

KMMMU: A Korean Massive Multi-discipline Multimodal Understanding Benchmark

Anonymous ACL submission

Abstract

The evaluation of expert-level multimodal capability remains constrained by a reliance on English-centric corpora and translated approximations that fail to capture the nuances of localized professional reasoning. To address this limitation, we introduce **KMMMU** which is a native Korean benchmark constructed exclusively from high-stakes professional and academic examinations. The dataset comprises 3,466 verified image-text pairs spanning diverse disciplines and task types. These questions necessitate the interpretation of information-dense visual inputs including technical diagrams and administrative tables where visual evidence is indispensable for the solution. While scaling and inference-time compute improve logical reasoning, our results suggest that they do not fully mitigate limitations in structural perception or cultural grounding. The observed disparity between processing text-rich documents and abstract diagrams points to ongoing challenges in structural visual reasoning. We hope that **KMMMU** serves as a useful diagnostic resource to address these fine-grained visual and institutional blind spots in future multimodal systems.

1 Introduction

Multimodal Large Language Models have demonstrated remarkable progress in visual recognition and document understanding (Alayrac et al., 2022; Kim et al., 2022; OpenAI, 2023; Team and Google, 2023; Li et al., 2023a; Liu et al., 2023; Bai et al., 2023; OpenAI, 2025; Bai et al., 2025; Google Cloud, 2025). However, current benchmarking methodologies lag behind the complex environments where these systems are increasingly deployed (Sun et al., 2024; Fu et al., 2024; Guan et al., 2024). Most existing evaluations rely on English-centric datasets or translated variants of general knowledge tasks (Li et al., 2023b; Yue et al., 2024). This approach understates two critical challenges in

real-world deployment which are expert-level reasoning under high-stakes assessment formats and deep localization that requires more than simple linguistic fluency (Liu et al., 2024; Li et al., 2023b; Lu et al., 2023; Huang et al., 2023; Son et al., 2025a). Consequently, high scores on current leaderboards often fail to reflect a model’s true proficiency in specific cultural and professional contexts (Zhou et al., 2023; Gudibande et al., 2023).

To contribute to addressing this gap, we introduce **KMMMU**, a native Korean benchmark designed to evaluate expert-level multimodal understanding. Unlike datasets derived from translation, **KMMMU** is constructed from sources including the Public Service Aptitude Test, National Technical Qualifications, academic Olympiads, and National Competency Standards in South Korea. These sources present problems that necessitate reasoning over information-dense visuals such as technical diagrams, administrative tables, and legal schematics. The dataset comprises 3,466 quality-assured questions curated through a hybrid pipeline of automated digitization and multi-stage human verification.

One distinguishing aspect of **KMMMU** is its emphasis on institutional grounding. We include a Korean-specific subset that targets domestic legal frameworks, administrative procedures, and cultural conventions, which helps separate general reasoning from localization demands. We also define a hard subset from consensus failures of strong baselines to focus analysis on challenging instances.

Evaluations of proprietary and open-source models suggest that scaling and inference-time reasoning improve some logical tasks but do not fully resolve perceptual challenges. Even advanced models struggle with structural visual parsing in diagrams and with grounded knowledge needed for localized expert problems. We hope **KMMMU** will serve as a diagnostic benchmark for culturally grounded and professionally relevant multimodal systems.

2 Related Works

Prior multimodal evaluation primarily addressed general visual question answering (Goyal et al., 2017; Hudson and Manning, 2019) or specific skills such as chart and document understanding (Masry et al., 2022; Methani et al., 2020; Mathew et al., 2021; Singh et al., 2019). Recent initiatives have introduced expert-oriented benchmarks to approximate professional problem solving (Lu et al., 2023, 2022; Zhang et al., 2024). We position KMMMU within this landscape to address the specific requirements of Korean expert settings (Table 1).

Global Expert Benchmarks MMMU (Yue et al., 2024) and its extensions like CMMMU (Zhang et al., 2024) and JMMMU (Onohara et al., 2025) evaluate broad reasoning across academic disciplines. While these benchmarks establish a rigorous standard for expert evaluation, they generally rely on globally oriented content and do not assess cultural or institutional grounding. KMMMU diverges from this approach by utilizing native sources from the Public Service Aptitude Test and national technical qualifications. This design ensures the evaluation captures the distinct legal and administrative conventions of South Korea rather than relying on translation.

Korean-Centric Benchmarks Existing Korean resources often dissociate professional expertise from multimodal capabilities. Benchmarks such as KOFFVQA (Kim and Jung, 2025) focus on general scene understanding but lack the complexity required to test advanced reasoning. Conversely, KMMLU (Son et al., 2025a) and KMMLU-Pro (Hong et al., 2025) provides a robust assessment of expert knowledge but remains restricted to a text-only format. Although the recent KO-VLM suite (Marker-Inc-Korea, 2024) improves realism through industrial document parsing, it prioritizes information extraction over complex logic. KMMMU bridges this gap by presenting high-level problems that demand multi-step quantitative reasoning on original visual formats. It integrates an explicit Korean-specific flag to quantify localization and serves as a complementary expert-level benchmark in the field.

3 The KMMMU Benchmark

3.1 Data Collection and Annotation

KMMMU is constructed from high stakes national examinations and academic competitions in

South Korea. We select these sources to organize problems that require advanced reasoning and domain knowledge rather than surface level pattern matching. The benchmark draws from four primary sources including PSAT, national technical qualifications, Olympiads, and NCS. Source specific collection scope and examples are provided in Appendix A.

Annotation and quality control. We combine automated extraction with manual verification. Technical qualification exam data are collected through web crawling, while other sources are digitized using the MinerU-2.5 OCR system (Niu et al., 2025). Five Korean annotators review each instance to verify the question text, associated images, and explanatory materials. Details on the verification interface and the filtering protocol are provided in Appendix B. The resulting dataset contains 3,466 curated questions.

3.2 Taxonomy and Dataset Composition

KMMMU is designed to evaluate expert-level multimodal understanding across diverse domains and problem formats. To support analysis beyond overall accuracy, we annotate each instance along four axes: subject, task, visual modality, and question format, together with a Korean-specific flag.

Labels and quality control. We assign each instance to one of 45 fine-grained subjects aggregated into 9 macro disciplines, and label the primary competency required (task), the dominant visual form, and the answer format. We mark items that require institutional or culturally grounded knowledge specific to South Korea as Korean-specific. We obtain initial subject and task labels using GEMINI-2.5-FLASH-LITE, and visual modality labels and the Korean-specific flag using GEMINI-2.5-FLASH, followed by human consolidation. We validate label quality via a manual audit of roughly 300 randomly sampled instances and manually verify all positive cases for the Korean-specific flag. Additional implementation details are provided in Appendix B.2.

Distributions. Figure 1 shows broad coverage across professional disciplines. KMMMU is dominated by structure-heavy inputs such as diagrams, document-style visuals, and tables (Table 2), and many instances contain in-image text that requires combining text extraction with multimodal reasoning. KMMMU also includes 299 Korean-

Benchmark	Year	Language	Subjects	Expert	Multimodal	Native KO	KO-Spec
<i>Global Expert Multimodal Benchmarks</i>							
MMMU (Yue et al., 2024)	2023	EN	30	✓	✓	✗	✗
CMMM (Zhang et al., 2024)	2024	ZH	39	✓	✓	✗	✗
JMMM (Onohara et al., 2025)	2025	JA	30	✓	✓	✗	✗
<i>Korean-Centric Benchmarks</i>							
KOFFVQA (Kim and Jung, 2025)	2023	KO	General	✗	✓	✓	Low
KMMLU (Son et al., 2025a)	2024	KO	45	✓	✗	✓	High
KMMM (Ours)	2025	KO	45	✓	✓	✓	High

Table 1: **Comparison of KMMM with the MMMU series and existing Korean-centric benchmarks.** *Expert* denotes college-level or professional expertise requirements. *Native KO* indicates construction from original Korean sources rather than translation. *KO-Spec* refers to the inclusion of South Korean-specific legal, administrative, or cultural context.

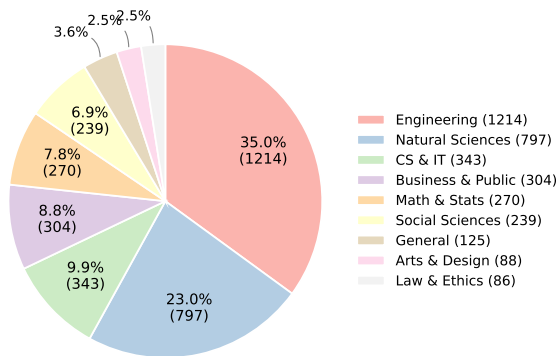


Figure 1: **Macro subject type distribution.**

Visual Modality	Count	%
Diagram	1,751	50.52
Text, Code & Documents	869	25.07
Tables	549	15.84
Charts & Plots	151	4.36
Mathematics	88	2.54
Maps	23	0.66
Symbols & Art	22	0.64
Photographs	13	0.38

Table 2: **Visual modality categories.**

specific items that target localization and institutional grounding, concentrated in domains such as law and public administration. Detailed task and question-format distributions are provided in Appendix D.

4 Experimental setup

4.1 Evaluated Models

We evaluate a diverse set of multimodal models covering both proprietary and open-source systems. The models are organized according to whether they employ explicit reasoning mechanisms during

inference.

Non-reasoning models. This group consists of instruction-tuned multimodal models that generate answers in a single pass without an explicit reasoning mode. Our evaluation covers open-source multilingual models including GEMMA-3 (Team et al., 2025) in the 4B, 12B, and 27B variants, QWEN3-VL (Bai et al., 2025) in the 2B, 4B, 8B, 32B, 30B-A3B, and 235B-A22B variants, as well as the LLAMA-4 (Meta, 2025b,a) family with the Scout and Maverick variants. We also include Korean-focused multimodal models, VARCO-VISION-2.0 (Cha et al., 2025) with 1.7B and 14B parameters and HYPERCLOVAX-SEED-VISION-3B (NAVER HyperCLOVAX, 2025).

Reasoning models. We additionally evaluate models that are designed to perform explicit multi-step reasoning prior to producing a final answer. This group consists of QWEN3-VL-THINKING (Bai et al., 2025) variants with 30B-A3B, 32B, and 235B-A22B parameters, as well as the Korean reasoning model KO-REASON-G3-12B (Son et al., 2025b). These models are configured to internally allocate additional computation for structured reasoning, which is relevant for problems requiring multi-step visual and domain-specific analysis.

Frontier models. We report results on this subset separately, including closed-source frontier models like GPT-5, CLAUDE-OPUS-4.5, GEMINI-3-PRO, etc (OpenAI, 2025; Anthropic, 2025; Google Cloud, 2025).

Hard subset. We define a hard subset to stress-test failure cases that persist under strong baselines.

