
VARCO-VISION-2.0 Technical Report

Young-rok Cha* Jeongho Ju* SunYoung Park
Jong-Hyeon Lee Younghyun Yu Youngjune Kim†

NC AI

{jaycha, jeongho, sun0park, leejh1230, zrohyun}@ncsoft.com
datajuny@gmail.com

Abstract

We introduce VARCO-VISION-2.0, an open-weight bilingual vision-language model (VLM) for Korean and English with improved capabilities compared to the previous model VARCO-VISION-14B. The model supports multi-image understanding for complex inputs such as documents, charts, and tables, and delivers layout-aware OCR by predicting both textual content and its spatial location. Trained with a four-stage curriculum with memory-efficient techniques, the model achieves enhanced multimodal alignment, while preserving core language abilities and improving safety via preference optimization. Extensive benchmark evaluations demonstrate strong spatial grounding and competitive results for both languages, with the 14B model achieving 8th place on the OpenCompass VLM leaderboard³ among models of comparable scale. Alongside the 14B-scale model, we release a 1.7B version optimized for on-device deployment. We believe these models advance the development of bilingual VLMs and their practical applications. Two variants of VARCO-VISION-2.0 are available at Hugging Face: a full-scale 14B model⁴ and a lightweight 1.7B model⁵.

1 Introduction

In December 2024, we released our first vision-language model (VLM), VARCO-VISION-14B [1]. While it achieved strong performance on many benchmarks compared to models of a similar size, it showed limitations in handling multi-image scenarios and Korean-localized tasks. To address these challenges, we present VARCO-VISION-2.0, a new Korean-English VLM designed to understand both images and text and respond to user prompts with greater accuracy and fidelity.

With support for multi-image inputs, the model can effectively process complex visual content, such as documents, tables, and charts. It exhibits robust comprehension in both Korean and English, with notable advancements in Korean language generation and cultural contextual understanding. It shows improved performance across benchmarks and greater usability in practical scenarios, including general Q&A, document parsing, and summarization. As of August 4, 2025, the 14B model ranks 8th on the OpenCompass VLM leaderboard among models with fewer than 20B parameters.

The training strategy follows a four-stage curriculum with memory-efficient techniques, yielding competitive results compared to other open-weight state-of-the-art (SoTA) models. VARCO-VISION-2.0 demonstrates strong capabilities in spatial grounding and real-world perception, delivering high-quality OCR with text localization and robust performance on text-only tasks. To enhance

*Equal contribution

†Corresponding author

³<https://rank.opencompass.org.cn/leaderboard-multimodal>

⁴<https://huggingface.co/NCISOFT/VARCO-VISION-2.0-14B>

⁵<https://huggingface.co/NCISOFT/VARCO-VISION-2.0-1.7B>

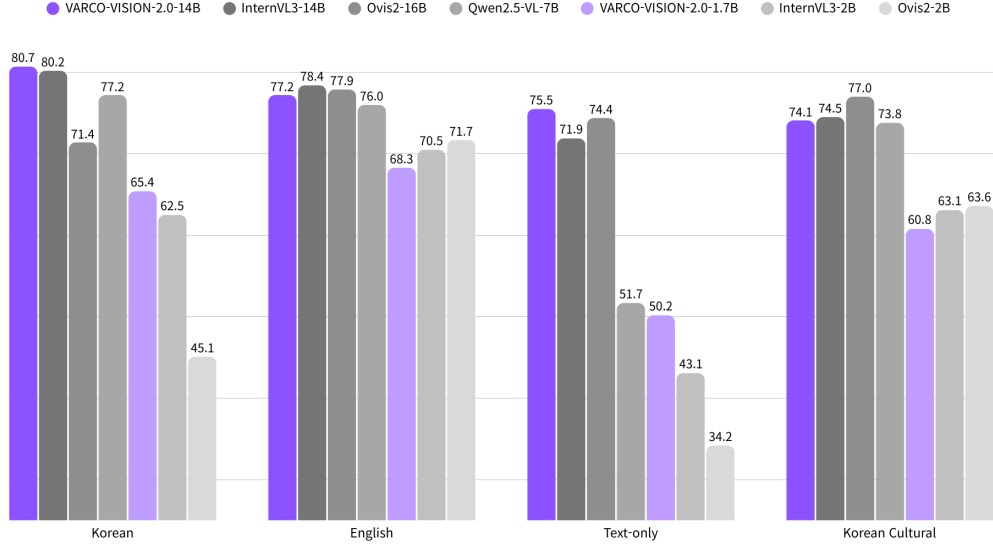


Figure 1: Average Performance of Models across Benchmark Categories.

accessibility, we additionally release a compact 1.7B variant, optimized for deployment on personal devices such as smartphones and PCs.

Major advancements include:

1. **Multi-image Understanding:** Support for multi-image inputs allows the model to analyze multiple images simultaneously and make more holistic and context-aware decisions.
2. **Korean Language Specialization:** The model is further specialized for Korean, with deeper understanding of Korean language, context, and culture. Korean text generation has been significantly improved, resulting in more natural, fluent, and accurate responses.
3. **OCR with Text Localization:** VARCO-VISION-2.0 can identify the position of the text and provide bounding boxes around it. This makes it especially useful for document understanding, signage interpretation, and structured visual data.
4. **Enhanced Safety:** The model now offers improved handling of harmful or sexually explicit content, ensuring safer and more reliable interactions.
5. **Two Model Types of VARCO-VISION-2.0:** We release both a 14B full-scale model and a 1.7B lightweight variant of VARCO-VISION-2.0 on Hugging Face, advancing bilingual VLM research and enabling practical on-device applications.

Together, these features position VARCO-VISION-2.0 as a practical, culturally adapted, open-access vision-language solution.

2 Training

2.1 Architecture

VARCO-VISION-2.0 is built on the LLaVA-OneVision [2] architecture, combining a large language model (LLM), a vision encoder, and a two-layer MLP connector that projects image features into the LLM’s embedding space. We implement the model following the Hugging Face Transformers standard, which ensures immediate compatibility with the transformers ecosystem and enables deployment via vLLM [3] for production-scale inference without requiring additional code adaptation.

We adopt Qwen3 [4] as the LLM and SigLIP2 [5] as the vision encoder. For VARCO-VISION-2.0, we employ SigLIP2 with a patch-16 configuration, replacing the SigLIP patch-14 setup used in VARCO-VISION-14B. Prior work [2] showed that token count has only a limited effect on overall

quality of VLMs, motivating our use of the patch-16 setting. The patch-16 setting reduces the number of visual tokens compared to patch-14; for example, a 384×384 input produces $24^2 = 576$ tokens with patch-16, considerably fewer than the $27^2 = 729$ tokens produced by patch-14. To process high-resolution images effectively, we apply the AnyRes strategy [2], which permits arbitrary input resolutions via tiled/cropped views with resolution-aware aggregation.

