

MiniCPM-V: A GPT-4V Level MLLM on Your Phone

Yuan Yao^{*†} Tianyu Yu^{*} Ao Zhang^{*} Chongyi Wang^{*} Junbo Cui^{*} Hongji Zhu^{*}
 Tianchi Cai^{*} Haoyu Li^{*} Weilin Zhao^{*} Zhihui He^{*} Qianyu Chen^{*} Huarong Zhou^{*}
 Zhensheng Zou^{*} Haoye Zhang^{*} Shengding Hu^{*} Zhi Zheng^{*} Jie Zhou^{*} Jie Cai^{*}
 Xu Han^{*} Guoyang Zeng^{*} Dahai Li^{*} Zhiyuan Liu^{*} Maosong Sun[‡]

MiniCPM-V Team, OpenBMB
 yaoyuanthu@gmail.com

<https://github.com/OpenBMB/MiniCPM-V>

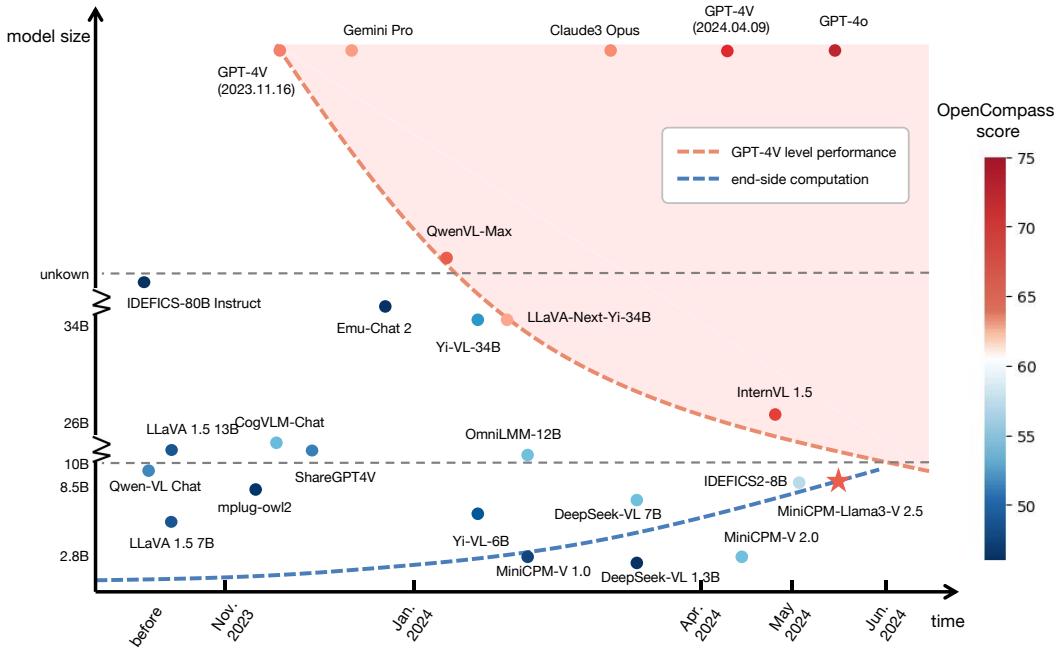


Figure 1: Moore’s Law for MLLM? Trends of MLLM development in terms of time (x-axis), model size (y-axis), and performance (color). The **red line** shows the decreasing model sizes for achieving GPT-4V level performance, while the **blue line** represents the growing end-side computation capacity. This jointly shows that GPT-4V level MLLMs deployed on end devices are becoming increasingly possible, unlocking a wider spectrum of real-world AI applications in the near future.

Abstract

The recent surge of Multimodal Large Language Models (MLLMs) has fundamentally reshaped the landscape of AI research and industry, shedding light on a promising path toward the next AI milestone. However, significant challenges remain preventing MLLMs from being practical in real-world applications. The most notable challenge comes from the huge cost of running an MLLM with a massive number of parameters and extensive computation. As a result, most MLLMs

* Core team, † Project lead, ‡ Corresponding author

need to be deployed on high-performing cloud servers, which greatly limits their application scopes such as mobile, offline, energy-sensitive, and privacy-protective scenarios. In this work, we present MiniCPM-V, a series of efficient MLLMs deployable on end-side devices. By integrating the latest MLLM techniques in architecture, pretraining and alignment, the latest MiniCPM-Llama3-V 2.5 has several notable features: (1) Strong performance, outperforming GPT-4V-1106, Gemini Pro and Claude 3 on OpenCompass, a comprehensive evaluation over 11 popular benchmarks, (2) strong OCR capability and 1.8M pixel high-resolution image perception at any aspect ratio, (3) trustworthy behavior with low hallucination rates, (4) multilingual support for 30+ languages, and (5) efficient deployment on mobile phones. More importantly, MiniCPM-V can be viewed as a representative example of a promising trend (Fig. 1): The model sizes for achieving usable (e.g., GPT-4V) level performance are rapidly decreasing, along with the fast growth of end-side computation capacity. This jointly shows that GPT-4V level MLLMs deployed on end devices are becoming increasingly possible, unlocking a wider spectrum of real-world AI applications in the near future.

1 Introduction

The rapid development of Multimodal Large Language Models (MLLMs) [83, 2, 68, 62, 102, 23, 7, 41, 52, 76, 11] have brought an impressive surge in multimodal capabilities in understanding, reasoning and interaction. This has not only fundamentally reshaped the landscape of AI research and industry, but also shed light on a promising path towards the next AI milestone. However, current MLLMs are still far from being practical in real-world applications. One of the most predominant challenges is that current MLLMs typically entail a massive number of parameters and impose heavy computational burdens. As a result, most MLLMs can only be deployed on high-performing cloud servers, leading to significant energy consumption and carbon emissions. This limitation significantly constrains the potential application scopes such as on mobile devices, energy-sensitive scenarios, offline scenarios without stable network connections, and privacy/security protective scenarios for both personal and industrial users.

In light of these limitations, there is a growing interest in exploring more efficient lightweight MLLMs [1, 68, 83, 11] that can run on end-side devices. End-side scenarios encompass a broader scope of equipment, including mobile phones, personal computers, vehicles and robotics, etc., which are ubiquitous in users' daily lives and experiencing rapid advancements in computation capacities. End-side MLLMs provide a promising solution towards more practical applications due to their broader usage scope, better computation efficiency, more robust offline behaviors, and better privacy/security protection.

However, developing capable end-side MLLMs is challenging due to significantly constrained parameter and inference computation budgets. As a result, more careful architecture designs and training recipes are required to fully unleash the potential of end-side MLLMs. In this work, we present MiniCPM-V, a series of efficient MLLMs deployable on end-side devices. The philosophy of MiniCPM-V is to achieve a good balance between performance and efficiency, a more important objective in real-world applications. To date in 2024, we have unveiled three models: (1) In February, we launched MiniCPM-V 1.0 2B, one of the first MLLMs designed for mobile phones. (2) In April, MiniCPM-V 2.0 2B was introduced, outperforming strong larger MLLMs such as Qwen-VL 9B [7], CogVLM 17B [102], and Yi-VL 34B [108]. This iteration also introduces support for high-resolution image input and exhibits promising OCR capabilities. (3) Most recently in May, we released MiniCPM-Llama3-V 2.5 8B, which outperforms strong GPT-4V-1106, Gemini Pro and Claude 3 on OpenCompass evaluation. Noteworthy features of this model include strong OCR capability, high-resolution image perception, trustworthy behavior, multilingual support, and efficient end-side deployment optimization.

More importantly, MiniCPM-V can be viewed as a representative example of a promising trend. Fig. 1 summarizes the recent development of MLLMs [1, 68, 57] in terms of performance, parameters and release time. We observe an interesting trend akin to Moore's Law [103] indicated by the red line: the sizes of models reaching the GPT-4V level performance are rapidly decreasing over time. This phenomenon could perhaps be called the *Moore's Law of MLLMs*. Simultaneously, the computational capacity of end-side devices such as phones and personal computers is steadily

increasing (qualitatively depicted by the blue line). The convergence of these two trends indicates usable (e.g., GPT-4V level) MLLMs deployable on end-side devices are soon within reach, opening up broader possibilities and benefiting more application scenarios in the near future. From a historical perspective of human technology development, this trend can also be viewed as human pursuit of miniaturization of state-of-the-art technologies, which have been repeatedly witnessed in other science and technology fields. For example, in aerospace, the latest SpaceX Raptor 2 rocket engine can achieve a strong thrust of 2,256 kN with a mass of 1.6 tons, whereas 20 years ago, the RD-0750 rocket engine could only achieve a thrust of 1,413 kN with a mass exceeding 4 tons [104].

MiniCPM-V Series Techniques. In this paper, we will take the latest MiniCPM-Llama3-V 2.5 as an example, and systematically introduce the notable features of MiniCPM-V series and the key techniques behind them:

- **Leading Performance.** MiniCPM-Llama3-V 2.5 achieves better performance than GPT-4V-1106, Gemini Pro and Claude 3 on OpenCompass collection, a comprehensive evaluation over 11 popular benchmarks. This is jointly contributed by its careful design in architecture, data and training recipes, which we will detail in the following.
- **Strong OCR Capability.** MiniCPM-Llama3-V 2.5 outperforms GPT-4V, Gemini Pro and Qwen-VL-Max on OCRBench. It also supports high-utility functions such as table-to-markdown conversion and full OCR content transcription. These are largely attributed to the 1.8M pixel high-resolution (e.g., 1344×1344) image perception technique across any aspect ratios [107].
- **Trustworthy Behavior.** Based on the RLAIF-V [112] and RLHF-V [111] techniques that align MLLM behaviors from AI/human feedback, MiniCPM-Llama3-V 2.5 exhibits more trustworthy behaviors, achieving lower hallucination rates than GPT-4V-1106 on Object HalBench.
- **Multilingual Support.** Inspired by the findings from VisCPM [41], the integration of multilingual LLM significantly alleviates the heavy reliance on multimodal training data in low-resource languages. Based on the foundation, a high-quality multilingual multimodal instruction tuning helps MiniCPM-Llama3-V 2.5 generalize its multimodal capabilities to more than 30 languages.
- **Efficient End-side Deployment.** We systematically integrate a suite of end-side optimization techniques, encompassing quantization, memory optimization, compilation optimization and NPU acceleration, enabling efficient deployment on end-side devices.

We hope MiniCPM-V series can serve as an example for unveiling the potential of end-side MLLMs, and help draw more attention to promote the research in this direction. Following Moore’s Law for MLLM, we believe there will be increasingly powerful end-side MLLMs with reduced sizes, bringing efficient, safe, and trustworthy AI services on devices soon.

The contribution of this work is summarized as follows: (1) We introduce and open-source MiniCPM-V, a series of efficient end-side MLLMs achieving a good balance between performance and efficiency. (2) We investigate key techniques driving MLLMs towards the performance-efficiency balance at scale, unveiling the potential of these techniques. (3) We summarize the trend of MLLM development in its Moore’s Law, and empirically instantiate the trend with representative examples of MiniCPM-V.

2 Related Works

Multimodal Large Language Models. The development of LLMs has significantly advanced the progress in MLLMs. Flamingo [4] first proposes to connect a pre-trained visual encoder with the Chinchilla 70B [40] LLM and demonstrate the MLLM’s zero-shot and few-shot ability across a series visual language tasks. After the appearance of ChatGPT, many open-source models including BLIP-2 [54], Kosmos-1 [43], MiniGPT-4 [121], LLaVA [62], and VPGTrans [117] are proposed. Among them, most are built upon existing pre-trained LLMs like Llama [100] and Vicuna [120], while Kosmos-1 tries to train the LLM from scratch. Later, researchers continue to extend the function scope of MLLMs and improve the visual perception capabilities. Kosmos-2 [80], CogVLM [102], Shikra [20], and NExT-Chat [116] further incorporate the localization capabilities to the MLLMs with either pix2seq paradigm or connecting with detection/segmentation models. Qwen-VL-Chat [7], Yi-VL [108], DeepSeek-VL [68], InternVL [23] and Intern-XComposer [31] pay more attention to improving the models’ capability with different techniques like high-resolution input, more training data, and better data ratio.

