

PolicyLLM: Towards Excellent Comprehension of Public Policy for Large Language Models

Anonymous ACL submission

Abstract

Large Language Models (LLMs) are increasingly integrated into real-world decision-making, including in the domain of public policy. Yet, their ability to comprehend and reason about policy-related content remains under-explored. To fill this gap, we present *Policy-Bench*, the first large-scale [cross-system benchmark \(US-China\)](#) evaluating policy comprehension, comprising 21K cases across a broad spectrum of policy areas, capturing the diversity and complexity of real-world governance. Following Bloom’s taxonomy, the benchmark assesses three core capabilities: (1) **Memorization**: factual recall of policy knowledge, (2) **Understanding**: conceptual and contextual reasoning, and (3) **Application**: problem-solving in real-life policy scenarios. Building on this benchmark, we further propose *PolicyMoE*, a domain-specialized Mixture-of-Experts (MoE) model with expert modules aligned to each cognitive level. The proposed models demonstrate stronger performance on application-oriented policy tasks than on memorization or conceptual understanding, and yields the highest accuracy on structured reasoning tasks. Our results reveal key limitations of current LLMs in policy understanding and suggest paths toward more reliable, policy-focused models.

1 Introduction

In recent years, Large Language Models (LLMs) (Vaswani et al., 2017; Touvron et al., 2023; Achiam et al., 2023) have achieved remarkable progress, demonstrating intelligent and superior performance in a wide range of natural language processing (NLP) tasks, including machine translation (Zhu et al., 2024), code generation (Svyatkovskiy et al., 2020; Chen et al., 2021), and article writing (Yuan et al., 2022). In parallel with these advancements, a growing body of research has focused on systematically benchmarking LLM capabilities across multiple cognitive dimensions, including language

understanding (Wang et al., 2024), reasoning (Xiang, 2023; Cobbe et al., 2021; Joshi et al., 2017), and knowledge acquisition (Yang et al., 2015).

Beyond traditional NLP tasks, LLMs are increasingly being deployed in high-stakes real-world decision-making contexts, such as education (Xiao et al., 2023), law (Fei et al., 2024; Guha et al., 2023; Zhou et al., 2024), healthcare (Tang et al., 2023) and public administration (Pesch, 2025). Among these, public policy stands out as particularly consequential: supporting policy analysis and generation requires not only factual knowledge, but also contextual reasoning and value-sensitive judgment (Hou et al., 2025). Missteps can have tangible social consequences—for example, a model that miscalculates rural funding allocations by relying on the wrong fiscal base may cause substantial under-allocation of resources. Ensuring that LLMs develop a reliable and nuanced understanding of policy content is therefore both a technical necessity and an ethical imperative.

Understanding and applying public policy presents unique challenges for LLMs. While the field’s interdisciplinary nature, contextual dependence, and linguistic complexity are well recognized, the central obstacles to advancing policy-aware AI can be framed as a three-tiered problem. 1) *The Evaluation Challenge: Lack of Rigorous Benchmarks*. There is currently no comprehensive benchmark to systematically assess the policy comprehension capabilities of LLMs. Without standardized evaluation frameworks, it is difficult to measure performance across skills ranging from factual recall to conceptual reasoning and practical application, hindering objective comparison and targeted improvement. 2) *The Diagnostic Challenge: Identifying Strengths and Weaknesses*. Aggregate metrics obscure where models succeed and fail. It remains unclear which cognitive abilities, policy domains, or linguistic contexts pose the greatest difficulties. Fine-grained diagnostic

analysis is therefore essential to pinpoint strengths and weaknesses. 3) *The Adaptation Challenge: Developing Specialized Models*. General-purpose LLMs often struggle with the distinct demands of policy tasks. A key challenge is how to adapt existing architectures to better handle the multifaceted requirements of policy analysis, thereby closing the gaps revealed through rigorous evaluation and diagnosis.

To rigorously evaluate the gap, we present *PolicyBench*, a [cross-system benchmark \(US-China\)](#) specifically designed to assess LLMs’ understanding of public policy in both China and the United States. *PolicyBench* encompasses a broad spectrum of policy domains and features meticulously crafted questions targeting three cognitive levels: memorization, understanding, and application. Through extensive experiments, we find that model performance improves steadily from memorization to application tasks, with LLMs showing particular strength in structured reasoning scenarios such as numerical calculation and scenario-based decision-making, while still facing challenges in abstract or ambiguous policy contexts and in handling Chinese policy texts. To further enhance LLMs’ policy-related reasoning, we propose *PolicyMoE*—a MoE model (Jacobs et al., 1991; Jordan and Jacobs, 1994) trained on policy-focused data (Kang et al., 2024). *PolicyMoE* integrates three specialized expert models, each excelling in distinct capabilities. Experimental results demonstrate that *PolicyMoE* significantly outperforms general-purpose LLMs on policy tasks.

Overall, our main contributions are as follows:

- ▷ We construct *PolicyBench*, a comprehensive [cross-system](#) benchmark for evaluating LLMs’ policy understanding across diverse domains—in both Chinese and US contexts.
- ▷ Through extensive experiments and human evaluation on *PolicyBench*, we uncover key findings on the strengths and limitations of LLMs in [cross-system](#) policy understanding.
- ▷ We propose *PolicyMoE*, an MoE model fine-tuned on *PolicyBench*, which achieves superior performance over strong baselines and underscores the potential of domain-adaptive pretraining for governance-related tasks.

2 The *PolicyBench*

In this section, we provide a detailed introduction to the design and construction principles of *Policy-*

Bench. To construct our benchmark, we focused on two of the world’s most significant yet distinct policy environments: mainland China (CN) and the United States federal government (US). This deliberate selection provides a high-contrast, [cross-system](#) testbed for evaluating an LLM’s core policy comprehension capabilities across different governance systems. While we acknowledge this does not encompass the full global policy frameworks, it establishes a critical and challenging baseline for this foundational area.

2.1 Policy Acquisition

In the process of policy collection, we initially gathered a broad set of Chinese and US policy documents and related materials. To ensure relevance and timeliness of the content, we applied a filtering process that removed outdated policies, duplicate entries, and documents not related to substantive policy content (e.g., *purely procedural notices or administrative logistics*), details in [Appendix A](#). After this curation step, we retained **721** Chinese policies and **1,890** supplementary Chinese policy materials (e.g., *official commentaries, media news, expert interviews, and public consultations*), as well as **603** US policies and **1,082** supplementary US materials:

- For Chinese policies, all documents were sourced exclusively from the Policy Document Repository of the State Council of China.¹ To categorize the policies, we first followed the organizational structure of the State Council, then refined it to better reflect the content distribution within the corpus. Based on this adapted structure, we grouped the policies into eight domains (e.g., *Public safety*). To retrieve relevant materials, we selected representative search terms—such as “Belt and Road Initiative” and “Double Reduction Policy”—based on “Hot Words” highlighted by major official media and platforms, see [Figure 6](#) for details. In addition, we collected supplementary materials including official interpretations, policy outcomes, and expert interviews from extended social media sources.
- For US policies: As there is no centralized repository for federal policies in the US, we collected policy documents from the official

¹<https://www.gov.cn/zhengce/zhengcewenjianku/>

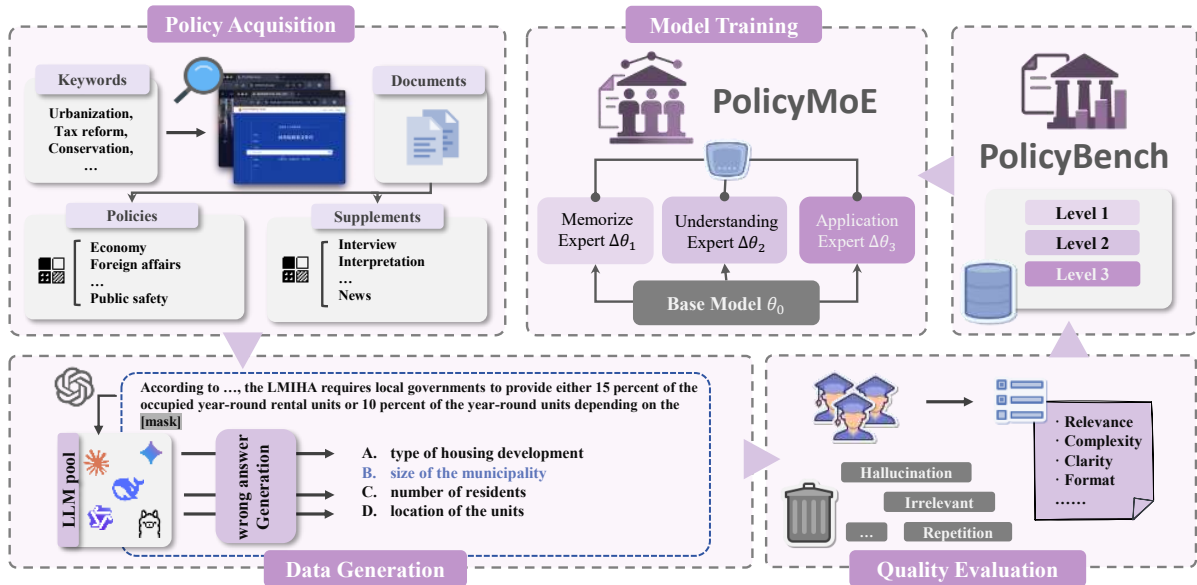


Figure 1: Three levels of evaluating LLM in *PolicyBench*.

websites of 12 US federal departments (Details in Table 5). Supplementary materials were gathered from authoritative news outlets such as *Reuters*, *Fox News*, etc.

All collected policies fall within the timeframe of 2000 to January 2025, with a primary focus on the most recent decade. All operations fully complied with ethical standards, and all data were lawfully sourced from open-access public databases.

Table 1: Task list of *PolicyBench* and respective numbers (“Mem” = Memorization, “Und” = Understanding, “App” = Application).

Level	ID	Task	Number	
			CN	US
Mem.	1-1	Article/Date Memorization	4,498	3,351
	1-2	Terminology Recognition	237	126
	1-3	Organization Identification	268	228
Und.	2-1	Idea Understanding	1,219	1,013
	2-2	Interest Understanding	1,145	996
	2-3	Institution Understanding	1,620	1,484
App.	3-1	Policy-Based Numerical Reasoning	918	452
	3-2	Scenario-Based Decision-Making	77	179
	3-3	Procedural/Institutional Implementation	320	391
	3-4	Policy Logic and Value Explanation	1,199	1,214

2.2 Dataset Curation

To facilitate fine-grained analysis of model capabilities, we categorize the benchmark into **10 task types**, each targeting a distinct subskill relevant to public policy comprehension. This taxonomy enables a comprehensive assessment of LLMs across a wide range of cognitive and policy domains.

The categorization is structured around a three-level hierarchy informed by Bloom’s Taxonomy of Educational Objectives (Krathwohl, 2002), a widely recognized framework in cognitive and educational psychology. Bloom’s taxonomy organizes cognitive skills from basic factual recall to deeper understanding and the application of knowledge in real-world contexts. Drawing on this structure, our benchmark defines three assessment levels:

- **Level 1: Memorization.** This level focuses on factual recall, such as memorizing publication dates, institutional actors, specific provisions, or technical terminology. Tasks at this level require minimal inference and primarily test a model’s ability to retrieve explicit information from policy documents.
- **Level 2: Understanding.** At this tier, we move beyond literal recall to examine a model’s capacity for conceptual understanding and contextualization. Guided by the **3I framework** from policy studies, which emphasizes Ideas, Interests, and Institutions (Hall and Taylor, 1996). Level 2 tasks probe how well a model can interpret underlying motivations, identify key stakeholders, and comprehend institutional logic within policies.
- **Level 3: Application.** The highest level assesses the model’s ability to apply policy knowledge in practical scenarios. Tasks at this tier require reasoning about hypothetical or real-world situations, evaluating implications, and suggesting appropriate actions based on policy content.

As summarized in Table 1, this hierarchical task design enables a fine-grained evaluation of pol-



Figure 2: Selected examples from PolicyBench spanning three levels and two languages.

