

Self-Verification is All You Need To Pass The Japanese Bar Examination

Anonymous ACL submission

Abstract

Despite rapid advances in large language models (LLMs), achieving reliable performance on highly professional and structured examinations remains a significant challenge. The Japanese bar examination is a particularly demanding benchmark, requiring not only advanced legal reasoning but also strict adherence to complex answer formats that involve joint evaluation of multiple propositions. While recent studies have reported improvements by decomposing such questions into simpler true-false judgments, these approaches have not been systematically evaluated under the original exam format and scoring scheme, leaving open the question of whether they truly capture exam-level competence. In this paper, we present a self-verification model trained on a newly constructed dataset that faithfully replicates the authentic format and evaluation scale of the exam. Our model is able to exceed the official passing score when evaluated on the actual exam scale, marking the first demonstration, to our knowledge, of an LLM passing the Japanese bar examination without altering its original question structure or scoring rules. We further conduct extensive comparisons with alternative strategies, including multi-agent inference and decomposition-based supervision, and find that these methods fail to achieve comparable performance. Our results highlight the importance of format-faithful supervision and consistency verification, and suggest that carefully designed single-model approaches can outperform more complex systems in high-stakes professional reasoning tasks. Our dataset and codes are publicly available.¹

1 Introduction

Large language models (LLMs) have demonstrated remarkable capabilities across a wide range of natural language processing tasks, including question

answering (Yue, 2025; Lehmann et al., 2024), summarization (Liu et al., 2023), and even domain-specific reasoning in areas such as mathematics (Shao et al., 2024; Yang et al., 2024) and programming (Jiang et al., 2024; El-Kishky et al., 2025). Nevertheless, their performance in highly professional domains, particularly law, remains uneven. Legal reasoning often requires precise interpretation of statutes, careful evaluation of multiple interacting conditions, and strict adherence to task-specific output formats, posing challenges that go beyond surface-level linguistic competence.

The Japanese bar examination (司法試験) represents one of the most demanding legal benchmarks in this regard. In addition to its substantive difficulty, its multiple choice exam (短答式) is characterized by a distinctive question format in which examinees must jointly evaluate multiple statements and select correct combinations under rigid answer constraints. Errors in even a single constituent can invalidate an otherwise plausible answer. As we show later in this paper, base LLMs perform poorly under this evaluation regime, highlighting a substantial gap between general language understanding and exam-level legal competence.

Recent work has sought to improve LLM performance on the Japanese bar examination through dataset construction and task reformulation. Notably, the Japanese Bar Exam Question Answering (JBE-QA) dataset (Cao et al., 2025) decomposes complex exam questions into collections of independent true-false judgments, thereby simplifying the learning problem and enabling more stable training. While such decomposition-based approaches have demonstrated promising results on their own benchmarks, they fundamentally alter the structure of the original exam and do not directly address whether models trained in this manner can succeed when confronted with authentic exam questions and scoring criteria.

In this work, rather than modifying the task to

¹To be available upon publication.

083 suit the model, we develop a consistency-verifying
084 fine-tuning strategy in which the model is trained to
085 generate an answer and subsequently verify its own
086 prediction in the context of the original question.
087 We also design a dataset that preserves the original
088 exam format. This answer-conditioned verification
089 step leverages the model’s strength as an evalua-
090 tor of candidate solutions and leads to substantial
091 performance gains without increasing model size
092 or relying on external tools. We also investigate
093 more complex inference strategies, including multi-
094 agent architectures, but find that they do not yield
095 improvements under realistic exam conditions.

096 Through experiments on the Japanese bar exam-
097 ination, we demonstrate that a single fine-tuned
098 model with self-verification trained on a format-
099 faithful dataset can surpass the official passing
100 threshold on the actual exam scale. Specifically,
101 our model obtains the score of 96 from 2024
102 Japanese bar exam whose passing score is 93.
103 These results suggest that careful alignment be-
104 tween supervision format and evaluation criteria is
105 crucial for advancing LLM performance in high-
106 stakes professional domains.

107 2 Related Work

108 Numerous benchmarks have been proposed to eval-
109 uate LLMs on legal reasoning tasks. In the United
110 States, prior work has shown that GPT-4 achieves
111 passing-level performance on simulated bar exam-
112 inations (Katz et al., 2024), substantially outper-
113 forming earlier models. Other studies have ex-
114 amined legal reasoning across jurisdictions using
115 datasets such as LegalBench (Guha et al., 2023),
116 LawBench (Fei et al., 2023), and national exam-
117 derived corpora in China (Li et al., 2024), Korea
118 (Kimyeeun et al., 2024), and Europe (Chlapanis
119 et al., 2024). These benchmarks cover a wide
120 range of tasks, including multiple-choice question
121 answering, statute interpretation, case outcome pre-
122 diction, and legal text entailment. While these ef-
123 forts demonstrate that modern LLMs can perform
124 competitively on legal tasks, many rely on task
125 reformulation, simplified supervision, or indirect
126 evaluation metrics.

