

Quality Assurance of LLM-generated Code: Addressing Non-Functional Quality Characteristics

Xin Sun*, Daniel Ståhl, Kristian Sandahl, Christoph Kessler

Department of Computer and Information Science, Linköping University, Sweden

Abstract

In recent years, LLMs have been widely integrated into software engineering workflows, supporting tasks like code generation. However, while these models often generate functionally correct outputs, we still lack a systematic understanding and evaluation of their non-functional qualities. Existing studies focus mainly on whether generated code passes the tests rather than whether it passes with quality. Guided by the ISO/IEC 25010 quality model, this study conducted three complementary investigations: a systematic review of 108 papers, two industry workshops with practitioners from multiple organizations, and an empirical analysis of patching real-world software issues using three LLMs. Motivated by insights from both the literature and practitioners, the empirical study examined the quality of generated patches on security, maintainability, and performance efficiency. Across the literature, we found that security and performance efficiency dominate academic attention, while maintainability and other qualities are understudied. In contrast, industry experts prioritize maintainability and readability, warning that generated code may accelerate the accumulation of technical debt. In our evaluation of functionally correct patches generated by three LLMs, improvements in one quality dimension often come at the cost of others. Runtime and memory results further show high variance across models and optimization strategies. Overall, our findings reveal a mismatch between academic focus, industry priorities, and model performance, highlighting the urgent need to integrate quality assurance mechanisms into LLM code gen-

*Corresponding author.

Email addresses: xin.sun@liu.se (Xin Sun), daniel.stahl@liu.se (Daniel Ståhl), kristian.sandahl@liu.se (Kristian Sandahl), christoph.kessler@liu.se (Christoph Kessler)

eration pipelines to ensure that future generated code not only *passes tests* but truly *passes with quality*.

Keywords: LLM-generated code, code generation, non-functional quality characteristics, maintainability, security, performance efficiency

1. Introduction

The prevalence of large language models for code (LLMs for code) has brought new changes to the field of software engineering (Chen et al., 2021)(Rozière et al., 2023). These models are tailored for understanding and processing code and can generate functionally correct outputs under appropriate prompts (Luo et al., 2024)(Hui et al., 2024)(Jiang et al., 2023). Currently, these powerful models have been integrated into the workflow of software engineering and are used by millions of developers worldwide. GitHub Copilot, released by GitHub and OpenAI, is an “AI pair programme” that is built on Codex, which is a family of code LLMs (Chen et al., 2021)(GitHub, 2025). Copilot can generate code across multiple programming languages from various types of prompts, such as natural language descriptions, function signatures, and surrounding code. In some cases, it can produce complete applications such as interactive websites or data pipelines from a single prompt. As these models continue to evolve rapidly, their potential and implications for software engineering remain in flux (The AI Digest, 2024).

Several carefully curated metrics and benchmarks have been developed to evaluate the generation capability of code LLMs Chen et al. (2021)Jimenez et al. (2024). Chen et al. (2021) released HumanEval, together with GitHub Copilot, to evaluate the ability of code LLMs to generate functionally correct code. SWE-bench (Jimenez et al., 2024) was introduced to assess the ability of LLMs to solve real-world software engineering tasks.

While these benchmarks provide a good basis for evaluating the ability of LLMs to generate functionally correct code, they offer limited support for assessing the non-functional quality characteristics (NFQCs) of generated code. According to the ISO/IEC 25010 quality model (International Organization for Standardization, 2023), software quality encompasses not only functional correctness but also NFQCs such as performance efficiency, maintainability, and security.

Due to the limited assessments of NFQCs in the current evaluation system, code LLMs may produce code with varying quality along different

NFQCs, which makes the generated code¹ unreliable and sometimes causes severe faults. For example, Pearce et al. (2022) studied the performance of GitHub Copilot in high-risk cybersecurity scenarios MITRE (2024). In their study of 1,689 programs generated across 89 scenarios, approximately 40% were found to contain code patterns considered vulnerable, illustrating the security weakness in LLM-generated code.

While prior work has begun to explore NFQCs in LLM-generated code, the field still lacks a study that not only synthesizes existing research but also incorporates insights from practitioners and experimental evidence of real-world software systems to identify key challenges and outline future directions.

To address this gap, we present a comprehensive study that investigates the NFQCs of LLM-generated code from multiple perspectives. The overall goal of this study is to understand challenges and practices in NFQCs for LLM-generated code, explore the existing issues in this field, and facilitate the integration of LLMs in the software development workflow. Guided by the ISO/IEC 25010 software quality model, the study adopted mixed methods to examine the NFQCs of generated code. We first conducted a systematic literature review to synthesize existing studies and highlight current challenges and future directions. To complement the literature findings, we organized two workshops with industry practitioners to capture real-world perspectives on quality assurance of LLM-generated code, with a particular focus on NFQCs. Finally, we performed an empirical study to examine the performance of generated patches along three NFQCs (maintainability, security, and performance efficiency), which were consistently identified as most important by both the literature and practitioners. Together, these three components provide a consolidated overview of the state of the field and offer insights for integrating NFQCs into the quality assurance process to improve the reliability of LLM-generated code in real-world software systems.

To guide our study, we defined three overarching **Research Goals (RGs)**, each addressing a different aspect. The research goals are as follows:

- **RG1:** Explore how NFQCs in generated code have been addressed in existing research, and identify research gaps in this area.

¹In this study, **generated code** refers to source code generated by LLMs, including both code snippets and code patches.

- **RG2:** Identify the real-world expectations and pain points that practitioners have when integrating generated code in software projects.
- **RG3:** Characterize the empirical behavior of generated code across key NFQCs identified in the previous research goals, and examine how feedback-driven strategies affect these qualities and the trade-offs among them in a real-world software engineering context.

For each research goal, we further formulated corresponding **Research Questions (RQs)** to guide the analysis, which are presented in Section 3.

The rest of the paper is organized as follows: Section 2 introduces the background on code generation and related work. Section 3 details our methodology, including the systematic literature review, the design of workshops, and the experimental setup. Section 4 presents our results, combining insights from the literature, industry discussion, and empirical evaluation. Additionally, we discuss the implications of our findings, highlighting challenges and opportunities for future research. We discuss the threats to validity in Section 5. Finally, Section 6 concludes the study and outlines directions for integrating NFQCs into the quality assurance of LLM-generated code in real-world software systems.

2. Background and Related Work

This section presents the background of LLMs for code generation and related works.

2.1. Code generation

Code generation has become a more important application scenario since LLMs were first proposed Manna and Waldinger (1971) Liventsev et al. (2023). Code generation aims to automatically generate code given the problem description, high-level specifications, or existing code. Before the advent of LLMs, code generation was limited to structured methods Hemel et al. (2008). These methods usually require detailed formal specifications and several steps for generation. LLMs changed this paradigm by learning from large code corpora, thus eliminating the need for formal specifications and allowing code generation from natural language prompts Luo et al. (2024).

The release of GitHub Copilot marked the beginning of integrating LLMs into daily software workflows. Since then, numerous code LLMs have emerged, demonstrating remarkable performance in generating functionally correct

and meaningful code across various programming languages Rozière et al. (2023)Hui et al. (2024)Jiang et al. (2023)Chen et al. (2021)Li et al. (2022)Li et al. (2023).

By 2025, many AI programming tools have emerged, and a new approach called *vibe coding* has gained prominence (Lovable, 2024)(Bolt, 2024)(Cursor, 2024)(Windsurf, 2024). In vibe coding, users do not directly engage with the code; instead, they remain entirely at the level of natural language, formulating prompts that instruct the code LLMs to generate, modify, and deploy components of the software. For example, Lovable² is an AI-driven platform that enables users to build full-stack applications using natural language prompts, allowing rapid prototyping and deployment without extensive coding knowledge. These advancements have significantly lowered the barrier to programming and accelerated software development processes.

2.2. Literature Reviews

A review by Wang and Chen (2023) examined LLM code generation from a different angle. By reviewing 20 existing studies, they demonstrated that while using LLMs to generate code has been widely studied, the evaluation of LLM-generated code had received relatively little attention. The study addressed functional requirements such as functional correctness, as well as NFQCs, including security, maintainability, and others. It also highlighted key limitations in current evaluation practices, such as the use of inadequate evaluation criteria, the absence of systematic and quantitative evaluation frameworks, and the lack of consideration for human involvement in the evaluation process.

Yang et al. (2024) presented a systematic literature review of 146 studies focusing on NFQCs of code LLMs. They identified and discussed six NFQCs beyond functional correctness: *robustness, security, privacy, explainability, efficiency* and *usability*. The study highlighted the vulnerability of LLM4Code systems to adversarial attacks, data poisoning, and privacy leaks, as well as challenges in explainability and usability. To address these issues, they proposed three complementary perspectives, which are data-centric, human-centric, and system-centric, for developing more reliable and effective LLM4Code systems in the future. This work provided a comprehensive overview for understanding the broader implications of adopting LLM4Code

²<https://lovable.dev/>

in software engineering. However, most of the studies included in Yang et al. (2024)'s literature review focus on evaluating the models themselves, such as the architecture, training data, or prompting strategies, rather than the quality of the code generated by LLMs. In practice, a model capable of producing high-quality code in principle does not guarantee that the actual outputs will meet high standards across NFQCs. It is therefore essential to acknowledge that, in real-world practice, systematic evaluation of generated code is required to ensure the overall quality of the software system. This motivates our study, which differs from prior work by shifting the perspective from the model to its outputs, and by examining to what extent the generated code satisfies NFQCs.

2.3. Empirical Studies

Beyond conceptual discussions and literature reviews, a growing body of studies has examined the NFQCs of LLM-generated code. Niu et al. (2024) conducted a comprehensive performance study using HumanEval and MBPP benchmarks and a set of programming problems from LeetCode, an online programming practice platform. They also investigated generating efficient code using prompt engineering methods. Their findings noted that the performance can depend heavily on how a prompt is given. If not instructed, LLMs might output a straightforward solution that ensures correctness but is not efficiency optimal.

Fu et al. (2025) found that approximately 27.3% of generated code contained security weaknesses, spanning 43 Common Weakness Enumeration categories, including critical issues like Insufficiently Random Values and Cross-site Scripting. The study also explored the effectiveness of Copilot in fixing these vulnerabilities, demonstrating that enhanced prompts could resolve up to 55.5% of issues.

In the area of maintainability, Liu et al. (2024a) applied static analysis tools such as Pylint to assess generated code and reported frequent code smells and issues, suggesting the need for human oversight.

Taken together, prior work reveals a growing effort towards understanding the NFQCs of generated code. Our study extends that effort by integrating evidence from research, practice, and empirical evaluation to provide a holistic view of NFQCs in a real-world context.

3. Methodology

This section describes the methodological framework and overall structure of the study. It is organized into three main parts. Subsection 3.1 outlines the structure of the systematic literature review. Subsection 3.2 presents the organization of the two industry workshops, detailing the participant selection, session design, and qualitative analysis approach. Subsection 3.3 summarizes the setup of the empirical evaluation, covering dataset selection, experiment configuration, evaluation metrics, and analysis workflow.

3.1. Literature Review

To be able to achieve RG1 in Section 1, we first conducted a literature review to gain enough understanding of the current state of research and potential issues in this field. Furthermore, the literature review sought to answer the following research questions:

- **RQ1.1:** What trends can be observed in research on NFQCs of generated code?
- **RQ1.2:** How have different NFQCs been examined in the existing studies, and what challenges and limitations have been identified?

Research Design. LLMs have been utilized across various domains. But in this literature review, we focus especially on code LLMs, which are LLMs tailored for code generation. Our review followed key elements of the PRISMA methodology (Page et al., 2021), including a structured *search strategy*, clearly defined *inclusion and exclusion criteria*, and the use of *snowballing* to identify additional relevant literature. We included papers that examine code LLMs and AI programming tools or platforms that integrate these models, such as GitHub Copilot. To be eligible, the selected papers must explore the NFQCs of the code generated by these tools or models. To facilitate consistent discussion and comparison, we map the diverse quality characteristics mentioned across studies to their corresponding terms in the ISO/IEC 25010 software product quality model, thereby standardizing terminology for our analysis.

Search Strategy. Currently, the popularity of code LLMs has caused a surge of studies in the research community. However, the specific focus on the quality of generated code, particularly with respect to NFQCs, is still relatively recent. In addition to peer-reviewed studies, many relevant and insightful preprints have emerged on platforms like arXiv.³ These preprints also provide valuable perspectives and contribute meaningfully to the ongoing research. Accordingly, our review also considers preprints in this field.

We conducted the search on **Scopus**, chosen for its extensive coverage of peer-reviewed publications and inclusion of some preprints. Our search query was formulated in two stages. First, we identified three core concepts: (1) LLMs, (2) code generation, and (3) software quality. Next, we broadened each concept using closely related terms to maximize recall while maintaining precision. For example, for *LLMs*, we included terms like ‘AI’ and ‘artificial intelligence’; for *quality*, we included ‘performance’ and ‘nonfunctional’.

- LLM; large language model; AI; artificial intelligence
- code generation
- quality; performance; non-functional

To ensure a high coverage of relevant papers, we deliberately avoided overly specific keywords related to individual NFQCs (*e.g.*, "robust", "efficient", "latency") in the keywords list, as these terms can vary widely across papers and might appear in different forms. Finally, we set our search query as follows:

```
(“LLM” OR “large language models” OR “AI” OR “Artificial intelligence” ) AND ( “code generation” ) AND ( “quality” OR “performance” OR “non-functional” )
```

The initial results contained many irrelevant studies, so we applied a series of filters to refine the search. Our inclusion and exclusion criteria are summarized in Table 1. Specifically, we excluded papers published before 2022, as the majority of studies in code generation prior to that year were

³<https://arxiv.org/>

based on traditional approaches rather than LLMs, which gained major traction with the release of ChatGPT in late 2022. We also restricted the subject area to *computer science* and *Engineering*, as these two categories in Scopus most accurately capture research related to LLMs and code generation. Additionally, we included peer-reviewed journal and conference papers, as well as preprints that were publicly available and relevant to the research topic, all written in English.

Table 1: The selection criteria of our literature review.

Inclusion criteria
<ul style="list-style-type: none">• The paper discusses one of the NFQCs or similar quality characteristics of the AI-generated code.• The paper proposes an approach, study or benchmark to evaluate one of the quality attributes of AI-generated code.• The paper proposes new metrics to evaluate one of the quality attributes of AI-generated code.
Exclusion criteria
<ul style="list-style-type: none">• Paper that is not written in English.• Paper that is not within the subject of Computer Science.• Paper that appeared before 2022.

Applying these filters, we obtained 471 publications and 776 preprints. We then performed manual screening based on the title, abstract, and keywords of each paper to determine whether it was relevant to the evaluation of LLM-generated code with respect to NFQCs. Each study was further assigned to one or more categories defined by the ISO/IEC 25010 standard. In addition, we applied forward and backward snowballing to identify further relevant literature that might not have been captured by the initial query, and this added three more studies in the dataset. Following this process, the final dataset comprises 106 papers, of which 38 are preprints. The structure of the applied search and selection process is shown in Figure 1.

