

From KMMLU-REDUX to PRO: A Professional Korean Benchmark Suite for LLM Evaluation

Anonymous ACL submission

Abstract

The development of Large Language Models (LLMs) requires robust benchmarks that encompass not only academic domains but also industrial fields to effectively evaluate their applicability in real-world scenarios. In this paper, we introduce two Korean expert-level benchmarks. **KMMLU-REDUX**, reconstructed from the existing KMMLU (Son et al., 2024a), consists of questions from the Korean National Technical Qualification exams, with critical errors removed to enhance reliability. **KMMLU-PRO** is based on Korean National Professional Licensure exams to reflect professional knowledge in Korea. Our experiments demonstrate that these benchmarks comprehensively represent industrial knowledge in Korea.

1 Introduction

As LLMs continue to achieve strong performance across a wide range of subjects (OpenAI et al., 2024b; Deepmind, 2024; DeepSeek-AI et al., 2025a; Research et al., 2025), the demand for comprehensive benchmarks has grown. MMLU (Hendrycks et al., 2021) is widely used for its broad coverage of general knowledge from elementary to college level. However, its publicly available online problems has raised concerns about reliability and potential data contamination (Gema et al., 2025; Vendrow et al., 2025; Zhao et al., 2024).

We identify similar issues in KMMLU (Son et al., 2024a), a widely used benchmark for evaluating Korean expert-level knowledge. The dataset was constructed by crawling websites that provide questions from various exams. We observe noisy samples, including problems that explicitly reveal the answer or non-existent reference, which can mislead performance evaluation. Additionally, we find evidence of contamination between the train and test splits, as well as with common web corpus such as FineWeb2 (Penedo et al., 2024).

Instead of collecting data from online sources directly, recent challenging benchmarks (Srivastava et al., 2023; Rein et al., 2024; Phan et al., 2025; Kazemi et al., 2025; Team et al., 2025b) have been constructed through problems and answers authored by human expert. Although this approach ensures high-quality, contamination-free benchmarks, it is costly to construct and maintain. The high construction cost hinders regular updates (White et al., 2025; Jain et al., 2025), leaving the benchmarks vulnerable to depreciation and potential contamination. Furthermore, existing benchmarks focus primarily on academic knowledge and evaluate models using a single overall score. They often overlook the practical applicability of models in industrial or professional contexts, specifically whether the models meet the standards for real-world professional certification.

We introduce two benchmarks to address the limitations above and reflect real-world professional knowledge. First, **KMMLU-REDUX**, a refined subset of KMMLU with 2,587 problems, is built through rigorous manual examination by authors to reduce errors and contamination within KMMLU. Furthermore, while the KMMLU contains problems from wide range of exams, even with high-school level, KMMLU-REDUX only selects Korean National Technical Qualification (KNTQ) exams as sources of the benchmark. The exams require applicants to have either a bachelor's degree or at least nine years experience in industrial field, thus making the benchmark more challenging.

Second, we build **KMMLU-PRO**, a new challenging benchmark, which consists of 2,822 problems from acquisition exams for Korean National Professional Licensure (KNPL), representing highly specialized professions in Korea. We include 14 professions from diverse domains. Unlike KMMLU that crawls websites, we collect data directly from the official source of each license. After that, human annotators manually examine it to

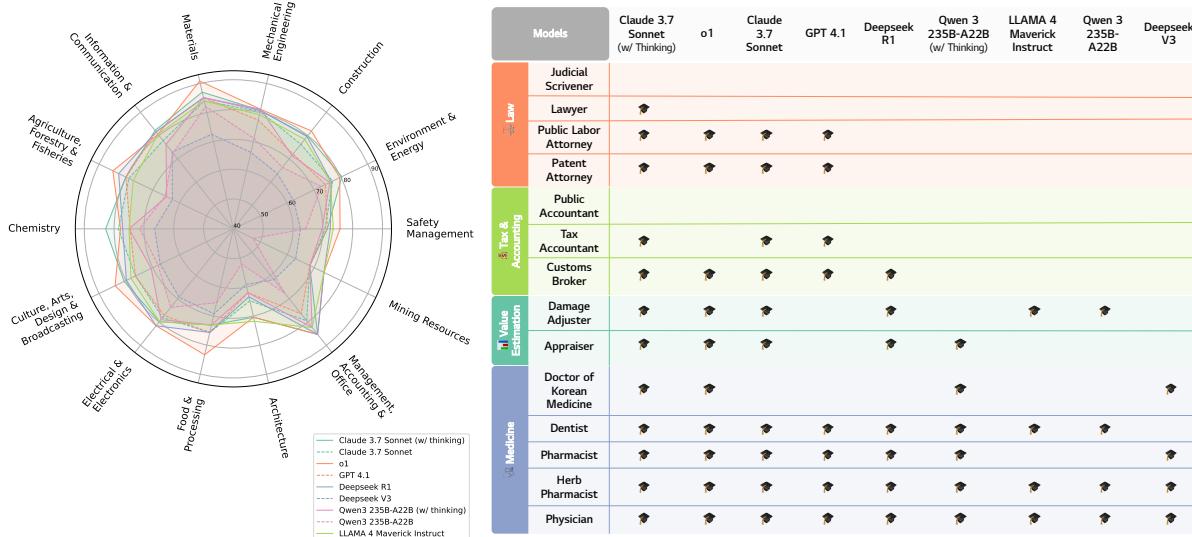


Figure 1: Model performance across industry domains on KMMLU-Redux (left) and profession acquisition status with 🎓 on KMMLU-Pro (right).

083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
avoid noises. KMMLU-PRO only includes exams held in the most recent year and would be updated annually with the latest exam to maintain long-term reliability and prevent contamination.

We conduct extensive evaluations of various LLMs on the two benchmarks. Our benchmarks, based on real-world exams, enable aligned analysis with industrial and professional qualifications, effectively revealing the practical strengths of each model. As shown in Figure 1, KMMLU-PRO assesses which professional license a model is capable of obtaining, while KMMLU-REDUX evaluates the breadth of industrial knowledge. The results show that several state-of-the-art models perform strongly in the medicine domain, passing most licenses, yet nearly fail in law-related licenses.

Moreover, we observe the significant performance gaps between our datasets and merely translated datasets like MMMLU (OpenAI, 2024), especially in domains such as laws, where in-depth knowledge of specific countries is required (see Section 6.1). We argue that these findings underscore the practicality of our benchmarks for assessing the capabilities of models in professional fields within Korea.

To summarize, our contributions are as follows:

- We improve the previous benchmark, KMMLU, to construct the refined and compact version of the benchmark, **KMMLU-REDUX**, by correcting various errors.
- We introduce **KMMLU-PRO**, a new bench-

mark designed to evaluate high-level professional knowledge in Korea. By imitating real-world license acquisition systems, KMMLU-PRO assesses the industrial practicality across various professions in Korea.

- We comprehensively analyze the results of two benchmarks, highlighting the importance of benchmarks specialized in Korea-specific professional knowledge.

2 KMMLU-REDUX

We first revisit KMMLU (Son et al., 2024a) to examine its quality. Building on these insights, we construct KMMLU-REDUX, a cleaned and compact version of the KMMLU. We carefully denoise against the KMMLU and increase difficulty.

2.1 Revisiting KMMLU

KMMLU plays a significant role in the NLP community as a de facto standard for evaluating LLMs on Korea-specific expert knowledge (Yoo et al., 2024; Research et al., 2024; Team et al., 2025a). The dataset was constructed by crawling websites¹ where various exam questions are uploaded by people online, spanning from high school to professional qualification exams. Upon closer inspection, we identify several noises and limitations, which can be categorized into three types: 1) duplication issues, 2) dataset errors, and 3) contamination.

Duplication Issue As the KMMLU crawls problems from hundreds of exams in Korea, we observe

¹<https://www.kinz.kr/>

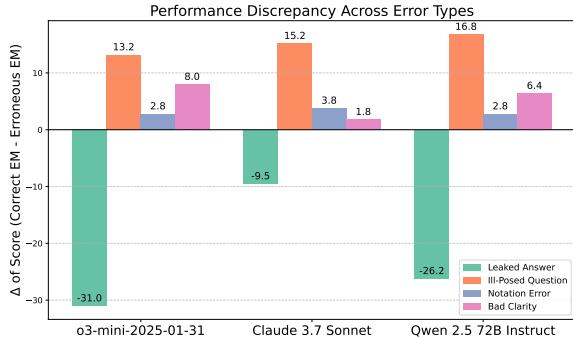


Figure 2: Performance differences in LLMs on the erroneous versus correct dataset. Leaked answer errors tend to overestimate model capabilities, while three other error types hinder LLMs to correctly predict true answers.

multiple duplicated questions across related exams. Notably, duplicated samples occur not only within the test set but also between the training and test sets. Using the Longest Common Sequence (LCS) algorithm to investigate overlaps, we could find 5.36% duplication within the test set and a 5.46% contamination between train and test set.

Dataset Errors Following Gema et al. (2025), we investigate the extent to which various error types appear in the KMMLU test set and their potential impact on LLM performance. We identify four representative error types: leaked answers, ill-posed questions, poor clarity, and notation errors. We then annotate the entire test set using GPT-4o² (OpenAI et al., 2024a). In total, we find that 7.66% of the data contains one of the errors mentioned above. We describe the details of the investigation of dataset errors in Appendix A.

In Figure 2, we observe significant discrepancies in the LLM’s performance between erroneous and correct samples. Notably, the performance of instances with leaked answers drops significantly when the LLMs are assessed on the clean dataset. Conversely, all models’ scores increase on ill-posed questions, underscoring their inability to identify the correct answer for poorly formulated questions.

Contamination The KMMLU was primarily sourced online, making them highly susceptible to contamination from web-crawled training corpora. When applying n-gram contamination detection (Lambert et al., 2025; Grattafiori et al., 2024) to the Korean subset of FineWeb2 (Penedo et al., 2024) and the KMMLU dataset, 1.88% of the data were flagged as contaminated.

²gpt-4o-2024-11-20

2.2 Dataset Construction

KMMLU consists of approximately 35k examples, making evaluation resource-intensive. To reduce this burden, we first restrict the scope to a subset of high-difficulty exams (Section 2.2.1). Furthermore, to ensure the reliability of the benchmark, we manually conduct a thorough examination of the dataset to eliminate errors (Section 2.2.2).

2.2.1 Filtering Non-Challenging Problems

Since KMMLU includes a variety of exams in Korea, spanning from high school to professional certification exams, we filter out the easier exams to build a more challenging benchmark. Specifically, we choose Korean National Technical Qualification (KNTQ) exams, which are primarily designed to assess practical technical competencies required in industrial field. The qualifications require applicants to have either a bachelor’s degree or at least nine years of professional experience. We adopt a collection of 100 KNTQ exams across the 14 domains in total. We only include the most recent exam for each qualification, thereby avoiding outdated knowledge being evaluated³.

To further discriminate simple problems, we leverage the performances of LLMs on the data (Wang et al., 2024; Zellers et al., 2019; Lee et al., 2023). By employing seven smaller LLMs⁴, we mark a data as *easy* if four or more models correctly predict its answer. Through this process, we remove 38.6% of the dataset.

2.2.2 Denoising

To remove noises, we follow the processes described in Section 2.1. Specifically, we first manually review the dataset to minimize errors. Next, we perform decontamination to prevent potential data leakage from pre-training corpora. Additionally, we detect inner duplication and against the training and test sets in KMMLU, thus finally remove all duplicates.

³We have compiled a list of all KNTQs exams alongside their most recent exam dates in Appendix B. Our aim is to make these available to LLM researchers and developers to help prevent data contamination.

⁴Llama 3.2 3B (Meta, 2024c), Qwen 2.5 3B (Qwen et al., 2025), Gemma 3 4B IT (Team, 2025a), Kanana Nano 2.1B Instruct (Team et al., 2025a), EXAONE 3.5 2.4B (Research et al., 2024), DeepSeek-R1-Distill-Qwen-1.5B (DeepSeek-AI et al., 2025a) EXAONE Deep 2.4B (Research et al., 2025), and Ko-R1-7B-v2.1 (OneLineAI).

Domain	Names of KNPLs	U.S. Equivalent	# of Instances
Law	Certified Judicial Scrivener	Paralegal, Legal Document Assistant, or Notary Public (no direct equivalent)	198
	Lawyer (Kim et al., 2024b)	Attorney-at-Law	150
	Certified Public Labor Attorney	Labor & Employment Lawyer (requires J.D. and bar admission; no separate certification in the U.S.)	239
	Certified Patent Attorney	Patent Attorney (JD + USPTO registration required)	109
Tax & Accounting	Certified Public Accountant (CPA)	Certified Public Accountant (CPA) – Exact Equivalent	208
	Certified Tax Accountant	Enrolled Agent (IRS) or CPA with Tax Specialization	238
	Certified Customs Broker	U.S. Customs Broker (licensed by U.S. Customs and Border Protection - CBP)	159
Value Estimation	Certified Damage Adjuster (CDA)	Claims Adjuster / Insurance Adjuster (state-licensed)	120
	Certified Appraiser	Certified Real Estate Appraiser (licensed at the state level)	196
Medicine	Doctor of Korean Medicine	Licensed Acupuncturist (L.Ac.) or Doctor of Acupuncture and Oriental Medicine (D.A.O.M.)	288
	Dentist (Kweon et al., 2024)	Doctor of Dental Surgery (D.D.S.) / Doctor of Dental Medicine (D.M.D.)	252
	Pharmacist (Kweon et al., 2024)	Doctor of Pharmacy (Pharm.D.)	271
	Herb Pharmacist	Herbalist (non-licensed or CAM-certified depending on state)	244
	Physician (Kweon et al., 2024)	Medical Doctor (M.D./D.O.)	150
Total			2822

Table 1: The list of National Professional Licenses (NPLs) used for KMMLU-PRO and their corresponding statistics. The names of NPLs are translated from those in Korea and we also report equivalent licences in U.S. We use KorMedMCQA (Kweon et al., 2024) for three licenses in the Medical category, and KBL (Kim et al., 2024b) for the bar exam of lawyer.

2.2.3 Final Statistics

For KMMLU-REDUX, we have collected 2,587 problems from 100 KNTQ exams. Among these, 596 problems are from exams that require over nine years of professional experience to acquire the qualification. To categorize the dataset into 14 domains, we follow the Korean Standard Industrial Classification (KSIC) published by Statistics Korea⁵ as the qualification system is primarily designed to align with industrial fields. Figure 6 in Appendix B illustrates the distribution of domain of KMMLU-REDUX.

3 KMMLU-PRO

A major challenge in building benchmarks from online sources is data contamination (Zhao et al., 2024; Jain et al., 2025; Roberts et al., 2024). While some studies address this by having experts manually create problems (Srivastava et al., 2023; Rein et al., 2024; Phan et al., 2025; Kazemi et al., 2025; Team et al., 2025b), this approach is costly and time-consuming. As an alternative, recent work explores periodic releases of fresh subsets. (White et al., 2025; Jain et al., 2025).

Motivated by these approaches, we focus on the Korean National Professional Licensure (KNPL) exams. These are high-stakes exams administered annually that pose a significant challenge. Unlike benchmarks crafted by a small set of experts (Rein et al., 2024; Phan et al., 2025), our approach leverages the well-established curricula of professional licensing systems, designed to assess real-world professional knowledge.

⁵<https://www.kostat.go.kr/>

3.1 Korean National Professional Licensures

We choose KNPL exams for main source of KMMLU-PRO. KNPL exams target high-level professionals, such as lawyers, accountants, or physicians, requiring advanced knowledge, critical reasoning, and ethical judgment. Among them, we select 14 KNPLs representing highly specialized and regulated professions in Korea (See Table 1 for the list of KNPLs and their equivalent licensure in U.S.). These licenses are legally mandated credentials required to practice in their respective domains. As such, they serve as institutionalized gateways to high-status occupations with significant entry barriers. Our evaluation simulates real-world assessment standards by incorporating official exam pass criteria, aligning model performance with human standards.