2.2 Training Strategies and Datasets

Our training pipeline follows a four-stage curriculum [6, 7] designed to progressively build multimodal capability. We construct a high-quality corpus by integrating English-dominant sources with carefully selected Korean additions. The corpus includes multi-image instruction data, culturally contextual content, and safety-critical scenarios (both curated and synthesized). This diverse composition ensures robust cross-modal learning, bilingual competence, and alignment with safety and localization goals. The model is trained on approximately 6.5 billion text tokens and 30.4 billion image tokens. A breakdown of token usage per stage is shown in Table 1, and key training hyperparameters are summarized in Table 2.

2.2.1 Stage 1. Feature Alignment Pre-training

To bridge the inherent mismatch between the separately pre-trained vision encoder and language model, we train only the MLP connector to project visual features into the language embedding space. Both the vision encoder and language model remain frozen, while paired image-caption data allow the MLP to learn a reliable alignment between the two modalities. All input images are uniformly resized to a fixed resolution prior to encoding, ensuring consistent visual representations and stable training dynamics.

At this stage, the model is trained to generate textual outputs from images alone, without relying on explicit text prompts. To reduce the impact of noisy image-caption pairs, we use a filtered dataset of real-world images paired with concise, well-formatted English descriptions of their key objects. The structured training dataset enables the model to generate sentences in a consistent output format and learn robust input-output alignment, resulting in strong image-to-text generation capabilities.

2.2.2 Stage 2. Basic Supervised Fine-tuning

From this stage onward, we use the AnyRes training strategy introduced in Section 2.1, which allows flexible handling of input images with varying resolutions. In Stage 2, images are processed at relatively low resolutions to reduce computational overhead, and all model components are jointly trained in single-image settings. To achieve effective instruction tuning across diverse downstream tasks, we focus on building a foundation of broad world knowledge and strong visual-textual understanding.

- **General:** We curate captioning datasets covering real-world images, charts, and tables, with a strong emphasis on quality improvement. In addition to selecting high-value open-source data, we re-caption the datasets in-house using VLMs to enhance accuracy, fluency, and consistency. This re-captioning process helps the model to better acquire knowledge about diverse objects, structural layouts, and reasoning over tabular or graphical information.
- **Text Recognition:** We construct a bilingual text-recognition dataset containing Korean and English text in diverse fonts and styles, composed of both collected samples and synthetically generated images with embedded text. Training prompts are formatted as standardized instructions to enhance robust bilingual text recognition and reading comprehension capabilities (e.g., “OCR this image section by section, from top to bottom, and left to right. Do not insert line breaks in the output text. If a word is split due to a line break in the image, use a space instead.”).

2.2.3 Stage 3. Advanced Supervised Fine-tuning

In Stage 3, we aim to train the model to handle more complex scenarios and improve spatial precision. In single-image settings, input images are processed at higher resolutions than in the previous phase, producing finer-grained visual representations. For multi-image scenarios, we adopt a fixed-size image representation strategy to keep token lengths manageable and ensure compatibility with the

model’s context window. This allows the model to scale to multi-image reasoning tasks without sacrificing visual fidelity. We expand the dataset to support instruction tuning for a wide range of image-based tasks.

- **General:** We construct an image-based QA dataset spanning diverse tasks. To strengthen the model’s bilingual capabilities, we regenerate Korean queries from the original English prompts and their corresponding answers, producing a high-quality Korean image-QA dataset. Furthermore, we employ human annotators to enrich the queries by appending task-specific prompts aligned with the target output format (e.g., “Answer the question using a single word or phrase.”, “Explain the solution process step by step and provide the final answer.”). Exposure to such datasets enhances the model’s instruction-following ability, which in turn increases its usability.
- **Document:** We use an in-house dataset designed for document-based question answering (QA), consisting of up to 12 Korean-English images per sample across a wide range of domains. Since generating QA pairs directly from images can often lead to hallucinations, we adopt two strategies that use text as a reference. The first is to collect the accompanying text when crawling document images and construct QA pairs based on the text. The second strategy is to create new QA pairs from the document text and generate corresponding synthetic images for each document using different templates. This pipeline allows us to minimize hallucination and significantly improve the model’s performance on multi-image document QA tasks.
- **Fine-grained:** Inspired by Kosmos-2 [8], we produce grounding and referring datasets. Grounding tasks require the model to identify the locations of objects mentioned in the user query within images. Queries should be annotated with the special `<gro>` token to prompt the model to perform grounding. In contrast, referring is to provide an appropriate, context-based answer based on the objects designated by the user. The locations of objects are represented in the following format:

`<obj>{object}</obj><bbox>{x1}, {y1}, {x2}, {y2}</bbox>`

where `<obj>` encloses the recognized object text and `<bbox>` specifies its bounding box as normalized coordinates $(x_{\min}, y_{\min}, x_{\max}, y_{\max})$ within $[0, 1]$.

- **OCR:** We further develop fine-grained OCR datasets in which inputs for the model are a single image and the `<ocr>` query. The model is required to detect and recognize every word in the image. Words are segmented at whitespace boundaries, ensuring consistent unit-level annotation. Each word is annotated using the format:

`<char>{word}</char><bbox>{x1}, {y1}, {x2}, {y2}</bbox>`

The recognized text is ordered based on y-coordinate clustering to imitate the human reading pattern, proceeding from top to bottom and left to right. This structured annotation lets the model to demonstrate precise character-level recognition and robust word-level understanding in both Korean and English.

2.2.4 Stage 4. Preference Optimization

The final stage employs Direct Preference Optimization (DPO) [9]. Unlike the previous model VARCO-VISION-14B, which updates only the LLM layers in Stage 4, VARCO-VISION-2.0 unfreezes the entire model for full end-to-end preference optimization. These improvements lead to more accurate grounding, richer visual understanding, and better alignment with user intent in practical applications. Preference optimization focuses on guiding the model towards safe, accurate, and culturally appropriate responses. We design training datasets for Stage 4 across three categories:

- **General:** Following the approach of MMPR-v1.2 [10], we develop an in-house dataset specifically for general-purpose images. The goal is to make the model generate precise and reliable answers without hallucination. This dataset serves as the backbone dataset for refining the model’s response quality across a broad spectrum of visual tasks.
- **Safety:** We build datasets targeting safety-critical scenarios, where the model needs to appropriately refuse harmful or unsafe queries, including cases where users submit misleading instructions or provide sensitive images.

- **Localization:** We emphasize cultural and contextual awareness by constructing datasets with Korean-specific images and enriched responses. For instance, an image of Gyeongbokgung Palace is not simply annotated as “a palace,” but as “Gyeongbokgung Palace in Seoul, the main royal palace of the Joseon dynasty.” This enables the model to provide culturally relevant and context-aware answers.

Table 1: Training data usage across stages.

Stage	Text tokens	Image tokens	Total tokens
1	4M	560M	564M
2	760M	4.7B	5.46B
3	5.7B	25B	30.7B
4	61M	93M	154M
Total	6.5B	30.4B	36.9B

Table 2: Hyperparameters.

