


Bidirectional Encoder Representations From Transformers (BERT)

September 6, 2024

Seoul National University

Young rae Cho

Background

A man with a beard and a blue cap is talking to Elmo. The man is wearing a colorful patterned shirt. Elmo is a red Muppet character. They are in a studio setting with a blue background.

Hey ELMo, what's the embedding of the word "stick"?

There are multiple possible embeddings! Use it in a sentence.

Oh, okay. Here:
"Let's stick to improvisation in this skit"

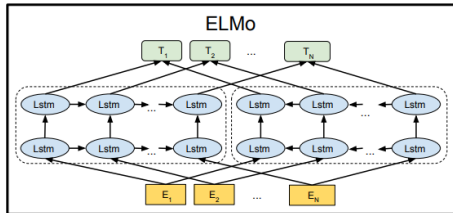
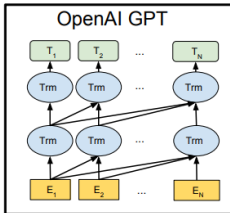
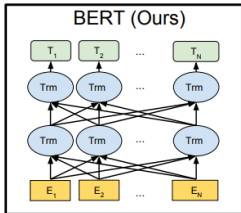
Oh in that case, the embedding is:
-0.02, -0.16, 0.12, -0.1etc

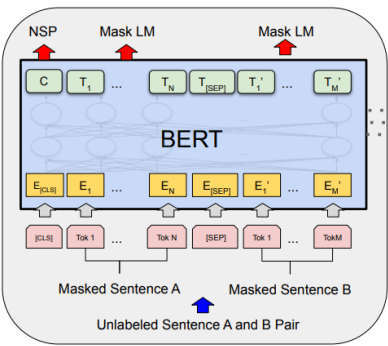
Background

BERT uses a bidirectional Transformer (encoder)

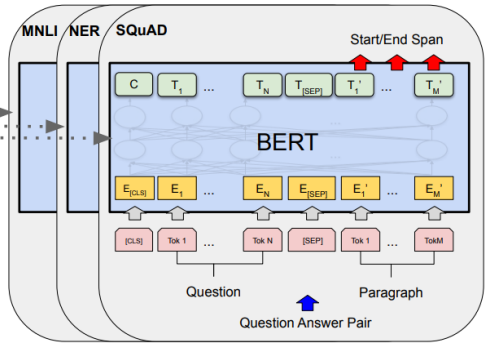
OpenAI GPT uses a left-to-right Transformer (decoder)

ELMo uses the concatenation of independently trained Left-to-right and right-to-left LSTMs.



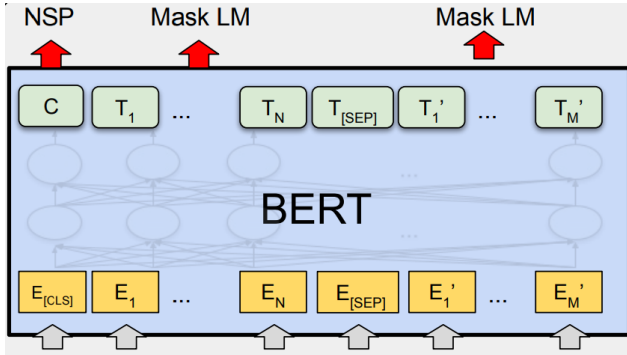


Pre-training



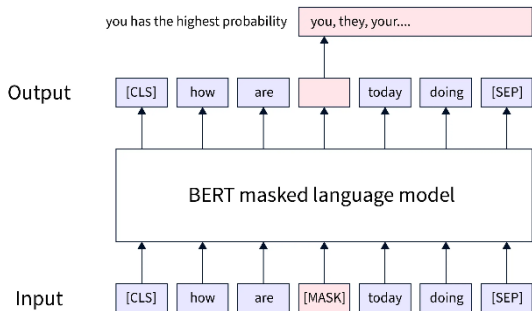
Fine-Tuning

Embedding



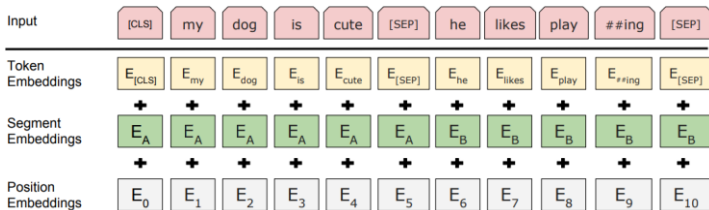
Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{##ing}$	$E_{[SEP]}$
Segment Embeddings	+	+	+	+	+	+	+	+	+	+	+
	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B
Position Embeddings	+	+	+	+	+	+	+	+	+	+	+
	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}

