# Deep NLP 3: Transformers and Attention

Oliver Guhr

# 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

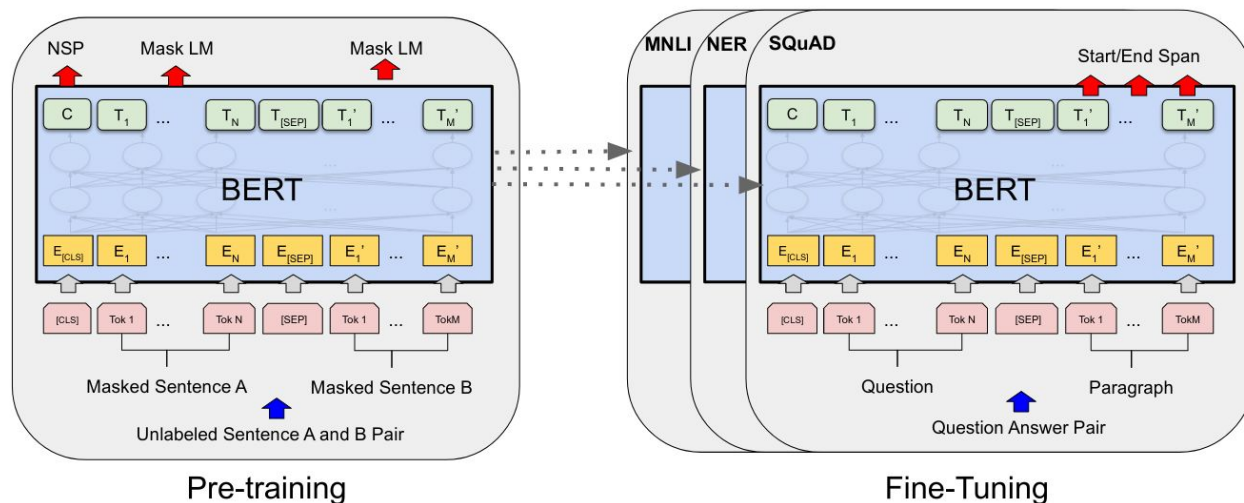B) pretraining a model with a labeled data and retraining it with an unsupervised task
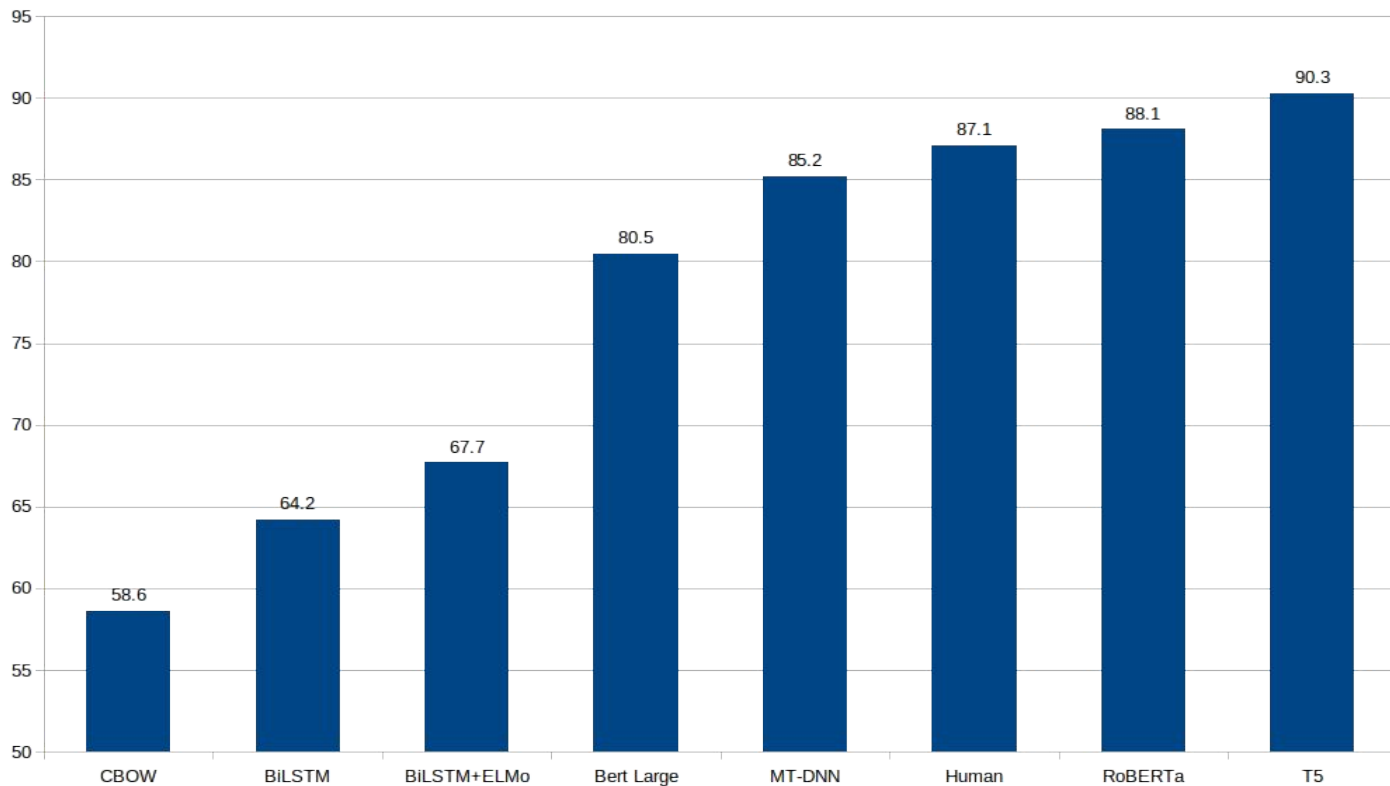
# Recap

