
LLaVAR: Enhanced Visual Instruction Tuning for Text-Rich Image Understanding

Yanzhe Zhang^{1*}, Ruiyi Zhang², Jiuxiang Gu², Yufan Zhou², Nedim Lipka²,
Diyi Yang³, Tong Sun²

¹Georgia Tech, ²Adobe Research, ³Stanford University

Abstract

Instruction tuning enhances the capability of Large Language Models (LLMs) to interact with humans. Furthermore, recent instruction-following datasets include images as visual input, collecting responses for image-based instructions. However, current visual instruction-tuned models cannot comprehend textual details within images well. This work enhances the current visual instruction tuning pipeline with text-rich images (*e.g.*, movie posters, book covers, etc.). Specifically, we first used publicly available OCR tools to collect results on 422K text-rich images from the LAION dataset. Furthermore, we prompt text-only GPT-4 with recognized text and image captions to generate 16K conversations, each containing question-answer pairs for text-rich images. By combining our collected data with previous multi-modal instruction-following data, our model, **LLaVAR**, substantially improves the capability of the LLaVA model on text-based VQA datasets (up to 20% accuracy improvement). The GPT-4-based instruction-following evaluation also demonstrates the improvement of our model on both natural images and text-rich images. Through qualitative analysis, LLaVAR shows promising interaction skills (*e.g.*, reasoning, writing, and elaboration) with humans based on the latest real-world online content that combines text and images. We make our code/data/models publicly available.

1 Introduction

Instruction tuning [1, 2] improves generalization to unseen tasks by formulating various tasks into instructions. Such open-ended question-answering capability fosters the recent chatbot boom since ChatGPT. Recently, visual instruction-tuned models [3–5] further augment conversation agents with visual encoders such as CLIP-ViT [6, 7], enabling human-agent interaction based on images. However, possibly due to the dominance of natural images in training data (*e.g.*, Conceptual Captions [8] and COCO [9]), they struggle with understanding texts within images [10]. However, textual understanding is integral to visual perception in everyday life.

Fortunately, tools such as Optical Character Recognition (OCR, 11) allow us to recognize text in images. One naive way to utilize this is to add recognized texts to the input of visual instruction-tuned models [12]. However, such approach significantly increases the computation (longer context lengths), and might not fully leverage the encoding capability of visual encoders. To this end, we propose to enhance the end-to-end visual instruction-tuned model by collecting instruction-following data that require understanding texts within images.

*Collaborations through Adobe University Gift Program.

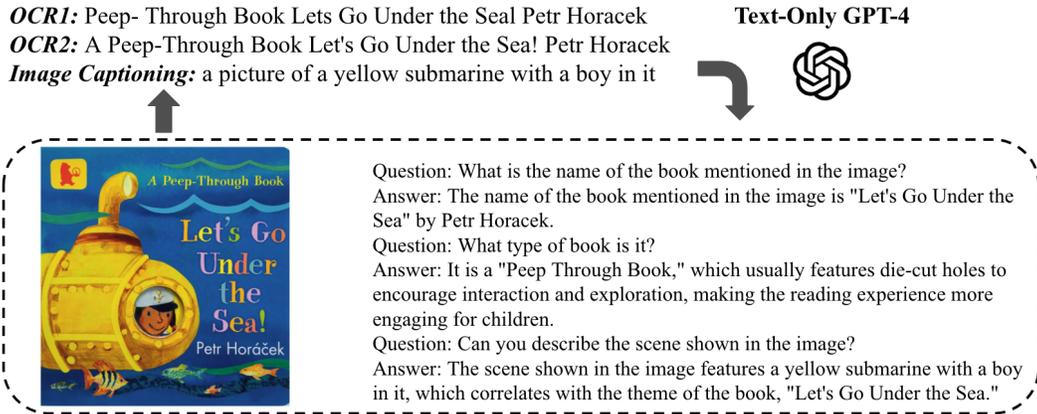


Figure 1: The process of collecting high-quality instruction-following data.

Specifically, we first collect 422K noisy instruction-following data using text-rich² images by combining manually written instructions (e.g., “Identify any text visible in the provided image.”) and the OCR results. Such large-scale noisy-aligned data effectively improve feature alignment between visual features and the language decoder. Furthermore, we prompt text-only GPT-4 [13] with OCR results and image captions to generate 16K conversations, where each conversation can be multiple turns of question & answer pairs, as high-quality instruction-following examples. This process requires GPT-4 to de-noise the OCR results and develop specific questions to create complex instructions based on the input (Figure 1).

To evaluate the effectiveness of the collected data, we use noisy and high-quality examples to augment the pretraining and fine-tuning stages of LLaVA [3] accordingly. We name our model **LLaVAR**, signifying the LLaVA (Large Language and Vision Assistant) that can **Read**. Compared to the original LLaVA, we also conducted experiments scaling the input resolution from 224² to 336² to better encode small textual details. Empirically, we report the results on four text-based VQA datasets following the evaluation protocol from Liu et al. [10]. Moreover, we apply GPT-4-based instruction-following evaluation to 30 natural images from COCO [9, 3] and 50 text-rich images from LAION [14]. We also provide qualitative analysis (e.g., on posters, website screenshots, and tweets) to test more complex instruction-following skills.

To sum up, our contributions are as follows:

- We collect 422K noisy instruction-following data and 16K high-quality instruction-following data. Both are shown to be effective in augmenting visual instruction tuning.
- Our model, LLaVAR, significantly enhances text understanding within images while slightly improving the model’s performance on natural images.
- The enhanced capability enables our model to provide end-to-end interactions based on various forms of online content that combine text and images.
- We open source the training and evaluation data together with the model checkpoints.

2 Related Work

Instruction Tuning Following natural language instructions is the key capability for an agent to interact with real-world users. Instruction tuning starts from collecting human-preferred feedback for human written instructions [1] or formulating multi-task training in a multi-task instruction-following manner [2, 15]. However, large, capable instruction-tuned models are usually closed-sourced and serve as commercial APIs only. Recently, Alpaca [16, 17], Vicuna [18], and Baize [19] start the trend of generating high-quality instruction-following data based on LLMs such as GPT-3.5 / ChatGPT /

²In this work, we use the phrase “text-rich images” to describe images with text in them, such as posters and book covers. In contrast, we refer to images without text as “natural images”.

GPT-4 and finetuning the open source LLaMA model [20]. However, evaluating the ability to follow instructions remains a challenge. While GPT-4 has demonstrated superior evaluation capabilities [21], there are still a number of concerns, such as bias toward response length [19] and lack of robustness to the order of examples [22]. Following Chiang et al. [18], Liu et al. [3], Dubois et al. [23], we use GPT-4-based instruction-following evaluation in this work.

Multimodal Instruction Tuning Recently, instruction tuning has been expanded to the multimodal setting, including image, video [24, 25], and audio [26, 27]. For image-based instruction tuning, MiniGPT-4 [28] employs ChatGPT to curate and improve detailed captions for high-quality instruction-following data. LLaVA [3] generates multimodal instruction-following data by prompting text-only GPT-4 with captions and object’s bounding boxes. LLaMA-Adapter [29, 12] uses COCO data for text-image feature alignment and utilizes textual data only for instruction tuning. mPLUG-owl [30] combines more than 1000M image-text pairs for pretraining and a 400K mixture of text-only/multimodal instruction-following data for finetuning. However, according to Liu et al. [10], most of these models struggle to accomplish tasks requiring OCR capability. InstructBLIP [31] transforms 13 vision-language tasks (including OCR-VQA [32]) into the instruction-following format for instruction tuning. Cream [33] applies multi-task learning that includes predicting masked texts in images. A more comprehensive survey can be found in Li et al. [34]. In this work, we select LLaVA as our baseline, which is the most data-efficient and powerful model, and demonstrate the effectiveness of our proposed pipeline.

3 Data Collection

Starting from the LAION-5B [14] dataset ³, our goal is only to keep images that are text-rich. Considering that documents usually contain plenty of text, we first obtained a binary classification dataset by combining natural images and document data. Subsequently, we trained an image classifier using a DiT [35]-base backbone, which was fine-tuned on the RVL-CDIP dataset [36]. Hopefully, such a classifier can predict whether an image contains text or not. We first build a subset by selecting images with a predicted probability greater than 0.8 while also satisfying $p(\text{watermark}) < 0.8$ and $p(\text{unsafe}) < 0.5$ ⁴. The derived subset is noisy due to the limitation of the classifier. To further clean up the data and incorporate human judgment, we randomly sampled 50K images and clustered them into 100 clusters based on CLIP-ViT-B/32 visual features. After inspecting the clustering results, we carefully select 14 clusters (see Figure 10 in the Appendix for examples) containing diverse text-rich images ranging from posters, covers, advertisements, infographics, educational materials, and logos. The cluster model is then used as the filter to collect images for constructing our instruction-following examples. As a reference, we provide a CLIP [7]-based categorization (see Appendix A for details.) to illustrate the distribution of images for both two types of data we collected in Figure 2. We summarize our collected data and LLaVA’s data in Table 1.

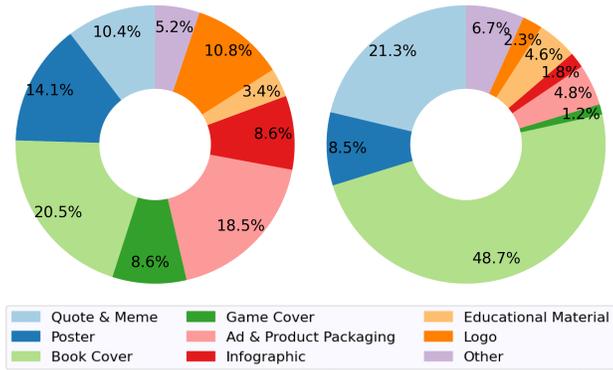


Figure 2: CLIP-based categorization of our collected images. The left refers to images used to collect noisy data, and the right refers to images used in the GPT-4 prompting. Both pie charts are based on 10K sampled images from the corresponding datasets.

