

# Deep NLP 3: Transformers and Attention

Oliver Guhr

Two large, solid orange parallelogram shapes are positioned at the bottom of the slide, overlapping each other. The one on the left is smaller and tilted more steeply, while the one on the right is larger and tilted less steeply, creating a sense of depth and modern design.

# Identify offensive language using Transformers



# Quiz Time!

# How is the pace of this course?

- A) too slow
- B) just right
- C) too fast

# What is unsupervised learning?

- A) compressing sparse into dense vectors
- B) learning based on example input-output pairs
- C) a different name for transfer learning
- D) an algorithm that learns patterns from untagged data

# What does the Distributional Hypothesis say?

- A) Words can be encoded in a vector space
- B) Words are described by their context words
- C) Words can be drawn on maps
- D) Similarity between words can be calculated using the euclidean distance

# Why do we need dense vector representations for texts?

- A) to efficiently compute neural networks
- B) to encode the relationships between words
- C) to create word clouds
- D) to pretrain neural networks

# Transfer learning for NLP works by:

- A) training a model with an unsupervised task and retraining it with labeled data
- B) pretraining a model with a labeled data and retraining it with an unsupervised task

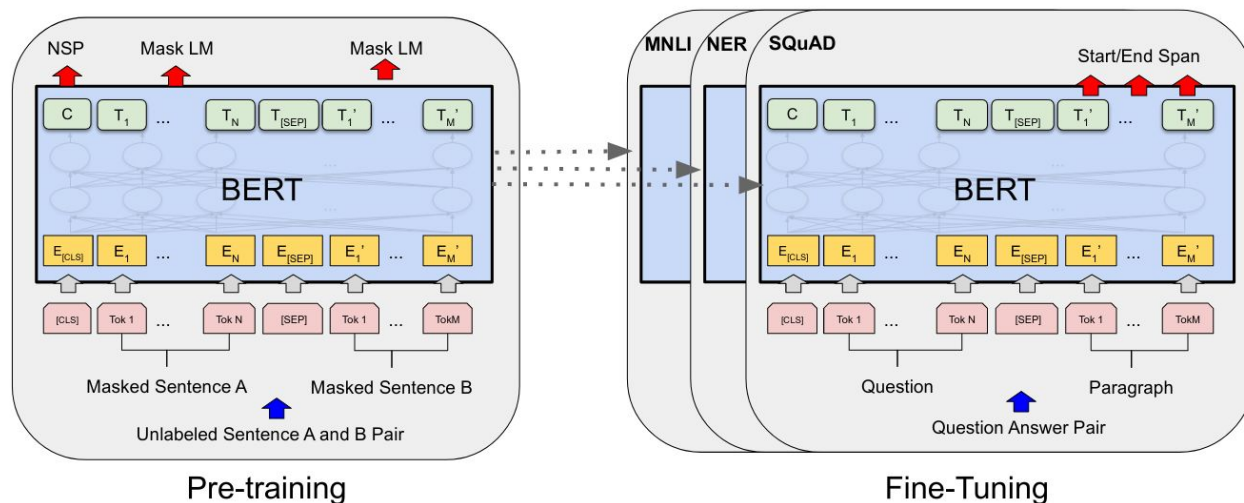


# Recap

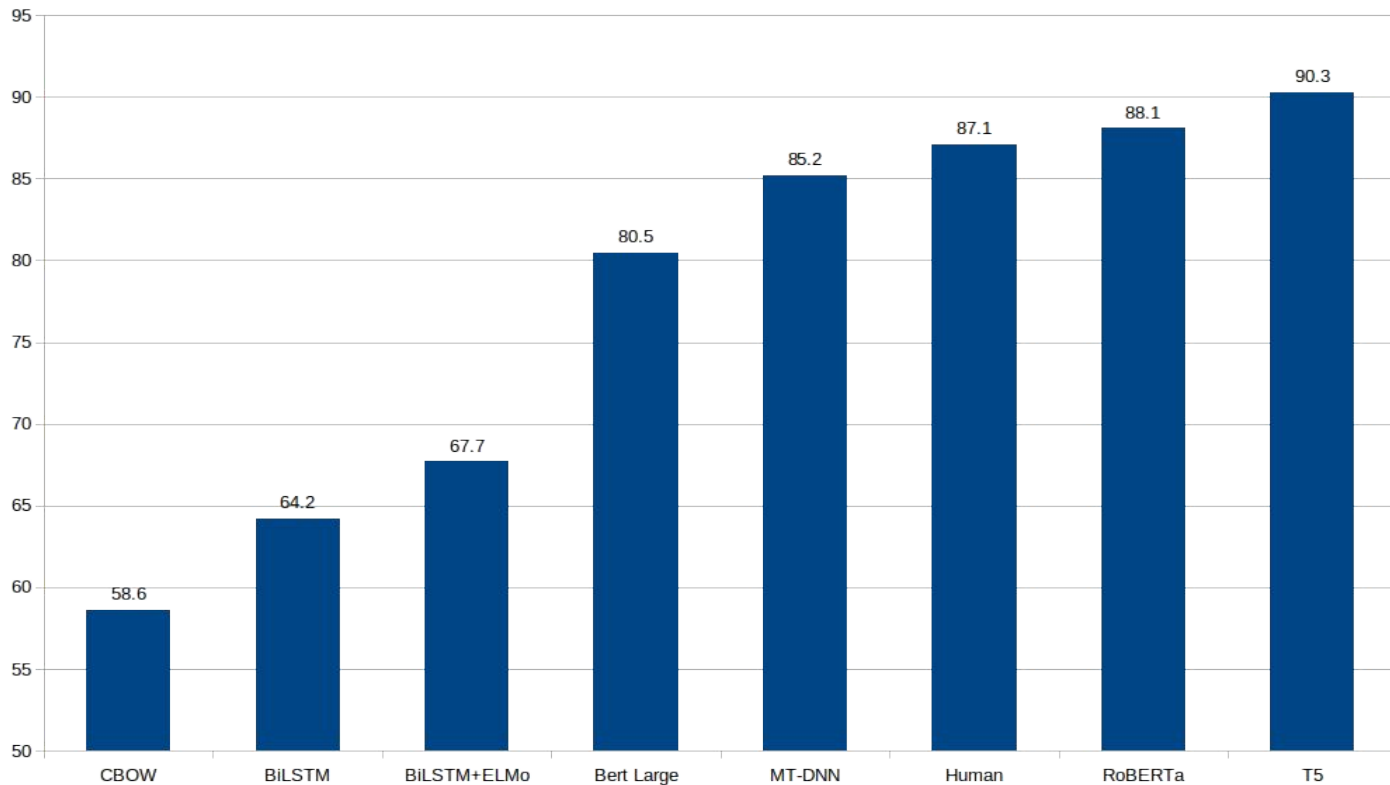
# Bert



- [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#)
- Paper by Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova
- Published in 2018
- improved the state-of-the-art in most important benchmarks

