# 005.11. ELECTRA

**ELECTRA (Efficiently Learning an Encoder that Classifies Token Replacements Accurately)** is a transformer-based pre-training method introduced in 2020 that improves the efficiency of pre-training language models compared to models like BERT. Unlike traditional masked language models (MLMs) like BERT, which focus on predicting missing tokens, ELECTRA introduces a new training objective that makes it more efficient while achieving comparable or better performance.

**ELECTRA** is a pre-training method that introduces a novel task of **Replaced Token Detection (RTD)**, where a discriminator learns to classify whether tokens in the input sequence have been replaced by a generator. This approach enables the model to learn from every token in the input sequence, making it much more efficient than traditional masked language models like BERT. ELECTRA achieves better performance with smaller models and less computational cost, making it a popular choice for pre-training large-scale language models.

## Key Concepts of ELECTRA:

### 1. Generator and Discriminator Framework

ELECTRA adopts a two-part architecture during pre-training:

- **Generator**: This part of the model generates corrupted (incorrect) tokens by replacing some of the tokens in the input sequence with alternative ones.
- **Discriminator**: Instead of trying to predict the missing tokens (like BERT), the discriminator's job is to classify whether each token in the input sequence is the original token or a replaced (corrupted) token generated by the generator.

This approach is somewhat similar to a **Generative Adversarial Network (GAN)**, where the generator creates examples, and the discriminator tries to identify whether an example is real or fake. In ELECTRA, the generator corrupts the input, and the discriminator detects the corruption.

### 2. Generator

The generator in ELECTRA is typically a small model (often a small variant of BERT), which masks a certain percentage of tokens in the input sequence (like in BERT) and replaces them with tokens predicted by the model. This is called **Masked Language Modeling (MLM)**.

- The generator learns to predict missing or corrupted tokens in the sequence.
- It is trained similarly to BERT, but it is typically smaller and less computationally intensive.

### 3. Discriminator

The discriminator is the core part of ELECTRA and performs the **Replaced Token Detection (RTD)** task. The input to the discriminator is the sequence of tokens where some tokens have been replaced by the generator. The discriminator's goal is to determine whether each token in the sequence is the original token or a replacement (fake) token.

- **Replaced Token Detection**: The discriminator processes the input sequence and outputs a binary classification for each token: 1 if the token is replaced, 0 if it is the original token.
- This task helps the model learn meaningful language representations more efficiently because the model processes every token in the input sequence, unlike BERT, where only the masked tokens are predicted.

### 4. Pre-Training Process

During pre-training, ELECTRA uses the following steps:

1. The **generator** corrupts the input sequence by replacing some tokens with predicted tokens.
2. The **discriminator** then attempts to classify each token in the sequence as either an original token or a replaced token.
3. Both the generator and discriminator are trained simultaneously, but more focus is placed on the discriminator, which becomes the pre-trained model after training.

The **discriminator** is the part of ELECTRA that gets fine-tuned for downstream tasks, such as text classification, sentiment analysis, or question answering.

### 5. Advantages of ELECTRA

- **Efficiency**: ELECTRA is much more efficient than BERT. While BERT only learns from the masked tokens (which are a small fraction of the total tokens), ELECTRA learns from every token in the input sequence, as the discriminator processes the entire sequence. This results in better use of the input data and faster convergence during training.
- **Smaller Model Size**: Since ELECTRA can achieve comparable or even better performance than BERT with smaller models, it is more computationally efficient. The smaller model size also makes it more practical for real-world applications where resources may be limited.
- **Better Sample Efficiency**: ELECTRA learns faster from a given amount of data compared to BERT, which means it can achieve strong performance using fewer computational resources and training time.

### 6. Pre-Training Objective

The main pre-training task for ELECTRA is **Replaced Token Detection (RTD)**, which differs from the **Masked Language Modeling (MLM)** used in BERT:

- In BERT's MLM, only a small percentage of tokens are masked and predicted, meaning the model learns from a small subset of tokens in each sequence.

- In ELECTRA's RTD, every token is classified as real or replaced, meaning the model learns from every token in the sequence, making the learning process more efficient.

The loss function for ELECTRA's pre-training is based on the discriminator's ability to correctly classify tokens as replaced or not.

### 7. Fine-Tuning

Once pre-training is complete, only the discriminator is kept, and the generator is discarded. The pre-trained discriminator can then be fine-tuned for various downstream tasks, such as:

- **Text Classification**: Classifying text into categories like sentiment (positive/negative).
- **Question Answering**: Predicting answers to questions based on input text.
- **Named Entity Recognition (NER)**: Identifying entities like names, places, or organizations in a text.

### 8. Comparison to BERT

- **Training Objective**: BERT uses **Masked Language Modeling (MLM)**, while ELECTRA uses **Replaced Token Detection (RTD)**.
- **Training Efficiency**: ELECTRA is more efficient because it learns from every token in the input sequence, while BERT learns from only the masked tokens.
- **Model Size**: ELECTRA can achieve better performance with smaller models compared to BERT, making it computationally cheaper.

Despite being smaller and more efficient, ELECTRA achieves state-of-the-art performance on several NLP benchmarks, demonstrating its effectiveness in learning high-quality language representations.

