



Group reading

NLP: From BERT to ELECTRA

Shasha Zhou 1.3

目录

CONTENTS

- 1 **Introduction of Natural Language Processing**
- 2 **BERT (Masked LM)**
- 3 **ELECTRA (Replace Token Detection)**
- 4 **Summary and Reflection**



Introduction of NLP

● NLP

	One Sequence	Multiple Sequences
One Class	Sentiment Classification Stance Detection Veracity Prediction...	NLI Search Engine Relation Extraction
Class for each Token	POS tagging Word Segmentation ...	
Copy from Input		Extractive QA
General Sequence	Abstractive Summarization Translation Grammar Correction...	General QA Chatbot

Introduction of NLP

● Word Vector

1. One-hot
2. Distributed Representation(word2vec)

$$w^{Aardvark} = \begin{vmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{vmatrix}$$

$$w^{zebra} = \begin{vmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{vmatrix}$$

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

Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova

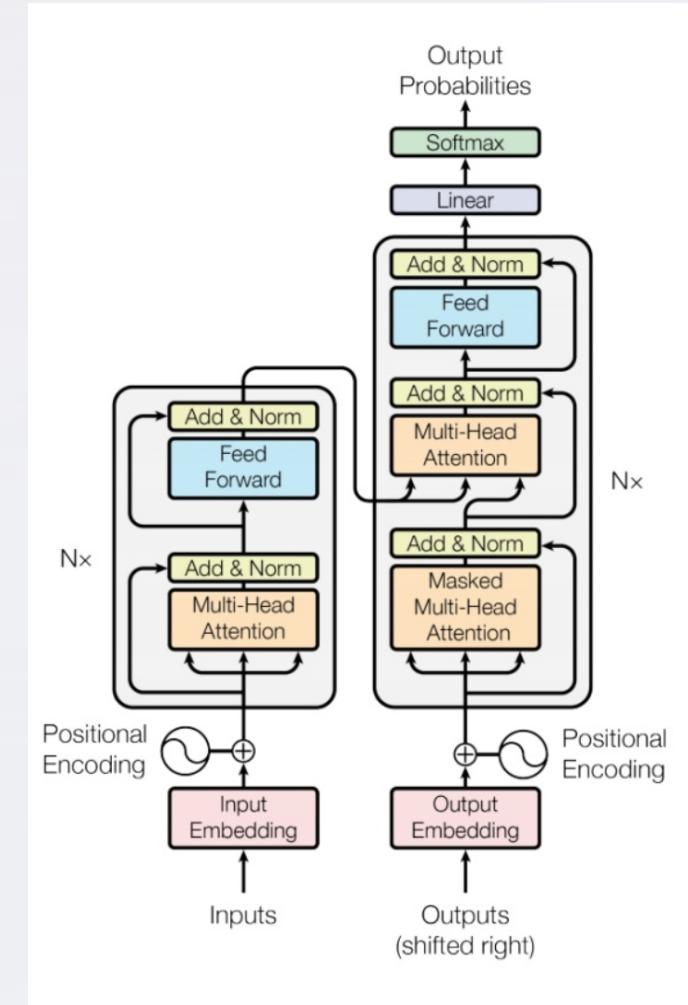
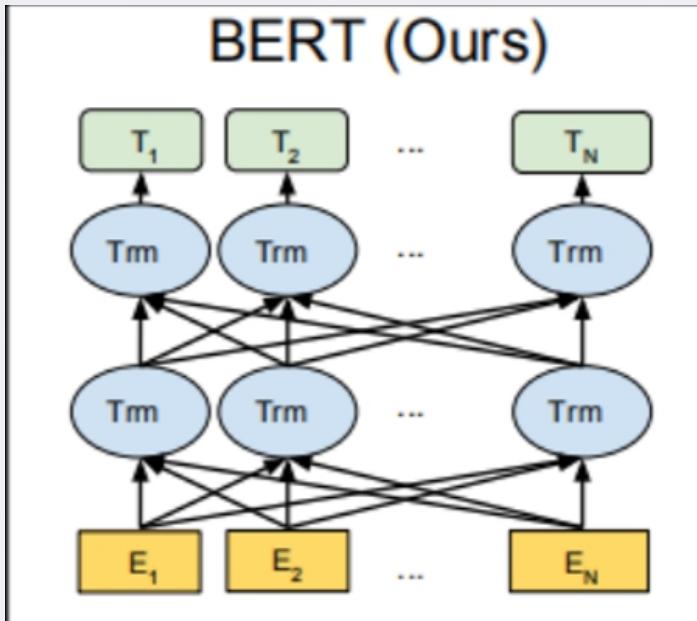
Google AI Language