# 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
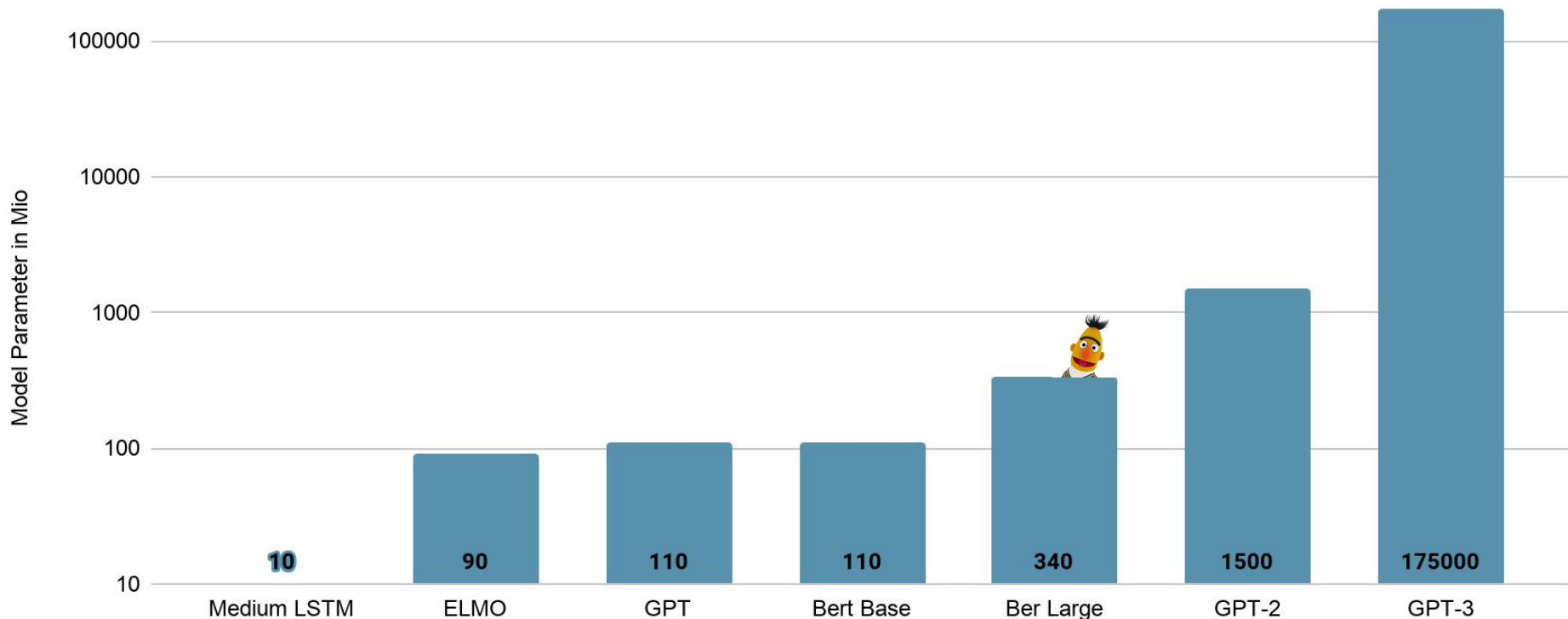- improved the state-of-the-art in most important benchmarks
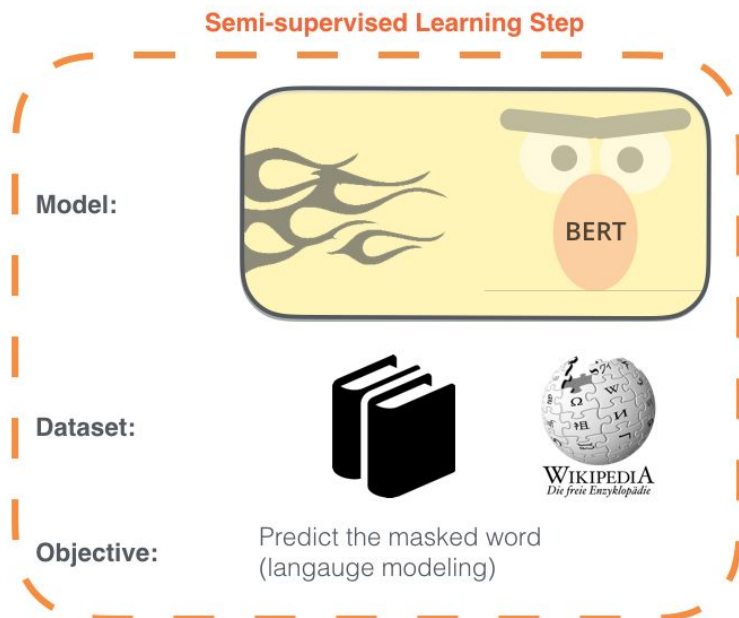
# GLUE Benchmark



