

FEDERAL UNIVERSITY OF SÃO CARLOS

Vitor Hugo Guilherme

**AN INITIAL INVESTIGATION OF CHATGPT UNIT TEST GENERATION
CAPABILITY**

São Carlos, São Paulo

2023

Vitor Hugo Guilherme

**AN INITIAL INVESTIGATION OF CHATGPT UNIT TEST GENERATION
CAPABILITY¹**

Undergraduate thesis presented as a partial requirement to obtain the title of Computer Engineer by the Federal University of São Carlos.

Supervisor: Prof. Auri Marcelo Rizzo Vincenzi

São Carlos, São Paulo

2023

¹ O conteúdo deste Trabalho de Conclusão e Curso é o mesmo do artigo sobre o mesmo tema aceito para publicação: Guilherme, Vitor Hugo, e Auri M. R. Vincenzi. "An initial investigation of ChatGPT unit test generation capability". Em 8th Brazilian Symposium on Systematic and Automated Software Testing -- SAST'2023. Campo Grande, MS: ACM Press, 2023.



FUNDAÇÃO UNIVERSIDADE FEDERAL DE SÃO CARLOS

DEPARTAMENTO DE COMPUTAÇÃO - DC/CCET

Rod. Washington Luís km 235 - SP-310, s/n - Bairro Monjolinho, São Carlos/SP, CEP 13565-905

Telefone: (16) 33518231 - <http://www.ufscar.br>

DP-TCC-FA nº 13/2023/DC/CCET

Graduação: Defesa Pública de Trabalho de Conclusão de Curso

Folha Aprovação (GDP-TCC-FA)

FOLHA DE APROVAÇÃO

VITOR HUGO GUILHERME

AN INITIAL INVESTIGATION OF CHATGPT UNIT TEST GENERATION CAPABILITY

Trabalho de Conclusão de Curso

Universidade Federal de São Carlos – Campus São Carlos

São Carlos, 29 de agosto de 2023

ASSINATURAS E CIÊNCIAS

Cargo/Função	Nome Completo
Orientador	Auri Marcelo Rizzo Vincenzi
Membro da Banca 1	André Takeshi Endo
Membro da Banca 2	Delano Medeiros Beder



Documento assinado eletronicamente por **André Takeshi Endo, Professor(a) do Ensino Superior**, em 29/08/2023, às 15:23, conforme horário oficial de Brasília, com fundamento no art. 6º, § 1º, do [Decreto nº 8.539, de 8 de outubro de 2015](#).



Documento assinado eletronicamente por **Delano Medeiros Beder, Professor(a) do Ensino Superior**, em 29/08/2023, às 15:28, conforme horário oficial de Brasília, com fundamento no art. 6º, § 1º, do [Decreto nº 8.539, de 8 de outubro de 2015](#).



Documento assinado eletronicamente por **Auri Marcelo Rizzo Vincenzi, Docente**, em 01/09/2023, às 19:14, conforme horário oficial de Brasília, com fundamento no art. 6º, § 1º, do [Decreto nº 8.539, de 8 de outubro de 2015](#).



A autenticidade deste documento pode ser conferida no site <https://sei.ufscar.br/autenticacao>, informando o código verificador **1167421** e o código CRC **F45B474C**.

ABSTRACT

Context: Software testing plays a crucial role in ensuring the quality of software, but developers often disregard it. The use of automated testing generation is pursued with the aim of reducing the consequences of overlooked test cases in a software project.

Problem: In the context of Java programs, several tools can completely automate generating unit test sets. Additionally, there are studies conducted to offer evidence regarding the quality of the generated test sets. However, it is worth noting that these tools rely on machine learning and other AI algorithms rather than incorporating the latest advancements in Large Language Models (LLMs). **Solution:** This work aims to evaluate the quality of Java unit tests generated by an OpenAI LLM algorithm, using metrics like code coverage and mutation test score. **Method:** For this study, 33 programs used by other researchers in the field of automated test generation were selected. This approach was employed to establish a baseline for comparison purposes. For each program, 33 unit test sets were generated automatically, without human interference, by changing Open AI API parameters. After executing each test set, metrics such as code coverage, mutation score, and success rate of test execution were collected to evaluate the efficiency and effectiveness of each set. **Summary of Results:** Our findings revealed that the OpenAI LLM test set demonstrated similar performance across all evaluated aspects compared to traditional automated Java test generation tools used in the previous research. These results are particularly remarkable considering the simplicity of the experiment and the fact that the generated test code did not undergo human analysis.

Keywords: software testing, experimental software engineering, automated test generation, coverage testing, mutation testing, testing tools

FIGURES LIST

Figure 1 - Experiment design diagram	21
Figure 2 - Prompt version 1 for test set generation	26
Figure 3 - Prompt version 2 for test set generation	29
Figure 4 - Example of JUnit test suite for Max program.	35

TABLES LIST

Table 1 - Static information of the Java programs	24
Table 2 - Tools version and purpose	26
Table 3 - Average Data for Each Temperature Parameter -- Prompt version 1	28
Table 4 - Average Data for Each Temperature Parameter -- Prompt version 2	30
Table 5 - Number of successful tests per temperature per project	31
Table 6 - All LLM test sets versus baseline test sets.	35

SUMMARY

1 INTRODUCTION	13
2 CONTEXT	13
3 BACKGROUND	15
3.1 SOFTWARE TESTING	15
3.2 TRADITIONAL AUTOMATIC TEST DATA GENERATION	16
3.3 LLM AND SOFTWARE ENGINEERING	16
3.4 LLM FOR AUTOMATIC TEST DATA GENERATION	17
4 RELATED WORK	18
5 EXPERIMENT DESIGN	20
5.1 LLM SELECTION	22
5.2 PROGRAM SELECTION	23
5.3 TOOLS SELECTION	25
6 DATA COLLECTION AND ANALYSIS	26
7 DISCUSSION	34
8 THREATS TO VALIDITY	37
9 CONCLUSION	38

1 INTRODUCTION

In recent years, the field of artificial intelligence (AI) has witnessed remarkable advancements, largely attributed to the emergence of Large Language Models (LLMs) like ChatGPT. These models have revolutionized various domains by demonstrating an exceptional ability to comprehend and generate human-like text. One significant area that has gained traction within the AI community is software testing, a critical process that ensures the quality and reliability of software systems. With the advent of LLMs, the potential to automate certain aspects of software testing has become a subject of exploration. This undergraduate thesis delves into the intriguing realm of software testing by focusing on the unit test generation capability of ChatGPT. The study's findings were presented as an article at *VIII Simpósio Brasileiro de Teste de Software Sistemático e Automatizado*² and are expounded upon in the subsequent chapters.

2 CONTEXT

Unit testing is an essential practice in software development aimed at ensuring the correctness and robustness of individual code units. These tests, typically written by developers, play a crucial role in identifying defects and validating the expected behavior of software components. DevOps pipelines are strongly based on the quality of unit tests. However, manually generating comprehensive unit tests can be a challenging and time-consuming task, often requiring significant effort and expertise. To address these challenges, researchers have explored automated approaches for test generation (ABDI; DEMEYER, 2022) (FERNANDES, 2022), leveraging advanced techniques and tools.