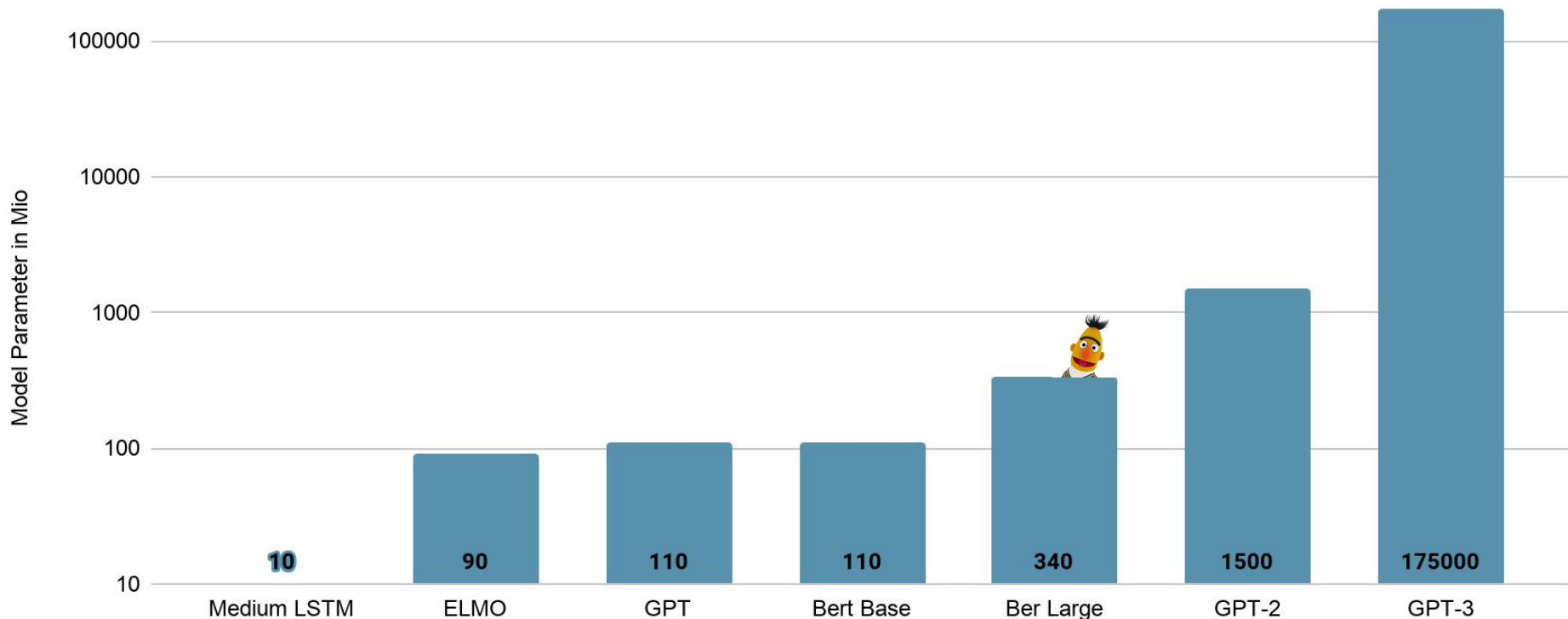


# GLUE Benchmark



GLUE Leaderboard: <https://gluebenchmark.com/leaderboard>

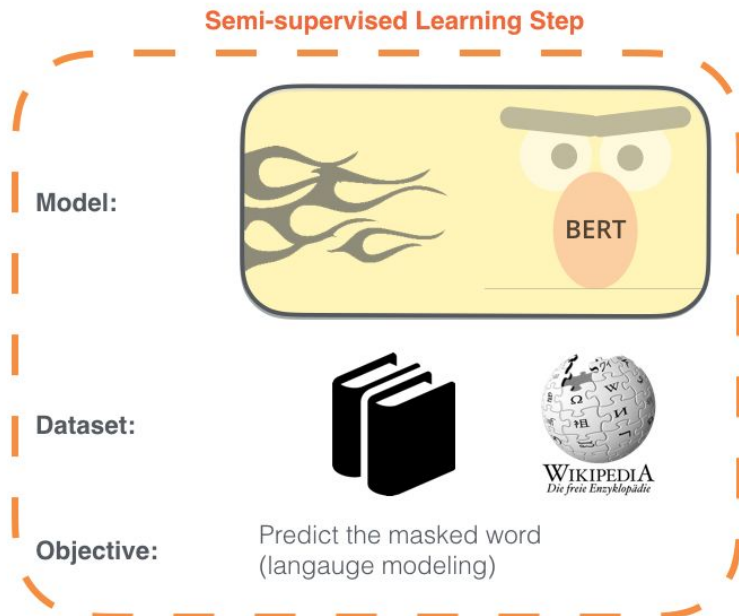
# How deep are these models?



# Bert

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



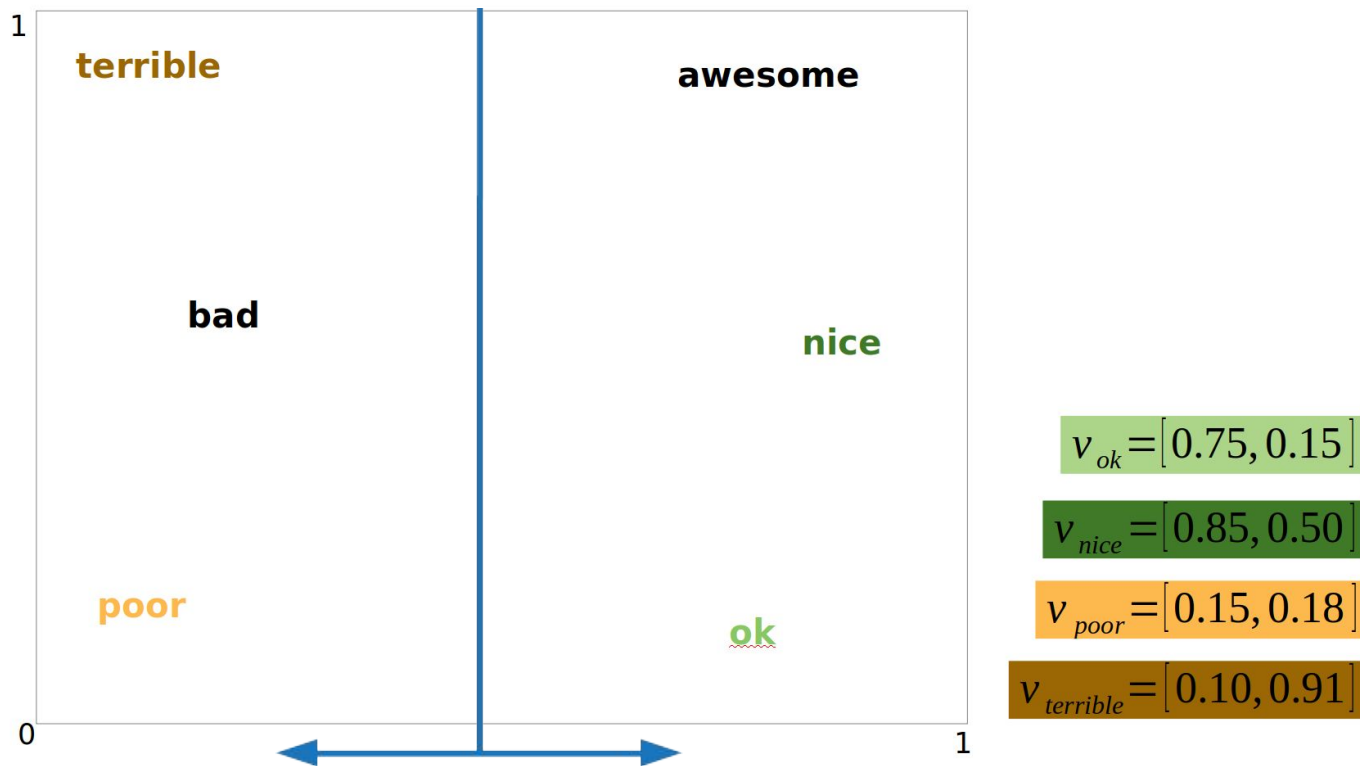
Words that occur in the same contexts tend to have similar meanings.

Harris (1954)

A word is characterized by the company it keeps.

Firth (1957)