127 Evaluating language models on the Japanese bar
128 examination has recently attracted increasing atten-
129 tion as a challenging benchmark for legal reasoning
130 in Japanese. The examination covers multiple le-
131 gal domains, including constitutional law, civil law,
132 and criminal law, and is characterized by questions

133 that require joint evaluation of multiple proposi-
134 tions under strict answer-format constraints. Earlier
135 Japanese legal NLP resources, such as the COLIEE
136 shared tasks (Thanh et al., 2020, 2021), primarily
137 focused on civil law and emphasized subtasks like
138 information retrieval or textual entailment, rather
139 than full exam-style question answering. More re-
140 cent efforts (Nguyen et al., 2025) have attempted
141 to construct datasets directly derived from the bar
142 examination, highlighting both the difficulty of the
143 content and the importance of handling the exam’s
144 unique structure. The most notable recent bench-
145 mark is the Japanese Bar Exam Question Answer-
146 ing (JBE-QA) dataset (Cao et al., 2025), which
147 reformulates past bar exam questions into collec-
148 tions of independent true/false statements. Each
149 original question is decomposed into multiple bi-
150 nary judgments, allowing models to be trained and
151 evaluated on simplified supervision signals. Using
152 this formulation, JBE-QA evaluates a wide range
153 of proprietary and open-source LLMs under zero-
154 shot and few-shot settings, and reports substantial
155 performance gains for state-of-the-art models, par-
156 ticularly when chain-of-thought prompting is en-
157 abled. While this decomposition strategy stabilizes
158 learning and evaluation, it fundamentally alters the
159 structure of the original exam. As a result, it re-
160 mains unclear whether models trained under this
161 paradigm can succeed when confronted with intact
162 exam questions that require reasoning over multiple
163 interacting propositions and adherence to the origi-
164 nal scoring rules. Our work directly addresses this
165 gap by evaluating models on the authentic exam
166 format and scale.

167 Recent work has explored multi-agent architec-
168 tures for legal reasoning (Zhang and Ashley, 2025;
169 Sun et al., 2024), in which multiple LLM agents
170 collaborate or debate to produce a final answer.
171 Such approaches have been applied to tasks like
172 legal argument generation and multi-step legal an-
173 alysis, often improving factual coverage or inter-
174 pretability. However, these systems introduce ad-
175 ditional complexity and inference cost, and their
176 effectiveness under strict exam-style evaluation re-
177 mains underexplored. On the other hand, consis-
178 tency verification (Patwardhan et al., 2024) and
179 reflection-based methods (Renze and Guven, 2024)
180 have been proposed as general techniques to im-
181 prove LLM reliability. These approaches encour-
182 age a model to evaluate or critique its own output,
183 leveraging the observation that LLMs are often
184 stronger evaluators than generators.

Table 1: Comparison between our dataset, and the JBE-QA dataset. While JBE-QA decomposes a single exam question into multiple independent true/false items, our dataset preserves the original joint decision structure. Despite having substantially fewer questions, fine-tuning on our dataset yields significantly higher performance on the actual exam.

Dataset	Question	Answer	#Questions
JBE-QA	憲法第31条の定める法定手続の保障は、直接には刑事手続に関するものであるが、行政手続にも及ぶと解すべき場合があり、その場合には行政処分相手方に常に事前の告知、弁解、防御の機会を与える必要がある。(Article 31 always requires prior notice and defense opportunity in administrative procedures.)	False	2,770
	憲法第35条は、住居、書類及び所持品について、侵入、搜索及び押収を受けることのない権利を規定しているが、この規定の保障対象には、住居、書類及び所持品に準ずる私的領域に侵入されることのない権利が含まれる。(Article 35 protects private domains equivalent to residences, papers, and effects.)	True	
	憲法第38条第1項は、自己が刑事上の責任を問われるおそれのある事項について供述を強要されないことを保障するものであり、氏名の供述も、これによって自己が刑事上の責任を問われるおそれがあることから、原則として保障が及ぶ。(Article 38(1) generally protects refusal to state one's name.)	False	
Ours (identical format as the original exam)	刑事手続上の権利に関する次のアからウまでの各記述について、最高裁判所の判例の趣旨に照らして、それぞれ正しい場合には1を、誤っている場合には2を選びなさい。 ア. 憲法第31条の定める法定手続の保障は、直接には刑事手続に関するものであるが、行政手続にも及ぶと解すべき場合があり、その場合には行政処分相手方に常に事前の告知、弁解、防御の機会を与える必要がある。 イ. 憲法第35条は、住居、書類及び所持品について、侵入、搜索及び押収を受けることのない権利を規定しているが、この規定の保障対象には、住居、書類及び所持品に準ずる私的領域に侵入されることのない権利が含まれる。 ウ. 憲法第38条第1項は、自己が刑事上の責任を問われるおそれのある事項について供述を強要されないことを保障するものであり、氏名の供述も、これによって自己が刑事上の責任を問われるおそれがあることから、原則として保障が及ぶ。 (Regarding the following statements (A) through (C) concerning rights in criminal procedure, select 1 if the statement is correct in light of Supreme Court precedents, and select 2 if it is incorrect. (A) Article 31 may extend to administrative procedures and always requires prior notice and opportunity to defend. (B) Article 35 protects not only residences, papers, and effects but also equivalent private domains. (C) Article 38(1) protects against compelled self-incrimination, and this protection generally applies to stating one's name.)	2,1,2	460

3 Method

3.1 Dataset Construction

We collect actual exam questions spanning 6 years (2019-2024) from the Japanese Ministry of Justice, where we separate 2024 (*Reiwa 6* or R6) as the test set. We construct a dataset that faithfully replicates the format and evaluation criteria of the Japanese bar examination. Unlike prior work that decomposes questions into independent true/false statements, each instance in our dataset corresponds to a complete exam question, including all constituent statements and the original answer choices. Answers are represented exactly as required in the exam, such as concatenated numeric labels indicating the correctness of each statement or indices corresponding to valid combinations.

Each question is annotated with its subject category (constitutional law(憲法), civil law(民法), or criminal law(刑法)), year of administration, and the points rewarded. This allows evaluation not only in terms of accuracy but also in terms of the official point-based scoring scheme used in the exam. We split the dataset by year, using earlier years (2019-2023, R1-R5) for training and reserving 2024 (R6) exam for evaluation, mirroring realistic exam preparation, where the examinees rely on the past exams for reference.