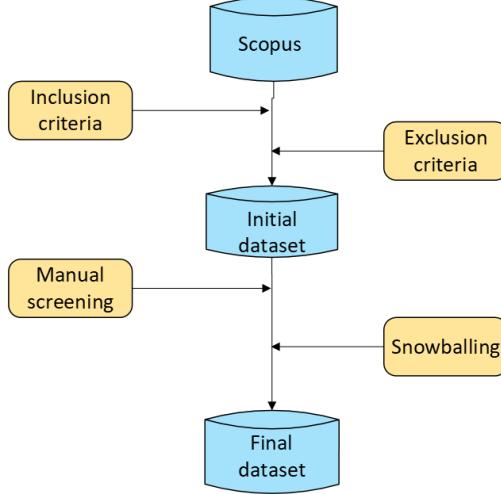


Figure 1: Search and selection process

3.2. Workshop Design and Execution

To complement and validate the findings of our literature review, we held two interactive workshops with industry experts from several organizations. These organizations are actively exploring the integration of LLMs into large-scale software systems and seeking to ensure the reliability of LLMs' outputs. The workshops are designed to answer the following research questions:

- **RQ2.1:** What NFQCs are considered most important by industry practitioners when evaluating generated code?
- **RQ2.2:** What challenges and risks do practitioners perceive when integrating generated code into existing development workflows?

Participants. We invited several groups of industry experts to participate in the workshops, with a total of 15 participants, including three of the authors. The participants are from seven different organizations in various industry areas, and most of the organizations are members of Software Center⁴ (SWC), which is an industry-academic collaborative network to accelerate industrial digitization and the adoption of novel approaches to software engineering.

⁴<https://www.software-center.se/>

The participants have a wide range of professional skills, including quality-in-design, system integration and testing, CI/CD architecture, and automotive software test-driven development. Their professional experience in software engineering ranges from 3 to over 13 years, ensuring both breadth and depth of knowledge.

Workshop Procedure. The workshops were held remotely through Microsoft Teams and divided into two sessions to accommodate different groups and foster focused discussion. In the first session, we met experts from a non-SWC organization, and later in the second session, we met SWC organization participants. During the workshop, we intentionally made the discussion open-ended, and questions were welcomed at any time to encourage spontaneous and candid input. Both sessions followed the same overall structure:

- **Presentation of our findings:** Each session began with a presentation of our literature review, highlighting the recent trends and the key findings related to NFQCs in LLM-generated code (see details in Section 4). The presentation acted as a shared foundation for discussion, providing participants with a common frame of reference for the subsequent workshop activities.
- **Participant sharing:** After our presentation, the participants were invited to present their ongoing projects, observed challenges, and useful practices related to this topic. This segment enabled different organizations (in the second workshop) to communicate with each other on the practical implementations and establish the synchronization between industry and academia.
- **Open discussion:** The discussion phase happened not only during the presentations, but also after the presentations. We did not follow a predetermined set of questions. Instead, it was facilitated in an exploratory manner, allowing participants to raise issues most relevant to their contexts and to reflect on both the presentation and their practical experience.

Data collection and analysis. During the two sessions, two of the authors independently took detailed notes during each session. After each workshop, the presentation slides were compiled, and all authors had meetings to consolidate the individual notes into a comprehensive workshop transcript. In

cases where discrepancies arose between the two notes, these were discussed and resolved through consensus.

For data analysis, we used the definitions from ISO/IEC 25010 to systematically examine and categorize the workshop results. Unclear cases were discussed until the authors reached an agreement. Finally, the categorized data were synthesized to identify patterns and to highlight which aspects of NFQCs were emphasized during the workshops.

3.3. Experimental Design

According to the findings from our literature review and workshops, we identified three NFQCs that are of primary concern in both academia and industry: maintainability, security, and performance efficiency. For performance efficiency, we further considered two measurable dimensions: runtime and memory.

Based on these insights and RG3, we designed an experiment to evaluate the NFQC performance of generated code and to answer the following questions:

- **RQ3.1:** How do LLMs perform with respect to NFQCs when generating code for real-world software engineering tasks?
- **RQ3.2:** Can incorporating static analysis feedback, such as results from CodeQL, into prompts improve the NFQC performance of generated code?
- **RQ3.3:** Does improving one NFQC result in trade-offs with other NFQCs?

To address these research questions, we designed a three-stage experimental pipeline, as summarized in Figure 2. In the **first stage (Baseline Evaluation)**, we generated patches for a benchmark dataset of real-world software issues Jimenez et al. (2024) and evaluated both the generated and the benchmark gold patches. This stage established baseline measurements of functional correctness, runtime, memory usage, and static analysis results related to security and maintainability. In the **second stage (Filter and Prompt Design)**, we identified a set of comparable, correctly resolved instances across models and used the corresponding baseline results to design targeted prompts for NFQC-specific optimization. In the **final stage (NFQC-specific Regeneration and Evaluation)**, we regenerated

the patches using the new prompts, evaluated them under the same conditions, and compared the results against the baseline to analyze improvements and trade-offs among different NFQCs.

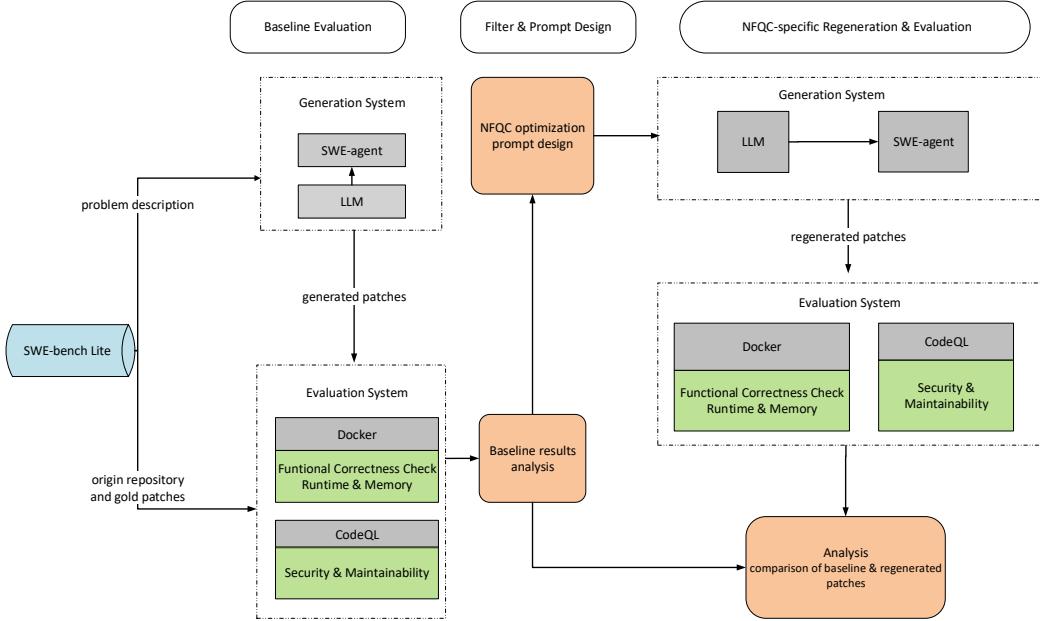


Figure 2: Overview of the experiment procedure. The workflow is divided into three main stages: (1) **Baseline Evaluation**, where initial patches were generated using SWE-agent and different LLMs, generated and benchmark gold patches are evaluated with Docker and CodeQL; (2) **Filter and Prompt Design**, where baseline results are analyzed and NFQC-specific prompts are constructed; and (3) **NFQC-specific Regeneration and Evaluation**, where patches are regenerated using NFQC-specific prompts, re-evaluated and compared against the baseline to analyze improvements and potential trade-offs.

Dataset. We selected *SWE-bench Lite* Jimenez et al. (2024), which is a subset of SWE-bench, as the evaluation benchmark for our experiment. SWE-bench Lite is a carefully curated subset of this benchmark, created to support faster iteration and model development. In our experiment, we used the *test* set of SWE-bench Lite, which consists of 300 issue-pull request pairs from 11 Python repositories. Each instance is evaluated by executing unit tests, using the post-pull request behavior as the reference.

Using SWE-bench Lite allows us to evaluate model performance on real-world software tasks with manageable resources. Although the benchmark

SWE-bench Lite Instance Structure

```
instance_id: Unique identifier of the instance.  
repo: GitHub repository name.  
base_commit: Commit hash before applying the fix.  
problem_statement: Title and description of the issue.  
hints_text: Optional hints or comments from the issue.  
test_patch: Unit test modifications introduced by the PR.  
patch: Reference solution patch.  
FAIL_TO_PASS: Tests that failed before and passed after.  
PASS_TO_PASS: Tests that passed both before and after.
```

Figure 3: SWE-bench Lite Instance Structure

itself focuses on functional correctness, its automated patch verification system provides a reliable foundation on which we build our NFQC evaluation pipeline. This enables us to analyze NFQCs of functionally correct patches without modifying the benchmark itself. Figure 3 shows the structure of each instance in SWE-bench Lite.

Model Selection. While LLMs deliver promising performance in generating basic functions and algorithms, their performance on the SWE-bench family benchmarks remains limited. As of August 20, 2025, Claude 4 Opus achieves a resolved rate of 67.60% on SWE-bench, and Claude 4 Sonnet reports 66.93% (submission on May 21, 2025) and 56.67% (submission on May 26, 2025) on SWE-bench and SWE-bench Lite, respectively.⁵

For model selection, we reviewed all the submissions on the SWE-bench Lite leaderboard and aimed to minimize confounding factors unrelated to the models themselves. To this end, we decided to use the SWE-agent framework as the generation pipeline. SWE-agent is an agent-based system designed to automate the generation process in LLM code generation. It enables LLMs to use tools, interact with the environment, and produce patches step by step, without human instruction. The use of SWE-agent also ensures that the generation process is fully driven by the LLM, thereby reducing potential biases introduced by manual intervention. This design makes comparisons

⁵<https://github.com/SWE-bench/experiments>

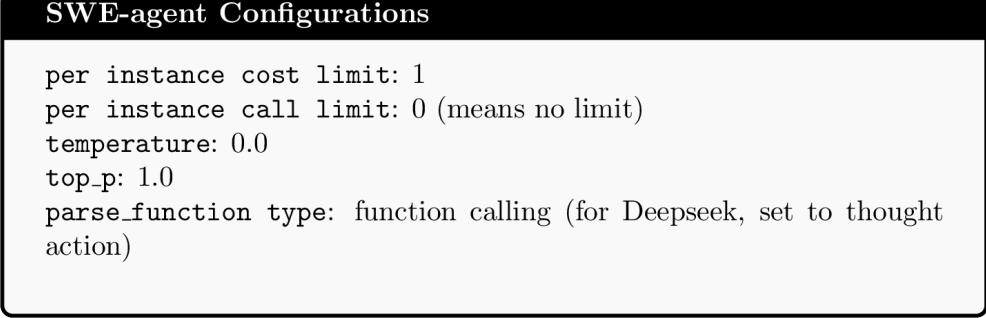
across models more fair and enhances the reliability of the evaluation. To examine how models differ in their NFQCs improvement potential, we selected three models that combine representativeness and contrastiveness:

- Claude 4 Sonnet: A mid-size model in the Claude 4 family. It demonstrates strong coding and reasoning abilities with enhanced instruction following and extended context support.
- DeepSeek-Reasoner: A reasoning-augmented model that generates chain-of-thoughts to improve the generation accuracy. It shows competitive performance on reasoning and code-related tasks compared to state-of-the-art closed-source models.
- GPT-4o: OpenAI’s LLM, released in 2024. It achieves performance on par with GPT-4 Turbo in code generation, while offering faster inference and lower API costs.

These models cover closed-source and open-source models, providing a more representative and generalizable comparison for evaluating NFQC in generated code.

Generation Strategy. In our experiment, we employed SWE-agent as the generation framework. Each model has its identical configuration parameters, and interacts with the environment through command-line tool use. The agent supports various stopping criteria and cost control strategies. In the first stage, we set a cost limit of \$1 per instance for each model, without restricting the number of retries. This configuration balanced computational cost with the ability to maximize resolved instances, which provides more baseline resolved patches for subsequent analysis. In the final stage, we limited the number of retries to 1, where the agent produced one patch per instance before moving to the next instance. This multi-round generation design more closely mirrors how humans interact with LLMs in practice Zhong et al. (2024). The configuration is shown in Figure 4.

Evaluation System. The evaluation system was designed to assess the generated patches across three NFQC dimensions: **performance efficiency** (including runtime and memory usage), **security**, and **Maintainability**, along with functional correctness as a prerequisite. Thus, we measured four metrics in the experiment. It consisted of two main components:



```

SWE-agent Configurations

per instance cost limit: 1
per instance call limit: 0 (means no limit)
temperature: 0.0
top_p: 1.0
parse_function type: function calling (for Deepseek, set to thought
action)

```

Figure 4: SWE-agent Configurations

The first component focused on functional verification and performance measurement. For each instance, we used Docker to create an isolated environment, cloned the corresponding repository, applied the generated patch, and executed the test cases provided by SWE-bench Lite. This step verified whether the patch was functionally correct and resolved the given issue. For patches that resolved the given issues, we further measured the total execution time required to run the entire test suite and the peak memory usage during the execution.

The second component performed static code analysis to evaluate security and maintainability. We employed CodeQL to perform holistic repository analysis rather than incremental scans, ensuring comprehensive detection of potential security and maintainability issues. We applied the official Python security and quality query suites released by GitHub for static analysis of source code. These suites are part of the standard CodeQL query packs and are publicly available in the CodeQL documentation. Each query in these suites is labeled with tags such as maintainability, correctness, security and performance. The description and related information were extracted from the official document, and we categorized the queries based on the tag of each query. In Table 2, Table 3, we list the most frequently triggered maintainability and security rules in our experiment.

The relative improvements of the metrics were calculated as:

$$\Delta\text{metric} = \frac{\text{metric}_{\text{before}} - \text{metric}_{\text{after}}}{\text{metric}_{\text{before}}} \quad (1)$$

For the static analysis results, we further derived a composite risk score to quantify potential quality risks in security and maintainability:

$$\text{risk_score} = \sum W_{\text{severity}} \times W_{\text{precision}} \times \text{trigger_count} \quad (2)$$

where `trigger_count` denotes the number of times a specific rule is triggered. The severity levels indicate the potential impact of the detected issues. Specifically, *Error* denotes critical flaws that may cause failures or security vulnerabilities. *Warning* denotes problematic or non-robust patterns, and *Recommendation* suggests minor improvements for better maintainability or readability. Based on the description of severity level, we assigned weights of 1.0, 0.6 and 0.2 to Error, Warning, and Recommendation, respectively. The $W_{\text{precision}}$ adjusts for the reliability of each rule in reporting true positives. We assigned 1.0, 0.7, and 0.4 to Very-high, High, and Medium. A higher risk score indicates greater potential quality risk introduced by the generated patch.

