# VARCO-VISION: Expanding Frontiers in Korean Vision-Language Models

**Jeongho Ju**[*]    **Daeyoung Kim**[*]    **SunYoung Park**[*]    **Youngjune Kim**[*†]

NC Research, NCSOFT

{jeongho, daeyoungk, sun0park, youngjune}@ncsoft.com

## Abstract

In this paper, we introduce an open-source Korean-English vision-language model (VLM), VARCO-VISION. We incorporate a step-by-step training strategy that allows a model learn both linguistic and visual information while preserving the backbone model's knowledge. Our model demonstrates outstanding performance in diverse settings requiring bilingual image-text understanding and generation abilities compared to models of similar size. VARCO-VISION is also capable of grounding, referring, and OCR, expanding its usage and potential applications for real-world scenarios. In addition to the model, we release five Korean evaluation datasets, including four closed-set and one open-set benchmarks. We anticipate that our milestone will broaden the opportunities for AI researchers aiming to train VLMs. VARCO-VISION is available at https://huggingface.co/NCSOFT/VARCO-VISION-14B.

## 1  Introduction

Recent advancements in Large Language Models (LLMs) have increased attention on handling multimodality, providing robust backbones for Vision-Language Models (VLMs). Incorporating high-performing LLMs in VLMs demonstrated significant improvements across various visual tasks requiring text understanding, reasoning, and generation abilities [1, 2, 4, 6, 7, 14, 16, 20, 24, 25, 27, 30, 31, 33, 34]. The development of Multimodal Large Language Models (MLLMs) can further widen the usage of AI models and enhance user experiences to a great extent. Accordingly, the AI community, including both academia and industry, is committing substantial time and resources to train MLLMs [18, 39, 41] and establish evaluation frameworks [8, 38].

Although numerous multimodal models and benchmark datasets are being rapidly developed, they focus primarily on major languages such as English and Chinese [1, 6, 7, 14, 16, 19, 24, 34]. On the other hand, only a handful of open-source and commercial MLLMs are available for low-resource languages. This may cause a heavy reliance on proprietary model APIs for users, instead of fostering a research environment. As of now, even in South Korea where a huge AI community exists, there is a limited selection of Korean-supported models and datasets. While open-source Korean datasets for simple vision-text tasks like Visual Question Answering (VQA) or Optical Character Recognition (OCR) can be found [12, 13], assessing the models' general performance remains challenging.

In this work, we present (1) a strong English-Korean VLM called VARCO-VISION-14B (VARCO-VISION for short) and (2) five Korean benchmarks. VARCO-VISION is trained for four distinct phases with the final preference optimization stage. To evaluate the model's overall comprehension and generation abilities, we translated three closed-set (MMBench [22], SEED [17], MMStar [5]) and one open-set (LLaVA-W [21]) English benchmarks, and human-validated the datasets to ensure quality.

---

[*]These authors contributed equally to this work. Authors are listed in alphabetical order.

[†]Corresponding author

Figure 1: VARCO-VISION Application Examples: Visual Question Answering (VQA), Optical Character Recognition (OCR), Referring, and Grounding. Our model excels at both Korean/English vision-text and text-only tasks. Please see B for more detailed examples.

The Korean closed-set benchmarks, K-MMBench, K-SEED, and K-MMStar, are multiple-choice question answering (MCQA) tasks, which allow objective evaluation of MLLMs. K-LLaVA-W is based on LLaVA-Bench-in-the-wild (LLaVA-W) dataset with LLM-based automatic evaluation, and measures the Korean generation skill of a model compared to a GPT [2] model. Aside from these four benchmarks, we introduce a novel closed-set Korean benchmark (K-DTCBench) on understanding documents, tables, and charts. This benchmark is designed from scratch to specifically assess the capability of VLMs to process diverse image types.

We report performances of our own model as well as other baselines on English and Korean benchmarks to test the models' bilingual proficiency and multimodal capabilities. The models were also evaluated on text-only benchmarks, as one of our primary goals is to develop a VLM with strong language proficiency. In an overall quantitative evaluation encompassing both multimodal and text-only schemes, VARCO-VISION-14B not only surpasses other models of similar size in performance but also achieves scores comparable to proprietary models. We observe that the preference optimization phase significantly increases the readability and fluency of model outputs, leading to improved performance in K-LLaVA-W and text-only benchmarks. Furthermore, our model demonstrates proficient grounding and referring capabilities, which indicate its substantial potential for practical application.

By releasing VARCO-VISION and five Korean vision-text benchmarks, we are looking forward to promoting a more open AI community and expanding opportunities for researchers. Our contributions can be summarized as follows:

1. **English-Korean Bilingual Model**: We release a powerful 14B bilingual vision-language model, VARCO-VISION, that outperforms other models of similar scale. Despite being a VLM, it achieves high scores on language benchmarks in both Korean and English, demonstrating strong language capabilities.

2. **High-quality Korean Multimodal Benchmarks**: Based on widely recognized English multimodal benchmarks, we devise four closed-set and one open-set Korean benchmarks to evaluate VLMs' bilingual proficiency.

3. **Gradual Four-step Training**: We train our model in four stages with different objectives, so that it can absorb visual and linguistic capabilities progressively without losing the pre-trained backbone models' prior knowledge. As a result, our model exhibits outstanding performance across various benchmarks.

4. **Grounding, Referring, and OCR**: VARCO-VISION is capable of grounding, referring, and OCR tasks in both Korean and English, showing its high potentials for applications in real-world scenarios.

## 2 Training

### 2.1 Model Architecture

VARCO-VISION-14B consists of three main components: a vision encoder, a projector, and a Large Language Model (LLM). We leverage Qwen-2.5-14B-Instruct [32] as the language foundation model and SigLIP [37] as the vision encoder. The overall model architecture and visual representation processing method follow LLaVA-OneVision [16]. In this work, we concentrate on training VARCO-VISION with single-image examples.

Special tokens are added in the tokenizer for specific usages, such as OCR, grounding, and referring. The added special tokens are:

- **<gro>** for grounding tasks
- **<ocr>** for OCR tasks
- **<char> </char>** for indicating a text phrase
- **<obj> </obj>** for indicating an object
- **** for representing a bounding box
- **<delim>** for representing multiple location points for one object or text

The specific examples for each of these tokens are illustrated in the Appendix B.

### 2.2 Training Strategy

Our training pipeline aims to teach the model to gradually absorb and integrate both visual and linguistic understanding capabilities. Extending LLaVA-OneVision's training framework, we train our model in four stages. To preserve the model's language proficiency throughout the training process, we incorporate text-only data from Stages 2 to 4.

**Stage 1. Feature Alignment Pre-training**: We optimize the randomly initialized MLP projection layers while keeping other components frozen. This pre-training stage lets the model learn the mapping between the vision encoder and LLM. We use image-caption pairs as the training dataset to facilitate basic alignment between the two modalities.

