

Cornell Bowers CIS

## Logistics

- HW2 is out
- We have a feedback form (due Friday February 23)
- Tuesday 10am office hours might change
- We will talking about projects on Thursday
  - HW3 will be a shorter

Cornell Bowers CIS

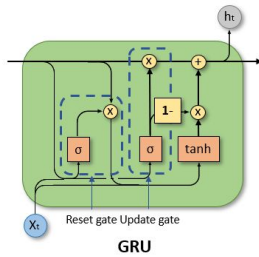
## Previously: Using LSTMs to solve sequence problems

- Process sequences one element at a time.
- Maintain a 'memory' (cell state) to capture information about previous steps.
- Mitigates the RNN vanishing gradient problem
- Suitable for time series, speech, text, and other sequential data.

Cornell Bowers CIS

## Bidirectional LSTM

## Gated recurrent units (GRUs)



$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

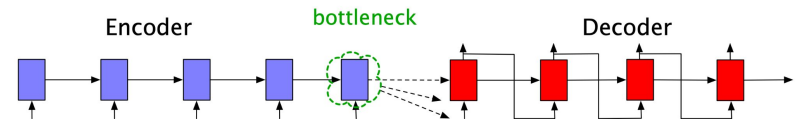
$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

$$\hat{h}_t = \phi(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$$

$$h_t = (1 - z_t) \odot \hat{h}_t + z_t \odot h_{t-1}$$

## Bottleneck Problem

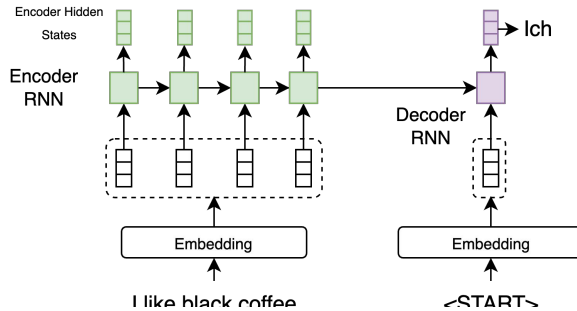
- All the information about the source sequence must be stored in a single vector
  - How to translate a long paragraph?
  - How to summarize long articles?



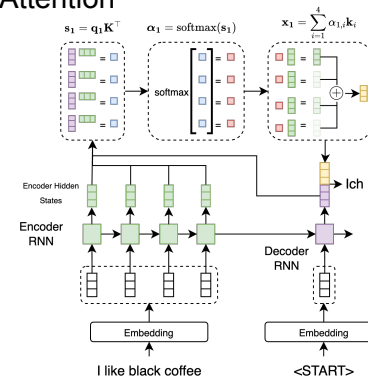
<https://web.stanford.edu/~jurafsky/slp3/>

## RNN for Machine Translation

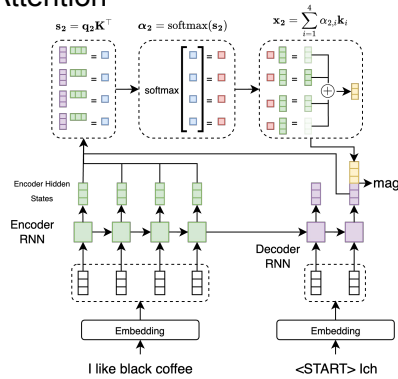
Would be nice if we could “look back” at previous hidden states



## RNN with Attention

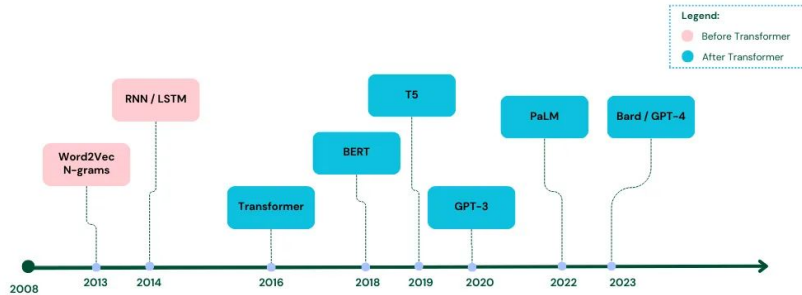


### RNN with Attention



Discuss: What are limitations of such sequence to sequence models? Hint: Think about runtime.

