

ASSIGNMENT 1

Software Tools and Techniques for CSE

Tejas Lohia

23110335

https://github.com/TejasLohia21/TejasLohia_Assessment1.git

TABLE OF CONTENTS

Contents

1	Introduction to Version Controlling, Git Workflows, and Actions	3
1.1	Introduction	3
1.2	Tools	3
1.3	Setup	3
1.4	Methodology and Execution	4
1.5	Results and Analysis	8
1.6	Discussion and Conclusion	10
2	Commit Message Rectification for Bug-Fixing Commits in the Wild	11
2.1	Introduction	11
2.2	Tools	11
2.3	Setup	12
2.4	Part A: Repository selection and Part B: Define Selection Criteria	12
2.5	Part C: Bug-Fixing Commit Identification:	14
2.6	Part D and E: Diff Extraction and Analyses and Rectification of the Message	17
2.6.1	Main Loop: CSV Processing and Analysis	20
2.7	Part F: Evaluation: Research Questions	24
2.7.1	Possible Reasons for high hit rate:	27
2.8	Discussion and Conclusion	29
3	Multi-Metric Bug Context Analysis and Agreement Detection in Bug-Fix Commits	31
3.1	Introduction	31
3.2	Tools	31
3.3	Setup	32
3.4	Part A and B: Compute and report baseline descriptive statistics	33
3.5	Part C: Structural Metrics with radon	39
3.6	Part D: Change Magnitude Metrics	43
3.7	Part E: Classification and Agreement	45
4	Exploration of different diff algorithms on Open-Source Repositories	49
4.1	Introduction	49
4.2	Tools	49

4.3	SETUP	50
4.4	Part A: Repository Selection and Part B: Define Selection Criteria	50
4.5	Part C: Run Software Tool on the Selected Repository	52
4.6	Part D: Compare Diff Outputs for Discrepancy Analysis	54
4.7	Part E: Report Final Dataset Statistics	55
4.8	Analysis and Conclusion	60
4.9	Part F: Which algorithm performed better	62

Chapter 1

Introduction to Version Controlling, Git Workflows, and Actions

1.1 Introduction

The aim of the laboratory session was to provide a hands-on introduction to version control using Git and GitHub. The activities focused on core concepts—initializing repositories, staging and committing changes, and synchronizing work—while also highlighting good practices for collaboration and project maintainance.

1.2 Tools

- **Programming Language:** Python 3.12.9 — Used for code and using pylint library.
- **Editor/IDE:** Visual Studio Code — Used for coding, debugging and execution.
- **Version Control:** Git and GitHub — Used to track and the changes in the code, and to improve maintainability of the codebases.
- **Linting Tool:** Pylint — Used to ensure that the code follows PEP8 standards.
- **Automation Platform:** GitHub Actions — used to automate the linting workflow and provide continuous integration (CI) feedback.
- **Virtual Environment:** venv — To prevent library version conflicts by isolating working environments.

1.3 Setup

To configure GitHub for this project, I had to setup git on my machine and Visual Studio Code (VS Code) as the code editor. For github setup I had to sign up using email ID and password;

To maintain isolated coding environment, created a new venv **lab1** with python version **3.12.9**.

1.4 Methodology and Execution

The following methodology was followed execute tasks.

Part A: Initial Compilation and Git Setup

1. Git initialization:

- Installed git on MacOS, and checked git version:

```
git --version
```

2. Version Control Initialization:

- Configured Git global username and email:

```
git config --global user.name "TejasLohia21"  
git config --global user.email 23110335@iitgn.ac.in
```

- Verified configurations:

```
git config --list  
git config user.name  
git config user.email
```

3. Repository Setup in a New Folder:

- A new folder named `testing_lab1` was created and entered:

```
cd testing_lab1
```

- A Git repository was initialized inside the folder:

```
git init
```

- A `README.md` file was created and initialized with content:

```
echo "# read_test" > README.md
```

- The file was staged and committed with a message:

```
git add README.md  
git commit -m "add readme"
```

4. Repository Setup in a New Folder:

- A new folder named `testing_lab1` was created and entered:

```
cd testing_lab1
```

- A Git repository was initialized inside the folder:

```
git init
```

- A `README.md` file was created and initialized with content:

```
echo "# read_test" > README.md
```

- The file was staged and committed with a message:

```
git add README.md  
git commit -m "add readme"
```

5. Pushing Code to GitHub:

- The local commits were pushed to the remote repository using:

```
git push -u origin main
```

- The `-u` flag sets the upstream branch so that subsequent pushes can simply use `git push`.

6. Checking Commit History:

- The commit history was viewed using:

```
git log
```

- Output:

```
commit 7e4d7da93c58274df903e73f653de9e9fe8e84fb (HEAD -> main)  
Author: TejasLohia21 <23110335@iitgn.ac.in>  
Date:   Sat Sep 6 01:00:00 2025 +0530  
  
add readme
```

Part C: Working with Remote Repositories

1. Connecting to GitHub:

- A new repository was created on GitHub named TejasLohialab1.
- The local repository was linked to GitHub using:

```
git remote add origin git@github.com:TejasLohia21/TejasLohialab1.git  
git branch -M main  
git push -u origin main
```

2. Pushing Changes to GitHub:

- The committed changes were pushed to GitHub using:

```
git push -u origin main
```

3. Cloning a Repository:

- An existing repository was cloned from GitHub to the local machine using:

```
git clone git@github.com:TejasLohia21/datascienc-HNSW.git
```

4. Pulling Changes:

- Updates from the remote repository were pulled using:

```
git pull origin main
```

Part D: Setting up Pylint Workflow with GitHub Actions

In this part, a continuous integration (CI) workflow was created using **GitHub Actions** to automatically check the Python code using `pylint`. The steps followed were:

1. **Python Script Creation:** A Python file `code.py` was created with more than 30 lines of code. The script implemented functions for factorial calculation, Fibonacci sequence generation, and prime number detection.
2. **Workflow Configuration:** A workflow file was created at the path:

```
.github/workflows/pylint.yml
```

The content of the workflow file is shown below:

```
name: Pylint Check

on: [push, pull_request]

jobs:
  lint:
    runs-on: ubuntu-latest
    steps:
      - name: Checkout repository
        uses: actions/checkout@v2

      - name: Set up Python
        uses: actions/setup-python@v2
        with:
          python-version: '3.12'

      - name: Install dependencies
        run: |
          python -m pip install --upgrade pip
          pip install pylint

      - name: Run pylint
        run: |
          pylint code.py
```

3. **Commit and Push:** The workflow file and the Python script were staged, committed, and pushed to the GitHub repository:

```
git add main.py .github/workflows/pylint.yml
git commit -m "Add Python script and pylint workflow"
git push
```

4. **Verification:** After pushing, the GitHub Actions workflow was triggered. The Python script was linted using `pylint`, and all errors were resolved until a green tick (✓) appeared, confirming successful execution.

This ensured that the code followed Python coding standards and passed linting checks automatically on every push.

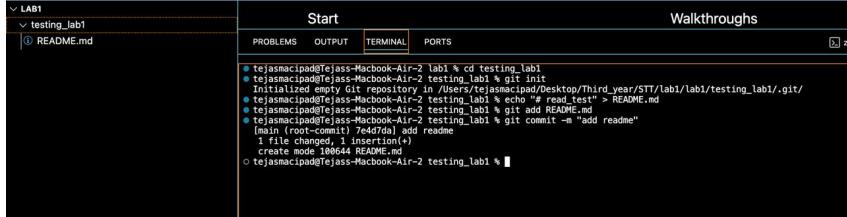
1.5 Results and Analysis

```
tejasmacipad@Tejass-Macbook-Air-2 lab1 % git --version
git version 2.39.5 (Apple Git-154)
tejasmacipad@Tejass-Macbook-Air-2 lab1 % git config --global user.name "TejasLohia21"
tejasmacipad@Tejass-Macbook-Air-2 lab1 % git config --global user.email 23110335@iitgn.ac.in
credential.helper=osxkeychain
init.defaultbranch=main
user.name=TejasLohia21
user.email=23110335@iitgn.ac.in
init.default=branch
init.defaultbranch=main
pull.rebase=false
core.excludesfile=/Users/tejasmacipad/.gitignore_global
filter.lfs.clean=git-lfs clean -- %f
filter.lfs.smudge=git-lfs smudge -- %f
filter.lfs.process=git-lfs filter-process
filter.lfs.required=true
core.repositoryformatversion=0
core.filemode=true
core.bare=false
core.logallrefupdates=true
core.ignorecase=true
core.precomposeunicode=true
remote.origin.url=git@github.com:TejasLohia21/lab1.git
remote.origin.fetch=refs/heads/*:refs/remotes/origin/*
branch.main.remote=origin
branch.main.merge=refs/heads/main
tejasmacipad@Tejass-Macbook-Air-2 lab1 % git config user.name
TejasLohia21
tejasmacipad@Tejass-Macbook-Air-2 lab1 % git config user.email
23110335@iitgn.ac.in
tejasmacipad@Tejass-Macbook-Air-2 lab1 %
```

Figure 1.1: Setting up git

Commands to check the version and verify initialization gave the following outputs.

- Git version initialized: 2.39.5 (Apple Git - 154)
- Git username: TejasLohia21
- Git user email: 23110335@iitgn.ac.in



```
LAB1
└── testing_lab1
    └── README.md

Start Walkthroughs
PROBLEMS OUTPUT TERMINAL PORTS

tejasmacipad@Tejass-Macbook-Air-2 lab1 % cd testing_lab1
tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % git init
Initialized empty Git repository in /Users/tejasmacipad/Desktop/Third_year/STT/lab1/lab1/testing_lab1/.git/
tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % touch README.md
tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % echo "# read test" > README.md
tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % git add README.md
[main (root-commit) 7e4d7da] add readme
1 file changed, 1 insertion(+)
create mode 100644 README.md
tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 %
```

Figure 1.2: Init git Repo and addition of README file

Upon committing the staged files in the main branch of the initialized local git repository:

```
[main (root-commit) 7e4d7da] add readme
1 file changed, 1 insertion(+)
create mode 100644 README.md
```

```

● tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % git log
commit 7e4d7da93c58274df903e73f653de9e9fe8e84fb (HEAD -> main)
Author: TejasLohia21 <23110335@iitgn.ac.in>
Date:   Sat Sep 6 01:00:00 2025 +0530

    add readme
○ tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 %

```

Figure 1.3: Checking commit history

Command git log generated the commit history along with Metadata.

```

Author: TejasLohia21 <23110335@iitgn.ac.in>
Date: Sat Sep 6 01:00:00 2025 +530

```

```
add readme
```

```

● tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % git remote add origin git@github.com:TejasLohia21/TejasLohialab1.git
● tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % git branch -M main
● tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % git push -u origin main
Enumerating objects: 3, done.
Counting objects: 100% (3/3), done.
Writing objects: 100% (3/3), 224 bytes | 224.00 KiB/s, done.
Total 3 (delta 0), reused 0 (delta 0), pack-reused 0
To github.com:TejasLohia21/TejasLohialab1.git
 * [new branch]      main -> main
branch 'main' set up to track 'origin/main'.
○ tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 %

```

Figure 1.4: Linking local repositories with github

Created a new repository on github named lab1. Using the command, the local repository was linked to Online Repo. The repository was then pushed to the online repository on github in the main branch.

```

● tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % git clone git@github.com:TejasLohia21/datasience-HNSW.git
Cloning into 'datasience-HNSW'...
remote: Enumerating objects: 89, done.
remote: Counting objects: 100% (89/89), done.
remote: Compressing objects: 100% (63/63), done.
remote: Total 89 (delta 34), reused 73 (delta 21), pack-reused 0 (from 0)
Receiving objects: 100% (89/89), 852.94 KiB | 948.00 KiB/s, done.
Resolving deltas: 100% (34/34), done.
● tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 % git pull origin main
From github.com:TejasLohia21/TejasLohialab1
 * branch      main    -> FETCH_HEAD
Already up to date.
○ tejasmacipad@Tejass-Macbook-Air-2 testing_lab1 %

```

Figure 1.5: Cloning existing repositories and initiating pull

To get a hands on experience of pulling repositories, I cloned an existing repository and pulled that. As the repository did not have any changes, the output was 'Already upto date'.