227	Specifically, we select examples that are answered	is insufficient and correctness depends on more re-	275
228	incorrectly by both GEMMA-3-27B and GEMINI-	liable grounding and reasoning.	276
229	2.5-FLASH-LITE. We report results on this subset		
230	separately, including closed-source frontier models,	Cross-subject variability. Macro subject break-	277
231	to better expose differences that are muted on the	down reveals that some domains are more stable	278
232	full split.	across model choices, while others exhibit larger	279
233		dispersion. This dispersion is consistent with dif-	280
234	4.2 Evaluation Protocols	ferences in localization and institutional grounding	281
235	For all models, All evaluations were conducted in	requirements, which are more prominent in BUSI-	282
236	a zero-shot setting. Response generation was also	NESS & PUBLIC and LAW & ETHICS. In contrast,	283
237	performed using either the officially recommended	subjects that are supported by structured visual evi-	284
238	decoding parameters provided by the model de-	dence tend to be less sensitive to the model choice.	285
239	velopers or the default settings when no explicit	We further quantify these patterns through macro	286
240	recommendations were available. All models were	task analysis in Appendix E.	287
241	evaluated using the same prompt template across		
242	all experiments, and no additional prompt-specific	5.2 Performance on Korean-Specific Content	288
243	tuning or parameter optimization was applied.	To test whether Korean-specific knowledge intro-	289
244	The primary metric is mean accuracy. We per-	duces an additional bottleneck, we stratify items by	290
245	formed 3 independent trials for each model, and	is_korean_specific and compare accuracy on	291
246	report the mean accuracy and the standard devia-	True versus False items. Appendix Table 18 re-	292
247	tion.	ports the counts and ratios of Korean-specific items	293
248		in the Full and Hard splits.	294
249	5 Results		
250		Localization as a remaining barrier. Across	295
251	5.1 Main Results	stronger models, accuracy drops on the Korean-	296
252	Table 3 reports accuracy by macro subject and over-	specific subset even when performance on non-	297
253	all average. Two patterns emerge. First, stronger	Korean-specific items is high. This suggests that	298
254	models benefit from scale, but scale alone is not	these items are not simply harder samples from the	299
255	sufficient to close the gap on KMMMU. Second,	same distribution, but require additional localiza-	300
256	enabling reasoning mode produces broad improve-	tion, including institutional conventions, adminis-	301
257	ments that are most salient when questions require	trative context, and culturally grounded constraints.	302
258	multi step computation or constraint satisfaction		
259	rather than direct recognition. This suggests that	Interpreting small or positive gaps. Some	303
260	KMMMU contains a substantial fraction of items	weaker models exhibit near-zero or even positive	304
261	where their limiting factor is not only perception but	differences between True and False. This does	305
262	also the ability to execute and verify intermediate	not necessarily indicate better localization. When	306
263	reasoning steps.	overall accuracy is low, both subsets can approach	307
264		a similar error rate, which compresses the gap and	308
265	Open-source Korean models. Korean focused	can occasionally reverse its sign. For this reason,	309
266	open source models remain behind the strongest	the stratified gap is most informative for models	310
267	multilingual open source baselines. This gap is	that perform reliably on the non-Korean-specific	311
268	consistent across subjects, which indicates that im-	subset. We provide per-model stratified results in	312
269	proved Korean coverage alone does not fully ad-	Appendix Table 21.	313
270	dress the benchmark demands. Instead, competitive		
271	performance appears to require both robust visual	5.3 Is LLM-as-a-Judge a Consistent	314
272	grounding and broad expert level knowledge across	Evaluator?	315
273	domains.	Because KMMMU includes both multiple-choice	316
274		and free-form questions, we use LLM-as-a-judge	317
	Closed-source Frontier models. We addition-	for scalable evaluation. We run a controlled study	318
	ally evaluate frontier models on the hard subset.	to assess prompt sensitivity under a fixed set of	319
	Results reported in Table 4 show clearer separa-	model outputs.	320
	tion among frontier models, suggesting that the sub-		
	set concentrates cases where shallow visual parsing		

Model	Arts & Design	Business & Public	CS & IT	Engineering	General	Law & Ethics	Math & Stats	Natural Sciences	Social Sciences	Overall Acc.
Open-Source Multilingual Non-Reasoning Models										
Gemma-3-4B-IT	25.00 _{4,10}	20.18 _{1,00}	11.76 _{1,35}	12.14 _{0,58}	10.40 _{1,60}	18.99 _{4,70}	2.84 _{0,21}	7.49 _{1,57}	11.58 _{0,64}	11.41 _{0,50}
Gemma-3-12B-IT	30.30 _{2,62}	25.99 _{2,48}	19.14 _{1,98}	15.60 _{1,08}	9.33 _{1,22}	31.01 _{1,78}	10.00 _{1,70}	14.09 _{1,50}	16.60 _{2,56}	16.68 _{0,42}
Gemma-3-27B-IT	22.73 _{4,10}	24.89 _{2,56}	20.31 _{3,22}	17.93 _{1,84}	14.13 _{4,41}	21.71 _{8,25}	15.19 _{2,06}	17.40 _{1,19}	19.53 _{3,25}	18.63 _{0,36}
Qwen3-VL-2B-IT	24.62 _{3,99}	27.19 _{1,98}	20.41 _{0,87}	13.95 _{0,48}	9.07 _{0,46}	24.03 _{0,67}	11.48 _{3,03}	11.59 _{0,38}	17.29 _{1,35}	15.59 _{0,23}
Qwen3-VL-4B-IT	23.11 _{1,74}	37.06 _{1,33}	25.95 _{2,49}	17.11 _{0,95}	13.60 _{1,39}	27.52 _{2,93}	15.31 _{1,40}	18.95 _{0,78}	23.71 _{2,72}	20.75 _{0,32}
Qwen3-VL-8B-IT	24.62 _{1,74}	40.24 _{1,81}	28.57 _{1,82}	21.33 _{1,65}	17.60 _{0,00}	31.01 _{3,74}	24.94 _{2,52}	23.96 _{1,07}	30.54 _{0,72}	25.42 _{0,53}
Qwen3-VL-32B-IT	27.27 _{3,01}	52.63 _{2,16}	36.73 _{1,05}	28.91 _{0,94}	24.27 _{1,85}	37.21 _{6,98}	48.52 _{4,55}	40.07 _{0,97}	46.72 _{2,38}	37.08 _{0,81}
Qwen3-VL-30B-A3B-IT	25.76 _{8,53}	50.00 _{2,30}	33.33 _{3,03}	26.44 _{1,57}	20.80 _{2,12}	31.78 _{4,08}	41.85 _{2,25}	38.31 _{2,13}	37.10 _{1,05}	33.77 _{0,44}
Qwen3-VL-235B-A22B-IT	23.48 _{2,37}	57.79 _{1,00}	36.73 _{1,46}	31.58 _{1,19}	28.27 _{3,95}	35.66 _{1,34}	55.06 _{1,83}	44.12 _{1,02}	49.37 _{3,27}	40.10 _{0,31}
Llama-4-Scout-17B-16E-IT	25.38 _{3,28}	52.19 _{0,68}	35.86 _{1,54}	30.20 _{1,26}	27.47 _{3,23}	32.17 _{0,67}	34.57 _{1,07}	32.12 _{0,50}	39.61 _{3,36}	33.95 _{0,54}
Llama-4-Maverick-17B-128E-IT	28.79 _{2,37}	56.91 _{2,05}	39.46 _{1,50}	33.66 _{1,67}	25.07 _{3,23}	37.98 _{2,93}	44.32 _{2,04}	39.98 _{1,01}	47.70 _{2,74}	39.20 _{0,38}
Open-Source Korean Non-Reasoning Models										
HyperCLOVAX-SEED-Vision-3B	21.59 _{2,27}	20.83 _{0,95}	18.46 _{0,17}	11.15 _{0,50}	12.00 _{1,39}	21.71 _{1,34}	10.25 _{1,40}	9.37 _{0,59}	16.32 _{0,84}	13.16 _{0,13}
VARCO-VISION-2.0-1.7B	18.18 _{7,10}	12.73 _{3,81}	15.06 _{1,38}	12.52 _{1,93}	6.67 _{2,01}	11.63 _{4,03}	5.93 _{3,21}	7.90 _{2,29}	11.16 _{2,72}	11.03 _{2,39}
VARCO-VISION-2.0-14B	24.24 _{6,46}	28.73 _{2,01}	26.63 _{2,37}	24.44 _{0,26}	13.87 _{2,31}	27.52 _{1,34}	15.19 _{0,74}	18.90 _{1,71}	25.80 _{3,14}	22.82 _{0,16}
Reasoning Models										
Qwen3-VL-30B-A3B-Thinking	34.09 _{15,75}	64.25 _{4,93}	49.76 _{10,21}	49.81 _{7,87}	50.40 _{9,09}	41.86 _{9,93}	76.05 _{5,91}	58.43 _{6,64}	67.78 _{6,54}	55.76 _{7,53}
Qwen3-VL-32B-Thinking	28.41 _{3,41}	62.83 _{2,37}	49.66 _{0,84}	47.50 _{0,81}	44.80 _{3,49}	36.82 _{3,55}	72.84 _{1,83}	56.92 _{0,51}	67.50 _{1,89}	53.73 _{0,49}
Qwen3-VL-235B-A22B-Thinking	23.86 _{3,41}	64.04 _{1,16}	45.87 _{1,02}	45.77 _{1,26}	49.60 _{3,67}	36.05 _{1,16}	75.43 _{0,93}	58.89 _{1,28}	65.27 _{1,45}	53.39 _{0,19}
KO-REASON-G3-12B-1009	30.68 _{1,97}	53.40 _{0,68}	47.42 _{2,15}	55.82 _{1,42}	53.07 _{1,22}	32.56 _{4,03}	76.67 _{1,28}	57.05 _{1,19}	63.46 _{1,69}	55.90 _{0,33}

Table 3: **Accuracy (%) on the full set by macro subject with overall accuracy.** Overall accuracy is averaged across all subjects. Mean accuracy is reported in percentage, with standard deviation shown as a subscript.

Model	Arts & Design	Business & Public	CS & IT	Engineering	General	Law & Ethics	Math & Stats	Natural Sciences	Social Sciences	Overall Acc.
Gemini-3-Pro	69.01 _{0,29}	57.84 _{3,40}	69.30 _{0,76}	64.10 _{1,27}	64.24 _{1,32}	63.46 _{1,92}	60.78 _{6,79}	82.99 _{0,59}	74.16 _{0,99}	72.13 _{1,64}
Gemini-3-Flash	60.27 _{0,15}	52.94 _{2,94}	61.40 _{2,74}	54.21 _{2,77}	55.93 _{1,32}	47.44 _{4,84}	54.90 _{3,40}	73.81 _{1,56}	65.17 _{2,46}	66.67 _{0,95}
GPT-5	51.50 _{0,88}	28.43 _{4,49}	51.75 _{0,76}	39.56 _{2,91}	43.51 _{2,25}	45.51 _{1,11}	37.25 _{3,40}	76.87 _{2,12}	59.55 _{4,54}	61.75 _{1,89}
GPT-5-Mini	42.80 _{0,33}	17.65 _{7,78}	53.95 _{4,56}	31.14 _{2,29}	31.65 _{2,22}	32.05 _{4,44}	29.41 _{10,19}	70.75 _{1,56}	52.81 _{9,99}	49.73 _{3,41}
Claude-Opus-4.5	42.58 _{0,63}	34.31 _{6,12}	56.14 _{3,31}	35.16 _{1,10}	32.87 _{1,17}	38.46 _{1,92}	35.29 _{5,88}	64.97 _{1,18}	45.82 _{1,81}	53.55 _{3,41}
Claude-Sonnet-4.5	34.00 _{0,83}	21.57 _{1,70}	55.70 _{5,32}	24.54 _{3,86}	25.02 _{1,13}	22.44 _{4,84}	41.18 _{5,88}	56.46 _{4,25}	36.20 _{2,49}	42.62 _{3,28}
Grok-4	42.42 _{1,62}	29.41 _{5,09}	42.11 _{2,63}	36.63 _{2,29}	37.63 _{3,09}	33.33 _{6,75}	19.61 _{3,40}	57.82 _{1,18}	49.69 _{0,78}	44.26 _{2,84}
Grok-4.1-Fast	32.54 _{1,10}	27.45 _{1,70}	29.82 _{3,31}	19.05 _{1,68}	28.85 _{1,28}	25.00 _{5,09}	17.65 _{5,88}	37.41 _{1,18}	45.57 _{2,19}	26.23 _{4,34}
Mistral-Large-3-675B-Instruct-2512	23.39 _{1,99}	22.55 _{6,79}	35.53 _{3,48}	21.25 _{1,68}	20.54 _{2,60}	19.87 _{2,94}	15.69 _{3,40}	24.49 _{2,04}	23.22 _{2,34}	32.79 _{4,92}

Table 4: **Accuracy (%) on the hard subset by macro subject with overall accuracy.** Overall accuracy is averaged across all subjects. Results are reported as mean accuracy, with standard deviation shown as a subscript.