In this work, we focus on evaluating the quality of Java unit tests generated by an OpenAI Large Language Model (LLM), which has demonstrated remarkable capabilities in generating tests across various domains (YUAN et al, 2023) (SCHWEIKL et al, 2022) (SIDDIQ et al, 2023). Our evaluation will utilize three key

² <https://cbsoft2023.ufms.br/sast>

quality parameters: code coverage, mutation score, and build and executing success rate of test sets. Code coverage quantifies the extent to which the tests exercise different parts of the code, indicating the thoroughness of the test suite. On the other hand, the mutation score measures the ability of the tests to detect and kill mutated versions of the code, providing insights into the fault-detection capability (ANDREWS; BRIAND; LABICHE, 2005). Finally, the build and success execution rate measures the reliability of the generated tests.

In order to conduct a thorough and comprehensive analysis, this study will compare the quality of the unit tests generated by the selected LLM with those produced by other prominent Java test generation tools, such as EvoSuite³. This comparative evaluation aims to determine if the LLMs can outperform state-of-the-art Java test generation tools and will leverage relevant data from (ARAUJO; VINCENZI, 2020) research to provide a meaningful benchmark for comparison. By assessing the effectiveness and performance of the LLMs against established tools, we can gain valuable insights into their capabilities and potential advantages in generating high-quality unit tests for Java programs.

Therefore, we can summarize these paper's contributions:

- To provide evidence of the quality of LLMs in generating unit test sets for Java programs with respect to their efficiency and efficacy;
- To evaluate the improvements a combination of test sets can achieve over individual test sets with respect to efficiency and efficacy;
- To collect data for supporting further comparison of different LLMs on generating Java unit test sets;
- To develop and to make available a set of artifacts for easing the experimentation for different sets of programs.

The structure of the rest of this paper is as follows: We outline the essential subjects for comprehension of this paper in section 3. In section 4, we touch on other studies that are related to ours and highlight the differences. The design of our experiment, along with our choices of programs and tools, is detailed in Section 5. In Section 6, we display the data we've gathered and the subsequent analysis. A discussion of the outcomes derived from the collected data is provided in Section 7.

³ <https://www.evosuite.org/>

We then discuss potential risks that may affect our experiment results in Section 8. Finally, in Section 9, we wrap up the paper by indicating possible future research directions informed by this study and the data collected.

3 BACKGROUND

This section will explain software testing, automatic test data generation, and large language models so that the rest of the paper can be understood.

3.1 SOFTWARE TESTING

In the sphere of software development, it is crucial to ensure the robustness and reliability of a program. A primary technique used for this goal is software testing, which is a systematic process that checks the functionality and accuracy of a software application. However, with software systems growing increasingly complex and versatile, covering a broadening range of use cases and inputs, makes software testing an arduous task.

To analyze the effectiveness of a test set, various criteria come into play, two of which are lines of coverage and mutation testing (ROPER, 1994). Code coverage entails analyzing the extent to which the test suite exercises the internal structure of a software product like its statements or conditions. The goal is to achieve complete or near-complete coverage to ensure that each statement or each branch in the code has been executed at least once during testing.

On the other hand, mutation testing evaluates the test suite ability to identify and “kill” mutated versions of the software (DeMILLO; LIPTON; SAYWARD, 1978). Mutation testing can be seen as a fault model representation (ANDREWS; BRIAND; LABICHE, 2005). These mutations involve making small syntax changes to the code to simulate potential faults. A successful mutation test is one in which the test suite effectively detects these mutations, highlighting its proficiency in identifying vulnerabilities and potential issues within the software. Both code coverage and mutation testing are used in this study as metrics to measure the reliability and

thoroughness of the automatically generated test sets. Moreover, these are traditional metrics used in other studies, like the one developed by Araujo and Vincenzi (2020) which we will use as a baseline.

3.2 TRADITIONAL AUTOMATIC TEST DATA GENERATION

The automatic generation of test data poses an undecidable problem from a computational perspective. While random testing or search-based strategies are commonly employed, other research has shown that the problem remains unsolved when using traditional tools that rely on these approaches (VINCENZI et al, 2016) (ARAUJO; VINCENZI, 2020). The shortcomings of traditional test data generators become apparent when attempting to achieve all testing objectives, such as complete code coverage or eliminating all mutants (ABDI; DEMEYER, 2022). Consequently, the pursuit of comprehensive and efficient test data generation techniques continues to be an ongoing challenge in the dynamic field of software testing.

Even considering the state-of-the-art unit testing generation for Java, (FRASER; ARCURI, 2014), the resultant test set reaches low mutation scores in traditional competitions (VOGL et al, 2021) (SCHWEIKL et al, 2022). Other tools have been discontinued, like Palus (ZHANG, 2011) and JTEExpert (SAKTI et al, 2015), which also employed search-based algorithms. And there are also tools that employ a random generation approach like Randoop (PACHECO; ERNST, 2007) but are still being improved.

3.3 LLM AND SOFTWARE ENGINEERING

LLMs, like ChatGPT⁴, are state-of-the-art language models based on the Transformer architecture (VASWANI et al, 2017). They are designed to process and understand human language, enabling machines to generate coherent and contextually relevant text. These models have been trained on vast amounts of language data, allowing them to capture intricate patterns and relationships in

⁴ <https://chat.openai.com/>

language usage. As a result, they demonstrate impressive capabilities in tasks such as text generation, translation, question-answering, and even software-related activities.

Ma et al (2023) conduct a comprehensive exploration of ChatGPT's applicability and its potential in the software engineering field. The authors examine various tasks, including code generation, code summarizing, bug detection, and code completion, to evaluate the performance of ChatGPT. Through a rigorous investigation and comparison with existing software engineering tools and techniques, the study reveals both the strengths and limitations of ChatGPT in different software engineering scenarios. The findings provide valuable insights into the capabilities of ChatGPT and offer guidance for leveraging its potential to improve software development practices while also highlighting areas where further advancements are needed.

White et al (2023) also explore the potential applications of ChatGPT in various software engineering tasks. The researchers introduce a collection of prompt patterns specifically designed to leverage ChatGPT's language generation capabilities for code quality improvement, refactoring, requirements elicitation, and software design tasks. Through experiments and case studies, they demonstrate the effectiveness of using ChatGPT with these prompt patterns in aiding developers and software engineers in their day-to-day activities. The article highlights the versatility of ChatGPT as a tool for supporting software engineering practices and fostering better code development and design.

3.4 LLM FOR AUTOMATIC TEST DATA GENERATION

The related work section explores the possibility of leveraging LLMs for automatic test data generation. These studies involved exploratory investigations into the use of LLMs for generating test data across different testing phases, ranging from unit testing to end-to-end testing. Notably, the context provided to the LLM was the only aspect that changed during these experiments.

In the case of unit testing, the LLM was presented with code snippets as input (LI et al, 2023) (YUAN et al, 2023) (SIDDIQ et al, 2023) (XIE et al, 2023). For instance, a prompt could be formulated as follows:

“Given the code snippet provided, please generate test cases to cover all possible scenarios and branches within the code.”

