# Bridging the Language Gap: Enhancing Multilingual Prompt-Based Code Generation in LLMs via Zero-Shot Cross-Lingual Transfer

Mingda Li[1,2], Abhijit Mishra[2], Utkarsh Mujumdar[2]
[1]Department of Statistics and Data Science, Yale University, New Haven, USA
[2]School of Information, University of Texas at Austin, Austin, USA
Email: {mingdali, abhijitmishra, utkarsh.mujumdar}@utexas.edu

*Abstract*—**The use of Large Language Models (LLMs) for program code generation has gained substantial attention, but their biases and limitations with non-English prompts challenge global inclusivity. This paper investigates the complexities of multilingual prompt-based code generation. Our evaluations of LLMs, including CODELLAMA and CODEGEMMA, reveal significant disparities in code quality for non-English prompts; we also demonstrate the inadequacy of simple approaches like prompt translation, bootstrapped data augmentation, and fine-tuning. To address this, we propose a zero-shot cross-lingual approach using a neural projection technique, integrating a cross-lingual encoder like LASER to map multilingual embeddings from it into the LLM's token space. This method requires training only on English data and scales effectively to other languages. Results on a translated and quality-checked *MBPP* dataset show substantial improvements in code quality. This research promotes a more inclusive code generation landscape by empowering LLMs with multilingual capabilities to support the diverse linguistic spectrum in programming.**

*Index Terms*—**Multilingual Artificial Intelligence, Code Generation, Machine Translation, Transformer**

## I. INTRODUCTION

The use of Large Language Models (LLMs) for code generation, such as generating Python programs from problem specifications, has gained substantial interest due to their effectiveness in handling complex language tasks [2]–[4]. This capability has led to the development of innovative applications like GitHub Copilot [5] and specialized LLMs such as CodeLLaMa [1], underscoring the growing importance of this field. Although LLMs are globally prevalent and proficient in processing multilingual inputs, they often exhibit biases against non-English prompts [6], [7], which is particularly evident in code generation. Figure 1 under *Baseline Output* illustrates how the quality of generated code diminishes as prompts shift from English to other languages, a disparity linked to the availability of data used in LLM training and fine-tuning. Ensuring LLMs deliver consistent quality across languages is crucial for fostering fair and inclusive code generation, especially as the global programming community is increasingly composed of non-English speakers. Data from coding platforms like The Competitive Programming Hall

of Fame[1] highlight the skew towards non-English-speaking regions. As the global population of coders grows, addressing these biases in LLMs is essential to promoting an equitable and accessible environment for developers worldwide.

This paper first tries to bridge the language gap in multilingual prompt-based code generation by enhancing LLMs through self-supervised fine-tuning. To begin with, we evaluate the performance of LLMs such as CODELLAMA and GPT-4 on English and five non-English languages, by translating and quality-checking a sanitized version of the MBPP dataset [4]. Significant disparities in code quality across languages are observed, even when using the same problem statements. Inspired by previous work on improving LLM performance in multilingual tasks with Chain-of-Thought (CoT) [8], [9], and using bootstrapping for multilingual data generation [10], we form strong baselines for code generation with CoT and fine-tuning on bootstrapped multilingual data. However, these approaches showed only marginal and inconsistent improvements.

To address data sparsity and limited multilingual exposure, we propose a novel approach: (a) using a pre-trained multilingual encoder like LASER [11] to encode multilingual inputs into a joint vector space; (b) projecting these embeddings into the LLM's input space and aligning them through training solely on English data; (c) employing this LASER→LLM pipeline at inference for zero-shot cross-lingual processing of non-English inputs. This method familiarizes models with multiple languages without needing additional external training data. Our evaluation shows improved code quality and reduced syntax and logical errors, as illustrated in Figure 1 under *Zero-shot Cross-Lingual Inference*. Our approach integrates seamlessly as a minor training step that does not require any expensive pretraining or fine-tuning, thus offering a promising way to enhance LLMs' multilingual capabilities in code generation.

The contributions of the paper are as follows:

1) We create a novel multilingual test dataset with quality-checked translations and new evaluation metrics.

---

[1]https://cphof.org/countries

| Language | Prompt | Baseline Output | Zero-shot Cross Lingual Output |
|---|---|---|---|
| English | Write a python function to convert complex numbers to polar coordinates. | ```def convert_to_polar(z):
    """
    Takes a complex number z and returns its polar form.
    """
    r = abs(z)
    theta = cmath.phase(z)
    return r, theta``` <br> **Complete and Correct Output** | ```def convert_to_polar(z):
    r = abs(z)
    theta = cmath.phase(z)
    return r, theta``` <br> **Complete and Correct Output** |
| Chinese | 编写一个 python 函数将复数转换为极坐标。<br>(Gloss: Bianxie yige python hanshu jiang fushu zhuanhuanwei jizuobiao。) | ```def complex_to_polar(z):
    # 请在这里补充代码
    pass``` <br> **Complete but Logically Incorrect Output** | ```def convert_to_polar(complex_number):
    real = complex_number.real
    imag = complex_number.imag
    modulus = math.sqrt(real**2 + imag**2)
    argument = math.atan2(imag, real)
    return modulus, argument``` <br> **Complete and Correct Output** |
| Spanish | Escribe una función de Python para convertir números complejos en coordenadas polares. | La función debe recibir un número complejo como parámetro y devolver un tupla con la magnitud y el ángulo del número complejo en radianes. Por ejemplo, si se llama a la función con el número complejo 3+4j, la función debe devolver el tupla (5, 0.9272952180016122). La función debe utilizar la función cmath.polar() para calcular la magnitud <br> **Incomplete and Noisy Output** | ```def convert_to_polar_coordinates(z):
    r = abs(z)
    theta = np.angle(z)
    return (r, theta)``` <br> **Complete and Correct Output** |

**Fig. 1:** Disparity in output code generated by *CodeLLaMa-Instruct* model [1] with 7B parameters for the same problem statement given in multiple languages

2) We introduce a scalable projection technique by integrating the LASER multilingual encoder with popular open-source LLMs like CodeLLaMa [1], CodeGemma [12], and Mistral [13] for zero-shot cross-lingual code generation.

