

# Self-Instruct Tuning with GPT-4 for Japanese Language Models: Performance Gains over Human Benchmarks

Assignee Research

May 30, 2026

## Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: What is the impact of self-instruct methods based on GPT-4 on the performance of Japanese language models compared to traditional human-annotated benchmarks, as measured by BLEU or ROUGE scores. Despite vision-language models' (VLMs) remarkable capabilities as versatile visual assistants, two substantial challenges persist within the existing VLM frameworks: (1) lacking task diversity in pretraining and visual instruction tuning, and (2) annotation error and bias in. 17 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Vision-Flan: Scaling Human-Labeled Tasks in Visual Instruction Tuning. Research question: What is the impact of self-instruct methods based on GPT-4 on the performance of Japanese language models compared to traditional human-annotated benchmarks, as measured by BLEU or ROUGE scores?.

## 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 6.8/10.

### **3 Results**

12 papers retrieved. 17 claims extracted; 1 independently verified. Quality review score: 6.8/10.

### **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
Vision-Flan was evaluated on MMbench, MME, MMMU, MM-Vet, LLaVA-Bench, POPE, CIFAR-10, CIFAR-100, MNIST, and miniImageNet	×	0.02
For MMbench, MME, MM-Vet, LLaVA-Bench, POPE, and MMMU, the study strictly followed official implementations.	×	0.02
Vicuna 1.5 13B was used to evaluate performance on MNIST and miniImageNet datasets.	×	0.06
Baselines compared in the study include BLIP-2, InstructBLIP, Shikra, LLaVA, Qwen-VL, Qwen-VL-Chat, and LLaVA-1.5.	×	0.02
VISION-FLAN BASE achieves state-of-the-art performance on MME, MM-Bench, and MMMU benchmarks.	×	0.10
VISION-FLAN BASE scores significantly lower on the LLaVA-Bench dataset compared to VLMs trained using GPT-4 synthesized	×	0.14
VISION-FLAN CHAT was tuned on 1,000 GPT-4 synthesized data instances.	✓	0.19
Replacing instruction-tuned MLPs with pre-trained MLPs from the pre-trained LLaVA model allows VISION-FLAN models to retrain	×	0.06
The VISION-FLAN dataset contains 1.6 million instances.	×	0.07
The VISION-FLAN dataset contains 196 distinct tasks.	×	0.06
The VISION-FLAN dataset is publicly available.	×	0.11
The LLaVA dataset contains 150,000 instances covering 3 task categories.	×	0.04
The SVIT dataset contains 4.2 million instances covering 4 task categories.	×	0.03
The MultiInstruct dataset contains 510,000 instances covering 62 tasks.	×	0.03
The majority of existing visual instruction tuning datasets compared were generated using proprietary language models such as	×	0.14
The VL-Qwen dataset is annotated by humans but remains inaccessible to the public.	×	0.01
The MultiInstruct dataset mainly focuses on visual grounding tasks and contains 29 tasks that do not involve region-spec	×	0.03

## References

- <http://arxiv.org/abs/2304.03277v1>
- <http://arxiv.org/abs/2402.11690v1>
- <http://arxiv.org/abs/2403.03690v1>