

BERT Fine-Tuning

Context

... right ...

... they were on the right ...



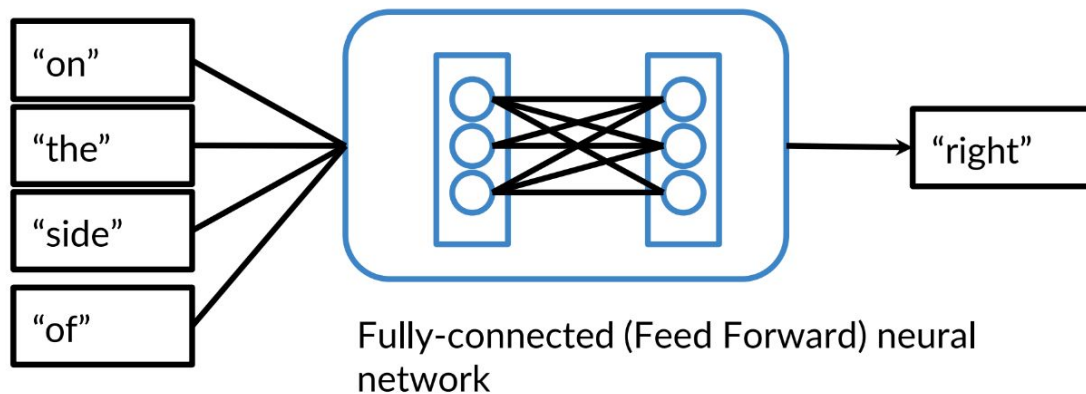
... they were on the right side of the street



Continuous Bag of Words

... they were on the right side of the street

Fixed window Fixed window



Need more context?

... they were on the right side of the street.



Fixed window Fixed window

... they were on the right side of history.

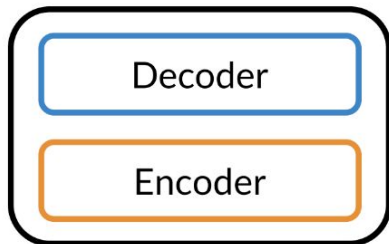
Use all context words

The legislators believed that they were on the right side of history, so they changed the law.



BERT

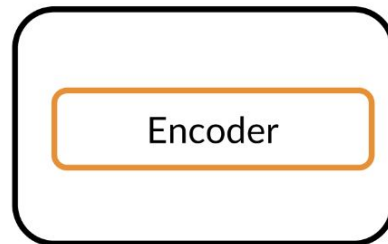
Transformer



GPT



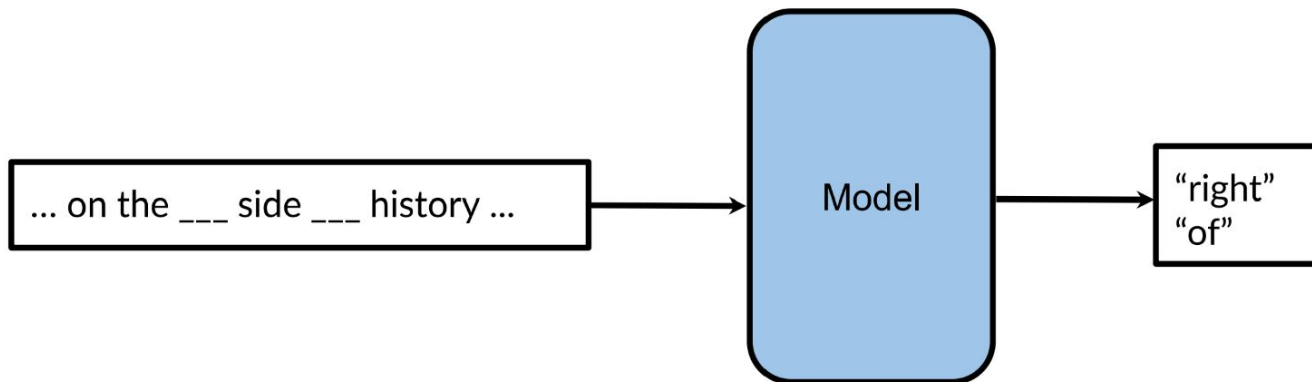
BERT



The legislators believed that they were on the ____ side of history, so they changed the law.

