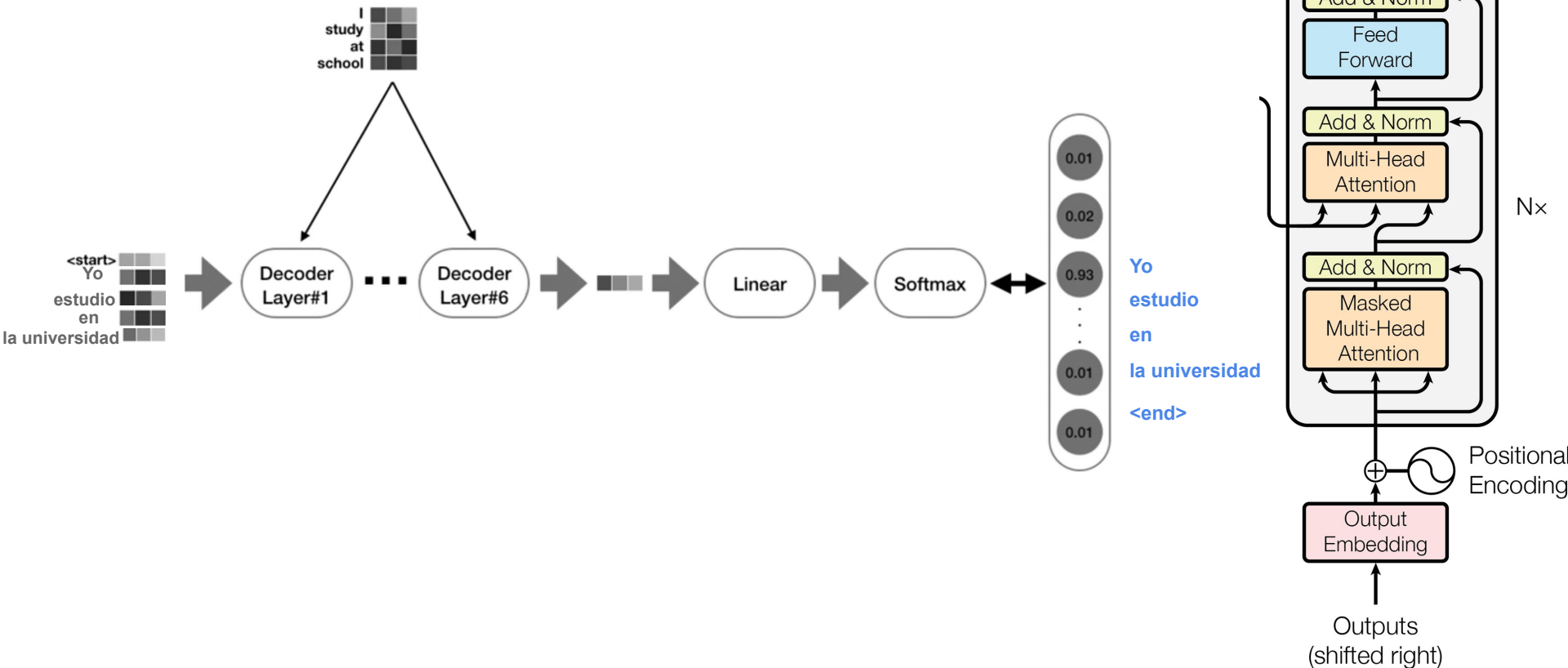
Softmax()

