# Large Language Models for Code Refactoring
## *PhD Thesis Project*

**Researcher:** Henrique Nunes
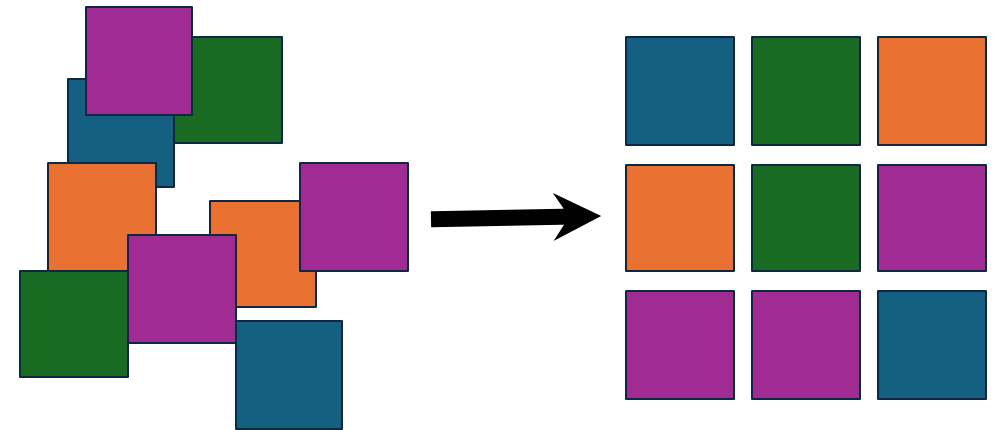
**Advisor & Coadvisor:** Eduardo Figueiredo & Tushar Sharma

UFMG
UNIVERSIDADE FEDERAL
DE MINAS GERAIS

lab·soft
software engineering laboratory

# Context of Our Work
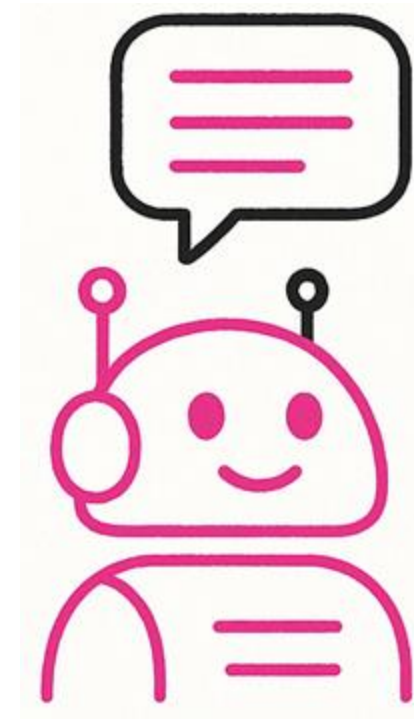*Why LLMs for Code Refactoring?*

# Refactoring

- Refactoring reduces code **complexity** and improves **maintainability**.

- Manual refactoring is costly and slow.

- Refactoring tools are semi-automated, inaccurate and require human intervention.

# Large Language Models (LLMs)

- LLMs are models trained on a large amount of data and have shown good capability to **write and understand natural language**.

