

Indian Institute of Technology
Department of Computer Science

CS202: Software Tools and Technologies Assessment 1

Jaskirat Singh Maskeen
(23110146 | jaskirat.maskeen@iitgn.ac.in)

Instructor: Prof. Shouvick Mondal

September 7, 2025

Contents

1	Introduction to VCS, Git Workflow and Actions	1
1.1	Introduction	1
1.2	Setup and Tools	1
1.3	Methodology and Execution	2
1.4	Results and Analysis	5
1.5	Discussion and Conclusion	7
	References	8
2	Commit Message Rectification for Bug-Fixing Commits	9
2.1	Introduction	9
2.2	Setup and Tools	9
2.3	Methodology and Execution	9
2.4	Results and Analysis	12
2.5	Discussion and Conclusion	15
	References	16
3	Bug Context Analysis and Agreement Detection	17
3.1	Introduction	17
3.2	Setup and Tools	17
3.3	Methodology and Execution	17
3.4	Results and Analysis	19
3.5	Discussion and Conclusion	22
	References	22
4	Different diff algorithms on Repositories	23
4.1	Introduction	23
4.2	Setup and Tools	23
4.3	Methodology and Execution	23
4.4	Results and Analysis	26
4.4.1	Private-gpt	26
4.4.2	Crawl4ai	27
4.4.3	Marker	27
4.5	Discussion and Conclusion	28
	References	28

Lab 1

Introduction to Version Controlling, Git Workflows, and Actions

[Github Repository for Lab 1](#)

1.1 Introduction

Git is a version control system that helps us manage the versions of our projects either locally or on some Git provider such as GitLab or GitHub. The primary advantage of using GitHub or any such provider is that we can collaborate with other users and share our code online with other potential users. This sort of framework supports active development and constant feedback, as anyone can create an Issue or fork our repository and provide pull requests to improve or fix issues with our code.

The main goal of this lab is to understand version control systems (VCS), set up and configure git, verify the installation, and run simple Git commands to push files.

1.2 Setup and Tools

I have Git, Python, and VSCode installed. Since I have already used GitHub before, I have it preconfigured and installed, so I will skip the screenshots for installation details. I will use Git CLI on the command prompt and GitHub Web Interface to create repositories and run workflows. I will also use pylint to perform static type checking. `git --version` gives `git version 2.45.2.windows.1` as the output. While logging into git using the CLI, we first have to sign in using username and password, and verify on GitHub Mobile, as I have 2FA enabled. I also confirm if the current credentials are set (Fig. 1.1). The `--global` flag is

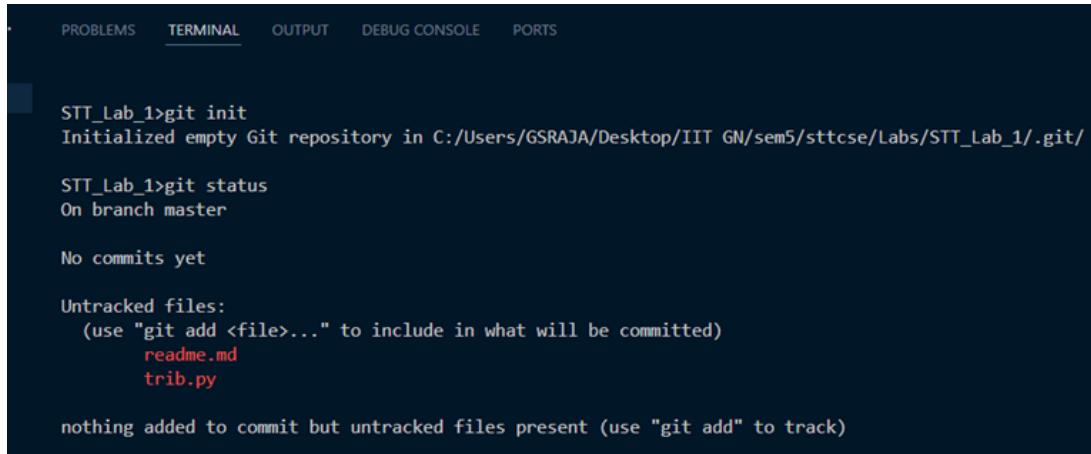
```
STT_Lab_1>git config --global user.name  
jsmaskeen  
  
STT_Lab_1>git config --global user.email  
jsmaskeen@gmail.com  
  
STT_Lab_1>|
```

Figure 1.1: Configuring the git CLI

used because I have two GitHub accounts on my laptop, and I want to verify the global configuration. These commands were looked through using <https://git-scm.com/docs/git-config>. My working directory is named `STT_Lab_1`.

1.3 Methodology and Execution

I first initialized a repository in my working directory using `git init`. I have put in a sample file named `trib.py`, which calculates the n^{th} Tribonacci number. I also added a `readme.md` file, which will get rendered on the link of this repository. I also plan to have a single repository on GitHub, with a new branch for each lab to avoid clutter. Hence, I will also rename the branch to `lab1`,



```
PROBLEMS TERMINAL OUTPUT DEBUG CONSOLE PORTS

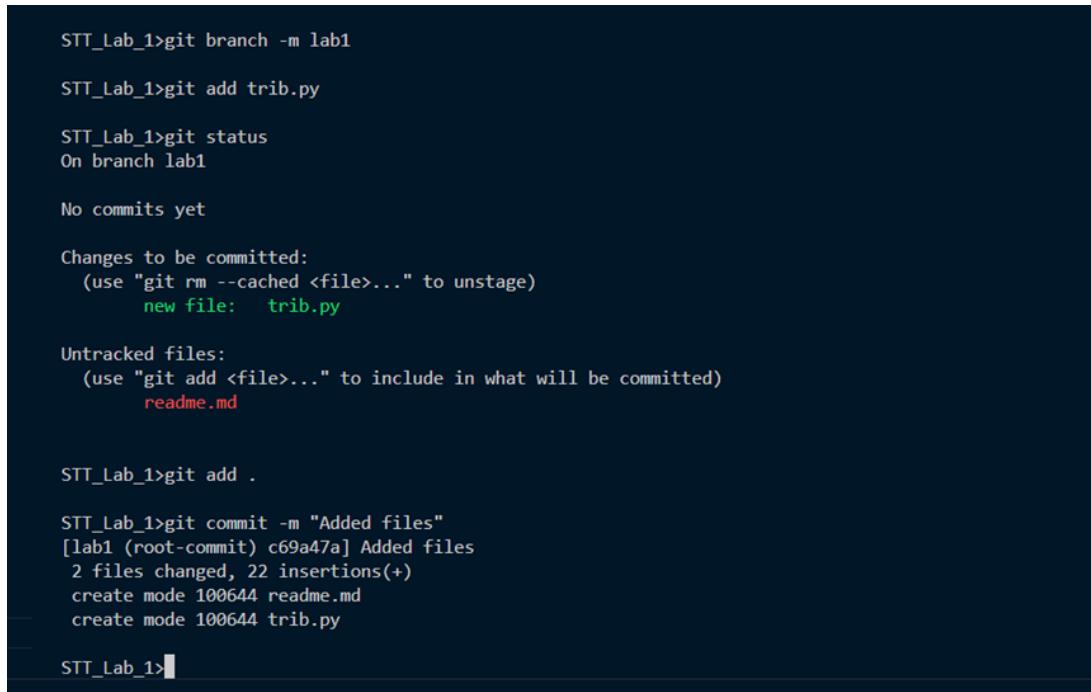
STT_Lab_1>git init
Initialized empty Git repository in C:/Users/GSRAJA/Desktop/IIT GN/sem5/sttcse/Labs/STT_Lab_1/.git/
STT_Lab_1>git status
On branch master

No commits yet

Untracked files:
  (use "git add <file>..." to include in what will be committed)
    README.md
    trib.py

nothing added to commit but untracked files present (use "git add" to track)
```

Figure 1.2: Initializing a git repository



```
STT_Lab_1>git branch -m lab1
STT_Lab_1>git add trib.py
STT_Lab_1>git status
On branch lab1

No commits yet

Changes to be committed:
  (use "git rm --cached <file>..." to unstage)
    new file:   trib.py

Untracked files:
  (use "git add <file>..." to include in what will be committed)
    README.md

STT_Lab_1>git add .
STT_Lab_1>git commit -m "Added files"
[lab1 (root-commit) c69a47a] Added files
  2 files changed, 22 insertions(+)
  create mode 100644 README.md
  create mode 100644 trib.py

STT_Lab_1>
```

Figure 1.3: Committing to `lab1` branch.

I can check the status of the repository using `git status` (Fig. 1.2), and add files one by one using `git add filename`, or all unstaged files (i.e., which have a delta), at once using `git add .`, I can commit these changes using the command `git commit -m "Added files"`, where "Added files" is a commit message (Fig. 1.3). `git log` is also a useful command to check the history of the commits being made. Now I create a GitHub repository, **CS202-STTCSE**, which will be the repository where I will upload all the labs, just on their respective branches. GitHub Link: <https://github.com/jsmaskeen/CS202-STTCSE/tree/lab1>.

I then added the local repository to GitHub using the commands,
`git remote add origin git@github.com:jsmaskeen/CS202-STTCSE.git`,
Followed by `git push -u origin lab1`.

```
STT_Lab_1>git log
commit c69a47ab4d509a597b546310d6f0222f5f2761d7 (HEAD -> lab1)
Author: js maskeen <jsmaskeen@gmail.com>
Date:   Sat Aug 9 20:32:36 2025 +0530

    Added files

STT_Lab_1>
```

Figure 1.4: git log command

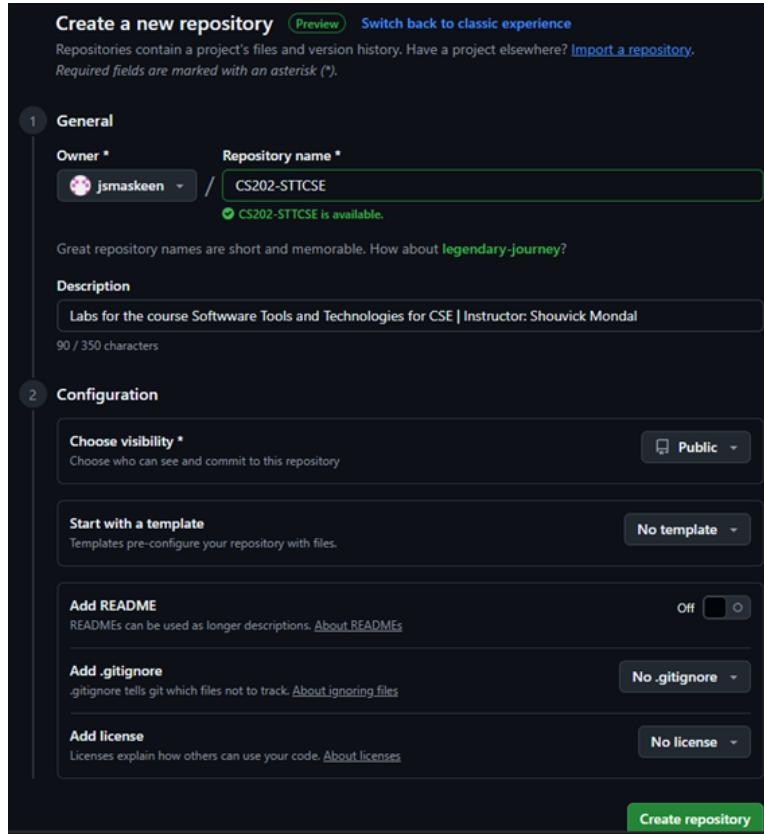


Figure 1.5: GitHub repository creation

```
STT_Lab_1>git remote add origin git@github.com:jsmaskeen/CS202-STTCSE.git
STT_Lab_1>git push -u origin lab1
Enumerating objects: 4, done.
Counting objects: 100% (4/4), done.
Delta compression using up to 16 threads
Compressing objects: 100% (4/4), done.
Writing objects: 100% (4/4), 535 bytes | 535.00 KiB/s, done.
Total 4 (delta 0), reused 0 (delta 0), pack-reused 0 (from 0)
To github.com:jsmaskeen/CS202-STTCSE.git
 * [new branch]      lab1 -> lab1
branch 'lab1' set up to track 'origin/lab1'.

STT_Lab_1>
```

Figure 1.6: Pushing the code to remote repository on GitHub

So, in this way, my local repository is now on GitHub. Now, if I do any commits and push them locally, they will reflect on GitHub as well. For example, I add another file, say **mergesort.py**. I will stage and commit the file, but not push it. Then I will clone this repository elsewhere, and to demonstrate the pull command, I will push the **mergesort.py** file after I have cloned the repository. Please see the sequence of commands carefully ((Figures 1.7, 1.8, 1.9, 1.10, and 1.11)).

