# ELECTRA

## ABSTRACT

## INTRODUCTION

Learning Bidirectional Representation    BERT                                         .
          BERT   Pre-training   MLM    NSP                                              .
MLM                (15%)   [MASK]                     (       Noise Token)                Encoder                 ,
Classifier(Linear Layer)          Noise Token          (Denoising;                        )
          Encoder   **Denoising AutoEncoder**              .
                                                                                                .
                                                  (Representation)              ,
        **Representation Learning**                  .

BRET    MLM task                        15%                                                            **BERT**
**Substantial Computing Cost(                )**              .


                15%                                                      (Original  Token  +  Noise  Token)
Noise  Token                              Loss                                        .
                                                              .

                              **Replaced Token Detection**              .
        MLM                                                    .
        **Computationally Efficient(              )**      .
                ,                                        .


Replaced  Token  Detection                      Masked  Token                      Synthetically  Generated  Token(SGT
Token)
          (Real;Not  Replaced)                      (Fake;Replaced)                                            .
SGT  Token      Generator                MLM                                  .
      MLM                                        .

                                            .


We  call  our  approach  ELECTRA  for  " Efficiently  Learning  an  Encoder  that  Classifies  Token  Replaceme
nts  Accurately."
                                                                    .

                                    (Replaced Token Detection)

BERT          SOTA

- 
- 

            .


## METHOD

**Replaced Token Detection Task**

Replaced Token Detection Task                                                    .
                    Generator    Discriminator                  ,  Generator                           [MASK]               SGT
         (MLM  Task)
Discriminator    SGT                                                 Fake  or  Real        (Binary  Classification  Task)
       .

                         Generator    Discriminator                                                        Representation  V
ector                     .

         Generator   Discriminator                              GAN                      .

1. Generator                              Noise Input

   Electra     Generator                                                        .

2. Generator    Discriminator

                                              (Maximum  Likelihood)                    .

3. Maximum Likelihood Estimation vs Adversarial Learning

   Kullback-Leibler  Divergence(          KL-Div)        Jensen-Shannon  Divergence(         JSD)
        .

   KL-Div

   KL-Div                                                          .

                                                     .
                                /                                               MLE(Maximum
   Likelihood  Estimation)                .

                (Entropy)        .
   ,  E(-log(p_real)  +  log(q_model))   ,  where  E()  is  Expectation  (i.e  KL(P||Q))
        (              )                                          .

   JSD

   JSD                                       /                          .
         ,                                                                   .
                           (          P,Q),  P      Q                            Q      P                      .
       1/2*KL(P||(P+Q)/2)+1/2*KL(Q||(P+Q)/2)                          .  (i.e  JSD(P||Q)  =  JSD(Q||P))

## Objective Function(Loss Function)

ELECTRA                    Generator   Discriminator   Cross Entropy                                    .

Generator              [MASK]                                                                    ,
Discriminator          SGT      Real/Fake                                         .

Generator   Discriminator              Maximum  Likelihood                    (MLE)      ,
                                                            .
                                                      KL-Div              .
Cross  Entropy     KL-Div              .
              One-Hot.(i.e  All  or  None)      ,  Cross  Entropy                  .(
      )

      ,

## Generator Loss

Generator                                          Representation Vector              .
                                                    .

$$p_G(x_t|\boldsymbol{x}) = \exp\left(e(x_t)^T h_G(\boldsymbol{x})_t\right) / \sum_{x'} \exp\left(e(x')^T h_G(\boldsymbol{x})_t\right)$$

x,  x'              input  token              .
h_G()     Generator              .
e()     Embedding              .

Generator            Look-up  Table              .
          ,              (        )                              Embedding  Matrix              .
Embedding                        Embedding  Vector
          Generator            Representation  Vector
    ,  Header  Network                Embedding  Maxtrix              .

h_G(x)   ,  vocab_size  x  len(embedding  vector)
e(x)   ,  1  x  len(embedding  vector)

      e(x)*h_G(x)^T   ,  1  x  vocab_size
              Embedding  Vector     Eembedding  Matrix(i.e  Look-Up  Table)                  ,
                                                .(                        ,                    .)

                                            Softmax              ,                    .

$$\mathcal{L}_{\text{MLM}}(\boldsymbol{x}, \theta_G) = \mathbb{E}\left(\sum_{i\in\boldsymbol{m}} -\log p_G\left(x_i|\boldsymbol{x}^{\text{masked}}\right)\right)$$

                                            ,                              .
          KL-Div                    One-Hot            KL-Div                    Cross  Entropy
      .

## Discriminator Loss

Discriminator              Fake/Real                        .                .

$$D(\boldsymbol{x}, t) = \text{sigmoid}(w^T h_D(\boldsymbol{x})_t)$$

$$\mathcal{L}_{\text{Disc}}(\boldsymbol{x}, \theta_D) = \mathbb{E}\left(\sum_{t=1}^{n} -\mathbb{1}(x_t^{\text{corrupt}} = x_t) \log D(\boldsymbol{x}^{\text{corrupt}}, t) - \mathbb{1}(x_t^{\text{corrupt}} \neq x_t) \log(1 - D(\boldsymbol{x}^{\text{corrupt}}, t))\right)$$

Discriminator                                    Generator          Softmax    CrossEntropy
(Binary)              .