- Studies on LLMs are growing exponentially in software engineering, yet their capability show **limitations** and still need to be explored.
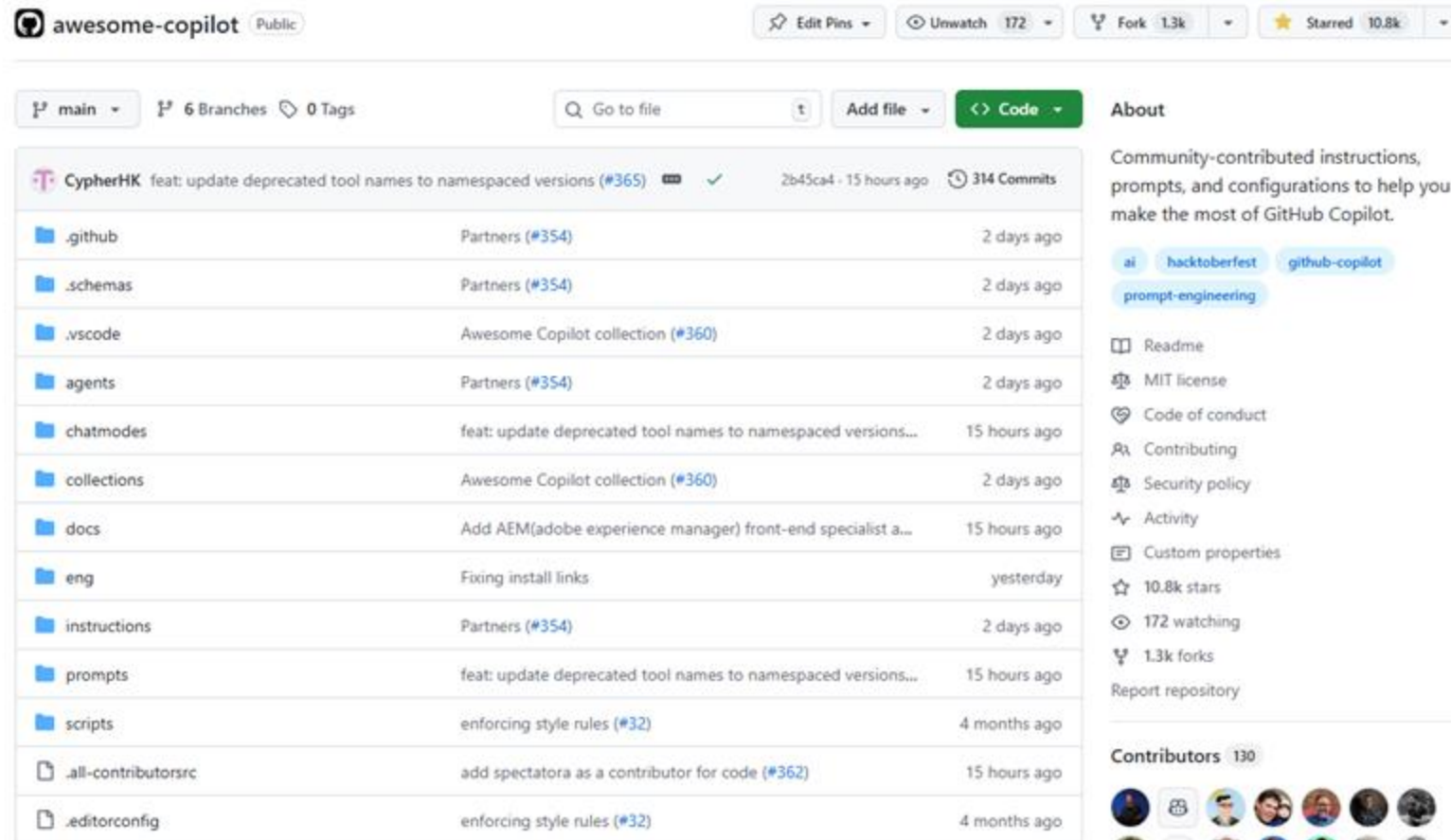
# Goals

- We aims to explore the extent to which variation in **prompt techniques** and their combinations affects **different refactoring types**.

- We intend to provide developers and tools with **practical insights** for creating refactoring prompts for real-world projects.

# Relevance

# Relevance

# Research Method



| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| Informal Literature Review | Machine Learning Study | LLM Refactoring Effectiveness Study | Refactoring Dataset Study | Refactoring Study |

# Contribution Progress

- Nunes, H. G., Santana, A., Figueiredo, E., & Costa, H. *Tuning code smell prediction models: A replication study*. In International Conference on Program Comprehension (**ICPC 2024**) - **Study 1**

- Nunes, H. G., Figueiredo, E., Rocha, L., Nadi S., Ferreira F., & Esteves, G. *Evaluating the effectiveness of llms in fixing maintainability issues in real-world projects*. In International Conference on Software Analysis, Evolution and Reengineering (**SANER 2025**) - **Study 2**

- Nunes, H. G., Sharma, T., Figueiredo, E. *MaRV: A Manually Validated Refactoring Dataset*. In International Conference on AI Foundation Models and Software Engineering (**FORGE 2025**) - **Study 3**

# Machine Learning to Detect Code Smells

*Study 1*

# Study Context & Goals

- **Replicate** the study by Cruz et al. (2023), which uses machine learning to detect code smells, but **with a dataset of modern systems**.