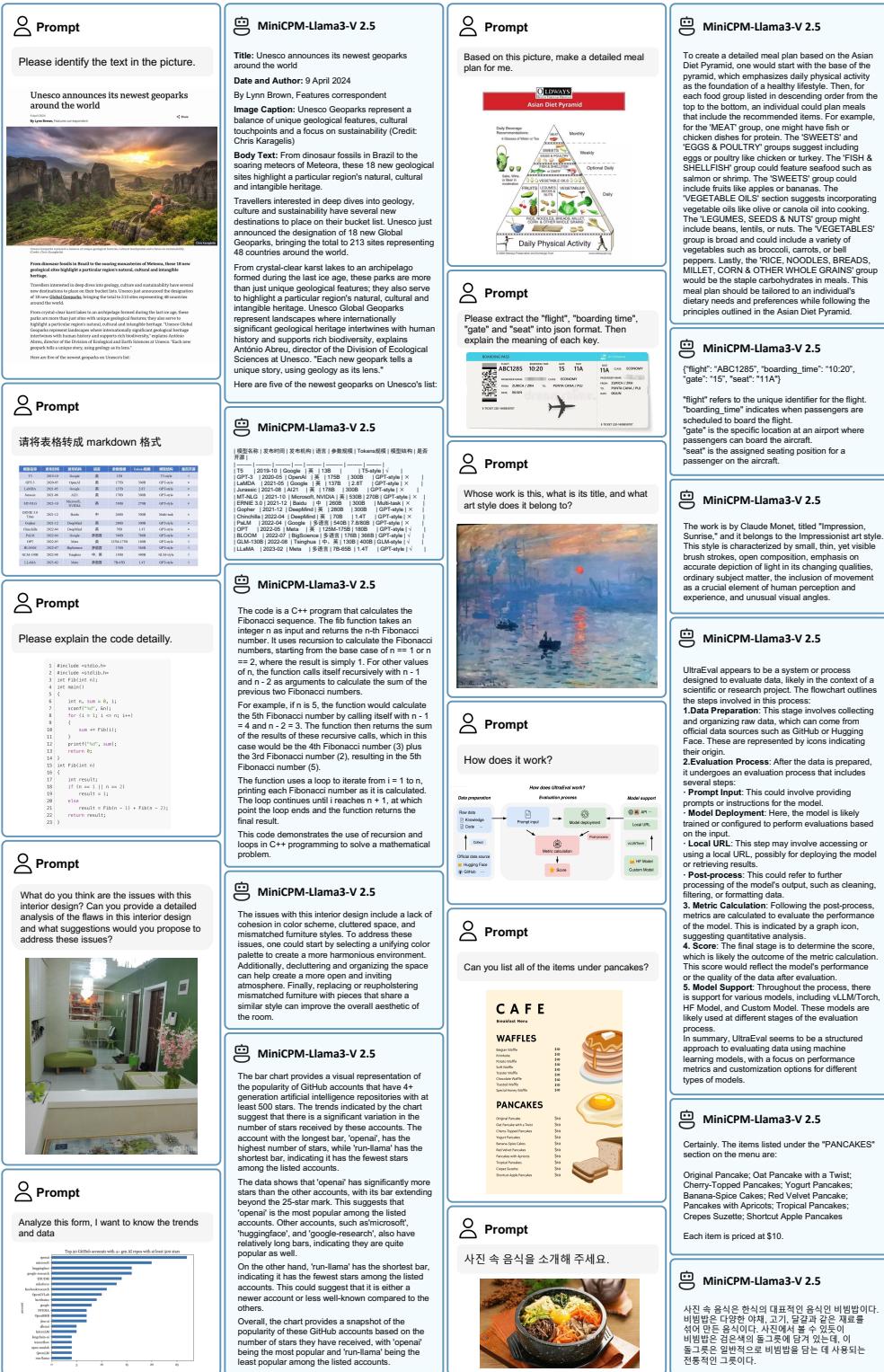


Figure 2: Qualitative results of MiniCPM-Llama3-V 2.5 in reading text in images, converting tables to markdown, performing complex reasoning, and multilingual interaction, etc. See here for the screen recordings of MiniCPM-Llama3-V 2.5 running on mobile phones.

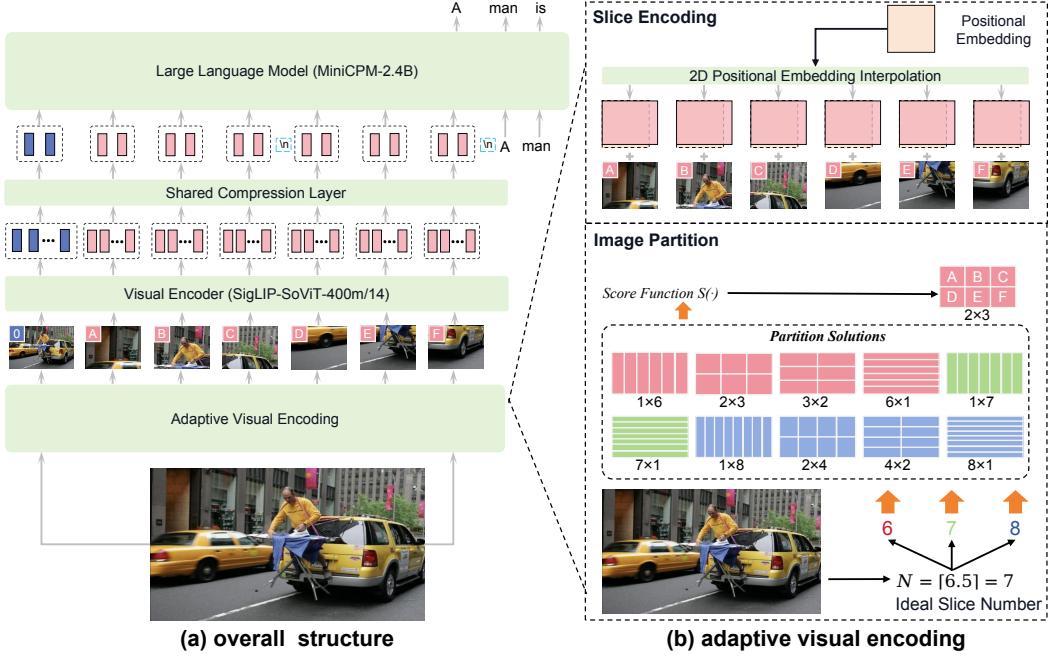


Figure 3: Model architecture. (a) **Overall structure** presents the architecture of the model including the visual encoder, shared compression layer, and LLM. (b) **Adaptive visual encoding** deals with high-resolution inputs in various aspect ratios.

End-side Multimodal Large Language Models. The huge number of parameters of MLLMs incurs prohibitively high computation costs in both training and deployment, greatly limiting the widespread applications. Recently, there has been a trend of building smaller LLMs with fewer parameters. The representative models are Phi [45], Gemma [8], MobileLLM [66], MiniCPM [42], etc. The moderate size of these models makes them applicable on end-side devices such as personal computers and even mobile phones. With optimized training strategies, end-side LLMs like MiniCPM 2B can achieve comparable performance with strong 7B models like Llama2-7B [101]. Similar trends have also been witnessed in MLLMs. For example, Mini-Gemini [57] and PaliGemma [11] are built based on Gemma 2B [8] and MobileVLM V2 is built based on MobileLlama [25]. However, the fewer-parameter nature of end-side MLLMs presents significant challenges to building a capable model. MiniCPM-V series aims to push forward the potential of end-side MLLMs by addressing the key bottleneck problems through careful designs in architecture, training, inference and deployment.

3 Model Architecture

In this section, we present the model architecture of MiniCPM-V, outlining the overall structure and the adaptive high-resolution visual encoding approach. The design philosophy of MiniCPM-V series is to achieve a good balance between performance and efficiency, a more practical objective for a broader scope of real-world applications, which is implemented in architecture design, training, inference, and deployment.

3.1 Overall Structure

The model comprises three key modules: the visual encoder, compression layer, and LLM. The input image is first encoded by a visual encoder, utilizing the adaptive visual encoding approach. Specifically, we employ SigLIP SoViT-400m/14 [115] as the visual encoder. The visual tokens are then compressed by the compression layer, which adopts a perceiver resampler structure with one layer cross-attention. Finally, the compressed visual tokens, along with the text input, are fed into the LLM for conditional text generation.

Dùng phiên bản SigLIP SoViT-400m/14 làm Visual Encoder, 1 thành phần quan trọng mà mình không hiểu là paper có ko nói k là compression layer có khống nén compress token nhg i mà trnhi u ki ut 384x16 chn 32 token ghi y

3.2 Adaptive Visual Encoding

Recently, there has been growing consensus on the fundamental role of visual encoding in MLLM performance [76, 68], especially for fine-grained capabilities such as OCR. For effectiveness, a good visual encoding strategy should both respect the raw aspect ratio of the input and preserve sufficient visual details (high resolution). For efficiency, the number of visual tokens from image encoding should be moderate to be affordable on end-side devices. To this end, we take advantage of the adaptive visual encoding method proposed by LLaVA-UHD [107].

Image Partition. To handle the high-resolution images with different aspect ratios, we divide images into slices, where each slice better matches ViT’s pre-training setting in terms of resolution and aspect ratio. Specifically, we first calculate the ideal number of slices based on the input image size. Given an image with resolution (W_I, H_I) and a ViT pre-trained on images with resolution (W_v, H_v) , we calculate the ideal slice number $N = \lceil \frac{W_I \times H_I}{W_v \times H_v} \rceil$. Then, we choose the combination of rows n and columns m from the set $\mathbb{C}_N = \{(m, n) | m \times n = N, m \in \mathbb{N}, n \in \mathbb{N}\}$. A good partition (m, n) should result in slices that match well with ViT’s pre-training setting. To achieve this, we use a score function to evaluate each potential partition:

$$S(m, n) = - \left| \log \frac{W_I/m}{H_I/n} - \log \frac{W_v}{H_v} \right|. \quad (1)$$

We select the partition with the highest score from all possible candidates:

$$m^*, n^* = \arg \max_{(m, n) \in \bar{\mathbb{C}}} S(m, n), \quad (2)$$

where $\bar{\mathbb{C}}$ is the possible (m, n) combinations with the product N . However, when N is a prime number, the feasible solutions can be limited to $(N, 1)$ and $(1, N)$. Therefore, we additionally introduce \mathbb{C}_{N-1} and \mathbb{C}_{N+1} , and set $\bar{\mathbb{C}} = \mathbb{C}_{N-1} \cup \mathbb{C}_N \cup \mathbb{C}_{N+1}$. In practice, we set $N < 10$, supporting 1.8 million pixels (e.g., 1344×1344 resolution) at most during encoding. Although we can encompass more image slices for higher resolutions, we purposely impose this resolution upper-bound, since it already well covers most real-world application scenarios, and the benefit of further increasing encoding resolution is marginal considering the performance and overhead.

