

Towards Vision Enhancing LLMs: Empowering Multimodal Knowledge Storage and Sharing in LLMs

Yunxin Li¹ Baotian Hu¹ Wei Wang² Xiaochun Cao² Min Zhang¹

Abstract

Recent advancements in multimodal large language models (MLLMs) have achieved significant multimodal generation capabilities, akin to GPT-4. These models predominantly map visual information into language representation space, leveraging the vast knowledge and powerful text generation abilities of LLMs to produce multimodal instruction-following responses. We could term this method as *LLMs for Vision* because of its employing LLMs for visual-language understanding, yet observe that these MLLMs neglect the potential of harnessing visual knowledge to enhance overall capabilities of LLMs, which could be regraded as *Vision Enhancing LLMs*. In this paper, we propose an approach called **MKS2**, aimed at enhancing LLMs through empowering **M**ultimodal **K**nowledge **S**torage and **S**haring in LLMs. Specifically, we introduce the Modular Visual Memory, a component integrated into the internal blocks of LLMs, designed to store open-world visual information efficiently. Additionally, we present a soft Mixtures-of-Multimodal Experts architecture in LLMs to invoke multimodal knowledge collaboration during generation. Our comprehensive experiments demonstrate that MKS2 substantially augments the reasoning capabilities of LLMs in contexts necessitating physical or commonsense knowledge. It also delivers competitive results on multimodal benchmarks.

1. Introduction

Recent advances (Yin et al., 2023; Driess et al., 2023; Li et al., 2023c; Ye et al., 2023) on Multimodal Large Language Models (MLLMs) have opened the eyes of text-only large language models (LLM, “blind” to visual informa-

¹School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen ²School of Cyber Science and Technology, Shenzhen Campus of Sun Yat-sen University. Correspondence to: Baotian Hu <hubaotian@hit.edu.cn>.

Preliminary work. Working in Progress.

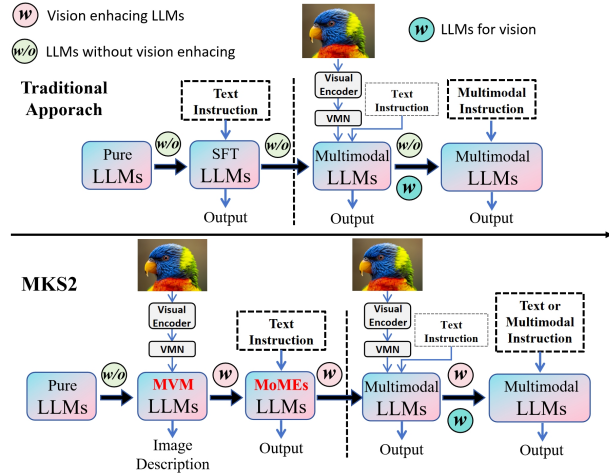


Figure 1. Comparisons between the proposed MKS2 and previous supervised fine-tuned (SFT) and multimodal LLMs. MKS2 focuses on improving LLMs with visual knowledge. VMN refers to the visual mapping network, transferring image encoding to the language space. MVM and MoMEs represent the proposed modular visual memory and soft mixtures-of-multimodal experts architecture in LLMs, respectively.

tion), allowing them to understand and process multimodal information, thereby promoting the further development of LLMs-centered Artificial General Intelligence (AGI). In this line of works such as MiniGPT-4 (Zhu et al., 2023), LLaVA (Liu et al., 2023a), and BLIP-2 (Li et al., 2023a), information outside language modality is usually aligned into language space, and then the rich knowledge stored in LLM and its powerful text generation capability are used to understand various multimodal information and generate the response corresponding to human instructions. They took a significant step towards constructing a multimodal large visual-language model similar to GPT-4 (OpenAI, 2023), contributing a lot of multimodal instruction-following data (Zhang et al., 2023; Liu et al., 2023a;b) and efficient multimodal fine-tuning technical (Ye et al., 2023; Zhu et al., 2023). These approaches concentrating on multimodal information understanding could be regarded as “*LLMs for vision*” because of mainly utilizing LLMs for

processing visual-language problems.

However, current MLLMs, pretrained and supervised finetuned (SFT) LLMs both overlook enhancing the ability of LLMs to tap into visual knowledge. Ideally, just as the human brain retains and utilizes visual information, MLLMs or LLMs should be equipped to store external visual information. In situations that require visual common sense, even in the absence of direct visual input, LLMs should be able to access this stored visual-language knowledge for combined reasoning. This goes beyond merely processing multimodal input, as “LLMs for Vision” depicted in Figure 1. Hence, we present a term “*Vision Enhancing LLMs*” to describe the desired capability for LLMs. Through this enhancement, large models would store and effectively draw upon multimodal knowledge and their knowledge base and reasoning capabilities would be enhanced.

To this end, We present **MKS2**, an innovative approach designed for empowering **Multimodal Knowledge Storage and Sharing** within LLM, consisting of two core stages: Visual Information Storage and Multimodal Knowledge Collaboration. In the first stage, we introduce Modular Visual Memory (MVM) in internal transformer blocks of LLMs to store visual information. Specifically, inspired by previous works (Kazemnejad et al., 2023; Wang et al., 2022b) focused on measuring parametric knowledge of pretrained language models and observing the knowledge storage role of feed-forward neural networks (FNN), we incorporate a two layers of FNN into each LLM block to build a lightweight visual memory. Subsequently, we employ a collection of image-text pairs to exclusively train and update MVM using two learning approaches: image-to-text generation and text-to-image retrieval. In both ways, soft image and token embeddings pass through the visual memory following attention calculations. These strategies empower LLMs to comprehend, translate, and store visual information in LLMs via a linguistic framework.