**Stage 2. Basic Supervised Fine-tuning**: All model layers are fully trained on six different tasks: basic instruction-following, OCR, grounding/referring, caption description, document/table/chart/mathematical contents, and text-only examples. The model can acquire fundamental vision-language capabilities from Stage 2.

**Stage 3. Advanced Supervised Fine-tuning**: Advanced supervised fine-tuning is similar to Stage 2, but differs in that training tasks demand much more complex problem-solving abilities. The model is trained to enhance its reasoning and instruction-following skills across a range of tasks, from detailed image analysis to multi-step reasoning.

**Stage 4. Preference Optimization**: During the final stage, we focus exclusively on training the LLM layers to improve response alignment and generation capabilities. By applying Direct Preference Optimization (DPO) [28], we refine various aspects of the model responses—including but not limited to consistency, safety, and task-specific performance.

## 3 Evaluation

### 3.1 Korean Evaluation Benchmarks

In this section, we explain how we devise five Korean multimodal benchmarks. This paper is the first to release open-source Korean evaluation benchmarks for general comprehension and Korean generation capabilities. Based on the widely recognized benchmarks, MMBench [22], SEED [17], MMStar [5], and LLaVA-W [21], we translate the datasets into Korean and curate the translated outputs to create high-quality benchmarks. K-DTCBench is a Korean benchmark with documents, tables, and charts that we constructed from scratch. Four of the datasets (K-MMBench, K-SEED, K-MMStar, and

K-DTCBench) are closed-set multiple-choice question answering tasks, and K-LLaVA-W is an open-set freeform answer generation task with automatic LLM evaluation.

### 3.1.1 Closed-set Dataset

For VLM evaluation, it is common to employ multiple-choice question answering benchmarks to test fundamental abilities in processing image-text information. The three selected English benchmarks (MMBench, SEED, and MMStar) are composed of various evaluation dimensions and were thoroughly curated, making them suitable candidates for reference benchmarks. When using these benchmarks, we utilize the latest GPT API [3] to translate the datasets and enhance the outputs with the help of human annotators. Post-editing was necessary to improve the fluency and accuracy of translation due to translationese and localization issues.

- **K-MMBench**[4]: The original English MMBench is comprised of 20 ability dimensions, such as identity reasoning, image emotion, and attribute recognition. We use questions and ground truth answers from the dev subset of English MMBench.

- **K-SEED**[5]: SEED-Bench consists of images and videos, and evaluate VLMs with regard to 12 dimensions. From English SEED-Bench, we select the first 20 percent of questions requiring images.

- **K-MMStar**[6]: MMStar is a vision-indispensable benchmark of 1500 image-oriented questions. We observe that there are unanswerable cases (e.g., multiple images required to answer the question but only have a single image, vague questions or options) in the original MMStar dataset. Thus, we modify or re-create the questions to ensure that they can be answered within a single image. The examples of K-MMStar can be found in the Appendix A.1.

- **K-DTCBench**[7]: K-DTCBench is a newly developed benchmark featuring both computer-generated and handwritten images of three different types (documents, tables, and charts), all written in Korean. It consists of 80 questions for each image type and two questions per image, summing up to 240 questions in total. This benchmark is designed to evaluate whether VLMs can process images in different formats and be applicable to diverse domains. All images are generated with made-up values and statements for evaluation purposes only. We scanned handwritten documents/tables/charts, or created digital objects with matplotlib library to build K-DTCBench. The proportions of digital and handwritten images are equal, each constituting 50%.

### 3.1.2 Open-set Dataset

While it is important for VLMs to predict correct answers for given questions, generating fluent outputs is also a significant task for models. However, to the best of our knowledge, there are currently no Korean benchmarks available for assessing the quality of models' generation capability. Therefore, we adopt LLaVA-Bench-in-the-wild (LLaVA-W) benchmark and translate-validate the benchmark as we did for the closed-set benchmarks.

LLaVA-W contains 24 images of various domains and 60 daily-life questions. Since our goal was to build a benchmark exclusively focused on Korean, we change the English texts in images into Korean for localization. Figure 4 shows the examples of LLaVA-W and K-LLaVA-W[8]. Given an image and a question related to the image, models need to generate open-ended answers. The captions of the images are only used during evaluation.

For evaluation, our benchmark follows the pipeline of LLaVA-W's LLM automatic evaluation [9], but with a little twist in the JudgeLLM prompts. We convert English prompts into Korean and added instruction details as shown in the Appendix A.3.2. Based on the provided caption, question, and

---

[3] gpt-4o-2024-08-06
[4] https://huggingface.co/datasets/NCSOFT/K-MMBench
[5] https://huggingface.co/datasets/NCSOFT/K-SEED
[6] https://huggingface.co/datasets/NCSOFT/K-MMStar
[7] https://huggingface.co/datasets/NCSOFT/K-DTCBench
[8] https://huggingface.co/datasets/NCSOFT/K-LLaVA-W
[9] https://github.com/EvolvingLMMs-Lab/lmms-eval/tree/main/lmms_eval/tasks/llava-in-the-wild

model's generated output, JudgeLLM measures the model's helpfulness, relevance, accuracy, level of detail, and Korean generation capability. The final K-LLaVA-W score of the target model is calculated as the ratio of the target model's JudgeLLM score to the baseline model's JudgeLLM score.[10]

## 3.2 Benchmark Results

We leverage four types of benchmarks to fully evaluate VARCO-VISION's performance on diverse dimensions. To assess the model's ability to understand and generate both languages, ten benchmarks—five each for Korean and English—are used for evaluation. We also use text-only benchmarks to test the effectiveness of our training strategy to incorporate text datasets throughout Stages 2 to 4. In addition to comparing with other VLMs, we evaluate our model on OCR tasks and compare its performance with that of OCR-focused models.

### 3.2.1 Korean Benchmarks

| | KOREAN BENCHMARKS | | | | |
|---|---|---|---|---|---|
| **Model** | **Image Understanding (MCQA)** | | | | **Generation (Image-based)** |
| | K-MMB (dev) | K-SEED | K-MMSTAR | K-DTCBench | K-LLaVA-W |
| **VARCO-VISION-14B** | **82.21** | **75.39** | **57.33** | **84.58** | **84.74** |
| **Pangea-7B [36]** | 71.64 | 73.34 | 35.00 | 48.33 | 69.70 |
| **Pixtral-12B [3]** | 57.47 | 46.44 | 23.93 | 27.50 | 82.00 |
| **Molmo-7B-D [7]** | 63.83 | 69.53 | 47.40 | 45.83 | 63.90 |
| **Qwen2-VL-7B-Instruct [33]** | 78.26 | 74.08 | 50.67 | 75.00 | 62.00 |
| **LLaVA-OneVision-7B [16]** | 76.28 | 73.21 | 54.00 | 52.91 | 48.80 |
| Qwen2-VL-72B-Instruct [33] | 84.27 | 78.25 | 63.53 | 88.75 | 97.40 |
| LLaVA-OneVision-72B [16] | 88.01 | 77.86 | 62.66 | 60.83 | 84.10 |
| GPT-4o-mini [2] | 74.48 | 73.30 | 42.33 | 74.58 | 101.90 |
| GPT-4V [2] | 77.92 | 71.66 | 35.20 | 47.50 | 98.90 |
| GPT-4o [2] | 81.01 | 76.98 | 56.20 | 85.80 | 103.90 |

Table 1: Model Comparison on Korean Benchmarks. The models in the first upper block are open-source models with similar scale, and the second block are open-source 72B models. The last block shows the performance of proprietary GPT models. We primarily compare VARCO-VISION with the models mentioned in the first block, as they are similar in size to our model.