This demonstrated the git pull as well as git clone commands. Now I have to create a pylint

```

PROBLEMS TERMINAL OUTPUT DEBUG CONSOLE PORTS

STT_Lab_1>git status
On branch lab1
Your branch is up to date with 'origin/lab1'.

Untracked files:
  (use "git add <file>..." to include in what will be committed)
    mergesort.py

nothing added to commit but untracked files present (use "git add" to track)

STT_Lab_1>git add .

STT_Lab_1>git commit -m "Added mergesort"
[lab1 114ceed] Added mergesort
 1 file changed, 56 insertions(+)
 create mode 100644 mergesort.py

```

Figure 1.7: Step 1: Stage the mergesort.py and commit

```

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\labs>cd ..

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse>git clone https://github.com/jsmaskeen/CS202-STTCSE.git
Cloning into 'CS202-STTCSE'...
remote: Enumerating objects: 4, done.
remote: Counting objects: 100% (4/4), done.
remote: Compressing objects: 100% (4/4), done.
remote: Total 4 (delta 0), reused 4 (delta 0), pack-reused 0 (from 0)
Receiving objects: 100% (4/4), done.

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse>cd CS202-STTCSE

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>git status
On branch lab1
Your branch is up to date with 'origin/lab1'.

nothing to commit, working tree clean

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>git log
commit c69a47ab4d509a5970546310d6f0222ff5f2761d7 (HEAD -> lab1, origin/lab1, origin/HEAD)
Author: jsmaskeen <jsmaskeen@gmail.com>
Date:   Sat Aug 9 20:32:36 2025 +0530

  Added files

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>tree /f
Folder PATH listing for volume Windows
Volume serial number is E6E3-610C
C:.
  README.md
  trib.py

No subfolders exist

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>

```

Figure 1.8: Step 2: Clone the repository elsewhere

```

STT_Lab_1>git push -u origin lab1
Enumerating objects: 4, done.
Counting objects: 100% (4/4), done.
Delta compression using up to 16 threads
Compressing objects: 100% (3/3), done.
Writing objects: 100% (3/3), 765 bytes | 765.00 KiB/s, done.
Total 3 (delta 0), reused 0 (delta 0), pack-reused 0 (from 0)
To github.com:jsmaskeen/CS202-STTCSE.git
  c69a47a..114ceed  lab1 -> lab1
branch 'lab1' set up to track 'origin/lab1'.

STT_Lab_1>[]

```

Figure 1.9: Step 3: Now I push to GitHub.

```

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>git fetch
remote: Enumerating objects: 4, done.
remote: Counting objects: 100% (4/4), done.
remote: Compressing objects: 100% (3/3), done.
remote: Total 3 (delta 0), reused 3 (delta 0), pack-reused 0 (from 0)
Unpacking objects: 100% (3/3), 745 bytes | 57.00 KiB/s, done.
From https://github.com/jsmaskeen/CS202-STTCSE
  c69a47a..114ceed  lab1      -> origin/lab1

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>git status
On branch lab1
Your branch is behind 'origin/lab1' by 1 commit, and can be fast-forwarded.
  (use "git pull" to update your local branch)

nothing to commit, working tree clean

```

Figure 1.10: Step 4: Check with the cloned repository (git fetch, git status)

```

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>tree /f
Folder PATH listing for volume Windows
Volume serial number is E6E3-610C
C:.
  README.md
  trib.py

No subfolders exist

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>git pull
Updating c69a47a..114ceed
Fast-forward
  mergesort.py | 56 ++++++-----+
  1 file changed, 56 insertions(+)
  create mode 100644 mergesort.py

C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\CS202-STTCSE>tree /f
Folder PATH listing for volume Windows
Volume serial number is E6E3-610C
C:.
  mergesort.py
  README.md
  trib.py

No subfolders exist

```

Figure 1.11: Step 5: Check the file structure before and after pulling. (git pull)

workflow, so I have deliberately kept some issues in the `mergesory.py` file, such as unused imports, unused

functions, and an unused while loop. I hope to see pylint detect all these issues.

To create a pylint workflow, I configured the GitHub action corresponding to pylint.

Once I clicked on configure, a **pylint.yml** file appeared in editing format under a new folder named workflows (Fig. 1.12). Once I committed and pushed this file on the GitHub UI (This means I will have to do git pull on local (Fig. 1.13), to stay up to date), the action ran and gave a few errors.

Whenever the action fails, we are greeted with a **X** (a red cross) (Fig. 1.14), and whenever it passes, we are greeted with a **✓** (a blue tick).

```

1 name: Pylint
2
3 on: [push]
4
5 jobs:
6   build:
7     runs-on: ubuntu-latest
8     strategy:
9       matrix:
10         python-version: ["3.8", "3.9", "3.10"]
11     steps:
12       - uses: actions/checkout@v4
13       - name: Set up Python ${{ matrix.python-version }}
14         uses: actions/setup-python@v3
15       - with:
16         python-version: ${{ matrix.python-version }}
17       - name: Install dependencies
18         run: |
19           python -m pip install --upgrade pip
20           pip install pylint
21       - name: Analysing the code with pylint
22         run: |
23           pylint $(git ls-files *.py)

```

Figure 1.12: Pylint workflow

```

PROBLEMS TERMINAL OUTPUT DEBUG CONSOLE PORTS

STT_Lab_1>git status
On branch lab1
Your branch is up to date with 'origin/lab1'.

remote: Enumerating objects: 6, done.
remote: Counting objects: 100% (6/6), done.
remote: Compressing objects: 100% (3/3), done.
remote: Total 5 (delta 0), reused 0 (delta 0), pack-reused 0 (from 0)
Unpacking objects: 100% (5/5), 1.34 KiB | 91.00 KiB/s, done.
From github.com:jsmaskeen/CS202-STTCSE
  413007e..f2ed760  lab1      -> origin/lab1

STT_Lab_1>git pull
Updating 413007e..f2ed760
Fast-forward
 .github/workflows/pylint.yml | 23 ++++++-----+
  1 file changed, 23 insertions(+)
 create mode 100644 .github/workflows/pylint.yml

STT_Lab_1>git status
On branch lab1
Your branch is up to date with 'origin/lab1'.

nothing to commit, working tree clean

STT_Lab_1>

```

Figure 1.13: Pulling on local to stay up to date

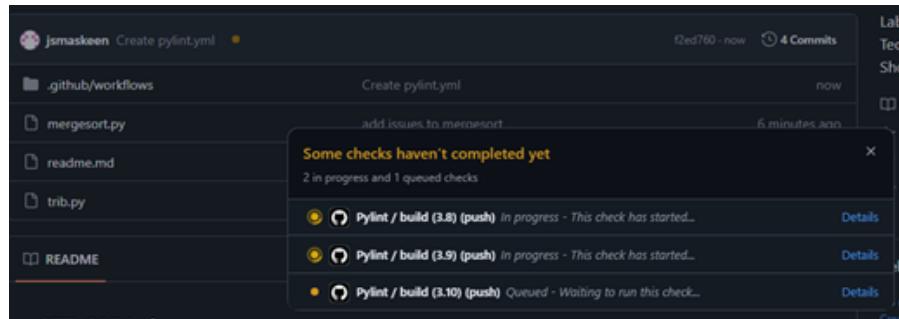


Figure 1.14: Workflow is running

1.4 Results and Analysis

Since I uploaded an erroneous file (which executes correctly, but has unused code, and should fail static checking) (Fig. 1.15).

Looking carefully at the details, we see multiple issues (Fig. 1.16, Fig. 1.17). However, the workflow was not able to detect a few issues, for example:

```
1 import random, pandas, numpy
2
3 def mergeSort(ls):
4     if len(ls) <= 1:
5         return ls
6     mid = len(ls) // 2
7     left_half = mergeSort(ls[:mid])
8     right_half = mergeSort(ls[mid:])
9     return merge(left_half, right_half)
10
11 def merge(left, Right):
12     mergeResult = []
13     i = 0; j=0
14     if len(left) > len(Right):
15         left, Right = Right, left
16     if len(left) > len(Right):
17         left, Right = Right, left
18     while i < len(left) and j < len(Right):
19         if left[i] < Right[j]:
20             mergeResult.append(left[i])
21             i += 1
22         else:
23             mergeResult.append(Right[j])
24             j += 1
25     while i < len(
26         left
27     ):
28         mergeResult.append(left[i])
29         i += 1
30
31     while j < len(Right):
32         mergeResult.append(Right[j])
33         j += 1
34
35     return mergeResult
36     print(mergeResult)
37
38 def unused_function():
39     for k in range(5):
40         print("Waterboy ?", k)
41     return "Firegirl"
42
43 def main():
44     nums = [random.randint(0, 100) for _ in range(15)]
45     print("Unsorted:", nums)
46     sorted_nums = mergeSort(nums)
47     print("Sorted:", sorted_nums)
48
49 main()
```

Figure 1.15: `mergesort.py` file (with errors)

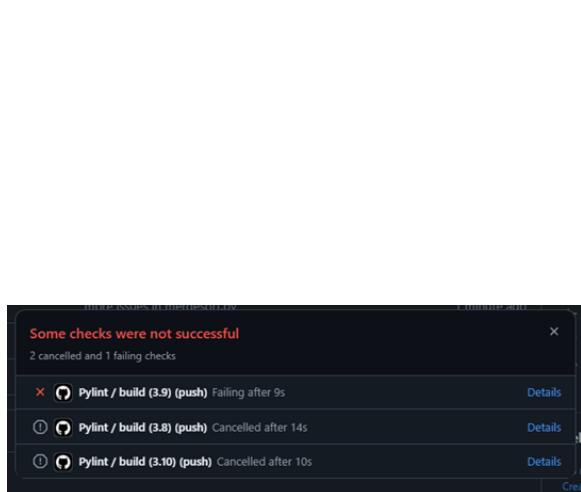


Figure 1.16: Failed Checks

```
Analysing the code with pylint

1 Run pylint $(git ls-files "*.py")  
2 ***** Module mergesort  
3 mergesort.py:27:0: C0301: Line too long (103/100) (line-too-long)  
4 mergesort.py:37:0: C0303: Trailing whitespace (trailing whitespace)  
5 mergesort.py:38:22: C0303: Trailing whitespace (trailing whitespace)  
6 mergesort.py:39:22: C0303: Trailing whitespace (trailing whitespace)  
7 mergesort.py:40:0: C0303: Trailing whitespace (trailing whitespace)  
8 mergesort.py:40:6: C0303: Trailing whitespace (trailing whitespace)  
9 mergesort.py:40:14: C0114: Missing module docstring (missing-module-docstring)  
10 mergesort.py:1:0: C0410: Multiple imports on one line (random, pandas, numpy) (multiple-imports)  
11 mergesort.py:1:0: E0401: Unable to import pandas (import-error)  
12 mergesort.py:1:0: E0401: Unable to import numpy (import-error)  
13 mergesort.py:1:0: C0101: Missing function or method docstring (missing-function-docstring)  
14 mergesort.py:1:0: C0101: Function name "mergeSort" doesn't conform to snake_case naming style (invalid-name)  
15 mergesort.py:1:0: C0101: Missing function or method docstring (missing-function-docstring)  
16 mergesort.py:11:0: C0101: Argument name "Right" doesn't conform to snake_case naming style (invalid-name)  
17 mergesort.py:12:4: C0403: Variable name "MergeResult" doesn't conform to snake_case naming style (invalid-name)  
18 mergesort.py:13:10: C0211: More than one statement on a single line (multiple-statements)  
19 mergesort.py:36:4: W0101: Unreachable code (unreachable)  
20 mergesort.py:38:0: C0101: Missing function or method docstring (missing-function-docstring)  
21 mergesort.py:42:0: C0101: Missing function or method docstring (missing-function-docstring)  
22 mergesort.py:11:0: W0611: Unused import pandas (unused-import)  
23 mergesort.py:11:0: W0611: Unused import numpy (unused-import)  
24 ***** Module trib  
25 trib.py:15:0: C0303: Trailing whitespace (trailing whitespace)  
26 trib.py:1:0: C0114: Missing module docstring (missing-module-docstring)  
27  
28 ***** Module tribr  
29 trib.py:15:0: C0303: Trailing whitespace (trailing whitespace)  
30 trib.py:1:0: C0114: Missing module docstring (missing-module-docstring)  
31 trib.py:1:0: C0114: Missing module docstring (missing-module-docstring)  
32  
33 ***** Module tribr  
34 Your code has been rated at 3.92/10  
35  
36 Error: Process completed with exit code 22.
```

Figure 1.17: Issues found by workflow.