# Word Vectors - Klassifikation



# Task One: Mask Words

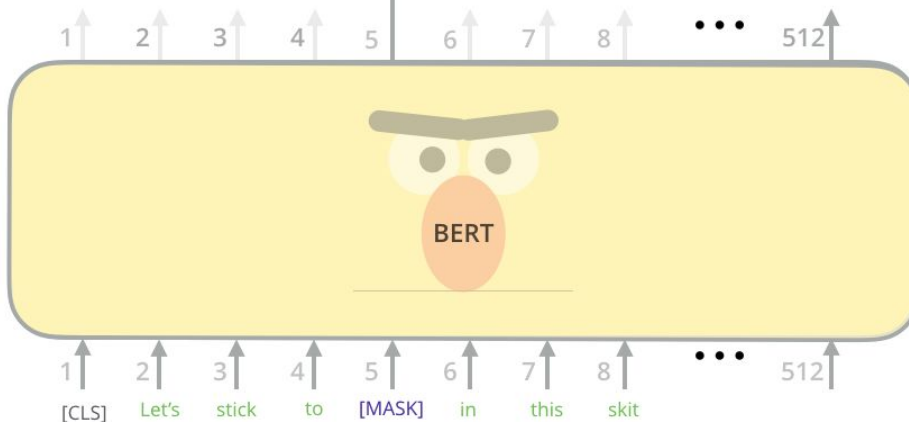


Use the output of the masked word's position to predict the masked word

Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva

FFNN + Softmax



Randomly mask  
15% of tokens

Input

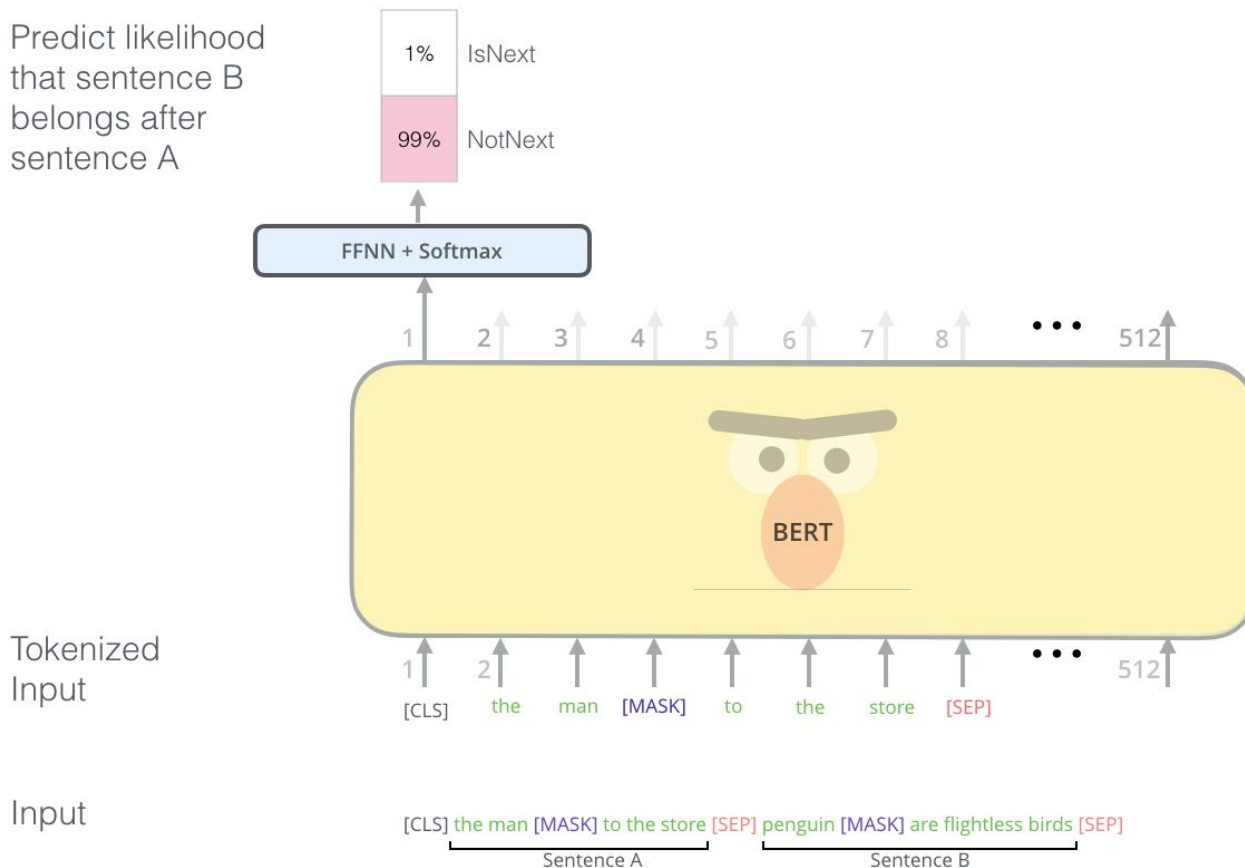
[CLS] Let's stick to improvisation in this skit



# Task Two: Next Sentence Prediction



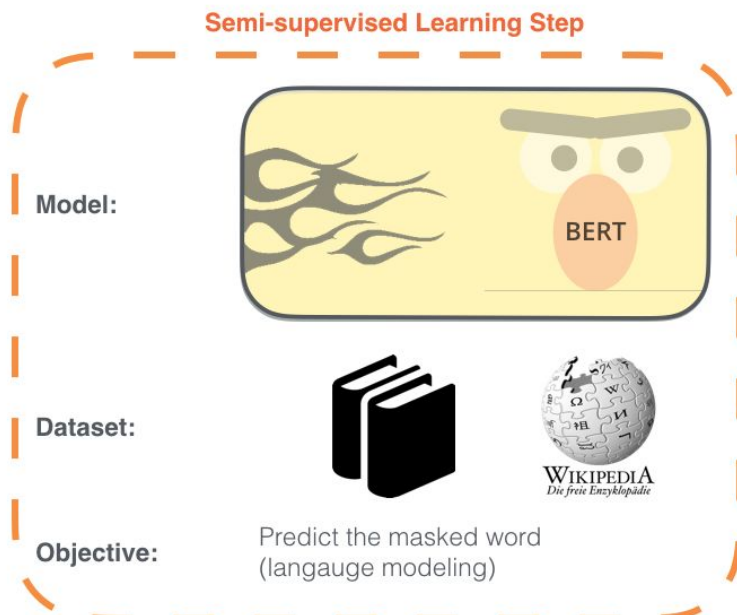
Predict likelihood  
that sentence B  
belongs after  
sentence A



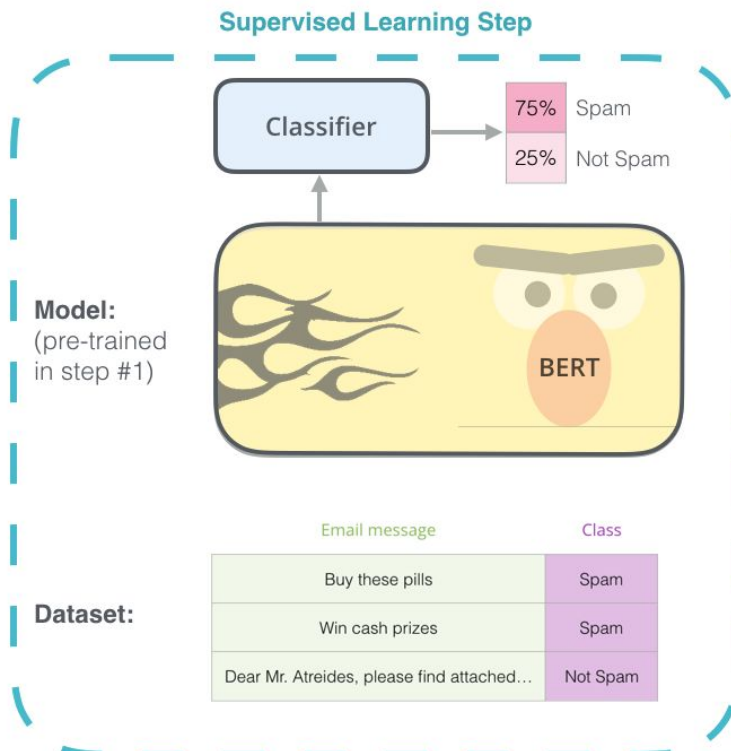
# Bert

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

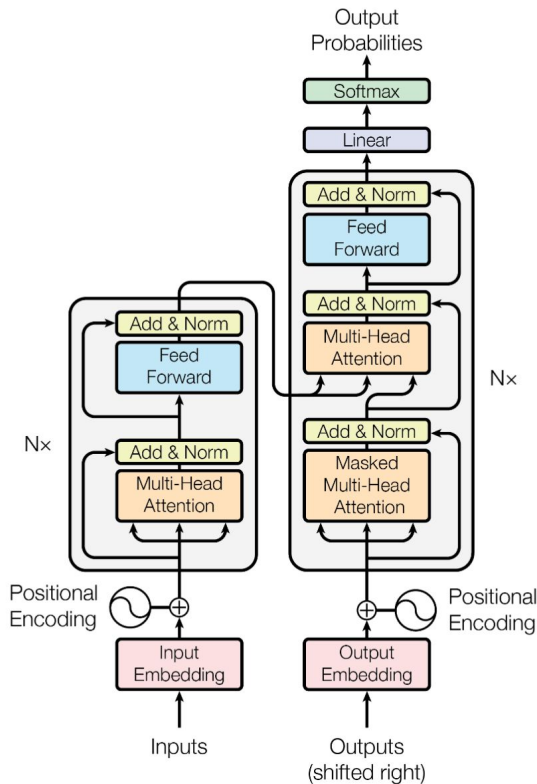


2 - **Supervised** training on a specific task with a labeled dataset.



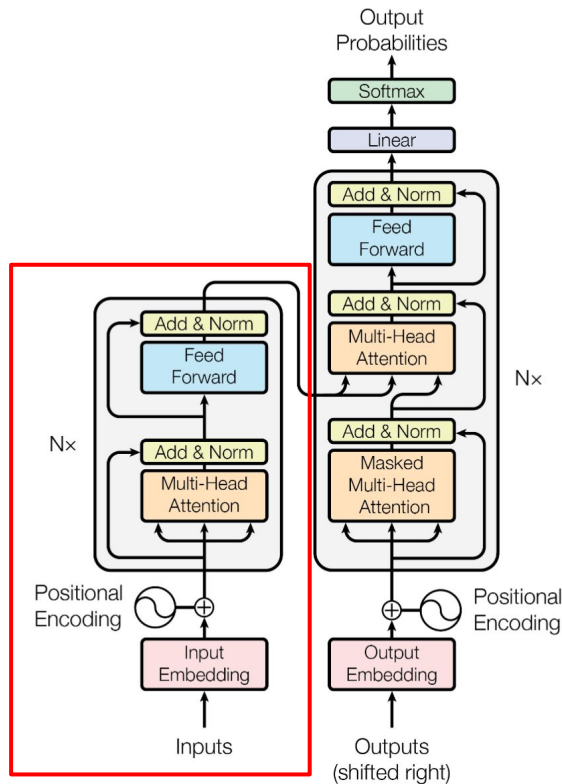
# How do Transformers work?