Model	Hyperparameter	Stage 1	Stage 2	Stage 3	Stage 4
14B	Trainable	MLP	Full Model	Full Model	Full Model
	Batch Size	128	128	128	128
	Context Length	1024	16384	16384	9216
	LR	1e-3	1e-5	1e-5	3e-7
	LR (Vision)	–	2e-6	2e-6	3e-7
	LR Schedule	cosine	cosine	cosine	constant
	Max. #Grids	1×1	2×2	6×6	6×6
	Max. #Tokens	576	(4+1)×576	(9+1)×576	(9+1)×576
1.7B	Trainable	MLP	Full Model	Full Model	Full Model
	Batch Size	128	128	128	128
	Context Length	1024	16384	16384	9216
	LR	1e-3	1e-5	1e-5	6e-7
	LR (Vision)	–	2e-6	2e-6	6e-7
	LR Schedule	cosine	cosine	cosine	constant
	Max. #Grids	1×1	2×2	6×6	6×6
	Max. #Tokens	576	(4+1)×576	(9+1)×576	(9+1)×576

2.3 Initialization Strategy for VARCO-VISION-2.0-1.7B

The VARCO-VISION-2.0-1.7B model shares the same overall architecture as its 14B counterpart, with the only difference being the use of Qwen3-1.7B as the language model in place of Qwen3-14B. Inspired by Progressive Scaling [11], we initialize the vision encoder of the 1.7B model with weights from the 14B model after Stage 3 training. This approach facilitates knowledge transfer from the larger model and accelerates convergence. The impact of this initialization strategy is further discussed in our ablation study (Section 3.2).

2.4 Model Merging

For the 14B variant, we adopt a merge-train-merge strategy to stabilize training and improve generalization of the model. Inspired by prior work on weight-space model averaging [12], we first merge multiple Stage 3 checkpoints to obtain a robust initializer for Stage 4 (preference optimization). After Stage 4 completes, we perform another round of checkpoint merging to produce the final model. This approach reduces checkpoint variance and aggregates distinct patterns learned across multiple checkpoints without introducing additional inference overhead.

In contrast, we do not apply model merging to the lightweight 1.7B variant, as trials on averaging multiple checkpoints with the 1.7B model have degraded validation performance. We find this is due to the lower effective dimensionality in smaller models, where parameter vectors are less likely

to be mutually orthogonal. This property of smaller models makes weight averaging particularly vulnerable to interference. The hypothesis is further supported by recent work that identifies parameter interference as a key limitation in merging low-dimensional models [13]. Therefore, we opt for a simpler strategy for the 1.7B model by selecting one single best-performing checkpoint as its final version.

2.5 Training Infrastructure and Efficiency Optimizations

We use a single compute node equipped with $8 \times \text{H100}$ GPUs. To alleviate memory bottlenecks and enable large-scale training within the constrained setup, we employ several efficiency-oriented strategies.

- **Parallelization and Memory Reduction:** We leverage Fully Sharded Data Parallel (FSDP) [14] as the primary parallelization strategy. FSDP shards not only model parameters but also optimizer states and gradients across GPUs, enabling substantial memory savings and efficient scaling. To further reduce memory overhead, we use activation checkpointing [15] and the 8-bit Adam optimizer [16], which minimize the memory footprint of intermediate activations and optimizer states, respectively.
- **Memory-Efficient Logit Computation:** As modern LLMs adopt extensive vocabulary sizes often exceeding 100K tokens, memory consumption during logit computation has emerged as a major bottleneck in training. VARCO-VISION-2.0 follows this trend with a vocabulary size of approximately 150K, making the logit tensor one of the dominant contributors to peak GPU memory usage. To mitigate this, we integrate the Liger kernel [17], which substantially reduces the memory overhead associated with logit computation.
- **CPU Offloading and Logit Chunking in DPO:** During DPO training, memory demands increase due to the simultaneous use of both the target and reference models. Thus, we apply CPU offloading, shifting portions of model storage to host memory at the cost of added communication latency. Although the Liger kernel proves its effectiveness during SFT, it offers limited benefits in the DPO setting. As a result, we implement a custom logit chunking strategy tailored to DPO, allowing us to manually partition the logit computation and better control peak memory usage.

Overall, training VARCO-VISION-2.0 requires approximately 700 hours of wall-clock time on a single $8 \times \text{H100}$ node, totaling around 5,600 GPU-hours across all stages.

3 Experiments

3.1 Evaluation

3.1.1 English Benchmarks

We evaluate VARCO-VISION-2.0 against SOTA open-weight VLMs, including InternVL3 [7], Ovis2 [18], and Qwen2.5-VL [6], across a broad set of widely adopted English benchmarks. The results are sourced from the OpenCompass VLM leaderboard⁶ if available, to ensure fair and consistent comparisons. Otherwise, evaluations are conducted using the VLMEvalKit [19].

Overall, VARCO-VISION-2.0 demonstrates performance competitive with leading open-weight models (Table 3). It shows strong perceptual capabilities and excels in spatial grounding tasks, particularly in benchmarks emphasizing real-world understanding and low-level visual features (e.g., RealWorldQA [20], SEEDBench_IMG [21], Q-Bench [22], A-Bench [23]). However, performance drops on tasks requiring complex reasoning, scientific knowledge, and document/OCR understanding (e.g., ScienceQA [24], DocVQA [25], TextVQA [26], OCRBench), suggesting clear directions for future improvement.

As presented in Table 4, the 1.7B variant delivers the best performance on RealWorldQA among lightweight models, indicating strong spatial perception in real-world contexts. On the other hand, its performance is relatively weaker on knowledge-intensive and document-oriented tasks such as

⁶<https://rank.opencompass.org.cn/leaderboard-multimodal>, https://huggingface.co/spaces/opencompass/open_vlm_leaderboard

ScienceQA and DocVQA. Although the average score is slightly lower than that of other lightweight models, it displays distinctive strengths in physical-world visual understanding.

3.1.2 Korean Benchmarks

To evaluate Korean language capabilities, we assess VARCO-VISION-2.0 on a set of publicly available Korean multimodal benchmarks: K-MMBench, K-MMStar, K-SEED, K-LLaVA-W, and K-DTCBench [1]. The same baseline models used in the English evaluation are included for consistency. All evaluations are conducted using the original Korean instructions provided in each benchmark.

In Table 5, VARCO-VISION-2.0 demonstrates the highest overall average among all compared models, proving its strong generalization across a wide range of Korean tasks. In particular, it shows substantial gains on K-LLaVA-W, indicating superior capabilities in Korean text generation and dialogue understanding. On the other hand, performance on K-DTCBench remains relatively weaker, likely due to the benchmark’s emphasis on structured data understanding involving documents, tables, and charts. This aligns with the domain-specific limitations in the English benchmarks.

We find that Ovis2-16B shows unusually low scores on K-MMStar, which appear to result from the model frequently failing to follow the expected response format. This suggests that the performance drop arises from challenges with Korean instruction formats rather than from limitations in language understanding. To validate the hypothesis, we re-evaluate all models using the English instruction provided by VLMEvalKit. The results are discussed in the Appendix A.