As a model trained with a strong emphasis on Korean linguistic ability, VARCO-VISION excels in all MCQA benchmarks compared to the models with similar size. In general comprehension tasks like K-MMBench and K-SEED, Qwen2-VL-7B-Instruct [33] and LLaVA-OneVision-7B [16] are next to our model, showing a slightly lower performance overall. While models have similar scores in K-MMBench, K-SEED, and K-MMStar, we observe significant variation in the performance of models in K-DTCBench while VARCO-VISION achieving dominant performance. This suggests that Korean documents, tables, and charts were not adequately trained by other open-source models. Moreover, VARCO-VISION reaches competitive performance to large-scale open-source and proprietary models, gaining higher scores than GPT-4o-mini and GPT-4V for MCQA tasks.

In K-LLaVA-W where fluent Korean generation is a key value, only two models among models under 20B, VARCO-VISION-14B and Pixtral-12B [3], obtain score over 80. On the other hand, Qwen2-VL-72B-Instruct [33] and GPT models achieve scores around 100, implying that model scale may affect the response quality in longer text generation.

---

[10]We used gpt-4o-2024-08-06 for both JudgeLLM and the baseline model. The original LLaVA-W provides GPT responses. However, instead of translating the responses in LLaVA-W, we ran the latest GPT model to obtain high-quality responses.

### 3.2.2 English Benchmarks

| ENGLISH BENCHMARKS | | | | | |
|---|---|---|---|---|---|
| **Model** | **Image Understanding (MCQA)** | | | | **OCR** |
| | MMBv1.1 (dev) | SEED (image) | MMStar (val) | MMMU (val) | OCRBench (test) |
| **VARCO-VISION-14B** | **84.28** | **76.73** | **63.33** | <u>51.33</u> | <u>820</u> |
| **Pangea-7B [36]** | 76.23 | 74.88 | 43.26 | 43.55 | 620 |
| **Pixtral-12B [3]** | 72.98 | 74.34 | 48.33 | 49.00 | 682 |
| **Molmo-7B-D [7]** | 72.05 | 74.36 | 52.73 | 45.30 | 708 |
| **Qwen2-VL-7B-Instruct [33]** | <u>80.95</u> | <u>76.45</u> | 60.00 | **54.10** | **866** |
| **LLaVA-OneVision-7B [16]** | 80.80 | 76.41 | <u>61.33</u> | 47.67 | 630 |
| Qwen2-VL-72B-Instruct [33] | 86.91 | 77.86 | 67.60 | 56.66 | 877 |
| LLaVA-OneVision-72B [16] | 85.44 | 77.43 | 65.33 | 56.80 | 741 |
| GPT-4o-mini [2] | 76.31 | 72.80 | 54.80 | 60.00 | 785 |
| GPT-4V [2] | 79.41 | 73.00 | 56.00 | 62.30 | 656 |
| GPT-4o [2] | 81.73 | 76.70 | 64.70 | 69.90 | 805 |

Table 2: Model Comparison on English Benchmarks. MMBench [22], SEED [17], MMStar [5], and MMMU [35] are multi-choice question answering tasks. MMBench and SEED are for comprehension evaluation, whereas MMStar is focused more on vision-indispensible reasoning. MMMU tests college-level subject knowledge of VLMs. OCRBench [23] is a specialized benchmark in OCR for VLMs, composed of 1000 question-answer pairs. The values in OCRBench indicate the number of questions correctly answered by models.

Our goal in training lies in boosting Korean and English proficiency, thus evaluating on English benchmarks was necessary to investigate our model's performance. In English MCQA benchmarks, VARCO-VISION gains higher performances over other models under 20B in all benchmarks, except for MMMU. However, we notice that all models generally demonstrate sufficient performance in English benchmarks in contrast to their performances in Korean benchmarks. The results suggest that training schemes for the majority of the models in Table 2 prioritized learning English. In OCRBench, our model shows superior performance over other VLMs, spotlighting the effectiveness of Stages 2 and 3, where the model is exposed to OCR, grounding, and referring tasks.

### 3.2.3 Text-only Benchmarks

| TEXT-ONLY BENCHMARKS | | | |
|---|---|---|---|
| **Model** | **Korean** | | **English** |
| | LogicKor | KoMT-Bench | MT-Bench |
| **VARCO-VISION-14B** | **8.69** | **8.39** | **8.80** |
| **Pangea-7B [36]** | 5.06 | 5.06 | 7.29 |
| **Pixtral-12B [3]** | <u>7.71</u> | <u>8.11</u> | <u>8.40</u> |
| **Molmo-7B-D [7]** | 2.64 | 3.58 | 6.93 |
| **Qwen2-VL-7B-Instruct [33]** | 4.62 | 4.54 | 7.13 |
| **LLaVA-OneVision-7B [16]** | 2.23 | 3.52 | 7.52 |
| Qwen2-VL-72B-Instruct [33] | 7.74 | 7.49 | 8.53 |
| LLaVA-OneVision-72B [16] | 8.22 | 7.87 | 8.78 |
| EXAONE 3.0 7.8B Inst.(LLM) [29] | 8.62 | 8.92 | 9.01 |
| GPT-4o-mini [2] | 9.14 | 8.88 | 9.09 |
| GPT-4V [2] | 8.66 | 9.25 | 9.41 |
| GPT-4o [2] | 9.57 | 9.24 | 9.30 |

Table 3: Model Performance on Korean and English Text-only Benchmarks. MT-Bench [40] is an English multi-turn dialogue benchmark, and KoMT-Bench [29] is built by translating MT-Bench. LogicKor[11] consists of multi-turn Korean dialogues across six categories.

We evaluate the models on text-only benchmarks to investigate whether our strategy of including text-only datasets during training stages contributed to textual understanding of VARCO-VISION. Two Korean language benchmarks (LogicKor[11] and KoMT-Bench [29]) and one English benchmark (MT-Bench [40]) are employed for text-only evaluation. Our model outperforms other VLMs in text-only benchmarks, even when compared to 72B open-source models. We believe that applying preference optimization in the final training phase boosted the overall quality of model responses, resulting in outstanding performance in all three text-only benchmarks. VARCO-VISION achieves scores comparable to that of EXAONE 3.0 7.8B, which is a Korean-English language model only trained on textual inputs. It reflects our model's capability to handle textual inputs as much as other bilingual language models.