# Attention is all you need



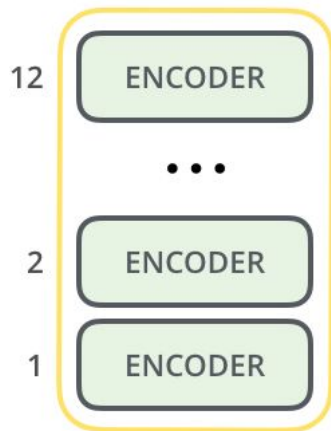
Attention Is All You Need, Vaswani et al. <https://arxiv.org/abs/1706.03762>

# Attention is all you need

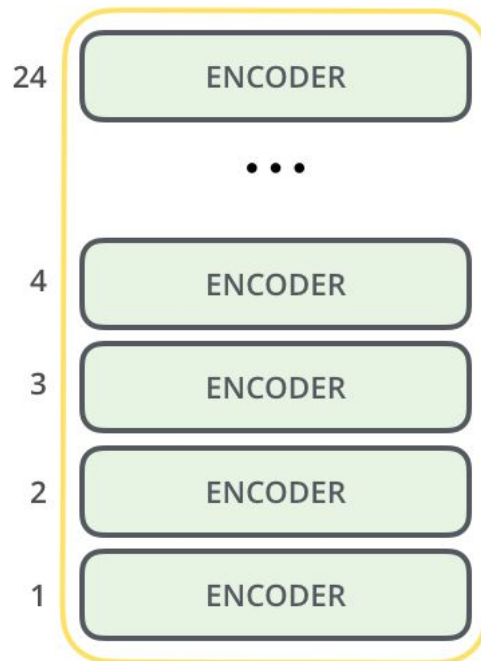


Attention Is All You Need, Vaswani et al. <https://arxiv.org/abs/1706.03762>

# How encoders work.

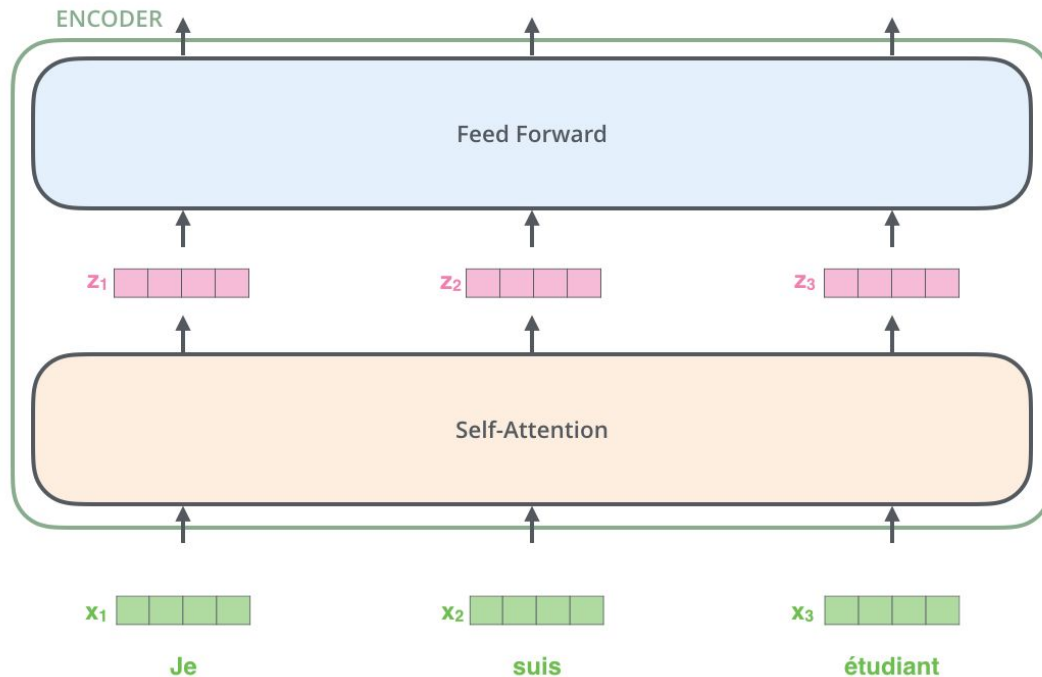


BERT<sub>BASE</sub>



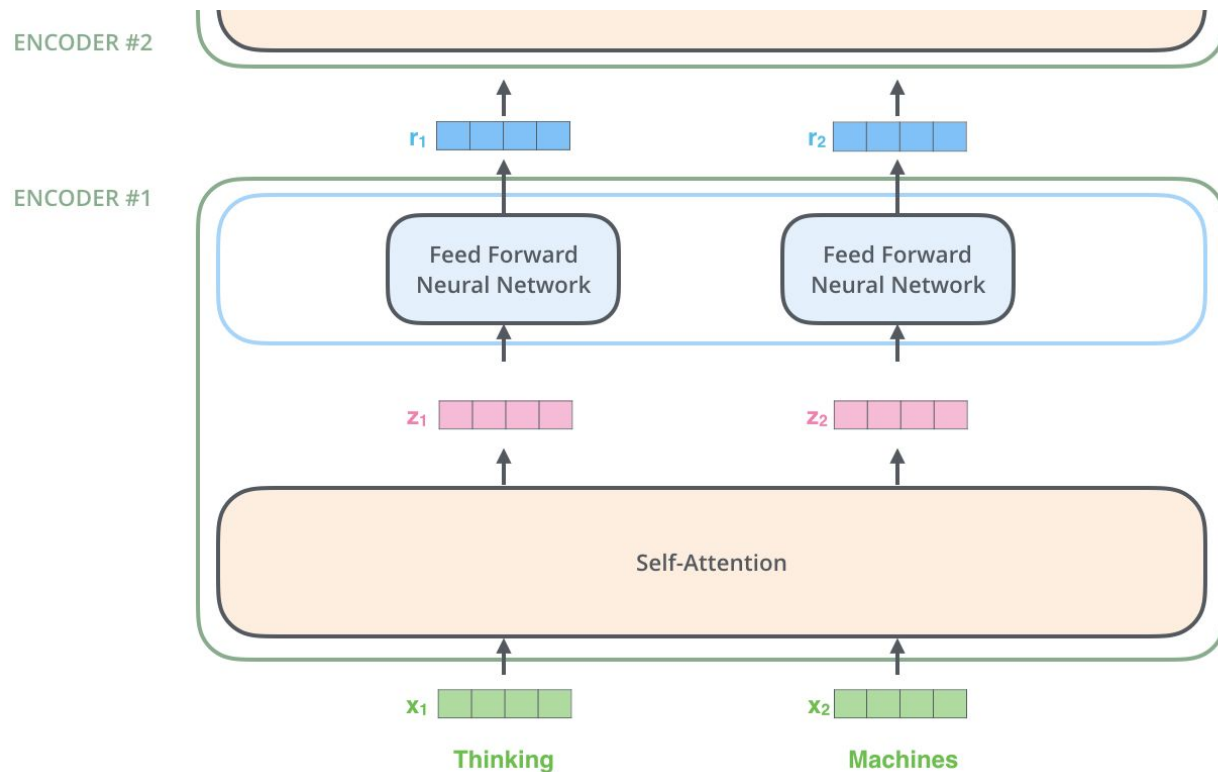
BERT<sub>LARGE</sub>

# Transformer Encoder



The Illustrated BERT, Jay Alammar: <http://http://jalammar.github.io/illustrated-transformer/>

# Transformer Encoder



The Illustrated BERT, Jay Alammar: <http://http://jalammar.github.io/illustrated-transformer/>



# What is self attention?

# Scaled dot product attention



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

Query Key Value

# Scaled dot product attention



$$\text{Attention}(\underbrace{Q}_{\text{Query}}, \underbrace{K}_{\text{Key}}, \underbrace{V}_{\text{Value}}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Take the current **word or token**, find the most similar **key** and return the corresponding **value**.



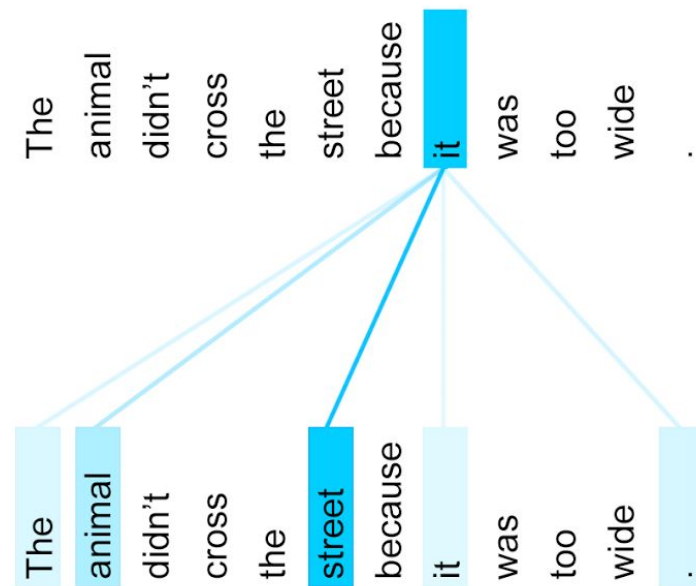
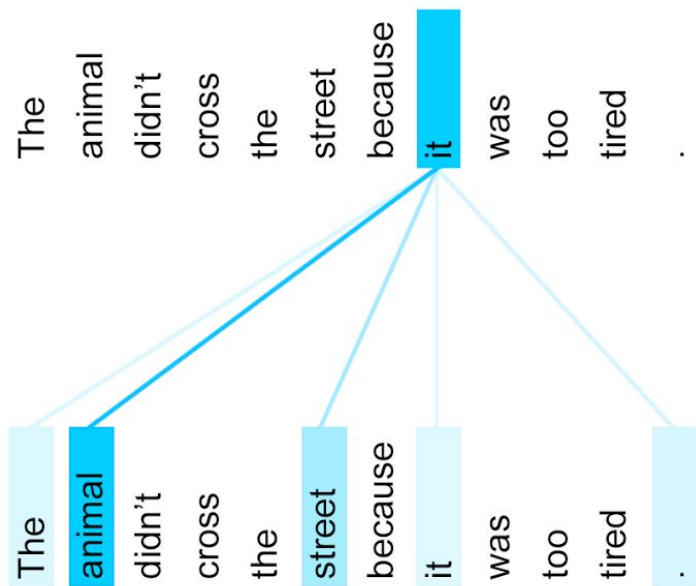
# What does Attention do?