NAACL-HLT, 2019

● Model Architecture

Multi-head attention

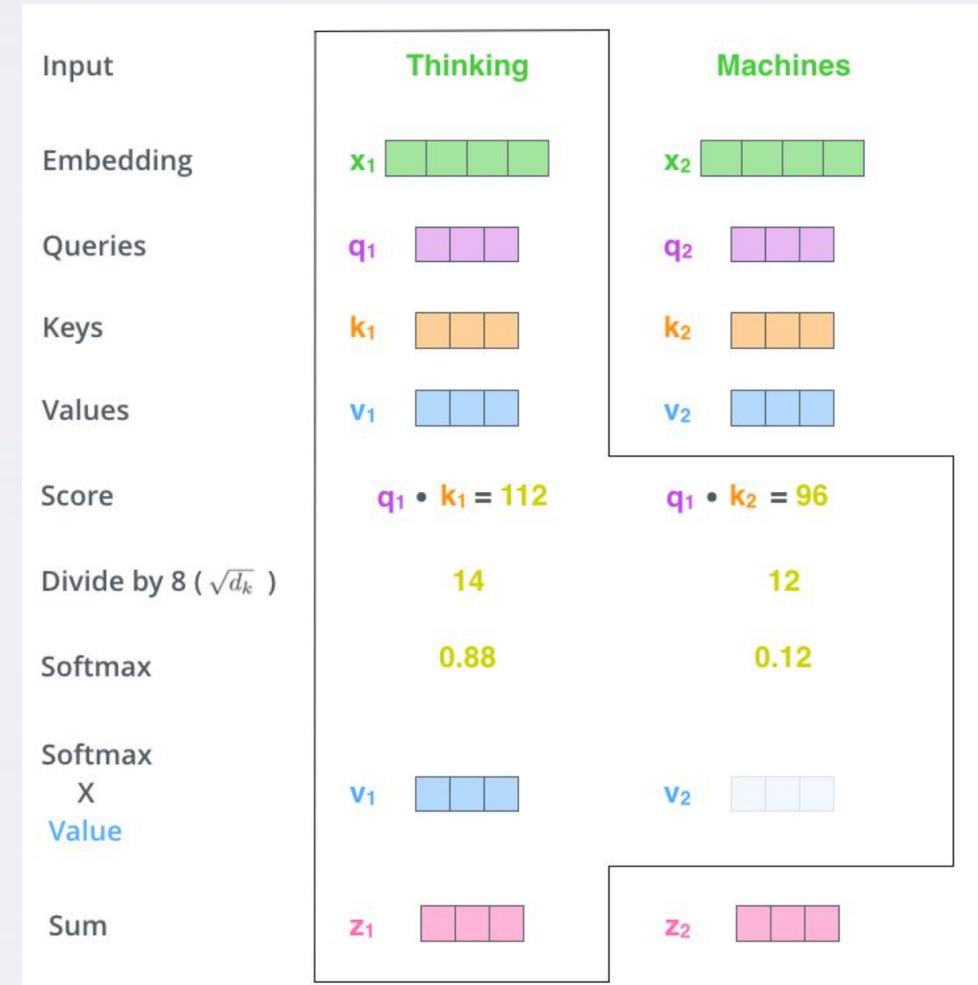
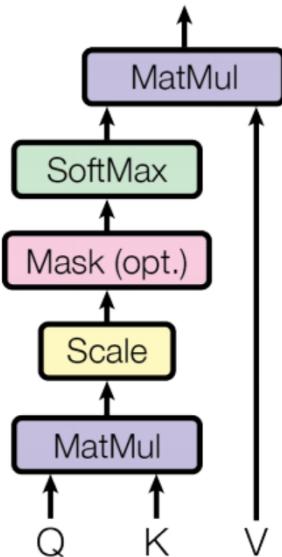
Feed forward