3) Our evaluation of the approach against Chain-of-Thought (CoT) and fine-tuning with multilingual bootstrapped data, highlights the strengths and limitations of each method.

We make our code[2] and multilingual evaluation data[3] publicly available for academic use.

## II. RELATED WORK

As AI technology advances, transformer-based LLMs like GPT [14], LLaMA [15], Mistral [13], and Gemma [16] have become prominent in research and applications. While pre-trained models such as LLaMA2 and Gemma are fine-tuned for code generation, their English-centric training data limits multilingual proficiency [17], [18]. Studies show these models face performance issues with non-English tasks [19], [20], and human supervision remains crucial for quality [21].

To address these gaps, Ahuja et al. [22] develop a multilingual benchmark for evaluating LLMs, revealing performance drops across languages. Tan et al. [23] and Huang et al. [24] suggest leveraging transfer learning and cross-lingual prompting to improve multilingual capabilities. Additional methods include language-specific pre-training [25] and consistency regularization for fine-tuning [26].

Optimizing prompts enhances multilingual LLM accuracy [27], [28], and CoT techniques improve code generation [29]. Fine-tuning with multilingual synthetic data, including translation and back-translation [30], [31], further refines LLMs [32]–[36]. Contrary to these popular approaches, we take an orthogonal route by using specialized multilingual encoders and lightweight projectors to bridge language gaps in popular LLMs. Our work is inspired by multimodal AI literature, integrating projection techniques for different modalities such as language, vision and speech [37]–[39].

## III. EXPERIMENTAL SETUP

This section details our experimental setup, including the creation of a multilingual benchmark dataset, the models evaluated, and the metrics used. Our focus is on Python code generation from multilingual prompts, though the methods and insights are applicable to other languages and contexts.

### A. Evaluation Dataset

To obtain datasets with multilingual prompts for code generation, we adapted the Mostly Basic Programming Problems (MBPP) dataset [4], specifically the sanitized version with its "test" split, containing 257 problems with solutions and three test cases each. We translate these prompts into five languages—Chinese-Simplified (*zh-cn*), Spanish (*es*), Japanese (*ja*), Russian (*ru*), and Hindi (*hi*)—using the Google Translate

---
[2]https://github.com/lmd0420/Multilingual_Code_Gen
[3]https://huggingface.co/datasets/AnonymousDoe/MBPP-Translated

| Translation | A1 | A2 | Agreement (%) |
|---|---|---|---|
| en_es | 0.94 | 0.96 | 89.69 |
| en_hi | 0.93 | 0.96 | 89.11 |
| en_ja | 0.93 | 0.96 | 89.88 |
| en_ru | 0.93 | 0.96 | 90.43 |
| en_zh-cn | 0.94 | 0.96 | 90.79 |

**TABLE I:** Human Evaluation of Translated Prompts. Two distinct bilingual speakers from MTurk rated translations with 1 (acceptable) or 0 (not acceptable) for each translation. A1 and A2 represent their average scores.

| Lang. Pair | Average Rating | St.Dev |
|---|---|---|
| en_hi | 4.88 | 0.40 |
| en_es | 4.90 | 0.48 |
| en_ru | 4.95 | 0.30 |
| en_zh-cn | 4.93 | 0.39 |
| en_ja | 4.87 | 0.55 |

**TABLE II:** Average Rating and Standard Deviation for Translation from English to Other Languages

API[4]. Our selection covers languages with similar characters to English (e.g., Spanish) as well as those with entirely different characters (e.g., Japanese), demonstrating robustness across different language families.

Translation quality was assessed by:

- (a) expert bilingual speakers via *Amazon Mechanical Turk*, who rate translations as acceptable or not. (Note: guidelines were provided for binary rating and consent was obtained to report the statistics in the paper.)
- (b) GPT-4, which rate translations on a scale of 1 to 5.

Human-evaluated results in table I show superior translation quality. Table II presents GPT-4's ratings, again indicating that translations are of high-quality with high mean scores and low standard deviations.

### B. Models Used for Evaluation

We consider three open source variants of instruction-tuned models for evaluation, namely CODELLAMA-7B[5], CODEGEMMA-7B[6] and MISTRAL-7B-V0.3[7]. These models are specialized versions of their base models to programming-related tasks. They have demonstrated greater efficacy at code generation, infilling, and debugging capabilities compared to the standard versions. Models are obtained from HuggingFace hub[8]. For benchmarking, we adopt GPT-4 as the reference system (*a.k.a Skyline*) because it has been extensively trained on various languages, and has proven effective in various tasks, including code generation.

---

[4]https://cloud.google.com/translate
[5]codellama/CodeLlama-7b-Instruct-hf
[6]google/codegemma-7b-it
[7]mistralai/Mistral-7B-Instruct-v0.3
[8]http://huggingface.co

### C. Inference Pipeline and Evaluation Metrics

We develop a pipeline to process task descriptions in various languages. The pipeline feeds these prompts into the models and variants described in Sections IV and V and stores the results. We aim to provide a solution for people with little programming experience who use a variety of languages, allowing them to access LLMs' coding capabilities through natural language. Thus, we choose to report our results on from-scratch code generation tasks in the zero-shot setting. Python code from the outputs is extracted using regular expressions. We then automatically identify function names in the code, replacing the function names in the MBPP test assertions with those from the model outputs, ensuring any function names would not affect the test results. Finally, we generate bash scripts to execute the extracted code and assertions, and measure the following metrics:

- **Total Error Rate (TotalER):** The ratio of code samples that fail at least one test case, to the total number of samples. Lower is better.
- **Logical Error Rate (LER):** The ratio of code samples that execute successfully but produce incorrect results, to the total number of samples. Lower is better.
- **Syntax Error Rate (SER):** The ratio of code samples containing syntax errors, to the total number of samples. Lower is better.
- **All Tests Passed Rate (ATPR):** The ratio of code samples that pass all given test cases, to the total number of samples. Higher is better.