Evaluation protocol. We grade hard outputs from three evaluated systems using two judge models with deterministic decoding (temperature 0). For each judge, we compare three prompt templates (P1–P3) and report accuracy on the full hard split ($N = 1,053$). Prompt templates and additional details are provided in Appendix G.

Prompt-induced variation. Table 5 shows that accuracy can vary noticeably across prompts even when evaluated outputs are held fixed. This effect is strongest for GEMINI-2.5-FLASH-LITE, indicating that in mixed-format settings, measured accuracy can be sensitive to prompt-specific parsing and extraction behavior. GEMINI-3-FLASH exhibits smaller variation, suggesting more stable grading under prompt changes. Prompts with explicit answer extraction can help when responses contain long reasoning traces, but they may also introduce additional failure modes when extraction is imperfect.

Takeaway. Overall, these results suggest that LLM-as-a-judge can be useful for mixed-format

evaluation, but prompt choice meaningfully affects reliability. In a small manual audit, P2 aligned most closely with human judgments, while P3 occasionally over-accepted answers when the extraction step was imperfect.

6 Error Analysis

To identify the cognitive bottlenecks of current MLLMs on KMMU, we categorized failure cases into five distinct types. *Knowledge Shortage* aggregates failures where the model lacks relevant domain rules or cannot correctly map the problem statement to the appropriate knowledge. *Reasoning Error* covers coherent attempts that follow an incorrect inference chain. *Visual Perception Error* captures misreading or missing visual evidence in diagrams, tables, or documents. Finally, *Incomplete Response* and *Hallucination* represent output validity issues including truncated answers or unsupported fabrications.

Evaluated Model Outputs (hard)	Gemini-3-Flash (Judge)				Gemini-2.5-Flash-Lite (Judge)			
	P1	P2	P3	Range	P1	P2	P3	Range
Llama-4-Maverick-17B-128E-Instruct	25.7%	26.3%	26.0%	0.6	11.2%	19.9%	30.3%	19.1
Qwen3-VL-30B-A3B-Thinking	50.3%	56.6%	60.5%	10.2	37.5%	28.8%	53.4%	24.6
Gemini-3-Pro	68.1%	68.9%	68.6%	0.8	38.5%	52.9%	76.6%	38.1

Table 5: **Prompt sensitivity of LLM-as-a-judge on mixed-format evaluation.** Accuracy is reported for three judge prompts: **P1** minimal binary verdict, **P2** binary verdict with answer-format matching examples, and **P3** parse-then-judge that explicitly extracts the final answer before emitting the verdict. We also report **Range** (max–min across P1–P3; percentage points) to summarize prompt-induced variation.

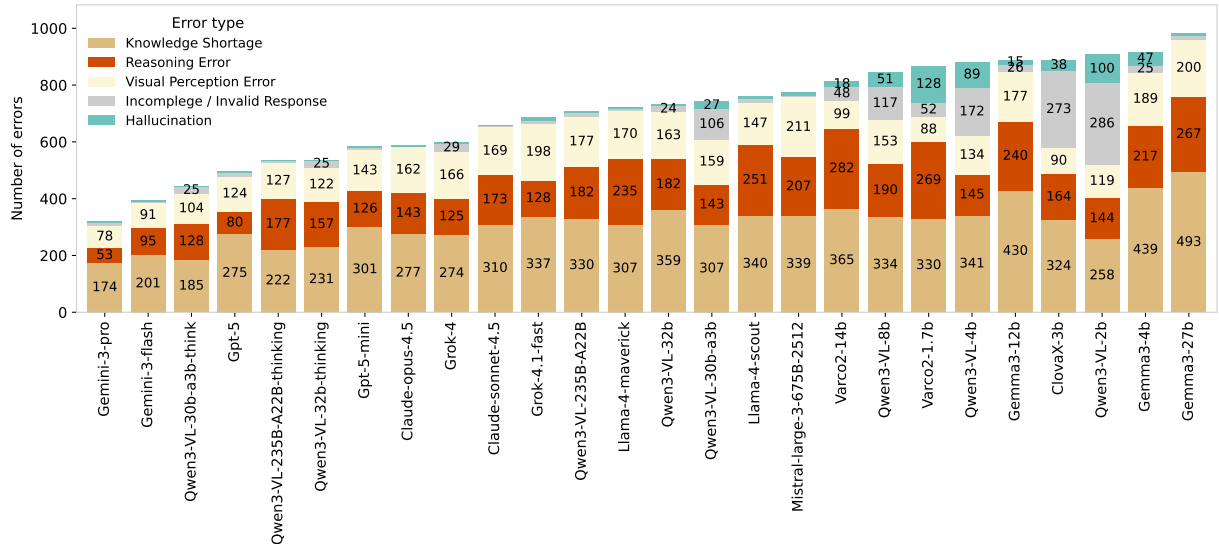


Figure 2: **Error composition on the KMMMU hard split.** We report the number of failures in each error category for each model.

6.1 The Shift in Cognitive Bottlenecks

Figure 2 illustrates that failure modes evolve as model capability increases. To determine which error types best predict performance ranking, we performed a regression analysis detailed in Appendix H.

For open source models, the primary bottleneck appears to be basic generation stability. Regression analysis identifies *Incomplete Response Rate* as the strongest predictor of ranking ($\beta = -0.09, p < 0.001$) followed closely by *Knowledge Shortage* ($\beta = -0.08, p < 0.001$). This indicates that many failures occur before the model can attempt complex multimodal reasoning. In this regime, *Visual Perception Error* has a smaller statistical impact ($\beta = -0.035$) which is consistent with a survivorship effect where instability and missing knowledge dominate first.

For frontier models, the landscape changes significantly. As output validity failures decrease, bottlenecks concentrate on substantive task difficulty. While *Knowledge Shortage* remains the top pre-

dictor ($\beta = -0.05, p < 0.001$) reflecting the expert domains of KMMMU, *Visual Perception Error* rises to become the second most critical bottleneck ($\beta = -0.045, p < 0.001$). This shift implies that once stability issues are resolved, fine grained visual understanding becomes the limiting factor for further gains.

6.2 Perceptual Bottlenecks and Structural Reasoning

To localize this perceptual limitation, we measured the visual perception error rate across different modalities as summarized in Table 6. The results reveal a stark contrast between text anchored modalities and structurally abstract ones. Models demonstrate high proficiency on text rich inputs such as *Text, Code & Documents* and *Tables* with minimal visual perception errors. This indicates that frontier models have effectively solved OCR based perception. However, performance degrades sharply on abstract visual inputs like *Charts & Plots* and *Diagrams*. Failures in these categories suggest that the

Visual Modality	Visual Perception Error (%)	Count (N)
Charts & Plots	39.11%	381
Diagrams	21.87%	4,386
Symbols & Art	19.67%	61
Mathematics	16.38%	232
Maps	11.63%	43
Photos	8.82%	34
Tables	5.48%	1,753
Text, Code & Documents	3.35%	2,360
Others	2.50%	40

Table 6: **Visual Perception Error Rate by Modality in Frontier Models.**

Model	Text, Code & Docs	Tables	Diagrams	Charts & Plots
Gemini-3-Pro	41.7%	18.8%	30.4%	27.9%
Gemini-3-Flash	48.9%	27.2%	40.4%	32.6%
GPT-5	66.7%	25.2%	49.1%	51.2%
Claude-4.5-Sonnet	73.1%	39.1%	70.8%	69.8%

Table 7: **Error rate comparison by visual modality among frontier models.** We report error rates on the hard split for major modalities.

challenge lies not in transcription but in preserving relational structure such as connectivity and layout dependent constraints.

This structural gap explains the performance advantage of GEMINI-3-PRO shown in Table 7. While most frontier models struggle with schematic inputs, GEMINI-3-PRO maintains significantly lower error rates on diagrams. This stability suggests stronger robustness to relational visual representations which effectively bridges the bottleneck between visual encoding and downstream reasoning.

6.3 Impact of Reasoning mode

We evaluated the impact of reasoning mode by comparing the QWEN3-VL series in standard instruction tuned (IT) form versus reasoning mode. Table 8 shows that reasoning mode yields large accuracy gains across model scales. Our analysis indicates that the thinking process primarily mitigates output stability issues and knowledge deficits. Comparing QWEN3-VL-30B variants reveals that incomplete responses are nearly eliminated and knowledge shortages are reduced. However, visual perception errors comprise a larger share of the remaining failures in reasoning models. This implies that while reasoning can refine logic and knowledge application, it cannot hallucinate correct visual details from a misperceived image embedding.

Base Model	Instruct (IT)	Thinking	Gain (Δ)
Qwen3-VL-30B-A3B	33.95%	55.76%	+21.81%
Qwen3-VL-32B	33.77%	60.28%	+26.51%
Qwen3-VL-235B	40.10%	58.87%	+18.77%

Table 8: **Accuracy gains from enabling reasoning mode in QWEN3-VL.** We compare instruction tuned models with their reasoning mode variants on the same evaluation setting.

6.4 Lessons from Error Patterns

Our error analysis points to three areas where progress appears most needed. First, the modality gap suggests that transcription oriented pretraining alone may not be sufficient for structural visual reasoning, motivating objectives that emphasize explicit structure extraction. Second, the persistence of knowledge shortages even in frontier models indicates that linguistic proficiency is not the same as institutional grounding, and that stronger coverage of Korean administrative and legal contexts may be important. Finally, the limited impact of reasoning mode on perceptual errors implies that long-form reasoning can drift away from the original visual evidence, suggesting a need for mechanisms that encourage iterative visual re-attention during inference.

7 Ablation Study

7.1 Evaluation of Image-Dependency

Model	Original	Text-only
Gemini-3-Flash	64.91 _{0.11}	37.14 _{0.48}
GPT-5-Mini	47.63 _{0.33}	23.18 _{0.29}

Table 9: **Ablation results for image dependency.** We compare the average accuracy and standard deviation (Acc_{std}) of models on the original multimodal dataset versus the text only baseline.

To verify the validity of KMMM as a multimodal benchmark, we evaluate the degree to which models rely on visual information to solve the problems. A robust visual reasoning benchmark should contain questions that are difficult or impossible to answer without access to the corresponding images. We conduct a text only baseline experiment where the models receive the textual question and options but the visual input is entirely omitted. Table 9 presents the performance gap between the original multimodal setting and the text only setting for GEMINI-3-FLASH and GPT-5-MINI. The results

indicate that the removal of images leads to a significant degradation in accuracy for both models. GEMINI-3-FLASH shows an average accuracy drop of 27.77% while GPT-5-MINI exhibits a decrease of 24.45%. This substantial performance decline demonstrates that KMMMU effectively captures the necessity of visual information. The models cannot rely solely on textual context or internal knowledge to solve a large majority of the tasks. The high performance gap confirms that our benchmark successfully measures the ability of models to integrate and reason over multimodal inputs.