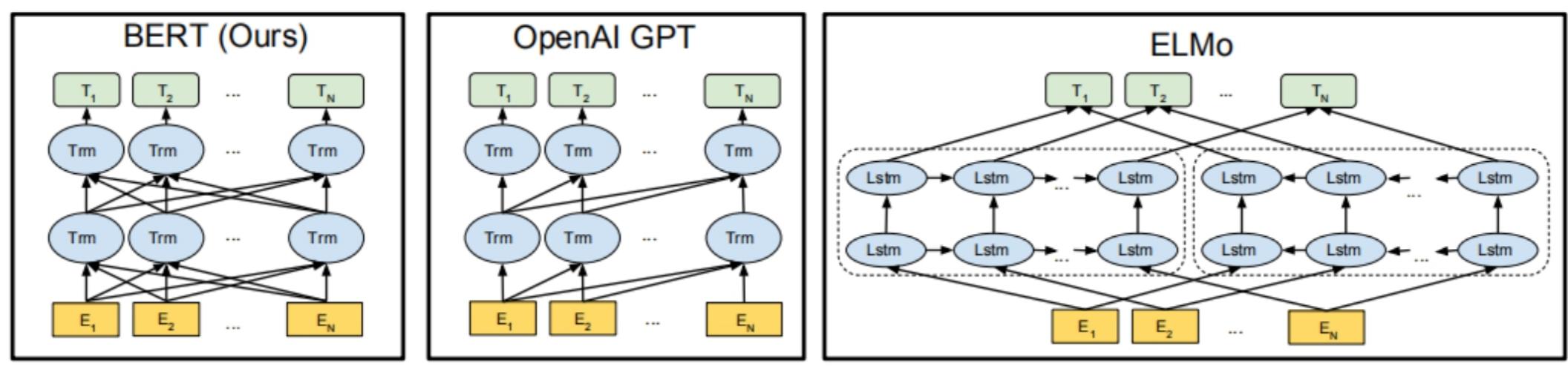
● Self-attention

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

Scaled Dot-Product Attention



- Model Architecture



- **Two steps in BERT**

1. Pre-Training

- Masked LM

- Next Sentence Prediction

2. Fine-Tuning

- **Two steps in BERT**

An example of Masked LM (15% masked):

my dog is hairy

- 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.

- Two steps in BERT

An example of Next Sentence Prediction(50% is next):

Input = [CLS] the man went to [MASK] store [SEP] he bought a gallon [MASK] milk [SEP]