### 3.2.4 OCR Benchmarks

| Model | OCR | | | Average |
|---|---|---|---|---|
| | CORD | ICDAR2013 | ICDAR2015 | |
| **VARCO-VISION-14B** | 82.69 | **94.42** | 72.95 | 83.35 |
| **EasyOCR**[12] | 79.56 | 84.97 | 57.90 | 74.14 |
| **Pororo [9]** | 78.73 | 84.29 | 64.65 | 75.89 |
| **PaddleOCR**[13] | **92.71** | 92.01 | **73.73** | **86.15** |
| CLOVA OCR[14] | 95.32 | 94.39 | 84.06 | 91.26 |

Table 4: OCR Benchmark Performance. EasyOCR, Pororo, and PaddleOCR are open-source OCR models. CLOVA OCR is a proprietary OCR model. PopEval [15] was used as the metric for all benchmarks.

We compare our model's OCR ability to other well-known OCR models using the PopEval [15] metric. In contrast to OCRBench, CORD [26], ICDAR2013 [10], and ICDAR2015 [11] are more complex tasks that require models to generate both textual content and spatial locations of texts within images. Despite the fact that other models in Table 4 are OCR expert models while VARCO-VISION is not, its performances across these benchmarks are remarkable. We find that VARCO-VISION holds the potential for use in a wide range of applications, not limited to specific tasks.

## 4 Discussion and Future Work

Throughout the process of model training and evaluation, we notice that benchmark performances do not fully reflect a model's true capability. MMBench[22] shows the possibility of choice biases in MLLMs that they tend to prefer a certain choice in MCQA tasks, and apply Circular Evaluation to regularize the problem. Nevertheless, the majority of multi-choice question answering tasks still follow naive evaluation techniques. In addition, there are a relatively small number of benchmarks aiming to evaluate the model's generation capability compared to benchmarks with short answers [19]. Although a model might achieve high performance in MCQA tasks, it may produce low-quality answers for tasks that require long responses. We seek significant future developments remain to be made in terms of MLLM benchmarks.

In this work, we focus on training VARCO-VISION on single-image scenarios, and the evaluation results show that our model excels in single-image benchmarks. However, MLLMs need to process multi-image scenarios (including videos) and sounds to be applicable in various domains. We plan to expand VARCO-VISION's modalities to video and audio in the near future, and broaden the horizon of our model's usage. Besides modality expansion, we are currently training an advanced, localized model with diverse images taken in Korea. With the improved capability to understand Korean culture,

---

[11]https://lk.instruct.kr
[12]https://github.com/JaidedAI/EasyOCR
[13]https://github.com/PaddlePaddle/PaddleOCR
[14]https://www.ncloud.com/product/aiService/ocr

the future model is expected to be applied in real-world tasks, such as multimodal search, multimodal retrieval-augmented generation, and visual agents.

## 5   Conclusion

We present an open-source Korean-English vision-language model, VARCO-VISION-14B, and five Korean benchmarks. Although extensive research and numerous models have been developed for MLLMs, this work is the first to release both bilingual VLM and benchmarks supporting Korean. Our model achieves remarkable results in both Korean and English benchmarks among other open-source models of similar scale. VARCO-VISION does not only excel in vision-text benchmarks, but also in text-only and OCR benchmarks. We find that this milestone may expand opportunities for many researchers and enable breakthroughs in training bilingual VLMs.

## Acknowledgements

## References

[1] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, 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] Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, et al. Pixtral 12b. *arXiv preprint arXiv:2410.07073*, 2024.

[4] Anthropic. Claude 3.5 sonnet, 2024. URL https://www.anthropic.com/news/claude-3-5-sonnet.

[5] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and Feng Zhao. Are we on the right way for evaluating large vision-language models? In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=evP9mxNNxJ.

[6] 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.

[7] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*, 2024.

[8] Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models, 2024. URL https://arxiv.org/abs/2407.11691.

[9] Hoon Heo, Hyunwoong Ko, Soohwan Kim, Gunsoo Han, Jiwoo Park, and Kyubyong Park. Pororo: Platform of neural models for natural language processing. https://github.com/kakaobrain/pororo, 2021.

[10] Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluis Gomez i Bigorda, Sergi Robles Mestre, Joan Mas, David Fernandez Mota, Jon Almazan Almazan, and Lluis Pere De Las Heras. Icdar 2013 robust reading competition. In *2013 12th international conference on document analysis and recognition*, pages 1484–1493. IEEE, 2013.

[11] Dimosthenis Karatzas, Lluis Gomez-Bigorda, Anguelos Nicolaou, Suman Ghosh, Andrew Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian Lu, et al. Icdar 2015 competition on robust reading. In *2015 13th international conference on document analysis and recognition (ICDAR)*, pages 1156–1160. IEEE, 2015.

[12] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. In *European Conference on Computer Vision (ECCV)*, 2022.

[13] Jin-hwa Kim, Soohyun Lim, Jaesun Park, and Hansu Cho. Korean Localization of Visual Question Answering for Blind People. In *AI for Social Good workshop at NeurIPS*, 2019.

[14] Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and better understanding vision-language models: insights and future directions. *arXiv preprint arXiv:2408.12637*, 2024.

[15] Hong-Seok Lee, Youngmin Yoon, Pil Hoon Jang, and Chankyu Choi. Popeval: A character-level approach to end-to-end evaluation compatible with word-level benchmark dataset. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 1207–1213. IEEE, 2019.

[16] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.

[17] Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench: Benchmarking multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13299–13308, 2024.

[18] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*, 2024.

[19] Jian Li, Weiheng Lu, Hao Fei, Meng Luo, Ming Dai, Min Xia, Yizhang Jin, Zhenye Gan, Ding Qi, Chaoyou Fu, Ying Tai, Wankou Yang, Yabiao Wang, and Chengjie Wang. A survey on benchmarks of multimodal large language models, 2024. URL https://arxiv.org/abs/2408.08632.

[20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 34892–34916. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/6dcf277ea32ce3288914faf369fe6de0-Paper-Conference.pdf.

[21] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.

[22] 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? In *European Conference on Computer Vision*, pages 216–233. Springer, 2025.

[23] Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng Yin, Cheng-lin Liu, Lianwen Jin, and Xiang Bai. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023.

[24] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024.

[25] Gen Luo, Yiyi Zhou, Yuxin Zhang, Xiawu Zheng, Xiaoshuai Sun, and Rongrong Ji. Feast your eyes: Mixture-of-resolution adaptation for multimodal large language models, 2024. URL https://arxiv.org/abs/2403.03003.

[26] Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk Lee. Cord: A consolidated receipt dataset for post-ocr parsing. In *Document Intelligence Workshop at Neural Information Processing Systems*, 2019.

[27] 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, 2023. URL https://arxiv.org/abs/2306.14824.

[28] 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. *Advances in Neural Information Processing Systems*, 36, 2024.

[29] LG Research, Soyoung An, Kyunghoon Bae, Eunbi Choi, Stanley Jungkyu Choi, Yemuk Choi, Seokhee Hong, Yeonjung Hong, Junwon Hwang, Hyojin Jeon, et al. Exaone 3.0 7.8 b instruction tuned language model. *arXiv preprint arXiv:2408.03541*, 2024.

