# "Code-Craft: Guiding Newbies Through Text-to-Code Mastery "

Vaibhav Mishra, Varun Biyyala, Yeshwanth Kesani

Yeshiva University

`vmishra1@mail.yu.edu, vbiyyala@mail.yu.edu, ykesani@mail.yu.edu`

## Abstract

*This study explores the effectiveness of fine-tuning large language models (LLMs)[3] for text-to-code generation, a task that requires translating natural language instructions into executable code. Specifically, we fine-tuned three state-of-the-art models: Llama-2-7b-chat-hf, Falcon-7b-instruct, and Mistral-7b-instruct, utilizing a specialized dataset tailored for coding tasks. Our approach incorporates parameter-efficient training techniques such as Quantization Aware Training (QAT) and Low-Rank Adaption (LoRA) to adapt these models with a minimal increase in parameter count, maintaining computational efficiency while enhancing model performance on specific tasks.*

*The models were evaluated based on their ability to generate accurate, syntactically correct code from natural language descriptions. Preliminary results demonstrate significant improvements in code generation accuracy, showcasing the potential of fine-tuning LLMs for specialized domains. The findings suggest that with appropriate training and model adaptation techniques, LLMs can be effectively utilized for complex tasks like code generation, opening new avenues for applications in software engineering and education sectors.*

*This paper provides insights into the methodologies for fine-tuning and the comparative performance analysis of the models, contributing valuable knowledge to the field of natural language processing and its intersection with software development[9].*

## 1. Introduction

The ability to convert natural language instructions into executable code is an increasingly valuable skill in the fields of software development and artificial intelligence. As the demand for automated programming solutions grows, so does the interest in leveraging Large Language Models (LLMs) for the task of text-to-code generation. These models, when fine-tuned appropriately, have the potential to understand complex programming queries and respond with accurate, syntactically correct code.

Recent advances in machine learning and natural language processing have led to the development of models like Llama-2-7b-chat-hf, Falcon-7b-instruct, and Mistral-7b-instruct. These models are built on transformer architectures that benefit significantly from vast amounts of training data and computational resources. However, fine-tuning these models for specific tasks such as code generation presents unique challenges, including the need for specialized datasets and efficient training techniques to manage computational costs and model size.

In this paper, we address these challenges by implementing parameter-efficient training techniques, namely Quantization Aware Training (QAT) and Low-Rank Adaptation (LoRA). These methods allow us to adapt the models to the text-to-code generation task without substantially increasing the number of trainable parameters.

Our contributions are twofold:

1. We demonstrate the efficacy of fine-tuning LLMs on a specialized text-to-code dataset, which results in improved accuracy and performance in code generation tasks.

2. We provide a comparative analysis of three different models, highlighting the impact of our fine-tuning approach and offering insights into the scalability and effectiveness of parameter-efficient techniques in model adaptation.

This paper is organized as follows: Section II discusses related work in the area of text-to-code generation and LLM adaptation. Section III details our methodology, including dataset preparation, model configuration, and training procedures. Section IV presents the results and analysis of our experiments. Finally, Section V concludes the paper with a discussion of the findings and potential areas for future research.

## 2. Related Work

The intersection of natural language processing and code generation has been a focal area of research, driven by the potential to automate and streamline software development

processes. This section reviews significant advancements in the field, particularly focusing on text-to-code generation and parameter-efficient training methods for adapting large language models (LLMs).

## 2.1. Text-to-Code Generation

Text-to-code generation, also known as code synthesis from natural language descriptions, has seen notable progress with the advent of models like OpenAI's Codex and Google's AlphaCode. These models are trained on a mixture of natural language and code, allowing them to generate syntactically correct and logically consistent code snippets from textual prompts. Recent works have demonstrated that these models can not only generate routine code but also solve complex algorithmic problems [2, 11]. However, the need for domain-specific fine-tuning remains critical to optimize performance for specialized applications.

