

JurisBench: A Deep Benchmark for Assessing Large Language Models in Professional Legal Practice

Anonymous ACL submission

Abstract

While Large Language Models (LLMs) have achieved high accuracy on isolated legal QA and "exam-style" benchmarks, their reliability in handling the interdependent, procedural workflows of real-world professional legal practice remains largely unproven. To address this gap, we introduce **JurisBench**, a vertical, depth-oriented benchmark designed to evaluate legal LLMs across the full lifecycle of Chinese civil litigation. JurisBench introduces a **Linear Depth Simulation** track that mirrors the cognitive workflow of professional judges through four sequential, dependency-aware phases: *Cause of Action* prediction, *Focus of Disputes* prediction, *Rationale of the Judgment* prediction and *Result of the Judgment* prediction. Experimental results from state-of-the-art LLMs reveal a stark "illusion of competence": while models excel in isolated generative tasks, their performance collapses in an end-to-end pipeline due to substantial error propagation. We identify **precise statutory grounding** as a persistent bottleneck, highlighting a critical gap between fluent linguistic output and practical judicial reliability. JurisBench provides a diagnostic framework for developing more robust, workflow-aware legal AI. JurisBench provides a principled framework and a diagnostic testbed for developing next-generation legal AI capable of professional-grade adjudication.

1 Introduction

While Large Language Models (LLMs) exhibit broad knowledge spanning diverse domains, they lack deep understanding of specialized tasks inherent to professional practice. Legal reasoning exemplifies this challenge.

The rapid advancements of LLMs do demonstrate their transformative potential in the legal field (Surden, 2018; Walters and Novak, 2021; Lee et al., 2023), and recent studies report strong performance on legal question answering and knowledge-

oriented tasks (Siino et al., 2025), but such results do not necessarily imply reliable performance in realistic judicial scenarios, where decisions are interdependent, procedural constraints are strict, and errors may propagate across multiple stages. Robust, practice-aligned evaluation benchmarks are therefore essential for accessing and understanding the actual capabilities and limitations of LLMs in real-world case-processing settings.

Meanwhile, the escalating volume of civil and commercial caseloads in China has placed immense pressure on judicial systems amplifying the demand for reliable automated tools (see Appendix A). However, the complexity of real-world adjudication means that superficial legal competence—often measured by accuracy on isolated or exam-style tasks (Guha et al., 2023; Fei et al., 2023; Dai et al., 2023; Li et al., 2024b; Fan et al., 2025)—is insufficient for end-to-end judicial workflows. This mismatch between *surface-level knowledge* and *procedural reasoning* underscores the need for benchmarks that assess legal reliability across the full lifecycle of a judicial process. Existing research, while providing valuable signals of linguistic competence, often focuses on "breadth-first" task coverage at the expense of "professional depth," offering limited insight into the procedural coherence required for case-level adjudication.

To address these limitations, this paper introduces **JurisBench**, a vertical, depth-oriented benchmark specifically tailored to the structured judicial workflow of the Chinese legal system. JurisBench prioritizes "professional depth" by simulating the cognitive workflow of judges (Ashley, 2017) through a four-phase **Linear Depth Simulation** pipeline: *Cause of Action* prediction, *Focus of Disputes* identification, *Rationale of the Judgment* generation, and *Result of the Judgment* determination (detailed in Section 3). By focusing on these interdependent tasks, JurisBench enables the systematic analysis of **reasoning stability** and **error**

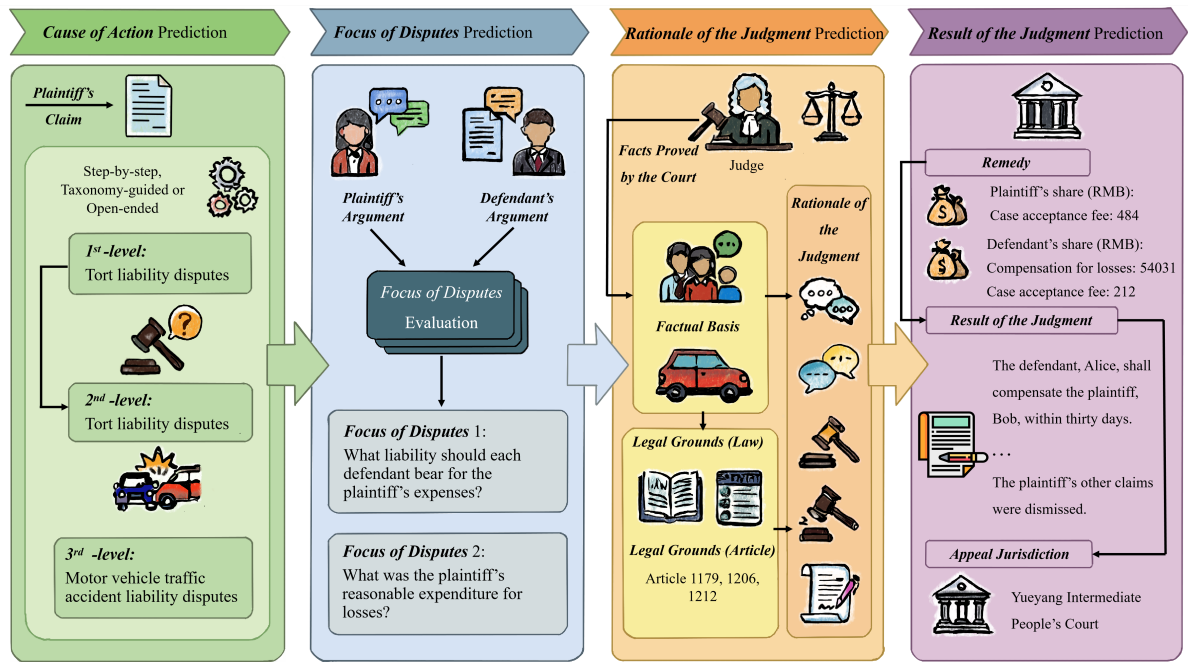


Figure 1: **An Overview of JurisBench Pipeline.** The figure illustrates the *Linear Depth Simulation* track, which models a structured judicial reasoning process for each case. The pipeline follows a decision chain: (1) Predict the *Cause of Action* based on the *Plaintiff’s Argument*; (2) Identify the *Focus of Disputes* by incorporating the *Defendant’s Argument*; (3) Derive the *Rationale of the Judgment* grounded in *Factual Basis* and *Legal Grounds*; and (4) Determine the *Result of the Judgment*. The logical justifications for this sequential workflow design, specifically the procedural dependencies where each phase serves as a prerequisite for the next, are elaborated in Appendix C.

propagation—critical dimensions that remain invisible in traditional QA-style evaluations in the field (Zheng et al., 2023; Prince Tritto and Ponce, 2025).

Benchmarking deep understanding necessarily differs from testing broad knowledge. This purpose requires a focus on specific case types with varying levels of difficulty. While JurisBench maintains broad coverage at the categorization level, its deep evaluation track focuses on *motor vehicle traffic accident liability disputes*—a representative archetype of civil litigation involving multi-party participation, factual causality analysis, and statutory grounding. This unified domain serves as a diagnostic testbed for probing professional-level reasoning depth that broader but shallower benchmarks are not designed to capture.

In summary, our contributions are as follows:

- We introduce **JurisBench**, the first deep, vertical, workflow-aligned benchmark that simulates the full lifecycle of Chinese civil litigation to evaluate the practical performance of LLMs in professional legal practice.
- We propose the **Linear Depth Simulation** track to analyze the stability of reasoning

chains and the impact of cascading error propagation across interdependent judicial stages.

- Through extensive experiments on 10+ state-of-the-art LLMs, we reveal a significant gap between surface-level legal knowledge and deep judicial reasoning capability, identifying accurate comprehension of statutory grounding as a persistent bottleneck.

2 Related Work

2.1 Legal Large Language Models

LLMs have been widely touted for their exceptional performance across diverse industries, with claims that they can match—or are poised to replace—domain experts (Bubeck et al., 2023; Co-manici et al., 2025). Such rhetoric has spurred researchers and startups to invest heavily in the training and development of professional domain LLMs, among which legal LLMs are particularly prominent, especially Chinese legal LLMs.

Legal LLMs have progressed from pre-trained models like LegalBERT (Chalkidis et al., 2020) and LexLM (Chalkidis et al., 2023) to generative systems mainly based on GPT (Walters and Novak, 2021; Lee et al., 2023; Surden, 2018) and

LLaMA (Touvron et al., 2023; Kassianik et al., 2025; Hossain et al., 2025; Prince Tritto and Ponce, 2025), with Chinese pioneers including ChatLaw (Cui et al., 2023), Lawyer LLaMA (Huang et al., 2023), and Lawyer GPT (Yao et al., 2024).

Technically, mainstream legal LLM adaptation relies on (i) **domain-adaptive pre-training** (Chalkidis et al., 2023), (ii) **task-driven supervised or instruction tuning** for legal QA and case analysis (Hendrycks et al., 2021; Shen et al., 2022), and (iii) **preference-based alignment** (e.g., RLHF) (Ouyang et al., 2022; Prince Tritto and Ponce, 2025). To ensure factual rigor, Retrieval-augmented generation (RAG) is therefore widely adopted to mitigate knowledge incompleteness (Lewis et al., 2020), with legal-specific pipelines further incorporating citation grounding and verification mechanisms (Qian et al., 2025; Zheng et al., 2025; Zhang et al., 2024).

However, even with this external knowledge support, a recurring reliability issue persists in practice-aligned workflows: error propagation across multi-stage reasoning chains (Surden, 2018; Prince Tritto and Ponce, 2025). Even with authoritative materials retrieved, legal rule application may remain inconsistent across procedurally coupled stages (Doyle and Tucker, 2025).

2.2 Benchmarks: From General-purpose to Legal-domain

General-purpose benchmarks, such as MMLU (Hendrycks et al., 2020) and AGIEval (Zhong et al., 2024), standardize LLM evaluation via closed-form tasks and outcome-oriented metrics. While effective for broad comparisons, they abstract away procedural structures and dependencies, limiting their utility for structurally constrained domains.

Existing legal benchmarks generally follow three paradigms: (i) **exam-style batteries** for broad automated scoring (e.g., LexGLUE (Chalkidis et al., 2022), LawBench (Fei et al., 2023)); (ii) **typology-driven sets** focused on reasoning types or verifiability (e.g., LegalBench (Guha et al., 2023), CitaLaw (Zhang et al., 2024)); and (iii) **interactive benchmarks** for tool-augmented workflows (e.g., LegalAgentBench (Li et al., 2024a)). Despite their diversity, these benchmarks predominantly treat tasks as independent units, inheriting an "exam-style" assumption that legal capability is a sum of isolated outcomes.

Consequently, they under-represent the proce-

dural interdependence of real-world adjudication, making cross-stage inconsistencies difficult to diagnose (Surden, 2018; Prince Tritto and Ponce, 2025; Posner and Saran, 2025). This gap motivates workflow-aligned benchmarks that evaluate procedural depth and error propagation, such as **JurisBench** proposed in this work.

3 JurisBench

This section details the systematic development of JurisBench, encompassing data curation, the mapping of sub-tasks to target judicial capabilities, complexity stratification, and Evaluation metrics, as summarized in Figure 2.

3.1 Overview

JurisBench operationalizes the end-to-end cognitive process of Chinese civil litigation, mapping case filing to final judgments using authentic judicial documents (**Supreme People’s Court of the PRC**). To ensure ethical compliance, all records were anonymized via placeholder-based de-identification. Given the high frequency of missing fields in raw judicial records, we applied systematic filtering to ensure structural integrity; only 5.23% of the original corpus met the criteria for inclusion (see Appendix B). A comprehensive semantic analysis and t-SNE visualization, provided in Appendix D, confirm the high diversity and task variance of the JurisBench dataset across its stratified complexity levels. The benchmark employs a *dual construction strategy* to evaluate models across two complementary dimensions (see Table 1).

Table 1: JurisBench dataset statistics across breadth and depth dimensions.

Evaluation Track	Complexity Phase	Number of Cases
Depth: Pipeline Evaluation (MVTAL disputes)	Simple	394
	Intermediate	394
	Complex	212
Breadth: Taxonomy Alignment (General 282 <i>Causes of Action</i>)	Phase 1-a	2,820
	Phase 1-b	

Breadth Evaluation of Cause of Action At the breadth level, the benchmark assesses the models’ systematic capability to navigate the wider Chinese civil law landscape. This evaluation set comprises 2,820 cases spanning 282 distinct *Causes of Action*. To investigate the impact of informational support, this track is bifurcated into two settings: **Phase 1-a**, where a hierarchical candidate list (legal

JurisBench: Dataset Construction and LLM Assessment Methodology

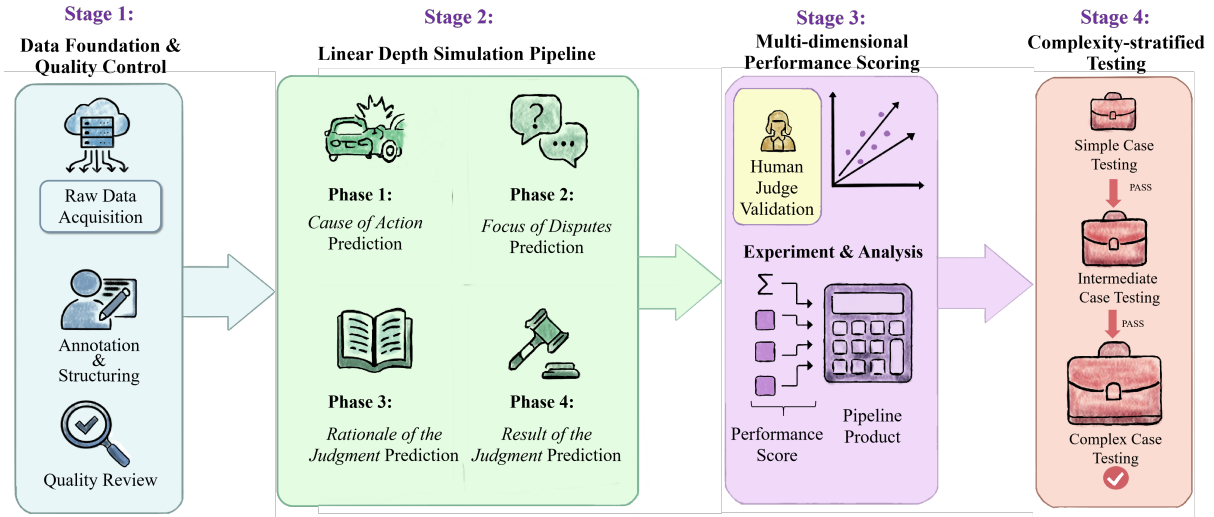


Figure 2: **Overall Methodology.** It transitions from a high-quality Data Foundation (Stage 1) to a core Linear Depth Simulation mirroring four judicial phases (Stage 2). Performance is then quantified using a Multi-dimensional Scoring system validated by human judges (Stage 3), followed by a Progressive Evaluation across three complexity-stratified tiers (Stage 4).

database) is provided within the context to simulate a restricted-selection scenario, and **Phase 1-b**, an open-ended prediction task requiring purely parameterized recall. These benchmarks determine whether models can accurately align factual narratives with rigid legal ontologies.