Additionally, we also observe the *Code Completion Rate* as a supplementary metric, which indicates the proportion of complete codes in model responses. A higher value represents a better result.
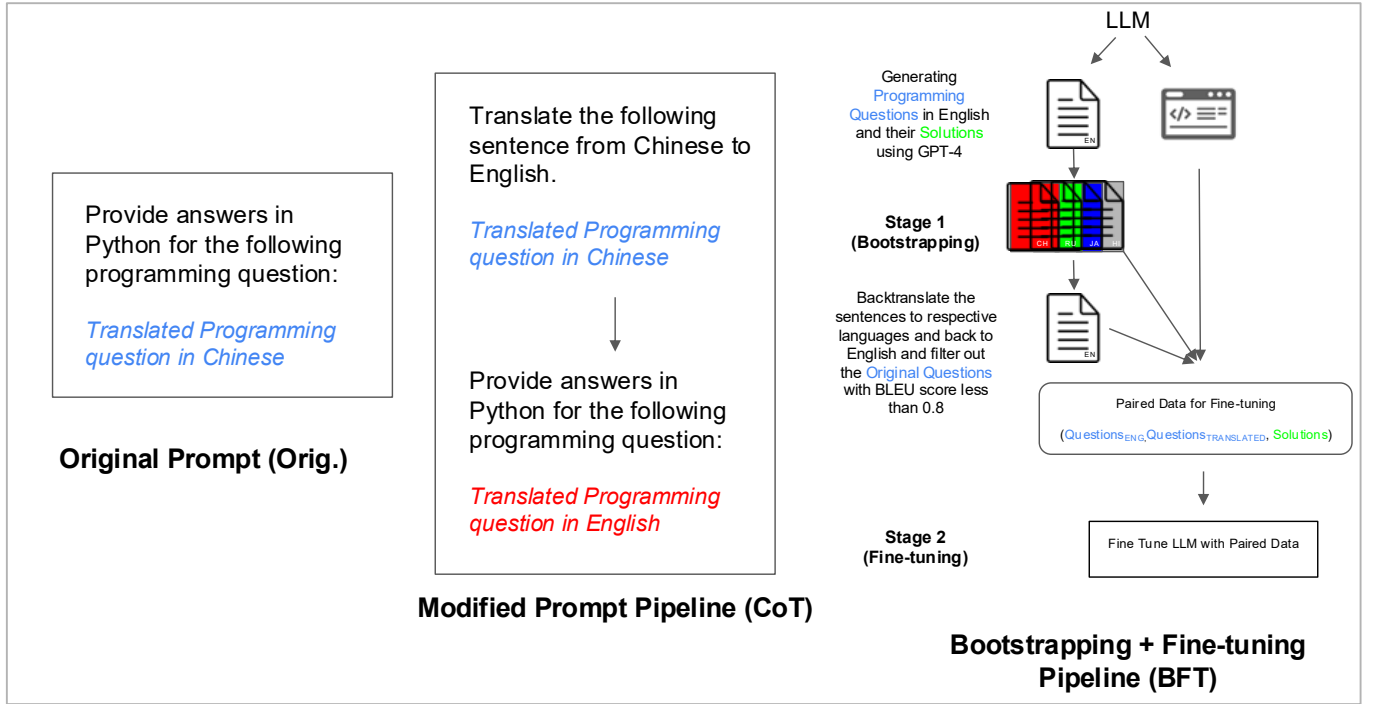
With this setup, we can now evaluate LLM code generation quality and propose mitigation strategies. Our approach and baselines are summarized in Figure 2, detailed in the following section.
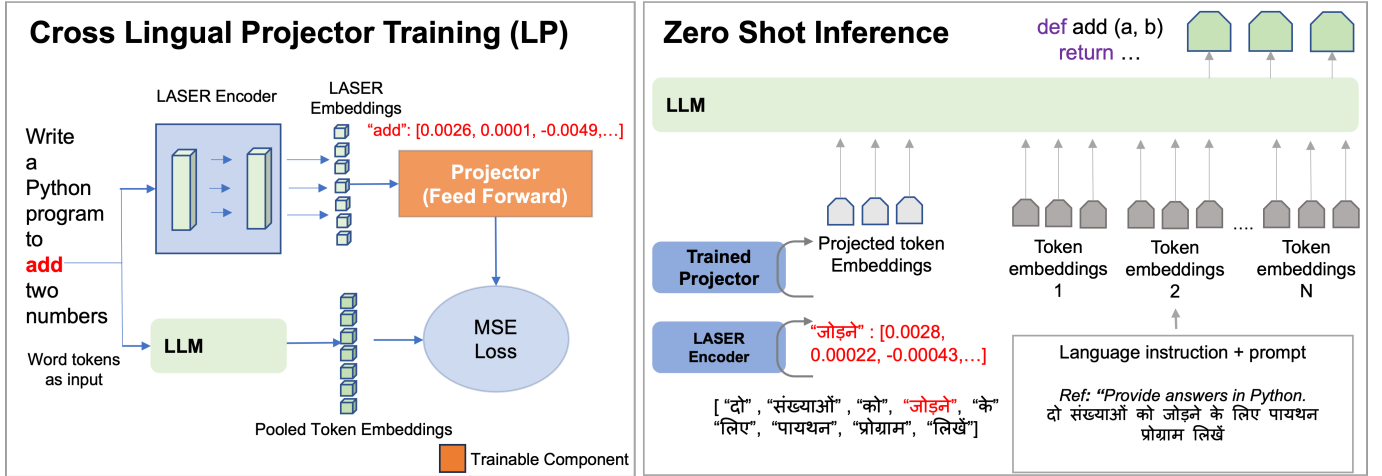
### IV. ISSUES WITH TRIVIAL BASELINES

Given that language models exhibit emergent capabilities and scale effectively across tasks and languages, efficient prompting and prompt tuning are generally preferred over costly training or fine-tuning that demands extensive data curation. Based on our experimental setup, we highlight the challenges LLMs face with multilingual code generation in conventional settings, providing a detailed analysis of existing models' performance and their limitations. Throughout this section, we will reference Table III for a comprehensive discussion of the results.

