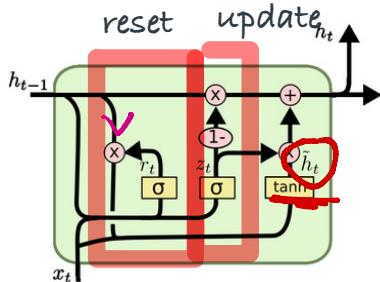


Agenda

1. GRU
2. BERT
3. BART
4. GPT
5. Vision Transformers



GRU



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

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



Reset Gate: Determine how to combine new input with the past hidden state, i.e. h_{t-1} with x_t

$h_t(1, 512)$ where 512 is the hidden state dimension

$x_t(1, 512)$ where 512 is embedding dimension of x_t

hidden^{last}
 $h_{t-1} + x_t$
 $(1024, 512)$

In Reset gate $\int [h_{t-1}; x_t] \rightarrow (1, 1024) \cdot w_r$

$(1, 1024) \times (1024, 512)$

$r_t = \sigma(1, 512)$

series of on/off

$r_t \odot h_{t-1}$

$\int [r_t \odot h_{t-1}; x_t] \odot w$

$(1, 512) \rightarrow (1, 1024)$

$(1024, 512) \rightarrow \text{---}$

value
 bus of 1

11

$$\text{candidate} = \tanh(1, 512)$$

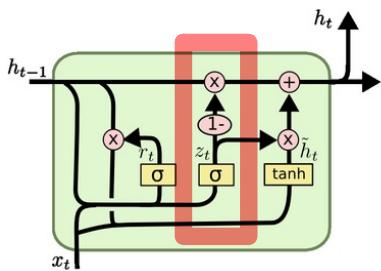
hidden state

$$= \tanh \left(\underbrace{\{ r_t \odot h_{t-1}; n_t \} \odot \omega}_{\text{determine how to combine}} \right)$$

determine how to
combine

h_{t-1} & n_t

update gate: Decide how much of past information to keep and how much of new information to add to hidden-state



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

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

$$\{ h_{t-1} : n_t \} \rightarrow (1, 1024) \cdot \omega_z$$

$$z_t = 6(1, 512) \quad \text{blue one}$$

$$h_t = \underbrace{z_t}_{\text{high}} \hat{h}_t + (1 - z_t) h_{t-1}$$

High value of z_t means low value for $(1 - z_t)$ which means forget more about past as $(1 - z_t)$ will be small