Depth Evaluation via Reasoning Pipeline. At the depth level, JurisBench focuses on a representative and structurally rich domain—Motor Vehicle Traffic Accident Liability (MVTAL) disputes. Evaluation is organized around 4 phases (see Table 2 for details), which is a dependency-aware evaluation pipeline that reflects how intermediate judicial decisions condition downstream reasoning, while parallel evaluation tracks are retained to assess standalone legal knowledge. The logical justifications for this sequential pipeline/workflow design, specifically the procedural dependencies where each phase serves as a prerequisite for the next, are elaborated in Appendix C. See Appendix K for a detailed example of the case data structure.

3.2 Task Taxonomy and Capability Mapping

To systematically analyze the cognitive demands on LLMs, we map the JurisBench phases into five distinct capability dimensions (detailed in Table 2). This taxonomy enables a granular diagnosis of legal reasoning: (1) **Semantic Taxonomy Alignment** (Phase 1) assesses the ability to map informal narratives to rigid legal ontologies; (2) **Adversarial**

Factual Reconstruction (Phase 2, 3-1) evaluates information synthesis from conflicting arguments to reconstruct a coherent Factual Basis; (3) **Statutory Symbolic Grounding** (Phase 3-2, 3-3) tests the precision of legal provision citations to mitigate hallucinations; (4) **Logical & Numerical Inference** (Phase 4-1, 4-3) represents System 2 reasoning for executing fixed algorithms and jurisdictional rules; and (5) **Constrained Legal Generation** (Phase 3-4, 4-2) evaluates the synthesis of logically consistent judicial rationales and results under strict procedural constraints.

3.3 Complexity Stratification and Evaluation Strategy

JurisBench evaluates legal reasoning under varying levels of case complexity using a three-tier stratification scheme: *Simple*, *Intermediate*, and *Complex*. Evaluation follows an simple-to-complex progression, in which models advance to higher-complexity subsets only after demonstrating adequate performance on simpler cases. This design reflects judicial practices for distinguishing difficult and complex cases and enables efficient identification of both baseline competence and expert-level reasoning limits.

The distinction between case complexity levels is grounded in the existence of baseline criteria for assessing adjudicative difficulty in judicial practice. Such criteria encompass multiple factors that af-

Table 2: **Detailed task design and evaluation metrics for JurisBench.** We define sub-tasks across four judicial phases, detailing their sequential input configurations, target outputs, and the metrics used to assess performance.

Index	Task Description	Output	Judicial Capability	Metric
Phase 1: Cause of Action Prediction				
<i>Input: Basic Case Info + Plaintiff's Argument</i>				
1	MVTAL disputes, Taxonomy-guided	<i>Cause of Action</i>	Semantic Taxonomy Alignment	Level-Acc.
1-a	Predict general 282 <i>Causes of Action</i> , Taxonomy-guided	<i>Cause of Action</i>	Semantic Taxonomy Alignment	Level-Acc.
1-b	Predict general 282 <i>Cause of Action</i> , Open-ended	<i>Cause of Action</i>	Semantic Taxonomy Alignment	Level-Acc.
Phase 2: Focus of Disputes Prediction				
<i>Input: Basic Case Info + Plaintiff's Argument + Defendant's Argument + Cause of Action</i>				
2	Evaluate each disputed issue	<i>Focus of Disputes</i>	Adversarial Factual Reconstruction	STS, EM
Phase 3: Rationale of the Judgment Prediction				
<i>Input: Basic Case Info + Plaintiff's Argument + Defendant's Argument + Cause of Action + Facts Proved by the Court</i>				
3-1	Extract useful <i>Factual Basis</i> for <i>Rationale of the Judgment</i>	<i>Factual Basis</i>	Adversarial Factual Reconstruction	STS
3-2	Retrieve case-relevant law provisions	<i>Legal Grounds (Law)</i>	Statutory Symbolic Grounding	F1-score
3-3	Predict exact legal provisions for <i>Rationale of the Judgment</i>	<i>Legal Grounds (Article)</i>	Statutory Symbolic Grounding	F1-score
3-4	Predict <i>Rationale of the Judgment</i> based on litigation request	<i>Rationale of the Judgment</i>	Constrained Legal Generation	STS
Phase 4: Result of the Judgment Prediction				
<i>Input: Basic Case Info + Plaintiff's Argument + Defendant's Argument + Cause of Action + Facts Proved by the Court + Rationale of the Judgment</i>				
4-1	Calculate individual costs	<i>Remedy</i>	Logical & Numerical Inference	Tol-EM
4-2	Predict individual <i>Result of the Judgment</i> based on litigation request	<i>Result of the Judgment</i>	Constrained Legal Generation	STS
4-3	Analyse <i>Appeal Jurisdiction</i>	<i>Appeal Jurisdiction</i>	Logical & Numerical Inference	EM
Note: Level-Acc.: Level-wise Top-3 Accuracy; STS: Semantic Textual Similarity; Tol-EM: $\pm 10\%$ Tolerance Exact Match; EM: Exact Match; F1: F1-score; (detailed in Section 3.4)				

fect both factual complexity and legal relationship complexity, including the difficulty of fact-finding and the structure of legal relations involved. These factors are not arbitrarily defined but are abstracted from statutory provisions and judicial interpretations that explicitly recognize their impact on trial organization and adjudicative effort (Supreme People's Court of the PRC, 2022c; National People's Congress of the PRC, 2021).

Based on aggregated complexity scores derived from these normatively grounded indicators, the Jenks Natural Breaks Optimization method (Jenks, 1967) is applied to partition cases into three difficulty levels with coherent legal characteristics and clear separations in adjudicative complexity. Detailed indicator definitions and scoring procedures are provided in Appendix E.

3.4 Evaluation Metrics

To accurately assess model performance across the diverse cognitive demands of the judicial process, we categorize our evaluation metrics into four functional groups, covering every subtask from initial filing to final judgment.

Classification and Exact Match. For the prediction of the *Cause of Action* (including the representative disputes in Phase 1, as well as Phase 1-a and Phase 1-b), we utilize **Level-wise Top-3**

Accuracy. Under this metric, a prediction at a specific hierarchical level is deemed correct if the ground-truth label is contained within the model's three highest-probability candidates. Given the hierarchical nature of the Chinese legal ontology, the evaluation for a specific sub-category (level $l + 1$) is conditioned on the ground-truth of its parent category (level l). This approach mirrors the deductive reasoning of judges and avoids penalizing valid high-level categorizations due to minor errors at the leaf nodes. For the analysis of *Appeal Jurisdiction* (Phase 4-3), we employ **Exact Match** to ensure the model correctly identifies the specific appellate court and its location.

Set-based Retrieval. For the prediction of *Legal Grounds*, which involves identifying both the relevant *Laws* (Phase 3-2) and specific *Articles* (Phase 3-3), we utilize the **F1-score**. Since a single case typically relies on multiple statutory references, this metric balances the precision of cited provisions with the recall of all mandatory legal grounds, effectively penalizing both the omission of critical laws and the hallucination of irrelevant ones.

Discretionary Numerical Tolerance. For the calculation of individual *Remedy* costs (Phase 4-1), we apply a $\pm 10\%$ **Tolerance Exact Match**. This metric recognizes the legally sanctioned "judgment margin" in Chinese civil law, where compen-

sation items such as "mental distress" or "reasonable expenses" are determined by locality-specific socio-economic statistics and judicial discretion rather than a unified national schedule (National People’s Congress of the PRC, 2020; Supreme People’s Court of the PRC, 2003). A prediction is deemed correct if it falls within a 10% interval of the reference amount.

Semantic Generative Evaluation. For narrative and complex reasoning tasks—including the identification of the *Focus of Disputes*, the extraction of the *Factual Basis*, the generation of the *Rationale of the Judgment*, and the determination of the final *Result of the Judgment*(Phase 2, 3-1, 3-4, 4-2)—surface-level n-gram metrics, such as ROUGE(Lin, 2004) and its series, are insufficient. To handle synonyms and logical sensitivity (e.g., "support" vs. "dismiss"), we implement a triple-strategy: (1) **Embedding-based Similarity**, calculating the cosine similarity between vector representations; (2) **LLM-as-a-Judge**, leveraging a high-capability model to evaluate factual consistency and logical coherence. An expert-led validation study confirmed that these automated proxies maintain a high correlation with manual ratings provided by professional judges (see Appendix G); (3) **Keyword-based Lexical Accuracy**, where students with legal backgrounds were tasked with annotating a set of mandatory professional keywords for each generative subtask (see Appendix G). The accuracy for these phases is then determined by the proportion of these essential keywords successfully recovered in the model’s generated output. While the primary results for generative phases in the main text are calculated using the first two metrics, we provide a complementary assessment based on professional keyword recovery in Appendix F.4.

4 Experiment

4.1 Experimental Setup

We evaluate 13 representative LLMs, encompassing state-of-the-art Multilingual models (Claude 4.5 (Anthropic, 2025), GPT-4o-mini (OpenAI, 2025), Gemini-3 Pro (DeepMind, 2025), Llama 3.3 70B (AI, 2024), and Grok 4 (xAI, 2025)) and Chinese-Oriented models (DeepSeek-R1 (Guo et al., 2025), Doubao-1.5-Pro (ByteDance, 2025), Kimi-K2 (Team et al., 2025), Qwen3-Instruct (Yang et al., 2025), and GLM-4.5 (Zeng et al., 2025)). Additionally, 3 Legal-Specific models are included to assess the impact of

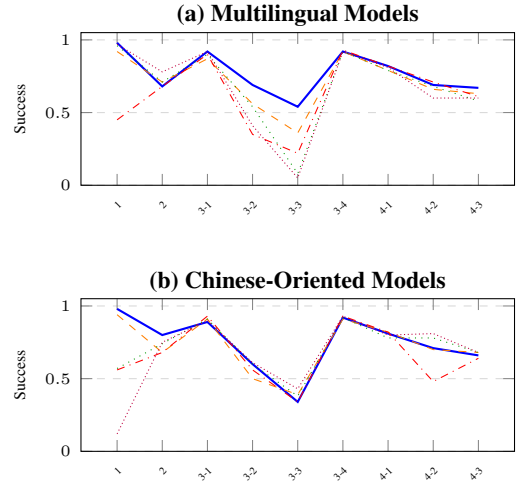


Figure 3: Parallel coordinates plot visualizing performance decay. The longitudinal axis represents accuracy across sequential pipeline phases. Each line represents a LLM under test; the blue lines represent Gemini-3 Pro and Qwen3-Instruct respectively.

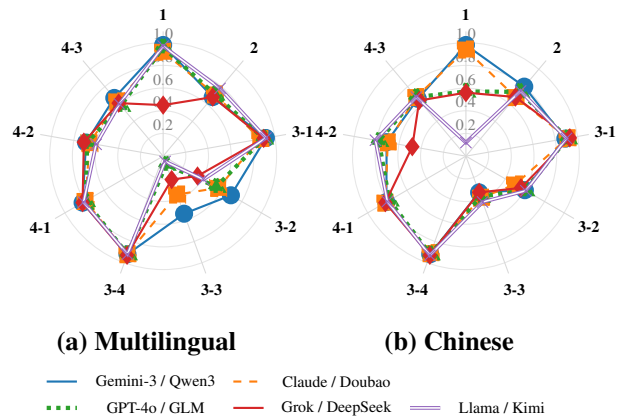


Figure 4: Capability radar charts across 9 sub-tasks of the 4 phases. Concentric rings indicate accuracy levels (0.2–1.0).

domain-specific fine-tuning (see Appendix F.3). The full evaluation prompts used in experiments are detailed in Appendix L.

4.2 Methodology & Main Results

Workflow-Oriented Modular Evaluation To precisely measure model capabilities at each judicial stage, we first adopt a *modular evaluation* strategy. For each pipeline phase, the model is provided with ground-truth (*gold*) inputs from all preceding stages (e.g., the true *Cause of Action* is provided as a premise for *Focus of Disputes* identification). Detailed results are shown in Table 3. This "Oracle-assisted" approach ensures that performance failures at a later stage are not confounded

Table 3: Experimental results (zero-shot) of 10 LLMs (*Simple* difficulty subset). **Bold** and underline indicate the best and second-best performance in each column, respectively. Cell background colors represent a performance heatmap, transitioning from **red** (lower performance) to **green** (higher performance). The "Product" column represents per-mille (‰) success rate with one decimal place.

Group	Model	Phase 1	Phase 2	Phase 3				Phase 4			Product
		1	2	3-1	3-2	3-3	3-4	4-1	4-2	4-3	
Multilingual	Gemini-3 Pro	0.982	0.684	<u>0.917</u>	0.689	0.535	0.919	0.820	0.693	0.669	79.3‰
	Claude 4.5	0.923	0.711	0.868	0.561	0.355	0.921	0.790	0.661	0.626	34.2‰
	GPT-4o-mini	0.957	0.709	0.901	0.542	0.083	0.919	0.794	<u>0.681</u>	<u>0.576</u>	7.9‰
	Grok 4	0.447	0.681	0.901	0.347	0.224	0.926	0.817	0.711	0.606	6.9‰
	Llama 3.3 70B	<u>0.970</u>	<u>0.782</u>	<u>0.917</u>	<u>0.413</u>	0.051	0.923	0.808	0.604	0.598	3.9‰
Chinese	Qwen3-Instruct	0.982	0.798	0.892	0.602	0.336	0.924	0.810	0.710	0.660	49.5‰
	Doubao-1.5-Pro	0.937	0.679	0.909	0.502	0.387	0.909	<u>0.818</u>	0.700	0.679	39.7‰
	GLM-4.5	0.571	0.738	0.912	0.604	0.383	0.922	0.776	<u>0.776</u>	0.675	33.3‰
	DeepSeek-R1	0.556	0.683	0.927	0.564	0.344	0.927	0.815	0.480	0.640	15.9‰
	Kimi-K2	0.117	0.750	0.904	<u>0.607</u>	<u>0.427</u>	0.930	0.800	0.807	0.684	8.4‰

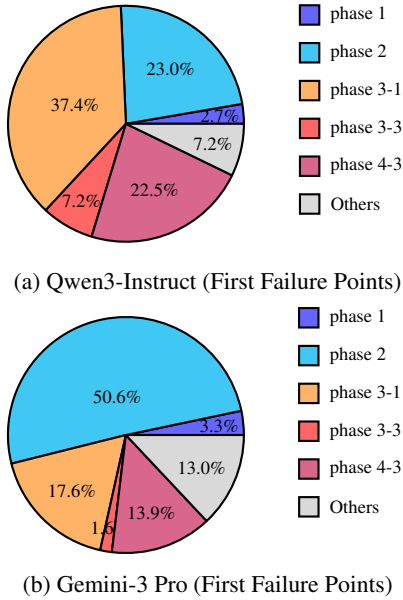


Figure 5: Distribution of the first point of failure in the unassisted pipeline.

by reasoning errors inherited from earlier phases, allowing for a granular diagnosis of specific capability bottlenecks (visualized in Figure 3 and 4).

Unassisted Simulation and Systemic Reliability

To investigate "true" end-to-end reliability—where errors are unassisted and cumulative—we transition from oracle assistance to a sequential, *unassisted simulation* on top-performing models (Gemini-3 Pro and Qwen3-Instruct). We enforce a *strict dependency chain*: for each individual case, if a model yields a score of 0 at any stage, it triggers a **Terminal Pipeline Failure**, rendering all downstream reasoning moot (see Figure 6).

