



**Hochschule für Technik  
und Wirtschaft Berlin**

University of Applied Sciences

*Modular Multi-Stage Agent for Bug Fixing - Analysis of Potentials and Limitations*

Abschlussarbeit

zur Erlangung des akademischen Grades

**Bachelor of Science (B.Sc.)**

an der

Hochschule für Technik und Wirtschaft (HTW) Berlin  
Fachbereich 4: Informatik, Kommunikation und Wirtschaft  
Studiengang *Internationale Medieninformatik*

1. Gutachter\_in: Prof. Dr. Gefei Zhang
2. Gutachter\_in: Stephan Lindauer

Eingereicht von Justin Gebert [s0583511]

22.07.2025

# Danksagung

[Text der Danksagung]

## Abstract

Generative AI technologies are reshaping software engineering practices by automating critical tasks, including code generation, debugging and program repair. Despite these advancements, existing Automated Program Repair (APR) systems frequently suffer from complexity, high computational demands, and missing integration within practical software development lifecycles, particularly Continuous Integration and Continuous Deployment (CI/CD) workflows. Such shortcomings often lead to frequent context switching, which negatively impacts developer productivity. In this thesis, we address these challenges by introducing a novel and lightweight APR system leveraging LLMs, explicitly designed for seamless integration into CI/CD pipelines deployed in budget constrained environments. Our containerized approach, developed with a strong emphasis on security and isolation, manages the complete bug-fixing lifecycle from issue creation on GitHub to the generation and validation of pull requests. By automating these processes end-to-end, the system significantly reduces manual intervention, streamlining developer workflows and enhancing overall productivity. We evaluate our APR system using the QuixBugs benchmark, a recognized dataset for testing APR methodologies. The experimental results indicate that our streamlined and cost-effective solution effectively repairs small-scale software bugs, demonstrating practical applicability within typical software development environments. The outcomes underscore the feasibility and advantages of integrating APR directly into real-world CI/CD pipelines. We also discuss limitations inherent in LLM-based solutions, such as accuracy and reliability issues and suggest future enhancement and research.

-add that this approach directly integrates into the development lifecycle reducing configuration ...

# Contents

<b>1. Introduction</b>	<b>1</b>
1.1. Background and Motivation . . . . .	1
1.2. Problem Statement . . . . .	1
1.3. Objectives and Research Questions . . . . .	1
<b>2. Background and Related Work</b>	<b>3</b>
2.1. Software Engineering . . . . .	3
2.1.1. code hosting platforms . . . . .	3
2.1.2. Software Development Lifecycle . . . . .	3
2.1.3. Continuous Integration and Continuous Deployment (CI/CD) . . . . .	3
2.2. Automated Programm Repair . . . . .	3
2.2.1. Evolution of Automated Program Repair . . . . .	3
2.3. LLMs in Software Engineering . . . . .	4
2.4. LLM-Based Tool Use and CI Context . . . . .	4
2.5. Related Work - Existing Systems . . . . .	4
<b>3. Requirements</b>	<b>6</b>
3.1. Functional Requirements . . . . .	6
3.2. Non-Functional Requirements . . . . .	6
<b>4. Methodology</b>	<b>9</b>
4.1. Preparation . . . . .	9
4.1.1. Dataset Selection . . . . .	9
4.1.2. Environment Setup . . . . .	9
4.1.3. Pipeline Architecture . . . . .	9
4.2. Evaluation . . . . .	9
<b>5. Implementation</b>	<b>11</b>
5.1. Implementation . . . . .	11
5.1.1. System Architecture . . . . .	11
5.1.2. System Components . . . . .	11
5.1.3. System Configuration . . . . .	11
<b>6. Results</b>	<b>12</b>
6.1. APR System Results . . . . .	12
6.2. Evaluation Results . . . . .	12
<b>7. Discussion</b>	<b>13</b>
7.1. Validity . . . . .	13
7.2. Potentials . . . . .	13

## *Contents*

7.3. Limitations . . . . .	13
7.4. Summary of Findings . . . . .	13
7.5. Lessons Learned . . . . .	13
7.6. Roadmap for Extensions . . . . .	13
<b>8. Conclusion</b>	<b>14</b>
<b>References</b>	<b>15</b>
<b>A. Appendix</b>	<b>16</b>
A.1. Quell-Code . . . . .	16
A.2. Tipps zum Schreiben Ihrer Abschlussarbeit . . . . .	16