Noisy Instruction-following Data Using the clustering model as a filter, we collect 422K deduplicated images that belong to the 14 preferred clusters. To balance the examples from different categories, we keep at most 52K examples for one cluster. We run all images through PaddleOCR ⁵.

³<https://huggingface.co/datasets/laion/laion-high-resolution>

⁴Both probabilities are from the LAION dataset’s metadata.

⁵<https://github.com/PaddlePaddle/PaddleOCR>

Data	Image	Instruction	# Conv	Avg Ins Len	Avg Res Len
LLaVA pretraining	CC3M	CC3M	595K	15.9	15.4
R _{pretraining} (Ours)	LAION	PaddleOCR	422K	17.2	48.8
LLaVA finetuning	COCO	GPT-4	158K	15.9	93.1
R _{finetuning} (Ours)	LAION	GPT-4	16K	15.1	40.5

Table 1: Summary of data statistics. R_{pretraining} and R_{finetuning} denote the additional pre-training / finetuning data we collected. The average instruction and response length are calculated after LLaMA tokenization.

Note that running OCR at the original resolution (e.g., 1024²) might recognize small fonts that are not visible by visual encoders like CLIP ViT (6, 7, resolution up to 336²). To ensure the recognition of visible fonts while maintaining OCR accuracy, we perform OCR on the image after downsampling (the short edge is resized to 384 pixels if longer than that.) to extract the text. Then, based on the geometric relationships between the recognized words, we merge them into paragraphs before concatenating them. As a robust instruction-following model should react similarly to instructions with similar meanings, we reword “Identify any text visible in the provided image.” into ten distinct instructions (Table 6 in Appendix). We then create a single-turn conversation for a given image by (i) randomly sampling an *input instruction* and (ii) using recognized texts as the desired *output response*. Such instruction-following data is noisy because of the relatively limited performance of OCR tools on diverse fonts and colorful backgrounds.

GPT-4-based Instruction-following Data Compared to high-quality instruction-following data, there are mainly two issues for the noisy data collected above. (i) Responses should contain organized sentences instead of raw OCR results with missing words and grammar errors. (ii) Instructions should be diverse, suitable and specific to the given image instead of monotonously asking for all visible texts. To address these issues, we follow Liu et al. [3] to generate instruction-following data by prompting text-only GPT-4 [13] with OCR results and captions.

It is challenging to prompt GPT-4 with fragmented OCR results in a few words to generate non-trivial instructions. To this end, we carefully select 4 of the 14 previously mentioned clusters (the 3rd, 4th, 6th and 9th clusters in Figure 10) to collect images with enough visible and coherent sentences. As shown in Figure 2, such filtering dramatically increases the percentage of book covers and quote images. We randomly selected 4K examples from each cluster (no overlap with images used for noisy instruction-following data), yielding a total of 16K images. Following prior work [16, 17, 3], we provide the visualization of verb-noun pairs for instructions generated by GPT-4 in Appendix Figure 11. For those instructions without a verb-noun pair, we demonstrate the frequency of objects being asked in Appendix Figure 12.

Furthermore, based on the system message and two in-context few-shot examples (shown in Appendix B), we ask GPT-4 to generate conversational data based on OCR results and image captions (Figure 1). The generated questions are used as *input instructions*, and answers are used as *output responses*. Concretely, for a given image, we first provide two OCR results from EasyOCR and PaddleOCR, which can complement each other. To illustrate visual elements other than texts within the image, we also provide the result of BLIP-2 image captioning [37]. To prevent the caption from focusing on the text, we use OCR bounding boxes to mask the text and then use the inpainting [38] to refill the mask before using generation captions. Note that captioning models might suffer from hallucinations [39]. We mention this unreliability in our system message and ask GPT-4 only to generate questions with sure answers. We leave the generation of more detailed captions [40, 41] for future work.

4 Model Architecture and Training

Architecture In most of our study, we use the same model architecture as LLaVA. For the visual encoder V , we use CLIP-ViT-L/14 for 224² resolution and CLIP-ViT-L/14-336 for 336² resolution. The grid features before the last transformer layer are then transformed into the word embedding space of the language decoder through a trainable projection matrix W . We use Vicuna-13B [18], a

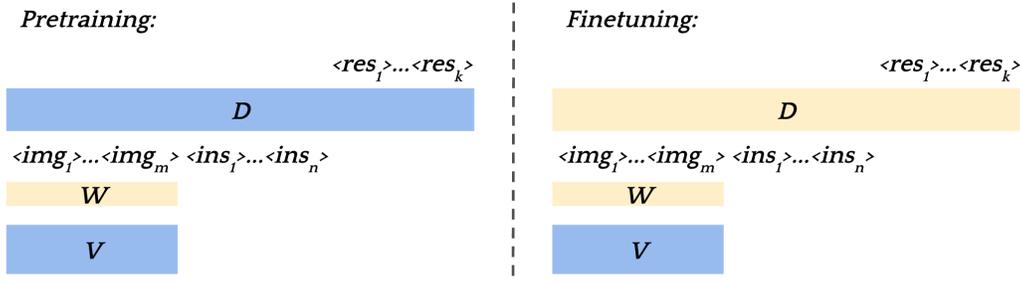


Figure 3: The model training process for the visual encoder V , projection matrix W , and language decoder D . **Blue blocks** denote frozen modules and **yellow blocks** denote trainable modules. The training input is image tokens ($\langle img \rangle$) and instruction tokens ($\langle ins \rangle$), while the target is response tokens ($\langle res \rangle$).

LLaMA-based [20] instruction-tuned language model, as the language decoder D except the ablation study in Table 4.

In Section 5.1 and Appendix H, we extend the current architecture by adding an extra high-resolution (high-res) visual encoder. Such a high-res encoder outputs thousands of patch features, which means that the transformed features and instruction tokens cannot fit in the context length of the language decoder. To this end, we propose to add cross-attention modules to the decoder, which attend to key-value pairs transformed from the high-res patch features.

Training We follow the two-stage training design of LLaVA (Figure 3). The training objectives of both stages are the same: generate **output responses** ($\langle res \rangle$) for the **input instructions** ($\langle ins \rangle$). The transformed image tokens ($\langle img \rangle$) are added before or after the first input instruction. **(i)** During the first pre-training stage, only the projection matrix W is trained for feature alignment. Since the decoder D is frozen, training tolerates noisy data. In the pre-training stage, we combine the 595K pre-training data from LLaVA with our 422K noisy instruction-following data. **(ii)** Both the projection matrix W and the language decoder D are trained during the finetuning stage, where we merge our 16K instruction-following data into the 158K instruction-following data from LLaVA as the training set. Note that the visual encoder is frozen throughout the training period, which might restrict text recognition performance, as CLIP is trained for general-purpose text-image alignment. Better choices of the visual encoder [42] or CLIP-ViT finetuning [30] may further benefit the visual understanding capability, which we leave for future work.

5 Experiments

We use the same training hyperparameters as LLaVA⁶, except that **(i)** We set the maximum sequence length to 1024 during pre-training. **(ii)** We first pad any given image to a square shape before resizing it to the desired input size, preventing some image content from cropping during preprocessing. For both resolutions (224^2 , 336^2), we reproduce the original LLaVA for a fair comparison. The GPT-4 model used in this work refers to the gpt-4-0314 version, while the cost to collect finetuning data is around \$300. The temperature used to sample GPT-4 is set to 1.0 for the generation of training data, 0.7 for the generation of evaluation data, and 0.2 for the evaluation based on GPT-4. All experiments are run on NVIDIA A100 80GB GPUs. During the evaluation, the temperature used to sample from our model is set at 0.9 for text-based VQA, 0.7 for GPT-4-based instruction-following evaluation, and 0.2 for other qualitative demonstrations.

5.1 Quantitative Analysis

Text-based VQA Following the evaluation protocol in Liu et al. [10], we evaluate the performance of LLaVAR on four text-based VQA datasets: ST-VQA [45], OCR-VQA [32], TextVQA [46], and DocVQA [47], representing various domains (see Appendix C for more details and Appendix E for

⁶<https://github.com/haotian-liu/LLaVA>

	Res.	ST-VQA	OCR-VQA	TextVQA	DocVQA
BLIP-2 [2023] †		21.7	30.7	32.2	4.9
OpenFlamingo [2023] †		19.3	27.8	29.1	5.1
MiniGPT4 [2023] †	224 ²	14.0	11.5	18.7	3.0
LLaVA [2023] †		22.1	11.4	28.9	4.5
mPLUG-Owl [2023] †		29.3	28.6	40.3	6.9
LLaVA ‡		24.3	10.8	31.0	5.2
LLaVAR	224 ²	30.2 (+5.9)	23.4 (+12.6)	39.5 (+8.5)	6.2 (+1.0)
LLaVA ‡		28.9	11.0	36.7	6.9
LLaVAR	336 ²	39.2 (+10.3)	23.8 (+12.8)	48.5 (+11.8)	11.6 (+4.7)

Table 2: Results (accuracy %) on text-based VQA. We use † to refer to the results obtained from Liu et al. [10] and ‡ to refer to our reproduced results. The accuracy metric used by Liu et al. [10] only counts for whether the ground truth appears in the response. For more metrics, please refer to Appendix D.

