

Benchmarking Large Language Models for Code Generation

**Sumedh Arun Patil¹, Vedant Hemant Pangudwale²,
Devansh Rakesh Rathi³, Prof. Rupali Kadu⁴**

^{1,2,3,4}Department of Electronics and Telecommunication Engineering, K.J. Somaiya Institute of Technology Mumbai, India

Abstract

As the landscape of software development continues to evolve, the need for efficient and innovative coding practices becomes increasingly apparent. This research endeavors to explore the effectiveness of Large Language Models (LLMs) in code generation, focusing on benchmarking their performance across various coding tasks. Leveraging advanced Natural Language Processing (NLP) techniques and deep learning architectures, our study investigates how LLMs, such as the `codellama-13b-instruct.Q5_K_S.gguf` engine, interpret and generate code from natural language instructions. With an emphasis on accuracy, efficiency, and user accessibility, our research seeks to shed light on the capabilities of LLMs in bridging the gap between human language and executable code. By evaluating factors such as model architecture, training data quality, and task complexity, we aim to provide insights into the potential of LLMs for revolutionizing the coding experience. Through meticulous benchmarking and analysis, this study aims to contribute to the advancement of LLM development and its applications in code generation, paving the way for more efficient and inclusive coding practices in the future.

Keywords: Benchmarking, LLM, Performance Evaluation, Code Generation, Training Data Quality.

1. INTRODUCTION

Large Language Models (LLMs) have emerged as transformative tools in natural language processing, demonstrating exceptional proficiency in tasks ranging from text generation to comprehension. Leveraging pre-trained transformer architectures, such as OpenAI's GPT series, LLMs have become increasingly instrumental in various domains, including software development, particularly in the realm of code generation.

The utilization of LLMs in code generation tasks signifies a paradigm shift in software engineering, promising automation and efficiency gains across diverse development processes. Notably, LLMs exhibit remarkable aptitude in understanding and synthesizing code snippets, enabling applications in unit testing, data cleaning, API calls, and beyond. This expanding role of LLMs in code generation underscores their potential to streamline software development workflows, reduce human effort, and enhance productivity. However, while the potential benefits of employing LLMs for code generation are evident, their efficacy across specific tasks remains contingent upon rigorous evaluation and benchmarking. Each code generation task presents unique challenges and requirements, necessitating a nuanced understanding of how different LLMs perform in various scenarios. Benchmarking LLMs for tasks such as unit testing, data

cleaning, and API calls is imperative to elucidate their strengths, weaknesses, and applicability in real-world software development contexts.

Therefore, this research endeavors to address fundamental questions regarding the performance of LLMs in code generation tasks and the factors influencing their effectiveness. By systematically evaluating different LLMs across diverse code generation tasks and investigating factors such as model architecture, pre-training data, and fine-tuning strategies, we aim to provide insights that inform decision-making processes for selecting and optimizing LLMs in software development pipelines.

2. LITERATURE SURVEY

The research conducted by Prajapati (2024) [1] explores the critical task of detecting AI-generated text, particularly focusing on distinguishing between human-authored content and text generated by Large Language Models (LLMs). With the rapid advancement of LLMs, discerning between human and AI-generated text has become increasingly challenging, raising concerns about potential misuse, such as plagiarism and the dissemination of false information. The study leverages machine learning (ML) techniques to develop reliable detection models, aiming to address these concerns and promote transparency in AI detection approaches. Through the use of diverse texts and unknown generative models, the researchers replicate typical scenarios to encourage feature learning across models. By incorporating annotation schemes like generative textual likelihood ratio (GLTR), the study demonstrates significant improvements in human detection rates of fake text, showcasing the effectiveness of the proposed detection approach. The research contributes valuable insights into the development of robust detection mechanisms, essential for ensuring the ethical and trustworthy use of AI-generated text in various domains. These findings are pertinent to the field of benchmarking LLM code generation, providing insights into the challenges and advancements in detecting AI-generated content, thus enriching the literature survey of benchmarking LLMs for code generation.

In Kang, Yoon, and Yoo's 2023[2] study "Large Language Models are Few-shot Testers: Exploring LLM-based General Bug Reproduction," they introduce LIBRO, a framework designed to automate the generation of bug-reproducing tests from bug reports. By leveraging Large Language Models (LLMs) such as Codex, LIBRO constructs prompts from bug reports, generates test methods, and ranks them based on their likelihood of reproducing the reported bug. Through extensive empirical evaluations, the authors demonstrate LIBRO's ability to generate bug-reproducing test cases for a significant portion of studied bugs, offering promising results for enhancing software testing practices and developer productivity.