For multimodal knowledge collaboration, we introduce a soft Mixtures-of-Multimodal Experts (MoMEs) architecture. This framework leverages specialized experts, including the Modular Visual Memory (Visual Expert) and the original MLPs (Textual Expert) in LLMs during the generation process. To efficiently achieve this, we freeze all parameters of LLMs, apply Low-Rank Adaption (LoRA (Hu et al., 2021)) to each expert module and facilitate information integration across LLM blocks through a token-level soft mixing approach. By doing so, the overall model becomes adept at accommodating both multimodal and text-modality information, enabling seamless collaboration across various input forms. During training, we collect a diverse set of instruction data, containing text-only instructions and image-text multimodal instruction-following data, to ensure the effectiveness of MoMEs in handling multimodal as well

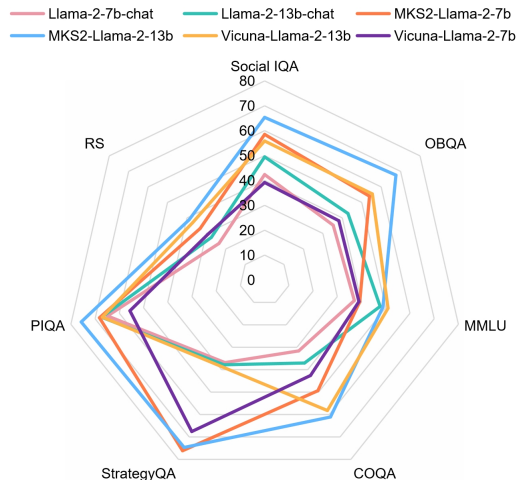


Figure 2. MKS2-Llama-2-13b achieves SOTA zero-shot performance on seven natural language reasoning tasks. It indicates achieving multimodal knowledge storage and share is effective for improving the overall capability of LLMs.

as text-only tasks.

To validate the effectiveness of our approach, we evaluate MKS2 on seven natural language processing (NLP) benchmarks and six image-text understanding datasets. Extensive experiment results indicate that MKS2 achieves superior performances on NLP tasks requiring physical or visual world knowledge, e.g, MKS2-Llama-2 significantly exceeds Llama-2-chat as shown in Figure 2. It also achieves competitive performances on image-text understanding scenarios compared to previous MLLMs.

Our main contributions can be summarized as follows:

- We introduce MKS2, a vision-enhanced learning framework for LLMs, designed for effective storage and sharing of multimodal knowledge. This framework efficiently handles both multimodal and text-only inputs.
- MKS2 demonstrates superior outcomes in knowledge-intensive tasks over traditional SFT LLMs and LLMs employing Reinforcement Learning from Human Feedback (RLHF).
- Ablation studies validate the efficacy of mixtures-of-multimodal-experts that incorporates a visual knowledge expert. This architecture distinctly improves the performance of LLMs beyond the capacities of conventional supervised fine-tuned LLMs.
- Our experiments indicate that multimodal instruction-following data further enhances LLMs’ performance in natural language reasoning tasks that require extensive commonsense.

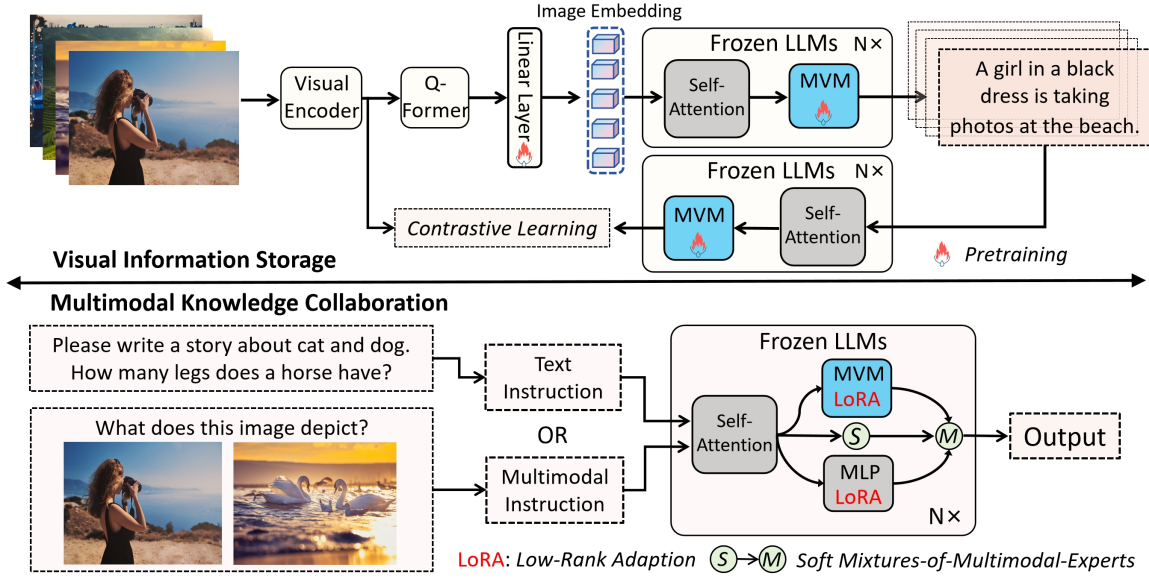


Figure 3. The overall work flow of MKS2. It realizes visual information storage and multimodal knowledge collaboration in LLMs. In the first stage, we introduce the modular visual memory (MVM, presented in blue blocks) and train it through language-centered learning strategies. We also present a soft mixtures-of-multimodal experts (MoMEs) architecture to accomplish multimodal knowledge collaboration during generation.

2. Preliminaries

We first review the supervised fine-tuning approach and recently proposed multimodal instruction-following tuning method for LLMs.

2.1. Supervised Fine-tuning

A pure pretrained LLM is fine-tuned on high-quality labeled datasets using token-level supervision to produce a Supervised Fine-Tuned model, dubbed as SFT-LLM. Common methods are using GPT-4 automatically constructed instruction data (Wang et al., 2022c) and manually annotated high-quality data from downstream tasks (Chung et al., 2022) to fine-tune pure LLMs. To reduce training costs, recent works present some efficient instruction-tuning approaches, e.g., LoRA (Hu et al., 2021), QLoRA (Detmers et al., 2023), etc. These SFT-LLMs are capable of generating human-like responses for various text-only instructions, having a profound impact on all walks of life.