Softmax Scores

1	0	0	0
0.37	0.62	0	0
0.26	0.31	0.43	0
0.21	0.26	0.26	0.26

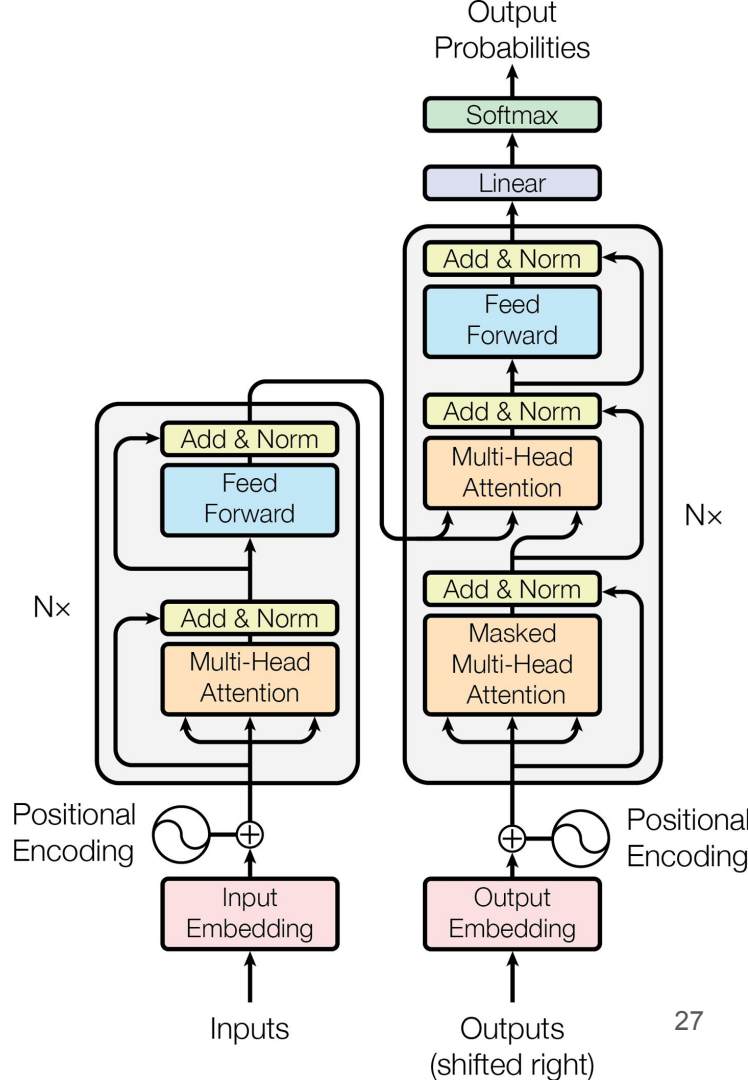


Output Layer

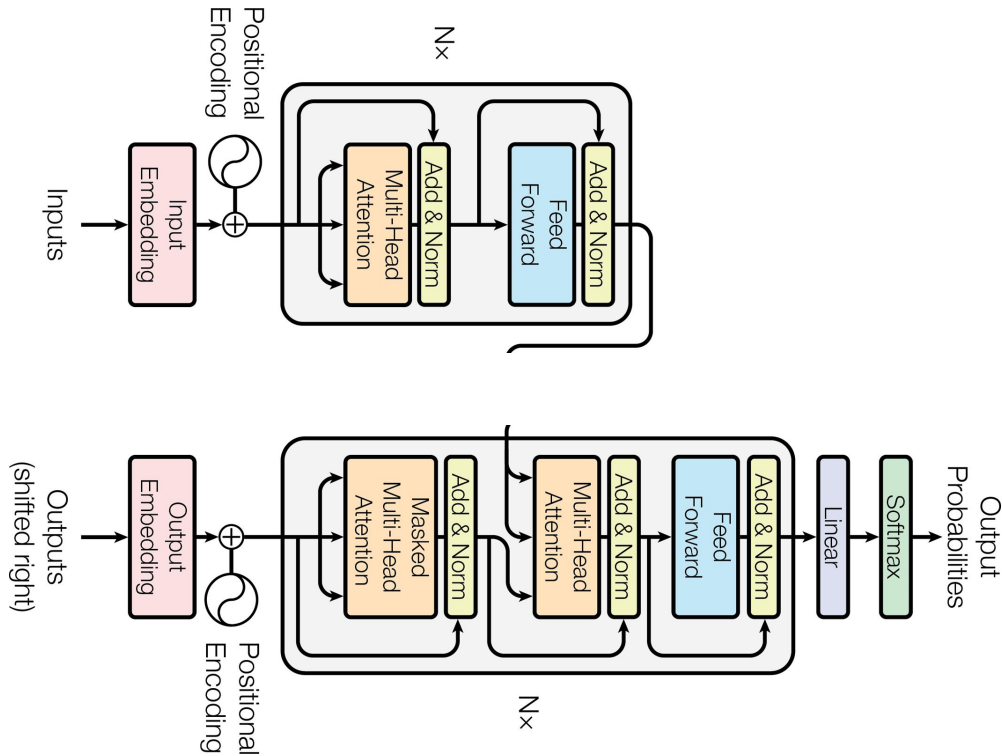


Properties of Transformers

- Does not suffer from short term memory! → **infinite** reference window → entire context!
- All-to-all comparisons, so each layer has $O(n^2)$ for sequence of length n . However, parallelism helps speedup training.
- Every output is a weighted sum of every input, and the weighting is a learned function.
- Because of this transformer architecture, the NLP industry can achieve unprecedented results.