229 icy comprehension across multiple cognitive levels.
 230 [Annotation guidelines and quality control procedures](#)
 231 [detailed in Appendix E, with formal task](#)
 232 [definitions are provided in Appendix F.](#)

233 2.3 Detail of Distractor Generation

234 To mitigate the bias that can arise when a single
 235 LLM generates all distractors, we employ a hetero-
 236 geneous *model pool* and harvest incorrect options
 237 in an iterative fashion. With details mathematically
 238 referred in Algorithm 1 and given a stem-answer
 239 pair $\langle q, a_{\text{gold}} \rangle$, we initially sample a model m
 240 from the pool and prompt it with the original question
 241 while *explicitly labelling* a_{gold} as a prior incorrect
 242 answer; the model is instructed to propose a new,
 243 plausible response that must differ from a_{gold} . If
 244 the resulting candidate d is neither redundant nor
 245 equal to the gold answer, it is added to the distractor
 246 set D . We iteratively resample additional models
 247 and repeat the procedure, supplying both the origi-
 248 nal question and the accumulated distractors, until
 249 the total number of distractors satisfies $|D| = k - 1$.
 250 Finally, we randomly permute $a_{\text{gold}} \cup D$ to form
 251 the list of k options.

252 3 The PolicyMoE

253 Specifically in policy domains, we propose *Policy-*
 254 *MoE*, a method that transforms the overall LLM
 255 into a compositional and modular system of experts
 256 with different expertise, drawing inspiration from
 257 (Kang et al., 2024). In this section, we present
 258 the details of *PolicyMoE*. By dividing the policy

259 domain expertise into three specialized expert mod-
 260 ules—Memory Policy, Understanding Policy, and
 261 Apply Policy—*PolicyMoE* creates a comprehen-
 262 sive system capable of handling diverse policy-
 263 related tasks with enhanced precision and effi-
 264 ciency.

265 3.1 Architecture Overview.

266 *PolicyMoE* follows the core principles of the MoE
 267 framework while specializing in policy domains:

- 268 • Expert Modules: Three dedicated expert models
 269 trained on specific policy-related capabilities:
 - 270 • Memory Expert: Specializes in recalling pol-
 271 icy facts, regulations, historical precedents,
 272 and exact policy language
 - 273 • Understanding Expert: Focuses on interpret-
 274 ing policy intent, analyzing implications, and
 275 explaining policy rationales.
 - 276 • Application Expert: Excels at applying poli-
 277 cies to specific scenarios, predicting outcomes,
 278 and recommending implementation strategies.
- 279 • Intelligent Router: A simple linear layer that is
 280 shared across all LoRA adapters, which efficiently
 281 analyzes input features to determine the most rele-
 282 vant policy domain expertise.

283 3.2 Constructing Expert Modules

284 **Specialization with LoRA:** Each expert module
 285 is specialized using Low-Rank Adaptation (LoRA)
 286 (Hu et al., 2022), which introduces lightweight,
 287 trainable parameters specific to each policy domain

Algorithm 1: Distractor Generation for Multiple-Choice Question Construction

Input: Question q , Correct Answer a_{gold} , LLM pool \mathcal{M} , Target number of choices $k = 4$
Output: Multiple-choice question with one correct answer and $k - 1$ distractors
Initialize distractor set $\mathcal{D} \leftarrow \emptyset$;
while $|\mathcal{D}| < k - 1$ **do**
 Sample a model $m \sim \mathcal{M}$;
 Construct prompt;;
 Include the question q ;
 Provide a_{gold} as a previous incorrect answer;
 Instruct model m to generate a new plausible answer that is **also incorrect**;
 Generate distractor candidate
 $d \leftarrow m(q, a_{\text{gold}}$ marked as incorrect);
 if $d \notin \mathcal{D}$ and $d \neq a_{\text{gold}}$ **then**
 Add d to \mathcal{D} ;
Randomly shuffle $a_{\text{gold}} \cup \mathcal{D}$ to form final options;
return $\{q, \text{options}, a_{\text{gold}}\}$;

288 while keeping the base LLM intact. The specialized
289 model $\Theta_{\text{spec},i}$ for each domain i is defined as:

$$290 \quad \Theta_{\text{spec},i} = \Theta_0 + \Delta\Theta_i$$

291 where $\Delta\Theta_i$ represents the LoRA parameters for
292 domain i . The forward pass for each expert module
293 is:

$$294 \quad h_i = \theta_0 x + \theta_{B_i} \theta_{A_i} x$$

295 Here, $\theta_{B_i} \in \mathbb{R}^{d \times \text{rank}}$ and $\theta_{A_i} \in \mathbb{R}^{\text{rank} \times k}$, with
296 $\text{rank} \ll \min(d, k)$.

297 3.3 Dynamic Integration of Experts

298 **Routing Mechanism:** A router module θ_r is intro-
299 duced to analyze each input token and route it to
300 the most appropriate expert module. The output h
301 for each input x is computed by combining the con-
302 tributions of the selected expert modules, weighted
303 by their relevance:

$$304 \quad h = \theta_0 x + \sum_{i=1}^n \alpha_i \Delta\theta_i x$$

305 where α represents the weights computed by the
306 router:

$$307 \quad \alpha = \text{top-k}(\text{softmax}(\theta_r x))$$

Algorithm 2: Inference Procedure of *PolyMoE*

Input: Input instruction x
Output: Final model response y
Step 1: Expert Modules Initialization
for each expert $i \in \{\text{Memory}, \text{Understanding}, \text{Application}\}$ **do**
 Load LoRA adapter $\Delta\Theta_i$ into base
 model Θ_0 ;
 Construct expert model
 $\Theta_{\text{spec},i} = \Theta_0 + \Delta\Theta_i$;
Step 2: Routing Decision
Compute routing score vector: $\mathbf{s} = \theta_r x$;
Compute expert weights: $\alpha = \text{softmax}(\mathbf{s})$;
Select top-1 expert index:
 $i^* = \arg \max(\alpha)$;
Step 3: Expert Inference
Use selected expert Θ_{spec,i^*} to generate
response:
 $y = \text{LM}_{\Theta_{\text{spec},i^*}}(x)$;
return y

The router is trained using the aggregated synthetic
308 data $D = \{D_i\}_{i=1}^n$ to learn optimal module selec-
309 tion for a given task: 310

$$311 \quad \mathcal{L}(\theta_r) = -\mathbb{E}_{(x,y) \sim D} [\log P_{\Theta_0}(y|x; \theta_r, \{\Delta\Theta_i\}_{i=1}^n)]$$

312 4 Experiments

313 4.1 Setup

314 **Models.** We select 11 representative state-
315 of-art models: GPT-4o (Achiam et al., 2023),
316 o4-mini, Claude-3.7-Sonnet (Anthropic,
317 2025), Claude-3.5-sonnet (Anthropic,
318 2024), Gemini-2.5-Flash (Google, 2025),
319 Gemini-2.0-Flash (Google, 2024), and the
320 open-source models: Gemma-3-27B (Team et al.,
321 2025), Qwen-QwQ-32B (QwQ, 2024), Llama-4
322 (Meta, 2025), Deepseek-V3 (Liu et al., 2024) and
323 Deepseek-R1 (Guo et al., 2025), details in Table 6.
324 **Scoring Mechanism.** *PolicyBench* adopts a
325 level-aware scoring framework tailored to ques-
326 tion type and format. **Levels 1–2** consist ex-
327 clusively of multiple-choice and true/false ques-
328 tions, which are evaluated using standard accu-
329 racy: $\text{Score} = \frac{\# \text{Correct Answers}}{\# \text{Total Questions}}$. **Level 3** includes both ob-
330 jective (multiple-choice, true/false) and subjective
331 (open-ended) questions. Open-ended responses are
332 scored on a 0–5 scale based on alignment with

Table 2: Performance (accuracy) of all models on different levels and regions. Gemini-2.5, Gemini-2.0, Claude-3.5 and Claude-3.7 denote Gemini-2.5-Flash, Gemini-2.0-Flash, Claude-3.5-Sonnet and Claude-3.7-sonnet respectively. Red and blue represent the highest and lowest scores in each row respectively.

Level	Region	GPT-4o	o4-mini	Gemini-2.5	Gemini-2.0	Claude-3.7	Claude-3.5	LLaMA-4	Gemma-3-27B	QwQ-32B	Deepseek-V3	Deepseek-R1
Level 1	CN	46.01%	45.93%	54.06%	47.87%	55.29%	53.77%	49.81%	41.75%	55.87%	48.61%	62.02%
	US	52.69%	54.90%	57.73%	53.71%	58.68%	58.76%	52.55%	49.91%	46.40%	50.12%	59.33%
Level 2	CN	56.34%	55.81%	60.57%	56.39%	60.47%	59.74%	56.56%	55.56%	59.79%	55.51%	62.92%
	US	63.40%	64.71%	64.91%	62.25%	68.23%	68.95%	61.17%	62.17%	57.71%	58.62%	65.37%
Level 3	CN	70.24%	79.49%	76.18%	73.80%	73.82%	72.83%	68.54%	71.51%	80.34%	72.33%	73.78%
	US	68.13%	77.00%	69.44%	66.55%	68.28%	68.47%	66.41%	68.37%	69.90%	69.39%	74.60%
AVERAGE		59.47%	62.97%	63.82%	60.10%	64.13%	63.75%	59.17%	58.21%	61.67%	59.10%	66.34%

Table 3: Average accuracy (%) of all models across Chinese and US (red and blue represent the highest and lowest in each row).

Region	GPT-4o	o4-mini	Gemini-2.5	Gemini-2.0	Claude-3.7	Claude-3.5	LLaMA-4	Gemma-3-27B	QwQ-32B	Deepseek-V3	Deepseek-R1
Chinese	57.53%	60.41%	63.60%	59.35%	63.19%	62.11%	58.30%	56.27%	65.33%	58.82%	62.24%
US	61.41%	65.54%	64.03%	60.84%	65.06%	65.39%	60.04%	60.15%	58.00%	59.38%	66.43%
Average	59.47%	62.98%	63.82%	60.10%	64.13%	63.75%	59.17%	58.21%	61.67%	59.10%	64.34%

reference answers. To ensure evaluation consistency, we adopt the **LLM-as-a-Judge** (Zheng et al., 2023): for each open-ended question, **two models are randomly sampled** from a pool of four state-of-the-art LLMs: o4-mini, gemini-2.5-flash, claude-3.7-sonnet, and Deepseek-R1, to serve as automated graders. Each grader compares the response to a reference answer and assigns a score according to predefined criteria. The final score for that question is computed as the average of the two model scores. **We analyze the potential bias in Appendix G.**

The overall Level 3 score is calculated as a weighted average across all question types:

$$\text{Score} = \frac{S_{mc} + S_{tf} + S_{oe}}{T_{mc} + T_{tf} + 5 \times T_{oe}},$$

where S_{mc} , S_{tf} , S_{oe} denote the cumulative scores for multiple-choice, true/false, and open-ended questions, respectively, and T_{mc} , T_{tf} , T_{oe} represent the corresponding question counts. The weighting reflects the maximum possible score (5) for open-ended responses. In addition, we conducted some experiments to demonstrate the robustness of LLM-as-a-Judge in Appendix H.

Training Setup. We initialize our MoE architecture using **Qwen2.5-7B-Instruct** as the base model, using bfloat16 precision. Expert modules are fine-tuned using LoRA (Hu et al., 2022) with a rank of 16, scaling factor $\alpha = 32$, and dropout rate of 0.05.

Training is conducted in two stages: expert training for 3 epochs and router training for 4 epochs. We use a batch size of 4 per device with a gradient accumulation step of 4, resulting in an effective batch size of 16. The learning rate is set to 5×10^{-5} throughout. The *PolicyBench* is partitioned into an 80/20 split for training and testing. To prevent data leakage, we perform a grouped split by policy, ensuring all questions from the same source document are kept in the same set. See subsection B.2 for more details.