7.2 Data Contamination Analysis

To assess the risk of memorization, we run a prefix completion test in which models receive the first 35% of a question and generate the remaining continuation. We apply this test to questions longer than 150 tokens and evaluate three frontier models under two settings. The first setting provides no additional metadata. The second setting provides the exam name and year as a potential memorization trigger.

We evaluate reconstructions using an LLM as a judge. We use GEMINI-3-FLASH to score exactness on a 0 to 100 scale based on lexical overlap, key detail preservation, and overall faithfulness to the reference question. We also report refusal and hallucination rates to characterize behavior under incomplete prompts.

Across all sources, GPT-5-MINI and GEMINI-3-PRO mostly refuse to complete the continuations, with refusal rates typically above 95%. This behavior is consistent with safety and policy training that discourages reproducing exam materials. GEMINI-3-FLASH attempts completions more often, but its exactness remains low and does not systematically increase in the hinted setting. These trends suggest that observed benchmark performance is unlikely to be explained by simple memorization. Full results by source, model, and hint setting are reported in Appendix Table 26.

8 Conclusion

We introduced KMMMU, a native Korean benchmark designed to evaluate expert-level multimodal understanding in real-world assessment settings. KMMMU is built from South Korean examinations and competitions and comprises 3,466 carefully verified questions with information-dense visuals, broad subject coverage, and rich annotations

for subject, task, and visual modality. To better disentangle general capability from localization, we additionally curated a Korean-specific subset that requires institutional grounding, and we constructed a hard subset to concentrate frontier-level failure cases.

Our evaluations across open-source, Korean-focused, and frontier multimodal models reveal a clear capability gradient, but also a consistent set of remaining bottlenecks. While larger models and reasoning mode generally improve accuracy on problems that demand multi-step computation or constraint satisfaction, they do not fully resolve perceptual challenges. Error analyses indicate that once output stability improves, failures increasingly concentrate on *structural* visual understanding and grounded domain knowledge, with diagrams and plots remaining particularly error-prone.

Taken together, these results position KMMMU as a diagnostic benchmark for the next stage of multimodal research: models that robustly preserve relational structure in complex visuals and that exhibit institutionally grounded reasoning in localized professional contexts. We hope KMMMU will serve as a reference point for developing culturally aware and professionally competent multimodal systems, and for evaluating progress beyond English-centric or translation-based testbeds.

Limitations

Coverage and representativeness. Although KMMMU spans many disciplines, it is not a comprehensive model of all real-world multimodal use cases. The benchmark is exam-centric and emphasizes information-dense, structure-heavy visuals (e.g., diagrams and documents), so performance may not directly transfer to everyday perception, interactive settings, or non-exam domains.

Annotation noise and taxonomy subjectivity. Subject, task, and modality labels are ultimately based on human consolidation of initial model-proposed annotations, and some categories admit ambiguous boundaries. While we audit a random subset and manually verify all Korean-specific positives, residual label noise is likely, especially for fine-grained subjects and multi-skill questions.

Evaluation noise for mixed-format answers. Because KMMMU includes both multiple-choice and free-form items, scalable evaluation relies on

562	LLM-as-a-judge, which can be sensitive to prompt design and answer formatting. Despite using deterministic decoding and spot-checks, some grading errors may remain, particularly when responses are verbose, underspecified, or unconventional in format.	
563		
564		
565		
566		
567		
568	References	
569	Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, and 1 others. 2022. Flamingo: a visual language model for few-shot learning. In <i>NeurIPS</i> .	
570		
571		
572	Anthropic. 2025. Introducing Claude Opus 4.5 . Anthropic Newsroom. Published: 2025-11-24. Accessed: 2026-01-05.	
573		
574		
575	Jinze Bai, Shuai Bai, Shusheng Yang, and 1 others. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. <i>arXiv preprint arXiv:2308.12966</i> .	
576		
577		
578		
579	Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025. Qwen3-vl technical report . Preprint, arXiv:2511.21631.	
580		
581		
582		
583		
584		
585		
586	Young-rok Cha, Jeongho Ju, SunYoung Park, Jong-Hyeon Lee, Younghyun Yu, and Youngjune Kim. 2025. Varco-vision-2.0 technical report. <i>arXiv preprint arXiv:2509.10105</i> .	
587		
588		
589		
590	Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasapat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. <i>arXiv preprint arXiv:2507.06261</i> .	
591		
592		
593		
594		
595		
596		
597	Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. 2024. Blink: Multimodal large language models can see but not perceive. In <i>European Conference on Computer Vision</i> , pages 148–166. Springer.	
598		
599		
600		
601		
602		
603	Google Cloud. 2025. Gemini 3 Pro (Preview) Generative AI on Vertex AI . Google Cloud Documentation. Release date: 2025-11-18. Accessed: 2026-01-05.	
604		
605		
606	Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 6904–6913.	
607		
608		
609		
610		
611		
612	Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen,	
613		
	Furong Huang, Yaser Yacoob, and 1 others. 2024. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14375–14385.	614 615 616 617 618 619
	Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. The false promise of imitating proprietary llms. <i>arXiv preprint arXiv:2305.15717</i> .	620 621 622 623
	Seokhee Hong, Sunkyoung Kim, Guijin Son, Soyeon Kim, Yeonjung Hong, and Jinsik Lee. 2025. From kmmlu-redux to kmmlu-pro: A professional korean benchmark suite for llm evaluation. <i>arXiv preprint arXiv:2507.08924</i> .	624 625 626 627 628
	Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, and 1 others. 2023. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. <i>Advances in Neural Information Processing Systems</i> , 36:62991–63010.	629 630 631 632 633 634 635
	Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 6700–6709.	636 637 638 639 640
	Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoon Yun, Dongyoon Han, and Seunghyun Park. 2022. Ocr-free document understanding transformer. In <i>European Conference on Computer Vision</i> , pages 498–517. Springer.	641 642 643 644 645 646
	Yoonshik Kim and Jaeyoon Jung. 2025. Koffvqa: An objectively evaluated free-form vqa benchmark for large vision-language models in the korean language. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 575–585.	647 648 649 650 651
	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>ICML</i> .	652 653 654 655
	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Evaluating object hallucination in large vision-language models. <i>arXiv preprint arXiv:2305.10355</i> .	656 657 658 659
	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In <i>NeurIPS</i> .	660 661
	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, and 1 others. 2024. Mmbench: Is your multi-modal model an all-around player? In <i>European conference on computer vision</i> , pages 216–233. Springer.	662 663 664 665 666 667

668	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. <i>arXiv preprint arXiv:2310.02255</i> .	japanese massive multi-discipline multimodal understanding benchmark for culture-aware evaluation. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 932–950.	723 724 725 726 727 728
674	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. <i>Advances in Neural Information Processing Systems</i> , 35:2507–2521.	OpenAI. 2023. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	729 730
675		OpenAI. 2025. Introducing GPT-5 . OpenAI. Published: 2025-08-07. Accessed: 2026-01-05.	731 732
676		Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 8317–8326.	733 734 735 736 737 738
677			
678			
679			
680	Marker-Inc-Korea. 2024. Ko-vlm benchmark: A comprehensive evaluation dataset for korean vision-language models. https://github.com/Marker-Inc-Korea/KO-VLM-Benchmark . Accessed: 2026-01-04.		
681			
682			
683			
684			
685	Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In <i>Findings of the association for computational linguistics: ACL 2022</i> , pages 2263–2279.	Guijin Son, Hanwool Lee, Sungdong Kim, Seungone Kim, Niklas Muennighoff, Taekyoon Choi, Cheonbok Park, Kang Min Yoo, and Stella Biderman. 2025a. Kmmlu: Measuring massive multitask language understanding in korean. In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 4076–4104.	739 740 741 742 743 744 745 746 747
686			
687			
688			
689			
690			
691	Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 2200–2209.	Guijin Son, Donghun Yang, Hitesh Laxmichand Patel, Amit Agarwal, Hyunwoo Ko, Chanuk Lim, Srikant Panda, Minhyuk Kim, Nikunj Drolia, Dasol Choi, and 1 others. 2025b. Pushing on multilingual reasoning models with language-mixed chain-of-thought. <i>arXiv preprint arXiv:2510.04230</i> .	748 749 750 751 752 753
692			
693			
694			
695			
696	Meta. 2025a. meta-llama/Llama-4-Maverick-17B-128E-Instruct . Hugging Face model card and weights. Model release date: 2025-04-05. Accessed: 2026-01-05.	Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. 2024. Scieval: A multi-level large language model evaluation benchmark for scientific research. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19053–19061.	754 755 756 757 758 759
697			
698			
699			
700	Meta. 2025b. meta-llama/Llama-4-Scout-17B-16E . Hugging Face model card and weights. Model release date: 2025-04-05. Accessed: 2026-01-05.	Gemini Team and Google. 2023. Gemini: A family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	760 761 762
701			
702			
703	Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar. 2020. Plotqa: Reasoning over scientific plots. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 1527–1536.	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. Gemma 3 technical report. <i>arXiv preprint arXiv:2503.19786</i> .	763 764 765 766 767
704			
705			
706			
707			
708	NAVER HyperCLOVAX. 2025. HyperCLOVAX-SEED-Vision-Instruct-3B . Hugging Face model card and weights. Model release date: 2025-04-24 (as stated in repository license). Accessed: 2026-01-05.	Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, and 1 others. 2025. InternV1.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. <i>arXiv preprint arXiv:2508.18265</i> .	768 769 770 771 772 773
709			
710			
711			
712	Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, Zhenjiang Jin, Guang Liang, Rui Zhang, Wenzheng Zhang, Yuan Qu, Zhifei Ren, Yuefeng Sun, Yuanhong Zheng, Dongsheng Ma, and 42 others. 2025. Mineru2.5: A decoupled vision-language model for efficient high-resolution document parsing . <i>Preprint</i> , arXiv:2509.22186.	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi.	774 775 776 777 778
713			
714			
715			
716			
717			
718			
719			
720	Shota Onohara, Atsuyuki Miyai, Yuki Imajuku, Kazuki Egashira, Jeonghun Baek, Xiang Yue, Graham Neubig, and Kiyoharu Aizawa. 2025. Jmmmu: A		
721			
722			