AI TIMELINE: 1947-2020

1947

Alan Turing talks about AI in London

1966

MIT releases the ELIZA chatbot

2011 2015

IBM Watson beats players on *Jeopardy!* AlphaGo beats Fan Hui

1940s

1950

Turing's papers on *Intelligent machines*

1950s

1960s

1970s

1980s

1990s

2000s

2010s

1997

Deep Blue beats Garry Kasparov in chess

2017 Google Transformer

2018 GPT-1 117M

2019 GPT-2 1.5B



2020



Jan

Google Meena 2.6B

May

GPT-3 175B

Jun

iGPT 6.8B

Sep

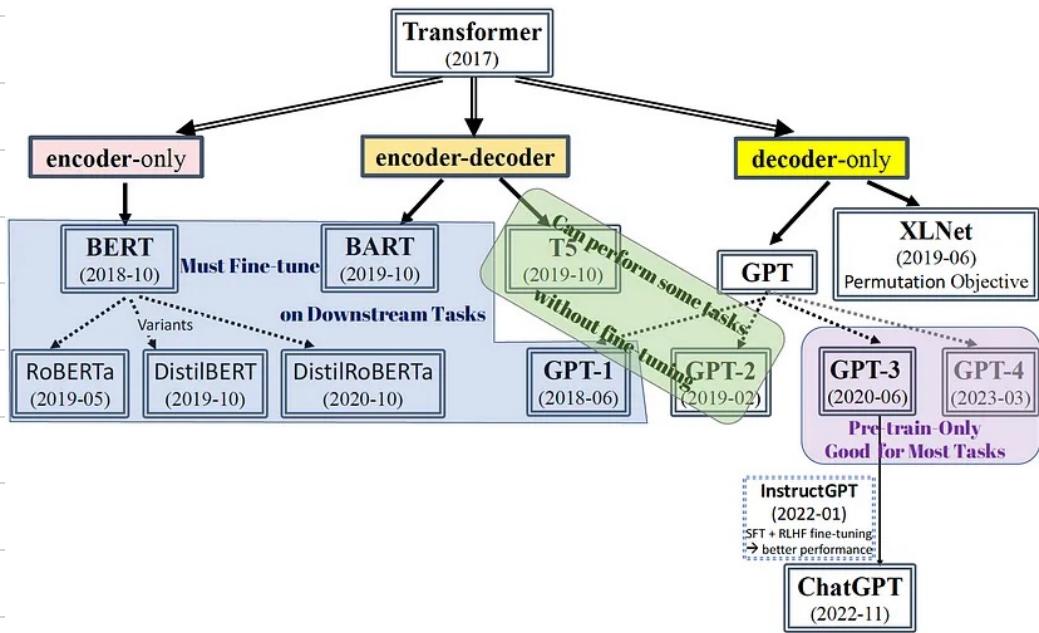
GPT-3 writes newspaper column

Selected highlights only. Alan D. Thompson. November 2021. <https://lifearchitect.ai/>



LifeArchitect.ai/timeline

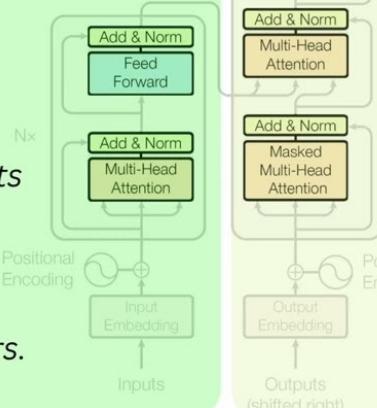
BERT - Bidirectional Encoder Representation with Transformers



BERT

Google

use transfer learning to **continue learning** from its existing data when adding user-specific tasks and layers.



GPT

OpenAI

decodes from its massive pre-trained embeddings to present output that matches user prompts. It

Figure 1: The Transformer - model architecture.

BERT is a stack of Encoders. We train BERT in two steps, "**pre-training**" and "**fine-tuning**".

During **pre-training**, the model is trained on a large dataset in an *unsupervised learning* task where the model is trained on an unlabelled dataset like the data from Wikipedia.

During **fine-tuning** the model is trained for downstream tasks like *classification, text generation, language translation, question-answering, and so forth*.

In the **pre-training** mode, BERT is trained on 2 specific tasks: **Masked Language Model** and **Next-sentence similarity prediction**.

First, let's talk about pre-training, how do you pre-train the stack of encoders

1. **Masked Language Modelling**
2. **Next Sentence Similarity Prediction**

The pre-training part, of BERT which includes above methods are **UNSUPERVISED learning**

We will discuss this each of above in elaborate detail

Masked Language Modelling

I could not visit the bank in the afternoon, as employees were on leave

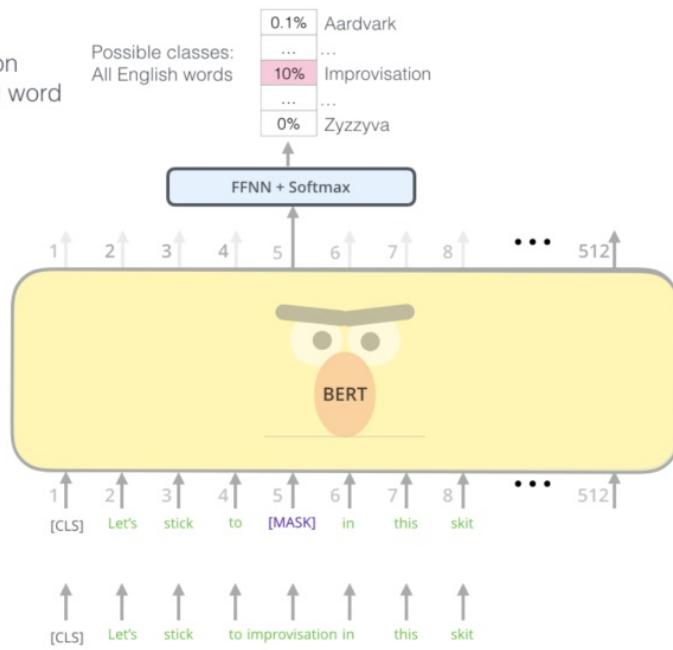
Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

| | |
|------|---------------|
| 0.1% | Aardvark |
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zzyzva |

FFNN + Softmax

Randomly mask 15% of tokens



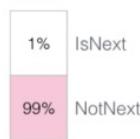
Input

Masked Language modelling uses - masked input like "the cat [MASK] the mouse" and BERT needs to predict ATE.