To reflect the overall coherence, we define the

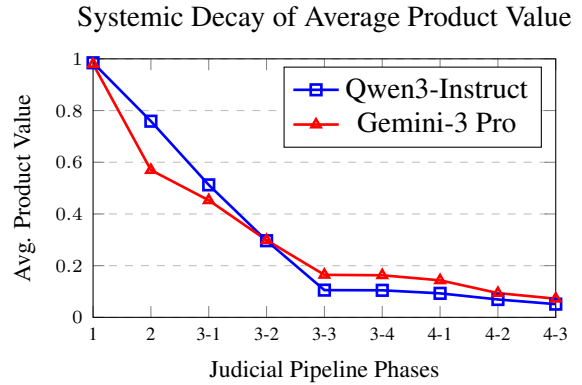


Figure 6: Systemic decay of the **Average Product** across the unassisted judicial pipeline. Plotted values represent the mean joint success rate calculated *only* for the subset of cases that have successfully traversed all preceding phases without a terminal failure. This visualization highlights the catastrophic impact of cascading errors on overall systemic reliability.

Product metric (‰) as the joint success rate. While the modular Product (Table 3) represents an *idealized upper bound*, the unassisted simulation tracks the *Effective Systemic Reliability* and the distribution of **First Failure Points** (Figure 5) to identify where the reasoning chain initially fractures.

Prompting Strategy and Scope All primary results reported in the main text are obtained under a **zero-shot** setting to assess the models' intrinsic legal reasoning capabilities. The impact of in-context learning via **one-shot** demonstrations, as well as the evaluation of Phase 1-a and 1-b, are detailed in Appendix F.

4.3 Analysis

Our analysis of the experimental results yields the following core findings:

The Bottleneck of Statutory Symbolic Grounding While most models exhibit robust performance in *Semantic Alignment* (Phase 1, avg. >0.9) and *Factual Reconstruction* (Phase 3-1, avg. >0.9), they struggle significantly with *Symbolic Grounding*. The capability footprints in Figure 4 reveal a universal "indentation" at Phase 3-3 (*Legal Grounds-Article Prediction*), where even Gemini-3 Pro achieves only 0.535. The parallel coordinates (Figure 3) present a stark "competence gap": LLMs excel at summarizing narratives but lack the "legal fidelity" required to anchor facts to specific statutory provisions. This acts as a **hard filter**: once a model fails this precise grounding, the subsequent accuracy in the pipeline decays drastically.

Modular Success vs. Systemic Fragility A striking discrepancy exists between modular scores and unassisted reliability. As shown in the simulation (Figure 6), systemic performance exhibits a catastrophic cascade effect. By the final phase (4-3), even the best-performing models achieve a **conditional pipeline success** of merely **5.12%** (Qwen3) and **7.22%** (Gemini-3 Pro).

Critically, these values represent a *conditional average*—calculated only from the subset of cases that have not yet encountered a terminal failure in preceding stages. As illustrated by the "reasoning fractures" in Figure 5, the vast majority of cases are filtered out much earlier; for instance, 50.6% of Gemini-3 Pro cases fracture at Phase 2. This implies that the *Effective Systemic Reliability*, factoring in the cumulative **Survival Rate**, is significantly lower than even these marginal percentages suggest. This evidence proves that judicial reasoning behaves as a "weakest-link" system: a single failure in early-stage information synthesis (as seen in Qwen3's 37.4% fracture at Phase 3-1) effectively nullifies any linguistic fluency achieved in later generative stages.

Thinking-Heavy Models vs. General-Purpose Stability A nuanced trade-off emerges between specialized reasoning and end-to-end stability. As shown in Table 3, while DeepSeek-R1 exhibits the highest stability in the initial extraction of the *Factual Basis* (Phase 3-1, 0.927), its final **Product** metric remains relatively low (15.9%), suggesting that general-purpose reinforcement learning for rea-

soning ("thinking" traces) does not automatically translate into the procedural rigor required for a multi-stage judicial workflow.

Furthermore, we observe a "specialization split" among Chinese-oriented models: for instance, Kimi-K2 shows the weakest performance in initial *Cause of Action* alignment (Phase 1, 0.117) but demonstrates a superior ability in final *Result of the Judgment* determination (Phase 4-2, 0.807). Despite these localized strengths, Gemini-3 Pro maintains a dominant lead in the final **Product** success rate (79.3%). This highlights a critical need for future legal AI to bridge the gap between "thinking" prowess and "procedural" consistency.

5 Conclusion

JurisBench introduces a vertical, workflow-aligned benchmark for evaluating LLMs in professional legal case-processing scenarios, emphasizing domain-specific depth rather than broad but shallow task coverage. By assessing model performance across interdependent judicial stages, the benchmark reveals that strong accuracy on isolated legal subtasks does not reliably translate into coherent end-to-end case handling, with early-stage errors often propagating and precise statutory grounding emerging as a persistent bottleneck. Although current models perform unsatisfactorily under this evaluation, such outcomes highlight the diagnostic value of JurisBench in exposing failure modes that remain largely invisible to exam-style or surface-level benchmarks. In this way, JurisBench provides a principled foundation for evaluating and guiding the development of legal LLMs toward deeper professional specialization and more reliable reasoning in real judicial contexts.

JurisBench not only provides a structured and operational framework for evaluating the real-world capabilities of LLMs in professional legal practice, but also carries broader methodological implications for the design of vertical domain benchmarks. This workflow-centric perspective offers a reusable paradigm for constructing benchmarks in other high-stakes, highly specialized domains—such as healthcare, financial regulation, and engineering decision-making—where surface-level competence is insufficient and reliability across the full decision process is critical.

525 Limitations

526 Despite the depth of JurisBench, several limitations
527 remain. First, **Domain Specificity**: while Motor
528 Vehicle Traffic Accident Liability (MVTAL) dis-
529 putes are structurally representative, they do not
530 encompass the full complexity of other civil do-
531 mains, such as intellectual property or corporate
532 liquidations, which may require different reason-
533 ing paradigms. Second, **Metric Sensitivity**: al-
534 though our embedding-based and LLM-aided met-
535 rics demonstrate high correlation with human ex-
536 pert ratings, automated proxies remain inherently
537 insufficient to capture every professional nuance
538 of judicial reasoning. There remains a resolution
539 gap between current automated evaluation frame-
540 works and the exhaustive logical rigor required in
541 a formal legal setting. Third, **Static Nature and**
542 **Contamination**: as a fixed benchmark, JurisBench
543 faces potential data contamination risks as LLMs
544 are increasingly trained on web-scale crawls. Addi-
545 tionally, the dynamic nature of judicial interpreta-
546 tions in China necessitates periodic updates to the
547 ground truth to remain legally valid. Finally, while
548 our study assesses both modular and unassisted
549 pipeline simulations to identify systemic bottle-
550 necks, it does not yet investigate fully autonomous
551 agentic workflows or more complex interactive rea-
552 soning architectures, which remain a promising
553 direction for future research.

554 Ethics Statement & Potential Risks

555 **Preventing Misinterpretation.** We explicitly
556 state that JurisBench is a *diagnostic tool* for iden-
557 tifying reasoning fractures in LLMs, not a certifi-
558 cation for autonomous adjudication. Our findings,
559 particularly the low "Product" success rates, un-
560 derscore that current LLMs are not ready for high-
561 stakes judicial decision-making. High benchmark
562 scores should not be equated with the professional
563 and moral judgment of a human judge.

564 **Bias and Privacy.** As JurisBench is derived from
565 authentic judicial documents, it may reflect his-
566 torical systemic biases. Models evaluated on this
567 data might inadvertently reinforce these biases;
568 thus, results should not be viewed as absolute leg-
569 al "correctness." Regarding privacy, although we
570 implemented rigorous anonymization for names
571 and identifiers, the unique combination of case
572 facts in open datasets poses a theoretical risk of
573 re-identification. We urge the community to use

this data strictly for research purposes and in com- 574
pliance with global data protection standards (e.g., 575
GDPR or local equivalents). 576

Jurisdictional Scope. JurisBench is tailored to 577
the Chinese civil law system. Users should exer- 578
cise extreme caution when applying its evaluation 579
logic to common law jurisdictions or other legal 580
frameworks, as reasoning patterns and procedural 581
requirements differ significantly. 582

References 583

- Meta AI. 2024. Introducing meta llama 3. <https://ai.meta.com/blog/meta-llama-3/>. Accessed: 584
2026-01-05. 585
586
- Anthropic. 2025. Claude sonnet 4.5 announce- 587
ment. <https://www.anthropic.com/news/claude-sonnet-4-5>. Accessed: 2026-01-05. 588
589
- Kevin D. Ashley. 2017. *Artificial Intelligence and Legal 590
Analytics: New Tools for Law Practice in the Digital 591
Age*. Cambridge University Press, Cambridge. 592
- Beijing Chaoyang District People’s Court. 2025. *Deep- 593
ening the reform of separating complex and simple 594
civil litigation procedures (reform case no. 180)*. 595
Supreme People’s Court Website, Accessed: 2025- 596
11-19. 597
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, 598
Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter 599
Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, and 600
1 others. 2023. Sparks of artificial general intelli- 601
gence: Early experiments with gpt-4. *arXiv preprint 602
arXiv:2303.12712*. 603
- ByteDance. 2025. Doubao-1.5-pro model docu- 604
mentation. https://seed.bytedance.com/zh/special/doubao_1_5_pro. Accessed: 2026-01-05. 605
606
- Ilias Chalkidis, Michalis Fergadiotis, Prodromos 607
Malakasiotis, Nikolaos Aletras, and Ion Androu- 608
sopoulos. 2020. Legal-bert: The muppets straight 609
out of law school. *arXiv preprint arXiv:2010.02559*. 610
- Ilias Chalkidis, Nicolas Garneau, Cătălina Goanță, 611
Daniel Katz, and Anders Søgaard. 2023. Lexfiles 612
and legallama: Facilitating english multinational 613
legal language model development. In *Proceedings 614
of the 61st Annual Meeting of the Association for 615
Computational Linguistics (Volume 1: Long Papers)*, 616
pages 15513–15535. 617
- Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael 618
Bommarito, Ion Androusoopoulos, Daniel Katz, and 619
Nikolaos Aletras. 2022. Lexglue: A benchmark 620
dataset for legal language understanding in english. 621
In *Proceedings of the 60th Annual Meeting of the 622
Association for Computational Linguistics (Volume 623
1: Long Papers)*, pages 4310–4330. 624

625	Daixuan Cheng, Shaohan Huang, and Furu Wei.	Deepseek-r1: Incentivizing reasoning capability in	681
626	2023. Adapting large language models to do-	llms via reinforcement learning. <i>arXiv preprint</i>	682
627	domains via reading comprehension. <i>arXiv preprint</i>	<i>arXiv:2501.12948</i> .	683
628	<i>arXiv:2309.09530</i> .		
629	Gheorghe Comanici, Eric Bieber, Mike Schaekermann,	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	684
630	Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	685
631	cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and	2020. Measuring massive multitask language under-	686
632	1 others. 2025. Gemini 2.5: Pushing the frontier with	standing. <i>arXiv preprint arXiv:2009.03300</i> .	687
633	advanced reasoning, multimodality, long context, and		
634	next generation agentic capabilities. <i>arXiv preprint</i>	Dan Hendrycks, Collin Burns, Anya Chen, and	688
635	<i>arXiv:2507.06261</i> .	Spencer Ball. 2021. Cuad: An expert-annotated	689
		nlp dataset for legal contract review. <i>arXiv preprint</i>	690
		<i>arXiv:2103.06268</i> .	691
636	Jiaxi Cui, Munan Ning, Zongjian Li, Bohua Chen, Yang	Sazzad Hossain, Touhidul Alam Seyam, Avijit Chowd-	692
637	Yan, Hao Li, Bin Ling, Yonghong Tian, and Li Yuan.	hury, Munis Xamidov, Rajib Ghose, and Abhijit	693
638	2023. Chatlaw: A multi-agent collaborative legal	Pathak. 2025. Fine-tuning llama 2 interference: a	694
639	assistant with knowledge graph enhanced mixture-	comparative study of language implementations for	695
640	of-experts large language model. <i>arXiv preprint</i>	optimal efficiency. <i>arXiv preprint arXiv:2502.01651</i> .	696
641	<i>arXiv:2306.16092</i> .		
642	Yongfu Dai, Duanyu Feng, Jimin Huang, Haochen Jia,	Quzhe Huang, Mingxu Tao, Chen Zhang, Zhenwei An,	697
643	Qianqian Xie, Yifang Zhang, Weiguang Han, Wei	Cong Jiang, Zhibin Chen, Zirui Wu, and Yansong	698
644	Tian, and Hao Wang. 2023. Laiw: A chinese legal	Feng. 2023. Lawyer llama technical report. <i>arXiv</i>	699
645	large language models benchmark. <i>arXiv preprint</i>	<i>preprint arXiv:2305.15062</i> .	700
646	<i>arXiv:2310.05620</i> .		
647	Google DeepMind. 2025. Gemini 3 pro. https://deepmind.google/models/gemini/pro/ . Ac-	George F. Jenks. 1967. The data model concept in	701
648	cessed: 2026-01-05.	statistical mapping. <i>International Yearbook of Car-</i>	702
649		<i>tography</i> , 7:186–190.	703
650	Wentao Deng, Jiahuan Pei, Keyi Kong, Zhe Chen, Furu	Paul Kassianik, Baturay Saglam, Alexander Chen,	704
651	Wei, Yujun Li, Zhaochun Ren, Zhumin Chen, and	Blaine Nelson, Anu Vellore, Massimo Aufiero, Fraser	705
652	Pengjie Ren. 2023. Syllogistic reasoning for legal	Burch, Dhruv Kedia, Avi Zohary, Sajana Weeraward-	706
653	judgment analysis. In <i>Proceedings of the 2023 con-</i>	hena, and 1 others. 2025. Llama-3.1-foundationai-	707
654	<i>ference on empirical methods in natural language</i>	securityllm-base-8b technical report. <i>arXiv preprint</i>	708
655	<i>processing</i> , pages 13997–14009.	<i>arXiv:2504.21039</i> .	709
656	Colin Doyle and Aaron D Tucker. 2025. If you give	Katherine Lee, A Feder Cooper, James Grimmelmann,	710
657	an llm a legal practice guide. In <i>Proceedings of</i>	and Daphne Ippolito. 2023. Ai and law: The next	711
658	<i>the 2025 Symposium on Computer Science and Law</i> ,	generation. <i>Available at SSRN 4580739</i> .	712
659	pages 194–205.		
660	Yu Fan, Jingwei Ni, Jakob Merane, Yang Tian, Yoan	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	713
661	Hermstrüwer, Yinya Huang, Mubashara Akhtar, Eti-	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	714
662	enne Salimbeni, Florian Geering, Oliver Dreyer,	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	715
663	and 1 others. 2025. Lexam: Benchmarking le-	täschel, and 1 others. 2020. Retrieval-augmented gen-	716
664	gal reasoning on 340 law exams. <i>arXiv preprint</i>	eration for knowledge-intensive nlp tasks. <i>Advances</i>	717
665	<i>arXiv:2505.12864</i> .	<i>in neural information processing systems</i> , 33:9459–	718
		9474.	719
666	Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou,	Haitao Li, Junjie Chen, Jingli Yang, Qingyao Ai, Wei	720
667	Zhuo Han, Songyang Zhang, Kai Chen, Zongwen	Jia, Youfeng Liu, Kai Lin, Yueyue Wu, Guozhi Yuan,	721
668	Shen, and Jidong Ge. 2023. Lawbench: Benchmark-	Yiran Hu, and 1 others. 2024a. Legalagentbench:	722
669	ing legal knowledge of large language models. <i>arXiv</i>	Evaluating llm agents in legal domain. <i>arXiv preprint</i>	723
670	<i>preprint arXiv:2309.16289</i> .	<i>arXiv:2412.17259</i> .	724
671	Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré,	Haitao Li, You Chen, Qingyao Ai, Yueyue Wu, Ruizhe	725
672	Adam Chilton, Alex Chohlas-Wood, Austin Peters,	Zhang, and Yiqun Liu. 2024b. Lexeval: A compre-	726
673	Brandon Waldon, Daniel Rockmore, Diego Zam-	hensive chinese legal benchmark for evaluating large	727
674	brano, and 1 others. 2023. Legalbench: A collabor-	language models. <i>Advances in Neural Information</i>	728
675	atively built benchmark for measuring legal reason-	<i>Processing Systems</i> , 37:25061–25094.	729
676	ing in large language models. <i>Advances in neural</i>	Chin-Yew Lin. 2004. Rouge: A package for automatic	730
677	<i>information processing systems</i> , 36:44123–44279.	evaluation of summaries. In <i>Text summarization</i>	731
		<i>branches out</i> , pages 74–81.	732
678	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao	Dong Liu and Yuan Gao. 2015. Shanghai courts pioneer	733
679	Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-	case weight coefficient assessment. <i>Wen Hui Bao</i> .	734
680	rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.		