2.2. Multimodal Instruction-Following Tuning

Compared to traditional visual-language models such as Oscar (Li et al., 2020), Flamingo (Alayrac et al., 2022), OFA (Wang et al., 2022a), etc, the multimodal instruction-following tuning approach explored extending the text-only instruction tuning in LLMs to multi-modality. These MLLMs applying LLMs as the multimodal information processor achieve impressive zero-shot performances on un-

seen tasks. Generally, as the traditional approach depicted in Figure 1, a frozen visual encoder (e.g., visual encoder of CLIP) is used to obtain the sequence representation of an image and a visual mapping network (VMN, a linear projection layer or Q-former from BLIP-2) projects the image encoding into soft image embeddings into the language space of LLMs. Then, we can utilize an efficient fine-tuning technical to allow LLMs to process multimodal information, thereby turning LLMs into MLLMs.

Formally, a multimodal image-text instruction sample could be expressed in the following triplet form, i.e., (I, T, R) , where I, T, R represent the input image, text description (about human demands or image-related premises), and ground-truth response, respectively. During training, the constructed MLLMs is forced to predict the next token of response via the autoregressive objective, which could be presented as:

$$\mathcal{L}(\theta) = - \sum_{i=1}^N \log P(R_i | I, T, R_{<i}; \theta), \quad (1)$$

where N is the length of response and θ refers to the training parameters in the whole framework.

In conclusion, we find that these two approaches ignore introducing visual knowledge to improve overall capabilities of LLMs for processing text-only tasks.

3. Methodology

In the following subsections, we will present the two stages of MKS2 in detail: Visual Information Storage and Multimodal Knowledge Collaboration.

3.1. Visual Information Storage

To realize visual information storage in LLMs, we propose injecting Modular Visual Memory (MVM) into internal blocks of LLMs and forcing MVM to memorize open-world visual information via language-centered learning strategies.

Modular Visual Memory (MVM). This module is two layers of feed-forward neural networks (FFN) and injected into each transformer block of LLMs. As the top part shown in Figure 3, the input image I is first projected into soft image embedding \mathbf{h}_I via the pretrained visual encoder, Q-former from BLIP-2, and a learnable linear layer. Take the first block as an example; the calculation process can be presented as follows:

$$\begin{aligned} \mathbf{h}_s^T &= \text{Self-Attention}(\mathbf{h}_I), \\ \mathbf{h}_F^T &= \mathbf{h}_s^T + \text{MVM}(\text{layernorm}(\mathbf{h}_s^T)), \end{aligned} \quad (2)$$

where Self-Attention is the original attention calculation in LLMs. We just inserted MVM inside the original LLMs and did not change other structures. All hidden states pass the MVM after gaining the output \mathbf{h}_s^T of Self-Attention, and we also set the overall size of visual memory by controlling the hidden state dimensions of FFN.

Language-Centered Learning Strategies. As we consider LLMs as analogs to the human brain, we have embarked on a groundbreaking endeavor to create the visual storage memory in LLMs. Our ultimate goal is to empower LLMs with the capability to comprehend a given image and conjure related visual scenarios based on textual input, akin to human cognition. To this end, we adopt two learning objects to train MVM with a large amount of image-text pairs. As shown in Figure 3, we allow LLMs to generate the language description of an image, which resembles understanding and translating an image like brain. Additionally, given a sentence with some visual objects, LLM should attach to the sentence-related image, which resembles imagination. Suppose that the short description (caption) of an input image I is D , the description generation loss is

$$\mathcal{L}_c = -\frac{1}{N} \sum_{i=1}^N l_c(\text{IMG}_i, D_i), \quad (3)$$

where N is the number of image-text pairs in a batch and l_c refers to the cross-entropy loss.

While retrieving the related image, we use the output hidden state h_e of the end token $\langle /s \rangle$ of input caption to match the image embedding. Concretely, we employ a learnable linear

layer to project it into the same dimension with image global encoding obtained by the visual encoder. Then we calculate the cosine similarity between them and minimize the InfoNCE loss for text-to-image (t2i) retrieval over a batch of N samples. The negatives are other irrelevant images in a batch. Hence, the total language-centered learning loss is

$$\begin{aligned} \mathcal{L}_{\text{Stage1}} &= \mathcal{L}_c + \mathcal{L}_{\text{t2i}}, \\ \mathcal{L}_{\text{t2i}} &= -\frac{1}{N} \sum_{i=1}^N \left(\log \frac{\exp(\text{sim}(D_i, \text{IMG}_i) / \tau)}{\sum_{j=1}^N \exp(\text{sim}(D_i, \text{IMG}_j) / \tau)} \right), \end{aligned} \quad (4)$$

where τ is a learnable temperature parameter. During training, we freeze all pretrained parameters of LLMs and only update MVM. In addition to the way of retrieving images to achieve visual information association, using image generation technology for joint training is also an alternative approach.

3.2. Multimodal Knowledge Collaboration

After gaining visual information storage inside LLMs, we need consider how to realize multimodal knowledge collaboration during generation. Regarding pretrained MVM and MLP in LLMs as visual and textual experts respectively, we propose a soft mixtures-of-multimodal experts (MoMEs) approach to achieve multimodal knowledge utilization at the token level.

Mixtures-of-Multimodal Experts (MoMEs). To speed up training process, as the bottom part shown in Figure 3, we freeze MVM and other parameters of LLMs, applying Low-Rank Adaption (Hu et al., 2021) (LoRA) to tow-modality experts: MVM and MLP. We denote the inputs tokens for one sequence inputted to MoMEs by $\mathbf{X} \in \mathbb{R}^{m \times d}$, where m is the number of tokens and d is their dimension. The computed process for visual and language knowledge expert could be given in

$$\begin{aligned} \mathbf{h}_{VE} &= \text{LoRA-MVM}(X), \\ \mathbf{h}_{TE} &= \text{LoRA-MLP}(X), \\ \text{LoRA}(W_0) &:= W_0 X + \Delta W X = W_0 x + B A X, \end{aligned} \quad (5)$$

where B, A are learnable parameters added for each pretrained weight of visual and textual experts. LoRA-MVM and $\text{LoRA-MLP}(X)$ represent original knowledge experts equipped with additional LoRA calculation. By doing so, the training process is efficient because of doing not update the overall parameters of experts.

