

# Mistral-Large-2 Code Generation Quality on MBPP: Human Evaluation vs. Ground Truth

Assignee Research

May 30, 2026

## Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How does the code generation quality of Mistral-Large-2 on MBPP benchmark compare to ground truth implementations when evaluated by human reviewers on functional correctness and code quality metrics. The creation of instruction data and evaluation benchmarks for serving Large language models often involves enormous human annotation. This issue becomes particularly pronounced when rapidly developing such resources for a non-English language like Japanese. 19 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Rapidly Developing High-quality Instruction Data and Evaluation Benchmark for Large Language Models with Minimal Human Effort: A Case Study on Japanese. Research question: How does the code generation quality of Mistral-Large-2 on MBPP benchmark compare to ground truth implementations when evaluated by human reviewers on functional correctness and code quality metrics?.

## 2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.7/10.

### **3 Results**

14 papers retrieved. 19 claims extracted; 1 independently verified. Quality review score: 5.7/10.

### **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
Japanese-Alpaca is fine-tuned with the instruction data translated from English Alpaca via GPT-3.5.	✓	0.20
The original English data is generated by outdated LLMs such as GPT-3.5.	×	0.07
The proposed method generates Japanese instruction with GPT-4 directly.	×	0.11
The original English manual seed instruction tasks are translated into Japanese and proofread by native Japanese speaker	×	0.07
The QA benchmark consists of 80 high-quality questions for evaluating Japanese LLMs.	×	0.12
The benchmark follows the reference-free evaluation manner, leveraging GPT-4 to assess the quality of the LLM answers.	×	0.07
The original questions data set is drawn from the English Vicuna benchmark and divided into 8 common question categories	×	0.06
The questions are designed to test instruction-following ability, covering common use cases and focusing on challenging	×	0.04
Each LLM being evaluated needs to answer all of those 80 Japanese questions.	×	0.05
GPT-4 performs pairwise comparisons of responses from different LLMs to ascertain which one performs better or yields co	×	0.02
GPT-4 directly assigns a score to an answer generated by LLMs, considering factors such as helpfulness, relevance, accur	×	0.04
LLaMA 7B MT Alpaca Self-instruct has a score of 2.36.	×	0.07
LLaMA2 7B MT Alpaca Self-instruct has a score of 5.71.	×	0.06
Open-calm 7B MT Alpaca Self-instruct has a score of 4.75.	×	0.05
LLaMA2 13B MT Alpaca Self-instruct has a score of 6.06.	×	0.07
LLaMA 7B MT Alpaca Self-instruct has a win-rate of 13.12.	×	0.09
LLaMA2 7B MT Alpaca Self-instruct has a win-rate of 46.25.	×	0.08
Open-calm 7B MT Alpaca Self-instruct has a win-rate of 34.37.	×	0.09
LLaMA2 13B MT Alpaca Self-instruct has a win-rate of 54.37.	×	0.14

## References

- <http://arxiv.org/abs/2403.03788v1>
- <http://arxiv.org/abs/2306.08568v2>
- <http://arxiv.org/abs/2403.03690v1>