779			
780		In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 9556–9567.	
781			
782	Ge Zhang, Xinrun Du, Bei Chen, Yiming Liang, Tongxu		
783	Luo, Tianyu Zheng, Kang Zhu, Yuyang Cheng,		
784	Chunpu Xu, Shuyue Guo, and 1 others. 2024. Cm-		
785	mmu: A chinese massive multi-discipline multi-		
786	modal understanding benchmark. <i>arXiv preprint</i>		
787	<i>arXiv:2401.11944</i> .		
788	Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen,		
789	Wayne Xin Zhao, Xu Chen, Yankai Lin, Ji-Rong		
790	Wen, and Jiawei Han. 2023. Don’t make your llm		
791	an evaluation benchmark cheater. <i>arXiv preprint</i>		
792	<i>arXiv:2311.01964</i> .		
793	A Data Sources and Collection Scope		
794	KMMM U is collected from four high stakes		
795	sources in South Korea. We summarize the col-		
796	lection scope for each source below.		
797	A.1 PSAT		
798	We annotate ten years of past examinations from		
799	civil service recruitment tracks. The PSAT includes		
800	Language Logic, Data Interpretation, and Situa-		
801	tional Judgment sections that assess logical reason-		
802	ing and information integration.		
803	A.2 National Technical Qualifications		
804	We collect fifteen years of questions from 252		
805	distinct certification exams, including Information		
806	Processing Engineer, Electric Engineer, and Fire		
807	Safety Manager. These exams cover a wide range		
808	of technical domains across industrial and engineer-		
809	ing fields.		
810	A.3 Olympiads		
811	To incorporate academically challenging reasoning		
812	problems, we gather ten years of Olympiad ques-		
813	tions spanning middle school, high school, and		
814	university levels. The collected problems focus pri-		
815	marily on mathematics and science.		
816	A.4 NCS		
817	We include three years of National Competency		
818	Standards examinations covering all ten compe-		
819	tency areas, such as Communication, Numeracy,		
820	and Problem Solving. These exams are used in		
821	recruitment for public sector organizations.		
822	B Annotation and Quality Control Details		
823	The construction of KMMM U uses a rigorous		
824	pipeline that combines automated processing with		
825	human verification to ensure high data fidelity.		
	B.1 Human Verification Interface		826
	We utilized a custom built annotation tool to verify		827
	and correct the output of the OCR pipeline. Raw		828
	data digitized by MinerU-2.5 (Niu et al., 2025) of-		829
	ten contained artifacts and formula errors. Figure 3		830
	shows the interface where five Korean annotators		831
	reviewed the parsed content against the original		832
	PDF source. Annotators were instructed to		833
	• Correct LaTeX formatting for mathematical		834
	formulas		835
	• Verify that image references in the text		836
	matched the cropped images		837
	• Discard questions where essential visual in-		838
	formation was illegible or missing		839
	B.2 Automatic Labeling and Taxonomy		840
	Consolidation		841
	We annotate several auxiliary attributes to sup-		842
	port analysis and stratified reporting, including		843
	subject, task, image type, and a Korean-specific		844
	flag. Initial subject and task labels are produced by		845
	GEMINI-2.5-FLASH-LITE. Image-type labels and		846
	the Korean-specific flag are produced by GEMINI-		847
	2.5-FLASH. For each labeling job, the model is		848
	given the question text and its associated image		849
	and outputs the most appropriate label.		850
	We use an open labeling step that does not con-		851
	strain predictions to a fixed label set. This reduces		852
	forced assignments when an instance does not		853
	cleanly match a predefined taxonomy. All label		854
	types are generated independently.		855
	Manual audit and consolidation We conduct a		856
	manual audit by randomly sampling around 300 in-		857
	stances and reviewing the assigned labels. Based on		858
	the audited outputs, we consolidate the subject tax-		859
	onomy through human curation into 45 sub-subject		860
	categories and 9 macro subject categories. Task la-		861
	belts are also consolidated through human review		862
	into a final set of 11 task types.		863
	Verification of Korean-specific cases Because		864
	false positives can inflate localization analyses, we		865
	manually verify all instances labeled as Korean-		866
	specific. We confirm that each positive case re-		867
	quires Korean-specific knowledge or context rather		868
	than general world knowledge expressed in Korean.		869

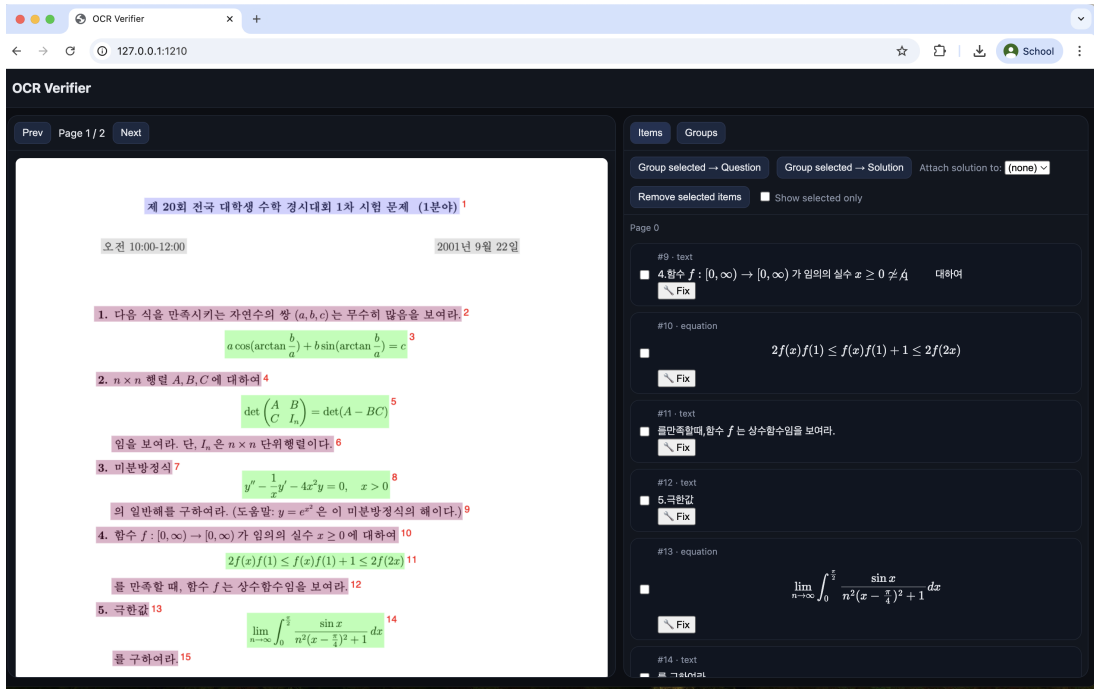


Figure 3: **Annotation tool interface used for OCR verification.** The tool displays the original PDF page on the left and the parsed text and images on the right, allowing annotators to correct OCR errors and validate image cropping in real time.

B.3 Adversarial Filtering Protocol

To ensure the benchmark evaluates frontier level capabilities, we implemented a multi stage adversarial filtering process designed to remove trivial samples that can be solved by existing mid tier models. We follow the rejection strategy below.

1. **De duplication.** We removed duplicate questions across different exam years using perceptual hashing for images and n gram overlap for text.
2. **Model based adversarial filtering.** We employed a cascade of multimodal models including INTERNVL-3.5-38B (Wang et al., 2025), GEMINI-2.5-FLASH-LITE, and GEMINI-2.5-FLASH (Comanici et al., 2025). If any of these models answered a question correctly with high confidence, the instance was excluded. This ensures that KMMMU focuses on the challenging tail of the distribution.

We initially collected approximately 68,000 raw samples. After removing invalid image links and duplicates, we apply multi stage adversarial filtering with frontier multimodal models to exclude trivial items and questions solvable without visual

reasoning. The resulting dataset contains 3,466 curated questions.

C Korean-Specific Context

To provide a concrete understanding of the KMMMU benchmark, we present a detailed visualization of a "Korean-Specific" instance in Figure 4. Standard multimodal benchmarks often rely on universal knowledge (e.g., physics, basic math) that is invariant across cultures. In contrast, KMMMU includes a dedicated subset of questions that require localized knowledge. Figure 4 shows a sample data card from the dataset. The input consists of an image containing a specific regulation text and the corresponding question. As shown in the translation panels, the model must interpret the visual text ("extraction area slope criteria") specifically within the context of the South Korean *Mountainous Districts Management Act* to identify the correct legal standards (Option 3). This necessitates that the model possesses not only optical character recognition capabilities but also grounded knowledge of Korean administrative laws.

D Detailed Dataset Statistics

In this section, we provide a granular breakdown of the dataset composition. Beyond the overview (Ta-

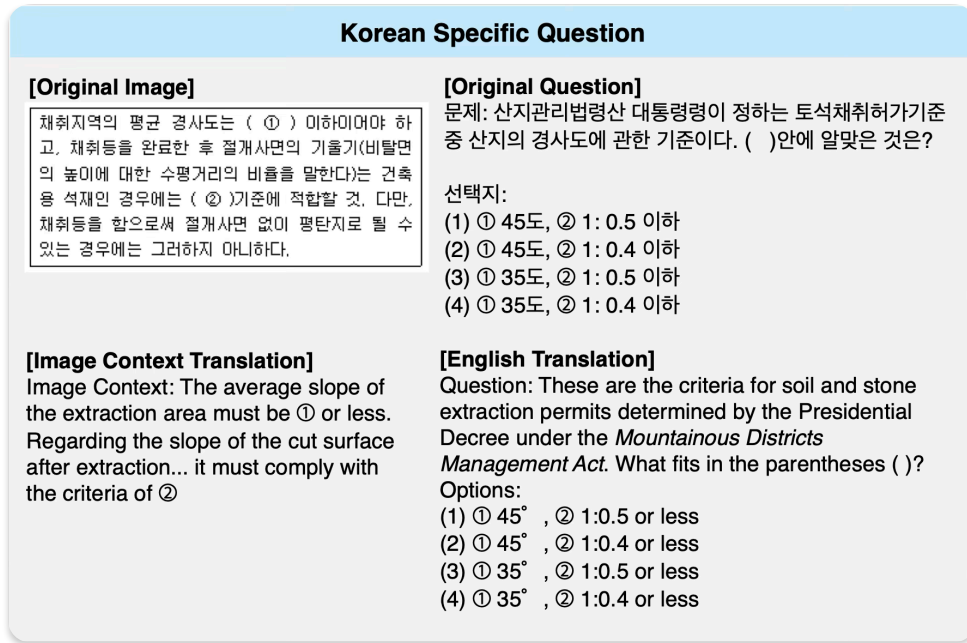


Figure 4: **Data Card for a Korean-Specific Question.** The figure aggregates the raw inputs and their translations. **[Original Image]** The original visual input containing a text-rich regulation box. **[Original Question]** The original question text in Korean. **[Translation]** English translations for both the visual context and the question. Correctly answering this question requires retrieving specific legal provisions regarding slope limits for soil extraction permits in South Korea, demonstrating the benchmark’s focus on localized expert knowledge.

ble 10), we report (i) the distribution of fine-grained subject categories (Table 11), (ii) the macro task distribution (Table 12), (iii) the question format distribution (Table 13), and (iv) the taxonomy of visual modalities (Table 14).

Statistic	Count
Total Questions	3,466
Macro Subject Categories	9
Sub-Subject Categories	45
Macro Visual Modality Types	8
Macro Task Types	11
Question Types	5
Questions with in-image texts	2,383 (68.75%)
Questions w/o in-image texts	1,083 (31.25%)
Korean-specific questions	299 (8.63%)

Table 10: **Dataset overview.** We report counts and percentages for key attributes such as in-image text and Korean-specific content.

D.1 Subject Category Distribution

Table 11 details the frequency of questions across 45 fine-grained subject categories. The distribution reflects the emphasis on STEM (Science, Technology, Engineering, and Mathematics) fields, with *Physics, Civil Engineering, and Mechanical Engineering* constituting the largest portions. This heavy tail in engineering disciplines ensures that KMMMU serves as a robust benchmark for technical domain expertise.

D.2 Macro Task Distribution

Table 12 summarizes the distribution of macro task labels. This breakdown clarifies what primary competency is most frequently required by KMMMU (e.g., concept application, data and visual interpretation, and quantitative reasoning), and supports task-wise analyses in later sections.