	ST-VQA	OCR-VQA	TextVQA	DocVQA
(1) LLaVA	28.9	11.0	36.7	6.9
(2) LLaVA + $R_{\text{pretraining}}$	36.7	26.1	46.5	9.6
(3) LLaVA + $R_{\text{finetuning}}$	34.1	21.6	43.6	9.5
(4) LLaVA + $C_{\text{pretraining}}$	35.4	27.0	45.6	9.2
(5) LLaVA + $N_{\text{finetuning}}$	34.1	25.9	43.3	10.2
(6) LLaVAR	39.2	23.8	48.5	11.6

Table 3: Ablation Study on pretraining/finetuning data. All results are from 336²-based models. $R_{\text{pretraining}}$ and $R_{\text{finetuning}}$ denote the extra pretraining/finetuning data we collected. $C_{\text{pretraining}}$ refers to using captions instead of OCR results as responses during pretraining. $N_{\text{finetuning}}$ refers to using written questions + raw OCR results instead of GPT-generated QA for finetuning.

more datasets). We present the results of the baseline models and our models in Table 2. Note that InstructBLIP includes OCR-VQA in its training sets, making it incomparable with our settings. In both resolution settings and all four datasets, LLaVAR substantially improves the LLaVA baseline, demonstrating that our collected data can bring about a robust improvement. Furthermore, the improvement is more significant in 336² resolution compared to 224², indicating that the collected data might bring a greater improvement at even higher resolutions. Our best model, 336²-based LLaVAR, performs best in 3 out of 4 evaluated datasets. Note that this is not a fair comparison. Some key factors include different language decoders, resolutions, and magnitudes of text-image training data. We provide more discussions on the comparison with mPLUG-Owl and the result of finetuning mPLUG-Owl using our data in Appendix F.

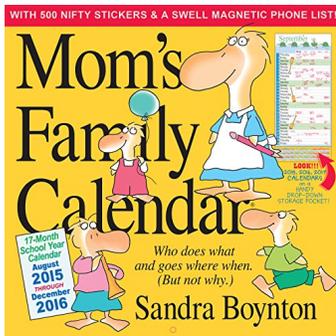
Ablation Study on pretraining/finetuning data We report the result in Table 3 and Figure 4.

(i) Based on variants (2) and (3), we find that the collected data can benefit the pretraining stage ($R_{\text{pretraining}}$) and finetuning stage ($R_{\text{finetuning}}$) separately while being complementary to each other in most cases⁷. More importantly, enhancing the pretraining stage alone achieves the second-best overall performance, indicating the potential to boost textual detail understanding without dependence on GPT-4-generated high-quality data. (ii) Using pretraining images, we obtain $C_{\text{pretraining}}$ by replacing the pretraining instructions with questions & captions, the same pattern as LLaVA. As variant (4) is not as good as (2), we can conclude that OCR is more advantageous than captions. (iii) We further validate the value of GPT-4 generated data by generating noisy finetuning data ($N_{\text{finetuning}}$), similar to pretraining data. Variant (5) achieves comparable accuracy as variant (3). However, as shown in Figure 4, such noisy finetuning data hurts the instruction-following capability: (5) responds with all recognized texts while ignoring the questions.

⁷Since the metric only consider the recall, it might favor variant (2)(4)(5) due to their longer outputs.

	CLIP Res.	Extra Enc.	$R_{\text{pretraining}}$	ST-VQA	OCR-VQA	TextVQA	DocVQA
(a)	224 ²	✗	Low	28.9	25.6	37.8	6.2
(b)	336 ²	✗	Low	37.4	31.0	45.7	11.4
(c)	224 ²	✗	High	28.9	24.9	35.8	6.2
(d)	336 ²	✗	High	36.9	30.6	46.5	11.7
(e)	224 ²	✓	Low	31.8	30.9	41.2	8.1
(f)	224 ²	✓	High	33.5	29.8	42.2	10.2
(g)	224 ²	✓	High (2×)	35.8	30.7	45.6	15.3

Table 4: Ablation Study on encoders and image resolution. “Extra Enc.” refers to whether we use Pix2Struct-base [44] as an extra high-res visual encoder. “ $R_{\text{pretraining}}$ ” refers to images used for pretraining (before preprocessing) and the collection of OCR results, where “Low” refers to the resized images (the same as the rest of this paper), “High” refers to the original high-res images, and “2×” means doubling the number of examples using the same criteria. All results from this table use Vicuna-7B, as adding the extra encoder requires more computation and memory.



Question: Who wrote this book?

- GT Sandra Boynton
- (1) The book was written by Sandra Boynton.
 - (2) The book was written by Sandra Boynton. It is a children’s book that features a family calendar with an accompanying activity guide. The book provides information (...) The cover of the book displays a picture of the family from the calendar, which helps to visually represent the topic of the book.
 - (3) The book is written by Sandra Byington.
 - (4) The book was written by Sandra Boynton.
 - (5) With 500+ Nity StickerS &aSwEllMagnetic Phone List! Moms 291573016171800 FAMILI (Pocket) Calendar (...) 2016 SANDRA BOYNTON
 - (6) The book was written by Sandra Boynton.

Figure 4: Ablation study based an example from OCR-VQA. GT refers to ground truth, and (1) - (6) are different model variants from Table 3. We replace the excessively long response with (...).

Ablation Study on encoders/image resolution While keeping finetuning data the same, we report the quantitative results of adding an extra visual encoder and varying the pretraining data in Table 4. (i) Take Pix2Struct-base as an example, we find that adding an extra high-res visual encoder with cross-attention indeed improves the performance ((g) vs. (a)), especially achieving the best zero-shot performance on DocVQA (15.3% accuracy). The performance gain on other datasets is relatively limited, probably due to the extra encoder we use being pretrained on web pages instead of natural images. On the other hand, the performance of (e) and (f) remains poor, without doubling the number of high-res examples in $R_{\text{pretraining}}$. Given the larger number of parameters initialized in the cross-attention module, they may be underfitting when trained on the same data as the projection matrix W (e.g., (e) vs. (b)), similar to the finding in Zeng et al. [48]. (ii) Considering (c) vs. (a) and (d) vs. (b), while the images are resized to the same size after preprocessing, high-res OCR results turn out to be not necessarily better than the low-resolution version, suggesting the capability of the visual encoder is almost saturated in (a) and (b). For more details and results on the extra high-res encoder, please refer to Appendix H.

GPT-4-based instruction-following evaluation We also report the GPT-4 evaluation results on instruction-following questions in Table 5. (i) **Natural Images:** 90 questions based on 30 COCO validation images from Liu et al. [3], including three aspects: conversation, detail description, and complex reasoning. This aims to test whether our collected data will hurt, maintain, or improve the model’s performance on natural images. First of all, using a higher resolution brings improvement (+2.9) in the performance of detail description, which is intuitive. Furthermore, LLaVAR achieves a better trade-off and increases the performance of all three aspects (+1.6 on average). More details are in Appendix J. (ii) **Text-Rich Images:** Similar to collecting the finetuning data, we leverage

	Res	Conversation	Detail	Complex	Read
LLaVA (Original)	224 ²	83.1	75.3	96.5	-
LLaVA	336 ²	83.9	78.2	95.3	87.9
LLaVAR	336 ²	84.5	78.9	96.5	91.7

Table 5: Relative scores (w.r.t. text-only GPT-4) for instruction-following questions, where the first three dimensions are based on natural images, the last dimension (“Read”) is based on text-rich images. In the first row, we show the original results (224²-based) fetched from Liu et al. [3]. We report our reproduced LLaVA on 336² resolution for a fair comparison.

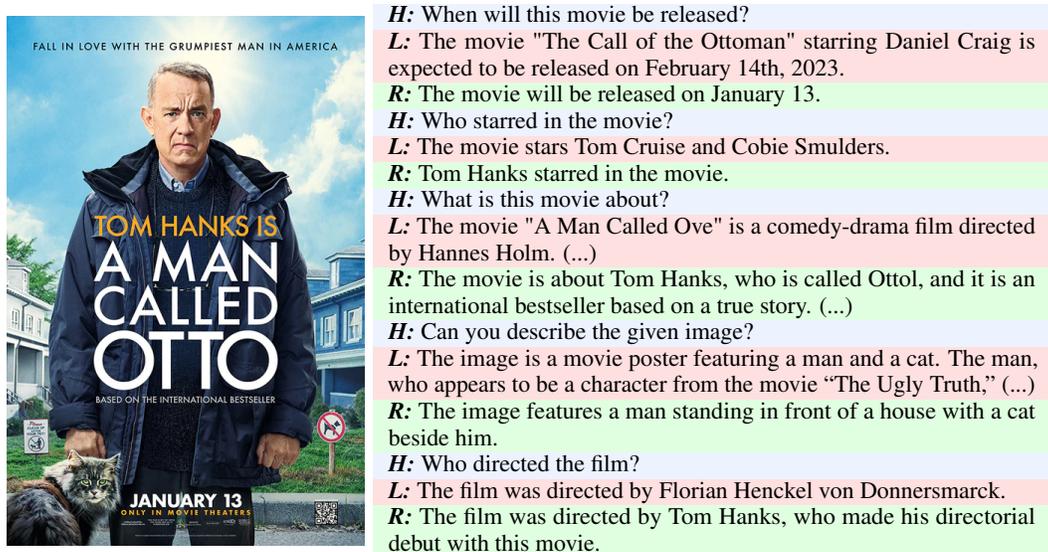


Figure 5: Qualitative Comparison. We color the human inputs (**H**) in blue, LLaVA’s responses (**L**) in pink, and LLaVAR’s responses (**R**) in green. Both are 336²-based models. We replace the excessively long response with (...).