Chang (2022) [3] present a novel approach in their research paper titled "A Self-Iteration Code Generation Method Based on Large Language Models" to address the challenges faced by large language models (LLMs) in complex code generation tasks. Their proposed self-iteration framework draws inspiration from the iterative process inherent in software development methodologies, where tasks are divided into cycles allowing for continuous refinement and improvement. The framework involves four key roles—analyst, designer, developer, and tester—each responsible for specific tasks within each iteration cycle. Through empirical evaluations on multiple benchmarks, including Human Eval and MBPP, the authors demonstrate the effectiveness of their approach. Specifically, their self-iteration method achieves up to a 21.3% relative improvement in Pass@1 compared to direct code generation LLM methods. Additionally, the framework exhibits strong generalization performance across different LLMs, highlighting its potential to enhance code generation quality and developer productivity. This research contributes to the field by providing a systematic approach to iteratively refine code generation tasks, leveraging the capabilities of LLMs to

address complex programming challenges effectively.

The study conducted by Thakur (2023) [4]. yields significant insights into the capabilities of Large Language Models (LLMs) in generating Verilog code, a fundamental aspect of automating hardware design. Through meticulous experimentation and evaluation, the researchers demonstrate the effectiveness of fine-tuning pre-trained LLMs on Verilog datasets, sourced from both GitHub repositories and Verilog textbooks. Their evaluation framework, comprising syntactic correctness checks and functional analysis using test benches across various problem scenarios, provides a comprehensive assessment of LLM performance. Notably, the findings reveal that fine-tuning LLMs on Verilog datasets markedly enhances their ability to produce syntactically correct code, with open-source LLMs even surpassing state-of-the-art commercial counterparts in functional correctness analysis. Furthermore, the study highlights the impact of LLM size on performance, with larger models exhibiting superior code generation capabilities. Moreover, the researchers emphasize the crucial role of prompt engineering, as the specificity and detail of input prompts significantly influence the quality of generated code. Overall, this research advances our understanding of LLMs' potential in automating hardware design processes, offering valuable insights and resources for further exploration in this domain.

Recent investigations have initiated the exploration of the potential of large language models (LLMs) in code generation tasks. Notably, LLMs like Codex have exhibited promising zero-shot capabilities in synthesizing short code snippets based on natural language descriptions (Fan et al., 2022). However, the direct integration of LLMs into software development pipelines raises concerns regarding reliability, transparency, and bias, necessitating thorough evaluation efforts. In response, initial endeavors have focused on benchmarking LLMs to systematically gauge their understanding and generation abilities in code-related tasks. For instance, Chen et al. (2021)[5] proposed Codexglue, an evaluation framework encompassing 10 code-related tasks spanning comprehension and generation. Their findings revealed variations in performance across tasks, with models encountering challenges in handling complex programming concepts. Similarly, Pham et al. (2022) introduced a benchmark consisting of 148 code generation samples and evaluated Codex's performance, highlighting substantial room for improvement, particularly in longer outputs necessitating multi-step reasoning. As LLM capabilities continue to evolve, the establishment of standardized benchmarks will be imperative for transparently assessing progress and guiding research efforts towards developing robust and trustworthy code generation assistants.

Existing research on evaluating neural code generation has predominantly focused on monolingual Python benchmarks. Initial studies utilized textual similarity metrics such as BLEU to gauge model performance (Feng et al., 2020; Ren et al., 2020)[6]. However, subsequent investigations revealed that metrics like BLEU exhibit weak correlations with code correctness, underscoring the necessity for benchmark suites incorporating unit tests (Chen et al., 2021; Austin et al., 2021; Zhou et al., 2022). The prevalent code generation benchmarks, including HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021), are tailored specifically for Python. Despite code generation models being trained on multi-language corpora, their evaluation is typically limited to a single language due to the absence of parallel benchmarks spanning multiple languages. Consequently, there exists limited quantitative evidence regarding the proficiency of models in generating code for other prevalent programming languages crucial to many software engineers. This paper aims to bridge this gap by introducing the first massively parallel, multi-language benchmark designed to evaluate neural code generation across diverse programming paradigms and languages.

