Harnessing PDF Data for Improving Japanese Large Multimodal Models

Jeonghun Baek[♠] Akiko Aizawa[♦] Kiyoharu Aizawa[♠]

◆The University of Tokyo ◇National Institute of Informatics baek@hal.t.u-tokyo.ac.jp

https://github.com/ku21fan/PDF-JLMM

Abstract

Large Multimodal Models (LMMs) have demonstrated strong performance in English, but their effectiveness in Japanese remains limited due to the lack of high-quality training data. Current Japanese LMMs often rely on translated English datasets, restricting their ability to capture Japan-specific cultural knowledge. To address this, we explore the potential of Japanese PDF data as a training resource, an area that remains largely underutilized. We introduce a fully automated pipeline that leverages pretrained models to extract image-text pairs from PDFs through layout analysis, OCR, and vision-language pairing, removing the need for manual annotation. Additionally, we construct instruction data from extracted imagetext pairs to enrich the training data. To evaluate the effectiveness of PDF-derived data, we train Japanese LMMs and assess their performance on the Japanese LMM Benchmark. Our results demonstrate substantial improvements, with performance gains ranging from 2.1% to 13.8% on Heron-Bench. Further analysis highlights the impact of PDF-derived data on various factors, such as model size and language models, reinforcing its value as a multimodal resource for Japanese LMMs.

1 Introduction

Large Multimodal Models (LMMs) have achieved high performance in English (OpenAI, 2024a; Team et al., 2024; Dubey et al., 2024; Yang et al., 2024), and their development is now expanding to other languages. Recently, several open-source Japanese LMMs have been released (Akiba et al., 2025; Tanahashi et al., 2023; Inoue et al., 2024; Lab, 2024; Sasagawa et al., 2025). While these models perform reasonably well, they still lag behind their English counterparts, partly due to the limited availability of Japanese training data.

Unlike English, where large-scale public imagetext pairs exist, Japanese models often rely on trans-



Figure 1: **Examples of PDF data.** We use various types of PDF data for Japanese LMM training.

lated English data (Shing and Akiba, 2023a,b; Tanahashi et al., 2023; Inoue et al., 2024). This means Japanese LMMs primarily learn content from western sources, lacking exposure to Japan-specific cultural knowledge.

To address this issue, we utilize Japanese PDF data to incorporate culturally relevant knowledge into LMM training. Unlike existing multimodal datasets, which are primarily web-based (Lin et al., 2014; Sharma et al., 2018; Schuhmann et al., 2022), PDFs contain a vast amount of valuable but underutilized information from books and documents. Despite this, research on integrating PDF data into LMM training remains limited. To our knowledge, no existing work has leveraged PDF data to enhance Japanese LMMs.

We investigate whether PDF data can effectively enhance Japanese LMMs. However, manually annotating large-scale PDFs is costly. To overcome this, we develop a fully automated pipeline, as

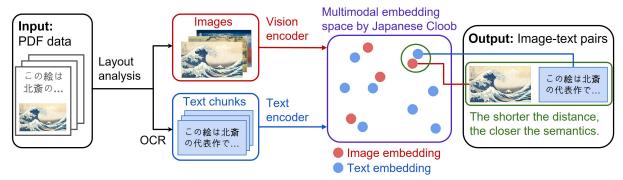


Figure 2: **An automated pipeline for extracting image-text pairs from PDFs.** It leverages pretrained models for layout analysis, OCR, and vision-language pairing.

shown in Figure 2. It extracts image-text pairs from PDFs using pretrained models for layout analysis, OCR, and vision-language pairing. Additionally, we generate instruction data from extracted imagetext pairs to enrich the training data.

We train Japanese LMMs with the PDF-derived data and evaluate them on the Japanese LMM Benchmark. Our results show that PDF-derived data significantly improves performance, with gains of 2.1% to 13.8% on Heron-Bench (Inoue et al., 2024). Furthermore, we conduct additional experiments to provide various insights into PDF-derived data's effectiveness. Our main contributions are as follows:

- We introduce a fully automated pipeline for extracting image-text pairs from PDFs using pretrained models, eliminating the need for manual annotation.
- We demonstrate that PDF-derived data can significantly improve Japanese LMM performance, achieving 2.1% to 13.8% gains on Heron-Bench.
- Through extensive experiments, we provide various insights into PDF-derived data's effectiveness. For example, we analyze the impact of PDF-derived data across different model sizes (3.8B, 8B, 14B) and evaluate the effectiveness of both image-text pairs and instruction data generated from images alone.

2 Related Work

2.1 Extracting Image-text Pairs from PDFs

Research on extracting images and their captions from PDFs, particularly scientific papers, has been actively explored (Clark and Divvala, 2015, 2016; Siegel et al., 2018; Naiman et al., 2022; Okamoto et al., 2023). These studies typically perform the layout analysis (Shen et al., 2021) to locate image

regions within a PDF, extract caption data from nearby text, and pair them together. When pairing, some approaches use distance-based matching, considering that caption text is generally closer to the corresponding image than other text (Okamoto et al., 2023).

However, to the best of our knowledge, no existing study has paired images with text other than the captions explicitly found in PDFs. We aim to extract image-text pairs from PDFs without being limited to captions. The closest existing work to our task is the identification of paragraphs that reference figures in scientific papers and summarizing their content to generate figure captions (Huang et al., 2023). However, this approach does not strictly pair images with non-caption text in PDFs, and its applicability is limited to scientific papers rather than general PDFs. In contrast, our work extends beyond scientific papers to cover a broader range of general PDFs.

2.2 Japanese LMM

Recently, Japanese large multimodal models (LMMs) have been emerging based on English LMMs. Proprietary LMMs have been improving their multilingual capabilities, achieving high performance in Japanese as well (OpenAI, 2023, 2024a; Anthropic, 2024; Team et al., 2024). Additionally, many open-source Japanese LMMs have been released (Shing and Akiba, 2023a,b; Akiba et al., 2025; Tanahashi et al., 2023; Inoue et al., 2024; Lab, 2024; Sasagawa et al., 2025).

Most open-source LMMs follow the LLaVA (Liu et al., 2023) approach, where a large language model (LLM) and a vision encoder are connected via a relatively shallow projector to form an LMM. For training, some use in-house Japanese data (Lab, 2024), while others rely on translated Japanese

data (Shing and Akiba, 2023a,b; Tanahashi et al., 2023; Inoue et al., 2024) and adopt a Japanese LLM as the base language model (Lab, 2024; Shing and Akiba, 2023a,b; Tanahashi et al., 2023; Inoue et al., 2024). This approach enables the development of Japanese LMMs with decent performance.

