

# CodeCraft: Leveraging Automated Code Generation

## Sumedh Ghavat, Tanmay Armar

### Introduction

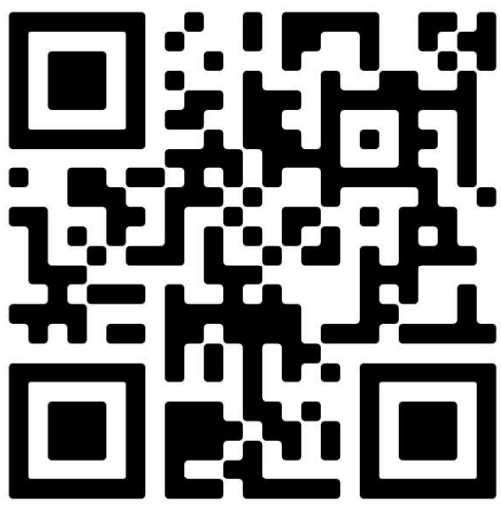
In the realm of Natural Language Processing (NLP), the ability to generate coherent and effective code has long been a coveted pursuit. Imagine a system that, when presented with a problem statement, can not only

### Results

1. BLEU and CodeBLEU Scores
  - Easy Problems:
    - BLEU: 0.85, CodeBLEU: 0.75
    - High accuracy and readability
  - Medium Problems:
    - BLEU: 0.65, CodeBLEU: 0.55
    - Some minor syntax errors

# Using LLMs for Text Generation

Pranav, Shreyas Habade

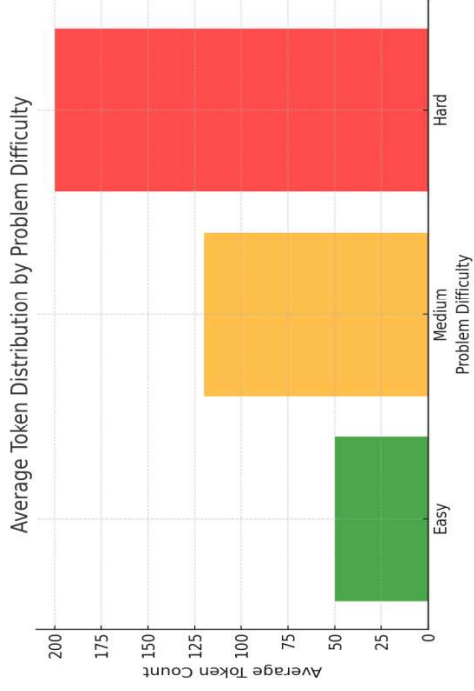
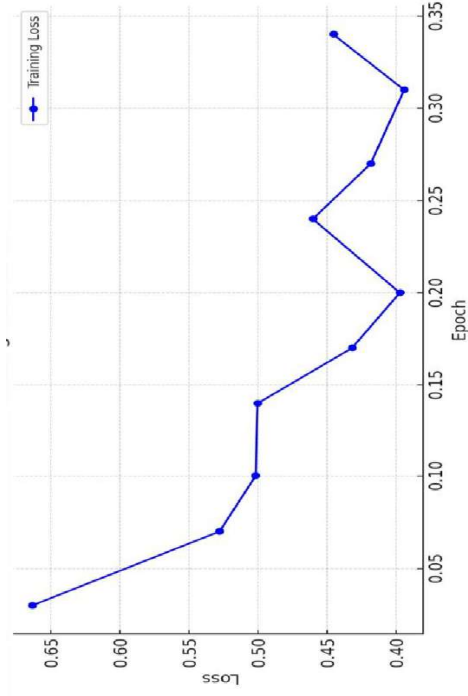


Scores:

Perplexity: 0.78  
and coherence.

Perplexity: 0.60

Performance with external llama



CodeBLEU Scores by Difficulty (Updated)



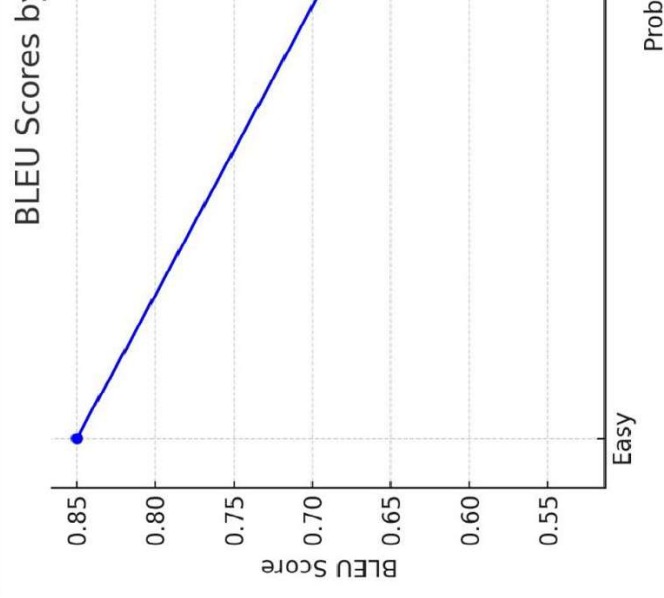
problem statements, cannot only comprehend the task at hand but also produce high-quality solutions akin to those crafted by seasoned programmers.