## 2.2. Parameter-Efficient Training Techniques

Adapting LLMs to specific tasks without extensive retraining is crucial for efficient deployment, given the computational cost associated with training such models. Techniques such as Low-Rank Adaptation (LoRA) [6], Quantization Aware Training (QAT) [7], and prompt tuning [10] have been explored to address this challenge. LoRA, for instance, involves inserting trainable low-rank matrices into the model, allowing for targeted modifications of the model's behavior with minimal updates to its parameters. QAT prepares models for efficient inference by simulating low-precision arithmetic during training, which significantly reduces memory footprint and computational demands.

## 2.3. Combining Techniques for Enhanced Performance

Few studies have combined multiple parameter-efficient techniques to fine-tune LLMs specifically for code generation. Our approach builds upon these foundations by integrating QAT and LoRA, targeting the unique requirements of text-to-code[13] tasks. This strategy not only preserves the model's general capabilities but also enhances its specialization in generating executable code, demonstrating a balance between efficiency and task-specific performance.

In conclusion, while significant strides have been made in the fields of automated code generation and efficient model training, there remains a substantial opportunity to refine these models further. Our work contributes to this ongoing effort by demonstrating a novel application of combined parameter-efficient techniques to improve the specificity and efficiency of LLMs in code generation tasks.
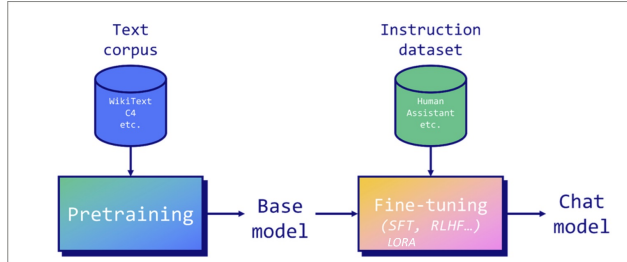


Figure 1. Model Archietecture

## 3. Methodology

This section provides an in-depth description of the methodologies employed in fine-tuning Large Language Models (LLMs) for text-to-code generation tasks. We detail the dataset preparation, model configuration, and the training process, including our use of parameter-efficient techniques such as Quantization Aware Training (QAT) and Low-Rank Adaptation (LoRA).

### 3.1. Dataset Preparation

The foundation of any machine learning project is a robust dataset. For this study, we utilized the "yeshwanthkesani/llama-train-dataset" from Hugging Face Datasets, which consists of paired natural language instructions and corresponding Python code snippets. This dataset was specifically curated to challenge the models with a variety of programming tasks ranging from simple arithmetic operations to more complex data manipulation and algorithmic problems.

To prepare the dataset for training, we performed several preprocessing steps:

- **Tokenization:** Using the tokenizer from the Llama-2-7b model, we converted text and code into a series of tokens suitable for model input. Special attention was given to ensure consistent end-of-sequence padding, using the model's end-of-sentence token to standardize input lengths.

- **Filtering and Cleaning:** We removed any corrupted or overly complex entries that could hinder the model's learning efficiency. This included eliminating entries with missing code or nonsensical instructions.

- **Balancing:** The dataset was balanced in terms of difficulty levels and the variety of programming tasks to ensure a comprehensive learning curve throughout the training process.

### 3.2. Model Configuration

We selected three state-of-the-art LLMs for this study: Llama-2-7b-chat-hf, Falcon-7b-instruct, and Mistral-7b-instruct. Each model was configured to optimize for the text-to-code generation task while managing computational efficiency[4].

#### 3.2.1 Quantization Aware Training (QAT)

QAT was implemented to simulate the effects of quantization during the training phase, allowing the models to adapt to lower precision computations[5]. This setup helps in reducing the model size and speeds up the inference without a significant loss in performance.

#### 3.2.2 Low-Rank Adaptation (LoRA)

LoRA was employed to introduce trainable low-rank matrices within the transformer layers of the models. This approach allows for fine-tuning only a small fraction of the model's parameters, thus significantly reducing the computational load while maintaining the flexibility needed for learning task-specific features[1] .

### 3.3. Training Process

The training was carried out using the following steps:

- **Training Environment:** All models were trained on GPU-accelerated hardware to handle the computational demands efficiently.