Unlike decomposed formulations that reduce

each statement to an independent true/false query, correctness in the actual exam depends on the joint evaluation of all statements. For example, an answer may require selecting a single option corresponding to a specific combination of statement truth values, or outputting a concatenated sequence of digits (e.g., “112”) aligned with multiple statements. An answer is considered incorrect if either the semantic judgment or the required format is violated. Table 1 contrasts the format of our dataset with that of JBE-QA.

3.2 Self-Verification

Fine-Tuning: Our approach combines supervised fine-tuning with answer-conditioned self-verification. During training, the model is fine-tuned to generate the correct exam-style answer given the full question, without decomposing the question into simpler subproblems. Given a question q_i consisting of a set of statements $\{s_{i1}, s_{i2}, \dots, s_{in}\}$ and a set of valid answer formats defined by the Japanese bar examination, the model is required to produce a single answer a_i that satisfies both semantic correctness and strict formatting constraints, where a_i may contain multiple integers as in the original exam.

Self-Verification: At inference time, we introduce a verification step in which the model re-evaluates its own predicted answer in the context of

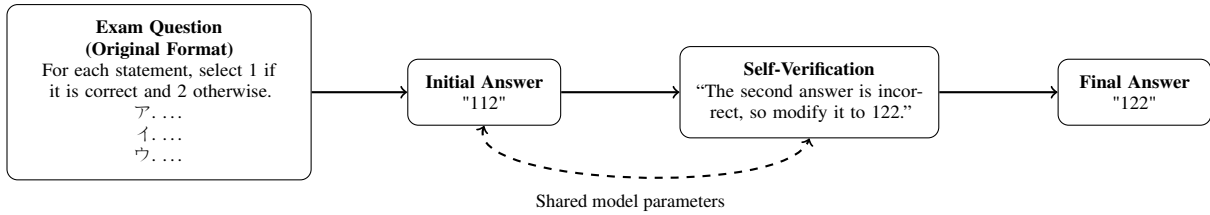


Figure 1: An overview of our method of self-verification with a shared model under the original exam format.

the original question. Importantly, this verification is performed by the same fine-tuned model, but under a different prompt that induces verification-oriented behavior.

Formally, let $f_{\theta}(q)$ denote the model’s initial prediction for question q . We then define a verification function $g_{\theta}(q, f_{\theta}(q))$, which produces a revised answer by assessing the consistency between the question and the initially predicted answer. The final output is given by $\hat{a} = g_{\theta}(q, f_{\theta}(q))$. Although f_{θ} and g_{θ} share the same parameters, they are instantiated with distinct prompts, where one encouraging answer generation and the other encouraging conservative correction. This procedure incurs only a single additional forward pass at inference time, while substantially improving robustness against formatting errors and local reasoning mistakes. Figure 1 describes the overall workflow of our approach.

Prompt Design: Table 2 summarizes the prompts used in our system. Answer format instruction prompts have been calibrated to enforce the strict formatting required, which eradicates the necessity for normalization schemes to account for various output formats that frequently occur when bar exam questions are asked. The verification prompt explicitly instructs the model to preserve the original answer unless a clear inconsistency is detected, which we find crucial for preventing unnecessary corrections.

4 Experiments

4.1 Experimental Setting

We evaluate our method on the Japanese bar examination questions from Reiwa 6 (R6), using questions from earlier years (R1–R5) exclusively for fine-tuning. We use GPT-4.1 (OpenAI, 2025) as our base model, and examine both zero-shot and few-shot setting. For few-shot setting, we chose 5 sample demonstrations from the training set. We also fine-tune separate GPT-4.1 models, with our dataset and with JBE-QA respectively. In order to

examine the effect of self-verification, we report the results obtained both with and without self-verification. The same exact prompts were used for all models. For each variant of the models, the experiments were repeated for 3 times.

We report exact-match accuracy as well as the official examination point score, which awards partial credit according to the exam’s grading rules. Partial credit scheme works as following; when 3 or more questions are grouped together with n points, getting one question wrong results in $n - 2$ points, whereas getting two or more questions wrong result in 0 point. For example, if 5 questions are grouped together with 4 points, getting 4 questions correct yields 2 points, but getting 3 questions correct results in 0 point, despite the accuracy being over 50%.

The exam consists of 50 points for constitutional law, 75 points for civil law, and 50 points for criminal law, adding up to 175 maximum points. In case of Reiwa 6 (R6) exam held in 2024, the actual passing score was 93. There is also an additional requirement that at least 40% of the points should be achieved in each law section.

4.2 Results & Analysis

Table 3 summarizes the results from the models examined.

Base models: Base model with zero-shot setting performs poorly, clearly suggesting that the legal knowledge on which the model has been pre-trained is insufficient for a reasonable performance on exam-level tasks. Few-shot setting hardly boosts the performance over zero-shot. While it is expected that providing a few samples would not significantly improve the legal knowledge, it also does not particularly seem to have helped in guiding exam-specific format. A clear performance boost is made by performing self-verification on base model. In fact, as we shall see, self-verification invariably boosts the performance regardless of model choice, demonstrating it is an efficient

Table 2: Prompts used for system role, answer format, and self-verification.