The LLM then utilized its language generation capabilities to produce comprehensive test data sets that catered to various testing scenarios. On the other hand, for end-to-end testing, the LLM was supplied with a description of the system's functional specifications (RIBEIRO, 2023) or a GUI (LI et al, 2023). The prompt may have asked the LLM to:

“Generate test cases that validate the entire system's functionality based on the provided functional specification.”

The results of these exploratory studies demonstrated the promising potential of LLMs in automating the test data generation process, streamlining testing efforts, and enhancing software quality. By tailoring the input context to the LLM's capabilities, it was possible to obtain effective test cases for different testing phases, further showcasing the versatility and adaptability of LLMs in software testing.

4 RELATED WORK

The field of test generation has witnessed significant advancements in recent years, with researchers exploring innovative approaches to automate the process and enhance software quality assurance practices. Among these emerging techniques, one notable area of exploration is the use of LLMs for test generation. This section provides a comprehensive overview of the existing literature investigating the application of these powerful tools in test generation.

Li et al (2023) introduce a novel approach to detecting failure-inducing test cases using ChatGPT. By leveraging the model's ability to understand natural

language and generate coherent responses, the researchers propose an interactive debugging technique that allows developers to converse with ChatGPT to identify test cases that are likely to trigger failures. Through experiments on real-world software projects, they demonstrate the effectiveness of their approach in improving fault localization and aiding in the debugging process, highlighting its potential to enhance software quality assurance practices.

Yuan et al (2023) explore the application of ChatGPT for automating unit test generation. The researchers evaluate the performance of ChatGPT in generating meaningful and effective unit tests by comparing them with existing test-generation tools. They also propose a novel approach to enhance ChatGPT's ability to generate high-quality unit tests by incorporating reinforcement learning techniques. Through rigorous experimentation and evaluation of various code bases, the authors demonstrate the potential of ChatGPT as a promising tool for automating the labor-intensive task of unit test generation, highlighting its ability to improve software testing efficiency and accuracy.

Siddiq et al (2023) investigate the efficacy of LLMs, specifically GPT-3, in the generation of unit tests for software programs. The study explores the ability of GPT-3 to understand the requirements of software functionalities and generate relevant test cases. The authors analyze the quality, diversity, and coverage of the generated unit tests through experiments conducted on real-world projects, comparing them with manually written tests. The findings highlight the potential of large language models in automated unit test generation but also reveal certain limitations and challenges that need to be addressed for more effective and reliable results. The research contributes to understanding the capabilities and limitations of large language models in the context of unit testing and provides insights for further advancements in this area.

Xie et al (2023) present ChatUniTest, an automated unit test generation tool that leverages ChatGPT. The tool allows developers to interact with ChatGPT in a conversational manner to generate unit tests for their code. By formulating test generation as a dialogue-based problem, developers can provide natural language prompts to ChatGPT, which then responds with relevant test case suggestions. The authors discuss the implementation details of ChatUniTest and evaluates its effectiveness through experiments on open-source projects. The results demonstrate

that ChatUniTest successfully generates meaningful unit tests, assisting developers in improving software quality and productivity. The study highlights the potential of ChatGPT in the context of automated unit test generation and presents an innovative approach for facilitating the software testing process.

Liu et al (2023) explore the application of GPT-3 for automated GUI testing in the context of mobile applications. The study proposes an innovative approach where GPT-3 is utilized as a conversational agent to interact with mobile apps and generate test cases. A series of experiments conducted on various real-world mobile apps demonstrate the feasibility of GPT-3 in performing human-like GUI testing. The approach achieves high code coverage and successfully detects critical issues, showcasing the potential of leveraging GPT-3 for efficient and effective automated GUI testing of mobile applications. The findings highlight the capabilities of GPT-3 in the domain of mobile app testing, opening avenues for further advancements in automated testing techniques.

Considering the studies carried out so far, the majority explores the use of ChatGPT in interactive mode. We intend to investigate the ChatGPT test generation capability fully automated, without human intervention, interacting, or correcting test cases, considering a possible scenario of no-touch testing (ABDI, 2022) (FERNANDES, 2022).

The next section presents our experiment design to perform this initial investigation.

5 EXPERIMENT DESIGN

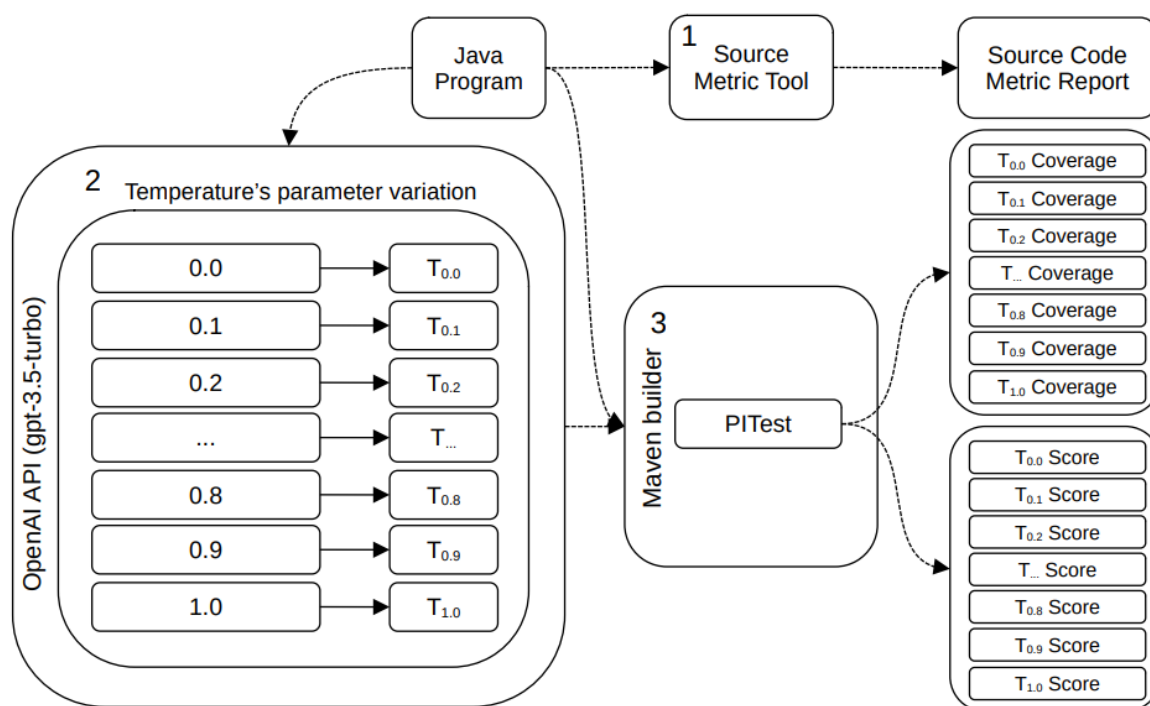
This paper presents an evaluation of the quality of automatically generated test sets by an LLM. A set of Java programs was carefully selected to accomplish this, and multiple JUnit⁵ test sets were generated using the LLM. We use the OpenAI API⁶ (Application Programming Interface) and develop a Python script for interacting with the model via API. The generated test sets will be evaluated based on code coverage, mutation score, and build and execution success rate using selected tools.

