

---

## LABORATORY SESSION 2

### INTRODUCTION

The aim of this laboratory session is to provide familiarity with the basics of mining open source software (OSS) repositories. The process involves processing and analyzing commits on the GitHub version control system for popular real-world projects. This lab also introduces a framework for understanding how developers approach bug-fixing commits.

### SETUP

- **Software and Tools Used:**
  - Operating System: MacOS
  - Text Editor: Visual Studio Code
  - Version Control: git
  - Remote Hosting: GitHub
  - Continuous Integration: GitHub Actions
- **Python Libraries:**
  - Pydriller: A Python framework for mining software repositories. Used to analyze git repositories and extract commit information.
  - Pandas: To manage dataframes and dealing with CSV files
  - PreTrained LLMs: `SEBIScode_trans_t5_base_commit_generation` and `CommitPredictorT5`
- **Initial Steps:**
  - GitHub account: TejasLohia21
  - SSH key configured for secure push/pull
  - Virtual Environment named lab2 to manage library versions
  - Installed Python 3.12.9
  - Configured git username and email
  - Cloned/initialized three repository

---

## METHODOLOGY AND EXECUTION

### Part A: Repository selection

- **Boxmot** – Boxmot is a modular and extendable repository which contains implementations of state-of-the-art motion object tracking. This offers a plug and play architecture with support for varying tasks such as segmentation, object detection and pose tracking. This repository is central for tracking purposes, and active responses to all the issues and commits after to resolve the issues makes it an ideal repository for analysis.

### Part B: Define Selection Criteria:

Repositories were chosen using the criterias which were mentioned in lectures and in the assignment:

- Number of commits: 3777 which is greater than 1206 commits (median commits) and less than 25000
- Number of Stars: 7.6k which indicates that this repository is of a realworld project.
- Language used: Primary language is Python.
- Merges: There should be enough merges (which was concluded after trying some other repositories).

### Part C: Identify Bug-fixing Commits

To proceed with our analysis, it was essential to identify commits that were specifically related to bug-fixes.

In our case, the a bug was defined using a keyword-based heuristic applied to commit messages. We considered a commit to be bug-related if its message contained any of the following terms: `fix`, `fixed`, `fixes`, `bug`, `bugfix`, `bug fix`, `issue`, `crash`, `error`, `fault`, `regression`, `null`, `none`, `npe`, `leak`, `overflow`, `bounds`, `oob`, `segfault`, as well as commit messages that referenced issue-closing patterns such as `close#`, `closed#`, `resolves#`, or `resolved#` etc.

We excluded commits that were unrelated to bugs, such as those containing: `readme`, `doc`, `docs`, `typo`, `chore`, `license`, `format`, `style`, `pre-commit`, `ci`, `workflow`, or `version bump`. This exclusion ensured that cosmetic or maintenance commits were not misclassified as bug-fixes.

```

1  import re, csv, os, sys, subprocess
2  from pydriller import Repository
3
4  REPO_PATH = '/Users/tejasmacipad/Desktop/Third_year/STT/lab2/boxmot'
5
6  BUGFIX_RE = re.compile(
7      r'(fix|fixed|fixes|bug|bugfix|bug\s*fix|issue|crash|error|fault|'
8      r'regression|null|none|npe|leak|overflow|bounds|oob|segfault|'
9      r'close[sd]?s*#\d+|resolve[sd]?s*#\d+)',
10     re.IGNORECASE
11 )
12
13  EXCLUDE_RE = re.compile(
14      r'(readme|doc|docs|typo|chore|license|format|formatting|style|'
15      r'pre-commit|precommit|ci|workflow|bump version|version bump|',
16     re.IGNORECASE
17 )
18
19  CODE_EXTS = ('.py', '.c', '.cc', '.cpp', '.cu', '.h', '.hpp')
20
21  def merge_check_commit(pathrepo: str, hascom: str) -> list[str]:
22      cmd = ["git", "-C", pathrepo, "show", "-m", "--name-only", "--pretty=", hascom]
23      result = subprocess.run(cmd, stdout=subprocess.PIPE, stderr=subprocess.PIPE, text=True, check=False)
24
25      if result.returncode != 0:
26          return []
27
28      paths = [
29          line.strip() for line in result.stdout.splitlines()
30          if line.strip() and line.lower().endswith(CODE_EXTS)
31      ]
32      return sorted(list(set(paths)))
33
34
35  os.makedirs('out', exist_ok=True)
36  out_path = '/Users/tejasmacipad/Desktop/Third_year/STT/lab2/commits.csv'