Purpose	Japanese Prompt	English Translation
System Role	あなたは日本の司法試験を受験する受験者である。	You are a test taker solving the Japanese bar examination.
Answer Format	<p>【回答形式の厳守】必ず「答えのみ」を出力せよ。理由・説明・記号は一切不要。</p> <p>1) 選択肢が番号で与えられている場合 (例: 1. アO イO ウO, 2. アO イO ウX...) → 正しい選択肢の番号のみ出力 (例: 2)</p> <p>2) 各記述 (ア・イ・ウ...) について1/2を答える問題の場合 → 数字列のみ出力 (例: 112)</p> <p>禁止:</p> <ul style="list-style-type: none"> - OOX - アO イO ウX - ア1 イ1 ウ2 - 説明文 	<p>【Strict answer format】 Output only the answer. Do not include any reasons, explanations, or symbols. 1) When the choices are given as numbered options (Example: 1. A○B○C○, 2. A○B○C× ...) → Output only the number of the correct option (Example: 2) 2) When each statement (A, B, C, ...) requires an answer of 1 or 2 → Output only the sequence of numbers (Example: 112) Prohibited:</p> <ul style="list-style-type: none"> - OOX - A○B○C× - A1 B1 C2 - Any explanatory text
Verification	<p>あなたは法律試験の答案を最終確認する役割である。以下の【問題】と【あなたの解答】を照らし合わせ、選択肢番号または数値の形式として最も正しい最終解答を一つだけ出力せよ。</p> <ul style="list-style-type: none"> ・ 問題文の条件に照らして明らかに誤っている場合のみ修正すること ・ 元の解答が正しい場合は、そのまま同じ解答を出力すること ・ 理由や説明は一切出力せず、最終的な数字のみを出力せよ 	<p>You are responsible for the final review of a law exam answer. Compare the following [Question] and [Your Answer], and output only one final answer in the form of a choice number or numeric value.</p> <ul style="list-style-type: none"> ・ Modify the answer only if it is clearly incorrect based on the question's conditions. ・ If the original answer is correct, output the same answer as is. ・ Do not include any reasons or explanations. Output only the final number.

Table 3: Performance comparison on the Japanese bar examination (Reiwa 6). Accuracy denotes exact-match answer accuracy, while scores follow the official exam grading scheme with partial credit. Average subject-wise scores are reported for constitutional law (憲法), civil law (民法), and criminal law (刑法).

Model	Accuracy	Exam Scale (Avg/Min/Max)	Const.	Civ.	Crim.
Passing Score for Examinees	N/A	93 (out of 175)	20	30	20
Base (Zero-Shot)	0.4036	67.0 / 65 / 68	8.0	32.0	27.0
Base (Few-Shot)	0.3896	68.3 / 63 / 71	8.0	33.3	27.0
Base (Few-Shot) + Self-Verification	0.4156	76.3 / 76 / 77	9.7	36.7	30.0
Fine-Tuned w/ JBE-QA	0.3766	64.0 / 62 / 66	8.0	30.0	26.0
Fine-Tuned w/ JBE-QA + Self-Verification	0.4226	80.7 / 78 / 82	21.0	32.7	27.0
Fine-Tuned w/ Ours	0.4675	92.3 / 91 / 93	20.3	42.0	30.0
Fine-Tuned w/ Ours + Self-Verification	0.4935	94.7 / 94 / 96	22.3	42.3	30.0
Multi-Agent (same model for all agents)	0.4026	75.7 / 74 / 79	19.3	30.7	25.7
Multi-Agent (separately fine-tuned models)	0.3969	71.0 / 66 / 77	12.7	34.7	25.6

model-agnostic technique.

JBE-QA: Despite being trained on a substantially larger dataset, the model fine-tuned on JBE-QA exhibits markedly poor performance on the actual Japanese bar examination, suggesting that improving the model’s legal knowledge does not automatically translate to performance boost in more complex tasks. While their decomposition strategy simplifies learning by isolating local factual judgments, it removes the requirement to reason over joint constraints among statements. The true/false reformulation introduces an implicit shift in task distribution, as the model is optimized for binary classification rather than constrained selection under combinatorial rules. As such, fine-tuning on decomposed propositions encourages a form of segmented knowledge representation that lacks mechanisms for re-composition at inference time.

These findings suggest that while proposition-level supervision may improve performance on simplified benchmarks, it does not necessarily transfer to evaluation settings that require holistic reason-

ing. For high-stakes professional examinations, preserving the native question format during training appears critical for enabling models to align local legal knowledge with global decision consistency.

As with the base models, the model fine-tuned with JBE-QA shows a significant performance boost with self-verification. This again reinforces that self-verification is an effective technique regardless of the model.

Ours: Fine-tuning with our dataset, with or without self-verification, clearly outperforms other approaches. Notably, fine-tuning with our dataset alone without self-verification already obtains scores around the passing score. Note that its performance gain cannot be attributed to format memorization or answer pattern learning, as the combinatorial space of such answers (e.g. "11221") makes correct prediction by memorizing answer frequencies or guessing implausible. Moreover, the training data does not provide decomposed supervision at the proposition level; the model must

internally reason over each constituent statement to produce a globally consistent output. This suggests that exposure to the authentic multi-proposition format during fine-tuning induces an ability to jointly assess multiple legal conditions within a single reasoning context. This is further supported by the fact that models trained on decomposed dataset fail to recover comparable performance when evaluated under the original exam structure, despite having access to substantially more training instances.

As with other model variants, we observe a further and consistent performance improvement by introducing self-verification at inference time. While the fine-tuned model already produces strong initial answers, the verification step allows the model to reassess the internal consistency of its own prediction against the original question, particularly for answers involving multiple propositions. This process is effective at correcting local inconsistencies, such as a single misjudged statement within an otherwise correct composite answer, which directly translates into higher exam scores under the official grading scheme. Importantly, self-verification does not introduce external supervision or additional training data, but instead leverages the model’s existing legal knowledge in a second-pass reasoning phase. The resulting improvement suggests that a non-trivial portion of errors made by strong fine-tuned models arise not from lack of legal knowledge, but from failures in global consistency, which self-verification is well suited to mitigate.