⁵ <https://junit.org/>

⁶ <https://openai.com/blog/openai-api>

The collected data will be summarized and analyzed using simple statistics to compare test sets generated by the LLMs with the ones generated by other automated test-generation tools. Figure 1 illustrates the experiment workflow and the steps involved in the evaluation process.

Figure 1: Experiment design diagram



Source: prepared by the author (2023)

Each dashed arrow indicates the input for the subsequent step. Initially, in the first step, we compute some static metrics from the Java source code using JavaNCSS⁷ metric tool.

In the second step, we provide the Python script with a personalized prompt with a program under testing and a “temperature” parameter, considering the OpenAI gpt-3.5-turbo model. This step results in the generation of 33 test sets per program, 3 for each temperature.

Subsequently, a program and its test sets are submitted to the PITest tool in the third step. The PITest tool generates all its mutants for the program under testing

⁷ <https://javancss.github.io/>

and executes each test set against its set of mutants, producing a comprehensive report that includes the mutation score and code coverage for each test set.

Below we comment on some decisions for experiment execution.

5.1 LLM SELECTION

The field of large language models has witnessed remarkable progress, with new ones being developed almost daily. OpenAI, one of the leading organizations in the domain, has been at the forefront of LLM development. In this paper, we choose to leverage the power of OpenAI's gpt-3.5-turbo model due to its availability as a free model and its association with ChatGPT, making it the most used model by final users.

It is worth mentioning that OpenAI had previously introduced a code generation-focused model named davinci-code. However, this model has been discontinued, making gpt-3.5-turbo the preferable option for code-related tasks in our study (OPENAI, 2023).

An important thing about OpenAI API is the temperature parameter. It is a feature that allows users to control the level of randomness and creativity in the generated text. The temperature value can be adjusted when using the API to influence the output's diversity and exploration.

Higher temperatures, such as 1.0, encourage more randomness in the generated text, resulting in imaginative and varied responses. On the other hand, a lower temperature, like 0.2, produces more focused and deterministic output, favoring predictable and conservative responses. By adjusting the temperature parameter, users can fine-tune the balance between generating creative and coherent text, enabling them to obtain the desired level of output for their specific application or task.

We conducted the experiment using the range of temperature values to investigate the variation in the results we will obtain. By exploring the entire spectrum of available temperature values, we aimed to identify the most suitable setting that would yield the best results for our specific test generation requirements. Because of the randomness of the model, especially with higher temperature values, we chose to generate 3 test sets for each temperature value and use the average results.

5.2 PROGRAM SELECTION

We decided to use the results from Araujo and Vincenzi (2020)'s work as a baseline. They used a set of 33 Java programs and conducted an experiment investigating the capability of four different automatic testing generators (EvoSuite, Palus, JTEExpert, and Randoop) for Java on covering code and killing mutants using PITest as the mutation tool.

Therefore, we selected the same set of programs to perform our experiment⁸. By comparing the test sets produced by GPT-Turbo-3.5 with those generated by these tools, our research aims to provide valuable insights into the effectiveness and efficacy of LLMs in automated unit test generation. Table 1 shows the selected programs and their characteristics.

⁸ We would like to thank Araujo and Vincenzi for making the set of programs, scripts, and data available to ease the comparison.

Table 1: Static information of the Java programs

ID	Program	#Classes	#Methods	NCSS	CCN	CCA	#Req	#Mut
1	Max	1	1	8	3	3,0	4	14
2	MaxMin1	1	1	13	4	4,0	8	21
3	MaxMin2	1	1	14	4	4,0	8	21
4	MaxMin3	1	1	32	9	9,0	16	61
5	Sort1	1	1	11	4	4,0	10	21
6	FibRec	1	1	8	2	2,0	6	12
7	FibIte	1	1	8	2	2,0	6	12
8	MaxMinRec	1	1	26	5	5,0	13	41
9	Mergesort	1	2	22	6	4,0	16	56
10	MultMatrixCost	1	1	18	6	6,0	14	75
11	ListArray	1	4	20	3	1,8	12	29
12	ListAutoRef	2	4	23	2	1,3	12	21
13	StackArray	1	5	20	3	1,8	12	27
14	StackAutoRef	2	5	27	3	1,4	17	27
15	QueueArray	1	5	24	3	2,0	19	40
16	QueueAutoRef	2	5	32	3	1,6	23	32
17	Sort2	2	7	74	6	3,4	49	141
18	HeapSort	1	9	59	5	2,7	40	116
19	PartialSorting	1	10	62	5	2,5	42	120
20	BinarySearch	1	4	32	8	3,5	21	55
21	BinaryTree	2	11	85	7	3,0	48	145
22	Hashing1	2	10	61	5	2,1	35	88
23	Hashing2	2	12	88	7	3,2	51	162
24	GraphMatAdj	1	9	60	5	2,9	42	134
25	GraphListAdj1	3	16	66	4	1,6	34	95
26	GraphListAdj2	2	14	88	6	2,2	51	113
27	DepthFirstSearch	3	16	65	4	1,6	33	94
28	BreadthFirstSearch	3	16	65	4	1,6	33	94
29	Graph	3	16	65	4	1,6	33	94
30	PrimAlg	1	5	40	7	2,6	31	71
31	ExactMatch	1	4	55	8	6,3	40	205
32	AproximateMatch	1	1	24	7	7,0	19	88
33	Identifier	1	3	30	9	7,7	22	114
Avg		1,5	6,1	40,2	4,9	3,3	24,8	73,9
SD		0,7	5,2	25,4	2,0	2,0	14,7	50,3

Source: adapted from Araujo and Vincenzi (2020)

For each program Araujo and Vincenzi (2020) computed the following metrics:

- Non-Commenting Source Statements (NCSS);
- Cyclomatic Complexity Number (CCN);
- Cyclomatic Complexity Average (CCA);
- Number of test cases on each test set: EvoSuite (E), JTEpert (J), Palus (P) and Randoop (R);

- Number of requirements demanded to cover statement coverage (Req); and
- Number of generated mutants considering all mutation operators available in PITest (Mut).

As can be observed, they are not complex programs but once we are working at unit testing levels, we understand that each program provides units with sufficient complexity, equivalent to units present in other real programs. In terms of lines of code, the average size is around 40, and cyclomatic complexity is around 4.9.

To simplify the experiment we make the assumption that all programs are correct, and any mutations on the source code generates an incorrect version (with the exception of equivalent mutants). Therefore, if a mutation breaks a test, we infer that the mutant is killed, avoiding the oracle problem.

5.3 TOOLS SELECTION

To evaluate the results, we used JUnit as a unit testing framework, which is widely recognized as the industry standard for testing and generating comprehensive reports. Another noteworthy aspect is the utilization of JUnit in the study of Araujo and Vincenzi (2020), which is a valuable reference point for comparing our results.

By employing the same testing framework, we establish a meaningful basis for comparison, enabling us to analyze and assess the effectiveness of our LLM-generated tests in relation to their findings. The same logic was used to select PITest⁹ as our mutation tool. All the mutation operators of PITest¹⁰ were selected to cover every possible change in the source code.

Therefore, Table 2 summarizes the tools and versions we used, which we kept the same as the ones adopted by Araujo and Vincenzi (2020) to minimize threats. We present the data we collected and some analysis we carried out so far.