As reported in Table 6, the 1.7B variant ranks highest among the lightweight models. It performs strongly on K-MMBench_DEV, K-SEED, and K-LLaVA-W, showcasing solid comprehension and generation capabilities in Korean. Nonetheless, its performance on K-DTCBench lags behind, reflecting challenges in structured visual reasoning once again.

The 1.7B model also underperforms on K-MMStar, primarily owing to frequent response-format violations under Korean instructions (e.g., generating free-form answers instead of selecting the correct choice letter). Notably, this issue does not occur with the 14B model. We propose two possible factors contributing to the 1.7B model’s low performance: (i) the absence of Korean multiple-choice data in the Stage 3 training corpus, which limits exposure to the task schema, and (ii) the greater capacity of the 14B model, which enables better generalization to unseen instruction formats, reducing such errors despite equivalent training data setup.

3.1.3 Text-only Benchmarks

We view VLMs as language models extended with visual understanding. From this perspective, our primary goal is to preserve the model’s linguistic foundation and demonstrate that VARCO-VISION maintains strong performance on text-only benchmarks. As shown in Tables 7 and 8, both the full-scale and lightweight variants of VARCO-VISION-2.0 achieve consistently high scores across various text-only tasks, including general knowledge, instruction following, and multi-turn reasoning. These results indicate that the language capabilities have been effectively maintained throughout our training phases.

We credit this performance largely to the strength of the underlying language foundation model, Qwen3. Furthermore, unlike other models that depend on massive amounts of multimodal data, VARCO-VISION-2.0 relies on more compact, high-quality training datasets and efficient learning strategies. This implies that strong performance does not necessarily require maximal training data. Instead, fine-grained data curation and optimized training strategies may be more critical for enhancing text-only capabilities.

We observe that some models show anomalously low performance on the MMLU benchmark [39] because of the response-format violations, as also seen in other multimodal benchmarks. These cases do not indicate significant deficiencies in the models’ knowledge or reasoning capabilities. We encourage readers to keep this in mind when interpreting benchmark results.

3.1.4 Korean Cultural Benchmarks

We also evaluate VARCO-VISION-2.0 on benchmarks assessing understanding of Korean culture and region-specific contexts (Tables 9 and 10). While the model shows improvements over its previous version, its performance remains less competitive compared to other leading models. We attribute

Table 3: English Benchmark Results (Large Models). Since no size-equivalent release of Qwen2.5-VL is available, we report Qwen2.5-VL-7B as a reference. VARCO-VISION is abbreviated as VV. Best in **bold**, runner-up underlined.

Benchmark	InternVL3-14B	Ovis2-16B	Qwen2.5-VL-7B	VV-1.0-14B	VV-2.0-14B
MMStar [27]	68.9	<u>67.2</u>	64.1	64.1	66.9
MMMU_VAL [28]	64.8	<u>60.7</u>	58.0	56.3	<u>61.9</u>
MathVista [29]	74.4	<u>73.7</u>	68.1	67.6	73.2
OCRBench	87.7	<u>87.9</u>	88.8	81.5	86.9
AI2D [30]	<u>86.0</u>	86.3	84.3	83.9	85.7
HallusionBench [31]	<u>55.9</u>	56.8	51.9	46.8	53.2
MMVet [32]	80.5	68.4	<u>69.7</u>	53.0	68.9
SEEDBench_IMG	77.5	<u>77.7</u>	<u>77.0</u>	76.6	78.0
LLaVABench [33]	84.4	93.0	<u>91.0</u>	72.8	90.2
RealWorldQA	69.8	<u>74.1</u>	68.4	72.5	74.6
POPE [34]	89.4	<u>87.5</u>	85.9	87.9	<u>89.2</u>
ScienceQA_TEST	98.6	95.2	89.0	<u>98.5</u>	<u>93.5</u>
SEEDBench2_Plus [35]	70.1	72.1	70.7	<u>69.9</u>	<u>71.9</u>
BLINK [36]	59.9	<u>59.0</u>	55.3	50.3	54.5
TextVQA_VAL	82.2	<u>83.0</u>	85.4	59.5	80.4
ChartQA_TEST [37]	87.8	79.1	80.6	34.2	<u>84.2</u>
Q-Bench1_VAL	76.5	<u>79.2</u>	78.2	74.7	79.9
A-Bench_VAL	76.3	79.6	75.4	74.4	<u>79.5</u>
DocVQA_TEST	94.1	<u>94.9</u>	95.7	77.8	90.9
InfoVQA_TEST [38]	83.6	<u>82.8</u>	82.6	60.2	80.4
Average	78.4	<u>77.9</u>	76.0	68.1	77.2

Table 4: English Benchmark Results (Lightweight Models).

Benchmark	InternVL3-2B	Ovis2-2B	VV-2.0-1.7B
MMStar	61.1	<u>56.7</u>	54.5
MMMU_VAL	48.7	<u>45.6</u>	44.1
MathVista	57.6	64.1	<u>61.1</u>
OCRBench	<u>83.1</u>	87.3	83.0
AI2D	<u>78.6</u>	82.7	76.0
HallusionBench	41.9	50.2	43.0
MMVet	67.0	<u>58.3</u>	<u>52.7</u>
SEEDBench_IMG	75.0	74.4	<u>74.5</u>
LLaVABench	72.1	<u>76.6</u>	77.3
RealWorldQA	65.1	<u>66.0</u>	66.8
POPE	90.1	87.8	<u>88.6</u>
ScienceQA_TEST	95.8	<u>91.2</u>	84.0
SEEDBench2_Plus	64.8	67.4	<u>66.9</u>
BLINK	53.1	<u>47.9</u>	47.2
TextVQA_VAL	<u>78.6</u>	80.0	77.0
ChartQA_TEST	<u>76.0</u>	81.4	75.7
Q-Bench1_VAL	71.9	76.3	<u>72.3</u>
A-Bench_VAL	<u>74.3</u>	76.2	<u>72.4</u>
DocVQA_TEST	<u>88.2</u>	91.9	83.5
InfoVQA_TEST	<u>66.9</u>	71.7	65.0
Average	<u>70.5</u>	71.7	68.3

Table 5: Korean Benchmark Results (Large Models).

Benchmark	InternVL3-14B	Ovis2-16B	Qwen2.5-VL-7B	VV-1.0-14B	VV-2.0-14B
K-MMBench_DEV	89.1	86.0	84.7	84.8	<u>87.7</u>
K-MMStar	64.9	29.7	49.3	58.8	<u>63.6</u>
K-SEED	78.2	73.2	75.7	75.4	<u>77.2</u>
K-LLaVA-W	80.9	86.3	<u>94.1</u>	83.1	96.5
K-DTCBench	87.9	81.7	82.1	<u>84.6</u>	78.3
Average	<u>80.2</u>	71.4	77.2	77.3	80.7

Table 6: Korean Benchmark Results (Lightweight Models).