Basically I can use entire text corpus of internet to train this

Next Sentence Similarity Prediction

Predict likelihood
that sentence B
belongs after
sentence A



FFNN + Softmax

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

Tokenized
Input

1 [CLS] 2 the 3 man 4 [MASK] 5 to 6 the 7 store 8 [SEP]



Input

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

Sentence A

Sentence B

CLS cat loves milk. cat hates dog.

Code From Scratch: <https://colab.research.google.com/drive/1V1n57ImzTLroU1K0cRX51x9XyQMe8aPD?usp=sharing>

Link to Beautiful Article: <https://neptune.ai/blog/bert-and-the-transformer-architecture>

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:



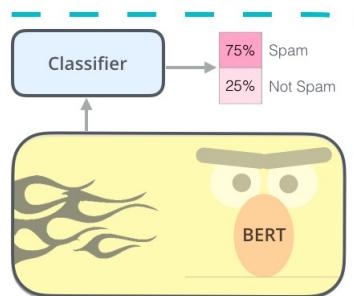
Predict the masked word
(language modeling)

Dataset:

Objective:

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

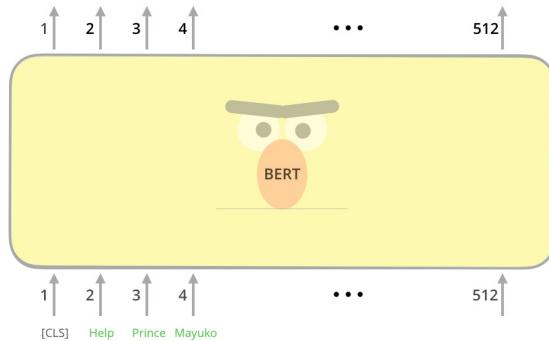
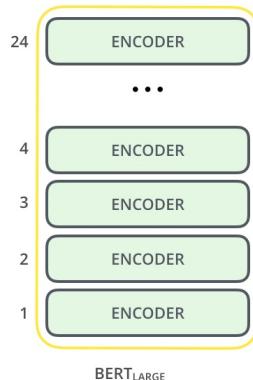
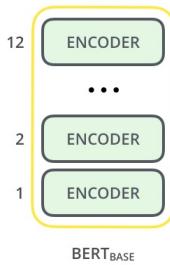
Supervised Learning Step

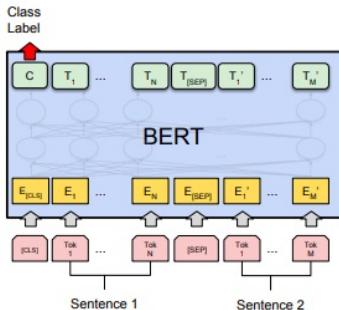


Model:
(pre-trained
in step #1)

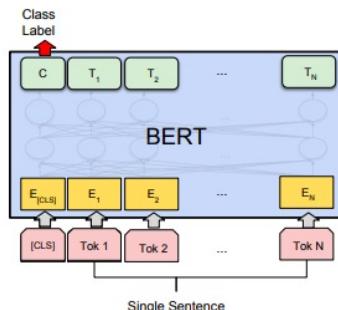
Dataset:

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

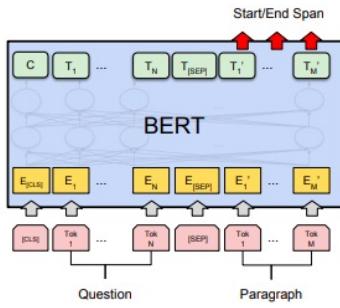




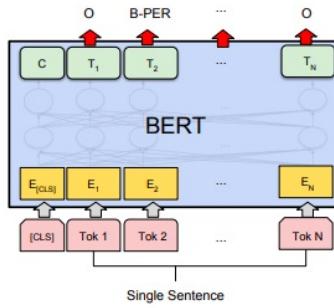
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(b) Single Sentence Classification Tasks:
SST-2, CoLA