Slice Encoding. Although image partitioning can ensure a good match between the slices and the ViT pre-training setting, each slice’s size is not precisely equal to (W_v, H_v) . To feed the slices into ViT, we first adjust each slice by resizing it proportionally so that the resultant area size matches ViT pre-training area size $W_v \times H_v$. This adjustment helps prevent a significant gap between the number of encoded patches and the ViT’s pre-training setting. Subsequently, we interpolate the ViT’s position embeddings to adapt to the slice’s ratio. This involves reshaping the ViT’s 1D embedding $P_1 \in \mathbb{R}^{Q \times l}$ back to its 2D format $P_2 \in \mathbb{R}^{q \times q \times l}$, where the number of position embeddings $Q = q \times q$. Then, we interpolate P_2 to fit the size of each slice via 2D interpolation. We also include the original image as an additional slice to provide holistic information about the entire image.

Token Compression. After visual encoding, each slice is encoded into 1,024 tokens, where 10 slices can yield over 10k tokens collectively. To manage this high token count, we employ a compression module comprising of one-layer cross-attention and a moderate number of queries, with 2D positions informed [7]. In practice, the visual tokens of each slice are compressed into 64 queries for MiniCPM V1&2 and 96 tokens for MiniCPM-Llama3-V 2.5 through this layer. Compared with other MLLMs with competitive performance, the significantly smaller number of visual tokens in MiniCPM-V series enables superior efficiency in terms of GPU memory consumption, inference speed, first-token latency and power consumption, making it more friendly to wider application scopes and communities.

Spatial Schema. To indicate each slice’s position relative to the whole image, inspired by [9], we additionally introduce a spatial schema. We first wrap tokens of each slice by two special tokens `<slice>` and `<\slice>`, and then employ a special token “\n” to separate slices from different rows.

Image partition khá n g i n c b n là nh s cchia thành nh i slices m i slice thi ph i chia sao cho trùng setting v ch t l ng nh và vát i trong mô hình ViT

4 Training

The model training consists of 3 phases: the pre-training phase, the supervised fine-tuning phase, and the RLAIF-V phase. We will introduce the training recipe in the following sections.

o nà y t clành nà y, theo công th c cross- attention, trong mô hình s cós n l ng query vector mà h c c t ng giá tr nhán v it ng giá tr vector key c a nh r tinh softmax và tính weightsum dùng giá tr l y t softmax nhán v it ng giá tr value c a nh Output là t ms l ng vector b ng v i query vector luôn nên m i nói là n g i m i s token c a nh i r t nh i u

4.1 Pre-training

In this phase, we utilize large-scale image-text pairs for MLLM pre-training. The primary goal of this phase is to align the visual modules (i.e., visual encoder and compression layer) with the input space of the LLM and learn foundational multimodal knowledge. The pre-training phase is further divided into 3 stages.

Stage-1. The role of stage-1 is to warm up the compression layer, primarily connecting the visual encoder and LLMs. (1) Trainable Modules. We randomly initialize the compression layer and train this module in stage-1, keeping other parameters frozen. The visual encoder’s resolution is set to 224×224 , which is the same as the visual encoder’s pre-training setting. (2) Data. To warm up the compression layer, we randomly select 200M data from the Image Captioning data in Table 1. Data cleaning is performed to remove image-text pairs with poor correlation and ill-formatted text data, ensuring the data quality.

Table 1: Pre-training data. The pre-training data consists of image captioning and OCR data in English and Chinese. LAION-2B-OCR is generated by applying OCR tools to LAION-2B images.

Category		Sources	Size
Image Captioning	English	COCO [59], VG [51], CC3M [89], CC12M [17] LAION-COCO [86], COYO [15], LAION-2B [86]	410M
	Chinese	AIC [105], LAION-2B-Chinese [86], WuKong [35] Zero-Chinese [106], etc.	110M
OCR+Knowledge	English	WIT [92], IDL [13], SynthText [37], SynthDoG-en [50] SynthDoG-zh [50], ArxivCap [55], etc.	39M
	Chinese	WIT [92], LAION-2B-OCR	11M

Stage-2. After the warm-up training of the compression layer, the role of stage-2 is to extend the input resolution of the pre-trained visual encoder. (1) Trainable Modules. In stage-2, we extend the image resolution from 224×224 to 448×448 . The whole visual encoder is trained, leaving other parameters frozen. (2) Data. To extend the pre-trained resolution, we additionally select 200M data from the Image Captioning data in Table 1.

Stage-3. After extending the primary input resolution of the visual encoder, we finally train the visual modules using the adaptive visual encoding strategy, which can further accommodate high-resolution inputs with any aspect ratio. (1) Trainable Modules. During the stage-3 training, both the compression layer and the visual encoder are trained to adapt to the language model embedding space. The LLM is kept frozen to avoid disruption from the relatively low-quality pre-training data. (2) Data. Different from the previous stages with only image captioning data, during the high-resolution pre-training stage, we additionally introduce OCR data to enhance the visual encoders’ OCR capability.

Caption Rewriting. Image-text pairs sourced from the Web [86, 15] can suffer from quality issues in the caption data, including non-fluent content, grammatical errors, and duplicated words. Such low-quality data can lead to unstable training dynamics. To address the issue, we introduce an auxiliary model for low-quality caption rewriting. The rewriting model takes the raw caption as input and is asked to convert it into a question-answer pair. The answer from this process is adopted as the updated caption. In practice, we leverage GPT-4 [14] to annotate a small number of seed samples, which are then used to fine-tune an LLM for the rewriting task.

Data Packing. Samples from different data sources usually have different lengths. The high variance of sample lengths across batches will lead to inefficiency in memory usage and the risk of out-of-memory (OOM) errors. To address the issue, we pack multiple samples into a single sequence with a fixed length. By truncating the last sample in the sequence, we ensure uniformity in sequence lengths, facilitating more consistent memory consumption and computational efficiency. Meanwhile, we modify the position ids and attention masks to avoid interference between different samples. In our experiments, the data packing strategy can bring 2~3 times acceleration in the pre-training phase.

Table 2: SFT data for MiniCPM-V series. Part-1&2 data are concatenated sequentially in the SFT phase. Part-1 focuses on bolstering basic recognition capabilities, while part-2 aims to enhance advanced capabilities in generating detailed responses and following human instructions.

Category	Sources	Size
Part-1	Short Caption Flickr-30K [81], COCO [59]	560K
	VQA FM-IQA [34], VGQA [51], IconQA [69], GQA [44], VQAv2 [6] CLEVR [46], VizWiz [38], Visual7W [122], COCO-QA [84]	1.4M
	Knowledge OKVQA [72], A-OKVQA [87], KVQA [88], ScienceQA [70]	60K
	Grounding RefCOCO [109]	570K
	Reasoning COMVINT [32], VCR [114], NLVR [94], LRV [60]	135K
	Math GeoQA [19], SMART-101 [24]	125K
	OCR DocVQA [74], TextVQA [91], OCR-VQA [77], ST-VQA [12], VisualMRC [96], DVQA [47] FigureQA [48], ChartQA [73], DeepForm [95], TabFact [22], InfographicsVQA [75] Kleister Charity [93], WikiTableQuestions [79], Real-CQA [3], AI2D [49], etc.	1.7M
Part-2	Chat FSVQA [90], Visual-Dialog [28]	780K
	Part-1 sample from Part-1 data	400K
	OCR DocVQA, TextVQA, OCR-VQA, VisualMRC, ChartQA, AI2D ArxivQA [56], LLaVAR [118], TextOCR-GPT4V [16], etc.	690K
	Instruct SVIT [119], LLaVA-Instruct-150K [62], UniMM-Chat [110], ShareGPT4V [21] LVIS [36], ALLaVA [18]	1.9M
Text-Only	Ultra-Chat [30], Alpaca [97], ShareGPT [120], BELLE [10]	-
	OpenOrca [58], OpenHermes [98], In-House-MiniCPM-SFT	-

Multilingual Generalization. Multimodal capability across multiple languages is essential for serving users from broader communities. Traditional solutions involve extensive multimodal data collection and cleaning, and training for the target languages. Fortunately, recent findings from VisCPM [41] have shown that the multimodal capabilities can be efficiently generalized across languages via a strong multilingual LLM pivot. This solution largely alleviates the heavy reliance on multimodal data in low-resource languages. In practice, we only pre-train our model on English and Chinese multimodal data, and then perform a lightweight but high-quality multilingual supervised fine-tuning to align to the target languages. Despite its simplicity, we find the resultant MiniCPM-Llama3-V 2.5 can achieve good performance in over 30 languages as compared with significantly larger MLLMs.

4.2 Supervised Fine-tuning

After learning foundational capabilities from pre-training, we perform supervised fine-tuning (SFT) on high-quality visual question answering datasets to further learn knowledge and interaction capability from human annotations.

Trainable Modules. Compared with the pre-training phase which mainly uses crawled data from the Web, the SFT phase mainly utilizes high-quality datasets annotated by either human lablers or strong models such as GPT-4. Therefore, we unlock all model parameters to better exploit the data and learn rich knowledge during SFT phase.

Data. Recent works [42, 83] show that data near the end of training plays a more important role in shaping the models’ capabilities and response styles. We categorize the SFT data into two parts. Part-1 focuses on bolstering the models’ basic recognition capabilities, while part-2 is tailored to enhance their capabilities in generating detailed responses and following human instructions. Specifically, part-1 data consists of the traditional QA/captioning datasets with relatively short response lengths, which helps enhance the model’s basic recognition capabilities. In comparison, part-2 encompasses datasets featuring long responses with complex interactions, either in text or multimodal context. During SFT, these two parts of data are concatenated and sequentially fed into the model. For

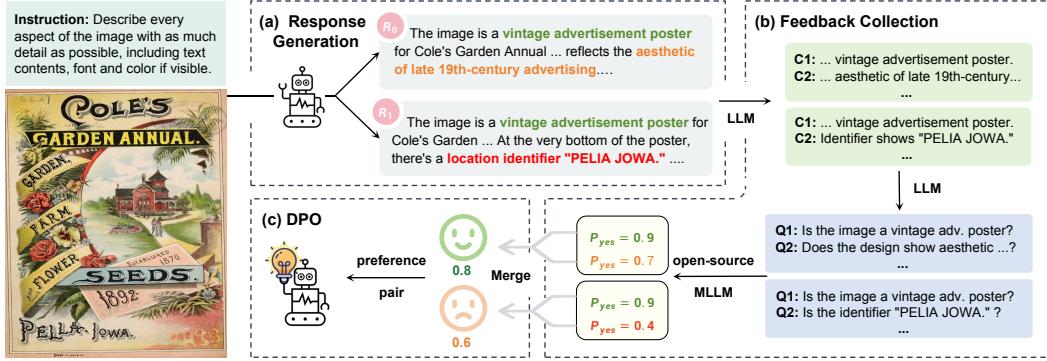


Figure 4: **RLAIF-V** framework for hallucination reduction. (1) **Response generation** produces multiple responses for an instruction using the policy model. (2) **Feedback collection** evaluates the correctness of each response in a divide-and-conquer fashion. (3) **DPO** optimizes the model on the preference dataset.

MiniCPM-Llama3-V 2.5, we integrate 2M data from the recent Cauldron dataset [52] for multimodal knowledge augmentation, and 90K multilingual data over 36 languages for boosting the multilingual conversation capability.

4.3 RLAIF-V

MLLMs are typically prone to hallucination problems, generating responses that are not factually grounded in the input image [111]. The issue greatly limits the wide application of MLLMs, especially in high-stakes scenarios, such as autonomous driving and assistance for visually impaired groups. To address the hallucination problem, we employ the recent RLAIF-V [112] approach (Fig. 4), where the key is to obtain scalable high-quality feedback from open-source models for direct preference optimization (DPO) [82].