50 text-rich images from LAION to collect instruction-following questions based on OCR results and human-annotated captions. We then collect responses from our trained model and use GPT-4 to calculate the relative score w.r.t GPT-4 responses. We add this as an extra dimension “**Read**” to Table 5, where our model demonstrates a significant (+3.8) improvement. The Appendix provides an example in Table 17.

5.2 Qualitative Analysis

We use a recent movie poster ⁸ to demonstrate the difference between LLaVA and LLaVAR when interacting with humans based on text-rich images. LLaVA, without augmenting textual understanding within images, suffers from hallucination when answering these questions. Some mentioned movies, like “A Man Called Ove” and “The Ugly Truth”, are real movies, suggesting that the language decoder is hallucinating its internal knowledge, while the visual encoder cannot encode helpful information. Alternatively, LLaVAR can correctly answer many of the provided questions with **faithful** information, which is clearly grounded in the image. However, some limitations remain, such as the spelling error “ottol” (We provide more statistics related to such spelling errors in Appendix I). Also, the final question asks for information that is not observable from the given poster, where an expected response should express such uncertainty instead of giving concrete answers. Nevertheless, neither model correctly answers the question.

⁸<https://www.imdb.com/title/tt7405458/>

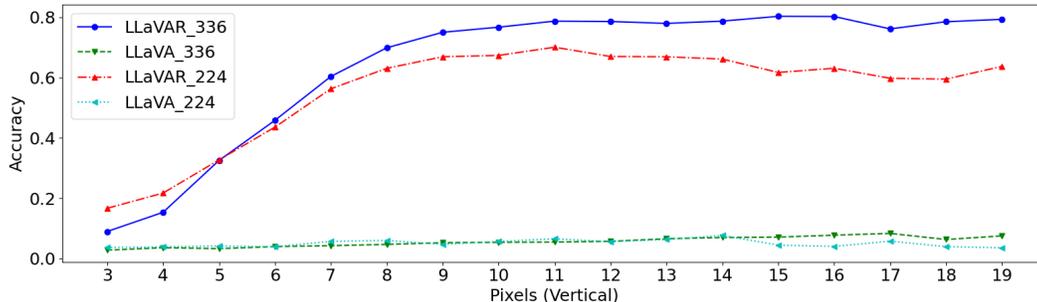


Figure 6: Case study of the recognizable font size, in which the x-axis refers to the height of ground truth answers in the image and the y-axis stands for the answer accuracy of models. We plot the results for both 224^2 -based models and 336^2 -based models.

5.3 Case Study: Recognizable Font Size

We first collect 825 examples from OCR-VQA, which have answers directly presented in the image and are detectable by the OCR tool. By rescaling the images, we test the model’s performance in answering these questions while the vertical heights of answers range from 3 pixels to 19 pixels. We report the result in Fig 6. **(i)** For the baseline model LLaVA, it struggles to provide correct answers in all scenarios, for both 224^2 -based and 336^2 -based versions. **(ii)** Our model LLaVAR achieves significantly better results in all scales. We observe a threshold for recognizable texts for both 224^2 -based and 336^2 -based versions as the accuracy sharply decreases when the height is smaller than 7 pixels. More interestingly, the 224^2 -based version achieves better performance on small texts with 3 pixels height while the 336^2 -based achieves better performance on large texts with more than 7 pixels height. We assume the extra training stage of CLIP 336^2 makes it better on the larger scale but worse on the smaller scale.

5.4 Transferred Instruction-following Capability

According to the dataset statistics (Table 1) and the visualization (Figure 11), our collected instruction-following data is not as diverse and substantial as LLaVA. This can be attributed to the relatively limited information given GPT-4 compared to five different human-written captions used in LLaVA. The content of text-rich images is also less diverse than that of natural images. While using more complex in-context examples can definitely stimulate generating more complicated instruction-following examples, it can also multiply the cost. In Appendix Figure 9, we demonstrate the transferred instruction-following capability of LLaVA, potentially from both the LLaVA data and the Vicuna backbone. While the extra data we add mainly focuses on understanding the visible texts within images, LLaVAR manages to build its reasoning, writing, and elaboration skills based on the top of its text recognition capability in an end-to-end manner. This allows users to interact with various online content based on simple screenshots.

6 Conclusion

In this work, we enhance visual instruction-tuned models in terms of their capability to read texts in images. Using text-rich images from the LAION dataset, we collect 422K noisy instruction-following examples using OCR results only and 16K high-quality instruction-following data based on text-only GPT-4. These two sets of data are leveraged to augment the pretraining stage and finetuning stage of LLaVA accordingly. Our model, LLaVAR, demonstrates superior performance in understanding texts within images and following human instructions on both prior benchmarks and real-world online content. Moreover, our analysis shows that the same augmented data is more effective with higher resolution. Additionally, using noisy instruction-following examples to augment pretraining essentially boosts the model performance without prompting GPT-4. For future work, we encourage exploration of **(i)** better image selection criteria or domain reweighting strategy [49] and **(ii)** more data-efficient and computation-efficient ways to augment instruction-following models with multimodal capability, especially in the high-res scenario.

References

- [1] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
- [2] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022.
- [3] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
- [4] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023.
- [5] Chunyuan Li. Large multimodal models: Notes on cvpr 2023 tutorial. *ArXiv*, abs/2306.14895, 2023.
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2020.
- [7] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.
- [8] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts, 2021.
- [9] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015.
- [10] Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Cheng lin Liu, Lianwen Jin, and Xiang Bai. On the hidden mystery of ocr in large multimodal models, 2023.
- [11] S. Mori, C.Y. Suen, and K. Yamamoto. Historical review of ocr research and development. *Proceedings of the IEEE*, 80(7):1029–1058, 1992. doi: 10.1109/5.156468.
- [12] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. Llama-adapter v2: Parameter-efficient visual instruction model, 2023.
- [13] OpenAI. Gpt-4 technical report, 2023.
- [14] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022.
- [15] Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krима Doshi, Maitreya Patel, Kuntal Kumar Pal, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Shailaja Keyur Sampat, Savan Doshi, Siddhartha Mishra, Sujan Reddy,

- Sumanta Patro, Tanay Dixit, Xudong Shen, Chitta Baral, Yejin Choi, Noah A. Smith, Hannaneh Hajishirzi, and Daniel Khashabi. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks, 2022.
- [16] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions, 2022.
- [17] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [18] 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. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- [19] Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data, 2023.
- [20] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- [21] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg evaluation using gpt-4 with better human alignment, 2023.
- [22] Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. Large language models are not fair evaluators, 2023.
- [23] Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback, 2023.
- [24] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.
- [25] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models, 2023.
- [26] Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jia-Bin Huang, Jinglin Liu, Yixiang Ren, Zhou Zhao, and Shinji Watanabe. Audiogpt: Understanding and generating speech, music, sound, and talking head. *ArXiv*, abs/2304.12995, 2023.
- [27] Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities, 2023.
- [28] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023.
- [29] Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention, 2023.
- [30] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. mplug-owl: Modularization empowers large language models with multimodality, 2023.
- [31] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.

- [32] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *2019 international conference on document analysis and recognition (ICDAR)*, pages 947–952. IEEE, 2019.
- [33] Geewook Kim, Hodong Lee, Daehee Kim, Haeji Jung, Sanghee Park, Yoonsik Kim, Sangdoon Yun, Taeho Kil, Bado Lee, and Seunghyun Park. Cream: Visually-situated natural language understanding with contrastive reading model and frozen large language models, 2023.
- [34] Chunyuan Li, Zhe Gan, Zhengyuan Yang, Jianwei Yang, Linjie Li, Lijuan Wang, and Jianfeng Gao. Multimodal foundation models: From specialists to general-purpose assistants, 2023.
- [35] Junlong Li, Yiheng Xu, Tengchao Lv, Lei Cui, Cha Zhang, and Furu Wei. Dit: Self-supervised pre-training for document image transformer. *Proceedings of the 30th ACM International Conference on Multimedia*, Oct 2022. doi: 10.1145/3503161.3547911. URL <http://dx.doi.org/10.1145/3503161.3547911>.
- [36] Adam W. Harley, Alex Ufkes, and Konstantinos G. Derpanis. Evaluation of deep convolutional nets for document image classification and retrieval, 2015.
- [37] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023.
- [38] Alexandru Telea. An image inpainting technique based on the fast marching method. *Journal of graphics tools*, 9(1):23–34, 2004.
- [39] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning, 2018.
- [40] Noam Rotstein, David Bensaid, Shaked Brody, Roy Ganz, and Ron Kimmel. Fusecap: Leveraging large language models to fuse visual data into enriched image captions, 2023.
- [41] Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A. Smith, and Jiebo Luo. Promptcap: Prompt-guided task-aware image captioning, 2022.
- [42] Michael Tschannen, Basil Mustafa, and Neil Houlsby. Clippo: Image-and-language understanding from pixels only, 2022.
- [43] Anas Awadalla, Irena Gao, Joshua Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo, March 2023. URL <https://doi.org/10.5281/zenodo.7733589>.
- [44] Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct: Screenshot parsing as pretraining for visual language understanding, 2022.
- [45] Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluís Gomez, Marçal Rusinol, Minesh Mathew, C.V. Jawahar, Ernest Valveny, and Dimosthenis Karatzas. Icdar 2019 competition on scene text visual question answering. *2019 International Conference on Document Analysis and Recognition (ICDAR)*, Sep 2019. doi: 10.1109/icdar.2019.00251. URL <http://dx.doi.org/10.1109/ICDAR.2019.00251>.
- [46] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2019. doi: 10.1109/cvpr.2019.00851. URL <http://dx.doi.org/10.1109/CVPR.2019.00851>.
- [47] Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images, 2020.
- [48] Yan Zeng, Hanbo Zhang, Jiani Zheng, Jiangnan Xia, Guoqiang Wei, Yang Wei, Yuchen Zhang, and Tao Kong. What matters in training a gpt4-style language model with multimodal inputs? *arXiv preprint arXiv:2307.02469*, 2023.

