

Fine-Tuning CodeT5 for Python Code Generation with Streamlit Deployment

Abstract

This paper presents a comprehensive approach to fine-tuning transformer models for automatic code generation from natural language queries. We implement an ensemble learning technique with k-fold cross-validation using the CodeT5 model as our foundation. The system achieves a ROUGE-L score of 39.14, BLEU score of 25.66, and an exact match score of 0.09 on our test dataset of Python code snippets. We further develop an interactive web application using Streamlit to demonstrate the practical utility of our model. This work contributes to the growing field of AI-assisted programming by enabling developers to generate code snippets through natural language descriptions.

1. Introduction

1.1 Motivation

Software development is a complex and time-consuming process, with programmers spending considerable time writing boilerplate code and implementing standard algorithms. Recent advances in natural language processing (NLP) and deep learning have enabled the development of systems that can generate code snippets from natural language descriptions, potentially increasing developer productivity and lowering the barrier to entry for programming.

1.2 Problem Statement

The primary challenge in code generation is accurately interpreting the intent behind natural language queries and translating them into syntactically correct and semantically appropriate code. Traditional approaches often fail to capture the nuances of both natural language and programming syntax, resulting in inaccurate or non-functional code generation.

1.3 Proposed Solution

We propose fine-tuning the CodeT5 transformer model on a dataset of Python problem-solution pairs. Our approach incorporates:

- K-fold cross-validation for robust model training
- Ensemble learning to improve prediction accuracy
- Interactive deployment via Streamlit for practical usage
- Temperature-based sampling for diverse code generation

2. Theoretical Background

2.1 Transformer Models

The Transformer architecture, introduced by Vaswani et al. in "Attention is All You Need," has revolutionized NLP tasks through its self-attention mechanism. Unlike recurrent neural networks, transformers process entire sequences simultaneously, allowing for better parallelization and improved learning of long-range dependencies.

The core components of transformer models include:

- **Multi-head attention:** Allows the model to focus on different parts of the input sequence
- **Position encodings:** Provide information about token positions
- **Feed-forward networks:** Process the contextual representations
- **Residual connections:** Help with gradient flow during training

2.2 Fine-Tuning Approach

Fine-tuning involves adapting a pre-trained model to a specific downstream task by continuing the training process on task-specific data. This approach leverages the knowledge embedded in the pre-trained model, requiring significantly less data and computational resources than training from scratch.

Parameter-efficient fine-tuning methods have emerged to reduce computational demands, including:

- Adapter layers
- LoRA (Low-Rank Adaptation)
- Selective fine-tuning of specific layers

2.3 CodeT5 Architecture

CodeT5 is a transformer-based encoder-decoder model specifically pre-trained on code-related tasks. It incorporates:

- **Identifier-aware pre-training:** Special handling for code identifiers
- **Dual-generation objectives:** Both understanding and generating code
- **Multi-programming language support:** Training on various languages

3. Methodology

3.1 Dataset Description

The dataset consists of approximately 3,300 problem-solution pairs, where each problem is a natural language description and each solution is a corresponding Python code snippet. The dataset covers diverse programming tasks including:

- NumPy array operations
- Data structure manipulations (lists, dictionaries)
- String processing
- Pandas dataframe operations
- Algorithm implementations (sorting, searching)
- Database operations

Sample problems include:

- "Write a NumPy program to repeat elements of an array"
- "Write a Python function to create and print a list where the values are square of numbers between 1 and 30"
- "Write a Python program to remove duplicates from a list of lists"

Data preprocessing steps included:

```
df = pd.read_csv("ProblemSolutionPythonV3.csv").dropna(subset=["Python Code"])
df = df.rename(columns={"Problem": "Query", "Python Code": "Code_Snippet"})
df["Query"] = df["Query"].str.lower()
df["Code_Snippet"] = df["Code_Snippet"].str.lower()
```

3.2 Model Architecture

We utilized the CodeT5-base model with approximately 220 million parameters. The model employs a transformer encoder-decoder architecture with:

- 12 layers in both encoder and decoder
- 12 attention heads
- 768-dimensional hidden states

3.3 Fine-Tuning Process

Our fine-tuning process involved:

1. **K-fold Cross-Validation (k=3):** Splitting the dataset into three parts, training on two and validating on one, rotating through all combinations
2. **Input Formatting:** Prefixing queries with "generate code: " to signal the task
3. **Tokenization:** Using the CodeT5 tokenizer with maximum length of 256 tokens
4. **Training Parameters:**

- Learning rate: 5e-5
- Batch size: 4
- Training epochs: 5
- Weight decay: 0.01
- AdamW optimizer

3.4 Ensemble Learning Approach

We created an ensemble of models from the k-fold training process, using a majority voting strategy for inference. This ensemble approach helps:

- Reduce overfitting
- Improve generalization
- Increase prediction robustness
- Mitigate the impact of random initialization

3.5 Evaluation Metrics

We evaluated our model using three complementary metrics:

- **ROUGE-L**: Measures the longest common subsequence between generated and reference code
- **BLEU**: Computes n-gram overlap between generated and reference code
- **Exact Match**: The percentage of generations that exactly match the reference

4. Implementation

4.1 Training Implementation

The training pipeline consists of:

```
# K-Fold Training + Ensemble
k = 3
kf = KFold(n_splits=k, shuffle=True, random_state=42)
models = []

for fold, (train_idx, val_idx) in enumerate(kf.split(df)):
    model = T5ForConditionalGeneration.from_pretrained(model_checkpoint)
    train_data = df.iloc[train_idx]
    val_data = df.iloc[val_idx]

    train_dataset = CodeGenDataset(train_data["Query"].tolist(),
                                   train_data["Code_Snippet"].tolist())
    val_dataset = CodeGenDataset(val_data["Query"].tolist(),
                                 val_data["Code_Snippet"].tolist())

    trainer = Trainer(
```

```

        model=model,
        args=training_args,
        train_dataset=train_dataset,
        eval_dataset=val_dataset,
        compute_metrics=compute_metrics,
    )

    trainer.train()
    models.append(model)

```

4.2 Inference Strategy

Our inference implementation uses temperature sampling to promote diversity in generated code:

```

def generate_code(query, max_length=256):
    input_text = "generate code: " + query.lower()
    inputs = tokenizer(input_text, return_tensors="pt").to(device)

    with torch.no_grad():
        generated_ids = model.generate(
            input_ids=inputs.input_ids,
            attention_mask=inputs.attention_mask,
            max_length=max_length,
            do_sample=True,
            temperature=0.7,
            top_k=50,
            top_p=0.9,
            num_return_sequences=1
        )

    output = tokenizer.decode(generated_ids[0], skip_special_tokens=True)
    return output.strip()

```

4.3 Streamlit Application

We developed a Streamlit application for easy interaction with our model:

```

def generate_code(query, max_length=256, num_beams=5):
    input_text = "generate code: " + query.lower()
    inputs = tokenizer(input_text, return_tensors="pt",
                       truncation=True, padding=True).to(device)

    with torch.no_grad():
        generated_ids = model.generate(
            input_ids=inputs.input_ids,
            attention_mask=inputs.attention_mask,
            max_length=max_length,
            do_sample=True,
            temperature=0.7,
            top_k=50,
            top_p=0.9,

```

```

        num_return_sequences=1
    )

    output = tokenizer.decode(generated_ids[0], skip_special_tokens=True)
    return output.strip()

```

The Streamlit interface includes:

- Text input for custom queries
- Example query selection dropdown
- Syntax-highlighted code output
- Error handling for model loading and generation

5. Evaluation Results

5.1 Quantitative Metrics

Our evaluation on the test dataset yielded the following results:

Metric	Score
ROUGE-L	39.14
BLEU Score	25.66
Exact Match	0.09

These metrics indicate that while our model captures some of the semantic intent of the queries (shown by the ROUGE-L and BLEU scores), achieving exact matches remains challenging due to the multitude of valid ways to implement the same functionality in Python.

5.2 Qualitative Analysis

The model demonstrates the ability to generate syntactically correct and functionally appropriate code for a variety of query types:

Query 1: "find smallest element in the array"

Generated Code:

```

arr = []
size = int(input("enter the size of the array: "))
print("enter the element of the array:")
for i in range(0, size):
    num = int(input())
    arr.append(num)
print("smallest element of the array is:", min(arr))

```

This code correctly:

- Creates an empty array
- Takes user input for the array size
- Prompts for and accepts array elements
- Uses Python's built-in `min()` function to find the smallest element
- Outputs the result with appropriate labeling

Query 2: "Write a NumPy program to convert an array to a float type."