⁹ <https://pitest.org/>

¹⁰ <https://pitest.org/quickstart/mutators/>

Table 2: Tools version and purpose

Tool	Version	Purpose
JavaNCSS	32.53	Static Metric Computation
PITest	1.3.2	Mutation and Coverage Testing
Maven	3.6.3	Application Builder
JUnit	4.12	Framework for Unit Testing
Python	3.7	Script language
Java	8	Programs language
LLM	gpt-3.5-turbo	OpenAI LLM for generating tests

Source: adapted from Araujo and Vincenzi (2020)

6 DATA COLLECTION AND ANALYSIS

The initial step involved creating a centralized repository storing all the selected programs, scripts, and experimental results. To achieve version control and facilitate seamless collaboration, we opted for GitHub as our hosting platform¹¹. Once we selected the programs from Araujo and Vincenzi (2020) the static metrics from their work were used without recomputing it. Table 1 presents such data about the programs. Subsequently, we proceeded with the test set generation using the *gpt-3.5-turbo* model. To accomplish this, we formulated a specific base prompt designed to request the model's assistance in generating JUnit unit tests tailored for a program. The first prompt version is as follows:

Figure 2: Prompt version 1 for test set generation

Generate test cases just for the {cut}
Java class in one Java class file with
imports using JUnit 4 and Java 8:
{code}

Source: prepared by the author (2023)

In the prompt above, {cut} is a variable that represents the name of the class under testing, and {code} is a variable containing the code of the class under testing and its dependencies.

¹¹ <https://github.com/vitor0x5/initial-investigation-chatgpt-unit-tests>

To automate the process, we developed a Python script (`generate-chatgpt.py`), which sends the request to the OpenAI API. Upon receiving the response from the API, the script removes any natural language comments that the LLM model added before or after the generated code. Additionally, the script ensured that the Java test class name matched the file name, following a pattern to enhance test data organization. As an output, the script generates 33 Java test classes for every selected program, with 3 test classes for each LLM temperature value, as mentioned in Section 5.1.

Then, with all tests generated for every program, we build and run them using Maven. To automate this process, we developed another Python script `compile-and-test-chatgpt.py`. However, at this stage, we encountered an issue where some tests generated by the model do not build successfully due to problems such as syntax errors and missing imports. The script discards any test set with failing test cases.

The script moves all test files to a directory outside the project, copies one test file at a time to the project's test directory, and then builds and runs the test for that specific file. In this manner, any build issues or errors in one test won't affect the others, ensuring a smoother and more effective testing process.

Finally, we developed the last Python script `reports-chatgpt.py` for extracting coverage and mutation score from PITest reports. It generates one CSV file for each Java program, including all test results that are executed successfully. Tables 3 and 4 present parts of the collected data.

Considering the first prompt version, presented in Figure 3, Table 3 presents average data for each temperature value we investigate.

Table 3: Average test data for each temperature parameter -- Prompt version 1

Temp.	# of Suc. Test	% of Suc.	AVG Cov.	AVG Score
0.0	37	37.4	83.0	51.3
0.1	37	37.4	85.7	51.6
0.2	37	37.4	84.6	52.3
0.3	38	38.4	86.1	53.9
0.4	39	39.4	86.9	53.4
0.5	35	35.4	88.3	53.6
0.6	52	52.5	83.9	54.4
0.7	36	36.4	88.9	54.8
0.8	45	45.5	87.6	54.2
0.9	42	42.4	81.5	49.2
1.0	41	41.4	81.8	52.9

Source: prepared by the author (2023)

Observe that the average results in Table 3 show that from a possible total of 99 test sets for each temperature (3 for each of 33 programs), the temperature most effective on generating successful test sets is 0.6. With this temperature, 52 out of 99 test sets run correctly with no errors, with a successful rate of 52.5%.

Table 3 also shows that quantity does not mean high-quality tests. Temperature 0.7 reaches 36.4% of the successful rate of test sets, around 16% less than temperature 0.6, but with 36 test sets, the average coverage and mutation score are the highest: 88.9 and 54.8, respectively.

Although we consider these results impressive due to the simplicity of the prompt, we analyzed the errors produced by the test sets and the parts of the source code not covered by the tests, and we tried to improve the prompt to mitigate some problems found. Figure 3 shows the prompt's second version.

Observe that in the prompt presented in Figure 3, {cut}, {clazz} and {code} have the same meaning, the name of the class under testing and the source code of the class under testing and its dependencies. We were more incisive regarding how we wanted the test set. Including mandatory dependencies, timeout, throws Exception, test set name, and also the calling of void methods and default constructors. We also enforce two testing criteria: decision coverage and boundary values.

Figure 3: Prompt version 2 for test set generation

I need functional test cases to cover all decisions in the methods of the class under testing.

All conditional expressions must assume true and false values.

Tests with Boundary Values are also mandatory. For numeric data, always use positive and negative values.

All tests must be in one Java class file.

Include all necessary imports.

It is mandatory to throws Exception in all test method declarations.

It is mandatory to include timeout=1000 in all @Test annotations.

It is mandatory a test for the default constructor.

Each method in the class under test must have at least one test case.

Even simple or void methods must have a test calling it with valid inputs.

@Test(expected= must be used only if the method under testing explicitly throws an exception.

Test must be in JUnit 4 framework format.

Test set heather package and import dependencies:

```
package ds;
import org.junit.Test;
import org.junit.Before;
import static org.junit.Assert.*;
import ds.*;
```

The class under testing is { clazz }.

The test class must be { cut }Test.

Class under testing

```
*****
{code}
```

Source: prepared by the author (2023)

After this prompt upgrade, we rerun all scripts to generate new test sets, check their quality, and measure coverage and mutation scores. Table 4 presents the

average data per temperature. The new prompt improved the test set's successful execution by more than 12%, observing temperature 0.2, 64 out of 99 test sets executed without failures, a successful rate of 64.6%. We also improved coverage and mutation scores to an average of 93.5% and 58.8%, respectively.

Table 4: Average Data for Each Temperature Parameter -- Prompt version 2

Temp.	# of Suc. Test	% of Suc.	AVG Cov.	AVG Score
0.0	61	61.6	93.5	58.8
0.1	59	59.6	93.4	57.4
0.2	64	64.6	90.7	57.4
0.3	63	63.6	91.2	57.7
0.4	59	59.6	92.0	57.8
0.5	55	55.6	93.3	57.7
0.6	63	63.6	88.0	55.9
0.7	54	54.5	89.9	55.4
0.8	55	55.6	88.6	55.3
0.9	54	54.5	85.8	54.1
1.0	61	61.6	87.7	54.1

Source: prepared by the author (2023)

Based on this data, we decided to detail the analysis per program and temperature to verify if each temperature has similar behavior for each program. Table 5 presents the data. The first thing we observed in the last two lines of the table is that, in general, the lower the temperature value, the greater the number of programs without successful test sets.

In the worst case, for temperature 0.0, 12 programs out of 33 (36,4%) have no test set running successfully. In the best case, temperature 1.0, 3 out of 33 programs (9.1%) have no test set running successfully. We tried to investigate the reasons, especially for these three programs, why it fails to generate successful runnable test sets. The general observation is that, for these specific programs, they define an *Item* interface and a *MyItem* class implementing the interface, but this class did not override *compareTo* and *equals* methods from *Object* class in Java. Nevertheless, ChatGPT seems to assume they are available for object comparison once several tests make use of object comparison, but they check reference equality and not object field contents, failing the test cases. This is the main reason all tests for programs 10, 18, and 19 have no test set available, independently of the temperature's parameter.