Response Generation. The first step of RLAIF-V is to generate multiple responses for a given instruction using the policy model. Specifically, given a model M waiting for alignment, we sample 10 responses $Y = \{y_1, y_2, \dots, y_n\}$ from M using sampling decoding with high temperatures. There are several benefits of using the policy model M for response generation: (1) Feedback collection and learning can better focus on trustworthiness, since different text styles from multiple MLLMs are avoided. (2) Feedback learning is more efficient since preference is directly collected on the distribution of the policy model.

Feedback Collection. Collecting high-quality feedback from open-source MLLMs can be challenging due to their typically weaker capabilities compared with proprietary models. To address the issue, RLAIF-V uses a divide-and-conquer strategy for response scoring. Specifically, each response y_i is divided into atomic claims $C_i = \{c_1, c_2, \dots, c_m\}$ using Llama-3 8B, where the correctness of atomic claims is much easier to evaluate. Then, we verify the claims by converting each claim to a yes/no question and employing an open-source MLLM to score each claim. In practice, we adopt OmniLMM 12B for MiniCPM-V 2.0 scoring and LLaVA-NeXT-Yi 34B for MiniCPM-Llama3-V 2.5 scoring. The final score s_i of the response y_i is given by $-n_{rej}$, where n_{rej} is the number of invalid atomic claims.

Direct Preference Optimization. After collecting the high-quality AI feedback, we perform preference learning via DPO method. The DPO algorithm requires training on preference pairs, where one sample y_w is preferred to the other one y_l . To compose the preference dataset, we randomly sample pairs from each response set $Y = \{y_1, y_2, \dots, y_n\}$, and determine (y_w, y_l) based on their relative scores. Finally, we construct a preference dataset consisting of 6K preference pairs from 3K unique images for preference learning.

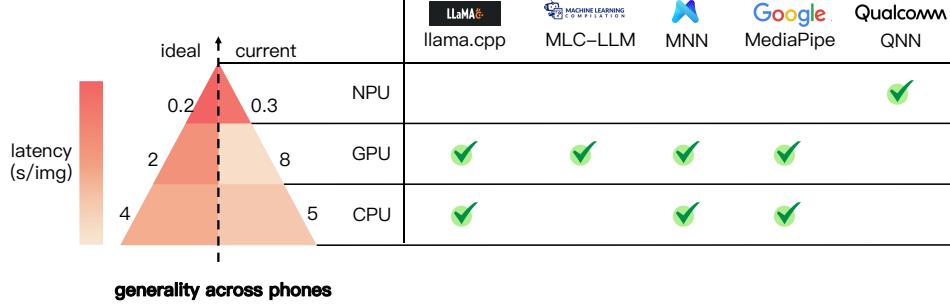


Figure 5: An overview of end-side deployment. Current end-side chips for computation acceleration include CPU, GPU and NPU. We list their generality, (estimated) ideal and current performance, and related deployment frameworks.

5 End-side Deployment

In this section, we investigate the deployment of MiniCPM-V on end-side devices. We first introduce the challenges, and then present the basic and advanced practices for end-side deployment. Finally, we analyze and discuss the evaluation results across different devices.

5.1 Challenges

End-side devices, such as smartphones and computers, often face resource limitations due to factors like heat dissipation, size constraints, and power consumption. We identify several key challenges of end-side deployment for MLLMs by comparing end-side devices with high-performance servers:

Memory Constraints. High-performance servers typically boast extensive memory capacities, often exceeding 100GB or even 1TB. In contrast, the memory available on mobile phones typically ranges from 12GB to 16GB, which can be insufficient for MLLM deployment.

CPU/GPU Speed Restriction. The overall processing speeds of CPUs in smartphones are notably slower. For instance, the Snapdragon 8 Gen3 features 8 CPU cores², whereas high-performance server like Intel Xeon Platinum 8580 has 60 CPU cores³. Similarly, mobile phone GPUs are not as powerful as server GPUs. For example, Qualcomm Adreno 750 only has 6 TFLOPS, while NVIDIA 4090 can reach 83 TFLOPS.

5.2 Basic Practice

To deploy the MLLM on end-side devices, we first employ quantization for reduced memory cost, and empirically investigate the deployment results on different frameworks.

Quantization. Quantization is a widely used technique to reduce memory consumption. The main idea of model quantization is to use a unified scaling factor to compress multiple weights into a narrower range, followed by discretization. This process is mathematically represented as:

$$w'_i = \text{round}\left(\frac{w_i}{s}\right), \forall 1 \leq i \leq n, \quad (3)$$

where w' denotes the quantized parameter and s signifies the calculated scale factor. The round function discretizes the quantized value.

For MiniCPM-Llama3-V 2.5, the fp16 version model typically demands 16~17G memory. We opt for the Q4_K_M mode 4-bit quantization strategy within GGML⁴ framework. This reduces the memory requirement to around 5G, which is friendly to mobile phone usage.

²https://docs.qualcomm.com/bundle/publicresource/87-71408-1_REV_E_Snapdragon_8_gen_3_Mobile_Platform_Product_Brief.pdf

³<https://www.intel.com/content/www/us/en/products/sku/237250/intel-xeon-platinum-8580-processor-300m-cache-2-00-ghz/specifications.html>

⁴<https://github.com/ggerganov/ggml>

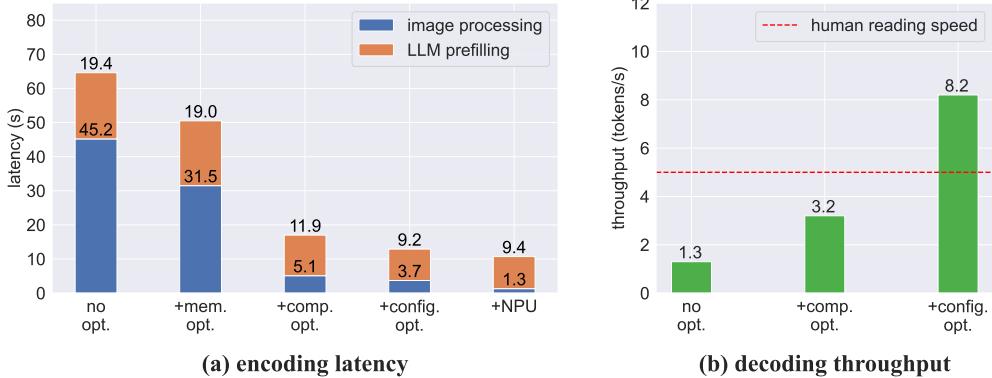


Figure 6: Ablation study on the Xiaomi 14 Pro (Snapdragon 8 Gen 3). We show the influence of different techniques on the (a) encoding latency and (b) decoding throughput. **No opt.**: non-optimized, **mem. opt.**: memory usage optimization, **comp. opt.**: compilation optimization, **config. opt.**: configuration optimization, **NPU**: NPU acceleration. Note that the encoding latency includes both model loading time and encoding time, which differs from Fig. 5’s encoding time only.

Deployment Framework. Several frameworks have been proposed for end-side deployment. Illustrated in Fig. 5, we make a thorough investigation of different frameworks for different chip types including CPU, GPU, and NPU.

Given the ubiquity of CPU usage across devices, we prioritize this chip type and opt for the llama.cpp [67] framework. Combining quantization and llama.cpp on Xiaomi 14 Pro (Snapdragon 8 Gen 3), the model achieves a text encoding latency of 64.2s and a text decoding speed of 1.3 tokens/s (as depicted in Fig. 6), which is still far from acceptable for users.

5.3 Advanced Practice

To enhance user experience, we investigate a series of advanced techniques including memory usage optimization, compilation optimization, configuration optimization, and NPU acceleration.

Memory Usage Optimization. Experimental results show that, without specific optimizations, image processing can be the bottleneck of the inference speed due to limited memory resources on mobile phones. To address the issue, we explore memory usage optimization strategies. Instead of loading both ViT and LLM simultaneously into memory, we adopt a sequential loading approach. Specifically, we first load ViT for visual encoding, followed by the LLM for visual and text token encoding. By releasing the large amount of memory occupied by LLM, we can prevent frequent paging (swapping in and out) during ViT encoding, thereby improving the program efficiency. This optimization technique, as illustrated in Fig. 6 (a), results in a notable reduction of image processing time from 45.2s to 31.5s.

Compilation Optimization. We find that directly compiling the models on the target devices can significantly improve the encoding latency and the decoding throughput. This can be attributed to better consistency between the compilation and target device instruction set architecture. As depicted in Fig. 6, this optimization endeavor yields promising results. Encoding latency shows a notable reduction from 50.5s to 17.0s, while decoding throughput experiences a significant boost from 1.3 tokens/s to 3.2 tokens/s.

Configuration Optimization. We observe that a single default configuration of the llama.cpp framework may not be optimal for diverse end-side devices. To maximize the inference speed, we devise an automatic parameter search algorithm that dynamically determines the most suitable configurations (e.g., computation allocation on different CPU cores). Through configuration optimization, we can achieve good improvements. Specifically, decoding throughput surged from 3.2 tokens/s to an impressive 8.2 tokens/s, surpassing the typical human reading speed.

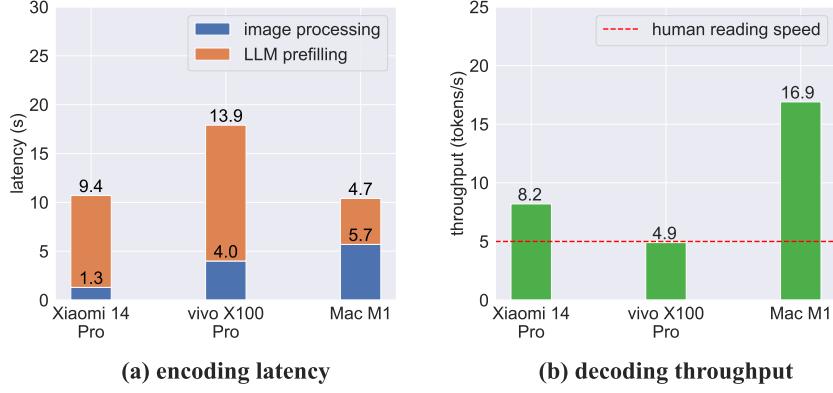


Figure 7: Results on different end-side devices. We show the (a) encoding latency and (b) decoding throughput across different device types. Xiaomi 14 Pro is the only device with NPU.

NPU Acceleration. The above techniques are mostly tailored for CPU deployment. Another promising avenue involves leveraging alternative chip types such as GPUs and NPUs. Despite the potential of GPU, we find in our experiments that current frameworks for mobile phone GPU are not optimized or compatible enough to exceed the results on CPU. As an alternative, we turn to NPUs (Neural Processing Units), which represent a novel class of specialized hardware introduced in recent years, specifically designed for accelerating AI applications. Some smartphones are already equipped with NPUs, which are recognized as better suited for addressing computation bottlenecks.

In practice, we primarily leverage NPUs to accelerate visual encoding. Specifically, we replace the backend framework of ViT to QNN, while retaining the llama.cpp backend for the LLM component. For mobile phones equipped with Qualcomm NPUs, this optimization results in a notable reduction in visual encoding time, decreasing from 3.7s to 1.3s, as illustrated in Fig. 6 (a).

5.4 Results.

Analysis. For a comprehensive assessment of MiniCPM-Llama3-V 2.5’s performance across various end-side devices, we present test results on Xiaomi 14 Pro (Snapdragon 8 Gen 3), vivo X100 Pro (Mediatek Dimensity 9300), and Macbook Pro (M1) in Fig. 7. Thanks to the deployment optimization techniques, MiniCPM-Llama3-V 2.5 can operate efficiently on both mobile phones and personal computers, delivering acceptable latency and throughput. For instance, leveraging NPU on Xiaomi 14 Pro enables it to achieve a similar encoding speed as the Mac M1. Furthermore, nearly all devices exhibit comparable or higher throughput compared with human reading speed.