Some models, such as Owen-VL (Bai et al., 2023), achieve strong performance in Japanese without using a Japanese LLM, instead leveraging a multilingual LLM. VILA-jp (Sasagawa et al., 2025) has further improved performance by utilizing interleaved data. MangaLMM (Baek et al., 2025) specifically targets Japanese manga and demonstrates the potential of domain-specific Japanese LMMs. However, no existing work has leveraged PDF data to enhance Japanese LMMs. To achieve higher performance, we utilize PDF data in our approach. Recently, several benchmarks have been introduced for evaluating Japanese LMMs, such as JDocQA (Onami et al., 2024), which focuses on Japanese document QA, and JMMMU (Onohara et al., 2025), which considers Japanese culture.

3 Harnessing PDF Data

We aim to enhance the performance of Japanese LMMs using PDF data. In this section, we describe the process of obtaining Japanese LMM training data from PDF data.

3.1 Collecting PDF data

The PDF dataset used in this study was collected from the Web based on URLs supplied by the Web ARchiving Project of the National Diet Library of Japan (National Diet Library). The total number of PDFs exceeds 51.38 million. However, we do not use all of these PDFs; instead, we select a subset through the following process. The PDF data includes a wide variety of document types, not only academic or scientific papers but also newsletters, magazines, reports, posters, advertisements, campaign materials, pamphlets, brochures, manuals, and books. Figure 1 presents examples of various types of PDF data.

3.2 Extracting Image-text Pairs

To create training data for LMMs, we extract imagetext pairs from PDF data. The overall process is illustrated in Figure 2.

Selecting PDFs that Contain Images. Before extracting images and text from PDFs, we first filter out PDFs that do not contain images. A significant

portion of the PDF data does not include any images, and many files contain only small logos or symbols rather than meaningful images.

After manually inspecting hundreds of PDFs, we observed that as the number of pages increases, PDFs tend to resemble books, where text dominates and images are scarce. To address this, we select only PDFs with five or fewer pages. Additionally, we found that images frequently appear on the first page of PDFs. If an image is absent on the first page, subsequent pages are often image-free as well. Thus, we extract only the first page from each selected PDF.

To detect whether a PDF contains images, we use a Python library. There are several libraries available for this task (Artifex Software Inc.; Belval, 2017; pdfminer.six). Among them, we choose PyMuPDF (Artifex Software Inc.), which is widely used and offers both high speed and accuracy. Using PyMuPDF, we identify PDFs that contain image data and filter out those without images. As a result, we select 200K PDFs. Since we only use the first page of each PDF, this amounts to a total of 200K PDF pages.

Extracting Image and Text through Layout Analysis and OCR. The PyMuPDF library used for PDF selection directly reads a PDF and extracts images and text stored within the file. However, this approach sometimes leads to issues. For instance, PyMuPDF may extract images that are invisible to the human eye within a PDF. Additionally, in some cases, it breaks down visible images and extracts them separately based on layout elements. For example, in some cases, the background and objects within an image are extracted separately.

To prevent such issues, we first convert a PDF into a JPEG image and then extract images and text from it. This ensures that only images visible to the human eye are extracted. To extract images and text, we use Surya (Paruchuri, 2024), a tool designed for PDF analysis based on image inputs.

Surya performs both layout analysis (Shen et al., 2021) and OCR. First, layout analysis identifies image and text regions within the PDF. Then, OCR is applied to the text regions to extract the text. Through this process, we obtain both images and text

Surya employs pretrained deep learning models for layout analysis and OCR, supporting over 90 languages, making it applicable to Japanese PDFs. However, their performance is not perfect. For ex-



(a) Paired text 恩田川沿いの桜並木(3月31日撮影) (Cherry tree-lined path along the Onda River (Photographed on March 31))

(b) PDF-style text

桜の花が咲き誇る風景は、日本の春を象徴する光景であり、多くの人々に愛されています。この時期には、花見を楽しむために多くの観光客が訪れ、自然の美しさを堪能します。

(The sight of cherry blossoms in full bloom is a symbol of spring in Japan and is cherished by many. During this season, numerous tourists visit to enjoy ...)

(c) Instruction data

質問: この画像には何が写っていますか?(Question: What is depicted in this image?)

回答: ごの画像には、桜の木が川の両側に並んでいる風景が写っています。桜の花は満開で、白や薄いピンクの花びらが見えます。空は青く、晴れた日... (Answer: This image depicts a scene where cherry trees line both sides of a river. The cherry blossoms are in full bloom, displaying white and pale pink...)

Figure 3: Example of an image with (a) paired text, (b) PDF-style text, and (c) instruction data

ample, despite processing Japanese text, it occasionally misidentifies characters as Hindi. Additionally, we filter out images detected by Surya if their width or height is less than 50 pixels, as many non-image elements were mistakenly classified as images. Further analysis of Surya's pretrained model performance is provided in §5.5.

Pairing Image and Text. For each image, we match the most semantically similar text. Specifically, we embed the image using a vision encoder and OCR-extracted text using a text encoder. We then compute the cosine similarity between the image and the text and select the text with the highest similarity as the paired text for the image. For pairing, we use Japanese-Cloob (Shing et al., 2022; Sawada et al., 2024), a pretrained vision-language model widely used in Japan (with 300K users last month), which follows a CLIP (Radford et al., 2021)-like architecture. Figure 3(a) presents an example of the image-text pair.

Filtering NSFW and PII data. We found that PDF data sometimes contains (1) NSFW (Not Safe For Work) content and (2) PII (Personally Identifiable Information). To filter out such data, we use GPT-4o-mini (gpt-4o-mini-2024-07-18) (OpenAI, 2024b) to detect NSFW and PII content and exclude any data falling into these categories. As a result, our dataset is largely free from such content and is expected to be safe for use.

Generating PDF-style text. Pretrained models for layout analysis, OCR, and pairing are expected to perform well on scientific papers, as this domain has been extensively studied. However, in our experiments, where we used a diverse range of PDFs beyond scientific papers, these pretrained models did not perform as effectively. As a result, the quality of image-text pairs obtained using pretrained models was not always high. Further details on this can be found in §5.5.

This raises the question: "What if we could extract image-text pairs from PDFs with higher accuracy?" To explore this, we define PDF-style text and generate it using GPT-40-mini, simulating an ideal paired text where each image is associated with a semantically relevant description. Unlike direct image captions, PDF-style text does not explicitly describe the image in detail but instead provides a brief, indirect explanation derived from the surrounding text in the PDF. For more details, refer to the prompt used for PDF-style text generation (Table A).

As shown in Figure 3(b), the generated PDF-style text resembles the type of sentences commonly found in PDFs while providing an indirect description of the image. The impact of training with PDF-style text is presented in Table 6.

3.3 Generating Instruction Data

Image-text pairs can be directly used for LMM training; however, their effectiveness was limited (see Table 6). Thus, instead of using them as they are, we followed the LLaVA (Liu et al., 2023) approach and generated instruction data using GPT.

Strictly speaking, our method differs slightly from LLaVA. At the time LLaVA was developed, GPT-4 (OpenAI, 2023) could not process image inputs, which likely explains why image data was not provided to GPT-4 during instruction generation. However, by the time of our study, GPT-40 had been released, enabling image recognition.