Another plausible explanation is that both format-specific fine-tuning and self-verification act as catalysts that elicit latent knowledge already present in the model. This may explain why performance improves substantially despite the relatively small size of our training data, which by itself is unlikely to contain all the legal knowledge required to pass the examination. Self-verification further amplifies this effect by encouraging the model to reassess and consolidate its own predictions, as evidenced by the consistent, model-agnostic performance gains observed when verification is applied.

In short, format-specific fine-tuning teaches the model how to exploit internal knowledge that would otherwise remain dormant, and self-verification further strengthens this elicitation process by promoting global consistency across multiple propositions.

Table 4 shows qualitative examples with the outputs from each model, along with the points

awarded to each output. Our model correctly answers both the single answer format and the composite answer format, which other models struggle to address. Improvements with self-verification can be seen in other models as well, where they receive partial credit. Note that there are instances where other models often produce incompatible outputs, such as the additional number of answers. Also, note that in many cases, 0 points are awarded even when the accuracy for the questions is over 50%. This reinforces the point that performing well on decomposed propositions does not automatically translate to equivalent performance on actual exam format and scale, and that it requires composite reasoning to demonstrate genuine success, rather than to simply claim competence based on simplified benchmarks.

5 Multi-Agent Reasoning

Multi-agent reasoning has demonstrated success in a number of complex tasks (Zhang et al., 2025), and has also been attempted in legal domain (Zhang and Ashley, 2025; Sun et al., 2024). To evaluate whether explicit decomposition into interacting agents improves performance on the Japanese bar examination, we implemented a multi-agent pipeline in which distinct agents are responsible for retrieval, verification, knowledge abstraction, and final answering. Unlike prior work that assumes homogeneous agents or informal collaboration, our implementation assigns clearly separated functional roles and evaluates both shared-model and independently fine-tuned agent configurations under the authentic exam grading scheme.

5.1 Multi-Agent Architecture

The pipeline consists of four sequential agents, loosely inspired by the architecture proposed by (Zhang et al., 2025);

- *Retriever Agent*: Given a test question q , the retriever agent selects a set of candidate past exam questions and answers from the training set (R1–R5) that it finds relevant to the question q .
- *Verifier Agent*: The verifier agent receives the test question and the retrieved candidates from the past exams, and filters them to retain only those deemed relevant. Its role is to discard superficially similar but legally irrelevant questions, thereby reducing noise before knowledge abstraction.

Table 4: Qualitative examples under the authentic Japanese bar examination format. For models, +V indicates that self-verification is performed. For each output, points awarded are displayed in parenthesis. Bold outputs indicate correct outputs and maximum points. Outputs that do not match the number of digits are actual mistakes by the models.