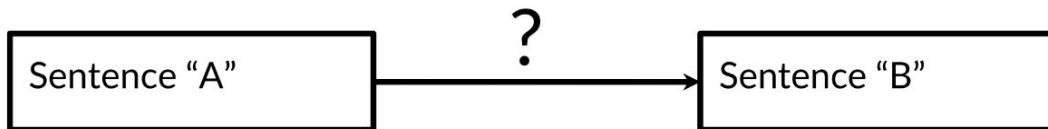
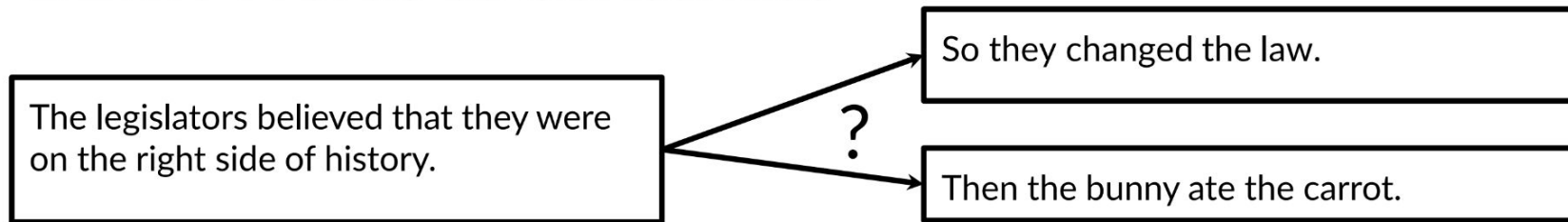
Bi-directional

Transformer + Bi-directional Context



Multi-Mask Language Modeling

BERT: Words to Sentences



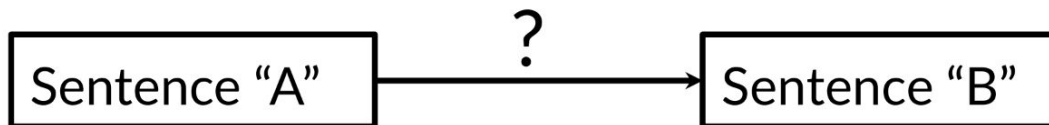
Next Sentence Prediction

BERT Pre-training Tasks

Multi-Mask Language Modeling



Next Sentence Prediction



Bidirectional Encoder Representations from Transformers (BERT)

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

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

Google AI Language

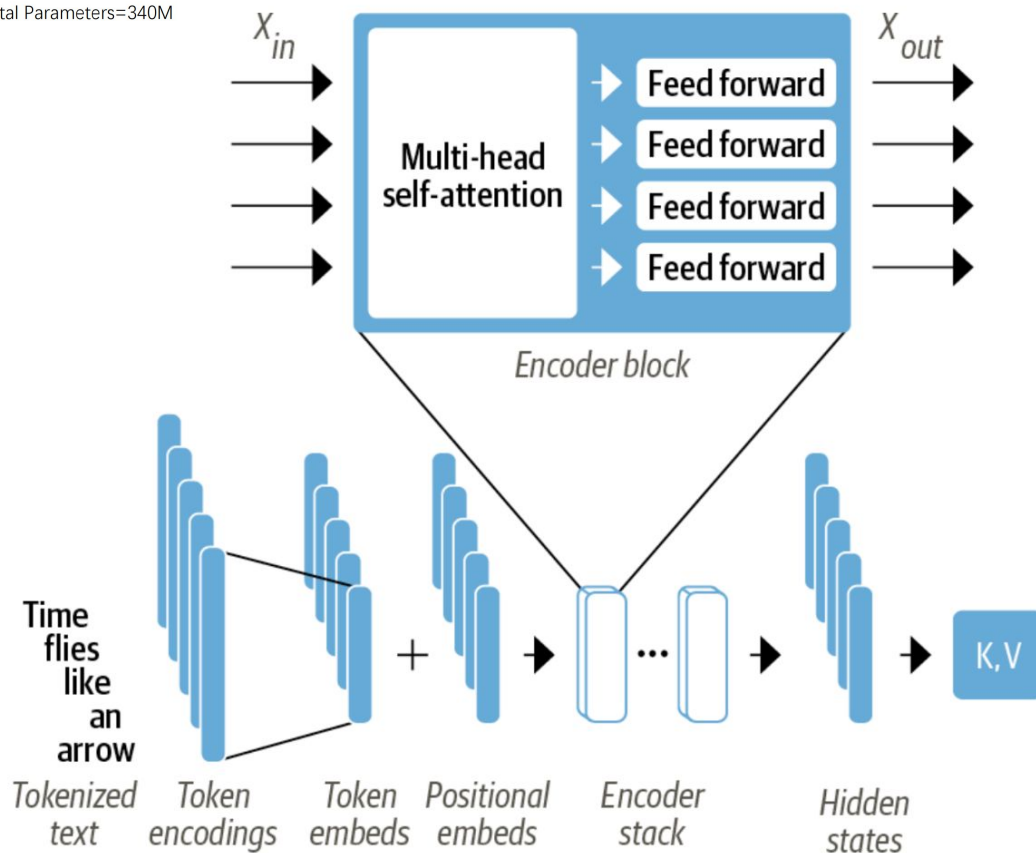
{jacobdevlin, mingweichang, kentonl, kristout}@google.com

