



**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 is reshaping software engineering practices by automating more tasks every day, 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. 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 Automated Bug Fixing 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

1. Introduction	1
2. Background and Related Work	2
2.1. Software Engineering	2
2.1.1. Software Development Lifecycle	2
2.1.2. Continuous Integration	3
2.1.3. Software Project Hosting Platforms	4
2.2. Generative Ai in Software Development	4
2.2.1. Generative AI and Large Language Models	4
2.2.2. Generative AI in Software Development	5
2.3. Automated Programm Repair	5
2.3.1. Evolution of Automated Program Repair	6
2.3.2. APR benchmarks	7
2.3.3. APR in CI Context	7
3. Methodology	8
3.1. Preparation	8
3.1.1. Dataset Selection	8
3.1.2. Environment Setup	8
3.1.3. Requirements Specification	9
3.2. Pipeline Implementation	9
3.3. Evaluation	9
4. Requirements	11
4.1. Functional Requirements	11
4.2. Non-Functional Requirements	11
5. Implementation	14
5.1. System Components	14
5.2. System Architecture	16
5.3. System Configuration	16
6. Results	17
6.1. Showcase of workflow	17
6.2. Evaluation Results	17
7. Discussion	18
7.1. Validity	18
7.2. Potentials	18
7.3. Limitations	18

Contents

7.4. Summary of Findings	18
7.5. Lessons Learned	18
7.6. Roadmap for Extensions	18
8. Conclusion	19
References	20
A. Appendix	23
A.1. Quell-Code	23
A.2. Tipps zum Schreiben Ihrer Abschlussarbeit	23

List of Figures

2.1. Agile Software Development Lifecycle	2
2.2. Continuous Integration Cycle	3
2.3. Simple Action	4

List of Tables

4.1. Functional requirements (F0–F8)	12
4.2. Non-Functional (N1–N5) requirements	13

Listings

1. Introduction

Generative AI is rapidly changing the software industry and how software is developed and maintained. The emergence of Large Language Models (LLMs), a subfield of Generative AI, has opened up new opportunities for enhancing and automating various domains of the software development lifecycle. Due to remarkable capabilities in understanding and generating code snippets, LLMs have become valuable tools for developers' everyday tasks such as requirement engineering, code generation, refactoring, and program repair [1, 2].

Despite these advances, bug fixing remains a resource-intensive task and is often perceived negatively [3]. It leads to frequent interruptions and context switching, which reduce developer productivity [4]. Bugs have direct impact on software quality by causing crashes, vulnerabilities or even data loss. [5] The process of bug fixing can be time-consuming and error-prone, leading to delays in software delivery and increased costs. In fact, according to CISQ: in 2022 alone poor software quality costs 2.41 trillion dollars only the US with at least 607 billion dollars spend on finding and fixing bugs [6].

Given the critical role of debugging and bug fixing in software development, Automated Program Repair (APR) has gained significant research interest. Typically, bug fixing involves multiple steps: bug reporting, localization, repair, and validation [7, 8, 9, 10, 11]. Recent research has shown that LLMs can effectively be used to enhance automated bug fixing, thereby introducing new standards in the APR world showing potential of making significant improvements in efficiency of the software development process [9, 12, 13, 14, 15, 16].

However existing APR approaches are often complex and require significant computational resources [17], making them less suitable for budget-constrained environments or individual developers. Additionally, the lack of integration with existing software development lifecycles and workflows limits their practical applicability in real-world development environments [18, 12].

Motivated by these challenges, this thesis explores the potential of integrating LLM based automated bug fixing within continuous integration and continuous deployment (CI/CD) pipelines. By leveraging the capabilities of LLMs, we aim to develop a cost-effective prototype for automated bug fixing that seamlessly integrates into existing software development workflows. Considering computational demands, complexity of integration and practical constraints we aim to provide insights into possibilities and limitations of our approach.

2. Background and Related Work

In this section we will provide an overview of the relevant background and context for this thesis. First introducing the software engineering lifecycle and the rising role of GenAI/LLMs in it. The Second part showcases the evolution and state of APR and explores existing approaches.

2.1. Software Engineering

In the following section introduces the software engineering lifecycle, the role of code hosting platforms, and the importance of Continuous Integration and Continuous Deployment (CI/CD) in modern software development.