Each MoE layer uses expert functions (shown in E.q. 5) applied on individual tokens, namely $\{f_i : \mathbb{R}^d \rightarrow \mathbb{R}^d\}_{1:2}$. Each expert will process p slots, and each slot has a corresponding d -dimensional vector of parameters. As $S \rightarrow M$ shown in Figure 3, the token-level combination for expert

outputs can be presented as

$$\begin{aligned} S &= \text{Softmax}(w_s X + b_s), \\ \mathbf{h}_M &= S_1 \mathbf{h}_{VE} + S_2 \mathbf{h}_{TE}, \\ \mathbf{h}_o &= \mathbf{h}_s^T + \mathbf{h}_M, \end{aligned} \quad (6)$$

where $S \in R^{X \times 2}$ and the final dimension is normalized with Softmax calculation. The output of each block in LLMs is denoted to \mathbf{h}_o .

3.3. Training

In the first stage, the size of used image-text pairs are about 2.3M from CC3M (Changpinyo et al., 2021), COCO Captioning (Chen et al., 2015), and Flickr-30k (Plummer et al., 2015). To achieve multimodal knowledge collaboration, as shown in Figure 1, we use text-only and image-text instruction-following data to train the overall architecture. The added modular visual memory and LLMs are frozen during training. We use widely-used instruction data including: high-quality natural language processing tasks from Flan-T5 (Chung et al., 2022), complex instruction-finetuning data from WizardLLMs (Xu et al., 2023), and multimodal instruction data LLaVAR (Zhang et al., 2023), which totally consists of 1.5M text-only and 166k image-text instruction tuning data.

4. Experiments

4.1. Datasets

Natural Language Processing Benchmarks. We use seven text-only downstream datasets to comprehensively evaluate MKS2, which consists of physical world knowledge-relevant datasets and basic ability assessment benchmark MMLU (Hendrycks et al., 2021). We use multiple choice question answering tasks that can benefit from visual knowledge: PIQA (Bisk et al., 2020) that requires physical commonsense reasoning, Commonsense QA (CSQA) (Talmor et al., 2019) for evaluating the commonsense reasoning capability of models, OpenBook QA (OBQA) (Mihaylov et al., 2018) that requires multi-step reasoning, use of additional common and commonsense knowledge, and rich text comprehension, RiddleSense (RS) (Lin et al., 2021) for complex understanding of figurative language and counterfactual reasoning skills, Social IQA (Sap et al., 2019) which focuses on physical or taxonomic knowledge for testing social commonsense intelligence, StrategyQA (Geva et al., 2021) that needs the reasoning steps should be inferred using a strategy.

Image-Text Understanding Benchmarks. To evaluate the multimodal capability of our proposed model, we introduce six classical Visual Question Answering (VQA) datasets: VQAv2 (Antol et al., 2015), OK-VQA (Marino et al., 2019), ST-VQA (Biten et al., 2019), OCR-VQA (Mishra et al., 2019), TextVQA (Singh et al.,

2019), and DocVQA (Mathew et al., 2021). VQAv2 is a classic open-world VQA dataset, containing more than 1 million samples. Scene Text Visual Question Answering (STVQA) consists of 31,000+ questions across 23,000+ images collected from various public datasets. The OCRVQA dataset includes more than 1 million question-answer pairs that cover over 207,000 book cover images. The TextVQA dataset consists of over 45,000 questions related to text on more than 28,000 images selected from specific categories of the OpenImages dataset. DocVQA is a comprehensive dataset comprising 12,767 document images with diverse types and content, accompanied by over 50,000 questions and answers. For datasets containing far more than 5000 image-question pairs, we selected the first 5000 pairs for test, similar to Liu et al. (2023c).

4.2. Comparing Models

The comparing models mainly comprise three types of open-source large language models Llama-2: SFT Llama-2, RLHF-tuned Llama-2, and recently proposed MLLMs. To verify the new vision enhancing supervised fine-tuning method MKS2, we present text-only instruction tuned variant Llama-2-7b-INST-LoRA and Vicuna-Llama-2 (Penedo et al., 2023), where we adopt the text instruction data identical to our approach and set $r = 16$ for LoRA to train Llama2-7b-INST. Hence, the training parameters of Llama-2-7b-INST-LoRA is similar to the proposed MKS2-Llama-2-7B, about 14M. RLHF-tuned models are language models that have been trained using a combination of human feedback and reinforcement learning techniques, achieving better performance to understand human instruction and generate high-quality responses. We mainly comparing with Llama-2-7b-chat and Llama-2-13b-chat variants released by Meta. Additionally, to evaluate the multimodal information processing capability of MKS2, we also introduce recently proposed MLLMs as baselines. Flamingo (Alayrac et al., 2022) and OFA (Wang et al., 2022a) are traditionally pretrained visual and language models, which have seen an amount of image-text pairs. BLIP-2 (Li et al., 2023a) is a widely-used visual and language models, achieving remarkable zero-shot performance on downstream image-text understanding tasks. MiniGPT-4 (Zhu et al., 2023), FROMAGE (Koh et al., 2023), mPLUG-Owl (Ye et al., 2023), LLaVR (Liu et al., 2023a) and InstructBLIP (Li et al., 2023a) are multimodal instruction tuned MLLMs, trained with enormous image-text instruction-following data.

4.3. Implementation Details

We take the pretrained Llama-2 version (Touvron et al., 2023) as the backbone of MKS2 and run all models with Adam Optimizer (Kingma & Ba, 2014) on 4 A100-80G GPUs with python environment. All models are trained and tested with Bfloat16 floating-point format. The dimension

Table 1. Zero-shot model performances on natural language processing benchmarks. Models with [†] indicate that their SFT or RLHF tuned data are unknown or unused in our work. “INST-LoRA” refers to applying the widely-used LoRA technical to fine-tune LLMs with same text-only instruction data. “Multimodal-SFT” represents the multimodal instruction-following data. “Avg.Score” shows the average evaluation score on the total tasks. **Bold** and underlined numbers refer to the best and second-best performance for comparative model variants of Llama-2-7b/13b, respectively.