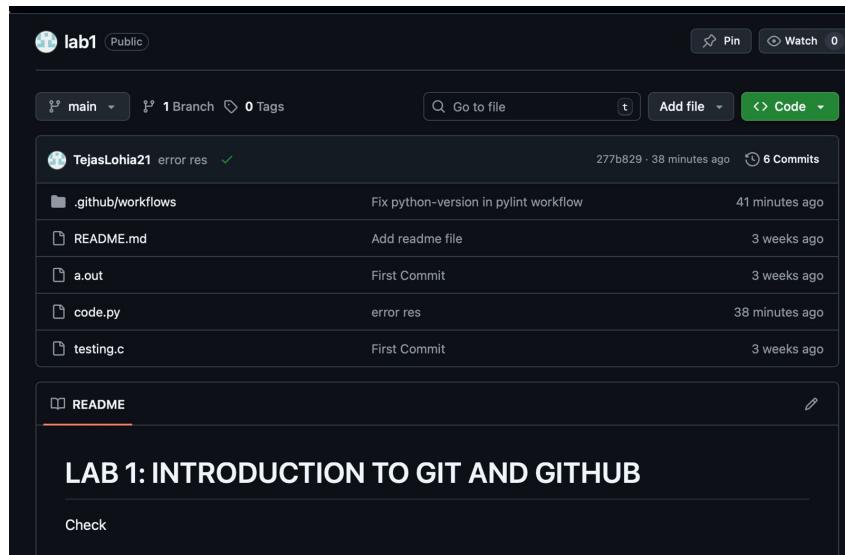


Figure 1.6: Pylint workflow

Initially there was a cross mark after pushing to the online repository. The error was because of no blank empty line in the end.

After rectifying this, there was a tick symbol in the github repository.

1.6 Discussion and Conclusion

This lab was quite important to understand and learn git in a systematic way. It introduced us to the very basics of git init, till using github workflows, making us familiar to use and maintain repositories using git and github.

Chapter 2

Commit Message Rectification for Bug-Fixing Commits in the Wild

2.1 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 also introduced about the quality of commit messaged in real-life projects which have huge userbases, and require active and effective maintainance. I was also introduced to different LLMs which could be used to classify the commits into various types of commits, as well as to understand the effectiveness of generating commit messages using LLMs. This task also made us to explore SEART search engine, which allows us to filter all the repositories on github based on various parameters.

2.2 Tools

- **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 library for mining software repositories, used to analyze git repositories and extract commit information.

- Pandas: Used for managing dataframes and working with CSV files.
- Pre-trained LLMs: SEBIS/code_trans_t5_base_commit_generation and CommitPredictorT5.

2.3 Setup

- GitHub account: TejasLohia21
- Created and activated a virtual environment named lab2 to manage library versions.
- Installed Python version 3.12.9.
- Cloned and initialized three repositories for experiments and version control.

2.4 Part A: Repository selection and Part B: Define Selection Criteria

METHODOLOGY AND EXECUTION

To choose a repository, I used criterias mentioned in the lecture to find real-life projects with active users. Used SEART search engine, to search for repositories which satisfied the criterias for filtering out repositories.

The screenshot shows the SEART search engine interface. At the top, there are links for Statistics, Publication, About, and a search icon. Below the header, there are several sections for filtering repositories:

- General**: Includes a search bar for "Search by keyword in name" with a dropdown for "Contains", and filters for "License" and "Has topic".
- History and Activity**: Filters for "Number of Commits" (1000 to 25000), "Number of Contributors" (min to max), "Number of Issues" (min to max), "Number of Pull Requests" (min to max), "Number of Branches" (min to max), and "Number of Releases" (min to max).
- Date-based Filters**: Filters for "Created Between" (09/07/2025 to 09/07/2025) and "Last Commit Between" (09/07/2025 to 09/07/2025).
- Popularity Filters**: Filters for "Number of Stars" (1000 to max), "Number of Watchers" (min to max), and "Number of Forks" (min to max).
- Size of codebase**: Filters for "Non Blank Lines" (min to max), "Code Lines" (min to max), and "Comment Lines" (min to max).
- Additional Filters**: Includes a "Sorting" section with dropdowns for "Name" and "Ascending", and a "Repository Characteristics" section with checkboxes for "Exclude Forks", "Only Forks", "Has License", "Has Open Issues", "Has Pull Requests", and "Has Wiki".

At the bottom center is a "Search" button.

Figure 2.1: SEART Search engine



Figure 2.2: BOXMOT repository

- **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 resolving them makes it an ideal repository for analysis.

Criteria for choosing this repository

This was chosen in reference to the funnel diagram in lecture 2 slides.

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).

I have previously also used this repository as this is central to implementations of all the state-of-the-art (**SOTA**) trackers benchmarked in research papers. This repository has a huge base of users, as well as there are active resolutions for the issues resolved, resulting in higher number of commits.

2.5 Part C: Bug-Fixing Commit Identification:

METHODOLOGY AND EXECUTION

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+)', re.IGNORECASE
10 )
11
12 EXCLUDE_RE = re.compile(
13     r'(readme|doc|docs|typo|chore|license|format|formatting|style|'
14     r'pre-commit|precommit|ci|workflow|bump version|version bump)', re.IGNORECASE
15 )
16
17
18 CODE_EXTS = ('.py', '.c', '.cc', '.cpp', '.cu', '.h', '.hpp')
19
20 def merge_check_commit(pathrepo: str, hascom: str) -> list[str]:
21     cmd = ['git', "-C", pathrepo, "show", "-m", "--name-only", "--pretty=", hascom]
22     result = subprocess.run(cmd, stdout=subprocess.PIPE, stderr=subprocess.PIPE, text=True, check=False)
23
24     if result.returncode != 0:
25         return []
26
27     paths = [
28         line.strip() for line in result.stdout.splitlines()
29         if line.strip() and line.lower().endswith(CODE_EXTS)
30     ]
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 2.3: 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,

```

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 2.4: Code for mine fixing.

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 and some other filtering. 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.

df.head()					
✓ 0.0s	Hash	Message	Hashes of parents	Is a merge commit?	List of modified files
0	b51eb87b4bce896550e647562a9526b34603eab9	Added possibility to write MOT compliant result...	718185a74e27c873833e840636f628153c73e7	False	track.py
1	21bc8b50d5727207b0fcfa440cb2c96dbdc542	Fixed video save bug	c41c8b5b514aa6487985e4ac50d2c070fe6564a06	False	track.py
2	5fb232162b2d7d9b6892c2e594192c217b31b3a45	increment ages when no detections fix	a0b0b0102c7956a9fb9d182246389c9b2c02923bf	False	track.py
3	6599919ba5723a7a3e5586952fa88212be483ae6	bugfix: writing of MOT16 compliant results	b117ddad1bdf8b0d494da9b73dcabe7f66f17fe	False	track.py
4	3b100c562af371b273b6679a72c53800f3afb69	fixed MOT compliant result bug, switched bbox_...	6599919ba5723a23e5586952fa8d812be483ae6	False	track.py

Figure 2.5: First five rows of the generated output.

Discussion and Conclusion

In this code, we implemented a pipeline that extracts all the commits, and respective Metadata and information is saved to a CSV file. The regular expression used to identify the commits related to bug-fixes is an efficient way of classifying commits, but lags in cases where commit messages might not explicitly contain keywords in the list. Merge commits were handled effectively using another function. This pipeline achieves the aim of storing the commit history in a CSV file for required for further analysis.

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

METHODOLOGY AND EXECUTION

In this, we work on a pipeline iterates through the csv file, and runs LLM models to classify the commits into fix-types. Another LLM was used to generate a rectified commit message. This task is fundamental to understand the quality of commit messages contributors commit with. Also, this analysis leads us to investigate the quality of the generated rectified messages by LLMs.

In the pipeline, we iterated through all the rows in the dataframe generated in the previous code and cleaned it to deal with abnormality in the strings.

To generate llm message, CommitPredictorT5 is instantiated, which is fed with code difference between current source code and previos codes.

Rectification

- 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 2000 to avoid hallucination of models which could have occured because of the prompt size.
- **Prompt designing:** Pileline implements a dynamic prompt based on the length of the source code.
 - Methodology: Models are expected to generate more relevant rectified commit message with more contexts. But, when models are provided with very long prompts, they tend to lose relevant information and often hallucinate.
 - Parameters in prompt: To provide better context, model is prompted with code difference, current source code, previos source code, human message and `llm_inference`.

- Reducing the prompt size: To avoid hallucination, we dynamically change prompt. Based on classification of the source code, based on its length being less than 2000, model changes prompt and does not pass source codes, if length > 2000.
- Expected improvement: Difference in source code provides crucial information, thus can't be sliced. Removal of codes might cause information lose, but would also prevent the model from hallucinating.

```

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 2.6: Code(Image 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 2.7: Code (Image 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 "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 MODEL_NAME = "mamiksik/CommitPredictorT5"
113 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
114 model = AutoModelForSeq2SeqLM.from_pretrained(MODEL_NAME)
115 model.to(device)
116 print("Model loaded.")
117
118 MODEL_NAME_2 = "SEBIS/code_trans_t5_base_commit_generation"
119 tokenizer_2 = T5Tokenizer.from_pretrained(MODEL_NAME_2, use_fast=False)
120 model_2 = AutoModelForSeq2SeqLM.from_pretrained(MODEL_NAME_2)
121 model_2.to(device)
122
123 print("Model and tokenizer for retification loaded successfully")
124
```

Figure 2.8: Code (Image 3).

```

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

Figure 2.9: Code (Image 4).