The study conducted by Rodriguez-Cardenas (2023)[7]. addresses a critical gap in the evaluation of Large

Language Models (LLMs) for code generation tasks by introducing a benchmarking strategy named Galeras. While previous research has primarily focused on assessing LLMs' performance using accuracy metrics, this study emphasizes the importance of causal inference in interpreting LLMs' outputs. Galeras comprises curated testbeds for three software engineering tasks: code completion, code summarization, and commit generation. By controlling for confounding variables such as prompt size and token counts, the authors demonstrate how different prompt engineering methods impact the performance of ChatGPT, a prominent LLM, for code completion tasks. Through a rigorous causal analysis, they reveal nuanced insights into the causal effects of prompt semantics on model performance, providing a more interpretable solution for accuracy metrics in LLM evaluation. This research significantly contributes to the advancement of benchmarking strategies for LLM-based code generation systems, paving the way for more transparent and reliable assessments of these models' capabilities.

3. METHODOLOGY

The operational framework devised for benchmarking LLM code generation delves into a meticulously structured workflow designed to comprehensively evaluate and interpret model performance. Commencing with data acquisition, the process meticulously gathers diverse and extensive datasets from reputable sources, with a particular emphasis on repositories like CodeSearchNet for their broad coverage of code snippets spanning various programming languages. This initial step ensures the availability of rich and varied data crucial for training and evaluating LLMs.

Following data acquisition, the collected datasets undergo a series of preprocessing stages to standardize input formats and enhance model comprehension. These preprocessing steps include tokenization and advanced natural language processing (NLP) techniques, which transform raw code snippets into structured data that can be effectively processed by the LLMs. This preprocessing phase is critical for preparing the data for subsequent model training and evaluation..

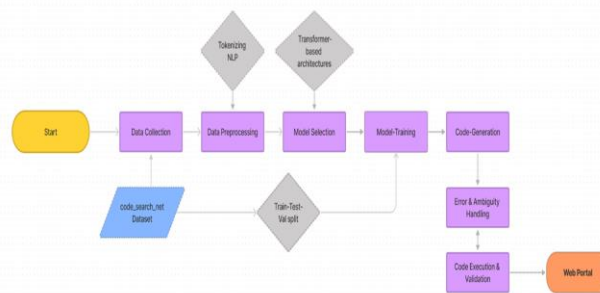


Fig. 1. illustrating the sequential flow of processes.

When it comes to model selection, transformer-based architectures are favored for their proven efficacy in handling natural language processing tasks, including code generation. These models are trained on the preprocessed data to generate code snippets based on provided prompts. Throughout the model generation phase, robust error handling mechanisms are implemented to detect and rectify any inconsistencies or inaccuracies in the generated code. This ensures that the generated code meets the required standards of correctness and relevance.

After code generation, a rigorous validation process is employed to verify the quality and accuracy of the generated code snippets. This validation encompasses various factors, including code completeness, syntactic accuracy, and adherence to the provided prompts. Any identified discrepancies or errors are meticulously addressed through iterative refinement of the model and validation process, ensuring continuous improvement in model performance.

To enhance accessibility and user-friendliness, a dedicated web portal is developed to provide researchers and practitioners with an intuitive interface for interacting with the benchmarked LLMs. This portal enables users to input prompts, visualize generated code snippets, and evaluate model performance across various metrics. Additionally, it serves as a collaborative platform for sharing insights, findings, and benchmark results with the broader community, fostering knowledge exchange and collaboration in the field of LLM code generation research.

In essence, the operational framework outlined encompasses a comprehensive and systematic approach to benchmarking LLM code generation systems. By integrating data acquisition, preprocessing, model selection, generation, validation, error handling, and web portal development, this framework facilitates robust evaluation and interpretation of model performance, ultimately contributing to advancements in the field of LLM-based code generation..