Models↓ Types →	COQA	StrategyQA	Social IQA	OBQA	PIQA	RS	MMLU	Avg.Score
Llama-2-13b-chat [†] (Touvron et al., 2023)	37.02	37.80	49.46	42.89	67.29	27.45	47.69	44.23
Vicuna-Llama-2-13b [†] (Chiang et al., 2023)	58.21	38.82	55.85	55.46	67.01	37.02	50.96	51.90
Llama-2-13b-INST-LoRA _{r=16}	57.68	63.73	63.80	58.6	71.98	38.51	46.70	57.28
MKS2-Llama-2-13b	62.10	<u>74.68</u>	65.71	67.6	76.11	41.03	<u>48.83</u>	62.30
w/o Multimodal-SFT	58.77	74.73	<u>64.56</u>	<u>60.6</u>	<u>75.03</u>	<u>38.74</u>	48.44	<u>60.12</u>
w/o Multimodal-SFT & MoMEs	54.81	68.21	62.25	54.0	67.95	35.08	46.50	55.54
Llama-2-7b-chat [†] (Touvron et al., 2023)	31.62	36.83	42.37	35.3	64.90	23.53	37.05	38.82
Vicuna-Llama-2-7b [†] (Chiang et al., 2023)	42.58	67.58	39.71	38.2	55.62	<u>29.32</u>	<u>38.94</u>	44.56
Llama-2-7b-INST-LoRA _{r=16}	41.93	74.10	54.65	39.4	53.42	27.01	38.68	47.02
MKS2-Llama-2-7b	49.38	<u>76.15</u>	58.51	54.0	68.19	33.20	39.27	54.10
w/o Multimodal-SFT	<u>44.06</u>	76.46	<u>57.72</u>	<u>50.5</u>	<u>67.10</u>	28.99	37.45	<u>51.84</u>
w/o Multimodal-SFT & MoMEs	42.84	70.46	55.42	37.0	60.71	25.27	37.71	47.06

of the middle layer of the inserted visual memory module is 1/4 of the hidden state size of LLMs. For Llama-2-7b, the total parameters of MVM is about 410 million. During visual information storage, we take the frozen visual encoder and Q-former from BLIP-2-FlanT5-xxl to obtain the image encoding, thus the length of soft image embedding is 32. Additionally, we set the initial learning rate to 1e-4 and train the model for about 2 epochs with warm up steps equaling 5000. The batch size is set to 32 with four-step gradient accumulation for single GPU device. While performing instruction-following learning, we set the batch size, r in LoRA to 3 and 8, respectively, and the max length of input is set to 1024. To tag the position of image embedding, we introduce two learnable tokens $\langle img-start \rangle$ and $\langle img-end \rangle$. Similar to Llama-2-chat, we add [INST] and [/INST] at the starting and ending of inputting text instruction, as like “[INST] Please write a short story about cat and dog [/INST]”. During generation, we set beam sizes to 1 and 4 for text-only and VQA tasks respectively.

4.4. Overall Performances

Performance of vision enhancing LLMs. We present zero-shot model performances on Table 1, aiming to evaluate the instruction-understanding and open-world problem solving abilities of LLMs. We observe that the proposed method MKS2-Llama-2-7B/13B achieves best performances on almost all evaluation datasets, especially on substantially suppressing Llama-2-7b/13b-chat. Compared to powerful Llama-2-7b-INST-LoRA of the same magnitude, MKS2-Llama-2-7b could gains by about 8% on CommensenseQA, 14.5% on OpenBookQA, 16% on PIQA, and 8.6% on RS, respectively. Hence, MKS2 is capable of markedly improv-

ing the overall performance on text-only tasks requiring physical world knowledge. Compared to Vicuna-Llama-2 models with all parameters of Llama-2 updating, our approach stands out by requiring being fine-tuned on only a small fraction of parameters ($< 0.2\%$ of LLM parameters) while still achieving superior performance on several tasks.

Competitive performances on multimodal benchmarks.

We also present the model’s zero-shot performance on the VQA dataset in Table 2. To gain suitable and robust image embeddings, we further fine-tuned the visual mapping network for one epoch and freeze all other parameters, which does not affect any text-only performance of LLMs. We can see that MKS2-Llama-2-7b could achieve comparative performances on open-world and scene text VQA datasets. It’s noteworthy that there was no discernible performance degradation when comparing MKS2 performance on multimodal tasks to that of the original large language model without visual enhancement. This implies that the addition of visual enhancement in LLMs did not lead to a loss in text-related knowledge, while performing LLMs for vision. Moreover, the incorporation of text-only data proves to be a valuable strategy for enhancing the model proficiency in answering open visual questions. While mixed instruction data leads to improvements in open-world question answering, it appears to have a detrimental effect on scene text recognition, possibly due to shifts in training data distribution. Further investigation into the fine-tuning process and data distributions can help optimize MKS2 performance across a wider range of tasks involving both text and images.

Table 2. Zero-shot performances on some multimodal datasets. “NumImg” represents the total number of images contained in the pretraining stage. The size of input image is always 224² for the following models. “[‡]” indicates that the corresponding model employs the training samples of following evaluation benchmarks such as VQAv2, OK-VQA, and OCR-VQA, leading to unfair comparison.

Models↓ Types →	NumImg	VQAv2	OK-VQA	STVQA	OCR-VQA	TextVQA	DocVQA
Flamingo (Alayrac et al., 2022)	>1B	49.2	41.2	19.3	27.8	29.0	5.0
MiniGPT-4 (Vicuna-7b) (Zhu et al., 2023)	5M	44.3	32.1	14.0	11.5	18.7	3.0
LLaVA (Vicuna-7b) (Liu et al., 2023a)	0.6M	53.5	47.4	22.1	11.4	28.9	4.5
LLaVAR (Vicuna-13b) (Zhang et al., 2023)	1M	-	-	30.2	23.4	39.5	6.2
OFA-Large (Wang et al., 2022a)	20M	40.2	19.3	-	-	-	-
FROMAGe (OPT-6.7b) (Koh et al., 2023)	3.3M	44.1	20.1	-	-	-	-
mPLUG-Owl [‡] (Ye et al., 2023)	11B	-	-	29.3	28.6	40.3	6.9
InstructBLIP [‡] (FlanT5XL) (Liu et al., 2023a)	129M	62.6	50.1	23.9	39.7	33.1	3.8
BLIP-2 (OPT-6.7b) (Li et al., 2023a)	129M	50.1	36.4	13.4	10.6	21.2	0.8
BLIP-2 (FlanT5XL) (Li et al., 2023a)	129M	42.8	25.6	15.8	26.6	25.2	2.9
BLIP-2 (FlanT5XXL-11B) (Li et al., 2023a)	129M	45.4	27.8	21.7	30.7	32.2	4.9
MKS2 -Llama-2-13b	2.3M	54.4	45.1	23.4	35.7	34.2	6.7
MKS2 -Llama-2-7b	2.3M	53.3	42.1	22.3	25.2	33.1	6.7
w/o Text-SFT	2.3M	50.2	40.8	21.5	36.5	34.2	7.4
w/o Text-SFT & MoMEs	2.3M	50.1	41.2	21.4	35.3	34.3	7.3