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.

```

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 2.10: Code (Image 5).

- **LLM-based commit message rectification** – To generate the rectified message, the diff and other parameters were passed to the model.

2.6.1 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.

Results and Analysis

Hash	Message	Filename	Source Code (.pyw)	Source Code (current)	DIF	LLM Inference	rectified message
11 9b0f512eacc2e07688a371c50e42982e95de	fix bad initial kf predictions for new objects.	deep_sort_pytorch/deep_sort/sorter/kalman_filter.py	# vim: expandtab:ts=4:sts=4:sw=4:et:ff=rump as g++	# vim: expandtab:ts=4:sts=4:sw=4:et:ff=rump as g++	-/+13,-13 +74,-53	add	Add note about data volume to enable_EQC
12 9b0f512eacc2e07688a371c50e42982e95de	fix bad initial kf predictions for new objects.	track.py	import bayesian_track.metrics	# import bayesian_track.metrics	-/-12/-1 +12/-7	update	Addition T-100 -17 to Changing
13 4531842bd4dd50e42982e95de	Fix frame ID M1 export bug	track.py	# limit the number of cups used by high performance	# limit the number of cups used by high performance	-/+15,-15 +15,-15	update	Add note about data volume to enable_EQC
16 2e62f49d02704999f33ca50e42982e95de	dt fix bug	track.py	# limit the number of cups used by high performance	# limit the number of cups used by high performance	-/+13,-13 +17,-17	unknown	from yolov5.models import model
17 f056f9a620765b0d3881527277e21866026	fix https://github.com/mikel-bromley/YOLOv5_D...	track.py	# limit the number of cups used by high performance	# limit the number of cups used by high performance	-/+205,-205 +205,-205	update	Add note about data volume to enable_EQC

Figure 2.11: First five rows of the generated output.

Commit 1

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

Commit 2

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

Commit 3

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

Commit 4

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

Commit 5

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

Discussion and Conclusion

In this pipeline, two LLMs were used to generation fix-type and rectified commit message. Many of the original commits, had same commit message for multiple modified files, resulting in loss of information for the exact change. Generated rectified messages are specific for file changes and are based on in depth context about the changes, thus making commit messages more relevant.

Initially CommitPredictorT5 was testeed to generate rectified commit message. Despite providing in depth context and specific prompt, it was not able to generate the rectified message properly, because of lack of fine-tuning. `SEBIS/code_trans_t5_base_commit_generation` was used to generate rectified messages, as it is trained for this specific task.

2.7 Part F: Evaluation: Research Questions

RQ1 (Developer eval.)

This question intends to investigate the precision of commit message and quantify hit rate.

To quantify the precision, pipeline uses an encoder **microsoft/codebert-base** to get embeddings for code difference and commit message. Then to obtain similarity, cosine distance is used and applied over a THRESHOLD to obtain HIT RATE.

- **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.

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.

```

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

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 2.12: Code for Similarity analysis.

- **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**.

RQ2 (LLM eval.)

This question intends to investigate the precision of fix-type and quantify hit rate.

To quantify the precision, pipeline uses an encoder **microsoft/codebert-base** to get embeddings for code difference and commit message. Then to obtain similarity, cosine distance is used and applied over a THRESHOLD to obtain HIT RATE.

- 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.

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

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)
A] ✓ 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 2.13: Semantic similar between fix type (LLM) embedding and commit message embedding.

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.

2.7.1 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.

RQ3 (Rectifier eval.)

This section intends to analyze the rectification in the commit message. We take the difference between the similarities between human commit message and rectified commit message, and then take Average for all the modifications.

- Aim: – This section requires to analyze the amount of rectification done, to generate 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.

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.

```

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()}")
print(f"Output path: {output_path}")

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: 0.9307055559414125
Saved CSV with rectification and hit rate to: /Users/tejasmacipad/Desktop/Third_year/STT/lab2/newmodelfile_with_rectification_codebert.csv

```

Figure 2.14: Code for rectification and improvement

- **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%

2.8 Discussion and Conclusion

This was an interesting investigation to understand the LLM outputs. I was surprised by the score obtained even with the high similarity obtained upon human commit message.

Reason for this could be the generic commit messages, which might match the code differences.

Rectified messages had an improvement but not as high as expected. Possible reason for this could be Initial high scores, but also when we have huge codes which are being embedded, it starts gettin generalized because of encoding of every line in the code combined into a single vector. Even if the embeddings did not show much of improvement, but checking the rectified commit messages manually, I found them to be much more relevant.

Chapter 3

Multi-Metric Bug Context Analysis and Agreement Detection in Bug-Fix Commits

3.1 Introduction

The aim of this laboratory was to give us an hands on experience to investigate the bug fix commits using structural metrix and similar measure between the previous and current versions of the source code. The structural matric was analyzed using methods of Maintainability Index(MI), Cyclomatic Complexity(CC), and Lines of Code(LOC) for each bugfix. To analyze the similarity we use semantic similarity and token similarity generated by models like codeBERT and BLEU. We analyzed the results based on code quality improvement (or changes) and the extend of changes. Based on this changes were analyzed to classify the commits into major bug and minor bug fixes, also the conflicting assessments between both the models.

3.2 Tools

- **Software and Tools Used:**
 - Operating System: macOS
 - Text Editor: Visual Studio Code
 - Version Control: Git
 - Remote Hosting: GitHub
 - Continuous Integration: GitHub Actions
- **Repository:**
 - `boxmot` – An open-source library that extends trackers by integrating state-of-

the-art multi-object tracking (MOT) algorithms. It provides implementations of popular trackers such as DeepSORT, StrongSORT, and ByteTrack, allowing seamless combination of object detection and tracking. The library is modular, easy to integrate with detection pipelines, and designed for real-time performance, making it suitable for research as well as production-level applications.

- **Python Libraries:**

- `pydriller` – Framework for mining software repositories and extracting commit information.
- `pandas` – For managing dataframes and extracting data from CSV files.
- `collections.Counter` – To create dictionaries with default values and count occurrences of elements.
- `matplotlib.pyplot` – For data visualization and plotting.
- `os` – Interacting with the operating system.
- `radon` – Code analysis tool to compute metrics like LOC, Maintainability Index, and Cyclomatic Complexity.
- `re` – Regular expressions for string pattern matching and text refinement.
- `subprocess` – Running external commands and interacting with system processes.
- `sentence_transformers` – For semantic similarity checking using pretrained transformer models.
- `torch` – PyTorch deep learning framework.
- `nltk` (BLEU) – Natural Language Toolkit, used for evaluating BLEU scores.

3.3 Setup

- GitHub account: `TejasLohia21`
- SSH key configured for secure push/pull
- Virtual Environment named `lab3` to manage library versions
- Configured git username and email
- Cloned/initialized three repositories

3.4 Part A and B: Compute and report baseline descriptive statistics

This part aims us analyze the fundamental properties of the Commits and repositories from the data generated in previous lab session.

Code analyses the number of commits and files, which give us an idea of the size of the repositories, and maintainance changes. It becomes easier to mantain, if the number of changes made in a single commit are less, as it provides better navigation and specific understanding to each of the change made. Thus we also analyze the number of file modifications per commit.

In the previous lab, we used CommitPredictorT5 to classify bug-fix commits into various types which allows us to interprete the changes required, or the type of errors that exist in the repository.

Its also important to know, which type of files have most bug-fixes commits, as this allows to deal differently with different types of files to handle efficiently. Some files might be critical and changes as well as bug-fixing in them might be a big deal, thus we also analyze the files which have most number of modifications

METHODOLOGY AND EXECUTION

Loaded the csv "diff.csv" generated in Lab2 using pandas into a dataframe "df"

- **Compute and report baseline descriptive statistics:**
 - Using nunique on the Hashes stored in the df to obtain the unique number of hashes
 - Using nunique on the Filename stored in the df to obtain the unique number of files.
 - Stored all the commits and the file counts for respective commits, and then reported the average number of modified files per commit.
 - To find the fix types from LLM, we use a dictionary to store the keywords used in commit messages and map them to six generic categories. We run through every commit message and if any word is in dictionary mapping, we return the type of that keyword.
 - Used OS and counter libraries to count the top ten most frequently modified files, and extensions.
- **Codes: Total number of commits and files**
 - **Explanation of Code for total commits and files:**

```

import pandas as pd
# Import module
# ✓ 0.0s

df = pd.read_csv("diff.csv")
# ✓ 0.2s

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

# df.head(2)
# ✓ 0.0s

      Hash           Message   Filename Source Code (prev) Source Code (current) Diff          LLM Inference rectified message
0  6d7eab7b40ce89650d4d47562a952635463ea9b  Added possibility to write MOT compliant track.py import argparse;import os;import ... import argparse;import os;import ... @@@ -114,6 +114,10 @@ def detectobj, update webcam detection
1  2fbcb850d87727079fbc2ca44fc52c98db5c542  Fixed video save bug  track.py  from yolov5.utils.datasets import LoadImages, ... from yolov5.utils.datasets import LoadImages, ... @@@ -45,11 +50,11 @@ def draw_bboxes(img, bbox, ...
Total number of commits
# ✓ 0.0s
Print("Total number of hashes are {len(df['Hash'])}, while the unique number hashes are {df['Hash'].nunique()}")
# ✓ 0.0s
Print("Total number of files are {len(df['FileCount'].unique())} which have changed in them")
# ✓ 0.0s
Print("Total number of hashes are 727, while the unique number hashes are 388")
# ✓ 0.0s
Print("Total number of files are 727, while the unique number of files which have been modified are 247")
# ✓ 0.0s

```

Figure 3.1: Loading of diff.csv file and finding total number of commits and files.

- * Loaded the row which stores hashes, and used .nunique to find the total unique hashes.
- * Loaded the row which stores Filename, and used .nunique to find the total unique files modified.

– Results:

- * Total number of hashes in the dataframe are **727**, while the unique number of hashes are **388**.
- * Total number of files in the dataframe are **727**, while the unique number of files which have been modified are **247**.

• Codes: Average number of files modified per commit

```

dict = {}
for index, row in df.iterrows():
    if row["Hash"] not in dict:
        dict[row["Hash"]] = 1
    else:
        dict[row["Hash"]] += 1

avg_value = 0

for _, count in dict.items():
    avg_value += count

avg_value /= len(dict)
print(f"Average number of modified files per commit: {avg_value}")
del dict

```

Figure 3.2: Average numbers of modified files per commit.

– Explanation of Code for Average file modifications per commit:

- * For every commit, maintain a dictionary which maps it to a list of files which have been modified.
- * Iterated thru' all the rows, as each row contained information about changes in distinct modifications in files.

- * Sum of all the values in dictionary divided by the length of dictionary yielded the average number of files modified per commit.

– **Results:**

- * Average number of files modified per commit are **1.8737**.

• **Codes: Distribution of number of fix types**

```

import pandas as pd
from collections import defaultdict

fix_types = {
    "add": ["add", "create"],
    "update": ["update", "upgrade", "refactor", "restructure", "clean", "optimize", "improve", "enhance", "performance"],
    "remove": ["remove", "delete", "drop"],
    "fix": ["fix", "bug", "error", "correct"],
    "docs": ["docs", "documentation", "readme"],
    "test": ["test", "tests", "config", "build"]
}
mappings = {s2: m for m, s in fix_types.items() for s2 in s}
for m, s in fix_types.items():
    for s2 in s:
        mappings[s2] = m

def classify(commit_message):
    words = commit_message.lower().split()
    for word in words:
        if word in mappings:
            return [mappings[word]]
    return []

def counttype(df, column):
    count_lst = defaultdict(int)
    for msg in df[column].dropna():
        typematch = classify(msg)
        for t in typematch:
            count_lst[t] += 1
    return dict(count_lst)

counts = counttype(df, "LLM Inference")
print("Fix type distribution:", counts)

Fix type distribution: {'update': 497, 'add': 191, 'fix': 37, 'remove': 1}

```

Figure 3.3: Distribution of number of fix types.