3.2 Dataset Collection and Annotation

In Korea, the government releases and manages the questions for KNPL acquisition exams. We directly download the PDF files from the government’s websites for each license and use GPT-4o (OpenAI et al., 2024a) for OCR parsing. As our dataset sources from the official PDFs, we can enhance the quality of the dataset, avoiding potential errors when collecting from online text (Team et al., 2025b). Since GPT-4o has difficulty handling tables and low-resolution PDFs, we employ human annotators to review parsed questions. When a problem contains an image, the annotators convert it into text that conveys the same meaning of the image, if possible⁶. Notably, our process remains relatively

⁶For further details, please see Appendix C

cost-efficient as it requires human annotators solely for error reviewing tasks, not full annotation.

We follow previous works, releasing a new set of questions periodically (White et al., 2025; Jain et al., 2025). Since the exams are conducted annually, we commit to collect and release questions from the exams held just before the current year.

3.3 Decontamination

We also adopt the same process outlined in Section 2.2.2 to ensure KMMLU-PRO is free from contamination. When conducting n-gram match between KMMLU-PRO, FineWeb2, and the training and validation sets of KMMLU, we did not find any contaminated examples. As a result, we can retain all 2,822 data points in the KMMLU-PRO, maintaining its contamination-free integrity.

4 Experiments

We select a diverse set of baseline models varying in size, multilingual capability, and reasoning ability. By default, we apply a zero-shot Chain-of-Thought (CoT) (Wei et al., 2022) prompt written in Korean. However, we observe that some models perform worse when prompted in Korean. For those models, we report results using the English prompt instead⁷. Our evaluation implementation is based on OpenAI’s simple-evals repository⁸. We use greedy decoding for non-reasoning models, while for reasoning models, we follow the settings in DeepSeek-AI et al. (2025a); Research et al. (2025), using a temerature of 0.6 and top-p (Holtzman et al., 2020) of 0.95. For more details, see Appendix D.

Metrics In addition to accuracy, our primary evaluation metric, we also report the number of licenses each LLM obtains in KMMLU-PRO. To better align with human evaluation standards, our procedure is desgned to mirror the official license acquisition criteria. However, for licenses that rely heavily on image-based questions, full replication is not possible. See Appendix D.3 for details.

5 Results

5.1 Main Results

Table 2 presents the overall performance on KMMLU-REDUX and KMMLU-PRO. The o1 model achieves the highest average accuracy

⁷For further details, please see Section 6.3.

⁸<https://github.com/openai/simple-evals>

(79.55), followed by Claude 3.7 with Thinking (78.49). Among open-weight models, DeepSeek’s R1 even outperforms many closed models. Models equipped with reasoning capabilities consistently perform better than their non-reasoning models, such as the Qwen3 series.

Beyond accuracy, we evaluate each model’s ability to obtain professional licenses under the official passing criteria for KMMLU-PRO. Specifically, in most licenses, a model must score at least 40% in each subject and achieve an overall average of 60% to pass. Claude 3.7 with Thinking obtains 12 out of 14 KNPL licenses, the highest among all evaluated models. In contrast, although the o1 model achieves higher accuracy, it qualifies for fewer licenses than Claude 3.7 with Thinking. This highlights the importance of balanced competence across subjects, as required in real-world certification exams.

5.2 Performance Across Industrial Domains in KMMLU-REDUX

Through KMMLU-REDUX, we assess the technical competencies of models across diverse industrial fields. As shown in Figure 1 (left), models with reasoning capabilities consistently outperform their non-reasoning counterparts across all domains. The improvement is particularly notable in Agriculture, Food & Processing, and Architecture. Despite these gains, all models continue to struggle in domains such as Mining Resources and Architecture, highlighting persistent challenges in underrepresented or highly specialized fields. The full results for all models across 14 domains are presented in Appendix E.3.

5.3 Professional License Acquisition Status on KMMLU-PRO

We provide a breakdown of the KMMLU-PRO results by analyzing which licenses are most frequently obtained by LLMs. Figure 1 (right) shows the pass rates of LLMs across 14 KNPLs. While many models obtain licenses in the medicine domain, most fail in Law and Tax & Accounting, with only DeepSeek R1 passing the Customs Broker exam among open-weight models. Notably, no models pass the Judicial Scrivener or Public Accountant exams. This trend becomes even more evident in the full results across a broader range of LLMs; the only licenses that relatively smaller models (<20B) are able to pass are in the medicine domain (see Table 8 in Appendix E.4).

Moreover, many LLMs fail to obtain licenses

	KMMLU-Redux Acc	KMMLU-Pro Acc	# of passed KNPLs	Avg. Acc (micro)
Open-weight Models				
Aya Expanse 32B (Dang et al., 2024)	33.05	31.26	0/14	32.12
Gemma 3 12B IT (Team, 2025a)	46.70	45.82	2/14	46.24
Phi-4 (14B) (Abdin et al., 2024)	49.75	45.32	1/14	47.44
EXAONE 3.5 32B Instruct (Research et al., 2024)	49.40	46.71	2/14	48.00
Mistral Small 3.1 Instruct (24B) (Mistral, 2025)	52.92	49.49	3/14	51.13
Gemma 3 27B IT (Team, 2025a)	54.04	51.03	2/14	52.47
Llama 3.3 70B Instruct (Meta, 2024a)	56.17	53.24	3/14	54.64
Qwen3-14B (Yang et al., 2025)	57.25	53.02	3/14	55.04
EXAONE Deep 32B (Research et al., 2025)	58.33	52.33	1/14	55.20
Qwen3-30B-A3B (Yang et al., 2025)	58.41	52.33	3/14	55.24
C4AI Command A (111B) (Cohere, 2025)	62.93	57.48	3/14	60.07
Qwen3-32B (Yang et al., 2025)	64.98	58.86	3/14	61.79
Llama-4-Scout-17B-16E-Instruct (Meta, 2025)	67.49	58.14	4/14	62.61
Qwen3-30B-A3B (w/thinking) (Yang et al., 2025)	65.25	60.52	3/14	62.78
Qwen3-14B (w/thinking) (Yang et al., 2025)	65.71	60.18	2/14	62.82
DeepSeek V3 (671B) (DeepSeek-AI et al., 2025b)	65.64	60.77	4/14	63.10
Qwen3-32B (w/thinking) (Yang et al., 2025)	68.77	61.14	3/14	64.79
QwQ 32B (Team, 2025b)	67.34	63.94	5/14	65.57
Qwen3-235B-A22B (Yang et al., 2025)	69.54	62.12	4/14	65.67
Qwen3-235B-A22B (w/thinking) (Yang et al., 2025)	74.49	68.22	6/14	71.22
Llama-4-Maverick-17B-128E-Instruct (Meta, 2025)	77.58	68.10	4/14	72.63
DeepSeek R1 (671B) (DeepSeek-AI et al., 2025a)	78.51	71.33	7/14	74.76
Closed Models				
GPT-4.1 mini (2025-04-14) (OpenAI, 2025a)	67.03	62.18	4/14	64.50
o3-mini (2025-01-31) (OpenAI, 2025c)	67.84	62.05	3/14	64.82
Grok-3-mini-beta (xAI, 2025)	71.47	65.08	5/14	68.14
Grok-3-beta (xAI, 2025)	72.90	68.37	7/14	70.54
o4-mini (2025-04-16) (OpenAI, 2025b)	75.80	69.65	6/14	72.59
GPT-4.1 (2025-04-14) (OpenAI, 2025a)	75.86	72.99	10/14	74.36
Claude 3.7 Sonnet (Anthropic, 2025)	76.88	74.52	10/14	75.65
o3 (OpenAI, 2025b)	79.92	73.60	9/14	76.62
Claude 3.7 Sonnet (w/thinking) (Anthropic, 2025)	79.36	77.70	12/14	78.49
o1 (OpenAI et al., 2024b)	81.14	78.09	10/14	79.55

Table 2: The main evaluation results of KMMLU-REDUX and KMMLU-PRO benchmarks on various LLMs. The gray-shaded models stand for reasoning models. The results of models with size < 10B are presented in Table 6 in Appendix E.2.

even when scoring above 60%; for example, o3-mini, Qwen3-235B-A22B, and Llama-4-Maverick score over 85% on the Pharmacist exam but still fail to qualify due to not meeting the threshold of law-related subject in the exam. These cases highlight the difficulty of acquiring region-specific domain knowledge, particularly in legal subjects governed by Korean law.

6 Analysis

6.1 The Importance of Locally Adapted Benchmarks

To highlight the importance of evaluation grounded in local context (Plaza et al., 2024; Singh et al., 2025), we compare category-level performance be-

tween benchmarks translated from English and KMMLU-PRO. Specifically, we focus on subjects related to law, accounting, and medicine,⁹ selecting relevant subjects from the Korean subset of MMMLU (OpenAI, 2024), as well as KMMLU-PRO.

As shown in Figure 3, the performance gap is relatively small in categories such as medicine, where domain knowledge is largely consistent across countries and cultures. In contrast, categories such as law, where substantial differences in content are expected, show a significantly larger gap. This sug-

⁹{professional_law, jurisprudence, international_law} for Law. {professional_accounting} for Accounting. {professional_medicine, clinical_knowledge, college_medicine, medical_genetics, anatomy} for Medicine.

Model Performance Comparison on {Medicine/Accounting/Law}-Related Problems

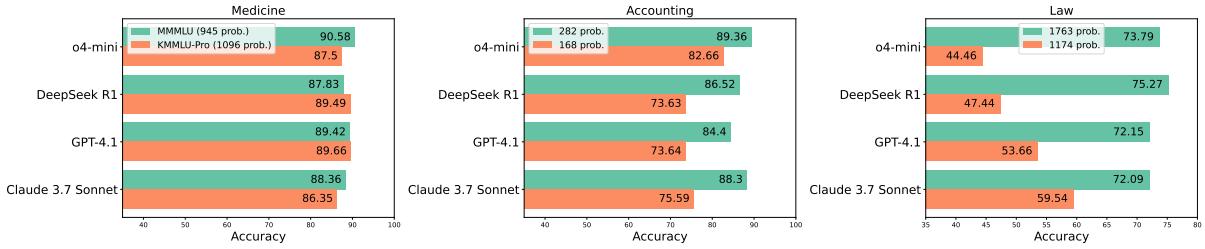


Figure 3: Performance of four LLMs on {Medical(left), Accounting(center), Law(right)}-relevant subsets from the MMMLU (Korean) (OpenAI, 2024) and KMMLU-PRO. While the discrepancies in scores are narrow in the medicine domain, they are wider in law-related problems, emphasizing the need for datasets that reflect real professional knowledge in Korea.

gests that MMMLU, which relies on direct translation of law questions based on U.S. standards, cannot adequately represent knowledge of Korean law. These findings highlight the importance of our dataset, which reflects authentic professional knowledge specific to the Korean context.

6.2 Impact of Reasoning Budget

Recent studies have shown that increasing reasoning efforts can enhance model performance (Muenninghoff et al., 2025; Anthropic, 2025; Qwen et al., 2025). To examine how reasoning *budget* affects performance on individual licenses in KMMLU-PRO, we conduct an experiment varying the number of tokens allocated to the reasoning path. We adopt Qwen3-32B (Yang et al., 2025) and Claude 3.7 Sonnet (Anthropic, 2025), for each budget $b \in B$, we generate $n = 4$ for Qwen3 and $n = 2$ for Claude.

As shown in Figure 4, we observe a positive correlation between the reasoning budget and the overall score of KMMLU-PRO for both models. However, this trend does not hold uniformly across all licenses. For example, we observe little to no improvement on the Judicial Scrivener and Herb Pharmacist licenses for both models, indicating that increased reasoning budget does not result in performance gains for certain licenses.

6.3 Impact of Prompt Language

Prompt language can significantly influence model behavior, raising concerns about consistency in multilingual settings (Wang et al., 2025; Zhang et al., 2025; Lai and Nissim, 2024). Since all exam questions in our datasets are written in Korean, using Korean prompts is a natural choice. However, we observe that some models perform worse when prompted in Korean. Table 3 presents the performance difference between English and Korean prompts. The Llama-4 model series exhibits

	KMMLU-REDUX			KMMLU-PRO		
	English	Korean	diff (%)	English	Korean	diff (%)
Qwen3-32B (w/ thinking)	68.77	69.08	+0.5%	61.14	60.66	-0.8%
o4-mini (2023-04-16)	75.80	76.17	+0.5%	69.65	69.10	-0.8%
Qwen3-235B-A22B (w/ thinking)	74.49	75.11	+0.8%	68.22	66.98	-1.8%
Grok-3-mini-beta	71.47	70.85	-0.9%	65.08	64.89	-0.3%
Qwen3-14B (w/ thinking)	65.71	65.40	-0.5%	60.18	59.48	-1.2%
EXAONE Deep 32B	58.33	56.17	-3.7%	52.33	52.19	-0.3%
DeepSeek R1 (671B)	78.51	75.38	-4.0%	71.33	70.62	-1.0%
QwQ 32B	67.34	62.66	-6.9%	63.94	59.95	-6.2%
EXAONE Deep 7.8B	44.99	40.82	-9.3%	41.53	38.98	-6.1%
Llama-4-Maverick-17B-128E-Instruct	77.58	72.52	-7.0%	68.10	57.15	-16.1%
Llama-4-Scout-17B-16E-Instruct	67.49	45.03	-33.3%	58.14	28.74	-50.6%

Table 3: Comparison results between English and Korean prompts of models whose main results are evaluated on English prompts. The *diff* values are the relative difference in scores between two prompts. The specific prompts used are detailed in Appendix D.1.

the most substantial drop in performance, while closed models such as Grok-3-mini-beta and o4-mini show minimal change.

7 Related Works

Reliability Issues of Benchmarks Recent works (Gema et al., 2025; Vendrow et al., 2025) have raised concerns about the reliability of LLM benchmarks due to dataset noise and contamination. MMLU-Redux (Gema et al., 2025) improved evaluation quality through systematic human reannotation, while GSM8K-Platinum (Vendrow et al., 2025) refined arithmetic benchmarks via automated and manual error detection. MMLU-CF (Zhao et al., 2024) prevents both unintentional and malicious contamination via sourcing diverse domains and question rewriting. LiveBench (White et al., 2025) and LiveCodeBench (Jain et al., 2025) adopted dynamic evaluation protocols with temporal cutoffs to prevent future leakage.