## We minimize the combined loss

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

Generator       Discriminator  Loss                                          .
lambda                     Discriminator                  ,
Generator    Task(          )    Discriminator    Task(          )    ,                              Entropy
.
Entropy                                                  ,
.
BERT       Vocab_size    30000                      ,  30000                                            ,
entropy                                  .
,  Discriminator                Weight        Hyper  Parameter                              .

Method                    ,

- Replaced  Token  Detection  Task                              .
- Maximun  Likelihood  Estimation                    .
- Discriminator  Loss                              .
.

## Experiment

### Experimental Setup

(          )

General  Language  Understanding  Evaluation(GLUE)      Stanford  Question  Answering  (SQuAD)    dataset
.

.
GLUE
[[https://huffon.github.io/2019/11/16/glue/]]
SQuAD
[[https://happygrammer.github.io/nlp/dataset/]]

BERT

GLUE
Simple  Linear  Classifier      Header  Network

SQuAD
XL-NET       Header  Network
BERT                          Header  Network(Independently  predict)              Jointly  predict              .
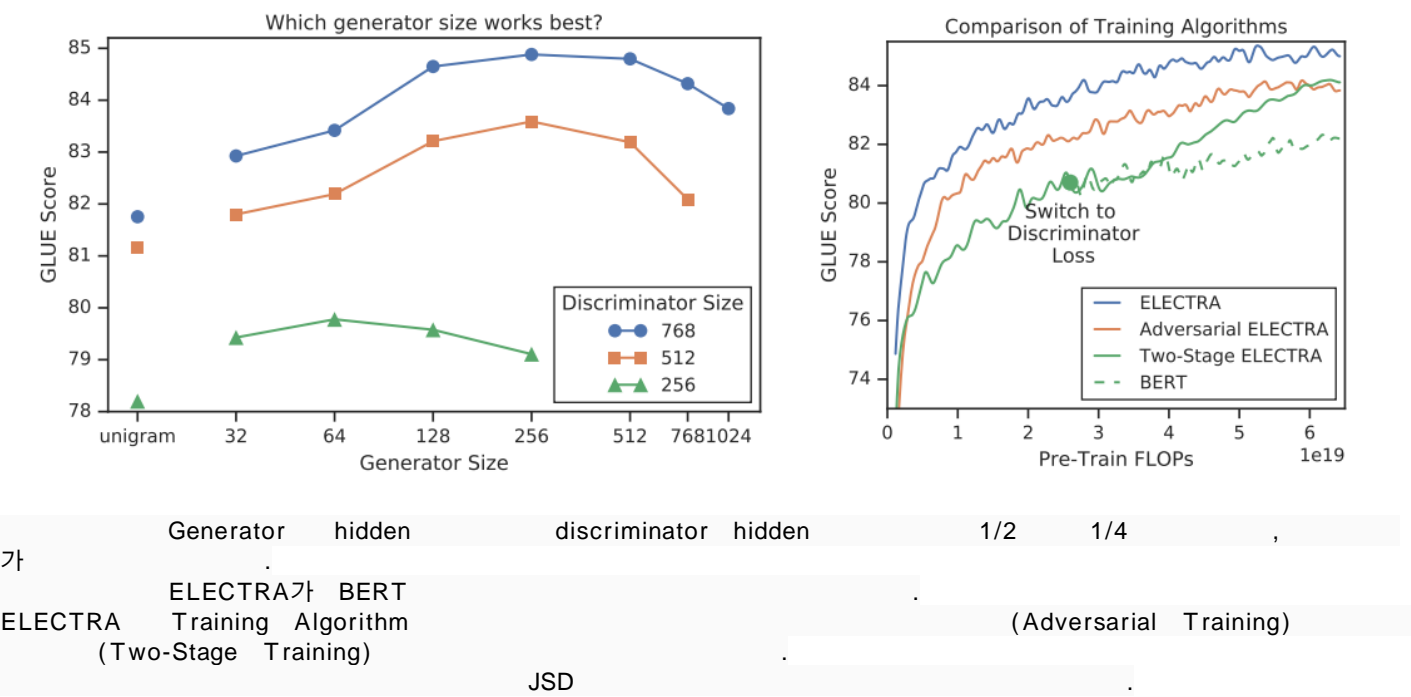
(Model Extentions)

ELECTRA                                        BERT-base                              .

.

1. Weight Sharing
   - Discriminator       Generator       Embedding Layer                            .(i.e Embedding Layer   Look-Up table       )

2. Smaller Generator
   - ELECTRA                    ,  Generator     Discriminator                                        BERT
         .
         BERT                                                                    .
         ,  Generator     Layer                                    .
   - Generator      Unigram  (Character  Level)
   -                              500k   step

3. Training Algorithm(Two-Step Training)
   - Generator     Discriminator                            Generator          n-step                   ,
                       Discriminator                          n-step          (    ,  Generator
                   )
   -                         .

