







Hierarchical Retrieval-Augmented Generation and Validation for Multimodal LLMs

D. Caffagni^{*}, F. Cocchi^{*}, N. Moratelli^{*}, S. Sarto^{*}, M. Cornia^{*}, L. Baraldi^{*}, R. Cucchiara[†],^{*}

* Università degli Studi di Modena e Reggio Emilia • Consiglio Nazionale delle Ricerche

Naples, 23/09/2024











Transversal Project on Vision, Language and Multimodal Challenges





















A Multimodal and Retrieval-Augmented Language Model

Main features

- The first Italian Multimodal Language Model, endowed with retrieval capabilities
- Retrieves and exploits multimodal knowledge from an external source
- Embeds input images via an MLP-based adapter
- Fine-tuned for Italian on translated visual conversation datasets

Proudly trained on the Leonardo Supercomputer thanks to an ISCRA-B grant.











A Multimodal and Retrieval-Augmented Language Model



- Long-tail concept understanding
- Visual reasoning capabilities
- Reading/OCR capabilities
- Q&A Capabilities











MORE: Multimodal Retrieval Augmented Vision-and-Language Mode	el	MORE: Multimod Augmented Vision-a	al Retrieval Ind-Language M	odel
Developed by AlmageLab (UNIMORE) within FAIR		Developed by AlmageLab (UNIM	ORE) within FAIR	
Che cosa è raffigurato in figura?	ia	Come si chiama il tipo di pasta raffig	urato in quest <mark>a</mark> immagine?	Invia
🖸 Rigenera 🛛 🕅 Pulisci		📓 Rigenera	🗑 Pulisci	









Hierarchical Retrieval for RAG

- Extending the model to incorporate **world-specific knowledge** (*e.g.* extracted from Wikipedia) and make the retrieval phase truly multimodal.
- We design a new model that integrates knowledge retrieved from an external knowledge base of documents through a **hierarchical retrieval pipeline**.

Downstream task:

- Knowledge-based VQA
 - For now, existing English benchmarks (*i.e.* Encyclopedic VQA and InfoSeek) are considered.













Hierarchical Retrieval for RAG

The visual encoder is employed to **provide the MLLM with visual context** and as a query to retrieve from an external knowledge base.











Hierarchical Retrieval for RAG

The visual encoder is employed to **provide the MLLM with visual context** and as a query to retrieve from an external knowledge base.











Hierarchical Retrieval for RAG

A **hierarchical retrieval module** is designed to first find the relevant document, using a similarity score between the CLIP-based embeddings extracted from the input image and the Wikipedia page title.











Hierarchical Retrieval for RAG

Then, the **most relevant passages** are retrieved inside the document computing similarities between **Contriever-based textual embeddings** extracted from each passage and the given question.











Hierarchical Retrieval for RAG

The retrieved passages are given as input to the multimodal LLM as additional input context, allowing the model to generate more specific answers.











Multimodal Self-reflection for RAG

- Effective strategies are needed to manage retrieved items and to improve CLIP-based models, which perform poorly in retrieving the most relevant document related to a given image.
- **Current focus:** Integration of **self-reflection and re-ranking techniques** inside the M-LLM to:
 - Decide if global visual features are sufficient or if fine-grained visual features are needed ("[FG]" token).
 - Decide when retrieval is needed, through the emission of a "[RET]" dedicated token.
 - Verify whether the retrieved knowledge is relevant or not to the given question (-> re-ranking of retrieved items).











Better Embedding spaces for RAG

- Most embedding spaces for multimodal RAG (i.e. CLIP) consider single-modality queries and values (e.g. images or text), limiting their encoding capabilities.
- **Current focus:** Design of embedding spaces for **RAG which support multimodal queries and documents** (e.g. image + question):
 - Textual features from the question guide the extraction of fine-grained features from the input image
 - Images are fused into the text of external documents, creating multi-modal retrievable items
 - Fusion between different modalities is done layer-wise and with learnable gates.











Evaluation of Multimodal LLMs

- Aligning machine-generated outputs with human judgment is of crucial importance.
- Previous work: training global image-text embedding spaces dedicated to evaluation.
- Current work: **BRIDGE**, a new learnable and reference-free image captioning metric that:
 - considers fine-grained visual features and maps them into dense multimodal vectors
 - employs multi-modal pseudo-captions built during the evaluation process to measure image-text correspondences at the detail level.
- SoTA on five human evaluation datasets.