4.2 Main Results

Performance improves from memorization to application. As shown in Table 2, models exhibit progressively higher accuracy from **Level 1** (memorization) to **Level 3** (application) across both Chinese and US settings. For instance, in the Chinese subset, average model scores range from **41.8%–62.0%** on Level 1 to **70.2%–80.3%** on Level 3. A similar trend is observed in the US subset. We posit that this phenomenon stems from the distinct capabilities emphasized during the different stages of LLM training. Level 1 tasks demand high-fidelity recall of specific facts (e.g., policy dates, exact terminology), which is a function of knowledge acquired during **pre-training**. While vast, this knowledge is stored implicitly, and its precise retrieval can be unreliable. In contrast, Level 2 (Understanding) and Level 3 (Application) tasks heavily rely on structured reasoning, contextual interpretation, and problem-solving—skills that

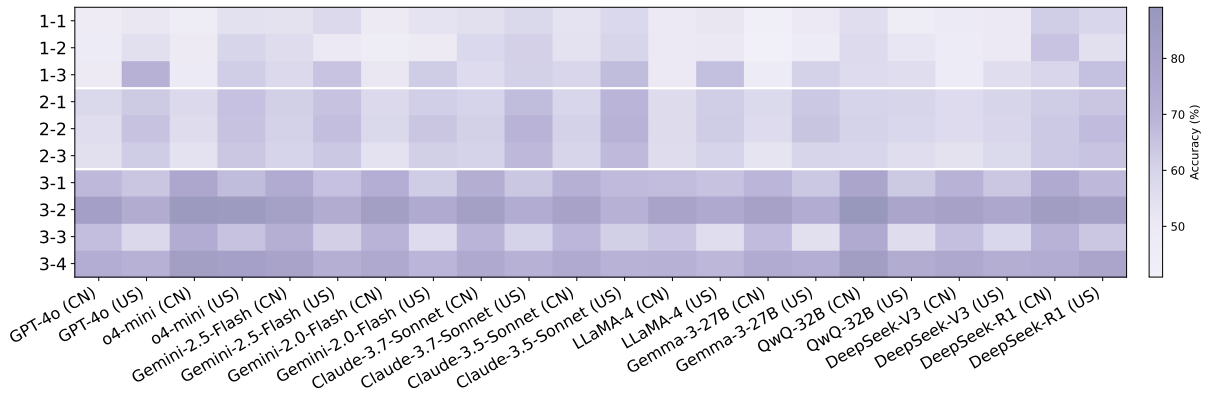


Figure 3: Model performance in 10 subtasks (ID and the specific task are shown in Table 1).

are explicitly and extensively honed during **post-training** (e.g., *instruction tuning, RLHF* (Ouyang et al., 2022)). Therefore, the models’ superior performance on these more complex tasks likely reflects a stronger alignment with the generalizable reasoning abilities optimized during fine-tuning, rather than a deeper mastery of the policy domain itself. *Among all evaluated models*, DeepSeek-R1 achieves the highest overall accuracy at **66.34%**, making it a strong candidate for practitioners seeking a general-purpose model for policy-related applications.

Models excel at structured reasoning tasks but falter on abstract concepts. As shown in the task-wise heatmap (Figure 3), models consistently achieve higher accuracy on specific application-oriented tasks, particularly **Policy-Based Numerical Reasoning** and **Scenario-Based Decision-Making**. For many models, accuracy in these categories exceeds **75%**, with some surpassing **80%**. These tasks typically involve concrete conditions, rule-based logic, or everyday reasoning scenarios—formats that closely align with the pre-training and instruction-following capabilities of large language models. In contrast, accuracy remains relatively lower on more abstract or ambiguous tasks such as **Ideas** and **Institutions**, which require understanding latent policy concepts or institutional relationships. This suggests that LLMs are better equipped to handle tasks with clear logical structures than those requiring interpretive or conceptual comprehension.

Models consistently perform better on US policy questions than Chinese ones. As shown in Table 3, most models achieve higher accuracy on US policy questions than on their Chinese counterparts. The overall average improves from **61.02%** (CN) to **62.39%** (US), with models like

o4-mini rising from **60.41%** to **65.54%**, and Claude-3.5-Sonnet from **62.11%** to **65.39%**. An exception is QwQ-32B, which performs notably better in Chinese (**65.33%**) than in English (**58.00%**). This overall trend may be attributed to the dominance of English in pretraining corpora, as well as the higher syntactic and semantic density of Chinese policy texts. These results highlight the need for more robust cross-lingual policy understanding in LLMs. Further more, we also conduct an error analysis (detailed in Appendix D).

Table 4: Qwen2.5-7B-Instruct performance across levels and regions before and after training (%)

Level	Region	Original	Training	Δ
Level 1	CN	36.85%	41.83%	$\uparrow 13.51\%$
	US	23.35%	35.43%	$\uparrow 51.73\%$
Level 2	CN	45.68%	47.02%	$\uparrow 2.93\%$
	US	42.31%	42.78%	$\uparrow 1.11\%$
Level 3	CN	64.73%	69.12%	$\uparrow 6.78\%$
	US	46.65%	57.48%	$\uparrow 23.22\%$

4.3 Results of PolicyMoE

Table 4 reports the performance of our model before and after fine-tuning with the *PolicyMoE* framework. The goal of this experiment is not to pursue state-of-the-art results with a moderately sized base model (Qwen2.5-7B-Instruct), but to demonstrate the efficacy of our expert specialization approach. Accordingly, we emphasize the **relative improvements** achieved. Notably, the fine-tuned 7B model not only shows substantial gains but also outperforms several larger baselines in Table 2, underscoring the effectiveness of domain-specific adaptation.

Across task levels, performance improves consistently after fine-tuning. The largest gain appears in US Level 1 tasks, where accuracy rises from

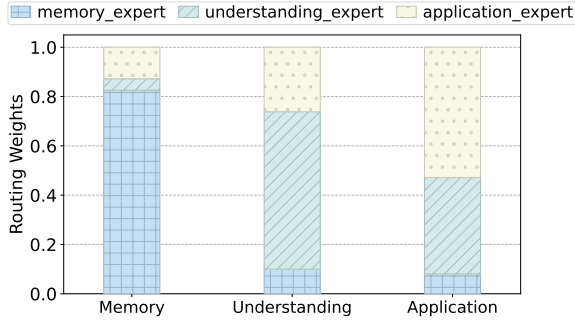


Figure 4: Routing distributions over three experts for each level.

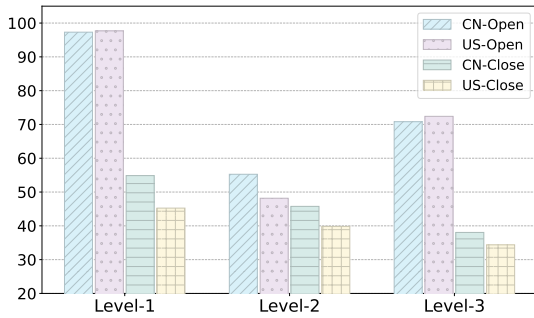


Figure 5: Human accuracy on three levels (%).

23.35% to 35.43%—a relative improvement of over 50%. China also records a 13.51% gain at the same level, highlighting the benefit of injecting structured domain knowledge. Improvements on Level 2, which emphasizes policy comprehension, are more modest (2.93% for China and 1.11% for the US), suggesting that higher-level reasoning is less sensitive to task-specific fine-tuning and may require advanced strategies such as chain-of-thought prompting (Wei et al., 2022) or richer supervision. By contrast, Level 3 tasks show more pronounced gains, with the US domain achieving a 23.22% relative improvement, indicating that the Application Expert effectively supports contextual reasoning and scenario-based decision-making.

Overall, *PolicyMoE* delivers clear benefits for both factual recall and applied reasoning. While improvements in abstract comprehension remain limited, the results point to a promising path toward specialized, capable models without relying on prohibitively large architectures. We also compare *PolicyMoE* with standard LoRA (Hu et al., 2022) in Appendix J.

4.4 *PolicyMoE* Router Analysis

To analyze the behavior of the router module, we randomly sample 10 questions from each cognitive level (Level 1–3) and compute the average expert

weights assigned by the router. This allows us to observe how the model distributes attention across three experts in response to different types of tasks.

As shown in Figure 4, router weight patterns show clear specialization for factual tasks, with memory weights peaking over 80% when the memory expert is selected. In contrast, understanding and application selections result in more distributed weights, reflecting the multi-dimensional nature of these tasks and the need for shared reasoning across modules.

4.5 Human Performance

To contextualize LLM performance, we established human baselines with 12 university students from the United States and China, distinct from the annotators. Participants, proficient in English or Mandarin but without policy expertise, each answered 100 randomly sampled *PolicyBench* questions under two conditions: **open-book** (with access to policy texts) and **closed-book** (relying on prior knowledge). As shown in Figure 5, the results indicate that:

- **Open- vs. closed-book gaps distinguish memory from reasoning.** The large gap at Level 1 reflects reliance on factual recall, while the minimal difference at Level 2 suggests a shift toward reasoning-intensive challenges.
- **Cross-linguistic consistency supports benchmark validity.** Comparable performance across languages within each setting indicates that task difficulty is not driven by language differences.
- **Non-monotonic accuracy across levels.** Accuracy drops at Level 2 and partially recovers at Level 3, suggesting that abstract reasoning without direct retrieval is most challenging, whereas Level 3 allows participants to combine reasoning with general knowledge.

5 Conclusion

We introduce *PolicyBench*, a **cross-system benchmark (US-China)** assessing LLM comprehension of public policy, which reveals that models perform well on recall and application tasks, they struggle with understanding questions involving policy intent and institutional reasoning. To bridge this gap, we propose *PolicyMoE*, a domain-specialized MoE model that achieves improved performance. Our findings underscore the need for targeted adaptation to support real-world policy analysis.

533 Limitations

534 **Geographic and systemic scope.** *PolicyBench*
535 currently covers only China and the United States.
536 While these two cases offer a strong contrast across
537 governance systems, they cannot fully represent the
538 diversity of global policy environments. Extending
539 to additional regions would improve generalizabil-
540 ity and cross-cultural robustness.

541 **Task diversity.** The benchmark mainly relies on
542 multiple-choice and true/false formats, with limited
543 coverage of open-ended tasks. Real-world policy
544 analysis, however, involves greater ambiguity, nu-
545 ance, and value-sensitive reasoning than what struc-
546 tured formats can capture. Expanding task types
547 would better reflect such challenges.

548 **Model adaptation and architecture.** *PolicyMoE*
549 shows clear improvements in factual recall and
550 applied reasoning, but its gains in abstract com-
551 prehension remain limited. Moreover, the current
552 router design can only select one expert per query,
553 whereas increasingly complex tasks may require
554 activating multiple experts simultaneously. Fu-
555 ture work could explore reasoning-focused models,
556 advanced prompting strategies, and more flexible
557 routing mechanisms that allocate experts based on
558 learned weight distributions.

559 Ethics Statement

560 All data used in this study were collected from pub-
561 licly accessible platforms in compliance with ethi-
562 cal and legal standards. No proprietary or private
563 materials were included. For the human evaluation
564 component, all participants were recruited on a vol-
565 untary basis and provided informed consent prior
566 to their involvement. They were clearly informed
567 of the research purpose and their rights, includ-
568 ing the right to withdraw at any stage. Finally, all
569 content generated by large language models was
570 carefully reviewed to ensure that it did not contain
571 sensitive, harmful, or inappropriate material. The
572 study adheres to responsible AI research practices
573 and aims to contribute to the safe and transparent
574 development of policy-focused benchmarks and
575 models.

576 References

577 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
578 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
579 Diogo Almeida, Janko Altenschmidt, Sam Altman,
580 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-
581 cal report. *arXiv preprint arXiv:2303.08774*.