Experimental Setup. Our experiment was conducted on a laptop with an AMD Ryzen 7 CPU and 16GB RAM. The operating system was Ubuntu 22.04 with Python 3.11. We used SWE-agent v1.1.0, the latest release available as of May 2025, along with SWE-bench.

The details of the three stages are as follows:

- **Baseline evaluation:** We first established baseline measurements for the original repositories and SWE-bench Lite gold patches across the selected NFQCs using the evaluation system described earlier. These measurements established the baseline results for subsequent comparison. After that, we generated the initial patches using the SWE-bench Lite benchmark and the SWE-agent framework. SWE-agent functioned as a unified framework that coordinated the patch generation process across different LLMs. Since the official SWE-bench Lite leaderboard already provides submissions using SWE-agent with Claude 4 Sonnet and GPT-4o, we reused the same configurations to reproduce the results. For DeepSeek-Reasoner, although no official submission was available, we used the same configuration to maintain fairness across models. Thus, in this stage, the configurations were identical across experiments, with the only difference being the LLM used.

After patch generation, the SWE-agent output the generated patches corresponding to the instances in the dataset. These patches were then evaluated through the evaluation system, performing the functional correctness check and the NFQC evaluation. An instance was

considered *resolved* if its generated patch successfully fixed the issue described in the prompt. For all resolved instances, we additionally measured the runtime and peak memory usage during the functional correctness check. The output from this stage included both the generated patches and the evaluation results of the original repositories with and without the generated patches applied.

- **Filter and prompt design:** Next, we identified the resolved instances for each model and derived their intersection, meaning we retained only those instances that were successfully resolved by all three models. This ensured that subsequent analyses compared the same issues under identical functional conditions. For these common instances, we analyzed the corresponding NFQC evaluation results from CodeQL and recorded the number and type of newly introduced issues. CodeQL rules were grouped into two categories: security and maintainability, and weighted by severity and precision levels according to Equation 2. Using these results, we designed NFQC optimization prompts to approximate how developers might provide targeted feedback to an LLM in practice. Security and maintainability prompts incorporated explicit issue descriptions extracted from CodeQL reports, while prompts for runtime and memory optimization specified the desired performance focus directly. The structure of the NFQC optimization prompts is shown in Listing 1.

Listing 1: NFQC-specific optimization prompt

```

1 PROMPT_TMPL_BASE = """Your previous patch is functionally correct.  

2     You must re-implement its behavior directly on the base commit,  

3     while focusing on improving **{focus}**. Always generate a  

4     unified diff against the ORIGINAL BASE COMMIT.  

5  

6 <issue>  

7 {problem_statement}  

8 </issue>  

9  

10 {prev_patch_block}  

11 {feedback_block}  

12 REQUIREMENTS:  

13 - Output **ONLY a unified diff** (git patch) with file paths  

14     relative to the repo root.  

15 - Do NOT include explanations, markdown fences, numbering, or shell  

16     commands.

```

```

12 - Keep changes minimal; do NOT modify tests.
13 {requirement_line}
14 - If a trade-off arises, prioritize functional correctness.
15 - Re-implement the behavior of the previous patch on the base commit
    , preserving functionality.
16
17 FORMAT:
18 - Start with lines like: diff --git a/<path> b/<path>
19 - Use standard unified diff hunks: @@ -old,+new @@
20 - Ensure the patch applies cleanly.
21
22 Now produce ONLY the patch:
23 """

```

- **NFQC-specific regeneration and evaluation:** In this stage, for each model and each instance in the intersection set, we applied the designed optimization prompts separately, producing four regenerated patches per instance per model. Every regenerated patch was evaluated with the same pipeline as in first stage: functional correctness in isolation, followed by measurement of total runtime and peak memory if resolved, and a holistic CodeQL scan to compute security and maintainability scores. We then compared these outcomes against the corresponding unoptimized results for the same instance and model to quantify the improvements and assess trade-offs across NFQCs. (As in the baseline, only patches that resolve the original issue were included in comparisons.)

4. Results and Discussion

In this section, we present the main results of our study, together with the discussion. The section is structured in three parts, aligned with the methodology: (1) literature review, (2) insights from industry workshops, and (3) our experiments. In each part, we provide both the results and corresponding discussion.

4.1. Literature Review Findings

Following the search strategy described in Subsection 3.1, 106 studies were finally included, of which 38 were preprints. This subsection presents the main trends and research focuses observed across these studies, as well as how individual NFQCs have been investigated and evaluated in prior work.

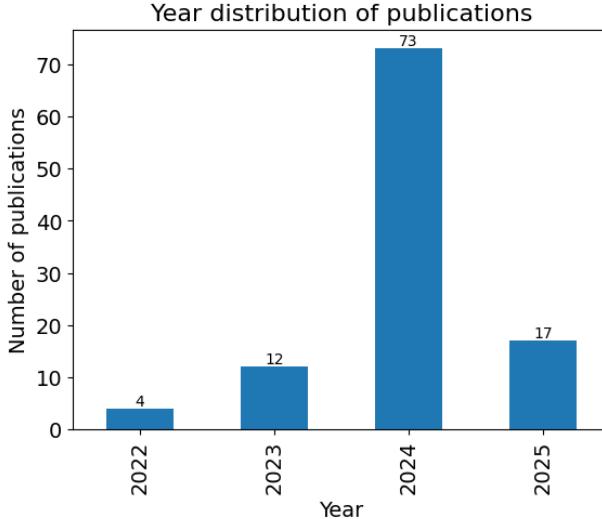


Figure 5: The year distribution of papers.

As shown in Figure 5, the publication trend shows a rapid increase in studies on NFQCs of generated code between 2022 and 2024. While only four papers appeared in 2022 and 12 in 2023, the number surged to 73 in 2024. This growth indicates that the research community has begun to recognize the importance of NFQCs in generated code, particularly after the release of LLM-based programming tools such as GitHub Copilot.

The word cloud in Figure 6 illustrates the thematic focus of the reviewed literature. Among all the reviewed papers, **security**, **performance efficiency**, **Maintainability**, and **readability** dominate the discourse, while concepts such as **reliability** and **robustness** appear less frequently. Also, we observed that different studies often use varying terminology to refer to the same quality attributes defined in ISO/IEC 25010. For example, we found that in studies addressing performance efficiency, terms such as *runtime*, *green code*, *memory management* were frequently used as proxies for performance efficiency. *Readability* and *understandability* are often used to denote aspects of maintainability.

As shown in Figure 7, the distribution of studies across ISO/IEC 25010 categories reveals that only four NFQCs were addressed in the literature. Among them, **security** attracted the most attention (41%), followed by **performance efficiency** (31.2%). **Maintainability** accounted for 20% of reviewed studies, whereas **reliability** received the least attention (7.5%).

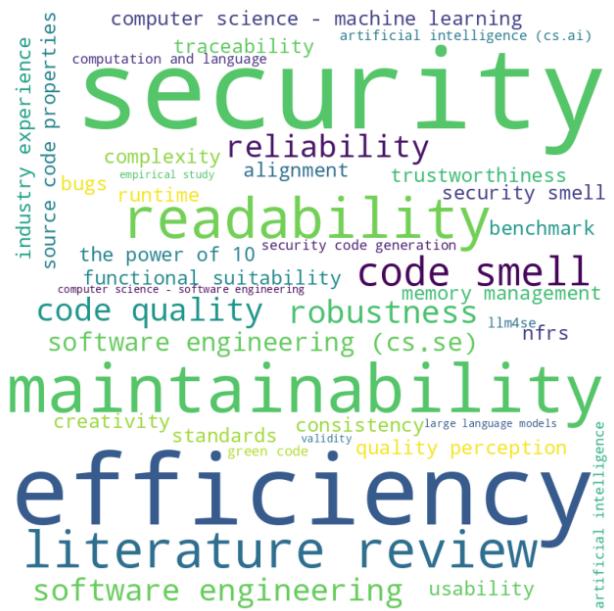


Figure 6: Word cloud of the identified literature.

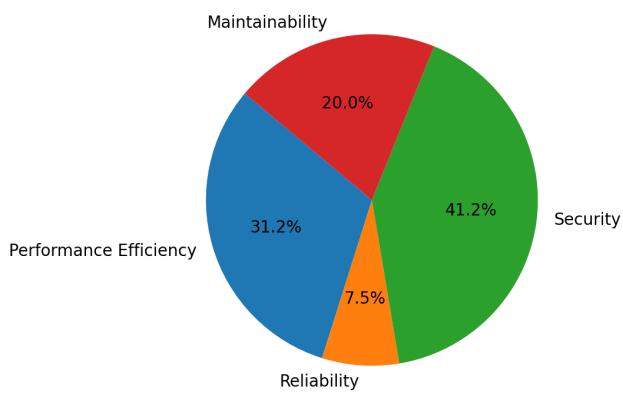


Figure 7: The distribution of papers by NFQCs.

These findings demonstrate both the concentration of research efforts on certain NFQCs and the relative neglect of others, suggesting the research gaps for future investigation.

4.2. Current research state

Security. Security is a major concern in LLM-generated code, as insecure code can lead to vulnerabilities in software systems and sometimes cause severe issues. Security always refers to the ability of the system to defend against attack patterns and protect information and data from inappropriate access International Organization for Standardization (2023).

A growing body of studies has observed that LLMs sometimes suggest code that is functionally correct but has security issues Perry et al. (2023)Khoury et al. (2023). Thus, many studies specifically evaluate whether the generated code contains security vulnerabilities or poor security practices Pearce et al. (2025)Elgedawy et al. (2024)Siddiq and Santos (2022)Asare et al. (2023). Pearce *et al.* Pearce et al. (2025) conducted an empirical study to investigate the security vulnerabilities in code generated by GitHub Copilot. They evaluated the generated code across three dimensions: diversity of weaknesses, diversity of prompts, and diversity of domains. Their analysis of 1,689 generated code reveals that approximately 40% of the generated code contains security vulnerabilities, with some high-risk weaknesses appearing frequently. Later, Siddiq and Santos Siddiq and Santos (2022) introduced SecurityEval, a dataset to evaluate the ability of LLMs to generate vulnerability-free code.

Among the studies, static code analysis is a common method for security evaluation. Tools like Bandit or CodeQL can detect common security issues and vulnerabilities in code. For example, Pearce *et al.* Pearce et al. (2025) used CodeQL to detect security vulnerabilities in the generated code. In Liu et al. (2024b), Liu *et al.* utilized CodeQL for vulnerability detection on the code snippets generated by ChatGPT. They found that 33% of the valid generated code snippets are vulnerable, and the majority of vulnerable code snippets are related to the MissingNullTest, which may lead to runtime exceptions or unexpected program crashes. Also, Siddiq *et al.* Siddiq et al. (2024) investigated the quality issues of ChatGPT-generated code using the DevGPT dataset, which is a curated dataset of developer-ChatGPT conversations encompassing prompts with ChatGPT’s responses, including code snippets. They found that the generated code contains security smells, such as hard-coded credentials.

Another dimension is comparing AI-generated code against humans in security. Asare *et al.* Asare et al. (2023) posed the question: “Is GitHub’s Copilot as bad as humans at introducing vulnerabilities in code?” They found that Copilot replicates the original vulnerable code about 33% of the time and replicates the fixed code at a rate of 25%. They concluded that Copilot is not as bad as humans at introducing vulnerabilities, but they also pointed out that it is risky to use Copilot to generate code and fix security bugs. Licorish *et al.* Licorish et al. (2025) evaluated the code written by GPT-4 and humans using several static analysis tools. They observed security flaws in code generated by both, and found that human-written code tends to include a great variety of problems, while GPT-4 code has more severe issues.

We also see works focusing on evaluating and improving LLMs’ performance on generating security code. Yetistiren *et al.* Yetistiren et al. (2023) assessed the security of code generated by several AI programming tools using SonarQube on the HumanEval benchmark, and found all tools generated code containing code smells. To support systematic evaluation in this area, several benchmarks were introduced, including SecurityEval Siddiq and Santos (2022), LLMSECEval Tony et al. (2023), and the AI Code Generators for Security benchmark Natella et al. (2024).

Empirical studies consistently report that LLM-generated code often contains security issues, even when the code is functionally correct. These findings suggest that LLMs do not reliably incorporate good coding practices by default.

Performance Efficiency. Performance efficiency refers to whether the code is efficient in using resources (such as CPU, memory, storage, and network) and performing its functions within a specific time International Organization for Standardization (2023). Traditional software development often includes algorithmic complexity, execution time, memory usage, and even energy consumption.

Among the studies, most of the work focus on measuring the runtime performance and memory usage of code produced by LLMs Niu et al. (2024) Singhal et al. (2024) Coignion et al. (2024a) Huang et al. (2024) Waghjale et al. (2024). Coignion *et al.* Coignion et al. (2024a) evaluated the efficiency of the code generated by 18 LLMs and measured the performance against human-crafted solutions using problems from the LeetCode platform, they found that on average, generated solutions were more efficient than human-written code.

Shypula et al. Shypula et al. (2024) introduced a framework for adapting LLMs for learning performance-improving code edits. With the carefully curated datasets and the proposed adaptation strategies, their enhanced model achieved a mean speedup of $6.86\times$, significantly higher than the average human programmer’s optimization.

Memory usage is also reported by several studies. Huang et al. (2024) measured the memory consumption of the generated code using their proposed benchmark. Their results show that sometimes the generated code consumes more memory due to using simpler but memory-hungry data structures.

Beyond runtime and memory, the energy efficiency of generated code has been a niche focus Vartziotis et al. (2024)Solovyeva et al. (2025). Vartziotis et al. (2024) evaluated the energy consumption of the generated code together with the run time and memory usage. Preliminary observations suggest that LLMs do not inherently optimize for energy use unless explicitly instructed.

In summary, the performance efficiency evaluation of generated code includes runtime, memory usage, and sometimes energy consumption, and often uses LeetCode problems as the dataset. The generated solutions, when correct, often have performance on par with human-crafted optimal solutions. With specific prompting techniques or frameworks, LLMs can generate high-performance code, in some cases outperforming human solutions. Also, the intersections with other qualities are observed, which is an issue that needs consideration.

Maintainability. Maintainability refers to how easily the code can be understood, modified by developers in the future International Organization for Standardization (2023). Generally, maintainability requires the code to be clear, readable, have good organization and structure, proper documentation, and be free from code smells. Maintainability is a critical aspect of software quality, especially for large codebases, as it directly impacts the cost and effort required for future modifications, debugging, and enhancements.