Therefore, we directly feed images to GPT-40-mini to generate instruction data. The paired text associated with each image is used as context information during instruction generation. For the prompt, we slightly modify the final part of the LLaVA prompt, ensuring that the responses are generated in Japanese. The actual prompt used is provided in Table B.

Through our experiments, we found that when

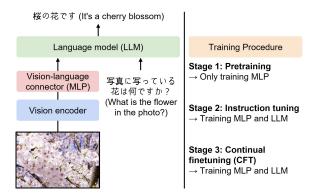


Figure 4: **Training LMM.** We use the LLaVA1.5 (Liu et al., 2024a) framework in our experiments. The training consists of 3 stages: pretraining, instruction tuning, and continual finetuning (CFT).

| Dataset | Stage | Count |
|-----------------------------|-------|-------|
| LLaVA-Pretrain-JA | 1 | 558K |
| LLaVA-v1.5-Instruct-620K-JA | 2 | 620K |
| Instruct-from-200K PDF | 3 | 362K |

Table 1: Details of the Japanese LMM training sets.

the paired text matched with an image is imperfect, it is more effective to generate instruction data using only the image, rather than including the paired text. Further details on this can be found in §5.6. Therefore, in our experiments, all instruction data, except for that used in Table 7, is generated using images only. From 200K PDFs, we generate a total of 362K instruction tuning data. Figure 3(c) presents an example of the instruction data.

4 Training Japanese LMM

Recently, various open-source Japanese LMMs have been released (Shing and Akiba, 2023a,b; Akiba et al., 2025; Tanahashi et al., 2023; Inoue et al., 2024; Lab, 2024; Sasagawa et al., 2025). Among them, we adopt the widely used LLaVA (Liu et al., 2023) framework to evaluate the effectiveness of PDF data, specifically using LLaVA1.5 (Liu et al., 2024a). Figure 4 illustrates the LLaVA1.5 framework. Most hyperparameters follow the original LLaVA1.5 settings, with a few modifications. We replace the vision encoder CLIP (clip-vit-largepatch14-336) (Radford et al., 2021) with SigLIP (siglip-so400m-patch14-384) (Zhai et al., 2023) and experiment with different large language models (LLMs) instead of Vicuna-7B (Chiang et al., 2023). Details on the LLM selection will be discussed later. For training, we employ LoRA (Hu et al., 2022) for parameter-efficient finetuning.

| Dataset | #Questions | #Images |
|-----------------------|------------|---------|
| JA-LLaVA-Bench (COCO) | 90 | 30 |
| JA-LLaVA-Bench (Wild) | 60 | 24 |
| Heron-Bench | 103 | 21 |

Table 2: Details of Japanese LMM evaluation sets.

4.1 Training Procedure

LLaVA1.5 training consists of two stages, aiming to integrate the pretrained vision encoder and pretrained LLM to create an LMM that effectively handles vision input.

Stage 1: Pretraining. In this stage, only the vision-language connector (MLP) is trained using imagetext pairs, linking the vision encoder with the LLM.

Stage 2: Instruction Tuning. Using visual instruction data, both the MLP and LLM are instructiontuned to improve multimodal understanding.

Stage 3: Continual Fine-Tuning (CFT). In our study, we introduce an additional CFT stage after Stages 1 and 2. Here, we perform CFT on both the MLP and LLM using PDF data.

Training Data. Table 1 presents the details of the training data. In Stage 1 (Pretraining) and Stage 2 (Instruction Tuning), we use the Japanese-translated version of the original LLaVA training data (Inoue et al., 2024), translated using DeepL (DeepL, 2017). Specifically, Stage 1 uses 558K samples from LLaVA-Pretrain-JA (Motors, 2024a), Stage 2 uses 620K samples from LLaVA-v1.5-Instruct-620K-JA (Motors, 2024b), and Stage 3 uses 362K PDF-derived samples created in §3.

Elapsed Time for Training. LLaVA1.5 is trained for one epoch per stage. Training LLaVA1.5 with a Llama3 8B-based LLM using four NVIDIA A100 GPUs took about 11 hours for Stage 1, 42 hours for Stage 2, and 19 hours for Stage 3.

4.2 LLM Selection

Japanese LLM. We train LLaVA1.5 using three well-known Llama3-8B-based Japanese LLMs: Suzume (suzume-llama-3-8B-japanese (Devine, 2024)), ELYZA (Llama-3-ELYZA-JP-8B (Hirakawa et al., 2024)), and Swallow (Llama-3-Swallow-8B-Instruct-v0.1 (Okazaki et al., 2024)) We build our main model by finetuning LLaVA1.5, using Swallow as the base Japanese LLM, with Japanese PDF data. We name this model PDF-JLMM (PDF-based Japanese Large Multimodal Model).

| | | JA-Ll | LaVA-l | Bench (CO | CO) | JA-L | LaVA | Bench (W | ild) | | Heron | 1-Bench | |
|---------|--------------------------|--------|--------|-----------|------|--------|------|----------|------|--------|-------|---------|------|
| | Method | Detail | Conv | Complex | Avg. | Detail | Conv | Complex | Avg. | Detail | Conv | Complex | Avg. |
| -U | GPT-4V | 88.0 | 84.7 | 97.5 | 90.1 | 89.9 | 93.1 | 99.1 | 94.1 | 83.3 | 77.5 | 78.3 | 79.7 |
| Bench | Claude 3 Opus | 67.9 | 70.4 | 95.3 | 77.9 | 87.4 | 73.4 | 94.6 | 85.1 | 74.5 | 68.4 | 77.7 | 73.6 |
| -Be | Gemini Pro | 65.3 | 83.0 | 75.4 | 74.6 | 61.7 | 84.2 | 84.2 | 76.7 | 55.6 | 64.3 | 64.0 | 61.3 |
| ron | LLaVA 1.6 7B | 60.3 | 83.7 | 60.6 | 68.2 | 36.7 | 44.7 | 53.4 | 44.9 | 30.9 | 37.3 | 31.0 | 33.1 |
| He | LLaVA 1.5 7B | 71.1 | 83.7 | 69.6 | 74.8 | 49.2 | 48.7 | 54.7 | 50.8 | 42.4 | 45.9 | 35.5 | 41.3 |
| from | Qwen-VL 7B | 78.0 | 81.0 | 82.2 | 80.4 | 55.9 | 49.7 | 56.4 | 54.0 | 46.3 | 50.6 | 52.3 | 49.7 |
| | Japanese StableVLM 7B | 18.9 | 54.8 | 24.1 | 32.6 | 26.0 | 24.8 | 29.2 | 26.7 | 25.2 | 51.2 | 37.8 | 38.1 |
| esults | EvoVLM-JP-v1 7B | 61.0 | 75.7 | 71.0 | 69.2 | 49.6 | 65.5 | 54.2 | 56.4 | 50.3 | 44.4 | 40.5 | 45.1 |
| esı | Heron BLIP v1 7B | 84.8 | 94.3 | 89.5 | 89.5 | 45.5 | 32.9 | 56.9 | 45.1 | 49.1 | 41.5 | 45.7 | 45.4 |
| Ä | Heron GIT 7B | 83.0 | 78.2 | 91.1 | 84.1 | 41.0 | 39.9 | 54.6 | 45.2 | 42.8 | 54.2 | 43.5 | 46.8 |
| lts | PDF-JLMM 8B | 81.6 | 95.9 | 87.1 | 88.2 | 69.3 | 50.6 | 77.6 | 65.8 | 70.1 | 62.3 | 65.0 | 65.8 |
| results | LLaVA1.5-Llama3 8B | 82.8 | 90.6 | 87.2 | 86.9 | 60.5 | 38.7 | 71.5 | 56.9 | 68.7 | 60.1 | 56.1 | 61.6 |
| II I | LLaVA1.5-Phi3-mini 3.8B | 77.3 | 89.8 | 83.7 | 83.6 | 66.7 | 36.5 | 66.8 | 56.7 | 61.6 | 61.2 | 48.5 | 57.1 |
| Our | LLaVA1.5-Phi3-medium 14B | 84.7 | 89.6 | 86.2 | 86.8 | 77.0 | 68.2 | 77.0 | 74.1 | 62.0 | 56.1 | 54.1 | 57.4 |