Anthropic. 2024. Anthropic claude-3.5-sonnet. 582
<https://www.anthropic.com/news/claude-3-5-sonnet>. 583
584

Anthropic. 2025. Anthropic claude-3.7-sonnet. 585
<https://www.anthropic.com/claude/sonnet>. 586

Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi, 587
Amanpreet Singh, Joseph Chee Chang, Kyle Lo, 588
Luca Soldaini, Sergey Feldman, Mike D’arcy, 589
David Wadden, Matt Latzke, Minyang Tian, Pan Ji, 590
Shengyan Liu, Hao Tong, Bohao Wu, Yanyu Xiong, 591
Luke Zettlemoyer, and 6 others. 2024. *Openscholar: 592*
Synthesizing scientific literature with retrieval- 593
augmented lms. *Preprint*, arXiv:2411.14199. 594

Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, 595
and Sung Ju Hwang. 2025. *Researchagent: Iterative 596*
research idea generation over scientific literature with 597
large language models. *Preprint*, arXiv:2404.07738. 598

Han Bao, Yue Huang, Yanbo Wang, Jiayi Ye, Xi- 599
angqi Wang, Xiuying Chen, Yue Zhao, Tianyi 600
Zhou, Mohamed Elhoseiny, and Xiangliang Zhang. 601
2024. *Autobench-v: Can large vision-language 602*
models benchmark themselves? *arXiv preprint 603*
arXiv:2410.21259. 604

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, 605
Henrique Ponde De Oliveira Pinto, Jared Kaplan, 606
Harri Edwards, Yuri Burda, Nicholas Joseph, Greg 607
Brockman, and 1 others. 2021. *Evaluating large 608*
language models trained on code. *arXiv preprint 609*
arXiv:2107.03374. 610

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 611
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 612
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 613
Nakano, and 1 others. 2021. *Training verifiers 614*
to solve math word problems. *arXiv preprint 615*
arXiv:2110.14168. 616

Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, 617
Dominic Culver, Rui Melo, Caio Corro, Andre F. T. 618
Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia 619
Morgado, and Michael Desa. 2024. *Saullm-7b: A 620*
pioneering large language model for law. *Preprint*, 621
arXiv:2403.03883. 622

Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, 623
Zhuo Han, Alan Huang, Songyang Zhang, Kai Chen, 624
Zhixin Yin, Zongwen Shen, and 1 others. 2024. *Law- 625*
bench: Benchmarking legal knowledge of large lan- 626
guage models. In *Proceedings of the 2024 Confer- 627*
ence on Empirical Methods in Natural Language 628
Processing, pages 7933–7962. 629

Google. 2024. Google gemini-2-flash. 630
[https://deepmind.google/technologies/gemini/flash- 631](https://deepmind.google/technologies/gemini/flash-lite/)
lite/. 632

Google. 2025. Google gemini-2.5-flash. 633
[https://deepmind.google/technologies/gemini/flash/ 634](https://deepmind.google/technologies/gemini/flash/)

635	Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré,	Nikos Karacapilidis, Evangelos Kalampokis, Nikolaos	690
636	Adam Chilton, Alex Chohlas-Wood, Austin Peters,	Giarelis, and Charalampos Mastrokostas. 2024. Gen-	691
637	Brandon Waldon, Daniel Rockmore, Diego Zam-	erative ai and public deliberation: A framework	692
638	brano, and 1 others. 2023. Legalbench: A collab-	for llm-augmented digital democracy. <i>Proceedings</i>	693
639	oratively built benchmark for measuring legal reason-	<i>http://ceur-ws.org ISSN</i> , 1613:0073.	694
640	ing in large language models. <i>Advances in Neural</i>		
641	<i>Information Processing Systems</i> , 36:44123–44279.		
642	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao	David R Krathwohl. 2002. A revision of bloom’s taxon-	695
643	Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-	omy: An overview. <i>Theory into practice</i> , 41(4):212–	696
644	rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.	218.	697
645	Deepseek-r1: Incentivizing reasoning capability in		
646	llms via reinforcement learning. <i>arXiv preprint</i>	Yueqing Liang, Liangwei Yang, Chen Wang, Congying	698
647	<i>arXiv:2501.12948</i> .	Xia, Rui Meng, Xiong Xiao Xu, Haoran Wang, Ali	699
		Payani, and Kai Shu. 2025. Benchmarking llms for	700
		political science: A united nations perspective. <i>arXiv</i>	701
		<i>preprint arXiv:2502.14122</i> .	702
648	Peter A Hall and Rosemary CR Taylor. 1996. Political		
649	science and the three new institutionalisms. <i>Political</i>	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	703
650	<i>studies</i> , 44(5):936–957.	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	704
		Deng, Chenyu Zhang, Chong Ruan, and 1 others.	705
651	Yichen He, Guanhua Huang, Peiyuan Feng, Yuan Lin,	2024. Deepseek-v3 technical report. <i>arXiv preprint</i>	706
652	Yuchen Zhang, Hang Li, and Weinan E. 2025. Pasa:	<i>arXiv:2412.19437</i> .	707
653	An llm agent for comprehensive academic paper		
654	search . <i>Preprint</i> , arXiv:2501.10120.	Meta. 2025. Meta llama-4.	708
		https://www.llama.com/models/llama-4/ .	709
655	Ce Hou, Fan Zhang, Yong Li, Haifeng Li, Gengchen		
656	Mai, Yuhao Kang, Ling Yao, Wenhao Yu, Yao Yao,	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	710
657	Song Gao, and 1 others. 2025. Urban sensing in the	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	711
658	era of large language models. <i>The Innovation</i> , 6(1).	Sandhini Agarwal, Katarina Slama, Alex Ray, and 1	712
		others. 2022. Training language models to follow in-	713
659	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	structions with human feedback. <i>Advances in neural</i>	714
660	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,	<i>information processing systems</i> , 35:27730–27744.	715
661	Weizhu Chen, and 1 others. 2022. Lora: Low-rank		
662	adaptation of large language models. <i>ICLR</i> , 1(2):3.	Paulina Jo Pesch. 2025. Potentials and challenges of	716
		large language models (llms) in the context of admin-	717
663	Robert A Jacobs, Michael I Jordan, Steven J Nowlan,	istrative decision-making. <i>European Journal of Risk</i>	718
664	and Geoffrey E Hinton. 1991. Adaptive mixtures of	<i>Regulation</i> , pages 1–20.	719
665	local experts. <i>Neural computation</i> , 3(1):79–87.		
		QwQ. 2024. Qwen-qwq-32b.	720
666	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,	https://qwenlm.github.io/blog/qwq-32b/ .	721
667	Hanyi Fang, and Peter Szolovits. 2021. What disease		
668	does this patient have? a large-scale open domain	Mehrdad Safaei and Justin Longo. 2024. The end of the	722
669	question answering dataset from medical exams. <i>Ap-</i>	policy analyst? testing the capability of artificial intel-	723
670	<i>plied Sciences</i> , 11(14):6421.	ligence to generate plausible, persuasive, and useful	724
		policy analysis. <i>Digital Government: Research and</i>	725
671	Michael I Jordan and Robert A Jacobs. 1994. Hierarchi-	<i>Practice</i> , 5(1):1–35.	726
672	cal mixtures of experts and the em algorithm. <i>Neural</i>		
673	<i>computation</i> , 6(2):181–214.	Jaromir Savelka. 2023. Unlocking practical applica-	727
		tions in legal domain: Evaluation of gpt for zero-shot	728
674	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke	semantic annotation of legal texts . In <i>Proceedings of</i>	729
675	Zettlemoyer. 2017. Triviaqa: A large scale distantly	<i>the Nineteenth International Conference on Artificial</i>	730
676	supervised challenge dataset for reading comprehen-	<i>Intelligence and Law</i> , ICAIL 2023, page 447–451.	731
677	sion. In <i>Proceedings of the 55th Annual Meeting of</i>	ACM.	732
678	<i>the Association for Computational Linguistics (Vol-</i>		
679	<i>ume 1: Long Papers)</i> , pages 1601–1611.	Alexey Svyatkovskiy, Shao Kun Deng, Shengyu Fu,	733
		and Neel Sundaresan. 2020. Intellicode compose:	734
680	Jiaju Kang, Puyu Han, Tian Zhang, and Luqi Gong.	Code generation using transformer. In <i>Proceedings</i>	735
681	2025. Polycysimeval: A benchmark for evaluat-	<i>of the 28th ACM joint meeting on European software</i>	736
682	ing policy outcomes through agent-based simulation.	<i>engineering conference and symposium on the founda-</i>	737
683	<i>arXiv preprint arXiv:2502.07853</i> .	<i>tions of software engineering</i> , pages 1433–1443.	738
684	Junmo Kang, Leonid Karlinsky, Hongyin Luo, Zhen	Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and	739
685	Wang, Jacob Hansen, James Glass, David Cox,	Xia Hu. 2023. Does synthetic data generation of	740
686	Rameswar Panda, Rogerio Feris, and Alan Ritter.	llms help clinical text mining? <i>arXiv preprint</i>	741
687	2024. Self-moe: Towards compositional large lan-	<i>arXiv:2303.04360</i> .	742
688	guage models with self-specialized experts. <i>arXiv</i>		
689	<i>preprint arXiv:2406.12034</i> .		

743	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. Gemma 3 technical report. <i>arXiv preprint arXiv:2503.19786</i> .	<i>of the 18th workshop on innovative use of NLP for building educational applications (BEA 2023)</i> , pages 610–625.	800
744			801
745			802
746			
747			
748	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	Yi Yang, Wen-tau Yih, and Christopher Meek. 2015. Wikiqa: A challenge dataset for open-domain question answering. In <i>Proceedings of the 2015 conference on empirical methods in natural language processing</i> , pages 2013–2018.	803
749			804
750			805
751			806
752			807
753			
754	George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, Ioannis Partalas, Matthias Zschunke, Michael R Alvers, Dirk Weissenborn, Anastasia Krithara, Sergios Petridis, Dimitris Polychronopoulos, and 1 others. 2015. An overview of the bioasq large-scale biomedical semantic indexing and question answering competition. <i>BMC bioinformatics</i> , 16:1–28.	Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. 2022. Wordcraft: story writing with large language models. In <i>Proceedings of the 27th International Conference on Intelligent User Interfaces</i> , pages 841–852.	808
755			809
756			810
757			811
758			812
759			
760			
761			
762	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	Jerrold H Zar. 2005. Spearman rank correlation. <i>Encyclopedia of biostatistics</i> , 7.	813
763			814
764			
765			
766			
767	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In <i>International Conference on Learning Representations</i> .	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in Neural Information Processing Systems</i> , 36:46595–46623.	815
768			816
769			817
770			818
771			819
772	Dongsheng Wang, Natraj Raman, Mathieu Sibue, Zhiqiang Ma, Petr Babkin, Simerjot Kaur, Yulong Pei, Armineh Nourbakhsh, and Xiaomo Liu. 2023. Docllm: A layout-aware generative language model for multimodal document understanding. <i>arXiv preprint arXiv:2401.00908</i> .	Zhi Zhou, Jiang-Xin Shi, Peng-Xiao Song, Xiaowen Yang, Yi-Xuan Jin, Lan-Zhe Guo, and Yu-Feng Li. 2024. Lawgpt: A chinese legal knowledge-enhanced large language model. <i>CoRR</i> .	821
773			822
774			823
775			824
776			
777			
778	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In <i>The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2024. Multilingual machine translation with large language models: Empirical results and analysis. In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 2765–2781.	825
779			826
780			827
781			828
782			829
783			830
784			
785	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.		
786			
787			
788			
789			
790			
791	Alice Xiang. 2023. Beyond the imitation game: Collaborative benchmark for measuring and extrapolating the capabilities of language models [co-authored]: What is the tao? <i>Transactions on Machine Learning Research</i> .		
792			
793			
794			
795			
796	Changrong Xiao, Sean Xin Xu, Kunpeng Zhang, Yufang Wang, and Lei Xia. 2023. Evaluating reading comprehension exercises generated by llms: A showcase of chatgpt in education applications. In <i>Proceedings</i>		
797			
798			
799			

831 A Details of Data Collection

832 This section details the curation of our dataset,
833 outlining the collection methodology for both pol-
834 icy documents and their supplementary materials,
835 the tools utilized in the process, and the filtering
836 pipeline applied to ensure data quality.

837 A.1 Data Collection Tools and Process

838 Our data was collected primarily between January
839 2015 and March 2025 using a hybrid approach to
840 navigate the structural inconsistencies of official
841 government websites.

842 The majority of documents were collected manu-
843 ally from the portals listed in Table 5. This manual
844 approach was essential for bypassing complex site
845 navigation and security mechanisms, and for en-
846 suring the correct retrieval of policy documents
847 with associated attachments, which challenge stan-
848 dard web scrapers. For a small subset of highly
849 structured content, such as news archives, we em-
850 ployed targeted automation scripts using Python’s
851 **Selenium** library to assist with batch-downloading.