– **Explanation of Code for Fix types distribution:**

- * As the commit messages with same fix type may have different keywords, initialized a dictionary fix_types which contains mapping from six keywords to list of words which contain frequently used keywords of that fix type.
- * Created another dictionary, 'mappings' which contains map from each of the word in list to the fix type it belongs to.
- * Classify function iterates through all the words in each commits, and returns the type using the dictionary mapping.
- * counttype function iterates through each of the commit and calls for the function classify and increments the count of the type returned by function classify.

– **Results:**

- * Distribution for the fix type is:
- * update – 497

- * add – 191
- * fix – 37
- * remove – 1

- **Codes: Most frequently modified filenames/extensions**

```

import os
import matplotlib.pyplot as plt
from collections import Counter

def extret(filename):
    return os.path.splitext(filename)[1] if "." in filename else "no_ext"

def freqmax(df, column="filename", top_n=10):
    files = df[column].dropna().tolist()
    file_counts = Counter(files)
    ext_counts = Counter(extret(f) for f in files)

    return file_counts.most_common(top_n), ext_counts.most_common(top_n)

file_counts, ext_counts = freqmax(df, column="Filename")

print("Top modified filenames:")
for f, c in file_counts:
    print(f"{f} → {c}")

print("\nTop modified extensions:")
for e, c in ext_counts:
    print(f"{e} → {c}")

```

Python

Figure 3.4: Most frequently modified filenames/extensions.

- **Explanation of Code for Most frequently modified filenames/extensions:**
 - * 'df' is passed to the function freqmax which extracts the column Filename to extract extensions using another function extret.
 - * Freqmax function applies the instantiates the counter object over files to find the top_n number of files, and extret function returns the extension for all the files.
 - * Another counter object is instantiated to create a count for all the extensions
- **Most modified extensions:**
 - * Distribution for the fix type is:
 - * .py – **699**
 - * .yml – **14**
 - * .yaml – **6**

* .toml – **2**

* .txt – **1**

* .gz – **1**

* .lock – **1**

* .md – **1**

– **Frequently modified files:**

* Distribution for the fix type is:

* tracking/val.py – **40**

* track.py – **32**

* boxmotutils__init__.py – **24**

* val.py – **14**

* boxmottrackersbotsortbot_sort.py – **14**

* boxmotappearancereid_export.py – **13**

* trackingevolve.py – **13**

* githubworkflowsci.yml – **11**

* boxmotappearancereid_multibackend.py – **11**

* boxmottrackersbyteTrack_tracker.py – **11**

Plots

- Fix_type distribution

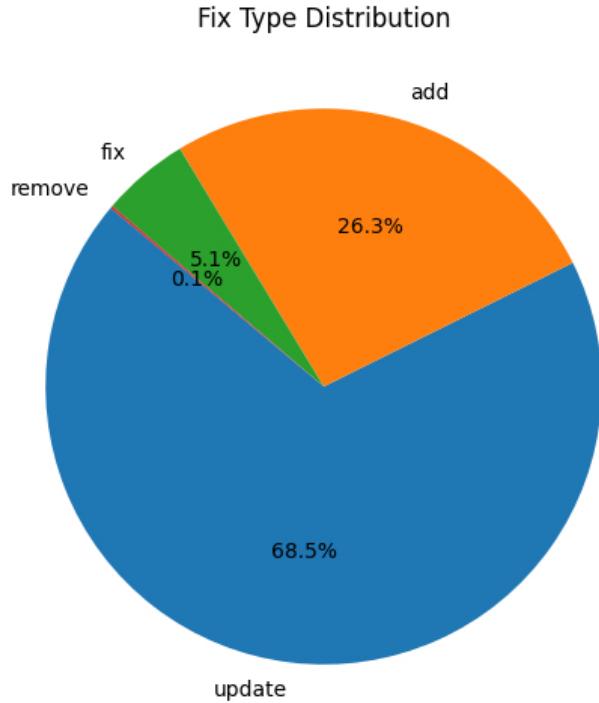


Figure 3.5: Fix_type distribution.

- frequently modified files

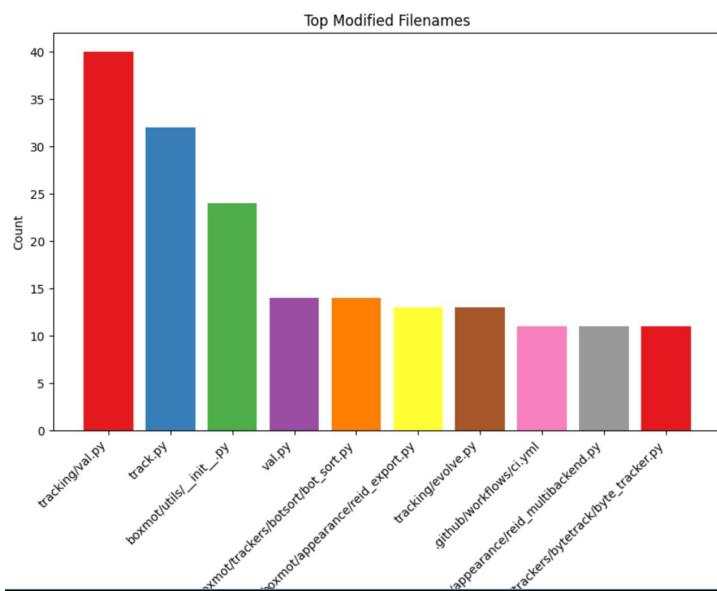


Figure 3.6: Top modified filenames.

- frequently modified extensions

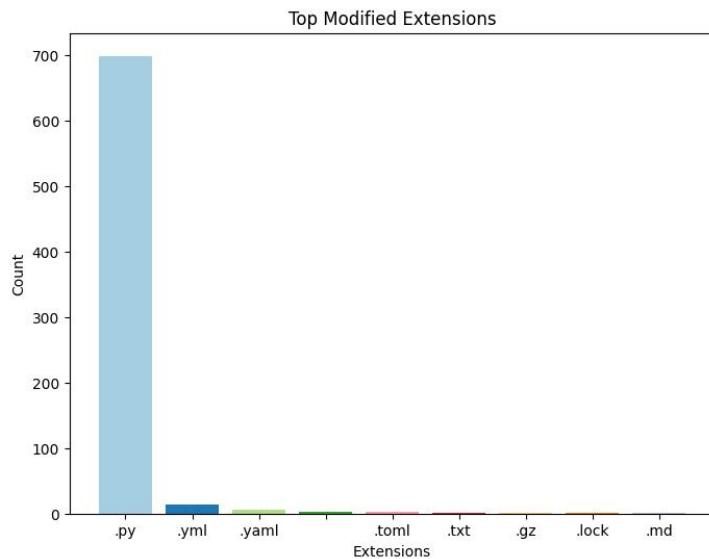


Figure 3.7: Top modified filenames.

Analysis and Conclusion

After running the code, we find that out of **727** rows in the CSV file, there were **388** commits and **247** unique files modified. This implies that same files were modified multiple times. Also, the average number of files modified per commit is **1.87**, which means that more than a single file, must be having same commit message, which results in loss of specific information in the commit message

Interestingly, highest number of fix types were of **update** follows by **add**. Most modified extension was obviously of **.py** format, which is obvious citing that python was the primary language. Also, most modified files were **val.py** and **track.py** which is because they contain the most central functions in the tracking pipelines.

3.5 Part C: Structural Metrics with radon

METHODOLOGY AND EXECUTION

Loaded the csv "diffprocessed.csv" generated in Lab2 using pandas into a dataframe "df"

- Three types of results are obtained in Structural metric.
 - Maintainability Index: Measure how well the code is structured in terms of easiness to understanding the code, modifications and extension of the code. It compresses static code properties on a scale of 0, 100 to obtain a score. This method uses parameters; Size of Code, Control flow Complexity and token level volume.

- Cyclomatic Complexity: Measures the Complexity of control flow of the code. It counts the number of independent execution path through a piece of code. The higher the Complexity, implies more the number of paths, making it difficult to test and debug.
- Lines of Code (LOC): It counts the number of lines present in the source code.
- Change in MI: MI tells whether the bug-fix commit made the code easier to interpretable.
- Change in Lines of Code: Indicates the Change in the lines of code, and this is the simplest parameter used for classification, as changes in few lines of code indicate minor changes.
- Cyclomatic Complexity: Increase in the CC indicates, that the number of paths have increased making the code complex to maintain.

- **Codes:**

```

import numpy as np
from radon.metrics import mi_visit, analyze
from radon.complexity import cc_visit

csv_path = "diff.csv"
df = pd.read_csv(csv_path)

def radonmet(source_code):
    if not isinstance(source_code, str) or not source_code.strip():
        return np.nan, np.nan, np.nan
    cleaned_code = source_code.encode('ascii', 'ignore').decode('utf-8')
    try:
        loc = analyze(cleaned_code).sloc
        mi = mi_visit(cleaned_code, multi=True)
        cc_results = cc_visit(cleaned_code)
        total_cc = sum(block.complexity for block in cc_results)
        return mi, total_cc, loc
    except Exception:
        return np.nan, np.nan, np.nan

metrics_before = []
metrics_after = []

for index, row in df.iterrows():
    code_before = str(row["Source Code (prev)"])
    code_after = str(row["Source Code (current)"])
    metrics_before.append(radonmet(code_before))
    metrics_after.append(radonmet(code_after))

```

Figure 3.8: Loading of diff.csv file and finding total number of commits and files.

```

df[['MI_Before', 'CC_Before', 'LOC_Before']] = metrics_before
df[['MI_After', 'CC_After', 'LOC_After']] = metrics_after

df['MI_Change'] = df['MI_After'] - df['MI_Before']
df['CC_Change'] = df['CC_After'] - df['CC_Before']
df['LOC_Change'] = df['LOC_After'] - df['LOC_Before']

display_cols = [
    'Hash', 'Filename', 'MI_Before', 'MI_After', 'MI_Change',
    'CC_Before', 'CC_After', 'CC_Change',
    'LOC_Before', 'LOC_After', 'LOC_Change'
]
if 'Filename' not in df.columns:
    display_cols.remove('Filename')

print(df[display_cols].to_string())
metric_cols = ['MI_Before', 'CC_Before', 'LOC_Before', 'MI_After', 'CC_After', 'LOC_After', 'MI_Change', 'CC_Change', 'LOC_Change']
cleaned_df = df.dropna(subset=metric_cols)
output_filename = 'diffprocessed.csv'
cleaned_df.to_csv(output_filename, index=False)

```

Figure 3.9: Loading of diff.csv file and finding total number of commits and files.