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.
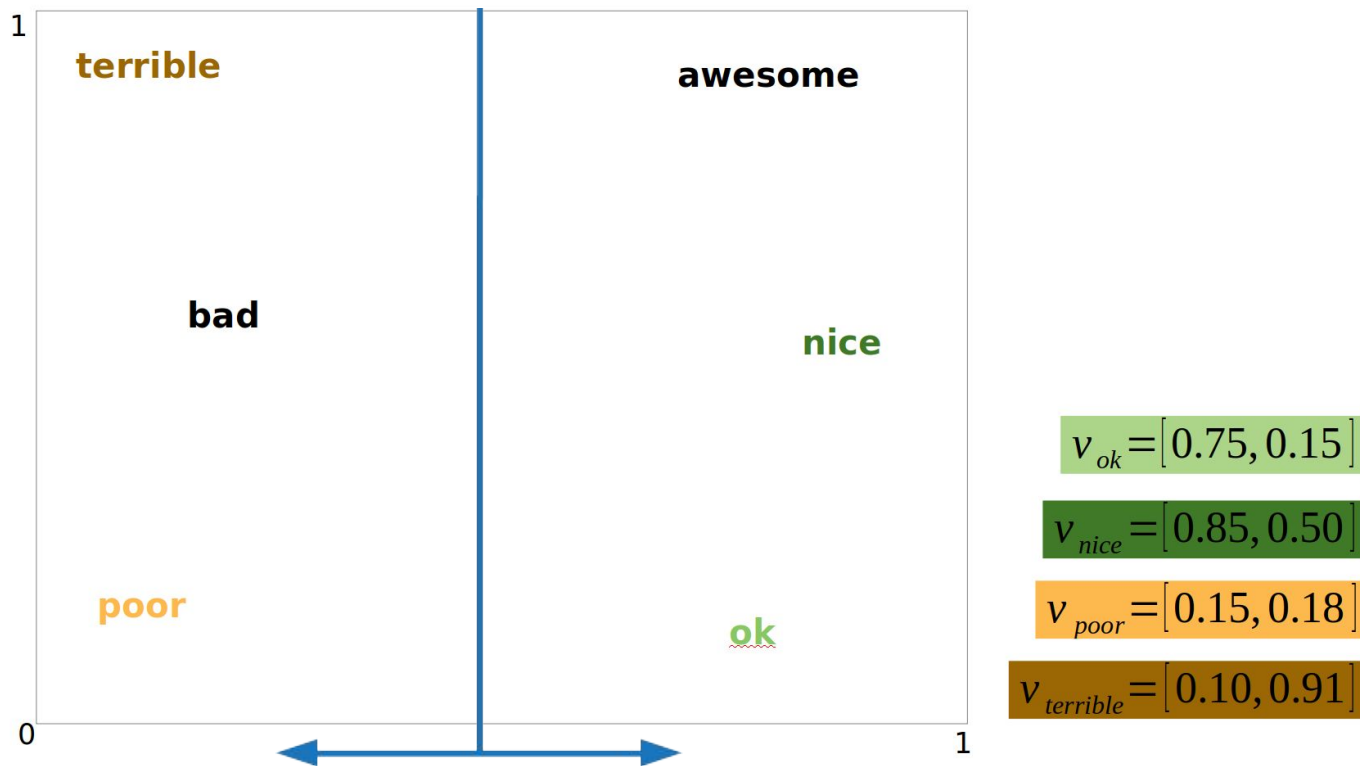


**Semi-supervised Learning Step**

**Model:** BERT

**Dataset:**

WIKIPEDIA
Die freie Enzyklopädie

**Objective:** Predict the masked word (langauge modeling)

# Distributional Hypothesis

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

terrible

awesome

bad

nice

$v_{ok} = [0.75, 0.15]$

$v_{nice} = [0.85, 0.50]$

poor

ok

$v_{poor} = [0.15, 0.18]$

$v_{terrible} = [0.10, 0.91]$

# 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% | Zyzzyva |

FFNN + Softmax

1 2 3 4 5 6 7 8 ... 512

BERT

Randomly mask
15% of tokens

1 2 3 4 5 6 7 8 ... 512

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

Input

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

# Task Two: Next Sentence Prediction

Predict likelihood that sentence B belongs after sentence A

1% IsNext

99% NotNext

FFNN + Softmax

1 2 3 4 5 6 7 8 ... 512

BERT

Tokenized Input

1 2 ... 512

[CLS] the man [MASK] to the store [SEP]

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]

Sentence A          Sentence B

# 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.

**Semi-supervised Learning Step**

**Model:**



BERT

**Dataset:**



WIKIPEDIA
*Die freie Enzyklopädie*

**Objective:** Predict the masked word (langauge modeling)

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

**Supervised Learning Step**

| Classifier | → | 75% Spam |
| | | 25% Not Spam |

**Model:**
(pre-trained in step #1)



BERT

**Dataset:**

| Email message | Class |
| --- | --- |
| Buy these pills | Spam |
| Win cash prizes | Spam |
| Dear Mr. Atreides, please find attached… | Not Spam |

# GPT
Generative Pretrained Transformer

# GPT - Pretraining



**Unsupervised Pre-training**

Input (features): a | robot | must

Correct output (label): obey

GPT-3 (under training)

Output (Prediction)

# Reinforcement Learning from Human Feedback (RLHF)

**Step One:**

Train a scoring model

**Prompts Dataset**

*Sample many prompts*

**Initial Language Model**

Lorem ipsum dolor
sit amet, consectet
adipiscing elit. Aen
Donec quam felis
vulputate eget, arc
Nam quam nunc
eros faucibus tincic
luctus pulvinar, her