4.5. Ablation Study and Analysis

Effects of MKS2. Comparing the experimental results of MKS2 w/o Multimodal-SFT, MKS2 w/o Multimodal-SFT & MoMEs, and Llama-2-7b-INST-LoRA in Table 1, we observed that the incorporation of visual information into the MVM positively impacted the model’s performance on common sense reasoning tasks. This demonstrates that MoMEs can leverage stored visual information to improve its understanding and reasoning abilities in various contexts. Additionally, the performances of MKS2 variants w/o Text-SFT and w/o Text-SFT + MoMEs on multimodal tasks further emphasizes its ability to access textual knowledge without significant compromise. The integration of visual information alongside textual data did not hinder the model’s capacity to extract and utilize textual knowledge effectively. This suggests that the core text-related capabilities of LLMs remain robust when dealing with multimodal inputs.

Impact of multimodal instruction data for MKS2 performance. Our experimental results in Table 1 highlights the impact of multimodal instruction data on MKS2’s performance. While it effectively enhances the accuracy in addressing questions related to physical world knowledge, it does not provide a substantial boost in solving intricate and complex problems that demand advanced reasoning and strategic thinking. These findings emphasize the importance of tailoring data and approaches to specific task requirements when leveraging multimodal data in large language models for optimal performance. Further research may uncover strategies to bridge this gap for complex problem-solving tasks.

Table 3. Ablation Experiments on five datasets requiring common sense. We explore various data and visual memory sizes to check model performances. All models are built upon Llama-2-7b, where MKS2[♣] signifies MKS2 w/o Multimodal-SFT. “-AM-BM” indicates that it includes a “A” size of visual memory and was trained using “B” number of image-text pairs.

Models↓ Types →	COQA	OBQA	PIQA	RS	Social IQA
Llama-2-7b-chat [†]	31.62	35.3	64.90	23.53	42.37
Vicuna-Llama-2-7b [†]	42.58	38.2	53.42	29.32	39.71
MKS2 [♣] -410M-2.3M	44.06	50.5	67.10	28.99	57.72
MKS2 [♣] -410M-12.3M	46.12	50.8	64.68↓	30.06	55.31↓
MKS2 [♣] -810M-2.3M	43.78	48.6	66.97	28.89	57.21
MKS2 [♣] -810M-12.3M	46.56	49.8	68.10	30.65	57.84

Could text-only instruction data be effective for enhancing multimodal performance while building MLLMs?

In the ablation experiments discussed in Table 2, text-only instruction data was found to enhance the performance of pre-trained LLMs on various open multimodal problems, particularly open-ended questions that demand a fusion of textual and visual information. However, the introduction of text-only instruction data may lead to a trade-off, as it could potentially diminish the LLMs to answer visual questions involving scene text recognition. These findings underscore the importance of striking a careful balance when incorporating additional data modalities into LLMs, with a keen consideration of task requirements and data characteristics. Such careful integration ensures that overall model performance is optimized without unintended consequences on

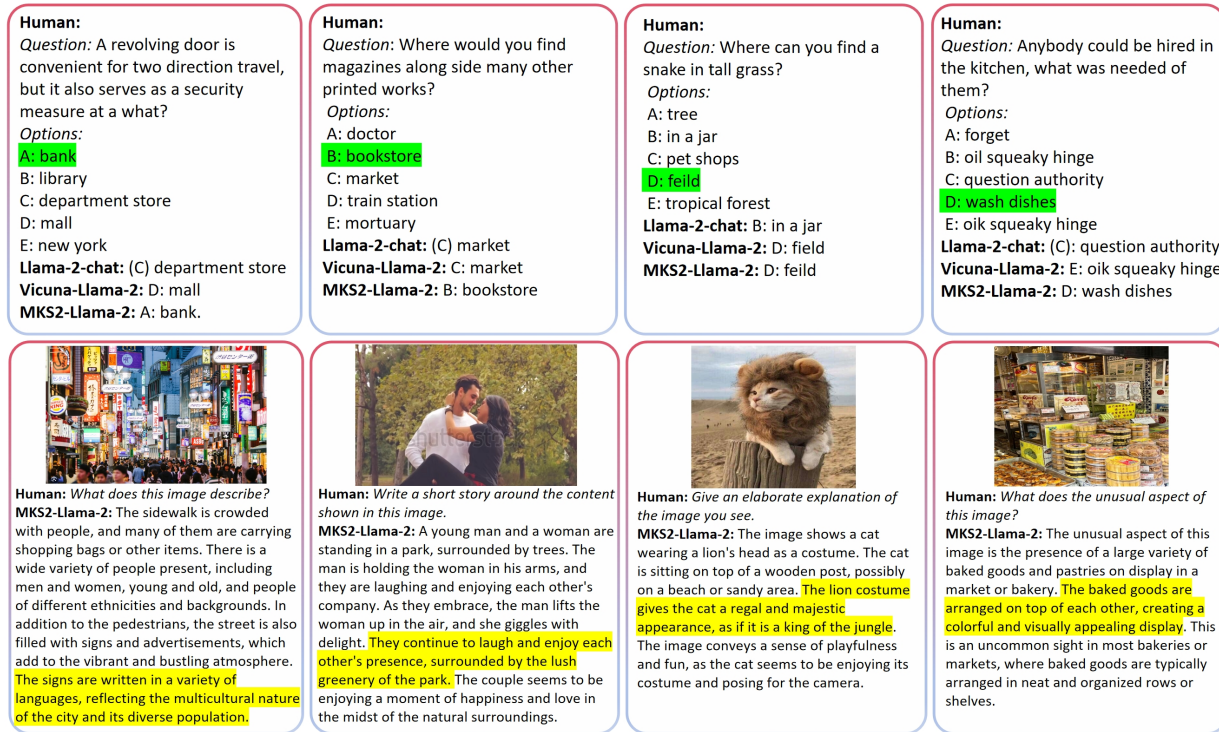


Figure 4. An illustration of cases generated by comparative models and MKS2. Green words refer to the correct answer in the natural question answering tasks. Yellow parts represent the interesting and correct content in the response.