Benchmark	InternVL3-2B	Ovis2-2B	VV-2.0-1.7B
K-MMBench_DEV	76.9	68.4	77.9
K-MMStar	50.1	10.9	<u>40.8</u>
K-SEED	<u>69.2</u>	34.5	70.7
K-LLaVA-W	47.6	<u>67.2</u>	73.5
K-DTCBench	68.8	44.6	<u>64.2</u>
Average	<u>62.5</u>	45.1	65.4

this performance gap to the relatively limited availability of training data related to Korean culture, which constrains the model’s capacity to understand cultural contexts. Expanding and diversifying this portion of the training corpus presents a promising direction for future work.

3.1.5 OCR Benchmarks

We evaluate VARCO-VISION-2.0’s OCR capabilities on three datasets: CORD [46], ICDAR2013 [47], and ICDAR2015 [48]. These benchmarks require not only accurate recognition of textual content but also precise spatial localization within images.

We find that upscaling input images to a minimum resolution of 2,304 pixels on the longer side (for smaller images) significantly improves the OCR performance. The SigLIP2 vision encoder [5] used in our model operates with an input resolution of 384, and training is performed with a maximum spatial grid size of 6×6 . This implies that the effective maximum resolution is $384 \times 6 = 2,304$, which corresponds to the most fine-grained tokenization supported by the model. Maximizing the number of visual tokens enables the model to capture visual signals in greater detail, and we therefore adopt this resolution for all reported results.

As presented in Table 11, VARCO-VISION-2.0 shows substantial gains over popular open-source OCR systems such as PaddleOCR [49] and EasyOCR [50], achieving notably higher accuracy across all benchmarks. When compared to CLOVA OCR [51]—a strong commercial OCR system—our model demonstrates competitive performance, outperforming it on CORD and ICDAR2013, and closely approaching its accuracy on ICDAR2015. These results are particularly noteworthy given that VARCO-VISION-2.0 is not explicitly trained as an OCR-specific model, which highlights its strong visual-textual alignment and potential for broader real-world applications beyond conventional OCR.

3.2 Ablation Study

Due to compute budget constraints, we only conduct ablations on the 1.7B variant using the eight main benchmarks from the OpenCompass VLM leaderboard. Unless otherwise specified, we follow the same training recipe as in the main setting and report average performance across benchmarks.

3.2.1 Vision Encoder Sharing

As described earlier, the 1.7B model initializes its vision encoder with weights from the 14B model trained up to Stage3. Compared to off-the-shelf SigLIP2 initialization, this shared-encoder initialization results in higher average performance (Table12). We also experiment with freezing the

Table 7: Text-only Benchmark Results (Large Models).

Benchmark	InternVL3-14B	Ovis2-16B	Qwen2.5-VL-7B	VV-1.0-14B	VV-2.0-14B
MMLU	78.5	<u>78.4</u>	4.6	4.9	77.9
MT-Bench [40]	<u>89.3</u>	85.9	80.7	87.7	89.8
KMMLU [41]	<u>51.4</u>	49.3	39.6	37.7	57.5
KoMT-Bench [42]	<u>70.1</u>	<u>79.1</u>	68.4	83.8	78.3
LogicKor [43]	70.0	<u>79.4</u>	65.5	86.7	74.0
Average	71.9	<u>74.4</u>	51.7	60.2	75.5

Table 8: Text-only Benchmark Results (Lightweight Models).

Benchmark	InternVL3-2B	Ovis2-2B	VV-2.0-1.7B
MMLU	59.9	12.9	<u>55.3</u>
MT-Bench	<u>62.8</u>	61.4	72.3
KMMLU	38.0	<u>31.1</u>	10.4
KoMT-Bench	29.1	<u>34.4</u>	59.1
LogicKor	25.6	<u>31.2</u>	53.7
Average	<u>43.1</u>	34.2	50.2

Table 9: Korean Cultural Benchmark Results (Large Models).

Benchmark	InternVL3-14B	Ovis2-16B	Qwen2.5-VL-7B	VV-1.0-14B	VV-2.0-14B
K-Viscuit [44]	71.7	77.0	70.9	69.3	<u>73.7</u>
PangeaBench (ko) [45]	77.2	<u>76.9</u>	76.6	67.6	74.5
Average	<u>74.5</u>	77.0	73.8	68.5	74.1

Table 10: Korean Cultural Benchmark Results (Lightweight Models).

Benchmark	InternVL3-2B	Ovis2-2B	VV-2.0-1.7B
K-Viscuit	<u>60.0</u>	64.1	57.7
PangeaBench (ko)	66.2	63.1	<u>63.8</u>
Average	<u>63.1</u>	63.6	60.8

vision encoder and training only the remaining components of the 1.7B model to evaluate whether this preserves the encoder’s original capabilities throughout the four-step training process. However, this approach leads to lower performance than full end-to-end training, highlighting the importance of continued joint optimization even when transferring from a strong encoder.

3.2.2 DPO Variant Experiments

Inspired by MPO [10], we also investigate a combined objective (DPO+SFT) as an alternative to standard DPO-only for preference optimization (Table 13). (Note: BCO is excluded from this round of experiments due to engineering constraints, such as logit chunking.) Across the eight main OpenCompass benchmarks, the combined objective setting does not outperform DPO-only setting. In fact, DPO-only model achieves marginally better average performance, though the difference may not be significant. A plausible explanation is that our preference dataset is not well aligned with the SFT-style objective, and the auxiliary loss term may introduce gradient interference that weakens the preference learning signal. Based on these findings, our recipe for preference optimization employs shared-encoder initialization, end-to-end training, and DPO objective.

Table 11: OCR benchmark Results (Recognition Accuracy).

Benchmark	CLOVA OCR	PaddleOCR	EasyOCR	VV-1.0-14B	VV-2.0-1.7B	VV-2.0-14B
CORD	93.9	91.4	77.8	81.9	<u>96.2</u>	97.1
ICDAR2013	94.4	92.0	85.0	94.4	95.9	<u>95.7</u>
ICDAR2015	84.1	73.7	57.9	73.5	73.7	<u>79.4</u>
Average	90.8	85.7	73.6	83.3	88.6	<u>90.7</u>

Table 12: Ablation on the 1.7B model: vision-encoder sharing/freezing. Abbrev.: NS = not shared; SH = shared; FRZ_{≤2/3} = freeze vision-encoder until Stage 2/3.

Setting	MMBv1.1	MMStar	AI2D	MMMU	MathVista	HallusionBench	OCRBench	MMVet	Avg.
NS	73.9	51.7	74.6	43.0	56.6	35.3	79.4	50.4	58.1
SH	74.5	54.6	75.8	44.2	61.8	41.7	82.2	49.3	60.5
FRZ _{≤2}	72.8	52.4	75.2	42.6	55.1	36.9	78.6	51.9	58.2
FRZ _{≤3}	72.7	50.7	73.8	42.7	57.4	36.4	76.5	50.8	57.6

3.3 Other Experiments

Extrapolation to Larger Grids and Token Lengths. As we find that the number of visual tokens and grid size affect the model’s OCR performance in Section 3.1.5, we further explore whether the model could extrapolate beyond its original training configuration. When increasing the maximum grid size to 8×8 and the token sequence length to $(16+1) \times 576$, the model exhibits improved performance on OCRBench [52] compared to the original setting. However, for OCR tasks requiring precise text localization, performance drops sharply, with accuracy occasionally reaching zero. This suggests that while content-level understanding extrapolates, spatial generalization remains limited.