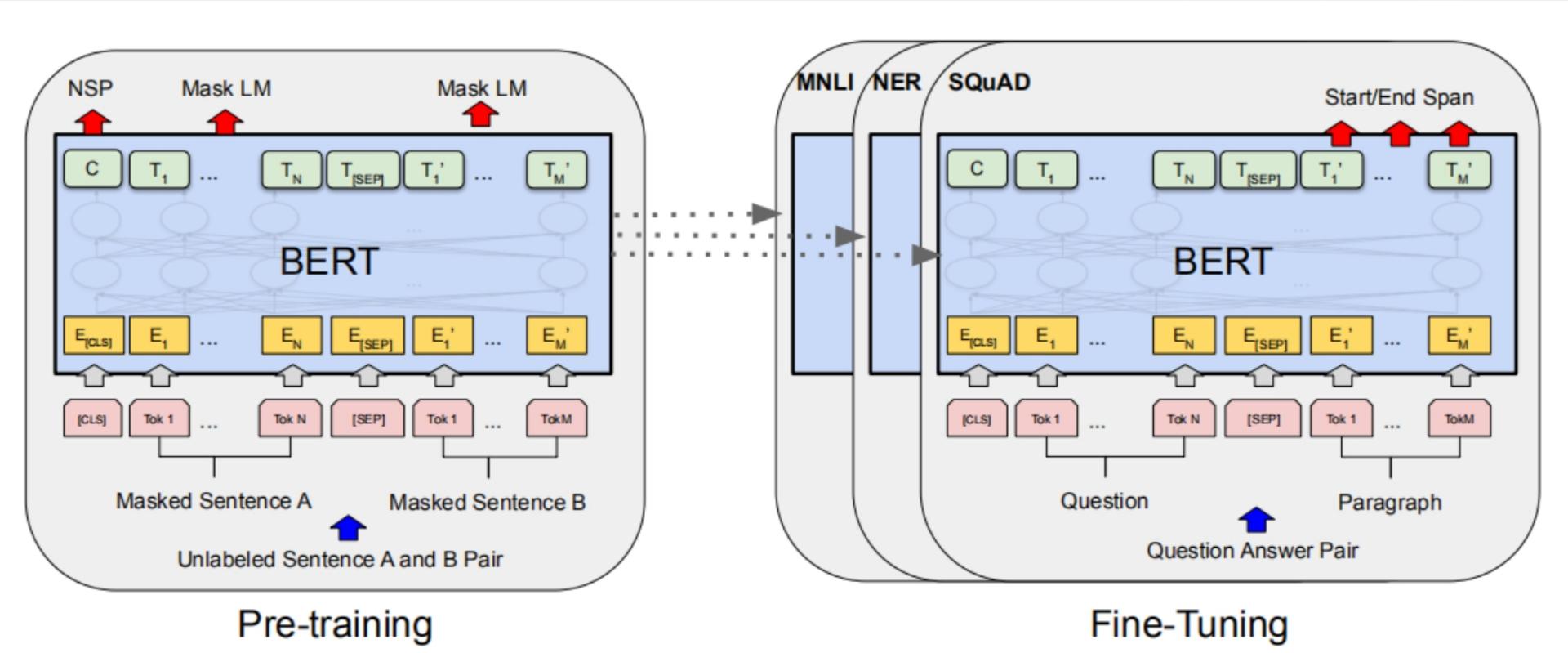
Label = IsNext

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

Label = NotNext

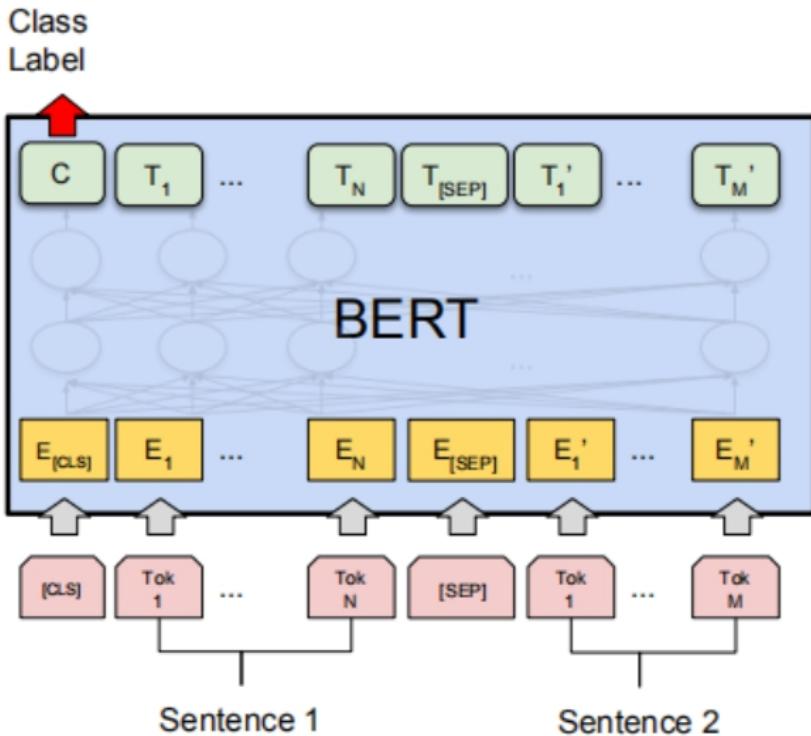
BERT

- Two steps in BERT

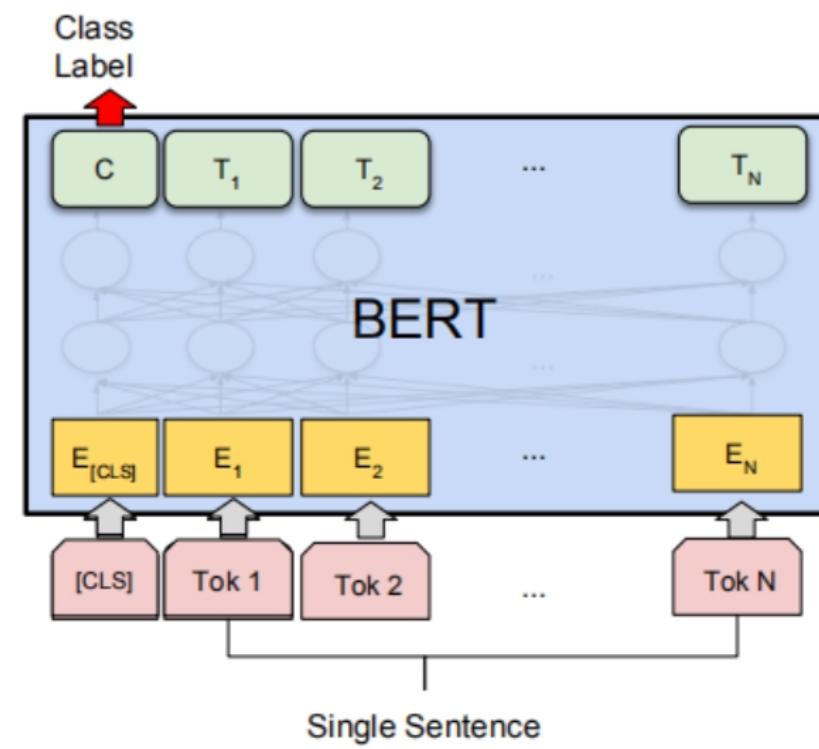


BERT

- Two steps in BERT



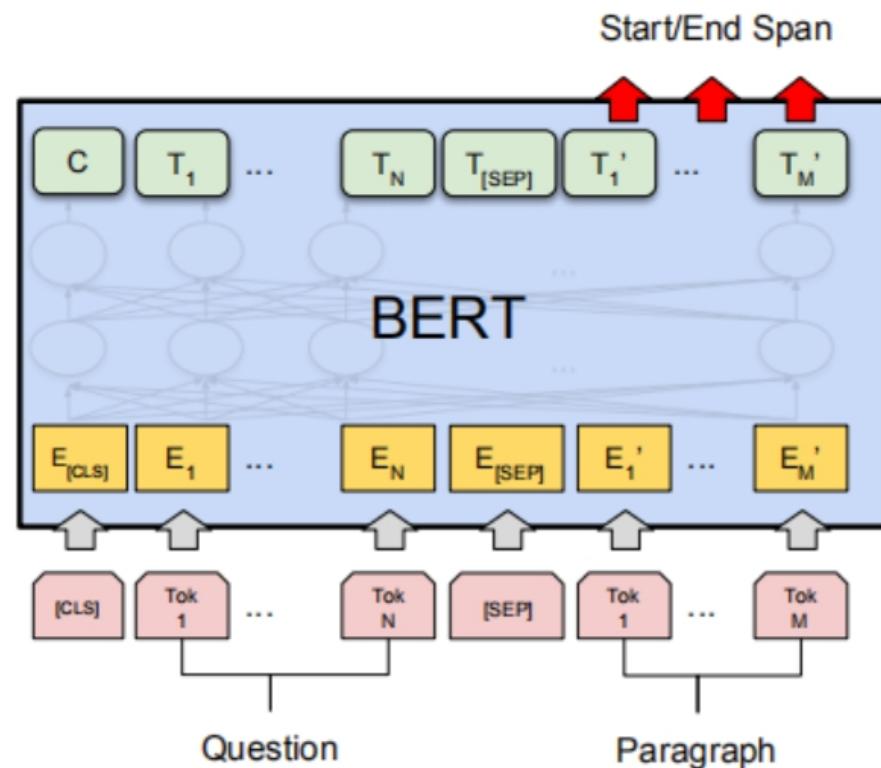
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



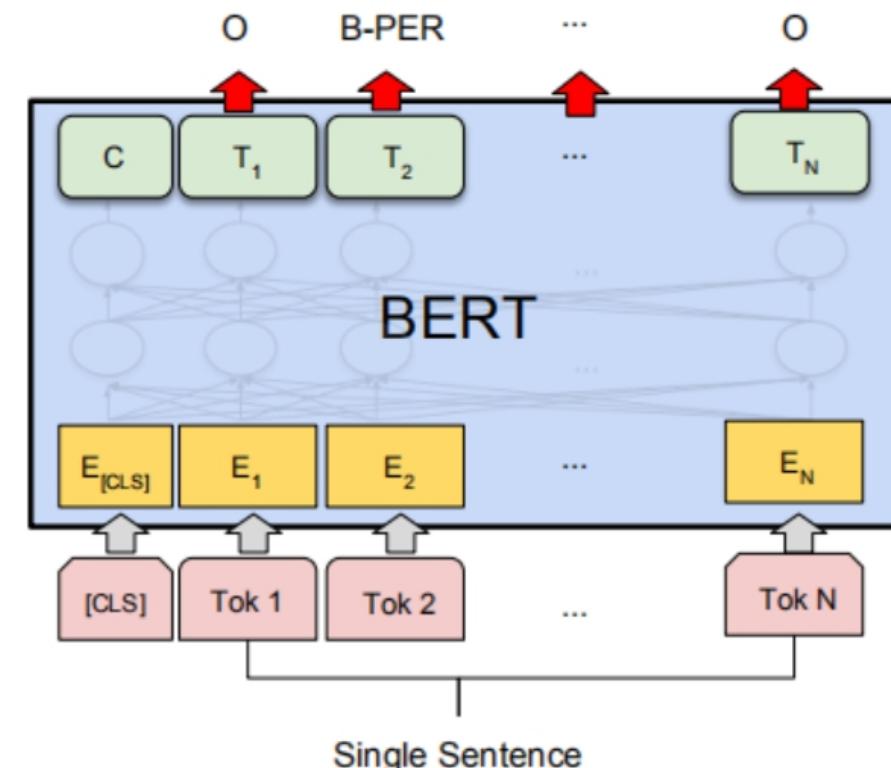
(b) Single Sentence Classification Tasks:
SST-2, CoLA

BERT

- Two steps in BERT



(c) Question Answering Tasks:
SQuAD v1.1



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

- Sentence Classification Tasks Results

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	79.8	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

● Question Answering Task Results

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
Published				
BiDAF+ELMo (Single)	-	85.6	-	85.8
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 1.1 results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	86.3	89.0	86.9	89.5
#1 Single - MIR-MRC (F-Net)	-	-	74.8	78.0
#2 Single - nlnet	-	-	74.2	77.1
Published				
unet (Ensemble)	-	-	71.4	74.9
SLQA+ (Single)	-	-	71.4	74.4
Ours				
BERT _{LARGE} (Single)	78.7	81.9	80.0	83.1

Table 3: SQuAD 2.0 results. We exclude entries that use BERT as one of their components.