Professional Benchmark With the rapid advancement of LLMs, more challenging benchmarks have become essential. GPQA (Rein et al., 2024) and SuperGPQA (Team et al., 2025b) assess graduate-level knowledge; MMLU-Pro (Wang

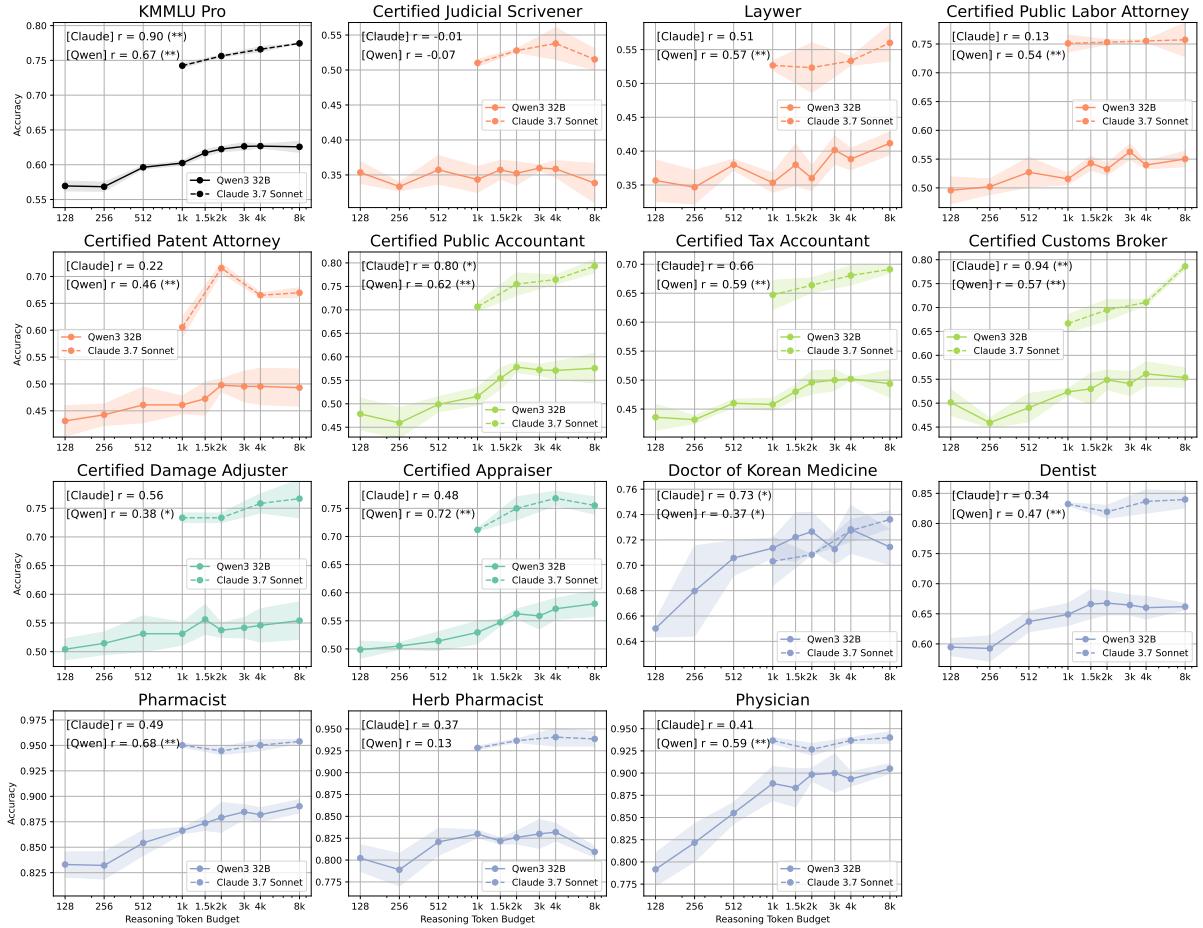


Figure 4: Reasoning budget results of Qwen3-32B and Claude 3.7 Sonnet on KMMLU-PRO. r indicates the Pearson Correlation coefficient value between the reasoning token budget and the accuracy. The responses are sampled multiple times for each thinking budget setting; $n = 4$ and $n = 2$, respectively, for Qwen and Claude.

et al., 2024) extends MMLU by increasing the share of college-level questions and expanding answer choices. Humanity’s Last Exam (Phan et al., 2025) introduces a frontier benchmark composed of manually authored, research-level questions.

Korean Benchmark While prior benchmarks focus primarily on English, recent efforts have produced Korean-specific evaluations (Son et al., 2024b; Kim et al., 2024a). Some rely on translated datasets (Park et al., 2024; Kim et al., 2025; OpenAI, 2024; Singh et al., 2024), often with human post-editing, but these lack regional context, institutional norms, and domain-specific fluency. In contrast, native Korean benchmarks such as KMMLU (Son et al., 2024a), KorMedMCQA (Kweon et al., 2024), and KBL (Kim et al., 2024b) address cultural and linguistic specificity. However, KMMLU (Son et al., 2024a) suffers from quality issues, including leaked answers, and KorMedMCQA (Kweon et al., 2024) and KBL (Kim

et al., 2024b) are limited to narrow domains.

In this work, we introduce KMMLU-REDUX and KMMLU-PRO, two contamination-free, industry grade benchmarks, providing a practical assessment of LLM capabilities in Korean industries.

8 Conclusion

We present two benchmarks constructed from real-world professional licensing exams, designed to reflect industrial domain knowledge and practical application standards. To ensure reliability, we identify and eliminate various sources of errors. Through extensive experiments, we evaluate the professional knowledge capabilities of LLMs across a wide range of domains. Our analysis further identifies key factors that influence performance, including region-specific knowledge, reasoning budget, and prompt language. We hope this work provides a foundation for more rigorous evaluation and continued advancement of real-world competence in language models.

504 Limitations

505 Our benchmarks are limited to text-only and
506 multiple-choice questions for text-only LLMs. It
507 restricts its coverage of real-world licensure
508 exams. Many real-word professional qualification
509 exams include non-textual modalities or require
510 constructed responses such as essay. Our bench-
511 mark cannot fully assess all aspects of professional
512 competence or reasoning required in such exams.
513 Expanding to multimodal inputs and open-ended
514 question formats is an important direction for fu-
515 ture work.

516 Ethical Statements

517 All data used in our benchmarks are either pub-
518 licly available or collected from official licensing
519 materials released by government or professional
520 institutions. For quality control, we hired human
521 annotators to review parsed questions from PDF;
522 they were compensated over the minimum wage in
523 Korea. Our benchmarks would be released under
524 CC-BY-NC-ND 4.0 license.

525 References

526 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien
527 Bubeck, Ronen Eldan, Suriya Gunasekar, Michael
528 Harrison, Russell J. Hewett, Mojan Javaheripi, Piero
529 Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li,
530 Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric
531 Price, Gustavo de Rosa, Olli Saarikivi, Adil Salim,
532 Shital Shah, Xin Wang, Rachel Ward, Yue Wu, Dingli
533 Yu, Cyril Zhang, and Yi Zhang. 2024. *Phi-4 technical*
534 *report*. Preprint, arXiv:2412.08905.

535 Anthropic. 2025. *Claude 3.7 sonnet* and *claude code*.

536 Cohere. 2025. *Model card for c4ai command a.*

537 John Dang, Shivalika Singh, Daniel D’souza, Arash
538 Ahmadian, Alejandro Salamanca, Madeline Smith,
539 Aidan Peppin, Sungjin Hong, Manoj Govindassamy,
540 Terrence Zhao, Sandra Kublik, Meor Amer, Viraat
541 Aryabumi, Jon Ander Campos, Yi-Chern Tan, Tom
542 Kocmi, Florian Strub, Nathan Grinsztajn, Yannis
543 Flet-Berliac, Acyr Locatelli, Hangyu Lin, Dwarak
544 Talupuru, Bharat Venkitesh, David Cairuz, Bowen
545 Yang, Tim Chung, Wei-Yin Ko, Sylvie Shang Shi,
546 Amir Shukayev, Sammie Bae, Aleksandra Piktus, Ro-
547 man Castagné, Felipe Cruz-Salinas, Eddie Kim, Lu-
548 cas Crawhall-Stein, Adrien Morisot, Sudip Roy, Phil
549 Blunsom, Ivan Zhang, Aidan Gomez, Nick Frosst,
550 Marzieh Fadaee, Beyza Ermis, Ahmet Üstün, and
551 Sara Hooker. 2024. *Aya expanse: Combining re-*
552 *search breakthroughs for a new multilingual frontier*.
553 Preprint, arXiv:2412.04261.

554 Google Deepmind. 2024. *Introducing gemini 2.0: our*
555 *new ai model for the agentic era*.

556 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,
557 Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
558 Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,
559 Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong
560 Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue,
561 Bingxuan Wang, Bochao Wu, Bei Feng, Chengda
562 Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang,
563 Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji,
564 Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo,
565 Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang,
566 Han Bao, Hanwei Xu, Haocheng Wang, Honghui
567 Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li,
568 Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang
569 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L.
570 Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai
571 Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai
572 Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong
573 Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan
574 Zhang, Minghua Zhang, Minghui Tang, Meng Li,
575 Miaojun Wang, Mingming Li, Ning Tian, Panpan
576 Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen,
577 Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan,
578 Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen,
579 Shanghao Lu, Shangyan Zhou, Shanhuang Chen,
580 Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng
581 Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing
582 Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun,
583 T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu,
584 Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao
585 Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan
586 Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin
587 Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li,
588 Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin,
589 Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang
590 Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang,
591 Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang
592 Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng
593 Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi,
594 Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang,
595 Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,
596 Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yu-
597 jia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You,
598 Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu,
599 Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu,
600 Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan,
601 Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean
602 Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao,
603 Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zi-
604 jia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song,
605 Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu
606 Zhang, and Zhen Zhang. 2025a. *Deepseek-r1: In-*
607 *centivizing reasoning capability in llms via reinforce-*
608 *ment learning*. Preprint, arXiv:2501.12948.

609 DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan
610 Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
611 Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai
612 Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie
613 Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo,
614 Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang,
615 Han Bao, Hanwei Xu, Haocheng Wang, Haowei

616 Zhang, Honghui Ding, Huajian Xin, Huazuo Gao,
 617 Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo,
 618 Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang
 619 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junx-
 620 iao Song, Kai Dong, Kai Hu, Kaige Gao, Kang
 621 Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong
 622 Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang,
 623 Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan
 624 Zhang, Minghua Zhang, Minghui Tang, Mingming
 625 Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng
 626 Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen,
 627 Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong
 628 Zhang, Ruijie Pan, Runji Wang, Runxin Xu, Ruoyu
 629 Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan
 630 Zhou, Shanhua Chen, Shaoqing Wu, Shengfeng
 631 Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang
 632 Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan,
 633 T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao,
 634 Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu,
 635 Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao
 636 Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao
 637 Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xia-
 638 oakkang Chen, Xiaokang Zhang, Xiaosha Chen, Xiao-
 639 tao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng,
 640 Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xin-
 641 nan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang,
 642 Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li,
 643 Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yan-
 644 hong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao
 645 Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu,
 646 Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong,
 647 Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yix-
 648 uan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,
 649 Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue
 650 Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan
 651 Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxi-
 652 ang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z.
 653 Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu,
 654 Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan
 655 Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhi-
 656 gang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu,
 657 Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu,
 658 Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi
 659 Gao, and Zizheng Pan. 2025b. Deepseek-v3 tech-
 660 nical report. Preprint, arXiv:2412.19437.

661 Aryo Pradipta Gema, Joshua Ong Jun Leang, Gi-
 662 won Hong, Alessio Devoto, Alberto Carlo Maria
 663 Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xi-
 664 aotang Du, Mohammad Reza Ghasemi Madani,
 665 Claire Barale, Robert McHardy, Joshua Harris, Jean
 666 Kaddour, Emile van Krieken, and Pasquale Min-
 667 ervini. 2025. Are we done with mmlu? Preprint,
 668 arXiv:2406.04127.

669 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
 670 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
 671 Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-
 672 ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh
 673 Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra,
 674 Archie Sravankumar, Artem Korenev, Arthur
 675 Hinsvark, Arun Rao, Aston Zhang, Aurelien Ro-
 676 driguez, Austen Gregerson, Ava Spataru, Baptiste
 677 Roziere, Bethany Biron, Binh Tang, Bobbie Chern,

678 Charlotte Caucheteux, Chaya Nayak, Chloe Bi,
 679 Chris Marra, Chris McConnell, Christian Keller,
 680 Christophe Touret, Chunyang Wu, Corinne Wong,
 681 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-
 682 lonsius, Daniel Song, Danielle Pintz, Danny Livshits,
 683 Danny Wyatt, David Esiobu, Dhruv Choudhary,
 684 Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
 685 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,
 686 Elina Lobanova, Emily Dinan, Eric Michael Smith,
 687 Filip Radenovic, Francisco Guzmán, Frank Zhang,
 688 Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis An-
 689 derson, Govind Thattai, Graeme Nail, Gregoire Mi-
 690 alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,
 691 Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
 692 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-
 693 han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
 694 Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,
 695 Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
 696 Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,
 697 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang,
 698 Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,
 699 Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-
 700 teng Jia, Kalyan Vasuden Alwala, Karthik Prasad,
 701 Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth
 702 Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,
 703 Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal
 704 Lakhota, Lauren Rantala-Yeary, Laurens van der
 705 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,
 706 Louis Martin, Lovish Madaan, Lubo Malo, Lukas
 707 Blecher, Lukas Landzaat, Luke de Oliveira, Madeline
 708 Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar
 709 Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew
 710 Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-
 711 badur, Mike Lewis, Min Si, Mitesh Kumar Singh,
 712 Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay
 713 Bashlykov, Nikolay Bogoychev, Niladri Chatterji,
 714 Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick
 715 Alrassy, Pengchuan Zhang, Pengwei Li, Petar Va-
 716 sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,
 717 Praveen Krishnan, Punit Singh Koura, Puxin Xu,
 718 Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj
 719 Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,
 720 Robert Stojnic, Roberta Raileanu, Rohan Maheswari,
 721 Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-
 722 nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan
 723 Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-
 724 hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-
 725 hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-
 726 ran Narang, Sharath Raparthi, Sheng Shen, Shengye
 727 Wan, Shruti Bhosale, Shun Zhang, Simon Van-
 728 denhende, Soumya Batra, Spencer Whitman, Sten
 729 Sootla, Stephane Collot, Suchin Gururangan, Syd-
 730 ney Borodinsky, Tamar Herman, Tara Fowler, Tarek
 731 Sheasha, Thomas Georgiou, Thomas Scialom, Tobias
 732 Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal
 733 Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh
 734 Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-
 735 ginie Do, Vish Vogeti, Vitor Albiero, Vladan Petro-
 736 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-
 737 ney Meers, Xavier Martinet, Xiaodong Wang, Xi-
 738 aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-
 739 feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-
 740 schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yi-
 741 wen Song, Yuchen Zhang, Yue Li, Yuning Mao,

742	Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Bader, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojaianzeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Nor-	man Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghatham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shahn Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models . Preprint, arXiv:2407.21783.	806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843
743	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding . In <i>International Conference on Learning Representations</i> .	844 845 846 847 848	
744	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration . In <i>International Conference on Learning Representations</i> .	849 850 851 852	
745	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fan-jia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2025. Live-codebench: Holistic and contamination free evaluation of large language models for code . In <i>The Thirteenth International Conference on Learning Representations</i> .	853 854 855 856 857 858 859	
746	Mehran Kazemi, Bahare Fatemi, Hritik Bansal, John Palowitch, Chrysovalantis Anastasiou, Sanket Vaibhav Mehta, Lalit K. Jain, Virginia Aglietti, Disha Jindal, Peter Chen, Nishanth Dikkala, Gladys Tyen, Xin Liu, Uri Shalit, Silvia Chiappa, Kate Olszewska, Yi Tay, Vinh Q. Tran, Quoc V. Le, and Orhan	860 861 862 863 864 865	
747			
748			
749			
750			
751			
752			
753			
754			
755			
756			
757			
758			
759			
760			
761			
762			
763			
764			
765			
766			
767			
768			
769			
770			
771			
772			
773			
774			
775			
776			
777			
778			
779			
780			
781			
782			
783			
784			
785			
786			
787			
788			
789			
790			
791			
792			
793			
794			
795			
796			
797			
798			
799			
800			
801			
802			
803			
804			
805			