## List of Figures

## List of Tables

3.1. Functional requirements (F0–F8) . . . . .	7
3.2. Non-Functional (N1–N5) requirements . . . . .	8

# Listings

5.1. Ein Beispiel: Hello World (Scala) . . . . .	11
--	----



# 1. Introduction

## 1.1. Background and Motivation

Generative AI is changing the software industry and how software is developed. The emergence of Large Language Models (LLMs) has opened up new possibilities for automating various aspects of software engineering, including code generation, debugging, and program repair. These models have demonstrated remarkable capabilities in understanding and generating code snippets, making them valuable tools for developers. Since debugging and fixing bugs are critical tasks in software development, the integration of LLMs into Automated Program Repair (APR) systems has gained significant attention. Recent studies have shown that LLMs can effectively assist in identifying and fixing bugs, thereby improving the efficiency of the software development process. However, the integration of these practices into existing software development workflows remains understudied.

existing approaches are often complex and require significant computational resources, making them less suitable for budget-constrained environments. Additionally, the lack of integration with Continuous Integration and Continuous Deployment (CI/CD) practices limits their practical applicability in real-world scenarios. The motivation behind this thesis is to explore the potential of LLMs in enhancing automated program repair (APR) systems, particularly in continuous integration and continuous deployment (CI/CD) environments. By leveraging the capabilities of LLMs, we aim to develop a more efficient and cost-effective approach to automated bug fixing that can seamlessly integrate into existing software development workflows.

while writing this paper lots of research and tools have been published clearly showing the importance of the topic and the need for further research in this area. The rapid advancements in LLMs and their applications in software engineering highlight the potential for significant improvements in automated program repair systems. As the software industry continues to evolve, it is crucial to explore innovative solutions that can enhance the efficiency and effectiveness of software development processes.

## 1.2. Problem Statement

modern APR systems require a lot of computational power budget and manual effort - studies suggest emphasizing connections with DevOps Procedures [6]

## 1.3. Objectives and Research Questions

The primary objective of this thesis is to investigate the feasibility and effectiveness of LLMs in APR within the Software Development Lifecycle (CI/CD pipelines) in budget

## *1. Introduction*

restrained environment. The research questions guiding this investigation include:

## 2. Background and Related Work

### 2.1. Software Engineering

Software engineering is a complex discipline consisting of constant

Continuous Integration and Continuous Deployment (CI/CD) is a software development practice that emphasizes frequent integration of code changes into a shared repository, followed by automated testing and deployment. This approach aims to enhance collaboration, reduce integration issues, and accelerate the delivery of high-quality software.

#### 2.1.1. code hosting platforms

github is where open source lives and development takes place. It is a platform that allows developers to host, share, and collaborate on code repositories. GitHub provides version control, issue tracking, and collaboration tools, making it a popular choice for open-source projects and software development teams.

allows for integration of systems like renovate

software development moves more and more towards loosely coupled multi-microservices which makes code repositories smaller but more numerous. This trend is driven by the need for flexibility, scalability, and faster development cycles. Smaller code repositories allow teams to work on specific components or services independently, reducing dependencies and enabling quicker iterations. This approach aligns with modern software development practices, such as microservices architecture and agile methodologies.

Bug fixing is a highly resource intensive task in software engineering, consisting of multiple stages.

APR:

this is where APR comes into play

#### 2.1.2. Software Development Lifecycle

#### 2.1.3. Continuous Integration and Continuous Deployment (CI/CD)

### 2.2. Automated Program Repair

Automated Program Repair (APR) helps developers fix bugs

#### 2.2.1. Evolution of Automated Program Repair

Earliest APR techniques were based on version control history, using the history to roll back to a previous version of the code part, where no issues were present. This

## 2. Background and Related Work

approach, while effective in some cases, often lacked the ability to preserve new features. (more like instant rollback) history based

Initial template based repair, apply predefined transformations to the code based on rules

The emergence of LLM based techniques LLM based APR techniques have demonstrated significant improvements over all other state of the art techniques, benefiting from their coding knowledge [hossainDeepDiveLarge2024]

Agent Based agent based system improve fixing abilities by providing LLMs the ability to interact with the code base and the environment, allowing them to plan their actions [8].

LLMs lay the groundwork of a new APR paradigm [1]

complex agent architectures produce good results especially paired with containerized environments. Emphasis on quality assurance and Devops practices [6]

modern APRs usually consist of multiple stages, including localization, repair, and validation. These stages are often implemented as separate modules, allowing for flexibility and modularity in the repair process [8].