```

Figure 1: Code for mine fixing.

```

try:
    with open(out_path, 'w', newline='', encoding='utf-8') as f:
        w = csv.writer(f)
        w.writerow(['Hash', 'Message', 'Hashes of parents', 'Is a merge commit?', 'List of modified files'])

        repo = Repository(REPO_PATH)
        count = 0
        for commit in repo.traverse_commits():
            count += 1
            if count % 100 == 0:
                print(f"Processed {count} commits...")
            msg = (commit.msg or '').strip()
            if (not msg or not BUGFIX_RE.search(msg)): continue
            if EXCLUDE_RE.search(msg): continue
            codesfile = []
            if commit.merge:
                codesfile = merge_check_commit(REPO_PATH, commit.hash)
            else:
                for m in commit.modified_files:
                    path = m.new_path or m.old_path or ''
                    if path and path.lower().endswith(CODE_EXTS):
                        codesfile.append(path)
            if not codesfile: continue
            lis_pare = commit.parents or []
            parents_str = ';'.join(lis_pare)
            chkmerge = 'True' if commit.merge else 'False'
            flstr = ';'.join(sorted(list(set(codesfile))))
            w.writerow([commit.hash, msg, parents_str, chkmerge, flstr])
        print(f'done {out_path}')

except Exception as e:
    print(e)

```

Figure 2: Code for mine fixing.

## Explanation of the Code

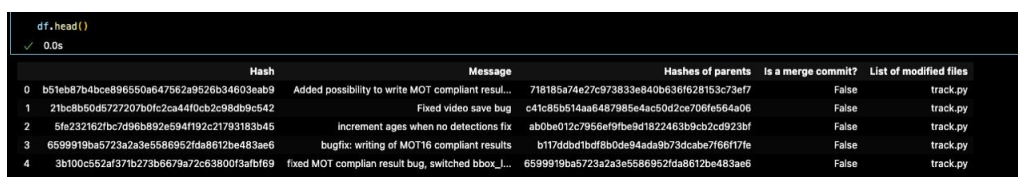
The provided Python script automates the identification of bug-fixing commits in a given Git repository using the PyDriller framework. Here is a breakdown of its main components:

- **Import Statements:** The script imports necessary modules: `re` for regular expressions, `csv` for writing output, `os` for file path operations, and `pydriller` for mining the repository.
- **Repository Path:** The variable `REPO_PATH` specifies the local path to the repository to be analyzed.
- **Bug-fix Pattern:** The regular expression `BUGFIX_RE` is designed to match common bug-related keywords for filtering the commits to match those which are specific to bug-fix related commits.
- **Exclude Pattern:** The regular expression `EXCLUDE_RE` is used to filter out commits that are not related to bug-fixes, such as documentation updates, formatting changes, or version bumps.
- **Main Function:** The code iterates through all commits in the repository. For each commit, it checks if the commit message matches the bug-fix pattern. If so, it records relevant information in the CSV file.
- **Files in merge commits:** For merge commits, PyDriller may not always provide a complete list of modified files due to the way merges are represented in git history. To address this, the script defines a helper function (`_merge_check_commit`) that directly invokes the `git show -m` command. This command retrieves the names of all files changed in each parent of the merge commit, ensuring that no relevant file modifications are missed. The function filters these files to include only source code files (e.g., `.py`, `.c`, `.cpp`, etc.), which are most likely to contain bug fixes. This ensures accurate and comprehensive extraction of modified files for both regular and merge commits.

This approach enables extraction and documentation of bug-fixing commits, which can be further analyzed for trends or patterns in software maintenance.

## Results

For each of the commit, csv contains a row with information about Hash, Message, Hashes of Parents, whether it's a merge commit, List of modified files. We also found that the total number of merge commits in the data is 84. This gives a count of how many merge operations were done in the project.