[1810.04805.pdf \(arxiv.org\)](#)

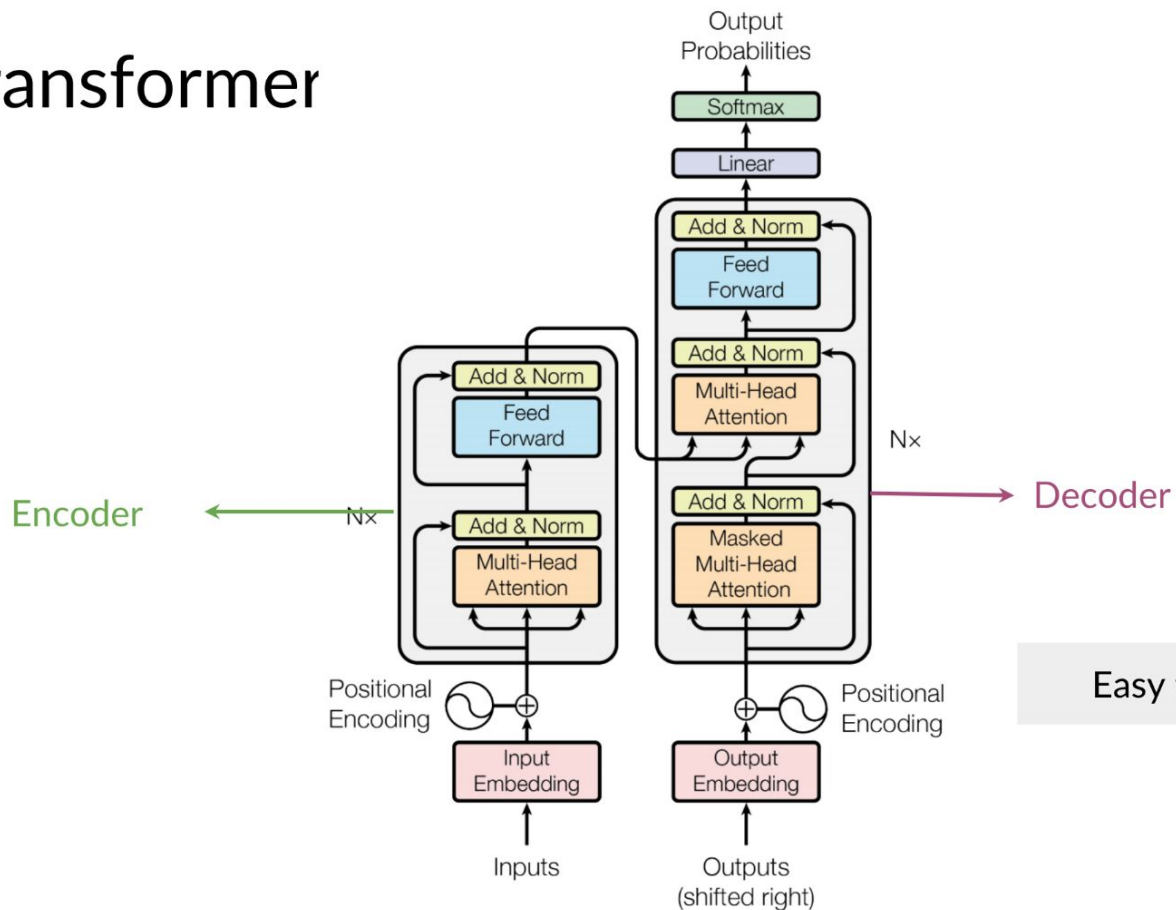
BERT

- A multi layer bidirectional transformer
- Positional embeddings
- BERT_base:
 - 12 layers (12 transformer blocks)
 - 12 attentions heads
 - 110 million parameters

- BERT_{BASE}: $N = 6$, $d_{\text{model}} = 512$, $h = 12$, Total Parameters=110M
- 4 cloud TPUs in Pod configuration (16 TPU chips total)
- BERT_{LARGE}: $N = 24$, $d_{\text{model}} = 1024$, $h = 16$, Total Parameters=340M
- 16 Cloud TPUs (64 TPU chips total)
- Each pretraining took 4 days to complete.