### 2.3. LLMs in Software Engineering

modern large language models have billions of parameters, are pre-trained on massive codebases which results in extraordinary capabilities in this area [1].

problems with LLMs are: Information leakage, hallucinations, and security issues

first LLMs now research is looking into developing and improving workflows leveraging LLMs [6].

problems we are looking into Agents using tools, LLMs + RAG,

### 2.4. LLM-Based Tool Use and CI Context

CI allows seamless integration... this way there is no harmful code executed on own machine, its encapsulated multiple times container in CI runner

### 2.5. Related Work - Existing Systems

end to end without LLMs Sapfix from Facebook. Fixing bugs in production environments lowering incidents mean time of recovery significantly [5]

FixAgent [3]

swe agent [8]

Agentless minimal system [7] claims existing systems are too complex and compute/-costs intensive. lacks the ability to control the decision planning. they use a minimalist approach using localization, repair and validation.

this approach inspired the thesis to test how a simple approach will perform in a real world scenario on a code hosting platform.

–point out areas where current systems see limitation

## *2. Background and Related Work*

– during the research for this thesis, full integrations were published by companies like OpenAI Codex and Github Copilot - but these are not open source

## 3. Requirements

- for this prototype we constructed Requirements which the system shall satisfy - we split into functional requirements: t.3.1 (EXPLAINATION), nonfunctional Requirements combined with security requirements

- these requirements allow for better planning and prioritisation during development. - the satisfaction of all the requirements will allow for evaluation of the integration into the software development lifecycle

### 3.1. Functional Requirements

### 3.2. Non-Functional Requirements

### 3. Requirements

ID	Title		Description	Verification
F0	Multi	Trig- ger	The Pipeline can be triggered: manually, scheudled via cron, or by GitHub issue creation/labeling.	Runs can be found for these triggers
F1	Issue	Gather- ing	Retrieve GitHub repository issues and filter them for correct state and configured labeles BUG.	gate logs list of fetched issues.
F2	Code	Check- out	Fetch the repository code into a fresh workspace and branch (via Docker mount).	After F1, workspace/ contains the correct source files.
F3	Issue	Local- ization	Use LLM to analyze the issue description and identify relevant files.	LLM output contains file paths with files that shall be edited.
F4	Fix	Genera- tion	Use LLM to edit the identified files.	LLM output contains adjusted content for the identified files.
F5	Change	Vali- dation	Run format, lint and relevant tests and capture pass/fail status.	Logs show build and test results.
F6	Iterative	Patch Gener- ation (retry logic)	If F4 reports failures, retry F4–F5 up to X times.	After retries, either F4 passes or fails with no further retries.
F7	Apply	Patch	Commit LLM-generated edits and generate patch.	Git history shows a new commit which is referencing the Github issue
F8	Result	Re- porting	Open a PR or post a comment on GitHub with the diff and summary metrics.	A PR or comment appears for each issue, showing diff and summary.
F9	Logs	and Metrics Col- lection	Provide log files and Metrics with fix-rate, attempt history, timings, token ussage.	A metrics file contains fields: issue, success, timings, stages.

**Table 3.1.:** *Functional requirements (F0–F8)*

### 3. Requirements

ID	Title	Description	Verification
N1	Containerized Execution	All agent code runs in CI runner in a Docker container to isolate it.	Workflow shows Docker container usage
N2	Configurability	User can specify issue labels, branches, attempts, LLM models via YAML.	Changing the config file alters agent behavior accordingly.
N3	Portability	The system can be deployed on any repository on GitHub.	???
N4	Reproducibility	Runs are deterministic given identical repo state and config.	Multiple runs on the same issue report similar metrics.
N5	Operability	The system provides logs and metrics for each run, including execution times, token usage and success rate.	Logs and metrics files are generated after each run, containing all relevant data.

**Table 3.2.:** *Non-Functional (N1–N5) requirements*



## 4. Methodology

The primary objective of this evaluation is to assess the potentials and limitations of our APR pipeline when integrated into a real-world software development lifecycle. We aim to answer the following research question (RQs) to evaluate the system’s capabilities and impact on the software development process: RQ1: Which potentials and limitations does the APR pipeline have ???

### 4.1. Preparation

#### 4.1.1. Dataset Selection

For the evaluation of the APR integration into the software development lifecycle, we selected the Quixbugs dataset [4] as our primary benchmark for testing APR integration. This dataset is well-suited for our purposes due to its focus on small-scale software bugs in Python. It consists of 40 algorithmic bugs each in one file consisting of a single erroneous line, each with a corresponding test for repair validation. Because these bugs were developed as challenging problems for developers [4], we can evaluate if our system can take over the complex fixing of small bugs without developer intervention to prevent context switching for developers.