A common approach to evaluating maintainability is to use static analyzers and metrics. Tools like SonarQube and CodeQL are frequently employed to measure the maintainability of generated code among the studies Nguyen and Nadi (2022)Zheng et al. (2024)Liu et al. (2024a)Yetistiren et al. (2023). Liu *et al.* Liu et al. (2024a) applied several static analysis tools to assess the code quality of ChatGPT-generated code. They found that nearly half of the generated code suffered from maintainability issues. The presence of main-

tainability issues in the generated code is also reported by Nguyen and Nadi (2022), who used SonarQube to evaluate the maintainability of the generated code. They found that Copilot generally produces low-complexity, easily understandable code. However, they warned that Copilot may produce bloated or overly verbose code, which needs oversight from developers.

Some studies also focus on the readability of generated code. Readability is a sub-factor of maintainability, and it refers to how easily a human can read and understand the code International Organization for Standardization (2023). Weyssow *et al.* Weyssow et al. (2024) evaluated the readability of code generated by LLMs using their proposed benchmark, *CodeUtra*. They found that code generated by GPT-4 and GPT-3.5 is more readable than that generated by open-source models, likely due to the well-formatted and structured training data. Zheng *et al.* Zheng et al. (2024) proposed the RACE benchmark to evaluate the quality of generated code across several dimensions, including readability. These studies show that LLMs still have a long way to go in producing maintainable code. The generated code is often correct, but not always maintainable. Besides readability, some studies also investigate the modularity and structure of the generated code. The study also reported that LLMs often generate code with poor modularity and structure, making it difficult for developers to understand and modify. Kang *et al.* Kang et al. (2024) questioned the traditional assumption that modularity improves code quality for LLM-generated code, finding that modular code does not consistently enhance performance and may not be favored by LLMs during generation.

Comparing the maintainability between generated code and human-crafted code is also a common approach. Licorish *et al.* Licorish et al. (2025) compared the maintainability of code written by GPT-4 and humans using several static analysis tools. They reported that human-written code was generally more readable and followed code standards better. Eltabakh *et al.* Eltabakh et al. (2024) reported a different finding using a larger and more diverse dataset. They found that the generated code is more structured, better documented, and easier to maintain, especially for simpler tasks. While human-written code is more flexible and adaptable, especially when addressing complex tasks, which require reasoning and creative thinking.

Overall, the studies show that LLMs can produce code that is easy to read and understand, but they often require human oversight to ensure maintainability. Some studies have attempted to have AI fix the maintainability issues in generated code Liu et al. (2024a). The iterative process of generating and

refining code can lead to improved maintainability, which suggests that LLMs are capable of producing maintainable code, but they do not inherently do so. The studies also show that the maintainability of generated code is not always better than human-written code, and it can depend on the task and the prompt used.

4.3. Discussion of literature review findings

Our analysis of the literature reveals that the research community has made some progress in evaluating the quality of generated code. However, several challenges and open gaps need to be addressed before NFQC evaluation in this context can become systematic and comparable across studies.

- **The definition of quality characteristics varies across studies.** The literature uses diverse terms when describing software quality characteristics. For example, maintainability is sometimes referred to as readability, and performance efficiency is called runtime or resource usage. The ISO/IEC 25010 standards provide a comprehensive framework for assessing software product quality, including both functional and non-functional characteristics. Currently, the research community has not adapted the standards to the context of LLM-generated code.
- **The evaluation of NFQCs is not comprehensive.** According to ISO/IEC 25010, there are several quality characteristics and sub-characteristics to evaluate the quality of a software at the code level. From the results of our literature review, currently, the research community has focused on a few quality characteristics, while giving less attention to others. The reason could be that the three quality characteristics are the primary concerns when developers assess the quality of their code. Also, other quality characteristics may not be easy to evaluate at the code level. For example, compatibility and flexibility are often evaluated at the system level. Nevertheless, a broader coverage of quality characteristics is essential for a holistic understanding of generated code quality.
- **A systematic and quantitative evaluation framework is needed.** Researchers measure the functional correctness of the generated code using accuracy or $pass@k$, which measures how often LLMs generate a correct response on its k_{th} attempt for a given problem. However, they employ diverse metrics and benchmarks to assess NFQCs. The use of

different metrics and benchmarks makes it difficult to compare results across models or studies. Overall, the field lacks a common methodological framework capable of quantitatively assessing multiple NFQCs alongside functional correctness at both code-level and system-level, which limits comparability and reproducibility.

- **The trade-offs between NFQCs should be explored.** According to our literature review, most studies investigate the NFQCs of the generated code, which is functionally correct. Some studies also reported the trade-offs between functional correctness and NFQCs. While such trade-offs are common in traditional software engineering (for example, between runtime and memory usage), LLMs introduced a new opportunity: these trade-offs can be influenced or even optimized through prompting strategies or other configurations. Understanding how LLMs balance or amplify these trade-offs is therefore essential for guiding their practical use in software development.

In summary, current research on NFQCs in generated code is inconsistent in terminology, selective in scope, and inconsistent in evaluation methods. Addressing these gaps will require unified definitions, standardized benchmarks, and comprehensive evaluation frameworks to enable a better understanding of generated code quality.

4.4. Workshop Findings

In this section, we summarize the primary insights from the workshops, and we identified several themes:

- **Not all ISO/IEC 25010 NFQCs are directly applicable when evaluating the generated code.** Participants highlighted that some characteristics, such as safety and compatibility, are typically assessed at the system level rather than the code level. By contrast, maintainability and security were recognized as characteristics that can be assessed from the generated code. Assessing whether a single piece of generated code satisfies system-level properties is generally not feasible.
- **The prioritization of NFQCs differed between academia and industry.** Although practitioners agree with the importance of performance efficiency and security, they emphasized that maintainability is particularly critical in practice, especially for large codebases. They

observed that the growing reliance on LLM-generated code may lead developers to overlook intricate code details, thereby increasing the effort required for future maintenance and repair. Practitioners were not satisfied with the current metrics used to measure maintainability. They asked for more reliable metrics that can reflect what they want.

- **Readability was identified as an essential code quality.** In addition to the ISO/IEC 25010 quality characteristics, many participants specifically highlighted readability. Although it is not formally defined as an NFQC in the ISO/IEC standard, it can be regarded as a subset of maintainability. Practitioners expressed a strong desire for tools that can assess and check the readability of the generated code. They emphasized that such support would help developers better evaluate the long-term quality of code, particularly in large projects where understanding and maintaining code are critical.
- **Participants discussed the issue of trade-offs among NFQCs.** They first emphasized the importance of functional correctness for generated code. Then they noted that the interactions among NFQCs themselves are less investigated, while existing studies focus on the trade-off between functional correctness and NFQCs. Questions were raised about whether NFQCs inherently exist in a zero-sum relationship and how developers might control or tune specific NFQCs, such as maintainability, security, or performance efficiency, when using LLM in development.

4.5. Discussion of Workshop findings

- **Scoping quality characteristics for evaluating generated code.** The difference between code-level and system-level NFQCs underscores the need to scope NFQC evaluation carefully. This difference also helps explain the imbalance observed in our literature review, where maintainability, security, and performance efficiency are most frequently examined. These three attributes are among the few that can be reliably assessed from code alone, as they have established automated metrics and tool support, such as static analysis and linters. In contrast, characteristics such as usability and compatibility depend on system context or user interaction and can not be reliably evaluated at the

code level. This observation also reflects a broader trend in current research, where both our empirical study and prior work primarily focus on these NFQCs due to their measurability at the code level.

- **Prioritized quality characteristics for generated code.** The emphasis on maintainability and readability from the practitioners highlights their concerns about maintaining software systems that incorporate generated code. As developers increasingly utilize LLMs for code generation and companies integrate these tools into their workflows, there is a higher propensity for developers to overlook intricate code details. Consequently, this oversight can significantly increase the time and effort required for future maintenance activities, including understanding the generated code, identifying potential issues, and devising solutions, ultimately contributing to the accumulation of technical debt. Currently, experienced engineers can still review generated code and identify potential vulnerabilities. However, as LLMs become more capable and produce increasingly complex code, the main challenge may shift from detecting flaws to understanding the code itself. In such a scenario, ensuring maintainability will become ever more critical.
- **The impact of trade-offs and enhancing specific quality characteristics.** A significant point of discussion centered on the complex nature and implications of trade-offs among the various quality attributes within generated code. While existing research has examined trade-offs between NFQCs and functional correctness Waghjale et al. (2024)Coignion et al. (2024b), a critical gap persists in understanding the intricate interactions among different NFQCs themselves and Functional Correctness. This lack of understanding gives rise to fundamental questions regarding the relationships between NFQCs: Do they inherently exist in a zero-sum relationship, where enhancing one attribute necessarily diminishes another? Or can methods and techniques be developed to navigate these interactions effectively, potentially enabling the simultaneous satisfaction of multiple NFQCs? Furthermore, what strategies can be employed to control and tune specific NFQCs, allowing developers to prioritize, for instance, security in security-critical contexts or optimize performance when that attribute is paramount, while managing the impact on other attributes?

Overall, the workshop discussions reveal three main insights. First, only a subset of NFQCs can be automatically and reliably measured at the code level, as there are established tools such as static analysis and linters. Second, practitioners emphasized that maintainability and readability are important aspects given their direct impact on long-term projects and complex systems. Third, participants highlighted the need to better understand and manage trade-offs between functional correctness and individual NFQCs. Future work should explore how multiple NFQCs interact and how to manage these trade-offs effectively. Together, these insights outline the current landscape of NFQC evaluation of generated code and motivate the empirical study presented in the following section.

Addressing the questions raised during the workshop regarding these NFQC trade-offs and interactions, the subsequent section of this paper details the establishment of our dataset and methodology specifically designed to investigate the interactions among NFQCs in generated code.

4.6. Experimental Results

This section presents the results of our empirical study, following the research questions and the stages of our experimental design.

Functional correctness of the generated patches. Figure 8 presents the comparative performance of the three models and the gold patches across 300 input instances. The results show clear differences among the three models across the two metrics: the **patch successfully applied rate (PSA)** and the **resolved rate**. The PSA indicates the proportion of generated patches that could be merged without errors, regardless of whether it compiles or passes the tests, and the resolved rate measures the proportion of patches that successfully resolved the issues described in the prompts. In all cases, their performance falls well below that of the gold patches, underscoring a substantial gap in functional correctness between LLM-generated patches and the curated reference patches produced by humans. GPT-4o has the highest PSA rate (75%) among the three models, but its resolved rate (16%) is the lowest. Claude-Sonnet-4 achieves a PSA rate of 58% together with the highest resolved rate of 36% among the three models. DeepSeek shows a more moderate performance, with a PSA rate of 46% and a resolved rate of 20%.

Among the failed attempts, we further inspected the agent logs to understand common failure causes. Notably, we observed that the agent’s retries

were largely devoted to resolving environment and shell command syntax errors (*e.g.*, import failures) instead of iteratively refining the patch logic, suggesting a weakness in context and tooling rather than code generation itself.

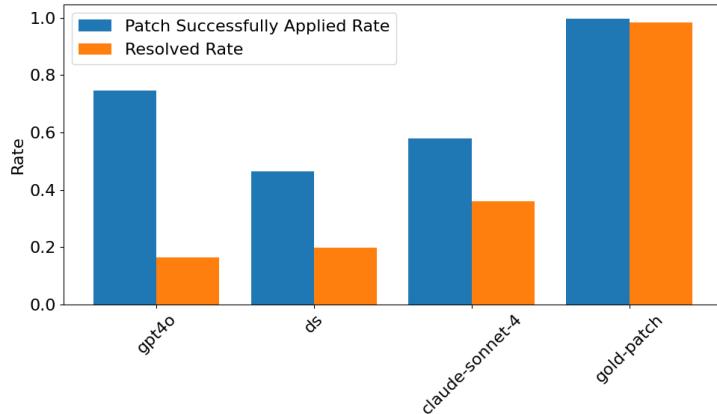


Figure 8: Patch generation results of different models in the baseline evaluation. The total number is 300.

Figure 9 reports functional correctness for patches regenerated with prompts tailored to specific NFQCs. Besides PSA and resolved rate, we included the submitted rate, which denotes the model provides a prediction for a instance. Under these NFQC-specific prompts, the functional correctness of the previously resolved patches decreased. The security prompt induces the biggest degradation in correctness, while memory and time have a lower impact on functional correctness. For GPT-4o, the submitted rate drops from 100% to 94% under the maintainability prompt, indicating that the model failed to generate patches for some instances. Under the security prompt, all three models show reduced submitted rates, with GPT-4o having the lowest resolved rate (12%). Claude-Sonnet-4 achieves the highest PSA and resolved rates under the maintainability, security, and time prompts, while DeepSeek-Reasoner achieves the best under the memory prompt, with both the PSA and the resolved rate reaching 94%.

RQ3.1: Baseline performance across NFQCs. To address RQ1, we first evaluated the resolved generated patches from the first stage in terms of four NFQCs: maintainability, security, runtime, and memory usage. To ensure

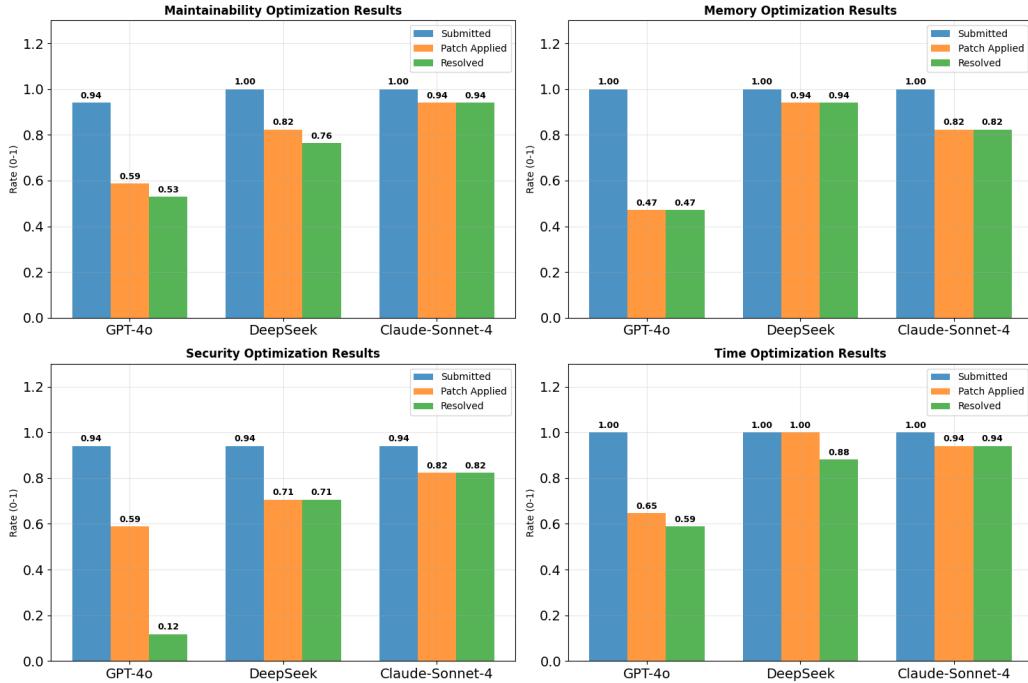


Figure 9: Patch generation results under different prompt strategies optimized for different NFQCs. Results are shown for three models across four NFQCs: (a) maintainability, (b) memory, (c) security, and (d) test runtime. For each NFQC, the bars illustrate the rate of submitted patches, the PSA, and the resolved rate.

comparability, the evaluation was conducted on the common subset of 17 resolved patches across the three models. Table 4 summarizes CodeQL results on the common subset of resolved patches for maintainability and security. Overall, the generated patches trigger substantially more CodeQL rules than the gold patches. When comparing the two dimensions, the number of triggered security rules is much lower than that of maintainability rules. Most rule hits for both dimensions fall under the "recommendation" severity level, but the maintainability checks also show considerable errors. For maintainability, each model triggers more than 150 errors. Among them, DeepSeek-Reasoner records the highest error count with 157, while GPT-4o shows the lowest counts across all three severity levels, with 152 errors, 377 recommendations, and 139 warnings. But for security, the largest number of rule hits appears in the recommendation severity level, where GPT-4o reaches 64. Across all three models, the generated patches trigger fewer than 10 rules at

the error and warning levels. Table 2 and Table 3 show the most triggered 5 rules of maintainability and security in our experiment. The description of each rule is extracted from the CodeQL document.