In Fig. 1.18 we have the same `if` loop repeated (which swaps left and right), which serves no function. So, it should be kept only once. [But static analysis cannot capture that]. It was also not able to detect the unused function. To fix the issues, I have edited the docstrings, multiple statements in a single line, and line too long issues.

```

def merge(left, right):
    """Merges two sorted lists."""
    merge_result = []
    i = 0
    j = 0
    if len(left) > len(right):
        left, right = right, left

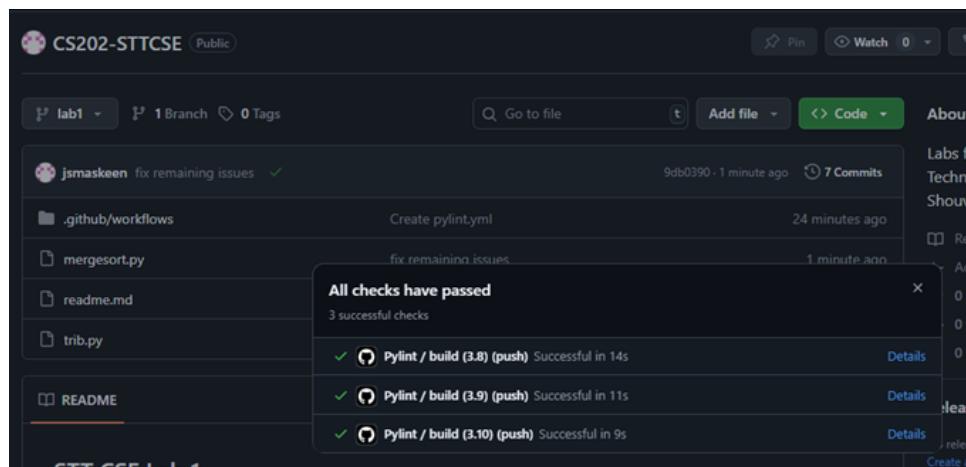
    if len(left) > len(right):
        left, right = right, left

    while i < len(left) and j < len(right):
        if left[i] < right[j]:
            merge_result.append(left[i])
            i += 1
        else:
            merge_result.append(right[j])
            j += 1
    while i < len(left):
        merge_result.append(left[i])
        i += 1
    while j < len(right):
        merge_result.append(right[j])
        j += 1
    return merge_result

```

Figure 1.18: Corrected Code

After pushing with these changes, we see that our tests have passed (Fig. 1.19).

Figure 1.19: All checks passed after fixing the `mergesort.py` file.

1.5 Discussion and Conclusion

This was a really interesting and useful lab session. Using git/GitHub will help us manage our code and collaborate with others; specifically, GitHub actions can configure our workflows, be it static checking or deploying a website after some post-processing. Maybe we write a library that requires a compilation, and various languages comprise the code base, so we can have a GitHub action to compile in the order we like and merge the files, to finally create an executable, for different OS. I did not face any difficulty in this lab; however, while configuring git, I had two accounts. I configured it with SSH tokens, so to set origin for a repository in my main account, I do `git remote add origin git@github.com:jsmaskeen/repo.git`, and for my other account, I do `git remote add origin git@github.com-unofficial:jsxxx/repo.git`.

Another thing I learnt, if we have ssh enabled, we can check if we are authenticated to GitHub via the command `ssh -T git@github.com`. Lastly, pylint was good, but it can only check issues statically. One mode of improvement can be to have two or more workflows, which also perform static checking, and they do not yield conflicting results, to find any issues missed by pylint (for example, the unused functions).

References

- [1] "Git documentation," <https://git-scm.com/docs>, accessed: 2025-08-08.
- [2] Stack Overflow, "How to verify github ssh key," <https://stackoverflow.com/questions/70371192/how-to-verify-github-ssh-key>, 2022, accessed: 2025-08-08.
- [3] "Lab assignment document," Google Classroom, accessed: 2025-08-08.

Lab 2

Commit Message Rectification for Bug-Fixing Commits in the Wild

[Github Repository for Lab 2](#)

Discussed with Siddhesh (23110347) and Aarsh (23110003)

2.1 Introduction

This lab is targeted at mining bugs in open source software repositories. The basic idea is that developers often fix certain bugs and issues throughout the development of the project. Our goal is to identify these bug-fixing commits using keyword analysis or regex (inspired by `minecpp`). Then, with these commits, I check the CommitPredictorT5 (CPT5) LLM model for its inference by giving it the diff between each of the modified files in the bug-fixing commit and the previous commit. Once this is done, I check if any rectification is needed. Basically, sometimes developers might fix multiple bugs in one go, and they might not write a descriptive enough commit message. So, I use `gpt-4o-mini` (rectifier), pass it the diff for a modified file, the original message, and get a modified commit message, had that file been committed independently. I then utilize the embedding model from Hugging Face, and try to classify the developer's commit message quality, the quality of the CPT5, and the quality of the rectifier. Algorithms are explained in the subsequent sections.

2.2 Setup and Tools

I am using `pydriller`, the `mamiksik/CommitPredictorT5` (CPT5) LLM model, `sentence-transformers/all-MiniLM-L12-v2` model for embedding generation, and `gpt-4o-mini` for the rectifier. I will also be using a virtual environment, so that my global packages do not interfere with `torch`, `transformers`, and `pydriller` libraries. and `python --version` gives Python 3.12.4, the `pydriller` version is shown in Fig. 2.1. SEART Search Engine <https://seart-ghs.si.usi.ch/> is used to search and filter the repository based on my criteria.

```
(venv) C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\Labs\STT_Lab_2>pip show pydriller
Name: PyDriller
Version: 2.8
```

Figure 2.1: Pydriller version

2.3 Methodology and Execution

I first chose the repository to mine the bugs. For this, I used the tool SEART Search Engine. I set a very simple criterion, I want to analyze the repositories that feature the creation of Discord bots, and they should have

- a minimum of 4000 starts
- a minimum of 1000 commits
- language should be Python
- discord in the title

The screenshot shows the GitHub search interface with the following parameters:

- General:** Search term: discord, License: discord, Contains: Python.
- History and Activity:** Number of Commits: 1000 to max, Number of Contributors: 10 to max, Created Between: dd-mm-yyyy to dd-mm-yyyy.
- Popularity Filters:** Number of Stars: 4000 to max, Number of Watchers: min to max, Number of Forks: min to max.
- Date-based Filters:** Last Commit Between: dd-mm-yyyy to dd-mm-yyyy.
- Additional Filters:** Size of codebase (Non Blank Lines: min to max), Repository Characteristics (Exclude Forks, Only Forks, Has Wild, Has License, Has Open Issues, Has Pull Requests).
- Sorting:** Name (Ascending).

A "Search" button is at the bottom.

Figure 2.2: Search parameters for finding the repository

The screenshot shows the GitHub search results page with two repository cards:

- cog-creators/red-discordbot**: Commits: 2900, Stars: 5216, Forks: 2389, Contributors: 174, Size: 82.02 KB, Last Commit: 2023-08-10.
- rapptz/discord.py**: Commits: 2943, Stars: 15579, Forks: 3878, Contributors: 159, Size: 21.72 KB, Last Commit: 2023-08-08.

Both cards have "Show More" links below them.

Figure 2.3: Search results

With this filtering, I was down to just two results. I will choose `rapptz/discord.py` as I have previously used this library to create Discord bots as well. I cloned the repository using
`git clone https://github.com/Rapptz/discord.py.git`.

Now I read the pydriller documentation, and wrote a custom script to go through the commits of the repository, and go through each modified file, and store the required details, as asked in the lab assignment.

The next step was to identify the bug-fixing commits. For this, I took inspiration from `minecapp` and from the lecture slides, and used the 52 keywords that were provided and the regex to search for in commit messages. I also have a word cloud function that gives me the word cloud of the words used in commit messages.

```
self.keywords = list(
    map(re.compile,
        [
            "fixed", "bug", "fixes", "fix", "fix", "fixed", "fixes", "crash", "solves", "resolves",
            "resolves", "issue", "issue", "regression", "fall back", "assertion", "covertify", "reproducible",
            "stack-wanted", "steps-wanted", "testcase", "failun", "fail", "npe", "npe", "except", "broken",
            "differential testing", "error", "hang", "hang", "test fix", "steps to reproduce", "crash",
            "assertion", "failure", "leak", "stack trace", "heap overflow", "freez", "problem", "problem",
            "overflow", "overflow", "avoid", "avoid", "workaround", "workaround", "break", "break",
            "stop", "stop"
        ],
    )
)
self.pattern = re.compile(
    r".*((solv(ed|es|eing))|(fix(s|es|ing|ed)?)((error|bug|issue)(s)?)).*"
)
```

Figure 2.4: 52 keywords for bug message identification

```
def gen_wordcloud(self):
    text = "\n".join(self.all_commits)
    self.word_cloud = WordCloud(
        width=1600, height=800, background_color="white"
    ).generate(text)
    return self.word_cloud

def bugfix_commits(self, count=-1):
    prev_curd_map: list[[list[Commit, Commit]]] = []
    i = 1
    while i < len(self.all_commits):
        current = self.all_commits[i]
        prev = self.all_commits[i - 1]
        commit_message = current.msg.lower()
        for kw in self.keywords:
            if (
                kw.search(commit_message)
                or self.pattern.search(commit_message)
                or kw.pattern in commit_message
            ):
                prev_curd_map.append([prev, current])
                break
        i += 1
    if count > -1:
        prev_curd_map = np.array(prev_curd_map)
        rows = np.random.choice(len(prev_curd_map), size=count, replace=False)
        prev_curd_map = prev_curd_map[rows]
    self.prev_curd_map = prev_curd_map
    return prev_curd_map
```

Figure 2.5: Wordcloud function

So, with this, I am able to find around 1875 commits, out of the total 5243 commits, which fix some bug(s). Due to resource and time constraints, I will only analyze 1397 bug-fix commits (out of 1875), which include 2405 modified files. The exact details for the bug-fixing commits are in the next section.

Rectifier is simply gpt-4o-mini running on a specific prompt, which is,

```

1 You are an assistant that produces a rectified commit message for a single modified
   file.
2 Input:
3 FILE: {}
4 ORIGINAL_COMMIT_MSG: "{}"
5 PARSED_DIFF:
6 {}
7
8 Produce a short rectified commit message [12-16 words] for this file maybe
   consisting of:
9 - Title
10 - One-line root cause
11 - One-line change summary
12
13 DO NOT OUTPUT ANYTHING ELSE and 8 TO 10 WORDS ONLY, starting filler messages etc.

```

Here, for each modified file, I fill in the corresponding details and save the output as the rectified message. The LLM inference is obtained in the following code snippet. The only preprocessing done on diff is to replace + with <add>, use <ide> as a prefix to unchanged lines, and to replace - with , so that it is easy for the CPT5 to understand.