Table 3: **Main results.** Our models trained with PDF data achieve high performance across all benchmarks except JA-LLaVA-Bench (COCO). Our LLaVA1.5-based models are named after their backbone LLMs (e.g., LLaVA1.5-Llama3). PDF-JLMM refers to LLaVA1.5 with Swallow as the backbone LLM. The best-performing open-source models are highlighted in **bold**.

General (Non-Japanese) LLM. To verify whether Japanese PDF data is effective in adapting a general (non-Japanese) LLM into a Japanese LMM, we use three non-Japanese LLMs: Llama3 (Llama-3-8B-instruct (Dubey et al., 2024)), Phi3-mini (Phi-3-mini-4k-instruct (Abdin et al., 2024), 3.8B parameters), and Phi3-medium (Phi-3-medium-4k-instruct (Abdin et al., 2024), 14B parameters).

5 Experiments and Analysis

5.1 Evaluation Metric

To evaluate Japanese LMMs, we adopt the evaluation method used in Heron-Bench (Inoue et al., 2024), a standard Japanese LMM benchmark. The authors of Heron-Bench provide three evaluation datasets: JA-LLaVA-Bench (COCO) and JA-LLaVA-Bench (Wild), which are Japanese translations of LLaVA-Bench(Liu et al., 2023), and Heron-Bench, specifically designed for Japanese evaluation. The details of them are presented in Table 2.

Heron-Bench follows the same score calculation method as LLaVA-Bench. First, GPT-4 (gpt-4-0125-preview) (OpenAI, 2023) generates reference answers using the question's context. Then, GPT-4 evaluates both LMM's answers and the reference answers using the LLM-as-a-judge approach (Zheng et al., 2024). The final score is calculated as the ratio (%) of the average score of the LMM's answers to the average score of GPT-4's reference answers. A score of 100% indicates performance on par with GPT-4, while scores above 100% suggest that the LMM outperforms GPT-4.

LMMs used for comparison. We use three proprietary LMMs—GPT-4V (OpenAI, 2023), Claude 3 Opus (Anthropic, 2024), and Gemini Pro (Team et al., 2024)—along with seven open-source LMMs—LLaVA 1.6 7B (Liu et al., 2024b), LLaVA 1.5 7B (Liu et al., 2023), Qwen-VL 7B (Bai et al., 2023), Japanese StableVLM 7B (Shing and Akiba, 2023b), EvoVLM-JP-v1-7B (Akiba et al., 2025), Heron BLIP v1 (620k) (Tanahashi et al., 2023), and Heron GIT (Inoue et al., 2024).

5.2 Main Result

Table 3 presents the results of our model trained on PDF-derived data. Compared to existing models from the Heron-Bench paper, our PDF-JLMM outperforms most open-source Japanese LMMs. It lags behind Heron BLIP v1 7B by only 1.3% on JA-LLaVA-Bench (COCO). On JA-LLaVA-Bench (Wild), it achieves a performance of 65.8%, surpassing the previous best of 56.4% by 9.4%. For Heron-Bench, it outperforms the previous best of 49.7% by a significant margin of 16.1%. Additionally, our LLaVA1.5-Llama3, LLaVA1.5-Phi3-mini, and LLaVA1.5-Phi3-medium also achieve higher performance than existing models. These results demonstrate that *utilizing PDF data effectively enhances model performance*.

In the following subsections, we conduct additional experiments on training with PDF-derived data and provide a component-wise analysis.

| LLM | Stage | L-COCO | L-Wild | Heron |
|-----------------|----------------------------------|--------|-------------|-------|
| | 1. Pretraining | 30.6 | 17.4 | 22.6 |
| 8B | 2. Instruction tuning | 84.0 | 59.8 | 54.7 |
| ≥ | 3. CFT on 50K PDF | 87.3 | 61.6 | 65.7 |
| Swallow | 3. CFT on 100K PDF | 88.2 | 65.8 | 65.8 |
| Sw | 3. CFT on 150K PDF | 88.1 | 65.5 | 63.8 |
| | 3. CFT on 200K PDF | 86.6 | 64.7 | 64.6 |
| | 1. Pretraining | 26.8 | 16.5 | 23.2 |
| 8B | 2. Instruction tuning | 84.4 | 57.0 | 54.8 |
| | 3. CFT on 50K PDF | 86.5 | 54.0 | 58.7 |
| Llama3 | 3. CFT on 100K PDF | 86.3 | 57.1 | 61.0 |
| Ϊ | 3. CFT on 150K PDF | 85.1 | 53.9 | 61.8 |
| | 3. CFT on 200K PDF | 86.9 | 56.9 | 61.6 |
| | 1. Pretraining | 21.8 | 13.5 | 18.4 |
| Phi3-mini 3.8B | 2. Instruction tuning | 82.7 | 50.6 | 43.3 |
| . <u>E</u> | 3. CFT on 50K PDF | 83.0 | 52.3 | 51.9 |
| . <u>Ē</u> | 3. CFT on 100K PDF | 83.3 | 52.0 | 53.0 |
| hi3 | 3. CFT on 150K PDF | 83.6 | 56.7 | 57.1 |
| | 3. CFT on 200K PDF | 83.8 | 48.2 | 54.3 |
| 4B | 1. Pretraining | 27.5 | 16.0 | 23.4 |
| n 1, | 2. Instruction tuning | 86.3 | 71.4 | 54.2 |
| li. | CFT on 50K PDF | 87.4 | 65.2 | 58.8 |
| Phi3-medium 14B | 3. CFT on 100K PDF | 85.9 | 66.8 | 56.3 |
| i3-r | 3. CFT on 150K PDF | 86.8 | 74.1 | 57.4 |
| Phi | 3. CFT on 200K PDF | 88.5 | 70.7 | 58.1 |

Table 4: Results for each stage and increasing amounts of PDF. CFT on PDF-derived data is effective. L-COCO, L-Wild, and Heron correspond to JA-LLaVA-Bench (COCO), JA-LLaVA-Bench (Wild), and Heron-Bench, respectively. The average value per benchmark is shown. For each LLM, the highest-performing stage and value are highlighted in **bold**.