### A. Baseline 1. Original Prompt

Here, each query in the dataset is passed through the pipeline, where the model generates response code, filtered from extraneous information such as code explanations, and executed using an automatically constructed bash script. The results are presented in first column of each section of Table III, with the following key observations:

**(a)** Baselines with direct prompting, Chain of Thought (CoT) and fine-tuning with bootstrapped data



**(b)** Our proposed approach based on cross lingual encoder and projector training and zero shot inference

**Fig. 2:** Explored approaches

GPT-4, recognized for its robustness and extensive engineering, reliably generates code across all language prompts, though with slightly varying error profiles - except for Hindi and Chinese. In contrast, open-source models like CodeLLaMa show more pronounced disparities between languages, with higher error rates and lower all-tests-passed rates compared to English. Notably, some models, such as CodeLLaMa-Instruct-7B, perform better in non-English languages like Spanish. This may seem unusual but aligns with findings from Chen et al. [40], which show that LLaMa 7B, when instruction-tuned for multilingual tasks, performs better in Spanish than English. Since CodeLLaMa is based on this instruction-tuned model, this could explain the atypical performance in Spanish. Over-

all, these results highlight a lack of consistency in code output quality as the language changes. We use the abbreviation *Orig.* to refer to this baseline henceforth.

### B. Chain-of-Thought with Back-Translation

Due to uneven language representation in LLM training datasets, achieving consistent results with direct prompting is challenging. A potential solution is to use back-translation: translate non-English prompts into English and use the English version as the query. This CoT approach involves translating the problem statement with the prompt: `Translate the sentence $PROBLEM from $TARGET-LANG to English`, then generating code outputs from the translated

| LLM | Lang | TotalER↓ | | | | LER↓ | | | | SER↓ | | | | ATPR↑ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Orig. | CoT | BFT | LP | Orig. | CoT | BFT | LP | Orig. | CoT | BFT | LP | Orig. | CoT | BFT | LP |
| GPT-4 (Skyline) | en | *58.37* | - | - | - | *10.9* | - | - | - | *47.47* | - | - | - | *41.63* | - | - | - |
| | es | *62.65* | - | - | - | *12.85* | - | - | - | *49.8* | - | - | - | *37.35* | - | - | - |
| | hi | *67.7* | - | - | - | *17.9* | - | - | - | *49.8* | - | - | - | *32.3* | - | - | - |
| | ja | *64.2* | - | - | - | *13.62* | - | - | - | *50.58* | - | - | - | *35.8* | - | - | - |
| | ru | *65.37* | - | - | - | *17.12* | - | - | - | *48.25* | - | - | - | *34.63* | - | - | - |
| | zh | *67.7* | - | - | - | *16.73* | - | - | - | *50.97* | - | - | - | *32.3* | - | - | - |
| Code LLaMa-7B | en | 87.16 | - | 82.1 | **75.49** | 63.04 | - | 28.79 | **22.57** | **24.12*** | - | 53.31 | 52.92 | 12.84 | - | 17.9 | **24.51** |
| | es | 79.77 | 91.83 | **81.71** | 81.71 | 28.8 | 56.81 | 26.07 | 24.9 | 50.97 | **35.02*** | 55.64 | 56.81 | **20.23** | 8.17 | 18.29 | 18.29 |
| | hi | 96.5 | 97.66 | 96.5 | **95.72** | 65.37 | 61.08 | 61.87 | 25.29 | **31.13** | 36.58 | 34.63 | 70.43 | 3.5 | 2.34 | 3.5 | **4.28** |
| | ja | 89.49 | 84.82 | 84.82 | **84.44** | 50.58 | 52.91 | 34.24 | 22.96 | 38.91 | **31.91** | 50.58 | 61.48 | 10.51 | 15.18 | 15.18 | **15.56** |
| | ru | **82.1** | 86.38 | 85.21 | 82.88 | 39.69 | 61.87 | 31.51 | 23.35 | 42.41 | **24.51** | 53.7 | 59.53 | 17.9 | 13.62 | 14.79 | **17.12** |
| | zh | 93.77 | 96.5 | 88.72 | **82.1** | 77.43 | 73.15 | 35.41 | 26.46 | **16.34** | 23.35 | 53.31 | 55.64 | 6.23 | 3.5 | 11.28 | **17.9** |
| Code Gemma-7B | en | 82.1 | - | 92.22 | **77.04** | 41.63 | - | 63.04 | 25.68 | 40.47 | - | **29.18** | 51.36 | 17.9 | - | 7.78 | **22.96** |
| | es | 86.38 | 89.1 | 91.05 | **77.82** | 47.86 | 42.02 | 57.59 | 24.51 | 38.52 | 47.08 | **33.46** | 53.31 | 13.62 | 10.9 | 8.95 | **22.18** |
| | hi | 89.49 | 91.05 | 94.16 | **81.71** | 49.41 | 50.51 | 74.71 | 29.18 | 40.08 | 40.47 | **19.45** | 52.53 | 10.51 | 8.95 | 5.84 | **18.29** |
| | ja | 83.66 | 90.27 | 91.05 | **79.77** | 38.91 | 44.75 | 50.58 | 24.13 | 44.75 | 45.52 | **40.47** | 55.64 | 16.34 | 9.73 | 8.95 | **20.23** |
| | ru | 85.99 | 88.72 | 89.1 | **77.04** | 42.41 | 48.25 | 59.53 | 25.68 | 43.58 | 40.47 | **29.57** | 51.36 | 14.01 | 11.28 | 10.9 | **22.96** |
| | zh | 84.82 | 86.38 | 93.0 | **79.38** | 39.68 | 48.64 | 62.26 | 28.02 | 45.14 | 37.74 | **30.74** | 51.36 | 15.18 | 13.62 | 7.0 | **20.62** |
| Mistral -7B-v0.3 | en | 85.21 | - | 92.61 | **83.27** | 35.41 | - | 28.41 | 27.24 | 49.8 | - | 64.2 | 56.03 | 14.79 | - | 7.39 | **16.73** |
| | es | 87.55 | 86.38 | 94.94 | **84.82** | 39.69 | 29.18 | 26.46 | 26.06 | 47.86 | 57.2 | 68.48 | 58.76 | 12.45 | 13.62 | 5.06 | **15.18** |
| | hi | 91.44 | **91.05** | 98.83 | 92.22 | 35.41 | 35.41 | **24.12** | 30.74 | 56.03 | **55.64** | 74.71 | 61.48 | 8.56 | **8.95** | 1.17 | 7.78 |
| | ja | 88.72 | **86.77** | 96.11 | 87.55 | 35.8 | 31.91 | 28.02 | 22.57 | **52.92** | 54.86 | 68.09 | 64.98 | 11.28 | **13.23** | 3.89 | 12.45 |
| | ru | 85.6 | **84.05** | 95.33 | 84.05 | 33.85 | 30.74 | 26.85 | 24.13 | **51.75** | 53.31 | 68.48 | 59.92 | 14.4 | **15.95** | 4.67 | 15.95 |
| | zh | 88.72 | 87.16 | 94.55 | **84.05** | 39.3 | 30.74 | **26.07** | 26.85 | **49.42** | 56.42 | 68.48 | 57.2 | 11.28 | 12.84 | 5.45 | **15.95** |