Now, to evaluate how good the LLM, the rectifier, and the developer are, I follow the following procedure:

For any commit C_i , let the commit message be o_i . This commit will have some modified files. So, for any modified file m_j , I first get the LLM Inference, l_j . Now, using the prompt described above, we obtain a rectified message, r_j .

Now, for each of these, o_i , l_j and, r_j We first get the embeddings, \tilde{o}_i , \tilde{l}_j and, \tilde{r}_j respectively. I consider the rectified message to be the common comparison factor.

To assess the developer's quality, I find the cosine similarity between the original commit message and the rectified message. Similarly, to assess the LLM, I find the cosine similarity between the LLM Inference and the rectified commit message.

So,

$$D_{(i,j)} = \frac{\tilde{o}_i \cdot \tilde{r}_j}{\|\tilde{o}_i\| \|\tilde{r}_j\|}, \quad L_j = \frac{\tilde{l}_j \cdot \tilde{r}_j}{\|\tilde{l}_j\| \|\tilde{r}_j\|}$$

where both these values are between -1 and 1 (1 denotes maximum alignment, and -1 denotes opposite alignment).

Finally, to evaluate how good the rectifier is, I consider the ideal scenario. Ideally, the developer should give a clear idea of all the changes made across all modified files in the commit message. So, assuming the developer is good, we can say, if \tilde{o}_i , and $\sum_j \tilde{r}_j$ are close, then rectifier correctly captures the meaning of the original commit message, so no rectification is needed. If they are far, then the developer is not able to correctly capture all the changes done in the modified files for that commit, so the rectifier is needed.

So, again choosing the cosine similarity as a quantitative metric,

$$R_i = \frac{\left(\sum_j \tilde{r}_j\right) \cdot \tilde{o}_i}{\left\|\sum_j \tilde{r}_j\right\| \|\tilde{o}_i\|}$$

After all the bug-fixing commits have been processed, I can plot them and set a threshold for each of these three values. From there, we can quantify the hit rate as,

$$\text{Developer hit rate} = \frac{|\{D_{(i,j)} : D_{(i,j)} \geq D_{\text{thresh}}\}|}{\sum_i \sum_j 1}$$

$$\text{LLM hit rate} = \frac{\sum_i |\{L_j : L_j \geq L_{\text{thresh}}\}|}{\sum_i \sum_j 1}$$

$$\text{Rectifier hit rate} = \frac{|\{R_i : R_i \leq R_{\text{thresh}}\}|}{\sum_i 1}$$

can then be multiplied by 100 to achieve hit rates as a percentage.

2.4 Results and Analysis

After letting the code analyze the discord.py repository (I had enabled checkpoint saves, after every 100 modified files are analyzed), and running the wordcloud function on commit messages, we see the results described in subsequent paragraphs.

```
ddf,df2 = await miner.analyze_all_pairs(verbose=True,max_workers=16,save_every=100)

[6] ⓘ 96m 30.4s
...
Resuming from saved state: 1004 files processed, 1004 entries.
Total modified-files (ordered): 3774. To process: 2770
Processed 1000/3774 files
Processed 1020/3774 files
Processed 1030/3774 files
Processed 1040/3774 files
Processed 1050/3774 files
Processed 1060/3774 files
Processed 1070/3774 files
Processed 1080/3774 files
Processed 1090/3774 files
Processed 1100/3774 files
Checkpoint saved at 1100 files -> checkpoints\pairs_1100.csv, checkpoints\commits_1100.csv
Processed 1110/3774 files
Processed 1120/3774 files
Processed 1130/3774 files
Processed 1140/3774 files
Processed 1150/3774 files
Processed 1160/3774 files
Processed 1170/3774 files
Processed 1180/3774 files
Processed 1190/3774 files
Processed 1200/3774 files
Checkpoint saved at 1200 files -> checkpoints\pairs_1200.csv, checkpoints\commits_1200.csv
Processed 1210/3774 files
...
Processed 1370/3774 files
Processed 1380/3774 files
Processed 1390/3774 files
Processed 1400/3774 files
```

Figure 2.6: Running of the commit analysis function

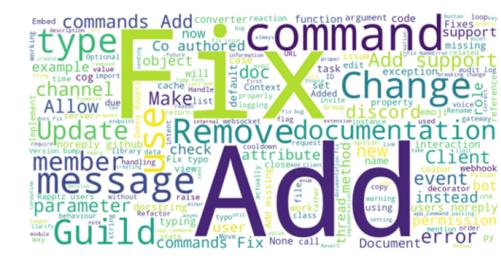


Figure 2.7: Wordcloud of commit messages.

Here I am sharing a subset of details of bug-fixing commits (refer to Table 2.1), the CSV file containing all of these commits is present on the [GitHub link](#)

Table 2.1: Sample of Bug-Fix Commits

Hash	Message	Hashes of parents	Is a merge commit ?	List of modified files
4c74523794b655b cff72ace92c7ea6f5 35a0e26a	Fix <i>versionadded</i> doc-strings in voice connect methods	cf031f71b91854b3 eedb16004412ff5 e05c9360	No	abc.py, voice_client.py
af311bff09f57ba3a 769cae4298559ea1 4b50d50	<i>Change .pot files to only contain at most 5 context lines</i> Crowdin did not like having many comments which caused the upload to fail.	78a026aae2da0ae 20669880e6c6384 57f7166419	No	message.pot_t, builder.py
b28a4a115e72977 a9091091983a121 3e91f8a5e9	Fix potentially stuck ratelimit buckets	ef06d7d9db2b868c 57f45fe1252d0554 90ad31a9	No	http.py
446c5029952d2b6 f466348b973888c0 1b0b2ef6e	Change lowercase detection to work with CJK languages. str.islower() does not properly work with characters in the Lo category, so CJK languages fail the check. Fix #7698	fd5dea4e340e69cd 2527114283d8c06 b2396689b	No	commands.py
400936df6989748c 2ec391778088a07 af355dbcfc	Fix type for content param in HTTP-Client.send_message	cdf46127ae2e1a18 3e5c94d9eb970bb ccbed64af	No	http.py
cdfccfbe74ba3680 58702b8035f1e637 8b245664	Fix description of user parameter in reaction remove event	53694724c13d079 1d77fcfd42564817 3771555a8	No	api.rst

One thing to note here is that none of the bug-fix commits are merge commits, even though in the repository, there are thousands of merged pull requests. The reason is that discord.py tries to maintain a linear commit tree, so it squashes the merges.

If I had used some other repository, then of course some of the bug fix commits would be merge commits, and they would have two parents.

We know every commit has at least one modified file, and for each modified file, I store the details (refer to Fig. 2.8) (again, this is a subset of the modified files; the entire CSV is on GitHub, moreover, I have pasted in links for the diff, source code (before) and (after), to display here, because they were very large. However in the CSV they are stored as plain text).

Hash	Message	Filename	Source Code (before)	Source Code (current)	Diff	LLM Inference (fix type)	Rectified Message
72e84a1b47049d4449873845f8bf89c302a033a	Change all email/password pair examples to use token.	docs\migrating.rst	https://katb.in/ozowedariro	https://katb.in/erixufiviju	use token instead of password		Update examples to use token instead of email/password.
f56dbb0379eabfe99c4fdeca1fbabc0bf5b0f6df	Fix Reaction not importing. Also bump version to v0.14.1	discord\init__.py	https://katb.in/gekoziwunos	https://katb.in/uukefijiwigi	update version numbers for discord 2.0		Fix Reaction import and update version to 0.14.1.
46bab822a11ba482be864f687d2bcef66c19005	Fix bug with editing messages over private messages.	discord\client.py	https://katb.in/ticaqodinoy	https://katb.in/afodubicula	fix logs from for private channels		Fix message editing in private channels. Handle <code>guild_id</code> for private messages. Update <code>edit_message</code> call accordingly.
73919fe1523d6b5e2e61d01b493a080a68766eb0	Documentation fixes and version bump.	discord\init__.py	https://katb.in/ecivevoxuvu	https://katb.in/wujupakinar	set version to 0.5.0		Bump version to 0.5.0. Update version and build numbers in <code>__init__.py</code> .
ab2512785b16d76a78f33e61a60ce90a2b225233	Handle VOICE_STATE_UPDATE <u>websocket</u> events. This adds a lot of new attributes into the Member class such as giving a <code>voice_channel</code> that the user is currently connected to. Initially there was a plan to have a <code>voice_members</code> attribute in the Channel class but this proved to be difficult when it came to actually removing users from the voice channel as the response would return <code>channel_id</code> as null. Fixes #16.	docs\api.rst	https://katb.in/apalezeholo	https://katb.in/amaxagapaza	https://katb.in/lopacuridig	add documentation for <u>on_voice_state_update</u>	Add <code>on_voice_state_update</code> function to handle voice state changes.

Figure 2.8: Table 2.1 extended with LLM inference and rectified message columns. Links are clickable in Table 2.2

Table 2.2: Source Code Comparison with Hashes

Hash	Source Code (before)	Source Code (current)	Diff
72e84a1b47049d4449873845f8bf89c302a033a	https://katb.in/ozowedariro	https://katb.in/botudarebeq	https://katb.in/erixufiviju
f56dbb0379eabfe99c4fdeca1fbabc0bf5b0f6df	https://katb.in/gekoziwunos	https://katb.in/uukefijiwigi	https://katb.in/gaxeletejucuh
46bab822a11ba482be864f687d2bcef66c19005	https://katb.in/ticaqodinoy	https://katb.in/afodubicula	https://katb.in/ubofowuzebo
73919fe1523d6b5e2e61d01b493a080a68766eb0	https://katb.in/ecivevoxuvu	https://katb.in/ihepcobiku	https://katb.in/wujupakinar
ab2512785b16d76a78f33e61a60ce90a2b225233	https://katb.in/apalezeholo	https://katb.in/amaxagapaza	https://katb.in/lopacuridig

I wrote functions to plot the three metrics, $D_{(i,j)}$, L_j , and R_i . Based on these I decide the values of the threshold for each of these metric (that is plotted as well).

$$L_{\text{thresh}} = \text{mean} \left(\bigcup_i \{ D_{(i,j)} : L_j > D_{(i,j)} \} \right) = 0.4388$$

[The threshold for LLM is the mean of similarity values for all the times when it outperforms the developer's original message.]

$$D_{\text{thresh}} = \text{mean} \left(\bigcup_i \{ L_j : D_{(i,j)} > L_j \} \right) = 0.3574$$

[The threshold for Developer is the mean of similarity values for all the times when Developer outperforms the LLM (CPT5)'s inference.]

$$R_{\text{thresh}} = \text{median}(R_i) - 1.3 \times \text{median}(|R_i - \text{median}(R_i)|) = 0.6156$$

[So, if the similarity between the rectified message for all modified files in a commit is more than 1.3 times the median absolute deviation away from the developer's commit message, then the rectifier is needed.]

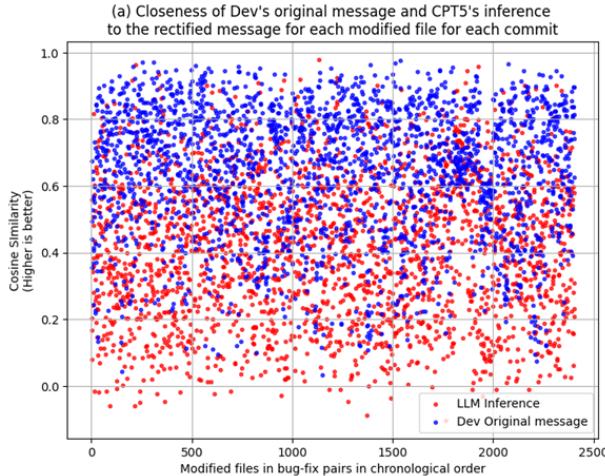


Figure 2.9: Dev's message, and CPT5's inference, as compared on rectified message.