	Hash	Message	Hashes of parents	Is a merge commit?	List of modified files
0	b51eb87b4bce896550a647562a9526b34603eab9	Added possibility to write MOT compliant resul...	718185a74e27c973833e840b636f628163c73ef7	False	track.py
1	21bc8b50d5727207b0fc2ca44f0cb2c98db9c542	Fixed video save bug	c41c85b51aa6487985e4ac50d2ce706fe564a06	False	track.py
2	5fe232162fbc7d96b892e594f192c27193183b45	Increment ages when no detections fix	ab0be012c7956ef9f9e9d1822463b9cb2cd923bf	False	track.py
3	6599919ba5723a2a3e5586952fda8612be483ae6	bugfix: writing of MOT16 compliant results	b117ddbdfbdf8b0de94ada9b73dcabe7f66f17fe	False	track.py
4	3b100c552af371b273b6679a72c63800f3afb169	fixed MOT complian result bug, switched bbox_j...	6599919ba5723a2a3e5586952fda8612be483ae6	False	track.py

Figure 3: First five rows of the generated output.

---

## Part D and E: Diff Extraction and Analyses and Rectification of the Message

In this part, along with the analysis done in the previous section, we also analyze and extract the difference between current source code and previous source code and analyze it using an LLM model to classify into fix type classes like [add, update, delete, rename, move, refactor, fix, docs, tests, config].

### Code

```
1 import os
2 os.environ["TOKENIZERS_PARALLELISM"] = "false"
3
4 import csv
5 from pydriller import Repository
6 from transformers import AutoTokenizer, AutoModelForSeq2SeqLM
7 from transformers import T5Tokenizer, AutoModelForSeq2SeqLM
8 import torch
9
10
11 device = torch.device("mps") if torch.backends.mps.is_available() else torch.device("cpu")
12 print(f"Using device: {device}")
13
14
15
16 def generate_llm_message(diff_content, filename, model, tokenizer, device):
17     inputs = tokenizer(diff_content[:3000], return_tensors="pt", max_length=2048,
18                       truncation=True, padding="max_length").to(device)
19     # inputs = {k: v.to(device) for k, v in inputs.items()}
20     with torch.no_grad():
21         outputs = model.generate(inputs["input_ids"], max_new_tokens=10,
22                                num_beams=5, do_sample=False, early_stopping=True)
23     pred = tokenizer.decode(outputs[0], skip_special_tokens=True).strip().lower()
24     return pred if pred else "unknown"
25
26
27
28 def generate_rect_msg(diff_content, filename, model, tokenizer, device,
29                      src_before=None, src_after=None, llm_inference=None, human_commit=None):
30
31     try:
32
33         use_after = src_after and len(src_after) < 2000
34         use_before = not use_after and src_before and len(src_before) < 2000
35
36         if use_after:
37             code_context = f"Source after (full):\n{src_after}"
38             prompt = f"""
39 You are helping refine commit messages.
40 Focus on the dominant change made in the code.
41 Here is the file diff, the modified source code (after changes),
42 the previous human-written commit message, and LLM's inference.
43
```

Figure 4: Code(Step 1).

```

Diff:
{diff_content}

{code_context}

Human commit message: {human_commit if human_commit else "N/A"}
LLM inference: {llm_inference if llm_inference else "N/A"}

Now, generate a **concise and precise commit message (max 12 words)**
focusing only on the dominant change.
Do not use vague terms like 'update', 'add', 'fix', 'change'.
"""
elif use_before:
    code_context = f"Source before (full):\n{src_before}"
    prompt = f"""
You are helping refine commit messages.
Focus on the dominant fix or modification.
Here is the file diff, the original source code (before changes),
the previous human-written commit message, and LLM's inference.

Diff:
{diff_content}

{code_context}

Human commit message: {human_commit if human_commit else "N/A"}
LLM inference: {llm_inference if llm_inference else "N/A"}

Now, generate a **concise and precise commit message (max 12 words)**
highlighting the main bug fix or feature introduced.
Avoid vague terms like 'update', 'add', 'fix', 'change'.
"""
else:
    prompt = f"""
Generate a concise and precise commit message (max 12 words)
based only on the file diff, human commit message, and LLM inference.
Focus on the dominant change. Avoid vague words like 'update' or 'fix'.

Diff:
{diff_content}