2.1.1. Software Development Lifecycle

Engineering Software is complex and including multiple stages. For structuring this work different Software Development Lifecycle Models have been introduced. Software Development Lifecycle Models evolve constantly to adapt to the changing needs of creating software. The most promising and widely used model is the Agile Software Development Lifecycle [19].

The Agile lifecycle brings an iterative approach to development, focusing on collaboration, feedback and adaptivity. The Goal frequent delivery of small functional features of software, allowing for continuous improvement and adaptation to changing requirements. Agile can be used with multiple frameworks like Scrum or Kanban but follows a similar approach. [19].

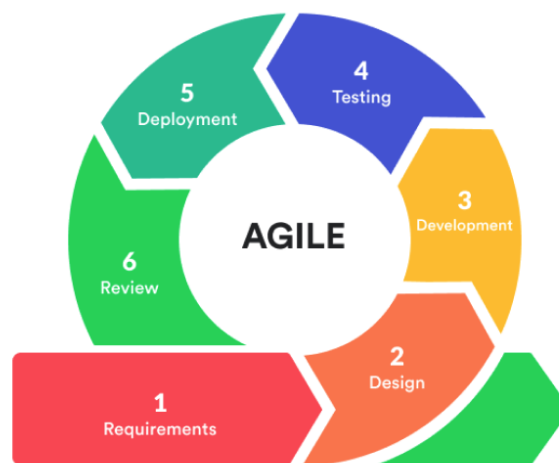


Figure 2.1.: *Agile Software Development Lifecycle*

2. Background and Related Work

A Agile Software Development Lifecycle iteration consists of several key stages like in Figure starting with planning phase where requirements for the iteration are gathered and prioritized.

Since agile focuses adaptivity arising bugs can alter iterations if prioritised and therefore slow down delivery of features. APR is supposed to help with this problem by accelerating the process of fixing bugs. —explain how bugs are handled

Software development is moving towards lightly coupled microservices which results in more repositories which are smaller in scale tailored towards a specialized domain. 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. With this trend developers work on multiple projects at the same time, which can lead to more interruptions and context switching when problems arise and priorities shift.

2.1.2. Continuous Integration

For accelerating the delivery of software in an iteration continuous integration has become a standard in agile software development.

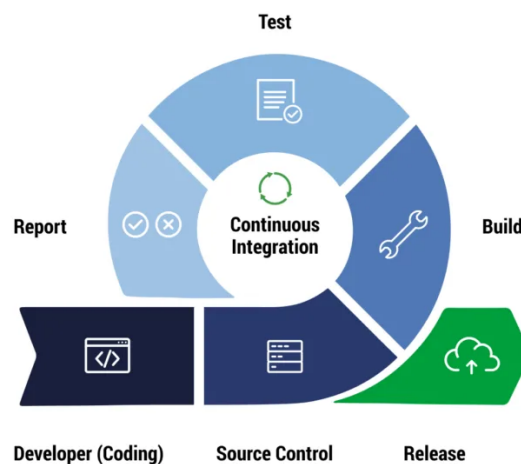


Figure 2.2.: *Continuous Integration Cycle*

Continuous Integration (CI) allows for frequent code integration into a code repository. CI can integrate steps like automated building and testing into the development resulting in rapid feedback right where the changes committed to the shared repository.

Although CI brings a lot of potential to development it can also have problems which can be long build durations and high maintenance.

CI supports aspects like fast delivery, fast feedback, enhanced collaboration which are critical for agile software development. [20]

2. Background and Related Work

2.1.3. Software Project Hosting Platforms

Software projects are hosted on platforms like Github or gitlab. With Github being the most popular and most used [**<empty citation>**] These platforms provides tools and feature for the complete software development lifecycle. Project hosting, verssion control, bug and issue tracking, project management, backups, collaboration, and documentation. [21]

GitHub has features like Issue tracking for requiremntns and planning with issues looking like this: [**githubdoc**] —image of issue has title, description, comments, labels and mroe assoicated information

also provides a manged soltution for integrating CI into reposiries by writing work-flows in YAML files called Github Actions. The pipeliens can run on github hosted runners or self hosted runner. A workflow can be triggered by one or more event. One or more jobscan be executed on a provided runner maschiene. A job can consist of multiple steps. [22]