**Generated text**

**Human Scoring**

**Train** on
{sample, reward} pairs

**Reward (Preference) Model**

text

$r_\theta$

**Outputs are ranked (relative, ELO, etc.)**

# Reinforcement Learning from Human Feedback (RLHF)

**Step Two:**

fine-tune the language model using the scoring model with RL

# 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_BASE

BERT_LARGE

# 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

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

**Query**  **Key**  **Value**

# Scaled dot product attention

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

**Query**  **Key**  **Value**

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

# Attention

# Attention

| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

# Attention



|  | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# Matrix Calculation

# Matrix Calculation

# Multi Head Attention



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

# Multi Head Attention

# Multi Head Attention



X

Thinking
Machines

ATTENTION HEAD #0

$Q_0$

$W_0^Q$

$K_0$

$W_0^K$

$V_0$

$W_0^V$

ATTENTION HEAD #1

$Q_1$

$W_1^Q$

$K_1$

$W_1^K$

$V_1$

$W_1^V$

# Multi Head Attention

HTW

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

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

X

$W^O$

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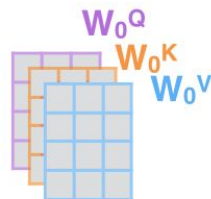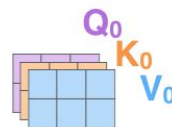
$Z$

=

1) This is our input sentence*

2) We embed each word*

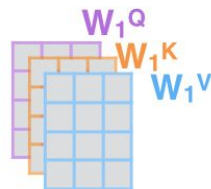3) Split into 8 heads. We multiply X or R with weight matrices

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

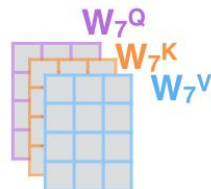5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer

X

Thinking Machines

$W_0^Q$
$W_0^K$
$W_0^V$

$Q_0$
$K_0$
$V_0$

$Z_0$

$W^O$

$W_1^Q$
$W_1^K$
$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

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
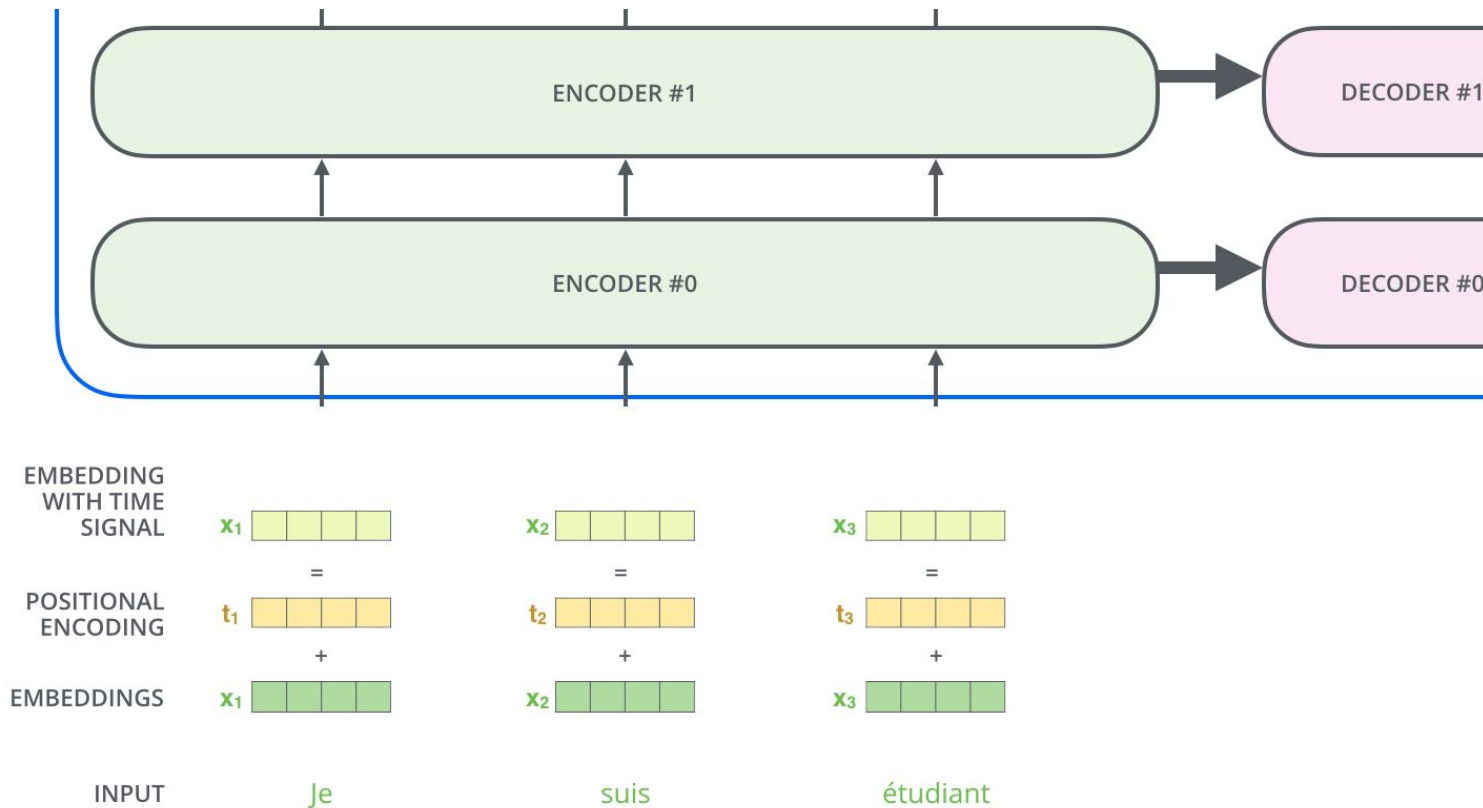
R

...