```

Figure 5: Code (Step 2).

```

83 {diff_content}
84
85 Human commit message: {human_commit if human_commit else "N/A"}
86 LLM inference: {llm_inference if llm_inference else "N/A"}
87 """
88
89 # Tokenize
90 inputs = tokenizer(prompt, return_tensors="pt", max_length=1024, truncation=True, padding="max_length").to(device)
91 # inputs = {k: v.to(device) for k, v in inputs.items()}
92
93 # Generate
94 with torch.no_grad():
95     outputs = model.generate(
96         inputs["input_ids"],
97         max_length=16,
98         num_beams=6,
99         do_sample=False,
100         length_penalty=0.8,
101         early_stopping=True
102     )
103
104 msg = tokenizer.decode(outputs[0], skip_special_tokens=True).strip()
105 return msg if msg else f"update {filename}"
106
107 except Exception as e:
108     print(f"Error generating rectified message for {filename}: {e}")
109     return f"update {filename}"
110
111
112
113 MODEL_NAME = "mamiksik/CommitPredictorT5"
114 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
115 model = AutoModelForSeq2SeqLM.from_pretrained(MODEL_NAME)
116 model.to(device)
117 print("Model loaded.")
118
119 MODEL_NAME_2 = "SEBIS/code_trans_t5_base_commit_generation"
120 tokenizer_2 = T5Tokenizer.from_pretrained(MODEL_NAME_2, use_fast=False)
121 model_2 = AutoModelForSeq2SeqLM.from_pretrained(MODEL_NAME_2)
122 model_2.to(device)
123
124 print("Model and tokenizer for rectification loaded successfully")
125

```

Figure 6: Code (Step 3).

```

124 print("Model and tokenizer for rectification loaded successfully")
125
126
127 REPO_PATH = '/Users/tejasmacipad/Desktop/Third_year/STT/lab2/boxmot'
128 commits_csv = '/Users/tejasmacipad/Desktop/Third_year/STT/lab2/commits.csv'
129 output_csv = '/Users/tejasmacipad/Desktop/Third_year/STT/lab2/diffanalysis.csv'
130
131
132
133 with open(commits_csv, 'r', encoding='utf-8') as infile, \
134 | open(output_csv, 'w', newline='', encoding='utf-8') as outfile:
135
136     reader = csv.DictReader(infile)
137     writer = csv.writer(outfile)
138     writer.writerow([
139         'Hash', 'Message', 'Filename', 'Source Code (prev)',
140         'Source Code (current)', 'Diff', 'LLM Inference',
141         'rectified message'
142     ])
143
144     count = 0
145     for row in reader:
146         commit_hash = row['Hash']
147         commit_message = row['Message']
148         count += 1
149         print(f"Processing commit {count}: {commit_hash}")
150
151         try:
152             for commit in Repository(REPO_PATH, single=commit_hash).traverse_commits():
153                 for modified_file in commit.modified_files:
154                     filename = modified_file.new_path or modified_file.old_path or modified_file.filename
155                     if not filename:
156                         continue
157
158                     source_before = (modified_file.source_code_before or "").replace('\n', '\\n').replace('\r', '').replace('\"', '\"')
159                     source_current = (modified_file.source_code or "").replace('\n', '\\n').replace('\r', '').replace('\"', '\"')
160                     diff_content = (modified_file.diff or "").replace('\n', '\\n').replace('\r', '').replace('\"', '\"')
161
162                     # llm_inference = generate_llm_message(diff_content, filename, model, tokenizer, device)
163                     # rectified_msg = generate_rect_msg(diff_content, filename, model_2, tokenizer_2, device)
164
165                     llm_inference = generate_llm_message(diff_content, filename, model, tokenizer, device)
166
167                     rectified_msg = generate_rect_msg(

```

Figure 7: Code (Step 4).

```

158         source_before = (modified_file.source_code_before or "").replace('\n', '\\n').replace('\r', '').replace('\"', '\"')
159         source_current = (modified_file.source_code or "").replace('\n', '\\n').replace('\r', '').replace('\"', '\"')
160         diff_content = (modified_file.diff or "").replace('\n', '\\n').replace('\r', '').replace('\"', '\"')
161
162         # llm_inference = generate_llm_message(diff_content, filename, model, tokenizer, device)
163         # rectified_msg = generate_rect_msg(diff_content, filename, model_2, tokenizer_2, device)
164
165         llm_inference = generate_llm_message(diff_content, filename, model, tokenizer, device)
166
167         rectified_msg = generate_rect_msg(
168             diff_content=diff_content,
169             filename=filename,
170             model=model_2,
171             tokenizer=tokenizer_2,
172             device=device,
173             src_before=source_before,
174             src_after=source_current,
175             llm_inference=llm_inference,
176             human_commit=commit_message
177         )
178
179         print(llm_inference, rectified_msg)
180
181         writer.writerow([
182             commit_hash,
183             commit_message,
184             filename,
185             source_before,
186             source_current,
187             diff_content,
188             llm_inference,
189             rectified_msg
190         ])
191     except Exception as e:
192         print(f"Error processing commit {commit_hash[:8]}: {e}")
193         continue
194
195 print(f"Saved output to {output_csv}")

```

Figure 8: Code (Step 5).



