# Tech Debt Transformer Pipeline
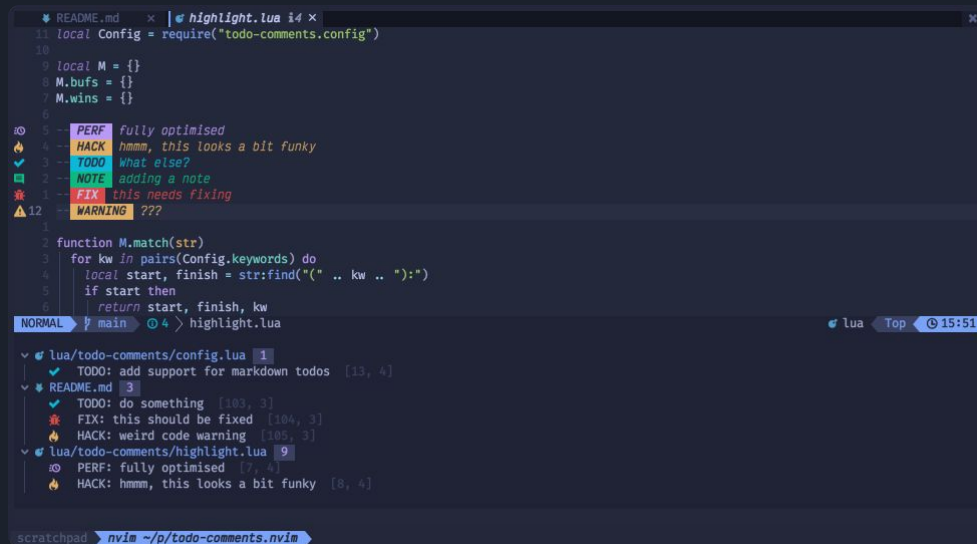
David Gao

# Recap: Transformers!

A good way to capture a lot of information!

Classification

Generation

# Recap: Technical Debt

Definition: "Implied cost of future reworking required when choosing an easy but limited solution instead of a better approach that could take more time" Wikipedia.

# Existing Solutions, Their Limitations, and Novelty

- Existing solutions usually utilize rules to find technical debt (more have started to use machine learning recently)
- Only one transformer based approach: finding technical debt in issue trackers. (Skryseth, Daniel, et al.)

- This study will be focused on directly looking at code from git commits to get insights into technical debt!

# Problem Statement

In open source projects, managing technical debt can be a challenging process. This study explores how fine-tuned transformer models can help improve the identification and classification of these issues.
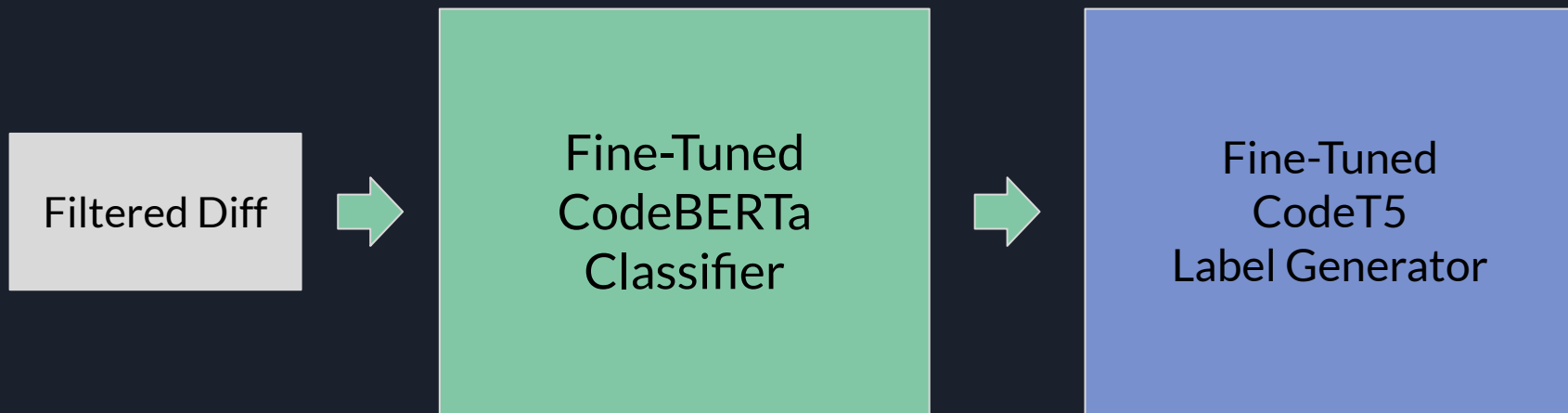
# Research Questions