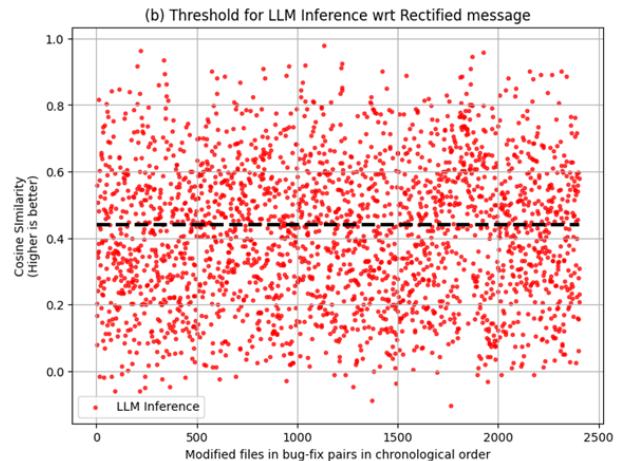


Figure 2.10: CPT5's inference, and threshold.

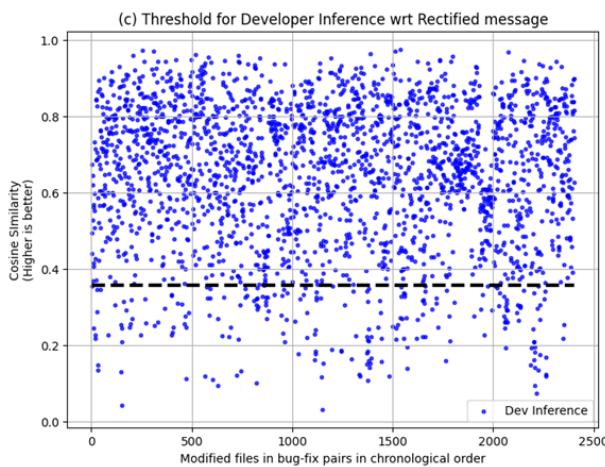


Figure 2.11: Dev's inference, and threshold.

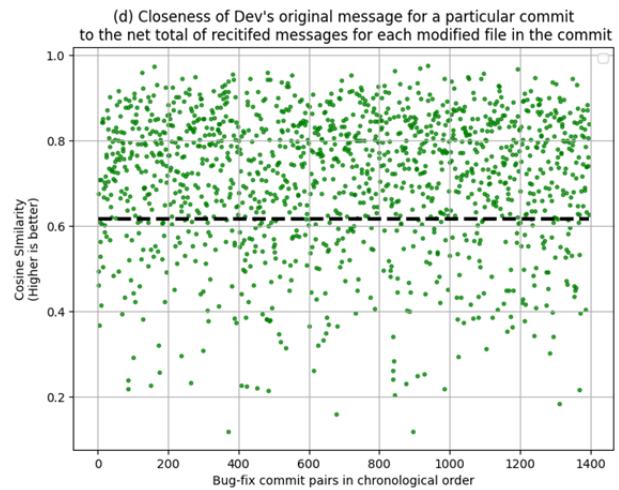


Figure 2.12: Rectified message, and threshold.

So, finally, the hit rates come out to be:

Developer hit rate = 90.686 %

LLM hit rate = 46.611%

Rectifier hit rate = 23.407%

2.5 Discussion and Conclusion

In this lab session, I built a tool to mine a GitHub repository, analyze the bug-fixing commits, and answer really interesting questions, such as, Do developers use a precise commit message? Does using the CPT5 LLM work well, and finally, does the rectifier I created work well?

The results indicate that 90% of the bug fix commits were clearly explained by the developer. However, there were cases where the developer committed multiple files in one commit but failed to list all the changes. That's where the rectifier comes into play. Out of the 1397 commits analyzed, 23.4% of the commit messages required rectification. As far as CPT5 LLM is considered, it only performed well about 46.6% of the time, which shows that a small local model like CPT5 does not generalize well. It is helpful to generate a commit message, but it requires the user's input as well. The major difficulty I faced during this lab was of the resources to compute the embeddings, and gpt-4o-mini. I also had to use async functions

to speed up the analysis of the repository.

For future work in this, I can weigh each modified file in a particular commit by the number of lines that have changed for that file. In this way, I will ensure that the commit message stays short in length and still effectively captures the major changes that have been made (major bug fixes). I also took inspiration to find the bug-fix commits from the minecpp library available on GitHub. I would like to explore that even more, and how it works with ASTs, and how that may outperform my code.

References

- [1] Puter, "Puter developer documentation," <https://developer.puter.com/>, accessed: 15 August 2025.
- [2] "Minecpp github repository," <https://github.com/SET-IITGN/MineCPP/>, SET-IITGN, accessed: 15 August 2025.
- [3] "Lab assignment document (on google classroom)," Google Classroom, SET-IITGN, accessed: 15 August 2025.
- [4] Wikipedia contributors, "Median absolute deviation — Wikipedia, the free encyclopedia," https://en.wikipedia.org/wiki/Median_absolute_deviation, 2025, accessed: 15 August 2025.
- [5] Contributors, "wordcloud: A little word cloud generator in python," <https://pypi.org/project/wordcloud/>, accessed: 15 August 2025.
- [6] D. Spadini, M. Aniche, and A. Bacchelli, *PyDriller: Python Framework for Mining Software Repositories*, 2018.
- [7] U. d. S. i. Software Institute, "Seart github search," <https://seart-ghs.si.usi.ch/>, accessed: 15 August 2025.

Lab 3

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

[Github Repository for Lab 3](#)

3.1 Introduction

This lab is a continuation of the previous lab. We add additional code quality metrics to the rectified messages database. We are instructed to use `radon`, to measure **Maintainability Index** (it is a function of Lines of Code, Cyclomatic Complexity, and other factors. It tells how easy to support and change the source code is), **Cyclomatic Complexity** (number of linearly independent paths through the code), and **Lines of Code** (I have excluded blank spaces, so basically, I have counted the number of source lines of code). Then we were to use `SacreBLEU` (a python library), and `CodeBERT` (from Hugging Face) to analyze the difference in BLEU and BERT scores of the source code after the commit and source code before the commit. This was then used to find a threshold, to find if there was a Major Fix or a Minor Fix in a modified file of any commit.

3.2 Setup and Tools

I am using `radon`, the `microsoft/codebert-base` model, `sacrebleu` model for embedding generation, and `rectified_commits.csv` (abbreviated as **rect database**) from the previous lab. I will also be using a virtual environment, so that my global packages do not interfere with `torch` and, `transformers` libraries. and `python --version` gives Python 3.12.4. Other versions are shown in the images.

```
C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\Labs\STT_Lab_3>pip show sacrebleu
Name: sacrebleu
Version: 2.5.1
Summary: Hassle-free computation of shareable, comparable, and reproducible BLEU, chrF, and TER scores
Home-page:
```

Figure 3.1: sacrebleu version

```
C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\Labs\STT_Lab_3>pip show radon
Name: radon
Version: 6.0.1
Summary: Code Metrics in Python
```

Figure 3.2: radon version

3.3 Methodology and Execution

I start with the trivial data analysis, such as finding the total number of commits and files. For total commits I just check the unique number of hashes in the rect database. And for total number of files, I check the unique filenames in the rect database. (One thing to note is, I do not store just the filename, I store the file path also). Since each row in the database is a file, to get average files per commit, we can simply get a frequency distribution of the hashes, and take the mean.

```
1 total_commits = len(modified_files_dataset['hash'].unique())
2 total_files = len(modified_files_dataset['filename'].unique())
3 average_files_per_commit = modified_files_dataset['hash'].value_counts().mean()
```

Next task is to get the distribution for the fix type, so, upon inspecting the LLM Inference, almost always the first word directly hints at what the fix type should be, So I make five categories of fixes, namely add, remove, refactor, fix, update, misc. And following is the first word assignment to each category:

```

1 term_to_category = {
2     "update": "update", "upgrade": "update", "improve": "update", "apply": "update",
3     "set": "update",
4     "use": "update", "simplify": "update", "replace": "update", "convert": "update",
5     "reorganize": "update", "fix": "fix", "escape": "fix", "ignore": "fix", "handle"
6     ": "fix",
7     "reset": "fix", "add": "add", "support": "add", "call": "add", "poll": "add",
8     "remove": "remove", "delete": "remove", "disconnect": "remove", "stop": "remove"
9     ",
10    "terminate": "remove", "rename": "refactor", "return": "refactor", "pass": "refactor",
11    "deprecate": "refactor", "send": "misc", "check": "misc", "raise": "misc"
12 }

```

With this I can easily check what the fix type of my commit is. In a similar manner I also find the frequency distribution for the file extensions, which tells me what are the extensions of most modified files. Now I proceed to updating the rect database, by first computing the structural metrics using radon. I run the radon's analyze and visit methods to find out the Lines of Code, Cyclomatic Complexity and Maintainability Index. One thing to note is, that incase a new file is created, or a file is deleted then one of the Source Code Before, or Source Code After values will be NaN. So to handle them, I simply return 0 as each of the metrics. These metrics are computed for both the Source Code Before commit and the Source Code After commit. Then the differences are appended to the rect database.

Moreover, radon only works for python files, so, for files which have codes other than what radon can analyze, the CC and the MI values are set to 0, however, the lines of code can still be computed.

Next task is to compute the Sematic Similarity (using CodeBERT) and the Token Similarity (using BLEU). So I load the CodeBert Model, generate embeddings for the Source Code Before commit and Source Code After Commit, and then find the cosine similarity between them. Since I was running this on my GPU, so I had to split the rect database rows in batches.

Similarly, I use the BLEU score to find the token similarity between Source Code Before and Source Code Current.

```

def radon_metrics(code):
    if isinstance(code,float):
        return 0,0,0
    try:
        raw = analyze(code)
        sloc = raw.sloc
        mi = mi_visit(code, True)
        cc_values = [r.complexity for r in cc_visit(code)]
        cc_total = int(sum(cc_values)) if cc_values else 0
    except SyntaxError:
        mi = 0
        cc_total = 0
        sloc = len(code.splitlines())
    return mi,cc_total,sloc

```

Figure 3.3: Radon metrics function

```

for i in range(0, len(before_code_list), batch_size):
    batch_before = before_code_list[i : i + batch_size]
    batch_current = current_code_list[i : i + batch_size]

    embeddings_before = generate_embeddings(batch_before)
    embeddings_current = generate_embeddings(batch_current)

    bert_cos_similarity = F.cosine_similarity(embeddings_before, embeddings_current)
    all_similarities.extend(bert_cos_similarity.cpu().numpy())

```

✓ 6m 59.7s

Figure 3.4: Computing similarity for BERT.

```

def compute_bleu(row):
    if pd.isna(row['source_code_before']) or pd.isna(row['source_code_current']):
        return np.nan
    before_ls = [row['source_code_before']]
    after = row['source_code_current']
    return sacrebleu.sentence_bleu(after, before_ls).score

```

Figure 3.5: Compute BLEU function

Then I plotted these values as a histogram with y axis in log scale, to figure out a good threshold, for which, I can classify the commit as a Minor Fix or a Major Fix for a particular modified file. Lastly, I

added one more column, Classes Agree, where if both Semantic Similarity is less than semantic threshold (Major Fix) and Token Similarity is less than token threshold (Major Fix), or vice-versa then both classes agree with each other.

3.4 Results and Analysis

The total number of commits are: 1402

The total (unique) number of modified files are: 146

The number of average files per commit is: 1.715

I categorize each row of the rect database, on the basis of first word of LLM Inference into six bug-fix type categories, and the histogram for the same is shown in Fig. 3.6.

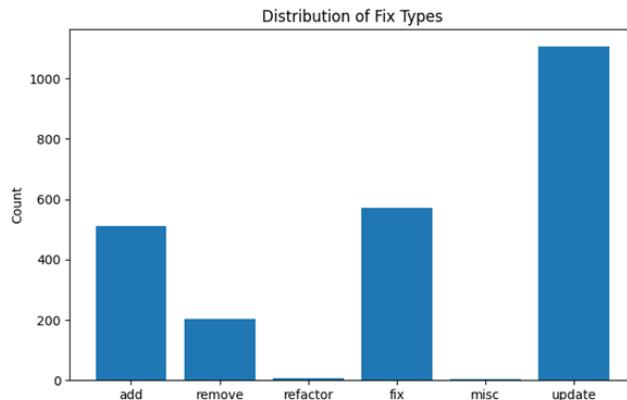


Figure 3.6: LLM inference fix types

This shows that most of the bug-fix commits are update in nature, that is, apart from fixing a bug in that commit, the developers update other code too. This category is followed by the fix category, where the developers explicitly fix a described issue. Next category is add, here developers, apart from fixing a bug, ad additional functionality. Lastly we have the remove category, here it may be the case that bug causing code is removed, or deprecated. The last two categories, refactor and misc, might be false positives, as the rect database might have some rows which are not bug-fix commits (because the developer might have used one of the 52 keywords in some context other than bug fixing (Refer to Lab 2)).