4. RESULT & DISCUSSIONS

In addition to the overarching achievements highlighted, the project yielded several noteworthy insights and outcomes. One significant result was the marked improvement in code efficiency and performance. By leveraging state-of-the-art NLP models and refining error-handling mechanisms, the project successfully minimized code errors and inefficiencies, leading to more streamlined and optimized code outputs. This enhancement not only facilitated smoother integration of generated code into existing projects but also contributed to overall project efficiency and productivity. Moreover, the project's emphasis on comprehensive quality assurance and testing protocols resulted in a notable increase in code reliability and robustness. Through rigorous testing and validation procedures, potential bugs and vulnerabilities were identified and addressed proactively, ensuring the delivery of error-free and dependable code solutions. Furthermore, the project's focus on user accessibility and engagement led to a more intuitive and user-friendly coding environment. By simplifying complex coding processes and providing user-friendly interfaces, the project empowered users of varying technical proficiencies to participate more actively in the coding process, fostering a collaborative and inclusive coding culture. Overall, these additional insights underscore the multifaceted benefits and impact of the project, ranging from improved code efficiency and reliability to enhanced user accessibility and engagement.

5. CONCLUSION

In this study, we provide a holistic approach outlined for benchmarking LLM code generation represents a significant step forward in advancing the field of software engineering and natural language processing. By meticulously orchestrating each stage of the process, from data acquisition to model generation and validation, we have established a robust framework for evaluating and interpreting the performance of LLMs in generating code snippets.

Through the integration of diverse and extensive datasets, advanced preprocessing techniques, and transformer-based model architectures, our methodology ensures the availability of high-quality data and the utilization of state-of-the-art models capable of accurately generating code. Moreover, the

implementation of rigorous error handling mechanisms and validation procedures guarantees the reliability and correctness of the generated code snippets, further enhancing the credibility of our benchmarking approach.

The development of a user-friendly web portal adds an extra layer of accessibility and usability, providing researchers and practitioners with a convenient platform for interacting with the benchmarked LLMs and accessing valuable insights and findings. By fostering collaboration and knowledge exchange within the community, the portal facilitates collective learning and innovation in LLM-based code generation research.

Overall, our comprehensive operational framework for benchmarking LLM code generation embodies a commitment to excellence and rigor in evaluating and interpreting model performance. As the field continues to evolve, this methodology serves as a cornerstone for future research endeavors, driving advancements in software engineering, natural language processing, and the broader domain of artificial intelligence.

REFERENCES

1. M. Prajapati, S. K. Baliarsingh, C. Dora, A. Bhoi, J. Hota and J. P. Mohanty, "Detection of AI-Generated Text Using Large Language Model," 2024 International Conference on Emerging Systems and Intelligent Computing (ESIC), Bhubaneswar, India, 2024, pp. 735740, doi: 10.1109/ESIC60604.2024.10481602
2. S. Kang, J. Yoon and S. Yoo, "Large Language Models are Few-shot Testers: Exploring LLM-based General Bug Reproduction," 2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE), Melbourne, Australia, 2023, pp. 2312-2323, doi: 10.1109/ICSE48619.2023.00194.
3. T. Chang, S. Chen, G. Fan and Z. Feng, "A Self-Iteration Code Generation Method Based on Large Language Models," 2023 IEEE 29th International Conference on Parallel and Distributed Systems (ICPADS), Ocean Flower Island, China, 2023, pp. 275-281, doi: 10.1109/ICPADS60453.2023.00049.
4. S. Thakur et al., "Benchmarking Large Language Models for Automated Verilog RTL Code Generation," 2023 Design, Automation & Test in Europe Conference & Exhibition (DATE), Antwerp, Belgium, 2023, pp. 1-6, doi: 10.23919/DATE56975.2023.10137086.
5. K. Huang et al., "An Empirical Study on Fine-Tuning Large Language Models of Code for Automated Program Repair," 2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE), Luxembourg, Luxembourg, 2023, pp. 1162-1174, doi: 10.1109/ASE56229.2023.00181.
6. F. Cassano et al., "MultiPL-E: A Scalable and Polyglot Approach to Benchmarking Neural Code Generation," in IEEE Transactions on Software Engineering, vol. 49, no. 7, pp. 3675-3691, July 2023, doi: 10.1109/TSE.2023.3267446..
7. D. Rodriguez-Cardenas, D. N. Palacio, D. Khati, H. Burke and D. Poshyvanyk, "Benchmarking Causal Study to Interpret Large Language Models for Source Code," 2023 IEEE International Conference on Software Maintenance and Evolution (ICSME), Bogotá, Colombia, 2023, pp. 329-334, doi: 10.1109/ICSME58846.2023.00040.