866	Firat. 2025. Big-bench extra hard. <i>Preprint</i> , arXiv:2502.19187.	923
867		924
868	Eunsu Kim, Juyoung Suk, Philhoon Oh, Haneul Yoo, James Thorne, and Alice Oh. 2024a. CLICK: A benchmark dataset of cultural and linguistic intelligence in Korean. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 3335–3346, Torino, Italia. ELRA and ICCL.	925
869		926
870		927
871		928
872		929
873		
874		
875		
876	Hyeonwoo Kim, Dahyun Kim, Jihoo Kim, Sukyung Lee, Yungi Kim, and Chanjun Park. 2025. Open ko-llm leaderboard2: Bridging foundational and practical evaluation for korean llms. <i>Preprint</i> , arXiv:2410.12445.	930
877		931
878		932
879		
880		
881	Yeeun Kim, Youngrok Choi, Eunkyoung Choi, JinHwan Choi, Hai Jin Park, and Wonseok Hwang. 2024b. Developing a pragmatic benchmark for assessing Korean legal language understanding in large language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 5573–5595, Miami, Florida, USA. Association for Computational Linguistics.	933
882		934
883		
884		
885		
886		
887		
888		
889	Sunjun Kweon, Byungjin Choi, Gyouk Chu, Junyeong Song, Daewon Hyeon, Sujin Gan, Jueon Kim, Minkyu Kim, Rae Woong Park, and Edward Choi. 2024. Ko-rmedmcqa: Multi-choice question answering benchmark for korean healthcare professional licensing examinations. <i>Preprint</i> , arXiv:2403.01469.	935
890		936
891		
892		
893		
894		
895	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	937
896		938
897		939
898		940
899		941
900		942
901		
902	Huiyuan Lai and Malvina Nissim. 2024. mCoT: Multilingual instruction tuning for reasoning consistency in language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12012–12026, Bangkok, Thailand. Association for Computational Linguistics.	943
903		944
904		
905		
906		
907		
908		
909	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. 2025. Tulu 3: Pushing frontiers in open language model post-training. <i>Preprint</i> , arXiv:2411.15124.	945
910		946
911		947
912		
913		
914		
915		
916		
917		
918		
919	Hwaran Lee, Seokhee Hong, Joonsuk Park, Takyoung Kim, Meeyoung Cha, Yejin Choi, Byoungpil Kim, Gunhee Kim, Eun-Ju Lee, Yong Lim, Alice Oh, Sangchul Park, and Jung-Woo Ha. 2023. SQuARE:	948
920		949
921		950
922		951
923	A large-scale dataset of sensitive questions and acceptable responses created through human-machine collaboration. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6692–6712, Toronto, Canada. Association for Computational Linguistics.	952
924		953
925		954
926		955
927		956
928		957
929		958
930	Meta. 2024a. The future of ai: Built with llama.	959
931	Meta. 2024b. Introducing llama 3.1: Our most capable models to date.	960
932	Meta. 2024c. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models.	961
933	Meta. 2025. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation.	962
934	Mistral. 2025. Mistral small 3.1.	963
935	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. <i>Preprint</i> , arXiv:2501.19393.	964
936	naver hyperclovax. 2025. Model card for hyperclovax-seed-text-instruct-1.5b.	965
937	OneLineAI. Ko-r1-7b-v2.1. https://huggingface.co/OLAIR/ko-r1-7b-v2.1 . Accessed: 26 March 2025.	966
938	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guaraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy,	967
939		968
940		969
941		970
942		971
943		972
944		973
945		974
946		975
947		976
948		977
949		978
950		979
951		
952		
953		
954		
955		
956		
957		
958		
959		
960		
961		
962		
963		
964		
965		
966		
967		
968		
969		
970		
971		
972		
973		
974		
975		
976		
977		
978		
979		

980	David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian O'Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lillian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shiron Wu, Shuaiqi Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. 2024a. Gpt-4o system card . Preprint, arXiv:2410.21276.	1044
981		1045
982		1046
983		1047
984		1048
985		1049
986		1050
987		1051
988		1052
989		1053
990		1054
991		1055
992		1056
993		1057
994		1058
995		1059
996		1060
997		1061
998		1062
999		1063
1000		1064
1001		1065
1002		1066
1003		1067
1004		1068
1005		1069
1006		1070
1007		1071
1008		1072
1009		1073
1010		1074
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		
1041		
1042		
1043		
OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Laguesi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya	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	Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Ji- acheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondracik, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingx- uan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Iz- mailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. 2024b. Openai o1 system card. Preprint, arXiv:2412.16720.	1167
1108	OpenAI. 2024. Multilingual massive multitask language understanding (mmmlu).	1168
1109	OpenAI. 2025a. Introducing gpt-4.1 in the api.	1169
1110	OpenAI. 2025b. Introducing openai o3 and o4-mini.	1170
1111	OpenAI. 2025c. Openai o3-mini.	1171
1112	Chanjun Park, Hyeyoung Kim, Dahyun Kim, SeongHwan Cho, Sanghoon Kim, Sukyung Lee, Yungi Kim, and Hwalsuk Lee. 2024. Open Ko-LLM leaderboard: Evaluating large language models in Korean with Ko-h5 benchmark. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3220–3234,	1172
1113	Bangkok, Thailand. Association for Computational Linguistics.	1173
1114	Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Martin Jaggi, Leandro von Werra, and Thomas Wolf. 2024. Fineweb2: A sparkling update with 1000s of languages.	1174
1115	Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, Michael Choi, Anish Agrawal, Arnav Chopra, Adam Khoja, Ryan Kim, Richard Ren, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, Tung Nguyen, Daron Anderson, Imad Ali Shah, Mikhail Doroshenko, Alun Cennyth Stokes, Mobeen Mahmood, Jaeho Lee, Oleksandr Pokutnyi, Oleg Iskra, Jessica P. Wang, Robert Gerbicz, John-Clark Levin, Serguei Popov, Fiona Feng, Steven Y. Feng, Haoran Zhao, Michael Yu, Varun Gangal, Chelsea Zou, Zihan Wang, Mstyslav Kazakov, Geoff Galgon, Johannes Schmitt, Alvaro Sanchez, Yongki Lee, Will Yeadon, Scott Sauers, Marc Roth, Chidozie Agu, Søren Riis, Fabian Giska, Saiteja Utpala, Antrell Cheatom, Zachary Giboney, Gashaw M. Goshu, Sarah-Jane Crowson, Mohinder Maheshbhai Naiya, Noah Burns, Lennart Finke, Zerui Cheng, Hyunwoo Park, Francesco Fournier-Facio, Jennifer Zampese, John B. Wydallis, Ryan G. Hoerr, Mark Nandor, Tim Gehrunger, Jiaqi Cai, Ben McCarty, Jungbae Nam, Edwin Taylor, Jun Jin, Gautier Abou Loume, Hangrui Cao, Alexis C Garretson, Damien Sileo, Qiuyu Ren, Doru Cojoc, Pavel Arkhipov, Usman Qazi, Aras Bacho, Lianghui Li, Sumeet Motwani, Christian Schroeder de Witt, Alexei Kopylov, Johannes Veith, Eric Singer, Paolo Rissone, Jaehyeok Jin, Jack Wei Lun Shi, Chris G. Willcocks, Ameya Prabhu, Longke Tang, Kevin Zhou, Emily de Oliveira Santos, Andrey Pupasov Maksimov, Edward Vendrow, Kengo Zenitani, Joshua Robinson, Aleksandar Mikov, Julien Guillod, Yuqi Li, Ben Pageler, Joshua Vendrow, Vladyslav Kuchkin, Pierre Marion, Denis Efremov, Jayson Lynch, Kaiqu Liang, Andrew Gritsevskiy, Dakotah Martinez, Nick Crispino, Dimitri Zvonkine, Natanael Wildner Fraga, Saeed Soori, Ori Press, Henry Tang, Julian Salazar, Sean R. Green, Lina Brüssel, Moon Twayana, Aymeric Dieuleveut, T. Ryan Rogers, Wenjin Zhang, Ross Finocchio, Bikun Li, Jinzhou Yang, Arun Rao, Gabriel Loiseau, Mikhail Kalinin, Marco Lukas, Ciprian Manolescu, Nate Stambaugh, Subrata Mishra, Ariel Ghislain Kemogne Kamdoum, Tad Hogg, Alvin Jin, Carlo Bosio, Gongbo Sun, Brian P Coppola, Haline Heidinger, Rafael Sayous, Stefan Ivanov, Joseph M Cavanagh, Jiawei Shen, Joseph Marvin Imperial, Philippe Schwaller, Shaipranesh Senthilkuma, Andres M Bran, Andres Algaba, Brecht Verbeken, Kelsey Van den Houte, Lynn Van Der Sypt, David Noever, Lisa Schut, Ilia Sucholutsky, Evgenii Zheltonozhskii, Qiaochu Yuan, Derek Lim, Richard Stanley, Shankar Sivarajan, Tong Yang, John Maar, Julian Wykowski, Martí Oller, Jennifer Sandlin, Anmol Sahu, Cesare Giulio Arditò, Yuzheng Hu, Felipe Meneguitti Dias, Tobias Kreiman, Kaivalya Rawal, Tobias Garcia Vilchis,	1175
1116		1176
1117		1177
1118		1178
1119		1179
1120		1180
1121		1181
1122		1182
1123		1183
1124		1184
1125		1185
1126		1186
1127		1187
1128		1188
1129		1189
1130		1190
1131		1191
1132		1192
1133		1193
1134		1194
1135		1195
1136		1196
1137		1197
1138		1198
1139		1199
1140		1200
1141		1201
1142		1202
1143		1203
1144		1204
1145		1205
1146		1206
1147		1207
1148		1208
1149		1209
1150		1210
1151		1211
1152		1212
1153		1213
1154		1214
1155		1215
1156		1216
1157		1217
1158		1218
1159		1219
1160		1220
1161		1221
1162		1222
1163		1223
1164		1224
1165		1225
1166		1226
1167		1227
1168		1228

1229	Yuexuan Zu, Martin Lackner, James Koppel, Jeremy	Tordera, George Balabanian, Earth Anderson, Lynna	1293
1230	Nguyen, Daniil S. Antonenko, Steffi Chern, Bingchen	Kvistad, Alejandro José Moyano, Hsiaoyun Mill-	1294
1231	Zhao, Pierrot Arsene, Sergey Ivanov, Rafal Poświaty,	iron, Ahmad Sakor, Murat Eron, Isaac C. McAl-	1295
1232	Chenguang Wang, Daofeng Li, Donato Crisostomi,	ister, Andrew Favre D. O., Shailesh Shah, Xiaox-	1296
1233	Ali Dehghan, Andrea Achilleos, John Arnold Am-	iang Zhou, Firuz Kamalov, Ronald Clark, Sher-	1297
1234	bay, Benjamin Myklebust, Archan Sen, David Per-	win Abdoli, Tim Santens, Harrison K Wang, Evan	1298
1235	rella, Nurdin Kaparov, Mark H Inlow, Allen Zang,	Chen, Alessandro Tomasiello, G. Bruno De Luca,	1299
1236	Kalyan Ramakrishnan, Daniil Orel, Vladislav Porit-	Shi-Zhuo Looi, Vinh-Kha Le, Noam Kolt, Niels	1300
1237	ski, Shalev Ben-David, Zachary Berger, Parker	Mündler, Avi Semler, Emma Rodman, Jacob Drori,	1301
1238	Whitfill, Michael Foster, Daniel Munro, Linh Ho,	Carl J Fossum, Luk Gloor, Milind Jagota, Ronak	1302
1239	Dan Bar Hava, Aleksey Kuchkin, Robert Lauff,	Pradeep, Honglu Fan, Tej Shah, Jonathan Eicher,	1303
1240	David Holmes, Frank Sommerhage, Anji Zhang,	Michael Chen, Kushal Thaman, William Merrill,	1304
1241	Richard Moat, Keith Schneider, Daniel Pyda, Zakayo	Moritz Firsching, Carter Harris, Stefan Ciobăcă,	1305
1242	Kazibwe, Mukhwinder Singh, Don Clarke, Dae Hyun	Jason Gross, Rohan Pandey, Ilya Gusev, Adam Jones,	1306
1243	Kim, Sara Fish, Veit Elser, Victor Efren Guadar-	Shashank Agnihotri, Pavel Zhelnov, Siranut Us-	1307
1244	rama Vilchis, Immo Klose, Christoph Demian,	awasutsakorn, Mohammadreza Mofayez, Alexander	1308
1245	Ujjwala Amantheswaran, Adam Zweiger, Guglielmo	Piperski, Marc Carauleanu, David K. Zhang,	1309
1246	Albani, Jeffery Li, Nicolas Daans, Maksim Radionov,	Kostiantyn Dobarskyi, Dylan Ler, Roman Leven-	1310
1247	Václav Rozhoň, Vincent Ginis, Ziqiao Ma, Chris-	tov, Ignat Soroko, Thorben Jansen, Scott Creighton,	1311
1248	tian Stump, Jacob Platnick, Volodymyr Nevirkovets,	Pascal Lauer, Joshua Duersch, Vage Taamazyan,	1312
1249	Luke Basler, Marco Piccardo, Niv Cohen, Viren-	Dario Beazzi, Wiktor Morak, Wenjie Ma, William	1313
1250	dra Singh, Josef Tkadlec, Paul Rosu, Alan Goldfarb,	Held, Tran Duc Huy, Ruicheng Xian, Armel Randy	1314
1251	Piotr Padlewski, Stanislaw Barzowski, Kyle Mont-	Zebaze, Mohanad Mohamed, Julian Noah Leser,	1315
1252	gomery, Aline Menezes, Arkil Patel, Zixuan Wang,	Michelle X Yuan, Laila Yacar, Johannes Lengler,	1316
1253	Jamie Tucker-Foltz, Jack Stade, Declan Grabb, Tom	Katarzyna Olszewska, Hossein Shahrtash, Edson	1317
1254	Goertzen, Fereshteh Kazemi, Jeremiah Milbauer, Ab-	Oliveira, Joseph W. Jackson, Daniel Espinosa Gon-	1318
1255	hishek Shukla, Hossam Elgnainy, Yan Carlos Leyva	zalez, Andy Zou, Muthu Chidambaram, Timothy	1319
1256	Labrador, Hao He, Ling Zhang, Alan Givré, Hew	Manik, Hector Haffenden, Dashiell Stander, Ali	1320
1257	Wolff, Gözdenur Demir, Muhammad Fayeaz Aziz,	Dasouqi, Alexander Shen, Emilien Duc, Bita Gol-	1321
1258	Younesse Kaddar, Ivar Ängquist, Yanxu Chen, El-	shani, David Stap, Mikalai Uzhou, Alina Borisovna	1322
1259	liott Thornley, Robin Zhang, Jiayi Pan, Antonio Ter-	Zhidkovskaya, Lukas Lewark, Miguel Orbegozo Ro-	1323
1260	pin, Niklas Muennighoff, Hailey Schoelkopf, Eric	driguez, Mátyás Vincze, Dustin Wehr, Colin Tang,	1324
1261	Zheng, Avishy Carmi, Jainam Shah, Ethan D. L.	Shaun Phillips, Fortuna Samuele, Jiang Muzhen,	1325
1262	Brown, Kelin Zhu, Max Bartolo, Richard Wheeler,	Fredrik Ekström, Angela Hammon, Oam Patel, Faraz	1326
1263	Andrew Ho, Shaul Barkan, Jiaqi Wang, Martin Ste-	Farhidi, George Medley, Forough Mohammadzadeh,	1327
1264	hberger, Egor Kretov, Peter Bradshaw, JP Heimo-	Madellene Peñaflor, Haile Kassahun, Alena Friedrich,	1328
1265	n, Kaustubh Sridhar, Zaki Hossain, Ido Akov, Yury	Claire Sparrow, Rayner Hernandez Perez, Taom	1329
1266	Makarychev, Joanna Tam, Hieu Hoang, David M.	Sakal, Omkar Dhamane, Ali Khajegili Mirabadi,	1330
1267	Cunningham, Vladimir Goryachev, Demosthenes Pa-	Eric Hallman, Kenchi Okutsu, Mike Battaglia, Mo-	1331
1268	tramanis, Michael Krause, Andrew Redenti, David	hammad Maghsoudimehrabani, Alon Amit, Dave	1332
1269	Aldous, Jesyin Lai, Shannon Coleman, Jiangnan Xu,	Hulbert, Roberto Pereira, Simon Weber, Handoko,	1333
1270	Sangwon Lee, Ilias Magoulas, Sandy Zhao, Ning	Anton Peristy, Stephen Malina, Samuel Albanie,	1334
1271	Tang, Michael K. Cohen, Micah Carroll, Orr Par-	Will Cai, Mustafa Mehkary, Rami Aly, Frank Rei-	1335
1272	adise, Jan Hendrik Kirchner, Stefan Steinerberger,	degeld, Anna-Katharina Dick, Cary Friday, Jasdeep	1336
1273	Maksym Ovchynnikov, Jason O. Matos, Adithya	Sidhu, Hassan Shapourian, Wanyoung Kim, Mariana	1337
1274	Shenoy, Michael Wang, Yuzhou Nie, Paolo Gior-	Costa, Hubeyb Gurdogan, Brian Weber, Harsh Ku-	1338
1275	dano, Philipp Petersen, Anna Sztyber-Betley, Paolo	mar, Tong Jiang, Arunim Agarwal, Chiara Ceconello,	1339
1276	Faraboschi, Robin Riblet, Jonathan Crozier, Shiv Ha-	Warren S. Vaz, Chao Zhuang, Haon Park, Andrew R.	1340
1277	lasyamani, Antonella Pinto, Shreyas Verma, Prashant	Tawfeek, Daattavya Aggarwal, Michael Kirchhof,	1341
1278	Joshi, Eli Meril, Zheng-Xin Yong, Allison Tee,	Linjie Dai, Evan Kim, Johan Ferret, Yuzhou Wang,	1342
1279	Jérémie Andréoletti, Orion Weller, Raghav Singhal,	Minghao Yan, Krzysztof Burdzy, Lixin Zhang, Anto-	1343
1280	Gang Zhang, Alexander Ivanov, Seri Khouri, Nils	nio Franca, Diana T. Pham, Kang Yong Loh,	1344
1281	Gustafsson, Hamid Mostaghimi, Kunvar Thaman,	Joshua Robinson, Abram Jackson, Shreen Gul, Gun-	1345
1282	Qijia Chen, Tran Quoc Khánh, Jacob Loader, Ste-	jan Chhablani, Zhehang Du, Adrian Cosma, Jesus	1346
1283	fano Cavalleri, Hannah Szlyk, Zachary Brown, Hi-	Colino, Colin White, Jacob Votava, Vladimir Vin-	1347
1284	manshu Narayan, Jonathan Roberts, William Alley,	nnykov, Ethan Delaney, Petr Spelda, Vit Stritecky,	1348
1285	Kunyang Sun, Ryan Stendall, Max Lamparth, Anka	Syed M. Shahid, Jean-Christophe Mourrat, Lavr	1349
1286	Reuel, Ting Wang, Hanmeng Xu, Pablo Hernández-	Vetoshkin, Koen Sponselee, Renas Bacho, Floren-	1350
1287	Cámar, Freddie Martin, Thomas Preu, Tomek Kor-	cia de la Rosa, Xiuyu Li, Guillaume Malod, Leon	1351
1288	bak, Marcus Abramovitch, Dominic Williamson, Ida	Lang, Julien Laurendeau, Dmitry Kazakov, Fatimah	1352
1289	Bosio, Ziye Chen, Biró Bálint, Eve J. Y. Lo, Maria	Adesanya, Julien Portier, Lawrence Hollom, Victor	1353
1290	Inês S. Nunes, Yibo Jiang, M Saiful Bari, Peyman	Souza, Yuchen Anna Zhou, Julien Degorre, Yiğit	1354
1291	Kassani, Zihao Wang, Behzad Ansarinejad, Yewen	Yalın, Gbenga Daniel Obikoya, Luca Arnaboldi, Rai,	1355
1292	Sun, Stephane Durand, Guillaume Douville, Daniel	Filippo Bigi, M. C. Boscá, Oleg Shumar, Kani-	1356