- [49] Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V. Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model pretraining, 2023.
- [50] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering, 2022.
- [51] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multi-modal chain-of-thought reasoning in language models, 2023.
- [52] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. *ArXiv*, abs/2304.09842, 2023.
- [53] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. Visual genome: Connecting language and vision using crowdsourced dense image annotations, 2016.
- [54] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- [55] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, Serge Belongie, Victor Gomes, Abhinav Gupta, Chen Sun, Gal Chechik, David Cai, Zheyun Feng, Dhyanesh Narayanan, and Kevin Murphy. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*, 2017.
- [56] Brian Kenji Iwana, Syed Tahseen Raza Rizvi, Sheraz Ahmed, Andreas Dengel, and Seiichi Uchida. Judging a book by its cover, 2016.
- [57] Anhar Risnumawan, Palaiiahankote Shivakumara, Chee Seng Chan, and Chew Lim Tan. A robust arbitrary text detection system for natural scene images. *Expert Systems with Applications*, 41(18):8027–8048, 2014.
- [58] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning, 2022.
- [59] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015.
- [60] Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In Jade Goldstein, Alon Lavie, Chin-Yew Lin, and Clare Voss, editors, *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics. URL <https://aclanthology.org/W05-0909>.
- [61] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL <https://aclanthology.org/W04-1013>.
- [62] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation, 2014.
- [63] Jianfeng Kuang, Wei Hua, Dingkan Liang, Mingkun Yang, Deqiang Jiang, Bo Ren, and Xiang Bai. Visual information extraction in the wild: Practical dataset and end-to-end solution, 2023.

Instructions
Identify any text visible in the image provided.
List all the text you can see in the given image.
Enumerate the words or sentences visible in the picture.
Describe any readable text present in the image.
Report any discernible text you see in the image.
Share any legible words or sentences visible in the picture.
Provide a list of texts observed in the provided image.
Note down any readable words or phrases shown in the photo.
Report on any text that can be clearly read in the image.
Mention any discernible and legible text present in the given picture.

Table 6: Ten instructions asking for OCR results.

A

CLIP-based categorization Based on the observation of selected clusters, we divide the images used into 8 categories. For each category, we use one or multiple words as labels.

- **Quote & Meme:** “quote”, “internet meme”.
- **Poster:** “movie poster”, “podcast poster”, “TV show poster”, “event poster”, “poster”.
- **Book Cover:** “book cover”, “magazine cover”.
- **Game Cover:** “game cover”.
- **Ad & Product Packaging:** “ad”, “advertisement”, “food packaging”, “product packaging”.
- **Infographic:** “chart”, “bar chart”, “pie chart”, “scatter plot”.
- **Educational Material:** “exam paper”, “quiz”, “certificate”, “book page”.
- **Logo:** “logo”.

For each word, we use the following templates to achieve embedding-space ensembling [7]:

- “a photo of a {}.”
- “a blurry photo of a {}.”
- “a black and white photo of a {}.”
- “a low contrast photo of a {}.”
- “a high contrast photo of a {}.”
- “a bad photo of a {}.”
- “a good photo of a {}.”
- “a photo of a small {}.”
- “a photo of a big {}.”

For each image, we calculate the similarity between the image and all words mentioned above using CLIP-ViT-L/14. If the highest similarity is less than 0.15, we then classify the image into **Other**, otherwise we classify into the “super class” (e.g., **Poster**) of the word (e.g., “movie poster”) with the highest similarity.

B

System Message adapted from Liu et al. [3].

You are an AI visual assistant, and you are seeing a single image. What you see is provided with two OCR results and one image caption describing the information within the same image you are looking at. Image captions might include hallucinations, while OCR results are more accurate. Answer all questions with definite answers 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 (e.g., the man, the sunset, the ocean.) and the texts contained in the image. 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 questions that cannot be answered confidently;

(3) DO NOT mention OCR or image caption in your questions and answers;

(4) DO NOT ask about information from captions while it looks unrelated to or contradicts OCR results.

Also include complex questions that are relevant to the content in the image, for example, asking about background knowledge of the texts in the image, asking to discuss about the design of 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.

Few-shot Examples

Human:

WHAT F YOUR BLESSINGS COME THROUGH RAINDROPS PDF Included LAURASTORY A DEVOTIONAL BASED ON THE GRAMMY AWARD WINNING SONG THAT HAS INSPIRED MILLION

WHAT IF YOUR BLESSINGS COME THROUGH RAINDROPS PDF Included LAURASTORY A DEVOTIONAL BASED ON THE GRAMMY AWARD WINNING SONG THAT HAS INSPIRED MILLION

a girl is standing in a field with a rainbow

GPT:

Question: What is the name of the devotional mentioned in the image?

Answer: The devotional is called "What If Your Blessings Come Through Raindrops" by Laura Story.

Question: What is special about the song mentioned in the image?

Answer: The song is special because it's a Grammy Award-winning song that has inspired millions of people."

Human:

*One of the hardest things in life to accept is a called third strike Robert Frost te
One of the hardest things in life to accept is a called third strike Robert Frost quote fancy
a close up of a baseball glove*

GPT:

Question: Why is the third strike hard to accept? Explain the quote to me.

Answer: The quote "One of the hardest things in life to accept is a called third strike" attributed to Robert Frost is a metaphor for life situations. In baseball, a called third strike is when the umpire determines that a pitch was a strike, but the batter did not swing. This can be frustrating for the batter because they missed an opportunity to hit the ball due to either a lack of confidence or misjudgment.

In life, this metaphor refers to difficult moments where we might miss opportunities or face setbacks because of our own inaction or hesitation. The quote suggests that it is

hard to accept these missed opportunities or challenges, especially when they result from our own choices or lack of action.

Method	Subject			Context Modality			Grade		Average
	NAT	SOC	LAN	TXT	IMG	NO	G1-6	G7-12	
Human [2022]	90.23	84.97	87.48	89.60	87.50	88.10	91.59	82.42	88.40
GPT-3.5 [2022]	74.64	69.74	76.00	74.44	67.28	77.42	76.80	68.89	73.97
GPT-3.5 w/ CoT [2022]	75.44	70.87	78.09	74.68	67.43	79.93	78.23	69.68	75.17
LLaMA-Adapter [2023]	84.37	88.30	84.36	83.72	80.32	86.90	85.83	84.05	85.19
MM-CoT _{Base} [2023]	87.52	77.17	85.82	87.88	82.90	86.83	84.65	85.37	84.91
MM-CoT _{Large} [2023]	95.91	82.00	90.82	95.26	88.80	92.89	92.44	90.31	91.68
LLaVA [2023]	90.36	95.95	88.00	89.49	88.00	90.66	90.93	90.90	90.92
LLaVA+GPT-4 [2023]	91.56	96.74	91.09	90.62	88.99	93.52	92.73	92.16	92.53
Chameleon (GPT-4) [2023]	89.83	74.13	89.82	88.27	77.64	92.13	88.03	83.72	86.54
LLaVAR	91.79	93.81	88.73	90.57	88.70	91.57	91.30	91.63	91.42

Table 7: Results (accuracy %) on Science QA dataset. All baseline results are from Liu et al. [3], Lu et al. [52]. The categories are denoted as NAT: natural science, SOC: social science, LAN: language science, TXT: text context, IMG: image context, NO: no context, G1-6: grades 1-6, G7-12: grades 7-12.

C

Details of evaluation datasets used in the main paper:

- ST-VQA [45] contains 31791 questions that require understanding the scene text, based on images from COCO [9], Visual Genome [53], ImageNet [54], etc.
- TextVQA [46] contains 45,336 questions that need reading and reasoning about the text in images to answer, based on images from OpenImages [55].
- OCR-VQA [32] contains more than 1 million questions asking about information from book cover images [56].
- DocVQA [47] contains 50000 questions based on document images.

Details of extra datasets in Appendix:

- CT80 [57] contains 80 images for curved text OCR evaluation. The formats of questions are: (1) "What is written in the image?" for English words. (2) "What is the number in the image?" for digit string.
- POIE [46] contains 3000 camera images collected from the Nutrition Facts label of products, together with 111,155 text instances. The format of questions is "What is {entity name} in the image?".
- ChartQA [58] includes 4,804 charts with 9608 human-written questions.

D

	Res.	METEOR	ROUGE-L	CIDEr
LLaVA	224 ²	7.0	8.2	15.3
LLaVAR		10.0	11.4	24.5
LLaVA	336 ²	8.4	9.9	19.1
LLaVAR		12.8	14.3	30.9

Table 8: Results on ST-VQA using text-matching metrics.

	Res.	METEOR	ROUGE-L	CIDEr
LLaVA	224 ²	8.7	10.5	12.2
LLaVAR		12.5	14.9	21.4
LLaVA	336 ²	9.9	12.1	15.3
LLaVAR		14.8	17.4	27.0

Table 9: Results on textVQA using text-matching metrics.

	Res.	METEOR	ROUGE-L	CIDEr
LLaVA	224 ²	0.2	0.1	0.0
LLaVAR		0.3	0.1	0.0
LLaVA	336 ²	0.3	0.1	0.0
LLaVAR		0.2	0.1	0.0

Table 10: Results on OCR-VQA using text-matching metrics.