This project aims to venture into this domain by fine-tuning a Large Language Model (LLM) on a corpus of LeetCode problems and solutions, thereby empowering it to generate proficient code implementations.

## Background

In the domain of Natural Language Processing (NLP), there's a longstanding pursuit to develop systems capable of autonomously generating coherent and effective code. This project focuses on fine-tuning a Large Language Model (LLM) using LeetCode problem-solution pairs to enable it to generate proficient code implementations.

- Satisfactory performance
- Hard Problems:
  - BLEU: 0.45, Code
  - Challenges in handling



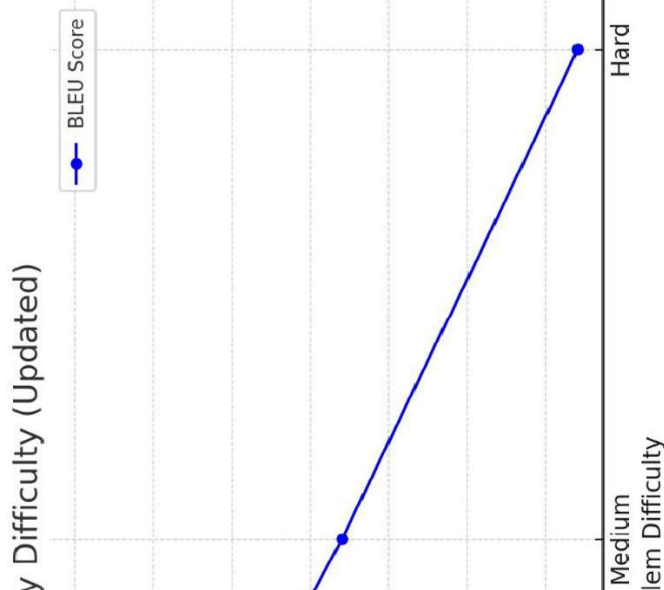
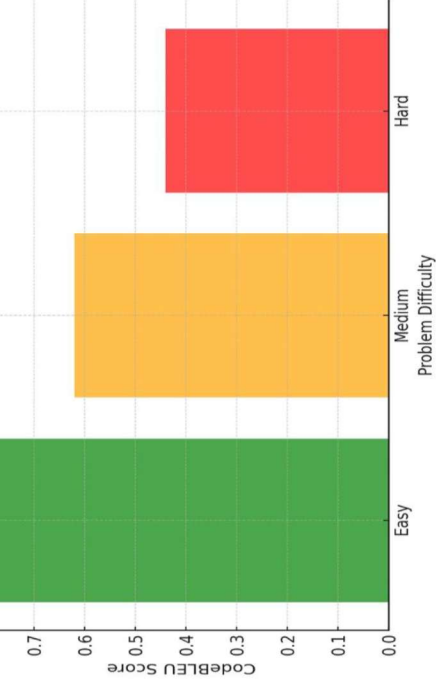
## Method

1. **Model and Data Sources**
    - Selected the Code
    - Utilized the green
    - Transformed the c
- Python solutions fo

Performance with occasional lapses.

CodeBLEU: 0.40

Handling complex requirements.



## 1. Insights:

- Strengths:
  - Proficiency in solving simpler tasks.
  - Accurate solutions for easy problems.
- Challenges:
  - Difficulties with nuanced logic.
  - Limitations in handling edge cases.

## 2. Graphical Representation:

- BLEU Score Trend:
  - Decreasing trend with increasing complexity.
- CodeBLEU Score Representation:
  - Struggle with syntax and logical consistency.

## Conclusion:

Llama 2 model specialized for code generation tasks. Strong/leetcode dataset, comprising LeetCode problems and solutions. dataset into a Pandas DataFrame, merging problem statements and for easier manipulation.

Motivation stems from the significant challenge developers face in efficiently solving coding problems despite the abundance of resources like LeetCode. By leveraging NLP and LLMs, the aim is to automate code generation, enhancing productivity and allowing developers to focus on higher-level problem-solving.