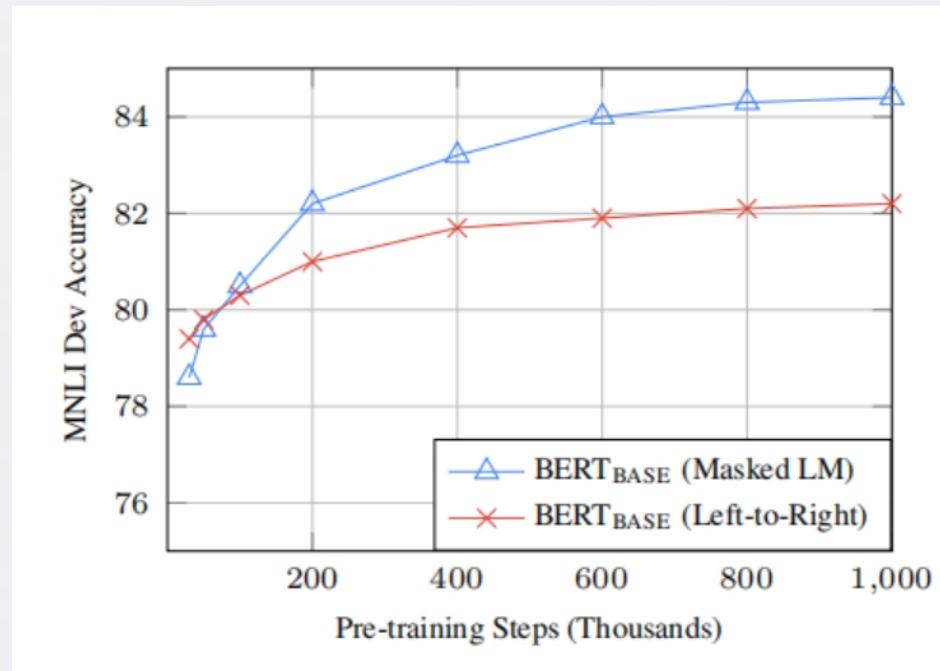
- **Situations With Adversarial Generations(SWAG) Results**

System	Dev	Test
ESIM+GloVe	51.9	52.7
ESIM+ELMo	59.1	59.2
OpenAI GPT	-	78.0
BERT _{BASE}	81.6	-
BERT _{LARGE}	86.6	86.3
Human (expert) [†]	-	85.0
Human (5 annotations) [†]	-	88.0

Table 4: SWAG Dev and Test accuracies. [†]Human performance is measured with 100 samples, as reported in the SWAG paper.

- The Shortcomings of BERT

1. More pre-training steps may be required for the model to converge.
2. BERT need a large amount of pre-training (128,000words/batch * 1,000,000 steps)

