# Transformers

**Akshitha Kumbam**

Acknowledgement - Yilun Kuang

10/16/2023

# Outline

- **Encoder-Only / Decoder-Only**
- **Huggingface Transformer - Provides pretrained models, tools and APIs for using transformer based architectures.**

**NYU**

# Encoder-Only | BERT

**Bidirectional Encoder Representations from Transformers (BERT) - Google.**

Model Architecture

- Transformer Encoder

What is Special About It

- Mask Language Modeling & Next Sentence Prediction
- Downstream task adaptation

**NYU**

BERT's primary goal is to create contextualized word embeddings, meaning it captures the meaning of a word based on its context in the sentence.
BERT's embeddings are dynamic and depend on the **words surrounding the target word**.

## How BERT Achieves This Goal:

To achieve contextualized embeddings, BERT is trained using two key tasks:
1. Next Sentence Prediction (NSP):
   a. What is NSP? The model is given two sentences, and it must predict whether the second sentence logically follows the first sentence or if it is a random sentence.
   b. Why? This task helps BERT understand relationships between sentences, which is useful for downstream tasks like question answering, text classification, or sentence similarity.
2. Masked Language Model (MLM):
   a. What is MLM? In this task, random words in the input sentence are masked (replaced with [MASK]), and the model is trained to predict those masked words based on the context provided by the other words in the sentence.
   b. Why? By doing this, BERT learns how words relate to each other in context. It helps the model understand how different words fit together, which is essential for producing rich, context-aware word representations.
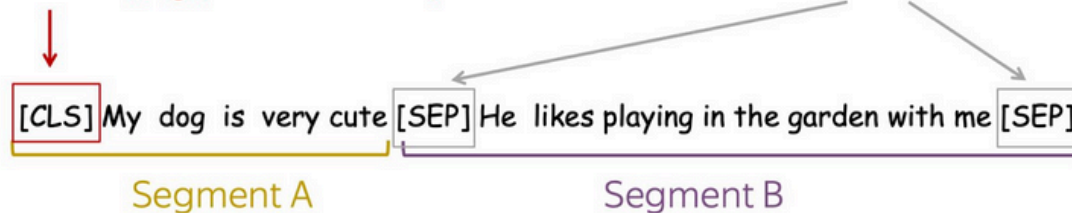
# BERT - Pretrain | Inputs

**Training Input:** 1) pairs of sentence; 2) [CLS] token; 3) [SEP] token

**[CLS]**: Special token
- Training time: predict if sentences are consecutive or not (Next Sentence Prediction /NSP objective)
- Test time: downstream tasks (e.g., classification)

[SEP]: Special token-separator

[CLS] My dog is very cute [SEP] He likes playing in the garden with me [SEP]

Segment A        Segment B

Training on pairs of sentences: either consecutive or random (50%/50%)

# Embedding

1) Token

2) Segment

3) Position

Training time: predict if sentences are consecutive (NSP objective)

Test time: classification

Model (Transformer encoder)

several layers

Input

positions
0 1 2 3 4 ...

segments
A A A A A A B B B B B

tokens
[CLS] My dog is ...

Training on pairs of sentences: either consecutive or random (50%/50%)

[CLS] My dog is cute [SEP] He is very fluffy [SEP]

Segment A          Segment B

NYU

# BERT - Pretrain | Inputs



Loss

Cross-entropy loss

Target          saw    grey        mat

Prediction      $P(*|...)$ $P(*|...)$    $P(*|...)$

several layers
(Transformer's encoder)

I  saw  a  grey  cat  on  a  mat  .  <eos>
   [MASK]  tea           mat

○ **[MASK]**,          ○ Random token,      ○ Original token,
   with p = 80%            with p = 10%          with p = 10%

At each training step:

- pick randomly 15% of tokens
- replace each of the chosen tokens with something
- predict original chosen tokens