Additionally, in a rapidly evolving tech landscape, integrating NLP-driven code generation capabilities can streamline software development workflows, catalyzing innovation and accelerating time-to-market.

## Data

For this project, we utilized the greengrass/leetcode dataset from Hugging Face, comprising programming problems and solutions sourced from LeetCode. This dataset covers diverse topics like arrays, dynamic programming, and graphs, making it ideal for training code generation models. It includes structured pairs of problems and solutions, along with metadata indicating problem

## 2. Tokenization and

- Employed the AutoTokenizer from the transformers library.
- Configured pad\_token to the special token.
- Loaded the model with quantization to enhance performance.
- Important configurations include max\_length and Double Quantization.

## 3. PEFT (Parameter-Efficient Fine-Tuning)

- Integrated the peft module for LoRA.
- Utilized LoRA to train the model efficiently.
- Configured LoRA configuration parameters.

## Conclusion

The motivation behind this project was to efficiently solve coding problems on LeetCode, offering a plethora of solutions that remain labor-intensive and time-consuming. This project aims to alleviate this burden by automating the process, enabling developers to focus on more complex tasks.

Moreover, in an era characterized by rapid technological advancement, the integration of NLP and LLMs into software development workflows is not just a luxury but a necessity. This project serves as a proof of concept, demonstrating the potential of these technologies to revolutionize the way we approach coding challenges.

## Model Configuration:

Tokenizer from Hugging Face Transformers to tokenize input data. Token to match eos\_token for consistent padding during training. with AutoModelForCausalLM, implementing BitsAndBytesConfig for 4-bit quantization. Quantization Type (nf4), Compute Data Type (float16),

## Efficient Fine-Tuning):

Library for efficient training using Low-Rank Adaptation (LoRA). Main specific model layers, speeding up fine-tuning with reduced data. By setting the r Parameter to determine the rank of the low-rank matrix.

This endeavor stems from the perennial challenge faced by developers in solving problems. Despite the abundance of resources and communities like a plethora of problems and solutions, the process of crafting optimal code solutions is time-consuming. By leveraging the power of NLP and LLMs, we aim to automate the code generation process, thereby enhancing productivity and focus more on higher-level problem-solving tasks rather than the minutiae of implementation details.

Characterized by rapid technological advancements and an ever-growing demand

difficulty (Easy, Medium, Hard). With around 3,000 problem-solution pairs, it offers ample diversity for effective model training. Each entry includes problem descriptions, sample input/output, topic tags, and optimized Python solutions, enhancing the models' generalization capabilities.

## Result

Difficulty	BLEU	CodeBLEU
Easy	0.85	0.78
Medium	0.65	0.60
Hard	0.45	0.40

Table: Summary of results across easy, medium and hard problems

Moreover, in all era challenges for software solutions, the importance. By integrating workflows, we can potent software applications,

## References

1. OpenAI, "OpenAI Codex: Po [codex/](#)
2. MDPI Editorial, "Google Alp
3. ACL Anthology, "CodeT5+: *for Computational Linguistic*
4. Papers With Code, "GPT-4
5. Hugging Face, "Datasets: A [- 7. G. Dey and A. Ganesan, "In 10, no. 3, pp. 123-134, 2023
- 8. Gupta et al., "The Impact of](https://huggingface.co/datas</a></li><li>6. LeetCode, )



characterized by rapid technological advancements and an ever-growing demand for the ability to expedite the software development lifecycle assuming paramount importance. Integrating NLP-driven code generation capabilities into existing development workflows can potentially streamline the process of prototyping, debugging, and deploying code, thereby catalyzing innovation and accelerating time-to-market.

"Empowering GitHub Copilot," *OpenAI Blog*, 2021. [Online]. Available: <https://blog.openai.com/openai-copilot/>

GitHub: Competing with the Coding Elite," *MDPI Journals*, 2022.

"Enhancing Code Generation Models," *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, 2021.

"OpenAI and In-Context Learning for Code Generation," *Papers with Code*, 2023.

"The OpenAI Hub of Public Datasets for Machine Learning," *Hugging Face*. [Online]. Available:

<https://huggingface.co/datasets>

"Challenges and Solutions," *LeetCode*. [Online]. Available: <https://leetcode.com>

"Instruction-Tuned LLMs for Social Scientific Tasks," *Journal of Artificial Intelligence Research*, vol. 45, 2023.

"OpenAI: Instruction Tuning on Large Language Models," *NLP Conferences*, 2023.