Results of other metrics The metric used for text-based VQA in the main paper is the standard practice in VQA benchmarks [59]. For STVQA and DocVQA, previous works use ANLS (Average Normalized Levenshtein Similarity) as the metric [45, 47], which calculates the average normalized edit distance and only works for supervised models trained to output short and precise answers. It works badly for instruction-following models that usually output long sequences instead of brief answers. For reference, we provide more text-matching metrics (METEOR, 60, ROUGE-L, 61, CIDEr, 62) to demonstrate the improvement of our model (Table 8, 9, 10, 11), which works well except for OCR-VQA. We assume these metrics are not valuable for OCR-VQA since the ground truth answers are usually too short.

E

Results on extra datasets In Table 12, we provide results on three extra datasets: CT80 (OCR, 57), POIE (Information Extraction, 63), and ChartQA [58]. We use the same VQA metric as other text-based VQA datasets. We observe similar trends as the main paper results: LLaVAR data significantly improves over the LLaVA baseline, usually more considerably in a higher resolution.

F

Comparison with mPLUG-Owl We find that LLaVAR usually performs similarly well with mPLUG-Owl in the same 224² resolution. We further clarify the setting differences between mPLUG-Owl and ours: mPLUG-Owl is trained on 1000M+ text-image pairs, while the original LLaVA is trained on about 0.6M text-image pairs. Our model, LLaVAR, is trained on about 1M text-image pairs. Within the same resolution, LLaVAR demonstrates a good performance with decent data efficiency.

We presume that training on large-scale non-OCR data improves OCR performance, as many of the captions in LAION datasets are equivalent to incomplete OCR results (Texts in an online image will sometimes appear in the captions). In the scale of our experiment, we observe similar improvement that just training on captions of text-rich images can help with text recognition capability: In Table 3, variant (4) is better than variant (1). However, training on captions only (variant (4)) is not as good as training on OCR-based data (variant (2)(6)), at least in the scale of our experiments.

Results of finetuning mPLUG-Owl To further validate the effectiveness of our collected data, we provide the results of finetuning mPLUG-Owl using our 16K GPT-4-based instruction-following data in Table 13. Though the mPLUG-Owl checkpoint is extensively trained on 1000M+ text-image pairs, we find that our data can boost performance in most cases, demonstrating the effectiveness of our data.

	Res.	METEOR	ROUGE-L	CIDEr
LLaVA	224 ²	3.8	4.8	6.3
LLaVAR		5.6	6.9	12.7
LLaVA	336 ²	4.6	5.6	8.7
LLaVAR		8.6	10.0	21.5

Table 11: Results on DocVQA using text-matching metrics.

	Res.	CT80	POIE	ChartQA
BLIP-2 [2023] †		80.9	2.5	7.2
OpenFlamingo [2023] †		67.7	2.1	9.1
MiniGPT4 [2023] †	224 ²	57.3	1.3	4.3
LLaVA [2023] †		61.1	2.1	7.3
mPLUG-Owl [2023] †		81.9	3.3	9.5
LLaVA ‡	224 ²	61.5	1.9	9.2
LLaVAR		81.6 (+20.1)	5.7(+3.8)	10.2 (+1.0)
LLaVA ‡	336 ²	64.9	2.5	10.2
LLaVAR		83.0 (+18.1)	8.7(+6.2)	13.5 (+3.3)

Table 12: Results (accuracy %) on three extra datasets: OCR, Information Extraction, and Chart Question Answering. We use † to refer to the results obtained from Liu et al. [10] and ‡ to refer to our reproduced results.

G

ScienceQA Results Starting from our pretrained LLaVAR (336²-based, without finetuning), we also report the results of further finetuning on the ScienceQA dataset [50] in Table 7, which is a multimodal multi-choice QA dataset covering diverse domains. Our motivation is that some images in this dataset contain text descriptions and tables that require textual understanding within images. The LLaVAR model finetuned on ScienceQA achieves an average accuracy of 91.42%, better than LLaVA (90.92%), while the most considerable improvement comes from natural science questions (+1.43%).

H

The original version of LLaVAR only supports up to 336² resolution, while our case study has also shown the threshold for the recognizable font size. Both suggest the difficulty of processing real-world high-res images without scaling and cutting. To this end, we test a dual visual encoder system for the high-res variant of LLaVAR, where a high-res visual encoder is added to work with the standard one. Ideally, the standard visual encoder extracts general, high-level information, while the high-res one specifically helps with detailed information.

Architecture A high-res visual encoder usually outputs thousands of visual features. Simply following LLaVA to feed the transformed visual features into the context of the language decoder is impractical, as the maximum sequence length of the language decoder is usually 2048/4096. To this end, we propose to handle high-res visual features by cross-attention module and standard visual features by feature transformation. We depict the proposed system in Figure 7.

Specifically, given a standard visual encoder V_1 , the extracted features are transformed into the word embedding space of the language decoder through a trainable projection matrix W . These transformed features are then concatenated with the word embeddings to build the input embeddings of the language decoder D .

$$\text{input_emb} = \text{concat}([\text{emb}(\langle \text{img}_1 \rangle), \dots, \text{emb}(\langle \text{img}_m \rangle), \text{emb}(\langle \text{ins}_1 \rangle), \dots, \text{emb}(\langle \text{ins}_n \rangle)]) \quad (1)$$

	ST-VQA	OCR-VQA	TextVQA	DocVQA	CT80	POIE	ChartQA
mPLUG-Owl	29.3	28.6	40.3	6.9	81.9	3.3	9.5
mPLUG-Owl _{ours}	29.6	31.2	40.8	7.0	84.7	3.7	10.2

Table 13: Results (accuracy %) of finetuning mPLUG-Owl. mPLUG-Owl_{ours} denotes mPLUG-Owl finetuned on our 16K GPT-4-based instruction-following data.

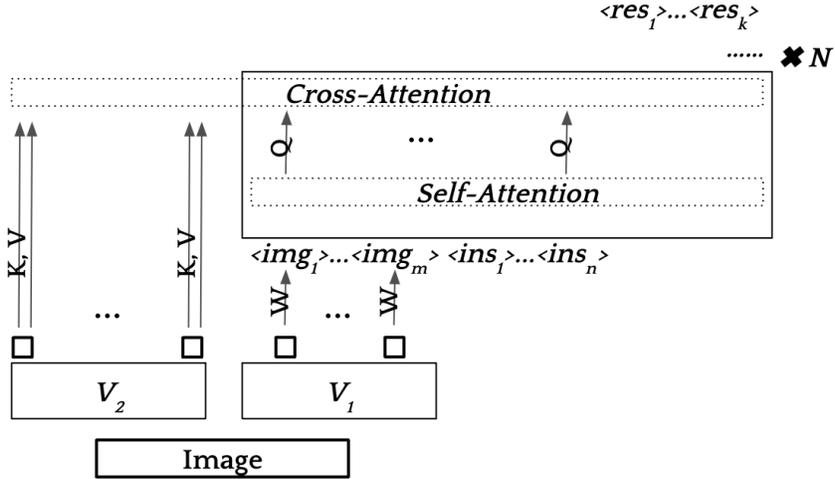


Figure 7: Illustration of the dual visual encoder system. Given an image, it is simultaneously processed by visual encoders V_1 and V_2 . V_1 features are transformed by transformation matrix W and directly used as input embeddings to the language model. For V_2 features, they are transformed by transformation matrix K and V and used as keys and values to calculate the cross attention in every transformer layer (assume there are N layers), which uses the transformed hidden states (through Q) from the self-attention module as queries. For the language decoder D , the input is image tokens ($\langle img \rangle$) and instruction tokens ($\langle ins \rangle$), while the target is response tokens ($\langle res \rangle$).

where I is the input image, V_1 denotes extracting the grid features before the last transformer layer. At the same time, we use the high-res visual encoder V_2 to extract high-res visual features, which are then transformed into keys/values as the inputs of the cross-attention module in transformer layers. Given h^j as the hidden state before the cross-attention module in layer j ,

$$\text{CrossAttention}(h, V_2, I) = \text{softmax}\left(\frac{Q^j h^j (K^j V_2(I))^T}{\sqrt{d}}\right) V^j V_2(I) \quad (2)$$

where Q^j, K^j, V^j denotes the query/key/value projection matrix in the j -th transformers layer. In practice, there is a pre-attention LayerNorm before calculating the attention and another output projection matrix O^j to project the aggregated values back to the hidden space.

As the pretrained language decoder D might only have self-attention modules, we manually add another cross-attention module after the original self-attention module in every transformer layer. Considering the random initialization of cross-attention modules might hurt the original language generation capability, we initialize the value projection matrix V^j as a zero matrix and the output projection matrix O^j as an identity matrix.

Implementation We use CLIP-ViT-L/14 as the standard visual encoder. For the high-resolution encoder, we test two models: (i) Pix2Struct-base [44] is a visual encoder trained on screenshot to HTML transformation. It supports up to 2048 patches with size 16^2 , equivalent to $1024 * 512$. (ii) ConcatCLIP refers to using 16 CLIP-ViT-L/14 models to encode the $4 * 4$ grids of images separately and then concatenate the extracted features together. In other words, it supports 896^2 resolution. We use Vicuna-7B as the language decoder for the high-res version of LLaVAR.

Training Only cross-attention modules and the projection matrix W are trained during pretraining, while visual encoders and the language decoder are frozen. Cross-attention modules, the projection matrix W , and the language decoder D are trained during finetuning.

Data To fully unlock the potential of the augmented visual encoder, we also double the number of pretraining examples using the same criteria mentioned in Section 3. This corresponds to the variant (g) in Table 4.