### Language Modelling History



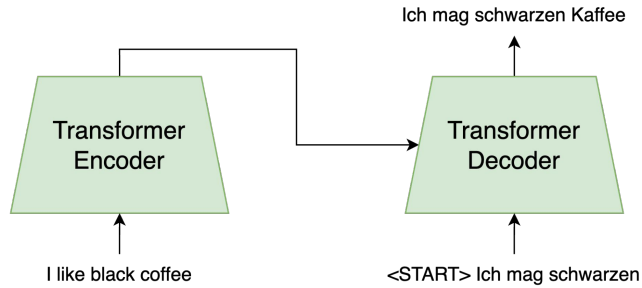
<https://medium.com/@kirudang/language-model-history-before-and-after-transformer-the-ai-revolution-bedc7948a130>

### Self-Attention

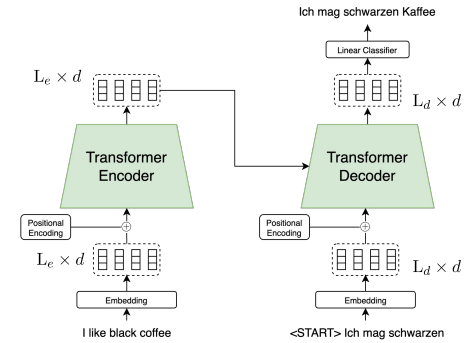
A **bat** flew out of the dugout, startling the baseball player and making him drop his **bat**.

## Transformer Architecture

Introduced for seq2seq tasks like Machine translation, summarization, question answering, etc.



## Transformer Architecture

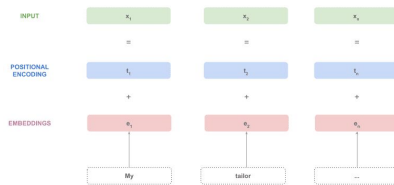


## Input Embeddings

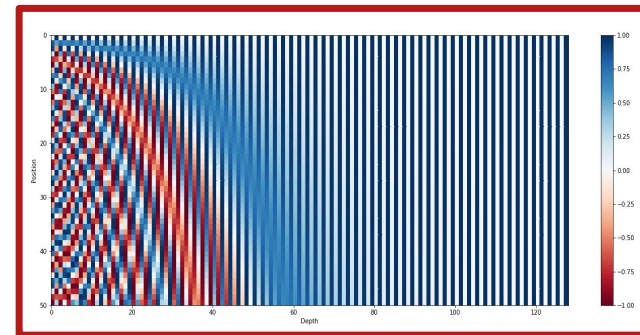
- Replace tokens with continuous vectors
- Made up of two components:
  - token embeddings
  - positional encoding depends on position and dimension index as follows:

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

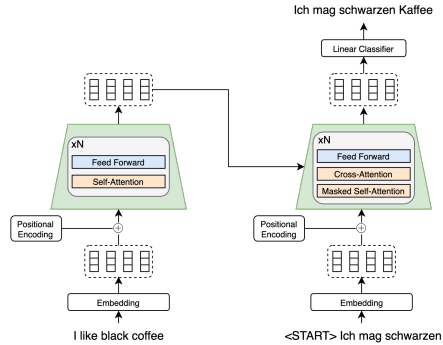
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$