(c) Question Answering Tasks:
SQuAD v1.1



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

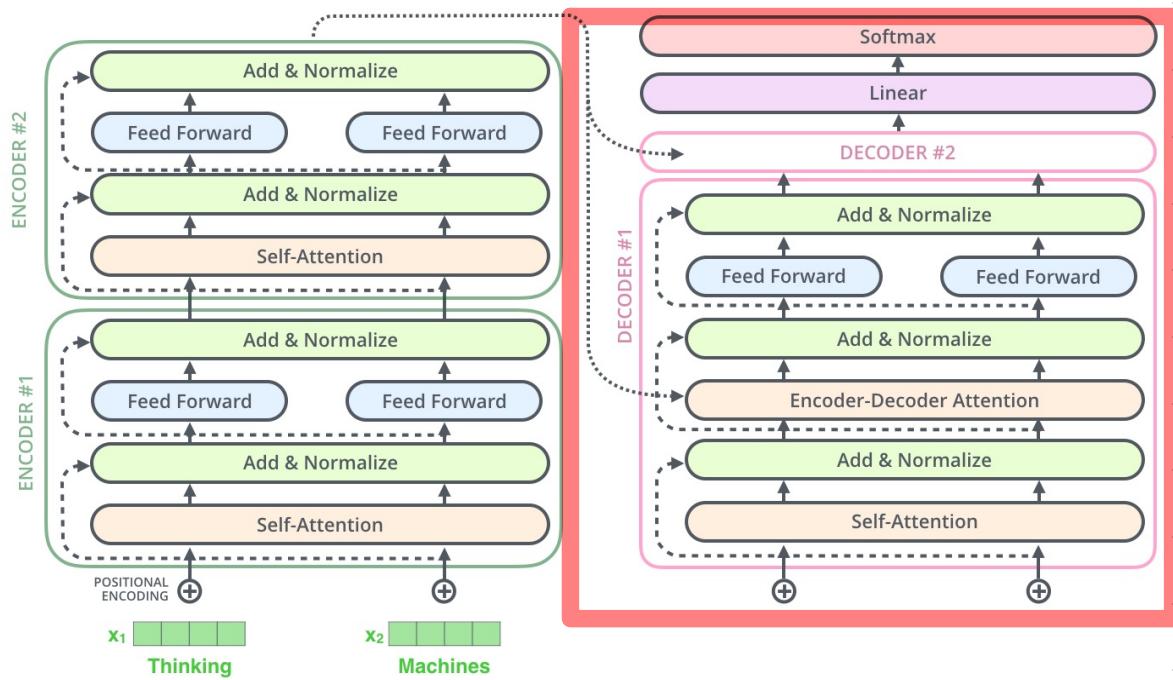
Handy vs

BART (Bidirectional and Auto-Regressive Transformers)

1. Token masking - a random subset of the input is replaced with [MASK] token, like in BERT
2. Token deletion - random tokens are deleted from the input. The model must decide which positions are missing (as the tokens are simply deleted and not replaced with anything else)
3. Text infilling - a number of text spans (length may vary) are each replaced with a single [MASK] token.
4. Sentence Permutation - the input is split based on period (.) and the sentences are shuffled.
5. Document Rotation - a token is chosen at random, and the sequence is rotated so that it starts with the chosen token.

| Corruption Scheme | Original Text | Corrupted Text | Explanation |
|----------------------|---------------|----------------|--|
| Token Masking | ABC.DE. | A._C._E. | Both B and D are masked with a single mask token for each. |
| Token Deletion | ABC.DE. | A.C.E. | Both B and D are deleted (and not replaced). |
| Text Infilling | ABC.DE. | A._D_E. | The span BC is replaced with a single mask token. A 0 length span is inserted between D and E. |
| Sentence Permutation | ABC.DE. | DE.ABC. | Split into sentences at periods (.) and shuffled. |
| Document Rotation | ABC.DE. | C.DE.AB | The sequence is rotated around C. |