Table 5: Number of successful tests per temperature per project

ID	Program	Temperature											
		0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0	All
1	Max	0	1	0	1	2	0	1	1	1	1	2	10
2	MaxMin1	2	1	2	3	3	2	2	2	0	3	2	22
3	MaxMin2	3	3	3	2	2	2	3	3	2	2	2	27
4	MaxMin3	0	0	0	0	0	0	0	0	0	0	1	1
5	Sort1	3	3	3	3	3	3	3	3	3	3	3	33
6	FibRec	0	0	1	0	1	2	2	2	2	3	2	15
7	FibIte	0	0	0	1	2	1	2	0	3	1	2	12
8	MaxMinRec	3	3	3	3	2	2	3	2	3	2	3	29
9	Mergesort	3	3	3	3	3	3	3	3	3	3	3	33
10	MultMatrixCost	0	0	0	0	0	0	0	0	0	0	0	0
11	ListArray	3	3	3	3	3	2	3	3	2	2	2	29
12	ListAutoRef	3	3	3	2	3	3	3	3	3	2	1	29
13	StackArray	3	3	3	3	3	3	3	3	3	3	2	32
14	StackAutoRef	3	3	3	3	2	3	3	2	3	2	3	30
15	QueueArray	0	0	1	2	0	2	3	3	3	2	2	18
16	QueueAutoRef	3	3	3	3	3	2	3	2	2	3	1	28
17	Sort2	0	0	1	0	0	0	0	0	0	1	1	3
18	HeapSort	0	0	0	0	0	0	0	0	0	0	0	0
19	PartialSorting	0	0	0	0	0	0	0	0	0	0	0	0
20	BinarySearch	3	2	3	2	3	2	2	3	0	0	2	22
21	BinaryTree	3	3	3	3	3	2	3	2	2	3	3	30
22	Hashing1	3	3	3	3	1	2	1	2	2	2	2	24
23	Hashing2	0	1	1	1	0	1	0	0	2	1	1	8
24	GraphMatAdj	3	3	3	3	3	3	3	0	2	2	3	28
25	GraphListAdj1	3	3	2	3	3	1	2	1	1	3	2	24
26	GraphListAdj2	3	3	3	3	3	3	3	2	3	3	2	31
27	DepthFirstSearch	3	2	3	3	2	2	3	2	2	0	2	24
28	BreadthFirstSearch	3	2	3	1	1	1	0	1	1	0	3	16
29	Graph	2	2	2	3	2	2	2	2	1	1	1	20
30	PrimAlg	0	0	0	0	0	0	1	1	0	0	1	3
31	ExactMatch	3	3	3	3	3	3	3	3	3	3	3	33
32	AproximateMatch	3	3	3	3	3	3	3	2	3	2	3	31
33	Identifier	0	0	0	0	0	0	0	1	0	1	1	3
# Successful Test		61	59	64	63	59	55	63	54	55	54	61	648
% Successful Test		61.6	59.6	64.6	63.6	59.6	55.6	63.6	54.5	55.6	54.5	61.6	59.5
# of Programs without test		12	10	8	8	9	8	8	8	9	8	3	3
% of Programs without test		36.4	30.3	24.2	24.2	27.3	24.2	24.2	24.2	27.3	24.2	9.1	9.1

Source: prepared by the author (2023)

Also inspired by Araujo and Vincenzi (2020), who observed that by merging test sets from EvoSuite, Palus, JTEExpert, and Randoop, the resultant merged test set performs better than any other individual test set in terms of coverage and mutation score, we decided to create a merged test set considering the test sets provided by different temperatures. Moreover, in our case, by merging all test sets, only 3 out of our 33 programs will remain without valid tests. The last column of Table 5 presents the number of valid tests for each program. Only for two programs (9 - Mergesort and

31 - ExactMatch) we got the maximum number of 33 valid tests, 3 for each different temperature value.

Then, we use the JUnit test suite to create a test suite corresponding to all successful test sets. Figure 4 presents an example of a JUnit test suite, considering the 10 successful test sets for program 1 - Max. We built a *ds.All.java* test suite file for each program and used it to collect coverage and mutation scores for all programs. The collected data is shown in Table 6.

Figure 4: Example of JUnit test suite for Max program.

```

1 package ds;
2 import org.junit.runner.RunWith;
3 import org.junit.runners.Suite;
4
5 @RunWith(Suite.class)
6 @Suite.SuiteClasses({MaxTest2.class,MaxTest5.class, MaxTest8.class,
7   MaxTest9.class, MaxTest10.class, MaxTest14.class, MaxTest18.class,
8   MaxTest20.class, MaxTest22.class, MaxTest27.class })
9 public class All { }

```

Source: prepared by the author (2023)

In the two last columns of Table 6, we present the best results that Araujo and Vincenzi (2020) obtained considering the merged test set in their experiment. We will refer to our merged test set as LLM Suite and Araujo and Vincenzi (2020)'s merged test set as Baseline Suite. We highlight in gray the cells with the best values with respect to the coverage or mutation score of each merged test set.

Regarding code coverage, LLM Suite did not reach Baseline Suite results in 6 out of 33 programs (10, 18, 19, 21, 23, and 26). As already mentioned, for three of these 6 programs (10, 18, and 19), ChatGPT was unable to create runnable without-fail tests, and we got zero coverage. For all the other programs, both suites covered all program source code. On average, LLM Suite coverage is 90.2% and Baseline Suite coverage is 99.5%. If we remove programs 19, 18, and 19 from the analysis, Baseline Suite keeps the same coverage of 99.5%, but LLM Suite coverage reaches 99.2%, almost the same.

The biggest surprise occurred with the mutation score. As can be observed, for 14 out of 33 programs, LLM Suite overcomes the mutation score of Baseline Suite, and in some cases, it improves by more than 20% the mutation score, like in programs 1, 20, and 30. On the other hand, Baseline Suite reaches better scores for 17 out of 33 programs, and for two programs, 6 and 7, we have a tie. On average, the average mutation score for Baseline Suite reaches 78.5%, and for LLM Suite, it is 70.5%. Again, removing programs 10, 18, and 19 from our analysis, we got very similar mutation scores of 77.6 and 79.5 for LLM and Baseline suites, respectively.

7 DISCUSSION

The idea for this work was just investigate the capability of LLM chats, ChatGPT in our experiment, on generating unit test sets but, when we got the first results from these interactions using the very simple prompt presented in Figure 2, we decided to investigate its potential with more emphasis.

The final results presented in Table 6 suggest that these prompts have a very good potential, if not to be used as a single way for unit testing generation, its combination in a coordinated way with traditional automatic testing generators can be very promising. Testing will always be a challenging activity, as many useful tools we have to automate this process better.

Prompts also show us huge flexibility in asking for test cases considering specific testing criteria or asking for test cases to reach a specific objective, like covering a specific statement or killing a specific mutant. In this work, we decided only on a standard predefined prompt, as shown in Figure 2, to use the generated unit testing fully automated, i.e., without interacting with the chat asking for additional testing or testing corrections.