Model	Exam(GT)	Base(ZS)	Base(FS)	Base+V	JBE-QA	JBE-QA+V	Ours	Ours+V	MA(Same)	MA(Sep)
Question 1	<p>憲法第22条に関する次のアからウまでの各記述について、それぞれ正しい場合には1を、誤っている場合には2を選びなさい。(For each of the following statements (A)-(C) concerning Article 22 of the Constitution, select 1 if correct and 2 if incorrect.)</p> <p>ア. 判例は、日本に適法に在留する外国人には、憲法上、その在留期間内において外国へ一時旅行する自由が保障されているものと解している。(Precedent holds that foreign nationals lawfully residing in Japan are constitutionally guaranteed the freedom to temporarily travel abroad during their period of stay.)</p> <p>イ. 居住・移動の自由は、複合的な性格を有する人権と解されており、広く知的な接触の機会を得るために不可欠であることから、精神的自由の要素も併せ持っている。(Freedom of residence and movement is understood as a right with a composite character, and because it is indispensable for obtaining broad opportunities for intellectual contact, it also possesses aspects of spiritual freedom.)</p> <p>ウ. 判例は、市営住宅の入居者が暴力団員であることが判明したときには当該住宅の明渡しを請求することができるとする条例の規定による居住の制限は、公共の福祉による必要かつ合理的なものであるから、この規定は憲法第22条第1項に違反しないと解している。(Precedent holds that an ordinance allowing eviction from public housing when a resident is found to be a gang member constitutes a necessary and reasonable restriction for the public welfare, and thus does not violate Article 22(1) of the Constitution.)</p>									
Model	Exam(GT)	Base(ZS)	Base(FS)	Base+V	JBE-QA	JBE-QA+V	Ours	Ours+V	MA(Same)	MA(Sep)
Output (Pts)	211(3)	121(0)	121(0)	221(1)	1112(0)	121(1)	211(3)	211(3)	122(0)	112(0)
Question 2	<p>相続人に関する次のアからオまでの各記述のうち、判例の趣旨に照らし誤っているものを組み合わせたものは、後記1から5までのうちどれか。(Which of the following combinations consists of statements that are incorrect in light of Supreme Court precedent?)</p> <p>ア. 被相続人の内縁の配偶者は、相続人となる。(A de facto spouse of the decedent becomes an heir.)</p> <p>イ. 被相続人が妻の懐胎中に死亡したときは、その後に出産した子は、相続人となる。(If the decedent dies while his wife is pregnant, the child subsequently born becomes an heir.)</p> <p>ウ. 被相続人Aと子Bが死亡し、その前後関係が不明な場合、Bの子Cは代襲相続する。(If both the decedent A and child B die and the order of death is unknown, B's child C succeeds by representation.)</p> <p>エ. 子Bが相続放棄した場合、Bの子Cは代襲相続する。(If child B renounces inheritance, B's child C succeeds by representation.)</p> <p>オ. 遺言書を破棄しても不当目的がなければ相続欠格に当たらない。(An heir who destroys a will without unjust intent does not become disqualified from inheritance.)</p>									
Model	Exam(GT)	Base(ZS)	Base(FS)	Base+V	JBE-QA	JBE-QA+V	Ours	Ours+V	MA(Same)	MA(Sep)
Output (Pts)	2(2)	3(0)	3(0)	2(2)	3(0)	3(0)	2(2)	2(2)	4(0)	1(0)
Question 3	<p>次のアからオまでの各記述を判例の立場に従って検討し、正しい場合には1を、誤っている場合には2を選びなさい。(Examine each of the following statements A-E according to judicial precedent; choose 1 if correct and 2 if incorrect.)</p> <p>ア. 甲は、宝くじの当せん金を得るため、外れた宝くじに印字された番号を当せん番号に改ざんした。この場合、甲に有印私文書変造罪が成立する。(For the purpose of obtaining lottery winnings, X altered the number printed on a losing lottery ticket to match the winning number. In this case, the crime of alteration of a private document with a seal is established.)</p> <p>イ. 甲は、事情を知らない乙に対し、偽造通貨を真正な通貨のように装って代金として交付し、乙から商品を購入した。この場合、甲に詐欺罪及び偽造通貨行使罪が成立し、両罪は観念的競合となる。(X handed counterfeit currency to Y, who was unaware of the falsity, as if it were genuine currency, and purchased goods from Y. In this case, fraud and uttering counterfeit currency are established, and the two crimes are in conceptual concurrence.)</p> <p>ウ. 甲は、乙から、乙の代わりにA大学の入学試験を受けてほしいと頼まれ、これを引き受け、乙に成り済まして同入学試験を受け、氏名欄に乙の氏名を記載し、乙名義で答案を作成した。この場合、甲に有印私文書偽造罪が成立する。(X agreed to take an entrance examination for University A on behalf of Y, impersonated Y, wrote Y's name, and prepared an answer sheet under Y's name. In this case, the crime of forgery of a private document with a seal is established.)</p> <p>エ. 甲は、行使の目的で、他人が振り出した額面100万円の小切手の金額欄に「0」を加え、額面1000万円の小切手に改ざんした。この場合、甲に有価証券偽造罪が成立する。(For the purpose of use, X added a zero to the amount field of a 1-million-yen check issued by another person, altering it into a 10-million-yen check. In this case, the crime of forgery of a valuable security is established.)</p> <p>オ. 甲は、乙から金銭の借入れとして1万円札10枚を受け取った際、それらの中に偽造の1万円札が含まれていることに気付かず、その後、偽造の1万円札の存在に気付いたが、行使の目的でそのまま保持した。この場合、甲に偽造通貨取得罪は成立しない。(X received ten 10,000-yen bills as a loan from Y without noticing that one was counterfeit; later noticing the counterfeit bill, X continued to possess it for the purpose of use. In this case, the crime of acquisition of counterfeit currency is not established.)</p>									
Model	Exam(GT)	Base(ZS)	Base(FS)	Base+V	JBE-QA	JBE-QA+V	Ours	Ours+V	MA(Same)	MA(Sep)
Output (Pts)	22121(4)	21222(0)	21222(0)	21122(0)	21212(0)	21211(0)	21121(2)	22121(4)	21222(0)	21212(0)

- **Knowledge Extraction Agent:** For each verified past question, the extraction agent abstracts generalizable legal principles from the question/answer pair. The agent is instructed to output only reusable criteria, conditions, or patterns in bullet-point form.
- **Final Reasoning Agent:** The reasoning agent receives the original test question together with the aggregated extracted knowledge and produces the final answer in the strict exam-required format. This agent is solely responsible for joint evaluation of all propositions and adherence to formatting constraints.

We evaluated two configurations of this architec-

ture. In the shared-model setting, all four agents were instantiated from the same model fine-tuned with our dataset as in the previous experiment, isolating the effect of role separation and interaction. In the independently fine-tuned setting, each agent was fine-tuned separately on the same training data but with role-specific prompts, with the goal of increasing functional specialization and diversity. Prompts for each agent are shown in Table 5. Answer format instruction is identical as the previous experiment.

5.2 Results & Analysis

As shown in Table 3, both configurations performed substantially worse than a single fine-tuned

Table 5: Role-specific prompts used in the multi-agent pipeline.

Agent	Prompt (English translation in parentheses)
Retriever	以下の問題に関連すると考えられる過去問とその回答を選択せよ。選択の基準は、扱われている法分野、論点、条文、または判例の種類が共通しているかどうかである。最大で数問まで選んでよい。 (Select past exam questions and answers that you consider relevant to the following problem. Relevance should be judged based on shared legal domain, issues, statutes, or types of precedents. You may select up to a few questions.)
Verifier	以下の問題に対して参考になる過去問と回答のみを選別してください。 (Select only the past exam questions and answers that are relevant to the following problem.)
Extractor	以下の問題と正解から、将来の類似問題に使える一般化可能な法的知識を抽出せよ。 (Extract generalizable legal knowledge from the following question and answer that can be reused for future similar problems.)
Reasoner	以下は関連する法的知識である。上記を踏まえて、次の問題に答えよ。 (Below is relevant legal knowledge. Based on it, answer the following question.)

model. The shared-model multi-agent system achieved an average score of 75.7 points, while the independently fine-tuned multi-agent system further degraded performance to 71.0 points, both of which are substantially below the passing score. These results fall far below the single-model approach, despite significantly increased inference complexity.