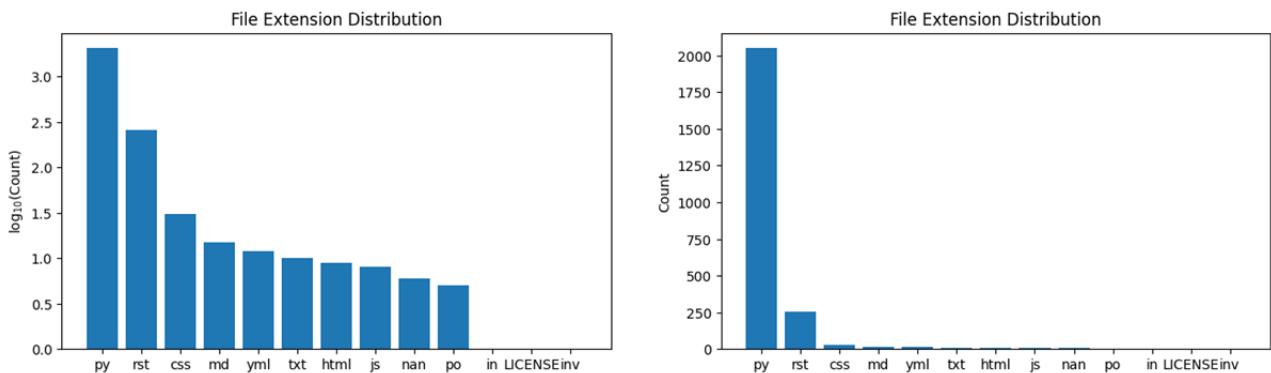


Figure 3.7: File extension distribution, with log scale and linear scale

The two graphs in Fig. 3.7, depict the same file extension distribution, the first graph is in log scale for y axis to promote visibility of the less frequent file extensions. The most common file extension for modified files, with over 2000 count is .py. which is preferred as the radon analysis works only with python files. The second most frequent extension is .rst, which is for restructured text, and it is used in the files which have documentation of the analyzed library.

After running the radon analysis for all the modified files, (a few of them will error out, as they have extension which is not .py), we obtain the subsequent plots.

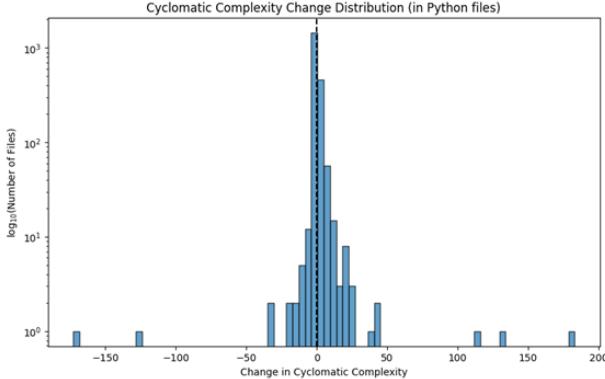


Figure 3.8: CC change distribution

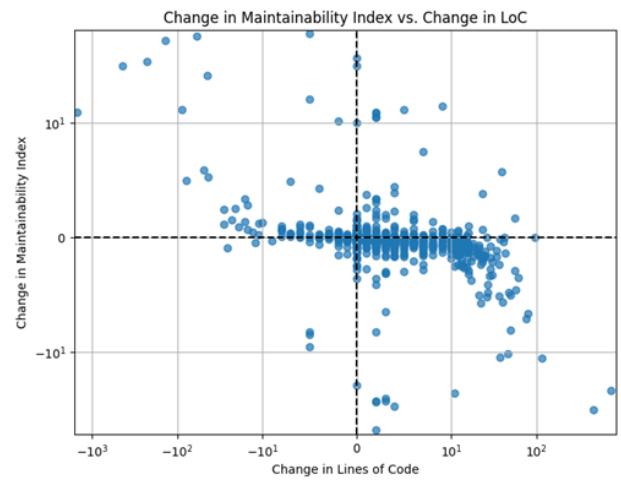


Figure 3.9: MI vs LOC distribution

Fig. 3.8 shows a almost gaussian like distribution (slightly skewed to right), almost equal number of files have an increased and a decreased complexity. Most of the files (1000) have no significant complexity change. But overall we can say that the complexity has slightly increased (distribution is skewed to right, mean of cc change is 0.73 which is more than 0).

Fig. 3.9 indicates how maintainability is affected by change in number of lines of code. Note that both the axes are in log scale.

We can observe a clear negative correlation with change in maintainability index and change in lines of code. As more code is introduced, the maintainability index goes down, and as code is removed, the maintainability index goes up.

Similar trend is observed if we plot change in maintainability index vs change in code complexity (Fig. 3.10).

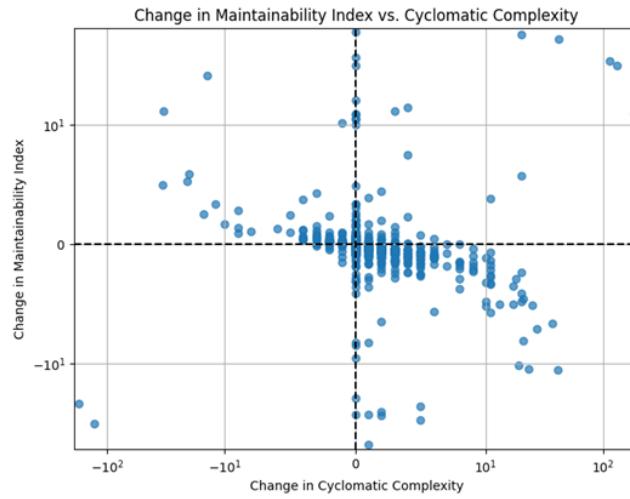


Figure 3.10: MI change vs CC distribution

Then I used CodeBERT and sacrebleu to compute and plot the BERT and BLEU scores, and on that basis decide the thresholds to classify a fix as Major or Minor. One thing to note is, the repository I chose is in continuous development, and is backwards compatible, So I do not think there would be many Major fixes, hence I will choose conservative thresholds, which favor Minor fixes.

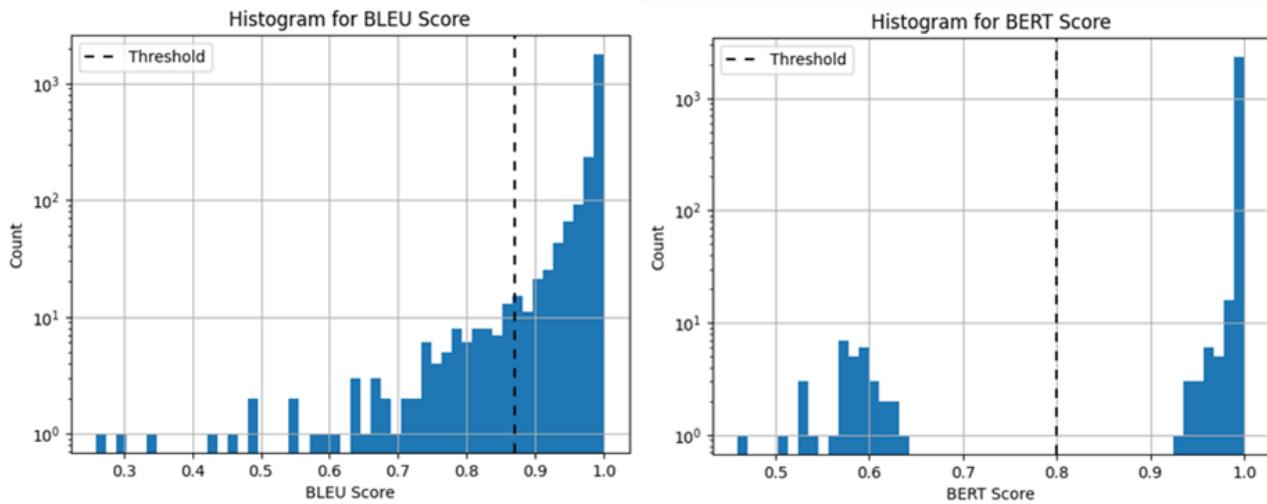


Figure 3.11: Threshold for BLEU and BERT

For BERT, we can clearly see a bimodal distribution (Fig. 3.11). This indicates a clear boundary for classifying a fix as a Major Fix or a Minor Fix. So for BERT I chose a threshold of 0.8. However, for BLEU there is no such boundary but the shape of histogram changes near 0.9, but a better way would be to set the threshold to the 4th quantile, which comes out to be 0.87. After these commits are classified into Major or Minor on the basis of BERT and BLEU scores, I make a confusion matrix, to see how they both agree on whether a commit is a major fix or a minor fix.

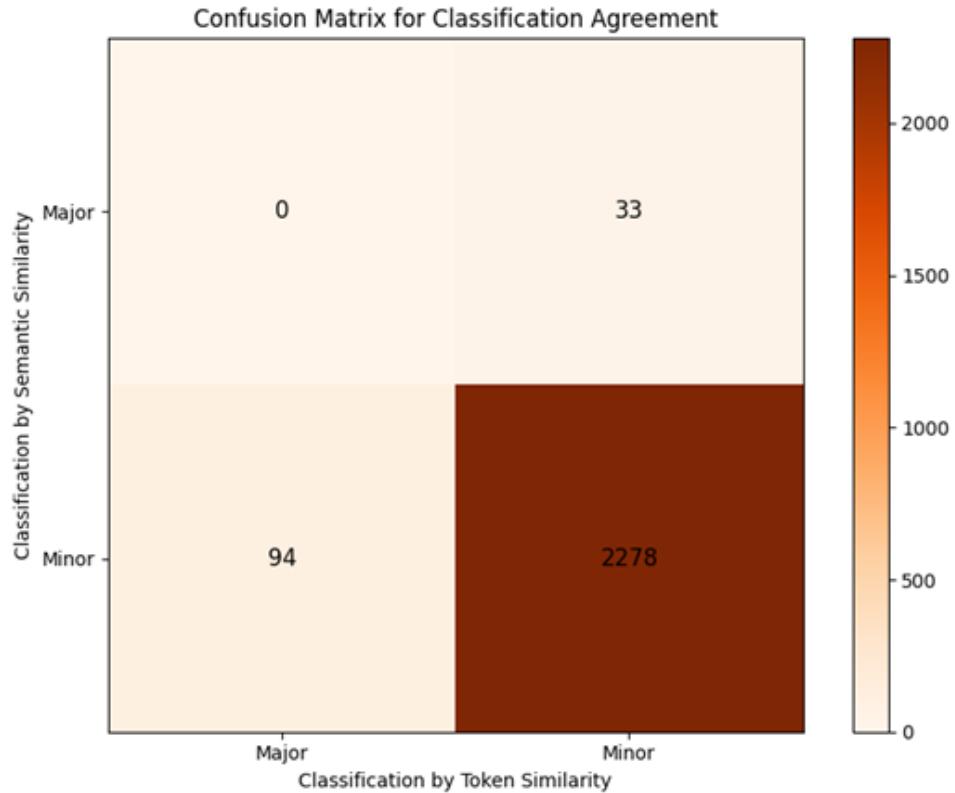


Figure 3.12: Confusion Matrix using BLEU and BERT

Fig. 3.12 shows that minor fixes are classified really well. The overall agreement rate is

$$\frac{2278}{2405} \times 100 = 94.72\%$$

The most common type of disagreement is CodeBERT says Minor Fix and BLEU says Major Fix. This

indicates that most of the text shares similar words (high token similarity) but their underlying meaning is different (low semantic similarity).