1. How well can a fine-tuned transformer model classify technical debt, given code from a commit?
2. How well can a fine-tuned transformer model generate labels for technical debt?

# Proposed Architecture

Filtered Diff

Fine-Tuned CodeBERTa Classifier

Fine-Tuned CodeT5 Label Generator

# The Technical Debt Dataset

- 78K commits from 33 Java Projects
- SZZ algorithm annotations
  - find the fix, trace backwards
- SonarQube
  - static code analysis, rule based

# Classification Task

How well can a fine-tuned transformer model classify technical debt, given code from a commit?

Metrics: Accuracy, F1 Score

# Preliminary Work

# Preliminary Work (continued)



davidgaofc / **TechDebtClassifier**  ♡ like  0

Text Classification   Transformers   TensorBoard   Safetensors   roberta   generated_from_trainer   Inference Endpoints

Model card   Files and versions   Training metrics   Community   Settings   Train ▾   Deploy ▾   Use in Transformers

Edit model card

## training

This model is a fine-tuned version of huggingface/CodeBERTa-small-v1 on an unknown dataset. It achieves the following results on the evaluation set:

Downloads last month
0

Safetensors   Model size  83.5M params   Tensor type  F32

```
trainer.evaluate()

[8645/8645 17:35]
{'eval_loss': 0.13908784091472626,
 'eval_accuracy': 0.9541031017280023,
 'eval_f1': 0.9492902966412665,
 'eval_runtime': 1056.6171,
 'eval_samples_per_second': 65.449,
 'eval_steps_per_second': 8.182}
```

$$\text{F1 Score} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}$$
$$= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

11

# Generation Task

How well can a fine-tuned transformer model generate labels for technical debt?

Metrics: BLEU, ROUGE

# Preliminary Work

# Preliminary Work (continued)

# Pipeline



```
@@ -0,0 +1,25 @@
+package org.apache.accumulo.server.test.randomwalk.shard;
+
+import java.util.Properties;
+import java.util.Random;
+import java.util.SortedSet;
+
+import org.apache.accumulo.server.test.randomwalk.State;
+import org.apache.accumulo.server.test.randomwalk.Test;
+import org.apache.hadoop.io.Text;
+
+
+public class Split extends Test {
+
+    @Override
+    public void visit(State state, Properties props) throws Exception {
+        String indexTableName = (String)state.get("indexTableName");
+        int numPartitions = (Integer)state.get("numPartitions");
+        Random rand = (Random) state.get("rand");
+
+        SortedSet<Text> splitSet = ShardFixture.genSplits(numPartitions,
+        log.debug("adding splits " + indexTableName);
+        state.getConnector().tableOperations().addSplits(indexTableName,
+    }
+
+}
```

output

[[{'label': 'LABEL_1', 'score': 0.9938503503799438}, {'label': 'LABEL_0', 'score': 0.00614961888641119}]]

output

[{'generated_text': "Rename this local variable to match the regular expression '^[a-z][a-zA"}]

# Questions/Comments?

# References

Valentina Lenarduzzi, Nyyti Saarimäki, and Davide Taibi. 2019. The Technical Debt Dataset. In Proceedings of the Fifteenth International Conference on Predictive Models and Data Analytics in Software Engineering (PROMISE'19). Association for Computing Machinery, New York, NY, USA, 2–11. https://doi.org/10.1145/3345629.3345630