   SMALL MODEL    LARGE MODEL                .

## SMALL MODEL

*As a goal of this work is to improve the efficiency of pre-training, we develop a small model that can be quickly trained on a single GPU.*

SMALL                                                                                    .



             Generator      hidden              discriminator   hidden             1/2       1/4                    ,
                  .
             ELECTRA     BERT                                                    .
ELECTRA     Training Algorithm                                            (Adversarial  Training)
      (Two-Stage  Training)                                    .
                                    JSD                                            .

| 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 |

.

1.                                           .
ELMO    GPT                                                    .
    ,  BERT-small                                      ,  ELECTRA-small                    .
2.
BERT                                        BERT                          .
3.  BERT
BERT-small/base        BERT                          .
4.  1  GPU
1  GPU

**LARGE MODEL**

*We train big ELECTRA models to measure the effectiveness of the replaced token detection pretraining task at the large scale of current state-of-the-art pre-trained Transformers.*

LARGE                                Replaced Token Detection Task        SOTA
              .

| Model | Train FLOPs | Params | SQuAD 1.1 dev | | SQuAD 2.0 dev | | SQuAD 2.0 test | |
|---|---|---|---|---|---|---|---|---|
| | | | EM | F1 | EM | F1 | EM | F1 |
| BERT-Base | 6.4e19 (0.09x) | 110M | 80.8 | 88.5 | – | – | – | – |
| BERT | 1.9e20 (0.27x) | 335M | 84.1 | 90.9 | 79.0 | 81.8 | 80.0 | 83.0 |
| SpanBERT | 7.1e20 (1x) | 335M | 88.8 | 94.6 | 85.7 | 88.7 | 85.7 | 88.7 |
| XLNet-Base | 6.6e19 (0.09x) | 117M | 81.3 | – | 78.5 | – | – | – |
| XLNet | 3.9e21 (5.4x) | 360M | **89.7** | **95.1** | 87.9 | **90.6** | 87.9 | 90.7 |
| RoBERTa-100K | 6.4e20 (0.90x) | 356M | – | 94.0 | – | 87.7 | – | – |
| RoBERTa-500K | 3.2e21 (4.5x) | 356M | 88.9 | 94.6 | 86.5 | 89.4 | 86.8 | 89.8 |
| ALBERT | 3.1e22 (44x) | 235M | 89.3 | 94.8 | 87.4 | 90.2 | 88.1 | 90.9 |
| BERT (ours) | 7.1e20 (1x) | 335M | 88.0 | 93.7 | 84.7 | 87.5 | – | – |
| ELECTRA-Base | 6.4e19 (0.09x) | 110M | 84.5 | 90.8 | 80.5 | 83.3 | – | – |
| ELECTRA-400K | 7.1e20 (1x) | 335M | 88.7 | 94.2 | 86.9 | 89.6 | – | – |
| ELECTRA-1.75M | 3.1e21 (4.4x) | 335M | **89.7** | 94.9 | **88.0** | **90.6** | **88.7** | **91.4** |

                                RoBERTa    ELECTRA              .
RoBERTa-500k    ELECTRA-400k              ,              ELECTRA    1/4                  R
oBERTa              .
ELECTRA-1.75M              RoBERTa-500k
          .
        ELECTRA
                    .

**EFFICIENCY ANALYSIS**

ELECTRA                                                              .

1.　　　　　　　　　　　Loss
2. Masked Language Model Task vs Replaced Token Detection Task
3. [MASK]

| Model | ELECTRA | All-Tokens MLM | Replace MLM | ELECTRA 15% | BERT |
|---|---|---|---|---|---|
| GLUE score | 85.0 | 84.3 | 82.4 | 82.4 | 82.2 |

Table 5: Compute-efficiency experiments (see text for details).

1.　　　　　　　　　　Loss　　　　　　　　　　　　　.
2. [MASK]　　　　　　　　　.
3. Replaced Token Detection Task　Masked Language Model Task

, ELECTRA　　　　　　　　　　　　　　.

| | | |
|---|---|---|
| Screen Shot 2022-05-01 at 11.49.40 PM.png | 55.4 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 12.41.59 AM.png | 15 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 12.42.36 AM.png | 14.9 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 12.42.36 AM.png | 14.9 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 12.56.44 AM.png | 14.9 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 1.01.26 AM.png | 9.81 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 1.05.57 AM.png | 20.2 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 1.11.42 AM.png | 19.9 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 1.43.02 AM.png | 106 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 1.49.15 AM.png | 107 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 2.04.31 AM.png | 115 KB | 2022-05-01 |
| Screen Shot 2022-05-02 at 2.38.35 AM.png | 30.6 KB | 2022-05-01 |