$W_7^Q$
$W_7^K$
$W_7^V$

...

$Q_7$
$K_7$
$V_7$

...

$Z_7$
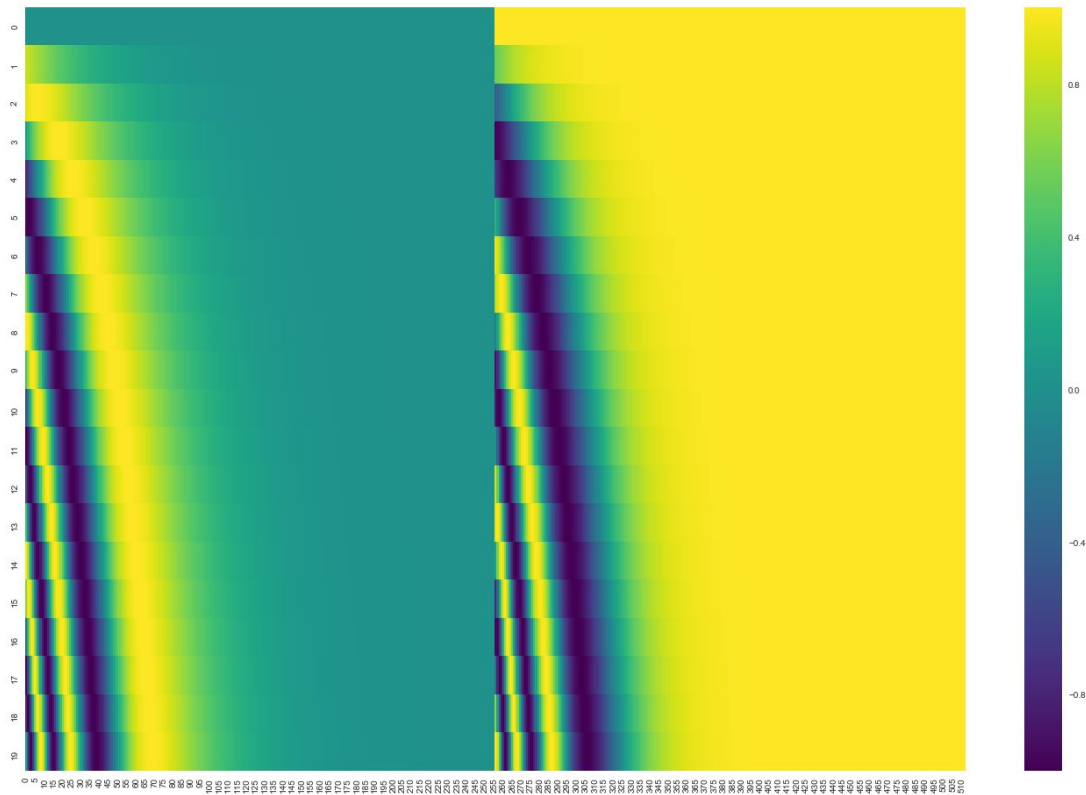
# Positional Encoding

# Positional Encoding
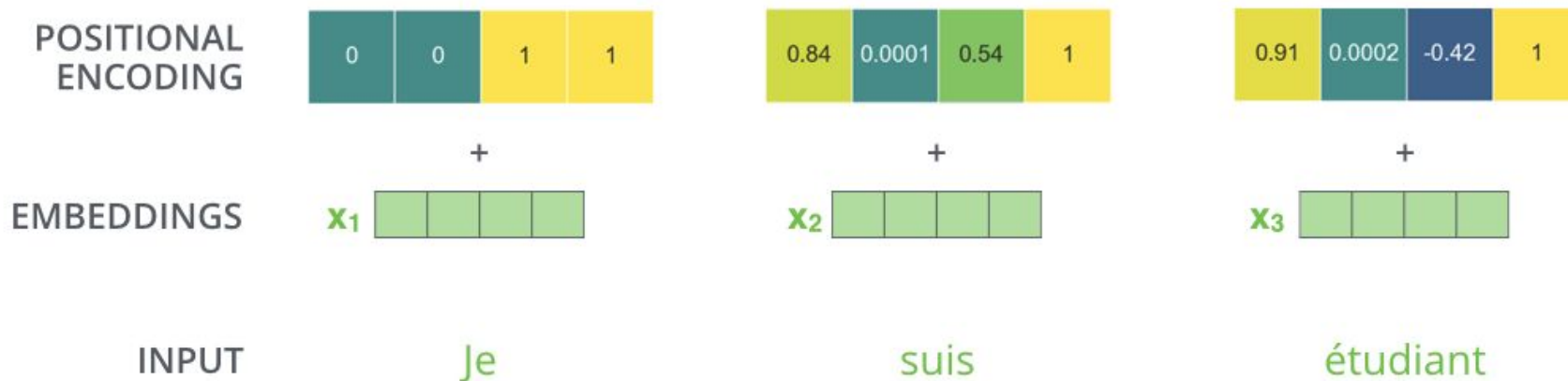


$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$
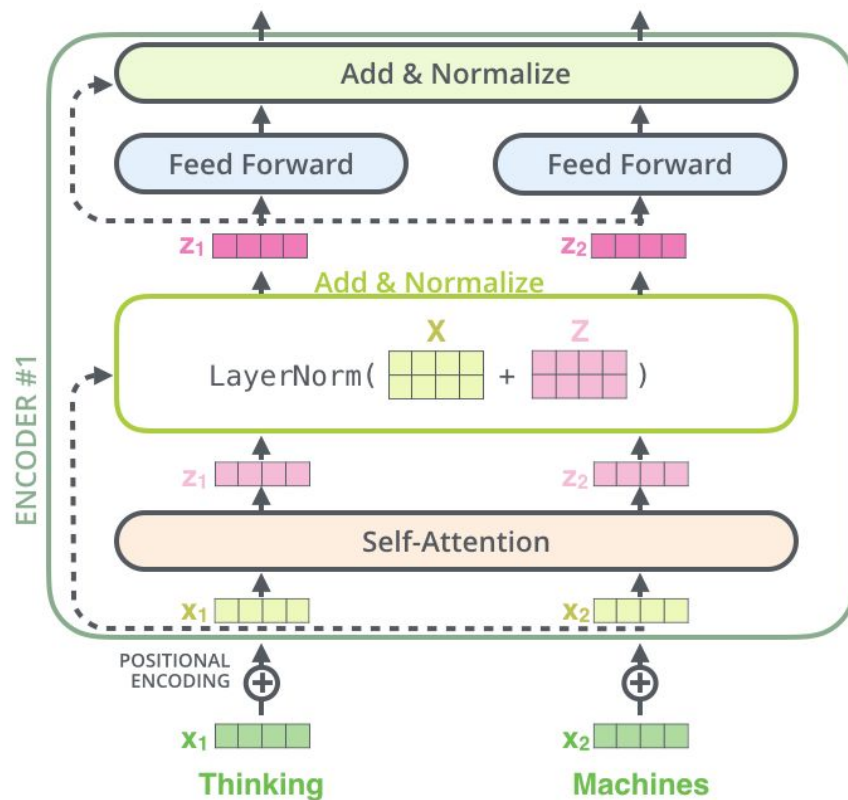