[30] DongJae Shin, HyeonSeok Lim, Inho Won, ChangSu Choi, Minjun Kim, SeungWoo Song, HanGyeol Yoo, SangMin Kim, and KyungTae Lim. X-LLaVA: Optimizing bilingual large vision-language alignment. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 2463–2473, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. findings-naacl.158. URL https://aclanthology.org/2024.findings-naacl.158.

[31] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

[32] Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.github.io/blog/qwen2.5/.

[33] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

[34] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024.

[35] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024.

[36] Xiang Yue, Yueqi Song, Akari Asai, Seungone Kim, Jean de Dieu Nyandwi, Simran Khanuja, Anjali Kantharuban, Lintang Sutawika, Sathyanarayanan Ramamoorthy, and Graham Neubig. Pangea: A fully open multilingual multimodal llm for 39 languages. *arXiv preprint arXiv:2410.16153*, 2024.

[37] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023.

[38] Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Lmms-eval: Reality check on the evaluation of large multimodal models, 2024. URL https://arxiv.org/abs/2407.12772.

[39] Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang, Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu, Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda Chen. Swift:a scalable lightweight infrastructure for fine-tuning, 2024. URL https://arxiv.org/abs/2408.05517.

[40] 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. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.

[41] Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024. Association for Computational Linguistics. URL `http://arxiv.org/abs/2403.13372`.

# A   Korean Multimodal Benchmarks

In this section, we provide examples of Korean vision-text benchmarks (K-MMStar, K-DTCBench, and K-LLaVA-W). Since these three benchmarks required more than simple translation, we elaborate on benchmark construction with specific examples.

## A.1   K-MMStar



| Type 1 | Type 2 | Type 3 |
| Simple Translation | Modification | Re-creation |

**MMStar**

Question:
Which option describe the object relationship in the image correctly?

Options:
A: The suitcase is on the book.
B: The suitcase is beneath the cat.
C: The suitcase is beneath the bed.
D: The suitcase is beneath the book.

Answer: A

**MMStar**

Question:
What is the predominant sport being played in the image?

Options:
A: football
B: basketball
C: soccer
D: baseball

Answer: A

**MMStar**

Question:
Which of the following graphs are not isomorphic?

Options:
A: <image 1>
B: <image 2>
C: <image 3>
D: nan

Answer: B

**K-MMStar**

Question:
이미지에서 물체들의 관계를
올바르게 설명하는 옵션은 무엇인가요?

Options:
A: 가방이 책 위에 있다.
B: 가방이 고양이 아래에 있다.
C: 가방이 침대 아래에 있다.
D: 가방이 책 아래에 있다.

Answer: A

**K-MMStar**

Question:
이미지에 표시된 주요 스포츠는 무엇입니까?

Options:
A: 축구
B: 농구
C: 미식축구
D: 야구

Answer: A

**K-MMStar**

Questions:
Graph 1에서 다음 중 이웃하는 두 꼭짓점 쌍이
아닌 것은 무엇입니까?

Options:
A: a와 h
B: b와 d
C: e와 f
D: f와 g

Answer: B

Figure 2: K-MMStar Example

K-MMStar has three different types of questions. We noticed that there are unanswerable or vague questions in the original MMStar, thus modified the question or created a new one. Figure 2 shows the examples of questions. Type 1 refers to cases where direct translations from English to Korean are sufficient for Korean version questions. In Type 2, the question from MMStar asks about which sport is being played in the image. However, the options include both 'football' and 'soccer,' either of which can be correct depending on British English or American English conventions. Therefore, we changed the options to "축구" (soccer) and "미식축구" (American football) for clarity. In Type 3 example, it has <image 1>, <image 2>, and <image 3> in the options but MMStar did not provide images corresponding to the options. Hence, we re-created a new question about the image.

## A.2   K-DTCBench

K-DTCBench is developed from scratch with synthetic images. Each image has a question and four options to choose from.

Figure 3: K-DTCBench Example

## A.3 K-LLaVA-W

### A.3.1 Example

In K-LLaVA-W, we changed the English text into Korean text for images with texts. If an original LLaVA-W image did not contain any text, we left it unchanged to preserve its authenticity. In Figure 4, we changed "MONDAY. JUST... MONDAY" into "월요일. 단지... 월요일".



Figure 4: K-LLaVA-W Example

### A.3.2 K-LLaVA-W JudgeLLM Prompt

[설명]
{caption}

[질문]
{question}

[어시스턴트 1]
{gpt_answer}
[어시스턴트 1 끝]

[어시스턴트 2]
{target_model_answer}
[어시스턴트 2 끝]

[System]
두 인공지능 어시스턴트의 성능을 [질문]에 대한 응답에 기반하여 평가하세요. 해당 [질문]은 특정 이미지를 보고 생성되었습니다. `유용성`, `관련성`, `정확성`, `세부 수준`, `한국어 생성능력`을 기준으로 응답을 평가하세요. 각각의 어시스턴트에게 1에서 10까지의 전반적인 점수를 부여하며, 높은 점수일수록 더 나은 전반적인 성능을 나타냅니다.

# 단계
1. 제공된 이미지 설명을 검토하세요.
2. 각 어시스턴트의 응답을 다음 기준으로 분석하세요:
   - `유용성`: 응답이 사용자의 질문을 얼마나 잘 해결하는가?
   - `관련성`: 응답이 사용자의 질문에 얼마나 적절한가?
   - `정확성`: 응답에서 제공한 정보가 얼마나 정확한가?
   - `세부 수준`: 응답이 과하지 않게 충분히 자세한가?
   - `한국어 생성능력`: 생성된 한국어 문장이 자연스럽고 문법적으로 올바른가?
3. 분석에 기반하여 각 어시스턴트에게 1에서 10까지의 점수를 부여하세요.
4. 두 점수를 공백으로 구분하여 한 줄로 제공하세요.
5. 점수에 대한 이유를 강조하면서 포괄적인 평가를 제공하고, 편견을 피하며 응답의 순서가 판단에 영향을 미치지 않도록 하세요.

# 출력 형식
- 첫 번째 줄: `어시스턴트1_점수 어시스턴트2_점수` (예: `8 9`)
- 두 번째 줄: `유용성`, `관련성`, `정확성`, `세부 수준`, `한국어 생성능력` 기준으로 점수를 설명하는 자세한 문단을 제공합니다.

