

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Paper Presentation

Isaac Kobby Anni

Computer Science
Bowling Green State University

November 18, 2024

- ① Background
- ② Introducing Transformer
- ③ BERT

1 Background

2 Introducing Transformer

3 BERT

Prerequisite

- Knowledge of recurrent neural networks and how they work.

Prerequisite

- Knowledge of recurrent neural networks and how they work.
- **Attention is all you need paper.**

Introduction

- RNNs like LSTM and GRU were established as the state of the art for sequence modeling.

Introduction

- RNNs like LSTM and GRU were established as the state of the art for sequence modeling.
- For problems like language modeling and machine translation.

Introduction

Recurrent Neural Networks (RNN)

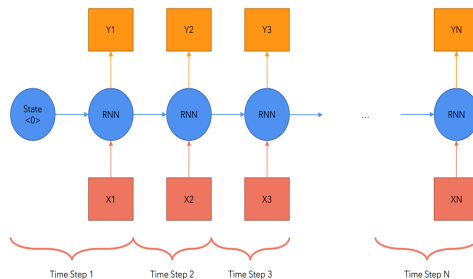


Figure 1: Workflow of RNN: *photo credit, Umar Jamil*

Introduction

- Problem as sequence length increases;

Introduction

- Problem as sequence length increases;
 - i **Slow computation**

Introduction

- Problem as sequence length increases;
 - i Slow computation
 - ii Varnishing and Exploding gradients

Introduction

- Problem as sequence length increases;
 - i Slow computation
 - ii Varnishing and Exploding gradients
 - iii Difficulty in accessing information from longer past time step

① Background

② Introducing Transformer
Attention is all you need.

③ BERT

① Background

② Introducing Transformer
Attention is all you need.

③ BERT

Transformer Architecture

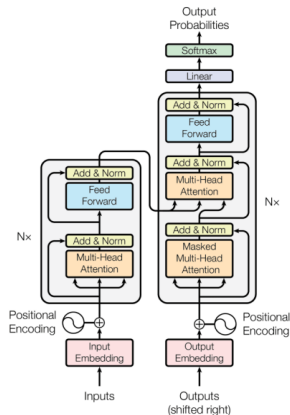


Figure 2: Transformer Model: *photo credit, Attention paper*

The Encoder

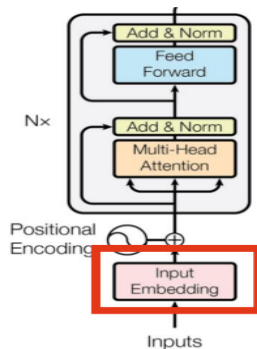


Figure 3: Transformer Model: *photo credit, Attention paper*

Input Embedding

Original sentence
(tokens)

YOUR

CAT

IS

A

LOVELY

CAT

Figure 4: Document Representation: *photo credit, Umar Jamil*

Input Embedding

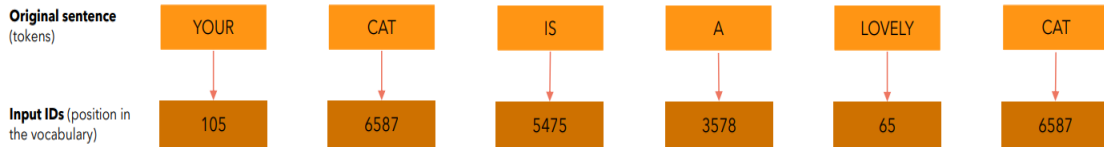


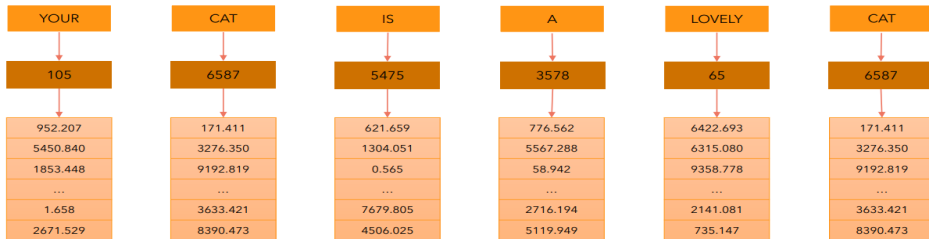
Figure 5: Document Representation

Input Embedding

Original sentence
(tokens)

Input IDs (position in
the vocabulary)