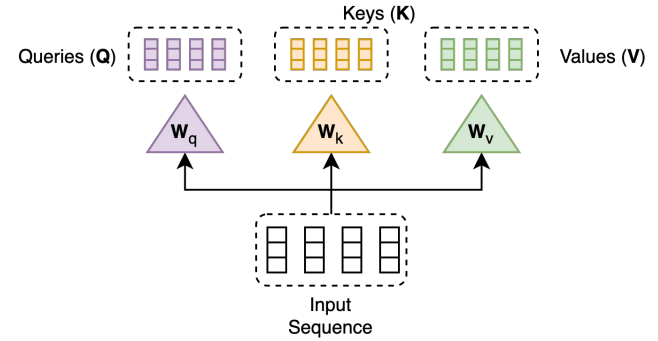
## Example Positional Encoding



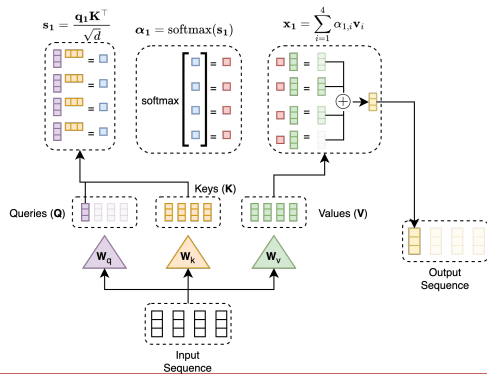
# Transformer Architecture



# Self-Attention



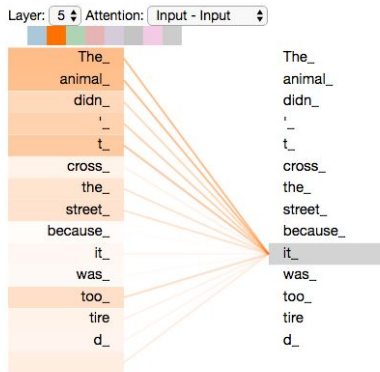
# Self-Attention



# Self-Attention: General Formula

$$Attention(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{d_k}\right)V$$

## Self-attention



## Discuss:

- Q, K, V are all (n x d) matrices. Consider have an input of shape b x n x d.
- What is the shape of  $QK^T$ ?
  - What does this matrix represent?
- What is the shape of the final output?
  - What does this matrix represent?

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{d_k}\right)V$$

## Multi-Head Attention

What if I want to pay attention to different things at the same time!?

**Content-based** This is my big red dog, Clifford.

**Description-based** This is my big red dog, Clifford.

**Reference-based** This is my big red dog, Clifford.

What's useful depends on the task. How do I pick what to do?

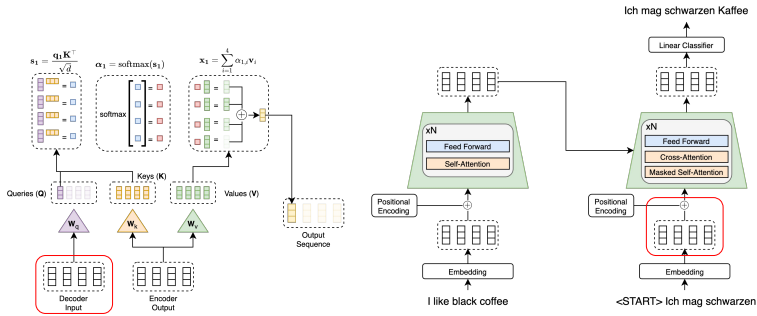
## Multi-Head Attention

- The Scaled Dot-Product Attention attends to one or few entries in the input key-value pairs.
- Idea: apply Scaled Dot-Product Attention multiple times on the linearly transformed inputs.

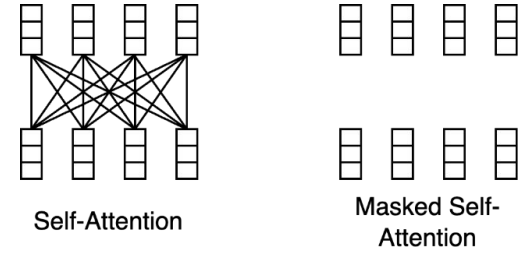
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where  $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