Persistence of Qwen3 Thinking Mode. Qwen3 [4] offers a special “thinking mode”, activated by the ‘/think’ flag appended to user queries. We examine if this language model’s inherent thinking capability persists even after VLM training phases. However, our experiments with VARCO-VISION-2.0 show no evidence that the thinking mode behavior has transferred to our model.

Prompted Reasoning Behavior. Although our training data includes only a small set of reasoning-style supervision tasks, we evaluate whether the model might benefit from step-by-step prompting strategies. Specifically, we test prompts such as “Let’s think step by step” [53]. We find that these reasoning cues do not lead to improved accuracy. Even when explicitly instructed to follow multi-step reasoning, the model often remains confident in its initial predictions and rarely changes its outputs. We provide several examples of the explicit reasoning instructions we tested in Appendix B.

4 Limitations

Instruction Robustness. Outputs are sensitive to superficial formatting changes (e.g., whitespace, newlines). This suggests an over-reliance on fixed prompt templates and insufficient exposure to diverse instruction formats during training.

Knowledge and Document Understanding. Despite strong performance in perception and spatial reasoning, the model underperforms on knowledge-intensive and document-centric tasks. We attribute this to limited inclusion of curated knowledge sources and a lack of robust layout-aware supervision during training.

Weakened Referring Capability. The model shows reduced performance on referring tasks compared to the previous VARCO-VISION model. Referring tasks require identifying and resolving references to specific objects or regions within images, and it is essential for embodied agents and multimodal systems that must interact with environments based on visual understanding. The current weakness may limit the model’s performance in practical applications involving precise object selection or instruction following grounded in visual context.

Table 13: Ablation on the 1.7B model: DPO variants.

Setting	MMBv1.1	MMStar	AI2D	MMMU	MathVista	HallusionBench	OCRBench	MMVet	Avg.
SFT	74.6	53.9	75.7	42.6	61.4	43.7	81.0	50.3	60.4
DPO	75.0	54.5	76.0	44.1	61.1	43.0	83.0	52.7	61.2
DPO+SFT	75.0	55.0	75.7	44.6	61.0	43.3	83.1	48.6	60.8

5 Future Work

Small-Model Distillation. Prior studies indicate that small models often benefit more from knowledge distillation than from direct supervised fine-tuning [54, 55, 56]. Based on this, we plan to explore teacher-student training regimes, using the larger model as a teacher for the smaller variant.

Reasoning Improvements. We look forward to enhancing multi-step reasoning capabilities by incentivizing reasoning behavior through reinforcement learning, as explored in recent studies such as DeepSeek-R1 [57]. This involves techniques like chain-of-thought distillation [58, 59] and preference-based fine-tuning [60].

Efficient Context Handling for High-Resolution and Long-Horizon Inputs. Extending the model’s context length is necessary to support high-resolution multi-image inputs and long-horizon video understanding. However, obtaining sufficient training data for very long contexts is challenging, and even when available, learning with such extended sequences greatly increases memory costs. To address this, we plan to adopt techniques such as YaRN [61] for efficient context extension and leverage sequence parallelism [62] to reduce memory overhead. In addition, we aim to reduce the number of visual tokens—currently about five times larger than that of text—by exploring strategies such as pixel unshuffling [63], which can improve overall context efficiency without sacrificing visual fidelity.

Extending to Video Modality and Beyond. A key goal is to enhance the model’s capacity for long-horizon video understanding through spatiotemporal encoders, including support for 3D free-viewpoint video [64] with controllable camera trajectories. Additionally, we aim to develop omni-modal models that incorporate modalities such as audio and speech, allowing the model to interpret diverse and complex signals.

Toward Embodied Multimodal Agents. The AI community is increasingly developing agents that both perceive and act, including interacting with GUIs and manipulating on-screen environments (e.g., Anthropic’s Claude 3.5 Sonnet ‘computer use’ feature [65]). Recent work, such as UI-TARS, shows that end-to-end GUI agents can integrate perception and action within a single model [66], highlighting a shift toward embodied agents that tightly couple vision, language, and action across digital and physical environments. With VARCO-VISION models, we hope to gradually contribute to the development of embodied agents.

Scaling. We plan to scale both model capacity and data volume to enhance generalization and robustness. This includes adopting efficient architectural designs, expanding the multimodal training corpus with higher-quality and more diverse supervision, and refining preference learning through improved reward signals and scalable algorithms. Together, these efforts would yield more capable, aligned, and deployment-ready multimodal models.

6 Conclusion

In this technical report, we present VARCO-VISION-2.0, a Korean-specialized, open-weight VLM available in two sizes (14B and 1.7B). Built on LLaVA-OneVision with a four-stage training curriculum, the models achieve competitive results across multiple tasks compared to leading open-weight VLMs. Their strengths include real-world perception, multi-image understanding, and high-fidelity OCR with bounding boxes, making them well-suited for practical applications in both Korean and English contexts. Text-only benchmark results further demonstrate that the VARCO-VISION series retains core language capabilities, preserving its foundation as a language model. While the 14B model ranks 8th on the VLM leaderboard among models of comparable size, the lightweight 1.7B model provides a practical option for on-device deployment. Overall, VARCO-VISION-2.0 offers a competitive, efficient, and culturally aware foundation for building practical multimodal systems.

References

- [1] Jeongho Ju, Daeyoung Kim, SunYoung Park, and Youngjune Kim. Varco-vision: Expanding frontiers in korean vision-language models, 2024. URL <https://arxiv.org/abs/2411.19103>.
- [2] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024. URL <https://arxiv.org/abs/2408.03326>.
- [3] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [4] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- [5] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, Olivier Hénaff, Jeremiah Harmsen, Andreas Steiner, and Xiaohua Zhai. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features, 2025. URL <https://arxiv.org/abs/2502.14786>.
- [6] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Huihui Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zhenru Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. URL <https://arxiv.org/abs/2502.13923>.
- [7] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Xuehui Wang, Yue Cao, Yangzhou Liu, Xingguang Wei, Hongjie Zhang, Haomin Wang, Weiye Xu, Hao Li, Jiahao Wang, Nianchen Deng, Songze Li, Yinan He, Tan Jiang, Jiapeng Luo, Yi Wang, Conghui He, Botian Shi, Xingcheng Zhang, Wenqi Shao, Junjun He, Yingting Xiong, Wenwen Qu, Peng Sun, Penglong Jiao, Han Lv, Lijun Wu, Kaipeng Zhang, Huipeng Deng, Jiaye Ge, Kai Chen, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models, 2025. URL <https://arxiv.org/abs/2504.10479>.
- [8] 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>.
- [9] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. URL <https://arxiv.org/abs/2305.18290>.
- [10] Weiyun Wang, Zhe Chen, Wenhai Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Jinguo Zhu, Xizhou Zhu, Lewei Lu, Yu Qiao, and Jifeng Dai. Enhancing the reasoning ability of multimodal large language models via mixed preference optimization, 2025. URL <https://arxiv.org/abs/2411.10442>.
- [11] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang, Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo, Jiahao Wang, Tan Jiang, Bo Wang, Conghui He, Botian Shi, Xingcheng Zhang, Han Lv, Yi Wang, Wenqi Shao, Pei Chu, Zhongying Tu, Tong He, Zhiyong Wu, Huipeng Deng, Jiaye Ge, Kai Chen, Kaipeng Zhang, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling, 2025. URL <https://arxiv.org/abs/2412.05271>.
- [12] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization, 2019. URL <https://arxiv.org/abs/1803.05407>.