- We use **seven traditional machine learning algorithms** to detect **four code smells**.

- We also evaluate how data resampling, feature selection, and polynomial feature techniques affect code smell predictions.

# Design of Study 1



1. Data Separation
2. Data Analysis
3. Models Parametrization
4. Models Comparision
5. Test on Unseen Data
6. Feature Engineering and Resample Study

# Study Result & Conclusion

- The performance of traditional ML algorithms for code smell detection are **limited**.

- ML techniques, especially resampling, improve prediction performance, but **not enough**.

- **This study indicated that we should explore for other solutions. LLMs?**

# Large Language Models for Code Refactoring

*Study 2*

# Study Context & Goals

- We evaluated the **effectiveness of LLMs for refactoring** 10 issue categories in real-world projects.

- Furthermore, the study assessed the most common **LLM failures**.

- We also conducted a human **readability evaluation**.

# Design of Study 2



① Collect and Build Projects → binaries

② Identify Maintainability Issues → 2,446 issues

③ Select Issues → 127 samples

④ Use LLMs to Fix Issues → 51 LLM-based solutions

⑤ Human Evaluation → readability results & LLM solutions

⑥ Evaluate Results

**lab·soft**
software engineering laboratory

# Effectiveness of the LLM

# Effectiveness of the LLM



Few-shot Llama fixed more maintainability issues (~45%)

# Effectiveness of the LLM



**All LLMs included compilation errors and caused test failures**

# LLM Failures (Hallucination)

Original

LLM

```java
private void put(DirectoryEntry entry,
    boolean overwriteExisting) {
    DirectoryEntry prev = null;
    [...]
    if (prev != null) {
        prev.next = entry;
    }
    [...]
}
```

```java
private void put(DirectoryEntry entry,
    boolean overwriteExisting) {
    [...]
    overwriteEntry(existingEntry, entry);
    [...]
}

private void overwriteEntry(DirectoryEntry existingEntry,
    DirectoryEntry newEntry) {
    [...]
    if (existingEntry.prev!= null) {
        existingEntry.prev.next = newEntry;
    }
    [...]
}
```

lab·soft
software engineering laboratory

# LLM Failures (Hallucination)

Original

LLM

Refactoring

```
private void put(DirectoryEntry entry,
    boolean overwriteExisting) {
    DirectoryEntry prev = null;

    if (prev != null) {
        prev.next = entry;
    }
    [...]
}
```

? 

```
private void put(DirectoryEntry entry,
    boolean overwriteExisting) {
    [...]
    overwriteEntry(existingEntry, entry);
    [...]
}
```

```
private void overwriteEntry(DirectoryEntry existingEntry,
    DirectoryEntry newEntry) {
    [...]
    if (existingEntry.prev != null) {
        existingEntry.prev.next = newEntry;
    }
    [...]
}
```

lab·soft
software engineering laboratory

# Human Evaluation of LLM Readability (RQ3)

Out of 51 method pairs, **~70%** considered LLM-refactored code more readable.

# Study Conclusion

- Despite **limited** effectiveness (below 50%), **few shot learning** presented the highest effectiveness (~40%), if compared with zero shot learning (~30%).

- Overall, participants considered LLM-refactored code **more readable**.

- **The results of this study motivated us to evaluate more deeply the use of prompt techniques for code refactoring.**

# Refactoring Dataset

*Study 3*

# Study Context & Goal

- We used RefactoringMiner to **collect refactorings** and conducted a human evaluation to **select the most representative examples**.

- We aim to produce a high-quality **manually validated dataset** of actual refactorings from open-source projects, **MaRV**.