GPT-S



let ↗

2 ↗

3 ↗

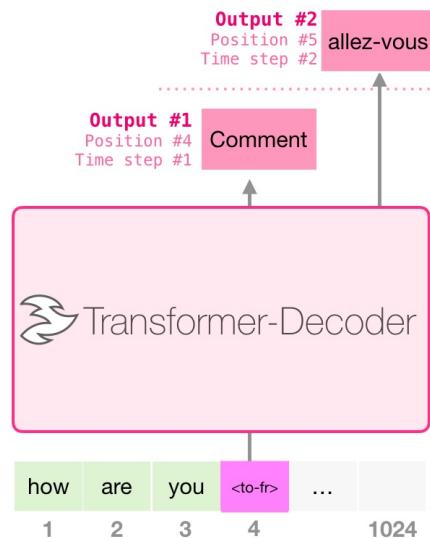
out¹ → 5

out¹ → suis

out¹ → étudie

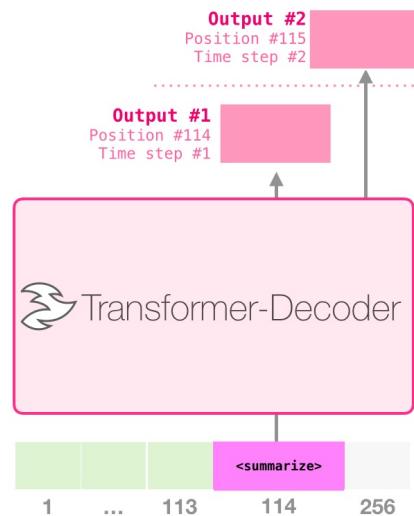
Training Dataset

| | | | | | | | |
|------|---------|---------|---------|---------|--------|---------|----------|
| I | am | a | student | <to-fr> | je | suis | étudiant |
| let | them | eat | cake | <to-fr> | Qu'ils | mangent | de |
| good | morning | <to-fr> | Bonjour | | | | |



Training Dataset

| | | |
|-------------------|-------------|--------------------|
| Article #1 tokens | <summarize> | Article #1 Summary |
| Article #2 tokens | <summarize> | Article #2 Summary |
| Article #3 tokens | <summarize> | Article #3 Summary |



very detailed article for GPT-2 <https://jalammar.github.io/illustrated-gpt2/>

GPT:

Generative Pre-trained Transformers are decoder-only transformer models. These models are pre-trained on massive amounts of data, such as books, web-pages, to generate contextually relevant and semantically coherent language.

GPT1 (117M parameters) was released in 2018 and was trained on [Common Crawl](#)[Links to an external site.](#) and the [BookCorpus](#)[Links to an external site.](#)

While GPT was great, it was prone to generating repetitive text and also failed to reason over multiple turns of dialogue.

GPT2:

GPT2 was released in 2019 and had 1.5B parameters. It was trained on much larger and more diverse dataset. It also struggled with tasks that required complex reasoning and understanding of context. It failed to maintain context and coherence over longer passages.

GPT3:

GPT3 was released in 2020 and has 175B parameters! It was trained on an even larger and diverse range of data sources, including BookCorpus, Common Crawl, and Wikipedia, among other things. The datasets comprise nearly a trillion words, allowing GPT3 to generate sophisticated responses on a wide range of NLP tasks, even without providing any prior example data.

One of the main improvements of GPT3 over its previous models is its ability to generate coherent text, write computer code, and even create art (SVG, etc.). Unlike the previous models, GPT3 understands the context of a given text and can generate appropriate responses which led to applications like content creation, chatbots, and language translation.

GPT3 had its flaws. The model can return biased, inaccurate, and inappropriate responses. This issue arises because it is trained on massive amounts of texts that possibly contain biased and inaccurate information.

GPT4:

GPT4 was released on March 14, 2023. It's a significant step up from GPT3. Rumors claim that GPT4 has 1.8 trillion parameters, which was estimated by the speed it was running. Some other [interestingLinks to an external site](#), facts/rumors: GPT4

has 1.8 trillion parameters

has 120 layers (10 times more than GPT3)

Has 16 Mixtures of Expert (MoE)