its ability to handle diverse multimodal challenges. Consequently, optimizing LLMs for multimodal tasks requires a nuanced approach tailored to the specific problem domain.

Impact of visual memory size and data scale. We introduce more image-text pairs from LAION-400M (Schuhmann et al., 2021) to analyse the impacts of vision memory and data sizes, and the experimental results are shown in Table 3. Our experimental results reveal that enlarging the size of the pre-trained image-text data leads to improvements in the MKS2 model’s performance across various downstream tasks. Additionally, we observe that expanding the visual storage module of the MKS2 model not only enhances performance but also contributes to greater stability in this enhancement, as the size of the image-text data increases. To summarize, these findings suggest a twofold strategy for optimizing the MKS2-based LLMs during the supervised fine-tuning phase: increasing the size of the image-text data and selecting a proportionately larger visual storage size. This approach is crucial for achieving more consistent and stable performance improvements.

4.6. Case Study

We present some cases in Figure 4 to further show the overall capability of MKS2-Llama-2. We can see that the proposed model achieve better performance while answering com-

monsense Q&A with physical knowledge. In addition, we also observe that the multimodal understanding capability of MKS2-Llama-2 is powerful, such as it could recognize the funny of the cat with a lion’s head. In addition, it can employ the relevant knowledge to enrich the response based on visual clues, e.g., the short story around the content shown in the second image.

5. Related Works

Visual knowledge enhanced methods. There is a long line of work on utilizing explicit visual information to improve the imaginative representation of language, thus promoting diverse generation capability of LLMs. Particularly, Jin et al. (2022) leverage visual knowledge in NLP tasks, developing multiple cross-model enhanced methods to improve the representation capability of pretrained language models. Some works (Shi et al., 2019; Lu et al., 2022; Li et al., 2023b) proposed to retrieve images corresponding to texts from the image corpus and use visual knowledge to improve the performance on the downstream tasks such as text completion (Zellers et al., 2019), story generation (Fan et al., 2018), and concept-to-text (Barzilay & Lapata, 2005). Recently, some researchers (Long et al., 2021; Yang et al., 2021; Zhu et al., 2022) proposed to utilize the powerful text-to-image technical to obtain the imagination represen-

tation of language and infuse them into the language model via the prefix-tuning way. In this paper, we present visual information storage in LLMs and achieve visual knowledge enhancing LLMs without explicitly inputting images to language models.

LLMs for vision. Recent works (Zhang et al., 2023; Zhu et al., 2023; Li et al., 2023a) towards multimodal LLMs focus on utilizing the extensive knowledge and language generation capabilities of LLMs to solve multimodal tasks (Li et al., 2023d), especially for visual understanding and reasoning. Firstly, these works usually map the visual information obtained by pretrained visual encoder into the representation space of LLMs, through a learnable linear projection layer (Merullo et al., 2023), MLP, or Q-Former (Li et al., 2023a). This stage is usually called feature alignment and only a few hundred thousand data may be needed to do a good job (Liu et al., 2023a). Afterwards, the initial MLLMs will be tuned via multimodal instruction-following data (Ye et al., 2023; Li et al., 2023c; Bai et al., 2023). At this stage, LLMs and projection layer are often tuned together only with multimodal instruction data. The commonly used large language model is the SFT LLM and the tuning approach adopts widely used lightweight LoRA (Hu et al., 2021). These works, however, rarely consider using visual knowledge to enhance the pure text processing capabilities of LLMs, thereby building a more robust LLM or MLLM.

6. Conclusion

In this paper, we present a new approach MKS2 that allows LLMs to memorize and employ visual information, achieving multimodal knowledge storage and collaboration in LLMs. MKS2 consists of modular visual memory and soft mixtures-of-multimodal experts, which are used to store visual information and realize multimodal knowledge collaboration, respectively. We conduct extensive experiments on many NLP and VQA tasks and the experimental results show that MKS2 is capable of enhancing the reasoning capability of LLMs and being used to solve multimodal problems.

References

- Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C. L., and Parikh, D. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015.
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- Barzilay, R. and Lapata, M. Collective content selection for concept-to-text generation. In *Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing*, pp. 331–338, 2005.
- Bisk, Y., Zellers, R., Bras, R. L., Gao, J., and Choi, Y. Piqa: Reasoning about physical commonsense in natural language. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- Biten, A. F., Tito, R., Mafla, A., Gomez, L., Rusinol, M., Valveny, E., Jawahar, C., and Karatzas, D. Scene text visual question answering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4291–4301, 2019.
- Changpinyo, S., Sharma, P., Ding, N., and Soricut, R. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3558–3568, 2021.
- Chen, X., Fang, H., Lin, T.-Y., Vedantam, R., Gupta, S., Dollár, P., and Zitnick, C. L. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J. E., Stoica, I., and Xing, E. P. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, E., Wang, X., Dehghani, M., Brahma, S., et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- Driess, D., Xia, F., Sajjadi, M. S., Lynch, C., Chowdhery, A., Ichter, B., Wahid, A., Tompson, J., Vuong, Q., Yu, T., et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- Fan, A., Lewis, M., and Dauphin, Y. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 889–898, Melbourne, Australia,