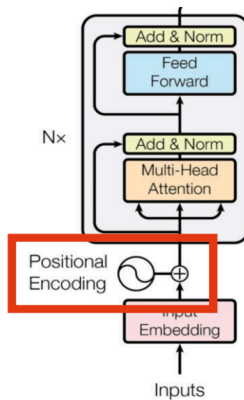
Embedding
(vector of size 512)



We define $d_{\text{model}} = 512$, which represents the size of the embedding vector of each word

Figure 6: Document Representation

Positional Embedding



Positional Embedding

- Addressing the sequential encoding like in RNNs;

Positional Embedding

- Addressing the sequential encoding like in RNNs;
i Position of each word in the document.

Positional Embedding

- Addressing the sequential encoding like in RNNs;
 - i Position of each word in the document.
 - ii Treating words that are close as "close" and distant as "distant"

Positional Embedding

- Addressing the sequential encoding like in RNNs;
 - i Position of each word in the document.
 - ii Treating words that are close as "close" and distant as "distant"
 - iii Position encoding portray patterns learned by the model.

Positional Embedding

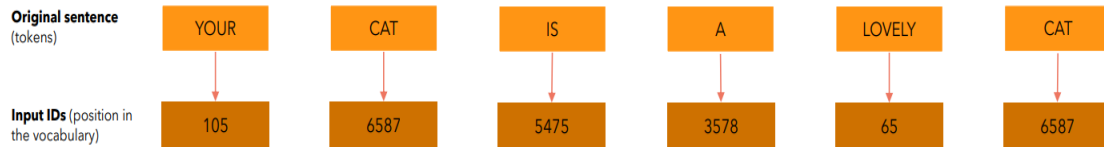


Figure 8: Document Representation

Positional Embedding

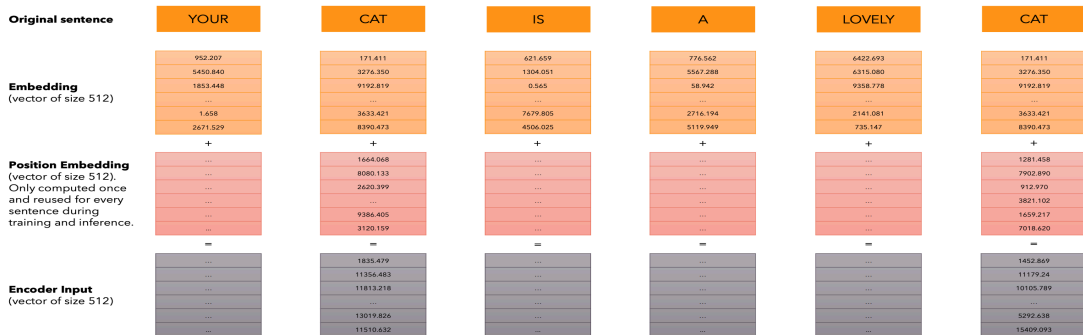


Figure 9: Document Representation

Positional Embedding Calculation

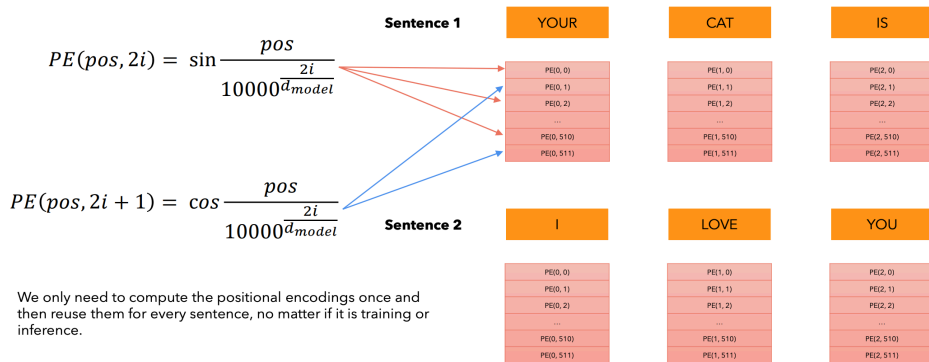
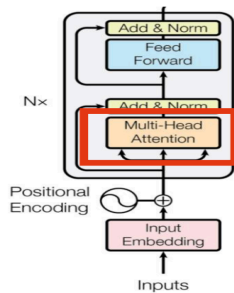


Figure 10: Document Representation

Attention Mechanism



Attention Mechanism

Self-Attention allows the model to relate words to each other.