1357	uar Bacho, Pierre Clavier, Gabriel Recchia, Mara Popescu, Nikita Shulga, Ngefor Mildred Tanwie, Dennis Peskoff, Thomas C. H. Lux, Ben Rank, Colin Ni, Matthew Brooks, Alesia Yakimchyk, Huanxu, Liu, Olle Häggström, Emil Verkama, Hans Gundlach, Leonor Brito-Santana, Brian Amaro, Vivek Vajipey, Rynaa Grover, Yiyang Fan, Gabriel Poesia Reis e Silva, Linwei Xin, Yosi Kratish, Jakub Łucki, Wending Li, Sivakanth Gopi, Andrea Caciolai, Justin Xu, Kevin Joseph Scaria, Freddie Vargus, Farzad Habibi, Long, Lian, Emanuele Rodolà, Jules Robins, Vincent Cheng, Tony Fruhauff, Brad Raynor, Hao Qi, Xi Jiang, Ben Segev, Jingxuan Fan, Sarah Martinson, Erik Y. Wang, Kaylie Hausknecht, Michael P. Brenner, Mao Mao, Xinyu Zhang, David Avagian, Es-hawn Jessica Scipio, Alon Ragoler, Justin Tan, Blake Sims, Rebeka Plecnik, Aaron Kirtland, Omer Faruk Bodur, D. P. Shinde, Zahra Adoul, Mohamed Zekry, Ali Karakoc, Tania C. B. Santos, Samir Sham-seldeen, Loukmene Karim, Anna Liakhovitskaia, Nate Resman, Nicholas Farina, Juan Carlos Gonzalez, Gabe Maayan, Sarah Hoback, Rodrigo De Oliveira Pena, Glen Sherman, Elizabeth Kelley, Hodjat Mar-iji, Rasoul Pouriamanesh, Wentao Wu, Sandra Mendoza, Ismail Alarab, Joshua Cole, Danyelle Ferreira, Bryan Johnson, Mohammad Safdari, Liangti Dai, Siriphap Arthornthurasuk, Alexey Pronin, Jing Fan, Angel Ramirez-Trinidad, Ashley Cartwright, Daphny Pottmaier, Omid Taheri, David Outevsky, Stanley Stepanic, Samuel Perry, Luke Askew, Raúl Adrián Huerta Rodríguez, Ali M. R. Minissi, Sam Ali, Ricardo Lorena, Krishnamurthy Iyer, Arshad Anil Fasiludeen, Sk Md Salauddin, Murat Islam, Juan Gonzalez, Josh Ducey, Maja Somrak, Vasilios Mavroudis, Eric Vergo, Juehang Qin, Benjamín Borbás, Eric Chu, Jack Lindsey, Anil Radhakrishnan, Antoine Jallon, I. M. J. McInnis, Pawan Kumar, Laxman Prasad Goswami, Daniel Bugas, Nasser Heydari, Ferenc Jeanplong, Archimedes Apronti, Abdallah Galal, Ng Ze-An, Ankit Singh, Joan of Arc Xavier, Kanu Priya Agarwal, Mohammed Berkani, Benedito Alves de Oliveira Junior, Dmitry Malishev, Nicolas Remy, Taylor D. Hartman, Tim Tarver, Stephen Mensah, Javier Gimenez, Roselynn Grace Montecillo, Russell Campbell, Asankhaya Sharma, Khalida Meer, Xavier Alapont, Deepakkumar Patil, Rajat Madeshwari, Abdelkader Dendane, Priti Shukla, Sergei Bogdanov, Sören Möller, Muhammad Rehan Siddiqi, Prajvi Saxena, Himanshu Gupta, Innocent Enyekwe, Ragavendran P V, Zienab EL-Wasif, Aleksandr Mak-sapetyan, Vivien Rossbach, Chris Harjadi, Mohsen Bahaloohoreh, Song Bian, John Lai, Justine Leon Uro, Greg Bateman, Mohamed Sayed, Ahmed Men-shawy, Darling Duclosel, Yashaswini Jain, Ashley Aaron, Murat Tirayioglu, Sheeshram Siddh, Keith Krenek, Alex Hoover, Joseph McGowan, Tejal Pat-wardhan, Summer Yue, Alexandr Wang, and Dan Hendrycks. 2025. Humanity’s last exam. Preprint, arXiv:2501.14249.	marks: is mmlu lost in translation? Preprint, arXiv:2406.17789.	1419
1358			1420
1359			1421
1360			1422
1361			1423
1362			1424
1363			1425
1364			1426
1365			1427
1366			1428
1367			1429
1368			1429
1369			1430
1370			1431
1371			1432
1372			
1373			
1374			
1375			
1376			
1377			
1378			
1379			
1380			
1381			
1382			
1383			
1384			
1385			
1386			
1387			
1388			
1389			
1390			
1391			
1392			
1393			
1394			
1395			
1396			
1397			
1398			
1399			
1400			
1401			
1402			
1403			
1404			
1405			
1406			
1407			
1408			
1409			
1410			
1411			
1412			
1413			
1414			
1415			
1416	Irene Plaza, Nina Melero, Cristina del Pozo, Javier Conde, Pedro Reviriego, Marina Mayor-Rocher, and María Grandury. 2024. Spanish and llm bench-		
1417			
1418			
	marks: is mmlu lost in translation? Preprint, arXiv:2406.17789.		
	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Bin yuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 technical report. Preprint, arXiv:2412.15115.		
	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. GPQA: A graduate-level google-proof q&a benchmark. In <i>First Conference on Language Modeling</i> .		
	LG AI Research, Soyoung An, Kyunghoon Bae, Eunbi Choi, Kibong Choi, Stanley Jungkyu Choi, Seokhee Hong, Junwon Hwang, Hyojin Jeon, Gerrard Jeongwon Jo, Hyunjik Jo, Jiyeon Jung, Yountae Jung, Hyosang Kim, Joonkee Kim, Seonghwan Kim, Soyeon Kim, Sunkyoung Kim, Yireun Kim, Yongil Kim, Youchul Kim, Edward Hwayoung Lee, Haeju Lee, Honglak Lee, Jinsik Lee, Kyungmin Lee, Woohyung Lim, Sangha Park, Sooyoun Park, Yongmin Park, Sihoon Yang, Heuiyean Yeen, and Hyeongu Yun. 2024. Exaone 3.5: Series of large language models for real-world use cases. Preprint, arXiv:2412.04862.		
	LG AI Research, Kyunghoon Bae, Eunbi Choi, Ki-bong Choi, Stanley Jungkyu Choi, Yemuk Choi, Seokhee Hong, Junwon Hwang, Hyojin Jeon, Ki-jeong Jeon, Gerrard Jeongwon Jo, Hyunjik Jo, Jiyeon Jung, Hyosang Kim, Joonkee Kim, Seonghwan Kim, Soyeon Kim, Sunkyoung Kim, Yireun Kim, Yongil Kim, Youchul Kim, Edward Hwayoung Lee, Haeju Lee, Honglak Lee, Jinsik Lee, Kyungmin Lee, Sangha Park, Yongmin Park, Sihoon Yang, Heuiyean Yeen, Sihyuk Yi, and Hyeongu Yun. 2025. Exaone deep: Reasoning enhanced language models. Preprint, arXiv:2503.12524.		
	Manley Roberts, Himanshu Thakur, Christine Herlihy, Colin White, and Samuel Dooley. 2024. To the cut-off... and beyond? a longitudinal perspective on LLM data contamination. In <i>The Twelfth International Conference on Learning Representations</i> .		
	Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I. Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond Ng, Shayne Longpre, Wei-Yin Ko, Sebastian Ruder, Madeline Smith, Antoine Bosselut, Alice Oh, Andre F. T. Martins, Leshem Choshen, Daphne Ippolito, Enzo Ferrante, Marzieh Fadaee, Beyza Ermis, and Sara Hooker. 2025. Global mmlu: Understanding and addressing cultural and linguistic biases in multi-lingual evaluation. Preprint, arXiv:2412.03304.		

1479	Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I. Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond Ng, Shayne Longpre, Wei-Yin Ko, Madeline Smith, Antoine Bosselut, Alice Oh, Andre F. T. Martins, Leshem Choshen, Daphne Ippolito, Enzo Ferrante, Marzieh Fadaee, Beyza Ermis, and Sara Hooker.	2024. <i>Global mmlu: Understanding and addressing cultural and linguistic biases in multilingual evaluation</i> . Preprint, arXiv:2412.03304.	1540
1480			1541
1481			1542
1482			1543
1483			1544
1484			1545
1485			1546
1486			1547
1487			1548
1488			1549
1489			1550
1490	Guojin Son, Hanwool Lee, Sungdong Kim, Seungone Kim, Niklas Muennighoff, Taekyoon Choi, Cheonbok Park, Kang Min Yoo, and Stella Biderman.	2024a. <i>Kmmlu: Measuring massive multi-task language understanding in korean</i> . Preprint, arXiv:2402.11548.	1551
1491			1552
1492			1553
1493			1554
1494			1555
1495			1556
1496	Guojin Son, Hanwool Lee, Suwan Kim, Huiseo Kim, Jae cheol Lee, Je Won Yeom, Jihyu Jung, Jung woo Kim, and Songseong Kim.	2024b. <i>HAE-RAE bench: Evaluation of Korean knowledge in language models</i> . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 7993–8007, Torino, Italia. ELRA and ICCL.	1557
1497			1558
1498			1559
1499			1560
1500			1561
1501			1562
1502			1563
1503			1564
1504	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Johan Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew Kyle Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mollokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekeci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, Cesar Ferri, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Christopher Waites, Christian Voigt, Christopher D Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, C. Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin,	Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodolà, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, François Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Xinyue Wang, Gonzalo Jaimovitch-Lopez, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Francis Anthony Shevlin, Hinrich Schuetze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jae-hoon Lee, Jaime Fernández Fisac, James B Simon, James Koppel, James Zheng, James Zou, Jan Kocon, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh Dhole, Kevin Gimpel, Kevin Omundi, Kory Wallace Mathewson, Kristen Chiaffullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros-Colón, Luke Metz, Lütfi Kerem Senel, Maarten Bosma, Maarten Sap, Maartje Ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramirez-Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael Andrew Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Świderski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdheh Gheini, Mukund Varma T, Nanyun Peng, Nathan Andrew Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Anto-	1559
1499			1560
1500			1561
1501			1562
1502			1563
1503			1564
1504			1565
1505			1566
1506			1567
1507			1568
1508			1569
1509			1570
1510			1571
1511			1572
1512			1573
1513			1574
1514			1575
1515			1576
1516			1577
1517			1578
1518			1579
1519			1580
1520			1581
1521			1582
1522			1583
1523			1584
1524			1585
1525			1586
1526			1587
1527			1588
1528			1589
1529			1590
1530			1591
1531			1592
1532			1593
1533			1594
1534			1595
1535			1596
1536			1597
1537			1598
1538			1599
1539			1600