– Explanation of the code:

- * Loaded the dataframe processed in the previous question to extract the source code before commit and the source code after commit.
- * Function `radonmet` which filters the code and asserts the code to be a python instance. Function calls for inbuilt functions from the radon library to extract Lines of Code (LOC), Maintainability Index (MI) and Cyclomatic Complexity (CC).
- * Code iterates through all the rows in the loaded dataframe, and uses the defined function `radonmet` to obtain Maintainability index, Cyclomatic Complexity and Lines of Code respectively for previous code and current code.
- * We call this function on both the source code before and source code after the commit and store them in lists.
- * In the main dataframe, six columns are added for Maintainability Index, Cyclomatic Complexity and Lines of Code before and after respectively.
- * Three more columns are appended in the dataframe which store the difference between MI, CC and LOC.
- * Because of the presence of some characters which the csv files are not able to store properly, radon returns errors which are handled as an exception in the `radonmet` function by returning a NaN value. These rows are filtered out to avoid processing errors further.

Results

CSV file was updated with addition of 9 columns

- Average of Maintainability Index before: 55.32

-
- Average of Maintainability Index after: 55.06
 - Average of Cyclomatic Complexity before: 40.46
 - Average of Cyclomatic Complexity after: 40.57
 - Average of LOC before: 178.80
 - Average of LOC after: 181.20
 - Average MI change: -0.26
 - Average CC change: +0.12
 - Average LOC change: +2.39

Analysis and Conclusion

The commits are intended to remove bugs and improve the structural metrics for the code. The pipeline here used inbuilt libraries to generate metric of MI, CC and LOC. Here the commits did not have any major improvement in these metrics, as they aim to bug fix without emphasizing much on metrics.

3.6 Part D: Change Magnitude Metrics

METHODOLOGY AND EXECUTION

Loaded the csv "diffprocessed.csv" generated in previous section.

- **Change Magnitude Metrics:**

- **Bilingual Evaluation Understud (BLEU)** – Its a metric designed for machine translation evaluation. BLEU provides an evaluation score to find textual similarity between two sequences.
- **CodeBERT** – A pre-trained model for source code and natural language, used to compute semantic similarity between two code snippets.
- **How its an Evaluation Method** – Similarity score between source code before and after the commit and tokens of codes give information about the extend of change in the files modified in the commit.
- **classification** – Similarity score between codes and token could be set as thresholds for classification of the commit into major and minor change category.

- **Codes:**

- **Explanation of the code:**

- * The CodeBERT model from HuggingFace's `transformers` library is used to compute semantic similarity between the `Source Code (prev)` and `Source Code (current)`.
 - * Each source code snippet is tokenized and passed through CodeBERT to generate embeddings.
 - * Cosine similarity is then calculated between the embeddings of the two code snippets, and this value is stored in a new column `CodeBERT_Similarity`.
 - * In cases where the code before or after the commit is missing (such as additions or deletions), a similarity score of 0.0 is assigned.
 - * Similarly, the BLEU model is applied as a complementary metric. The `BLEU_Similarity` column is appended to the dataframe to capture token-level similarity.
 - * The code iterates over each row, splits the source code into tokens, and applies the `sentence_bleu` function from NLTK with smoothing to compute a BLEU score.
 - * If either the previous code or the current code is missing, a BLEU score of 0.0 is assigned to handle such cases gracefully.

- * Thus, the two models provide different perspectives: CodeBERT for semantic similarity (meaning of code), and BLEU for token similarity (surface-level code overlap).

```

import pandas as pd
import nltk
from nltk.translate.bleu_score import sentence_bleu, SmoothingFunction
nltk.download("punkt")

csv_path = "diffprocessed.csv"
df = pd.read_csv(csv_path)

df["BLEU_Similarity"] = 0.0

smooth_fn = SmoothingFunction().method1

for i, row in df.iterrows():
    prev_code = str(row["Source Code (prev)"])
    curr_code = str(row["Source Code (current)"])

    prev_tokens = prev_code.split()
    curr_tokens = curr_code.split()

    if prev_tokens and curr_tokens:
        bleu = sentence_bleu([prev_tokens], curr_tokens, smoothing_function=smooth_fn)
    else:
        bleu = 0.0

    df.at[i, "BLEU_Similarity"] = bleu

    if i % 25 == 0:
        print(f"{i}/{len(df)} processed...")

out_csv = "diffprocessed.csv"
df.to_csv(out_csv, index=False)
print("done ->", out_csv)

```

Figure 3.10: Applications of BLUE Model over source code before and after.

```

import pandas as pd
import torch
from transformers import RobertaTokenizer, RobertaModel
from sklearn.metrics.pairwise import cosine_similarity

csv_path = "diffprocessed.csv"
df = pd.read_csv(csv_path)

tokenizer = RobertaTokenizer.from_pretrained("microsoft/codebert-base")
model = RobertaModel.from_pretrained("microsoft/codebert-base")
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
model.to(device)
model.eval()

df["CodeBERT_Similarity"] = 0.0

def get_codebert_embedding(text):
    inputs = tokenizer(text, return_tensors="pt", truncation=True, padding=True, max_length=512)
    inputs = {k: v.to(device) for k, v in inputs.items()}
    with torch.no_grad():
        outputs = model(**inputs)
        embedding = outputs.last_hidden_state[:, 0, :].cpu().numpy()
    return embedding

for i, row in df.iterrows():
    prev_code = str(row["Source Code (prev)"])
    curr_code = str(row["Source Code (current)"])
    if prev_code.strip() and curr_code.strip():
        prev_emb = get_codebert_embedding(prev_code)
        curr_emb = get_codebert_embedding(curr_code)
        similarity = cosine_similarity(prev_emb, curr_emb)[0][0]
    else:
        similarity = 0.0
    df.at[i, "CodeBERT_Similarity"] = similarity
    if i % 25 == 0:
        print(f"{i}/{len(df)} processed...")

out_csv = "diffprocessed.csv"
df.to_csv(out_csv, index=False)

```

Figure 3.11: Applications of CodeBERT Model over source code before and after.

Results

- Average BLEU_similarity score is 92.76%.
- Average CodeBERT_similarity between source code tokens 99.95%.

3.7 Part E: Classification and Agreement

METHODOLOGY AND EXECUTION

Loaded the csv `diffprocessed.csv` generated in previous section.

- **Classification and Agreement:**
 - Classification of commits is important for maintenance and risk assessments.
 - **Classification:** Commits can be classified as a major commit, which implies significant changes, while a minor commit, which implies small changes.
 - Commits which are classified as major commit need more testing for maintenance purposes, thus making this important.
 - **Method:** Once the similarity values between the previous code and the current code is obtained, certain thresholds need to be set for classification tasks.
 - To find an approximate THRESHOLD value, code plots the similarity vs frequency histogram and chooses the optimal value of 0.97 for BLEU score and CodeBERT model.
 - The `Classes_Agree` column is introduced to check whether the classifications from BLEU similarity and CodeBERT similarity are consistent.
 - If both methods give the same label, it is marked as “YES”, otherwise it is marked as “NO”, thus helping to analyze the alignment between semantic and token-level classifications.

- **Codes:**

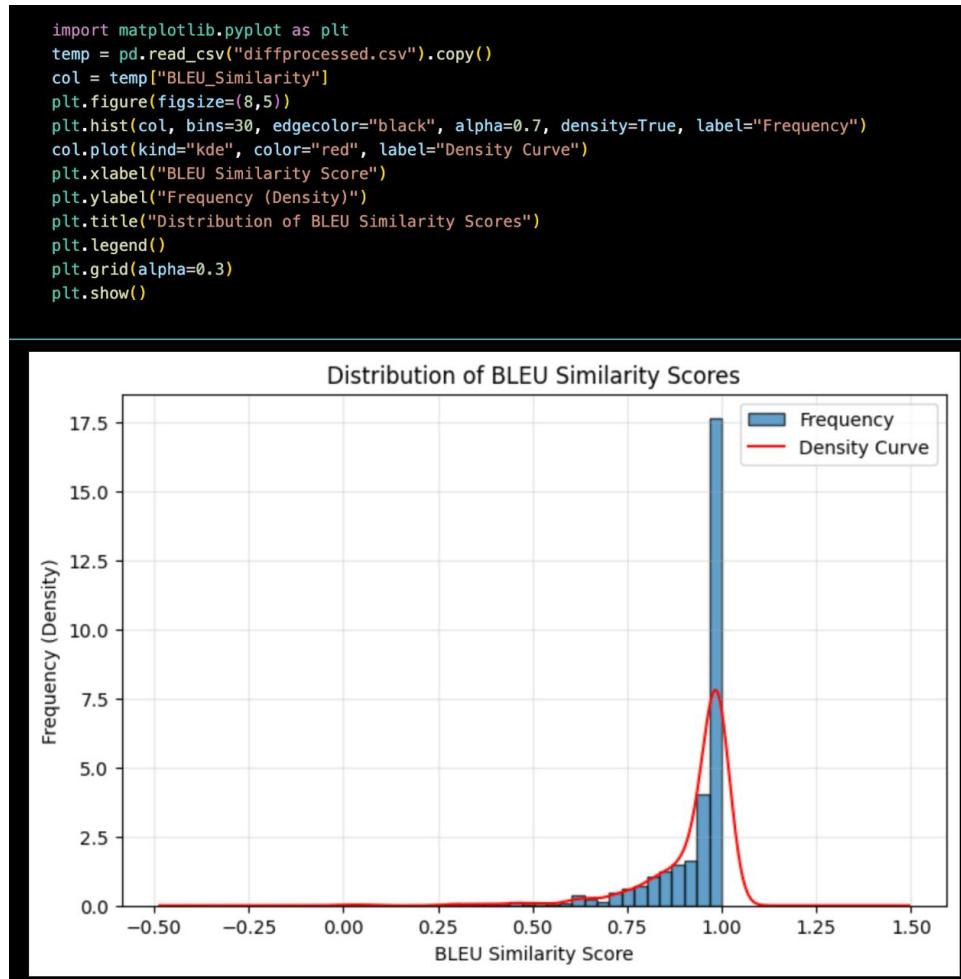


Figure 3.12: histogram plot on BLEU similarity score

```
import pandas as pd
df = pd.read_csv("diffprocessed.csv").copy()
cutoff = 0.97
df["Semantic_class"] = df["BLEU_Similarity"].apply(lambda x: "High" if x <= cutoff else "Low")
df["Token_class"] = df["CodeBERT_Similarity"].apply(lambda x: "High" if x <= cutoff else "Low")
df.to_csv("diffprocessed.csv", index=False)
```

Python

Figure 3.13: Classification based on THRESHOLD

- **Explanation of the code:**

- Plotted the Distribution histogram for the BLEU scored and set threshold of 0.97.
- Based on the threshold value, Semantic and Token_class columns are generated.
- These columns are binary based on classification at THRESHOLD values.
- If classification in both the values are equal, it is classified as Agree.