852 A.2 Data Filtering and Curation Pipeline

853 To ensure the quality and relevance of our final
854 dataset, all collected documents underwent a rigor-
855 ous multi-stage filtering and curation pipeline. The
856 objective was to create a clean, non-redundant, and
857 substantive corpus for constructing **PolicyBench**.
858 The pipeline consisted of the following sequential
859 steps:

- 860 1. **Duplicate Removal:** The first step was to
861 eliminate redundant files. We identified and
862 removed duplicates by checking for high simi-
863 larity in document titles and textual content.
864 Initially, documents with identical or near-
865 identical titles were flagged, after which their
866 content overlap was assessed. Documents
867 with a high degree of textual similarity were
868 considered duplicates, and all but the most
869 complete version were discarded.
- 870 2. **Substantive Content Filtering:** Next, we fil-
871 tered out documents that were not substantive
872 policy texts. A document was classified as
873 "non-substantive" and excluded if it met any
874 of the following criteria:
 - 875 • It was purely administrative or procedu-
876 ral (e.g., public meeting announcements,
877 personnel appointment notices, holiday
878 schedules).

- 879 • It was a table of contents, an index, or
880 a cover page without the corresponding
881 full document.
- 882 • The document’s title contained
883 keywords from a predefined ex-
884 clusion list, such as "Notice of
885 Public Hearing", "Personnel
886 Appointments", "Weekly Agenda", or
887 "Annual Report Summary".

- 888 3. **Temporal and Relevancy Filtering:** Finally,
889 we applied a filter to remove documents that
890 were considered "outdated" or irrelevant to the
891 contemporary policy landscape. A policy was
892 flagged and removed if:

- 893 • It was explicitly superseded by a more
894 recent version or subsequent legislation
895 from the same issuing authority.
- 896 • It was promulgated before the year 2000,
897 which we established as the historical
898 cutoff for our benchmark to maintain
899 modern relevance.

900 This structured pipeline allowed us to distill the
901 large volume of collected data into the high-quality,
902 curated set of 721 Chinese policies and 603 US
903 policies that form the foundation of **PolicyBench**.

904 B Details of Experiment Setting

905 B.1 Details of **PolicyBench**

906 **Level-1: Memorization Task.** Level-1 tasks are
907 designed to assess factual recall. We begin by us-
908 ing large language models to automatically gen-
909 erate cloze-style and true/false questions. Cloze
910 questions are created by masking factual elements
911 in policy texts—such as *dates*, *legal terms*, *orga-*
912 *nization names*, and *key definitions*—that reflect
913 domain-specific knowledge and span various pol-
914 icy areas. These cloze items are then transformed
915 into multiple-choice questions. To construct high-
916 quality distractors, we prompt multiple LLMs inde-
917 pendently to generate alternative options inspired
918 by (Bao et al., 2024). This multi-model strategy re-
919 duces single-model bias and increases the plausibil-
920 ity and diversity of distractors, thereby enhancing
921 the benchmark’s robustness.

922 **Level-2: Understanding Task.** Level-2 tasks
923 evaluate a model’s ability to comprehend the deeper
924 meaning and context of policy content. We prompt
925 LLMs to analyze policies using the 3I framework
926 from policy studies, which highlights three core

dimensions: *Ideas* (underlying beliefs and values), *Interests* (stakeholders involved), and *Institutions* (rules and structures guiding implementation) (Hall and Taylor, 1996). Based on these structured analyses, we generate question-answer pairs that probe a model’s understanding of policy motivations, actors, and institutional dynamics. The question generation process parallels that of Level-1: we first construct cloze-style prompts grounded in 3I insights, then convert them into multiple-choice format with distractors generated via multiple LLMs to improve quality and difficulty.

Level-3: Application Task. Level-3 tasks focus on practical reasoning and real-world contextual adaptation. To build these tasks, we draw on supplementary policy materials (*e.g.*, *official commentaries*, *media coverage*, *expert interviews*, and *public consultations*) to develop realistic scenarios where a policy might be applied. Based on these scenarios, we recruit students with relevant academic backgrounds to manually craft questions that require reasoning about a policy’s implications, suitability, or potential outcomes in novel contexts. This manual, context-driven approach ensures that Level-3 tasks closely mirror real-world decision-making challenges and reflect authentic policy discourse.

Expert-Led Question Design for Levels 2 and 3. To ensure high cognitive alignment and domain validity, the construction of **Level 2** (Understanding) and **Level 3** (Application) items was strictly led and executed by senior domain experts. The core writing team consisted of five senior Ph.D. candidates specializing in Public Policy, Law, and Computational Social Science. While three undergraduate assistants provided support for data formatting and preliminary cleaning, the substantive generation and validation of reasoning logic were exclusively performed by the senior doctoral researchers to ensure rigor.

Level 2 questions are designed to assess the model’s ability to comprehend policy intent, stakeholder interests, and institutional logic, following the 3I framework (Ideas, Interests, Institutions). These questions emphasize abstraction and interpretation rather than factual recall. **Level 3** questions focus on practical reasoning, including scenario-based decision-making, numerical calculations, and value-driven trade-offs, often grounded in real-world policy contexts.

All questions are constructed to be clear, faithful to source policies, and cognitively representative

of their designated levels (see Figure 2 for some examples). To ensure quality and consistency, the resulting items undergo a human evaluation process (for details, see Appendix E).

B.2 Details of PolicyMoE

Our MoE architecture is initialized using Qwen2.5-7B-Instruct as the base model with bfloat16 precision. Expert modules are fine-tuned using LoRA with rank $r = 16$, scaling factor $\alpha = 32$, and dropout rate of 0.05, targeting the attention and MLP projection layers. Expert training is conducted for 3 epochs with a per-device batch size of 4 and gradient accumulation of 4 steps, resulting in an effective batch size of 16. The learning rate is set to 5×10^{-5} with weight decay of 0.01. Router training follows for 2 epochs with batch size 8 and learning rate 1×10^{-4} , using a two-layer MLP architecture with LayerNorm and GELU activation. The dataset is split into 80% training and 20% testing, with maximum sequence lengths of 2048 tokens for expert training and 512 tokens for router training across the three domains: memory, comprehension, and application.

C Related Works

LLM Benchmarks in Social Science. The rapid advancements in LLMs have enabled their application in social science research. The ability to efficiently comprehend long textual inputs and generate human-like responses makes them favorable for text-intensive tasks, like academic paper searching (He et al., 2025) (Baek et al., 2025) (Asai et al., 2024), document analysis (Wang et al., 2023) and legal research (Colombo et al., 2024) (Savelka, 2023). Such ability underscores the need for tailored evaluation methods that measure their performance in these domains. Knowledge-intensive benchmarks, such as MMLU (Wang et al., 2024) and TriviaQA (Joshi et al., 2017), test models on a wide range of subjects, from STEM fields to humanities, shedding light on the recalling and reasoning ability of LLMs to handle complex queries in real-world scenarios.

Domain-specific benchmarks now probe specialized knowledge: BioASQ (Tsatsaronis et al., 2015) tests biomedical QA, MedQA (Jin et al., 2021) evaluates clinical reasoning. Legal NLP has advanced through LawBench (Fei et al., 2024) for Chinese statutory analysis and LegalBench (Guha et al., 2023) for Anglo-American jurisprudence

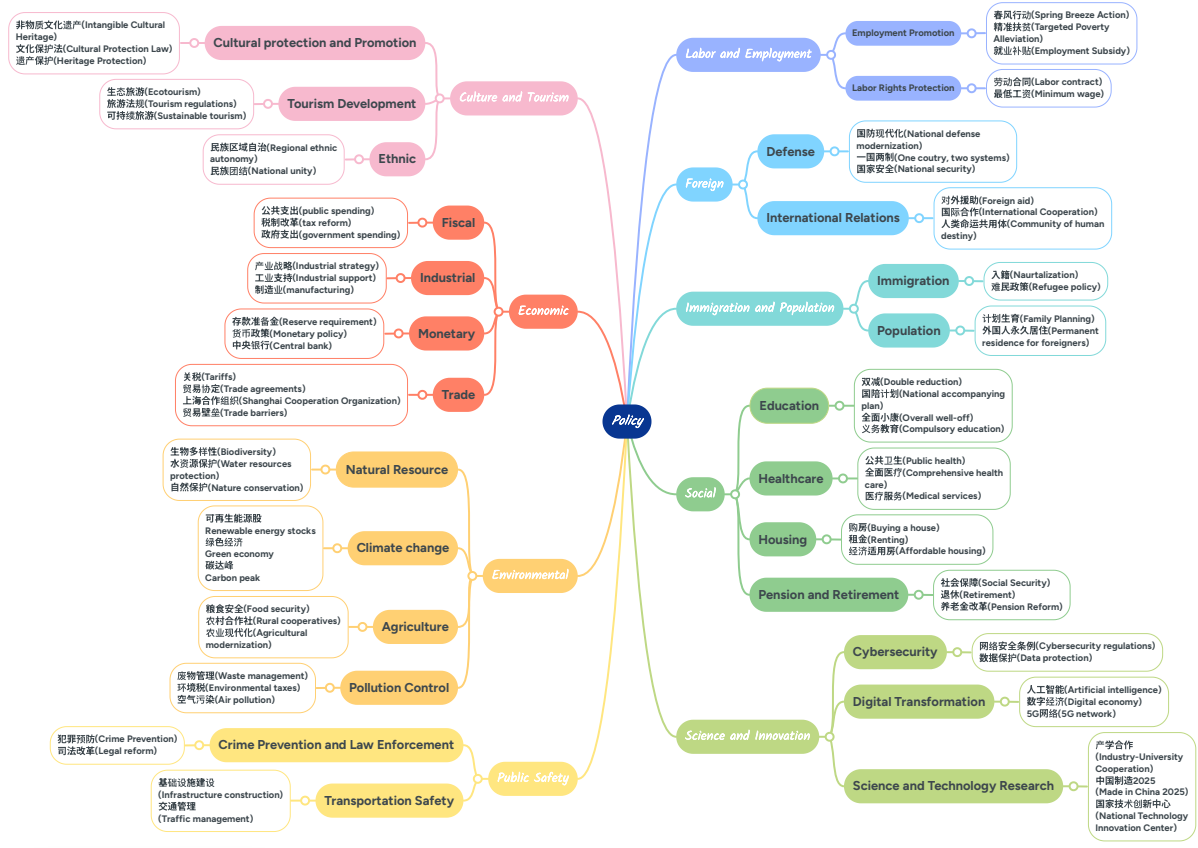


Figure 6: Categories of Chinese Policy Documents and Representative Keywords (Partial List).

tasks. Recent works have begun to address LLM evaluation in the policy domain. For instance, PolicySimEval (Kang et al., 2025) assesses policy outcomes through simulation, while UNBench (Liang et al., 2025) evaluates performance on political science tasks from a UN perspective. These benchmarks differ from PolicyBench, which focuses on the fine-grained comprehension of policy texts across broad domestic domains and governmental systems. In parallel, studies by Safaei et al. (Safaei and Longo, 2024) and Karacapilidis et al. (Karacapilidis et al., 2024) explore the application of LLMs for policy generation and deliberation. While these efforts target the 'output' capabilities of LLMs, our work addresses a complementary and foundational gap: evaluating the model's 'input' capability to precisely understand policy language, a crucial prerequisite for any reliable downstream application.

D Error Study

To better understand the limitations of current LLMs on PolicyBench, we conducted a qualitative error analysis by sampling representative failure cases from each level, as illustrated in Figure 8.

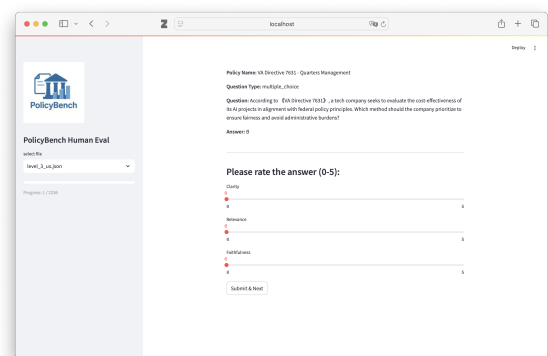


Figure 7: Screenshot of human evaluation interface.