Generated Code:

```
import numpy as np
x = np.array([1.5, 1.4, 2.5, 3.6, 7.6, 8.5, 2.7, 9.0, 10])
print("original array:")
print(x)
print("\nconvert the array to a float type:")
print(np.floor(x.astype(float)))
```

This code correctly:

- Imports the NumPy library
- Creates a sample array with decimal values
- Prints the original array
- Converts the array to float type using `astype(float)`
- Applies the `floor()` function to round down values
- Formats the output with descriptive labels

Both examples demonstrate the model's ability to:

1. Understand the programming task intent
2. Generate syntactically correct Python code
3. Apply appropriate code structure for the problem
4. Use relevant Python libraries and functions
5. Provide clear output formatting

The model shows proficiency in handling different types of operations including basic array manipulation, user input handling, and NumPy-specific functions. The code structure follows common Python conventions with appropriate variable naming and logical flow.

6. Discussion

6.1 Limitations

Several limitations were observed in our approach:

1. **Low Exact Match Score:** The exact match score of 0.09 indicates that achieving precisely the same code as the reference is challenging, though this metric is strict and doesn't account for alternative correct implementations.
2. **Variable Code Quality:** The quality of generated code varies depending on the complexity and specificity of the query. More ambiguous queries may result in less precise implementations.
3. **Domain Specificity:** The model performs better on common programming patterns but struggles with domain-specific or complex algorithms that are underrepresented in the training data.

6.2 Comparison with Alternative Approaches

Alternative approaches to code generation include:

- **Rule-based Systems:** More predictable but less flexible
- **Retrieval-based Methods:** Higher accuracy but limited to seen examples
- **Full Model Pre-training:** Better performance but requires massive computational resources

Our fine-tuning approach offers a balanced compromise, leveraging pre-trained knowledge while adapting to our specific task with reasonable computational requirements.

6.3 Future Improvements

Potential improvements to our system include:

1. **Extended Dataset:** Expanding the training data with more diverse problem-solution pairs
2. **Code Execution Validation:** Incorporating runtime validation to ensure functional correctness
3. **Multi-language Support:** Extending the model to generate code in multiple programming languages
4. **Parameter-Efficient Fine-Tuning:** Implementing techniques like LoRA to reduce computational requirements
5. **Human Feedback Loop:** Incorporating user feedback to continuously improve the model

7. Conclusion

This paper presented a comprehensive approach to fine-tuning transformer models for code snippet generation. We demonstrated that the CodeT5 model can be effectively fine-tuned on a dataset of approximately 3,300 Python problem-solution pairs, achieving reasonable performance metrics. The integration with Streamlit provides a practical interface for developers to leverage this technology.

Our work contributes to the field of AI-assisted programming by:

1. Documenting a complete pipeline for fine-tuning code generation models
2. Implementing an ensemble approach with k-fold cross-validation
3. Demonstrating practical deployment through a web interface
4. Providing a detailed analysis of performance and limitations

Future work will focus on expanding the dataset, improving evaluation metrics, and incorporating more sophisticated fine-tuning techniques to enhance performance.

References

1. Wang et al., "CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation," EMNLP 2021.
2. Wolf et al., "HuggingFace's Transformers: State-of-the-art Natural Language Processing," arXiv 2020.
3. Vaswani et al., "Attention Is All You Need," NeurIPS 2017.
4. "Fine-tuning - Hugging Face," <https://huggingface.co/docs/transformers/en/training>.
5. "Code Generation Using Natural Language Processing," International Journal of Research Publication and Reviews, Vol (5), Issue (6), June (2024).
6. "Fine-Tuning Transformers with MLflow for Enhanced Model Management," MLflow Documentation, 2023.
7. "Transformer Models: NLP's New Powerhouse," Data Science Dojo.
8. "Exploring Parameter-Efficient Fine-Tuning Techniques for Code Generation," ACM Digital Library, 2024.

Appendix

Figure 1: Transformer Architecture Diagram

Figure 2: Streamlit Application Interface

Figure 3: Training and Evaluation Pipeline

Figure 4: Example Code Generations

✱✱