# Example in Python

Here's an example of how to apply **ELECTRA** for an NLP task like text classification using Hugging Face's `transformers` library. In this example, we'll use a pre-trained ELECTRA model for sentiment analysis on a dataset like IMDb.

## Requirements

You'll need to install the following libraries:

```
pip install transformers torch datasets
```

Python Code: ELECTRA for Sentiment Analysis

```
import torch
from transformers import ElectraTokenizer,
ElectraForSequenceClassification, AdamW
from datasets import load_dataset
```

```python
from torch.utils.data import DataLoader
from transformers import get_scheduler
from tqdm.auto import tqdm

# Load the IMDb dataset
dataset = load_dataset("imdb")

# Split the dataset into train and test sets
train_dataset = dataset['train']
test_dataset = dataset['test']

# Load the pre-trained ELECTRA tokenizer and model
tokenizer = ElectraTokenizer.from_pretrained(
    'google/electra-small-discriminator')
model = ElectraForSequenceClassification.from_pretrained(
    'google/electra-small-discriminator', num_labels=2)

# Tokenization function
def tokenize_function(examples):
    return tokenizer(
        examples['text'], padding="max_length", truncation=True
    )

# Tokenize the datasets
train_dataset = train_dataset.map(tokenize_function, batched=True)
test_dataset = test_dataset.map(tokenize_function, batched=True)

# Set the format for PyTorch (remove non-input columns)
train_dataset.set_format(
    type="torch", columns=["input_ids", "attention_mask",
"label"])
test_dataset.set_format(
    type="torch", columns=["input_ids", "attention_mask",
"label"])

# Create data loaders
train_dataloader = DataLoader(
    train_dataset, shuffle=True, batch_size=8)
test_dataloader = DataLoader(test_dataset, batch_size=8)

# Move the model to GPU if available
device = torch.device("cuda") if torch.cuda.is_available() else
torch.device("cpu")
model.to(device)

# Define optimizer and learning rate scheduler
optimizer = AdamW(model.parameters(), lr=5e-5)
num_epochs = 3
```

```python
num_training_steps = num_epochs * len(train_dataloader)
lr_scheduler = get_scheduler(
    name="linear",
    optimizer=optimizer,
    num_warmup_steps=0,
    num_training_steps=num_training_steps
)

# Training loop
progress_bar = tqdm(range(num_training_steps))
model.train()

for epoch in range(num_epochs):
    for batch in train_dataloader:
        batch = {k: v.to(device) for k, v in batch.items()}
        outputs = model(**batch)
        loss = outputs.loss
        loss.backward()

        optimizer.step()
        lr_scheduler.step()
        optimizer.zero_grad()
        progress_bar.update(1)

# Evaluation loop
model.eval()
accuracy = 0
num_eval_steps = 0

for batch in test_dataloader:
    batch = {k: v.to(device) for k, v in batch.items()}
    with torch.no_grad():
        outputs = model(**batch)

    logits = outputs.logits
    predictions = torch.argmax(logits, dim=-1)
    accuracy += (predictions == batch['label']).sum().item()
    num_eval_steps += len(batch['label'])

accuracy = accuracy / len(test_dataset)
print(f"Test Accuracy: {accuracy:.4f}")

# Save the fine-tuned model
model.save_pretrained("./electra-sentiment-model")
```

## Key Components of the Code:

1. **Data Loading and Tokenization**:
   - We load the **IMDb dataset** using the Hugging Face `datasets` library. IMDb is a popular dataset for binary sentiment classification (positive/negative).
   - The **ELECTRA tokenizer** is used to tokenize the input text. This tokenizer processes the input text into token IDs, adds padding, and truncates longer sequences.
2. **Model Definition**:
   - We use **ELECTRA**'s discriminator (`ElectraForSequenceClassification`) from the Hugging Face `transformers` library. The model is fine-tuned for binary classification (positive/negative sentiment), and we set `num_labels=2` since IMDb is a binary classification task.
   - The discriminator detects whether each token has been replaced, making ELECTRA particularly effective for tasks like text classification.
3. **Training Setup**:
   - The model is fine-tuned using the **AdamW optimizer** with a linear learning rate scheduler. We specify a learning rate of `5e-5`, and the training runs for `3` epochs.
   - The model's parameters are updated using backpropagation during each training step.
4. **Evaluation**:
   - After training, the model is evaluated on the test set. Accuracy is computed by comparing the model's predictions with the actual labels in the test data.
   - The model is put in evaluation mode (`model.eval()`), and predictions are made without updating the model's weights.
5. **Saving the Model**:
   - After training, the fine-tuned model is saved using `model.save_pretrained()`. This allows you to reuse the model later for inference or further fine-tuning.

## Output:

The model will print the **test accuracy** after training, showing how well the fine-tuned ELECTRA model performed on the sentiment classification task. For example, it might output something like:

```
Test Accuracy: 0.9456
```

## Notes:

- **Fine-Tuning ELECTRA**: In this example, we're fine-tuning the ELECTRA discriminator for text classification. ELECTRA learns more efficiently than models like BERT because it processes every token in the sequence, not just masked tokens.

- **Other Use Cases**: You can fine-tune ELECTRA for various other tasks such as question answering, named entity recognition, or sentence pair classification by adjusting the task-specific head (classification layer) and dataset.

This example demonstrates how to use **ELECTRA** for **sentiment analysis**, and the same approach can be applied to other NLP tasks.