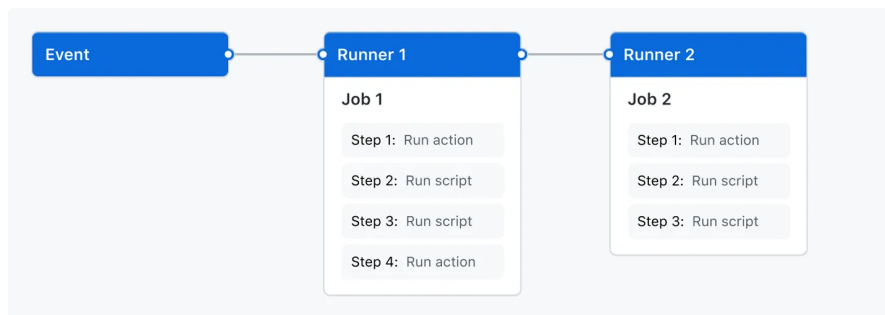


Figure 2.3.: Simple Action

Since constant feedback and reviews of code play an imporant role in agile workflows github also provides a pull request feature. Pull Requests allow for proposal of changes to the codebase, with an integrated review process which allows for collaboration and review before changes are integrated into the production codebase. Code changes are displayed in a diff format allowing reviews to see and dig into the changes made. This process is essential for maintaining code quality and ensuring that changes are validated before being merged. [**githubdoc**]

2.2. Generative Ai in Software Development

2.2.1. Generative AI and Large Language Models

Gen ai is subfield of Ai

LLms work using human langague data

relatively new advancement are AI Agents which let LLMs interact with the environ-ment and plan their actions

2.2.2. Generative AI in Software Development

Generative AI is reshaping software development by automating various tasks. Modern large language models have billions of parameters, are pre-trained on massive codebases which results in extraordinary capabilities in this area [18]. Tools like Github Copilot, OpenAI Codex, and ChatGPT have become popular in the software development community, providing developers with AI-powered code suggestions and completions for different tasks [23]. These tools get applied in various stages of the software development lifecycle, including requirement engineering, code generation, debugging, refactoring, and testing [bargavmallampatiRoleGenerativeAI2025, 1, 2]. LLMs can generate code in different programming languages. With python being the most supported [empty citation] By using LLMs to enhance these tasks development cycle times can be reduced by up to 30 percent [23, 24]. Furthermore these tools have positive impacts like improving developer satisfaction and reducing cognitive load [24].

Although Generative AI gets adopted really quickly in many areas of software development this transition still faces challenges and limitations. LLMs have issues working on tasks that are outside their scope of training or require specific domain knowledge [1]. Additionally LLMs have limited context windows, which can lead to challenges when working with large codebases or complex projects where context windows are too small for true contextual or requirements understanding [23]. When generating code LLMs can produce incorrect or insecure code, which can lead to further bugs and vulnerabilities in the software [1]. Generating content or code with LLMs based on training data raises questions about ownership, reproducibility and intellectual property rights. When it comes to security generating code using prompts is vulnerable to prompt injection, where unintended instructions are injected at some point and can lead to production of harmful code [empty citation]

[23, 1]. When generating code based on training data intellectual property issues arise [25]

recently research and companies are looking into developing solutions which integrate LLMs into existing software development practices and workflows [IntroducingCodex, 2, 26, 25].

recently more attention is given to integration of AI/ML into CI/CD [27]

2.3. Automated Program Repair

Automated Program Repair (APR) is software that helps detect and repair bugs in code with minimal human intervention. This field has also benefited of the rapid advancements in AI.

APR systems are supposed to take over the process of fixing bugs therefore making more time for developers to focus on more relevant work. [1]

Specific bugs can be fixed using a resulting patch from an APR system. For creating working patches APR takes a 3 stage approach: First localizing the bug. Then repairing the bug, in the end validation decides where the bug will be passed on.[empty citation]

2. Background and Related Work

In this section we will provide an overview of the evolution of APR, related work, and the current state of APR systems.