In this simple case we consider the sequence length **seq** = 6 and **d_{model}** = **d_k** = 512.

The matrices **Q**, **K** and **V** are just the input sentence.

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

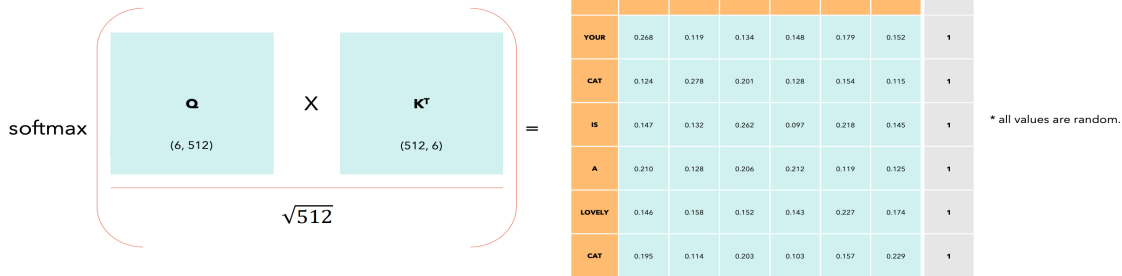


Figure 12: Self Attention computation

Attention Mechanism

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

(6, 6)

×

V

(6, 512)

=

Attention

(6, 512)

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

Each row in this matrix captures not only the meaning (given by the embedding) or the position in the sentence (represented by the positional encodings) but also each word's interaction with other words.

Figure 13: Self Attention computation

Properties of Self Attention

- No requires no parameter(s)

Properties of Self Attention

- No requires no parameter(s)
- Values along diagonals are higher.

Properties of Self Attention

- No requires no parameter(s)
- Values along diagonals are higher.
- Can switch off word interactions. etc

Self-Attention Computation

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}$$

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

Attention Mechanism

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

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

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

Figure 14: Multi-Head Attention

Multi-Head Attention Mechanism - DECODER

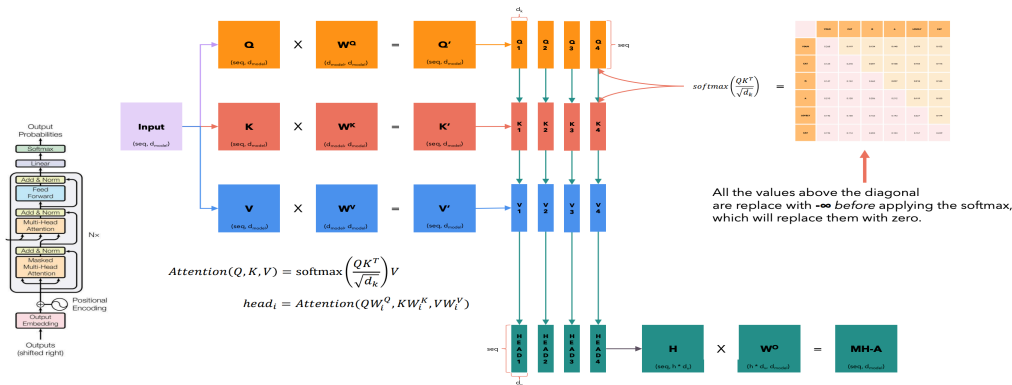


Figure 15: Masked Multi-head Attention computation

1 Background

2 Introducing Transformer

3 BERT

Motivation

- Existing models like OpenAI GPT were unidirectional.

Motivation

- Existing models like OpenAI GPT were unidirectional.
- Understanding context in both directions was a challenge.

Goal

- To build a model that fully leverages context from both directions for improved language understanding.

Model Architecture

- Architecture is made of layers of **ENCODER**'s of the transformer model.