Our results highlight that multi-agent systems do not automatically outperform strong single-model baselines. On the contrary, in tightly constrained tasks such as the Japanese bar examination, distributing reasoning across agents can be harmful, as errors introduced by individual agents tend to propagate and compound throughout the pipeline. In particular, abstraction and verification stages may introduce subtle inconsistencies that the final reasoning agent must reconcile under strict formatting and joint-consistency constraints.

We further find that increasing agent diversity through independent fine-tuning does not improve performance and instead exacerbates coordination failures. While independently trained agents exhibit greater surface-level variation, they lack a shared representational space that would allow their outputs to be reliably integrated. In contrast, shared representations appear to be crucial for effective agentic behavior in multi-agent setting, especially in settings where success depends on maintaining global consistency across multiple interdependent propositions.

6 Conclusion & Future Work

In this paper, we presented the first large language model system to achieve a passing score on the

Japanese bar examination when evaluated under its original question format and official grading scale. Our results demonstrate that preserving the exam’s multi-proposition structure during dataset construction and fine-tuning is essential. A single fine-tuned model, augmented with lightweight self-verification, is sufficient to meet the passing threshold without resorting to question decomposition or external supervision, despite a relatively small size of training data.

In contrast, approaches based on decomposed propositions or multi-agent deliberation fail to achieve comparable performance when assessed under realistic exam conditions, despite increased architectural complexity. These findings indicate that success on the Japanese bar examination depends less on distributing reasoning across components than on maintaining global consistency over tightly coupled propositions. The additional gains from self-verification further underscore that many remaining errors arise from coherence failures rather than missing legal knowledge. Overall, both format-specific fine-tuning and self-verification act as effective catalysts to extract the latent knowledge present in the model that would not be exploited otherwise.

As future work, we aim to extend this approach to the free-response (論文式) portion of the exam, which requires structured legal argumentation rather than discrete answer selection. More broadly, our results highlight the importance of evaluating legal reasoning systems under authentic task formats and caution against drawing conclusions from benchmarks that substantially simplify the structure and evaluation criteria.

565 Limitations

566 This work focuses exclusively on the multiple-
567 choice (短答式) portion of the Japanese bar ex-
568 amination and does not address the free-response
569 (論文式) component, which requires structured leg-
570 al argumentation, citation control, and extended
571 reasoning. As such, our results should not be inter-
572 preted as demonstrating comprehensive legal rea-
573 soning ability or qualification-level competence. In
574 addition, although our dataset faithfully preserves
575 the original exam format, its size is relatively small
576 compared to large-scale legal benchmarks, and per-
577 formance gains rely on the presence of substantial
578 prior knowledge in the underlying base model. Our
579 approach therefore assumes access to a strong pre-
580 trained language model and may not generalize to
581 weaker or domain-mismatched models. Finally,
582 while self-verification improves robustness under
583 the exam’s grading scheme, it does not guarantee
584 correctness in cases where the model’s internal leg-
585 al knowledge is itself incorrect or outdated.

586 Ethical Statement

587 This study does not involve human subjects, per-
588 sonal data, or sensitive private information. All
589 exam questions used in our dataset are publicly
590 available past questions from the Japanese bar ex-
591 amination, used solely for research and evaluation
592 purposes. Our results should not be construed as
593 endorsing the use of language models as a substi-
594 tute for formal legal education, professional train-
595 ing, or legal advice. In particular, passing or near-
596 passing performance on an exam-style benchmark
597 does not imply real-world legal competence or eth-
598 ical judgment. We emphasize that any deployment
599 of such models in legal contexts must be accom-
600 panied by appropriate human oversight and clear
601 disclosure of limitations. We caution against over-
602 interpreting benchmark success without evaluation
603 under authentic task formats and grading schemes,
604 a concern that our work explicitly aims to highlight
605 rather than exacerbate.

606 References

607 Zhihan Cao, Fumihito Nishino, Hiroaki Yamada,
608 Nguyen Ha Thanh, Yusuke Miyao, and Ken Satoh.
609 2025. *Jbe-qa: Japanese bar exam qa dataset for as-*
610 *sessing legal domain knowledge.*

611 Odysseas S. Chlapanis, Dimitrios Galanis, and Ion An-
612 droustopoulos. 2024. *Lar-echr: A new legal argu-*

ment reasoning task and dataset for cases of the euro- 613
pean court of human rights. *ArXiv*, abs/2410.13352. 614

Ahmed El-Kishky, Alexander Wei, Andre Saraiva, Bo- 615
rys Minaev, Daniel Selsam, David Dohan, Francis 616
Song, Hunter Lightman, Ignasi Clavera, Jakub W. Pa- 617
chocki, Jerry Tworek, Lorenz Kuhn, Lukasz Kaiser, 618
Mark Chen, Max Schwarzer, Mostafa Rohaninejad, 619
Nat McAleese, o3 contributors, Oleg Murk, and 5 620
others. 2025. *Competitive programming with large*
reasoning models. *ArXiv*, abs/2502.06807. 621
622

Zhiwei Fei, Xiaoyu Shen, D. Zhu, Fengzhe Zhou, Zhuo 623
Han, Songyang Zhang, Kai Chen, Zongwen Shen, 624
and Jidong Ge. 2023. *Lawbench: Benchmarking*
legal knowledge of large language models. *ArXiv*,
abs/2309.16289. 625
626
627

Neel Guha, Julian Nyarko, Daniel E. Ho, Christo- 628
pher Ré, Adam Chilton, Aditya Narayana, Alex 629
Chohlas-Wood, Austin M. K. Peters, Brandon Wal- 630
don, Daniel N. Rockmore, Diego A. Zambrano, 631
Dmitry Talisman, Enam Hoque, Faiz Surani, Frank 632
Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai 633
Porat, Jason Hegland, and 21 others. 2023. *Legal-*
bench: A collaboratively built benchmark for measur-
ing legal reasoning in large language models. *ArXiv*,
abs/2308.11462. 634
635
636
637

Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and 638
Sunghun Kim. 2024. *A survey on large language*
models for code generation. *ACM Transactions on*
Software Engineering and Methodology. 639
640
641

Daniel Martin Katz, Michael James Bommarito, Shang 642
Gao, and Pablo Arredondo. 2024. *Gpt-4 passes*
the bar exam. *Philosophical transactions. Series*
A, Mathematical, physical, and engineering sciences,
382. 643
644
645
646

Kimyeeun Kimyeeun, Choi Youngrok, Eunkyung Choi, 647
Jinhwan Choi, Hai Jin Park, and Wonseok Hwang.
2024. *Developing a pragmatic benchmark for as-*
sessing korean legal language understanding in large
language models. In *Conference on Empirical Meth-*
ods in Natural Language Processing. 648
649
650
651
652

Jens Lehmann, Antonello Meloni, Enrico Motta, 653
Francesco Osborne, Diego Reforgiato Recupero, An- 654
gelo Salatino, Sahar Vahdati, TU ScaDS.AI, Dresden,
and De. 2024. *Large language models for scientific*
question answering: An extensive analysis of the
sciq benchmark. In *Extended Semantic Web Confer-*
ence. 655
656
657
658
659

Haitao Li, You Chen, Qingyao Ai, Yueyue Wu, Ruizhe 660
Zhang, and Yiqun Liu. 2024. *Lexeval: A compre-*
hensive chinese legal benchmark for evaluating large
language models. *ArXiv*, abs/2409.20288. 661
662
663

Yixin Liu, Alexander R. Fabbri, Pengfei Liu, 664
Dragomir R. Radev, and Arman Cohan. 2023. *On*
learning to summarize with large language models
as references. In *North American Chapter of the*
Association for Computational Linguistics. 665
666
667
668

669 Hoang-Trung Nguyen, Tan-Minh Nguyen, Xuan-Bach
670 Le, Tuan-Kiet Le, Khanh-Huyen Nguyen, Ha Thanh
671 Nguyen, Thi-Hai-Yen Vuong, and Le-Minh Nguyen.
672 2025. [Nowj@coliee 2025: A multi-stage framework
673 integrating embedding models and large language
674 models for legal retrieval and entailment.](#) *ArXiv*,
675 [abs/2509.08025](#).

676 OpenAI. 2025. Introducing gpt-4.1 in the api. <https://openai.com/index/gpt-4-1/>.
677

678 Aditya Patwardhan, Vivek Vaidya, and Ashish Kundu.
679 2024. [Automated consistency analysis of llms.](#) *2024
680 IEEE 6th International Conference on Trust, Privacy
681 and Security in Intelligent Systems, and Applications
682 (TPS-ISA)*, pages 118–127.

683 Matthew Renze and Erhan Guven. 2024. [Self-reflection
684 in large language model agents: Effects on problem-
685 solving performance.](#) *2024 2nd International Con-
686 ference on Foundation and Large Language Models
687 (FLLM)*, pages 516–525.

688 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Jun-
689 Mei Song, Mingchuan Zhang, Y. K. Li, Yu Wu, and
690 Daya Guo. 2024. [Deepseekmath: Pushing the limits
691 of mathematical reasoning in open language models.](#)
692 *ArXiv*, [abs/2402.03300](#).

693 Jingyun Sun, Chengxiao Dai, Zhongze Luo, Yangbo
694 Chang, and Yang Li. 2024. [Lawluo: A multi-agent
695 collaborative framework for multi-round chinese le-
696 gal consultation.](#)

697 Nguyen Ha Thanh, Phuong Minh Nguyen, Thi-Hai-Yen
698 Vuong, Quan Minh Bui, Chau Nguyen, Binh Dang,
699 Vu Tran, Minh Le Nguyen, and Ken Satoh. 2021.
700 [Jnlp team: Deep learning approaches for legal pro-
701 cessing tasks in coliee 2021.](#) *ArXiv*, [abs/2106.13405](#).

702 Nguyen Ha Thanh, Hai-Yen Thi Vuong, Phuong Minh
703 Nguyen, Binh Dang, Quan Minh Bui, Sinh Trong
704 Vu, Chau Nguyen, Vu Tran, Ken Satoh, and Minh Le
705 Nguyen. 2020. [Jnlp team: Deep learning for legal
706 processing in coliee 2020.](#) *ArXiv*, [abs/2011.08071](#).

707 An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,
708 Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-
709 hong Tu, Jingren Zhou, Junyang Lin, Keming Lu,
710 Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang
711 Ren, and Zhenru Zhang. 2024. [Qwen2.5-math techni-
712 cal report: Toward mathematical expert model via
713 self-improvement.](#) *ArXiv*, [abs/2409.12122](#).

714 Murong Yue. 2025. [A survey of large language
715 model agents for question answering.](#) *ArXiv*,
716 [abs/2503.19213](#).

717 Li Zhang and Kevin Ashley. 2025. [Mitigating manipu-
718 lation and enhancing persuasion: A reflective multi-
719 agent approach for legal argument generation.](#) *ArXiv*,
720 [abs/2506.02992](#).

721 Qizheng Zhang, Changran Hu, Shubhangi Upasani,
722 Boyuan Ma, Fenglu Hong, Vamsidhar Reddy Kama-
723 nuru, Jay Rainton, Chen Wu, Mengmeng Ji, Hanchen
Li, Urmish Thakker, James Zou, and Kunle Oluko-
tun. 2025. [Agentic context engineering: Evolving
contexts for self-improving language models.](#) *ArXiv*,
[abs/2510.04618](#).