# BERT - Pretrain | Objective

**Next Sentence Prediction (NSP)**

**Input:** [CLS] the man went to [MASK] store [SEP] he bought a gallon [MASK] milk [SEP]
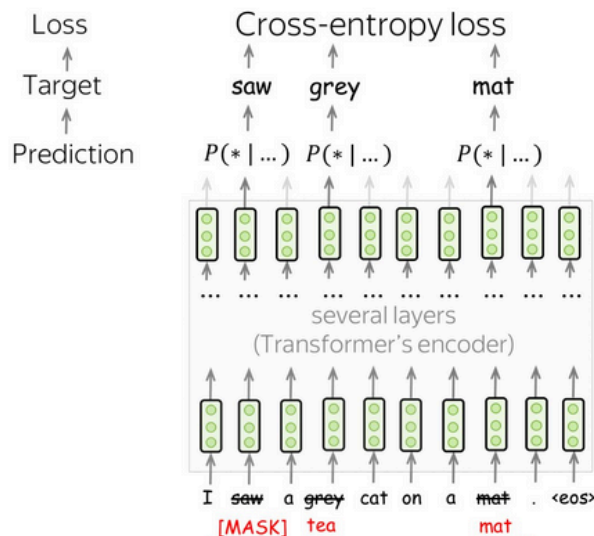**Label:** isNext

**Input:** [CLS] the man went to [MASK] store [SEP] penguin [MASK] are flight ##less birds [SEP]

**Label:** notNext

# BERT - Pretrain | Objective

**Masked Language Modeling (MLM)**

Output layer of MLM - Output Linear Layer: Generates logits for vocabulary.

**During training:**

Masked tokens: Predicted tokens are compared to original tokens.
Loss calculation: Cross-entropy loss calculated for each masked token.
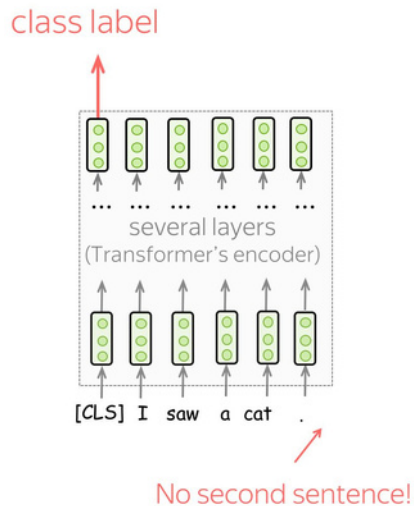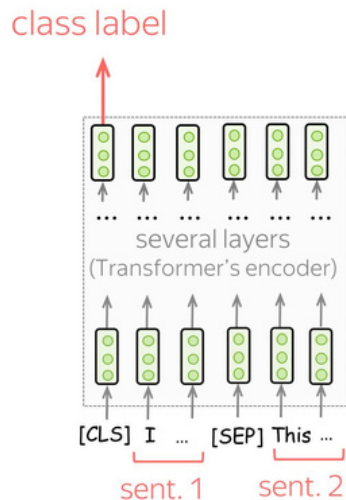Backpropagation: Gradients computed to update model weights.

**Benefits:**

1. Contextual understanding: BERT learns to represent tokens in context.
2. Improved language modeling: Better handling of out-of-vocabulary tokens.
3. Transfer learning: Fine-tuning benefits from pre-trained MLM.

# BERT - Finetune | Tasks
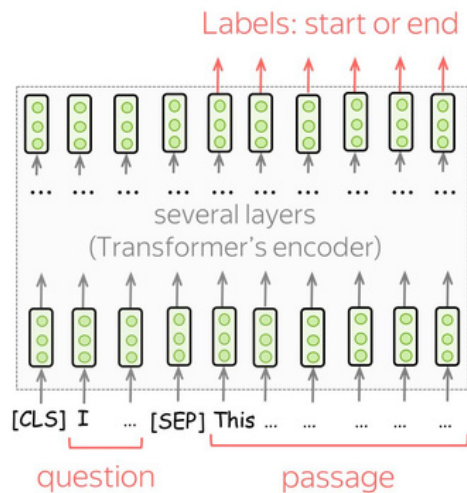
**Single sentence classification**
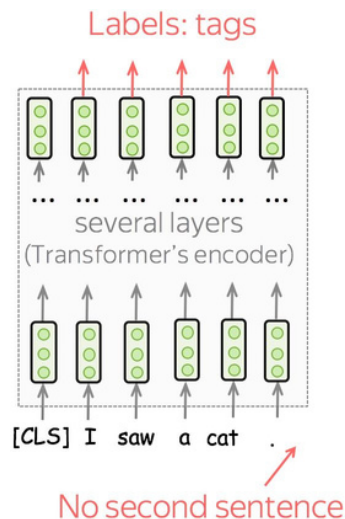