Compared to other APR benchmarks like SWE-Bench [2] Quixbugs is relatively small which accelerates setup and development.

If achieved I will add SWE-Bench later [2]

#### 4.1.2. Environment Setup

To mirror realistic software development environment, we prepared the Quixbugs dataset by creating a GitHub repository. This repository serves as the basis for the bug fixing process, allowing the system to interact with the codebase and perform repairs. The repository contains only relevant files and folders required for the bug fixing process, ensuring a clean environment for the system to operate in.

We automatically generated a GitHub issue for each bug, using a consistent template that captures just the Title of the Problem. These issues serve as the entry points to our APR pipeline.

#### 4.1.3. Pipeline Architecture

### 4.2. Evaluation

In this section, we describe how we measure the effectiveness and performance of our APR pipeline when integrated into a real-world CI process, using our QuixBugs

#### 4. Methodology

repository as

For Evaluation we will focus on several key metrics to assess the system's performance and abilities in repairing software bugs. These metrics will provide insights into the system's efficiency, reliability, and overall impact on the software development lifecycle. The following metrics will be used:

The following metrics are automatically collected for each run of the APR pipeline:

- **Repair Success Rate**: Calculate the percentage of successfully repaired bugs out of the total number of bugs attempted by using test results.
- **Number of Attempts**: Track the number of attempts made by the system to repair each bug.
- **Overall Execution Time in CI/CD**: Evaluate the time taken for the system to execute within a CI/CD pipeline, providing insights into its performance in real-world development environments.
- **Execution Time of Dockerized Agent**: Measure the time taken by the Containerized agent to execute the repair process, which helps in understanding the efficiency of the containerized environment and help evaluate the computer overhead of CICD.
- **Token Usage**: Monitor the number of tokens used by the LLM during the repair process, which can help in understanding the cost and efficiency of the model's usage.
- **Cost per Issue**: Calculate the cost associated with repairing each bug, considering factors such as resource usage, execution time, and any additional overhead.

what I evaluate: script execution time + CICD overhead one issue vs multiple issue times model vs model metrics costs attempts vs no attempts in all these categories

## 5. Implementation

```
1 object HelloWorld {  
2   def main(args: Array[String]): Unit = {  
3     println("Hello , world!")  
4   }  
5 }
```

**Listing 5.1:** *Ein Beispiel: Hello World (Scala)*

### 5.1. Implementation

Here we break down the implementation of the system into its core components, following the design and methodology outlined in the previous sections. The full code is attached in the ?? appendix.

#### 5.1.1. System Architecture

IMAGE of Figma diagram

#### 5.1.2. System Components

ci pipeline agent core

#### 5.1.3. System Configuration

secrets that need to be added: LLM provider API key, GITHUB TOKEN need to enable github actions permission to create pull requests in repo settings  
explain fields:

---

<sup>0</sup>During development we followed the paradigms introduced in clean code. [martin2008]).

## 6. Results

[Beschreibung der Ergebnisse aus allen voran gegangenen Kapiteln sowie der zuvor generierten Ergebnisartefakte mit Bewertung, wie diese einzuordnen sind]

### 6.1. APR System Results

### 6.2. Evaluation Results

- describe and showcase execution times and repair success rates - show a log?
  - with small models and attempt loop makes small models pass the whole benchmark?

## 7. Discussion

### 7.1. Validity

- quixbugs a small dataset, not representative of real world software development - only python, not representative of real world software development - but shows the potential of applying llm based agents in a real world cicd environment

### 7.2. Potentials

- can take over small tasks in encapsulated environment without intervention
  - small models can solve more problems with retrying with feedback

### 7.3. Limitations

- potential is highly dependant on model - github actions from github have a lot of computational noise - workflow runs on every issue and therefor has some ci minute overhead this could be solved by using a github app which replies on webhook events
  - SECURITY ISSUE: Prompt injection in issue: CICD makes this a bit safer?

### 7.4. Summary of Findings

### 7.5. Lessons Learned

- ai is a fast moving field with a lot of noise

### 7.6. Roadmap for Extensions

- Service Accounts for better and more transparent integration - try out complex agent architectures and compare metrics and results - try out more complex bug fixing tasks - SWE bench