**TABLE III:** Comprehensive comparison of different models across multiple languages and configurations. TotalER: Total Error Rate, LER: Logical Error Rate, SER: Syntax Error Rate, ATPR: All Test Passed Rate. *Orig:* Directly Querying LLMs, *CoT*: Chain of Thought with Translation, *BFT*: Fine tuning on Bootstrapped Multilingual Data, *LP* (Our approach): Fine tuning on Multilingual Projection with LASER Encoders

prompt. Our experiments, detailed in the second column of Table III, show that back-translation does not significantly improve results. In some cases, it even reduce performance, as indicated by lower ATPR scores. Qualitative analysis suggests that models struggle with non-canonical language representations and topic drift, despite the translations not being of poor quality. We use the abbreviation *CoT* to refer to this baseline henceforth.

### C. Bootstrapping Multilingual Data and Fine Tuning

Fine-tuning pre-trained models is effective for many NLP tasks but is often resource-intensive, requiring costly and time-consuming task-specific labeled data. Instead of manually creating such data for multiple languages - designing prompts, validating answers, and translating while preserving semantic meaning - we use a bootstrapping approach. In this method, we utilize a powerful LLM like ChatGPT to generate English programming problems and their answers. These problems are then translated into target languages and back-translated into English. We assess the similarity of translations using the BLEU score [41], retaining translations that meet a quality threshold (e.g., 0.8) to create new training data. This method preserves text quality in target languages and allows the model to validate its translations, as detailed in Algorithm 1.

After bootstrapping data for all target languages, we shuffle and use it to fine-tune the LLMs with a single A100 GPU. Models are quantized to FP16 and fine-tuned using parameter-efficient techniques, including low-rank adaptation [42]. We set the temperature to 0.8 for consistency and use two epochs.

---

**Algorithm 1** Bootstrap Training Data

```
1: function BOOTSTRAPDATA(LLM,Lang)
2:     n ← number of attempts
3:     threshold ← 0.9
4:     Initialize a query set Q ← {}
5:     Initialize training data TD ← {}
6:     squery ← "Generate 100 python problems"
7:     trquery ← "Translate from English into Lang"
8:     btrquery ← "Translate from Lang into English"
9:     for i ← 1 to n do
10:        q ← LLM(squery)
11:        Push q into query set Q
12:    end for
13:    for q in Q do
14:        a ← LLM(q)
15:        Push a into answer set A
16:        Push < q, a > into TD
17:    end for
18:    for q, a in TD do
19:        t ← LLM(trquery, q)
20:        bt ← LLM(btrquery, t)
21:        score ← BLEU(t,bt)
22:        if score > threshold then
23:            Push < t, a > into TD
24:        end if
25:    end for
26:    return TD
27: end function
```

As shown in the third column of Table III, while bootstrapping with ChatGPT reduces syntax errors, it also increases hallucinations, leading to lower test pass rates and higher total errors. This suggests that the model, although producing more complete code, struggles with accuracy and reliability. We use the abbreviation *BFT* to refer to this baseline henceforth.

## V. Our Approach: Projection-Based Zero-Shot Transfer

Our approach focuses on avoiding the use of in-language training data, which can be costly and impractical. Instead, we utilize an intermediate, lightweight method that relies on abundant English data and the LASER multilingual encoder [11], which provides joint embeddings for over 100 languages. In this setup, the LASER encoder preprocesses and embeds input tokens before passing them to the LLM, which then operates on these embeddings rather than raw input IDs. This method enables efficient language scaling, as similar meanings are represented consistently across languages (e.g., the English token "add" and its Hindi counterpart "JoDaNe" are embedded similarly, as shown in Figure 2 part (b)).

Two key challenges arise with this approach: (A) differing tokenization between the multilingual encoder and the LLM, and (B) the LLM's unfamiliarity with the multilingual embeddings. To address (A), we use word tokens and extract mean-pooled embeddings from subwords using tokenizers such as NLTK [9] for space-sparated lanague inputs, Jieba [10] for Chinese, and Janome[11] for Japanese. We then train a projector to align these embeddings. For a given word token, we compute the LLM's subword embeddings ($\hat{H}_{llm}$) through max pooling, and the multilingual embeddings ($H_{laser}$) from the LASER encoder. The projector, with learnable parameters $\mathbf{W}_{llm}$ and $\mathbf{b}_{llm}$, is defined as:

$$\mathbf{H}_{llm} = \mathbf{W}_{llm} \cdot \mathbf{H}_{laser} + \mathbf{b}_{llm}$$

The model is trained by minimizing the Mean Squared Error (MSE) between $\hat{H}_{llm}$ and $\mathbf{H}_{llm}$:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^{N} \left\| \hat{H}_{llm}^{i} - \mathbf{H}_{llm}^{i} \right\|^2$$

where $N$ is the number of word tokens. Training utilizes English tokens from the MBPP dataset, which includes 120 examples. We train the projector for 200 epochs on a single consumer-grade NVIDIA RTX-4060; the entire run finishes in under one hour. In contrast, LoRA fine-tuning or re-pretraining on extra multilingual data typically requires at least several GPU hours or even days, highlighting an orders-of-magnitude efficiency advantage for our lightweight projector training.