Table 5: Sources for Policies and Supplementary Materials (By Language)

Language & Content Type	Primary Collection Websites / Platforms
US Sources	
Official Policies	https://www.transportation.gov https://www.hhs.gov https://www.va.gov https://www.commerce.gov https://www.usda.gov https://www.energy.gov https://www.doi.gov https://www.ed.gov https://www.treasury.gov https://www.state.gov https://www.dhs.gov https://www.hud.gov
Supplementary Materials	https://www.cnn.com/ https://www.foxnews.com https://www.reuters.com
Chinese Sources	
Official Policies	https://www.gov.cn/zhengce/zhengcewenjianku/
Supplementary Materials	https://www.xuexi.cn/ http://www.people.com.cn/

Level1:Memorization

Question-1: "According to HHS Authorizes Five States to Provide Historic Health Care Coverage for People Transitioning out of Incarceration, the states of Illinois, Kentucky, Oregon, Utah, and Vermont join states such as ___ in pursuing an innovative option for expanding access to health care."
Option: A: 'Massachusetts', B: 'New Mexico', C: 'California', D: 'Arizona'
Classification: 'Organization Identification'
A

04-mini: Illinois, Kentucky, Oregon, Utah, and Vermont join states such as New Mexico in pursuing an innovative option for expanding access to health care. Therefore, the answer is B
Answer: True

Question-2: 根据《国务院关于在海南博鳌乐城国际医疗旅游先行区暂时调整实施《中华人民共和国药品管理法实施条例》有关规定的决定》,海南省人民政府要加强医疗机构使用药品的___管理。
Option: A: "审批", B: "监督", C: "监管", D: "质量"
Classification: "条款/日期记忆"
B

QwQ-32B: 根据国务院的决定,海南博鳌乐城先行区在调整实施《药品管理法实施条例》时,需要强化药品使用的管理。政策调整通常伴随着监管职责的加强,特别是在简化审批流程后,后续监管和风险控制变得更加重要。"监管"一词在政府文件中常指动态的、全面的管理和监督职责,而"监督"更侧重于检查和督促。"质量"则特指产品特性。因此,最符合政策语境的选项是加强"监管"。因此,答案是C
Answer: C

Level3:Application

Question-5: 根据关于调整完善土地出让收入使用范围优先支持乡村振兴的意见,某市土地出让收入为1000亿元,但土地出让收益仅为200亿元。若该市选择按土地出让收入比例计提资金,则当年用于农业农村的资金至少应达到多少?
Option: A: "80亿元", B: "100亿元", C: "200亿元", D: "500亿元"
Classification: "政策计算类推理"
Explanation: 根据《意见》,若按土地出让收入计提,则用于农业农村的资金占比应逐步达到10%以上。因此,该市当年用于农业农村的资金至少为:1000亿元 × 8% = 80亿元。
E

QwQ-32B: 根据《关于调整完善土地出让收入使用范围优先支持乡村振兴的意见》,若某市选择按土地出让收入比例计提资金用于农业农村,到2025年需达到的比例不低于10%。题目中土地出让收入为1000亿元,因此至少应计提1000 × 10% = 100亿元。选项B对应100亿元。
Answer: B

Question-6: "According to (VA Handbook 5383), a volunteer appointed to a Testing Designated Position (TDP) by the VA Center for Development and Civic Engagement refuses to sign the Release of Information (ROI) form. What action should the VA take?"
Option: A: "Allow the volunteer to continue duties without testing"; B: "Terminate the volunteer's position immediately"; C: "Subject the volunteer to disciplinary action, including removal"; D: "Refer the volunteer to EAP counseling."
Classification: "Procedural/Institutional Implementation"
Explanation: "The policy states volunteers in TDPs must sign ROI forms. Refusal subjects them to disciplinary action (Section 8a: 'Employees refusing to sign the ROI will be subject to disciplinary action, up to and including removal')"
F

Gemini-2.0-Flash: "According to VA Handbook 5383, volunteers in TDPs are subject to drug testing. Refusal to sign the Release of Information form is considered a refusal to test. Refusal to test should result in the termination of the volunteer's position. Therefore the correct action is to terminate the volunteer's position immediately."
Answer: B

Question-7: 根据科技部关于支持建设国家合成生物技术创新中心的高,某跨国公司企业计划与创新中心合作开发新型生物基材料。根据政策,双方合作中需遵循哪些制度设计原则?
Classification: "流程/制度执行"
Explanation: "治理机制: 接受理事会领导下的中心主任负责制(第四段); 利益分配: 构建'科研与市场紧密结合、利益共享、风险共担'的创新共同体(第四段); 资源共享: 企业需参与开放共享科研设施(第四段); 国际合作: 需符合'链接全球资源的创新网络'导向(第四段)"
G

GPT-4o: 根据科技部关于支持建设国家合成生物技术创新中心的政策,跨国公司企业与创新中心的合作需遵循以下制度设计原则: 1. **创新导向**:: 确保合作项目符合科技创新的总体要求,推动技术迭代和产业升级。2. **公平合作**:: 保障合作双方在技术开发和成果转化中的权利和义务平等,促进资源互补和优势互补。3. **知识产权保障**:: 严格界定知识产权归属,明确合作产生的知识产权归属,保护各方合法权益。4. **合规管理**:: 严格遵守国家法律法规,遵守国家安全、环境保护和伦理等相关法律法规,确保技术开发过程的安全性和合规性。5. **透明公开**:: 合作过程保持透明,定期通报和评估项目进展,确保各方知情权和监督权,促进合作的健康发展。这些原则确保合作项目科学性、合法性和可持续性。
Score: 1.5/5

Level2:Understanding

Question-3: INVESTING IN AMERICA Biden-Harris Administration Announces \$2.4 Billion in New Rail Projects to Improve Safety, Strengthen Supply Chains and Grow America's Freight and Passenger Rail Networks
Option: true/false
Classification: "Institutions"
Explanation: "The policy states USDOT and FRA have pushed freight railroads to provide guaranteed paid sick leave to all of their workers, aligning with institutional safety and labor reforms."
C

Gemma-3: "The Consolidated Rail Infrastructure and Safety Improvements (CRISI) program, funded through the Bipartisan Infrastructure Law, focuses on rail safety projects. While improved worker safety is a goal, the specific mention of 'expanding paid sick leave mandates' for freight railroads workers isn't a direct component outlined in the official announcements regarding CRISI funding or FRA safety initiatives. The primary focus is on infrastructure improvements, technology deployment, and addressing specific safety risks associated with rail operations. The issue of paid sick leave for railroad workers has been a point of negotiation between unions and freight rail companies but isn't specifically funded or mandated through this particular program."
Answer: False

Question-4: 根据国务院办公厅关于促进金融租赁行业健康发展的指导意见,金融租赁行业被赋予推动国家战略的功能,以下哪项最能反映政策背后的意识形态导向?
Option: A: "服务'一带一路'建设,促进国际产能合作"; B: "支持供给侧结构性改革,推动产业转型升级"; C: "促进绿色金融发展,支持环保和可持续发展项目"; D: "支持实体经济发展,助力中小企业融资创新"
Classification: "Ideas"
Explanation: "政策第三部分将金融租赁与国家重大战略结合,强调'推动'一带一路'、京津冀协同发展等战略"。选项A体现了通过金融工具服务国家发展战略目标的理念,属于顶层设计中的意识形态表达。其余选项为具体执行层面内容。
D

Claude-3.7-Sonnet: 或更准确地说,分析这个过程中各选项与国家战略和意识形态导向的关系。国务院关于促进金融租赁行业健康发展的指导意见中,确实赋予了金融租赁行业推动国家战略的功能,从选项来看,选项A"服务'一带一路'和国际产能合作,体现的是中国对外开放和国际合作的战略,这是外向型发展。选项B"支持供给侧结构性改革,推动产业转型升级,这体现了供给侧结构性改革的核心理念,即通过调整经济结构和生产要素配置,提高供给体系的质量和效率,增强经济的内生动力。选项C"促进绿色金融发展,支持环保和可持续发展项目,这体现了绿色金融发展的核心理念,即通过金融手段支持环保和可持续发展,实现经济效益和社会效益的统一。选项D"支持实体经济发展,助力中小企业融资创新,这体现了支持实体经济发展的核心理念,即通过金融手段支持实体经济的发展,特别是助力中小企业融资创新,增强经济的活力和韧性。因此,最能反映政策背后的意识形态导向的是选项A,因为它体现了中国对外开放和国际合作的战略,这是外向型发展的核心理念。
Answer: B

Figure 8: Representative error cases across three cognitive levels in PolicyBench.

Table 6: Overview of models evaluated or used as judges in this study.

Model	Developer	Open-source	Version	Role
GPT-4o	OpenAI	✘	/	Evaluated model
o4-mini		✘	2025-04-16	Evaluated model; Judge (Level 3)
Gemini-2.5-Flash	Google DeepMind	✘	preview-04-17	Evaluated model; Judge (Level 3)
Gemini-2.0-Flash		✘	/	Evaluated model
Gemma-3-27B		✔	27B Instruct	Evaluated model
Claude-3.7-Sonnet	Anthropic	✘	20250219	Evaluated model; Judge (Level 3)
Claude-3.5-Sonnet		✘	20241022	Evaluated model
LLaMA-4	Meta	✔	maverick-instruct	Evaluated model
QwQ-32B	Alibaba (Qwen Team)	✔	/	Evaluated model
DeepSeek-V3	DeepSeek AI	✔	/	Evaluated model
DeepSeek-R1		✔	/	Evaluated model; Judge (Level 3)

Level 1 (Memorization): Factual Confusions. The most common errors involve incorrect recall or misidentification, often triggered by the presence of distractors with strong semantic similarity. For example, when asked to identify US states that previously adopted a healthcare policy, the model incorrectly selected “New Mexico” instead of “California”—likely due to co-occurrence bias in training data. In Chinese policy questions, LLMs sometimes confuse regulatory terms with subtle distinctions (e.g., “Jiān Dū” vs. “Jiān Guǎn”), indicating limited sensitivity to nuanced legal language.

Level 2 (Understanding): Misinterpretations. Errors at this stage largely stem from misreading policy intent, ideological framing, or institutional logic. Models tend to misread underlying motivations or over-rely on surface cues. For instance, a model incorrectly identified “Supply-side structural reform” as the core ideological signal of a financial leasing policy, despite explicit references to national strategies like the Belt and Road Initiative. Similarly, when analyzing US labor protection clauses, the model missed the correct interpretation of sick leave provisions, favoring superficial summaries over deeper institutional mandates described in the text.

Level 3 (Application): Reasoning and Procedural Failures. Models frequently struggled with quantitative reasoning, procedural interpretation, and hallucinated conclusions. For example, a model failed to compute the required rural funding quota from land sale revenue, despite having all necessary information. In another case, a model prema-

turely recommended immediate termination for a volunteer’s non-compliance, overlooking the step-wise disciplinary procedures mandated by policy. For open-ended tasks, models sometimes hallucinate plausible-sounding but unsupported principles, rather than extracting the specific cooperation principles explicitly mentioned in the document.

Overall, these error patterns reveal that LLMs struggle across recall, understanding, and real-world reasoning, especially with legal nuance and institutional complexity. Addressing these issues may require targeted supervision or retrieval-augmented approaches.

E Data Quality & Human Evaluation

E.1 Annotation Team & Methodology

To ensure the rigorousness of *PolicyBench*, our annotation team was structured hierarchically. The core annotation, quality control, and validation phases were led by **five senior Ph.D. candidates** specializing in Public Policy, Law, and Computational Social Science. Three undergraduate assistants provided support for preliminary data formatting and cleaning but did not perform substantive labeling or validation tasks.

E.2 General Quality Evaluation

We conducted a comprehensive human evaluation between March 2025 and April 2025. Each generated question was independently evaluated along three key dimensions on a 1–5 Likert scale. We sampled **1,500** (Level 1), **1,000** (Level 2), and **500** (Level 3) questions balanced across languages.

Table 7: The detailed introduction of 10 dimensions.

Dimension	ID	Definition
Article/Date Memorization	1-1	Tests memory of specific articles, dates, numbers, etc. in the policy.
Terminology Recognition	1-2	Tests ability to recognize and understand policy terminology.
Organization Identification	1-3	Tests ability to identify organizations mentioned in the policy.
Idea Understanding	2-1	Examines the ideological foundation and value orientation behind the policy.
Interest Understanding	2-2	Assesses the identification and analysis of key stakeholders affected by or involved in the policy.
Institution Understanding	2-3	Evaluates understanding of the formal and informal rules, organizations, and mechanisms shaping policy implementation.
Policy-Based Numerical Reasoning	3-1	Perform simple mathematical reasoning or calculation based on the numerical provisions in the policy text.
Scenario-Based Decision-Making	3-2	Based on specific scenarios, determine how the parties should make decisions or choose the most appropriate approach based on the policy.
Procedural/Institutional Implementation	3-3	Examine the understanding and memory of the specific operational procedures, implementation steps, and institutional regulations in the policy.
Policy Logic and Value Explanation	3-4	Focus on the background motivation, target value and logical structure of policy making.