735	Laurens van der Maaten and Geoffrey Hinton. 2008.	State Council of the People’s Republic of China. 2007.	790
736	Visualizing data using t-sne. <i>Journal of machine</i>	Measures for the payment of litigation costs. Ad-	791
737	<i>learning research</i> , 9(Nov):2579–2605.	ministrative regulations governing litigation fees; as	792
		amended.	793
738	National People’s Congress of the PRC. 2020. Civil	Supreme People’s Court of the PRC. China judgments	794
739	code of the people’s republic of china, article 1179.	online. https://wenshu.court.gov.cn/ .	795
740	National People’s Congress of the PRC. 2021. Civil pro-	Supreme People’s Court of the PRC. 2003. Interpre-	796
741	cedure law of the people’s republic of china (article	tation on several issues concerning the application	797
742	160). Legislative Law.	of law in the trial of personal injury compensation	798
		cases.	799
743	OpenAI. 2025. Gpt-4o mini: Advancing cost-efficient	Supreme People’s Court of the PRC. 2016a. 2015 na-	800
744	intelligence. https://openai.com/index/	tional judicial statistical bulletin of chinese courts.	801
745	gpt-4o-mini-advancing-cost-efficient-intelligence/	Accessed: 2025-11-01.	802
746	Accessed: 2026-01-05.	Supreme People’s Court of the PRC. 2016b. Provisions	803
		of the supreme people’s court on the drafting of ju-	804
747	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	dicial documents. Judicial provisions regulating the	805
748	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	structure and content of court judgments.	806
749	Sandhini Agarwal, Katarina Slama, Alex Ray, and 1	Supreme People’s Court of the PRC. 2016c. <i>Work re-</i>	807
750	others. 2022. Training language models to follow in-	port of the supreme people’s court – delivered at the	808
751	structions with human feedback. <i>Advances in neural</i>	fourth session of the 12th national people’s congress	809
752	<i>information processing systems</i> , 35:27730–27744.	on march 13, 2016. Accessed: 2025-11-01.	810
		Supreme People’s Court of the PRC. 2017a. 2016 na-	811
753	Eric A Posner and Shivam Saran. 2025. Judge ai: As-	tional judicial statistical bulletin of chinese courts.	812
754	sessing large language models in judicial decision-	Accessed: 2025-11-01.	813
755	making. <i>University of Chicago Coase-Sandor Insti-</i>	Supreme People’s Court of the PRC. 2017b. <i>Work</i>	814
756	<i>tute for Law & Economics Research Paper</i> , (2503).	report of the supreme people’s court – delivered at	815
		the fifth session of the 12th national people’s congress	816
757	Philippe Prince Tritto and Hiram Ponce. 2025. Assess-	on march 12, 2017. Accessed: 2025-11-01.	817
758	ing ai-generated legal reasoning: A benchmark for	Supreme People’s Court of the PRC. 2018a. 2017 na-	818
759	legal text quality from literature review. In <i>Mexi-</i>	tional judicial statistical bulletin of chinese courts.	819
760	<i>cican Congress on Artificial Intelligence</i> , pages 54–68.	Accessed: 2025-11-01.	820
761	Springer.	Supreme People’s Court of the PRC. 2018b. <i>Work</i>	821
		report of the supreme people’s court – delivered at	822
762	Haosheng Qian, Yixing Fan, Jiafeng Guo, Ruqing	the first session of the 13th national people’s congress	823
763	Zhang, Qi Chen, Dawei Yin, and Xueqi Cheng.	on march 9, 2018. Accessed: 2025-11-01.	824
764	2025. Vericite: Towards reliable citations in retrieval-	Supreme People’s Court of the PRC. 2019a. 2018 na-	825
765	augmented generation via rigorous verification. In	tional judicial statistical bulletin of chinese courts.	826
766	<i>Proceedings of the 2025 Annual International ACM</i>	Accessed: 2025-11-01.	827
767	<i>SIGIR Conference on Research and Development</i>	Supreme People’s Court of the PRC. 2019b. Provisions	828
768	<i>in Information Retrieval in the Asia Pacific Region</i> ,	of the supreme people’s court on evidence in civil	829
769	pages 47–54.	litigation. Judicial provisions governing evidence in	830
		civil litigation.	831
770	Gerard Salton and Christopher Buckley. 1988. Term-	Supreme People’s Court of the PRC. 2019c. <i>Work re-</i>	832
771	weighting approaches in automatic text retrieval. <i>In-</i>	port of the supreme people’s court – delivered at the	833
772	<i>formation processing & management</i> , 24(5):513–	second session of the 13th national people’s congress	834
773	523.	on march 12, 2019. Accessed: 2025-11-01.	835
774	Shandong High People’s Court. 2025. Characteristics	Supreme People’s Court of the PRC. 2020a. 2019 na-	836
775	of road traffic disputes and governance suggestions.	tional judicial statistical bulletin of chinese courts.	837
		Accessed: 2025-11-01.	838
776	Zejiang Shen, Kyle Lo, Lauren Yu, Nathan Dahlberg,	Supreme People’s Court of the PRC. 2020b. Provisions	839
777	Margo Schlanger, and Doug Downey. 2022. Multi-	on the causes of action of civil cases. Judicial pro-	840
778	lexsum: Real-world summaries of civil rights law-	visions establishing the classification system of civil	841
779	suits at multiple granularities. <i>Advances in Neural</i>	causes of action.	842
780	<i>Information Processing Systems</i> , 35:13158–13173.		
781	Marco Siino, Mariana Falco, Daniele Croce, and Paolo		
782	Rosso. 2025. Exploring llms applications in law:		
783	A literature review on current legal nlp approaches.		
784	<i>IEEE Access</i> .		
785	Standing Committee of the National People’s Congress		
786	of the PRC. 2023. Civil procedure law of the people’s		
787	republic of china (2023 amendment). Adopted by		
788	the Standing Committee of the National People’s		
789	Congress; effective January 1, 2024.		

843	Supreme People’s Court of the PRC. 2020c. Work report of the supreme people’s court – delivered at the third session of the 13th national people’s congress on may 25, 2020 . Accessed: 2025-11-01.	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> .	894
844			895
845			896
846			897
847	Supreme People’s Court of the PRC. 2021a. 2020 national judicial statistical bulletin of chinese courts . Accessed: 2025-11-01.	Robert Walters and Marko Novak. 2021. Artificial intelligence and law. In <i>Cyber security, artificial intelligence, data protection & the law</i> , pages 39–69. Springer.	898
848			899
849			900
850	Supreme People’s Court of the PRC. 2021b. Work report of the supreme people’s court – delivered at the fourth session of the 13th national people’s congress on march 8, 2021 . Accessed: 2025-11-01.	Lan Wang and Sufang Qiu. 2019. Measurement of judges’ workload: Econometric models and sichuan experience. <i>Journal of Shanghai Jiao Tong University (Philosophy and Social Sciences)</i> , 27(6):61–73.	901
851			902
852			903
853			904
854	Supreme People’s Court of the PRC. 2022a. 2021 national judicial statistical bulletin of chinese courts . Accessed: 2025-11-01.	xAI. 2025. Grok 4 announcement. https://x.ai/news/grok-4 . Accessed: 2026-01-05.	905
855			906
856			907
857	Supreme People’s Court of the PRC. 2022b. Interpretation of the supreme people’s court on the application of the civil procedure law of the people’s republic of china. Judicial interpretation issued by the Supreme People’s Court.	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	908
858			909
859			910
860			911
861			912
862	Supreme People’s Court of the PRC. 2022c. Interpretation of the supreme people’s court on the application of the civil procedure law of the people’s republic of china (article 257). Judicial Interpretation.	Shunyu Yao, Qingqing Ke, Qiwei Wang, Kangtong Li, and Jie Hu. 2024. Lawyer gpt: A legal large language model with enhanced domain knowledge and reasoning capabilities. In <i>Proceedings of the 2024 3rd International Symposium on Robotics, Artificial Intelligence and Information Engineering</i> , pages 108–112.	913
863			914
864			915
865			916
866	Supreme People’s Court of the PRC. 2022d. Work report of the supreme people’s court – delivered at the fifth session of the 13th national people’s congress on march 8, 2022 . Accessed: 2025-11-01.	Shengbin Yue, Wei Chen, Siyuan Wang, Bingxuan Li, Chenchen Shen, Shujun Liu, Yuxuan Zhou, Yao Xiao, Song Yun, Xuanjing Huang, and 1 others. 2023. Disc-lawllm: Fine-tuning large language models for intelligent legal services. <i>arXiv preprint arXiv:2309.11325</i> .	917
867			918
868			919
869			920
870	Supreme People’s Court of the PRC. 2023a. 2022 national judicial statistical bulletin of chinese courts . Accessed: 2025-11-01.	Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, and 1 others. 2025. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models. <i>arXiv preprint arXiv:2508.06471</i> .	921
871			922
872			923
873	Supreme People’s Court of the PRC. 2023b. Work report of the supreme people’s court – delivered at the first session of the 14th national people’s congress on march 7, 2023 . Accessed: 2025-11-01.	Kepu Zhang, Weijie Yu, Sunhao Dai, and Jun Xu. 2024. Citalaw: Enhancing llm with citations in legal domain. <i>arXiv preprint arXiv:2412.14556</i> .	924
874			925
875			926
876			927
877	Supreme People’s Court of the PRC. 2024a. 2023 national judicial statistical bulletin of chinese courts . Accessed: 2025-11-01.	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.	928
878			929
879			930
880	Supreme People’s Court of the PRC. 2024b. Work report of the supreme people’s court – delivered at the second session of the 14th national people’s congress on march 8, 2024 . Accessed: 2025-11-01.	Lucia Zheng, Neel Guha, Javokhir Arifov, Sarah Zhang, Michal Skreta, Christopher D Manning, Peter Henderson, and Daniel E Ho. 2025. A reasoning-focused legal retrieval benchmark. In <i>Proceedings of the 2025 Symposium on Computer Science and Law</i> , pages 169–193.	931
881			932
882			933
883			934
884	Supreme People’s Court of the PRC. 2025. 2024 national judicial statistical bulletin of chinese courts . Accessed: 2025-11-01.		935
885			936
886			937
887	Harry Surden. 2018. Artificial intelligence and law: An overview. <i>Ga. St. UL Rev.</i> , 35:1305.		938
888			939
889	Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, and 1 others. 2025. Kimi k2: Open agentic intelligence. <i>arXiv preprint arXiv:2507.20534</i> .		940
890			941
891			942
892			943
893			944

948 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang,
949 Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen,
950 and Nan Duan. 2024. Agieval: A human-centric
951 benchmark for evaluating foundation models. In
952 *Findings of the Association for Computational Lin-*
953 *guistics: NAACL 2024*, pages 2299–2314.

A Increasing Caseload in China

Data extracted from the annual *Work Report of the Supreme People’s Court* and the *National Judicial Statistical Bulletin of Chinese Courts* reveals a consistent and significant upward trend in the number of cases handled by Chinese courts over the past decade. This growth has been accompanied by severe pressure arising from the "judge–case conflict," namely the structural imbalance between the number of judges and the rapidly increasing caseload.

From 2015 to 2019, the number of first-instance civil cases rose steadily from 6.228 million to 9.393 million, reflecting the expansion of disputes in civil domains such as family, labor, and consumer protection (Supreme People’s Court of the PRC, 2016c, 2017b, 2018b, 2019c, 2020c).

Table 4: First-instance Civil Cases in China (2015–2019).

Year	Number of Cases (10,000 cases)
2015	62.28
2016	67.38
2017	82.95
2018	90.17
2019	93.93

Note: Data for years after 2019 is not available in the Work Report of the Supreme People’s Court. Data labeled as "first-instance civil cases" in the National Judicial Statistical Bulletin of Chinese Courts actually refers to civil and commercial cases and therefore cannot be used as a direct supplement.

For first-instance civil and commercial cases, which cover a broader range of business-related disputes, the increase was even more pronounced. The caseload grew from 9.575 million in 2015 to 18.237 million in 2024, nearly doubling over a ten-year period (Supreme People’s Court of the PRC, 2016a, 2017a, 2018a, 2019a, 2020a, 2021a, 2022a, 2023a, 2024a, 2025).

The rapid expansion of case volume has directly translated into a heavier workload for judges. The annual average number of cases handled per judge increased from 187 in 2017 to 357 in 2023. Intermediate statistics indicate a persistent upward trajectory, with averages of 225 cases in 2020, 238 in 2021, and 242 in 2022 (Supreme People’s Court of the PRC, 2021b, 2022d, 2023b, 2024b).

The Supreme People’s Court has explicitly ac-

Table 5: First-instance Civil and Commercial Cases in China (2015–2024)

Year	Number of Cases (10,000 cases)
2015	95.75
2016	107.64
2017	116.51
2018	124.35
2019	139.30
2020	133.06
2021	157.46
2022	161.14
2023	174.77
2024	182.37

Table 6: Annual Average Cases Handled per Judge in China (2017, 2020–2023).

Year	Average Cases per Judge
2017	187
2020	225
2021	238
2022	242
2023	357

knowledged that "the judge–case" conflict has become increasingly prominent" underscoring the mounting pressure faced by the judicial workforce (Supreme People’s Court of the PRC, 2024b).

B Data Construction and Filtering

Table 7 presents detailed statistics of field-level completeness after structured information extraction from the raw legal documents. Each row corresponds to a specific field required for constructing or characterizing cases in JurisBench, and the missing rate reflects the proportion of documents in which the corresponding field cannot be reliably extracted. All raw legal documents used in this study are sourced from judicial decisions published in 2024 or later.