During inference, tokens are first word-tokenized and embedded using LASER, then projected, and finally input to

[9]https://www.nltk.org
[10]https://github.com/fxsjy/jieba
[11]https://mocobeta.github.io/janome/en/

the LLM for multilingual processing without requiring in-language data. To enhance performance and align with baselines, we also concatenate system prompt embeddings with the original programming prompt embeddings. Notably, LASER embeddings are of size 1024, while LLM embeddings are typically 4096 or larger, necessitating a 4-fold upsampling. We achieve this using two linear projection layers as outlined in the above equations. We use the abbreviation *LP* to explain this system henceforth.

## VI. Results and Discussions

Table III presents the overall performance models and variants discussed in sections IV and V. Our observations indicate that across all metrics, our proposed model consistently reduces the performance gap between English and non-English languages, as reflected in the differences and deviations. This improvement is particularly evident when comparing the direct querying setup (Orig.) with our multilingual projector-based variant (LP), where deviations from English are generally smaller. It is also worth noting that the projection method we proposed completely eliminates the need for external data. Therefore, we compared it with methods like bootstrapping from the model itself followed by finetuning (BFT) and Chain-of-Thought (CoT), which also do not rely on external data or models. We explore the details of each metric below.

### A. Total Error Rate (TotalER)

TotalER is an important metric that quantifies the overall error rate of the generated code. Our proposed method, LP, consistently achieves the lowest TotalER across nearly all languages and models, demonstrating its effectiveness. For example, with the CodeLLaMa-7B model, LP significantly reduces the TotalER to 75.49 for English (en) and 82.1 for Chinese (zh), outperforming the Orig. and other methods. This improvement is especially pronounced in languages with complex syntax and morphology, such as Hindi (hi) and Russian (ru), where LP reduces the TotalER by over 10% in some cases compared to the original model. Even in cases where LP is the second-best, its performance is very close to the top-performing method, highlighting its reliability. In contrast, BFT, a strong trivial baseline, tends to increase the TotalER due to hallucinations, as observed in our data analysis, despite slightly improving the all test cases passed metric.

### B. Logical Error Rate (LER)

LER is a critical component of the total error, measuring the proportion of code samples that execute without errors but produce incorrect results. A lower LER indicates a model's ability to generate logically sound code, making it a key metric for evaluating performance. It's important to note that we classify a logical error not only when no valid code is generated but also when any of the test cases fail.

Our approach, LP, consistently outperforms other methods in terms of LER, with only a few exceptions where the difference is marginal and still better than other candidates. For instance, with the CodeGemma-7B model, LP achieved an
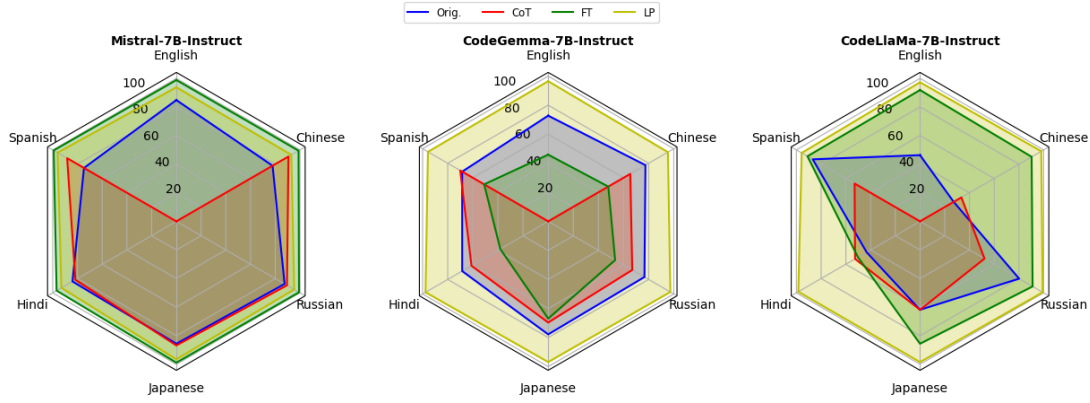
**Fig. 3:** Code Completion Rate (CCR) for Models and Languages, with LP represented by perfect polygons, always above 90%

LER of 25.68 for English, significantly lower than the 41.63 in Orig and 63.04 in BFT. This trend is also evident in other languages, such as Spanish (es) and Japanese (ja), where LP substantially reduces LER, underscoring its effectiveness in ensuring logical correctness across multilingual scenarios.

*C. Syntax Error Rate (SER)*

SER is a component of total error and indicates the proportion of code samples that contain syntax errors. A lower SER reflects the model's ability to generate syntactically correct code. Our overall observations with respect to this metric is models like ours that often produce code than omitting it (which is indicated by the higher code completion rate) are more prone to syntax error due to the high recall. While syntax error solving is a crucial step in program debugging, we believe such a form of error is slightly easier for beginners to solve than logical errors. Thus, given that LP consistently achieves the lowest LER across all languages and models, we consider this demonstrates the LP's proficiency in helping with generating error-free code, particularly in linguistically diverse contexts.

*D. All Test Passed Rate (ATPR)*

ATPR measures the proportion of code samples that successfully pass all given test cases. A higher ATPR signifies greater reliability of the generated code, making it a crucial metric. Our observations show that LP consistently outperforms other methods in terms of ATPR across most cases. Also, the improvement in English performance may be due to LASER slightly altering the way text enters the model. However, there are exceptions with the Mistral-7B-v0.3 model in a few languages. This model, being more recent, benefits from enhanced multilingual capabilities due to its diverse pretraining datasets and extended vocabulary. Overall, ATPR improvements are consistent across other languages, highlighting LP's superior performance in generating reliable and functional code.

Our observations using Multilingual Projections with LASER encoders reveal that LP not only reduces errors but also enhances the logical correctness and reliability of the generated code, establishing it as the leading approach for multilingual Python code generation. Additionally, we analyze the Code Completion Rate (CCR) to assess the robustness of these models in generating meaningful code rather than nonsensical explanations across languages. LP consistently outperforms other variants in this regard, as shown in the spider graph in Figure 3. This graph illustrates LP's strong performance in producing complete code across all languages. Notably, the shapes representing LP in the graph are perfect polygons, reflecting its consistent behavior and reliability across different languages.