Table 2: The most triggered 5 maintainability rules.

Maintainability	Severity	Description
cyclic-import	R	Module forms part of an import cycle, thereby indirectly importing itself.
unsafe-cyclic-import	E	Module uses member of cyclically imported module, which can lead to failure at import time.
conflicting-attributes	W	When a class subclasses multiple base classes and more than one base class defines the same attribute, attribute overriding may result in unexpected behavior by instances of this class.
uninitialized-local-variable	E	Using a local variable before it is initialized causes an <code>UnboundLocalError</code> .
mixed-returns	R	Mixing implicit and explicit returns indicates a likely error as implicit returns always return ‘None’.

Table 5 reports the performance of the unoptimized patches with respect to test runtime and memory usage. All statistics are computed over the common subset of 17 resolved patches. Compared with the gold patches, generated patches generally run more slowly and consume more memory under baseline prompts. For example, GPT-4o records an average runtime of 23.37 seconds and a memory usage of 44.54 MB, both of which are clearly higher than the gold patches, which are 14.26 seconds and 23.50 MB. The patches generated by DeepSeek-Reasoner and Claude-Sonnet-4 show similar trends, with maximum runtimes reaching 53.67 seconds and 51.21 seconds, respectively.

To further examine the performance of each resolved patch across the four NFQCs, we computed weighted scores for maintainability and security, where higher values indicate higher risk. As shown in Figure 11, the gold patches consistently outperform the generated patches in the runtime, security and maintainability across nearly all instances. In terms of memory

Table 3: The most triggered 5 security rules.

Security	Severity	Description
unused-global-variable	R	Global variable is defined but not used.
empty-except	R	Except doesn't do anything and has no comment.
ineffectual-statement	R	A statement has no effect.
unused-local-variable	R	Local variable is defined but not used.
clear-text-storage-sensitive-data	E	Sensitive information stored without encryption or hashing can expose it to an attacker.

 Table 4: Results of CodeQL analysis on generated patches in terms of maintainability and security after the first stage. The table presents the number of rules triggered for each severity level across the generated patches and the gold patches. The severity level follows the order: **Error** > **Warning** > **Recommendation**.

Patch	Maintainability			Security		
	E	W	R	E	W	R
GPT-4o	152	139	377	8	2	64
DeepSeek	157	154	383	9	2	59
Claude-Sonnet-4	154	159	378	7	2	58
Gold-patch	4	154	3	2	0	32

usage, the gold patch generally consumes less memory, although in one case, the *mwaskom__seaborn-3190*, it exceeds the generated patches. The results also reveal substantial heterogeneity across instances: some instances exhibit markedly higher risks in specific NFQCs. For example, the *pydata__xarray-5131* instance shows disproportionately high weighted scores for both security and maintainability, suggesting that generated patches for this case are particularly problematic. Overall, while the gold patch maintains stable performance across instances, generated patches show uneven quality and tend to concentrate weakness on specific cases.

RQ3.2: Effect of NFQC-specific Prompt. To address RQ2 and RQ3, we applied NFQC-specific prompts to the patches generated in the first stage and

Table 5: Test runtime (in seconds) and memory usage (in MB) of unoptimized patches generated by three LLMs compared with the gold patches. For each model, the mean, standard deviation, minimum, and maximum values are reported across all 17 instances from the common subset.

Model	Test Runtime (s)				Memory (MB)			
	mean	std	min	max	mean	std	min	max
Gold-patch	14.26	12.29	3.51	46.28	23.50	3.84	20.04	37.25
GPT-4o	23.27	19.75	5.64	68.77	44.54	25.79	8.79	111.09
DeepSeek	17.68	14.09	4.42	53.67	24.99	7.22	9.25	39.25
Claude-Sonnet-4	16.72	14.72	4.39	51.21	27.27	9.91	8.20	48.61

re-evaluated the regenerated patches with respect to both functional correctness and the targeted NFQCs. To illustrate the effects of these prompts, we present different heatmaps that show the changes across all NFQCs under each optimization prompt.

Figure 12 reports the results for GPT-4o. Overall, applying NFQC-specific prompts to the initially resolved patches often compromises functional correctness, reducing the number of resolved instances. In particular, the security optimization substantially degrades correctness, with only two patches remaining *resolved* after regeneration. Maintainability optimization yields moderate gains in maintainability scores. Security optimization does not improve security scores and even introduces additional security issues. Runtime optimization consistently improves execution time but frequently introduces additional memory overhead (Figure 12c). Finally, memory optimization reduces memory usage in most instances, although the *pydata__xarray-5131* exhibits higher memory usage.

For DeepSeek-Reasoner, the proportion of resolved patches after applying NFQC-specific prompts remains comparatively high, with 13/17 for maintainability, 15/17 for runtime, 12/17 for security, and 16/17 for memory. This indicates that, unlike GPT-4o, DeepSeek-Reasoner preserves the functional correctness of generated patches more robustly under NFQC-specific prompts. Maintainability optimization manages to improve the maintainability score only in one case (*pydata__xarray-5131*), but increases the maintainability scores of other instances. Security optimization does not lead to clear improvements, and even increases the security scores of some instances. For time optimization, 10 instances have a smaller runtime, but other instances even get worse performance. The results of memory optimization

vary a lot. Some instances get improvements, but there are several instances get larger memory usages. For *mwaskom__seaborn-3190*, which is labeled ‘N/A’, the memory optimization introduces new security issues to the previous security issue-free patch.

For Claude-Sonnet-4, most patches retained functional correctness, with 16/17 instances resolved under the maintainability and runtime optimizations, 14/17 under the security and memory optimizations. Across the four NFQC optimizations, maintainability scores generally deteriorate, with only two instances showing improvements, indicating that optimization often at the cost of maintainability. Security optimization similarly tends to increase the number of security issues, although a few instances experience reductions. In contrast, runtime optimization reduces execution time for most instances with improvements exceeding 15%. The memory optimization likewise lowers memory usage in all cases. Overall, the results suggest that while runtime and memory optimizations provide tangible performance gains, maintainability and security optimizations tend to compromise NFQC qualities rather than enhance them.

To better capture the difference before and after optimization, we provide box plots (see Figure 10) of test runtime and memory usage. For runtime, DeepSeek-Reasoner and Claude-Sonnet 4 achieve lower median values across all four NFQC-specific optimizations compared with both the gold patch and baseline evaluations. In contrast, for GPT-4o, only the maintainability and memory optimization result in lower median values, while the security and time optimizations yield higher median values. In terms of dispersion, the interquartile ranges (IQRs) of GPT-4o and Claude-Sonnet 4 under optimizations are narrower than their respective baselines, reflecting a more stable runtime after optimization. Conversely, DeepSeek-Reasoner shows broader IQRs in all optimizations, suggesting increased variability and less predictable runtime performance.

For memory usage, the models demonstrate distinct patterns. For Claude-Sonnet 4, all four optimizations lead to lower medians and smaller IQRs compared to the baseline, indicating a consistent improvement in both central tendency and stability. For DeepSeek-Reasoner, the maintainability, security, and runtime optimizations slightly reduce both the median and dispersion, whereas the memory optimization results in a higher median and a wider IQR, implying potential overhead when explicitly targeting memory usage. For GPT-4o, both maintainability and memory optimizations decrease the median and reduce dispersion. The security optimization also lowers the

median, but its distribution becomes right-skewed with substantially larger IQRs. The time optimization, however, increases the median value but reduces the dispersion, indicating a trade-off between runtime and memory usage.

RQ3.3: Trade-offs across NFQCs. As shown in Figure 12, Figure 13 and Figure 14, optimizing for a specific NFQC often introduces trade-offs in others.

For GPT-4o, maintainability optimization often leads to improvements in both security and runtime, with security scores reaching 100% improvement for most instances. It also reduces memory usage in many patches. In contrast, the security optimization produces too few resolved instances to support a general conclusion. Runtime optimization tends to degrade the other three NFQCs, with particularly pronounced declines in maintainability and memory. For example, in the case of *django_django-13933*, runtime optimization results in a 761.77% performance decrease in memory usage compared to the baseline generated patches. Compared with runtime optimization, memory optimization doesn't improve memory usage consistently across all resolved instances, but it enhances the performance of the other NFQCs on some instances.

For DeepSeek-Reasoner, maintainability is the most vulnerable NFQC. Across the four optimizations, maintainability optimization leads to an increase in maintainability issues for most instances, while security and memory optimizations also reduce maintainability performance. Notably, the *pydata_xarray-5131* shows improvements in maintainability under the three NFQC optimizations except security. Runtime optimization produces more mixed effects. It improves maintainability, security, and memory usage for the majority of instances, but a few cases (*e.g.*, *djangodjango-14238*) exhibit clear regressions.

Claude-Sonnet-4 displays a similar pattern. Maintainability is again the most affected NFQC, with nearly all optimizations causing significant degradation in maintainability for most instances, except for *mwaskom_seaborn-3190* and *pydata_xarray-5131*, which consistently improved by all optimizations. Security optimization is particularly detrimental, with maintainability performance declining by as much as 719.05%. Unlike DeepSeek-Reasoner, security optimization occasionally produces gains in runtime and memory usage. In addition, runtime and memory optimization do not exhibit strong trade-offs; instead, they often improve both runtime and memory

Table 6: Number of malformed patches generated by each model under different optimization prompts (out of 17 patches per optimization). A malformed patch refers to a generated patch that could not be successfully applied to the repository or failed to comply with the unified diff format.

Model / Optimization	Maintainability	Security	Time	Memory
GPT-4o	3	4	3	7
DeepSeek-Reasoner	2	0	0	0
Claude-Sonnet 4	0	1	1	3

performance simultaneously.

4.7. Discussion of experiment results

The results demonstrate that current code LLMs face substantial limitations when addressing complex engineering tasks. By simulating how developers realistically use LLMs in their workflows, our evaluation provides insights into their actual applicability in software development. As discussed in Subsection 4.6, many of the agent retries were spent on resolving environment and command issues rather than the patch logic. These issues reveal a weakness in environment configuration, dependency resolution, and command invocation. In an agentic and robust system, the ability to use tools effectively is an integral part of the overall capability; therefore, such failures indicate limitations at the process level, even when the code generation itself is adequate. Moreover, in the final stage, many patches failed to apply due to context mismatches (see Table 6 and Figure 15), highlighting persistent difficulties in adhering to the unified diff format despite clear instructions in the prompts.

The comparison across models reveals different patterns. GPT-4o achieved the highest PSA rate in the first evaluation, yet its resolved rate was the lowest, suggesting that many of its patches were syntactically valid but functionally ineffective. After introducing NFQC optimizations, both its application and resolution rates declined, showing limited robustness in balancing functional correctness and NFQCs. In contrast, Claude-Sonnet-4 exhibited a lower PSA in the first evaluation but achieved a higher resolved rate. After the NFQC optimization, Claude-Sonnet-4 achieved the highest PSA and resolved rates under maintainability, security, and runtime optimization, while DeepSeek-Reasoner performed best under memory optimization, with both PSA and resolved rate reaching 94%. These results suggest a specialization

of strengths among models across different optimization objectives, rather than a single model outperforming others universally.

A further challenge lies in the lack of NFQC awareness. Under baseline prompts, all models produced patches that were substantially inferior to the gold patches in maintainability, security, runtime, and memory usage. This suggests that agent-based generation pipelines, while capable of producing functionally correct patches, implicitly prioritize “passing the test” rather than “passing with quality.” Such a lack of NFQC awareness constitutes a fundamental barrier to the integration of LLMs into automated workflows. It is not simply that models are unaware of best practices; their decision mechanisms do not treat NFQCs as primary objectives. Since LLMs are trained to predict the most likely next token, they capture superficial patterns of code rather than underlying engineering principles. Without an intrinsic cost function for code quality, all test-passing solutions are treated as equivalent. As a result, directly applying these patches risks introducing technical debt: the code may be functionally correct but structurally fragile and suboptimal. This limitation reflects the current model architectures and training paradigms rather than implementation errors.

Our baseline analysis further showed that LLM-generated patches triggered substantially more maintainability rules, especially high-severity violations, than the gold patches. For example, in Table 4, GPT-4o triggers 152 maintainability errors compared only four in the gold patch. This pattern was reinforced by results from DeepSeek-Reasoner and Claude-Sonnet-4, confirming the fragility of maintainability. Importantly, our scanning covered the entire codebase, meaning that we measured not only issues introduced by the patches themselves but also how the patches interacted with the existing code, thereby exposing underlying quality problems. This difference stems from the nature of rule enforcement: maintainability rules (*e.g.*, excessive function length, code duplication) are global assessments of repository state and can be triggered once a threshold is crossed, even without introducing new defects. In contrast, security rules are only triggered when a patch introduces new insecure operations, which is less likely in localized modifications. Consequently, the surge in maintainability warnings reflects not only new issues but also the way LLM patches align with the existing quality baseline, often amplifying structural weaknesses. It also highlights the necessity of integrating static analysis tools into generation pipelines to provide immediate feedback and mitigate the replication and amplification of undesirable code structures.

The NFQC optimization results also revealed trade-offs between different NFQCs. For example, GPT-4o’s runtime optimization reduced execution time but significantly increased memory usage, while Claude’s runtime and memory optimizations were consistently accompanied by declines in maintainability. More critically, we observed negative optimization, where models degraded the very dimension they were instructed to improve. DeepSeek-Reasoner’s security optimization did not yield stable gains and, for several instances, increased the security score (worse), while Claude’s maintainability and security optimizations each degraded their respective targets. These findings suggest that current models lack a reliable internal representation of abstract rules relating to security and sustainability. When prompted with such rules and their triggered position, models often misinterpreted the intended objective and produced regressions. This indicates that current NFQC optimization is not yet trustworthy. Any optimization attempt must be accompanied by targeted validation, as improvements cannot be assumed to generalize across tasks.