	ST-VQA	OCR-VQA	TextVQA	DocVQA
Pix2Struct + LLaVA	21.9	11.8	28.7	4.4
Pix2Struct + LLaVAR	35.8 (+13.9)	30.7 (+18.9)	45.6 (+16.9)	15.3 (+10.9)
ConcatCLIP + LLaVA	23.1	14.2	30.5	5.1
ConcatCLIP + LLaVAR	42.1 (+19.0)	30.8 (+16.8)	52.1 (+21.6)	18.5 (+13.4)

Table 14: Additional results on the dual visual encoder system.

Discussion We report the performance of augmented architecture, using either LLaVA or LLaVAR data in Table 14. By comparing the relative improvement in Table 2 and 14, we find that higher-resolution models benefit more from our collected data, suggesting our data is underutilized in the original LLaVA architecture.

I

	Res.	Correct %	Partially Correct %
LLaVA	224 ²	1.6%	8.7%
LLaVAR		6.8%	22.8%
LLaVA	336 ²	2.2%	11.2%
LLaVAR		9.0%	26.8%

Table 15: Statistics of correct answers and partially correct answers on OCR-VQA.

The impact of OCR spelling errors We study such OCR errors by studying 1673 examples from OCR-VQA, which have ground truth answers with more than ten characters. We (i) define “correct” as the ground truth answers that are exactly in the predictions, and (ii) define “partially correct” as there exists a substring in the prediction that has high enough similarity with the ground truth but not the same. Specifically, we look at all substrings with the same length of the ground truth in the prediction to calculate ANLS (Average Normalized Levenshtein Similarity) and regard the prediction as “partially correct” if the highest ANLS is greater or equal to 0.5 but smaller than 1.

We find that many predictions can be considered partially correct, indicating the actual performance of tested models is better than the reported accuracy numbers. However, the percentage of partially correct predictions is highly correlated with the percentage of correct predictions. Therefore, we believe that the current metrics can effectively compare the performance of different models.

J

Ablation Study on Instruction-following Evaluation To provide a more detailed understanding of LLaVAR’s data, we provide the ablation study results on 90 instruction-following questions from LLaVA [3]. We find that including pretraining data improves the conversation capability, probably because longer training data leads to generating longer responses (Table 1), which is favored by GPT-4 [19]. On the other hand, including only finetuning data hurts the conversation capability but increases complex reasoning. Combining pretraining and finetuning data improves the trade-off between conversation and complex reasoning. Such results prove that incorporating our data will at least not harm the performance of interacting with natural images.

	Conversation	Detail	Complex
LLaVA	83.6	78.1	95.2
LLaVA + R _{pretrain}	86.7	79.3	95.1
LLaVA + R _{finetune}	79.0	79.4	98.0
LLaVAR	84.4	78.9	96.6

Table 16: Ablation Study on GPT-4-based Instruction-following Evaluation from LLaVA [3].

	Question: Based on the title and the image on the cover, what can be inferred about the content of "Bo's Lasting Lessons" and its potential target audience?	Relative Score
GPT-4 (Text-Only)	From the title and image, it can be inferred that "Bo's Lasting Lessons" is a book about leadership lessons taught by a legendary coach, likely in the context of sports. The image of men running suggests that the book might be focused on a specific sport, such as football or track and field. The target audience could be individuals interested in coaching, sports, or leadership in general, looking to learn from the experiences and teachings of a successful coach.	-
LLaVA	The book "Bo's Lasting Lessons: The Timeless Teachings of Bo Schembechler" is written by Bo Schembechler and Jim Harshaw. Given the title and the image on the cover, which shows a football coach running across the field with his players in pursuit, it can be inferred that the content of the book is likely focused on football and the strategies, techniques, or principles that can help players improve their performance and achieve success in the sport. The potential target audience for this book could be football players, coaches, and fans who are interested in learning more about the game and the experiences and insights of a seasoned coach like Bo Schembechler.	75%
LLaVAR	"Bo's Lasting Lessons" appears to be a book that shares the legacy coach Bo Schembechler's teachings on leadership and timeless fundamentals. The image on the cover, featuring a football coach leading his team, further supports this inference. The book might be targeted towards sports enthusiasts, particularly fans of football, and people interested in learning about leadership principles and fundamentals from a renowned coach's perspective. The audiobook, read by John H. Mayer and featuring a foreword and conversation with John U. Bacon, suggests that it could also appeal to audiences who prefer listening to books or want to gain more insights from the author and other experts.	129%

Table 17: An example for our collected evaluation data, corresponding to the "Read" dimension of GPT-4-based instruction-following evaluation. The corresponding image is Figure 8 in Appendix.

BO'S LASTING LESSONS

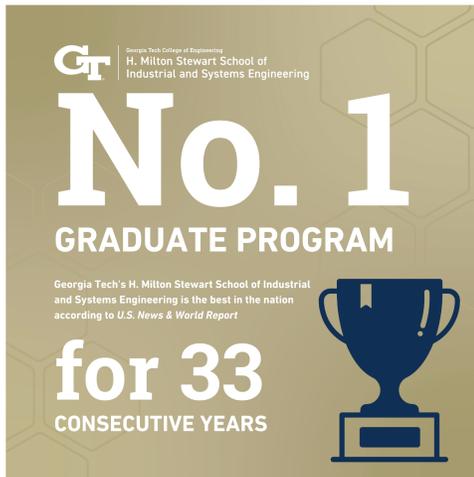
THE LEGENDARY COACH TEACHES THE TIMELESS FUNDAMENTALS OF LEADERSHIP



**BO SCHEMBECHLER
AND JOHN U. BACON**

READ BY JOHN H. MAYER
FEATURING A FOREWORD AND CONVERSATION WITH JOHN U. BACON

Figure 8: An example for the Read dimension of GPT-4-based instruction-following evaluation.



If I want to study optimization, is Georgia Tech a good place for me?



Yes, Georgia Tech is a good place for you if you want to study optimization. The image mentions that Georgia Tech's H. Milton Stewart School of Industrial and Systems Engineering is ranked as the best in the nation according to the U.S. News & World Report for 33 consecutive years. This indicates that the school has a strong reputation and offers quality education in the field of optimization and industrial engineering.



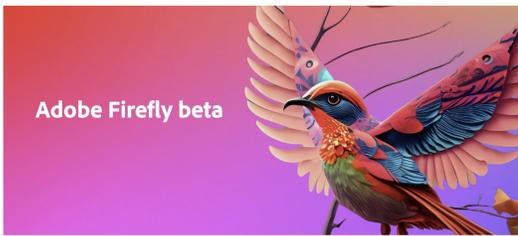
ADobe SUMMIT

Bringing Generative AI into Creative Cloud with Adobe Firefly



David Wedhwani
03-21-2023

in f @



Images generated using Adobe Firefly

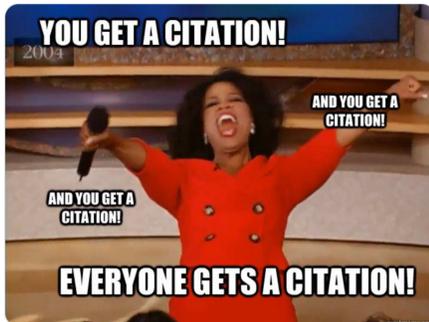


Jia-Bin Huang @jbhuang0604 · 1h
How to cite papers?

Citing papers properly

- gives credit where credit's due,
- provides supporting evidence of your claim, and
- presents an organized view of related work.

Sharing some tips I found useful. 📖



1 3 19 2,427

Pretend you are the author, write the next paragraph of this article.



Adobe Firefly is an exciting new tool that allows users to bring generative AI into their creative workflow. By integrating Firefly with Adobe's popular Creative Cloud, users can unlock new levels of creativity and automation, streamlining their design processes and enabling them to explore more innovative ideas. With Firefly, designers can explore data-driven design, real-time collaborative editing, and seamless integration with other Adobe products, making the entire creative process more efficient and intuitive.



What is this tweet about?



This tweet is about sharing tips on how to cite papers properly, giving credit where credit is due, providing support for evidence of your claim, and presenting an organized view of related work. It also humorously emphasizes the importance of citations by mentioning that "everyone gets a citation" in a playful manner.



By the way, who is the woman in this meme? Any background information?



The woman in the meme is Oprah Winfrey, a famous television host, author, and philanthropist. She is known for her influential book club and her ability to make people feel good about themselves. In this meme, she is seen enthusiastically endorsing the concept of citing sources in academic work.



Figure 9: Transferred instruction-following capability of LLaVAR.