The results show that the overall completeness of extracted legal documents is limited. While several meta-level attributes, such as case name and *Cause of Action*, exhibit negligible missing rates, many core litigation fields suffer from substantial absence. In particular, high-level *Rationale of the Judgment* components, including *Focus of Disputes* and *Facts Proved by the Court*, are missing in a large portion of documents, reflecting both the heterogeneity of judicial writing styles and the inherent difficulty of automatically identifying

1021 fine-grained legal reasoning elements. As a conse- 1070
1022 quence, only a relatively small subset of documents 1071
1023 contains all required fields and can be considered 1072
1024 fully usable for benchmark evaluation. 1073

1025 It is important to note that several fields listed un- 1074
1026 der *Basic Case Info*, such as *Trial Procedure*, *Case* 1075
1027 *Number*, and *Closing Date*, are not directly used as 1076
1028 inputs or targets in JurisBench tasks. Nevertheless, 1077
1029 we treat these attributes as mandatory metadata and 1078
1030 exclude documents in which they are missing. This 1079
1031 design choice is motivated by long-term bench- 1080
1032 mark reliability: missing basic case information 1081
1033 may introduce inconsistencies in dataset statistics, 1082
1034 hinder reproducibility, and complicate future exten- 1083
1035 sions involving temporal analysis, procedural strat- 1084
1036 ification, or cross-case aggregation. By enforcing 1085
1037 completeness even for non-task fields, we ensure 1086
1038 that all retained cases are well-formed, uniquely 1087
1039 identifiable, and suitable for systematic statistical 1088
1040 analysis. 1089

1041 Overall, this filtering process results in 1090
1042 $N_{usable} = 3,204$ fully structured cases, yielding 1091
1043 a ratio of approximately 19.13 incomplete docu- 1092
1044 ments for every usable one. (We ultimately se- 1093
1045 lected 1000 representative cases as data for pipeline 1094
1046 testing.) Although this substantially reduces the 1095
1047 dataset size, it provides a strong guarantee on data 1096
1048 integrity and supports reliable, interpretable, and 1097
1049 extensible evaluation of large language models in 1098
1050 the legal domain. 1099

1051 C Normative Foundations of the 1100 1052 JurisBench Pipeline 1101

1053 The JurisBench pipeline is normatively grounded 1102
1054 in the procedural framework of Chinese civil liti- 1103
1055 gation. All directed arrows in the benchmark cor- 1104
1056 respond to legally mandated or institutionally es- 1105
1057 tablished dependencies in judicial practice. Rather 1106
1058 than representing transitions between isolated case 1107
1059 elements, these arrows denote ordered transitions 1108
1060 between sub-benchmarks, each of which operates 1109
1061 on a structured and composite input set. We distin- 1110
1062 guish between **inter-benchmark (macro) transi-** 1111
1063 **tions**, which connect adjacent sub-benchmarks in 1112
1064 the pipeline, and **intra-benchmark (micro) tran-** 1113
1065 **sitions**, which decompose reasoning steps within a 1114
1066 single sub-benchmark. 1115

1067 C.1 Inter-Benchmark (Macro) Transitions 1116

1068 **Phase 1 (Cause of Action) → Phase 2 (Focus of** 1117
1069 **Disputes).** The transition from Phase 1 to Phase 2 1118

1070 reflects the procedural dependency between legal 1071
1072 characterization and dispute identification in Chi- 1073
1074 nese civil litigation. Phase 1 determines the *Cause* 1074
1075 *of Action* based on a composite input consisting of 1075
1076 *Basic Case Info* and *Plaintiff’s Argument*, corre- 1076
1077 sponding to the statutory case-filing stage. Under 1077
1078 the Civil Procedure Law of the People’s Republic 1078
1079 of China, a civil action must be initiated through 1079
1080 a written complaint specifying claims, facts, and 1080
1081 reasons, which constitute the legally required ba- 1081
1082 sis for determining the legal nature of the dispute 1082
1083 at docketing ([Standing Committee of the National](#) 1083
1084 [People’s Congress of the PRC, 2023](#)). This de- 1084
1085 termination governs subsequent trial organization 1085
1086 and adjudication, including judicial division assign- 1086
1087 ment. 1087

1088 Phase 2 operates on an expanded input set that 1088
1089 includes *Basic Case Info*, *Plaintiff’s Argument*, *De-* 1089
1090 *fendant’s Argument*, and the *Cause of Action* output 1090
1091 from Phase 1, and evaluates the identification of *Fo-* 1091
1092 *cus of Disputes* under adversarial conditions. *Judi-* 1092
1093 *cial Interpretations* authorize courts to summarize 1093
1094 dispute foci during pretrial proceedings and require 1094
1095 adjudication to proceed around such foci ([Supreme](#) 1095
1096 [People’s Court of the PRC, 2022b, 2019b](#)). Since 1096
1097 dispute identification must be framed within the 1097
1098 legally established nature of the case, the arrow 1098
1099 from Phase 1 to Phase 2 denotes the normative 1099
1100 requirement that *Focus of Disputes* identification 1100
1101 be conditioned on a prior determination of *Cause* 1101
1102 *of Action*, rather than treated as an unconstrained 1102
1103 issue-extraction task. 1103

1104 **Phase 2 (Focus of Disputes) → Phase 3 (Ra-** 1102
1105 **tionale of the Judgment).** The transition from 1103
1106 Phase 2 to Phase 3 corresponds to the statutory or- 1104
1107 dering between dispute identification and *Rationale* 1105
1108 *of the Judgment*. Phase 2 produces *Focus of Dis-* 1106
1109 *putes* based on a composite input integrating *Basic* 1107
1110 *Case Info*, *Plaintiff’s Argument*, *Defendant’s Argu-* 1108
1111 *ment*, and *Cause of Action*, reflecting the adversar- 1109
1112 ial identification of contested issues. Civil adjudi- 1110
1113 cation is required to proceed around disputed facts 1111
1114 and issues, and judgments must explicitly address 1112
1115 these disputes through reasoned analysis ([Standing](#) 1113
1116 [Committee of the National People’s Congress of](#) 1114
1117 [the PRC, 2023](#)). 1115

1116 Phase 3 evaluates *Rationale of the Judgment* un- 1116
1117 der a further expanded information set consisting of 1117
1118 *Basic Case Info*, *Plaintiff’s Argument*, *Defendant’s* 1118
1119 *Argument*, *Cause of Action*, and *Facts Proved by the* 1119
1120 *Court*. Although *Focus of Disputes* is not treated 1120

Table 7: **Data Filtering Ratio Statistics.** This table reports field-level missing statistics during the structured extraction of raw legal documents, illustrating substantial variation in completeness across different types of information. After enforcing strict completeness requirements, only 3,204 cases contain all required fields and are retained as usable data.

Category	Field Name	Missing	Total	Missing Rate (%)
Meta Information	<i>Meta Information</i>	0	64,504	0.00
	<i>Party</i>	11,736	64,504	18.19
Core Litigation	<i>Plaintiff’s Argument</i>	812	64,504	1.26
	<i>Defendant’s Argument</i>	8,573	64,504	13.29
	<i>Focus of Disputes</i>	55,480	64,504	86.01
	<i>Facts Proved by the Court</i>	18,474	64,504	28.64
	<i>Rationale of the Judgment</i>	3,894	64,504	6.04
	<i>Result of the Judgment</i>	1,068	64,504	1.66
Basic Case Info	<i>Court of Acceptance</i>	505	64,504	0.78
	<i>Trial Procedure</i>	107	64,504	0.17
	<i>Case Name</i>	0	64,504	0.00
	<i>Case Number</i>	1,253	64,504	1.94
	<i>Cause of Action</i>	0	64,504	0.00
	<i>Closing Date</i>	845	64,504	1.31

Note: Ratio of incomplete to usable documents ($N_{usable} = 3,204$) is 19.13:1.

as an explicit input variable, it is institutionally embedded in the task design as a binding contextual constraint. The arrow from Phase 2 to Phase 3 therefore represents the legal requirement that *Rationale of the Judgment* respond to previously identified dispute foci rather than being generated independently of the adversarial issue structure.

Phase 3 (*Rationale of the Judgment*) → Phase 4 (*Result of the Judgment*). The transition from Phase 3 to Phase 4 reflects the formal structure of civil judgments. Phase 3 assesses the model’s ability to generate *Rationale of the Judgment* grounded in established facts and applicable law, based on a composite input consisting of *Basic Case Info*, *Plaintiff’s Argument*, *Defendant’s Argument*, *Cause of Action*, and *Facts Proved by the Court*.

Phase 4 operates on an augmented input set that incorporates the *Rationale of the Judgment* output of Phase 3 in addition to *Basic Case Info*, *Plaintiff’s Argument*, *Defendant’s Argument*, *Cause of Action*, and *Facts Proved by the Court*, and evaluates whether the model can derive *Result of the Judgment* that is procedurally and substantively consistent with the preceding reasoning. Statutory provisions require civil judgments to include both the reasoning and the adjudicative outcome, rendering *Result of the Judgment* the logical culmination of *Rationale of the Judgment* ([Standing Committee of the National People’s Congress of the PRC,](#)

2023). Accordingly, the arrow from Phase 3 to Phase 4 denotes the legally mandated dependency between *Rationale of the Judgment* and *Result of the Judgment*.

C.2 Intra-Benchmark (Micro) Transitions

Internal Structure of Phase 1 (*Cause of Action*).

Within Phase 1, the hierarchical progression across levels of *Cause of Action* reflects the multi-level classification system established by the Provisions on the *Causes of Action of Civil Cases* ([Supreme People’s Court of the PRC, 2020b](#)). This internal structure supports standardized case filing and adjudication by progressively refining the legal characterization of disputes, while remaining anchored in the same composite input of *Basic Case Info* and *Plaintiff’s Argument*.

Internal Structure of Phase 2 (*Focus of Disputes*).

Within Phase 2, *Focus of Disputes* identification presupposes the joint consideration of *Plaintiff’s Argument* and *Defendant’s Argument* as part of its composite input. This reflects the adversarial structure of civil litigation, under which courts must hear arguments from both parties before determining contested issues ([Standing Committee of the National People’s Congress of the PRC, 2023](#)). The internal reasoning of Phase 2 therefore evaluates the synthesis and reconciliation of competing factual claims rather than unilateral issue extraction.

Internal Structure of Phase 3 (*Rationale of the Judgment*). Phase 3 decomposes *Rationale of the Judgment* into multiple normative steps, including grounding in *Facts Proved by the Court*, application of applicable law, citation of specific legal Articles, and synthesis of the written rationale. This internal decomposition is normatively supported by statutory principles requiring adjudication to be based on established facts and governed by law, as well as judicial drafting rules mandating explicit articulation of facts, *Legal Grounds*, and reasoning with accurate statutory citations ([Standing Committee of the National People’s Congress of the PRC, 2023](#); [Supreme People’s Court of the PRC, 2016b](#)).

Internal Structure of Phase 4 (*Result of the Judgment*). Within Phase 4, the determination of litigation costs precedes the assessment of appeal jurisdiction, reflecting auxiliary procedural regulations governing litigation fees and appellate rights ([State Council of the People’s Republic of China, 2007](#); [Standing Committee of the National People’s Congress of the PRC, 2023](#)). These internal steps jointly constitute the adjudicative outcome and are integrated into the finalized *Result of the Judgment* in accordance with statutory rules on the composition and effectiveness of civil judgments ([Standing Committee of the National People’s Congress of the PRC, 2023](#)).

D Semantic Diversity and Dataset Distribution

To evaluate the representational capacity and semantic diversity of the **JurisBench** dataset, we performed a high-dimensional feature analysis followed by manifold learning-based visualization. Specifically, we first extracted representative textual features from the judicial documents using a *Term Frequency-Inverse Document Frequency* (TF-IDF) vectorization scheme ([Salton and Buckley, 1988](#)). We incorporated both unigrams and bigrams to capture fine-grained legal terminology and structural markers. Given the high dimensionality of the resulting feature space, we applied Principal Component Analysis (PCA) to reduce the latent space to 50 dimensions, followed by *t*-distributed Stochastic Neighbor Embedding (*t*-SNE) ([Maaten and Hinton, 2008](#)) for non-linear dimensionality reduction into a two-dimensional visual manifold.

The semantic distribution across the three difficulty tiers—Simple ($n = 394$), Intermediate ($n = 394$), and Complex ($n = 212$)—is illustrated

in Figure 7. The visualization reveals an extensive and relatively uniform distribution across the latent space, avoiding narrow or isolated clusters. This suggests that **JurisBench** covers a broad and diverse semantic landscape within the legal domain, rather than being restricted to a few repetitive scenarios.

Quantitatively, our diversity analysis yields an average pairwise distance of 1.3603 ($\sigma = 0.0631$), which is remarkably close to the theoretical maximum distance of 1.414 in a normalized TF-IDF vector space. This high degree of sparsity and semantic spread indicates that the dataset possesses low redundancy and high task variance. Notably, the intermingling of simple, intermediate, and complex cases in the *t*-SNE plot demonstrates that case complexity in our benchmark is driven by procedural depth and reasoning logic rather than mere lexical variations or superficial keyword distributions, thereby validating the structural integrity of our difficulty stratification.

E Complexity Indicators and Stratification Procedure

The complexity stratification of JurisBench test cases is grounded in both statutory standards and empirical studies on case difficulty in Chinese judicial practice. Each case is quantitatively evaluated along five core dimensions: *Parties Involved*, *Claims*, *Focus of Disputes*, *Rationale of the Judgment*, and *Result of the Judgment*.

In the *Parties Involved* dimension, multi-party participation is treated as a primary driver of complexity, as it often entails intertwined legal relationships and challenges in liability allocation ([Supreme People’s Court of the PRC, 2022c](#); [Shandong High People’s Court, 2025](#); [Beijing Chaoyang District People’s Court, 2025](#)). In the *Claims* and *Focus of Disputes* dimensions, compound claims and intensive disputes—such as disagreements over appraisal conclusions or jurisdictional objections—are indicative of increased trial difficulty. This aligns with procedural standards that restrict simplified procedures to cases with clear rights and limited disputes ([National People’s Congress of the PRC, 2021](#)).