We do not think LLMs will solve all testing problems automatically. We believe a good automated testing strategy now gained important support from LLMs. Our intention is to observe the limits of LLM for unit testing generation. If some important testing requirement is missing, having time and people available for testing, it is possible to develop specialized prompts to solve and generate specific test cases with human support to check and correct possible mistakes. This is especially true

once it is difficult to maintain software testing generators. For instance, considering the ones used by Araujo and Vincenzi (2020), two of them (Palus and JTEExpert) are not available or did not work with new versions of Java.

On the other hand, LLMs just need a huge amount of data to work and can be easily personalized to meet different testing objectives. A possible alternative to improve the LLM capabilities, considering Java programs, for instance, is to use EvoSuite to start the test set generation and, later, to provide to the LMM the source code of the class under testing and also a previously generated EvoSuite test sets. In this way, we suppose the prompt can better understand the test case style, which may reduce the test case failures generated by LLMs.

Table 6: All LLM test sets versus baseline test sets.

ID	LLM Suite		Baseline Suite ⁹	
	Coverage	Score	Coverage	Score
1	100.0	85.7	100.0	64.3
2	100.0	85.7	100.0	83.8
3	100.0	85.7	100.0	84.3
4	100.0	64.5	100.0	79.8
5	100.0	80.0	100.0	78.5
6	100.0	100.0	100.0	100.0
7	100.0	100.0	100.0	100.0
8	100.0	83.3	100.0	83.1
9	100.0	96.4	100.0	95.5
10	0.0	0.0	100.0	45.6
11	100.0	93.5	100.0	78.1
12	100.0	87.0	100.0	83.9
13	100.0	96.8	100.0	81.3
14	100.0	93.5	100.0	85.5
15	100.0	93.0	100.0	90.5
16	100.0	67.6	100.0	82.4
17	100.0	57.6	100.0	68.8
18	0.0	0.0	100.0	73.9
19	0.0	0.0	100.0	84.4
20	100.0	94.7	100.0	70.4
21	81.3	62.5	88.8	93.8
22	100.0	57.3	100.0	93.9
23	98.1	63.3	100.0	86.6
24	100.0	73.4	96.2	69.0
25	100.0	78.9	100.0	81.9
26	98.0	77.2	99.2	78.1
27	100.0	78.7	100.0	81.0
28	100.0	78.7	100.0	82.1
29	100.0	78.7	100.0	81.4
30	100.0	71.1	100.0	42.4
31	100.0	39.5	100.0	58.6
32	100.0	38.4	100.0	51.6
33	100.0	64.0	100.0	75.8
AVG	90.2	70.5	99.5	78.5
SD	29.2	27.5	2.0	13.9
AVG**	99.2	77.6	99.5	79.5
SD**	3.4	16.5	2.1	13.1

** - AVG and SD removing zeros.

Source: prepared by the author (2023)

Another point is that we decided to provide the class under testing source code to LLM and asked it to generate tests for the entire class. Although, we believe if you ask for a test only for a specific method inside a class, we will get better results once the scope is reduced, and LLM will create more tests for each specific method.

Finally, we explore a single OpenAI API model called gpt-3.5-turbo, but OpenAI offers a variety of models, each with different capabilities. Deciding which one is more suitable for each situation demands additional experimentation. Moreover, there are also a lot of new LLMs available like Bing¹², Bard¹³ and LLaMa¹⁴ which may also demand more investigation with respect to their capacity on automatic generating unit testing for specific languages.

8 THREATS TO VALIDITY

There are several potential threats in this paper. One possible threat is sampling bias, which means that the selection of programs and tools used in the experiment may not accurately represent the entire software development landscape. This could lead to biased results that may not be applicable to other contexts. To minimize this threat, we tried to use tools and programs already explored in other experiments. Moreover, especially for the automated test generator, at least EvoSuite is a tool used in a vast number of experiments both in academia (VOGL et al, 2021) (SCHWEIKL et al, 2022) and in industry (FRASER; ARCURI, 2014) and is also integrated into professionals' integrated development environments (ARCURI; FRASER, 2016).

Another threat is the limited generalizability of the findings. The study's conclusions may only be relevant to a specific set of programs and tools and may not be applicable to different scenarios. Additionally, there is a risk of measurement bias, where the metrics used to measure the effectiveness of the generated test data may not fully capture its quality and comprehensiveness. In this way, we manually revise the Python scripts and check the collected data for some programs to ensure the information is accurate. Coverage and mutation score are traditional metrics for

¹² <https://www.bing.com>

¹³ <https://bard.google.com>

¹⁴ <https://labs.perplexity.ai/>

evaluating software testing quality. Mutation is confirmed to be an excellent fault model to evaluate the quality of test sets (ANDREWS; BRIAND; LABICHE, 2005) (JIA; HARMAN, 2011) (JUST et al, 2014). Although we work with Java programs on this initial investigation, other studies in the course also explore the LLM test generation capabilities for programs written in other languages like Python and C, for instance.

In Section 2.4.2, we presume the correctness of all programs, as they are basic and known algorithms. However, there remains the possibility of bugs that could potentially result in inaccurate mutation scores and failure to run correct tests.

The use of large language models for automatic test data generation may have limitations or biases that could impact the quality and comprehensiveness of the generated test sets. Using baseline results obtained from traditional automated test case generators (ARAUJO;VINCENZI, 2020) to confront the results obtained from test sets generated from LLM aims to minimize this threat. Moreover, we only used an LLM engine and model in this experiment, which may not represent the results for other LLMs or models. We intend to extend the experiment for a large number of programs, LLM engines, and models in further studies.

9 CONCLUSION

In this work, we presented an initial investigation of the use of OpenAI API, considering the LLM named *gpt-3.5-turbo*, for unit test generation in a fully automated way, i.e., with no human interaction for test case correction after prompt return. The idea was to detect to which extent the test cases will run directly, with no errors, for testing a set of Java programs.

Basically, we developed a prompt to ask test sets via API, only varying the code of the class under testing and the "temperature" parameter of the *gpt-3.5-turbo* model. We asked for three test sets for each one of the eleven different temperature values (0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0) for each program, resulting, in the best case, in a total of 33 test sets per program.

Our results show that not for all temperatures the API was able to produce useful test sets which run automatically with no error without human intervention. In

this way, we discarded these test sets in our experiment. For 3 out of 33 programs, the model was not able to generate useful test sets for any temperature, especially due to the non-overriding of traditional Java methods for object comparison like *equals()* and *compareTo()* for the application under testing.

We observed interesting results by keeping only test sets that run automatically and comparing our results with those obtained by other researchers that used traditional automated test set generators Araujo and Vincenzi (2020). We considered that, besides the simplicity of the prompt, asking for testing to the LLM, the results in terms of code coverage were very similar to the ones obtained in the baseline. Moreover, with respect to mutation score, we observed complementary aspects between LLM Suite and Baseline Suite. They complement each other.

Further work intends to investigate the best way to use a traditional automated testing generator together with LLM prompts to obtain better results than when using isolated tools.