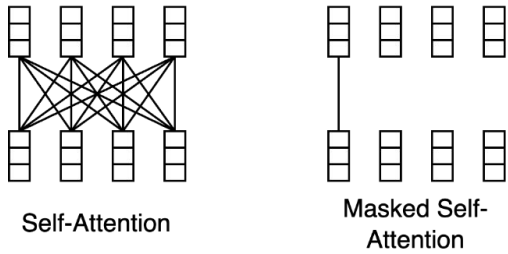
### Cross-Attention



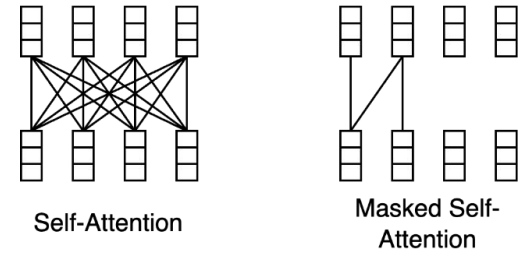
### Self-Attention vs. Masked Self-Attention



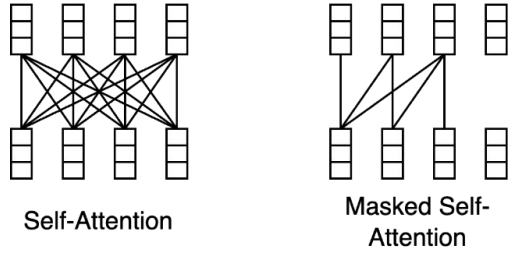
### Self-Attention vs. Masked Self-Attention



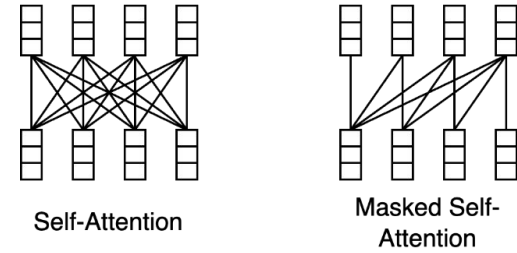
### Self-Attention vs. Masked Self-Attention



## Self-Attention vs. Masked Self-Attention



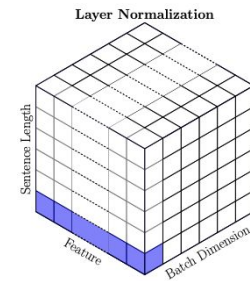
## Self-Attention vs. Masked Self-Attention



## Point-wise Feed-forward Networks

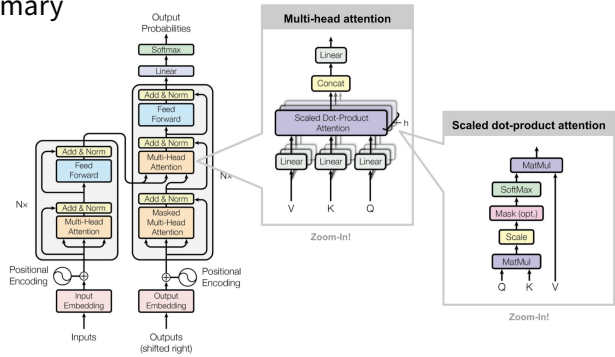
- Purpose
  - Applies non-linear transformations to the output of the attention layer
- Equation
  - $FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$ 
    - where  $W$  and  $b$  are learned weights and biases
- These FFN is applied separately to each position

## Layer Norm



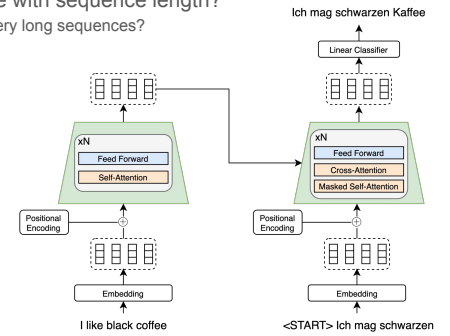


### Summary



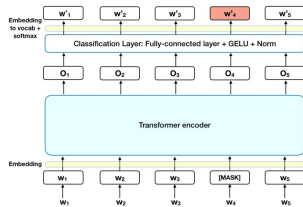
### Discuss:

- How does the transformer scale with sequence length?
  - Any problems with applying it to very long sequences?