$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

# Positional Encoding

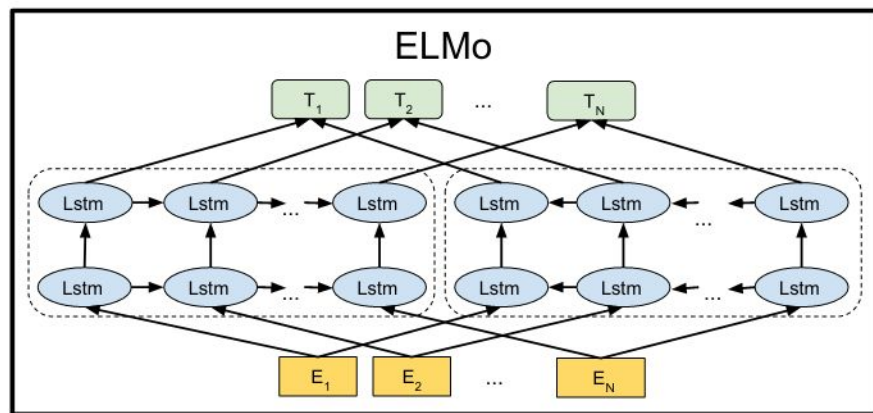For embedding with a dimensionality of 4 the encodings look like this:



POSITIONAL ENCODING

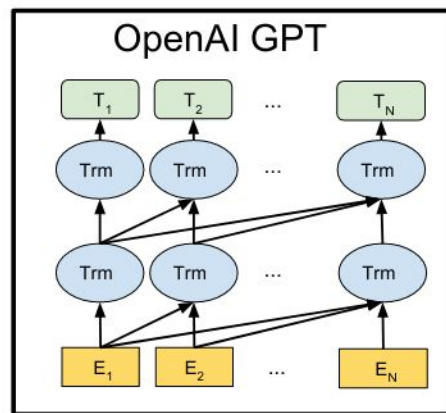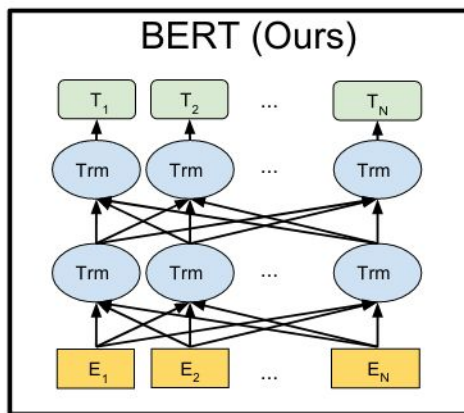| 0 | 0 | 1 | 1 |