# 주의사항
- 평가 시 잠재적 편견을 방지하여 객관성을 확보하세요.
- 분석과 설명에서 일관성과 명확성을 유지하세요.

Figure 5: K-LLaVA-W Evaluation Prompt. We translated the LLaVA-W prompts and added specific guidelines in the JudgeLLM prompt.
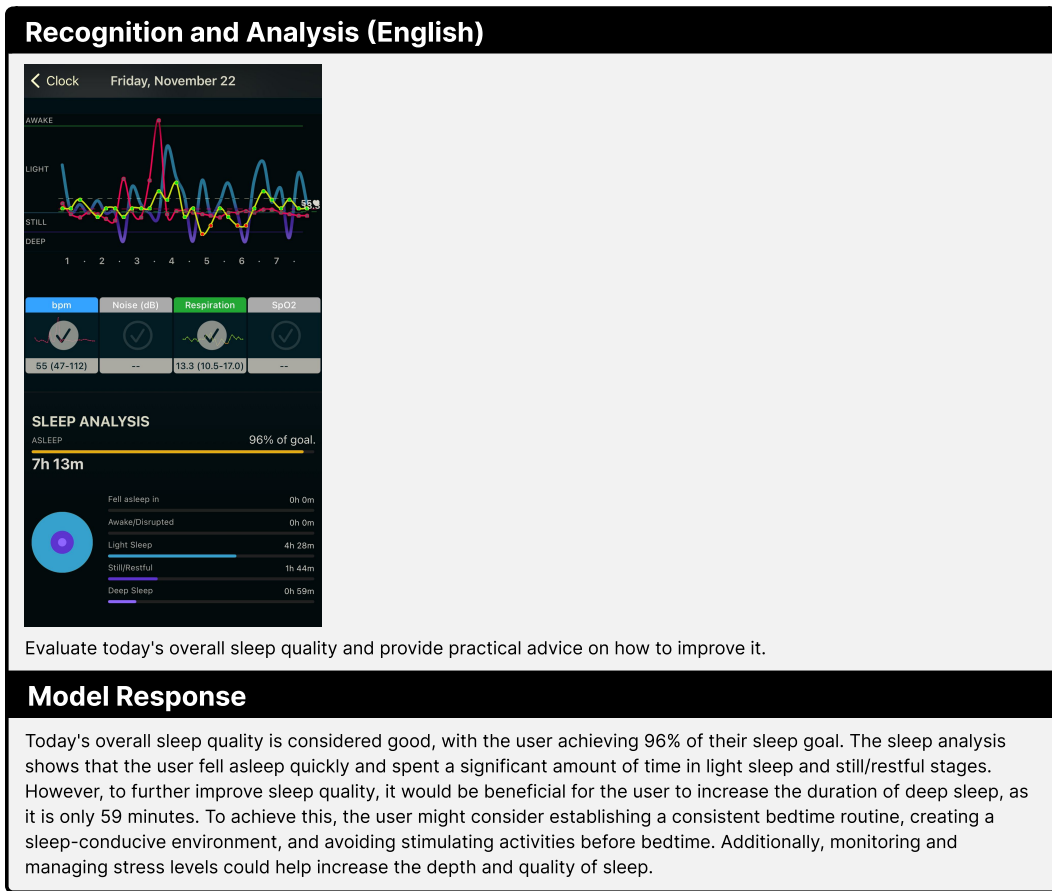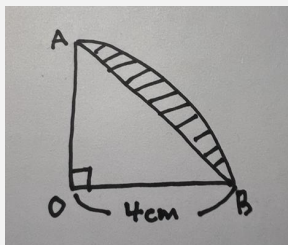
# B Application Examples



Figure 6: Text Recognition and Analysis Example (English)

## Mathematical Reasoning (Korean)



이 그림에서 색칠한 부분의 넓이를 구하세요.
(Translation: Calculate the area of the shaded region in this figure.)

## Model Response

이 문제를 해결하기 위해 다음 단계를 따라가겠습니다:

1. 전체 도형을 분석합니다:
- 중심각이 90°인 부채꼴
- 반지름이 4cm인 원의 일부

2. 부채꼴의 넓이를 계산합니다:
- 부채꼴의 면적 공식: (중심각/360) × π × r²
- 여기서 r = 4cm, 중심각 = 90°
- 계산: (90/360) × π × 4² = 4π cm²

3. 삼각형의 넓이를 계산합니다:
- 밑변 = 4cm, 높이 = 4cm
- 삼각형의 넓이 공식: (밑변 × 높이) / 2
- 계산: (4 × 4) / 2 = 8 cm²

4. 색칠한 부분의 넓이를 구합니다:
- 부채꼴 넓이 - 삼각형 넓이
- 4π cm² - 8 cm²

5. 최종 답을 계산합니다:
- 4π - 8 ≈ 4 × 3.14 - 8 = 12.56 - 8 = 4.56 cm²

따라서 색칠한 부분의 넓이는 약 4.56 cm²입니다.

To solve this problem, we will follow these steps:

1. Analyze the overall shape:
- A sector with a central angle of 90°
- Part of a circle with a radius of 4 cm

2. Calculate the area of the sector:
- Formula for the area of a sector: (central angle / 360) × π × r²
- Here, r = 4 cm, central angle = 90°
- Calculation: (90 / 360) × π × 4² = 4π cm²

3. Calculate the area of the triangle:
- Base = 4 cm, height = 4 cm
- Formula for the area of a triangle: (base × height) / 2
- Calculation: (4 × 4) / 2 = 8 cm²

4. Find the area of the shaded region:
- Sector area - Triangle area
- 4π cm² - 8 cm²

5. Calculate the final answer:
- 4π - 8 ≈ 4 × 3.14 - 8 = 12.56 - 8 = 4.56 cm²

Therefore, the area of the shaded region is approximately 4.56 cm².

Figure 7: Mathematical Reasoning Example (Korean)

## Information Extraction and Calculation (Korean)



치즈버거 한개와 탄산음료 한개를 주문했습니다. 총 얼마를 내야할까요? 화폐 단위는 천원 입니다.
(Translation: I ordered one cheeseburger and one soda. How much should I pay in total? The unit of currency is 1,000 won.)

### Model Response

총 금액을 계산하는 방법은 다음과 같습니다:

1. 치즈버거의 가격 확인:
- 치즈버거 (CHEESEBURGER) 가격: 14.9천원

2. 탄산음료의 가격 확인:
- 탄산음료 (SODA) 가격: 3.9천원

3. 두 메뉴의 가격 합산:
- 14.9천원 + 3.9천원 = 18.8천원

따라서, 치즈버거 한 개와 탄산음료 한 개를 주문할 경우 총 18.8천원을 내야 합니다.

Here's how to calculate the total amount:

1. Check the price of cheeseburger:
- CHEESEBURGER price: 14.9K won

2. Check the price of soda:
- SODA price: 3.9K won

3. Sum up the prices of both items:
-14.9K won + 3.9K won = 18.8K won

Therefore, when ordering one cheeseburger and one soda, you need to pay 18.8K won in total.

Figure 8: Information Extraction and Calculation Example (Korean)

## Grounding (Korean)

<gro>
이미지를 자세히 묘사해줘. (Translation: Describe the image in detail.)