### BERT (Bidirectional Encoder Representations from Transformers)

- Bidirectional Context
- Pre-trained on the language, and then fine-tuned



### BERT - Input Representation

Input:

- Use 30,000 WordPiece vocabulary on input.
- Each token is sum of three embeddings

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	#ring	[SEP]
Token Embeddings	$E_{[CLS]}$	$E_{my}$	$E_{dog}$	$E_{is}$	$E_{cute}$	$E_{[SEP]}$	$E_{he}$	$E_{likes}$	$E_{play}$	$E_{ring}$	$E_{[SEP]}$
Sentence Embedding	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_A$	$E_B$	$E_B$	$E_B$	$E_B$	$E_B$
Transformer Positional Embedding	$E_0$	$E_1$	$E_2$	$E_3$	$E_4$	$E_5$	$E_6$	$E_7$	$E_8$	$E_9$	$E_{10}$

## Training

- Masked Language Modelling
  - Mask out k% of the input words, and then predict the masked words
  - the man went to the store to [MASK] a [MASK] of milk
  - What can you use as a loss function?
- Next sentence prediction
  - To learn relationships between sentences, predict whether Sentence B is actual sentence that proceeds Sentence A, or a random sentence

**Sentence A** = The man went to the store.  
**Sentence B** = He bought a gallon of milk.  
**Label** = IsNextSentence

**Sentence A** = The man went to the store.  
**Sentence B** = Penguins are flightless.  
**Label** = NotNextSentence

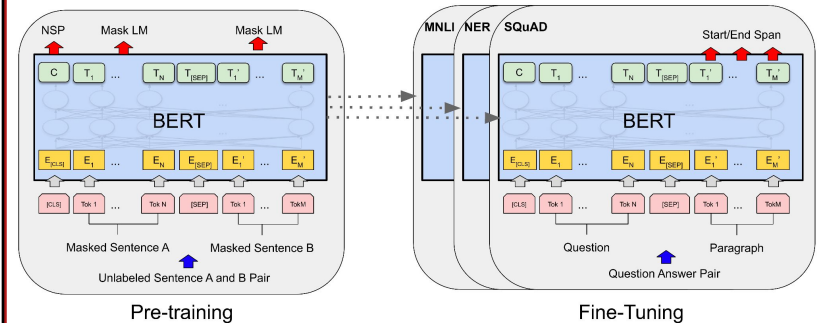
## Model Details

- Data: Wikipedia (2.5B words) + BookCorpus (800M words)
- Batch Size: 131,072 words (1024 sequences \* 128 length or 256 sequences \* 512 length)
- Training Time: 1M steps (~40 epochs)
- Optimizer: AdamW, 1e-4 learning rate, linear decay
- BERT-Base: 12-layer, 768-hidden, 12-head, 110M params
- BERT-Large: 24-layer, 1024-hidden, 16-head, 340M params
- Trained on 4x4 or 8x8 TPU slice for 4 days

## Demo

<https://huggingface.co/google-bert/bert-large-cased?text=Paris+is+the+capital+of+%5BMASK%5D.>

## Pre-training to Fine-tuning Pipeline



## Cornell Bowers CIS

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

<https://arxiv.org/pdf/1810.04805.pdf>

## Cornell Bowers CIS

### Self-supervised Learning

- Labels are generated automatically, no human labeling process
- Benefits
  - Scales well
  - Cost-Efficient
  - Flexible
- Challenges
  - Larger datasets are required
  - More compute is necessary

## Cornell Bowers CIS

### Review

- LSTMs/GRUs are recurrent
- Self-attention can effectively replace recurrence in sequence-to-sequence models
- Transformers use self-attention and are parallelizable
- Pre-training using self-supervised learning help train large models that learn very good representations