In the *Rationale of the Judgment* and *Result of the Judgment* dimensions, the length and structural depth of judgment documents are used as proxy variables for judicial effort and reasoning complexity. Prior empirical studies and case-weight reform

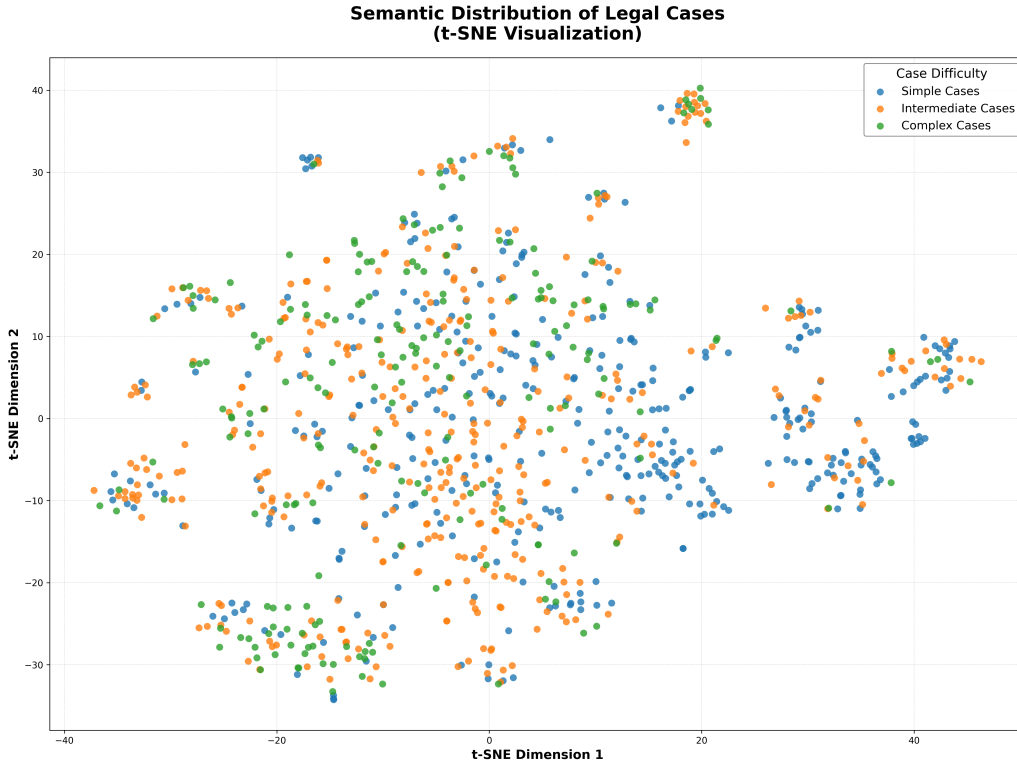


Figure 7: Semantic distribution of legal cases in JurisBench. The t -SNE visualization shows the high-dimensional embedding space of cases across three difficulty levels. The broad spread and high average pairwise distance (1.3603) reflect the semantic diversity and low redundancy of the dataset.

practices have shown that more extensive reasoning is strongly correlated with increased fact-finding difficulty and legal analysis depth (Liu and Gao, 2015; Wang and Qiu, 2019).

After scoring cases across all dimensions, the Jenks Natural Breaks Optimization method (Jenks, 1967) is applied to identify two optimal breakpoints in the complexity score distribution. This method minimizes within-group variance while maximizing between-group variance, enabling an objective partition of cases into *Simple*, *Intermediate*, and *Complex* subsets with coherent legal characteristics and clearly differentiated difficulty levels.

F Additional Experiments

F.1 General Legal Taxonomy Alignment (Phase 1-a & 1-b)

To verify the generalizability of LLMs across the full spectrum of Chinese civil law, we evaluated 10 representative models on a broad dataset comprising 2,820 cases across 282 *Causes of Action*. This experiment investigates the models’ performance when provided with **hierarchical taxonomy guidance** (Phase 1-a) versus their purely **parameterized recall** in a open-ended setting (Phase 1-b).

Table 8: Results of General Legal Taxonomy Alignment on 282 *Causes of Action* (2,820 cases). Values in **bold** and underline indicate the best and second-best performance per column.

Model	Phase 1-a \uparrow	Phase 1-b \uparrow	Δ (Drop %)
<i>Multilingual Models</i>			
Gemini-3 Pro	0.7128	0.3430	-51.88%
Claude 4.5	0.6071	<u>0.0699</u>	-88.49%
GPT-4o-mini	0.5032	0.0011	-99.78%
Grok 4	0.6114	0.0096	-98.43%
Llama 3.3 70B	<u>0.6233</u>	0.0024	-99.62%
<i>Chinese-Oriented Models</i>			
Qwen3-Instruct	0.3766	0.0167	-95.57%
Doubao-1.5-Pro	0.6170	0.0216	-96.50%
GLM-4.5	0.4043	0.0085	-97.90%
DeepSeek-R1	0.6092	0.0621	-89.81%
Kimi-K2	0.6227	0.0043	-99.31%

The comparative results in Table 8 yield the following critical insights:

- **The Open-Ended Collapse:** We observe a catastrophic performance degradation from Phase 1-a (taxonomy-guided) to Phase 1-b (open-ended) across almost all tested models. Excluding Gemini-3, the average accuracy drop exceeds 90%. This suggests that

while LLMs can effectively utilize a provided taxonomy as a reference for recognition, they lack the internal parameterized knowledge required to generate precise, hierarchical legal terms from scratch. In the open-ended setting, models frequently hallucinate non-existent causes of action or fail to conform to the standard 282-*Cause of Action* nomenclature.

- **Gemini-3’s Outlier Performance:** Gemini-3 Pro exhibits a significant advantage in the open-ended setting, achieving 0.3430 accuracy—nearly an order of magnitude higher than its peers. This indicates a superior alignment with Chinese legal taxonomies during its pre-training or instruction-tuning phase, enabling it to maintain structural fidelity without external taxonomy guidance.
- **Recall vs. Recognition:** Most Chinese-oriented models (e.g., Doubao-1.5-Pro, Kimi-K2) show strong "recognition" capabilities in Phase 1-a, matching the performance of top-tier global models like Llama 3.3 70B. However, their "recall" in Phase 1-b remains fragile. This confirms that localized pre-training primarily enhances the ability to process and select from domain-specific context rather than solving the "term recall" bottleneck in isolation.

F.2 One-Shot Performance Analysis

Overall, we observe that one-shot prompting exerts a relatively marginal influence across the majority of sub-benchmarks, with performance deltas frequently remaining near zero. To further investigate the impact of In-Context Learning (ICL) on the JurisBench pipeline, we conducted one-shot experiments on a randomly sampled subset of 100 cases. For each sub-benchmark, a single representative legal case with its corresponding "gold" reasoning and results was provided in the prompt as a demonstration.

As illustrated in Table 9, while the overall shifts are limited, the introduction of a single demonstration yields localized but significant gains for certain models. Notably, multilingual models like Llama 3.3 70B and Grok 4 exhibit substantial improvements in structural alignment (Phase 1) and judgment result determination (Phase 4-2), suggesting that ICL effectively helps these models adapt to the specific formatting requirements of Chinese

judicial documents. However, we also observe a "sensitivity trap" in certain Chinese-oriented models, such as Qwen3-Instruct, which showed a sharp decline in Phase 3-4. This suggests that while one-shot prompting can clarify task constraints, it may also introduce biased priors that interfere with the internal legal reasoning of models already heavily aligned with domestic legal corpora.

F.3 Experiment of Law-specific Models

Preliminary Evaluation on Domain-Specific LLMs. We conducted a pilot study on several representative legal-specific models, including DISC-LawLLM (Yue et al., 2023), AdaptLLM/Law-Chat (Cheng et al., 2023), and Fuzi-Mingcha (Deng et al., 2023), using a subset of 100 cases of the Simple difficulty subset. The results revealed significant challenges in these models’ stability across our multi-phase tasks. While they could partially complete basic tasks (e.g., Phase 1, 2, and 4-3), more than 50% of the outputs in remaining tasks exhibited severe **unfaithful generation**. Specifically, we observed frequent **instruction following failures** and symptoms of **model drift**, where models optimized for specific legal corpora failed to adapt to the structured, multi-step reasoning required by our benchmark.

Furthermore, these models often suffered from **textual degeneration** and semantic collapse. For instance, in Phase 3-2 (Legal Grounds-Law), Fuzi-Mingcha frequently failed to distinguish between legal entities and case evidence. Instead of outputting standardized legal categories, it generated a disorganized list of trial-related artifacts such as “xxx Hospital Discharge Records” and “Inpatient Bills” within the predicted legal name fields. Due to this high rate of unreliable outputs and the inability to maintain logical consistency, we opted to exclude these models from the full-scale large-scale evaluation to ensure the integrity of our comparative analysis, focusing instead on models that demonstrate sufficient functional stability.

F.4 Keyword-based Evaluation of Generative Tasks

To further validate the substantive legal precision of model outputs beyond linguistic fluency, we introduced a **Keyword-based Lexical Accuracy** metric for the four generative phases: *Focus of Disputes* (Phase 2), *Factual Basis* (Phase 3-1), *Rationale of the Judgment* (Phase 3-4), and *Result of the Judgment* (Phase 4-2).

Table 9: Experimental results (one-shot) of 10 LLMs.

Group	Model	P1	P2	Phase 3: <i>Rationale of the Judgment</i>				Phase 4: <i>Result of the Judgment</i>		
		1	2	3-1	3-2	3-3	3-4	4-1	4-2	4-3
Multilingual	Gemini-3 Pro	+0.000	+0.085	-0.006	+0.007	+0.036	+0.021	-0.002	+0.047	+0.124
	Claude 4.5	+0.209	+0.060	+0.025	+0.043	+0.101	+0.006	+0.016	+0.011	+0.014
	GPT-4o-mini	+0.010	-0.009	+0.006	+0.007	+0.003	+0.017	+0.005	+0.034	+0.006
	Grok 4	+0.267	+0.244	-0.019	+0.188	-0.114	+0.036	-0.026	+0.254	+0.030
	Llama 3.3 70B	+0.400	-0.031	-0.018	+0.176	+0.023	+0.035	-0.047	+0.446	-0.104
Chinese	Qwen3-Instruct	+0.000	+0.070	-0.100	+0.373	+0.027	-0.593	+0.183	+0.067	+0.050
	Doubao-1.5-Pro	-0.031	+0.046	+0.005	+0.175	+0.013	+0.049	-0.021	+0.255	+0.124
	GLM-4.5	+0.188	-0.027	-0.003	+0.045	+0.017	+0.029	+0.025	+0.204	+0.094
	DeepSeek-R1	+0.030	+0.082	-0.025	+0.004	-0.296	-0.011	+0.044	+0.055	+0.064
	Kimi-K2	-0.099	+0.016	-0.035	+0.003	-0.342	-0.011	+0.005	+0.017	+0.093

Comparative Results. Table 10 presents the comparison between the original semantic metrics (reported in the main text) and the keyword-based lexical metrics.

Key Observations. The comparison reveals a sharp **Fluency-Precision Gap**: while models achieve high semantic scores (> 0.9), keyword recovery collapses in complex reasoning phases (≈ 0.2). Notably, Chinese-oriented models consistently outperform their multilingual counterparts in keyword recall, demonstrating higher *terminological fidelity*. While multilingual models often maintain semantic correctness, they tend to employ general descriptions rather than the precise Chinese legal nomenclature required for professional adjudication. Qwen3-Instruct’s consistent lead highlights the efficacy of localized pre-training for professional grounding.

G Human Subjects, Annotation Process, and Compensation

This study involves human participation in dataset construction and model evaluation. All procedures were designed to ensure objectivity, role separation, and auditability, while minimizing subjective bias.

G.1 Compensation and Labor Remuneration

All participants received fixed compensation independent of task outcomes or model performance. Student annotators were compensated at a rate of 100 RMB per person, reflecting standard research assistance rates. Judicial professionals were compensated at 200 RMB per person, reflecting their professional expertise and time commitment. No performance-based incentives were provided, and compensation was not contingent on specific annotation content or evaluation results.

G.2 Dataset Construction (Student Annotators)

A total of 9 student annotators (demographics in Table 11) were divided into two functionally independent groups to ensure data integrity through a "construction-verification" pipeline.

Roles and Process. Six students formed the **Construction Group**, responsible for anonymizing personal identifiers, structuring raw judgments into the JurisBench schema, and screening outcome-relevant keywords (e.g., liability attribution, causal relations). Three students formed the **Verification Group**, who independently reviewed the structured outputs for consistency and completeness. Discrepancies were resolved through predefined consistency rules via a double-blind process.

Guidelines. Annotators were strictly prohibited from introducing interpretative judgments or reconstructed reasoning. Structured fields were required to reflect explicit content from the source judgment. Anonymization was limited to removing personal identifiers without altering legally relevant facts.

G.3 Human Evaluation and Validation (Judicial Professionals)

Fifteen judicial professionals (demographics in Table 12) participated in evaluating model performance and validating our automated metrics.

Evaluation Setting. Judicial professionals acted as judicial assistants under a strict blind-review protocol. They were not informed whether an answer was produced by a human or an LLM. All answers were presented in a uniform format to eliminate stylistic or metadata cues. Evaluators assessed model predictions for the *Focus of Disputes* task under two conditions: with and without access

Table 10: Comparison between Original Semantic Metrics (Main) and Keyword-based Metrics (KW) across generative phases. “Main” refers to Embedding-based Similarity & LLM-as-a-Judge scores, while “KW” indicates the keyword recovery rate. All scores are normalized to $[0, 1]$. Values in **bold** and underline indicate the best and second-best performance in each **KW** column.

Group	Model	Phase 2 (Focus)		Phase 3-1 (Facts)		Phase 3-4 (Reasoning)		Phase 4-2 (Result)	
		Main	KW	Main	KW	Main	KW	Main	KW
Multiling.	Gemini-3 Pro	0.684	0.459	0.917	0.558	0.919	0.202	0.693	0.269
	Claude 4.5	0.644	0.401	0.901	0.455	0.938	0.155	0.446	0.159
	GPT-4o-mini	0.624	0.419	0.897	0.502	0.909	0.174	0.593	0.209
	Grok 4	0.681	0.521	0.901	0.424	0.926	0.088	0.711	0.192
	Llama 3.3 70B	0.782	0.559	0.917	0.540	0.923	0.165	0.604	0.151
Chinese	Qwen3-Instruct	0.798	0.656	0.892	0.639	0.924	0.286	0.710	<u>0.225</u>
	Doubao-1.5-Pro	0.679	0.554	0.909	<u>0.620</u>	0.909	0.178	0.700	0.195
	GLM-4.5	0.738	<u>0.577</u>	0.912	0.556	0.922	<u>0.223</u>	<u>0.776</u>	0.209
	DeepSeek-R1	0.683	0.507	0.927	0.596	<u>0.927</u>	0.216	0.480	0.191
	Kimi-K2	0.750	0.571	<u>0.904</u>	0.493	0.930	0.205	0.807	0.219

Table 11: Demographic Information of Student Annotators

Student ID	Age	Gender	Work Location
Student 1	23	Male	Beijing
Student 2	23	Male	Beijing
Student 3	25	Female	Beijing
Student 4	28	Male	Beijing
Student 5	28	Male	Beijing
Student 6	28	Female	Beijing
Student 7	35	Female	Beijing
Student 8	29	Male	Beijing
Student 9	24	Male	Beijing

Table 12: Demographic Information of Judicial Professionals

Judge ID	Age	Gender	Work Location
Judge 1	42	Female	Beijing
Judge 2	36	Female	Beijing
Judge 3	33	Female	Beijing
Judge 4	36	Female	Beijing
Judge 5	45	Female	Beijing
Judge 6	45	Male	Beijing
Judge 7	27	Male	Shandong
Judge 8	37	Female	Shandong
Judge 9	32	Female	Shandong
Judge 10	27	Female	Shandong
Judge 11	36	Female	Shandong
Judge 12	28	Male	Shandong
Judge 13	29	Male	Shandong
Judge 14	34	Male	Shandong
Judge 15	28	Male	Shandong

to the "gold" reference extracted from the original document.

Judicial Evaluation Protocol and Human-Metric Alignment. To validate our automated *LLM-as-a-Judge* framework (which provides continuous scores in $[0, 1]$), we conducted a human-in-the-loop study on Phase 2 (*Focus of Disputes*). We adopted a discrete seven-point ordinal scale for judges to better mirror practical judicial decision-making: **0**: Completely incorrect; **1**: Extensive errors, unusable; **2**: Substantial errors, no meaningful assistance; **3**: Approximately half correct; **4**: Mostly correct, requiring minor corrections; **5**: Sufficiently accurate for adjudication; **5+**: Functionally equivalent to high-quality professional output.