2.3.1. Evolution of Automated Program Repair

We have seen multiple transformations in the field of Automated Program Repair (APR) over the years. This evolution of APR can be categorized into key stages, each marked by significant advancements in techniques and methodologies.

Traditional Approaches:

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

Template based systems relied on predefined template for transformations of commonly known bugs. Templates applied predefined transformations to the code based on fixed rules. This Approach had limited the flexibility and adaptability in a quickly transforming software landscape. [2]

Search based repair,

Semantic based repair,

One of the most outstanding system is Getafix develop and deployed at Meta [28]

Nevertheless traditional systems face significant limitations in scalability and adaptability, struggling to generalize to new scenarios or unseen bugs, or to adapt to evolving codebases. They often required extensive computational resources or manual effort. [2]

The emerge of llm based APR: LLM based APR techniques have demonstrated significant improvements over all other state of the art techniques, benefiting from their coding knowledge [29]. For that reason LLMs lay the groundwork of a new APR paradigm [18]. Common LLMs used for APR include GPT-4, ChatGPT, Codex, CodeLlama, DeepSeek-Coder, and CodeT5 [1, 30, 31].

Using LLMs different Paradigms have emerged and are being actively researched. These paradigms include:

Retrieval-Augmented Approaches repair bugs with the help of retrieving relevant context during the repair process. This approach allows adding external knowledge to the repair process, enhancing the LLM's ability to understand and fix bugs [1, 30].

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. These frameworks reconstruct the cognitive processes using multiple Agents that can generate code with the help of multi-step reasoning, usage of tools for with Environments and Tools [empty citation]. Examples for that are SWE-Agent [yangSWEAgentComputerInterfaces2024], FixAgent [8], MarsCodeAgent [12], GitHub Copilot.

complex agent architectures produce good results especially paired with containerized environments. [2]

Interactive approaches make use of LLMs dialogue capabilities to equip patch validation with instant feedback. This feedback is used to refine the generated patches trying to achieve better results. This process takes a more dynamic and iterative approach. [15]

2. Background and Related Work

Agentless systems are recent a push towards more lightweight solution, focusing on simplicity and efficiency. These approaches aim to reduce the complexity of APR systems while maintaining effectiveness in bug fixing [9]. Furthermore this approach provides clear rails to the LLMS improving the transparency of the bug fixing approach taken. These Systems have achieved promising results with low costs [9]

Common problems currently faced by state of the art APR system are: Existing system are overly complex with limited transparency and control over the bug fixing process.[puvvadiCodingAgentsComprehensive202, 9, 1] Repairing bugs takes a lot of computational resources and is time intensive therefore producing significant costs [14, 2] Repairing bugs is done on benchmarks or in controlled set up environments and not integrated into real world software development workflows [32, 2]

2.3.2. APR benchmarks

Popular are Quixbugs small bugs in python , Defects4J for java programs and the hardest: SWE Bench based on real world github issues in python repositories

2.3.3. APR in 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

3. Methodology

The primary objective of this thesis 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 to evaluate the system’s capabilities and impact on the software development process:

What are the potentials and limitations of integrating an LLM-based automated bug-fixing pipeline into a CI/CD workflow, as measured by repair success rate, end-to-end execution time, and API cost, and how does this integration impact the overall software development lifecycle?

For awnsering this question we streamlines this process into three phases Preparation, Implementation / Application and Evaluation.

3.1. Preparation

3.1.1. Dataset Selection

For the evaluation of the APR integration into the software development lifecycle, we selected the Quixbugs dataset [33] 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 algortithmic bugs each in one file consisting of a single erroinums line , each with a correspodng tests for repair validation. Because these bugs where developed as challending problems for developers [33], 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 [34] Quixbugs is relaivly small which accelerates setup and development.

if archieived I will ad swe bench lite later [34]

3.1.2. Environment Setup

To mirrow 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.

3.1.3. Requirements Specification

3.2. Pipeline Implementation

The Automated Bug Fixing Pipeline was developed using iterative prototyping and testing, with a focus on simplicity and modularity. Starting with constructing the functional and non functional requirements ?? we build the following System:

—IMAGE of HIGH LEVEL System Architecture here

When the system is in place in a target repository. An issue with the default or configured bug label can be created. This trigger trigger will spin up a github action runner which executes the APR pipeline. This Pipeline talks to the configured LLM Api (google and openai) to localize and fix the buggy files. With the supplied file edits the code is validated and tested. When validation passes the a pull request with the fiels changes is autoamatically open on the repository, linking the issue and providing details about the repair process. In case of a unsucessfull repair the failure is reported to the issue.