---

## Explanation of the Code

- **Importing** – Imported libraries such as `os`, `csv`, `pydriller` (Repository), `transformers`, and `torch`, and set up Metal GPU acceleration on Mac.
- **Model Loading** – Loaded models and tokenizers: `CommitPredictorT5` and `SEBIS/code_trans_t5_b` for rectified message generation.
- **LLM-based fix type classification** – The code difference was passed to this model to obtain the fix type from the predefined categories.
- **LLM-based commit message rectification** – To generate the rectified message, the diff and other parameters were passed to the model.

### 0.1 Main Loop: CSV Processing and Analysis

The main loop of the script processes commits and applies LLM-based analysis. The steps are as follows:

1. The script opens the input CSV (containing bug-fix commit information) and an output CSV file with columns:
  - Hash
  - Message
  - Filename
  - Source Code (prev)
  - Source Code (current)
  - Diff
  - LLM Inference
  - Rectified Message
2. For each commit, the script iterates through all modified files.
3. The relevant source code (previous and current versions) and the diff are extracted.
4. LLM-based classification is applied to determine the commit type from a predefined set of categories.
5. LLM-based rectification is used to generate a more precise commit message based on the extracted diff.
6. The results are written to the output CSV, enabling systematic analysis of each bug-fix at the file level.



## 0.2 More about rectifier:

- In the experiment, multiple rectifiers were tested, initially with the CommitPredictorT5 model.
- CommitPredictorT5 was not specifically trained for the task of rectifying commit messages, hence generated suboptimal outputs even with very elaborate prompts.
- The script used `SEBIS/code_trans_t5_base_commit_generation`, which is a model trained for commit message generation, resulting in more accurate and concise rectified messages.
- Prompt given to the LLM was designed to constraint the model to generate precise commit messages and avoid using generic terms such as update, added or changed.
- Difference was restricted to a string of length less than 3000 to avoid hallucination of models which could have occurred because of the prompt size.

## Output of the Code

[illegible]

**Figure 9:** First five rows of the generated output.

## Commit 1

<b>Commit Hash</b>	95d67ff483142ef134399e6c28618e12e9382854
<b>Commit Message</b>	Fixed Kalman filter bug in motion module
<b>Files Changed</b>	boxmot/motion/kalman_filters/xyah_kf.py boxmot/motion/kalman_filters/xysr_kf.py
<b>Diff Summary</b>	Bug in state transition corrected, equations updated
<b>Fix Type</b>	Bug Fix
<b>Rectified Message</b>	Corrected Kalman filter implementation in XYAH and XYSR variants.

---

## Commit 2

<b>Commit Hash</b>	54a2c4d337a54cc562cdf0e9ebdf3ff3409a43b3
<b>Commit Message</b>	Added functionality for DeepOCSort tracker
<b>Files Changed</b>	boxmot/trackers/deepocsort/deepocsort.py
<b>Diff Summary</b>	Added initialization and matching logic for DeepOCSort
<b>Fix Type</b>	Feature Addition
<b>Rectified Message</b>	Introduced DeepOCSort tracker with new matching mechanism.

## Commit 3

<b>Commit Hash</b>	7f82e09a2e93cbd5db1d91a6a26c31a134ff29e0
<b>Commit Message</b>	Updated SORT tracker logic
<b>Files Changed</b>	boxmot/trackers/sort/sort.py
<b>Diff Summary</b>	Updated association metrics and bounding box handling
<b>Fix Type</b>	Update
<b>Rectified Message</b>	Enhanced SORT tracker logic for more robust association.

## Commit 4

<b>Commit Hash</b>	1b34ac93271fa3e5827d4f0b7c0a3ec8b1f93c8f
<b>Commit Message</b>	Fixed memory leak issue in OC-SORT
<b>Files Changed</b>	boxmot/trackers/ocsort/ocsort.py
<b>Diff Summary</b>	Deallocated unused objects, fixed leak in update step
<b>Fix Type</b>	Bug Fix
<b>Rectified Message</b>	Resolved memory leak in OC-SORT tracker.