- **Results:**

```

import pandas as pd
df = pd.read_csv("diffprocessed.csv")
df["Classes_Agree"] = df.apply(
    lambda row: "YES" if row["Semantic_class"] == row["Token_class"] else "NO",
    axis=1
)
# df.to_csv("diffprocessed.csv", index=False)

print(df[["Semantic_class", "Token_class", "Classes_Agree"]].head())
print(f"Number of commits which are major according to CodeBERT Model {(df['Token_class'] == 'High').sum()}")
print(f"Number of commits which are major according to Bleu Tokens {(df['Semantic_class'] == 'High').sum()}")
print(f"Number of commits where methods agree {(df['Classes_Agree'] == 'YES').sum()}")


✓ 0.1s
Python

```

Semantic_class	Token_class	Classes_Agree	
0	High	Low	NO
1	High	Low	NO
2	Low	Low	YES
3	Low	Low	YES
4	High	Low	NO

Number of commits which are major according to CodeBERT Model 37
Number of commits which are major according to Bleu Tokens 291
Number of commits where methods agree 341

Figure 3.14: Model Agree

- Number of commits which are major according to CodeBERT Model 41, which is low because of the THRESHOLD equal to bleu threshold.
- Number of commits which are major according to Bleu Tokens 291
- Number of commits where methods agree 341.

Analysis and Conclusion

Analysis of the extent of semantic changes are very important. In repositories with multiple collaborators, its important to have checks on major fixes so that people are aware of changes made, which could be problematic to compile with other codes. In the section C, we used fundamental methods to find the extent of changes such as LOC, MI and CC. In the section D and E, we intended to find the extent of change semantically, which may detect changes in the functions and their outputs, which might get missed in the previous section if they are structurally same.

For this we use two models BLEU which is based on machine translation evaluation techniques, while CodeBERT is used to compute semantic similarity between two code snippets. We also contrast the outcomes of these models on the same code differences.

We run these models, to find cosine similarity in the embeddings generated by both the models. CodeBERT gave a very high average similarity value.

To use generated similarity index, we need to set some uniform threshold which could classify the changes into Major and Minor bug-fixes.

To set a THRESHOLD, we used the frequency vs similarity histogram. CodeBERT was not considered because, it had a very high similarity value because of which it might not have much of information gain.

Plots showed that the histogram had highest frequencies at somewhere around 0.97, therefore this was set as the THRESHOLD, for classification for both the models to maintain uniformity.

Commits where both the models agreed were 341, which is low because of the high similarity values obtained by CodeBERT model which made it classify most of the changes as minor fixes.

Chapter 4

Exploration of different diff algorithms on Open-Source Repositories

4.1 Introduction

The aim of the laboratory session was to provide gain a hands-on experience on how different algorithms like Myers and Histogram behave on real-world open source repositories. It also made us explore difference between code artifacts and non-code artifacts in the context of software development. It also aims to make us understand the difference in comparisons when the code ignores whitespace and blank lines. Importantly, how discrepancies can highlight the strengths of diff algorithms. After the lab, we aim to be proficient in analyzing diff outputs, generate datasets, visualize mismatches using python plots.

4.2 Tools

- **Software Used:**
 - Operating System: macOS
 - Text Editor: Visual Studio Code
 - Version Control: git
 - Remote Hosting: GitHub
 - Continuous Integration: GitHub Actions
- **Repositories:**
 - `butterknife` – A popular view binding library for Android that helps reduce boilerplate code when working with UI components.

-
- `retrofit` – A type-safe HTTP client for Android and Java, widely used for making API requests and handling network operations.
 - `glide` – An image loading and caching library for Android, optimized for smooth scrolling and efficient resource usage.
- **Python Libraries:**
 - Pydriller: A Python framework for mining software repositories. Used to analyze git repositories and extract commit information.
 - Pandas: To manage dataframes and deal with CSV files.

4.3 SETUP

- GitHub account: TejasLohia21
- SSH key configured for secure push/pull
- Virtual Environment named `lab4` to manage library versions
- Configured git username and email
- Cloned/initialized three repositories

4.4 Part A: Repository Selection and Part B: Define Selection Criteria

- `butterknife` – A popular view binding library for Android that helps reduce boilerplate code when working with UI components.
- `retrofit` – A type-safe HTTP client for Android and Java, widely used for making API requests and handling network operations.
- `glide` – An image loading and caching library for Android, optimized for smooth scrolling and efficient resource usage.

Methodology and Execution

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

- `butterknife` – A popular view binding library for Android that helps reduce boilerplate code when working with UI components.
 - **Commits:** The repository has 1016 commits, which is within the acceptable range for analysis.

-
- **Stars:** The repository has 28.8k stars, indicating strong community usage and real-world relevance.
 - **Myers Algorithm:** Since Git uses Myers as the default diff algorithm, we can directly apply `git diff -diff-algorithm=myers` on Butterknife.
 - **Histogram Algorithm:** Git also supports the Histogram algorithm, which works with the Butterknife repository using `git diff -diff-algorithm=histogram`.
 - **Discrepancy Analysis:** With over a thousand commits and frequent UI-related code changes, Butterknife provides meaningful opportunities to compare Myers and Histogram diffs.
 - **Required Files:** The repository includes `butterknife/` source code, a dedicated `tests/` suite, along with `README.md` and `LICENSE`, thus satisfying all the required file-type conditions.
 - **Conclusion:** Butterknife is one of the academically researched repositories for discrepancy analysis in both the Myers and Histogram algorithms.
 - **requests** – A popular Python library for making HTTP requests in a simple and human-friendly way.
 - **Commits:** The repository has 6372 commits, which is in the range of [1206, 25000].
 - **Stars:** The repository has 53.2k stars, showing that it is widely used and well-established.
 - **Myers Algorithm:** Since Git uses Myers as the default diff algorithm, we can directly apply `git diff -diff-algorithm=myers` on the Requests repository.
 - **Histogram Algorithm:** Git also supports the Histogram algorithm, which works with this repository using `git diff -diff-algorithm=histogram`.
 - **Discrepancy Analysis:** With thousands of commits and frequent code updates, the Requests repository provides sufficient changes to compare Myers and Histogram diffs effectively.
 - **Required Files:** The repository contains the `requests/` source directory, a large `tests/` suite, along with `README.md` and `LICENSE`, fulfilling the file-type requirements.
 - **glide** – A fast and efficient image loading and caching library for Android, widely used for handling media in mobile applications.
 - **Commits:** The repository has 3047 commits, which lies in the range of [1206, 25000], making it suitable for analysis.

-
- **Stars:** The repository has 34.9k stars, showing its strong adoption and relevance in real-world projects.
 - **Myers Algorithm:** Since Git uses Myers as the default diff algorithm, we can directly apply `git diff -diff-algorithm=myers` on the Glide repository.
 - **Histogram Algorithm:** Git also supports the Histogram algorithm, which works seamlessly with this repository using `git diff -diff-algorithm=histogram`.
 - **Discrepancy Analysis:** With over three thousand commits and consistent updates, the Glide repository provides enough modifications to effectively compare Myers and Histogram diffs.
 - **Required Files:** The repository includes the `glide/` source code, a dedicated `tests/` suite, along with `README.md` and `LICENSE`, satisfying the file-type requirements.

4.5 Part C: Run Software Tool on the Selected Repository

Introduction

This part requires to calculate the difference between the current source code and previous source code. This initially sounded a simple task, as we would exactly know in which lines some changes have been taken place. Surprisingly, a lot of research has been done to find the code differences. In this task we use two algorithms to find the code differences in the cloned repositories.

Myer's Algorithm

This algorithm uses Longest Common Subsequence method to find the differences in the codes. The algorithm efficiently searches diagonals using a greedy strategy and stores only frontier points, so it's very memory-efficient.

- Produces a minimal diff (shortest possible sequence of edits).
- Runs in $O(ND)$ time, where N = number of lines and D = edit distance.
- Git uses it by default because it balances accuracy and performance.

Myer's Algorithm

Designed for better readability of diffs in code, especially when the same line appears many times.

- First, it builds a histogram of line occurrences in both files.

- It then focuses on unique lines (lines that appear only once or very rarely) because they are more likely to be real anchors for matching.
- It tries to align these unique lines first, then expands outward to find context around them.
- Often produces a more human-friendly diff compared to Myers, though not always the minimal edit script.

Methodology and Execution

We use existing libraries to perform the task of extracting code differences.

```

local_repo_paths = [
    "/Users/tejaswicaiad/Desktop/Third_year/STT/lab4/butterknife",
    "/Users/tejaswicaiad/Desktop/Third_year/STT/lab4/retrifit",
    "/Users/tejaswicaiad/Desktop/Third_year/STT/lab4/glide"
]
output_csv_filename = "diff_analysis_results.csv"

csv_header = [
    "repository_name", "old_file_path", "new_file_path", "commit_SHA",
    "parent_commit_SHA", "commit_message", "diff_myers", "diff_hist", "Discrepancy"
]

with open(output_csv_filename, 'w', newline='') as csvfile:
    csv_writer = csv.writer(csvfile)
    csv_writer.writerow(csv_header)

    five_mb = 5 * 1024 * 1024 # measure the growth of csv file (preventive measure so that it does not become too big to process)
    print(f"Five MB = {five_mb} MB")

    for path in local_repo_paths: #looping each of the repository
        print(f"\nProcessing repository: {os.path.basename(path)}\n")
        try:
            for commit in Repository(path).traverse_commits(): #traversing across each of the repo
                if not commit.parents:
                    continue
                parent_SHA = commit.parents[0] # taking parent of the current commit, [0] used to avoid issues in case of merging commit

                for mod in commit.modified_files:
                    if mod.filename.endswith(('.png', '.jpg', '.jpeg', '.gif', '.zip')): #ignoring the processing on the files which are in the binary format
                        continue

                    filepatchdiff = mod.new_path #getting the path of the current committed commit, we take or with mod.old_path to avoid exception in case where the file might have been deleted later on
                    if not filepatchdiff:
                        continue

                    common_args = ["git", "-C", path, "diff", parent_SHA, commit.hash, "--"]
                    myers_proc = subprocess.run(common_args, capture_output=True, text=True, encoding='utf-8', errors='ignore')
                    hist_proc = subprocess.run(common_args + ["--diff-algorithm=histogram"], capture_output=True, text=True, encoding='utf-8', errors='ignore')

                    myers_diff_text = myers_proc.stdout #diff according to the myers algorithm
                    hist_diff_text = hist_proc.stdout #diff according to the hist_diff_text algorithm

                    norm_myers = "\n".join([re.sub(r'\n+', '\n', line) for line in myers_diff_text.splitlines()]) # normalization as mentioned in the question
                    norm_hist = "\n".join([re.sub(r'\n+', '\n', line) for line in hist_diff_text.splitlines()]) # normalization as mentioned in the question

                    discrepancy_found = "Yes" if norm_myers != norm_hist else "No" # manual check, if they exactly overlap then no discrepancy which is the majority case, else "Yes"

                    csv_writer.writerow([#writing the extracted content to the csv file
                        os.path.basename(path), mod.old_path, mod.new_path,
                        commit.message, parent_SHA, commit.hash,
                        myers_diff_text, hist_diff_text, discrepancy_found
                    ])

                    current_size = os.path.getsize(output_csv_filename)
                    if current_size >= printat:
                        size_in_mb = current_size / (1024 * 1024)
                        print(f"\nCSV file size has reached ({size_in_mb:.2f} MB.)")
                        printat += five_mb

        except Exception as error:
            print(f"\nAn error processing ({os.path.basename(path)}): ({error})")
            continue

print("\nAnalysis finished. Results are in ({output_csv_filename})")

```

Figure 4.1: Central Code to extract diff.