| 0.84 | 0.0001 | 0.54 | 1 |

| 0.91 | 0.0002 | -0.42 | 1 |

+

EMBEDDINGS  $x_1$  $x_2$  $x_3$

INPUT    Je        suis        étudiant

# 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
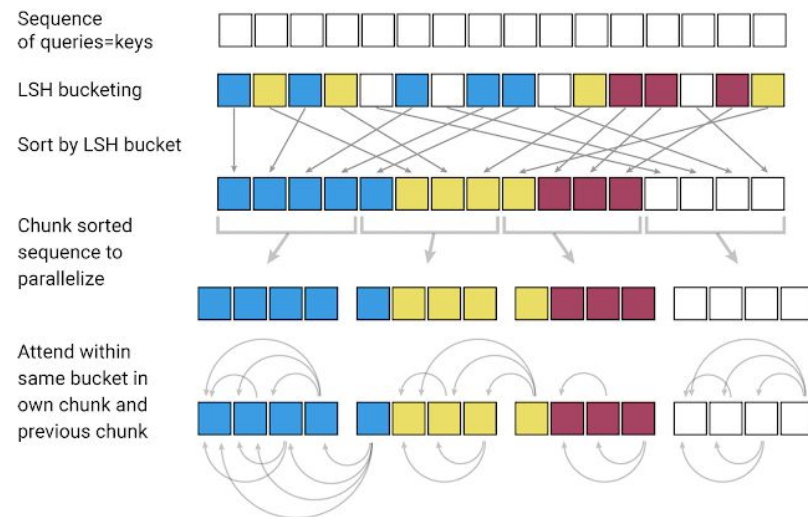


Source Bert Paper: https://arxiv.org/pdf/1810.04805.pdf

# 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, [\w*]former


- More Information
  - Paper by Kitaev, Kaiser and Levskya
  - Google AI Blog Post
  - Video Introduction
  - Background Info



Sequence of queries=keys

LSH bucketing

Sort by LSH bucket

Chunk sorted sequence to parallelize

Attend within same bucket in own chunk and previous chunk

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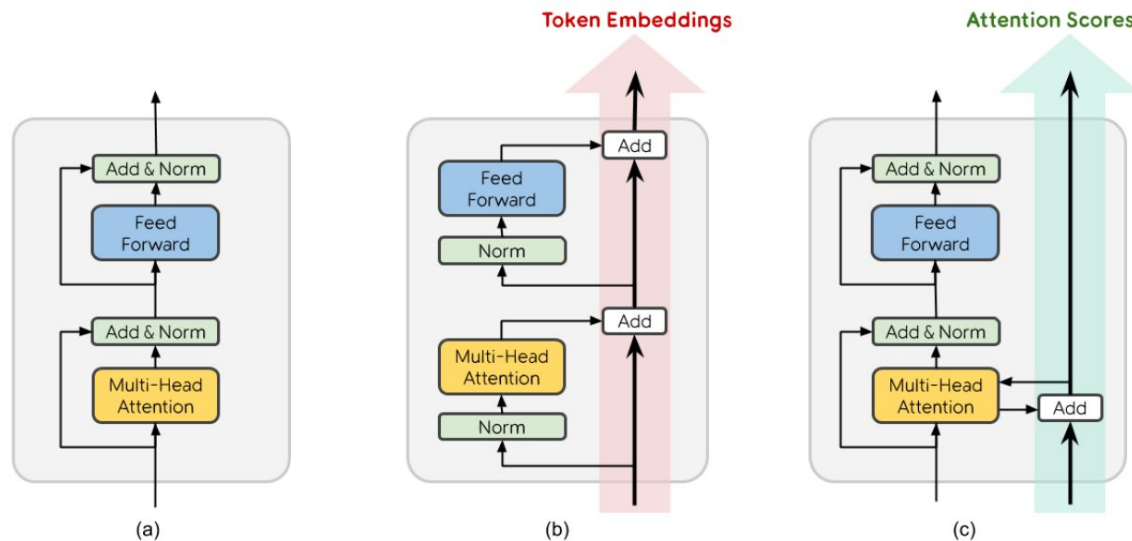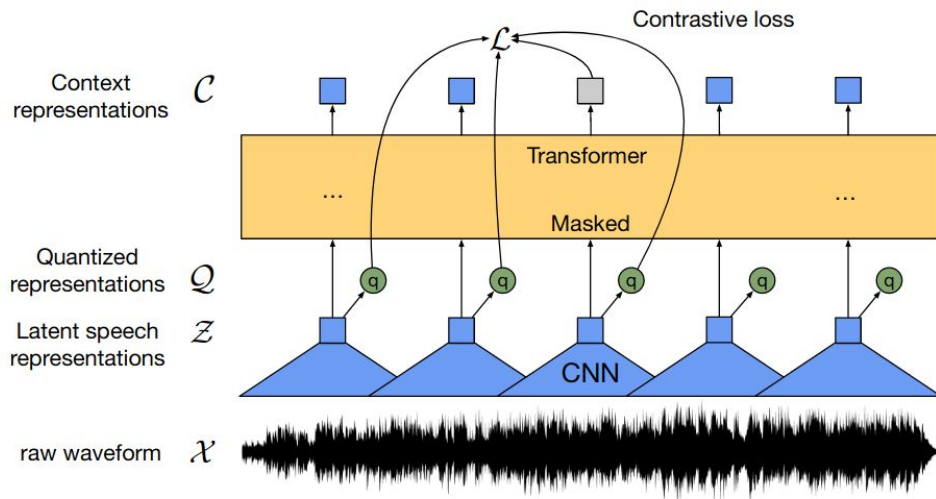- Paper by Ruining He, Anirudh Ravula, Bhargav Kanagal, Joshua Ainslie