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

The  
animal  
didn't  
cross  
the  
street  
because  
it  
was  
too  
wide  
.

# What does Attention do?



The encoder self-attention distribution for the word “it” from the 5th to the 6th layer of a Transformer trained on English to French translation (one of eight attention heads).

Source: <https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

# Attention

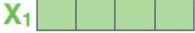


Input

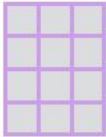
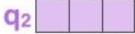
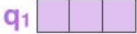
Thinking

Machines

Embedding

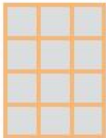
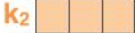
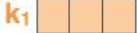


Queries



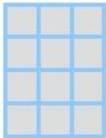
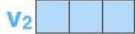
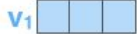
$W^Q$

Keys



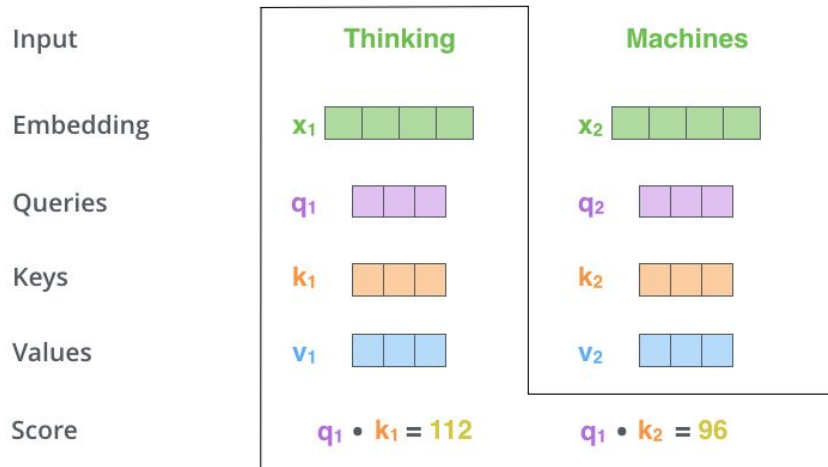
$W^K$

Values



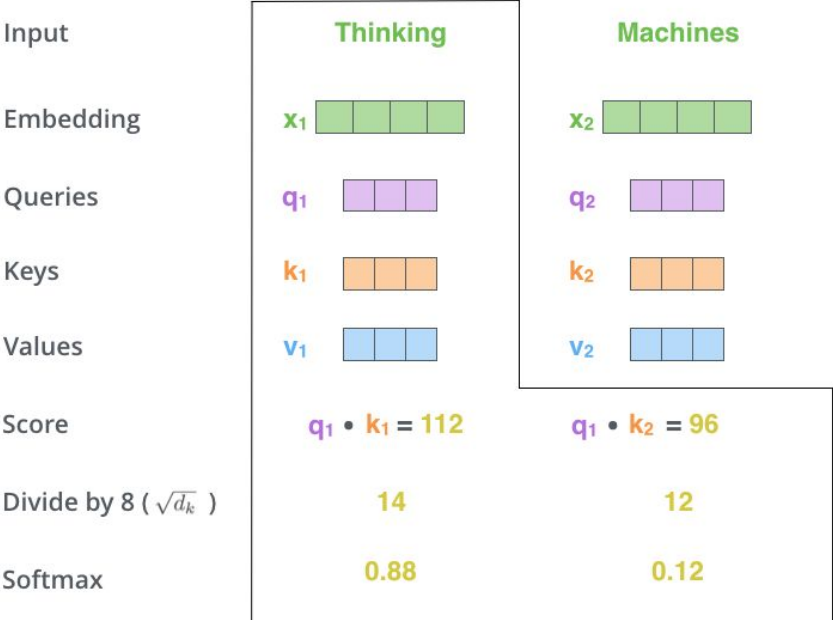
$W^V$

# Attention

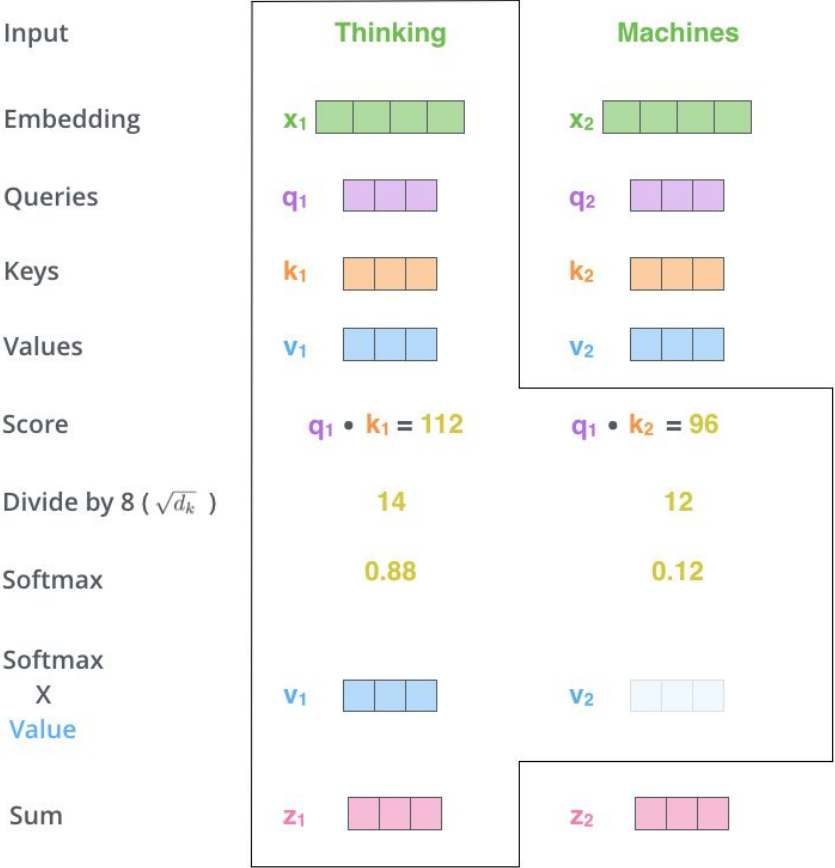




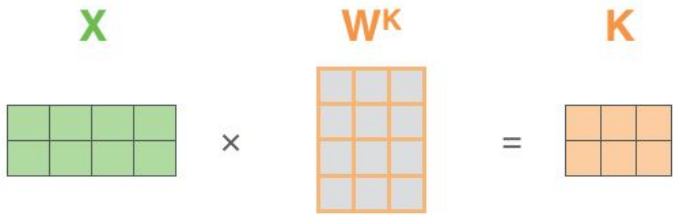
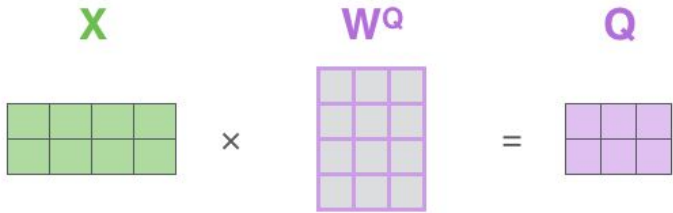
# Attention