Transformer Quest



Encoder:

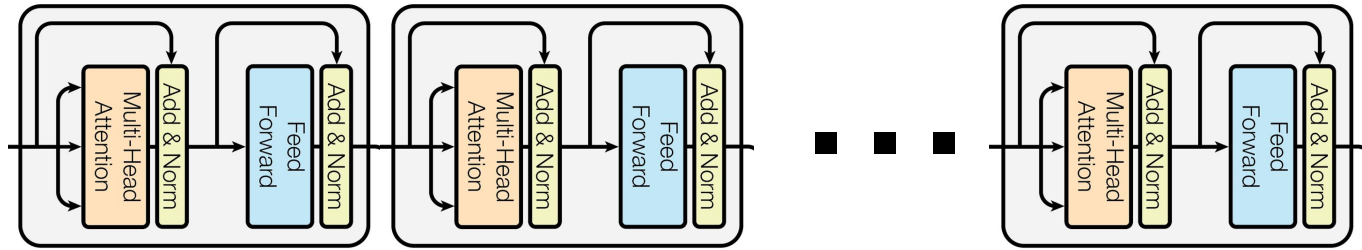
How to map the semantic context?

What is language?

Decoder:

How to map words from one language to another?

Stacking Transformer Encoders



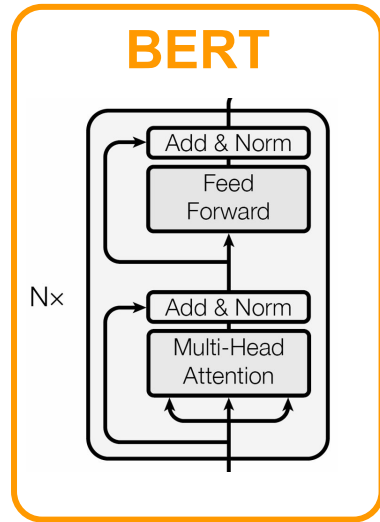
BERT: Bidirectional Encoder Representation from Transformers

BERT



No, not this Bert!

What is BERT?



BERT is a large transformer masked language model.

24 layers, 1024 nodes, 16
attention heads, 340M params

a new family of neural network
architectures based on attention

a kind of contextual word embedding
and can be used as input embedding

a pre-trained statistical model to predict
the probability of a sentence or phrase.

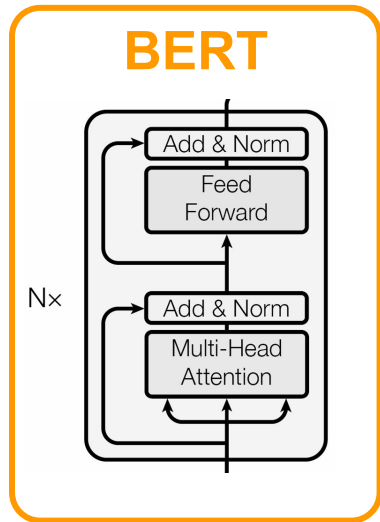
Fun fact: BERT was trained
on 16 TPUs for 4 days!

BERT Applications and Training

Solve more tasks than the regular Transformer:

- Neural Machine Translation (same as Transformers)
- Question Answering
- Sentiment Analysis
- Text Summarization

→ **BERT tries to understand language!**

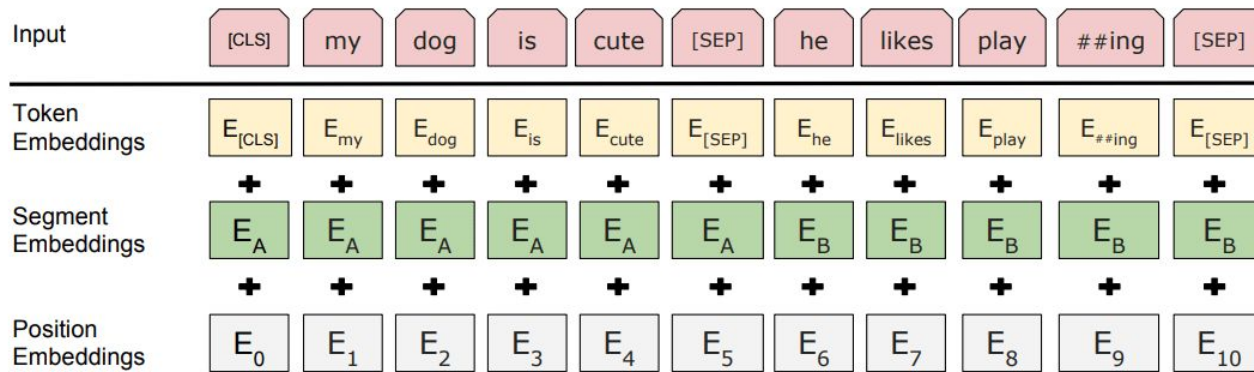


BERT has a two-phase training:

1. Pretrain to understand **language**
2. Finetune to learn a **specific task**

Pre-training - Input representation + metadata

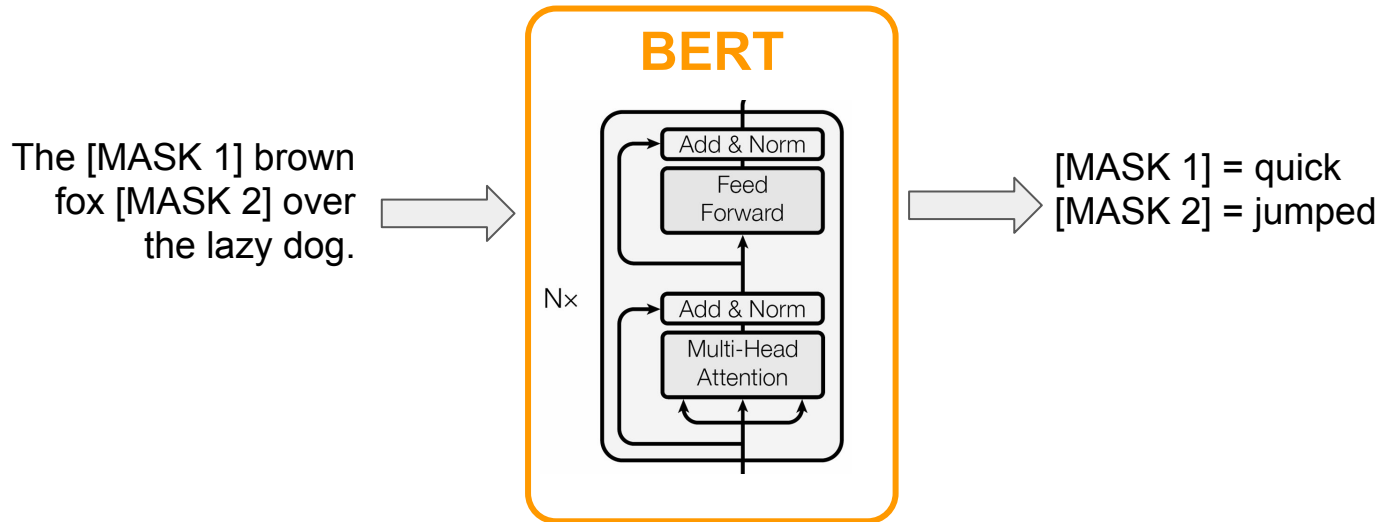
- BERT is trained on the **WikiText** dataset of over **100 million tokens** extracted from the set of verified Good and Featured articles on Wikipedia.



- Token embeddings:** A [CLS] at beginning, [SEP] at end of each sentence.
- Segment embeddings:** A marker distinguish sentences (A or B or ...)
- Positional embeddings:** A encoding to indicate input's position in the sequence

Pre-training - Mask Language Modeling (MLM)

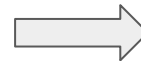
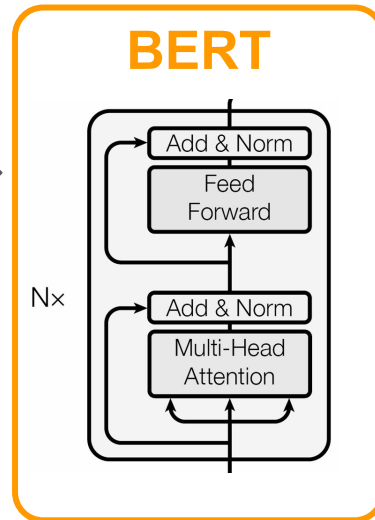
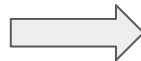
- Randomly maskout 15% of the words under [MASK] tokens
- **Use supervised learning** to predict only the masked words using the context of non-masked words.
- To avoid the model only try to predict [MASK] tokens, replace 10% of [MASK] tokens with a random one, and another 10% are left unchanged.



Pre-training - Next Sentence Prediction (NSP)

- Insteads of the next word, BERT tries to predict the **entire sentence**.
- To understand relationship between sentences separated by [SEP] token, BERT tries to determine the next sentence given a sentence with masks.
- **Use Supervised Learning:** Feed 2 sentences at a time, half the time the 2nd one comes next after the 1st one, another half it is a random one from corpus.