Was trained on 13 trillion tokens

Was trained on GPT3 data + twitter, youtube, reddit, and more books

Took \$63M to train

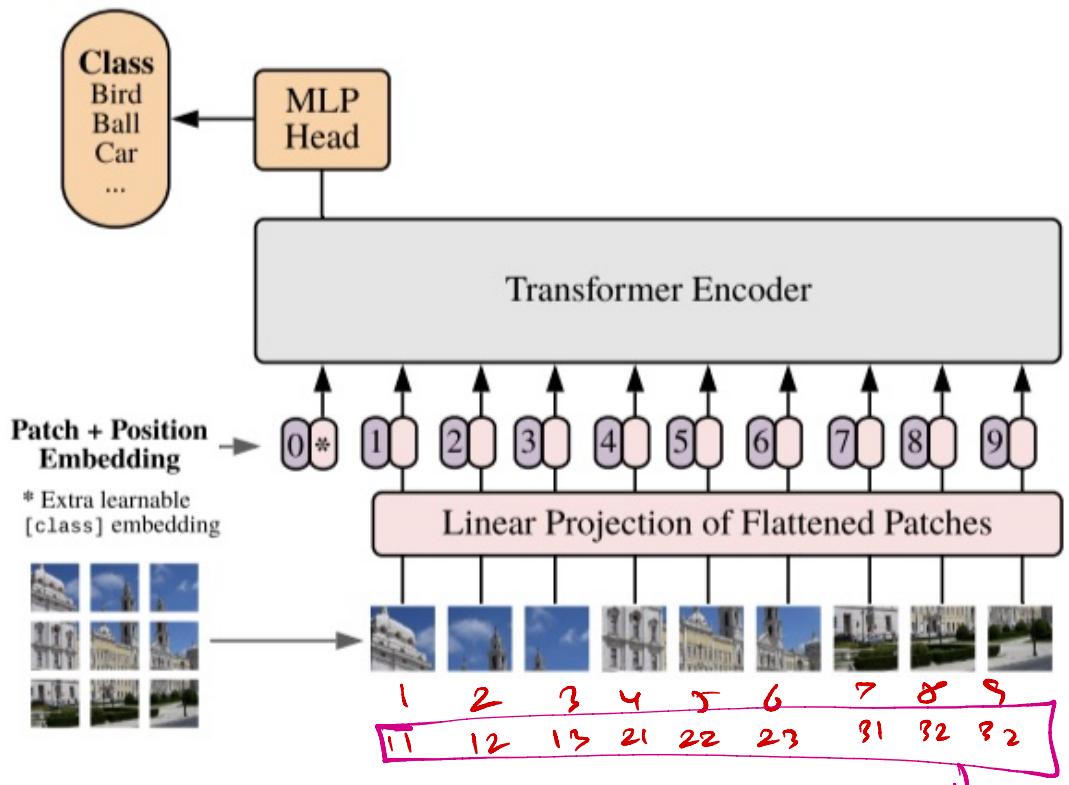
```
[example] an input that says "search" [toCode] Class App extends React Component... </div> } } }  
[example] a button that says "I'm feeling lucky" [toCode] Class App extends React Component...  
[example] an input that says "enter a todo" [toCode]
```