## 8. Conclusion

## References

- [1] Zhi Chen, Wei Ma, and Lingxiao Jiang. *Unveiling Pitfalls: Understanding Why AI-driven Code Agents Fail at GitHub Issue Resolution*. Mar. 2025. DOI: 10.48550/arXiv.2503.12374. arXiv: 2503.12374 [cs]. (Visited on 03/24/2025).
- [2] Carlos E. Jimenez et al. *SWE-bench: Can Language Models Resolve Real-World GitHub Issues?* Nov. 2024. DOI: 10.48550/arXiv.2310.06770. arXiv: 2310.06770 [cs]. (Visited on 03/06/2025).
- [3] Cheryl Lee et al. *A Unified Debugging Approach via LLM-Based Multi-Agent Synergy*. Oct. 2024. DOI: 10.48550/arXiv.2404.17153. arXiv: 2404.17153 [cs]. (Visited on 03/06/2025).
- [4] Derrick Lin et al. “QuixBugs: A Multi-Lingual Program Repair Benchmark Set Based on the Quixey Challenge”. In: *Proceedings Companion of the 2017 ACM SIGPLAN International Conference on Systems, Programming, Languages, and Applications: Software for Humanity*. Vancouver BC Canada: ACM, Oct. 2017, pp. 55–56. ISBN: 978-1-4503-5514-8. DOI: 10.1145/3135932.3135941. (Visited on 04/17/2025).
- [5] Alexandru Marginean et al. “SapFix: Automated End-to-End Repair at Scale”. In: *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. Montreal, QC, Canada: IEEE, May 2019, pp. 269–278. ISBN: 978-1-7281-1760-7. DOI: 10.1109/ICSE-SEIP.2019.00039. (Visited on 03/06/2025).
- [6] Meghana Puvvadi et al. “Coding Agents: A Comprehensive Survey of Automated Bug Fixing Systems and Benchmarks”. In: *2025 IEEE 14th International Conference on Communication Systems and Network Technologies (CSNT)*. Mar. 2025, pp. 680–686. DOI: 10.1109/CSNT64827.2025.10968728. (Visited on 04/27/2025).
- [7] Chunqiu Steven Xia et al. *Agentless: Demystifying LLM-based Software Engineering Agents*. Oct. 2024. DOI: 10.48550/arXiv.2407.01489. arXiv: 2407.01489 [cs]. (Visited on 04/24/2025).
- [8] John Yang et al. *SWE-agent: Agent-Computer Interfaces Enable Automated Software Engineering*. Nov. 2024. DOI: 10.48550/arXiv.2405.15793. arXiv: 2405.15793 [cs]. (Visited on 04/20/2025).

# A. Appendix

## A.1. Quell-Code

## A.2. Tipps zum Schreiben Ihrer Abschlussarbeit

- Achten Sie auf eine neutrale, fachliche Sprache. Keine „Ich“-Form.
- Zitieren Sie zitierfähige und -würdige Quellen (z.B. wissenschaftliche Artikel und Fachbücher; nach Möglichkeit keine Blogs und keinesfalls Wikipedia<sup>1</sup>).
- Zitieren Sie korrekt und homogen.
- Verwenden Sie keine Fußnoten für die Literaturangaben.
- Recherchieren Sie ausführlich den Stand der Wissenschaft und Technik.
- Achten Sie auf die Qualität der Ausarbeitung (z.B. auf Rechtschreibung).
- Informieren Sie sich ggf. vorab darüber, wie man wissenschaftlich arbeitet bzw. schreibt:
  - Mittels Fachliteratur<sup>2</sup>, oder
  - Beim Lernzentrum<sup>3</sup>.
- Nutzen Sie L<sup>A</sup>T<sub>E</sub>X<sup>4</sup>.

---

<sup>1</sup>Wikipedia selbst empfiehlt, von der Zitation von Wikipedia-Inhalten im akademischen Umfeld Abstand zu nehmen [wikipedia2019].

<sup>2</sup>Z.B. [balzert2011], [franck2013]

<sup>3</sup>Weitere Informationen zum Schreibcoaching finden sich hier: <https://www.htw-berlin.de/studium/lernzentrum/studierende/schreibcoaching/>; letzter Zugriff: 13 VI 19.

<sup>4</sup>Kein Support bei Installation, Nutzung und Anpassung allfälliger L<sup>A</sup>T<sub>E</sub>X-Templates!



## Eidesstattliche Versicherung

Hiermit versichere ich an Eides statt durch meine Unterschrift, dass ich die vorstehende Arbeit selbstständig und ohne fremde Hilfe angefertigt und alle Stellen, die ich wörtlich oder annähernd wörtlich aus Veröffentlichungen entnommen habe, als solche kenntlich gemacht habe, mich auch keiner anderen als der angegebenen Literatur oder sonstiger Hilfsmittel bedient habe. Die Arbeit hat in dieser oder ähnlicher Form noch keiner anderen Prüfungsbehörde vorgelegen.

---

Datum, Ort, Unterschrift