The code shown in Figure 4.1 is responsible for extracting the differences between file versions in the selected repositories. It leverages the Pydriller library to iterate through each commit and retrieve the modified files. For each file, the code applies both the Myers and Histogram diff algorithms using Git commands, capturing the output for further analysis.

Key steps in the code:

- Store all the paths to the repositories in `local_repo_paths`.
- Define the output `csv_filename` and also define the `csv_header` for that csv file which contains all the columns asked in the part (c) of the assignment.

-
- Algorithm then starts iterating through each repository path, and also tracks the size of growing csv file.
 - Iterating through every commit of that repository which initially fetches the hash value of the commit.
 - Extraction of parent_sha (We only take the first parent in case of a merge commit). Following this, we iterate through each of the modified files in that commit.
 - Defining common_args which stores commands to be passed on to the subprocess.run() as arguments.
 - We obtain the Myers diff and Histogram diff outputs by executing the respective Git commands and capturing their output and apply normalization
 - Defining discrepancy_found if the outputs of the two algorithms differ and commit the output to the CSV file

This approach allows for systematic extraction and comparison of diff outputs, facilitating the analysis of discrepancies between the Myers and Histogram algorithms across real-world code changes.

4.6 Part D: Compare Diff Outputs for Discrepancy Analysis

METHODOLOGY and EXECUTION

After processing all the repositories we record the output of both the algorithms (Myers and Histogram) in a CSV file.

1. **Normalization:** To ensure fairness in comparison, whitespace and blank lines were ignored.
2. **Comparison Logic:** The normalized Myers output and Histogram output were compared line by line.
 - If both were identical, the entry was labeled **No** discrepancy.
 - If they differed, the entry was labeled **Yes** discrepancy.
3. **CSV Output:** A dedicated column named **Discrepancy** was added to the CSV file. This column stores either **Yes** or **No** for each commit-file pair.

4.7 Part E: Report Final Dataset Statistics

Introduction

After collecting and storing the diff results from all three repositories, we performed an exploratory analysis of the final dataset. The analysis focused on quantifying discrepancies and categorizing them based on file types.

Methodology and Execution

1. **File Categorization:** Each modified file was categorized into one of the following groups: *Source code*, *Test code*, *README*, *LICENSE*, or *Other*. This categorization was performed by checking the file extensions and naming patterns (e.g., files containing “test” were labeled as Test code).
2. **Mismatch Detection by Category:** For each category, the total number of files and the number of files with discrepancies between the Myers and Histogram algorithms were computed. This enabled us to analyze which file types were more prone to mismatches.

Results

1. **Statistical Summary:** The following statistics were derived from the analysis:
 - Total files analyzed across all repositories: **31,344**.
 - Overall mismatches detected: **1,108**, corresponding to approximately **3.54%** of the total files.
 - Breakdown of mismatches by file category:
 - **Source code:** 14,515 files, 616 mismatches (**4.24%**).
 - **Test code:** 9,282 files, 428 mismatches (**4.61%**).
 - **README files:** 381 files, 12 mismatches (**3.15%**).
 - **LICENSE files:** 17 files, 0 mismatches (**0%**).
 - **Other files:** 7,149 files, 52 mismatches (**0.73%**).
 - **File Extension Summary (Top 10):**

-
- **.java:** 23,281 files, 1,040 mismatches (4.47%)
 - **.kt:** 350 files, 14 mismatches (4.00%)
 - **.md:** 541 files, 15 mismatches (2.77%)
 - **Other:** 142 files, 3 mismatches (2.11%)
 - **.json:** 70 files, 1 mismatch (1.43%)
 - **.gradle:** 1,410 files, 14 mismatches (0.99%)
 - **.yml:** 150 files, 1 mismatch (0.67%)
 - **.iml:** 152 files, 1 mismatch (0.66%)
 - **.xml:** 2,773 files, 15 mismatches (0.54%)
 - **.properties:** 587 files, 2 mismatches (0.34%)

2. Top 10 Files Causing Mismatches:

- library/src/main/java/com/bumptech/glide/load/resource.bitmap/Downsampler.java – 15 mismatches
- README.md – 12 mismatches
- library/src/com/bumptech/glide/Glide.java – 12 mismatches
- butterknife-compiler/src/main/java/butterknife/compiler/BindingClass.java – 11 mismatches
- library/src/main/java/com/bumptech/glide/load/engine/DecodeJob.java – 9 mismatches
- library/src/main/java/com/bumptech/glide/request/target/ViewTarget.java – 8 mismatches
- library/src/main/java/com/bumptech/glide/load/engine/EngineJob.java – 7 mismatches
- library/src/main/java/com/bumptech/glide/RequestBuilder.java – 7 mismatches
- butterknife-compiler/src/main/java/butterknife/compiler/ButterKnifeProcessor.java – 6 mismatches
- samples/giphy/src/main/java/com/bumptech/glide/samples/giphy/MainActivity.java – 6 mismatches

3. Commit-Level Summary:

- **Average mismatches per commit:** 0.03

-
- **Top 10 commits with most mismatches:**
 - f389e91ccecac6ddf736bbf1a4346782609eb034 - 229 mismatches
 - f7a6d65cf7c1a41908dd48e0dab68ee5b881387e - 49 mismatches
 - 8f8a1600826fd042abae9251cbad063dee5144b2 - 36 mismatches
 - c375a2fbf594bdf422c45a1395a65823141a8bd5 - 35 mismatches
 - 0dcc33fe6957657b910484599c499430cdf3461c - 32 mismatches
 - ee71a6858dfddcc4650f6ff4d71112911bdf6ca7 - 31 mismatches
 - 6f1d71715361d68384afa0cf5fc9ad0e898b3697 - 18 mismatches
 - ed20643fb94d4e17f4cdb3699a6d83621408dd34 - 15 mismatches
 - 91cb86225967b1e12cdab52d2d3ea35afc2b8c9e - 12 mismatches
 - eac2c5f7190d148e72341ea289b69cae1ae2a866 - 12 mismatches

4. Repository-Level Summary:

- **glide:** 20,105 files, 820 mismatches (4.08%)
- **butterknife:** 3,120 files, 109 mismatches (3.49%)
- **retrofit:** 8,119 files, 179 mismatches (2.20%)

5. **Visualization:** To better understand the dataset, we generated multiple plots summarizing files, mismatches, and repository behavior:

- **Total Files by Category:**

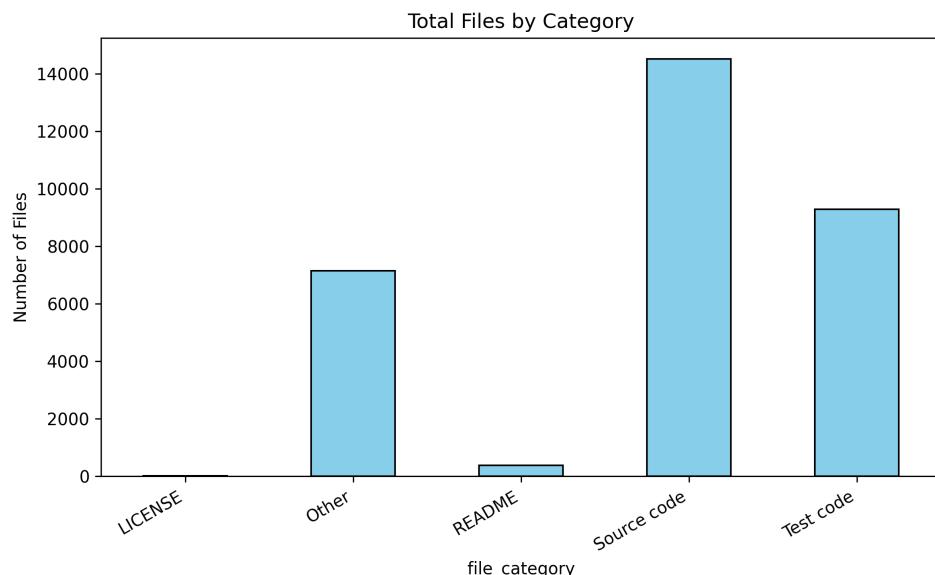


Figure 4.2: Bar chart showing the total number of files per file category.

- **Mismatches by File Category:**

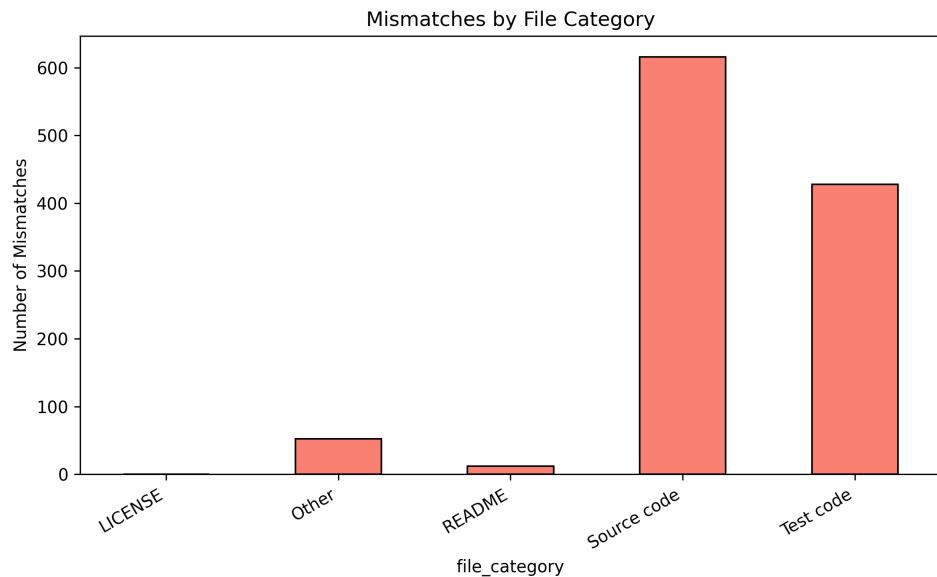


Figure 4.3: Bar chart showing the number of mismatches for each file category.

- **Mismatch Percentage by File Category:**

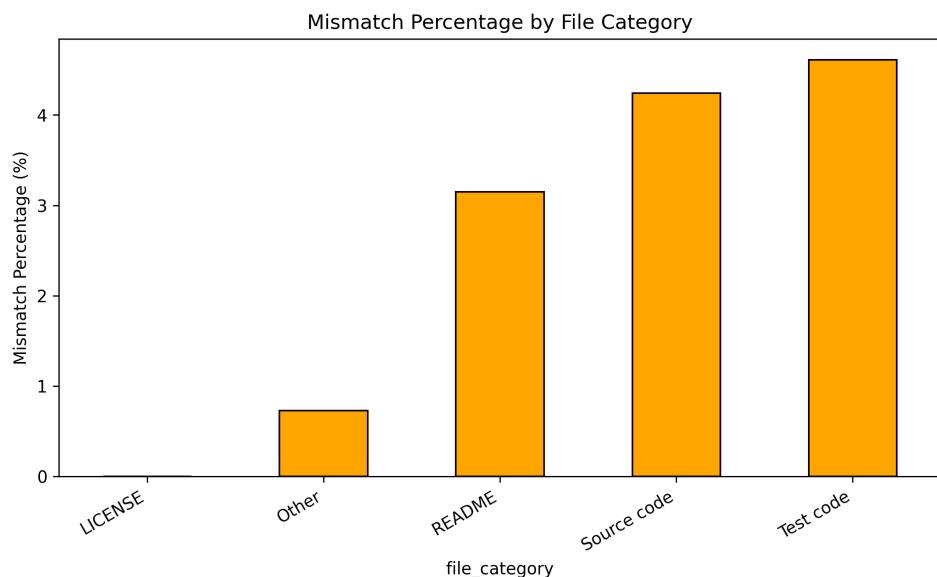


Figure 4.4: Bar chart displaying the mismatch percentage per file category.