Consistency and Validity. To ensure objectivity, we selected the answers of the tested LLMs of five representative cases for a cross-validation study, where each case was evaluated by five independent judges (25 total assessments). Evaluators

focused on substantive legal alignment rather than surface-level fluency. We then mapped the continuous LLM ratings onto this ordinal scale and performed a concordance analysis. The results demonstrate high consistency between the judges’ consensus and the automated scores, confirming that our framework effectively serves as a reliable proxy for professional-level judicial reasoning. Judges treated the reference *Focus of Disputes* as the normative baseline and assessed each case in isolation to prevent cross-referencing bias.

1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510

1511 **H Written Guidelines for Student**
1512 **Annotators**

1513 This section provides the complete and operational
1514 written guidelines for student annotators. The ob-
1515 jective of this task is to transform raw judicial doc-
1516 uments into a high-fidelity and structured dataset
1517 through two primary processes: **Data Structuring**
1518 and **Outcome-Relevant Keyword Labeling**.

1519 **H.1 Role Definition and Workflow**

1520 Student annotators act as technical data extractors.
1521 Your objective is to perform a high-fidelity map-
1522 ping of the judicial text into a structured schema.

- 1523 • **Verbatim Extraction:** You must use the exact
1524 wording from the judgment. You are strictly
1525 prohibited from performing any paraphrasing
1526 or summarizing.
- 1527 • **No Linguistic Cleaning:** If the phrasing of
1528 the judge is awkward or grammatically imper-
1529 fect, you must preserve it exactly as it appears.
- 1530 • **Mental Firewall:** You must disregard your
1531 personal knowledge of the law. If a fact is
1532 missing in the text, it must be recorded as *Not*
1533 *Specified*.

1534 **H.2 Task I: Verbatim Data Structuring**

1535 Annotators must decompose the judgment into
1536 seven specific fields using a **Key-Value pair** for-
1537 mat. The *Key* represents the attribute of the infor-
1538 mation segment, and the *Value* is the verbatim text
1539 extracted from the document.

1540 **The Seven Mandatory Fields:**

- 1541 1. **Basic Case Info Information:** This includes
1542 metadata such as the Case Name, the *Cause of*
1543 *Action*, the *Court of Acceptance*, the *Parties*
1544 and their Procedural Roles, the Procedural
1545 Posture, and the Date of Case Closure.
- 1546 2. **Plaintiff's Arguments:** This includes the
1547 claims, the factual grounds, and the litigation
1548 requests explicitly raised by the plaintiff or
1549 plaintiffs.
- 1550 3. **Defendant's Arguments:** This includes the
1551 responses, the defenses, the objections, or the
1552 admissions explicitly raised by the defendant
1553 or defendants.
- 1554 4. **Facts Proved by the Court:** This consists of
1555 the objective facts and evidence determined
1556 and verified by the court.

5. **Focus of Disputes:** This identifies the core
1557 points of contention or the disputed issues as
1558 defined by the court. 1559

6. **Rationale of the Judgment:** This covers the
1560 reasoned analysis of the court, the application
1561 of law to the established facts, and the legal
1562 logic. 1563

7. **Result of the Judgment:** This is the final ad-
1564 judicative outcome and the specific orders of
1565 the court. 1566

Structuring Rules: 1567

- **Format Requirement:** Every entry must be
1568 represented as: Attribute_Name: "Exact text
1569 from the document". 1570
- **Strict Copying:** You must not change a sin-
1571 gle character. If a segment of text does not
1572 logically fit into any structured attribute, you
1573 must abandon the structuring for that specific
1574 segment rather than forcedly adapting the con-
1575 tent. 1576
- **Missing Information:** If a field is not explic-
1577 itly stated in the judgment, it must be recorded
1578 as *Not Specified*. 1579

H.3 Task II: Outcome-Relevant Keyword 1580
Labeling 1581

1582 Keywords must be labeled within three specific
1583 fields: **Focus of Disputes**, **Rationale of the Judg-**
1584 **ment**, and **Result of the Judgment**. Annotators
1585 must identify all information that could materially
1586 influence the adjudicative outcome of the case.

Labeling Categories and Examples: 1587

1588 Annotators must identify and label the following
1589 categories:

- **Legal Terms:** Doctrinal concepts and termi-
1590 nology, for example, *tort* or *liability*. 1591
- **Legal Statutes and Articles:** Specific le-
1592 gal citations, for example, *Civil Code Article*
1593 *1165*. 1594
- **Factual Elements:** 1595
 - **Verbs:** Actions describing the core
1596 events, for example, *crash* or *breach*. 1597
 - **Nouns:** Key objects, entities, or loca-
1598 tions, for example, *hospital* or *contract*. 1599

1600 – **Amounts:** Specific monetary values,
1601 percentages, or numbers, for example,
1602 *50,000 Renminbi* or *10 percent*.

1603 • **Involved Personnel:** This includes all parties
1604 involved in the litigation and identified
1605 individuals relevant to the facts.

1606 **Negative Constraints:** To ensure the purity of
1607 the data, do not label words that represent personal
1608 writing habits or logical transitions that do not hold
1609 substantive legal weight, such as *through*, *therefore*,
1610 *accordingly*, or *moreover*.

1611 I One-Page Annotation Checklist

ONE-PAGE ANNOTATION CHECKLIST

[PHASE ONE: STRUCTURING AUDIT]

- Is the data organized into the seven required fields?
- Is every entry formatted as a Key-Value pair?
- Does the extracted text match the original document character for character?
- Have you avoided any paraphrasing or summarizing?
- Have you cleaned the text of digital artifacts (e.g., page numbers, headers)?

[PHASE TWO: KEYWORD AUDIT]

- Have keywords been labeled only in:
 - * Focus of Disputes
 - * Rationale of the Judgment
 - * Result of the Judgment
- Have you included all parties involved and identified personnel?
- Have you included all monetary amounts and specific legal articles?
- Have you excluded stylistic transition words (e.g., *therefore*)?

[FIELD-SPECIFIC REFERENCE GUIDE]

Base Information
Extract from the document header and initial paragraph.
(No keyword labeling required)

Plaintiff's Arguments
Extract from the section starting with the claims of the plaintiff.
(No keyword labeling required)

Defendant's Arguments
Extract from the section starting with the defenses of the defendant.
(No keyword labeling required)

Facts Proved by the Court
Extract from the section describing the established facts.

(No keyword labeling required)

Focus of Disputes
Extract from the judicial summary of disputed issues.
Labeling focus: legal terms, disputed actions, specific amounts.

Rationale of the Judgment
Extract from the section containing the opinion of the court.
Labeling focus: legal statutes, legal doctrines, key nouns.

Result of the Judgment
Extract from the final adjudicative orders.
Labeling focus: final amounts, party names, final outcomes.

1613

J Verification Principles and Review Procedure

1614

1615

Role of Verification Annotators:

1616

Verification annotators act as rule-based auditors rather than secondary annotators. Prior to review, they must study and understand all annotation guidelines and treat them as the sole normative standard.

1617

1618

1619

1620

1621

Scope of Verification:

1622

Verification requires checking:

1623

- Completeness of required fields

1624

- Consistency across structured sections

1625

- Fidelity to the original judicial text

1626

- Compliance with all annotation prohibitions

1627

Verification annotators must not optimize wording or introduce legal judgment.

1628

1629

Error Identification Criteria:

1630

An annotation is considered problematic if required information is missing, unsupported by the source text, improperly inferred, or placed in an incorrect field.

1631

1632

1633

1634

Error Logging and Correction Protocol:

1635

When an issue is identified, verification annotators must record:

1636

1637

- Case ID

1638

- Affected Field

1639

- Problem Description

1640

- Rule Violated

1641

- Revised Content

1642

1612

Standard Error Log Template:

- Case ID:
- Affected Field:
- Problem Description:
- Rule Violated:
- Revised Content:

All corrections must strictly follow the original annotation rules and remain fully traceable to the source text.

Resolution Principle:

Corrections must be made strictly according to written rules. When multiple revisions are possible, the version that minimizes interpretation and maximizes textual fidelity must be selected.

Auditability Statement:

The structured error logs ensure that all annotation revisions are auditable, transparent, and reproducible, enabling independent assessment of dataset construction quality.

K One Case Example of JurisBench

The test dataset for the LLMs is in Chinese. Here it is translated into English.

```
1 {
2   "case_id": 8,
3   "case_name": "Spring Breeze Center v. Sun Xiuying;
   Xiyang Insurance Co., Ltd. -- Civil Judgment (
   First Instance) in a Motor Vehicle Traffic
   Accident Liability Dispute",
4   "court_of_acceptance": "Tianjin Hebei District People's
   Court",
5   "acceptance_date": "2025-08-18",
6   "parties_involved": {
7     "plaintiff": ["Spring Breeze Center"],
8     "defendants": ["Sun Xiuying", "Xiyang Insurance Co.,
   Ltd."]
9   },
10  "causes_of_action": [
11    "Tort liability dispute",
12    "Tort liability dispute",
13    "Motor vehicle traffic accident liability dispute"
14  ],
15  "plaintiff's_argument": {
16    "full_plaintiff_argument": "The plaintiff, Spring
   Breeze Center, requested: (1) that the
   defendants be ordered to reimburse the plaintiff
   for the advanced emergency medical expenses in
   the amount of RMB 79,923.51; and (2) that the
   defendants bear the litigation costs of the
   present case. Facts and reasons: On November 11,
   2024 at 15:47, the defendant Sun Xiuying was
   driving vehicle plate Tianjin A Gxx and
   proceeded south along the second traffic lane of
   Yuhong Road in Hebei District. Approaching the
   Wenzhou Gate, she entered the intersection on a
   red light. At the same time, Chen Chunli, who
   had not lawfully obtained a motor vehicle
   driving license, was riding an unregistered
   motorcycle bearing a suspended plate reading \"
   Tianjin xxx\" without wearing a helmet and while
   intoxicated (blood alcohol content measured at
   50.5 mg/100 ml), travelling from east to west
   into the same intersection. The front of the
   vehicle driven by Sun Xiuying contacted the
   right side of the motorcycle driven by Chen
   Chunli, causing personal injury to Chen Chunli
   and damage to both vehicles. Determinations
   assigned primary fault to Sun Xiuying and
```

```
secondary fault to Chen Chunli. The vehicle
Tianjin A Gxx is registered to the defendant Sun
Xiuying and was insured with both compulsory
traffic insurance and commercial third-party
liability insurance with a limit of RMB 3,000,00
0. Upon application by Chen Chunli, the
plaintiff advanced medical expenses of RMB 114,1
76.45 on her behalf.",
17 "plaintiff's_argument_list": {
18   "plaintiff's_argument_1": "Order the defendants to
   reimburse the plaintiff for the advanced
   emergency medical expenses of RMB 79,923.51",
19   "plaintiff's_argument_2": "Litigation costs to be
   borne by the defendants"
20 }
21 },
22 "defendant's_argument": {
23   "full_defendant's_argument": "Defendant Sun Xiuying
   did not file a defense. Defendant Xiyang
   Insurance Co., Ltd. contended that the vehicle
   Tianjin A Gxx was insured with the respondent
   for compulsory traffic insurance and commercial
   third-party liability insurance of RMB 3,000,000
   ; the accident occurred during the insurance
   period; the compulsory traffic insurance medical
   expense limit of RMB 18,000 has been exhausted;
   the driver of Tianjin A Gxx bears primary
   responsibility and should compensate according
   to the fault allocation; and the insurer
   disputed liability for the litigation costs.",
24   "defendant's_argument_list": {
25     "defendant's_argument_1": "The vehicle Tianjin A
   Gxx was insured with the respondent for
   compulsory traffic insurance and commercial
   third-party liability coverage of RMB 3,000,00
   0; the accident occurred during the policy
   period; the compulsory traffic insurance's
   medical expense limit of RMB 18,000 has been
   exhausted; the driver bears primary
   responsibility and should compensate according
   to the fault proportion; the insurer does not
   agree to bear litigation costs."
26   }
27 },
28 "focus_of_disputes": {
29   "full_focus_of_disputes": "The disputed issues in
   this case are: (1) what payment obligations
   should each defendant bear regarding the
   plaintiff's reasonable and lawful expenditures?
   (2) the amount of the plaintiff's reasonable and
   lawful expenditures.",
30   "focus_of_disputes_list": {
31     "focus_of_disputes_1": {
32       "content": "What payment obligations should each
   defendant bear for the plaintiff's
   reasonable and lawful expenditures?",
33       "keywords": "defendant, plaintiff, reasonable and
   lawful, expenditures, bear, payment,
   liability"
34     },
35     "focus_of_disputes_2": {
36       "content": "What is the amount of the plaintiff's
   reasonable and lawful expenditures?",
37       "keywords": "plaintiff, reasonable and lawful,
   expenditures, amount"
38     }
39   }
40   "keywords": "issues in dispute, defendant, plaintiff,
   reasonable and lawful, expenditures, bear,
   payment, liability, amount"
41 },
42 "facts_proved_by_the_court": "Based on the parties'
   statements and evidence reviewed and confirmed by
   the court, the court finds: On November 11, 2024
   at 15:47, the defendant Sun Xiuying drove vehicle
   Tianjin A Gxx south on the second lane of Yuhong
   Road in Hebei District and entered the
   intersection near Wenzhou Gate at a red light; at
   that time, Chen Chunli, who had not lawfully
   obtained a driving license and was not wearing a
   helmet and was intoxicated (BAC 50.5 mg/100 ml),
   was riding an unregistered motorcycle with a
   suspended plate into the intersection from east to
   west. The front of Sun's vehicle collided with
   the right side of Chen's motorcycle, causing
   injury to Chen and damage to both vehicles. The
   Tianjin Traffic Police Hebei Detachment determined
   that Sun bears primary responsibility, Chen bears
   secondary responsibility. After the accident,
   Chen was sent to a Tianjin hospital for treatment,
   and Spring Breeze Center advanced medical
   expenses of RMB 114,176.45 for Chen. The vehicle
   Tianjin A Gxx is owned by defendant Sun Xiuying,
   and the vehicle was insured with defendant Xiyang
   Insurance Co., Ltd. with compulsory traffic
   insurance and commercial third-party liability
   insurance for RMB 3,000,000; the accident occurred
   during the insurance period, but the compulsory
```

```

43     traffic insurance medical expense limit of RMB 18,
44     000 has been exhausted.",
    "rationale_of_the_judgment": {
    "full_rationale_of_the_judgment": "The court holds
    that the traffic accident determination issued
    by the traffic authority is accurate and
    sufficient, and no party raised objections; the
    court accepts it. In this traffic accident,
    defendant Sun Xiuying bears primary
    responsibility while Chen bears secondary
    responsibility; the compulsory traffic insurance
    medical expense limit for the vehicle has been
    exhausted. Therefore, 70% of the plaintiff's
    reasonable and lawful expenditures arising from
    this traffic accident should be advanced by
    defendant Xiyang Insurance Co., Ltd. within the
    scope of the commercial third-party liability
    policy, and any remaining shortfall shall be
    borne by the tortfeasor Sun Xiuying. The
    plaintiff's submitted application, undertaking,
    medical expense invoices, hospitalization fee
    lists, and bank remittance documentation from
    China Construction Bank confirm that it did
    advance medical expenses of RMB 114,176.45 for
    Chen; the court so finds. The amount of RMB 79,9
    23.51 shall be paid by defendant Xiyang
    Insurance Co., Ltd. to Spring Breeze Center.",
    "rationale_of_the_judgment_list": {
    "rationale_of_the_judgment_1": {
    "text": "The traffic accident determination by
    the traffic authority is accurate and
    sufficient; no party objected; the court
    accepts it. Defendant Sun bears primary
    responsibility and Chen secondary
    responsibility; the compulsory insurance
    medical expense limit has been exhausted.
    Therefore, 70% of the plaintiff's reasonable
    and lawful expenditures arising from the
    accident should be advanced by the defendant
    insurer under the commercial third-party
    liability policy; any shortfall remains the
    tortfeasor's responsibility.",
    "keywords": "traffic accident, traffic authority,
    accident determination, accuracy,
    sufficiency, parties, no objection, court,
    acceptance, defendant, liability allocation,
    primary responsibility, secondary
    responsibility, vehicle, compulsory traffic
    insurance, medical expense, policy limit,
    exhausted, plaintiff, reasonable and lawful,
    expenditures, insurer, commercial third-
    party liability, advance payment, shortfall,
    tortfeasor, indemnity"
    },
    "rationale_of_the_judgment_2": {
    "text": "Spring Breeze Center provided proofs
    demonstrating that it advanced medical
    expenses of RMB 114,176.45 for Chen; the
    court so finds. Accordingly, defendant
    Xiyang Insurance Co., Ltd. shall pay RMB 79,
    923.51 to Spring Breeze Center.",
    "keywords": "Spring Breeze Center, submissions,
    application, undertaking, medical invoices,
    hospitalization fee lists, bank remittance,
    proof, treatment, advanced medical expenses,
    court confirmation, amount, insurer payment"
    }
    },
    "rationale_of_the_judgment_2": {
    "text": "Spring Breeze Center provided proofs
    demonstrating that it advanced medical
    expenses of RMB 114,176.45 for Chen; the
    court so finds. Accordingly, defendant
    Xiyang Insurance Co., Ltd. shall pay RMB 79,
    923.51 to Spring Breeze Center.",
    "keywords": "Spring Breeze Center, submissions,
    application, undertaking, medical invoices,
    hospitalization fee lists, bank remittance,
    proof, treatment, advanced medical expenses,
    court confirmation, amount, insurer payment"
    }
    },
    "factual_basis": "The traffic authority's accident
    determination is accurate. Defendant Sun bears
    primary responsibility. Spring Breeze Center
    advanced medical expenses of RMB 114,176.45.",
    "legal_grounds": [
    { "PRC Civil Code": ["Article 11179", "Article 1208
    ", "Article 1213"] },
    { "PRC Road Traffic Safety Law": ["Article 76"] },
    { "PRC Civil Procedure Law": ["Article 147"] },
    { "PRC Road Traffic Safety Law": ["Article 76"] },
    { "Supreme People's Court Interpretation on Road
    Traffic Accident Damage Compensation": ["
    Article 13"] },
    { "Supreme People's Court Interpretation on the
    Application of the PRCL": [ "Article 90" ] }
    ],
    "keywords": "court, traffic accident, traffic
    authority, accident determination, accuracy,
    sufficiency, parties, no objection, acceptance,
    defendant, liability allocation, medical expense
    limit, exhausted, Spring Breeze Center,
    advanced medical expenses, insurer, indemnity"
    },
    "result_of_the_judgment": {
    "full_result_of_the_judgment": "Within fifteen days
    after this judgment becomes effective, defendant
    Xiyang Insurance Co., Ltd. shall pay Spring
    Breeze Center RMB 79,923.51. If a monetary
    obligation is not performed within the period