As shown in Table 8, the dataset exhibits consistently high quality, with average scores exceeding **96%** across all dimensions.

E.3 Inter-Annotator Agreement (IAA)

To validate the reliability of the human evaluation, we performed a double-blind annotation on a random subset of 600 items (20% of the evaluation set). We utilized *Cohen’s Kappa* (κ) to measure agreement beyond chance, calculated as:

$$\kappa = \frac{P_o - P_e}{1 - P_e} \quad (1)$$

where P_o is the observed agreement and P_e is the expected agreement by chance.

Calculation Method:

- For **Bloom-level classification** (Nominal), we calculated standard unweighted κ .
- For **Quality dimensions** (Ordinal 1–5), we **binarized** the scores into “Accept” (Score ≥ 4) and “Reject” (Score ≤ 3) to strictly measure the consistency of inclusion criteria.

As shown in Table 9, we achieved strong agreement ($\kappa > 0.80$) across all dimensions. Disagreements were resolved via a two-stage adjudication process involving a senior expert.

E.4 Expert Validation Study

To further verify the dataset’s validity against a professional standard, we conducted an additional expert validation study with four external experts

(**2 Ph.D. candidates in Public Policy, 1 Professor in Computational Social Science, and 1 Senior Sociology Researcher**). They evaluated a stratified sample of 120 questions.

1. Data Validity: Experts assessed the correctness of the Gold Answers and the appropriateness of the content. Agreement for Bloom-level labeling was calculated using *Krippendorff’s* α , which is robust for multiple raters:

$$\alpha = 1 - \frac{D_o}{D_e} \quad (2)$$

where D_o is the observed disagreement and D_e is the expected disagreement by chance. The results in Table 10 confirm the dataset’s high quality.

2. Expert Performance Baseline: To establish a human ceiling, the experts answered the questions under an **open-book setting**, simulating realistic policy analysis workflows. As shown in Table 11, experts significantly outperform models in deeper understanding tasks (L2/L3), validating the benchmark’s difficulty gradient.

F Formal Definition of the Policy Tasks

To address reviewer concerns regarding the lack of a formal task definition, we provide here a unified formulation of the *Policy Comprehension Task*, followed by fine-grained instantiations corresponding to the three cognitive levels and ten sub-tasks, details in Table 7.

Table 8: Human evaluation results across levels and regions.

Level	Region	Clarity	Relevance	Faithfulness	Avg.
Level 1	CN	99.20%	98.01%	98.78%	98.66%
	US	99.15%	98.25%	97.94%	98.45%
Level 2	CN	98.86%	98.00%	98.12%	98.32%
	US	98.76%	97.60%	97.22%	97.86%
Level 3	CN	99.16%	96.02%	96.14%	97.11%
	US	97.66%	95.64%	96.32%	96.54%

Table 9: Inter-Annotator Agreement statistics. Quality scores were binarized for κ calculation.

Annotation Dimension	Kappa (κ)	Agreement Level
Bloom Classification (L1/L2/L3)	0.856	Strong
Clarity (Pass/Fail)	0.841	Strong
Relevance (Pass/Fail)	0.877	Strong
Faithfulness (Pass/Fail)	0.827	Strong

Table 10: Expert Validation Results ($N = 120$).

Metric	Score
Correctness Verification	96.1%
Content Validity (Rated ‘‘Appropriate’’)	94.5%
Bloom-Label Agreement (Krippendorff’s α)	0.86

Table 11: Human Expert Performance (Open-Book) vs. Models.

Level	Expert Accuracy	Comparison vs. SOTA LLM
Level 1 (Memorization)	82.3%	Comparable
Level 2 (Understanding)	88.4%	Significant Gap
Level 3 (Application)	90.1%	Significant Gap

F.1 General Formulation

At a high level, we model policy comprehension as a conditional question-answering problem grounded in policy documents. Let C denote a policy context, which may consist of a full policy document, a section, or a clause describing rules, actors, and institutional mechanisms. Let Q_k denote a query designed to probe a specific cognitive level $k \in \{1, 2, 3\}$, corresponding to memorization, understanding, and application, respectively.

Given (C, Q_k) , a language model f_θ is expected to produce an answer \hat{A} that is consistent with the policy content and satisfies the cognitive requirement implied by k :

$$\hat{A} = \arg \max_{A \in \mathcal{A}} P(A | C, Q_k; \theta), \quad (3)$$

where \mathcal{A} denotes the answer space, which can be either a finite set of options (for multiple-choice questions) or free-form text (for open-ended ques-

tions).

F.2 Task Instantiation by Cognitive Level

We instantiate the policy comprehension task into **10 concrete sub-tasks**, organized into three cognitive levels.

Level 1: Memorization (Factual Retrieval) *Objective: Retrieve explicit facts or entities directly stated in the policy text C .*

I-1 Article / Date Memorization Recall specific temporal markers or citation identifiers.

Example: ‘‘According to the *U.S.-Philippines Civil Nuclear Cooperation Agreement*, the Agreement entered into force on [Mask].’’

Answer: July 2.

I-2 Terminology Recognition Identify the definition of a domain-specific term as stated in the text.

Example: ‘‘Civil nuclear cooperation agreements provide a legal framework for exports of [Mask].’’

Answer: material, equipment, and components.

I-3 Organization Identification Identify organizational hierarchy or institutional affiliation.

Example: ‘‘According to the *National Behavioral Health Workforce Career Navigator*, SAMHSA is an agency within [Mask].’’

Answer: U.S. Department of Health and Human Services (HHS).

Level 2: Understanding (Conceptual Interpretation) *Objective: Map explicit policy text to implicit concepts under the ‘‘3I’’ framework (Ideas, Interests, Institutions).*

1223	2-1 Idea Understanding	Infer the underlying goal, ideology, or strategic intent.	3-4 Policy Logic and Value Explanation	Explain conflicts or trade-offs between policy objectives.	1268
1224					1269
1225		<i>Example:</i> “The <i>Big Data Project (BDP)</i> aims to democratize access to environmental data primarily through collaboration with whom?”		<i>Example:</i> “Scott Turner’s goal to reduce reliance on government aid conflicts with which aspect of DC’s housing strategy?”	1271
1226					1272
1227					1273
1228					1274
1229		<i>Answer:</i> Cloud Service Providers.		<i>Answer:</i> Funding for emergency rental assistance.	1275
1230	2-2 Interest Understanding	Identify stakeholders and their eligibility, benefits, or losses.			1276
1231				This structured definition clarifies both the scope and the granularity of policy comprehension capabilities evaluated by <i>PolicyBench</i> .	1277
1232		<i>Example:</i> “Are Tribal entities eligible to apply for 2501 Program grants according to the USDA announcement?”			1278
1233					1279
1234		<i>Answer:</i> True.	G Multi-Examiner Bias Analysis		1280
1235				To address critical concerns regarding “ <i>Model-Speak</i> ” (stylistic cues) and “ <i>Familiarity Bias</i> ” (models favoring their own generation patterns), we conducted a comprehensive Multi-Examiner Sensitivity Analysis . This study empirically validates that our heterogeneous examiner pool serves as a necessary safeguard against evaluation artifacts.	1281
1236	2-3 Institution Understanding	Comprehend rules, categories, or allocation mechanisms.			1282
1237					1283
1238		<i>Example:</i> “Which category under the <i>Clean Energy</i> program allocates capacity to projects ensuring 50% of benefits go to low-income households?”			1284
1239					1285
1240		<i>Answer:</i> Category 4.	G.1 Experimental Setup		1286
1241				We designed a controlled experiment to decouple the “Examiner” (Question Generator) from the “Examinee” (Evaluated Model).	1287
1242					1288
1243	Level 3: Application (Scenario Reasoning)	<i>Objective:</i> Apply policy rules to a novel or hypothetical scenario.		1. The Examiner Pool: We utilized three distinct top-tier models to generate questions and distractors, ensuring coverage across different model families:	1289
1244					1290
1245					1291
1246	3-1 Policy-Based Numerical Reasoning	Execute numerical calculations derived from policy formulas.			1292
1247					1293
1248					1294
1249		<i>Example:</i> “If a \$2 million FEMA project has an 80% federal cost share and a 1.91% cost increase, how much does the local government pay?”			1295
1250					1296
1251		<i>Answer:</i> \$7,640.			1297
1252					1298
1253					1299
1254	3-2 Scenario-Based Decision Making	Determine the compliant action in a simulated situation.			1300
1255					1301
1256		<i>Example:</i> “If the Colonial Pipeline were subject to the new safety rule during a leak, which action would be required?”			1302
1257					1303
1258		<i>Answer:</i> Stage response personnel in pre-defined zones.			1304
1259					1305
1260					1306
1261	3-3 Procedural / Institutional Implementation	Validate procedural conditions or sequences.			1307
1262					1308
1263		<i>Example:</i> “What is a necessary condition for the transfer of nuclear reactors under the U.S.-Thailand 123 Agreement?”			1309
1264					1310
1265		<i>Answer:</i> Commitment to nonproliferation standards.			1311
1266					1312
1267					1312

1313
1314
1315

1316

1317
1318
1319
1320

1321

1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354

1355

1356
1357
1358
1359
1360

- **LOEO (Leave-One-Examiner-Out):** Questions generated by the remaining two examiners (e.g., *Wo-GPT*).

G.2 Full Leaderboard Sensitivity Results

Table 12 presents the complete performance matrix. The results demonstrate significant score variations under single-examiner conditions, confirming the necessity of our Consensus Baseline.

G.3 Analysis of Biases

1. Self-Scoring Bias (Familiarity). Models consistently perform differently on questions they generated themselves, creating a "Familiarity Bonus" or "Penalty". As shown in Table 13, relying on a single generator creates severe distortions:

Insight: The data reveals divergent biases. GPT-4o benefits from "Familiarity Bias" (+7%), likely exploiting its own stylistic patterns. Conversely, Claude exhibits "Self-Strictness" (-15%), penalizing its own generation logic. **PolicyBench’s consensus approach effectively averages out these extremes**, anchoring scores to a neutral ground truth.

2. Mitigating Model-Speak: Leaderboard Stability. To determine if "Model-Speak" (stylistic tells) compromises the validity of the rankings, we analyzed the **Spearman Rank Correlation** (ρ) (Zar, 2005) between the Baseline leaderboard and other conditions.

3. External Model Robustness. For models outside the generator pool, like **Llama-4**, the Baseline provides the most stable evaluation. Llama-4’s score varies from 82.0% to 89.0% across single examiners. The Baseline (82.0%) successfully anchors it to a consensus difficulty, filtering out examiner-specific noise.

Conclusion: The low correlation of the GPT-Only set ($\rho = 0.34$) proves that single-source benchmarks are fundamentally biased. The multi-examiner design in *PolicyBench* is a necessary mechanism to ensure that high performance reflects genuine *Policy Comprehension* rather than *Stylistic Alignment*.

H LLM-as-a-Judge Validation

To ensure the reliability and validity of our automated evaluation pipeline, we conducted a two-fold validation process: assessing *stability* (consistency across runs) and *alignment* (accuracy against human experts).

H.1 Evaluation Stability Analysis

Evaluation Stability. To rigorously assess the stability of our LLM-as-a-Judge pipeline and address concerns about the stochastic nature of model outputs, we conducted a multi-run, multi-judge evaluation analysis. We randomly sampled 10 question-answer pairs from the Level 3 open-ended test set. The scoring for each sample was performed as follows:

- **Initial Scoring:** For each evaluation run, two distinct LLM judges were sampled from our pool to score the response on a 0–5 scale (in 0.5 increments). The score for the run was the average of these two ratings.
- **Discrepancy Resolution:** If the scores from the two initial judges differed by more than 1.0 point, a third tie-breaker judge was invoked to provide an additional score. In such cases, the final score for the run was the average of all three judges’ ratings. This protocol ensures robustness against outlier judgments.
- **Repetition:** This entire scoring process was conducted three independent times for each of the 10 cases.