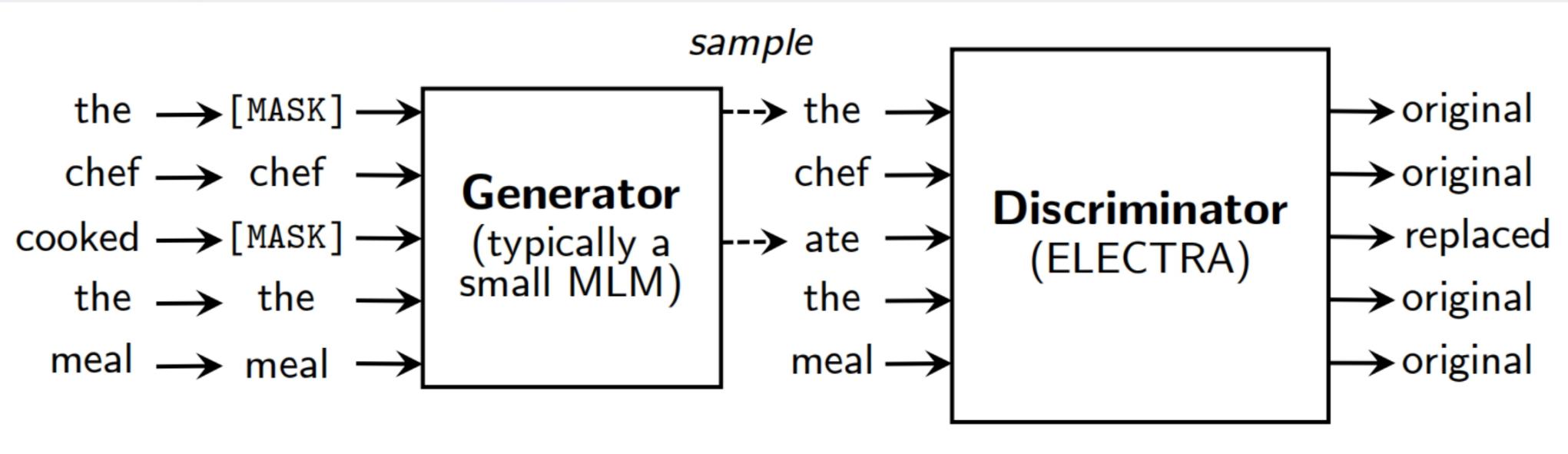


ELECTRA: PRE-TRAINING TEXT ENCODERS AS DISCRIMINATORS RATHER THAN GENERATORS

Kevin Clark, Minh-Thang Luong, Quoc V. Le

ICLR2020

ELECTRA



$$m_i \sim \text{unif}\{1, n\} \text{ for } i = 1 \text{ to } k$$

$$\hat{x}_i \sim p_G(x_i | \mathbf{x}^{\text{masked}}) \text{ for } i \in \mathbf{m}$$

$$\mathbf{x}^{\text{masked}} = \text{REPLACE}(\mathbf{x}, \mathbf{m}, [\text{MASK}])$$

$$\mathbf{x}^{\text{corrupt}} = \text{REPLACE}(\mathbf{x}, \mathbf{m}, \hat{\mathbf{x}})$$

- Loss Function

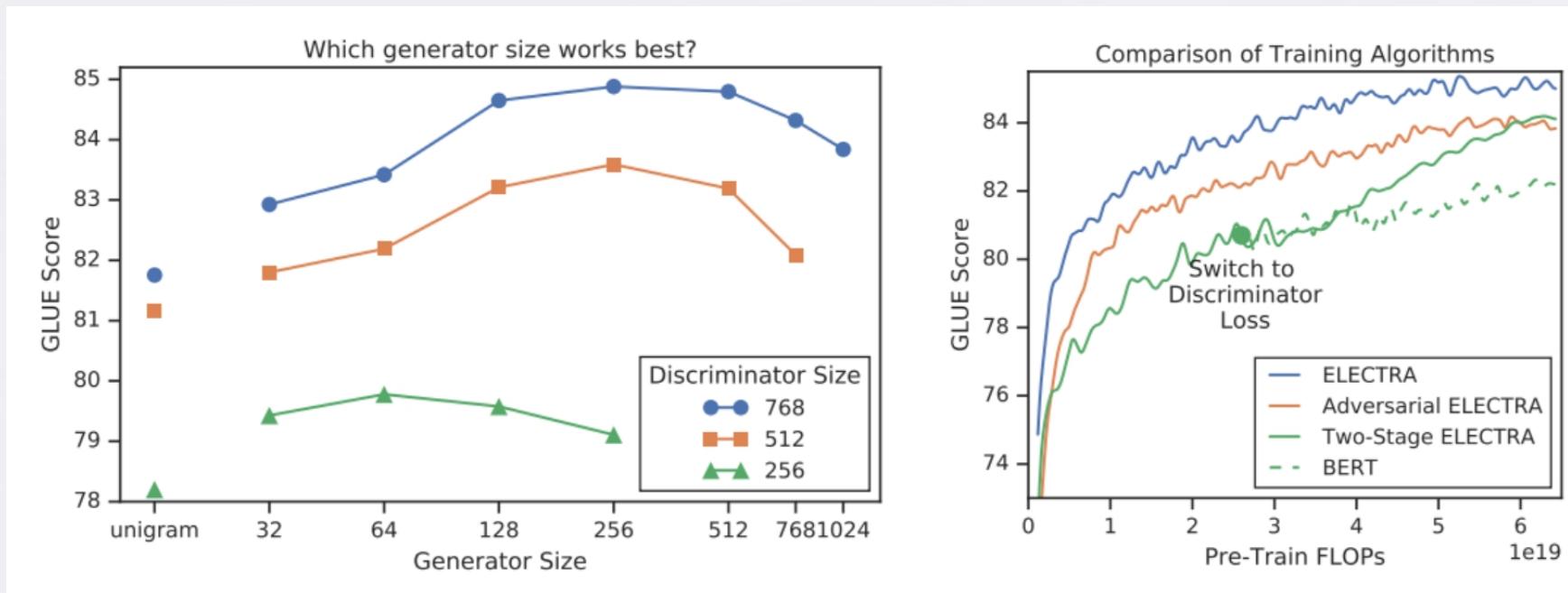
$$\begin{aligned}\mathcal{L}_{\text{MLM}}(\boldsymbol{x}, \theta_G) &= \mathbb{E} \left(\sum_{i \in \mathbf{m}} -\log p_G(x_i | \boldsymbol{x}^{\text{masked}}) \right) \\ \mathcal{L}_{\text{Disc}}(\boldsymbol{x}, \theta_D) &= \mathbb{E} \left(\sum_{t=1}^n -\mathbb{1}(x_t^{\text{corrupt}} = x_t) \log D(\boldsymbol{x}^{\text{corrupt}}, t) - \mathbb{1}(x_t^{\text{corrupt}} \neq x_t) \log(1 - D(\boldsymbol{x}^{\text{corrupt}}, t)) \right)\end{aligned}$$