[CLS] the man went to [MASK] store [SEP]
 he bought a gallon [MASK] milk [SEP]
 : [CLS] the man [MASK] to the store [SEP]
 penguin [MASK] are flight ##less birds [SEP]



IsNext

NotNext

Fine tuning BERT

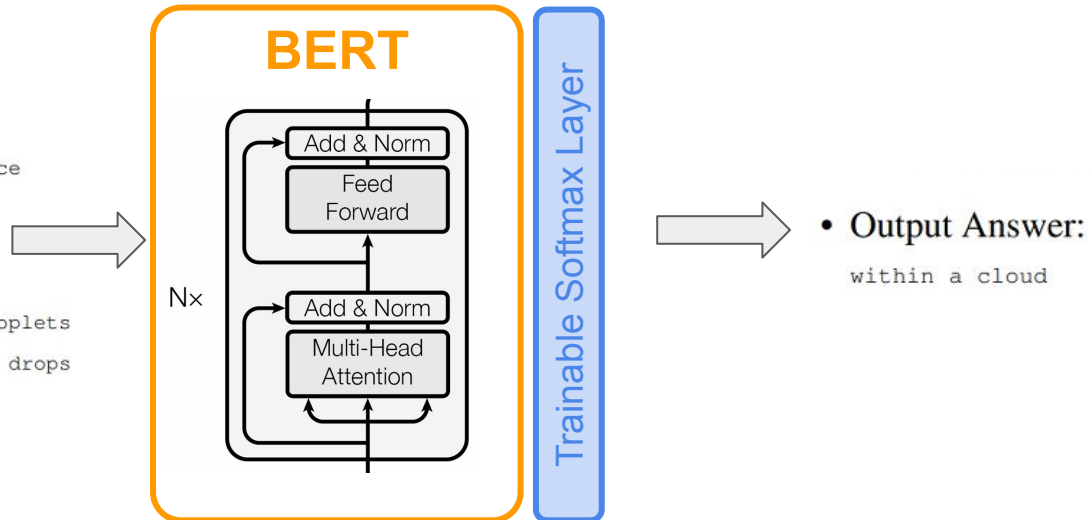
- After importing a pre-trained model, add a single layer on top of the core model for a specific task.
- **On a question answering task:** upon receiving a question as input, the goal is to identify the right answer from some corpus → learn two extra vectors **<start>** and **<end>** that marks the beginning and the end of the answer.

Input Question:

Where do water droplets collide with ice crystals to form precipitation?

Input Paragraph:

... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. ...



Fine tuning BERT (in 12 lines of code)

```
# Load dataset, tokenizer, model from pretrained
model/vocabulary

tokenizer = BertTokenizer.from_pretrained('bert-base-cased')
model = TFBertForSequenceClassification.from_pretrained('bert- #Pre-train BERT base model with cased text
base-cased')
data = tensorflow_datasets.load('glue/mrpc') #Microsoft Research Paraphrase Corpus dataset

# Prepare dataset for GLUE as a tf.data.Dataset instance
train_dataset = glue_convert_examples_to_features(data['train'],
                                                tokenizer, max_length=128, task='mrpc')
valid_dataset = glue_convert_examples_to_features(data['validation'],
                                                tokenizer, max_length=128, task='mrpc')
train_dataset = train_dataset.shuffle(100).batch(32).repeat(2)
valid_dataset = valid_dataset.batch(64) #Use operations from TensorFlow Dataset class

# Prepare training: Compile tf.keras model with optimizer, loss and
learning rate schedule
optimizer = tf.keras.optimizers.Adam(learning_rate=3e-5, epsilon=1e-08,
                                     clipnorm=1.0)

loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
metric = tf.keras.metrics.SparseCategoricalAccuracy('accuracy')
model.compile(optimizer=optimizer, loss=loss, metrics=[metric]) #Similar setup as other neural nets

# Train and evaluate using tf.keras.Model.fit()
model.fit(train_dataset, epochs=2, steps_per_epoch=115,
        validation_data=valid_dataset, validation_steps=7) #Train the top layer
```

BERT benefits and drawbacks

- + Transferable model: can be used as input to smaller, task specific models
 - + With successful fine tuning, can have very good accuracy
 - + Pretrained BERT models available in 100+ languages
-
- Big and slow to train (but you got the pretrain model)
 - Expensive (only if you want to train from scratch)
 - Needs to be fine-tuned for downstream tasks (this can be finicky)



Summary

- **Transfer Learning** works! Pretrain models can be fine-tuned for new tasks
- **Attention mechanism** revolutionizes machine translation and NLP
- **Transformers** and **BERT** are easier to train, more efficient neural networks
- **NLP Models:** Bag-of-words → RNNs / LSTM → Transformers → BERT → ?
- For neat tools on NLP, check out the **Hugging Face** project (huggingface.co)!



Bonus Content

Visual Attention Model

- Attention mechanism was applied in image captioning using **visual attention**
- **Idea:** a convolutional neural network first processes the image and outputs some feature maps, then a decoder RNN with attention generates the caption, one word at a time.
- At each time step, the decoder use attention model to focus on just the right part of the image to generate a word in the caption.