1604	nio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter W Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan Le Bras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Russ Salakhutdinov, Ryan Andrew Chi, Seungjae Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel Stern Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima Shammie Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven Piantadosi, Stuart Shieber, Summer Mishnerghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsunori Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Venkatesh Ramasesh, vinay uday prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Sophie Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>Transactions on Machine Learning Research</i> . Featured Certification.	1665
1656	Google Team. 2025a. Gemma 3 technical report.	1666
1657	Kanana LLM Team, Yunju Bak, Hojin Lee, Minho Ryu, Jiyeon Ham, Seungjae Jung, Daniel Wontae Nam, Taegyeong Eo, Donghun Lee, Doohae Jung, Boseop Kim, Nayeon Kim, Jaesun Park, Hyunho Kim, Hyunwoong Ko, Changmin Lee, Kyoung-Woon On, Seuye Baeg, Junrae Cho, Sunghee Jung, Jieun Kang, EungGyun Kim, Eunhwa Kim, Byeongil Ko, Daniel Lee, Minchul Lee, Miok Lee, Shinbok Lee, and Gaeun Seo. 2025a. Kanana: Compute-efficient bilingual language models. <i>Preprint</i> , arXiv:2502.18934.	1667
1658	M-A-P Team, Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, Kang Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, Chujie Zheng, Kaixing Deng, Shuyue Guo, Shian Jia, Sichao Jiang, Yiyan Liao, Rui Li, Qinrui Li, Sirun Li, Yizhi Li, Yunwen Li, Dehua Ma, Yuansheng Ni, Haoran Que, Qiyo Wang, Zhoufutu Wen, Siwei Wu, Tianshun Xing, Ming Xu, Zhenzhu Yang, Zekun Moore Wang, Junting Zhou, Yuelin Bai, Xingyuan Bu, Chenglin Cai, Liang Chen, Yifan Chen, Chengtuo Cheng, Tianhao Cheng, Keyi Ding, Siming Huang, Yun Huang, Yaoru Li, Yizhe Li, Zhaoqun Li, Tianhao Liang, Chengdong Lin, Hongquan Lin, Yinghao Ma, Zhongyuan Peng, Zifan Peng, Qige Qi, Shi Qiu, Xingwei Qu, Yizhou Tan, Zili Wang, Chenqing Wang, Hao Wang, Yiya Wang, Yubo Wang, Jiajun Xu, Kexin Yang, Ruibin Yuan, Yuanhao Yue, Tianyang Zhan, Chun Zhang, Jingyang Zhang, Xiyue Zhang, Xingjian Zhang, Yue Zhang, Yongchi Zhao, Xiangyu Zheng, Chenghua Zhong, Yang Gao, Zhoujun Li, Dayiheng Liu, Qian Liu, Tianyu Liu, Shiwen Ni, Junran Peng, Yujia Qin, Wenbo Su, Guoyin Wang, Shi Wang, Jian Yang, Min Yang, Meng Cao, Xiang Yue, Zhaoxiang Zhang, Wangchunshu Zhou, Jiaheng Liu, Qunshu Lin, Wenhao Huang, and Ge Zhang. 2025b. Supergpqa: Scaling llm evaluation across 285 graduate disciplines. <i>Preprint</i> , arXiv:2502.14739.	1668
1659	Qwen Team. 2025b. Qwq-32b: Embracing the power of reinforcement learning.	1669
1660	Joshua Vendrow, Edward Vendrow, Sara Beery, and Aleksander Madry. 2025. Do large language model benchmarks test reliability? <i>Preprint</i> , arXiv:2502.03461.	1670
1661	Yiming Wang, Pei Zhang, Jialong Tang, Haoran Wei, Baosong Yang, Rui Wang, Chenshu Sun, Feitong Sun, Jiran Zhang, Junxuan Wu, Qiqian Cang, Yichang Zhang, Fei Huang, Junyang Lin, Fei Huang, and Jingen Zhou. 2025. Polymath: Evaluating mathematical reasoning in multilingual contexts. <i>Preprint</i> , arXiv:2504.18428.	1671
1662	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024. MMLU-pro: A more robust and challenging multi-task language understanding benchmark. In <i>The Thirty-eighth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	1672
1663	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems</i> .	1673
1664	Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Benjamin Feuer, Siddhartha Jain, Ravid Shwartz-Ziv,	1674

1723	Neel Jain, Khalid Saifullah, Sreemanti Dey, Shubh-Agrawal, Sandeep Singh Sandha, Siddartha Venkat Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah Goldblum. 2025. Livebench: A challenging, contamination-limited LLM benchmark. In <i>The Thirteenth International Conference on Learning Representations</i> .	1784
1724		1785
1725		1786
1726		1787
1727		1788
1728		1789
1729		1790
1730	xAI. 2025. Grok 3 beta — the age of reasoning agents .	1791
1731		1792
1732		1793
1733	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. 2025. Qwen3 technical report . Preprint, arXiv:2505.09388.	1794
1734		1795
1735		1796
1736		1797
1737		1798
1738		1799
1739		1800
1740		1801
1741		1802
1742		1803
1743		1804
1744		1805
1745		1806
1746		1807
1747		1808
1748	Kang Min Yoo, Jaegeun Han, Sookyo In, Heewon Jeon, Jisu Jeong, Jaewook Kang, Hyunwook Kim, Kyung-Min Kim, Munhyong Kim, Sungju Kim, Donghyun Kwak, Hanock Kwak, Se Jung Kwon, Bado Lee, Dongsoo Lee, Gichang Lee, Jooho Lee, Baeseong Park, Seongjin Shin, Joon-sang Yu, Seolki Baek, Sumin Byeon, Eungsup Cho, Dooseok Choe, Jeesung Han, Youngkyun Jin, Hyein Jun, Jaeseung Jung, Chanwoong Kim, Jinhong Kim, Jinuk Kim, Dokyeong Lee, Dongwook Park, Jeong Min Sohn, Sujung Han, Jiae Heo, Sungju Hong, Mina Jeon, Hyunhoon Jung, Jungeun Jung, Wangkyo Jung, Chungjoon Kim, Hyeri Kim, Jonghyun Kim, Min Young Kim, Soeun Lee, Joonhee Park, Jieun Shin, Sojin Yang, Jungsoon Yoon, Hwaran Lee, Sanghwan Bae, Jeehwan Cha, Karl Gylleus, Donghoon Ham, Mihak Hong, Youngki Hong, Yunki Hong, Dahyun Jang, Hyojun Jeon, Yujin Jeon, Yeji Jeong, Myungeun Ji, Yeguk Jin, Chansong Jo, Shinyoung Joo, Seunghwan Jung, Adrian Jungmyung Kim, Byoung Hoon Kim, Hyomin Kim, Jungwhan Kim, Minkyung Kim, Minseung Kim, Sungdong Kim, Yonghee Kim, Youngjun Kim, Youngkwan Kim, Donghyeon Ko, Dughyun Lee, Ha Young Lee, Jaehong Lee, Jieun Lee, Jonghyun Lee, Jongjin Lee, Min Young Lee, Yehbin Lee, Taehong Min, Yuri Min, Kiyo Moon, Hyangnam Oh, Jaesun Park, Kyuyon Park, Younghun Park, Hanbae Seo, Seunghyun Seo, Mihyun Sim, Gyubin Son, Matt Yeo, Kyung Hoon Yeom, Wonjoon Yoo, Myungin You, Doheon Ahn, Homin Ahn, Joohee Ahn, Seongmin Ahn, Chanwoo An, Hyeryun An, Junho An, Sang-Min An, Boram Byun, Eunbin Byun, Jongho Cha, Minji Chang, Seung-gyu Chang, Haesong Cho, Youngdo Cho, Dalnim Choi, Daseul Choi, Hyoseok Choi, Minseong Choi,	1809
1749		1810
1750		1811
1751		1812
1752		1813
1753		1814
1754		1815
1755		1816
1756		1817
1757		1818
1758		1819
1759		1820
1760		1821
1761		1822
1762		1823
1763		1824
1764		1825
1765		1826
1766		1827
1767		1828
1768		1829
1769		1830
1770		1831
1771		1832
1772		1833
1773		1834
1774		1835
1775		1836
1776		1837
1777		1838
1778		1839
1779		1840
1780		1841
1781		1842
1782		1843
1783		1844
		1845
		1846
		1847

1848	Lee, Seung Jin Lee, Suhyeon Lee, Yeonjae Lee, Yesol Lee, Youngbeom Lee, Yujin Lee, Shaodong Li, Tianyu Liu, Seong-Eun Moon, Taehong Moon, Max-Lasse Nihlenramstroem, Wonseok Oh, Yuri Oh, Hongbeen Park, Hyekyung Park, Jaeho Park, Nohil Park, Sangjin Park, Jiwon Ryu, Miru Ryu, Simo Ryu, Ahreum Seo, Hee Seo, Kangdeok Seo, Jamin Shin, Seungyoun Shin, Heetae Sin, Jiangping Wang, Lei Wang, Ning Xiang, Longxiang Xiao, Jing Xu, Seonyeong Yi, Haanju Yoo, Haneul Yoo, Hwanhee Yoo, Liang Yu, Youngjae Yu, Weijie Yuan, Bo Zeng, Qian Zhou, Kyunghyun Cho, Jung-Woo Ha, Joonsuk Park, Jihyun Hwang, Hyoung Jo Kwon, Soonyong Kwon, Jungyeon Lee, Seungho Lee, Seonghyeon Lim, Hyunkyoung Noh, Seungho Choi, Sang-Woo Lee, Jung Hwa Lim, and Nako Sung. 2024. <i>Hyperclova x technical report</i> . Preprint, arXiv:2404.01954.	1889
1849		1890
1850		1891
1851		1892
1852		1893
1853		1894
1854		1895
1855		1896
1856		1897
1857		1898
1858		1899
1859		1900
1860		
1861		
1862		
1863		
1864		
1865	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. <i>HellaSwag: Can a machine really finish your sentence?</i> In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4791–4800, Florence, Italy. Association for Computational Linguistics.	1901
1866		1902
1867		
1868		
1869		
1870		
1871	Yidan Zhang, Yu Wan, Boyi Deng, Baosong Yang, Haoran Wei, Fei Huang, Bowen Yu, Junyang Lin, Fei Huang, and Jingren Zhou. 2025. <i>P-mmeval: A parallel multilingual multitask benchmark for consistent evaluation of llms</i> . Preprint, arXiv:2411.09116.	1903
1872		1904
1873		
1874		
1875		
1876	Qihao Zhao, Yangyu Huang, Tengchao Lv, Lei Cui, Qinzheng Sun, Shaoguang Mao, Xin Zhang, Ying Xin, Qiufeng Yin, Scarlett Li, and Furu Wei. 2024. <i>Mmlu-cf: A contamination-free multi-task language understanding benchmark</i> . Preprint, arXiv:2412.15194.	1905
1877		1906
1878		
1879		
1880		
1881		
1882	Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. 2024. <i>SGLang: Efficient execution of structured language model programs</i> . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	1907
1883		1908
1884		
1885		
1886		
1887		
1888		

A Details of KMMLU-REDUX Errors

A.1 Error Statistics

We define a set of error types based on recurring issues observed during our analysis of the KMMLU. Then, we find out the number of data instances per each error type as shown in Table 4. To identify leaked answer cases, we apply rule-based filtering using string-overlap heuristics. For other error types, including notation errors, bad clarity, and ill-posed questions, we leverage GPT-4o to assist in annotation. Also, we provide the prompt used for annotation in Figure 5.

- **Ill-posed Question** : The question lacks critical references or contextual information.
- **Leaked Answer** : The ground truth is explicitly stated within the question itself.
- **Notation Error** : Errors in mathematical expressions or chemical equations.
- **Bad Clarity** : The data itself is unclear and contains grammatical errors.

Error Type	# of Questions (Ratio)
Ill-posed Question	512 (1.46 %)
Leaked Answer	42 (0.11%)
Notation Error	846 (2.42%)
Bad Clarity	1284 (3.67%)

Table 4: Statistics of the error types

A.2 Example of Error Types

Table 10 provides examples of the error types in 2.2 identified in KMMLU. Each example illustrates a specific issue that affects benchmark reliability.

B Korean National Technical Qualification List of KMMLU-REDUX

Table 9 presents the list of Korean National Technical Qualifications (KNTQs) included in our benchmark along with their corresponding official exam dates. To facilitate analysis, we categorize the 100 NTQs into Korean Standard Industrial Classification (KSIC) where mapping to the exam (see Figure 6). This categorization enables structured evaluation across diverse domains and better reflects the real-world industrial fields.

GPT-4o Error Annotation Prompt

You are a helpful assistant that annotates error types in questions.

- Ill-posed question stands for which question miss the critical reference information to solve the question (such as table, formular, image, etc.).
- Notation correct stands for which question has correct math and chemical notation. Criteria of correctness whether the notation is impact to solve the question. (e.g. m_2 is incorrect, but \hat{m}_2 is correct. 센티미터 is correct.)
- Grammar correct stands for which question has correct grammar and spelling. Criteria of correctness whether the notation is impact to understand the question. (e.g. 상 담 기 법 보 다는 상 담 자 의 인 간 적 자 질 과 진 솔 한 태 도 를 중 시 한 다 . is grammatically incorrect due to spacing error.)

Please check the following question whether it has error types.

Check ill-posed question True/False.

Check notation correct True/False.

Check grammar correct True/False.

Then, if there is any error, please explain the reason.

Question

{ {question} }

Figure 5: The prompt is used for error type annotation. Each sample is annotated as an error if the respective field returns True. The 'Grammar Correct' field is used to detect 'Bad Clarity' cases.

C Details of KMMLU-PRO Annotation

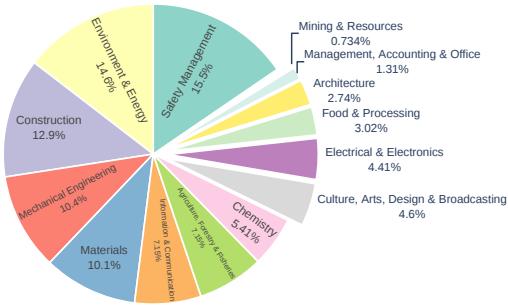


Figure 6: Domain distribution of problems in KMMLU-REDUX. The total size of the dataset is 2,587.

C.1 Annotation Pipeline

We first conduct OCR parsing with GPT-4o on PDF files of KNPL acquisition exams. With the parsed data, the main tasks for human annotators are: 1) reviewing parsing errors, 2) converting tables into latex format, and 3) converting images into text which conveys same meaning. If it is impossible to convert an image into text, we remove it. For the cases where multiple answers are allowed, commonly due to the ambiguity of question itself, we discard them.

As we leverage the official PDF files managed and controlled by the government, we can guarantee the correctness of answer label. This help us to save costs because we do not need experts for annotations nor the answer relabeling to avoid risk of data from online (Gema et al., 2025; Team et al., 2025b). Before annotation, we explained the context of benchmark construction about the Korean professional license exams to human annotators. We present annotation instructions in Fig 7.

1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945

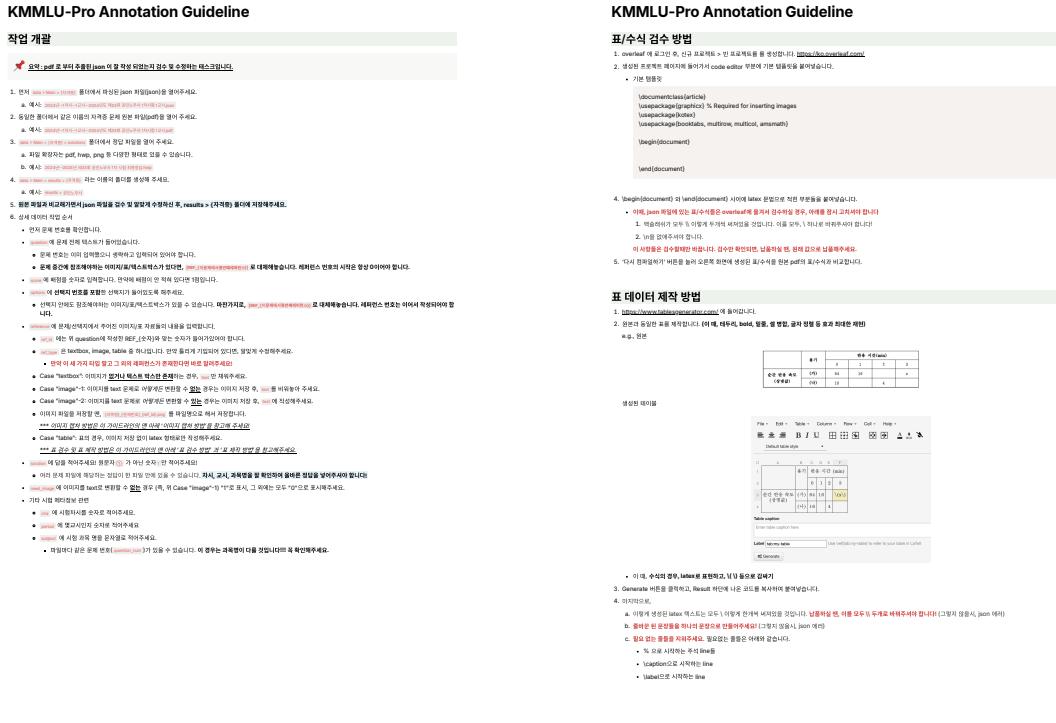


Figure 7: Excerpt from the annotation guidelines for converting PDF documents into structured text. We carefully instruct LaTeX table formatting.

