

# Multi-Task Pre-Training on Diverse Programming Languages Enhances Few-Shot Code Generation

Assignee Research

June 7, 2026

## Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does the multi-task pre-training on diverse programming languages impact the few-shot learning performance of LLMs on complex code generation tasks beyond HumanEval, such as MBPP or APPS, when. 11 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Multi-Programming Language Ensemble for Code Generation in Large Language Model. Research question: How does the multi-task pre-training on diverse programming languages impact the few-shot learning performance of LLMs on complex code generation tasks beyond HumanEval, such as MBPP or APPS, when compared to single-language training?.

## 2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

## 3 Results

12 papers retrieved. 11 claims extracted; 2 independently verified. Quality review score: 4.5/10.

## 4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

## 5 Extracted Claims

Claim	Verified	Confidence
HumanEval is designed for text-to-code (Python) generation tasks where the input is a brief passage describing the intention	×	0.04
HumanEval-plus extends the HumanEval dataset by incorporating a large number of additional valid unit test cases to rigorously	×	0.05
Pass@1 measures the percentage of tasks for which the model’s top output passes all hidden test cases.	×	0.02
We conducted experiments using both proprietary and open-source LLMs: GPT3.5-turbo, GPT-4o-mini, GPT-4o, Claude-Sonnet-3	×	0.02
The performance results of each method on the HumanEval and HumanEval-plus benchmarks are presented in Tables 1 and 2, respectively	×	0.11
The MPLE framework iteratively refines code by leveraging the strengths of different programming languages, reducing language-specific	✓	0.24
The MPLE framework integrates reflection algorithms and MCTS to enhance the overall robustness and accuracy of the generated code	×	0.08
The code generation task is formulated as a triplet (Q, Tv, Th), where Q is the code task description, Tv represents the test cases, and Th	×	0.05
The LLM is provided with Q and Tv to generate an initial program, P0, which is then refined iteratively to produce a sequence of programs	×	0.02
The final output program, P, is evaluated on the hidden test cases Th to verify its correctness.	×	0.01
The MPLE framework is designed to utilize the multi-language capabilities of LLMs to improve code generation.	✓	0.22

## References

- <http://arxiv.org/abs/2409.04114v1>
- <http://arxiv.org/abs/2003.04390v4>
- <http://arxiv.org/abs/2308.10783v2>