Task #1: Masked LM

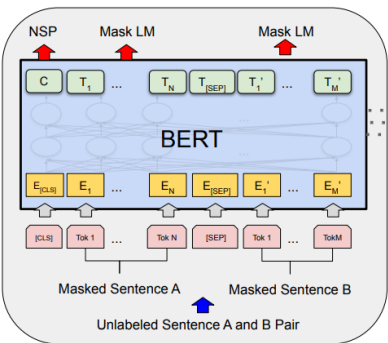


- 80% of the time: Replace the word with the [MASK] token, e.g., my dog is hairy → my dog is [MASK]
- 10% of the time: Replace the word with a random word, e.g., my dog is hairy → my dog is apple
- 10% of the time: Keep the word unchanged, e.g., my dog is hairy → my dog is hairy. The purpose of this is to bias the representation towards the actual observed word.

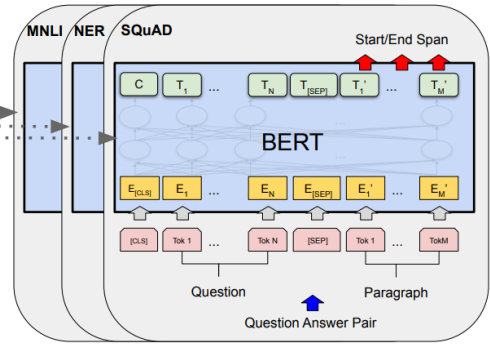
Task #2: Next Sentence Prediction (NSP)



Sentence 1	Sentence 2	Next Sentence?
I have a class	I will be back by 6	✓
I have a class	Zebra is a animal	✗

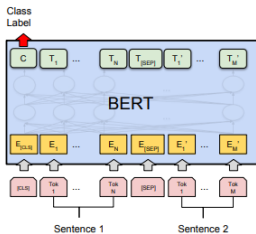


Pre-training

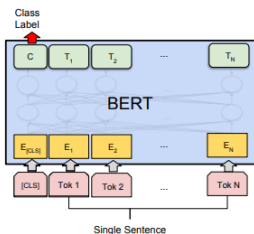


Fine-Tuning

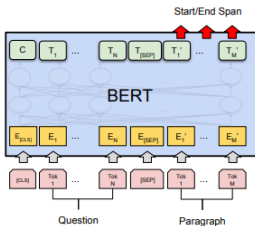
Fine-Tuning examples



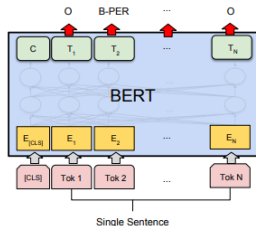
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



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



(c) Question Answering Tasks:
SQuAD v1.1



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

Effect of Pre-training Tasks

Tasks	Dev Set				
	MNLI-m (Acc)	QNLI (Acc)	MRPC (Acc)	SST-2 (Acc)	SQuAD (F1)
BERT _{BASE}	84.4	88.4	86.7	92.7	88.5
No NSP	83.9	84.9	86.5	92.6	87.9
LTR & No NSP	82.1	84.3	77.5	92.1	77.8
+ BiLSTM	82.1	84.1	75.7	91.6	84.9

Table 5: Ablation over the pre-training tasks using the BERT_{BASE} architecture. “No NSP” is trained without the next sentence prediction task. “LTR & No NSP” is trained as a left-to-right LM without the next sentence prediction, like OpenAI GPT. “+ BiLSTM” adds a randomly initialized BiLSTM on top of the “LTR + No NSP” model during fine-tuning.

Effect of Model Size

	Hyperparams				Dev Set Accuracy		
	#L	#H	#A	LM (ppl)	MNLI-m	MRPC	SST-2
Base	3	768	12	5.84	77.9	79.8	88.4
	6	768	3	5.24	80.6	82.2	90.7
	6	768	12	4.68	81.9	84.8	91.3
	12	768	12	3.99	84.4	86.7	92.9
	12	1024	16	3.54	85.7	86.9	93.3
Large	24	1024	16	3.23	86.6	87.8	93.7

Table 6: Ablation over BERT model size. #L = the number of layers; #H = hidden size; #A = number of attention heads. “LM (ppl)” is the masked LM perplexity of held-out training data.

BERT

Only Encoder

Bidirectional LM

Fine - tuning

GPT

Only Decoder

Left to Right LM

No Fine - tuning

Pre-training and Fine tuning model

Bidirectional model

State-of-the-art (SOTA)