Category	Demographics	Counts
Gender	Female	23
	Male	-
Age	20s	13
	30s	8
	40s	2
Academic background	Bachelor's degrees	20
	Master's degrees	3

Table 5: Demographics of the human annotators for KMMLU-PRO.

C.2 Annotators Demographics

The detailed demographics of annotators are presented in Table 5. The average hourly wage is 8.83 U.S. dollars, which is higher than the legal minimum wage at the time of hiring in South Korea.

D Evaluation Setup

D.1 Evaluation Prompts

The figure 8 and figure 9 present the prompt for the evaluation written in English and Korean, respectively. For the English prompt, we use the regex expression of `r"(?i)Answer[^A-E]*[^A-E]*([A-E])"`. For the Korean prompt

we use `r"정답[^A-E]*:[^A-E]*([A-E])"`. The regex expressions for the flexible parsing are `r"Answer[^A-E]*([A-E])|([A-E])\)"` and `r"정답[^A-E]*([A-E])|([A-E])\)"`, respectively.

D.2 Inference Engines

Excluding closed models, the main inference engine we use is SGLang (Zheng et al., 2024). However, for the Gemma 3 series and Mistral Small 3.1 Instruct models, we adopt vLLM (Kwon et al., 2023) due to their incompatibility with SGLang at the time of the paper writing. Evaluating open-sourced models was conducted over four days using sixteen NVIDIA A100 GPUs. For closed models, we accessed them via API calls, which incurred a total cost of approximately 4,000 USD.

D.3 License Acquisition Criteria for KMMLU-PRO

We follow the official scoring criteria used for each license examination. All licensing exams in KMMLU-PRO are composed of multiple subjects. Candidates are typically required to score at least 40% in every subject and achieve an average score of at least 60%, except for the cases of the Certified

English Prompt

Answer the following multiple choice question. The last line of your response should be of the following format: 'Answer: \$LETTER' (without quotes) where LETTER is one of ABCD. Think step by step before answering.

{ {question} }

- A) {{option A}}
- B) {{option B}}
- C) {{option C}}
- D) {{option D}}

Figure 8: The English prompt used for evaluating LLMs on our KMMLU-REDUX and KMMLU-PRO. This prompt is exactly same with the prompt used for Multiple-Choices Question Answering (MCQA) in OpenAI's simple-evals repository. The number of options is adjusted for each problems.

Korean Prompt

다음 문제에 대해 정답을 고르세요. 당신의 최종 정답은 ABCD 중 하나이고, "정답:" 뒤에 와야 합니다. 정답을 고르기 전에 차근차근 생각하고 추론하세요.

{ {question} }

- A) {{option A}}
- B) {{option B}}
- C) {{option C}}
- D) {{option D}}

Figure 9: The Korean prompt used for evaluating LLMs on our KMMLU-REDUX and KMMLU-PRO. This prompt is translated version of the English prompt. The number of options is adjusted for each problems.

Judicial Scrivener and the Lawyer. For the Certified Judicial Scrivener exam, candidates need only score at least 40% in each subject, with no requirement regarding the average. The Lawyer license exam uses a relative grading system where only a certain proportion of top-scoring candidates pass. The usual cut-off point is approximately 54.22 (900 out of 1660), which we use as the passing threshold.

It is also important to note that our evaluation benchmark is text-based, not multi-modal. Therefore, we exclude questions that include images. In addition, candidates often need to go through multiple exam stages to obtain a license, with the later stages, such as the second or third, containing descriptive questions. However, since we only collect multiple-choice questions, the descriptive problems

are excluded. Lastly, we exclude questions with multiple answers introduced by the ambiguity of the question.

E Detailed Results

E.1 KMMLU vs KMMLU-REDUX

Figure 10 illustrates the performance of LLMs on KMMLU and KMMLU-REDUX. The LLMs' performances on KMMLU-REDUX is lower than on KMMLU, due to our filtration process which aims to retain only challenging problems from KMMLU (see Section 2.2.1). Despite this decrease, there is a near-perfect monotonic association between the results, with a Spearman's rank correlation coefficient (ρ) of 0.995, suggesting they are highly correlated.

1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012

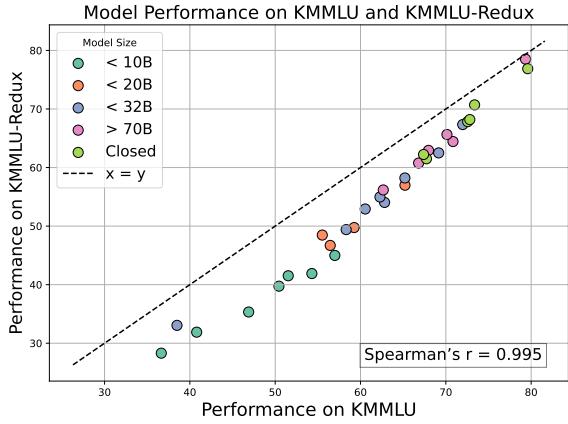


Figure 10: Performance of LLMs on KMMLU and KMMLU-REDUX. A high ρ value indicates a strong correlation between the results of the two benchmarks.

E.2 The Results for Smaller (<10B) Models

The Table 6 presents the results of KMMLU-PRO and KMMLU-REDUX for smaller (<10B) models. Since many of the *tiny* models in this table were used to construct KMMLU-REDUX through adversarial filtration, their KMMLU-REDUX scores are biased. Nevertheless, as shown by the results for larger models, models equipped with dense reasoning capabilities usually outperform their counterparts without reasoning (e.g., Qwen3-8B with and without “thinking”).

E.3 Breakdown of Results of KMMLU-REDUX

The Table 7 presents the breakdown results for 14 categories in KMMLU-REDUX across various LLMs.

E.4 Breakdown of Results of KMMLU-PRO

The Table 8 shows the breakdown results for all NPLs in KMMLU-PRO across various LLMs. While the models relatively easily pass licensing in the medicine domain, they struggle in the Law and Tax&Accounting domains.

	KMMLU-Redux Acc	KMMLU-Pro Acc	KMMLU-Pro # of passed KNPLs	Avg. Acc (micro)
DeepSeek-R1-Distill-Qwen-1.5B (DeepSeek-AI et al., 2025a)	21.30*	20.55	0/14	20.91
Llama 3.2 3B Instruct (Meta, 2024c)	17.59*	25.53	0/14	21.73
Gemma 3 4B IT (Team, 2025a)	25.09*	32.86	0/14	29.14
Qwen 2.5 3B Instruct (Qwen et al., 2025)	24.74*	33.27	0/14	29.19
Qwen3-1.7B (Yang et al., 2025)	28.99	30.42	0/14	29.74
Kanana Nano 2.1B Instruct (Team et al., 2025a)	27.25*	32.60	0/14	30.04
Aya Expanse 8B (Dang et al., 2024)	28.30	31.65	0/14	30.05
EXAONE Deep 2.4B (Research et al., 2025)	28.84*	32.08	0/14	30.53
EXAONE 3.5 2.4B Instruct (Research et al., 2024)	27.72*	34.49	0/14	31.25
HyperCLOVAX-SEED-Text-Instruct-1.5B (naver hyperclovax, 2025)	33.94	30.13	0/14	31.95
Llama 3.1 8B Instruct (Meta, 2024b)	31.89	33.81	0/14	32.89
Qwen3-1.7B (w/ thinking) (Yang et al., 2025)	37.80	38.27	1/14	38.05
Ko-R1-7B-v2.1 (OneLineAI)	41.94	38.70	1/14	40.25
EXAONE 3.5 7.8B Instruct (Research et al., 2024)	41.90	41.71	0/14	41.80
EXAONE Deep 7.8B (Research et al., 2025)	44.99	41.53	0/14	43.18
Qwen3-8B (Yang et al., 2025)	49.25	46.92	1/14	48.03
Qwen3-8B (w/ thinking) (Yang et al., 2025)	58.79	55.27	3/14	56.95

Table 6: The main evaluation results of KMMLU-REDUX and KMMLU-PRO benchmarks on smaller (< 10B) LLMs. The gray-shaded models are the dense-reasoning models. The KMMLU-REDUX scores with * are biased as these models are used for the dataset filtration (Section 2.2.1).

Domain		Safety Management	Environment & Energy	Construction	Mechanical Engineering	Materials	Information & Communication	Agriculture, Forestry & Fisheries	Chemistry	Culture, Arts, Design & Broadcasting	Electrical & Electronics	Food & Processing	Architecture	Management, Accounting & Office	Mining & Resources	Avg.
Opensourced Models																
< 5B	DeepSeek-R1-Distill-Qwen-1.5B	20.50	17.77	24.32	20.00	24.43	26.49	17.84	20.71	19.33	20.18	15.38	18.31	20.59	31.58	21.30*
	Llama 3.2 3B Instruct	13.00	12.20	13.81	14.81	11.07	15.68	16.76	11.43	17.65	7.89	16.67	11.27	5.88	15.79	17.59*
	Qwen 2.5 3B Instruct	16.50	17.24	15.32	17.41	15.27	20.54	18.38	11.43	26.05	16.67	11.54	15.49	14.71	10.53	24.74*
	Gemma 3.4B IT	23.00	23.87	22.52	25.56	20.23	28.65	32.43	25.71	21.01	27.19	21.79	21.13	14.71	10.53	25.09*
	Kanana Nano 2.1B Instruct	24.75	27.32	25.83	26.67	29.77	22.16	26.49	23.57	25.21	28.07	28.21	28.17	26.47	10.53	27.25*
	EXAONE 3.5 2.4B Instruct	26.75	26.79	22.22	28.15	29.39	28.65	29.19	32.86	31.93	28.95	23.08	33.80	26.47	31.58	27.72*
< 10B	EXAONE Deep 2.4B	27.50	30.77	23.72	28.52	32.06	30.27	22.16	40.00	25.21	34.21	25.64	18.31	17.65	15.79	28.84*
	Qwen3-1.7B	30.50	27.06	27.93	36.67	32.82	29.73	23.24	30.00	23.53	28.95	25.64	23.94	20.59	15.79	28.99
	HyperCLOVAX-SEED-Text-Instruct-1.5B	28.00	25.99	30.93	32.59	33.21	29.19	32.97	22.86	35.29	28.07	32.05	32.39	17.65	26.32	33.94
	Qwen3-1.7B (w/thinking)	37.25	35.54	33.03	39.63	40.08	41.62	36.76	40.71	35.29	46.49	42.31	29.58	41.18	42.11	37.80
	Aya Expans 8B	25.00	20.42	23.72	21.11	24.05	29.19	29.19	21.43	23.53	22.81	20.51	30.99	26.47	21.05	28.30
	Llama 3.1 8B Instruct	30.50	25.99	23.12	28.15	30.53	33.51	32.97	22.14	35.29	24.56	29.49	26.76	29.41	31.58	31.89
< 20B	EXAONE 3.5 7.8B Instruct	41.75	37.40	41.14	44.81	42.37	48.11	41.62	40.00	47.06	42.11	38.46	46.48	38.24	21.05	41.90
	Ko-R1-7B-v2.10	42.50	37.93	39.94	41.48	38.55	50.81	37.30	52.14	49.58	50.00	33.33	33.80	50.00	36.84	41.94
	EXAONE Deep 7.8B	40.75	45.09	40.54	49.26	46.56	51.35	31.89	57.14	42.02	55.26	47.44	39.44	44.12	26.32	44.99
	Qwen3-8B	47.00	48.01	43.84	49.63	56.11	56.22	41.62	58.57	47.06	60.53	46.15	43.66	50.00	31.58	49.25
	Qwen3-8B (w/thinking)	55.75	57.82	53.15	67.78	64.89	62.16	45.95	65.71	57.14	66.67	56.41	53.52	58.82	63.16	58.79
	Gemma 3.12B IT	44.50	42.18	46.25	44.81	46.95	50.81	51.35	50.71	57.98	48.25	46.15	39.44	50.00	42.11	46.70
< 32B	Phi-4 (14B)	48.50	47.75	48.05	52.59	49.62	52.43	44.86	58.57	54.62	56.14	50.00	36.62	52.94	36.84	49.75
	Qwen3-14B	52.75	55.70	51.95	63.33	61.45	65.95	47.57	61.43	57.98	67.54	52.56	59.15	61.76	47.37	57.25
	Qwen3-14B (w/thinking)	64.50	66.58	60.36	70.00	72.90	68.65	54.05	73.57	63.87	75.44	51.28	57.75	70.59	68.42	65.71
	Aya Expanse 32B	30.50	32.10	27.03	38.15	35.50	39.46	37.84	28.57	33.61	35.09	32.05	35.21	26.47	21.05	33.05
	EXAONE 3.5 32B Instruct	46.75	45.89	46.25	46.30	55.34	54.05	55.14	47.14	57.98	50.00	50.00	50.70	55.88	31.58	49.40
	Mistral Small 3.1 Instruct (24B)	45.00	50.40	51.95	55.93	59.54	58.38	51.35	58.57	56.30	58.77	47.44	52.11	58.82	31.58	52.92
< 70B	Gemma 3.27B IT	49.50	51.19	47.45	57.41	58.02	62.70	49.73	60.71	59.66	62.28	56.41	47.89	67.65	31.58	54.04
	EXAONE Deep 32B	53.25	62.07	53.75	66.30	59.54	67.03	53.51	62.86	57.14	60.53	56.41	45.07	58.82	21.05	58.33
	Qwen3-30B-A3B	54.25	57.56	54.65	58.52	64.50	61.08	51.89	65.00	62.18	68.42	58.97	56.34	58.82	52.63	58.41
	Qwen3-32B	59.50	64.19	60.66	68.52	72.52	69.19	58.38	68.57	74.79	72.81	58.97	53.52	70.59	63.16	64.98
	Qwen3-30B-A3B (w/thinking)	63.00	64.19	62.46	69.63	70.23	69.19	56.76	69.29	68.91	72.81	53.85	56.34	70.59	68.42	65.25
	QwQ 32B	61.25	67.64	68.17	71.11	74.05	72.97	58.38	71.43	72.27	71.05	53.85	56.34	79.41	52.63	67.34
	Qwen3-32B (w/thinking)	64.50	68.97	66.97	74.07	75.57	70.81	62.70	74.29	68.91	72.81	64.10	53.52	73.53	57.89	68.77
Closed Models	Llama 3.3 70B Instruct	52.25	54.11	52.85	62.59	59.92	64.32	52.43	50.71	60.50	60.53	55.13	53.52	64.71	36.84	56.17
	C4AI Command A (111B)	56.75	66.58	63.06	63.33	64.89	67.57	62.16	60.71	68.91	64.04	62.82	59.15	55.88	47.37	62.93
	DeepSeek V3 (671B)	62.50	62.33	63.66	66.30	72.52	72.97	62.70	66.43	67.23	69.30	69.23	59.15	61.76	63.16	65.64
	Llama-4-Scout-17B-16E-Instruct	64.00	67.64	65.77	69.63	77.10	71.35	61.08	69.29	66.39	71.93	62.82	57.75	64.71	57.89	67.49
	Qwen3-235B-A22B	64.25	72.68	66.37	71.11	82.06	72.43	65.41	71.43	68.07	73.68	65.38	52.11	67.65	47.37	69.54
	Qwen3-235B-A22B (w/thinking)	69.75	75.86	71.47	81.11	85.11	76.22	64.86	75.00	70.59	78.95	73.08	61.97	82.35	68.42	74.49
Closed Models	Llama-4-Maverick-17B-128E-Instruct	73.50	76.13	78.38	79.26	83.97	80.54	77.30	75.00	78.15	79.82	73.08	71.83	82.35	73.68	77.58
	DeepSeek R1 (671B)	72.50	76.39	79.58	80.74	85.11	81.62	82.70	77.14	79.83	81.58	75.64	63.38	85.29	73.68	78.51
	GPT-4.1 mini (2024-04-14)	61.00	64.99	63.66	69.26	74.81	72.43	63.78	73.57	70.59	66.67	75.64	56.34	73.53	57.89	67.03
	o3-mini (2025-01-31)	63.75	66.05	63.96	72.96	73.66	75.14	63.78	69.29	72.27	75.44	66.67	46.48	76.47	57.89	67.84
	Grok-3-mini-beta	69.75	69.76	69.07	73.33	83.97	77.30	65.95	70.71	68.07	73.68	61.54	59.15	76.47	73.68	71.47
	Grok-3-beta	70.00	70.56	68.47	76.30	83.97	76.76	69.19	75.71	74.79	69.30	70.51	71.83	73.53	57.89	72.90
Closed Models	o4-mini (2025-04-16)	67.75	73.74	76.28	78.52	82.06	81.08	76.76	78.57	80.67	79.82	71.79	61.97	79.41	78.95	75.80
	GPT-4.1	71.50	74.27	72.07	77.41	85.11	81.62	80.00	74.65	76.47	77.19	75.64	61.97	76.47	68.42	75.86
	Claude 3.7 Sonnet	70.75	76.92	76.28	80.37	83.97	77.84	78.92	78.57	76.47	78.07	75.64	64.79	79.41	68.42	76.88
	Claude 3.7 Sonnet (w/thinking)	71.50	80.37	79.88	80.74	87.02	82.16	80.00	82.86	80.67	79.82	73.08	70.42	85.29	68.42	79.36
	o3	77.75	78.25	77.78	80.37	85.50	78.92	83.24	78.17	84.03	85.09	82.05	66.20	85.29	78.95	79.92
	o1	75.75	79.58	81.98	81.11	90.84	80.54	84.86	77.86	84.03	81.58	83.33	70.42	85.29	73.68	81.14

Table 7: The break down results for 14 categories in KMMLU-REDUX. The gray-shaded models are the dense-reasoning models. The scores with * are biased as these models are used for the dataset filtration (Section 2.2.1).