Discussion. Upon analyzing Fig. 7, it becomes evident that the current computation bottleneck primarily stems from LLM prefilling, which mainly involves encoding image and text tokens for LLM inference. Promising research directions involve developing more efficient visual encoding methods with fewer visual tokens, and better leveraging GPU/NPU acceleration for LLM encoding. With increasing attention to end-side MLLMs and the rapid advancement of GPU/NPU acceleration techniques, we believe that real-time interaction with end-side MLLMs can be reached soon.

6 Experiments

In this section, we perform a comprehensive evaluation of MiniCPM-V series.

6.1 MiniCPM-V Series

We have released 3 models in the MiniCPM-V series, including MiniCPM-V 1.0, MiniCPM-V 2.0, and MiniCPM-Llama3-V 2.5. As shown in Table 3, MiniCPM-V 1.0 is trained with the pre-training stage1&2 and SFT without using the adaptive visual encoding and RLAIF-V. For MiniCPM-V 2.0, we include all of the training stages and the adaptive visual encoding strategy to further improve performance. In MiniCPM-Llama3-V 2.5, Llama3-Instruct 8B is adopted as the base LLM.

Table 3: The MiniCPM-V series, with key components and configurations. AR.: aspect ratio.

Model	Base LLM	Resolution			AR.	Pre-training	SFT	Alignment
MiniCPM-V 1.0	MiniCPM 2B	0.2M pixel (i.e., 448 × 448)			Fixed	stage-1&2	part1+2	No
MiniCPM-V 2.0	MiniCPM 2B	1.8M pixel (e.g., 1344 × 1344)			Any	stage-1&2&3	part1+2	RLHF-V
MiniCPM-Llama3-V 2.5	Llama3-Instruct 8B	1.8M pixel (e.g., 1344 × 1344)			Any	stage-1&2&3	part1+2	RLAIF-V

Table 4: Experimental results on general multimodal benchmarks. RW QA: RealWorldQA, Obj HalBench (Res./Men.) : Object HalBench with response/mention-level hallucination rates, *: our tested results with official checkpoints. The best open-source results are highlighted in **bold**.

Model	Size	Open-Compass	MME	MMB test (en)	MMB test (cn)	MMMU val	Math-Vista	LLaVA Bench	RW QA	Obj HalBench (Res./Men.) ↓
Proprietary										
GPT-4V (2023.11.06)	-	63.5	1771.5	77.0	74.4	53.8	47.8	93.1	63.0	13.6 / 7.3*
Gemini Pro	-	62.9	2148.9	73.6	74.3	48.9	45.8	79.9	60.4	-
Claude 3 Opus	-	57.7	1586.8	63.3	59.2	54.9	45.8	73.9	48.4	-
Open-source										
DeepSeek-VL-1.3B	1.7B	46.2	1531.6	66.4	62.9	33.8	29.4	51.1	49.7	16.7 / 9.6*
Mini-Gemini	2.2B	-	1653.0	-	-	31.7	-	-	-	-
Yi-VL-6B	6.7B	48.9	1915.1	68.4	66.6	40.3	28.8	51.9	53.5	19.4 / 11.7*
Qwen-VL-Chat	9.6B	51.6	1860.0	61.8	56.3	37.0	33.8	67.7	49.3	43.8 / 20.0*
Yi-VL-34B	34B	52.2	2050.2	72.4	70.7	45.1	30.7	62.3	54.8	20.7 / 14.0*
Phi-3-vision-128k-instruct	4.2B	-	-	-	-	40.4	44.5	64.2*	58.8*	-
XTuner-Llama-3-8B-v1.1	8.4B	53.3	1818.0	71.7	63.2	39.2	40.0	69.2	-	-
CogVLM-Chat	17B	54.2	1736.6	65.8	55.9	37.3	34.7	73.9	60.3	26.4 / 12.6*
Bunny-Llama-3-8B	8.4B	54.3	1920.3	77.0	73.9	41.3	31.5	61.2	58.8	-
DeepSeek-VL-7B	7.3B	54.6	1765.4	73.8	71.4	38.3	36.8	77.8	54.2	11.4 / 6.5*
LLaVA-NeXT-Llama3-8B	8.4B	-	1971.5	-	-	41.7	-	80.1	60.0	-
Idefics2	8.0B	57.2	1847.6	75.7	68.6	45.2	52.2	49.1	60.7	-
Cambrian-8B	8.3B	58.8	1802.9	74.6	67.9	41.8	47.0	71.0	60.0	-
CogVLM2-19B-Chat	19B	62.3	1869.5	73.9	69.8	42.6	38.6	83.0	62.9	-
LLaVA-NeXT-Yi-34B	34B	62.7	2006.5	81.1	79.0	48.8	40.4	81.8	66.0	-
Cambrian-34B	34B	64.9	2049.9	80.4	79.2	50.4	50.3	82.0	67.1	-
MiniCPM-V 1.0	2.8B	47.5	1650.2	64.1	62.6	38.3	28.9	51.3	51.2	21.6 / 11.5
MiniCPM-V 2.0	2.8B	54.5	1808.6	69.1	66.5	38.2	38.7	69.2	55.8	14.5 / 7.8
MiniCPM-Llama3-V 2.5	8.5B	65.1	2024.6	77.2	74.2	45.8	54.3	86.7	63.5	10.3 / 5.0

6.2 Experiment Settings

Benchmarks. We perform a comprehensive evaluation on popular benchmarks covering visual question answering, multimodal conversation, knowledge and reasoning, OCR, and hallucination. (1) **General benchmarks.** We adopt OpenCompass [26] as the general evaluation indicator, which is a comprehensive collection over 11 popular multimodal benchmarks, including MME [33], MM-Bench [63], MMMU [113], MathVista [71], LLaVA Bench [62], etc. We also report the results on RealWorldQA for real-world spatial understanding capabilities. (2) **OCR benchmarks.** We adopt three widely used benchmarks for OCR capability evaluation, including including OCR-Bench [64], TextVQA [91] and DocVQA [74]. (3) **Hallucination benchmarks.** We also include Object HalBench [85, 111] to evaluate the trustworthiness of the models.

Baselines. We compare with strong baselines in different series: For open-source models, we compare with strong models including Yi-VL-6B/34B [108], Qwen-VL-Chat [7], DeepSeek-VL-7B [68], TextMonkey [65], CogVLM-Chat-17B [102], CogVLM2-Llama3-19B [102], Idefics2-8B [52], Bunny-Llama-3-8B [39], XTuner-Llama-3-8B-v1.1 [27], LLaVA-NeXT-Llama-3-8B [53], Cambrian-8B/34B [99], LLaVA-NeXT-Yi-34B [61], DeepSeek-VL-1.3B [68], MobileVLM V2 [25], Mini-Gemini [57] and Phi-3-Vision-128k-instruct [1]. For proprietary models, we compare with GPT-4V-1106 [2], Gemini-Pro [83] and Claude 3 Opus [5].

6.3 Experimental Results

Main Results on General Multimodal Benchmarks. From the experimental results in Table 4, we have the following observations: (1) MiniCPM-Llama3-V 2.5 outperforms strong open-source models by a notable margin. For instance, MiniCPM-Llama3-V 2.5 surpasses the recent strong Idefics2-8B by 7.9 points on the OpenCompass benchmark, with similar model sizes. It also achieves better results than significantly larger models such as Cambrian-34B, LLaVA-NeXT-Yi-34B, Yi-VL-34B

Table 5: Results on OCR benchmarks. *: our tested results with official checkpoints. The best results are marked in **bold**.

Model	Size	OCRBench	TextVQA val	DocVQA test
Proprietary				
Gemini Pro	-	680	74.6	88.1
GPT-4V (2023.11.06)	-	645	78.0	88.4
Open-source				
Yi-VL-6B	6.7B	290	45.5*	17.1*
Yi-VL-34B	34B	290	43.4*	16.9*
Mini-Gemini	2.2B	-	56.2	34.2*
MobileVLM V2	3.1B	-	57.5	19.4*
DeepSeek-VL-1.3B	1.7B	413	58.4*	37.9*
Qwen-VL-Chat	9.6B	488	61.5	62.6
DeepSeek-VL-7B	7.3B	435	64.7*	47.0*
CogVLM-Chat	17.4B	590	70.4	33.3*
TextMonkey	9.7B	558	64.3	66.7
Idefics2	8.0B	-	73.0	74.0
Phi-3-vision-128k-instruct	4.2B	639*	70.9	-
MiniCPM-V 1.0	2.8B	366	60.6	38.2
MiniCPM-V 2.0	2.8B	605	74.1	71.9
MiniCPM-Llama3-V 2.5	8.5B	725	76.6	84.8

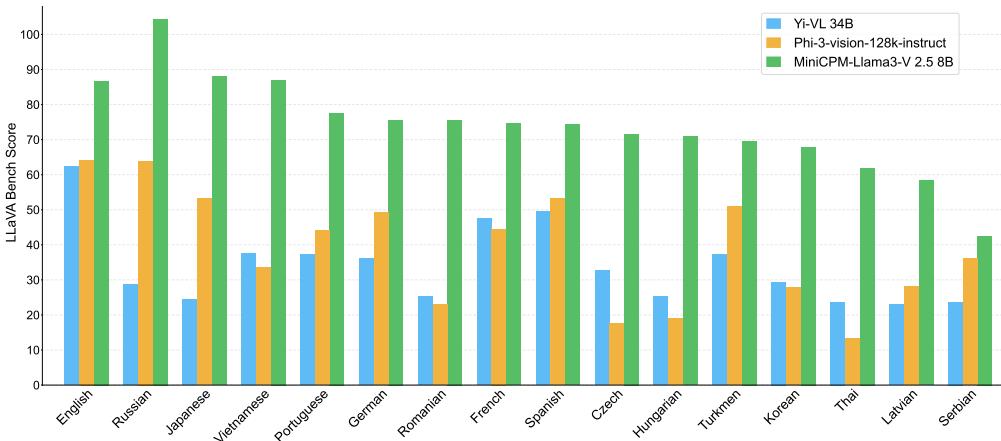


Figure 8: Experimental results on multilingual multimodal interaction. We translate LLaVA Bench reference response into different languages, and adopt GPT-4-Turbo for evaluation. Here we randomly select 15 languages for visualization.

and CogVLM2-Llama3-19B. (2) Compared with powerful proprietary models, such as GPT-4V-1106 and Gemini Pro, MiniCPM-Llama3-V 2.5 achieves better performance on the OpenCompass benchmark with significantly fewer parameters. In addition, MiniCPM-Llama3-V 2.5 also achieves lower hallucination rates than GPT-4V-1106 on Object HalBench, indicating its trustworthiness for real-world applications. (3) The smaller MiniCPM-V 2.0 with 2B parameters achieves significantly better performance compared with other 2B~3B models, and is even comparable with Llama3-based 8B MLLMs such as Bunny-Llama-3-8B. In summary, the results show that MiniCPM-V series achieves a good balance between performance and efficiency, making it more friendly for broader communities and applications.

Results on OCR Benchmarks. MiniCPM-V models also show strong OCR capabilities, including scene-text, document and screenshot understanding. As shown in Table 5, MiniCPM-Llama3-V 2.5 outperforms open-source MLLMs ranging 1.7B~34B on OCRBench, TextVQA, and DocVQA, and even performs comparably to proprietary models such as GPT-4V-1106 and Gemini Pro.