- July 2018. Association for Computational Linguistics. URL <https://aclanthology.org/P18-1082>.
- Geva, M., Khashabi, D., Segal, E., Khot, T., Roth, D., and Berant, J. Did Aristotle Use a Laptop? A Question Answering Benchmark with Implicit Reasoning Strategies. *Transactions of the Association for Computational Linguistics (ACL)*, 2021.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Jin, W., Lee, D.-H., Zhu, C., Pujara, J., and Ren, X. Leveraging visual knowledge in language tasks: An empirical study on intermediate pre-training for cross-modal knowledge transfer. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2750–2762, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- Kazemnejad, A., Rezagholizadeh, M., Parthasarathi, P., and Chandar, S. Measuring the knowledge acquisition-utilization gap in pretrained language models. *arXiv preprint arXiv:2305.14775*, 2023.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Koh, J. Y., Salakhutdinov, R., and Fried, D. Grounding language models to images for multimodal generation. *arXiv preprint arXiv:2301.13823*, 2023.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *ICML*, 2023a.
- Li, X., Yin, X., Li, C., Zhang, P., Hu, X., Zhang, L., Wang, L., Hu, H., Dong, L., Wei, F., et al. Oscar: Object-semantic aligned pre-training for vision-language tasks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*, pp. 121–137. Springer, 2020.
- Li, Y., Hu, B., Chen, X., Ding, Y., Ma, L., and Zhang, M. A multi-modal context reasoning approach for conditional inference on joint textual and visual clues. *ACL*, 2023b.
- Li, Y., Hu, B., Chen, X., Ma, L., and Zhang, M. Lmeyer: An interactive perception network for large language models. *arXiv preprint arXiv:2305.03701*, 2023c.
- Li, Y., Wang, L., Hu, B., Chen, X., Zhong, W., Lyu, C., and Zhang, M. A comprehensive evaluation of gpt-4v on knowledge-intensive visual question answering. *arXiv preprint arXiv:2311.07536*, 2023d.
- Lin, B. Y., Wu, Z., Yang, Y., Lee, D.-H., and Ren, X. RiddleSense: Reasoning about riddle questions featuring linguistic creativity and commonsense knowledge. *arXiv preprint arXiv:2101.00376*, 2021.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023a.
- Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023b.
- Liu, Y., Li, Z., Li, H., Yu, W., Huang, M., Peng, D., Liu, M., Chen, M., Li, C., Jin, L., et al. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023c.
- Long, Q., Wang, M., and Li, L. Generative imagination elevates machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5738–5748, Online, June 2021. Association for Computational Linguistics. URL <https://aclanthology.org/2021.naacl-main.457>.
- Lu, Y., Zhu, W., Wang, X. E., Eckstein, M., and Wang, W. Y. Imagination-augmented natural language understanding. *NACCL*, 2022.
- Marino, K., Rastegari, M., Farhadi, A., and Mottaghi, R. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pp. 3195–3204, 2019.
- Mathew, M., Karatzas, D., and Jawahar, C. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 2200–2209, 2021.
- Merullo, J., Castricato, L., Eickhoff, C., and Pavlick, E. Linearly mapping from image to text space. *ICLR*, 2023.
- Mihaylov, T., Clark, P., Khot, T., and Sabharwal, A. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.

- Mishra, A., Shekhar, S., Singh, A. K., and Chakraborty, A. Ocr-vqa: Visual question answering by reading text in images. In *ICDAR*, 2019.
- OpenAI. Gpt-4 technical report. <https://arxiv.org/abs/2303.08774>, 2023.
- Penedo, G., Malartic, Q., Hesslow, D., Cojocaru, R., Cappelli, A., Alobeidli, H., Pannier, B., Almazrouei, E., and Launay, J. The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*, 2023. URL <https://arxiv.org/abs/2306.01116>.
- Plummer, B. A., Wang, L., Cervantes, C. M., Caicedo, J. C., Hockenmaier, J., and Lazebnik, S. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pp. 2641–2649, 2015.
- Sap, M., Rashkin, H., Chen, D., LeBras, R., and Choi, Y. Socialiqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., and Komatsuzaki, A. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- Shi, H., Mao, J., Gimpel, K., and Livescu, K. Visually grounded neural syntax acquisition. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1842–1861, Florence, Italy, July 2019. Association for Computational Linguistics. URL <https://aclanthology.org/P19-1180>.
- Singh, A., Natarjan, V., Shah, M., Jiang, Y., Chen, X., Batra, D., Parikh, D., and Rohrbach, M. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8317–8326, 2019.
- Talmor, A., Herzig, J., Lourie, N., and Berant, J. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421>.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Wang, P., Yang, A., Men, R., Lin, J., Bai, S., Li, Z., Ma, J., Zhou, C., Zhou, J., and Yang, H. Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. *arXiv preprint arXiv:2202.03052*, 2022a.
- Wang, X., Wen, K., Zhang, Z., Hou, L., Liu, Z., and Li, J. Finding skill neurons in pre-trained transformer-based language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 11132–11152, Abu Dhabi, United Arab Emirates, December 2022b. Association for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.765>.
- Wang, Y., Kordi, Y., Mishra, S., Liu, A., Smith, N. A., Khoshabi, D., and Hajishirzi, H. Self-instruct: Aligning language model with self generated instructions, 2022c.
- Xu, C., Sun, Q., Zheng, K., Geng, X., Zhao, P., Feng, J., Tao, C., and Jiang, D. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023.
- Yang, Z., Wu, W., Hu, H., Xu, C., Wang, W., and Li, Z. Open domain dialogue generation with latent images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 14239–14247, 2021.
- Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- Yin, S., Fu, C., Zhao, S., Li, K., Sun, X., Xu, T., and Chen, E. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023.
- Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. URL <https://aclanthology.org/P19-1472>.
- Zhang, Y., Zhang, R., Gu, J., Zhou, Y., Lipka, N., Yang, D., and Sun, T. Llavav: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023.
- Zhu, D., Chen, J., Shen, X., Li, X., and Elhoseiny, M. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

Zhu, W., Yan, A., Lu, Y., Xu, W., Wang, X. E., Eckstein, M., and Wang, W. Y. Visualize before you write: Imagination-guided open-ended text generation. *arXiv preprint arXiv:2210.03765*, 2022.