3.3. 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 repository as a bases.

For Evaluation we will focus on several key metrics to assess the system's performance and abilites 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 are automatically collected and calculated 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. A sucessfull repair is defined as a bug passasing all tests associated with it.

Number of Attempts: Track the number of attempts made by the system to repair each bug with a maximum of 3 attempts.

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.

TODO should I split overall execution and stages into seperate metrics? **Execution Times of Dockerized Agent:** Measure the time taken by the Containerized agent and its stages to execute the repair process and the execution times of individual stages.

With this CICD overhead can be calculated and bottlenecks can be identified

Token Usage: Monitor the number of tokens used by the LLM during the repair process, which can help understanding the cost of repairing and issue and the relation between token usage and repair success.

Cost per Issue: Calculate the cost associated with repairing each bug, considering factors such as resource usage, execution time, and any additional overhead.

3. Methodology

explain how these metrics are collected and calculated, including any tools or scripts used to automate the process. explain how reapiir sucess rate is determined

ask zhang wether i need to include the evaluation of swe bench lite from the paper or if a ref is enough

explain statistical methods used to analyze the collected data, such as averages, medians, and standard deviations. This will help in understanding the variability and reliability of the results.

explain how resuLTS are calcuated for model comparisson

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

4. 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

4.1. Functional Requirements

4.2. Non-Functional Requirements

4. 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 usage.	A metrics file contains fields: issue, success, timings, stages.

Table 4.1.: *Functional requirements (F0–F8)*

4. 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 4.2.: *Non-Functional (N1–N5) requirements*

5. Implementation

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

the goal was to create a system which not only fixes bugs but is also portable /deployable across different repositories and configurable to some extent.

System Overview:

The system Consists of two main components: The APR core which hold the core logic for the repair process

The CI/CD Pipeline which put everything together and integrates the github entry-point of the target repository with the core logic of the agent.

Configuration Layer the programs behaviour can be altered by adjusting a configuration. A default YAML configuration is in place which allows for controlling: - labels - workdir - branches - Attempts - LLM models addiotnally these is the possiblity to over-write these values for a single repository using a dedicated 'bugfix.yml' configratuon which needs to be placed at the root of the repository.

5.1. System Components

Agent Core:

The agent core is written in python and dockerized so its slim and portable. the agent core contains the main bugfixing logic. To start fixing bugs it needs the following envrioment: The repo where to fix files on - provided using docker volumne mount the follwing Envrionment Variables: - Github Token - Github Repo - LLM API key - ISSUE TO PROCESS - Github repository

With this envriomnment set the system loops over all issues which are fethced from the envrionment varaibles.

For each issue the main APR logic is executed. This main logic consists of tools and stages: —IMAGE here First the worksapce (a new branch based on configraution naming) and a repair context is set up. the context is hing of the program its needed at every step and works as the main data structure for the APR system like memory. — JSON of Context init here: A stage uses the context to perform a specific task in the bug fixing process and returns the context with its added context. The stages are: Localize, Fix, Build, Test

With a repaired workspace the repairprocess start with the localization stage. This stage construct a prompt for the LLM to localize the bug in the codebase. This prompt makes use of a contrsucted hierachy of the repositories fielstrcutre and the issue description. The LLM is expected to return a list of files and lines where the bug is located. — PROMPT

5. Implementation

The results are stored in the context.

With the the locilized files in the context the fix stages constructs a prompt with the file content and the issue description. the llm can return code or no changes needed. — Prompt these edits are then applied to the files in the workspace. The context is updated with the new file content.

To ensure the changes are properly formatted the the build stages formats the code using the black formatter and lints the python code to ensure maintainable code.