Cross-model comparison further illustrates these differences. GPT-4o pursued a more aggressive strategy, delivering large improvements in some NFQCs (*e.g.*, under maintainability optimization, security scores for several instances improved up to 100%, see Figure 12a), but often at the cost of functional correctness. In contrast, DeepSeek-Reasoner and Claude-Sonnet-4 adopted more conservative approaches. They maintained higher rates of functional correctness across most optimizations but produced unstable and frequently negative results in NFQCs, particularly maintainability. This conservatism may stem from their stronger alignment toward correctness during training, which encouraged caution in logical code changes.

Finally, two particular instances stand out in our results. *Pydata__xarray-5131* showed consistent improvements in maintainability and security scores across models and optimizations, suggesting that instances with severe initial deficiencies and relatively lenient functional constraints allow for modifications without breaking tests. By contrast, *mwaskom__seaborn-3190* illustrated the inherent conflict among NFQCs. Across several optimizations, functional correctness was compromised, suggesting that its code was in a fragile equilibrium where any modification risked undermining both security and functionality.

5. Threats to Validity

This section presents the potential issues that may affect the validity and reproducibility of our findings.

5.1. Internal Validity

For the literature review, our search strategy followed key elements of the PRISMA guidelines, but there is still a risk that relevant studies have been missed due to database coverage limitations or the keywords. To mitigate the risk, we applied both forward and backward snowballing and found several relevant studies.

For the workshops, two authors independently took notes during the workshops; individual interpretations may have introduced subjective bias and affected the objectivity of the notes. To mitigate this, we conducted joint discussions after each workshop to reconcile the two records, ensuring a consistent interpretation of the practitioners' views.

In our experiments, we combined the prompts provided in the SWE-bench Lite benchmark with the CodeQL reports to design NFQC optimization prompts. This procedure mitigates subjective bias to some extent, but prompt sensitivity remains a threat. Minor changes in the prompts could affect LLM behavior and result in different generated patches. In addition, different hyperparameter settings, such as temperature, may influence the stochasticity and diversity of model outputs.

Second, our experiments were conducted on a single workstation (AMD Ryzen 7, 16GB RAM, Ubuntu 22.04). Although the evaluation environment was kept consistent across all models and stages, background system loads could have introduced small variations in runtime and memory measurements. These potential fluctuations were not explicitly recorded.

Third, we used three LLMs in the experiment. But the models update frequently during the experiment period, which could lead to model drift and different behaviors in the future. This is a known limitation when evaluating LLMs.

Finally, to ensure comparability across models, we restricted our analysis to instances successfully solved by three models. However, it may create a bias towards simpler tasks, potentially skewing the observed trade-offs among NFQCs.

5.2. Construct Validity

For the literature review, we mapped diverse quality characteristics from different studies to the ISO/IEC 25010 quality model; there is a risk that we may not always understand and reflect the authors' original intent during the mapping.

Individual understanding of concepts such as readability and security by participants might also have varied, leading to potential inconsistencies in the way in which some discussions were understood. To mitigate this, we reviewed and discussed our notes to ensure a consistent understanding of their perspectives.

The metrics used in the experiment represent practical but partial approximations of NFQCs. Maintainability and security were quantified through risk scores calculated from CodeQL rule counts, while performance efficiency was assessed through runtime and memory usage. Although these measures provide systematic and reproducible indicators, they capture only part of the underlying constructs.

5.3. External Validity

The literature review included studies up to early 2025. As the field evolves rapidly, newer studies may reveal additional findings that differ from ours.

Most of the workshop participants are Software Center companies, and there is a risk that their perspectives may not fully represent the broader software industry.

The generalizability of our results is limited by the dataset and models used. SWE-bench Lite was derived from real software repositories in GitHub, but it represents a relatively small subset of projects. After filtering for functional correctness, the number of usable instances further decreased to 17. This limits the applicability of the findings to instances solvable by the code LLMs. Additionally, different models may exhibit different behaviors. As code LLMs evolve rapidly, new models may show improved ability in delivering high-quality code. However, we believe that the observed interaction patterns and trade-off tendencies provide valuable insights that are likely to remain relevant.

6. Conclusion and Future Work

6.1. Conclusion

Our study systematically examined the performance of LLM-generated code with respect to NFQCs in real-world software engineering tasks. By combining evidence from a literature review, industry workshops, and empirical experiments, we identified both the current progress and the limitations of LLMs in addressing NFQCs and outlined key challenges that future research and practice need to address.

The literature review shows that existing research has mainly focused on security, performance efficiency, maintainability, and reliability, while other NFQCs have received comparatively little attention. At the same time, different studies often employ inconsistent terminology for the same NFQC, which hampers cross-study comparisons and highlights the absence of a unified framework.

Secondly, insights from workshops highlight the differences between academic and industrial perspectives. While academic research tends to emphasize security and performance efficiency, practitioners show stronger concern for maintainability, particularly in the context of large software systems and long-term maintenance.

Our experiment confirms that current LLMs in code generation lack awareness of NFQCs, with their default objective being to “pass the test” rather than to “pass with quality”. Applying LLM-generated patches in projects without validation entails the risk of introducing and amplifying technical debt.

Overall, this study highlights the significant limitations and trade-offs of LLMs with respect to NFQCs and underscores the importance of closer collaboration between academia and industry. We argue that only through the combined efforts of unified frameworks, context-aware optimization, and robust verification can the gap between functional correctness and NFQC quality be progressively narrowed, thus laying the foundation for the safe and sustainable application of AI-generated code in real-world software engineering practice.

6.2. Future Work

To move the generated code from “passing the tests” to “passing with quality”, future work should connect evaluation, modeling, and tooling. In future work, we are going to integrate strong verification mechanisms into the

generation pipeline. We will embed static analysis and automated detectors into generation and optimization so that NFQCs are monitored in real time, used as gates for candidate patches, and returned as structured feedback for revision. By addressing NFQCs before deployment, we aim to reduce technical debt accumulation and improve the reliability of LLM-generated code in software systems.

7. Acknowledgment

This work was funded by the Software Center project 61 and the Vinnova Competence Center for Continuous Digitalization 2023-00546. The authors are indebted to Willem Meijer for fruitful discussions.

References

- Asare, O., Nagappan, M., Asokan, N., 2023. Is github’s copilot as bad as humans at introducing vulnerabilities in code? *Empir. Softw. Eng.* 28, 129. URL: <https://doi.org/10.1007/s10664-023-10380-1>, doi:10.1007/S10664-023-10380-1.
- Bolt, 2024. Bolt – Your AI-powered development agent. URL: <https://bolt.new>. accessed: May 13, 2025.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H.P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F.P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W.H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A.N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., Zaremba, W., 2021. Evaluating large language models trained on code. *CoRR* abs/2107.03374. URL: <https://arxiv.org/abs/2107.03374>, arXiv:2107.03374.
- Coignion, T., Quinton, C., Rouvoy, R., 2024a. A performance study of llm-generated code on leetcode, in: Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering, EASE

- 2024, Salerno, Italy, June 18-21, 2024, ACM. pp. 79–89. URL: <https://doi.org/10.1145/3661167.3661221>, doi:10.1145/3661167.3661221.
- Coignion, T., Quinton, C., Rouvoy, R., 2024b. A Performance Study of LLM-Generated Code on Leetcode, in: Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering, ACM, Salerno Italy. pp. 79–89. doi:10.1145/3661167.3661221.
- Cursor, 2024. Cursor: The AI coding assistant. URL: <https://www.cursor.com>. accessed: May 13, 2025.
- Elgedawy, R., Sadik, J., Dutta, S., Gautam, A., Georgiou, K., Gholamrezae, F., Ji, F., Lim, K., Liu, Q., Ruoti, S., 2024. Ocassionally Secure: A Comparative Analysis of Code Generation Assistants. doi:10.48550/arXiv.2402.00689, arXiv:2402.00689.
- Eltabakh, T.M., Nabil Soudi, N., Shawky, D., 2024. Quality of AI-generated vs. human-generated code, in: 2024 34th International Conference on Computer Theory and Applications (ICCTA), pp. 200–205. doi:10.1109/ICCTA64612.2024.10974782.
- Fu, Y., Liang, P., et al., 2025. Security weaknesses of copilot-generated code in github projects: An empirical study. ACM Trans. Softw. Eng. Methodol. URL: <https://doi.org/10.1145/3716848>, doi:10.1145/3716848. just Accepted.
- GitHub, 2025. Github copilot: AI code completion tool. URL: <https://github.com/features/copilot>. accessed: 2025-04-20.
- Hemel, Z., Kats, L.C.L., Visser, E., 2008. Code generation by model transformation, in: Vallecillo, A., Gray, J., Pierantonio, A. (Eds.), Theory and Practice of Model Transformations, Springer Berlin Heidelberg, Berlin, Heidelberg. pp. 183–198.
- Huang, D., Qing, Y., Shang, W., Cui, H., Zhang, J., 2024. Effibench: Benchmarking the efficiency of automatically generated code, in: Globersons, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J.M., Zhang, C. (Eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 -

15, 2024. URL: http://papers.nips.cc/paper_files/paper/2024/hash/15807b6e09d691fe5e96cdecde6d7b80-Abstract-Datasets_and_Benchmarks_Track.html.

Hui, B., Yang, J., Cui, Z., Yang, J., Liu, D., Zhang, L., Liu, T., Zhang, J., Yu, B., Dang, K., Yang, A., Men, R., Huang, F., Ren, X., Ren, X., Zhou, J., Lin, J., 2024. Qwen2.5-coder technical report. CoRR abs/2409.12186. URL: <https://doi.org/10.48550/arXiv.2409.12186>, doi:10.48550/ARXIV.2409.12186, arXiv:2409.12186.

International Organization for Standardization, 2023. ISO/IEC 25010:2023 Systems and software engineering — Systems and software Quality Requirements and Evaluation (SQuaRE) — Product quality model. <https://www.iso.org/standard/78176.html>. Accessed: 2025-04-19.

Jiang, A.Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D.S., de Las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L.R., Lachaux, M., Stock, P., Scao, T.L., Lavril, T., Wang, T., Lacroix, T., Sayed, W.E., 2023. Mistral 7b. CoRR abs/2310.06825. URL: <https://doi.org/10.48550/arXiv.2310.06825>, doi:10.48550/ARXIV.2310.06825, arXiv:2310.06825.

Jimenez, C.E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., Narasimhan, K.R., 2024. Swe-bench: Can language models resolve real-world github issues?, in: The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024, OpenReview.net. URL: <https://openreview.net/forum?id=VTF8yNQM66>.

Kang, D., Seo, K.J., Kim, T., 2024. Revisiting the impact of pursuing modularity for code generation, in: Al-Onaizan, Y., Bansal, M., Chen, Y. (Eds.), Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024, Association for Computational Linguistics. pp. 11561–11571. URL: <https://aclanthology.org/2024.findings-emnlp.676>.

Khoury, R., Avila, A.R., Brunelle, J., Camara, B.M., 2023. How secure is code generated by chatgpt?, in: IEEE International Conference on Systems, Man, and Cybernetics, SMC 2023, Honolulu, Oahu, HI, USA, October 1-4, 2023, IEEE. pp. 2445–2451. URL: <https://doi.org/10.1109/SMC53992.2023.10394237>, doi:10.1109/SMC53992.2023.10394237.

- Li, R., Allal, L.B., Zi, Y., Muennighoff, N., Kocetkov, D., Mou, C., et al., 2023. Starcoder: may the source be with you! *Trans. Mach. Learn. Res.* 2023. URL: <https://openreview.net/forum?id=KoF0g41haE>.
- Li, Y., Choi, D.H., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., et al., 2022. Competition-level code generation with alphacode. *Science* 378, 1092–1097. URL: <https://www.science.org/doi/abs/10.1126/science.abq1158>, doi:10.1126/science.abq1158, arXiv:<https://www.science.org/doi/pdf/10.1126/science.abq1158>.
- Licorish, S.A., Bajpai, A., Arora, C., Wang, F., Tantithamthavorn, C., 2025. Comparing human and LLM generated code: The jury is still out! *CoRR* abs/2501.16857. URL: <https://doi.org/10.48550/arXiv.2501.16857>, doi:10.48550/ARXIV.2501.16857, arXiv:2501.16857.
- Liu, Y., Le-Cong, T., Widjyasaari, R., Tantithamthavorn, C., Li, L., Le, X.B.D., Lo, D., 2024a. Refining ChatGPT-Generated Code: Characterizing and Mitigating Code Quality Issues. *ACM Trans. Softw. Eng. Methodol.* 33, 116:1–116:26. doi:10.1145/3643674.
- Liu, Z., Tang, Y., Luo, X., Zhou, Y., Zhang, L.F., 2024b. No need to lift a finger anymore? assessing the quality of code generation by chatgpt. *IEEE Trans. Software Eng.* 50, 1548–1584. URL: <https://doi.org/10.1109/TSE.2024.3392499>, doi:10.1109/TSE.2024.3392499.
- Liventsev, V., Grishina, A., Härmä, A., Moonen, L., 2023. Fully autonomous programming with large language models, in: Silva, S., Paquete, L. (Eds.), *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO 2023, Lisbon, Portugal, July 15-19, 2023, ACM.* pp. 1146–1155. URL: <https://doi.org/10.1145/3583131.3590481>, doi:10.1145/3583131.3590481.
- Lovable, 2024. Build apps with an AI engineer. URL: <https://lovable.dev>. accessed: May 13, 2025.
- Luo, Z., Xu, C., Zhao, P., Sun, Q., Geng, X., Hu, W., Tao, C., Ma, J., Lin, Q., Jiang, D., 2024. Wizardcoder: Empowering code large language models with evol-instruct, in: *The Twelfth International Conference on Learning Representations.* URL: <https://openreview.net/forum?id=UnUwSIgK5W>.