5.3 Is Using PDF-derived Data Effective?

Table 4 presents the results of LLaVA1.5 training for each LLM across different stages. Overall, continual fine-tuning (CFT) on PDF-derived data is effective. Immediately after Stage 1 (Pretraining), all LLMs exhibit relatively low performance. After Stage 2 (Instruction tuning), Heron-Bench performance improves by approximately 20%–30%. After Stage 3 (CFT on PDF data), Heron-Bench performance further increases by at least 2.1% (Phi3-medium) and up to 13.8% (Phi3-mini).

For Stage 3, we also show experiments with increasing amounts of PDF data, where every additional 50K PDFs adds approximately 90K new instruction data. Figure 5 demonstrates the performance improvement on Heron-Bench when applying CFT on PDF data, highlighting the effectiveness of PDF data. However, performance does not always increase as the amount of PDF data grows. Performance does not improve in direct proportion to the amount of data, indicating that scaling is not

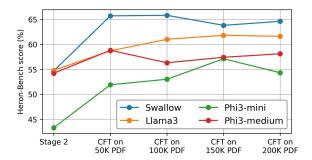


Figure 5: **Performance on Heron-Bench with CFT on PDF.** CFT improves performance. However, more data does not always lead to continuous performance gains.

| Method | L-COCO | L-Wild | Heron |
|---------------------|--------|--------|-------|
| LLaVA1.5-Suzume 8B | 83.1 | 55.4 | 51.9 |
| LLaVA1.5-ELYZA 8B | 81.8 | 55.1 | 53.5 |
| LLaVA1.5-Swallow 8B | 84.0 | 59.8 | 54.7 |

Table 5: **LLaVA1.5** results for each Japanese LLM. The results are from training up to Stage 1 and 2.

strictly linear.

Additionally, Japanese PDF data is also effective in adapting a general (non-Japanese) LLM into a Japanese LMM. To verify this, we conducted the same experiment using Llama3 and Phi3, both non-Japanese LLMs, as described in §4.2. As shown in Table 3 and Figure 5, PDF data is effective for training these models as Japanese LMMs.

Moreover, PDF data is effective across various model sizes. It benefits not only models based on Llama3 8B but also Phi3-mini (3.8B) and Phi3-medium (14B). However, in the case of Phi3, even at a larger model size of 14B, its performance lags behind the 8B Llama3-based model. This may be due to its limited Japanese vocabulary. Phi3 contains only 837 Japanese vocabulary tokens. Considering that commonly used kanji in Japanese amount to around 2,000 (Wikipedia contributors), 837 vocabulary tokens are likely insufficient for effectively handling Japanese. As a result, despite having a larger model size of 14B, its performance remains comparable to that of the 8B Llama3 model.

5.4 Which Japanese LLM Performs Best?

We train LLaVA1.5 up to Stage 1 and 2 using the three commonly used Japanese LLMs introduced in §4.2 (Suzume, ELYZA, and Swallow). Table 5 presents the results. We find that Swallow achieves the best performance across all three benchmarks, and therefore, we select Swallow as the Japanese LLM for this study.

| Method | L-COCO | L-Wild | Heron |
|----------------|--------|--------|-------|
| Stages 1 and 2 | 84.0 | 59.8 | 54.7 |
| Top 1 | 77.0 | 37.4 | 40.0 |
| Top 3 | 72.3 | 31.3 | 34.7 |
| Top 5 | 62.6 | 26.6 | 22.8 |
| Neighbor | 81.7 | 39.9 | 46.4 |
| PDF-style text | 81.5 | 56.5 | 65.5 |

Table 6: **Results from training with raw image-text** pairs. When performing CFT using only image-text pairs, the overall performance is lower than the baseline.

5.5 What If We Use Raw Image-Text Pairs Without Generating Instruction Data?

Instead of generating instruction data from extracted image-text pairs, an LMM can be trained directly on them. We train PDF-JLMM with 50K PDF data using extracted image-text pairs. Table 6 presents the results.

The results show that performance generally decreases compared to LLaVA1.5 trained only up to Stages 1 and 2. In the table, Top 1, 3, 5 refer to results obtained using the top 1, 3, or 5 texts ranked by cosine similarity when pairing with an image. Neighbor denotes the setting where, in addition to the Top 1 text, one preceding and one following text from the same PDF are also included. Even when using Top 3, Top 5, or Neighbor, performance is still lower than that of LLaVA1.5 trained only up to Stage 2. This suggests that image-text pairs extracted solely using pretrained models are not effective as training data.

There are two major possible reasons for this performance drop: (1) Many PDFs inherently contain little text that directly describes the images. (2) The limitations of pretrained models. The quality of image-text pairs may have been low due to insufficient performance in layout analysis, OCR, and pairing tasks. Since there is no annotation data for these PDFs, it is difficult to rigorously determine which factor is the primary cause. However, upon manually inspecting hundreds of examples, we observed that the OCR-extracted text was often inaccurate. Many extracted texts contained unintended line breaks, leading to broken words or sentences, and complex kanji characters were frequently misrecognized.

These results led us to question: "What if imagetext pairs were accurately extracted from PDFs?" To explore this, we generated an ideal paired text for each image using GPT-40-mini. This approach,

| Data source for instruction | L-COCO | L-Wild | Heron |
|-----------------------------|--------|--------|-------|
| Image | 87.3 | 61.6 | 65.7 |
| Image and paired text | 87.5 | 60.0 | 63.9 |
| Image and PDF-style text | 87.2 | 63.0 | 64.0 |

Table 7: Comparison of data sources for generating instruction data. PDF-JLMM is used in experiments.

referred to as PDF-style text, is described in §3.2.

Using PDF-style text significantly outperforms raw image-text pairs. Compared to Stages 1 and 2, it achieves a 10.8% improvement on Heron-Bench. These findings suggest that if text data is properly extracted from PDFs, training solely on image-text pairs could lead to further performance gains.

5.6 Is Paired Text Effective for Generating Instruction Data?

When generating instruction data, we use paired text matched with images as context. However, the effectiveness of this paired text remains uncertain. Since the text extracted from PDFs is often imperfect and contains noise, its usefulness for training may be uncertain. To investigate this, we generate instruction data using different data sources derived from 50K PDFs and present the results in Table 7. We compare three cases: using only images, images with paired text, and images with PDF-style text.