Multilingual Multimodal Capability. Based on the multilingual multimodal generalization approach from VisCPM, MiniCPM-Llama3-V 2.5 extends its multimodal capability to over 30 languages. As shown in Fig. 8, MiniCPM-Llama3-V 2.5 can outperform Yi-VL 34B and Phi-3-vision-

Table 6: The influence of RLAIF-V. The MLLM for the ablation is MiniCPM-Llama3-V 2.5.

Method	Open-Compass	MME	MMB dev(en)	MMB dev(zh)	MMMU val	Math-Vista	LLaVA Bench	Object HalBench
w/o RLAIF-V	64.5	2039.8	77.7	73.5	46.2	54.1	85.4	86.9 / 93.6
w RLAIF-V	65.1	2024.6	77.2	74.2	45.8	54.3	86.7	89.7 / 95.0

Table 7: The influence of multilingual generalization. We use 90k multilingual data (< 0.5% SFT data) for post SFT training, and show the performance changes below. ML: multilingual training.

Method	French	German	Portuguese	Spanish	Czech	Hungarian	Japanese	Korean	Thai
w/o ML	46.4	22.8	53.0	29.0	26.5	20.6	13.8	13.7	14.4
w ML	72.7	76.5	83.8	73.9	71.6	70.9	88.0	67.9	61.9

128k-instruct on the multilingual LLaVA Bench. The promising multilingual multimodal capability makes MiniCPM-Llama3-V 2.5 useful in serving larger groups with various languages.

Comparison with Other Llama-3 based Models. From experimental results in Table 4, we can observe that: (1) MiniCPM-Llama3-V 2.5 outperforms other Llama-3 based models by a large margin. For example, compared with the strong LLaVA-NeXT-Llama-3-8B, MiniCPM-Llama3-V 2.5 consistently achieves better results on all benchmarks. (2) Moreover, it is worth noting that MiniCPM-Llama3-V 2.5 requires significantly less inference computation. For example, the visual token number range of MiniCPM-Llama3-V 2.5 is (96, 960), which is lower than LLaVA-NeXT-Llama-3-8B’s (1728, 2880). This can be important especially for real-world end-side applications in terms of inference speed, first-token latency, memory usage, and power consumption.

6.4 Ablation Study

We perform an ablation study on components of MiniCPM-Llama3-V 2.5, including RLAIF-V and multilingual training.

Influence of RLAIF-V. From the results in Table 6, we can observe that RLAIF-V effectively reduces the hallucination rates of the base model on both response level and mention level. This makes the model more trustworthy in behaviors. Importantly, the hallucination reduction does not sacrifice the general capabilities. In contrast, RLAIF-V further improves the overall performance on OpenCompass by 0.6 points on an average of 11 benchmarks.

Multilingual Generalization. We investigate the necessity and effectiveness of the multilingual generalization technique. As shown in Table 8, we can see over 25 point improvement in all languages when using less than 0.5% multilingual SFT data. The results show that the multilingual generalization method can effectively improve multilingual capability with good data and computation efficiency. In addition, we notice that the performance improvement is uneven across languages. We hypothesize that the improvement extent might be related to multiple factors like the base LLM’s ability of the given language. We leave more systematical exploration for future works.

6.5 Case Study

We provide a more intuitive understanding of MiniCPM-Llama3-V 2.5 capabilities in the case study.

OCR Capability. MiniCPM-Llama3-V 2.5 shows strong OCR capabilities in real-world scenarios. Illustrated in Fig. 2, the model accurately transcribes English articles from screenshots into plain text, converts tables containing both English and Chinese content into Markdown format, comprehends code logic, and provides reasonable plans based on image content.

Any Aspect-ratio High-resolution Input. A standout feature of the MiniCPM-Llama3-V 2.5 is its capability to handle high-resolution input with extreme aspect ratios. As depicted in Fig. 9, the model well processes input with an aspect ratio of 10:1, accurately recognizing fine-grained article

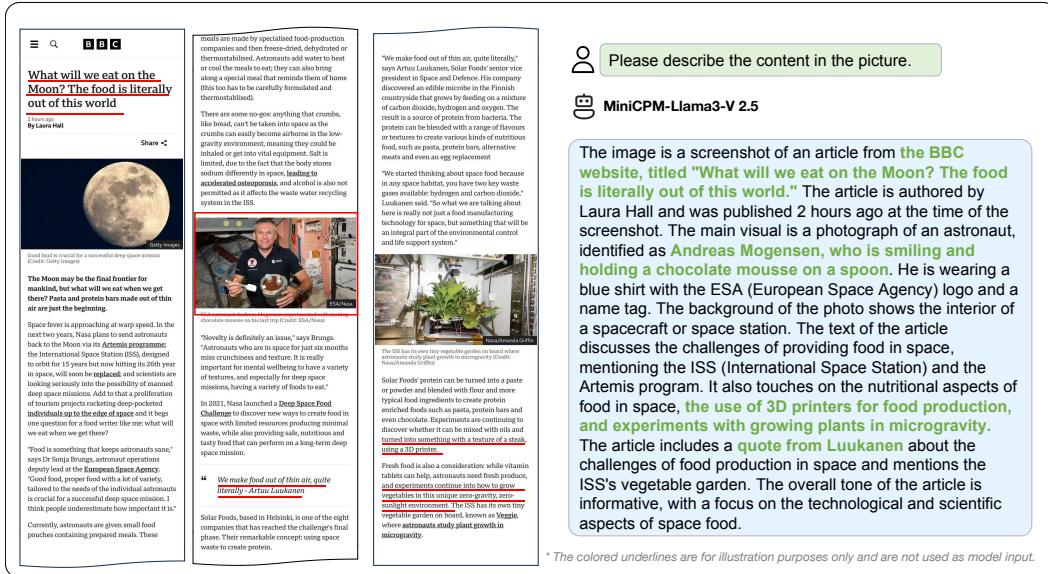


Figure 9: The example of MiniCPM-Llama3-V 2.5 for the input image with extreme aspect ratio.

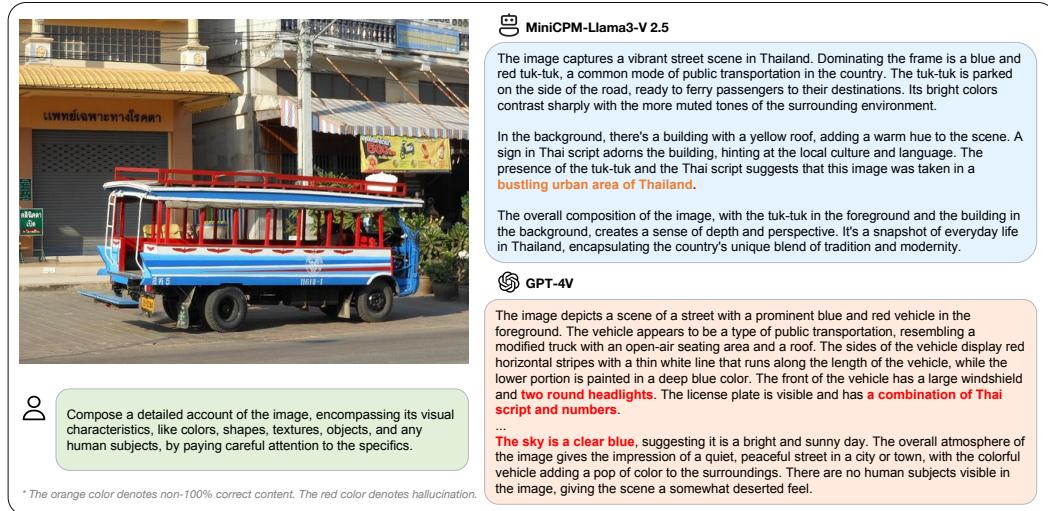


Figure 10: Comparison between MiniCPM-Llama3-V 2.5 and GPT-4V on hallucinations.

contents. Interestingly, the model can also interpret images within images, correctly describing the central image as a man “smiling and holding a chocolate mousse.”

Multilingual Multimodal Capability. Benefiting from the multilingual multimodal generalization approach from VisCPM [41], MiniCPM-Llama3-V 2.5 exhibits multilingual proficiency, generalizing across more than 30 languages. Fig. 11 showcases multimodal conversations in German, French, Japanese, Korean, and Spanish, showing good knowledge of language-specific cultures.

Trustworthy Behavior. Based on RLAIF-V, MiniCPM-Llama3-V 2.5 ensures more trustworthy responses with lower hallucination rates. As demonstrated in Fig. 10, the model’s responses exhibit fewer hallucinations as compared with powerful GPT-4V, showing its promising reliability and trustworthiness in real-world scenarios.

7 Conclusion

Contributions. In this work, we introduce the MiniCPM-V series models as a primary exploration into powerful end-side MLLMs. Thanks to techniques such as adaptive visual encoding, multilingual generalization, and the RLAIF-V method, MiniCPM-Llama3-V 2.5 can achieve GPT-4V level performance with significantly fewer parameters. With various end-side optimization techniques, this model ensures an acceptable user experience on mobile phones.

Limitations. Despite promising performance, there remain several limitations with current MiniCPM-V models. (1) Capability Depth. there is still plenty of room for improvement in enhancing multimodal understanding capability and inference efficiency. (2) Capability Width. In addition to image modality, it's promising to expand MLLM capabilities to encompass other modalities, such as video and audio, etc., where GPT-4o [78] and Google Astra [29] have given good examples.

In addition to MLLM capabilities, end-side deployment also presents unique challenges. The inference speed and latency are still far from good enough and the model service can be limited by the battery capacity. In addition, previous efforts on chips and deployment frameworks mainly target CNNs and LSTMs, which can be sub-optimal for MLLMs. Tailored efforts to MLLMs can bring ample room for improvement.

Future Works. Considering the current limitations and the promising future of end-side MLLMs, we anticipate increasing efforts from both academia and industry in enhancing model capabilities in terms of depth and width, and improving smartphone chips and deployment frameworks. We believe that simultaneous advancements in model capability and end-side device capacity will lead to end-side applications providing a satisfying user experience in the near future.

References

- [1] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] Saleem Ahmed, Bhavin Jawade, Shubham Pandey, Srirangaraj Setlur, and Venu Govindaraju. RealCQA: Scientific chart question answering as a test-bed for first-order logic. In *ICDAR*, pages 66–83. Springer, 2023.
- [4] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: A visual language model for few-shot learning. *NeurIPS*, 35:23716–23736, 2022.
- [5] Anthropic. Introducing the next generation of Claude, 2024. URL <https://www.anthropic.com/news/clause-3-family>.
- [6] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. VQA: Visual question answering. In *ICCV*, pages 2425–2433, 2015.
- [7] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- [8] Jeanine Banks and Tris Warkentin. Gemma: Introducing new state-of-the-art open models. <https://blog.google/technology/developers/gemma-open-models/>, 2024.
- [9] Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Soman, , and Sagnak TasÄşlar. Introducing our multimodal models. adept.ai/blog/fuyu-8b. 2023.
- [10] BELLEGroup. BELLE: Be everyone’s large language model engine. <https://github.com/LianJiaTech/BELLE>, 2023.

