## *Paper Dissection: Electra: Pre-training Text Encoders As Discriminators Rather Than Generators*

Yunfan Zhou, Lu Li

# 1. Introduction

**1) Overview**: This paper introduced a *replaced token detection* task *to* address the substantial computing resources problem in Masked Language Model(MLM). Their innovative solution involves two parts: Firstly, replace tokens with alternatives sampled from a smaller generator network. Secondly, training a discriminative model to predict the replaced tokens. This dual strategy not only tackles the computationally intensive nature of MLMs but also yields exceptional performance on downstream tasks, surpassing benchmarks set by renowned models like GPT, RoBERTa, and XLNet.
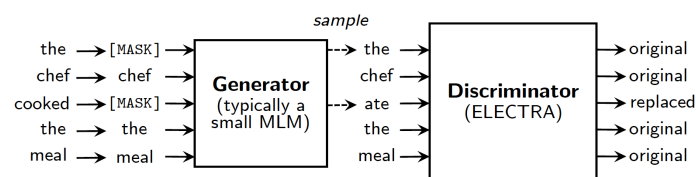
**2) Problem the paper aims to address:** A substantial compute cost occurs in current text representation methods such as BERT due to a small proportion (typically 15%) of the unlabeled sequences to train the network.

**3) The interest in solving this problem:** Reducing the compute cost can make developing and applying pre-trained text encoders more accessible to researchers and practitioners with less access to computing resources.

# 2. Methodology and Innovation

## A. Machine Learning Methods and their role

The structure of the model is inspired by the GAN(generative adversarial network) model, which contains a generator and a discriminator. The generator produces challenging negative samples using fake tokens to replace the masked token in sequences based on maximum likelihood method to confuse the discriminator. The discriminator(*the ELECTRA* model) reads all the output samples from the generator and predicts the replaced negative tokens from the original tokens.



*Picture 1*: the structure of replaced token detection, include a generator(usually a small MLM)) and a discriminator(the ELECTRA model)

**Pre-training:** They consider a small MLM model as a generator which generates "negative samples" to replace masked tokens in sequences, and trained jointly with the discriminator ELECTRA.

**Fine-tuning:** After pre-training, they throw out the generator and only fine-tune the discriminator (the ELECTRA model) on downstream tasks.

## B. Main Innovation

### 1) Primary innovation or contribution of the paper

The main innovation of their work is the discriminator, the ELECTRA model, learns from all input tokens instead of just the small masked-out subset, making it more computationally efficient. At the same time, achieving high performance on downstream tasks combined with other networks with comparatively less training time and computing resources.

### 2) Methods address the stated problem

In the *replaced token detection* structure, they introduced a generator as an intermediate to produce the full sequence samples to train the discriminator, joining with the ELECTRA model to utilize all the input tokens to raise the usage ratio of the sequences compared with the masked language models to make it compute-efficient.

## 3. Key Findings

### A. Summarize the key finding

### 1) ELECTRA-Small

To achieve the goal of this work, which is to improve the efficiency of pre-training, the study initially developed a small model that can be quickly trained on a single GPU. This ELECTRA-Small model starts with BERT-Base hyperparameters, but has a shorter sequence length and smaller token embeddings. The results showed that ELECTRA-Small outperformed BERT-Small models in terms of GLUE score while using fewer steps, indicating high efficiency and performance. Remarkably, ELECTRA-Small not only surpasses BERT-Small but also outperforms the much larger GPT model, despite requiring less compute and fewer parameters.

### 2) ELECTRA - Large

To measure the effectiveness of the replaced token detection pre-training task on a large scale, the study trained two versions of ELECTRA-Large based on different steps it used: ELECTRA-400K and ELECTRA-1.75M. In terms of GLUE score, ELECTRA-400K performs comparably to RoBERTa and XLNet but with significantly less compute, while ELECTRA-1.75M outperforms them in most GLUE tasks. Results on SQuAD show that ELECTRA models perform better than those based on masked-language modeling

### 3) Efficiency Analysis

To investigate the factor that influence the efficiency and performance of ELECTRA models, the study designed three variants: ELECTRA 15% (loss from only 15% masked tokens), Replace MLM (tokens replaced by a generator model), and All-Tokens MLM (predicts all tokens' identities, using tokens from a generator). The results turn out that ELECTR's strong performance is largely due to its learning from all input tokens, rather than just a subset.

### B. Novel insights

The replaced token detection method, which involves discriminating between real and fake tokens across the entire input sequence, is an innovation that leads to higher efficiency and performance. Also, the approach designed to investigate ELECTRA's gains is intriguing. Initially, the results challenge the prior assumption. The study then examines three aspects: the loss function, the replacement method, and predictions for all tokens, to identify the most influential factors. Such insights are valuable for designing our own experiments and offer guidance on enhancing efficiency when training BERT-based models with limited computing resources.

## 4. Possible problems

Despite its efficiency in smaller models, the ELECTRA-Large may still require substantial computational resources, which could be a limiting factor in resource-constrained environments. Also, the need for long periods of training time(steps) to achieve higher performance could be a possible problem. Finally, while ELECTRA demonstrates excellent results on standard benchmarks, its adaptability and generalizability across a wide range of real-world scenarios and diverse types of language data require further exploration

## 5. Conclusion

### A. Summarize the main points discussed in the paper

The study introduces a new self-supervised task for language representation learning, which is replacing token detection. This approach involves training a text encoder to differentiate between actual input tokens and high-quality negative samples generated by a generator network. This approach is more compute-efficient than traditional masked language modeling and leads to improved performance in benchmark tasks.

### B. Possible directions for future research

The findings of this work focus on making pre-trained text encoders more accessible, particularly for those with limited computational resources. As the study mentions, applying this model to multilingual data is considered a promising direction. However, the generalizability of the ELECTRA model needs further exploration.

## 6. Reference

Overview of GAN structure: https://developers.google.com/machine-learning/gan/gan_structure