The key takeaway from our experiment is that instruction data generated using only images performed best. This suggests that even extracting only image data from PDFs is valuable, and when paired text is inaccurate, using images alone can yield better performance. Thus, for experiments using 200K PDFs, we generated instruction data using only images.

From these results, we conclude that paired text containing noise is not effective. However, the results also show that PDF-style text improves performance compared to paired text. This implies that if text data is more accurately extracted from PDFs, the quality of instruction data generated from image-text pairs can also improve.

5.7 Comparison with the Translated Datasets

We compare the PDF-derived data with Japanese-translated versions of existing English instruction data. The 620K instruction data (LLaVA-v1.5-Instruct-620K-JA) used in Stage 2 is considered a different Japanese multimodal dataset, constructed by translating English instruction data. We compare it with our Instruct-from-200K PDF dataset (362K

| Data source for Stage 2 | L-COCO | L-Wild | Heron |
|-----------------------------|--------|--------|-------|
| LLaVA-v1.5-Instruct-620K-JA | 84.0 | 59.8 | 54.7 |
| Instruct-from-200K PDF | 88.1 | 72.7 | 70.0 |

Table 8: Comparison with the translated data, LLaVA-v1.5-Instruct-620K-JA.

instruction data). Instead of using LLaVA-v1.5-Instruct-620K-JA in Stage 2, we train the model with our Instruct-from-200K PDF dataset. Table 8 shows the results. LLaVA1.5-Swallow 8B trained with the PDF-derived data outperforms the one trained with LLaVA-v1.5-Instruct-620K-JA, even with fewer samples. This result suggests that the performance gains are not merely due to large-scale instruction tuning, but rather from leveraging instruction data constructed from Japanese PDFs containing culturally grounded content.

We also compare the PDF-derived data with other datasets that we translated into Japanese. Specifically, we translated Vision-Flan (Xu et al., 2024) (186K instruction data) and Image-Textualization (Pi et al., 2024) (99.6K instruction data)—subsets of LLaVA OneVision (Li et al., 2024)—into Japanese using GPT-40-mini. For translation, we use the prompt in Table C. We use these datasets for CFT in Stage 3. Table 9 shows the results. The results show that simply translating effective English data into Japanese may even lead to performance degradation on Japanese benchmarks. This underscores the importance of using Japanese PDF data that contains authentic, culturally grounded content.

Our approach differs from translated English datasets in that we extract images containing authentic Japanese content directly from Japanese PDFs. For example, as shown in Figure A of the qualitative analysis in the supplementary materials, the model trained only up to Stage 2 (before incorporating our PDF-derived data) describes "cherry blossoms"—a culturally symbolic flower in Japan—as "white flowers", whereas the model further trained on PDF-derived data correctly refers to it as "cherry blossoms". We believe this illustrates a key distinction from translated data, as our method helps the model learn cultural concepts grounded in native Japanese contexts rather than relying on translations.

6 Conclusion

We explore the use of Japanese PDF data to enhance LMM training and develop a fully automated

| Data source for Stage 3 | L-COCO | L-Wild | Heron |
|---------------------------|--------|--------|-------|
| 100K PDF (181K inst.) | 88.2 | 65.8 | 65.8 |
| Vision-Flan (VF) | 37.8 | 29.0 | 40.2 |
| Image-Textualization (IT) | 81.6 | 43.3 | 46.3 |
| VF + IT | 56.7 | 35.5 | 39.9 |

Table 9: Comparison with the translated subset of LLaVA-OneVision. LLaVA1.5-Swallow is used in experiments.

pipeline for extracting image-text pairs. Our experiments show significant performance gains by incorporating PDF-derived data, with up to 13.8% improvement on Heron-Bench. Further analysis confirms the effectiveness of PDF-derived data across different model sizes and its potential to complement existing multimodal datasets. These findings provide valuable insights into leveraging PDF data for LMM training and highlight its promise as a multimodal resource. While our focus is on Japanese, we believe this approach is applicable to other languages as well, and we hope it fosters further research toward improving LMMs across diverse languages.

Limitations

Following LLaVA's approach of generating instruction data using GPT, we used GPT-40-mini to generate instruction data from PDF data. While this approach is effective, it is dependent on GPT. To generate LMM training data without relying on GPT, high-quality image-text pair data is essential. Achieving this requires improving the performance of text extraction models for PDFs.

Currently, scaling beyond 100K PDFs has been challenging. For future work, we plan to investigate the underlying causes—whether the bottleneck lies in data quality, model capacity, suboptimal training settings (e.g., learning rate, number of epochs), or the limitations of existing test data. By addressing these factors, we believe that scaling to larger PDF datasets will become more feasible.

Acknowledgments

We acknowledge the support of the Research and Development Center for Large Language Models at the National Institute of Informatics for providing the PDF dataset. In particular, we would like to thank Professor Daisuke Kawahara for his guidance regarding copyright and permissions. This work was supported by JSPS KAKENHI Grant Number

24K23882 and by the NVIDIA Academic Grant Program.

References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. arXiv:2404.14219.
- Takuya Akiba, Makoto Shing, Yujin Tang, Qi Sun, and David Ha. 2025. Evolutionary optimization of model merging recipes. *Nature Machine Intelligence*.
- Anthropic. 2024. Claude 3 family. Accessed: 2025-02-15.
- Artifex Software Inc. PyMuPDF (pymupdf). https://github.com/pymupdf/PyMuPDF. Accessed: February 15, 2025.
- Jeonghun Baek, Kazuki Egashira, Shota Onohara, Atsuyuki Miyai, Yuki Imajuku, Hikaru Ikuta, and Kiyoharu Aizawa. 2025. Mangavqa and mangalmm: A benchmark and specialized model for multimodal manga understanding. *arXiv*:2505.20298.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- Belval. 2017. pdf2image: Convert pdf to image. https://github.com/Belval/pdf2image. Accessed: February 15, 2025.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. Accessed: 2025-02-15.
- Christopher Clark and Santosh Divvala. 2016. Pdffigures 2.0: Mining figures from research papers. In *Proceedings of the 16th ACM/IEEE-CS on Joint Conference on Digital Libraries*.
- Christopher Andreas Clark and Santosh Divvala. 2015. Looking beyond text: Extracting figures, tables and captions from computer science papers. In *Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*.
- DeepL. 2017. DeepL Translator. https://www.deepl.com/en/translator. Accessed: February 15, 2025.
- Peter Devine. 2024. Tagengo: A multilingual chat dataset. *Proceedings of the Fourth Workshop on Multilingual Representation Learning*.