### Model Response

이미지에는 <obj>아이</obj><bbox>0.431, 0.382, 0.597, 0.607</bbox>가 앉아 있는 <obj>카트</obj><bbox>0.022, 0.22, 0.626, 0.906</bbox>와 <obj>당나귀</obj><bbox>0.477, 0.413, 0.88, 0.974</bbox>가 있습니다. <obj>카트</obj><bbox>0.022, 0.22, 0.626, 0.906</bbox>는 나무로 만들어져 있으며, 햇빛을 피할 수 있도록 지붕이 달려 있습니다. <obj>아이</obj><bbox>0.431, 0.382, 0.597, 0.607</bbox>는 파란색 옷을 입고 카트에 앉아 있으며, 카트는 <obj>당나귀</obj><bbox>0.477, 0.413, 0.88, 0.974</bbox>에 의해 끌리고 있습니다. 배경에는 해질녘의 따뜻한 하늘과 멀리 보이는 건물들이 있습니다.

(Translation: In the image, there is a <obj>child</obj><bbox>0.431, 0.382, 0.597, 0.607</bbox> sitting in a <obj>cart</obj><bbox>0.022, 0.22, 0.626, 0.906</bbox> and a <obj>donkey</obj><bbox>0.477, 0.413, 0.88, 0.974</bbox>. The <obj>cart</obj><bbox>0.022, 0.22, 0.626, 0.906</bbox> is made of wood and has a roof to provide shade from the sunlight. The <obj>child</obj><bbox>0.431, 0.382, 0.597, 0.607</bbox> is wearing blue clothes and sitting in the cart, which is being pulled by the <obj>donkey</obj><bbox>0.477, 0.413, 0.88, 0.974</bbox>. In the background, there is a warm sunset sky and buildings visible in the distance.)
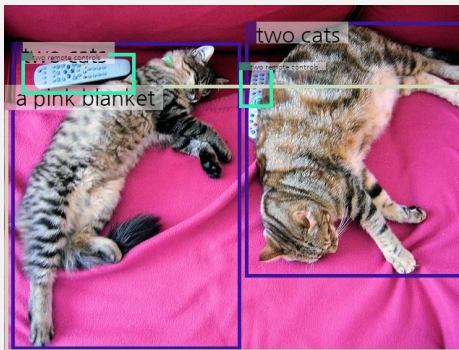
Figure 9: Grounding Example (Korean)

**Grounding (English)**

<gro>
Describe the image in detail.

**Model Response**

The image shows <obj>two cats</obj><bbox>0.521, 0.049, 0.997, 0.783<delim>0.016, 0.108, 0.512, 0.99</bbox> lying on <obj>a pink blanket</obj><bbox>0.002, 0.231, 0.999, 0.999</bbox>. The cat on the left is lying on its side with its head resting on the blanket and its body stretched out. The cat on the right is lying on its back with its paws stretched out and its head turned to the side. Both cats appear relaxed and comfortable. There are also <obj>two remote controls</obj><bbox>0.039, 0.138, 0.283, 0.257<delim>0.508, 0.166, 0.581, 0.295</bbox> placed near the cats, one on each side of them.

Figure 10: Grounding Example (English)

## Referring (Korean)



<char>이곳에</char><bbox>0.162, 0.076, 0.71, 0.275</bbox> 써있는 글자를 읽어주세요.
그리고 무엇을 하는 곳인지 설명해주세요.
(Translation: Read the letters <char>here</char><bbox>0.162, 0.076, 0.71, 0.275</bbox>.
And explain what the place does.)

## Model Response

이곳에 써있는 글자는 "STARBUCKS"입니다. 이는 세계적으로 유명한 커피 전문점 체인인 스타벅스의 로고입니다. 스타벅스는 커피, 차, 아이스드 음료와 같은 다양한 음료를 판매하며, 과자, 샌드위치 등의 간식과 함께 편안한 분위기에서 고객들이 쉴 수 있도록 공간을 제공합니다.
스타벅스는 한국뿐만 아니라 전 세계적으로 많은 매장을 운영하고 있으며, 친구들과 만나거나 혼자서 책을 읽거나 일할 수 있는 인기 있는 장소로 자리 잡았습니다. 이미지에서는 스타벅스 매장 밖을 보여주고 있으며, 유리창 너머로 매장 내부의 모습도 엿볼 수 있습니다.
(Translation: The text written here is "STARBUCKS". This is the logo of Starbucks, a globally famous coffee shop chain. Starbucks sells various beverages such as coffee, tea, and iced drinks, and provides a comfortable space where customers can relax with snacks like pastries and sandwiches.
Starbucks operates many stores not only in Korea but also worldwide, and has established itself as a popular place where people can meet friends or read books and work alone. The image shows the exterior of a Starbucks store, and through the glass windows, you can also glimpse the interior of the store.)

Figure 11: Referring Example (Korean)

**OCR**



<ocr>

**Model Response**



<char>SPEED</char><bbox>0.364, 0.301, 0.638, 0.375</bbox>
<char>LIMIT</char><bbox>0.362, 0.379, 0.638, 0.454</bbox>
<char>20</char><bbox>0.337, 0.538, 0.661, 0.737</bbox>

Figure 12: OCR Example

**OCR**

<ocr>

**Model Response**

<char>백범로</char><bbox>0.172, 0.265, 0.328, 0.34</bbox>
<char>124번길</char><bbox>0.349, 0.265, 0.512, 0.34</bbox>
<char>Baekbeom-ro</char><bbox>0.171, 0.335, 0.432, 0.391</bbox>
<char>124</char><bbox>0.444, 0.34, 0.508, 0.391</bbox>
<char>만수주공아파트</char><bbox>0.109, 0.528, 0.335, 0.594</bbox>
<char>시흥</char><bbox>0.443, 0.516, 0.522, 0.578</bbox>
...

Figure 13: OCR Example

## Summarization (English)

**British rock band Oasis announce a comeback in reunion tour in 2025**

By Zahid Mahmood, Catherine Nicholls and Alli Rosenbloom, CNN

3 minute read · Updated 3:44 PM EDT, Tue August 27, 2024



Liam Gallagher (left) and Noel Gallagher (right) of British band Oasis. Simon Emmett/Fear PR/PA

(CNN) — British rock band Oasis have announced a comeback in a reunion tour in 2025, according to the band's website.

The tour begins on July 4 at Cardiff Principality Stadium in Wales and the band will play a total of 14 dates across the UK and Ireland.

In a post on X by Oasis frontmen Noel and Liam Gallagher, the brothers said: "This is it, this is happening."

Tickets go on sale at 9 a.m. local time in the UK (4 a.m. ET) on Saturday August 31.

"There has been no great revelatory moment that has ignited the reunion - just the gradual realization that the time is right," said the band in a statement, adding that they are also planning to play more dates outside Europe later next year.

The announcement comes just two days shy of the 30th anniversary of Oasis' 1994 debut studio album "Definitely, Maybe."

The album included popular tracks "Live Forever," "Rock 'n' Roll Star" and "Supersonic" and marked the beginning of Oasis' remarkably fast rise to fame.

Brothers Noel and Liam Gallagher have had a notoriously fraught relationship since they disbanded Oasis in 2009, often trading barbs in the media.