The Transformer



Formalizing the input

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}



Tokenization

DistilBertTokenizer

101	1037	17453	14726	19379	12758	2006	2293	102
-----	------	-------	-------	-------	-------	------	------	-----

↑ 3) substitute tokens with their ids

[CLS]	a	visually	stunning	rum	##ination	on	love	[SEP]
-------	---	----------	----------	-----	-----------	----	------	-------

↑ 2) Add [CLS] and [SEP] tokens

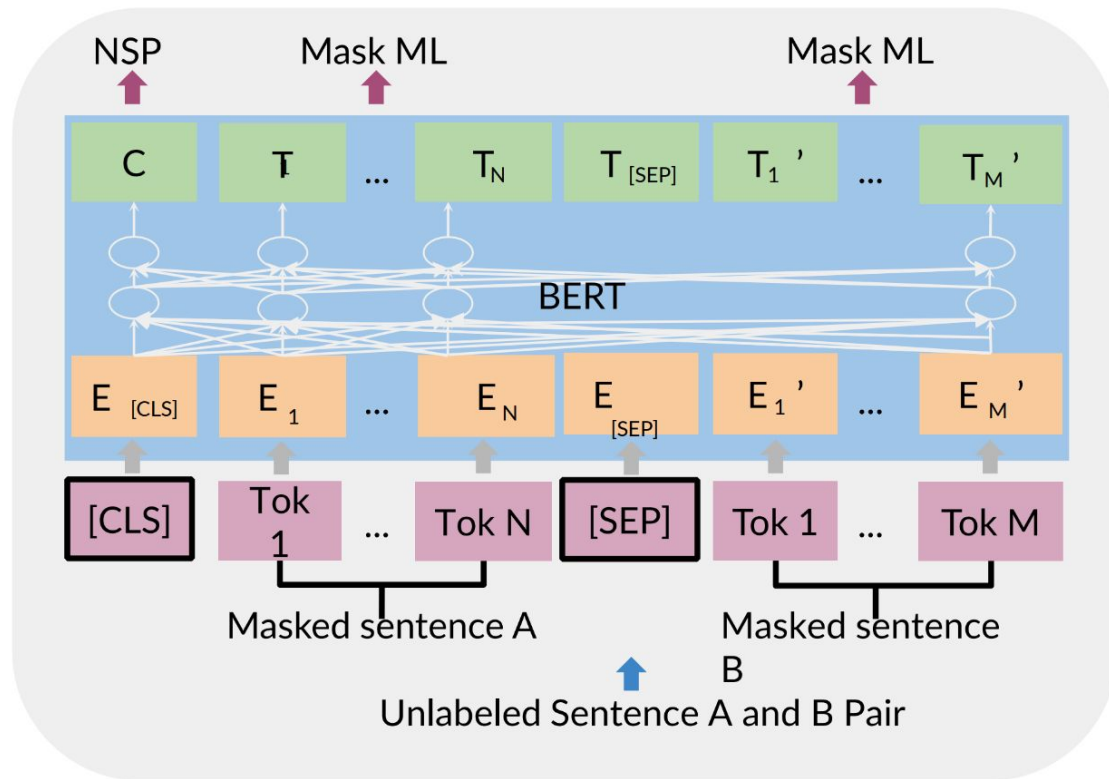
a	visually	stunning	rum	##ination	on	love
---	----------	----------	-----	-----------	----	------

↑ 1) Break words into tokens

↑ Tokenize

“a visually stunning rumination on love”

Visualizing the output

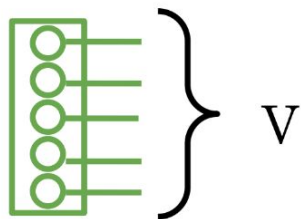


- **[CLS]**: a special classification symbol added in front of every input
- **[SEP]**: a special separator token