GPT-3

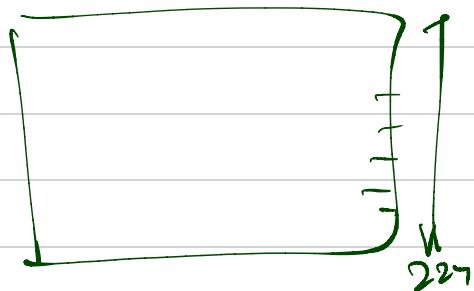


Vision Transformers



Input Image $\rightarrow 224 \times 224 \times 3$

14×14



14×14

196

$196 \text{ Patch} \rightarrow 16 \times 16 \times 3$

$16 \times 16 \times 3$

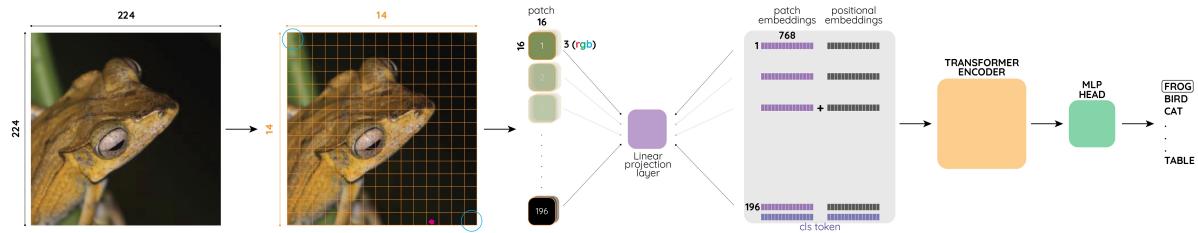
$224 / 16$
 $= 14$

$16 \times 16 \times 3$

$\rightarrow 196 \times 768 + PE$

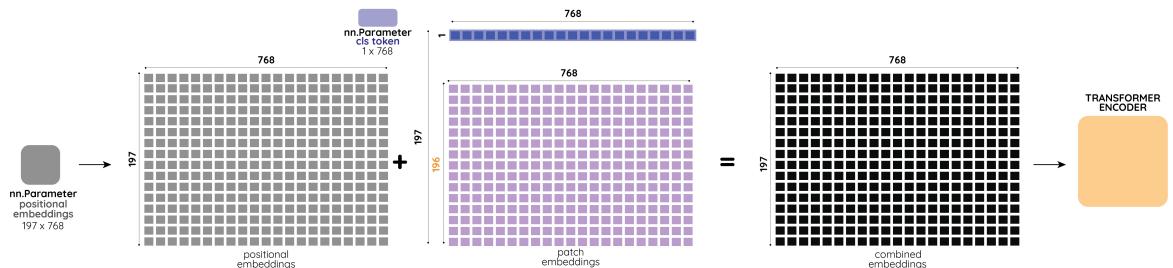
The overall architecture can be described easily in five simple steps:

1. Split an input image into patches
2. Get linear embeddings (representations) from each patch referred to as Patch Embeddings
3. Add positional embeddings and a [CLS] token to each of the Patch Embeddings.
4. Pass through a Transformer Encoder and get the output values for each of the [CLS] tokens.
5. Pass the representations of [CLS] tokens through an MLP Head to get final class predictions.



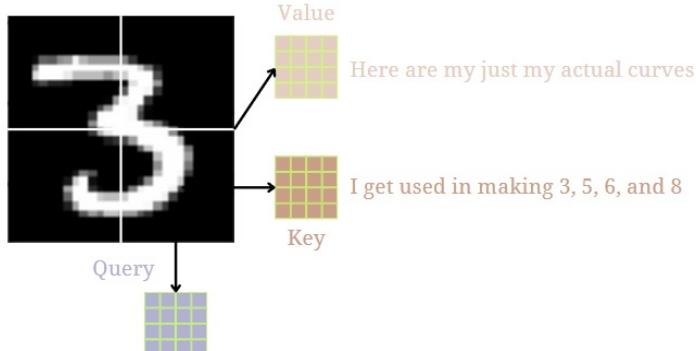
As an overall method, from the paper

We split an image into fixed-size patches, linearly embed each of them, add positional embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable "classification token" to the sequence.



Ques

Key: Value



I am looking for curved edges that can help me find 3

$$\text{Stride} = 32$$

$$\begin{array}{c}
 32 \times 82 \times 3 \\
 \text{Stride} = 32 \\
 224 \times 224 \times 3 \\
 \text{Stride} = 16 \\
 \text{Result} = 16 \times 16 \times 3 \times 768 \\
 \text{Stride} = 16
 \end{array}$$