- Manna, Z., Waldinger, R.J., 1971. Toward automatic program synthesis. Commun. ACM 14, 151–165. URL: <https://doi.org/10.1145/362566.362568>, doi:10.1145/362566.362568.
- MITRE, 2024. CWE Top 25 Most Dangerous Software Weaknesses. URL: https://cwe.mitre.org/top25/archive/2024/2024_cwe_top25.html. accessed: May 13, 2025.
- Natella, R., Liguori, P., Improta, C., Cukic, B., Cotroneo, D., 2024. AI code generators for security: Friend or foe? IEEE Secur. Priv. 22, 73–81. URL: <https://doi.org/10.1109/MSEC.2024.3355713>, doi:10.1109/MSEC.2024.3355713.
- Nguyen, N., Nadi, S., 2022. An empirical evaluation of github copilot’s code suggestions, in: 19th IEEE/ACM International Conference on Mining Software Repositories, MSR 2022, Pittsburgh, PA, USA, May 23-24, 2022, ACM. pp. 1–5. URL: <https://doi.org/10.1145/3524842.3528470>, doi:10.1145/3524842.3528470.
- Niu, C., Zhang, T., Li, C., Luo, B., Ng, V., 2024. On evaluating the efficiency of source code generated by llms, in: Lo, D., Xia, X., Penta, M.D., Hu, X. (Eds.), Proceedings of the 2024 IEEE/ACM First International Conference on AI Foundation Models and Software Engineering, FORGE 2024, Lisbon, Portugal, 14 April 2024, ACM. pp. 103–107. URL: <https://doi.org/10.1145/3650105.3652295>, doi:10.1145/3650105.3652295.
- Page, M.J., McKenzie, J.E., Bossuyt, P.M., Boutron, I., Hoffmann, T.C., Mulrow, C.D., Shamseer, L., Tetzlaff, J.M., Akl, E.A., Brennan, S.E., Chou, R., Glanville, J., Grimshaw, J.M., Hróbjartsson, A., Lalu, M.M., Li, T., Loder, E.W., Mayo-Wilson, E., McDonald, S., McGuinness, L.A., Stewart, L.A., Thomas, J., Tricco, A.C., Welch, V.A., Whiting, P., Moher, D., 2021. The PRISMA 2020 statement: An updated guideline for reporting systematic reviews. BMJ (Clinical research ed.) 372. doi:10.1136/bmj.n71, arXiv:<https://www.bmjjournals.com/content/372/bmj.n71.full.pdf>.
- Pearce, H., Ahmad, B., Tan, B., Dolan-Gavitt, B., Karri, R., 2022. Asleep at the Keyboard? Assessing the Security of GitHub Copilot’s Code Contributions, in: 2022 IEEE Symposium on Security and Privacy (SP), IEEE Computer Society. pp. 754–768. doi:10.1109/SP46214.2022.9833571.

- Pearce, H., Ahmad, B., Tan, B., Dolan-Gavitt, B., Karri, R., 2025. Asleep at the keyboard? assessing the security of github copilot's code contributions. *Commun. ACM* 68, 96–105. URL: <https://doi.org/10.1145/3610721>, doi:10.1145/3610721.
- Perry, N., Srivastava, M., Kumar, D., Boneh, D., 2023. Do users write more insecure code with AI assistants?, in: Meng, W., Jensen, C.D., Cremers, C., Kirda, E. (Eds.), *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security, CCS 2023, Copenhagen, Denmark, November 26-30, 2023*, ACM. pp. 2785–2799. URL: <https://doi.org/10.1145/3576915.3623157>, doi:10.1145/3576915.3623157.
- Rozière, B., Gehring, J., Gloeckle, F., Sootla, S., Gat, I., Tan, X.E., Adi, Y., Liu, J., Remez, T., Rapin, J., Kozhevnikov, A., Evtimov, I., Bitton, J., Bhatt, M., Canton-Ferrer, C., Grattafiori, A., Xiong, W., Défossez, A., Copet, J., Azhar, F., Touvron, H., Martin, L., Usunier, N., Scialom, T., Synnaeve, G., 2023. Code llama: Open foundation models for code. *CoRR* abs/2308.12950. URL: <https://doi.org/10.48550/arXiv.2308.12950>, doi:10.48550/ARXIV.2308.12950, arXiv:2308.12950.
- Shypula, A., Madaan, A., Zeng, Y., Alon, U., Gardner, J.R., Yang, Y., Hashemi, M., Neubig, G., Ranganathan, P., Bastani, O., Yazdanbakhsh, A., 2024. Learning performance-improving code edits, in: *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024, OpenReview.net*. URL: <https://openreview.net/forum?id=ix7rLVHXyY>.
- Siddiq, M.L., Roney, L., Zhang, J., Santos, J.C.S., 2024. Quality assessment of chatgpt generated code and their use by developers, in: Spinellis, D., Bacchelli, A., Constantinou, E. (Eds.), *21st IEEE/ACM International Conference on Mining Software Repositories, MSR 2024, Lisbon, Portugal, April 15-16, 2024*, ACM. pp. 152–156. URL: <https://doi.org/10.1145/3643991.3645071>, doi:10.1145/3643991.3645071.
- Siddiq, M.L., Santos, J.C.S., 2022. Securityeval dataset: mining vulnerability examples to evaluate machine learning-based code generation techniques, Association for Computing Machinery, New York, NY, USA. p. 29–33. URL: <https://doi.org/10.1145/3549035.3561184>, doi:10.1145/3549035.3561184.

- Singhal, M., Aggarwal, T., Awasthi, A., Natarajan, N., Kanade, A., 2024. NoFunEval: Funny How Code LMs Falter on Requirements Beyond Functional Correctness. URL: <http://arxiv.org/abs/2401.15963>, doi:10.48550/arXiv.2401.15963. arXiv:2401.15963 [cs].
- Solovyeva, L., Weidmann, S., Castor, F., 2025. AI-powered, but power-hungry? energy efficiency of llm-generated code, in: IEEE/ACM Second International Conference on AI Foundation Models and Software Engineering, Forge@ICSE 2025, Ottawa, ON, Canada, April 27-28, 2025, IEEE. pp. 49–60. URL: <https://doi.org/10.1109/Forge66646.2025.00012>, doi:10.1109/FORGE66646.2025.00012.
- The AI Digest, 2024. Time Horizons: How Far Can We See in AI? URL: <https://theaidigest.org/time-horizons>. accessed: May 12, 2025.
- Tony, C., Mutus, M., Ferreyra, N.E.D., Scandariato, R., 2023. Llmseceval: A dataset of natural language prompts for security evaluations, in: 20th IEEE/ACM International Conference on Mining Software Repositories, MSR 2023, Melbourne, Australia, May 15-16, 2023, IEEE. pp. 588–592. URL: <https://doi.org/10.1109/MSR59073.2023.00084>, doi:10.1109/MSR59073.2023.00084.
- Vartziotis, T., Dellatolas, I., Dasoulas, G., Schmidt, M., Schneider, F., Hoffmann, T., Kotsopoulos, S., Keckeisen, M., 2024. Learn to code sustainably: An empirical study on green code generation, in: LLM4CODE@ICSE, pp. 30–37. URL: <https://doi.org/10.1145/3643795.3648394>, doi:10.1145/3643795.3648394.
- Waghjale, S., Veerendranath, V., Wang, Z., Fried, D., 2024. ECCO: can we improve model-generated code efficiency without sacrificing functional correctness?, in: Al-Onaizan, Y., Bansal, M., Chen, Y. (Eds.), Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024, Association for Computational Linguistics. pp. 15362–15376. URL: <https://aclanthology.org/2024.emnlp-main.859>.
- Wang, J., Chen, Y., 2023. A review on code generation with llms: Application and evaluation, in: 2023 IEEE International Conference on Medical Artificial Intelligence (MedAI), pp. 284–289. doi:10.1109/MedAI59581.2023.00044.

- Weyssow, M., Kamanda, A., Sahraoui, H.A., 2024. Codeultrafeedback: An llm-as-a-judge dataset for aligning large language models to coding preferences. CoRR abs/2403.09032. URL: <https://doi.org/10.48550/arXiv.2403.09032>, doi:10.48550/ARXIV.2403.09032, arXiv:2403.09032.
- Windsurf, 2024. Windsurf Editor: The AI-native IDE. URL: <https://windsurf.com>. accessed: May 13, 2025.
- Yang, Z., Sun, Z., Yue, T.Z., Devanbu, P., Lo, D., 2024. Robustness, security, privacy, explainability, efficiency, and usability of large language models for code. URL: <https://arxiv.org/abs/2403.07506>, arXiv:2403.07506.
- Yetistiren, B., Özsoy, I., Ayerdem, M., Tüzin, E., 2023. Evaluating the code quality of ai-assisted code generation tools: An empirical study on github copilot, amazon codewhisperer, and chatgpt. CoRR abs/2304.10778. URL: <https://doi.org/10.48550/arXiv.2304.10778>, doi:10.48550/ARXIV.2304.10778, arXiv:2304.10778.
- Zheng, J., Cao, B., Ma, Z., Pan, R., Lin, H., Lu, Y., Han, X., Sun, L., 2024. Beyond correctness: Benchmarking multi-dimensional code generation for large language models. CoRR abs/2407.11470. URL: <https://doi.org/10.48550/arXiv.2407.11470>, doi:10.48550/ARXIV.2407.11470, arXiv:2407.11470.
- Zhong, L., Wang, Z., Shang, J., 2024. Debug like a human: A large language model debugger via verifying runtime execution step by step, in: Ku, L., Martins, A., Srikumar, V. (Eds.), Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024, Association for Computational Linguistics. pp. 851–870. URL: <https://doi.org/10.18653/v1/2024.findings-acl.49>, doi:10.18653/V1/2024.FINDINGS-ACL.49.

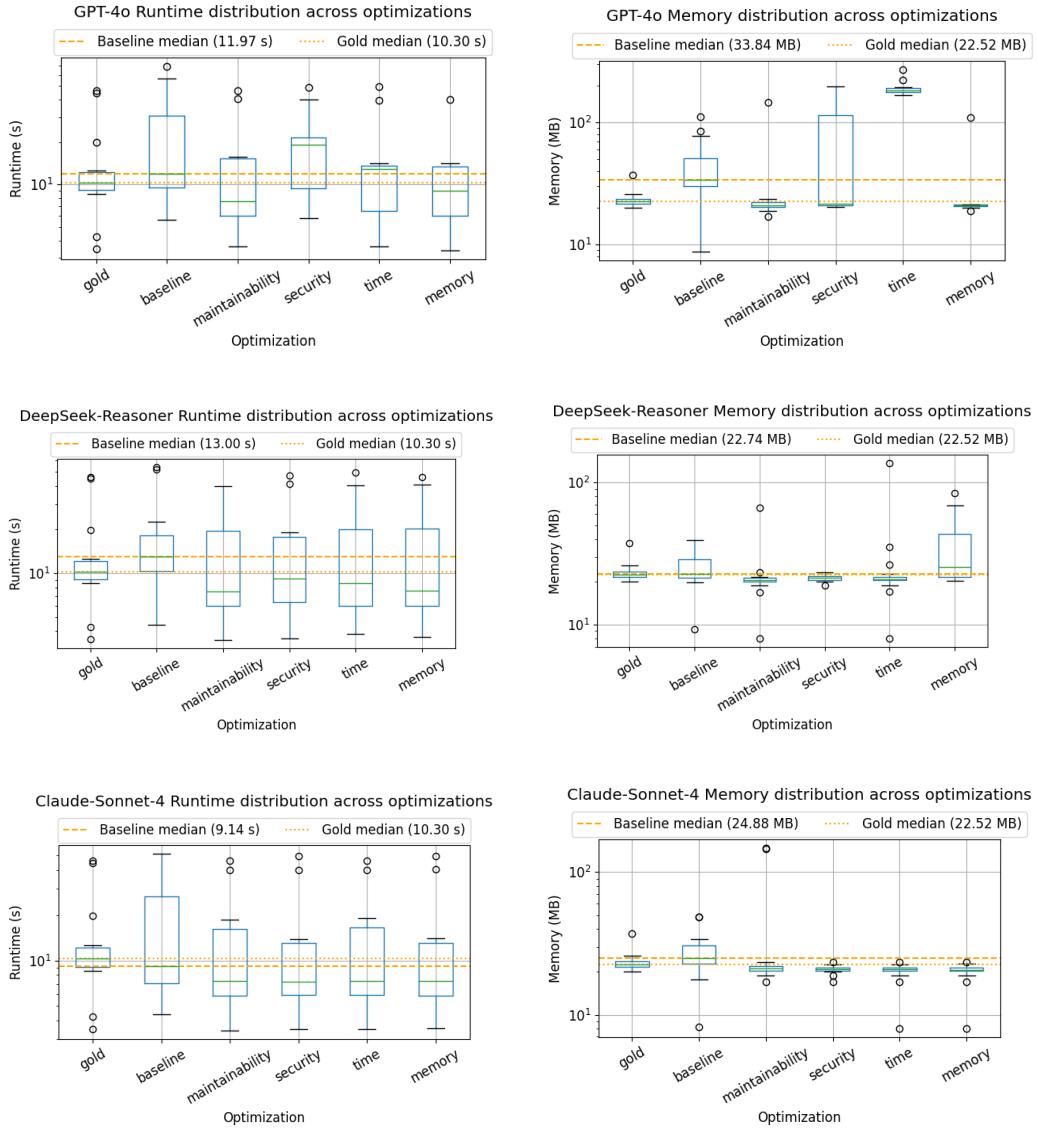


Figure 10: Boxplots of test runtime (s) and memory usage (MB) on a logarithmic scale for patches generated by GPT-4o, DeepSeek-Reasoner, and Claude-Sonnet-4, compared with the **baseline** results of baseline evaluation of each model. The boxes summarize the distribution across resolved instances, with dashed and dotted lines indicating the median values of the baseline evaluation and gold patches, respectively.

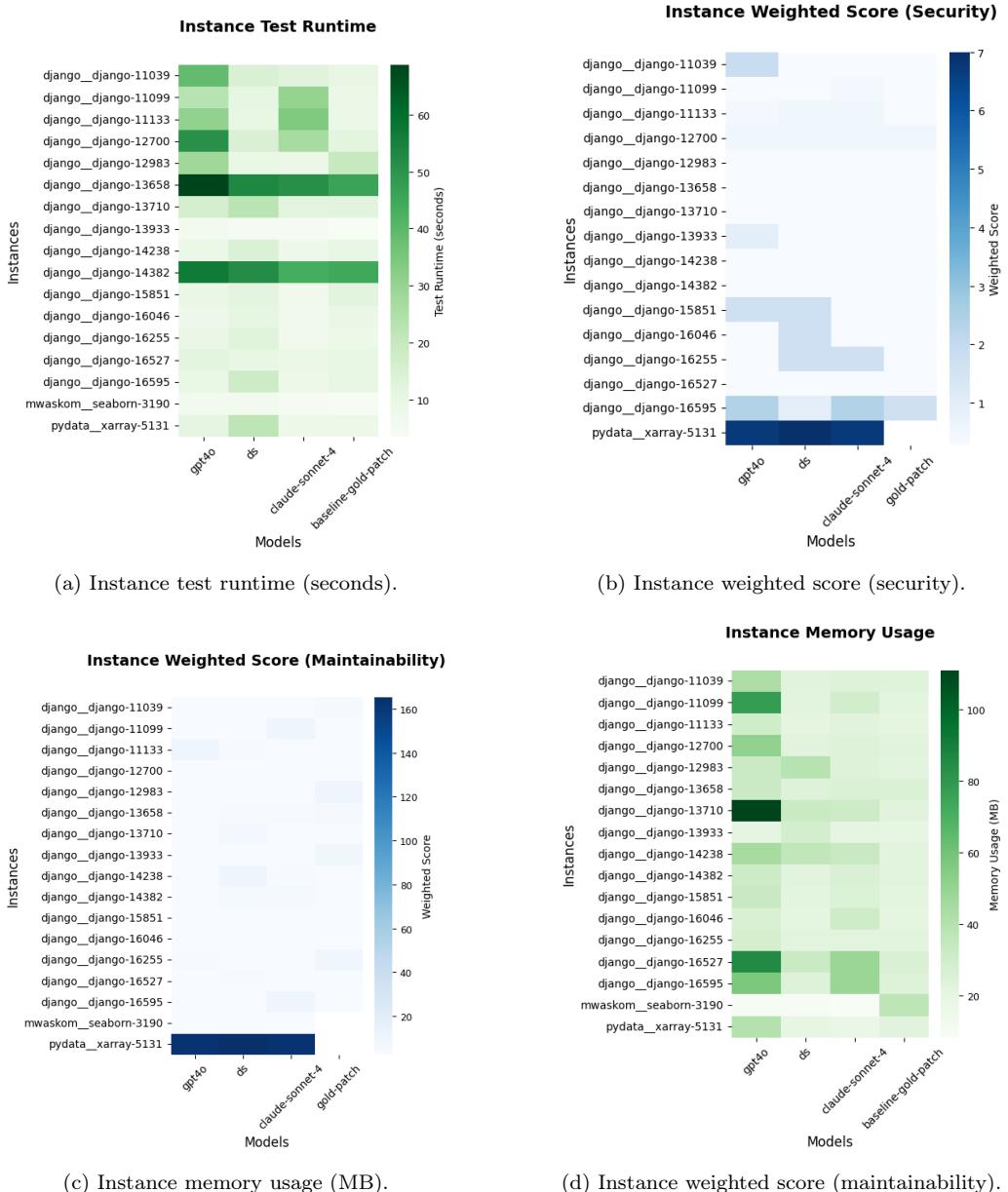


Figure 11: Per-instance comparison across models (GPT-4o, DeepSeek, Claude-Sonnet-4, and the gold patch). Panels show (a) test runtime, (b) weighted security score, (c) memory usage, and (d) weighted maintainability score. Darker colors indicate higher values; instances are listed on the vertical axis and models on the horizontal axis.