BERT Objective

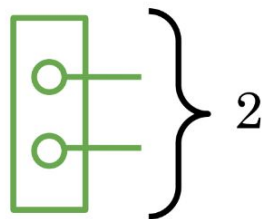
Objective 1:
Multi-Mask LM

Loss: Cross Entropy Loss



Objective 2:
Next Sentence Prediction

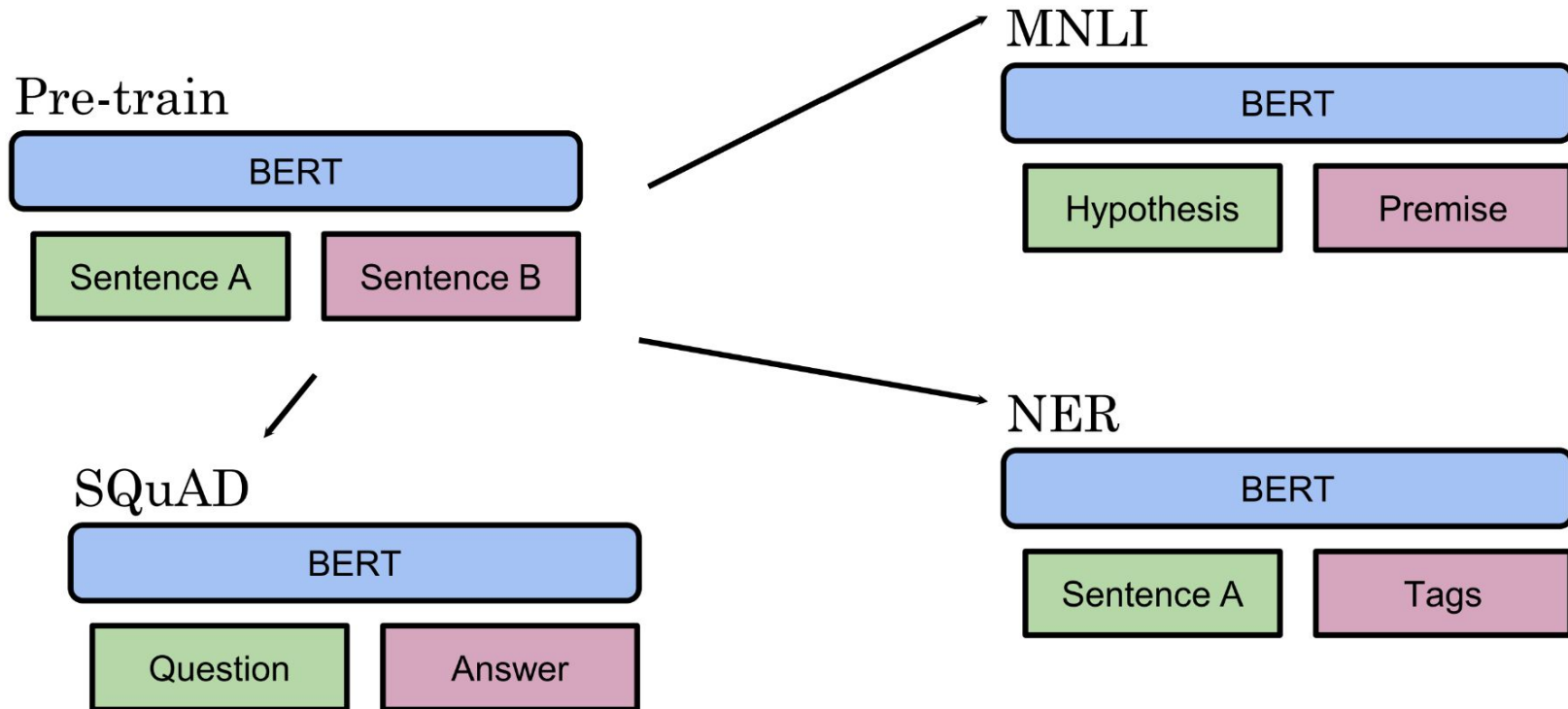
Loss: Binary Loss

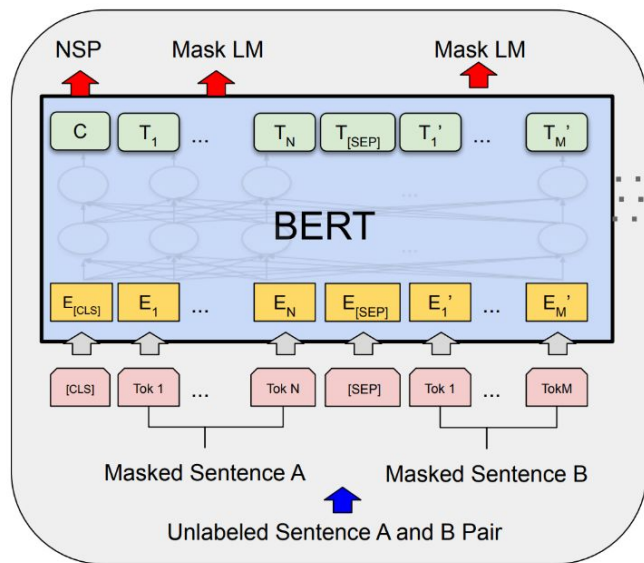


- 80% of the time: Replace the word with the [MASK] token, e.g., my dog is hairy → my dog is [MASK]
- 10% of the time: Replace the word with a random word, e.g., my dog is hairy → my dog is apple