# Design of Study 3

# Manual Evaluation Results

| | Votes | Count | Total |
|---|---|---|---|
| Consensus | [disagree, disagree]<br>[agree, agree] | 84<br>321 | 405 |
| Conflict | [disagree, agree]<br>[I don't know, disagree]<br>[I don't know, agree] | 217<br>27<br>41 | 285 |
| Other | [I don't know, I don't know] | 3 | 3 |
| | | | 693 |

# Manual Evaluation Results

| | Votes | Count | Total |
|---|---|---|---|
| Consensus | [disagree, disagree]<br>**[agree, agree]** | 84<br>**321** | 405 |
| Conflict | [disagree, agree]<br>[I don't know, disagree]<br>[I don't know, agree] | 217<br>27<br>41 | 285 |
| Other | [I don't know, I don't know] | 3 | 3 |
| | | | 693 |

**Both participants agreed in 321 cases (~46%). We considered this cases for MaRV**

# Study Conclusion

- **We aim to use MaRV to support several prompting techniques (e.g., few-shot prompting).**

# Agenda
*Next Steps*

# Future Steps

1. Improve MaRV Dataset

2. Select LLMs and Prompting Techniques

3. Define the Prompt Design

4. Execute and Validate LLM-Based Refactorings

5. Evaluate Results and Consolidate in a Paper

6. Extend our LLM Study

7. Write and Submit a Journal Paper

8. Write and Present Thesis

**lab·soft**
software engineering laboratory

# Improve MaRV

# Prompt Evaluation

- We aim to evaluate **different prompt techniques** and their variations.

- We will define several prompt designs, **varying and combining techniques**.

- We want to identify **which prompts are most effectiveness** for each different refactoring type.



**Prompt Template**

**Natural Language Instruction**

# Task:

You are an expert in refactoring Java methods. Apply [refactoring] to improve readability.

[constraints list]
[procedure steps]
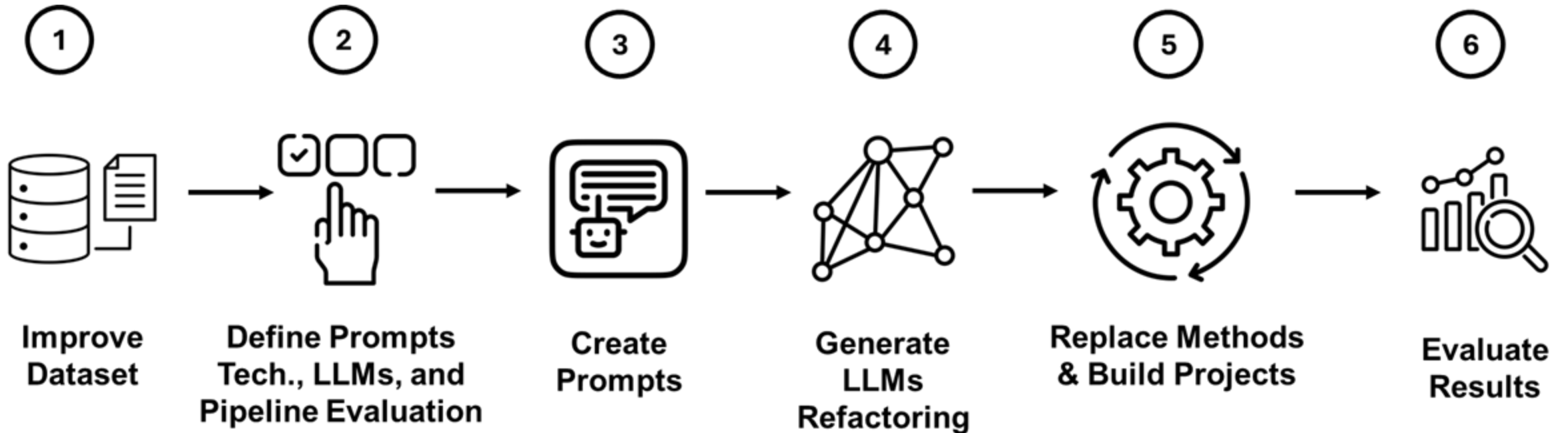
Below are 3 examples (input to output). Complete the 4th.

**Refactoring Demonstrations**

# Input:
[```example before commit```]
# Output:
[```example after commit```]

[two more demonstrations...]

**Method to be Refactored**

# Input:
[```input method```]
# Output:

# Next Steps Overview



1. Improve Dataset
2. Define Prompts Tech., LLMs, and Pipeline Evaluation
3. Create Prompts
4. Generate LLMs Refactoring
5. Replace Methods & Build Projects
6. Evaluate Results

# Agenda

# Thank you!
*Henrique Nunes | henrique.mg.bh@gmail.com*

lab-soft
software engineering laboratory