D.3 Question Format Distribution

Because KMMMU contains both multiple-choice and free-form items, the answer format affects eval-

945	uation difficulty and failure modes. Table 13 reports	993
946	the distribution of question formats in the bench-	994
947	mark.	995
948	D.4 Visual Modality Taxonomy	996
949	Table 14 presents our hierarchical visual taxonomy.	997
950	We grouped fine-grained visual types into 8 super-	998
951	categories. <i>Diagrams</i> (including circuit, mechan-	999
952	ical, and structural diagrams) are the most dominant	1000
953	modality (49.9%), reflecting the dataset’s focus on	1001
954	professional schematics interpretation. To provide	1002
955	visual references for these categories, we display	1003
956	representative examples in Figure 5.	
957	D.5 Question Type Taxonomy	
958	Table 15 summarizes the distribution of answer for-	1004
959	formats within each macro subject. This table clarifies	1005
960	which subjects are dominated by multiple choice	
961	items versus numerical or descriptive responses.	1006
962	It also provides context for interpreting task-wise	1007
963	performance, since answer format affects both eval-	1008
964	uation difficulty and failure modes.	1009
965	D.6 Analysis of the Hard Subset	1010
966	To further analyze model performance on the most	1011
967	challenging instances, we identified a "Hard Sub-	1012
968	set" ($N = 1,053$). This subset consists of ques-	1013
969	tions that require complex multi-step reasoning	
970	and were frequently answered incorrectly during	1014
971	our pilot evaluation with open-source models. Ta-	1015
972	ble 17 shows the distribution of tasks within this	1016
973	hard subset. Notably, <i>Data and Visual Interpreta-</i>	1017
974	<i>tion</i> and <i>Concept Application</i> dominate this subset,	1018
975	suggesting that the primary difficulty for current	1019
976	models lies not in simple knowledge recall, but in	1020
977	the synthesis of visual data with domain concepts.	1021
978	We also report the prevalence of Korean-specific	1022
979	items by split in Table 18. The Hard set contains	1023
980	104 Korean-specific questions (9.92%), which is	
981	slightly higher than the Full set (8.63%), indicat-	1024
982	ing that localization-heavy items are somewhat over-	1025
983	represented among challenging instances.	1026
984	E Macro-task Performance Breakdown	1027
985	This appendix provides extended quantitative re-	1028
986	sults that are omitted from the main text due to	1029
987	space constraints. We report macro-task perfor-	1030
988	mance breakdowns and dataset composition statis-	1031
989	tics to support replication and alternative aggrega-	1032
990	tions.	1033
991	Table 19 reports per-model accuracy by macro	1034
992	task on the full split. This breakdown complements	1035
		1036
		1037
		1038
		1039
		1040
		1004
		1005
		1006
		1007
		1008
		1009
		1010
		1011
		1012
		1013
		1014
		1015
		1016
		1017
		1018
		1019
		1020
		1021
		1022
		1023
		1024
		1025
		1026
		1027
		1028
		1029
		1030
		1031
		1032
		1033
		1034
		1035
		1036
		1037
		1038
		1039
		1040

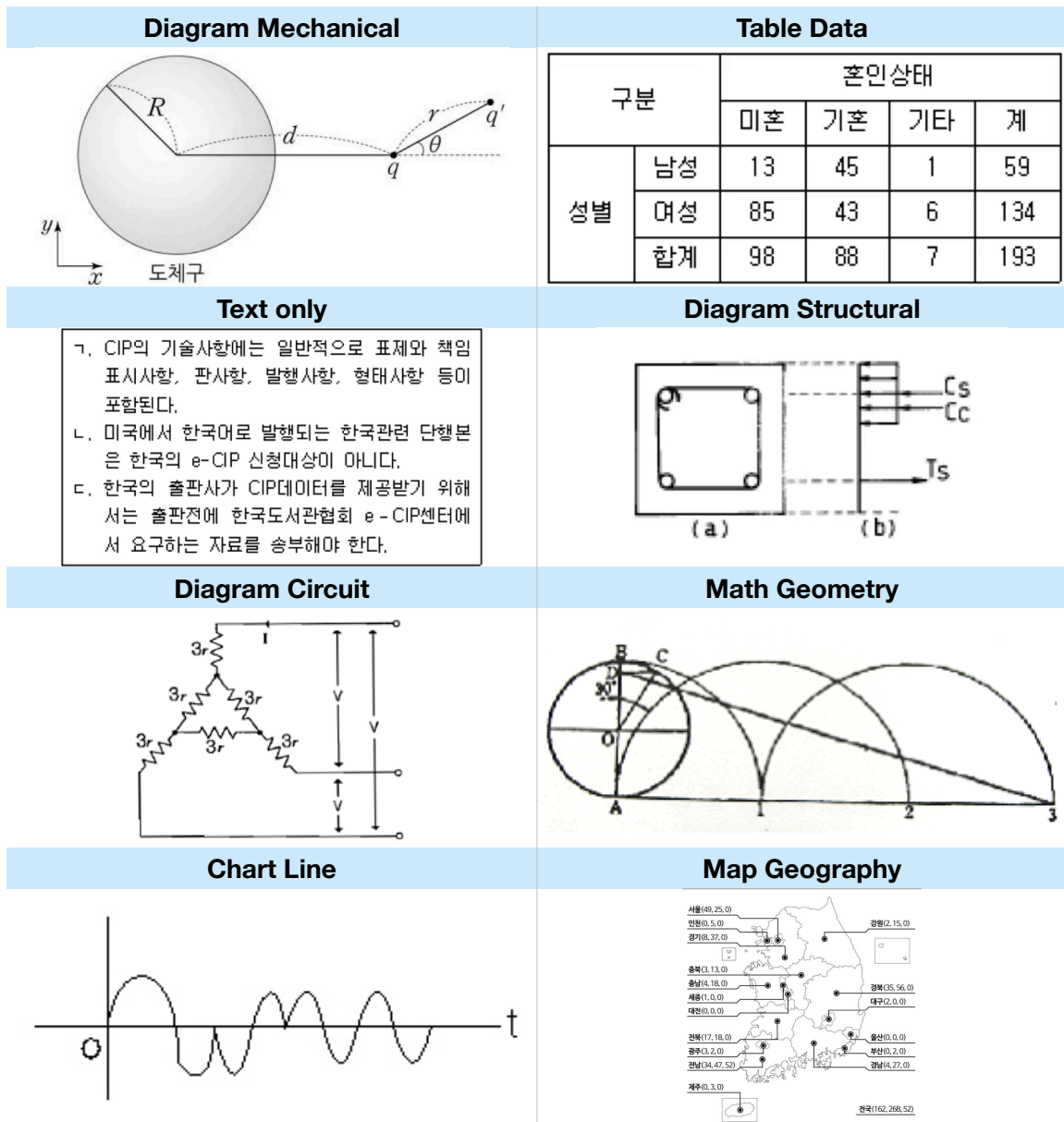


Figure 5: **Representative examples for each visual modality category.** KMMMU requires models to interpret various distinct visual formats, including specialized engineering diagrams and South Korean geographic maps.

G Additional Analysis on Prompt Sensitivity of LLM-as-a-Judge

We conduct a lightweight manual check to compare judge prompts against human verification. P2 matches human judgments most consistently. While P3 can reduce some format-related ambiguity by extracting a final answer, we observe cases where imperfect extraction leads to overly permissive verdicts, suggesting that explicit parsing can introduce additional failure modes.

Prompt-induced variation. Table 5 shows that measured accuracy can vary substantially across prompts even when the evaluated outputs are held

fixed. The effect is most pronounced for GEMINI-2.5-FLASH-LITE, indicating sensitivity to prompt-specific parsing and extraction behavior in mixed-format grading. In contrast, GEMINI-3-FLASH exhibits smaller variation, suggesting more stable grading under prompt changes.

Implications for evaluation. These results highlight that reported accuracy reflects not only model behavior but also the reliability of the grading protocol. Accordingly, prompt choice should be validated with spot-checks when grading mixed-format outputs.

1066
1067
1068
1069
1070
1071
1072

1073

1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115

G.1 LLM-as-a-Judge Prompt Templates

We adopt an LLM-as-a-judge protocol compatible with mixed-format questions. Table 23 lists the three templates used in our experiments, and the main text quantifies their prompt sensitivity under fixed evaluated outputs. Table 23 lists the three prompt templates used in our experiments.

H Detailed Regression Analysis

To rigorously quantify the impact of distinct failure modes on overall performance, we conducted an Ordinary Least Squares (OLS) regression analysis. The dependent variable was the overall accuracy of each model, while the independent variables were the error rates for the five error categories categorized in Section 6. We standardized all independent variables to unit variance to allow for a direct comparison of the resulting β coefficients. A negative coefficient indicates that an increase in that specific error type correlates with a decrease in overall accuracy, identifying it as a performance bottleneck.

Table 24 presents the standardized coefficients for both open-source and frontier model groups. For open-source models, the analysis identifies *Incomplete Response* as the most statistically significant predictor of performance ranking ($\beta = -0.088, p < 0.001$). This is closely followed by *Knowledge Shortage* ($\beta = -0.081, p < 0.001$). These figures suggest that the primary constraints for these models lie in generation stability and domain knowledge availability. Conversely, *Visual Perception Error* exhibits a relatively smaller impact ($\beta = -0.035$), consistent with a survivorship bias where models frequently fail due to instruction following or knowledge deficits before visual perception becomes the deciding factor.

In the frontier model regime, the regression coefficients indicate a fundamental shift in critical challenges. The impact of *Incomplete Response* diminishes to a negligible level ($\beta = -0.007$), confirming that generation stability issues are largely resolved. However, *Knowledge Shortage* remains the dominant predictor ($\beta = -0.052$) reflecting the high difficulty of the expert domains in KM-MMU. Most notably, *Visual Perception Error* rises in relative importance to become the second most critical bottleneck ($\beta = -0.045$). This statistical evidence supports the conclusion that expert-level multimodal proficiency is increasingly bounded by fine-grained visual understanding once foundational capabilities are established.

I Per Model Error Rates by Image Sub Category

This section provides a fine grained view of model behavior by reporting error rates for frequent image sub categories on the hard split. We include sub categories with at least 10 samples to reduce instability from small counts. For each sub category, we report the number of evaluated instances, the number of errors, and the resulting error rate. This breakdown is intended to support qualitative interpretation of the main text rather than to rank models on rare categories.

J Data Contamination Analysis Details

To assess the risk of memorization, we conduct a prefix completion test in which a model is given the first 35% of a question and asked to reconstruct the remaining 65%. We run this experiment on questions longer than 150 tokens to ensure that the continuation is non-trivial. We evaluate three frontier models under two settings. The first is a no-hint setting. The second is a hinted setting that provides the exam name and year, which can serve as a potential trigger for recall.

Evaluation with LLM-as-a-Judge We score reconstructions using an LLM-as-a-Judge framework based on GEMINI-3-FLASH. The judge compares the reconstructed continuation against the reference question and outputs an exactness score on a 0 to 100 scale. The rubric emphasizes lexical overlap, preservation of key details such as numbers, entities, and constraints, and overall faithfulness to the reference. The judge is instructed to be strict and to award high scores only when the reconstruction closely matches the original text. In addition, the judge flags whether the model refused and whether it hallucinated unsupported details.

Results and interpretation Table 26 reports mean exactness, median exactness, refusal rate, and hallucination rate for each source and hint setting. Two patterns are consistent across sources. First, GPT-5-MINI and GEMINI-3-PRO exhibit very high refusal rates, often above 95%, and near-zero exactness in both settings. This behavior is consistent with training that discourages reproducing exam materials, which reduces the practical risk of direct leakage. Second, GEMINI-3-FLASH attempts completions more frequently and achieves higher mean exactness, but the scores remain low overall and do not reliably increase in the hinted

1116
1117

1118
1119
1120
1121
1122
1123
1124
1125
1126
1127

1128

1129
1130
1131
1132
1133
1134
1135
1136
1137
1138

1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151

1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164

1165 setting. In some sources, hinted prompts reduce
 1166 exactness. Together, the absence of consistently
 1167 high exactness scores and the lack of systematic
 1168 improvement under hints suggest that benchmark
 1169 performance is unlikely to be explained by simple
 1170 memorization.