The model's focus before producing the word "frisbee"

MegatronLM

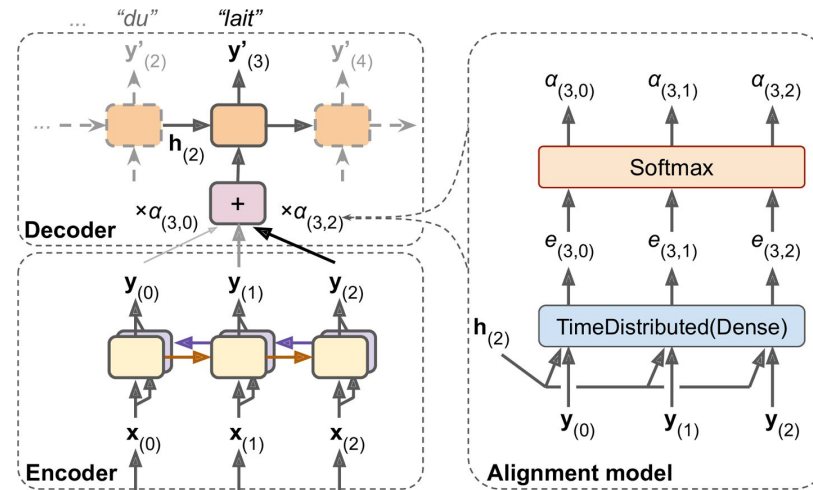
- Massive Transformer Language Model from NVIDIA
- 8.3B parameters
- Trained on 512 GPUs for 9 days (\$450K on EC2)
- Produced way to much green gas, but it is **reusable!**

XLM - RoBERTa

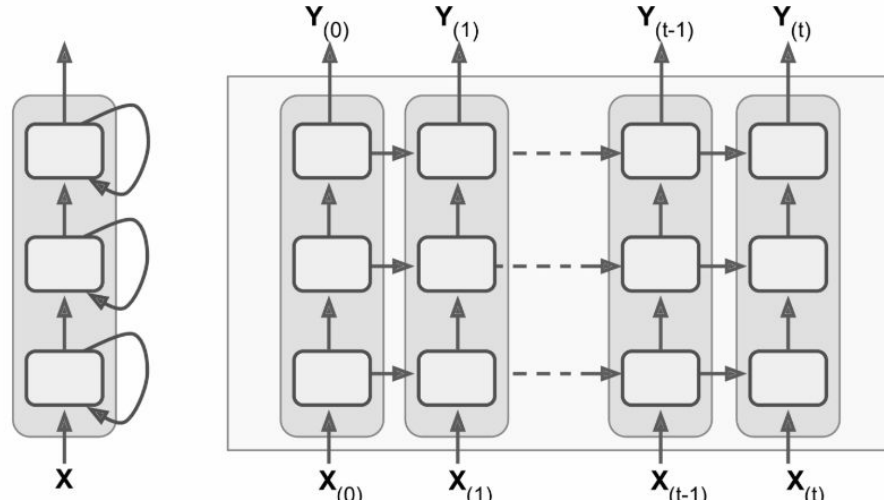
- Trained on 2.5TB of data is 100 languages
- Leverage other existing work

Intuitions of an Attention model

- Instead of sending the final hidden state to the decoder, **the encoder** send all outputs
- At each time step, **the decoder** computes a weighted sum $\alpha_{(t,i)}$ of encoder output i^{th} at time step t to determine which *words* it will focus on at that time step
- Weight $\alpha_{(t,i)}$ are generated by a small neural net called **an alignment model**, which is jointly trained
- Alignment model outputs a score for each **encoder output** to measure *how well each output is aligned with the decoder's previous hidden state*.