## VII. CONCLUSIONS AND FUTURE WORK

In this paper, we demonstrate the significant potential of Large Language Models to bridge language gaps and promote inclusivity in multilingual prompt-based code generation. While LLMs exhibit promising capabilities across various languages, their performance can be inconsistent, particularly with non-English prompts. Our comprehensive analysis and evaluation using a revised benchmark dataset revealed both strengths and limitations in multilingual code generation, highlighting areas needing improvement.

We showcase the limited effectiveness of traditional Chain-of-Thought method, as well as bootstrapping multilingual training data and fine-tuning LLMs to enhance code generation quality across multiple languages. Our zero-shot cross-lingual transfer approach, utilizing projected embeddings, proves effective, as evidenced by improved ATPR and reduced TotalER values. A single-layer projector already narrows most of the cross-lingual gap to within a 2-percentage-point range; this suggests that more complex projector architectures are unnecessary, as their potential accuracy gains are marginal while their memory and compute overheads would grow substantially.

Our lightweight method eliminates the need for extensive external multilingual data, maximizing the model's potential internally, and can be easily extended to other multilingual encoder + LLM combinations. Because the projector operates purely in embedding space, exactly the same training recipe can transfer seamlessly from the original LSTM-based LASER

encoder to the newer transformer-based LASER3 [43] variant, underscoring the generality of our approach.

Future work will expand this approach to include more languages, diverse prompt patterns, and programming languages beyond Python. Our findings underscore the importance of advancing these techniques to enhance LLM adaptability and utility for a global audience, stressing the need for ongoing efforts to improve their effectiveness and versatility in diverse linguistic contexts.

## VIII. Limitations

A major limitation of this work lies in the reliance on word tokenization and pooled token embeddings, which introduces external dependencies and may not scale effectively to extremely low-resource languages where tokenizers are not readily available. Furthermore, the sequence of projected embeddings from the target language can significantly differ from the canonical English order, potentially hindering the model's ability to fully leverage these embeddings. This misalignment could contribute to the generation of hallucinatory and erroneous outputs. To address this issue, some degree of fine-tuning of LLMs with denoising objectives may be necessary.

Additionally, our study focuses solely on generating code from scratch and does not cover code-filling scenarios, which is another important aspect that warrants future exploration. Due to resource constraints, the scope of this study has been limited to Python, but expanding the approach to encompass other general-purpose and special-purpose programming languages, as well as more benchmarks [44], [45] is essential for broader applicability.

## Acknowledgment

## References

[1] B. Roziere, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, T. Remez, J. Rapin *et al.*, "Code llama: open foundation models for code," arXiv preprint arXiv:2308.12950, 2023.

[2] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, "A survey of large language models," arXiv preprint arXiv:2303.18223, 2023.

[3] L. Gao, A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, and G. Neubig, "PAL: program-aided language models," in *Int. Conf. Machine Learning*, 2023, pp. 10 764–10 799.

[4] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le *et al.*, "Program synthesis with large language models," arXiv preprint arXiv:2108.07732, 2021.

[5] B. Yetistiren, I. Ozsoy, and E. Tuzun, "Assessing the quality of github copilot's code generation," in *Proc. 18th Int. Conf. Predictive Models and Data Analytics in Software Engineering*, 2022, pp. 62–71.

[6] Z. Talat, A. Névéol, S. Biderman, M. Clinciu, M. Dey, S. Longpre, S. Luccioni, M. Masoud, M. Mitchell, D. Radev *et al.*, "You reap what you sow: on the challenges of bias evaluation under multilingual settings," in *Proc. BigScience Episode #5 - Workshop on Challenges & Perspectives in Creating Large Language Models*, 2022, pp. 26–41.

[7] M. Choudhury and A. Deshpande, "How linguistically fair are multilingual pre-trained language models?" in *Proc. AAAI Conf. Artificial Intelligence*, vol. 35, 2021, pp. 12 710–12 718.

[8] F. Shi, M. Suzgun, M. Freitag, X. Wang, S. Srivats, S. Vosoughi, H. W. Chung, Y. Tay, S. Ruder, D. Zhou *et al.*, "Language models are multilingual chain-of-thought reasoners," arXiv preprint arXiv:2210.03057, 2022.

[9] L. Qin, Q. Chen, F. Wei, S. Huang, and W. Che, "Cross-lingual prompting: improving zero-shot chain-of-thought reasoning across languages," arXiv preprint arXiv:2310.14799, 2023.

[10] A. Awasthi, N. Gupta, B. Samanta, S. Dave, S. Sarawagi, and P. Talukdar, "Bootstrapping multilingual semantic parsers using large language models," arXiv preprint arXiv:2210.07313, 2022.

[11] M. Artetxe and H. Schwenk, "Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond," *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 597–610, 2019.

[12] C. Team, "Codegemma: open code models based on gemma," arXiv preprint arXiv:2406.11409, 2024.

[13] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier *et al.*, "Mistral 7b," arXiv preprint arXiv:2310.06825, 2023.

[14] OpenAI, "GPT-4 technical report," arXiv preprint arXiv:2303.08774, 2023.

[15] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, "Llama 2: open foundation and fine-tuned chat models," arXiv preprint arXiv:2307.09288, 2023.

[16] G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Rivière, M. S. Kale, J. Love *et al.*, "Gemma: open models based on gemini research and technology," arXiv preprint arXiv:2403.08295, 2024.

[17] V. D. Lai, N. T. Ngo, A. P. B. Veyseh, H. Man, F. Dernoncourt, T. Bui, and T. H. Nguyen, "Chatgpt beyond english: towards a comprehensive evaluation of large language models in multilingual learning," arXiv preprint arXiv:2304.05613, 2023.

[18] C. Akiki, G. Pistilli, M. Mieskes, M. Gallé, T. Wolf, S. Ilić, and Y. Jernite, "Bigscience: a case study in the social construction of a multilingual large language model," arXiv preprint arXiv:2212.04960, 2022.

[19] F. Shi, M. Suzgun, M. Freitag, X. Wang, S. Srivats, S. Vosoughi, H. W. Chung, Y. Tay, S. Ruder, D. Zhou *et al.*, "Language models are multilingual chain-of-thought reasoners," arXiv preprint arXiv:2210.03057, 2022.