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. arXiv:2407.21783.
- Masato Hirakawa, Shintaro Horie, Tomoaki Nakamura, Daisuke Oba, Sam Passaglia, and Akira Sasaki. 2024. elyza/llama-3-elyza-jp-8b. Accessed: 2025-02-14.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *ICLR*.
- Chieh-Yang Huang, Ting-Yao Hsu, Ryan Rossi, Ani Nenkova, Sungchul Kim, Gromit Yeuk-Yin Chan, Eunyee Koh, Clyde Lee Giles, and Ting-Hao'Kenneth' Huang. 2023. Summaries as captions: Generating figure captions for scientific documents with automated text summarization. In *International Natural Language Generation Conference*.
- Yuichi Inoue, Kento Sasaki, Yuma Ochi, Kazuki Fujii, Kotaro Tanahashi, and Yu Yamaguchi. 2024. Heronbench: A benchmark for evaluating vision language models in japanese. In CVPR The 3rd Workshop on Computer Vision in the Wild.
- CyberAgent AI Lab. 2024. Llava-calm2-siglip. https://huggingface.co/cyberagent/llava-calm2-siglip. Accessed: 2025-02-15.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. Llavaonevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In ECCV.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In *CVPR*.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In *NeurIPS*.
- Turing Motors. 2024a. Llava-pretrain-ja. https://huggingface.co/datasets/turing-motors/LLaVA-Pretrain-JA. Accessed: February 15, 2025.
- Turing Motors. 2024b. Llava-v1.5-instruct-620k-ja. https://huggingface.co/datasets/turing-motors/LLaVA-v1.5-Instruct-620K-JA. Accessed: February 15, 2025.

- Jill P Naiman, Peter KG Williams, and Alyssa Goodman. 2022. Figure and figure caption extraction for mixed raster and vector pdfs: digitization of astronomical literature with ocr features. In *International Conference on Theory and Practice of Digital Libraries*.
- National Diet Library. Web Archiving Project (WARP). Accessed: 2025-02-18.
- Yamato Okamoto, Haruto Toyonaga, Yoshihisa Ijiri, and Hirokatsu Kataoka. 2023. Constructing imagetext pair dataset from books. In *ICCV Workshop on Towards the Next Generation of Computer Vision Datasets*.
- Naoaki Okazaki, Kakeru Hattori, Hirai Shota, Hiroki Iida, Masanari Ohi, Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Rio Yokota, and Sakae Mizuki. 2024. Building a large japanese web corpus for large language models. In *Proceedings of the First Conference on Language Modeling*.
- Eri Onami, Shuhei Kurita, Taiki Miyanishi, and Taro Watanabe. 2024. JDocQA: Japanese document question answering dataset for generative language models. In *LREC-COLING*.
- Shota Onohara, Atsuyuki Miyai, Yuki Imajuku, Kazuki Egashira, Jeonghun Baek, Xiang Yue, Graham Neubig, and Kiyoharu Aizawa. 2025. Jmmmu: A japanese massive multi-discipline multimodal understanding benchmark for culture-aware evaluation. In *NAACL*.
- OpenAI. 2023. GPT-4 Technical Report. Accessed: 2025-02-15.
- OpenAI. 2024a. Gpt-4o. Accessed: 2025-02-15.
- OpenAI. 2024b. Gpt-40 mini: Advancing cost-efficient intelligence.
- Vik Paruchuri. 2024. Surya: A pdf processing library. https://github.com/VikParuchuri/surya. Accessed: February 15, 2025.
- pdfminer.six. pdfminer.six: Extract text from pdf documents. https://github.com/pdfminer/pdfminer.six. Accessed: February 16, 2025.
- Renjie Pi, Jianshu Zhang, Jipeng Zhang, Rui Pan, Zhekai Chen, and Tong Zhang. 2024. Image textualization: An automatic framework for creating accurate and detailed image descriptions. In *NeurIPS*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*.
- Keito Sasagawa, Koki Maeda, Issa Sugiura, Shuhei Kurita, Naoaki Okazaki, and Daisuke Kawahara. 2025. Constructing multimodal datasets from scratch for rapid development of a japanese visual language model. In *NAACL demo*.

- Kei Sawada, Tianyu Zhao, Makoto Shing, Kentaro Mitsui, Akio Kaga, Yukiya Hono, Toshiaki Wakatsuki, and Koh Mitsuda. 2024. Release of pre-trained models for the Japanese language. In *LREC-COLING* 2024
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation imagetext models. *NeurIPS*.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*.
- Zejiang Shen, Ruochen Zhang, Melissa Dell, Benjamin Charles Germain Lee, Jacob Carlson, and Weining Li. 2021. Layoutparser: A unified toolkit for deep learning based document image analysis. In *ICDAR*.
- Makoto Shing and Takuya Akiba. 2023a. Japanese instructblip alpha.
- Makoto Shing and Takuya Akiba. 2023b. Japanese stable vlm.
- Makoto Shing, Tianyu Zhao, and Kei Sawada. 2022. rinna/japanese-cloob-vit-b-16. Accessed: February 15, 2025.
- Noah Siegel, Nicholas Lourie, Russell Power, and Waleed Ammar. 2018. Extracting scientific figures with distantly supervised neural networks. In *Proceedings of the 18th ACM/IEEE on joint conference on digital libraries*.
- Kotaro Tanahashi, Yuichi Inoue, and Yu Yamaguchi. 2023. Heron blip japanese stablelm base 7b.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv* preprint arXiv:2403.05530.
- Wikipedia contributors. Jōyō kanji. https://en.wikipedia.org/wiki/J%C5%8Dy%C5%8D_kanji. [Online; accessed February 15, 2025].
- Zhiyang Xu, Chao Feng, Rulin Shao, Trevor Ashby, Ying Shen, Di Jin, Yu Cheng, Qifan Wang, and Lifu Huang. 2024. Vision-flan: Scaling human-labeled tasks in visual instruction tuning. In *ACL findings*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv:2412.15115*.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In *ICCV*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. In *NeurIPS*.

A Prompts Used in Our Experiments

This section provides the prompts used in our experiments. Table A, Table B, and Table C show the prompts used to generate PDF-style text, instruction data, and Japanese translations of English datasets, respectively.

B Qualitative Analysis

Figures A, B, and C present qualitative analyses. Each figure includes an image, a question, the reference answer from GPT-4, the response from LLaVA1.5-Swallow trained up to stages 1 and 2, and the response from LLaVA1.5-Swallow further trained with stage 3 (CFT on PDF). The results show that performance improves after training up to stage 3, demonstrating the effectiveness of CFT using PDF-derived data.

You are an AI visual assistant, and you are seeing a single image. Generate a passage that resembles text commonly found in PDF documents and is relevant to the given image. The provided image is extracted from a PDF, but no additional context, such as the document's text or structure, is available.

PDF-style text generally has the following characteristics:

- 1. No explicit captions or minimal captions: Instead of directly describing the image, related text may naturally integrate into the document's content.
- 2. Indirect descriptions: The text does not explicitly reference the image but provides supporting information that the image complements.