Next up the test stage runs tests for the fixes files and attaches these results to the context.

if tests do not pass the system will record a new attempt and start over from the state of the fix stage. Additionally a feedback is generated using the previous attempts code and results from the context which gets added to the fix prompt.

if the validation passes or the maximum number of attempts is reached the system will either report and unsucessfull repair to the issue or continue to the application steps

on sucessfull repair the follwing steps are executed: the fileschanges are committed and a diff file is generated. the changes are pushed to the remote repo usign the github token

a pull request is opened with a description of the changes and a link to the issue and more details about tests and some metrics like the number of attempts and the tokens used for the repair process.

During execution the core logs its actions, which can be used for debugging and analysis. Furthermore it collects metrics such as the number of attempts, execution time, and token usage, which are essential for evaluating the performance of the APR system. Logs and metrics are saved as .log and json files at the end

The agent core is designed to be modular and extensible, allowing for future enhancements and additional stages or tools to be integrated as needed. It is also designed to be lightweight, ensuring that it can run efficiently within the constraints of a CI/CD environment.

CI/CD Pipeline:

The Pipeline is written in YAML acoording to the Github CICD standard. ?? We will use runners hosted by Github which takes awayt the overhead of managing our own runners but comes at the cost of unknown performance and avaiablity It is made up of the Triggers and 3 Jobs: 'gate', 'skipped' and 'bugfix'. The triggers are set up first and will allow the execution of the APR process. Triggers can resutls in two types of runs: processing of all bug issues in correct state on manual execution request of the workflow ("workflow_dispatch") or scheudled execution ("cron"). processing a single issue: when an issue is openend and label with the configured labels ("issue_opened and issue_labeled ") or when extra information is added or edited on an issue in form of a comment.

The trigger event information gets passed as envrionment varaibles to the next job "gate" which is responsible for evaluating if the issue should be processed or skipped. This job checks the labels and resolves the issue state to determine if the issue is relevant for the APR process. If no issues pass this gate the job "skipped" is executed, which simply logs that no issues were found to process and exits the workflow.

5. Implementation

When the gate outputs issues that should be processed the job “bugfix” is executed. This job checks out the current repository and mounts it as a volume to the agent core container. Additionally it sets the necessary environment variables mentioned at ?? . For the agent to work permissions are set on the job level to allow the agent to edit repository content, create pull requests and write issues.

For giving access to the agent cores logs and metrics the job provides the logs directory as an artifact which is available after the workflow run is completed.

for all of this to run the following environment secrets need to be defined in the repo:

5.2. System Architecture

IMAGE of Figma diagram

5.3. 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:

The full implementation is listed in Appendix ??

6. Results

In the following section we will showcase the resulting workflow of our prototype and the evaluation results for the Quixbugs benchmark.

6.1. Showcase of workflow

To integrate the APR system into a repository living on GitHub we need to move the pipeline with its filter script to the dedicated github action workflow directory.

—IMAGE OF WORKFLOW IN PLACE

When the workflow is in place the APR system is ready to go. Optionally its default behavior can be altered by adding a configuration file (called: bugfix.yml) to the root of the repository. —EXAMPLE OF CONFIG

Now when an issue is created and labeled with the default label "bug" (or a custom label defined in the configuration file) the APR system will be triggered and start the bug fixing process. Manual triggering is also possible by using the "Run workflow" button in the GitHub actions tab of the repository. —IMAGE OF ISSUE BEING CREATED OR MANUALLY TRIGGERED

After the workflow is triggered and relevant issues are found the APR system start as a run of the GitHub action workflow. —IMAGE OF RUN BEING STARTED

When the automatic bug fixing process has completed there are two possible outcomes: Pull Request with patch for bug and link to the issue. —IMAGE OF PULL REQUEST

or a comment on the issue that bug fixing failed after all attempts. —IMAGE OF COMMENT ON ISSUE

After the Workflow finishes metrics and logs are available for download in the action. (during a run logs are live streamed in the Workflow run) —IMAGE OF LOGS

6.2. Evaluation Results

comparing different LLMS Models: google OpenAI

Single Issue Processing: CI overhead

Multi Issue Processing: CI overhead

Also include attempt loop

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 - this concept is applicable to other python repositories
 - ?? - accelerate bug fixing - lets developers focus on more complex tasks - therefore enhance software reliability and maintainability

7.3. Limitations

- github actions from github have a lot of computational noise - workflow runs on every issue and therefore 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?
 - its limit to small issues