# Attention



# Matrix Calculation



# Matrix Calculation



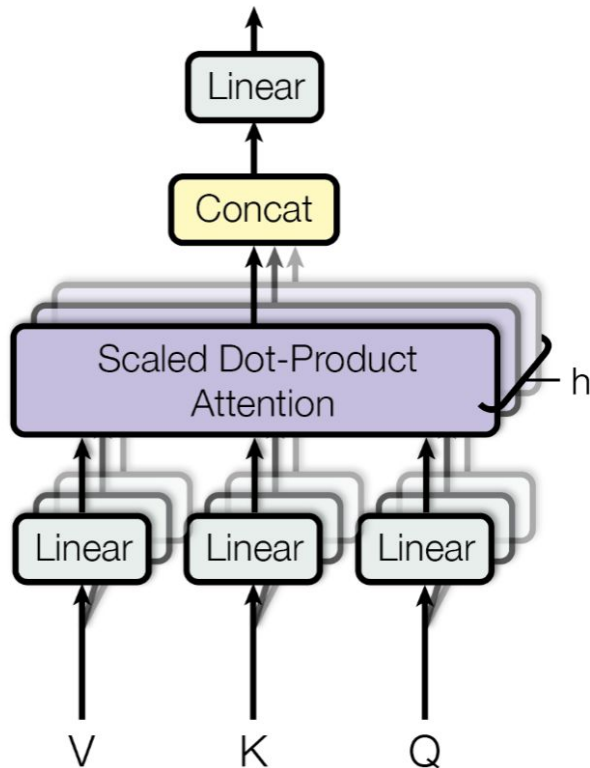
$$X \times W^Q = Q$$

$$X \times W^K = K$$

$$X \times W^V = V$$

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

# Multi Head Attention



Attention Is All You Need, Vaswani et al. <https://arxiv.org/abs/1706.03762>

# Multi Head Attention

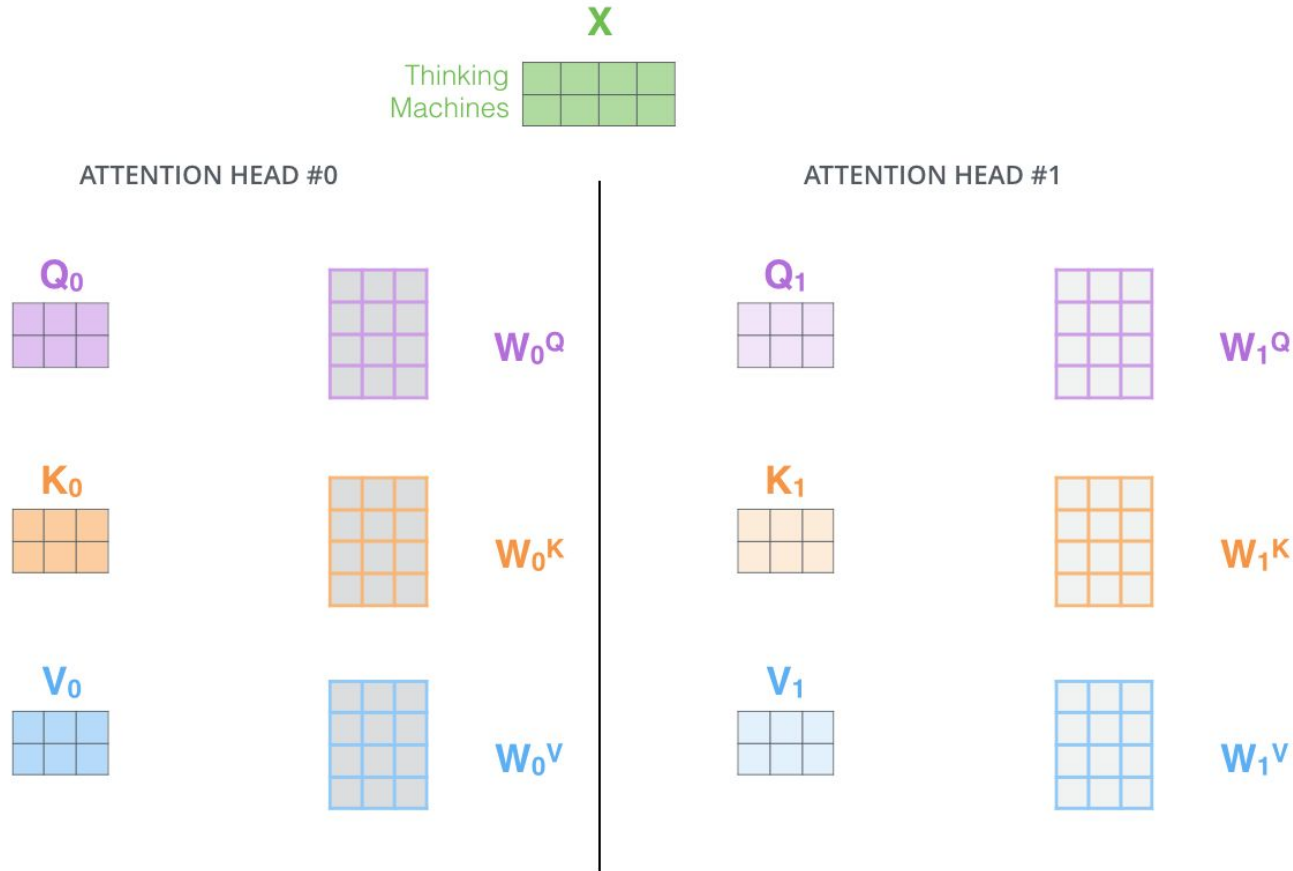


$$X \times W^Q = Q$$
A diagram illustrating the calculation of the Query matrix Q. It shows a green 2x4 matrix X multiplied by a purple 4x4 matrix W<sup>Q</sup> to produce a purple 2x4 matrix Q. The matrices are represented as grids of colored squares.

$$X \times W^K = K$$
A diagram illustrating the calculation of the Key matrix K. It shows a green 2x4 matrix X multiplied by an orange 4x4 matrix W<sup>K</sup> to produce an orange 2x4 matrix K. The matrices are represented as grids of colored squares.

$$X \times W^V = V$$
A diagram illustrating the calculation of the Value matrix V. It shows a green 2x4 matrix X multiplied by a blue 4x4 matrix W<sup>V</sup> to produce a blue 2x4 matrix V. The matrices are represented as grids of colored squares.

# Multi Head Attention



# Multi Head Attention

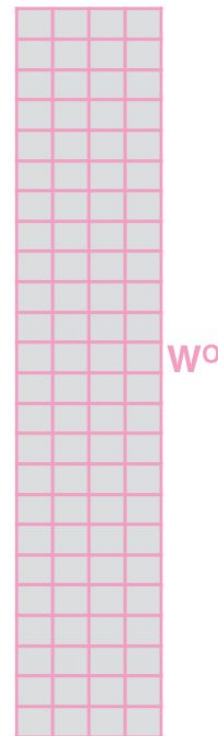


1) Concatenate all the attention heads



2) Multiply with a weight matrix  $W^O$  that was trained jointly with the model

$\times$



3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN





# Multi Head Attention



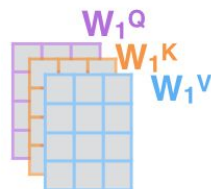
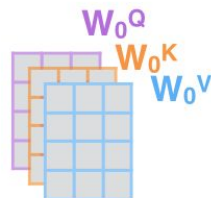
1) This is our input sentence\*

Thinking  
Machines

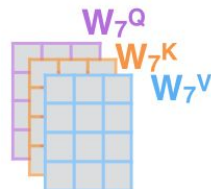
2) We embed each word\*



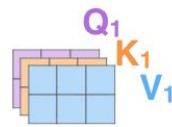
3) Split into 8 heads.  
We multiply  $X$  or  $R$  with weight matrices



...



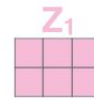
4) Calculate attention using the resulting  $Q/K/V$  matrices



...



5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer



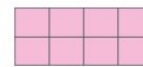
...