- [11] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. PaliGemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- [12] Ali Furkan Biten, Rubén Tito, Andres Mafla, Lluis Gomez, Marçal Rusinol, Ernest Valveny, CV Jawahar, and Dimosthenis Karatzas. Scene text visual question answering. In *CVPR*, pages 4291–4301, 2019.
- [13] Ali Furkan Biten, Rubén Tito, Lluis Gomez, Ernest Valveny, and Dimosthenis Karatzas. OCR-IDL: OCR annotations for industry document library dataset. In *ECCV*, pages 241–252. Springer, 2022.
- [14] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [15] Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. COYO-700M: Image-text pair dataset. <https://github.com/kakaobrain/coyo-dataset>, 2022.
- [16] Jimmy Carter. TextOCR-GPT4V. <https://huggingface.co/datasets/jimmycarter/textocr-gpt4v>, 2024.
- [17] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, pages 3558–3568, 2021.
- [18] Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. ALLaVA: Harnessing GPT4V-synthesized data for a lite vision-language model. *arXiv preprint arXiv:2402.11684*, 2024.
- [19] Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric P Xing, and Liang Lin. GeoQA: A geometric question answering benchmark towards multimodal numerical reasoning. *arXiv preprint arXiv:2105.14517*, 2021.
- [20] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal LLM’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023.
- [21] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. ShareGPT4V: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.
- [22] Wenhui Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. TabFact: A large-scale dataset for table-based fact verification. *arXiv preprint arXiv:1909.02164*, 2019.
- [23] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to GPT-4V? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.
- [24] Anoop Cherian, Kuan-Chuan Peng, Suhas Lohit, Kevin A Smith, and Joshua B Tenenbaum. Are deep neural networks smarter than second graders? In *CVPR*, pages 10834–10844, 2023.
- [25] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. MobileVLM: A fast, reproducible and strong vision language assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023.
- [26] OpenCompass Contributors. OpenCompass: A universal evaluation platform for foundation models. <https://github.com/open-compass/opencompass>, 2023.
- [27] XTuner Contributors. XTuner: A toolkit for efficiently fine-tuning LLM. <https://github.com/InternLM/xtuner>, 2023.
- [28] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. Visual Dialog. In *CVPR*, pages 326–335, 2017.
- [29] Google Deepmind. Project Astra, 2024. URL <https://deepmind.google/technologies/gemini/project-astra/>.
- [30] Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*, 2023.

- [31] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, et al. Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd. *arXiv preprint arXiv:2404.06512*, 2024.
- [32] Yifan Du, Hangyu Guo, Kun Zhou, Wayne Xin Zhao, Jinpeng Wang, Chuyuan Wang, Mingchen Cai, Ruihua Song, and Ji-Rong Wen. What makes for good visual instructions? Synthesizing complex visual reasoning instructions for visual instruction tuning. *arXiv preprint arXiv:2311.01487*, 2023.
- [33] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. MME: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- [34] Haoyuan Gao, Junhua Mao, Jie Zhou, Zhiheng Huang, Lei Wang, and Wei Xu. Are you talking to a machine? dataset and methods for multilingual image question. *NeurIPS*, 28, 2015.
- [35] Jiaxi Gu, Xiaojun Meng, Guansong Lu, Lu Hou, Niu Minzhe, Xiaodan Liang, Lewei Yao, Runhui Huang, Wei Zhang, Xin Jiang, et al. Wukong: A 100 million large-scale Chinese cross-modal pre-training benchmark. *NeurIPS*, 35:26418–26431, 2022.
- [36] Agrim Gupta, Piotr Dollar, and Ross Girshick. LVIS: A dataset for large vocabulary instance segmentation. In *CVPR*, pages 5356–5364, 2019.
- [37] Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Synthetic data for text localisation in natural images. In *CVPR*, pages 2315–2324, 2016.
- [38] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. VizWiz Grand Challenge: Answering visual questions from blind people. In *CVPR*, pages 3608–3617, 2018.
- [39] Muyang He, Yexin Liu, Boya Wu, Jianhao Yuan, Yueze Wang, Tiejun Huang, and Bo Zhao. Efficient multimodal learning from data-centric perspective. *arXiv preprint arXiv:2402.11530*, 2024.
- [40] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [41] Jinyi Hu, Yuan Yao, Chongyi Wang, Shan Wang, Yinxu Pan, Qianyu Chen, Tianyu Yu, Hanghao Wu, Yue Zhao, Haoye Zhang, et al. Large multilingual models pivot zero-shot multimodal learning across languages. *arXiv preprint arXiv:2308.12038*, 2023.
- [42] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, et al. MiniCPM: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv:2404.06395*, 2024.
- [43] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Barun Patra, et al. Language is not all you need: Aligning perception with language models. *NeurIPS*, 36, 2024.
- [44] Drew A Hudson and Christopher D Manning. GQA: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, pages 6700–6709, 2019.
- [45] Mojtaba Javaheripi and Sébastien Bubeck. Phi-2: The surprising power of small language models. <https://www.microsoft.com/en-us/research/blog/phi-2-the-surprising-power-of-small-language-models/>, 2023.
- [46] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, pages 2901–2910, 2017.
- [47] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. DVQA: Understanding data visualizations via question answering. In *CVPR*, pages 5648–5656, 2018.
- [48] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. FigureQA: An annotated figure dataset for visual reasoning. *arXiv preprint arXiv:1710.07300*, 2017.
- [49] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *ECCV*, pages 235–251. Springer, 2016.

- [50] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, JinYeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyo Han, and Seunghyun Park. OCR-free document understanding transformer. In *ECCV*, 2022.
- [51] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual Genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 123:32–73, 2017.
- [52] Hugo Laurencon, Leo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *arXiv preprint arXiv:2405.02246*, 2024.
- [53] Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. LLaVA-NeXT: Stronger LLMs supercharge multimodal capabilities in the wild, 2024. URL <https://llava-vl.github.io/blog/2024-05-10-llava-next-stronger-llms/>.
- [54] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *ICML*, pages 19730–19742, 2023.
- [55] Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. *arXiv preprint arXiv:2403.00231*, 2024.
- [56] Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. *arXiv preprint arXiv:2403.00231*, 2024.
- [57] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-Gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024.
- [58] Wing Lian, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". OpenOrca: An open dataset of GPT augmented FLAN reasoning traces. <https://huggingface.co/Open-Orca/OpenOrca>, 2023.
- [59] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014.
- [60] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023.
- [61] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. LLaVA-NeXT: Improved reasoning, OCR, and world knowledge, January 2024. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.
- [62] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 36, 2024.
- [63] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. MMBench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023.
- [64] Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezheng Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Lianwen Jin, et al. On the hidden mystery of OCR in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023.
- [65] Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai. TextMonkey: An OCR-free large multimodal model for understanding document. *arXiv preprint arXiv:2403.04473*, 2024.
- [66] Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen Lai, Yuandong Tian, Igor Fedorov, Yunyang Xiong, Ernie Chang, Yangyang Shi, Raghu Ram Krishnamoorthi, et al. MobileLLM: Optimizing sub-billion parameter language models for on-device use cases. *arXiv preprint arXiv:2402.14905*, 2024.
- [67] llama.cpp Group. llama.cpp. <https://github.com/ggerganov/llama.cpp>, 2023.
- [68] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. DeepSeek-VL: Towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024.

- [69] Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. IconQA: A new benchmark for abstract diagram understanding and visual language reasoning. *arXiv preprint arXiv:2110.13214*, 2021.
- [70] Pan Lu, Swaroop Mishra, Tanglin 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. *NeurIPS*, 35:2507–2521, 2022.
- [71] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. MathVista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- [72] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. OK-VQA: A visual question answering benchmark requiring external knowledge. In *CVPR*, pages 3195–3204, 2019.
- [73] 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. *arXiv preprint arXiv:2203.10244*, 2022.
- [74] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. DocVQA: A dataset for VQA on document images. In *WACV*, pages 2200–2209, 2021.
- [75] Minesh Mathew, Viraj Bagal, Rubén Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. InfographicVQA. In *WACV*, pages 1697–1706, 2022.
- [76] Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. MM1: Methods, analysis & insights from multimodal LLM pre-training. *arXiv preprint arXiv:2403.09611*, 2024.
- [77] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. OCR-VQA: Visual question answering by reading text in images. In *ICDAR*, pages 947–952, 2019.
- [78] OpenAI. Hello GPT-4o, 2024. URL <https://openai.com/index/hello-gpt-4o/>.
- [79] Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. *arXiv preprint arXiv:1508.00305*, 2015.
- [80] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.
- [81] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k Entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *ICCV*, pages 2641–2649, 2015.
- [82] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *NeurIPS*, 36, 2024.
- [83] Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittweiser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [84] Mengye Ren, Ryan Kiros, and Richard Zemel. Exploring models and data for image question answering. *NeurIPS*, 28, 2015.
- [85] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. *arXiv preprint arXiv:1809.02156*, 2018.
- [86] 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. *NeurIPS*, 35:25278–25294, 2022.
- [87] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-OKVQA: A benchmark for visual question answering using world knowledge. In *ECCV*, pages 146–162. Springer, 2022.
- [88] Sanket Shah, Anand Mishra, Naganand Yadati, and Partha Pratim Talukdar. KVQA: Knowledge-aware visual question answering. In *AAAI*, volume 33, pages 8876–8884, 2019.
- [89] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual Captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, pages 2556–2565, 2018.

- [90] Andrew Shin, Yoshitaka Ushiku, and Tatsuya Harada. The color of the cat is gray: 1 million full-sentences visual question answering (FSVQA). *arXiv preprint arXiv:1609.06657*, 2016.
- [91] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards VQA models that can read. In *CVPR*, pages 8317–8326, 2019.
- [92] Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. WIT: Wikipedia-based image text dataset for multimodal multilingual machine learning. In *SIGIR*, pages 2443–2449, 2021.
- [93] Tomasz Stanisławek, Filip Graliński, Anna Wróblewska, Dawid Lipiński, Agnieszka Kaliska, Paulina Rosalska, Bartosz Topolski, and Przemysław Biecek. Kleister: Key information extraction datasets involving long documents with complex layouts. In *ICDAR*, pages 564–579. Springer, 2021.
- [94] Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. A corpus of natural language for visual reasoning. In *ACL*, pages 217–223, 2017.
- [95] Stacey Svetlichnaya. DeepForm: Understand Structured Documents at Scale — wandb.ai. https://wandb.ai/stacey/\deepform_v1/reports/DeepForm-Understand-Structured-Documents/-at-Scale--Vm1ldzoy0DQ3Njg, 2020.
- [96] Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. VisualMRC: Machine reading comprehension on document images. In *AAAI*, volume 35, pages 13878–13888, 2021.
- [97] 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.
- [98] Teknium. OpenHermes 2.5: An open dataset of synthetic data for generalist LLM assistants, 2023. URL <https://huggingface.co/datasets/teknium/OpenHermes-2.5>.
- [99] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal LLMs. *arXiv preprint arXiv:2406.16860*, 2024.
- [100] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [101] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [102] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenqi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. CogVLM: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.
- [103] Wikipedia. Moore’s law, 2001. URL https://en.wikipedia.org/wiki/Moore%27s_law.
- [104] Wikipedia. Thrust-to-weight ratio, 2024. URL https://en.wikipedia.org/wiki/Thrust-to-weight_ratio.
- [105] Jiahong Wu, He Zheng, Bo Zhao, Yixin Li, Baoming Yan, Rui Liang, Wenjia Wang, Shipei Zhou, Guosen Lin, Yanwei Fu, et al. AI Challenger: A large-scale dataset for going deeper in image understanding. *arXiv preprint arXiv:1711.06475*, 2017.
- [106] Chunyu Xie, Jincheng Li, and Baochang Zhang. ZERO: A large-scale Chinese cross-modal benchmark with a new vision-language framework. 2022.
- [107] Ruyi Xu, Yuan Yao, Zonghao Guo, Junbo Cui, Zanlin Ni, Chunjiang Ge, Tat-Seng Chua, Zhiyuan Liu, Maosong Sun, and Gao Huang. LLaVA-UHD: An LMM perceiving any aspect ratio and high-resolution images. *arXiv preprint arXiv:2403.11703*, 2024.
- [108] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open Foundation Models by 01.AI. *arXiv preprint arXiv:2403.04652*, 2024.
- [109] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *ECCV*, pages 69–85. Springer, 2016.