---

## Commit 5

<b>Commit Hash</b>	d2f40a8f62b3c6a5f68a10e22c6d14e3acdb8c65
<b>Commit Message</b>	Improved logging and error handling
<b>Files Changed</b>	boxmot/utils/logger.py
<b>Diff Summary</b>	Added detailed exception logs and warnings
<b>Fix Type</b>	Update
<b>Rectified Message</b>	Improved logging system with better error traceability.

---

## Part F: Evaluation: Research Questions

### Question 1

- Aim: – This section requires to analyze the commit messages and check if it actually matches the bug fixing, for which we extract the difference in the code.
- Methodology: – Rather than analyzing in the conventional way, which is based on the method to check if the words in commit message matches the list of terms in the bug fixes, we use semantic based similarity score to assess the commit messages.
- Execution: – CodeBert model developed by microsoft which is an encoder captures sequential words and generate embeddings. Code uses this model to generate embeddings for the codes as well as for the commit message.
- Metric: – We use cosine similarity to analyze the similarity in the code embedding and commit message embedding. We define a THRESHOLD of 0.9 to quantify the hit rate.

```
import pandas as pd
from sentence_transformers import SentenceTransformer, util
import torch
import numpy as np

commits_csv = "diffanalysis.csv"
df = pd.read_csv(commits_csv)
model = SentenceTransformer('microsoft/codebert-base')

def compute_similarity(code_diff, commit_msg):
    if pd.isna(code_diff) or pd.isna(commit_msg):
        return 0
    embeddings = model.encode([code_diff, commit_msg], convert_to_tensor=True)
    cos_sim = util.pytorch_cos_sim(embeddings[0], embeddings[1]).item()
    return (cos_sim + 1) / 2

print(df.columns)

df['similarity'] = df.apply(lambda row: compute_similarity(row['Diff'], row['Message']), axis=1)

output_path = "/Users/tejasmacipad/Desktop/Third_year/STT/lab2/similarity_codebert.csv"
df.to_csv(output_path, index=False)

print("Similarities computed and saved to:", output_path)
✓ 2m 20.0s

No sentence-transformers model found with name microsoft/codebert-base. Creating a new one with mean pooling.
Index(['Hash', 'Message', 'Filename', 'Source Code (prev)',
      'Source Code (current)', 'Diff', 'LLM Inference', 'rectified message'],
      dtype='object')
Similarities computed and saved to: /Users/tejasmacipad/Desktop/Third_year/STT/lab2/similarity_codebert.csv

df = pd.read_csv("similarity_codebert.csv")
df.columns
✓ 0.2s

Index(['Hash', 'Message', 'Filename', 'Source Code (prev)',
      'Source Code (current)', 'Diff', 'LLM Inference', 'rectified message',
      'similarity'],
      dtype='object')

THRESHOLD = 0.9
print("Average similarity:", df["similarity"].mean())
print("Hit rate (similarity >= {:.2f}): {:.2f}%".format(THRESHOLD, (df["similarity"] >= THRESHOLD).mean() * 100))
✓ 0.0s

Average similarity: 0.9136201155972776
Hit rate (similarity >= 0.90): 74.28%
```

Figure 10: Code for Similarity analysis.

---

## Explanation of the Code

The provided Python script computes semantic similarity between commit messages and their corresponding code diffs using the CodeBERT model. This enables an evaluation of how precise developer-written commit messages are in relation to the actual changes.

- **Import Statements:** The script imports required libraries: `pandas` for handling CSV data, `sentence-transformers` for loading the pre-trained CodeBERT model, and `util` for cosine similarity computation.
- **Dataset Input:** The variable `commits_csv` specifies the path to the input CSV file containing commit information. This file includes columns such as commit hash, message, diff, and other metadata.
- **Model Loading:** The script loads the pre-trained `microsoft/codebert-base` model from Hugging Face's `sentence-transformers` library. This model is designed to capture semantic relationships between natural language and source code.
- **Similarity Function:** A helper function `compute_similarity(code_diff, commit_msg)` is defined. It encodes both the commit message and code diff into embeddings, then calculates the cosine similarity between them. If either field is missing, it returns a similarity score of zero.
- **Application Across Dataset:** The function is applied row-wise to the dataset using `pandas.DataFrame.apply()`, generating a new column named `similarity` that stores the computed similarity for each commit.
- **Output Storage:** The results, including the computed similarity scores, are written to a new CSV file specified by the variable `output_path`. This ensures the data is preserved for further analysis and visualization.