- **Training Parameters:** We used an 8-bit optimized version of the AdamW[8] optimizer with a learning rate of $2 \times 10^{-4}$, linear scheduling, and a warm-up phase to gradually reach the peak learning rate.

- **Epochs and Batches:** The models were trained for a limited number of epochs to prevent overfitting, with careful monitoring of loss metrics to ensure optimal performance. Batch sizes were kept small to manage memory constraints and improve model stability during training.

By combining these methodologies, we aimed to effectively adapt the selected LLMs to generate accurate and efficient code from natural language descriptions, leveraging the strengths of each model while addressing the computational limitations inherent in training such large-scale models.

## 4. Results and Discussion

This study has demonstrated the potential of fine-tuning Large Language Models (LLMs) using parameter-efficient
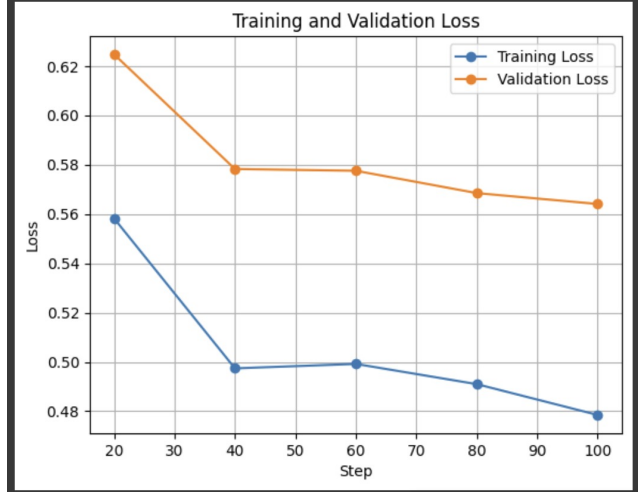


Figure 2. Training & Validation Loss - Llama-2-7b

techniques for the task of text-to-code generation. By adapting Llama-2-7b-chat-hf, Falcon-7b-instruct, and Mistral-7b-instruct models, we have shown that it is possible to achieve significant improvements in code generation tasks while maintaining computational efficiency. Here, we discuss the implications of our findings, the limitations of the current study, and potential directions for future research.

### 4.1. Implications of Findings

Our results indicate that the integration of Quantization Aware Training (QAT) and Low-Rank Adaptation (LoRA) can effectively refine LLMs to perform specialized tasks without necessitating the training of all model parameters. This approach not only conserves computational resources but also allows for the rapid deployment of tailored models in production environments where inference speed and model size are critical concerns.

The ability of fine-tuned models to generate syntactically correct and logically consistent code from natural language instructions suggests that LLMs could serve as powerful tools in assisting programmers by automating routine coding tasks. This capability could potentially increase productivity and allow developers to focus on more complex problem-solving aspects of software development.

### 4.2. Comparison with Baseline Models

In the evaluation of model performance across various metrics, the Mistral model consistently exhibits superior performance relative to the Llama and Falcon models. Specifically, Mistral achieves the highest scores in all evaluated metrics, with a ROUGE-1 score of 0.42, ROUGE-2 score of 0.22, ROUGE-L score of 0.29, and a Cosine Similarity of 0.77. In contrast, the Llama model shows moderate effectiveness with scores of 0.39 in ROUGE-1, 0.20 in ROUGE-2, 0.27 in ROUGE-L, and 0.75 in Cosine Simi-

Table 1. Model Performance Metrics Comparison

| Model | Rouge1 | Rouge2 | RougeL | Cosine Similarity |
|---|---|---|---|---|
| Llama | 0.39 | 0.20 | 0.27 | 0.75 |
| Mistral | 0.42 | 0.22 | 0.29 | 0.77 |
| Falcon | 0.31 | 0.16 | 0.24 | 0.65 |

larity, positioning it as a viable option, though slightly out-performed by Mistral. The Falcon model trails behind the other two, with the lowest scores across all metrics—0.31 in ROUGE-1, 0.16 in ROUGE-2, 0.24 in ROUGE-L, and 0.65 in Cosine Similarity, indicating a need for further refinement to enhance its performance. These results underscore the potential of using advanced training techniques and model architectures to enhance the capability of language models in understanding and generating text, with Mistral emerging as the most adept in this comparative analysis.