Domain	Law				Tax & Accounting				Value Estimation				Medical				
Names of KNPLs	Judicial Scrivener	Lawyer	Public Labor Attorney	Patent Attorney	Public Accountant (CPA)	Tax Accountant	Customs Broker	Damage Adjuster (CDA)	Appraiser	Doctor of Korean Medicine	Dentist	Pharmacist	Herb Pharmacist	Physician	Avg.	# of passed NPLs	
Opensourced Models																	
< 5B	DeepSeek-R1-Distill-Qwen-1.5B Llama 3.2 3B Instruct HyperCLOVAX-SEED-Text-Instruct-1.5B Qwen3-1.7B	15.66 27.27 24.24 19.19	13.33 24.67 30.13 20.67	24.27 25.94 26.61 30.96	25.69 16.35 25.48 32.11	23.56 23.95 29.83 21.63	18.07 23.95 30.19 21.01	18.24 23.9 30.19 23.39	20.00 23.33 41.67 30.0	21.94 19.39 25.00 29.08	20.49 23.26 30.90 30.56	18.49 26.45 27.31 30.32	23.25 36.53 36.16 46.49	23.77 30.74 35.66 46.72	18.67 28.67 33.33 33.33	20.55 25.53 30.13 30.42	0/14
< 10B	EXAONE Deep 2.4B Kanana Nano 2.1B Instruct Gemma 3 4B IT Qwen 2.5 3B Instruct EXAONE 3.5 2.4B Instruct Qwen3-1.7B (w/thinking)	16.67 19.19 23.23 20.20 28.28 25.25	21.33 30.13 32.22 35.15 24.67 25.33	31.38 22.02 26.61 24.77 32.22 31.8	36.70 23.08 26.47 30.19 27.4 34.86	25.00 23.08 26.47 31.09 27.4 25.48	23.53 22.08 26.47 30.19 27.67 29.56	32.08 30.19 27.67 30.0 27.04 29.56	37.50 28.33 29.17 35.00 35.00 34.17	30.61 26.53 23.47 37.85 31.63 36.22	30.21 42.01 37.15 37.85 36.81 44.79	29.25 31.61 36.56 35.05 35.91 38.06	52.77 52.40 49.08 47.97 53.14 61.99	44.67 44.67 43.03 52.46 43.44 59.43	32.00 38.67 36.00 34.67 38.67 44.67	30.27 32.60 32.86 33.27 34.49 38.27	0/14
< 20B	Aya Expans 8B Llama 3.1 8B Instruct Ko-R1-7B-v2.1 EXAONE Deep 7.8B EXAONE 3.5 7.8B Instruct Qwen3-8B Qwen3-8B (w/thinking)	20.20 28.28 40.59 26.26 28.79 29.29 27.78	20.00 30.96 28.85 38.91 38.49 34.00 36.0	30.54 30.96 38.33 42.20 35.10 46.03 49.37	22.02 23.85 39.8 42.20 31.09 35.78 41.28	23.56 21.15 35.22 32.21 32.21 33.17 48.08	18.49 23.53 33.94 38.99 31.09 35.71 41.28	37.11 33.33 26.47 38.99 37.74 34.59 49.69	22.45 26.02 42.21 40.31 33.16 40.82 53.06	36.11 35.42 29.29 42.36 44.10 50.69 61.81	29.25 33.76 48.67 44.73 43.44 46.02 54.19	50.18 54.61 46.24 64.94 67.53 71.96 83.03	46.31 50.82 42.00 53.69 56.56 59.33 76.23	42.00 30.67 38.70 41.53 41.71 46.92 75.33	31.65 33.81 38.70 41.53 41.71 46.92 55.27	0/14	
< 32B	Phi-4 (14B) Gemma 3 12B IT Qwen3-14B Qwen3-14B (w/thinking)	33.33 28.79 32.83 30.81	34.00 23.33 30.67 36.67	41.00 40.59 39.75 48.95	36.70 39.45 39.75 47.71	37.50 34.62 48.56 51.26	37.82 38.66 38.24 52.2	42.77 42.77 42.14 52.2	44.17 43.37 49.17 47.50	43.37 43.37 50.00 54.59	45.49 52.08 57.64 65.28	41.29 49.46 55.91 61.94	72.69 62.30 80.74 87.82	57.38 71.59 59.33 80.67	51.33 45.82 53.02 59.48	45.32 45.82 53.02 59.48	1/14
> 70B	Aya Expans 32B EXAONE 3.5 32B Instruct Mistral Small 3.1 Instruct (24B) Gemma 3 27B IT Qwen3-30B-A3B EXAONE Deep 32B Qwen3-32B Qwen3-30B-A3B (w/thinking) Qwen3-32B (w/thinking) QwQ 32B	24.75 33.33 27.27 29.80 31.31 30.30 33.33 35.35 33.33 35.35	18.67 26.67 33.43 31.33 43.51 46.00 53.14 50.21 55.23 47.12	24.27 43.10 43.51 44.35 33.94 47.28 53.14 45.87 55.23 43.12	32.11 38.53 39.45 38.46 39.90 35.78 43.12 48.62 50.00 44.04	20.67 32.35 38.66 42.38 39.90 45.19 57.69 59.21 55.29 61.06	21.01 42.77 42.77 40.88 39.90 35.71 43.12 48.62 50.00 56.72	27.67 42.77 42.77 44.39 42.77 44.03 53.46 50.44 54.09 56.60	34.17 52.50 47.50 43.33 49.17 47.96 53.33 50.00 57.50 66.84	23.98 41.33 41.33 44.39 49.17 47.96 48.98 48.98 55.61 66.84	32.64 52.08 52.08 52.08 56.60 53.47 61.46 70.49 70.14 71.53	31.83 58.49 58.49 58.49 53.98 52.26 60.86 59.78 66.88 62.24	53.14 64.95 81.18 88.19 84.43 55.91 83.03 77.87 72.00 84.47	45.49 49.89 49.89 51.88 56.60 52.26 85.24 84.84 84.00 87.82 80.74 82.67	31.26 49.49 51.03 58.86 52.33 52.33 58.86 60.52 61.14 63.94	0/14	
Closed Models	GPT-4.1 mini (2024-04-14) o3-mini (2025-01-31) Grok-3-mini-beta Grok-3-beta o4-mini (2025-04-16) GPT-4.1 o3 Claude 3.7 Sonnet Claude 3.7 Sonnet (w/thinking) o1	38.38 38.38 37.37 42.42 37.37 47.47 45.96 55.56 50.00 54.55	32.00 38.67 34.67 41.33 37.37 50.00 41.33 53.33 56.00 49.33	56.90 51.46 57.74 59.00 59.00 48.62 52.33 54.81 55.96 61.51	56.88 47.71 48.62 55.96 49.54 67.30 71.13 49.54 59.56 57.80	54.81 60.10 67.79 63.94 69.23 66.83 71.13 69.23 61.06 61.51	48.32 47.06 49.58 57.98 55.04 59.24 57.80 57.80 56.72 57.80	50.31 50.31 54.51 62.89 66.67 67.30 70.44 67.92 56.60 61.51	50.00 52.50 54.17 54.17 59.17 66.67 71.94 71.43 71.30 57.23	55.10 57.65 64.80 64.80 73.43 68.37 71.94 77.43 80.44 78.06	67.36 64.24 66.32 70.49 76.74 80.21 70.83 76.34 70.51 83.33	72.47 76.56 73.98 81.17 82.15 82.80 81.51 84.46 92.21 94.10	90.77 88.93 91.14 87.82 93.36 92.62 92.21 93.33 93.03 95.49	62.18 62.05 65.08 68.37 69.65 72.99 73.60 74.52 77.70	4/14 3/14 5/14 7/14 6/14 10/14 9/14 10/14 6/14 7/14		

Table 8: The break down results for all KNPLs in KMMLU-PRO. The gray-shaded models are the dense-reasoning models. The blue scores indicate that the LLM obtain the license. The details for the pass criteria of each license are described in Appendix D.3.

Year	Tests
2005	Master Craftsman Construction Equipment Maintenance, Master Craftsman Building General Work, Master Craftsman Precious Metal Processing, Master Craftsman Confectionary Making, Master Craftsman Casting, Master Craftsman Sheet-Metal & Boiler Making
2008	Master Craftsman Architectural Carpentry, Master Craftsman Surface Treatment
2010	Master Craftsman Railway Vehicles Maintenance
2014	Master Craftsman Welding
2016	Master Craftsman Steel Making
2017	Master Craftsman Metal Material, Master Craftman Metal Mould, Master Craftsman Cook
2018	Master Craftsman Gas, Master Craftsman Machinery Maintenance, Master Craftsman Plumbing, Master Craftsman Rolling, Master Craftsman Energy Management, Master Craftsman Hazardous Material, Master Craftsman Motor Vehicles Maintenance, Master Craftsman Electricity, Master Craftman Electronics, Master Craftsman Iron Making
2020	Engineer Radio Electronic Communication, Engineer Floral Design
2021	Engineer Construction Equipment, Engineer Machinery Design, Engineer Agricultural Health and Safety, Engineer Leak Nondestructive Testing, Engineer Radiation Nondestructive Testing, Engineer Biology Classification—Animal, Engineer Aquaculture, Engineer Visual Communication Design, Engineer Eddy Current Nondestructive Testing, Engineer Welding, Engineer Biomedical, Engineer Magnetic Nondestructive Testing, Engineer Electric Railway, Engineer Computer, Engineer Concrete, Engineer Explosives Handling
2022	Engineer Gas, Engineer Construction Safety, Engineer Construction Material Testing, Engineer Architecture, Engineer Building Facilities, Engineer Air-Conditioning Refrigerating Machinery, Engineer Transportation, Engineer Metal, Engineer Meteorology, Engineer Air Pollution Environmental, Engineer Urban Planning, Engineer Bioprocess, Engineer Forest, Engineer Industrial Safety, Engineer Industrial Hygiene Management, Engineer Plant Maintenance, Engineer Fire Protection System—Mechanical, Engineer Fire Protection System—Electrical, Engineer Noise & Vibration, Engineer Water Pollution Environmental, Engineer Elevator, Engineer Plant Protection, Engineer Food Processing Safety, New and Renewable Energy Equipment (Photovoltaic) Engineer, Engineer Interior Architecture, Engineer Energy Management, Engineer Greenhouse Gas Management, Engineer Organic Agriculture, Engineer Ergonomics, Engineer General Machinery, Engineer Motor Vehicles Maintenance, Engineer in Nature Environment and Ecological Restoration, Engineer Electric Work, Engineer Electricity, Engineer Computer System Application, Engineer Electronics, Engineer Information Processing, Engineer Landscape Architecture, Engineer Seeds, Engineer Cadastral Surveying, Engineer Railroad Signal Apparatus, Engineer Ultrasonic Nondestructive Testing, Engineer Livestock, Engineer Surveying Geo-Spatial Information, Engineer Penetrante Nondestructive Testing, Engineer Colorist, Engineer Civil Engineering, Engineer Soil Environment, Master Craftsman Telecommunication Apparatus, Engineer Wastes Treatment, Engineer Quality Management, Engineer Ocean Environment, Engineer Chemical Industry, Fire Investigation & Evaluation Engineer, Engineer Chemical Analysis
2023	Engineer Radio Telecommunication Equipment, Engineer Broadcasting Communication, Engineer Information Communication

Table 9: Redux Years and National Qualification Test Additions

Error Type	Examples
Ill-posed Question	<p>Category: Political science and sociology A, B에 대한 설명으로 옳은 것만을 <보기>에서 고르면? <i>What is the correct explanation about A, B in <Reference>?</i></p>
Leaked Answer	<p>Category: Ecology 산복수로에서 쌓기공작물의 높이가 3m이고, 수로깊이가 1m일 때 수로받이의 근사적 길이는? (문제 오류로 현재 복원중입니다. 보기 내용을 아시는 분들께서는 오류 신고를 통하여 보기 작성 부탁 드립니다. 정답은 3번입니다.) <i>What is the approximate length of the culvert if the pile is 3 meters high and the channel is 1 meter deep? (This is currently being restored due to a question error. If you know the referebce, please report the error. The correct answer is 3.)</i></p>
Notation Error	<p>Category: Math 다항식 $x^{2017}-1$을 x^2-x로 나누었을 때의 나머지를 $R(x)$라 할 때, $R(2017)$의 값은? <i>If the remainder of the polynomial $x^{2017}-1$ divided by x^2-x is called $R(x)$, what is the value of $R(2017)$?</i></p>
Bad Clarity	<p>Category: Education 정신분석 상담과 행동주의 상담의 공통점에 해당하는 것은? A. 상담과정에서 과거 경험보다 미래 경험을 중시한다. B. 상 담 기 법 보다는 상 담 자 의 인 간 적 자 질 과 진 솔 한 태 도 를 중 시 한 다 . C. 인간의 행동을 인과적 관계로 해석하는 결정론적 관점을 가진 다 . D. 비합리적 신념을 인식하고 수정하는 논박 과정을 중시한 다 . <i>Which is a common feature of psychoanalytic counseling and behavioral counseling? A. Emphasizes future experiences over past experiences in the counseling process. B. Prioritizes the counselor's human qualities and sincerity over counseling techniques. C. Interprets human behavior through a deterministic perspective based on causal relationships. D. Focuses on the disputation process to recognize and modify irrational beliefs.</i></p>

Table 10: Examples of error types in KMLLU. Each example demonstrates a specific issue that impacts the reliability of the benchmark. Gray text represents translation of the examples in English