- [13] Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. Ties-merging: Resolving interference when merging models, 2023. URL <https://arxiv.org/abs/2306.01708>.
- [14] Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, Alban Desmaison, Can Balioglu, Pritam Damania, Bernard Nguyen, Geeta Chauhan, Yuchen Hao, Ajit Mathews, and Shen Li. Pytorch fsdp: Experiences on scaling fully sharded data parallel, 2023. URL <https://arxiv.org/abs/2304.11277>.
- [15] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost, 2016. URL <https://arxiv.org/abs/1604.06174>.
- [16] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise quantization. arXiv:2110.02861, 2021.
- [17] Pin-Lun Hsu, Yun Dai, Vignesh Kothapalli, Qingquan Song, Shao Tang, Siyu Zhu, Steven Shimizu, Shivam Sahni, Haowen Ning, and Yanning Chen. Liger kernel: Efficient triton kernels for llm training, 2025. URL <https://arxiv.org/abs/2410.10989>.
- [18] Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. Ovis: Structural embedding alignment for multimodal large language model. *arXiv:2405.20797*, 2024.
- [19] Haodong Duan, Xinyu Fang, Junming Yang, Xiangyu Zhao, Yuxuan Qiao, Mo Li, Amit Agarwal, Zhe Chen, Lin Chen, Yuan Liu, Yubo Ma, Hailong Sun, Yifan Zhang, Shiyin Lu, Tack Hwa Wong, Weiyun Wang, Peiheng Zhou, Xiaozhe Li, Chaoyou Fu, Junbo Cui, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models, 2025. URL <https://arxiv.org/abs/2407.11691>.
- [20] xAI Corp. Grok-1.5 vision preview: Connecting the digital and physical worlds with our first multimodal model. <https://x.ai/blog/grok-1.5v>, April 2024.
- [21] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension, 2023. URL <https://arxiv.org/abs/2307.16125>.
- [22] Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li, Wenxiu Sun, Qiong Yan, Guangtao Zhai, and Weisi Lin. Q-bench: A benchmark for general-purpose foundation models on low-level vision, 2024. URL <https://arxiv.org/abs/2309.14181>.
- [23] Zicheng Zhang, Haoning Wu, Chunyi Li, Yingjie Zhou, Wei Sun, Xiongkuo Min, Zijian Chen, Xiaohong Liu, Weisi Lin, and Guangtao Zhai. A-bench: Are llms masters at evaluating ai-generated images?, 2025. URL <https://arxiv.org/abs/2406.03070>.
- [24] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering, 2022. URL <https://arxiv.org/abs/2209.09513>.
- [25] Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images, 2021. URL <https://arxiv.org/abs/2007.00398>.
- [26] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read, 2019. URL <https://arxiv.org/abs/1904.08920>.
- [27] 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?, 2024. URL <https://arxiv.org/abs/2403.20330>.
- [28] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi, 2024. URL <https://arxiv.org/abs/2311.16502>.
- [29] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts, 2024. URL <https://arxiv.org/abs/2310.02255>.
- [30] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images, 2016. URL <https://arxiv.org/abs/1603.07396>.

- [31] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. Hallusionbench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models, 2024. URL <https://arxiv.org/abs/2310.14566>.
- [32] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities, 2024. URL <https://arxiv.org/abs/2308.02490>.
- [33] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. URL <https://arxiv.org/abs/2304.08485>.
- [34] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models, 2023. URL <https://arxiv.org/abs/2305.10355>.
- [35] Bohao Li, Yuying Ge, Yi Chen, Yixiao Ge, Ruimao Zhang, and Ying Shan. Seed-bench-2-plus: Benchmarking multimodal large language models with text-rich visual comprehension, 2024. URL <https://arxiv.org/abs/2404.16790>.
- [36] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A. Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive, 2024. URL <https://arxiv.org/abs/2404.12390>.
- [37] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning, 2022. URL <https://arxiv.org/abs/2203.10244>.
- [38] Minesh Mathew, Viraj Bagal, Rubèn Pérez Tito, Dimosthenis Karatzas, Ernest Valveny, and C. V Jawahar. Infographicvqa, 2021. URL <https://arxiv.org/abs/2104.12756>.
- [39] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding, 2021. URL <https://arxiv.org/abs/2009.03300>.
- [40] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL <https://arxiv.org/abs/2306.05685>.
- [41] Guijin Son, Hanwool Lee, Sungdong Kim, Seungone Kim, Niklas Muennighoff, Taekyoon Choi, Cheonbok Park, Kang Min Yoo, and Stella Biderman. Kmmlu: Measuring massive multitask language understanding in korean, 2024. URL <https://arxiv.org/abs/2402.11548>.
- [42] LG AI Research, :, Soyoung An, Kyunghoon Bae, Eunbi Choi, Stanley Jungkyu Choi, Yemuk Choi, Seokhee Hong, Yeonjung Hong, Junwon Hwang, Hyojin Jeon, Gerrard Jeongwon Jo, Hyunjik Jo, Jiyeon Jung, Yountae Jung, Euisoon Kim, Hyosang Kim, Joonkee Kim, Seonghwan Kim, Soyeon Kim, Sunkyoung Kim, Yireun Kim, Youchul Kim, Edward Hwayoung Lee, Haeju Lee, Honglak Lee, Jinsik Lee, Kyungmin Lee, Moontae Lee, Seungjun Lee, Woohyung Lim, Sangha Park, Sooyoun Park, Yongmin Park, Boseong Seo, Sihoon Yang, Heuiyeen Yeen, Kyungjae Yoo, and Hyeongu Yun. Exaone 3.0 7.8b instruction tuned language model, 2024. URL <https://arxiv.org/abs/2408.03541>.
- [43] Jeonghwan Park. Logickor, 2024. URL <https://github.com/instructkr/LogicKor>.
- [44] ChaeHun Park, Yujin Baek, Jaeseok Kim, Yu-Jung Heo, Du-Seong Chang, and Jaegul Choo. Evaluating visual and cultural interpretation: The k-viscuit benchmark with human-vlm collaboration, 2025. URL <https://arxiv.org/abs/2406.16469>.
- [45] 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, 2025. URL <https://arxiv.org/abs/2410.16153>.
- [46] Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaehung Surh, Minjoon Seo, and Hwalsuk Lee. Cord: A consolidated receipt dataset for post-ocr parsing, 2019.
- [47] Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluís Gomez i Bigorda, Sergi Robles Mestre, Joan Mas, David Fernandez Mota, Jon Almazàn Almazàn, and Lluís Pere de las Heras. Icdar 2013 robust reading competition. In *2013 12th International Conference on Document Analysis and Recognition*, pages 1484–1493, 2013. doi: 10.1109/ICDAR.2013.221.