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] Xinyi Hou et al. *Large Language Models for Software Engineering: A Systematic Literature Review*. Apr. 2024. DOI: 10.48550/arXiv.2308.10620. arXiv: 2308.10620 [cs]. (Visited on 03/06/2025).
- [2] 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).
- [3] Emily Winter et al. "How Do Developers Really Feel About Bug Fixing? Directions for Automatic Program Repair". In: *IEEE Transactions on Software Engineering* 49.4 (Apr. 2023), pp. 1823–1841. ISSN: 1939-3520. DOI: 10.1109/TSE.2022.3194188. (Visited on 06/24/2025).
- [4] Bogdan Vasilescu et al. "The Sky Is Not the Limit: Multitasking across GitHub Projects". In: *Proceedings of the 38th International Conference on Software Engineering*. Austin Texas: ACM, May 2016, pp. 994–1005. ISBN: 978-1-4503-3900-1. DOI: 10.1145/2884781.2884875. (Visited on 06/24/2025).
- [5] Norbert Tihanyi et al. *A New Era in Software Security: Towards Self-Healing Software via Large Language Models and Formal Verification*. June 2024. DOI: 10.48550/arXiv.2305.14752. arXiv: 2305.14752 [cs]. (Visited on 06/26/2025).
- [6] *Cost of Poor Software Quality in the U.S.: A 2022 Report*. (Visited on 06/24/2025).
- [7] Feng Zhang et al. "An Empirical Study on Factors Impacting Bug Fixing Time". In: *2012 19th Working Conference on Reverse Engineering*. Oct. 2012, pp. 225–234. DOI: 10.1109/WCRE.2012.32. (Visited on 06/24/2025).
- [8] 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).
- [9] 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).
- [10] Yuwei Zhang et al. "PATCH: Empowering Large Language Model with Programmer-Intent Guidance and Collaborative-Behavior Simulation for Automatic Bug Fixing". In: *ACM Transactions on Software Engineering and Methodology* (Feb. 2025), p. 3718739. ISSN: 1049-331X, 1557-7392. DOI: 10.1145/3718739. (Visited on 03/24/2025).
- [11] Jialin Wang and Zhihua Duan. *Empirical Research on Utilizing LLM-based Agents for Automated Bug Fixing via LangGraph*. Jan. 2025. DOI: 10.33774/coe-2025-jbpg6. (Visited on 03/12/2025).

References

- [12] Yizhou Liu et al. *MarsCode Agent: AI-native Automated Bug Fixing*. Sept. 2024. DOI: 10.48550/arXiv.2409.00899. arXiv: 2409.00899 [cs]. (Visited on 03/06/2025).
- [13] 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).
- [14] Dominik Sobania et al. "An Analysis of the Automatic Bug Fixing Performance of ChatGPT". In: *2023 IEEE/ACM International Workshop on Automated Program Repair (APR)*. May 2023, pp. 23–30. DOI: 10.1109/APR59189.2023.00012. (Visited on 03/06/2025).
- [15] Chunqiu Steven Xia and Lingming Zhang. "Automated Program Repair via Conversation: Fixing 162 out of 337 Bugs for \$0.42 Each Using ChatGPT". In: *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. Vienna Austria: ACM, Sept. 2024, pp. 819–831. ISBN: 979-8-4007-0612-7. DOI: 10.1145/3650212.3680323. (Visited on 05/12/2025).
- [16] Haichuan Hu et al. *Can GPT-O1 Kill All Bugs? An Evaluation of GPT-Family LLMs on QuixBugs*. Dec. 2024. DOI: 10.48550/arXiv.2409.10033. arXiv: 2409.10033 [cs]. (Visited on 04/15/2025).
- [17] Pat Rondon et al. *Evaluating Agent-based Program Repair at Google*. Jan. 2025. DOI: 10.48550/arXiv.2501.07531. arXiv: 2501.07531 [cs]. (Visited on 03/24/2025).
- [18] 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).
- [19] Nayan B. Ruparelia. "Software Development Lifecycle Models". In: *ACM SIGSOFT Software Engineering Notes* 35.3 (May 2010), pp. 8–13. ISSN: 0163-5948. DOI: 10.1145/1764810.1764814. (Visited on 06/25/2025).
- [20] Vincent Ugwueze and Joseph Chukwunweike. "Continuous Integration and Deployment Strategies for Streamlined DevOps in Software Engineering and Application Delivery". In: *International Journal of Computer Applications Technology and Research* (Jan. 2024), pp. 1–24. DOI: 10.7753/IJCATR1401.1001.
- [21] Pekka Abrahamsson et al. *Agile Software Development Methods: Review and Analysis*. Sept. 2017. DOI: 10.48550/arXiv.1709.08439. arXiv: 1709.08439 [cs]. (Visited on 06/25/2025).
- [22] *About Workflows*. https://docs-internal.github.com/_next/data/9uQSGns-DWbCy3Cy8blUA/en/freepro-team%40latest/actions/concepts/workflows-and-actions/about-workflows.json?versionId=free-pro-team%40latest&productId=actions&restPage=concepts&restPage=workflows-and-actions&restPage=about-workflows. (Visited on 06/25/2025).
- [23] Bhargav Mallampati. "The Role of Generative AI in Software Development: Will It Replace Developers?" In: *World Journal of Advanced Research and Reviews* 26.1 (Apr. 2025), pp. 2972–2977. ISSN: 25819615. DOI: 10.30574/wjarr.2025.26.1.1387. (Visited on 06/25/2025).
- [24] Eirini Kalliamvakou. *Research: Quantifying GitHub Copilot's Impact on Developer Productivity and Happiness*. Sept. 2022. (Visited on 06/26/2025).