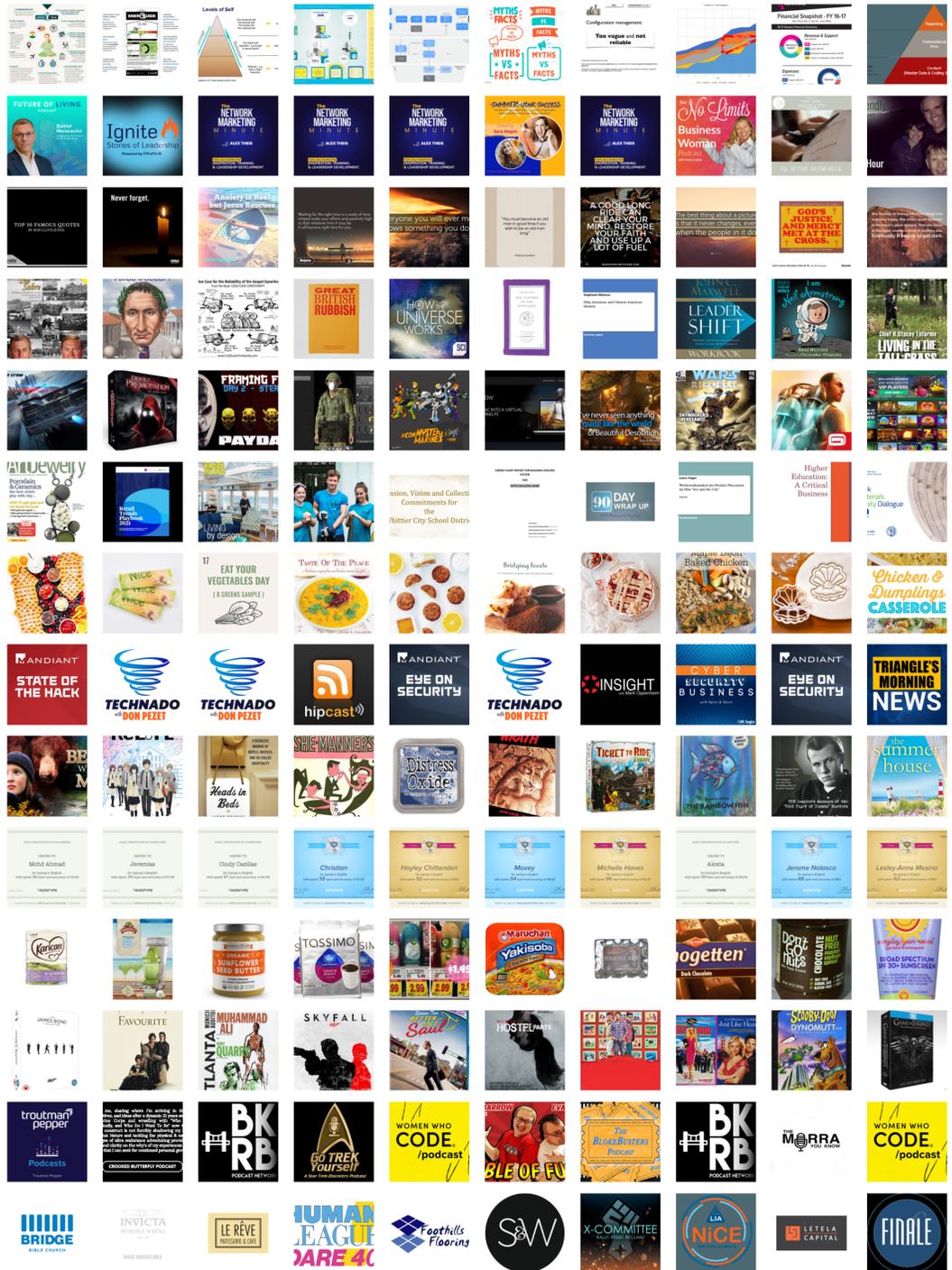


Figure 10: All 14 clusters we selected as text-rich images. Each row corresponds to one cluster, where we show ten randomly sampled examples before de-duplication.

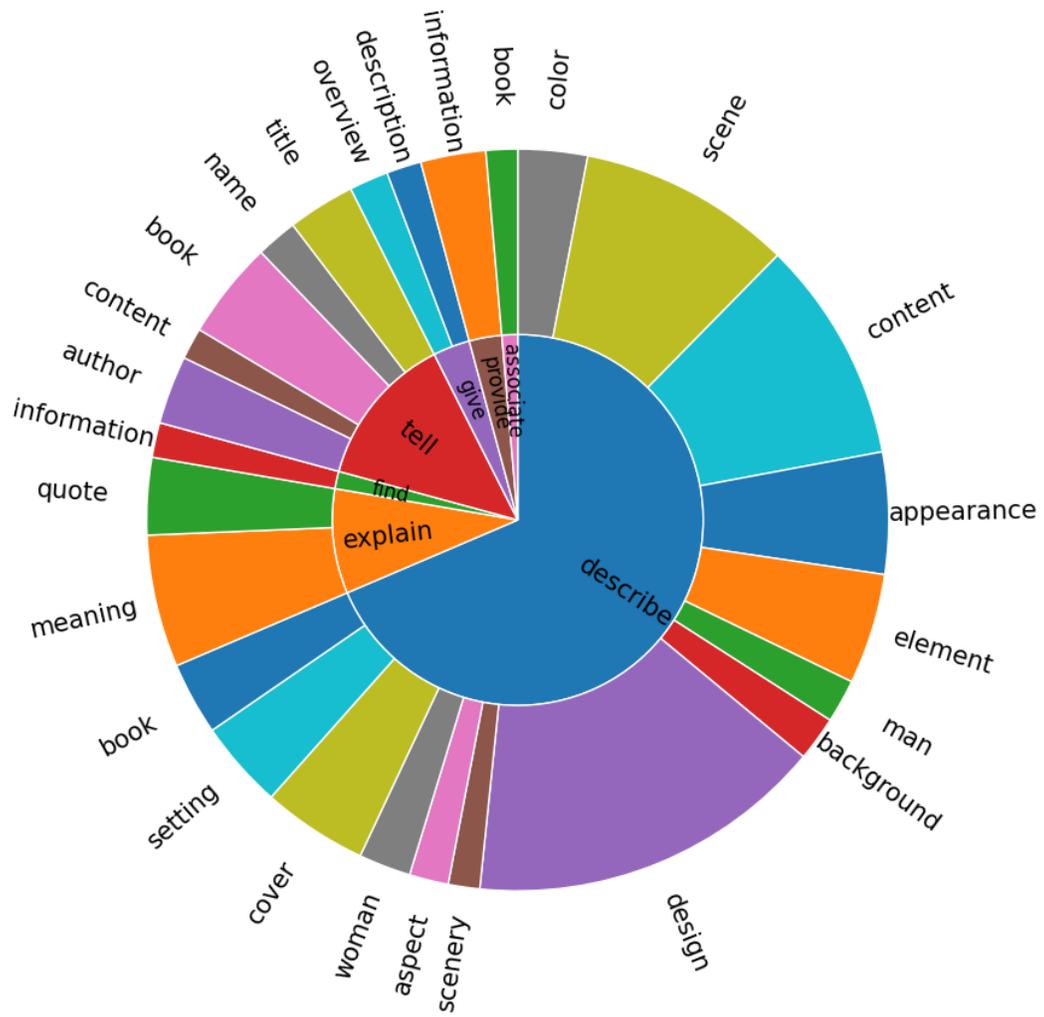


Figure 11: Visualization of collected instructions.

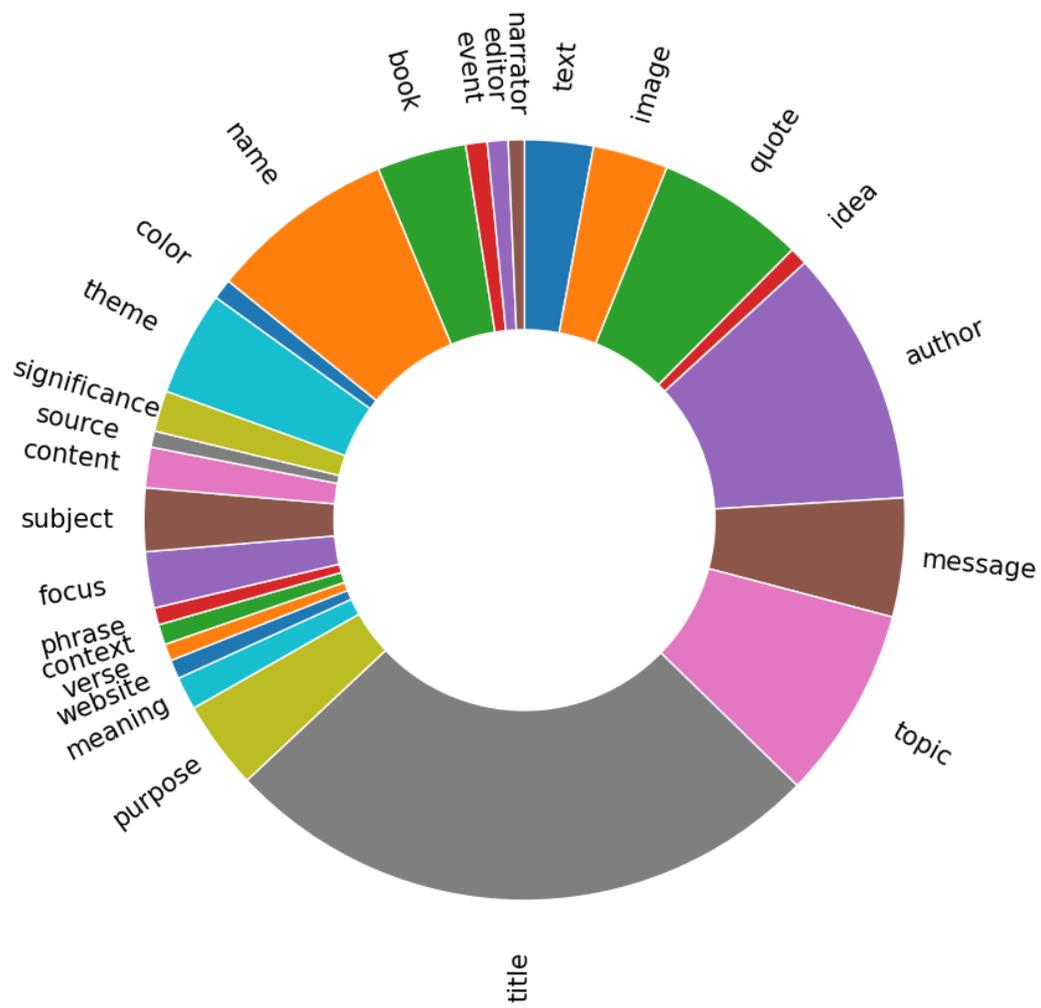


Figure 12: Visualization of collected instructions.