

Instruction tuning large language models (LLMs) using machine-generated instruction-following data has improved zero-shot capabilities on new tasks, but the idea is less explored in the multimodal field.

We use language-only GPT-4 to generate language-image instruction-following data. By using such generated data, we introduce LLaVA:

Assistant, an end-to-end trained large multimodal model that connects a vision encoder and LLM for

general-purpose visual and language understanding. Our

early experiments show that LLaVA demonstrates impressive multimodal chat abilities, sometimes exhibiting

In this paper, we present the first attempt to

generate multimodal

following data. By

such generated data, we

Large Language and Vision

end-to-end trained large multimodal

connects a vision encoder and LLM for

visual and language understanding. Our