- [110] Tianyu Yu, Jinyi Hu, Yuan Yao, Haoye Zhang, Yue Zhao, Chongyi Wang, Shan Wang, Yinxiu Pan, Jiao Xue, Dahai Li, et al. Reformulating vision-language foundation models and datasets towards universal multimodal assistants. *arXiv preprint arXiv:2310.00653*, 2023.
- [111] Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, et al. RLHF-V: Towards trustworthy MLLMs via behavior alignment from fine-grained correctional human feedback. In *Proceedings of CVPR*, pages 13807–13816, 2024.
- [112] Tianyu Yu, Haoye Zhang, Yuan Yao, Yunkai Dang, Da Chen, Xiaoman Lu, Ganqu Cui, Taiwen He, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. RLAIF-V: Aligning MLLMs through open-source AI feedback for super GPT-4V trustworthiness. *arXiv preprint arXiv:2405.17220*, 2024.
- [113] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert AGI. pages 9556–9567, 2024.
- [114] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *CVPR*, 2019.
- [115] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of ICCV*, pages 11975–11986, 2023.
- [116] Ao Zhang, Yuan Yao, Wei Ji, Zhiyuan Liu, and Tat-Seng Chua. NExT-Chat: An LMM for chat, detection and segmentation. *arXiv preprint arXiv:2311.04498*, 2023.
- [117] Ao Zhang, Hao Fei, Yuan Yao, Wei Ji, Li Li, Zhiyuan Liu, and Tat-Seng Chua. VPGTrans: Transfer visual prompt generator across LLMs. *NeurIPS*, 36, 2024.
- [118] Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. LLaVAR: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023.
- [119] Bo Zhao, Boya Wu, and Tiejun Huang. SVIT: Scaling up visual instruction tuning. *arXiv preprint arXiv:2307.04087*, 2023.
- [120] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging LLM-as-a-judge with MT-Bench and Chatbot Arena. *NeurIPS*, 36, 2024.
- [121] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- [122] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7W: Grounded question answering in images. In *CVPR*, 2016.

Wo ist dieses Bild? Bitte beschreiben Sie dieses Bild im Detail.

MiniCPM-Llama3-V 2.5

Dieses Bild erfasst den majestätischen Mount Everest, den höchsten Gipfel der Welt. Der Gipfel ist in einem goldenen Licht getaucht, das von der untergehenden Sonne reflektiert wird, was einen kontrastreichen Anblick gegen den klaren blauen Himmel bildet. Der Gipfel selbst ist mit Schnee bedeckt, was seine Höhe und Pracht noch unterstreicht. Die umliegenden Berge sind in Schatten getaucht, was Tiefe in die Szene bringt. Das Bild wurde von einer niedrigen Perspektive aufgenommen, was dem Betrachter den Eindruck vermittelt, als stünde er direkt vor dem Gipfel. Diese Perspektive betont die beeindruckende Höhe des Mount Everest und seine dominante Präsenz in der Landschaft. Das Bild ist ein Zeugnis für die Schönheit und Macht der Natur.

Que font les gens sur la photo ?

MiniCPM-Llama3-V 2.5

Les personnes dans la photo semblent apprécier une performance de musique en plein air. La présence d'un homme jouant de la guitare suggère qu'il pourrait être un musicien qui joue pour le public dans un environnement urbain.

この画像はどのアニメから来たものですか? 詳しく説明してください

MiniCPM-Llama3-V 2.5

この画像は日本の漫画家藤本弘が創作したキャラクター「ドラえもん」から来ています。ドラえもんは1980年代に始まり、多くの作品で人気を集めました。彼らは未来の猫で、時空旅行を可能にする能力があります。彼らは通常、主人公の野沢直樹と彼の家族と一緒に、科学技術、冒険、道徳の価値観のテーマに焦点を当てたストーリーを描いています。

사진 속 음식 만드는 법을 알려주세요.

MiniCPM-Llama3-V 2.5

사진 속 음식은 전통적인 한국 음식인 비빔밥이다. 비빔밥은 소스와 채소, 고기, 새우, 해면, 약이 들어 있는 나무 접시에 담아 만든다. 이 사진은 비빔밥을 준비하는 단계를 보여준다. 접시에는 비빔밥의 주요 성분인 다진고구마, 무, 피망, 오이, 미역, 소스를 포함한 비빔밥이 담겨 있다. 접시 위에 노른자 한 개가 놓여 있는데, 이는 비빔밥에 노른자를 더하는 데 사용될 것이다. 비빔밥은 일반적으로 접시에 담아 먹지만, 접시에서 접시로 옮기거나 접시에서 접시로 옮길 수도 있다. 비빔밥은 일반적으로 매운 소스와 함께 먹으며, 매운 소스는 매운 고추, 양념장, 간장과 같은 여러 가지 재료로 만들어진다.

por favor presenta este deporte

MiniCPM-Llama3-V 2.5

El deporte representado en la imagen es el toro, también conocido como tauromaquia. Es un espectáculo que se origina en España y otros países de habla hispana, que involucra a un torero o matador que enfrenta un toro en una arena rodeada por una valla. El torero utiliza una muleta, un pañuelo de cuero con una cinta de oro o plata, para manejar al toro. La indumentaria tradicional del torero, conocida como traje de luces, suele incluir un chaquetín con un patrón de flores doradas sobre fondo azul, calzones rojos, medias rojas y botas negras, todo lo cual se ve en la imagen. Este espectáculo ha sido objeto de debate sobre cuestiones éticas relacionadas con los animales y su tratamiento, y se celebra principalmente en festivales culturales, como la Feria de San Fermín en Pamplona, España.

Figure 11: Multilingual results of MiniCPM-Llama3-V 2.5. We select German, French, Japanese, Korean and Spanish as examples for demonstration.

DO 描述画面内容

MiniCPM-V 2.0

这幅图片描绘了一个繁忙的城市街道场景。在前景中，有一家名为“动画咖啡”(Animate Cafe)的商店，它有一个蓝色的招牌。这家店位于一栋砖砌建筑的前面，这栋建筑上装饰着各种广告牌和标志。在街道上，有几辆摩托车和一辆汽车停着或行驶着，表明这是一个活跃的交通环境。可以看到几个人，其中一些正在行走，而另一些则坐在摩托车上。在背景中，可以看到其他商店的招牌，包括一家名为“诚思不动产”的房地产公司，以及一家名为“全家便利店”(FamilyMart)的便利商店。这些商店的存在表明该地区可能是一个商业和住宅混合区。总体而言，这幅图像捕捉到了一个充满活力的城市环境，展示了各种商店、交通和行人，营造出一种繁忙的都市氛围。

*The colored boxes are for illustration purposes only and not used as model input.

DO Please explain in detail what this picture is about.

MiniCPM-V 2.0

The image appears to be a screenshot from a news article about the bioluminescence of the Sea of Stars on the Maldives. The article discusses the unique phenomenon where bioluminescent plankton illuminate the water at night, creating a mesmerizing light show. The author of the article, Sarah Harvey, provides a personal account of her experience witnessing this natural spectacle for the first time.

The article also touches on the scientific aspects of the phenomenon, explaining that it is caused by a chemical reaction between plankton and seawater, and it's one of the most beautiful attractions in the Maldives. The author notes that while there are no specific places to see this bioluminescence, it can be found in various locations around the Maldives.

The image includes a photo of the Sea of Stars, which is described as one of the highest-rated attractions in the Maldives. It's also mentioned that the bioluminescence of the Sea of Stars is similar to that of the Milky Way, which adds to its allure.

The article concludes with a call to action, inviting readers to choose an island with minimal light pollution to fully experience the bioluminescence of the Sea of Stars. It also suggests visiting during the Southwest Monsoon season from April to June for the most exciting experiences.

Overall, the image is a promotional piece for the Sea of Stars, highlighting its natural beauty and the unique experience it offers to visitors.

*The colored underlines are for illustration purposes only and not used as model input.

DO 描述画面内容

MiniCPM-V 2.0

The image captures a bustling scene at a gas station. Dominating the foreground is a white truck, its letters "TA" emblazoned on its side, parked under the shelter of a large, white canopy. The canopy, with its sleek lines, contrasts against the backdrop of the sky. Behind the truck, several large semi-trucks are parked, their metallic bodies gleaming under the clear sky. Among these, a red truck stands out, adding a splash of color to the scene. The ground is marked with yellow lines, guiding the flow of vehicles and ensuring safety. In the background, power lines stretch across the sky, a testament to the city's infrastructure. Trees dot the landscape, their green leaves adding a touch of nature to the urban setting. The entire scene is a harmonious blend of man-made structures and natural elements.

GPT-4V

This image features a busy fuel station during daylight. The sky is clear, hinting at a sunny day, and the natural light casts soft shadows on the ground. In the forefront, there's a white pickup truck with red and purple stripes and a logo on the door, parked at an angle. It has a utilitarian design with a cargo bed filled with assorted items, including what appears to be a red toolbox. Behind the pickup, under a large, flat-roofed shelter, multiple large trucks are parked in a row. Each truck has a different color scheme; one is predominantly blue, another is red, and a third has a white cab pulling a trailer with red and blue stripes. The shelter's roof is supported by thin pillars, creating an open space that looks clean and functional. To the left, there are fuel pumps with bright red housings and white tops, contrasting with the duller colors of the concrete ground. Various safety and direction signs are visible, such as a 'No Smoking' sign and arrows for traffic flow. There's an absence of people in the scene, which, combined with the stationary vehicles, gives the image a static feel. The overall appearance is one of an ordinary, mundane moment in a commercial or industrial area, captured in a photograph.

* The hallucinated output is marked in red color..

Figure 12: Qualitative results of MiniCPM-V 2.0. (1) Case 1 shows a scene-text understanding example on the street. (2) Case 2 shows an example of image understanding with an extreme aspect ratio. (3) Case 3 compares MiniCPM-V 2.0 and GPT-4V on detail captioning. The hallucinated outputs are marked in red color.

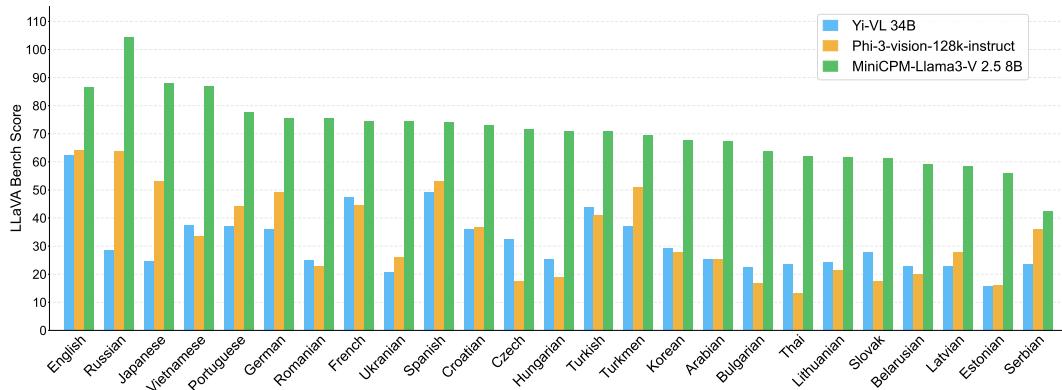


Figure 13: Multilingual results. Comparison between MiniCPM-Llama3-V 2.5, Yi-VL 34B, and Phi-3-vision-128k-instruct on more languages.