For cases where BLEU says Minor Fix and CodeBERT says Major Fix, here the words are different but the meaning remains the same.

Both CodeBERT and BLEU do not agree on even a single Major Fix. This shows the fundamental difference between checking similarity using both the models.

One captures the meanings and the other captures the semantics. However, it could also be the case that the repository (python library) I chose, discord.py, has only incremental changes, and there are no significant breaking changes (It can be confirmed that this library is backwards compatible).

3.5 Discussion and Conclusion

In this lab session, I extended the previous lab's dataset, and added structural metrics to better analyze the bug fix commits. I also capture the code complexity change, maintainability index change, and source lines change using the radon library. Next I use CodeBERT and sacrebleu to compute BERT and BLEU scores, and apply thresholds derived from the patterns in data, and classify each bug-fix pair as a Major Fix or a Minor Fix. Lastly, I plot the confusion matrix and analyze where both metrics agree or disagree.

The results indicate that, as the number of lines of code increase, the maintainability index decreases which is in turn an effect of increase in code complexity. Note that these structural metrics were only available for python files. Between BERT and BLEU, the former more accurately captures the semantics and gives a clear distinct threshold, below which I call the commits as Major fix. This is not the case with BLEU. It is hard to capture the notion of Major fix in BLEU, unless we know the developer's philosophy that there will always be incremental changes, and no major code change between commits. Lastly both these agree well on Minor Fixes, but do not generalize to Major Fixes, as same code can be written with less or different tokens (leading to low BLEU) but having similar semantic meaning (leading to high BERT).

I did not face much difficulty in this lab, as most of the tedious work of dataset creation was done in the previous lab. I would, however, like to run this analysis on some other repository to see how well these techniques generalize.

References

- [1] Radon-Developers, "Radon documentation," <https://radon.readthedocs.io/en/latest/intro.html>, accessed: 15 August 2025.
- [2] Shouwick-Mondal, "Lab assignment document (on google classroom)," Google Classroom, accessed: 15 August 2025.
- [3] Matt-Post, "Sacrebleu github repository," <https://github.com/mjpost/sacrebleu/blob/master/README.md>, accessed: 15 August 2025.
- [4] Microsoft, "Codebert github repository," <https://github.com/microsoft/CodeBERT/blob/master/README.md>, accessed: 15 August 2025.

Lab 4

Exploration of different diff algorithms on Open-Source Repositories

[Github Repository for Lab 4](#)

4.1 Introduction

In this lab, we explored how different diff algorithms in Git, specifically **Myers** and **Histogram**, produce different outputs when applied to real-world open-source repositories. Diff algorithms play an important role in version control systems such as git, as they determine how changes in **code** and **other** files (markdown files, test files and, license file) are represented. Even small differences in the algorithm can impact how developers interpret modifications. For this lab, I selected three open-source repositories and used `pydriller` to analyze their commit histories. Each repository was partitioned into code artifacts and non-code artifacts. The purpose of separation is to study if diff algorithms behave differently depending on the type of file being analyzed. I also heavily refer to the paper by Yusuf Nugroho, on different diff algorithms in Git. I generated both diff outputs (Myers and Histogram) for all modified files in all commits, and compared them to identify mismatches. Lastly the mismatches were analyzed across file type categories to understand the effectiveness of both the diff algorithms.

4.2 Setup and Tools

I am using `pydriller` to obtain data of commit history, histogram diff and myers diff for the chosen repositories, `pandas` library for dataframe processing, `matplotlib` for plotting. In my environment, `python --version` gives Python 3.12.4, the `pydriller` version is shown in Fig. 4.1. SEART Search Engine <https://seart-ghs.si.usi.ch/> is used to search and filter the repositories based on my criteria.

```
C:\Users\GSRAJA\Desktop\IIT GN\sem5\sttcse\Labs\STT_Lab_4>pip show pydriller
Name: PyDriller
Version: 2.8
Summary: Framework for MSR
Home-page: https://github.com/ishepard/pydriller
Author: Davide Spadini
Author-email: spadini.davide@gmail.com
License: Apache License
```

Figure 4.1: Pydriller version

4.3 Methodology and Execution

I first searched through the SEART search engine, three repositories, which are mainly based on python, to use for analysis during this lab. I chose the following repositories: `unclecode/crawl4ai`, `datalab-to/marker`,

and `zylon-ai/private-gpt`. Fig. 4.2 shows the search criteria I used, I ensured the following things:

- Major language should be python.
- It should have a minimum of 20000 stars.
- It has a license file.
- It has a wiki (having a wiki increases the chances that the repository has a `tests` folder).

The screenshot shows the SEART search engine interface with several sections for filtering repositories:

- General**: Includes a search bar for "Search by keyword in name" (Contains dropdown) and a "Python" filter under "Uses Label".
- License**: Has topic filter.
- History and Activity**: Filters for Number of Commits (min, max), Number of Contributors (min, max), Number of Issues (min, max), Number of Pull Requests (min, max), Number of Branches (min, max), and Number of Releases (min, max).
- Date-based Filters**: Created Between (dd-mm-yyyy to dd-mm-yyyy) and Last Commit Between (dd-mm-yyyy to dd-mm-yyyy).
- Popularity Filters**: Number of Stars (20000, max), Number of Watchers (min, max), and Number of Forks (min, max).
- Size of codebase ⓘ**: Non Blank Lines (min, max), Code Lines (min, max), and Comment Lines (min, max).
- Additional Filters**: Sorting (Name, Ascending), Repository Characteristics (Exclude Forks, Only Forks, Has License, Has Wiki, Has Open Issues, Has Pull Requests), and a "Search" button.

Figure 4.2: SEART search engine (criteria for searching)

I use the following procedure to extract information from each of the three repositories:

1. Iterate through the commits in the repository, with the parameters of `skip_whitespaces=True`, and `zip(Repository(..., histogram_diff=True), Repository(..., histogram_diff=False))` to access both the myers diff and histogram diff in same loop.
2. If the length of modified files in this commit is 0, then continue, as this would be a merge commit.
3. For each file in the modified files, store the following details:

```

1  {
2      "old_file_path": mf.new_path,
3      "new_file_path": mf.new_path,
4      "hash": commit.hash,
5      "parents_hash": commit.parents,
6      "message": commit.msg,
7      "diff_myers": mf.diff,
8      "diff_hist": mf_hist.diff,
9      "myers_added_lines": mf.diff_parsed["added"],
10     "myers_deleted_lines": mf.diff_parsed["deleted"],
11     "hist_added_lines": mf_hist.diff_parsed["added"],
12     "hist_deleted_lines": mf_hist.diff_parsed["deleted"],
13 }
14

```

To check if Myers and Histogram diff resulted in the same or the different outputs, I use the same methodology as described in the paper [1]. Fig. 4.3 shows the screenshot from that paper.

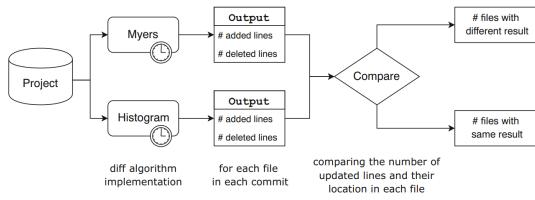


Fig.11 Overview of the metrics collection procedure

Figure 4.3: Screenshot from [1].

Using this, we say there is no discrepancy if #lines added by Myers = #lines added by Histogram as well as #lines removed by Myers = #lines removed by Histogram.

All the three repositories which I chose have LICENSE file, few markdown files, and test files under a folder name tests/. So I make bitmasks for these conditions on path of the modified file, and with that I can filter out the **code** and the **non code** files (Fig. 4.4).

```

def statistics(df):
    md_files_mask = df['new_file_path'].apply(lambda x: False if x == None else x.endswith('.md'))
    license_files_mask = df['new_file_path'].apply(lambda x: False if x == None else x.lower() == 'license')
    tests_files_mask = df['new_file_path'].apply(lambda x: False if x == None else x.startswith('tests\\'))
    source_files_mask = ~(md_files_mask | license_files_mask | tests_files_mask)
    return {
        'md':df[md_files_mask]['Discrepancy'].value_counts(),
        'license':df[license_files_mask]['Discrepancy'].value_counts(),
        'tests':df[tests_files_mask]['Discrepancy'].value_counts(),
        'source':df[source_files_mask]['Discrepancy'].value_counts()
    }

```

Figure 4.4: Bitmasks for code and non code files.

To evaluate programmatically which diff algorithm is better, I again refer to the paper given [1],

- 1 The diff results might show different change region with a contiguous list of deleted and added lines that is called as a change hunk (Ray et al., 2015).
- 2
- 3 We expect that a set of changing operations done by developers can be represented by change hunks.

Using this ideology, I evaluate both diff algorithms by treating contiguous added or deleted lines as change hunks. Here I assume that a developer's single logical edit often appears as one hunk.

1. For each modified file and for each algorithm (Myers and Histogram) I extract the change hunks separately for added lines and for deleted lines.
2. Ignore 1-line hunks when computing statistics, because one-line differences are common noise and do not reliably represent an operation.
3. Compute:
 - (a) The median length of the change hunks (combined over added and deleted hunks)
 - (b) The average number of change hunks approximated by,

$$\frac{\text{\# lines changed}}{\text{median hunk length}}$$

where # lines changed is computed separately for added and deleted lines and then averaged.

4. The average-number measure estimates how many logical hunks the changed lines are split into.
5. If $(\text{median hunk length})^{\text{algo1}} < (\text{median hunk length})^{\text{algo2}}$, then algo1 captures smaller hunks (more fine grained developer operations).
6. Else If there is a tie, then I use the maximum hunk length rather than the median.

7. I prefer the algorithm that overall produces smaller median hunk lengths and higher average number of hunks.

4.4 Results and Analysis

4.4.1 Private-gpt

Fig. 4.5 shows, in a tabular format, the extracted details for each commit.

	old_file_path	new_file_path	hash	parents_hash	message	diff_myers	diff_hist
24	chroma_preference.py	chroma_preference.py	8c6a81a07f9c800d53f62a33f5ae3b5247a22a6	[60226598b6bd46f905a952523a2974980ade0ca]	Fix Disable Chroma Telemetry\n\rOpt-outs of a... chromadb.config import_...	@@ -0.0 +1.11 @@-0\n\rfrom chromadb.config import_...	@@ -0.0 +1.11 @@-0\n\rfrom chromadb.config import_...
413	fern/fern.config.json	fern/fenr.config.json	f339f7608c5b1b7592b5075ac340df0ce40da7f	[ff7e2bc9d9da7a694a6e940313261f888be902122]	Move Docs to Fern (#1257)	@@ -0.0 +1.4 @@-0\n\rfrom "organization"; "priv... "organization": "priv... "organization": "priv...	@@ -0.0 +1.4 @@-0\n\rfrom "organization"; "priv... "organization": "priv... "organization": "priv...
507	settings.yaml	settings.yaml	56af625d71afc621f4830d07e0797aaea2e193862	[b7ca7d35a0c181bf79acaa00d1e66e50c18f97]	Fix the parallel ingestion mode, and make it a...	@@ -1.3 +1.6 @@-0\n#\tThe default configuration ...	@@ -1.3 +1.6 @@-0\n#\tThe default configuration ...
748	settings.yaml	settings.yaml	966af4771dbe5c3fd5f54b5fdff732407895c4	[d13029a046f6e19e8e65be3acd96365c738df2]	fix(settings); enable cors by default so it wi...	@@ -5.7 +5.7 @@-0\nserver\nenv_name: '\$APP_E...	@@ -5.7 +5.7 @@-0\nserver\nenv_name: '\$APP_E...
793	fern\docs\pages\installation\concepts.mdx	fern\docs\pages\installation\concepts.mdx	4523a30c8f004aac7a7e224671e2c45ec0b973	[01b7ccdd0648be032846647c9a184925d3682f612]	feat(docs): update documentation and fix previ...	@@ -8.18 +8.25 @@-0 It supports a variety of LLM...	@@ -8.18 +8.25 @@-0 It supports a variety of LLM...