Model Architecture

- Architecture is made of layers of **ENCODER**'s of the transformer model.
- **BERT-base:**

Model Architecture

- Architecture is made of layers of **ENCODER**'s of the transformer model.
- **BERT-base:**
 - i 12 encoder layers

Model Architecture

- Architecture is made of layers of **ENCODER**'s of the transformer model.
- **BERT-base:**
 - i 12 encoder layers
 - ii 12 attention heads

Model Architecture

- Architecture is made of layers of **ENCODER**'s of the transformer model.
- **BERT-base:**
 - i 12 encoder layers
 - ii 12 attention heads
 - iii 3072 hidden size of feed-forward layer

Model Architecture

- Architecture is made of layers of **ENCODER**'s of the transformer model.
- **BERT-base:**
 - i 12 encoder layers
 - ii 12 attention heads
 - iii 3072 hidden size of feed-forward layer
 - iv 768 embedding size

Model Architecture

- **BERT-large:**

Model Architecture

- **BERT-large:**
 - i 24 encoder layers

Model Architecture

- **BERT-large:**
 - i 24 encoder layers
 - ii 16 attention heads

Model Architecture

- **BERT-large:**
 - i 24 encoder layers
 - ii 16 attention heads
 - iii 4096 hidden size of feed-forward layer

Model Architecture

- **BERT-large:**
 - i 24 encoder layers
 - ii 16 attention heads
 - iii 4096 hidden size of feed-forward layer
 - iv 1024 embedding size

BERT Training

Two ways of training;

1 Pre-training

BERT Training

Two ways of training;

1 Pre-training

2 Fine-tuning

BERT Training

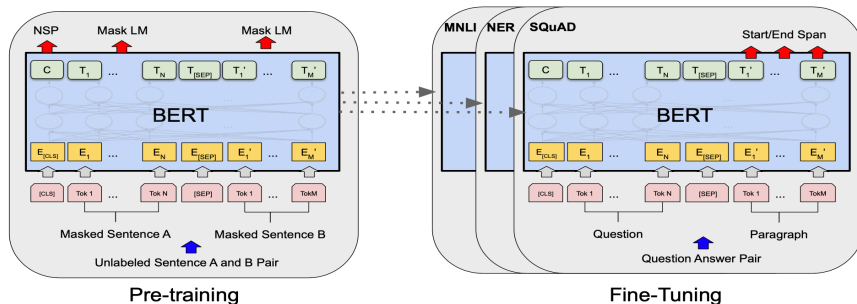


Figure 16: BERT Training

BERT: Pre-training tasks

Masked Language Modeling (MLM)

- Randomly masks 15% of the tokens.

BERT: Pre-training tasks

Masked Language Modeling (MLM)

- Randomly masks 15% of the tokens.
- Predicts the masked tokens using bidirectional context.

Masked Language Modeling (MLM)

SENTENCE: The **cat** sat on the mat

80% of the time [MASK] token → The [MASK] sat on the mat.

10% of the time random token → The **ant sat on the mat.**

10% of the time unchanged token → The **cat sat on the mat.**

Figure 17: BERT: MLM Pre-training

Masked Language Modeling (MLM)

Masked Language Model (MLM): training

Target (1 token):

capital

Loss

Run **backpropagation** to update the weights

Output (14 tokens):

TK1 TK2 TK3 TK4 TK5 TK6 TK7 TK8 TK9 TK10 TK11 TK12 TK13 TK14



Input (14 tokens):

Rome is the [mask] of Italy, which is why it hosts many government buildings.

Figure 18: BERT: MLM Pre-training

BERT: Pre-training tasks

Next Sentence Prediction (NSP)

- Determines if a second sentence follows logically from the first.

BERT: Pre-training tasks

Next Sentence Prediction (NSP)

- Determines if a second sentence follows logically from the first.
- Sentence pair tasks; (sentence A, sentence B)

BERT: Pre-training tasks

Next Sentence Prediction (NSP)

- Determines if a second sentence follows logically from the first.
- Sentence pair tasks; (sentence A, sentence B)
- 50% of the B is the actual sentence, 50% is a random sentence.

BERT: Pre-training tasks

Next Sentence Prediction (NSP)

- Determines if a second sentence follows logically from the first.
- Sentence pair tasks; (sentence A, sentence B)
- 50% of the B is the actual sentence, 50% is a random sentence.
- **Classification setup; *IsNEXT, NotNEXT***

Next Sentence Prediction (NSP)