- 10% of the time: Keep the word unchanged, e.g., my dog is hairy → my dog is hairy. The purpose of this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

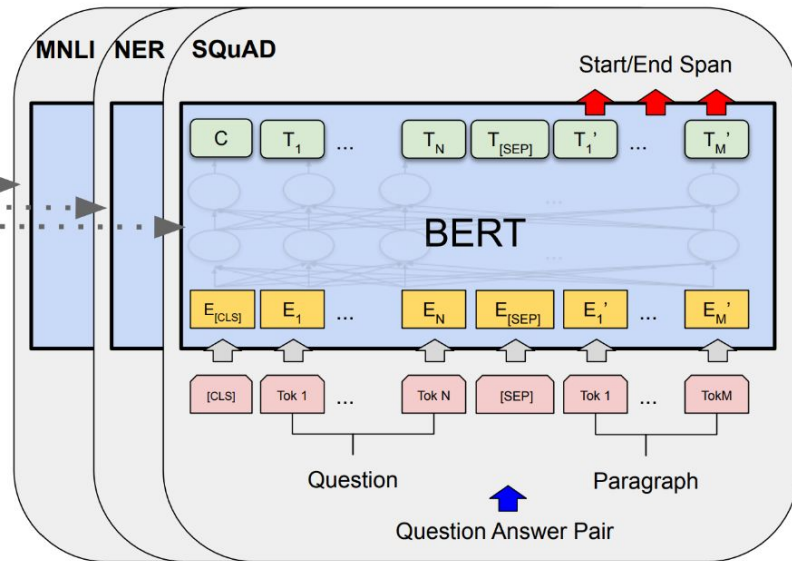
Fine-tuning BERT

Fine-tuning BERT: Outline



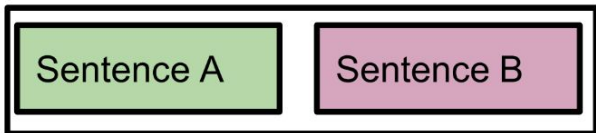


Pre-training



Fine-Tuning

Summary



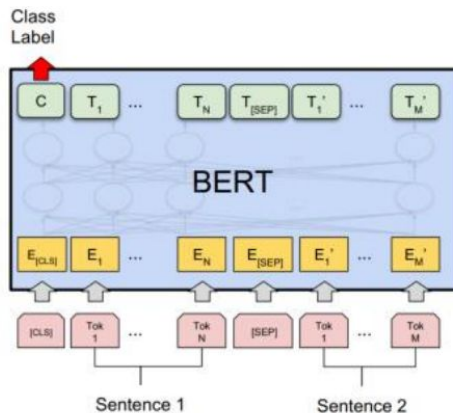
⋮

Fine-tuning with BERT

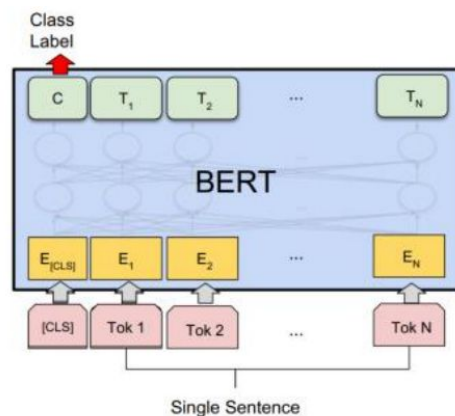
- Context vector C : Take the final hidden state corresponding to the first token in the input: [CLS].
- Transform to a probability distribution of the class labels:

$$P = \text{softmax}(CW^T)$$

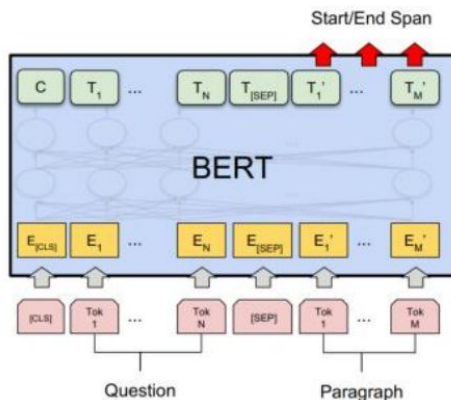
- **Batch size:** 16, 32
- **Learning rate (Adam):** 5e-5, 3e-5, 2e-5
- **Number of epochs:** 3, 4



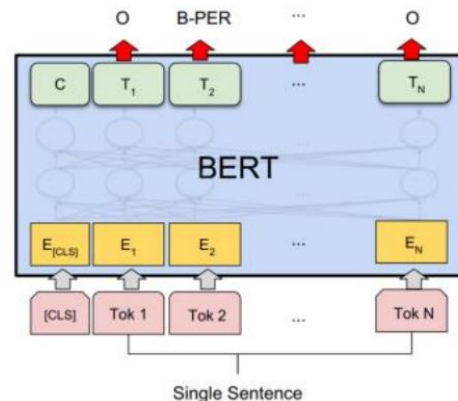
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



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



(c) Question Answering Tasks:
SQuAD v1.1



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

Evaluation for BERT: GLUE

- General Language Understanding Evaluation (**GLUE**) benchmark: Standard split of data to train, validation, test, where labels for the test set is only held in the server.
- Sentence pair tasks
 - **MNLI**, Multi-Genre Natural Language Inference
 - **QQP**, Quora Question Pairs
 - **QNLI**, Question Natural Language Inference
 - **STS-B** The Semantic Textual Similarity Benchmark
 - **MRPC** Microsoft Research Paraphrase Corpus
 - **RTE** Recognizing Textual Entailment
 - **WNLI** Winograd NLI is a small natural language inference dataset
- Single sentence classification
 - **SST-2** The Stanford Sentiment Treebank
 - **CoLA** The Corpus of Linguistic Acceptability

Evaluation for BERT: GLUE

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
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	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Table 1: GLUE Test results, scored by the GLUE evaluation server. The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set. OpenAI GPT = (L=12, H=768, A=12); BERT_{BASE} = (L=12, H=768, A=12); BERT_{LARGE} = (L=24, H=1024, A=16). BERT and OpenAI GPT are single-model, single task. All results obtained from <https://gluebenchmark.com/leaderboard> and <https://blog.openai.com/language-unsupervised/>.

Evaluation on SQUAD

- The Stanford Question Answering Dataset (SQuAD) is a collection of 100k crowdsourced question/answer pairs.

- Input Question:**

Where do water droplets collide with ice crystals to form precipitation?

- Input Paragraph:**

... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. ...

- Output Answer:**

within a cloud

System	Dev		Test	
	EM	F1	EM	F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
#1 Single - nlnet	-	-	83.5	90.1
#2 Single - QANet	-	-	82.5	89.3
Published				
BiDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQA)	84.2	91.1	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

Table in (Devlin *et al.*, 2018)

Evaluation on Named Entity Recognition

- The CoNLL 2003 Named Entity Recognition (NER) dataset. This dataset consists of 200k training words which have been annotated as **Person**, **Organization**, **Location**, **Miscellaneous**, or **Other** (non-named entity).

Jim Hen ##son was a puppet ##eer
I-PER I-PER X O O O X

System	Dev F1	Test F1
ELMo+BiLSTM+CRF	95.7	92.2
CVT+Multi (Clark et al., 2018)	-	92.6
BERT _{BASE}	96.4	92.4
BERT _{LARGE}	96.6	92.8

Table 3: CoNLL-2003 Named Entity Recognition results. The hyperparameters were selected using the Dev set, and the reported Dev and Test scores are averaged over 5 random restarts using those hyperparameters.