$$\min_{\theta_G, \theta_D} \sum_{\boldsymbol{x} \in \mathcal{X}} \mathcal{L}_{\text{MLM}}(\boldsymbol{x}, \theta_G) + \lambda \mathcal{L}_{\text{Disc}}(\boldsymbol{x}, \theta_D)$$

- **Differences Between ELECTRA And GAN**

1. if the generator happens to generate the correct token, that token is considered “real” instead of “fake”;
2. the generator is trained with maximum likelihood rather than being trained adversarially to fool the discriminator;
3. we do not supply the generator with a noise vector as input, as is typical with a GAN.

Result



- Result

Model	Train / Infer FLOPs	Speedup	Params	Train Time + Hardware	GLUE
ELMo	3.3e18 / 2.6e10	19x / 1.2x	96M	14d on 3 GTX 1080 GPUs	71.2
GPT	4.0e19 / 3.0e10	1.6x / 0.97x	117M	25d on 8 P6000 GPUs	78.8
BERT-Small	1.4e18 / 3.7e9	45x / 8x	14M	4d on 1 V100 GPU	75.1
BERT-Base	6.4e19 / 2.9e10	1x / 1x	110M	4d on 16 TPUv3s	82.2
ELECTRA-Small	1.4e18 / 3.7e9	45x / 8x	14M	4d on 1 V100 GPU	79.9
50% trained	7.1e17 / 3.7e9	90x / 8x	14M	2d on 1 V100 GPU	79.0
25% trained	3.6e17 / 3.7e9	181x / 8x	14M	1d on 1 V100 GPU	77.7
12.5% trained	1.8e17 / 3.7e9	361x / 8x	14M	12h on 1 V100 GPU	76.0
6.25% trained	8.9e16 / 3.7e9	722x / 8x	14M	6h on 1 V100 GPU	74.1
ELECTRA-Base	6.4e19 / 2.9e10	1x / 1x	110M	4d on 16 TPUv3s	85.1

- Result

Model	Train FLOPs	Params	CoLA	SST	MRPC	STS	QQP	MNLI	QNLI	RTE	Avg.
BERT	1.9e20 (0.27x)	335M	60.6	93.2	88.0	90.0	91.3	86.6	92.3	70.4	84.0
RoBERTa-100K	6.4e20 (0.90x)	356M	66.1	95.6	91.4	92.2	92.0	89.3	94.0	82.7	87.9
RoBERTa-500K	3.2e21 (4.5x)	356M	68.0	96.4	90.9	92.1	92.2	90.2	94.7	86.6	88.9
XLNet	3.9e21 (5.4x)	360M	69.0	97.0	90.8	92.2	92.3	90.8	94.9	85.9	89.1
BERT (ours)	7.1e20 (1x)	335M	67.0	95.9	89.1	91.2	91.5	89.6	93.5	79.5	87.2
ELECTRA-400K	7.1e20 (1x)	335M	69.3	96.0	90.6	92.1	92.4	90.5	94.5	86.8	89.0
ELECTRA-1.75M	3.1e21 (4.4x)	335M	69.1	96.9	90.8	92.6	92.4	90.9	95.0	88.0	89.5

Summary

1. BERT uses bidirectional pre-training and masked LM to achieve good training effect in 11 NLP tasks.
2. ELECTRA inspired by GAN, and propose a more sample-efficient pre-training task called replaced token detection.
3. Compared to MLM, ELECTRA's pre-training objective is more compute-efficient and results is better.

Reflection

1. The discriminator in ELECTRA is binary, it maybe not suit difficult NLP tasks.
2. ELECTRA needs to train two small ‘BERT’, so that it may need more memory.