$W^O$

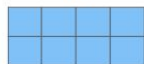


$Z$

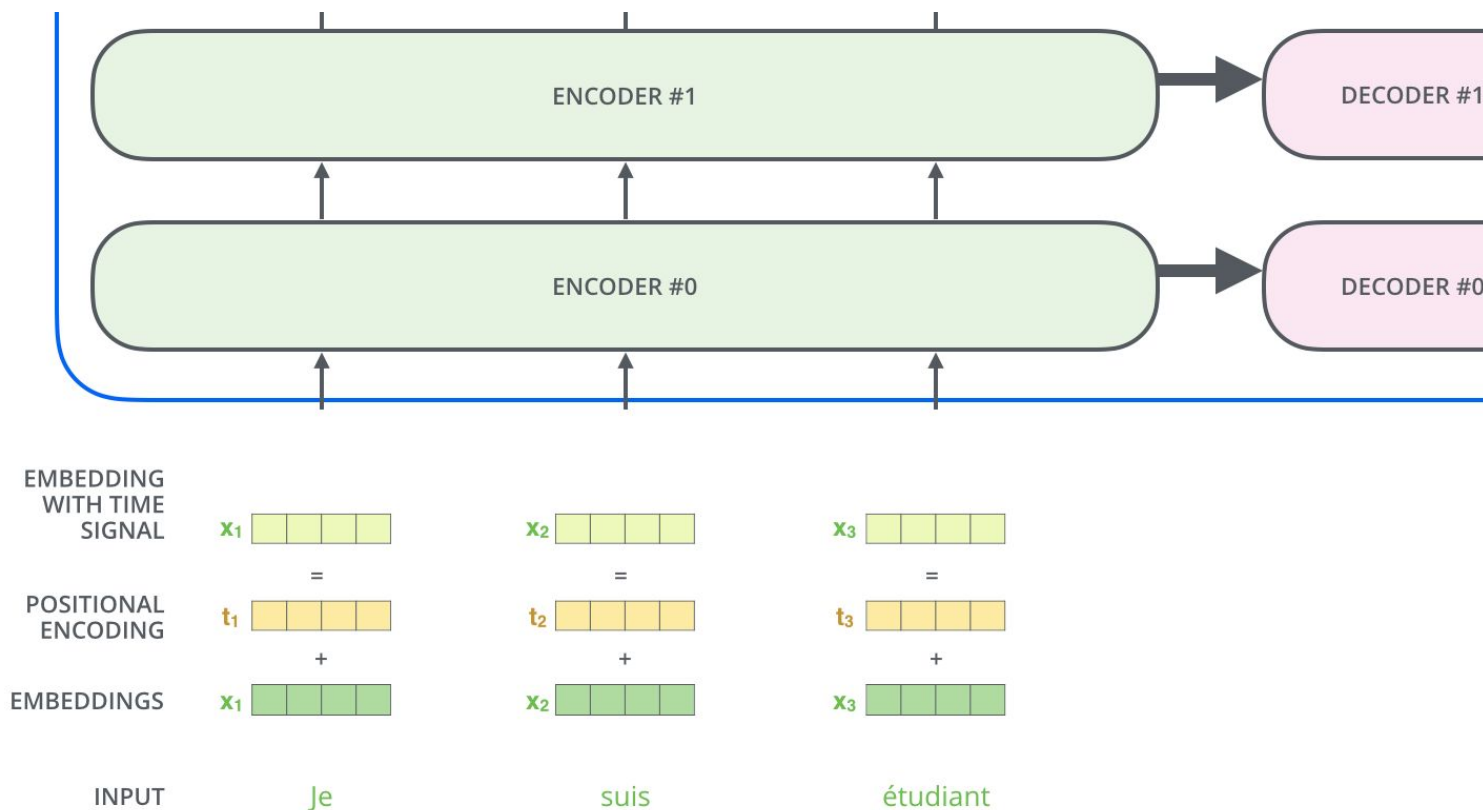


\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

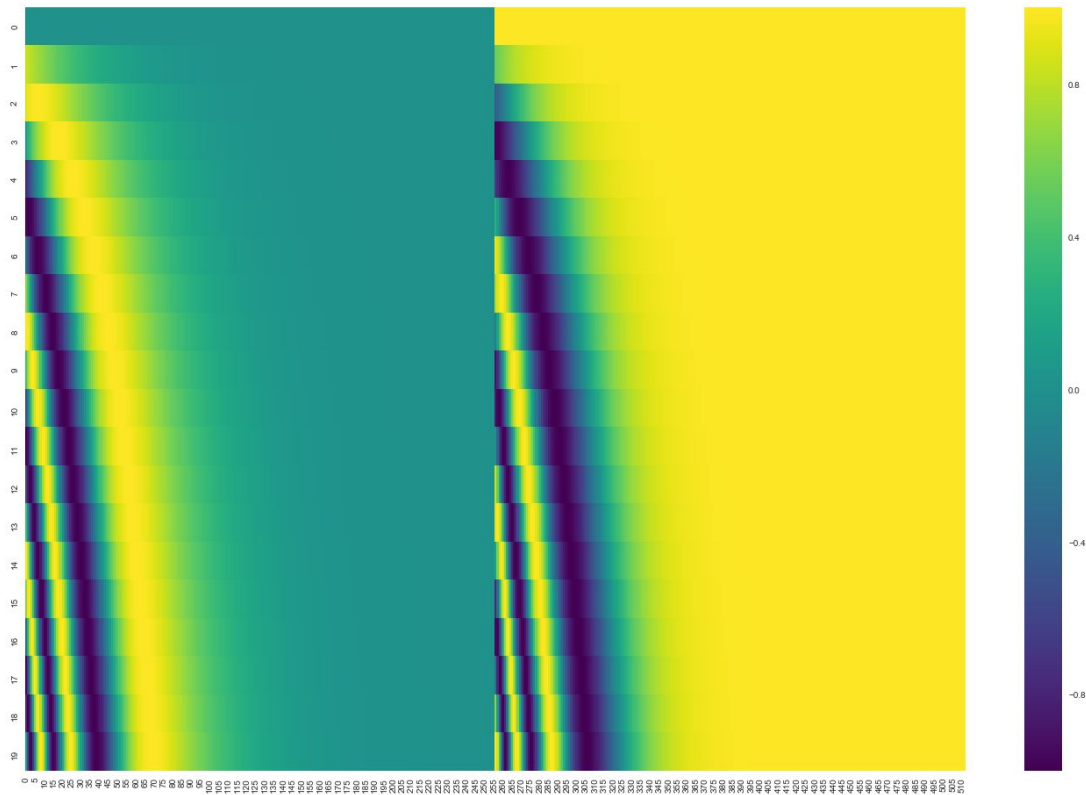
$R$



# Positional Encoding



# Positional Encoding



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

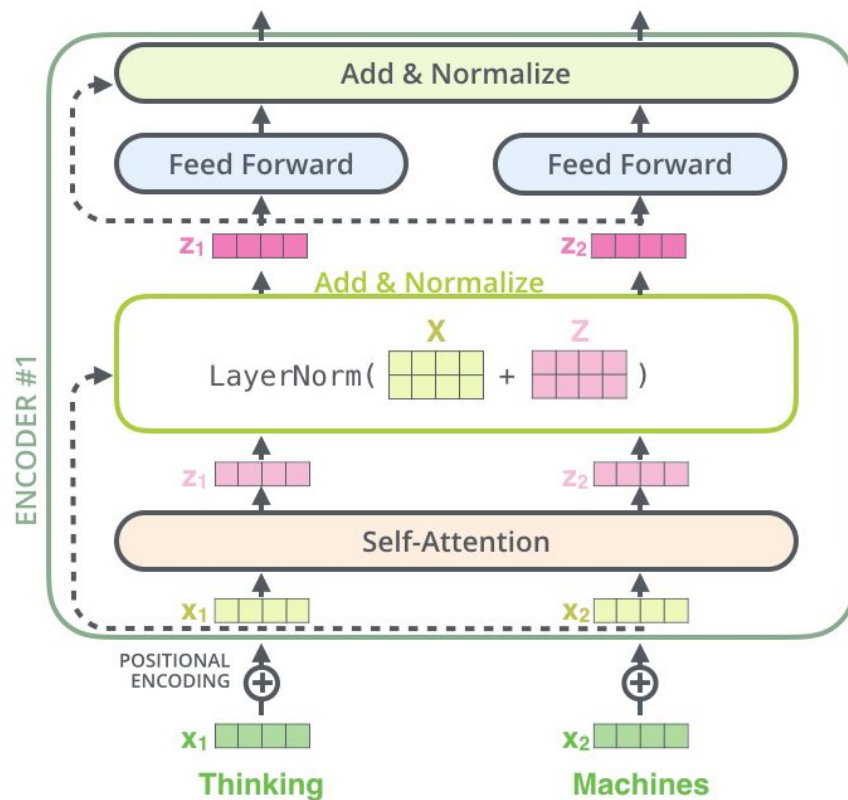
# Positional Encoding



For embedding with a dimensionality of 4 the encodings look like this:



# Add and Normalize

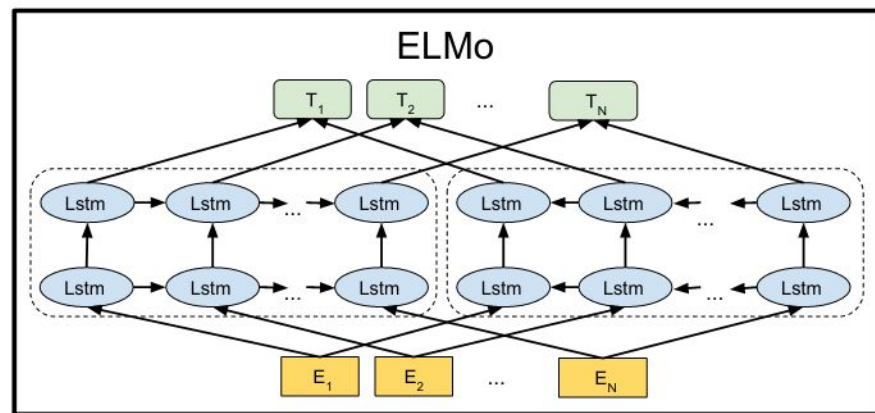
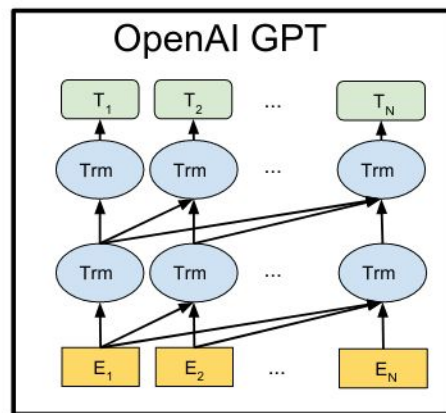
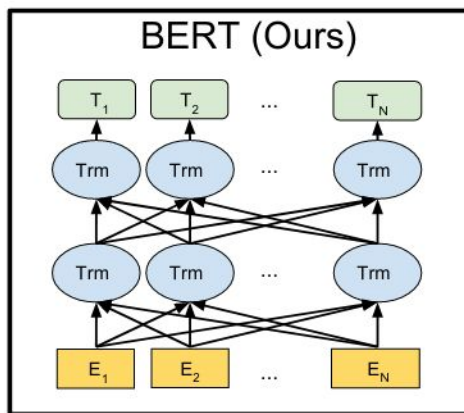


Layer Normalization Lei Ba et al. <https://arxiv.org/abs/1607.06450>

# Transformers vs LSTMs



- Can we build something similar using LSTMs?
  - Yes, its called ELMo



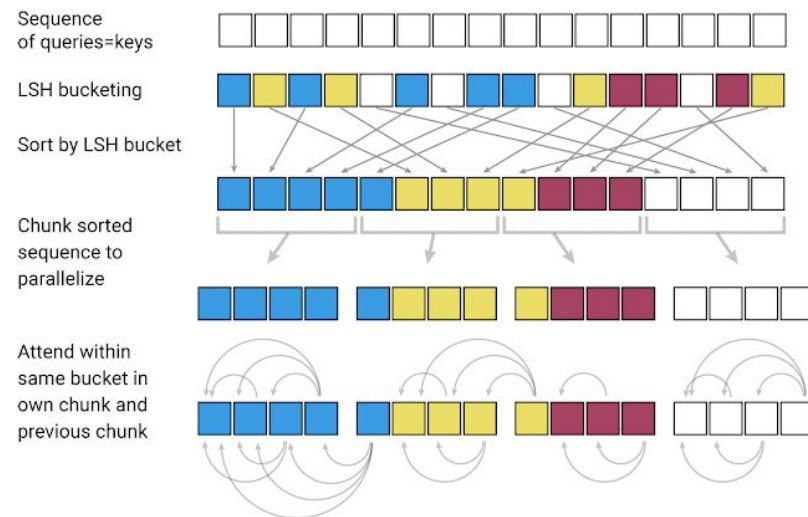
Source Bert Paper: <https://arxiv.org/pdf/1810.04805.pdf>