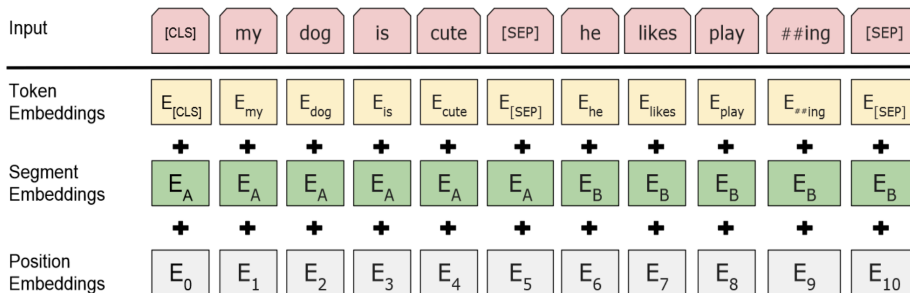


Figure 19: BERT: NSP Pre-training

Next Sentence Prediction (NSP)

Next Sentence Prediction (NSP): training

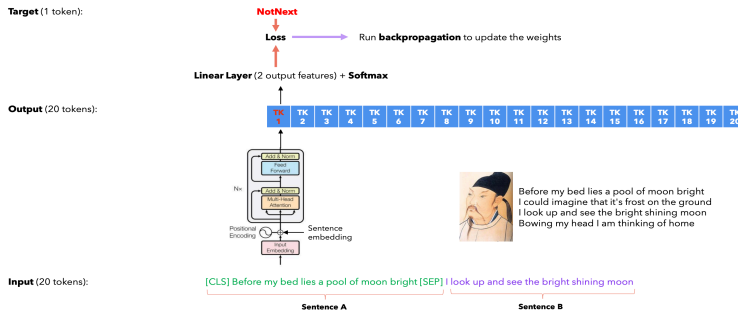


Figure 20: BERT: NSP Pre-training

BERT: Fine-tuning tasks

Text Classification Task

Text Classification: training

My router's led is not working, I tried changing the power socket but still nothing.

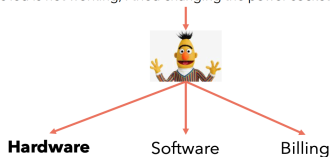


Figure 21: BERT: Fine-tuning for classification task

Text Classification Task

Text Classification: training

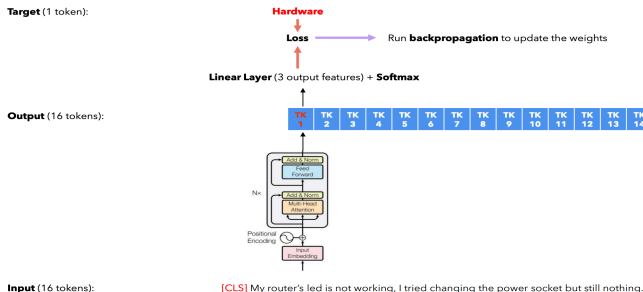


Figure 22: BERT: Fine-tuning for classification task

BERT: Fine-tuning tasks

Question Answering Task

Question Answering: sentence A and B

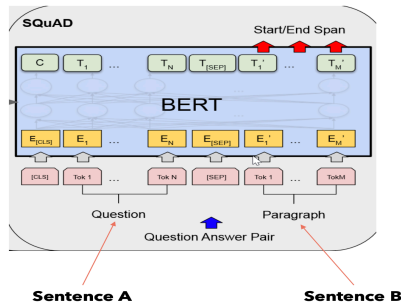


Figure 23: BERT: Fine-tuning for question answering task

Question Answering Task

Question Answering: **start** and **end** positions

Target (1 token):

start=TK10, end=TK10

Loss

Run **backpropagation** to update the weights

Linear Layer (2 output features) + **Softmax**

Output (27 tokens):

TK 1 TK 2 TK 3 ... TK 7 TK 8 TK 9 TK 10 TK 11 TK 12 TK 13 ... TK 22 TK 23 TK 24 TK 25 TK 26 TK 27



Input (27 tokens):

[CLS] What is the fashion capital of China? [SEP] Shanghai is a City in China, it is also a financial center, its fashion capital and industrial city.

Figure 24: BERT: Fine-tuning for question answering task

Experiment & Results

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
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 _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

Figure 25: Experiment and results

Limitations

- Computationally expensive.

Limitations

- Computationally expensive.
- Requires large datasets.

Limitations

- Computationally expensive.
- Requires large datasets.
- Limited understanding of rare or out-of-distribution words

Conclusions

- BERT revolutionized NLP with bidirectional transformers.

Conclusions

- BERT revolutionized NLP with bidirectional transformers.
- Future improvements focus on efficiency (e.g., DistilBERT) and domain-specific adaptations (e.g., BioBERT).

Thank You