We then calculated the mean and standard deviation of the three final scores for each case. The results, summarized in Table 15, demonstrate high consistency, with standard deviations remaining exceptionally low. This indicates that our multi-judge protocol effectively mitigates run-to-run variance and produces reliable evaluations.

H.2 Human-Model Alignment

Stability is a necessary but insufficient condition for validity. To demonstrate that our LLM-as-a-Judge pipeline accurately reflects expert judgment, we conducted a **Human-Model Alignment Study**. **Experimental Setup.** We sampled a subset of open-ended responses from Level 3 tasks. These responses were independently scored by our LLM Judge pipeline and by senior human experts (Ph.D. candidates in Public Policy) using the exact same rubric.

Quantitative Results. As shown in Table 16, the LLM Judge demonstrates strong alignment with human experts across three key metrics:

- **Pearson Correlation** ($r = 0.87$): Indicates a strong positive linear relationship between model and human scores.

Table 12: Complete Multi-Examiner Sensitivity Analysis. Scores represent accuracy (%). "Baseline" denotes the standard PolicyBench setting (3-Examiner Consensus). "Wo-X" denotes the Leave-One-Examiner-Out setting. The data reveals that relying on a single examiner (e.g., GPT-Only or Claude-Only) leads to drastic score fluctuations compared to the stable Baseline.

Model	Baseline (Consensus)	Single-Examiner Setting			LOEO Setting			Avg.
		Claude-Only	GPT-Only	Qwen-Only	Wo-Claude	Wo-GPT	Wo-Qwen	
Qwen-3	89.0	86.0	88.0	92.0	92.0	90.0	88.0	89.29
Llama-4	82.0	84.0	88.0	89.0	88.0	87.0	85.0	86.14
Qwen-2.5	83.0	81.0	85.0	90.0	83.0	86.0	87.0	85.00
GPT-4o-mini	66.0	59.0	74.0	85.0	74.0	71.0	67.0	82.67
GPT-4o	75.0	71.0	82.0	85.0	85.0	83.0	78.0	79.86
Claude-4-Sonnet	84.0	69.0	49.0	89.0	71.0	52.0	85.0	71.29
Claude-4-Haiku	60.0	47.0	81.0	86.0	59.0	48.0	75.0	91.20

Table 13: Analysis of Self-Scoring Bias. The results show that single-examiner benchmarks suffer from extreme variance (from +7 inflation to -15 deflation).

Model	Baseline Score	Own-Gen Score	Δ	Bias Type
GPT-4o	75.0%	82.0%	+7.0%	<i>Self-Leniency / Pattern Matching</i>
Qwen-3	89.0%	92.0%	+3.0%	<i>Moderate Inflation</i>
Claude-4-Sonnet	84.0%	69.0%	-15.0%	<i>Self-Strictness / Hyper-Critical</i>

Table 14: Leaderboard Stability Analysis. High correlation in LOEO conditions confirms the robustness of the ranking system, while single-examiner rankings (especially GPT-Only) are highly unstable.

Comparison Pair	Spearman ρ	Kendall τ	Interpretation
Baseline vs. Wo-Qwen	0.901	0.781	<i>Highly Stable</i>
Baseline vs. Wo-GPT	0.607	0.524	<i>Moderately Stable</i>
Baseline vs. GPT-Only	0.342	0.293	<i>Unstable / Biased</i>

- **Mean Absolute Error (MAE = 0.42):** On average, the model’s score deviates from the human score by less than 0.5 points (the smallest scoring increment).
- **Agreement Rate (94%):** In 94% of cases, the model’s score fell within an acceptable margin (≤ 1.0 point) of the expert score.

Mitigating Subjectivity via Rubric-Based Scoring. The high alignment is primarily attributed to our **Point-Based Rubric** design (detailed in Appendix G). Unlike generic "quality" assessments which can be subjective, our prompt explicitly directs the model to verify the presence of specific *Key Points* derived from the reference answer. The model assigns partial credit based on these matched points rather than an abstract "feeling" of quality. This structured approach significantly reduces hallucination and subjectivity, ensuring the judge acts as an objective verifier.

I Correlation with External Benchmarks

To demonstrate that *PolicyBench* evaluates a distinct capability orthogonal to general reasoning or pure legal knowledge, we analyzed the performance correlation between *PolicyBench* and two representative benchmarks: **MMLU-Pro** (Wang et al., 2024) (General Reasoning) and **LegalBench** (Guha et al., 2023) (Legal Reasoning).

Key Findings:

- **Negative Correlation with General Reasoning** ($r \approx -0.69$): Surprisingly, models with top-tier general reasoning (e.g., DeepSeek-V3) do not necessarily excel in policy comprehension. This suggests that policy analysis involves specific logic (e.g., institutional constraints) that general benchmarks fail to capture.
- **No Correlation with Legal Reasoning** ($r \approx -0.07$): The near-zero correlation with LegalBench indicates that *understanding policy* (Ideas, Interests) is fundamentally different from *applying law* (Statutes). *PolicyBench* fills this critical gap in the evaluation landscape.

J Comparing PolicyMoE with a Standard LoRA

To validate the architectural contribution of *PolicyMoE* and demonstrate its superiority over standard parameter-efficient fine-tuning, we conducted

Table 15: Stability analysis of the final LLM-as-a-Judge scores. The low standard deviation across three independent runs for each case demonstrates the reliability of our multi-judge evaluation protocol.

Case ID	Final Score (Run 1)	Final Score (Run 2)	Final Score (Run 3)	Mean	Std. Dev.
Case 1	4.00	4.50	4.00	4.17	0.29
Case 2	4.00	4.00	4.25	4.08	0.12
Case 3	3.00	3.25	3.00	3.08	0.12
Case 4	4.33	4.50	4.25	4.36	0.10
Case 5	2.25	2.25	2.50	2.33	0.12
Case 6	5.00	5.00	5.00	5.00	0.00
Case 7	3.67	4.00	3.50	3.72	0.25
Case 8	1.50	1.25	1.00	1.25	0.20
Case 9	4.00	3.75	4.00	3.92	0.12
Case 10	4.67	5.00	4.50	4.72	0.25

Table 16: Human–Model alignment statistics. High correlation, low error, and strong agreement indicate that the LLM judge closely matches expert evaluation.

Metric	Value
Pearson Correlation (r)	0.87
Mean Absolute Error (MAE)	0.42
Agreement Rate	94%

a controlled ablation study. We compared our *PolicyMoE* architecture against a well-tuned **Standard LoRA** baseline.

J.1 Experimental Setup

To ensure a fair comparison, both models were trained on the same stratified subset of the PolicyBench training data using Qwen2.5-7B-Instruct as the backbone.

- **Standard LoRA:** Fine-tuned using a single LoRA adapter (Rank=16, Alpha=32) applied to all target modules, treating all levels of tasks as a unified objective.
- **PolicyMoE (Ours):** Fine-tuned using our routing architecture with three specialized experts (Memory, Understanding, Application), utilizing the same data and hyperparameters.

J.2 Results & Analysis

As shown in Table 18, *PolicyMoE* outperforms the Standard LoRA baseline across all metrics, with an average accuracy improvement of **+2.68%**.

Mitigating Task Interference. These results suggest that the MoE architecture helps alleviate *task interference* arising from heterogeneous cognitive objectives.

- **Decoupling Memorization and Reasoning:** Standard LoRA must encode both exact factual recall (Level 1) and flexible scenario reasoning (Level 3) within a single low-rank adaptation, which can lead to competing gradient signals.
- **Improved Performance on Level 1:** By routing Level 1 queries to a dedicated **Memory Expert**, *PolicyMoE* achieves a **+3.81%** improvement over Standard LoRA, suggesting that specialized experts help preserve precise factual representations.
- **Implication:** These findings indicate that *PolicyMoE* provides a structurally meaningful extension over standard LoRA when adapting LLMs to multi-level policy comprehension tasks.

K Selected Policy Samples

This section showcases examples of the policy titles from our dataset, from both China and the United States. A partial list is provided in Figure 9.

L Prompt Template

This section shows the 3 key prompts used for data curation. **Prompt 1** converts clean policy text into cloze-style questions. **Prompt 2** instructs LLMs to generate incorrect answers, which serve as the distractors for multiple-choice questions. **Prompt 3** employs an LLM-as-a-judge methodology to evaluate the models’ responses to the questions. Prompts used for processing policies from different countries were designed to correspond to the respective national language. Only the English prompts are presented here. Prompts were translated using Google Translate to ensure consistency between the two languages. All translated content was sub-



Figure 9: Collected Policies (Part).

Table 17: Performance comparison and correlation analysis across benchmarks.

Model	MMLU-Pro	LegalBench	PolicyBench (Avg)	PolicyBench (L2)
DeepSeek-V3	81.9%	80.1%	59.10%	57.68%
GPT-4o	80.3%	79.8%	59.47%	56.08%
Claude-3.7-Sonnet	80.3%	78.1%	64.13%	58.48%
Gemini-2.5-Flash	77.9%	81.7%	63.82%	64.06%
Pearson Correlation (r)	-0.69	-0.07	1.00	–

Table 18: Performance comparison between the base model, Standard LoRA, and *PolicyMoE* on a controlled training subset. *PolicyMoE* achieves larger gains on Memorization (L1) and Application (L3) tasks.

Model	Level 1 (Memorization)	Level 2 (Understanding)	Level 3 (Application)	Average
Base (Qwen2.5-7B)	30.10%	44.00%	55.69%	43.26%
Standard LoRA	34.82%	44.51%	59.45%	46.26%
<i>PolicyMoE</i> (Ours)	38.63%	44.90%	63.30%	48.94%
Δ (<i>MoE</i> vs. <i>LoRA</i>)	+3.81%	+0.39%	+3.85%	+2.68%

sequently reviewed and tested by human annotators
to avoid potential semantic inconsistencies.

1512
1513

Prompt 1: Level-1 cloze-style Generation

You are a policy expert and your task is to generate questions based on the given policy text.

- Strictly follow the given material to generate questions, do not fabricate content.
- Generate 5-10 fill-in-the-blank questions and 3-8 true or false questions based on the length of the policy.
- Answer to the questions should be clear and precise, avoid ambiguous answers.
- Answer to the questions should preferably be a single word or phrase.
- If the answer is a false judgment question, Avoid altering, adding, or deleting the original text when the answer is not an incorrect judgment question.
- Don't generate questions related to non-key information such as file numbers.
- Each question should start with "According toXXX", please provide the full name of the policy.
- Questions should not be limited to a single aspect (such as time), and should be diversified.

policy title: {policy}

policy content: {policy_content}

Prompt 2: Distractor generation

You are an ai assistant tasked with answering policy-related questions.

- Answer the questions based on your knowledge.
- Please note that some incorrect answers are provided below. You must not make the same mistakes,
- Your answer needs to be semantically distinct from the given incorrect answer.
- Don't say you can't see the image, just answer based on your knowledge.
- Don't generate overly lengthy answers, keep them concise and to the point.
- The answer you generate needs to be factually different from the given incorrect answer.
- Try to use straightforward words instead of being too abstract or vague.

question:{question}

wrong answers:{wrong_answer}

Prompt 3: LLM-as-a-Judge

You are an expert evaluator. Your task is to score the following open-ended answer based on a reference answer and scoring criteria. Follow these rules carefully:

1. For calculation or factual questions where the result must be precise (e.g., math, unit conversion, logical problems), if the final answer is incorrect, the score should be 0, regardless of the explanation.

2. For general questions (e.g., reasoning, explanation, analysis), the reference answer includes multiple key points.

- Compare the given answer with the reference key points.
- For each matched key point, assign partial credit proportionally.
- If the answer includes correct but unlisted points (beyond the reference answer), you may award partial credit with explanation.

3. Provide a score from 0 to 5. Generally:

- 5 = Completely correct and well explained
- 4 = Mostly correct, with minor issues
- 3 = Partially correct, some key points missing or wrong
- 2 = Mostly incorrect but with small redeeming aspects
- 1 = Barely relevant or correct
- 0 = Completely wrong or irrelevant

4. In your reasoning, clearly list:

- Which points in the reference answer are matched
- Any extra correct points beyond the reference
- Justify any deductions

5. Be strict but fair. Do not be lenient.

—
Question: {question}

Reference Answer:

{reference_answer_with_point_marks}

User Answer: {user_answer}

—
Now output:

Score: X

Reasoning: ...