- **Top 10 File Extensions by Mismatch Percentage:**
- **Mismatch Percentage by Repository:**
- **Distribution of Mismatches per Commit:**
- **Top 10 Commits by Number of Mismatches:**

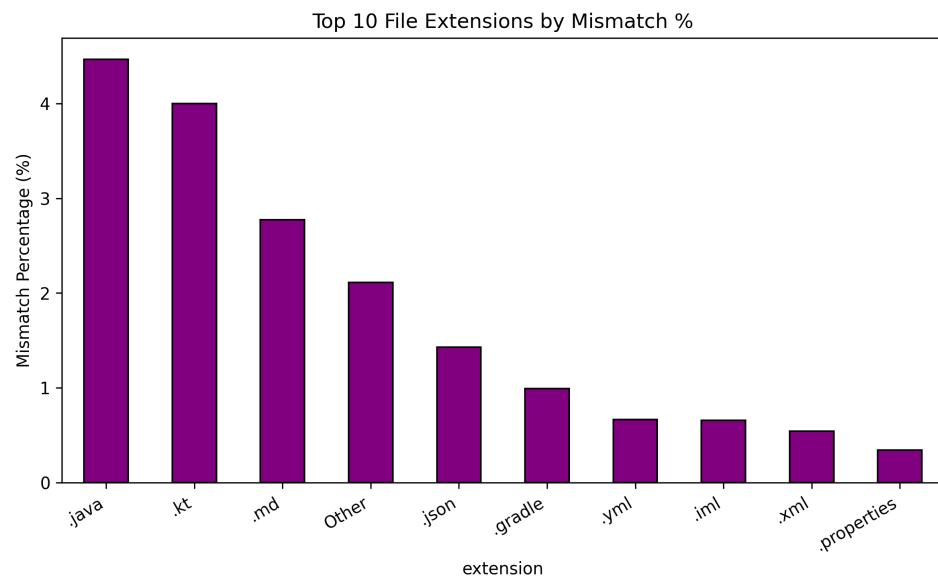


Figure 4.5: Bar chart highlighting the top 10 file extensions with the highest mismatch percentages.

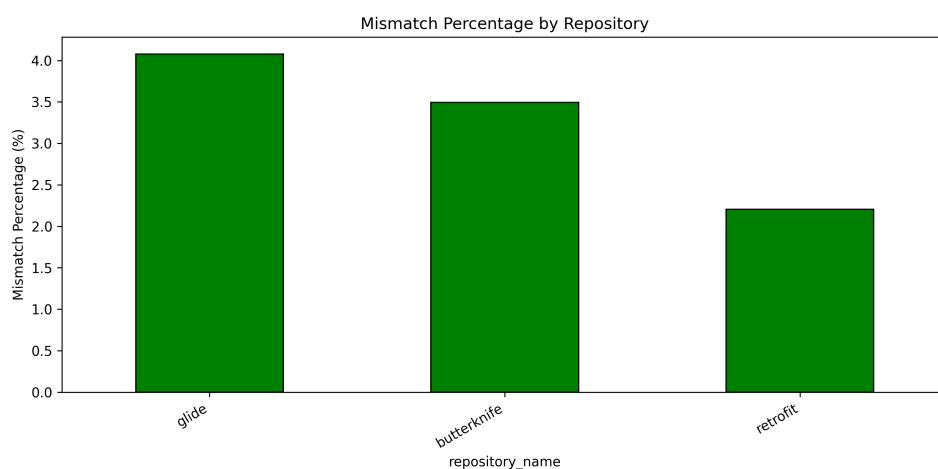


Figure 4.6: Bar chart showing mismatch percentages across repositories.

These plots and statistics prove an understanding of the dataset's structure and the nature of the mismatches present.

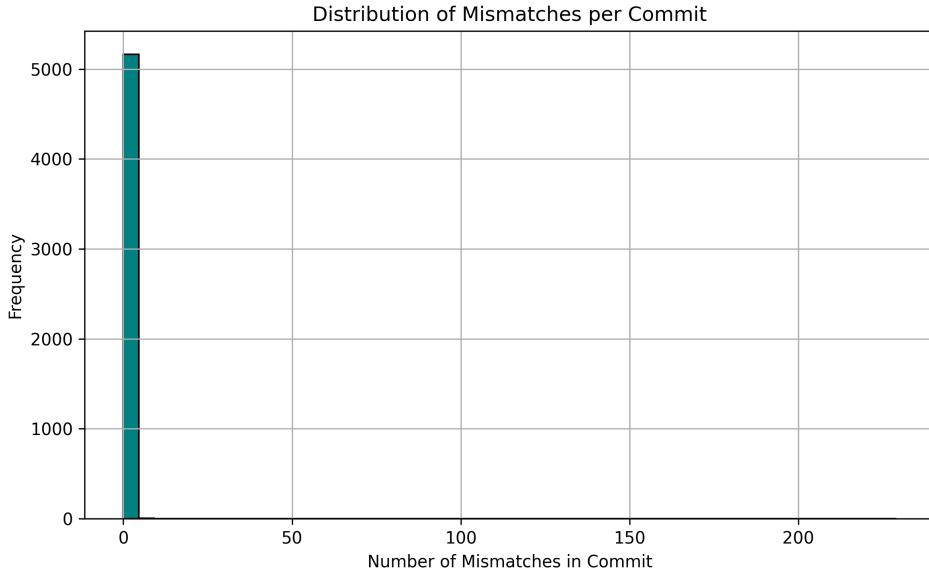


Figure 4.7: Histogram illustrating the distribution of mismatches per commit.

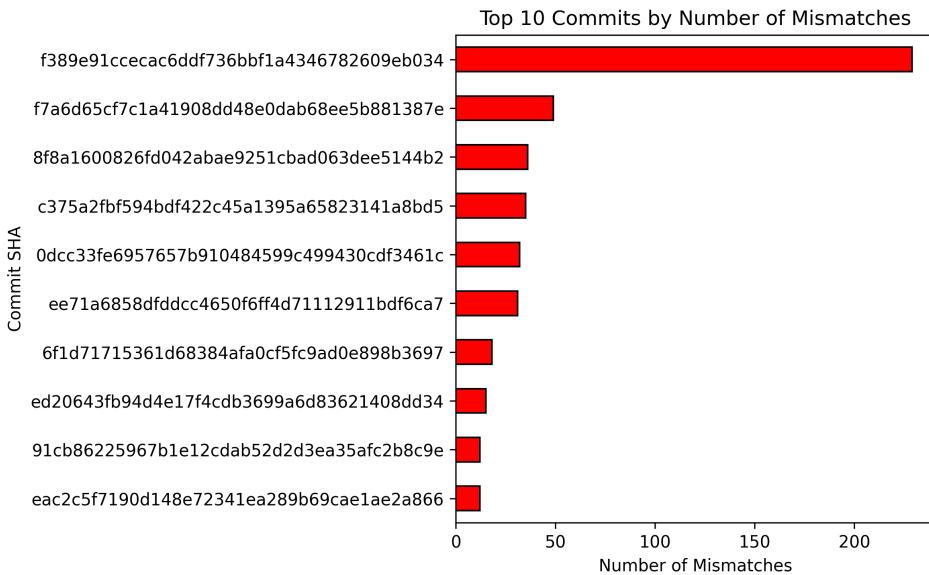


Figure 4.8: Horizontal bar chart displaying the top 10 commits responsible for the most mismatches.

4.8 Analysis and Conclusion

Analysis

The exploratory analysis revealed that mismatches between Myers and Histogram algorithms were relatively rare overall, with only **3.54%** of the 31,344 files affected. However, their distribution across file categories and repositories highlights certain patterns:

- **File Category Trends:**

-
- **Source code (4.24%)** and **Test code (4.61%)** showed the highest mismatch rates, suggesting that more complex files with frequent structural edits (e.g., function modifications, assertions) are more sensitive to algorithmic differences.
 - **README (3.15%)** and **LICENSE (0%)** files showed very few mismatches, indicating that simpler, more static text rarely produces discrepancies.
 - **Other files (0.73%)** had minimal mismatches, largely due to their structured formats (e.g., configuration files) where line changes are easier to track.
- **Extension-Level Insights:**
 - **Java files (.java)** dominated both in total count and mismatches (~23k files, 4.47%), reflecting their heavy presence in the repositories.
 - **Kotlin (.kt)** and **Markdown (.md)** files showed moderate mismatch percentages (~4% and 2.7% respectively).
 - Configuration-related extensions such as **.yml**, **.xml**, and **.properties** had very low mismatch rates, consistent with their line-oriented, less-structured nature.

- **Repository-Level Trends:**
 - The **glide repository** had the highest mismatch percentage (4.08%), reflecting its large codebase and frequent file modifications.
 - **butterknife (3.49%)** and **retrofit (2.20%)** showed relatively fewer mismatches, likely due to smaller or more modular project structures.
- **Commit-Level Observations:** Although the average mismatches per commit was only **0.03**, a few commits accounted for a disproportionate number of mismatches. For example, commit **f389e91cceac6ddf736bbf1a4346782609eb034** alone had 229 mismatches, highlighting how large-scale refactorings or bulk updates drive discrepancies.

Conclusion

The comparative analysis of Myers and Histogram algorithms demonstrated that while both generally agree, certain types of files and edits cause them to diverge. Source code and test files are particularly prone to mismatches, whereas configuration and license files remain stable.

At the repository level, larger and more actively developed projects such as **glide** exhibited higher mismatch rates. Similarly, individual commits involving large structural changes were significant contributors to observed differences.

Overall, these findings emphasize that the choice of diff algorithm can meaningfully impact results in code analysis tasks, especially in repositories with frequent, complex

modifications. Future work could extend this study by correlating mismatch frequency with developer activity, project size, and specific refactoring patterns.

4.9 Part F: Which algorithm performed better

- **Metric Definition** – We can define multiple metrics to evaluate the performance of the algorithms, such as number of mismatches, percentage of mismatches, runtime efficiency (Compute required).
- **Run Model to find various metrics** – Execute the models on the dataset to obtain the defined metrics for each algorithm.
- **Conclusive analysis** – Just based on the results on a single repository, we cannot conclude which algorithm is better. However, based on the results from all three repositories, we can perform Statistical tests to compare the algorithms more rigorously to prove/disprove any hypothesis.
- **Script to run these tests** – Automating the script to report the result which includes to run the model to get metrics, run statistical tests(optional), and report the one which scores better.