Subject Category	Count
Physics	474
Civil & Structural Engineering	332
Mechanical Engineering	223
Electrical & Electronics Engineering	221
Computer Science	164
Statistics & Probability	149
Mathematics	121
Business Administration & Management	104
Environmental Science & Engineering	99
Geography & Spatial Studies	95
Manufacturing & Production Engineering	90
Public Administration & Policy	89
Design & Visual Arts	88
Chemical Engineering	76
Chemistry	75
Law & Legal Studies	74
Materials Science & Metallurgy	69
Information Technology & Systems	69
General Knowledge & Interdisciplinary	67
Industrial & Systems Engineering	64
Architecture & Urban Studies	57
Human Resources & Organizational Studies	56
Linguistics & Language Studies	54
Biology	43
Safety, Risk & Reliability Engineering	42
Data Science & Analytics	41
Finance & Accounting	40
Sociology & Social Sciences	36
Economics	33
Earth & Geological Sciences	33
Library, Archival & Information Science	30
Transportation & Logistics	29
Agriculture & Life Sciences	27
Software & Programming	27
Astronomy & Space Science	26
Psychology	22
Cognitive & Behavioral Sciences	20
Biomedical & Health Sciences	20
Communication & Media Studies	17
Education & Pedagogy	16
Marketing & Consumer Studies	15
Ethics & Philosophy	12
Artificial Intelligence & Machine Learning	12
Energy & Power Systems	11
Sports & Exercise Science	4

Table 11: Distribution of Subject Categories in KM-MMU

Macro Task	Count	%
Concept Application	948	27.35
Data & Visual Interpretation	852	24.58
Quantitative Reasoning	516	14.89
Knowledge Recall	347	10.01
Language Understanding	183	5.28
Geometry & Spatial Reasoning	155	4.47
Information Extraction & Classification	134	3.87
Logic & Rule Reasoning	119	3.43
Optimization & Planning	107	3.09
Code & Debugging	58	1.67
Other	47	1.36

Table 12: Macro task distribution.

Question Type	Count	%
Multiple Choice (Single Answer)	2,851	82.25
Multiple Choice (Multiple Answers)	347	10.01
Short Answer (Numerical / Calculation)	176	5.08
Short Answer (Text)	69	2.00
Descriptive	23	0.66

Table 13: **Question type distribution.**

Super-category	Sub-category	Count
Diagrams (1,731)	diagram_mechanical	726
	diagram_structural	391
	diagram_circuit	330
	diagram_flowchart	103
	diagram_network	67
	diagram_geologic	43
	diagram_chemical	42
Text, Code & Documents (869)	diagram_biological	26
	text_only	456
	text_blank	276
	document_scan	105
Tables (549)	text_code	32
	table_data	508
Mathematics (88)	table_blank	40
	math_geometry	65
Charts & Plots (151)	math_function	23
	chart_line	109
	chart_bar	18
	chart_scatter	13
	chart_pie	5
Maps (23)	Others (mixed, box, radar)	6
	map_geographic	8
	map_topographic	7
	map_weather	5
Symbols & Art (22)	map_astronomy	3
	symbol_pictogram	16
Photographs (13)	painting	8
	photo_general	11
	Others (astronomy, organism)	2

Table 14: **KMMMU Dataset Visual Modality Taxonomy.** We organize fine-grained visual sub-categories into a set of semantically coherent super-categories to support structured analysis of multimodal reasoning. The dataset is rich in technical diagrams and document-based visuals.

Macro Subject	Multiple Choice (Single Answer)	Multiple Choice (Multiple Answers)	Short Answer (Text)	Descriptive	Short Answer (Numerical)
Engineering	1142	39	5	2	26
Natural Sciences	486	170	15	14	112
CS & IT	315	13	6	2	7
Business & Public	265	31	4	0	4
Math & Stats	175	64	9	1	21
Social Sciences	206	28	4	0	1
General	90	1	25	4	5
Arts & Design	88	0	0	0	0
Law & Ethics	84	1	1	0	0

Table 15: **Dataset attributes by macro subject and answer type.**

Table 16: Macro Subject Distribution in the KMMMU Hard Set

Macro Subject	Count
Engineering	357
Natural Sciences	267
Math & Stats	98
CS & IT	91
Business & Public Admin	76
Social Sciences	61
General/Interdisciplinary	52
Arts & Design	34
Law & Ethics	17
Total	1,053

Macro Task	Count	%
Data & Visual Interpretation	299	28.40
Concept Application	274	26.02
Quantitative Reasoning	123	11.68
Knowledge Recall	112	10.63
Language Understanding	67	6.36
Logic & Rule Reasoning	44	4.17
Optimization & Planning	39	3.70
Information Extraction & Classification	35	3.32
Geometry & Spatial Reasoning	33	3.13
Code & Debugging	15	1.42
Other	12	1.13

Table 17: **Macro task distribution on the hard subset.** Percentages are computed over the hard subset ($N = 1,053$).

Split	Korean Specific	Not Korean Specific	Ratio (%)
Full	299	3167	8.63
Hard	104	949	9.92

Table 18: **Number of Korean-specific questions each in Full and Hard sets.**

Model	Concept Application	Data & Visual Interpretation	Quantitative Reasoning	Knowledge Recall	Language Understanding	Geometry & Spatial Reasoning	Info Extraction & Classification	Logic & Rule Reasoning	Optimization & Planning	Code & Debugging	Other	Overall Acc.
clovaX-3b	13.470,64	11.110,38	8.140,67	17.480,33	19.671,89	7.741,71	23.132,89	20.451,94	13.082,47	15.521,72	9.932,46	13.160,13
gemma3-12b	20.250,11	16.741,36	8.850,62	22.380,60	13.300,63	9.250,74	18.912,62	22.973,97	14.641,95	14.943,59	9.221,23	16.680,42
gemma3-27b	20.960,34	18.540,41	13.181,21	24.881,92	14.392,28	10.111,49	23.383,11	18.492,22	16.201,95	22.996,97	18.442,46	18.630,36
gemma3-4b	14.240,02	8.841,42	5.490,99	21.231,09	14.032,70	4.520,65	13.180,43	12.042,57	6.851,43	12.074,56	8.513,69	11.410,50
koreason	47.781,11	60.331,23	76.682,20	38.421,30	36.072,89	75.480,65	46.771,14	51.543,88	66.360,00	47.701,99	75.894,43	55.900,33
llama-4-maverick	37.060,22	43.700,24	41.541,29	39.871,09	30.601,13	32.261,29	44.781,29	35.010,97	38.012,70	39.661,72	23.403,69	39.200,38
llama-4-scout	31.721,64	35.800,89	34.241,43	36.121,44	26.411,38	26.411,38	34.731,40	34.585,84	28.744,98	32.624,91	33.921,54	33.921,54
qwen3-2b	17.580,70	14.551,24	9.111,03	22.191,04	15.482,21	8.391,94	26.371,14	21.852,22	12.461,83	9.773,59	8.510,00	15.590,23
qwen3-30b-a3b	30.591,92	39.320,47	33.910,67	33.530,73	28.230,32	17.850,37	45.773,02	35.574,31	39.562,70	28.743,54	25.532,13	33.770,44
qwen3-30b-a3b-think	47.151,12	63.037,43	69.128,62	44.196,98	38.980,05	66.880,65	50.754,89	54.341,16	70.728,08	45.981,22	60.990,83	55.767,53
qwen3-32b	33.300,76	43.151,30	39.531,08	35.161,80	27.142,07	29.462,98	41.794,16	37.542,95	41.124,07	36.784,34	31.216,14	37.080,81
qwen3-32b-thinking	44.592,09	61.460,95	67.120,40	39.581,48	37.160,95	67.533,31	51.990,86	50.422,52	68.543,00	48.854,34	60.284,43	53.730,49
qwen3-4b	21.520,46	19.290,95	15.700,70	29.392,29	19.850,83	10.540,99	32.091,29	25.772,70	16.201,08	31.617,18	12.773,69	20.750,32
qwen3-4b	25.071,39	25.272,66	21.061,07	32.762,04	23.501,07	16.990,37	36.572,99	31.371,28	23.993,89	28.742,63	16.631,23	25.420,53
qwen3-VL-235B-A22B-thinking	45.250,18	61.270,41	65.962,50	36.702,58	39.530,01	64.522,23	48.012,62	53.500,28	71.960,93	46.551,72	58.873,25	53.390,19
qwen3-VL-235B-A22B-Instruct	34.771,68	47.810,82	45.541,51	35.640,88	31.880,01	28.392,92	44.531,55	39.781,75	45.172,35	35.062,63	34.756,84	40.100,31
varco2-1.7b	13.992,35	9.000,18	5.751,48	15.852,59	11.666,11	10.322,92	14.180,00	10.085,51	9.035,63	14.371,00	4.260,00	11.032,39
varco2-14b	24.300,60	19.480,31	21.961,48	26.133,21	20.223,94	24.950,99	30.852,15	26.053,85	16.200,54	27.591,72	19.862,46	22.820,16

Table 19: Accuracy (%) by macro task with overall accuracy. Values are reported as mean_{std} across three runs.

Model	Concept Application	Data & Visual Interpretation	Quantitative Reasoning	Knowledge Recall	Language Understanding	Geometry & Spatial Reasoning	Info Extraction & Classification	Logic & Rule Reasoning	Optimization & Planning	Code & Debugging	Other	Overall Acc.
Gemini-3-Pro	67.403,04	74.691,84	71.821,24	59.521,86	70.151,49	57.583,03	60.951,65	62.881,31	78.631,48	60.006,67	75.000,00	69.010,29
Gemini-3-Flash	61.070,56	64.552,74	65.851,73	58.332,73	49.254,48	44.444,63	53.337,19	50.764,73	64.963,92	48.893,85	55.561,23	60.270,15
GPT-5	48.302,43	61.430,97	58.541,41	36.903,14	42.791,72	31.313,50	39.057,19	46.973,47	67.521,48	42.223,85	58.338,33	51.500,88
GPT-5-mini	40.271,80	53.071,17	49.591,63	25.891,55	33.833,31	23.233,50	28.578,57	35.613,72	61.544,44	40.000,67	47.229,62	42.800,33
Claude-Opus-4.5	38.690,97	51.511,53	46.882,05	30.651,36	38.812,99	28.281,75	35.243,30	34.096,82	53.852,56	51.117,70	41.673,83	42.500,63
Claude-Sonnet-4.5	27.131,69	43.920,39	38.212,44	28.574,46	24.332,28	23.234,63	31.434,95	29.552,27	38.466,78	46.667,67	27.789,62	34.000,83
Grok-4	39.783,81	47.710,70	53.123,67	33.327,73	31.844,80	29.291,75	38.104,36	35.615,72	48.725,13	48.893,85	50.000,43	42.421,62
Grok-4.1-fast	35.403,12	30.210,51	47.700,47	23.210,89	23.885,38	31.311,75	26.673,30	21.211,31	41.022,56	24.443,85	50.000,33	32.541,10
Mistral-Large-3-675B-Instruct-2512	22.141,87	24.302,04	23.582,93	23.213,09	20.405,65	22.229,74	28.572,86	18.941,31	29.066,45	24.443,85	30.564,81	23.391,99

Table 20: Accuracy (%) on the hard subset by macro task. Values are reported as mean_{std} over three runs.