outlet shape Convolutn

$14 \times 14 \times 768$

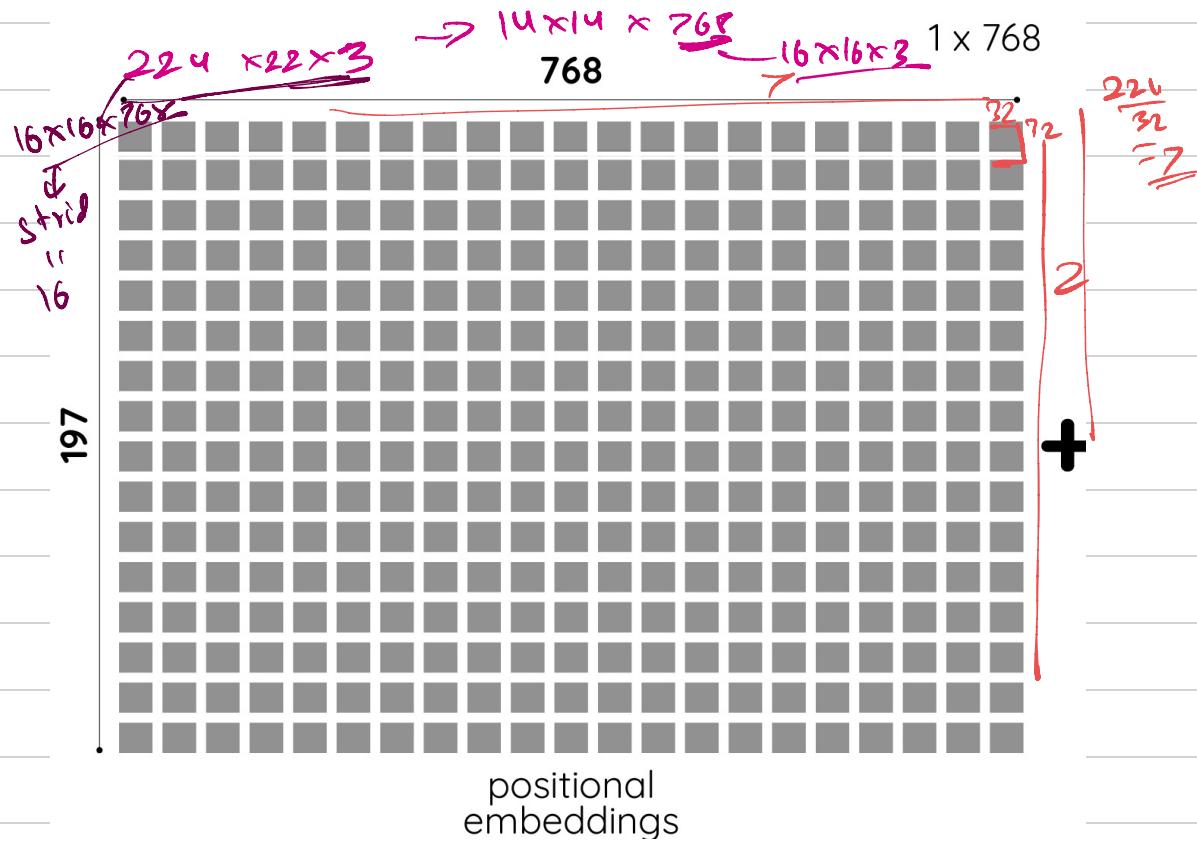
$$\frac{f \cdot P + 2P - K}{\text{Stride}} + 1$$

$$[(W-K+2P)/S]+1$$

$$\frac{224 + 2P - 16}{16} + 1$$

$$\frac{208 + 2 \times 0 - 16}{16} + 1$$

$$= \frac{13}{1} + 1 = 14$$



$I_{max} \rightarrow 224 \times 224 \times 3 \rightarrow \underbrace{n}_{49} \rightarrow \underbrace{32 \times 32 \times 3}_{11 \text{ 2048}} \rightarrow \underbrace{32 \times 32 \times 3}_{10} \rightarrow \text{Platmnd out}$

$$224 \times 224 \times 3 \quad \textcircled{=} \quad 32 \times 32 \times 3 \times 3072$$

Kernel

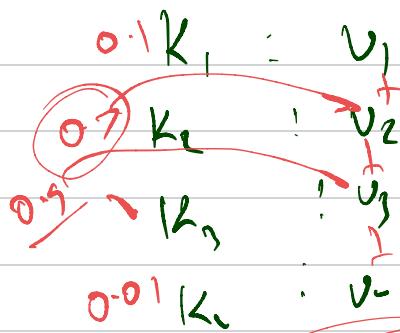
$$> 7 \times 7 \times 3072$$

$\Rightarrow 49 \times 3072$

$$\frac{224 - 2 \times 0 + 32}{32} + 1 \Rightarrow \underline{6+1} = 7$$

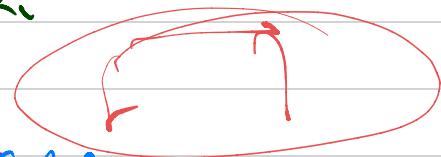
$Q_1 =$

wk



Best day of your life

dent



$v \xrightarrow{wv}$

Scddy day of in

dink \xrightarrow{gnP}

Scddy day of in