```

```

    specified in this judgment, interest on the
    overdue debt shall be doubled in accordance with
    Article 264 of the PRC Civil Procedure Law. The
    case acceptance fee of RMB 299.6 shall be borne
    by defendant Sun Xiuying. If dissatisfied, an
    appeal may be filed within fifteen days from the
    date of delivery of the judgment to this court;
    the appeal is to be heard by the Tianjin No.2
    Intermediate People's Court.",
    "keywords": "judgment effective, defendant, insurer
    payment, plaintiff, appeal jurisdiction, case
    fee",
    "result_of_the_judgment_list": {
    "result_of_the_judgment_1": {
    "text": "Within fifteen days after this judgment
    becomes effective, defendant Xiyang
    Insurance Co., Ltd. shall pay Spring Breeze
    Center RMB 79,923.51",
    "keywords": "judgment effective, defendant,
    Xiyang Insurance Co., Ltd., payment,
    plaintiff, Spring Breeze Center"
    },
    "result_of_the_judgment_2": {
    "text": "If the monetary obligation is not
    performed within the period specified in
    this judgment, interest on the overdue debt
    shall be doubled in accordance with Article
    264 of the PRC Civil Procedure Law.",
    "keywords": "judgment, period, performance,
    monetary obligation, PRC Civil Procedure Law
    , Article 264, double interest, overdue
    period, debt interest"
    },
    "result_of_the_judgment_3": {
    "text": "The case acceptance fee of RMB 299.6
    shall be borne by defendant Sun Xiuying.",
    "keywords": "case acceptance fee, defendant,
    borne"
    },
    "result_of_the_judgment_4": {
    "text": "If dissatisfied with this judgment, an
    appeal may be lodged within fifteen days
    from the date of delivery of the judgment;
    file an appeal with this court and submit
    copies according to the number of opposing
    parties; the appeal shall be heard by the
    Tianjin No.2 Intermediate People's Court.",
    "keywords": "appeal, dissatisfied, judgment,
    delivery, file appeal, copies, opposing
    parties, Tianjin No.2 Intermediate People's
    Court"
    }
    }
    },
    "Remedy": {
    "plaintiff_share": 0,
    "defendant_share": { "advanced_expenses_RMB": 79923
    .51, "fee_RMB": 299.6 }
    },
    "appeal_jurisdiction": "Tianjin No.2 Intermediate
    People's Court"
    }
    }

```

1668

L Input Prompts

1669

The input for the LLMs is in Chinese. Here it is translated into English.

1670

1671

```

1 # =====
2 Test 1: Multi-level cause of action prediction
3 =====
4 # 1-1: Step-by-step cause of action prediction with
5 knowledge base
6 # 1-2: Step-by-step cause of action prediction without
7 knowledge base
8 PT_test1 = '''You are a legal expert. Goal: Return the
9 top K most likely candidates for the cause of action
10 for the given case at the specified level (sorted
11 from high to low probability).
12
13 Return format requirement: Output only a JSON array (
14 brackets), containing string candidates. For example
15 : ["Contract Dispute", "Debt Dispute"]. Do not
16 output any explanation, reasoning, or other text.
17
18 Current level: Level {level}; please return no more than
19 {top_k} candidates.
20
21 {kb_candidates_text}
22 {one_shot_text}

```

1672

```

15 {reveal_path_text}
16 {case_name_text}
17 Input text: {case_text}
18
19 Now, please output only the top {top_k} candidate causes
    of action for Level {level} (JSON array), and do not
    output other text.
20
21
22
23 # =====
24 Test 2: Prediction of the focus of disputes
    =====
25
26 # 2-1: Predict each focus of disputes (item-by-item
    matching)
27 PT_test2_dispute = '''You are a legal expert. Please list
    all the focus of disputes of this case one by one
    based on the case information below.
28
29 Requirements:
30 - Please list the focus of disputes of this case one by
    one, with each focus of disputes on a separate line.
31 - The focus of disputes should accurately reflect the
    core disputed issues between the plaintiff and the
    defendant in this case.
32 - Must be listed in the format of "Focus of Disputes 1:",
    "Focus of Disputes 2:", etc., with each focus of
    disputes occupying one line.
33 - Output format example:
34 Focus of Disputes 1: The issue of liability
    determination in this traffic accident
35 Focus of Disputes 2: The losses of the plaintiff,
    Natural Person A, caused by this traffic accident
36 Focus of Disputes 3: The assumption of civil liability
    by the defendants, Natural Person D and Natural
    Person E
37 - Do not output any explanation, reasons, or other text.
38
39 {one_shot_text}
40
41 Case Information:
42 Court of Acceptance: {court}
43 Date of Conclusion: {date}
44 Party: {parties}
45 Cause of action: {case_cause}
46 Plaintiff's argument: {plaintiff}
47 Defendant's argument: {defendant}
48
49 Please list the focus of disputes one by one according to
    the requirements:'''
50
51
52
53 # =====
54 Test 3: Prediction of the rationale of the judgment
    =====
55
56 # 3-1: Predict factual basis
57 PT_test3_fact = '''You are a legal expert. Please
    generate the factual basis section of the rationale
    of the judgment based on the case information below.
58
59 Requirements:
60 - The factual basis should state the facts of the case
    comprehensively and objectively.
61 - It should include the court's factual findings
    regarding the focus of disputes.
62 - Relevant evidence materials should be cited.
63 - The language should be rigorous, standardized, and
    conform to the writing requirements of legal
    documents.
64 - Directly output the text of the factual basis, do not
    output titles like "Factual Basis:" or other
    explanatory text.
65
66 {one_shot_text}
67
68 Case Information:
69 Court of Acceptance: {court}
70 Date of Conclusion: {date}
71 Party: {parties}
72 Cause of action: {case_cause}
73 Plaintiff's argument: {plaintiff}
74 Defendant's argument: {defendant}
75 The focus of disputes: {dispute_focus}
76 Facts found by the court: {court_findings}
77
78 Please generate the factual basis:'''
79
80 # 3-2: Predict legal grounds (Law Names)
81 PT_test3_law_names = '''You are a legal expert. Please
    list all the names of the laws that should be cited
    in the rationale of the judgment of this case, based
    on the case information below.
82
83 Requirements:
84 - Only output the law names, in the format: "Law Name".
85 - One law name per line.

```

```

86 - You must use book title marks ("") to wrap the law name
    , otherwise it will not be counted in the accuracy
    rate.
87 - Sort by importance.
88 - Do not output any explanation or other text.
89
90 {one_shot_text}
91
92 Case Information:
93 Court of Acceptance: {court}
94 Date of Conclusion: {date}
95 Party: {parties}
96 Cause of action: {case_cause}
97 Plaintiff's argument: {plaintiff}
98 Defendant's argument: {defendant}
99 The focus of disputes: {dispute_focus}
100 Facts found by the court: {court_findings}
101
102 Please strictly follow the requirements to list the names
    of the laws that should be cited (one per line,
    must use book title marks):'''
103
104 # 3-3: Predict legal Articles (for each Law)
105 PT_test3_legal_articles = '''You are a legal expert.
    Please list the specific Articles of the involved
    law applicable to this case, based on the case
    information below and the identified name of the law
    involved in the judgment.
106
107 Requirements:
108 - Only list the specific Article numbers, in the format:
    Article XX.
109 - Each Article on a separate line.
110 - Arrange in the order of the Articles.
111 - Do not output any explanation or other text.
112
113 {one_shot_text}
114
115 Case Information:
116 Cause of action: {case_cause}
117 The focus of disputes: {dispute_focus}
118 Facts found by the court: {court_findings}
119
120 Law Name: {legal_name}
121
122 Please list the specific applicable Articles in this law
    (one per line):'''
123
124 # 3-4: Predict list of rationale of the judgment based on
    araguments
125 PT_test3_reasoning_list = '''You are a legal expert.
    Please generate the corresponding rationale of the
    judgment for each of the plaintiff's arguments
    separately, based on the case information and the
    plaintiff's arguments below.
126
127 Requirements:
128 - For each of the arguments, separately explain the
    reasons why the court supports or does not support
    it.
129 - Each rationale of the judgment must be listed in the
    format of "Rationale of the judgment 1:", "Rationale
    of the judgment 2:", etc., with each one occupying
    a paragraph.
130 - The language should be rigorous, standardized, and
    conform to the writing requirements of legal
    documents.
131 - Do not output other titles or explanations, directly
    output the content of the rationale of the judgment.
132
133 {one_shot_text}
134
135 Case Information:
136 Court of Acceptance: {court}
137 Date of Conclusion: {date}
138 Party: {parties}
139 Cause of action: {case_cause}
140 Plaintiff's claims: {litigation_requests}
141 Defendant's argument: {defendant}
142 The focus of disputes: {dispute_focus}
143 Facts found by the court: {court_findings}
144
145 Please generate the corresponding rationale of the
    judgment for each claim:'''
146
147
148
149
150 # =====
151 Test 4: Prediction of the result of the judgment
    =====
152
153 # 4-1: Predict remedy
154 PT_test4_fee = '''You are a legal expert. Please predict
    the sharing of remedies based on the case
    information and result of the judgment below.
155
156 Requirements:

```

```

157 - The expense items to be predicted include: Plaintiff
      share, Defendant share (there may be multiple
158 - Output only one JSON object, in the following format (
      strictly observe this format):
159   [{"Plaintiff share": {"Case acceptance fee and
      property preservation fee (Unit: RMB)": Amount}},
      "Defendant share": {"Compensation (Unit: RMB)":
      Amount, "Case acceptance fee and property
160 - If there are multiple defendants in the result of the
      judgment, the "Defendant bears" field needs to list
      them separately by defendant name.
161 - The amount unit is RMB, retain two decimal places or
      integer.
162 - Do not output any explanation or other text, only
      output the JSON object.
163
164 {one_shot_text}
165
166 Case Information:
167 Court of Acceptance: {court}
168 Party: {parties}
169 Cause of action: {case_cause}
170 Result of the judgment: {result}
171
172 Please output the expense prediction (JSON format):'''
173
174 # 4-2: Predict items of result of the judgment based on
      arguments
175 PT_test4_result_list = '''You are a legal expert. Please
      generate the corresponding result of the judgment
      for each of the plaintiff's claims separately, based
      on the case information and the plaintiff's claims
      below.
176
177 Requirements:
178 - For each argument, separately state the content of the
      court's judgment.
179 - Each result of the judgment must be listed in the
      format of "Result of the judgment 1:", "Result of
      the judgment 2:", etc., with each one occupying a
      paragraph.
180 - The language should be rigorous, standardized, and
      conform to the writing requirements of legal
      documents.
181 - Do not output other titles or explanations, directly
      output the content of the result of the judgment.
182
183 {one_shot_text}
184
185 Case Information:
186 Court of Acceptance: {court}
187 Date of Conclusion: {date}
188 Party: {parties}
189 Cause of action: {case_cause}
190 Plaintiff's argument: {litigation_requests}
191 Defendant's argument: {defendant}
192 The focus of disputes: {dispute_focus}
193 Facts proved by the court: {court_findings}
194 Rationale of the judgment: {reasoning}
195
196 Please generate the corresponding result of the judgment
      for each argument:'''
197
198 # 4-3: Predict appeals jurisdiction
199 PT_test4_appeal = '''You are a legal expert. Please
      predict the court of appeals jurisdiction for this
      case based on the case information below.
200
201 Requirements:
202 - Only output the full name of the court of appeals
      jurisdiction, for example: Jiangmen Intermediate
      People's Court of Guangdong Province.
203 - Do not output any explanation or other text.
204
205 {one_shot_text}
206
207 Case Information:
208 Court of Acceptance: {court}
209 Cause of action: {case_cause}
210
211 Please output the name of the court of appeals
      jurisdiction:'''

```