Model	#True	Acc _{True} (%)	#False	Acc _{False} (%)	Δ (pp)
KO-REASON-G3-12B-1009	299	32.66 _{2,04}	3167	58.01 _{0,17}	-25.35
Qwen3-32B-Thinking	299	31.77 _{1,16}	3167	55.83 _{0,45}	-24.05
Qwen3-VL-235B-A22B-Thinking	299	32.44 _{0,58}	3167	55.35 _{0,30}	-22.91
Qwen3-VL-30B-A3B-Thinking	299	35.45 _{7,91}	3167	57.62 _{7,50}	-22.17
Qwen3-VL-235B-A22B-IT	299	26.98 _{1,65}	3167	41.34 _{0,26}	-14.36
Llama-4-Maverick-17B-128E-IT	299	28.32 _{0,39}	3167	40.29 _{0,40}	-11.98
Qwen3-VL-32B-IT	299	26.31 _{3,35}	3167	38.12 _{0,82}	-11.81
Qwen3-VL-30B-A3B-IT	299	23.52 _{0,19}	3167	34.73 _{0,45}	-11.21
Llama-4-Scout-17B-16E-IT	299	27.20 _{0,39}	3167	34.69 _{0,62}	-7.49
VARCO-VISION-2.0-14B	299	22.30 _{1,58}	3167	22.99 _{0,07}	-0.70
Qwen3-VL-2B-IT	299	15.27 _{1,51}	3167	15.66 _{0,38}	-0.39
Qwen3-VL-4B-IT	299	20.62 _{2,92}	3167	20.86 _{0,34}	-0.24
Gemma-3-27B-IT	299	18.84 _{3,58}	3167	18.66 _{0,15}	+0.18
Qwen3-VL-8B-IT	299	25.98 _{3,37}	3167	25.46 _{0,52}	+0.51
VARCO-VISION-2.0-1.7B	299	11.71 _{3,48}	3167	11.00 _{2,29}	+0.71
Gemma-3-12B-IT	299	21.85 _{0,70}	3167	16.25 _{0,45}	+5.60
HyperCLOVAX-SEED-Vision-3B	299	18.73 _{1,86}	3167	12.69 _{0,30}	+6.04
Gemma-3-4B-IT	299	17.28 _{1,72}	3167	10.88 _{0,64}	+6.40

Table 21: Per-model accuracy on Korean-specific content (Full split). Values are reported in percentage using the notation mean_{std} across 3 runs. Δ is in percentage points (pp): $\Delta = Acc_{True} - Acc_{False}$, and models are sorted by Δ in ascending order (more negative indicates larger degradation on Korean-specific items).

Model	#True	Acc _{True}	#False	Acc _{False}	Δ (pp)
Gemini-3-Pro	104	58.01 _{1,47}	949	70.09 _{0,53}	-12.08
Gemini-3-Flash	104	52.88 _{1,67}	949	61.02 _{0,28}	-8.13
GPT-5	104	33.01 _{1,11}	949	53.50 _{0,92}	-20.48
GPT-5-mini	104	31.41 _{1,11}	949	43.96 _{0,28}	-12.55
Grok-4	104	30.45 _{2,00}	949	43.64 _{1,84}	-13.20
Grok-4.1-fast	104	19.55 _{1,11}	949	33.83 _{1,13}	-14.28
Claude-Opus-4.5	104	35.58 _{2,54}	949	43.36 _{0,43}	-7.78
Claude-Sonnet-4.5	104	32.37 _{2,22}	949	34.22 _{0,74}	-1.84
Mistral-Large-3-675B-Instruct-2512	104	26.28 _{2,42}	949	22.95 _{1,95}	+3.33

Table 22: Per-model accuracy on Korean-specific content (Hard split). Mean \pm std across 3 runs. Δ is in percentage points (pp): $\Delta = Acc_{True} - Acc_{False}$.

#1	<p>You will be given a question, gold answer, and a model response. Decide whether the model response matches the gold answer.</p> <ul style="list-style-type: none"> - Output EXACTLY ONE LINE (first line only): [TRUE] or [FALSE] Do not output anything else.
#2	<p>You will be given a question, gold answer, and a model response. Decide whether the model response matches the gold answer.</p> <p>Step 1) Determine question type:</p> <ul style="list-style-type: none"> - If the question contains explicit options/choices (e.g., ①②③④⑤, A/B/C/D, 1-5, ㄱ-ㄷ), treat it as MULTIPLE-CHOICE. - Otherwise treat it as FREE-FORM. <p>Step 2) For MULTIPLE-CHOICE:</p> <ul style="list-style-type: none"> - Extract the model's FINAL selected option only. - Normalize equivalent formats: <ul style="list-style-type: none"> - ①②③④⑤ == 1 2 3 4 5 - "4번", "정답은 ④", "answer: (4)" all mean option 4 - "A", "(A)", "Option A" all mean option A <p>Step 3) For FREE-FORM:</p> <ul style="list-style-type: none"> - Accept paraphrases that preserve the same meaning. - The response must contain the same final conclusion as the gold answer. <p>=====</p> <ul style="list-style-type: none"> - Output EXACTLY ONE LINE: [TRUE] or [FALSE] Do not output anything else.
#3	<p>You will be given a question, gold answer, and a model response. Decide whether the model response matches the gold answer.</p> <p>First, decide question type:</p> <ul style="list-style-type: none"> - MULTIPLE-CHOICE if explicit options exist, else FREE-FORM. <p>Second, extract the model's final answer:</p> <ul style="list-style-type: none"> - Extract the model's FINAL answer and normalize formats (①=1, "4번"=4, etc.). - Mark TRUE when the normalized final option equals the gold option. - If the response is partially correct, missing a required item, or contains any contradiction, output [FALSE]. - Do not reward verbosity. Do not infer missing information. <p>- Output EXACTLY ONE LINE: Final Answer: extracted final answer, Decision: [TRUE] or [FALSE] Do not output anything else.</p>

Table 23: Judge prompt templates used for model evaluation.

Error Type	Open-Source Models (β)	Frontier Models (β)
Incomplete Response	-0.088 ($p < 0.001$)	-0.007 ($p < 0.001$)
Knowledge Shortage	-0.081 ($p < 0.001$)	-0.052 ($p < 0.001$)
Reasoning Error	-0.052 ($p < 0.001$)	-0.042 ($p < 0.001$)
Visual Perception	-0.035 ($p < 0.001$)	-0.045 ($p < 0.001$)
Hallucination	-0.037 ($p < 0.001$)	-0.003 ($p = 0.002$)

Table 24: **OLS Regression Results: Predictors of Model Accuracy.** Values represent standardized beta coefficients (β). A larger negative value indicates that the error type is a stronger cause of performance degradation.

Image Sub-Category	Model	Count	Error	Error Rate (%)
chart_line	Gemini-3-Pro	32	11	34.38
	Gemini-3-Flash	32	10	31.25
	GPT-5	32	16	50.00
	Claude-Opus-4.5	32	20	62.50
	Grok-4	32	23	71.88
diagram_circuit	Gemini-3-Pro	109	41	37.61
	Gemini-3-Flash	109	49	44.95
	GPT-5	109	55	50.46
	Claude-Opus-4.5	109	77	70.64
	Grok-4	109	73	66.97
diagram_flowchart	Gemini-3-Pro	22	9	40.91
	Gemini-3-Flash	22	8	36.36
	GPT-5	22	12	54.55
	Claude-Opus-4.5	22	14	63.64
	Grok-4	22	13	59.09
diagram_mechanical	Gemini-3-Pro	242	59	24.38
	Gemini-3-Flash	242	88	36.36
	GPT-5	242	104	42.98
	Claude-Opus-4.5	242	135	55.79
	Grok-4	242	119	49.17
diagram_network	Gemini-3-Pro	19	8	42.11
	Gemini-3-Flash	19	10	52.63
	GPT-5	19	10	52.63
	Claude-Opus-4.5	19	12	63.16
	Grok-4	19	13	68.42
diagram_structural	Gemini-3-Pro	70	24	34.29
	Gemini-3-Flash	70	28	40.00
	GPT-5	70	44	62.86
	Claude-Opus-4.5	70	50	71.43
	Grok-4	70	40	57.14
document_scan	Gemini-3-Pro	19	4	21.05
	Gemini-3-Flash	19	5	26.32
	GPT-5	19	6	31.58
	Claude-Opus-4.5	19	10	52.63
	Grok-4	19	9	47.37
table_data	Gemini-3-Pro	188	32	17.02
	Gemini-3-Flash	188	50	26.60
	GPT-5	188	46	24.47
	Claude-Opus-4.5	188	64	34.04
	Grok-4	188	87	46.28
text_blank	Gemini-3-Pro	104	49	47.12
	Gemini-3-Flash	104	54	51.92
	GPT-5	104	84	80.77
	Claude-Opus-4.5	104	82	78.85
	Grok-4	104	80	76.92
text_only	Gemini-3-Pro	117	51	43.59
	Gemini-3-Flash	117	60	51.28
	GPT-5	117	75	64.10
	Claude-Opus-4.5	117	73	62.39
	Grok-4	117	85	72.65

Table 25: **Per model error rates by image sub category on the hard split.** We report the number of evaluated instances, the number of errors, and the resulting error rate. Only sub categories with at least 10 instances are shown to reduce noise from small counts.

Source	Model	Hinted	Mean Acc	Std	Median	Refusal %	Halluc %
PSAT	Gemini-3-Flash	False	14.20	13.39	15.0	17.53	76.10
		True	13.14	14.62	10.0	24.70	68.13
	Gemini-3-Pro	False	0.34	2.05	0.0	95.60	3.20
		True	0.46	2.03	0.0	94.78	3.21
	GPT-5-Mini	False	0.00	0.00	0.0	99.60	0.00
		True	0.00	0.00	0.0	100.00	0.00
NCS	Gemini-3-Flash	False	10.05	12.68	5.0	29.08	50.20
		True	6.54	11.02	0.0	50.20	35.57
	Gemini-3-Pro	False	0.18	2.54	0.0	99.21	0.40
		True	0.04	0.63	0.0	100.00	0.00
	GPT-5-Mini	False	0.65	3.99	0.0	95.65	1.19
		True	0.45	2.01	0.0	100.00	1.19
National Technical Qualification	Gemini-3-Flash	False	15.70	17.33	10.0	19.66	61.54
		True	18.58	16.97	15.0	15.61	76.84
	Gemini-3-Pro	False	0.33	2.31	0.0	93.86	4.44
		True	2.03	6.31	0.0	79.15	16.92
	GPT-5-Mini	False	0.45	4.13	0.0	97.96	0.51
		True	0.25	1.72	0.0	99.83	0.17
Olympiads	Gemini-3-Flash	False	10.48	12.45	5.0	23.12	55.58
		True	10.91	14.73	5.0	37.56	49.48
	Gemini-3-Pro	False	0.13	1.34	0.0	97.92	1.56
		True	0.38	2.46	0.0	97.41	2.33
	GPT-5-Mini	False	0.52	3.12	0.0	94.56	2.85
		True	0.18	1.72	0.0	98.96	1.04

Table 26: **Data contamination analysis via prefix completion.** We provide the first 35% of the question and measure the Exactness Score (0 to 100). Refusal and Hallucination rates are reported in percentages. Results are compared between the No-hint and Hinted settings across four data sources.