Trustworthiness and Safety

- Models trained on large-scale data can generate inappropriate content and lead to the development of unsafe behavior.
- We aim to **make Vision-and-Language models safer** by removing their sensitivity to NSFW concepts.
- By fine-tuning Llama2 to convert between safe and unsafe sentences, we created our ViSU dataset made of quadruplets of real safe text/image pairs and generated unsafe pairs.
- We then employ the ViSU dataset to **fine-tune the CLIP space by redirecting unsafe content while preserving the CLIP structure**. Applying our method to downstream tasks like retrieval and generation improves the safety of retrieved and generated content by enhancing vulnerable user's protection.
- Current focus: hyperbolic spaces for safety preservation.











LLaVA-MORE: Enhancing Visual Instruction Tuning with LLaMA 3.1

- The first LLaVa-based family of models that integrates LLaMa 3.1 as core language model.
- We train and release both pre-training (V&L feature alignment) and instruction fine-tuning stages.
- With a variety of visual backbones and multi-resolution encoding methodologies (OpenAI CLIP, SigLIP, S2).
- Available on Github and Huggingface: https://github.com/aimagelab/LLaVA-MORE













Finally...

Check out our survey on Multimodal LLMs, that:

- Revises the prevalent choices for **vision encoders** and adapter modules that equip LLMs with cross-modal capabilities, and gives an overview of the **training process** and data used.
- Explores the **range of tasks** addressed by MLLMs, including visual grounding and image generation and editing.
- Offers a comparative perspective on the **performance** and hardware requirements of existing MLLMs.

The Revolution of Multimodal Large Language Models: A Survey

Davide Caffagni1*, Federico Cocchi1,2*, Luca Barsellotti1*, Nicholas Moratelli1*, Sara Sarto1*, Lorenzo Baraldi2*, Lorenzo Baraldi1, Marcella Cornia1, and Rita Cucchiara1,3 ¹University of Modena and Reggio Emilia, Italy ²University of Pisa, Italy ³IIT-CNR, Italy ¹{name.surname}@unimore.it ²{name.surname}@phd.unipi.it

Abstract

Connecting text and visual modalities plays an essential role in generative intelligence. For this reason, inspired by the success of large language models, significant research efforts are being devoted to the development of Multimodal Large Language Models (MLLMs). These models can seamlessly integrate visual and textual modalities, both as input and output, while providing a dialogue-based interface and instruction-following capabilities. In this paper, we provide a comprehensive review of recent visual-based MLLMs, analyzing their architectural choices, multimodal alignment strategies, and training techniques. We also conduct a detailed analysis of these models across a wide range of tasks, including visual grounding, image generation and editing, visual understanding, and domain-specific applications. Additionally, we compile and describe training datasets and evaluation benchmarks, conducting comparisons among existing models in terms of performance and computational requirements. Overall, this survey offers a comprehensive overview of the current state of the art, laying the groundwork for future MLLMs.

1 Introduction

embedding spaces (Radford et al., 2021).



Figure 1: General architecture of Multimodal Large Language Models (MLLMs), composed of a visual encoder, a language model, and an adapter module that connects visual inputs to the textual space.

to broaden the scope of these models to encompass multiple modalities, both as inputs and outputs. This expansion has led to the development of cutting-edge models such as GPT-4V (Achiam et al., 2023) and Gemini (Anil et al., 2023), showcasing state-of-the-art performance.

The development of Multimodal Large Lan-The introduction of the attention operation and the guage Models (MLLMs) entails merging single-Transformer architecture (Vaswani et al., 2017) has modality architectures for vision and language, enabled the creation of models capable of handling various modalities on an increasingly large scale. establishing effective connections between them This advancement is largely attributed to the verthrough vision-to-language adapters, and devising satility of the operator and the adaptability of the innovative training approaches. These methodolo architecture. Initially, this breakthrough was levergies are crucial for ensuring modality alignment aged for language-specific models (Devlin et al., and the ability to follow instructions accurately. 2018: Brown et al., 2020) but quickly extended to In a context marked by the rapid release of new support diverse modalities (Li et al., 2019; Lu et al., models, our goal is to offer an exhaustive overview 2019) and facilitate their integration within unified of the MLLM landscape, with a focus on models exploiting the visual modality. This overview The surge in sophisticated Large Language serves as both an update on the current state and a Models (LLMs), particularly their capacity for source of inspiration for future developments. We in-context learning, has encouraged researchers identify three core aspects that define these models:

D. Caffagni, F. Cocchi, L. Barsellotti, N. Moratelli, S. Sarto, L. Baraldi, L. Baraldi, M. Cornia, R. Cucchiara "The Revolution of Multimodal Large Language Models: A Survey" ACL, 2024 (Findings)