- [48] Dimosthenis Karatzas, Lluís Gomez-Bigorda, Angelos Nicolaou, Suman Ghosh, Andrew Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian Lu, Faisal Shafait, Seiichi Uchida, and Ernest Valveny. Icdar 2015 competition on robust reading. In *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, pages 1156–1160, 2015. doi: 10.1109/ICDAR.2015.7333942.
- [49] PaddlePaddle Authors. Paddleocr, awesome multilingual ocr toolkits based on paddlepaddle. <https://github.com/PaddlePaddle/PaddleOCR>, 2020.
- [50] Jaidev AI. Easyocr: Ready-to-use ocr with 80+ languages. <https://github.com/JaidevAI/EasyOCR>, 2024.
- [51] Naver Cloud. Clova ocr service. <https://www.ncloud.com/product/aiService/ocr>, 2022.
- [52] Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12), December 2024. ISSN 1869-1919. doi: 10.1007/s11432-024-4235-6. URL <http://dx.doi.org/10.1007/s11432-024-4235-6>.
- [53] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners, 2023. URL <https://arxiv.org/abs/2205.11916>.
- [54] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015. URL <https://arxiv.org/abs/1503.02531>.
- [55] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, 2020. URL <https://arxiv.org/abs/1910.01108>.
- [56] Zhiyuan Fang, Jianfeng Wang, Xiaowei Hu, Lijuan Wang, Yezhou Yang, and Zicheng Liu. Compressing visual-linguistic model via knowledge distillation, 2021. URL <https://arxiv.org/abs/2104.02096>.
- [57] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaoqun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- [58] Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve, 2022. URL <https://arxiv.org/abs/2210.11610>.
- [59] Lucie Charlotte Magister, Jonathan Mallinson, Jakub Adamek, Eric Malmi, and Aliaksei Severyn. Teaching small language models to reason, 2023. URL <https://arxiv.org/abs/2212.08410>.

- [60] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL <https://arxiv.org/abs/2203.02155>.
- [61] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context window extension of large language models, 2023. URL <https://arxiv.org/abs/2309.00071>.
- [62] Hao Liu, Matei Zaharia, and Pieter Abbeel. Ringattention with blockwise transformers for near-infinite context. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=WsRHpHH4s0>.
- [63] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P. Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network, 2016. URL <https://arxiv.org/abs/1609.05158>.
- [64] Joel Carranza, Christian Theobalt, Marcus A. Magnor, and Hans-Peter Seidel. Free-viewpoint video of human actors. *ACM Trans. Graph.*, 22(3):569–577, July 2003. ISSN 0730-0301. doi: 10.1145/882262.882309. URL <https://doi.org/10.1145/882262.882309>.
- [65] Anthropic. Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku. <https://www.anthropic.com/news/3-5-models-and-computer-use>, October 2024.
- [66] Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, Wanjun Zhong, Kuanye Li, Jiale Yang, Yu Miao, Woyu Lin, Longxiang Liu, Xu Jiang, Qianli Ma, Jingyu Li, Xiaojun Xiao, Kai Cai, Chuang Li, Yaowei Zheng, Chaolin Jin, Chen Li, Xiao Zhou, Minchao Wang, Haoli Chen, Zhaojian Li, Haihua Yang, Haifeng Liu, Feng Lin, Tao Peng, Xin Liu, and Guang Shi. Ui-tars: Pioneering automated gui interaction with native agents, 2025. URL <https://arxiv.org/abs/2501.12326>.

A Korean Benchmarks under English Instructions

In this additional experiment, we re-run the evaluations in Table 5 using English instructions instead of the original Korean ones. This setting is designed to exclude the effects of Korean-specific instructions and directly assess models’ underlying Korean understanding ability. As shown in the results, Ovis2-16B performs even better than InternVL3-14B under this setup (Table 14), in contrast to the outcomes with Korean instructions. This observation suggests that benchmark scores obtained under a single instruction style may not fully represent a model’s overall language understanding ability. It also highlights that many open-weight models still lack robustness to diverse instruction formats.

Table 14: Korean benchmarks (large models, English instruction setting).

Benchmark	InternVL3-14B	Ovis2-16B	Qwen2.5-VL-7B	VV-2.0-14B
K-MMBench_DEV	83.2	83.2	78.7	82.4
K-MMStar	66.2	64.3	61.7	<u>65.1</u>
K-SEED	77.7	78.4	76.0	<u>78.2</u>
K-LLaVA-W	70.2	<u>80.9</u>	79.8	93.5
K-DTCBench	<u>88.8</u>	85.4	90.4	81.3
Average	75.8	<u>76.8</u>	75.7	79.2

B Explicit Reasoning Instructions

We test several explicit reasoning instructions to examine whether prompting could encourage the model to refine its answers. Below we list three representative examples.

Listing 1: Reasoning Prompt 1

```

First, provide a direct answer to the question based on your first
impression or intuition.
Then, describe the image in relation to the question.
Provide a more informed answer based on your description.
If the answer may be flawed or could be improved, critique it.
Finally, revise and present the final version of your answer if
necessary.

Use the following format. Include only the parts that are relevant or
necessary:

Question:
Direct Answer:
Description:
Description-based Answer: xxx
Critique: (if applicable)
Final Answer: (if applicable)

Question:

```

Listing 2: Reasoning Prompt 2

```

Step #1: Suggest at least two answer candidates. Think out loud.
Candidates can be single words, short phrases, or full sentences
Step #2: Explain which candidate makes the most sense overall.
Step #3: Restate the question in your own words while keeping its
original intent and think carefully about its core meaning before
answering.
Step #4: Provide your final answer in detail - give your best guess
even if you're unsure.

Use the following format:

```

```
Answer Candidates:
- Thinking process for candidate 1: xxx
  Candidate 1: xxx
- Thinking process for candidate 2: xxx
  Candidate 2: xxx
...
Comprehensive Analysis: xxx
Restated Question: xxx
Final Answer: xxx

Question:
```

Listing 3: Reasoning Prompt 3

```
Step #1: Provide a quick guess.
Step #2: Write down your thinking as you work toward the answer.
Step #3: Provide your answer based on your thought.
Step #4: Review and adjust your answer if needed.
Step #5: Provide your final answer in free form.

Use the following format:

Quick Guess: ...
Thinking process:
...
Answer: ...
Review:
...
Free form Answer:
...

Question:
```