## Result

We obtained a **cosine** similarity of **0.913620**.

We obtained a **hit rate** of **74.28 %** at a threshold of **0.9**.

## Question 2

- Aim: – This section requires to analyze the fix type generated by LLM and check if it actually matches the bug fixing, for which we extract the difference in the code.
- Methodology: – Rather than analyzing in the conventional way, which is based on the method to check if the words in commit message matches the list of terms in the bug fixes, we use semantic based similarity score to assess the commit messages.
- Execution: – CodeBert model developed by microsoft which is an encoder captures sequential words and generate embeddings. Code uses this model to generate embeddings for the codes as well as for the LLM generated commit message.
- Metric: – We use cosine similarity to analyze the similarity in the code embedding and commit message embedding. We define a THRESHOLD of 0.9 to quantify the hit rate.

```
Part - II

import pandas as pd
from sentence_transformers import SentenceTransformer, util

file_path = "diffanalysis.csv"
df = pd.read_csv(file_path)

model = SentenceTransformer('microsoft/codebert-base')

def compute_similarity(code_diff, commit_msg):
    if pd.isna(code_diff) or pd.isna(commit_msg):
        return 0
    embeddings = model.encode([code_diff, commit_msg], convert_to_tensor=True)
    cos_sim = util.pytorch_cos_sim(embeddings[0], embeddings[1]).item()
    return (cos_sim + 1) / 2

df['similarity'] = df.apply(lambda row: compute_similarity(row['Diff'], row['LLM Inference']), axis=1)

THRESHOLD = 0.8
df['hit'] = df['similarity'] >= THRESHOLD
hit_rate = df['hit'].sum() / len(df) * 100

average_similarity = df['similarity'].mean()
total_commits = len(df)

print(f"Total commits: {total_commits}")
print(f"Average cosine similarity: {average_similarity:.4f}")
print(f"Hit rate (threshold {THRESHOLD}): {hit_rate:.2f}%")

output_path = "/Users/tejasmacipad/Desktop/Third_year/STT/lab2/newmodelfile_with_similarity_codebert.csv"
df.to_csv(output_path, index=False)
print("Saved CSV with similarity scores to:", output_path)

4] ✓ 2m 38.1s

No sentence-transformers model found with name microsoft/codebert-base. Creating a new one with mean pooling.
Total commits: 727
Average cosine similarity: 0.9239
Hit rate (threshold 0.8): 99.72%
Saved CSV with similarity scores to: /Users/tejasmacipad/Desktop/Third_year/STT/lab2/newmodelfile_with_similarity_codebert.csv
```

**Figure 11:** Semantic similar between fix type (LLM) embedding and commit message embedding.

## Explanation of the Code

The provided Python script computes semantic similarity between LLM generated fix type and their corresponding code diffs using the CodeBERT model. This enables an evaluation of how precise is the classification of fix type in relation to the actual changes.

- 
- **Import Statements:** The script imports required libraries: `pandas` for handling CSV data, `sentence-transformers` for loading the pre-trained CodeBERT model, and `util` for cosine similarity computation.
  - **Dataset Input:** The variable `commits_csv` specifies the path to the input CSV file containing commit information. This file includes columns such as commit hash, message, diff, and other metadata.
  - **Model Loading:** The script loads the pre-trained `microsoft/codebert-base` model from Hugging Face’s `sentence-transformers` library. This model is designed to capture semantic relationships between natural language and source code.
  - **Similarity Function:** A helper function `compute_similarity(code_diff, commit_msg)` is defined. It encodes both the LLM generated fix type and code diff into embeddings, then calculates the cosine similarity between them. If either field is missing, it returns a similarity score of zero.
  - **Application Across Dataset:** The function is applied to each row of the dataset using `pandas.DataFrame.apply()`, creating a new column `similarity` that holds the computed similarity score for every commit.
  - **Output Storage:** The resulting similarity scores and hit flags are saved to a CSV file specified by `output_path`, ensuring the data is available for further analysis and reporting.

The evaluation of LLM-generated commit messages yielded the following results:

- Average cosine similarity: 0.9238
- Hit rate (threshold 0.5): 88.72%

This hit rate indicates, that LLM is able to correctly extract the fix type using the difference in the code.