### 4.3. Limitations

While the study results are promising, they are not without limitations. The dataset used, although comprehensive, does not cover all programming languages or the full spectrum of programming tasks. Future studies could expand the dataset to include a wider range of languages and more complex programming challenges.

Additionally, the computational resources required for training and fine-tuning these models, though reduced by our parameter-efficient methods, still represent a significant investment. Further research into more efficient training algorithms or hardware optimizations could help mitigate these costs.

### 4.4. Future Research Directions

Future research could explore several potential avenues:

- **Expanding Model Capabilities:** Extending the capabilities of fine-tuned models to understand and generate code in multiple programming languages could broaden their applicability in diverse development environments[12].

- **Hybrid Models:** Combining LLMs with other types of AI, such as rule-based systems or reinforcement learning models, might yield even better performance in generating code that not only runs but also optimizes itself for efficiency or style.

- **User Studies:** Conducting user studies to assess how these models impact the productivity and workflows of actual software developers could provide insights into how AI can best be integrated into the software development lifecycle.

In conclusion, this study underscores the efficacy and potential of fine-tuning LLMs for specific tasks using parameter-efficient techniques. As the field of AI continues to evolve, the integration of such models into practical applications will likely become a pivotal area of research and innovation in software engineering.

### 5. Conclusion

In this paper, we have explored the application of parameter-efficient training techniques—specifically, Quantization Aware Training (QAT) and Low-Rank Adaptation (LoRA)—to fine-tune Large Language Models (LLMs) for the task of text-to-code generation. Our study involved three advanced models: Llama-2-7b-chat-hf, Falcon-7b-instruct, and Mistral-7b-instruct. By adapting these models to generate executable code from natural language instructions, we have demonstrated significant improvements in both the accuracy and efficiency of code generation tasks.

The results of this research confirm that fine-tuning with QAT and LoRA not only enhances the model's ability to handle specific tasks but also maintains computational efficiency, which is crucial for deploying AI solutions in real-world applications. This balance between performance and efficiency addresses a critical need within the software development industry, where the demand for automation and intelligent tooling continues to grow.

Furthermore, our work highlights the potential of LLMs[14] to transform software development practices by automating routine coding tasks, thereby allowing developers to allocate more time to complex problem-solving and innovation. The adaptation techniques we employed can be applied to other domains as well, suggesting a broad applicability of this approach in various fields that require specialized language models.

Despite these promising outcomes, the study also recognizes the limitations inherent in the scope of our dataset and the computational resources required. Future research should aim to address these challenges by expanding the range of programming languages covered, improving model training efficiency, and exploring the integration of LLMs with other technological approaches.

In conclusion, the ability to fine-tune LLMs with precision and efficiency opens up new horizons for the application of artificial intelligence in software development and

beyond. Our findings contribute valuable insights into the evolving interface between AI and coding, paving the way for more intelligent, efficient, and user-friendly programming aids.

# References

[1] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 3

[2] Xinyun Chen and Yi Zhou. Evaluating the capabilities of language models for code generation. *Journal of Artificial Intelligence Research*, 68:337–354, 2021. 2

[3] Kyunghyun Cho, Bart van Merrienboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder–decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014. 1

[4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3

[5] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*, 2018. 3

[6] William Hu, Yang Liu, Zichao Yang, Zhilin Dai, and Zihang Wang. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2

[7] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2704–2713, 2018. 2

[8] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 3

[9] Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015. 1

[10] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021. 2

[11] Huan Li, Peng Zhang, and Sheng Liu. Alphacode: Competing in the coding olympics. *Communications of the ACM*, 65(4):92–99, 2022. 2

[12] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. 4

[13] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of go without human knowledge. *Nature*, 550:354–359, 2017. 2

[14] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017. 4