# Future...

# Reformer: The Efficient Transformer



- Improved efficiency of the attention algorithm
- **context windows of 1 million words** on a 16GB GPU (Transformer 512 Token)
- Main Contribution
  - locality-sensitive-hashing (LSH)
  - reversible residual layers
- **Similar ideas:**
  - Longformer, Linformer,  $[\text{w}^*]$ former
- More Information
  - [Paper by Kitaev, Kaiser and Levskya](#)
  - [Google AI Blog Post](#)
  - [Video Introduction](#)
  - [Background Info](#)





# RealFormer: Transformer Likes Residual Attention



- **Resnets idea** but for Transformers: Residual connections for attention values
- Improves overall results but not by much
- [Paper by Ruining He, Anirudh Ravula, Bhargav Kanagal, Joshua Ainslie](#)

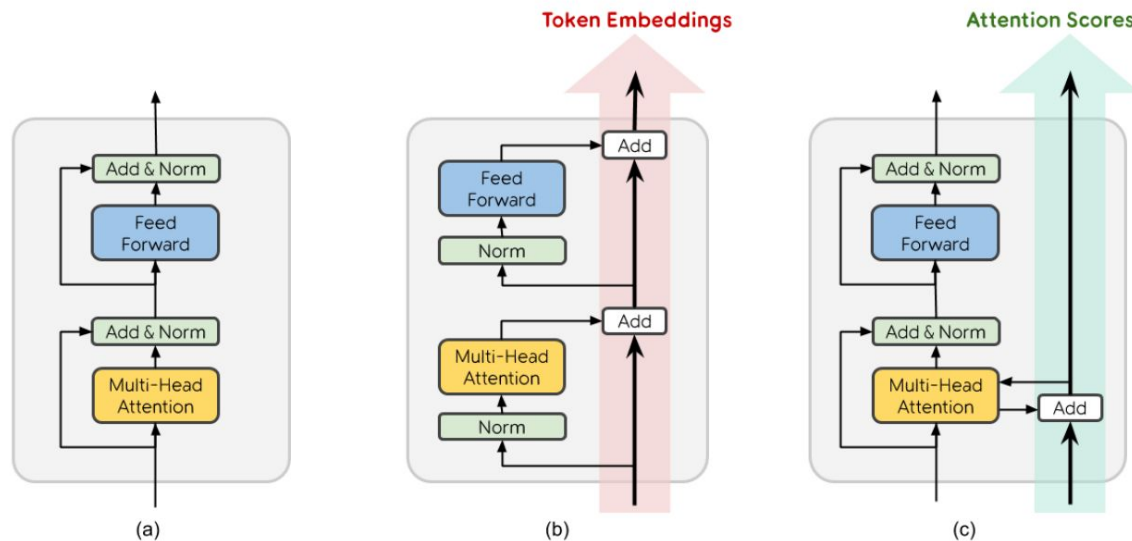
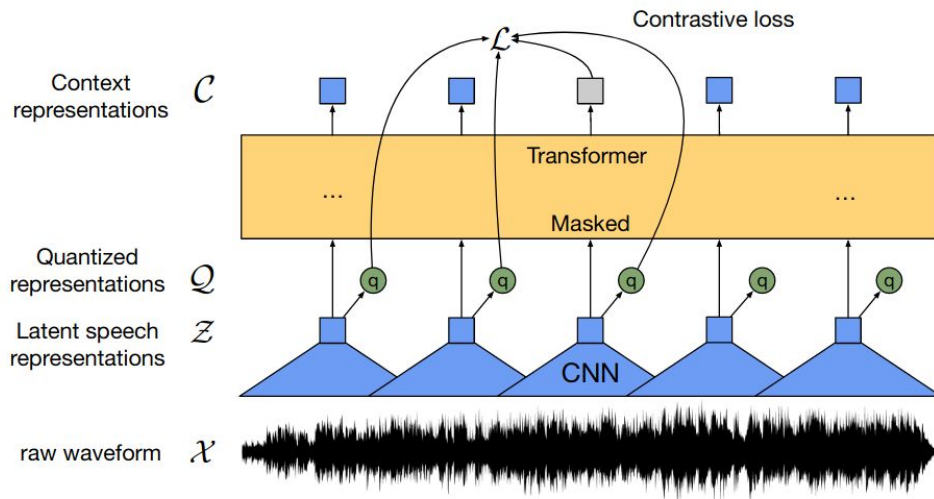


Figure 1: Comparison of different Transformer layers: (a) The prevalent Post-LN layer used by (e.g.) BERT; (b) Pre-LN layer used by (e.g.) GPT-2 that creates a “direct” path to propagate token embeddings; (c) Our RealFormer layer that creates a “direct” path to propagate attention scores (by adding a simple skip edge on top of (a)).

# Automatic Speech Recognition

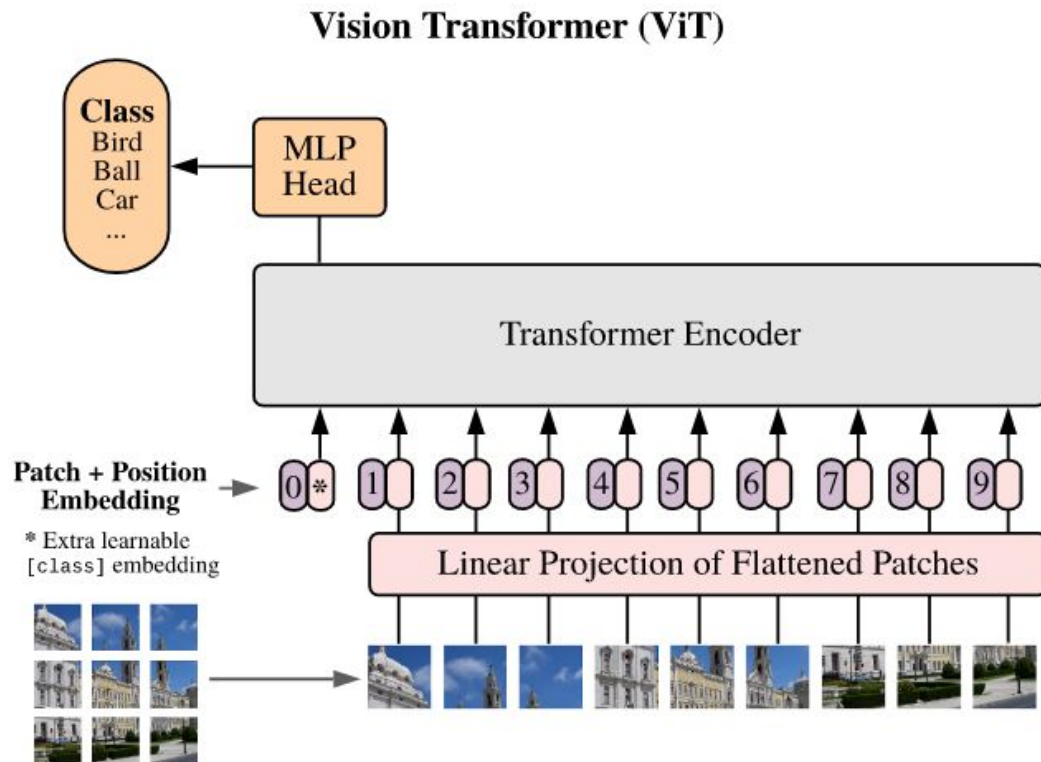


- wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations
- Key Ideas:
  - **CNN and Transformer based end to end model for speech recognition**
  - uses a novel **pretraining schema to learn for unlabeled audio data**
- outperforms the previous state of the art while using 100 times less labeled data
- can achieve good accuracy with very little data
- By Alexei Baevski, Henry Zhou, Abdelrahman Mohamed and Michael Auli



# An Image is Worth 16x16 Words

- Imagenet and CIFAR with transformers
  - 88.55% on ImageNet,
  - 90.72% on ImageNet-Real,
  - 94.55% on CIFAR-100
- Paper by [Dosovitskiy et al.](#)
- Other approaches to vision tasks
  - [Taming Transformers for High-Resolution Image Synthesis](#)



# Sources

- Paper
  - [Attention is all you need. Vaswani et al.](#)
  - [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding Devlin et al.](#)
  - [Reformer: The Efficient Transformer Kitev et al.](#)
- Good Read
  - [Jay Alammars The Illustrated Transformer](#)
  - [Jay Alammars The Illustrated BERT](#)
- Conference Talk:
  - [Attention is all you need attentional neural network models by Łukasz Kaiser](#)