Implementing a deep RNN / LSTM model



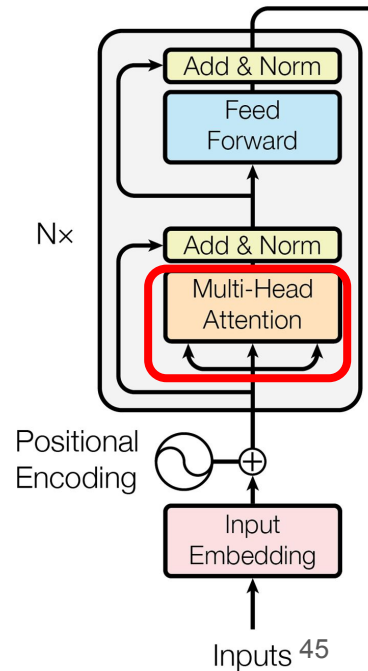
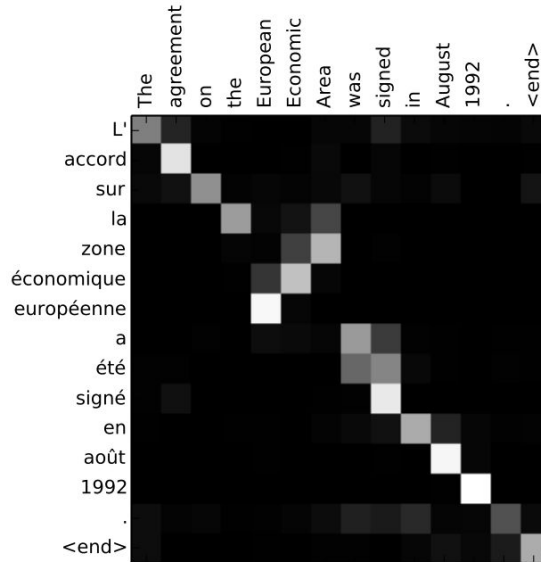
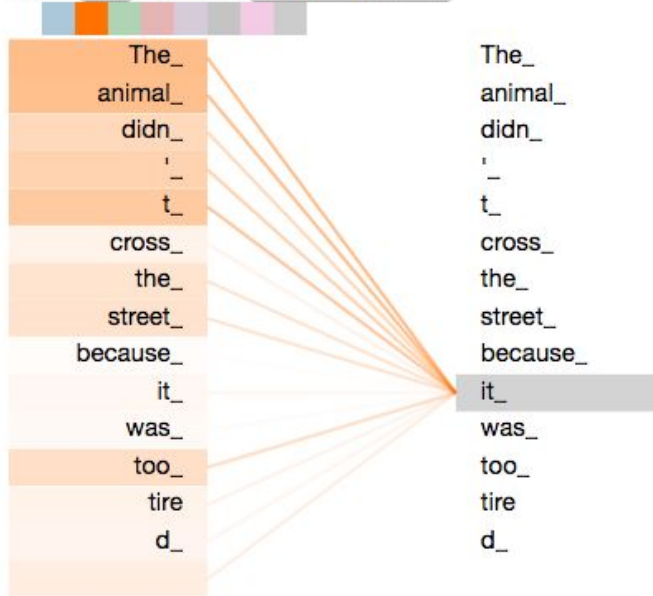
```
model = keras.models.Sequential([
    keras.layers.LSTM(20, return_sequences=True, input_shape=[None, 1]),
    keras.layers.LSTM(20, return_sequences=True),
    keras.layers.TimeDistributed(keras.layers.Dense(10))
])
```

Multi-headed Attention - Self-attention

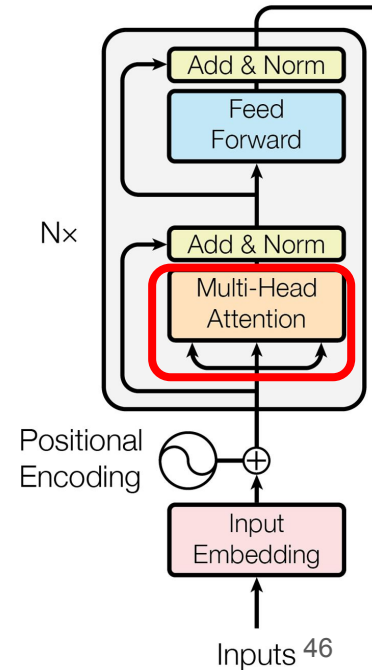
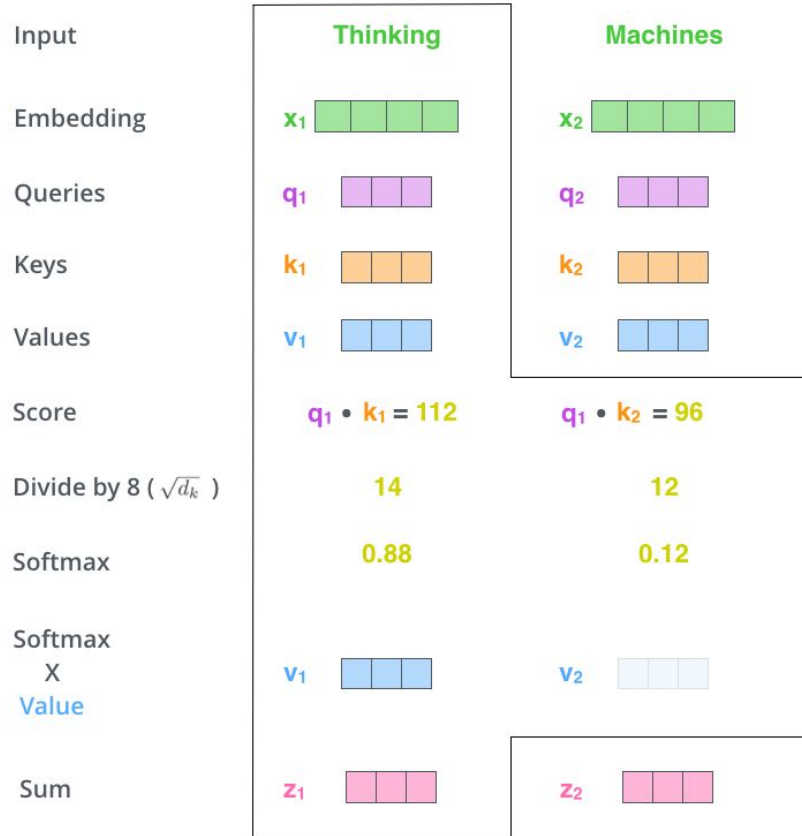
First, let's discuss Self-attention! Consider this sentence:

"The animal didn't cross the street because it was too tired"

Layer: 5 Attention: Input - Input

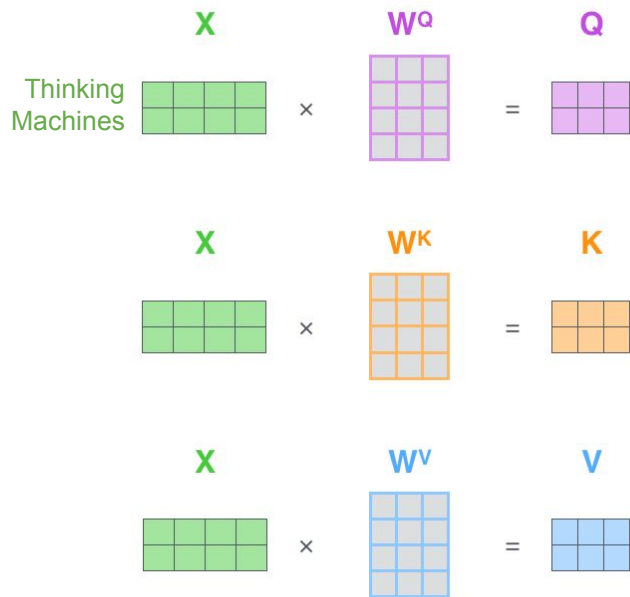


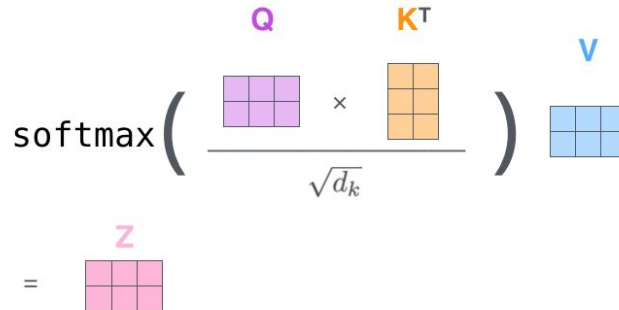
Multi-headed Attention - Self-attention

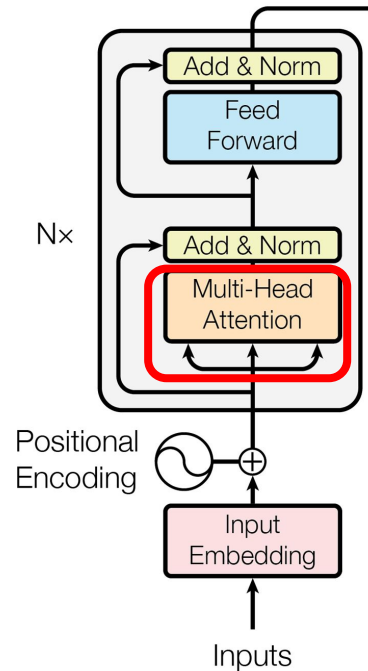


Multi-headed Attention - Matrix calculation

Compute attention weight for each input (how each word should attend to other words in the sequence)



$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) \times V = Z$$




Multi-headed Attention - Put it together

1) This is our input sentence*

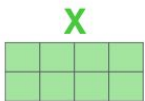
2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

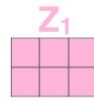
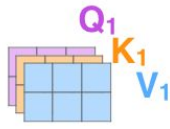
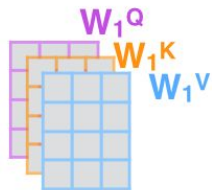
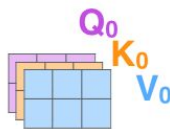
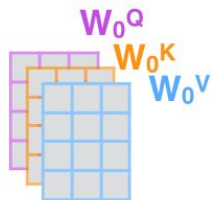
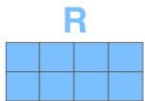
4) Calculate attention using the resulting $Q/K/V$ matrices

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



...

...

...