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)).
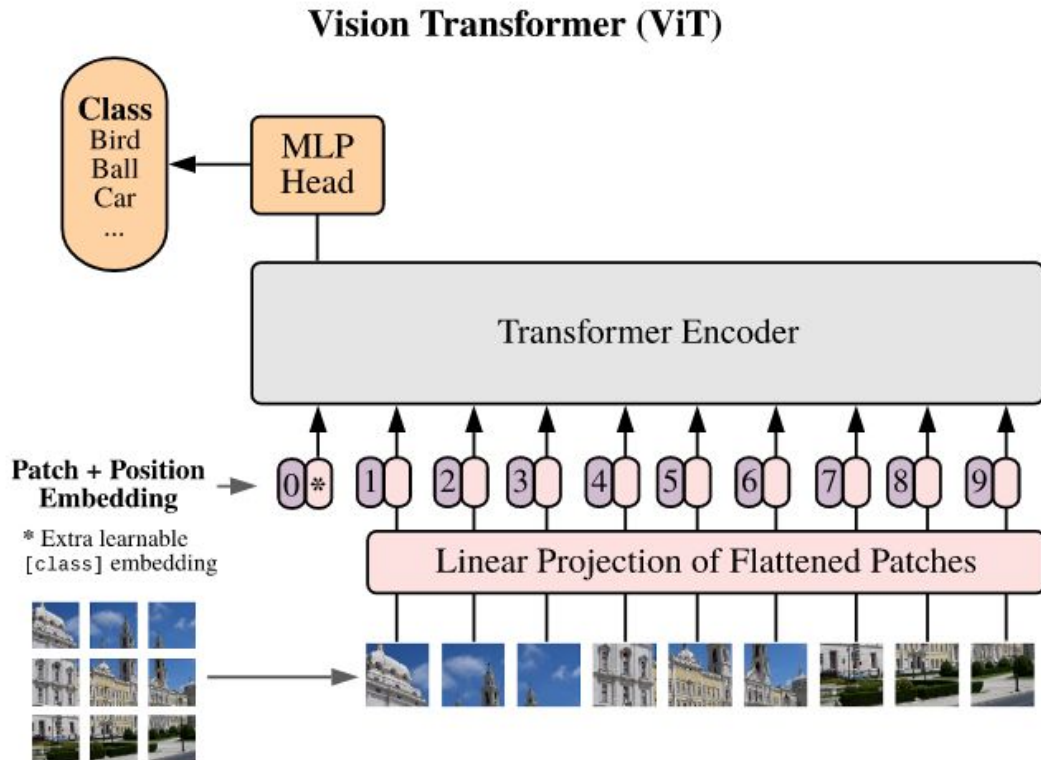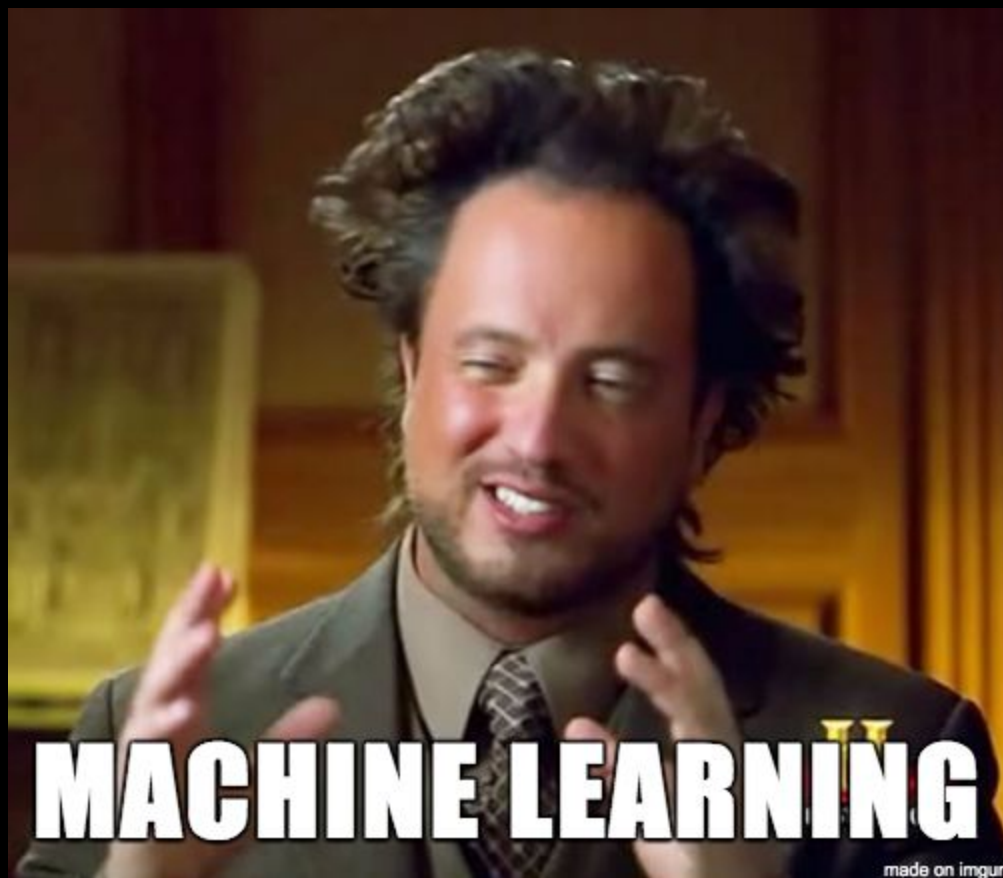
# Automatic Speech Recognition

- wav2vec 2.0: A Framework for Self-SupervisedLearning 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



**Vision Transformer (ViT)**

MACHINE LEARNING

made on imgur

I'M NOT SAYING IT'S MAGIC

BUT MAGIC

makeameme.org

# Sources

# Transformer

- 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

# Stanford Transformers Seminar