**Sentence Pair Classification**

# BERT - Finetune | Tasks

**Question Answering**



**Single sentence tagging**

# Decoder-Only | GPT

**Generative Pretrained Transformer (GPT) - OpenAI**

Model Architecture

- Transformer **Decoder**

What is Special About It

1. Autoregressive Language Modeling - Predict next token given previous tokens.

2. Downstream task adaptation.

Mathematically:
P(token_n | token_1, token_2, ..., token_{n-1})

Key features:
- **Unsupervised Learning**: Trained on large corpus without labeled data.
- **Generative Capabilities**: Generates coherent text.
- **Transfer Learning**: Fine-tunes for downstream tasks.

**Downstream Task Adaptation**
GPT's adaptation:
- **Language Translation**: Generate translations.
- **Text Summarization**: Summarize long texts.
- **Chatbots**: Engage in conversations.

# GPT - Pretrain | Inputs

**Training Input:** 1) sentences; 2) [PAD] / [EOS] token

Sentences: Sequential text data.

[PAD] token: Padding token for fixed-length sequences.

[EOS] token: End-of-sequence token.

**Example:**

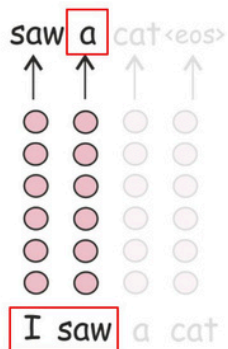My dog is very cute. He likes playing in the garden with me [EOS] [PAD] ... [PAD]

**Embedding:**

Token embeddings + positional embeddings

# GPT - Pretrain | Objective
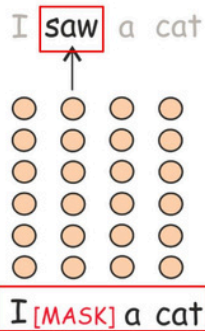
**Autoregressive Language Modeling**

## Language Modeling

- Target: next token
- Prediction: $P(* | \text{I saw})$



saw **a** cat <eos>

**I saw** a cat

left-to-right, does
not see future

## Masked Language Modeling

- Target: current token (the true one)
- Prediction: $P(* | \text{I [MASK] a cat})$



I **saw** a cat

**I** [MASK] **a cat**

sees the whole text, but
something is corrupted

$P(\text{token\_n} | \text{token\_1}, \text{token\_2}, ..., \text{token\_{n-1}})$

$P(\text{token\_i} | \text{token\_1}, ..., \text{token\_{i-1}}, \text{[MASK]}, \text{token\_{i+1}}, ...)$

# GPT - Finetune | NLU Tasks

**GPT for Natural Language Understanding (NLU) Tasks**

- You can train an additional linear head on top of the final transformer block activation vector.

Source: Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.

# Encoder-Decoder | T5

**Text-to-Text Transfer Transformer (T5)** Model Architecture

- Transformer Encoder & Decoder

What is Special About It

- pretraining on multi-task mixture of unsupervised and supervised tasks (converted to Text-to-Text format)

# Resources

- Encoder-Only (BERT)
  - https://github.com/JonasGeiping/cramming
- Decoder-Only (GPT)
  - https://github.com/karpathy/nanoGPT
- Encoder-Decoder (T5)
  - https://github.com/PiotrNawrot/nanoT5

# Summary

Encoder-Only (BERT)

Decoder-Only (GPT)

Encoder-Decoder (T5)

# **Acknowledgement**