Moreover, this initial investigation raised more questions than produced answers. To answer the raised questions more experimentation is necessary. A few of them are:

- Do the other OpenAI models produce similar or complementary results?
- Does the language used in the prompt influence the results?
- Does the language of the product under testing influence the results?
- How do other LLMs prompts perform the automation of unit testing generation?
- Does the LLM perform better by asking testing for a method instead of a class?

REFERÊNCIAS

ABDI, Mehrdad; DEMEYER, Serge. Steps towards zero-touch mutation testing in Pharo. In: 21st Belgium-Netherlands Software Evolution Workshop – BENEVOL'2022 (CEUR Workshop Proceedings, Vol. 1). Mons, 2022.

ANDREWS, J. H.; BRIAND, L. C.; LABICHE, Y. Is mutation an appropriate tool for testing experiments?. In: XXVII International Conference on Software Engineering – ICSE'05. ACM Press, St. Louis, MO, USA, 2005. p. 402-411. Disponível em: <https://doi.org/10.1145/1062455.1062530>.

ARAUJO, Filipe Santos; VINCENZI, Auri. How far are we from testing a program in a completely automated way, considering the mutation testing criterion at unit level?. In: Anais do Simpósio Brasileiro de Qualidade de Software (SBQS). SBC, 2020. p. 151-159. Disponível em: <https://doi.org/10.1145/3439961.3439977>.

ARCURI, Andrea; CAMPOS, José; FRASER, Gordon. Unit Test Generation During Software Development: EvoSuite Plugins for Maven, IntelliJ and Jenkins. In: 2016 IEEE International Conference on Software Testing, Verification and Validation (ICST). 2016. p. 401-408. Disponível em: <https://doi.org/10.1109/ICST.2016.44>.

DeMILLO, R. A.; LIPTON, R. J.; SAYWARD, F. G. Hints on Test Data Selection: Help for the Practicing Programmer. IEEE Computer, v. 11, n. 4, p. 34-43, April 1978. Disponível em: <https://doi.org/10.1109/C-M.1978.218136>.

FERNANDES, Leo et al. Put Your Hands In The Air! Reducing Manual Effort in Mutation Testing. In: Proceedings of the XXXVI Brazilian Symposium on Software Engineering (SBES '22, Vol. 1). Association for Computing Machinery, New York, NY, USA, 2022. p. 198-207. Disponível em: <https://doi.org/10.1145/3555228.3555233>.

FRASER, Gordon; ARCURI, Andrea. EvoSuite: automatic test suite generation for object-oriented software. In: Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering (ESEC/FSE '11). ACM, Szeged, Hungary, 2011. p. 416-419. Disponível em: <https://doi.org/10.1145/2025113.2025179>.

FRASER, Gordon; ARCURI, Andrea. A Large-Scale Evaluation of Automated Unit Test Generation Using EvoSuite. ACM Transactions on Software Engineering and Methodology, v. 24, n. 2, Dec. 2014. Disponível em: <https://doi.org/10.1145/2685612>.

JIA, Yue; HARMAN, Mark. An Analysis and Survey of the Development of Mutation Testing. IEEE Transactions on Software Engineering, v. 37, n. 5, Sept. 2011, p. 649-678. Disponível em: <https://doi.org/10.1109/TSE.2010.62>.

JUST, René et al. Are mutants a valid substitute for real faults in software testing?. In: Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2014, Vol. 1). Association for Computing

Machinery, Hong Kong, China, 2014. p. 654-665. Disponível em: <https://doi.org/10.1145/2635868.2635929>.

LI, Tsz-On et al. Finding Failure-Inducing Test Cases with ChatGPT. 2023.

LIU, Zhe et al. Chatting with GPT-3 for Zero-Shot Human-Like Mobile Automated GUI Testing. 2023. Disponível em: <https://doi.org/10.48550/arXiv.2305.09434>.

MA, Wei et al. The Scope of ChatGPT in Software Engineering: A Thorough Investigation. 2023. Disponível em: <https://doi.org/10.48550/arXiv.2305.12138>.

OpenAI. OpenAI GTP-3.5 Models Documentation. July 2023. Disponível em: <https://platform.openai.com/docs/models/gpt-3-5>.

PACHECO, Carlos; ERNST, Michael D. Randoop: Feedback-directed Random Testing for Java. In: Companion to the 22Nd ACM SIGPLAN Conference on Object-oriented Programming Systems and Applications Companion (OOPSLA '07). ACM, 2007. p. 815-816. Disponível em: <https://doi.org/10.1145/1297846.1297902>.

RIBEIRO, Marco Tulio. Testing Language Models (and Prompts) Like We Test Software. May 2023. Disponível em: <https://towardsdatascience.com/testing-large-language-models-like-we-test-software-92745d28a359>.

ROPER, M. Software Testing. McGrall Hill, 1994.

SAKTI, Abdelilah et al. JTEExpert at the Third Unit Testing Tool Competition. 2015. p. 52-55. Disponível em: <https://doi.org/10.1109/SBST.2015.20>.

SCHWEIKL, Sebastian et al. EvoSuite at the SBST 2022 Tool Competition. In: Proceedings of the 15th Workshop on Search-Based Software Testing (SBST '22). Association for Computing Machinery, New York, NY, USA, 2022. p. 33-34. Disponível em: <https://doi.org/10.1145/3526072.3527526>.

SIDDIQ, Mohammed Latif et al. Exploring the Effectiveness of Large Language Models in Generating Unit Tests. 2023. Disponível em: <https://doi.org/10.48550/arXiv.2305.00418>.

VASWANI, Ashish et al. Attention is All You Need. In: Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17). Curran Associates Inc., Red Hook, NY, USA, 2017. p. 6000-6010.

VINCENZI, Auri M. R. et al. The Complementary Aspect of Automatically and Manually Generated Test Case Sets. In: Proceedings of the 7th International Workshop on Automating Test Case Design, Selection, and Evaluation (A-TEST 2016, Vol. 1). ACM, 2016. p. 23-30. Disponível em: <https://doi.org/10.1145/2994291.2994295>. Event-place: Seattle, WA, USA.

VOGL, Sebastian et al. EvoSuite at the SBST 2021 Tool Competition. In: 2021 IEEE/ACM 14th International Workshop on Search-Based Software Testing (SBST). IEEE, 2021. p. 28-29.

WHITE, Jules et al. ChatGPT Prompt Patterns for Improving Code Quality, Refactoring, Requirements Elicitation, and Software Design. 2023. Disponível em: <https://doi.org/10.48550/arXiv.2303.07839>.

XIE, Zhuokui et al. ChatUniTest: a ChatGPT-based automated unit test generation tool. 2023. Disponível em: <https://doi.org/10.48550/arXiv.2305.04764>.

YUAN, Zhiqiang et al. No More Manual Tests? Evaluating and Improving ChatGPT for Unit Test Generation. 2023. Disponível em: <https://doi.org/10.48550/arXiv.2305.04207>.

ZHANG, Sai. Palus: A Hybrid Automated Test Generation Tool for Java. In: Proceedings of the 33rd International Conference on Software Engineering (ICSE'11). Association for Computing Machinery, 2011. p. 1182-1184. Disponível em: <https://doi.org/10.1145/1985793.1986036>. Event-place: Waikiki, Honolulu, HI, USA.