[20] B. A. Becker, P. Denny, J. Finnie-Ansley, A. Luxton-Reilly, J. Prather, and E. A. Santos, "Programming is hard-or at least it used to be: educational opportunities and challenges of ai code generation," in *Proc. 54th ACM Tech. Symp. Comput. Sci. Educ. V. 1*, 2023, pp. 500–506.

[21] S. Sarsa, P. Denny, A. Hellas, and J. Leinonen, "Automatic generation of programming exercises and code explanations using large language models," in *Proc. 2022 ACM Conf. Int. Comput. Educ. Res.*, 2022, pp. 27–43.

[22] K. Ahuja, H. Diddee, R. Hada, M. Ochieng, K. Ramesh, P. Jain, A. Nambi, T. Ganu, S. Segal, M. Axmed *et al.*, "MEGA: multilingual evaluation of generative AI," arXiv preprint arXiv:2303.12528, 2023.

[23] L. Tan and O. Golovneva, "Evaluating cross-lingual transfer learning approaches in multilingual conversational agent models," arXiv preprint arXiv:2012.03864, 2020.

[24] H. Huang, T. Tang, D. Zhang, W. X. Zhao, T. Song, Y. Xia, and F. Wei, "Not all languages are created equal in llms: improving multilingual capability by cross-lingual-thought prompting," arXiv preprint arXiv:2305.07004, 2023.

[25] J. Pfeiffer, N. Goyal, X. Lin, X. Li, J. Cross, S. Riedel, and M. Artetxe, "Lifting the curse of multilinguality by pre-training modular transformers," in *Proc. 2022 Conf. North American Chapter Assoc. Comput. Linguistics: Human Lang. Technol.*, 2022, pp. 3479–3495.

[26] B. Zheng, L. Dong, S. Huang, W. Wang, Z. Chi, S. Singhal, W. Che, T. Liu, X. Song, and F. Wei, "Consistency regularization for cross-lingual fine-tuning," in *Proc. 59th Annu. Meet. Assoc. Comput. Linguistics*, 2021, pp. 3403–3417.

[27] M. Zhao and H. Schütze, "Discrete and soft prompting for multilingual models," arXiv preprint arXiv:2109.03630, 2021.

[28] L. Huang, S. Ma, D. Zhang, F. Wei, and H. Wang, "Zero-shot cross-lingual transfer of prompt-based tuning with a unified multilingual prompt," arXiv preprint arXiv:2202.11451, 2022.

[29] Y. Ma, Y. Yu, S. Li, Y. Jiang, Y. Guo, Y. Zhang, Y. Xie, and X. Liao, "Bridging code semantic and llms: semantic chain-of-thought prompting for code generation," arXiv preprint arXiv:2310.10698, 2023.

[30] R. Sennrich, B. Haddow, and A. Birch, "Improving neural machine translation models with monolingual data," arXiv preprint arXiv:1511.06709, 2015.

[31] C. D. V. Hoang, P. Koehn, G. Haffari, and T. Cohn, "Iterative back-translation for neural machine translation," in *2nd Workshop Neural Machine Translation and Generation*, 2018, pp. 18–24.

[32] X. Li, P. Yu, C. Zhou, T. Schick, L. Zettlemoyer, O. Levy, J. Weston, and M. Lewis, "Self-alignment with instruction backtranslation," arXiv preprint arXiv:2308.06259, 2023.

[33] Y. Chai, S. Wang, C. Pang, Y. Sun, H. Tian, and H. Wu, "ERNIE-Code: beyond english-centric cross-lingual pretraining for programming languages," arXiv preprint arXiv:2212.06742, 2023.

[34] T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, M. Gallé *et al.*, "BLOOM: a 176b-parameter open-access multilingual language model," arXiv preprint arXiv:2211.05100, 2023.

[35] T. Nakamura, M. Mishra, S. Tedeschi, Y. Chai, J. T. Stillerman, F. Friedrich, P. Yadav, T. Laud, V. M. Chien, T. Y. Zhuo *et al.*, "Aurora-m: the first open source multilingual language model red-teamed according to the u.s. executive order," arXiv preprint arXiv:2404.00399, 2024.

[36] S. Zhang, C. Gao, W. Zhu, J. Chen, X. Huang, X. Han, J. Feng, C. Deng, and S. Huang, "Getting more from less: large language models are good spontaneous multilingual learners," arXiv preprint arXiv:2405.13816, 2024.

[37] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," in *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[38] Y. Fathullah, C. Wu, E. Lakomkin, J. Jia, Y. Shangguan, K. Li, J. Guo, W. Xiong, J. Mahadeokar, O. Kalinli *et al.*, "Prompting large language models with speech recognition abilities," in *ICASSP 2024-2024 IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP)*, 2024, pp. 13 351–13 355.

[39] L. Beyer, A. Steiner, A. S. Pinto, A. Kolesnikov, X. Wang, D. Salz, M. Neumann, I. Alabdulmohsin, M. Tschannen, E. Bugliarello *et al.*, "Paligemma: a versatile 3b vlm for transfer," arXiv preprint arXiv:2407.07726, 2024.

[40] P. Chen, S. Ji, N. Bogoychev, A. Kutuzov, B. Haddow, and K. Heafield, "Monolingual or multilingual instruction tuning: which makes a better alpaca," arXiv preprint arXiv:2309.08958, 2024.

[41] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proc. 40th Annu. Meet. Assoc. Comput. Linguistics*, 2002, pp. 311–318.

[42] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: low-rank adaptation of large language models," arXiv preprint arXiv:2106.09685, 2021.

[43] K. Heffernan, O. Çelebi, and H. Schwenk, "Bitext mining using distilled sentence representations for low-resource languages," 2022. [Online]. Available: https://arxiv.org/abs/2205.12654

[44] Z. Wang, S. Zhou, D. Fried, and G. Neubig, "Execution-based evaluation for open-domain code generation," arXiv preprint arXiv:2212.10481, 2023.

[45] Q. Peng, Y. Chai, and X. Li, "HumanEval-XL: a multilingual code generation benchmark for cross-lingual natural language generalization," arXiv preprint arXiv:2402.16694, 2024.