Figure 4.5: Details for each modified file in all commits.

All the table images, I have added in this document, are a sample of the actual csv files which are available on the [Github Link](#). After processing the discrepancies, we make a new csv with the discrepancy column, as shown in Fig. 4.6.

old_file_path	new_file_path	hash	parents_hash	message	diff_myers	diff_hist	Discrepancy
30 requirements.txt	requirements.txt	85528db7431d57a2bb2ca6e3e8d0fb37a66bf8d2	[918b384e38f34a709c196ccfa81bcb67e10de0e]	Update langchain to 0.16.0\n[n]Tested @ [v]n[Rele...	@@ -1.4 +1.4 @@-[v]n[Rele...	@@ -1.4 +1.4 @@-[v]n[Langcha...	False
70 ingest.py	ingest.py	7844553ca1fd39aa5962795f2209e503522212b1	[f42046c5ec0322e02e545cd43d96e09d2f4287]	Implement a way of ingestting more documents\n[v]M...	@@ -24.6 +24.16 @@ from langchain.docstore.doc...	@@ -2.6 +24.16 @@ from langchain.docstore.doc...	False
113 ingest.py	ingest.py	28537b6a84ce51487f477e2619d1c177e0f399b	[b1057afdf8f65fdb0e4160abd84626e08271]	Better error message if .env is empty/does not...	@@ -24.11 +24.12 @@ from langchain.text_split...	@@ -24.11 +24.13 @@ from langchain.text_split...	True
147 README.md	README.md	51c638758f68032bd971ed7d80fa5e284f185de	[78d1ef44adea1b72235a4cb603bbf0e4d9033d10]	Next version of PrivateGPT (#1077)\n[v]n# Docker...	@@ -1.152 +1.158 @@ [v]n# privateGPT\n[v]n# Ask quest...	@@ -1.152 +1.158 @@ [v]n# privateGPT\n[v]n# Ask quest...	True
159 poetry.lock	poetry.lock	51c638758f68032bd971ed7d80fa5e284f185de	[78d1ef44adea1b72235a4cb603bbf0e4d9033d10]	Next version of PrivateGPT (#1077)\n[v]n# Docker...	@@ -1.99 +1.110 @@ [v]n# This file is automatica...	@@ -1.99 +1.110 @@ [v]n# This file is automatica...	True

Figure 4.6: Sample of commits, after processing the discrepancies

Overall we see that, 1.7% of the modified files have a discrepancy (Fig. 4.7). Fig. 4.8 shows the discrepancy across the markdown, license, tests, and the source code files.

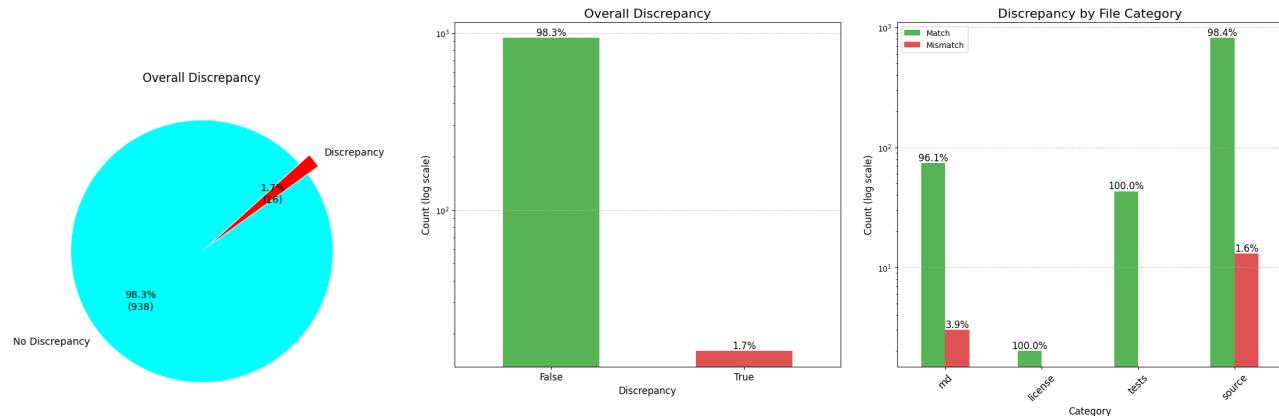


Figure 4.7:
Pie chart for
overall dis-
crepancy.

Figure 4.8: Bar graph showing discrepancy across categories of files.

```
1 #Mismatched for Source Code files : 13  
2 #Mismatched for Test Code files : 0  
3 #Mismatched for README files : 3  
4 #Mismatched for LICENSE files : 2
```

Comparing both the diff algorithms based on the approach described before, we get the following results.

- 1 Hist is better in 12 instances (1.8%)
- 2 Myers is better in 11 instances (1.6%)
- 3 ties 646 (96.6%)

4.4.2 Crawl4ai

I will proceed with the graphs and diff algorithm comparison, as the csv file for details of commits and discrepancy analysis is available on GitHub.

- 1 #Mismatches for Source Code files : 77
- 2 #Mismatches for Test Code files : 2
- 3 #Mismatches for README files : 56
- 4 #Mismatches for LICENSE files : 2

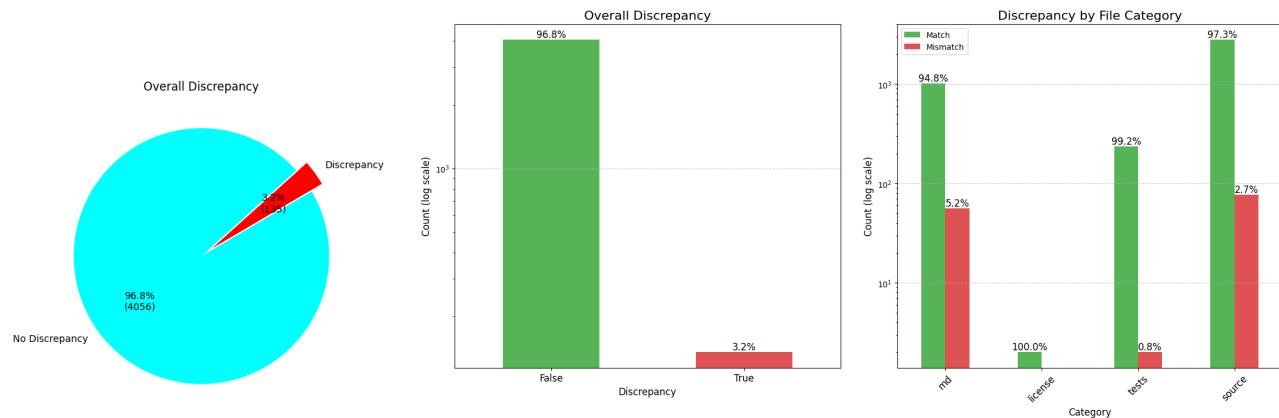


Figure 4.9:
Pie chart for
overall dis-
crepancy.

Figure 4.10: Bar graph showing discrepancy across cate-
gories of files.

We observe that, 3.2% of the modified files have a discrepancy (Fig. 4.9). Fig. 4.10 shows the discrepancy across the markdown, license, tests, and the source code files.

Comparing both the diff algorithms based on the approach described before, we get the following results.

- 1 Hist is better in 116 instances (3.8%)
- 2 Myers is better in 79 instances (2.6%)
- 3 ties 2878 (93.7%)

4.4.3 Marker

I will proceed with the graphs and diff algorithm comparison, as the csv file for details of commits and discrepancy analysis is available on GitHub.

- 1 #Mismatches for Source Code files : 73
- 2 #Mismatches for Test Code files : 2
- 3 #Mismatches for README files : 15
- 4 #Mismatches for LICENSE files : 2

We observe that, 2.1% of the modified files have a discrepancy (Fig. 4.11). Fig. 4.12 shows the discrepancy across the markdown, license, tests, and the source code files.

Comparing both the diff algorithms based on the approach described before, we get the following results.

- 1 Hist is better in 71 instances (2.7%)
- 2 Myers is better in 62 instances (2.3%)
- 3 ties 2516 (95.0%)

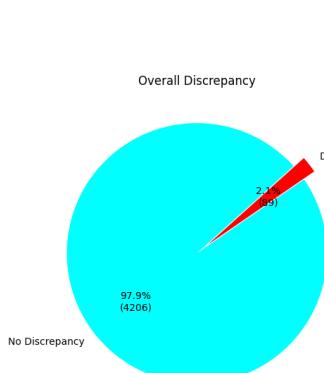


Figure 4.11:
Pie chart for
overall dis-
crepancy.

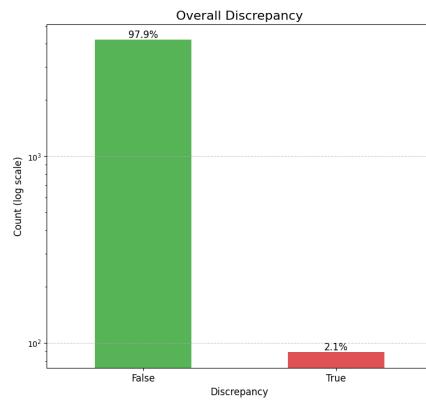
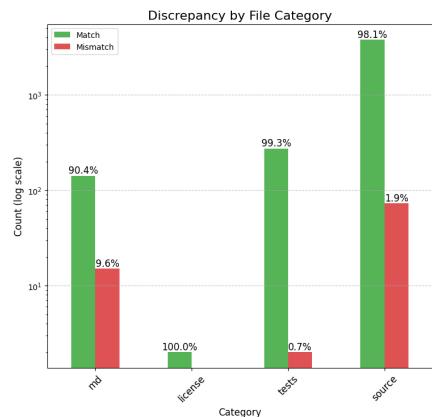


Figure 4.12: Bar graph showing discrepancy across cate-
gories of files.



4.5 Discussion and Conclusion

In this lab I compared Myers and Histogram diff algorithms on three real-world repositories to see which algorithm produces diff that better reflect developers' edits. I have assumed (motivated by [1]) that a developer's single logical change tends to appear as a contiguous run of changed lines (a change hunk). I simply find the median hunk length and the average number of code hunks, for both Myers and Histogram diff algorithms, and whichever has smaller median hunks and larger number of hunks is treated as better aligned with developer's intent.

The three repositories show low overall discrepancy rates between the two algorithms: 1.7%, 3.2% and 2.1% of modified files had different added/deleted counts, respectively, which is consistent with the [1]'s findings (1.7% to 8.2%).

Overall, while comparing the diff algorithms, Histogram wins more often than Myers. Across the three datasets I observed 199 Histogram-wins versus 152 Myers-wins (the remainder were ties). That imbalance supports the conclusion from [1] that Histogram is generally better at producing hunks that match developer edits in practice, although the effect size is small. This shows that the choice between Myers and Histogram usually does not change the results, since they agree most of the time. When discrepancies do occur, Histogram is often better. However, there are limitations: I assumed that developers' logical edits are contiguous. That assumption fails when developers interleave unrelated edits in the same commit or when edits are semantically equivalent (for example, `[x for x in range(10)]` vs. `list(range(10))`). We could combine the methodology from the previous lab (refer to Lab 3) for detecting semantic equivalence with our current approach to determine the better diff algorithm.

References

- [1] Y. S. Nugroho, H. Hata, and K. Matsumoto, "How different are different diff algorithms in git?" *Empirical Software Engineering*, vol. 25, no. 1, pp. 790–823, Jan 2020. [Online]. Available: <https://doi.org/10.1007/s10664-019-09772-z>
- [2] Shouvick-Mondal, "Lab assignment document (on google classroom)," Google Classroom, accessed: 28 August 2025.
- [3] "Git documentation," <https://git-scm.com/docs>, accessed: 2025-08-28.
- [4] D. Spadini, M. Aniche, and A. Bacchelli, *PyDriller: Python Framework for Mining Software Repositories*, 2018.
- [5] U. d. S. i. Software Institute, "Seart github search," <https://seart-ghs.si.usi.ch/>, accessed: 28 August 2025.