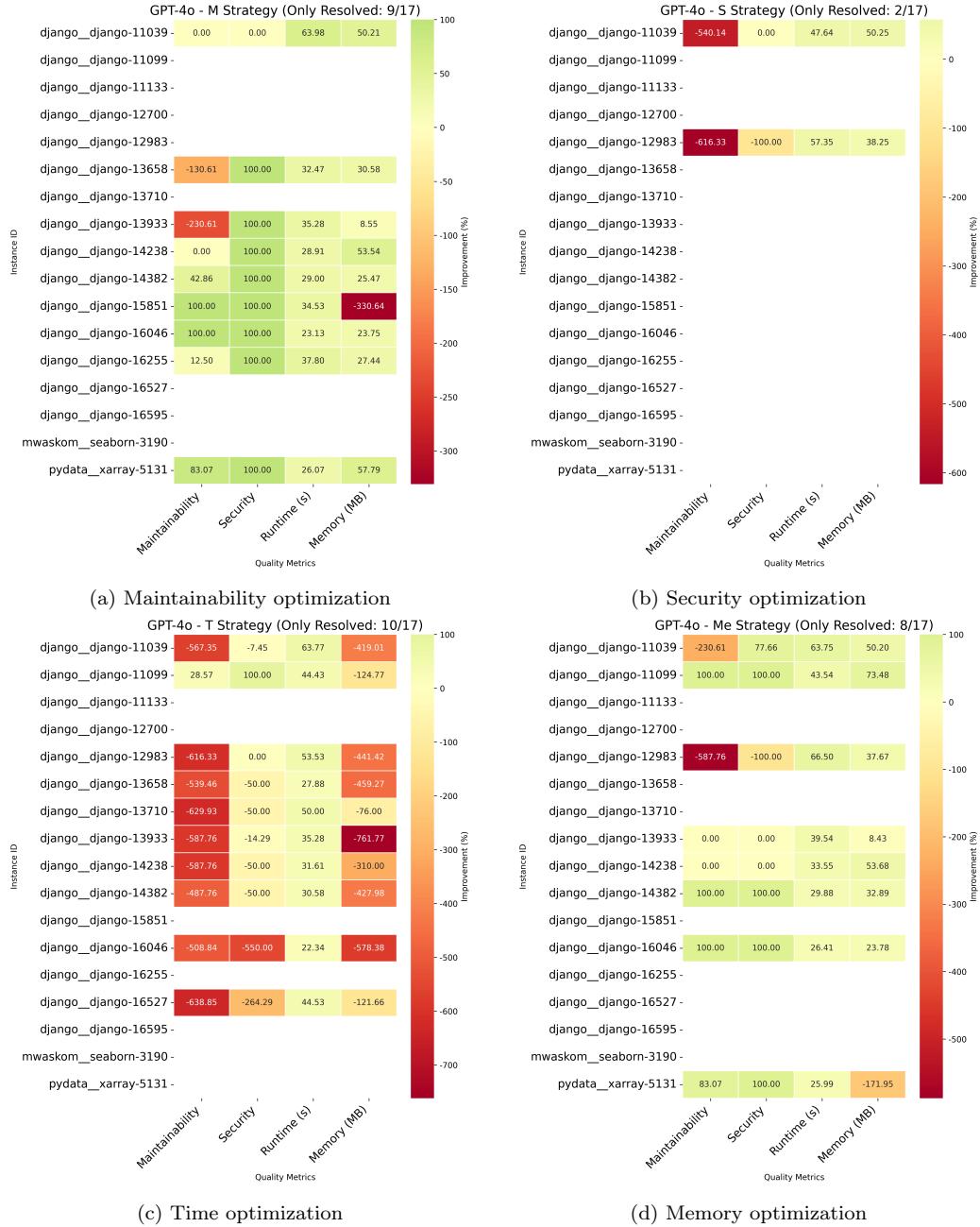


Figure 12: Difference heatmaps for GPT-4o under four NFQC-specific optimizations. Each heatmap shows the change in metric values between the first-stage resolved patches and the regenerated patches. Green color indicates improvement, while red indicates deterioration.

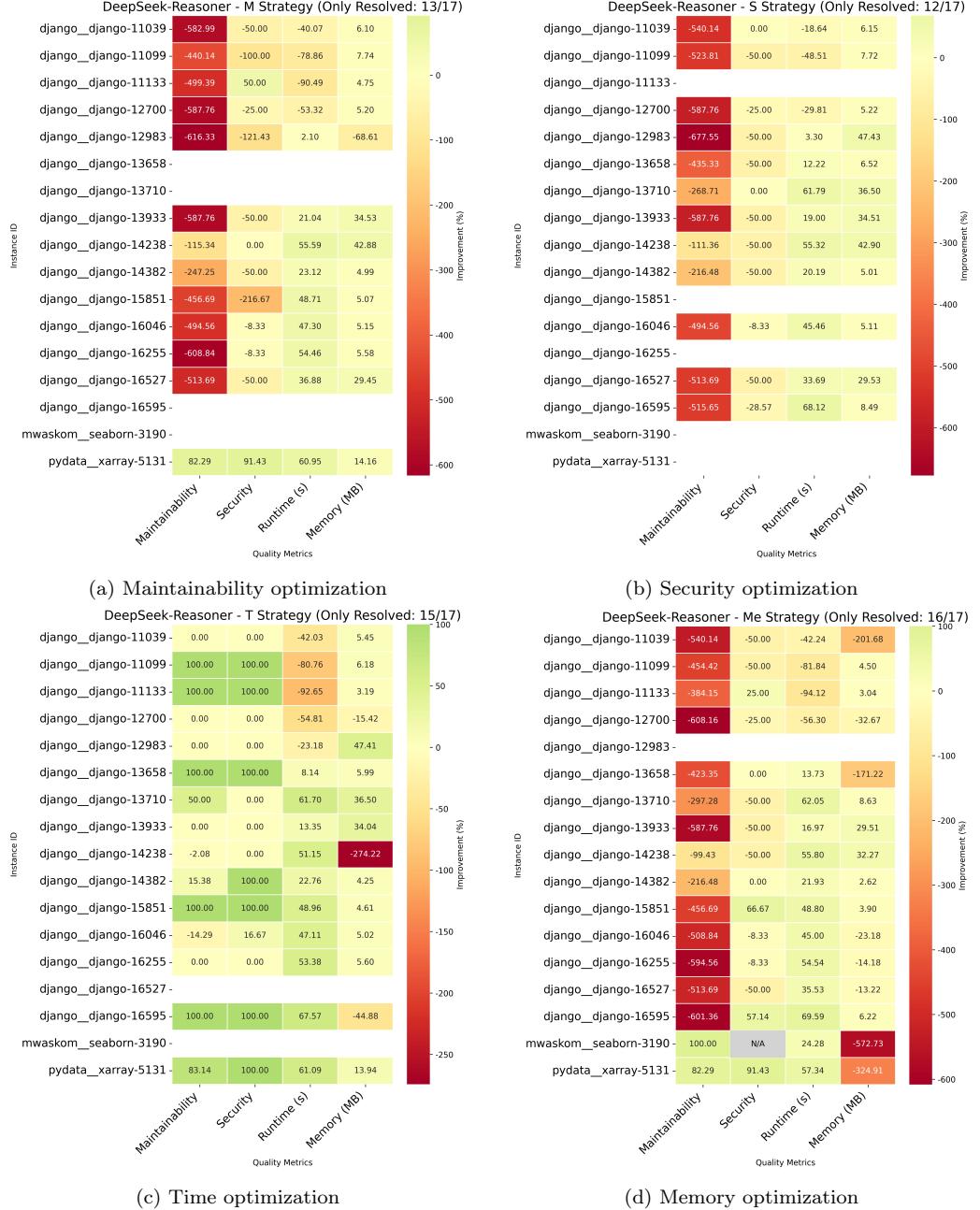


Figure 13: Difference heatmaps for DeepSeek-Reasoner under four NFQC-specific optimizations. Each heatmap shows the change in metric values between the first stage resolved patches and the regenerated patches. The gray color means the value is zero in the first evaluation, but the optimization increases the value.

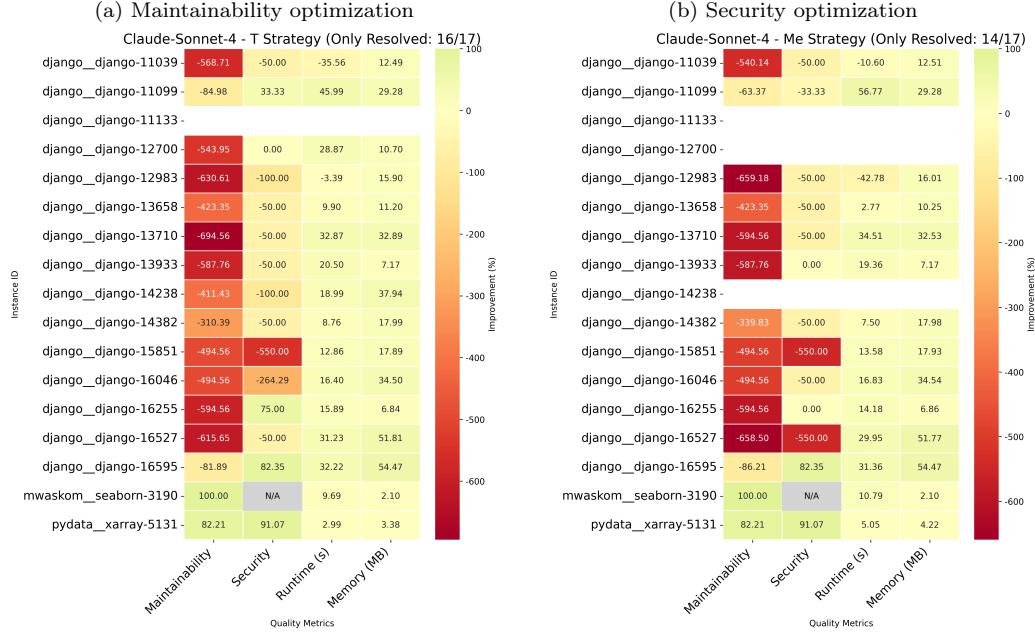
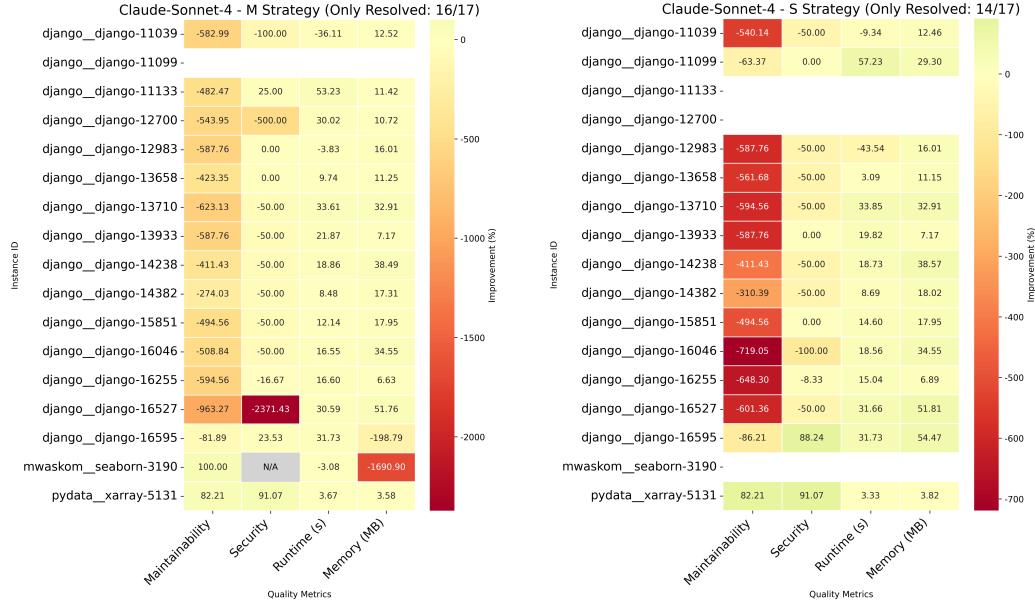


Figure 14: Difference heatmaps for Claude-Sonnet-4 under four NFQC-specific optimizations. Each heatmap shows the change in metric values between the first stage resolved patches and the regenerated patches.

```
INFO - Creating container for django_django-13710...
INFO - Container for django_django-13710 created...
INFO - Container for django_django-13710 started...
INFO - Intermediate patch for django_django-13710 written to logs
      now applying to container...
INFO - Failed to apply patch to container: git apply --verbose
INFO - Failed to apply patch to container: git apply --verbose --
      reject
INFO - Failed to apply patch to container: patch --batch --fuzz=5 -
      p1 -i
INFO - >>>> Patch Apply Failed:
patch: **** malformed patch at line 15:      def media(self):
patching file django/contrib/admin/options.py
```

Figure 15: Example of malformed patch error during patch application.