Summarize this news article in lyric style.

## Model Response

In the land of rock and roll, where Oasis once shone,
A reunion tour is announced, a comeback to be known.
Liam and Noel, brothers in fame, their bond now mended,
Their music will resound, in 2025, they'll be attended.

From Cardiff's stadium, they'll play fourteen nights,
Across the UK and Ireland, with fans in delight.
Tickets sold on August 31st, at 9 am local time,
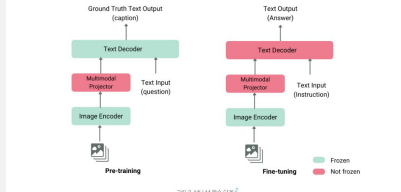A chance for fans to relive memories, of a band so divine.

Their debut album, "Definitely Maybe," a classic of old,
With tracks like "Live Forever" and "Supersonic," untold.
Their rise to fame was fast, a legacy that will last,
Now, after years apart, they're back, to make hearts restart.

Figure 14: Summarization Example (English)

## Text Recognition (Korean)



**MLLM 학습 단계**

MLLM의 학습 과정은 크게 사전 학습 단계(Pre-training)와 미세 조정 단계(Fine-tuning)로 구분됩니다. 아래 그림은 **학습 단계별 모델의 구조**를 표현한 것입니다. 모델의 구조를 간략히 설명드리면 먼저 이미지를 입력으로 받는 사전 학습된 이미지 인코더가 있습니다. 프로젝터는 이미지 인코더에서 출력된 벡터를 텍스트 디코더가 입력으로 받을 수 있는 차원으로 변환해 주는 역할을 합니다. 텍스트 디코더는 사전 학습된 LLM을 의미합니다. 그리고 텍스트는 LLM의 임베딩 계층을 통해서 벡터로 변환됩니다. 아래 그림에서 녹은색은 Not frozen을 의미하며 모델의 Weight를 학습 가능한 상태로 설정했음을 나타내고, 반대로 녹색은 Frozen을 의미하며 모델의 Weight를 학습 불가능한 상태로 두었다는 뜻입니다. 이제 각 단계를 살펴보겠습니다.

*그림 2. MLLM 학습 단계*

먼저 **사전 학습 단계**에서는 이미지-텍스트 입력(질문)을 바탕으로 텍스트 출력(답변)을 생성합니다. 이 단계에서 주목할 점은 이미지 인코더와 텍스트 디코더가 Frozen 상태, 즉 학습되지 않고, 프로젝터만 학습된다는 것입니다. 한편 MLLM 학습 방법론은 아직 표준이 정립되지 않았습니다. 따라서 일부 학습 방법론에서는 사전 학습 때 입력 텍스트 없이 이미지만 입력하고 캡션을 출력하기도 합니다.

**미세 조정 단계**에서는 텍스트 입력은 지시 사항이 들어가도록 하고, 텍스트 출력은 그에 맞는 답변이 생성되도록 합니다. 이전 단계와 다른 모델 구조별 학습 가능 여부는 텍스트 디코더를 Not frozen으로 설정한다는 점입니다. 학습 방법론에 따라 이미지 인코더도 Not frozen 상태로 두어 모든 파라미터를 학습 가능한 상태로 두기도 합니다. 이 외에도 다양한 방법론이 있으며 최근 들어 더욱 빠르게 발전하고 있습니다.

이처럼 사전 학습과 미세 조정이라는 두 단계를 거치면서, MLLM은 이미지를 이해하고 관련된 질문에 답할 수 있는 능력을 갖추게 됩니다. 이러한 MLLM 기술은 인간이 시각적 정보와 언어 정보를 통합하여 처리하는 방식과 유사합니다. 이는 인간-AI 상호작용의 새로운 장을 열어, 다양한 산업 분야에 변화를 불러올 것입니다.

이 게시물에 있는 텍스트를 모두 알려줄래?
(Translation: Can you tell me all the text in this post?)

## Model Response

이미지에 있는 텍스트를 모두 알려드리겠습니다:

---

MLLM 학습 단계

MLLM의 학습 과정은 크게 사전 학습 단계(Pre-training)와 미세 조정 단계(Fine-tuning)로 구분됩니다. 아래 그림은 학습 단계별 모델의 구조를 표현한 것입니다. 모델의 구조를 간략히 설명드리면 먼저 이미지를 입력으로 받는 사전 학습된 이미지 인코더가 있습니다. 프로젝터는 이미지 인코더에서 출력된 벡터를 텍스트 디코더가 입력으로 받아 번역해 주는 역할을 합니다. 텍스트 디코더는 사전 학습된 LLM을 의미합니다. 그리고 텍스트는 LLM의 임베딩 계층을 통해서 벡터로 변환됩니다. 아래 그림에서는 녹은색은 Not frozen을 의미하며 모델의 Weight를 학습 가능한 상태로 설정했음을 나타내고, 반대로 녹색은 Frozen을 의미하며 모델의 Weight를 학습 불가능한 상태로 두었다는 뜻입니다. 이제 각 단계를 살펴보겠습니다.

Ground Truth Text Output (caption)
Text Output (Answer)

Text Decoder
Multimodal Projector
Text Input (question)
Image Encoder

Frozen

Pre-training

Text Decoder
Multimodal Projector
Text Input (Instruction)
Image Encoder

Not frozen

Fine-tuning

그림 2. MLLM 학습 단계

먼저 사전 학습 단계에서는 이미지-텍스트 입력(질문)을 바탕으로 텍스트 출력(답변)을 생성합니다. 이 단계에서 주목할 점은 이미지 인코더와 텍스트 디코더가 Frozen 상태, 즉 학습되지 않고, 프로젝터만 학습된다는 것입니다. 한편 MLLM 학습 방법론은 아직 표준이 정립되지 않았습니다. 따라서 일부 학습 방법론에서는 사전 학습 때 입력 텍스트 없이 이미지만 입력하고 캡션을 출력하기도 합니다.

미세 조정 단계에서는 텍스트 입력은 지시 사항이 들어가도록 하고, 텍스트 출력은 그에 맞는 답변이 생성되도록 합니다. 이전 단계와 다른 모델 구조별 학습 가능 여부는 텍스트 디코더를 Not frozen으로 설정한다는 점입니다. 학습 방법론에 따라 이미지 인코더도 Not frozen 상태로 두어 모든 파라미터를 학습 가능한 상태로 두기도 합니다. 이 외에도 다양한 방법론이 있으며 최근 들어 더욱 빠르게 발전하고 있습니다.

이처럼 사전 학습과 미세 조정이라는 두 단계를 거치면서 MLLM은 이미지를 이해하고 관련된 질문에 답할 수 있는 능력을 갖추게 됩니다. 이러한 MLLM 기술은 인간이 시각적 정보와 언어 정보를 통합하여 처리하는 방식과 유사합니다. 이는 인간-AI 상호작용의 새로운 장을 열어, 다양한 산업 분야에 변화를 불러올 것입니다.

Figure 15: Text Recognition Example (Korean)