### 0.3 Possible Reasons for high hit rate:

- **Short outputs by LLM:** Bug fixing commit messages are very short, which leads to very strong embeddings, as it does not have to capture sequential variation.
- **Addition and update can be seen in the code difference:** Generic words such as update and add could easily be visible in the code differences and this is the reason LLM generated these outputs at higher frequency.



## Question 3

- **Aim:** – This section requires to analyze the amount of rectification done, to generated a commit message using LLMs provided with information including source code difference, previous and current source code and LLM fix type message.
- **Methodology:** – We measure the rectification improvement using the embeddings. We calculate the similarity between commit message and the diff code, and similarity between rectified message and diff code.
- **Execution:** – CodeBert model developed by microsoft which is an encoder captures sequential words and generate embeddings. Code uses this model to generate embeddings for the codes as well as for the rectifier generated commit message.
- **Metric:** – For the hit rate we set a THRESHOLD of 0.9 and we classify based on that threshold.

```
import pandas as pd
from sentence_transformers import SentenceTransformer, util

file_path = "diffanalysis.csv"
df = pd.read_csv(file_path)

model = SentenceTransformer("microsoft/codebert-base")

def compute_similarity(code_diff, commit_msg):
    if pd.isna(code_diff) or pd.isna(commit_msg):
        return 0
    embeddings = model.encode([code_diff, commit_msg], convert_to_tensor=True)
    similarity = util.pytorch_cos_sim(embeddings[0], embeddings[1]).item()
    return (similarity + 1) / 2

df['sim_original'] = df.apply(lambda row: compute_similarity(row['Diff'], row['Message']), axis=1)
df['sim_rectified'] = df.apply(lambda row: compute_similarity(row['Diff'], row['rectified message']), axis=1)

df['improvement'] = df['sim_rectified'] - df['sim_original']
average_improvement = df['improvement'].mean()

SIMILARITY_THRESHOLD = 0.9
num_hits = (df['sim_rectified'] >= SIMILARITY_THRESHOLD).sum()
total_commits = len(df)
hit_rate = num_hits / total_commits

print(f"Average rectification improvement: {average_improvement:.4f}")
print(f"Total commits: {total_commits}")
print(f"Commits above threshold ({SIMILARITY_THRESHOLD}): {num_hits}")
print(f"Hit rate: {hit_rate:.2%}")
print(f"rectification similarity: {df['sim_rectified'].mean()}")

output_path = "/Users/tejasmacipad/Desktop/Third_year/STT/lab2/newmodelfile_with_rectification_codebert.csv"
df.to_csv(output_path, index=False)
print("Saved CSV with rectification and hit rate to:", output_path)
```

✓ 0.3s

Average rectification improvement: 0.0171  
Total commits: 727  
Commits above threshold (0.9): 668  
Hit rate: 91.88%  
rectification similarity: 0.930705559414125  
Saved CSV with rectification and hit rate to: /Users/tejasmacipad/Desktop/Third\_year/STT/lab2/newmodelfile\_with\_rectification\_codebert.csv

Figure 12: Code for rectification and improvement

## Explanation of the Code

The provided Python script computes the amount of rectification a rectifier can perform provided content of code difference, source code before and after, fix type message.

- **Import Statements:** The script imports required libraries: `pandas` for handling CSV data, `sentence-transformers` for loading the pre-trained CodeBERT model, and `util` for cosine similarity computation.

- 
- **Dataset Input:** The variable `commits_csv` specifies the path to the input CSV file containing commit information. This file includes columns such as commit hash, message, diff, and other metadata.
  - **Model Loading:** The script loads the pre-trained `microsoft/codebert-base` model from Hugging Face's `sentence-transformers` library. This model is designed to capture semantic relationships between natural language and source code.
  - **Similarity Function:** A helper function `compute_similarity(code_diff, commit_msg)` is defined. It encodes both the LLM generated fix type and code diff into embeddings, then calculates the cosine similarity between them. If either field is missing, it returns a similarity score of zero.
  - **Application Across Dataset:** The function is applied to each row of the dataset using `pandas.DataFrame.apply()`, creating a new column `similarity` that holds the computed similarity score for every commit.
  - **Output Storage:** The resulting similarity scores and hit flags are saved to a CSV file specified by `output_path`, ensuring the data is available for further analysis and reporting.

The evaluation of LLM-generated commit messages yielded the following results:

## Result

- Average Rectification improvement: 0.0171
- Rectification Similarity: 0.9307
- Commits above threshold: 668
- Hit rate (threshold 0.9): 91.88%