References

- [25] Jaakko Sauvola et al. “Future of Software Development with Generative AI”. In: *Automated Software Engineering* 31.1 (May 2024), p. 26. ISSN: 0928-8910, 1573-7535. DOI: 10.1007/s10515-024-00426-z. (Visited on 06/25/2025).
- [26] Thomas Dohmke. *GitHub Copilot: Meet the New Coding Agent*. May 2025. (Visited on 06/26/2025).
- [27] Abdul Sajid Mohammed et al. “AI-Driven Continuous Integration and Continuous Deployment in Software Engineering”. In: *2024 2nd International Conference on Disruptive Technologies (ICDT)*. Mar. 2024, pp. 531–536. DOI: 10.1109/ICDT61202.2024.10489475. (Visited on 04/23/2025).
- [28] Johannes Bader et al. “Getafix: Learning to Fix Bugs Automatically”. In: *Proceedings of the ACM on Programming Languages* 3.OOPSLA (Oct. 2019), pp. 1–27. ISSN: 2475-1421. DOI: 10.1145/3360585. (Visited on 03/06/2025).
- [29] Soneya Binta Hossain et al. “A Deep Dive into Large Language Models for Automated Bug Localization and Repair”. In: *Proceedings of the ACM on Software Engineering* 1.FSE (July 2024), pp. 1471–1493. ISSN: 2994-970X. DOI: 10.1145/3660773. (Visited on 03/13/2025).
- [30] Xin Yin et al. “ThinkRepair: Self-Directed Automated Program Repair”. In: *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. Vienna Austria: ACM, Sept. 2024, pp. 1274–1286. ISBN: 979-8-4007-0612-7. DOI: 10.1145/3650212.3680359. (Visited on 03/12/2025).
- [31] Avinash Anand et al. *A Comprehensive Survey of AI-Driven Advancements and Techniques in Automated Program Repair and Code Generation*. Nov. 2024. DOI: 10.48550/arXiv.2411.07586. arXiv: 2411.07586 [cs]. (Visited on 06/01/2025).
- [32] Fairuz Nower Meem, Justin Smith, and Brittany Johnson. “Exploring Experiences with Automated Program Repair in Practice”. In: *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. Lisbon Portugal: ACM, Apr. 2024, pp. 1–11. ISBN: 979-8-4007-0217-4. DOI: 10.1145/3597503.3639182. (Visited on 06/26/2025).
- [33] 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).
- [34] 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).

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 Wikipedia¹).
- 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 Fachliteratur², oder
 - Beim Lernzentrum³.
- Nutzen Sie L^AT_EX⁴.

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

²Z.B. [balzert2011], [franck2013]

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

⁴Kein Support bei Installation, Nutzung und Anpassung allfälliger L^AT_EX-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