To keep the text concise, generate only 1 to 2 sentences per image, ensuring it aligns with common PDF writing styles.

You must respond in Japanese.

Table A: **Prompt for generating PDF-style text.** An image is provided to GPT-4o-mini along with this prompt.

You are an AI visual assistant, and you are seeing a single image. What you see are provided within several sentences, describing the same image you are looking at. Answer all questions as you are seeing the image.

Design a conversation between you and a person asking about this photo. The answers should be in a tone that a visual AI assistant is seeing the image and answering the question. Ask diverse questions and give corresponding answers.

Include questions asking about the visual content of the image, including the object types, counting the objects, object actions, object locations, relative positions between objects, etc. Only include questions that have definite answers:

- (1) one can see the content in the image that the question asks about and can answer confidently;
- (2) one can determine confidently from the image that it is not in the image.

Do not ask any question that cannot be answered confidently.

Also include complex questions that are relevant to the content in the image, for example, asking about background knowledge of the objects in the image, asking to discuss about events happening in the image, etc. Again, do not ask about uncertain details.

Provide detailed answers when answering complex questions. For example, give detailed examples or reasoning steps to make the content more convincing and well-organized. You can include multiple paragraphs if necessary.

You must use Japanese all the time.

When creating a question, start with '質問:'.

When creating a response, start with '回答:'.

After finishing a question or response, always separate them with '\n\n'.

Table B: **Prompt to generate instruction data.** An image and paired text are provided to GPT-4o-mini along with this prompt.

Given a JSON array of objects, each with a 'from' and 'value' field, translate only the English text inside the 'value' field into Japanese. Keep the special token <image>\in the 'value' field unchanged. Do not change the overall structure of the JSON. Translate all English content in the 'value' field, even if it is a single word. Output only the translated JSON data and nothing else. Make sure the output is a valid JSON array that can be parsed with json.loads(). Do not include any text before or after the JSON.

Table C: **Prompt for translating English datasets into Japanese datasets.** Vision-Flan (Xu et al., 2024) and Image-Textualization (Pi et al., 2024), subsets of LLaVA OneVision (Li et al., 2024), are translated into Japanese using GPT-40-mini.

| Image | |
|--|---|
| Question | 写真に写っている花は何ですか?(What is the flower in the photo?) |
| GPT-4 (reference answer) | 写真に写っている花は桜です。満開の状態で、春の青空の下に美しく咲いています。これらの桜は、季節の象徴として、新しい始まりを感じさせます。 (The flowers in the photo are cherry blossoms. They are in full bloom, beautifully blossoming under the spring blue sky. These cherry blossoms symbolize the season and evoke a sense of new beginnings.) |
| LLaVA1.5-Swallow (stage 1 and 2) score 6.0 | 写真は、木の枝に咲く <mark>白い花</mark> を特徴としています。 (The photo features <mark>white flowers</mark> blooming on the branches of a tree.) |
| LLaVA1.5-Swallow (stage 3, CFT on PDF) score 7.0 | 写真に写っている花は桜の花です。桜は日本の春の象徴であり、毎年多くの人々が楽しむ花です。 (The flowers in the photo are cherry blossoms. Cherry blossoms are a symbol of spring in Japan and are enjoyed by many people every year.) |

Figure A: **Qualitative analysis on Heron-Bench.** Correct parts of the responses are highlighted in green, while incorrect parts are marked in red.

| Image | |
|--|---|
| Question | この写真に写っている鳥居の普通とは違うところはどこですか? (What is different about the torii in this photo compared to a typical one?) |
| GPT-4 (reference answer) | この写真に写っている鳥居の普通とは違うところは、水に浮かんでいる点です。通常、鳥居は地上に設置されますが、この鳥居は潮の満ち引きによって海上に現れる特徴があります。また、その背後に広がる緑豊かな山々とのコントラストが美しい景色を作り出しています。 (The difference in the torii in this photo compared to a typical one is that it appears to be floating on the water. Normally, torii gates are installed on land, but this one emerges above the sea depending on the tide. Additionally, the contrast with the lush green mountains in the background creates a beautiful scenic view.) |
| LLaVA1.5-Swallow (stage 1 and 2) score 4.0 | この写真の鳥居の珍しい点は、水面から突き出ていることです。通常、鳥居は陸上にあり、水域の上に建てられることはありません。鳥居は水面から突き出ているため、水面下の建物や構造物にアクセスするためのユニークな方法を提供し、興味深い視覚的要素を加えています。 (The unique feature of the torii in this photo is that it rises from the water. Typically, torii gates are built on land and are not constructed over bodies of water. Since this torii protrudes from the water, it provides a unique way to access structures beneath the surface and adds an interesting visual element.) |
| LLaVA1.5-Swallow (stage 3, CFT on PDF) score 9.0 | この写真に写っている鳥居は、通常の鳥居とは異なり、水の中に立っているため、非常に珍しい光景です。一般的な鳥居は地上にあり、水の上に立っていることはありません。このため、水面に映る鳥居の姿が特に印象的です。 (The torii in this photo is quite rare, as it stands in the water, unlike typical torii gates, which are usually built on land. Since torii gates are not normally positioned over water, the reflection of the gate on the water's surface is particularly striking.) |

Figure B: Another example of qualitative analysis on Heron-Bench.

| Image | |
|--|--|
| Question | 写真に写っているのはどこの国の料理が教えてください。 (Please tell me which country's cuisine is shown in the photo.) |
| GPT-4 (reference answer) | 写真に写っているのは日本の料理です。天ぷらは日本の伝統的な揚げ物料理で、野菜や海鮮を衣で包んで揚げたものです。この写真には、特にエビの天ぷらといくつかの野菜の天ぷらが盛り付けられており、日本料理の一例として認識されます。 (The food shown in the photo is Japanese cuisine. Tempura is a traditional Japanese deep-fried dish, where vegetables and seafood are coated in batter and fried. This photo specifically features shrimp tempura along with several types of vegetable tempura, making it a representative example of Japanese cuisine.) |
| LLaVA1.5-Swallow (stage 1 and 2) score 7.0 | 写真は日本の料理で、エビの天ぶらが皿の上に置かれています。 (The photo shows Japanese cuisine, with shrimp tempura placed on a plate.) [Comment]: It mentions only 'shrimp,' which is not sufficient to fully describe the picture. |
| LLaVA1.5-Swallow (stage 3, CFT on PDF) score 8.0 | 写真に写っている料理は日本の揚げ物の一種である天ぷらです。天ぷらは、野菜や魚などの食材を薄衣で揚げたもので、特にエビや野菜がよく使われます。 (The dish shown in the photo is tempura, a type of Japanese fried food. Tempura consists of ingredients such as vegetables and fish coated in a light batter and deep-fried, with shrimp and vegetables being commonly used.) |

Figure C: Further qualitative analysis on Heron-Bench.