

# ScreenshotLegalBench: A Multimodal Benchmark for Legal Evidence Understanding in Chat Screenshots

Anonymous EMNLP submission

## Abstract

Chat screenshots from platforms such as WeChat are increasingly used as legal evidence in Chinese civil litigation. However, their informal layout, multimodal nature, and lack of structure pose significant challenges for automated understanding. We introduce **ScreenshotLegalBench**, the first large-scale multimodal benchmark for *Legal Screenshot Evidence Understanding (LSEU)*. It supports two key tasks: (1) structured key information extraction (KIE) and (2) legal visual question answering (VQA). The dataset contains over 4,600 chat screenshots annotated with 145,044 structured labels, a 143-image evaluation set with 2,678 verified annotations, and 1,176 VQA instances covering evidence relevance, format validity, and legal reasoning. Among these, 106 cases involve multi-image cause-of-action scenarios. We benchmark several open-source vision-language models (VLMs), including InternVL and Qwen-VL families. Experimental results show that current VLMs struggle with layout interpretation and domain-specific reasoning, despite instruction tuning. **ScreenshotLegalBench** offers a novel and scalable resource at the intersection of vision, language, and law, enabling future research on multimodal legal document understanding in real-world settings. The dataset and code are soon publicly available at Github.

## 1 Introduction

In recent years, multimodal large models have emerged as a major focus in AI research, demonstrating impressive performance on tasks that require integrating image and text modalities. Within the legal domain, such models are increasingly expected to automate the preprocessing of complex and loosely structured case materials, particularly in civil proceedings where parties are required to submit supporting evidence in diverse digital formats. Among these, WeChat chat screenshots have become a common form of statutory evidence in

China, as officially recognized in Art.116 of the Supreme People’s Court’s Interpretation on the Civil Procedure Law.

However, analyzing chat screenshots remains a deeply manual task in judicial practice. Legal practitioners must verify speaker identities, reconstruct conversation sequences, determine the legal relevance of each message, and evaluate whether the screenshot satisfies evidentiary requirements. Unlike formal documents, chat screenshots are informal, heterogeneous, and visually irregular—mixing text, images, emojis, file attachments, and transfer records within complex UI layouts. The absence of structured representations, coupled with the multimodal and noisy nature of the content, poses significant obstacles to automation. As a result, lawyers and judges must perform labor-intensive and error-prone manual reviews, which reduces efficiency and increases the risk of oversight.

To address this gap, we propose the task of *Legal Screenshot Evidence Understanding (LSEU)*, which aims to extract structured legal information and assess evidentiary attributes from real-world, multimodal inputs. Unlike generic visual question answering (VQA) or document understanding tasks, LSEU presents unique challenges: (1) informal and irregular visual layouts, (2) field-level legal structuring grounded in procedural norms, and (3) domain-specific reasoning required for admissibility assessment and case interpretation. As illustrated in Appendix Figure 3, legal practitioners typically process such evidence in five stages: relevance screening, timeline reconstruction, factual extraction, legal mapping, and litigation strategy development. Our benchmark focuses on three of these stages, namely structured perception, factual extraction, and legal mapping, which collectively capture the core components of real-world judicial workflows.

However, existing vision-language models strug-

gle to address LSEU effectively due to three key challenges: (1) the multi-modal and informal nature of chat screenshots, (2) the need for structured field extraction aligned with legal interpretation, and (3) the domain-specific reasoning required to infer case types or assess evidentiary admissibility.

To tackle these challenges, we decompose the LSEU task into two core components: **structured perception**, which detects and extracts legally relevant information (e.g., who said what, when, and with what legal implication), and **semantic reasoning**, which classifies legal relevance, determines evidentiary status, and infers preliminary case hypotheses. This decomposition reflects how legal professionals process chat evidence—first segmenting and organizing visual information, then reasoning over the structured content. **We are guided by the following research questions:**

- *RQ1: Can large vision-language models accurately extract legally structured information from noisy, multimodal chat screenshots?*
- *RQ2: Does a two-stage modeling pipeline that combines structured perception with downstream semantic reasoning outperform end-to-end baselines in legal classification and case reasoning tasks?*

To answer these questions, we present **ScreenshotLegalBench**, the first publicly available benchmark designed for multimodal legal evidence understanding in chat-based scenarios. The dataset includes over 4,600 WeChat chat screenshots annotated for KIE, as well as more than 1,100 VQA pairs targeting legal classification, evidence validity, and cause-of-action reasoning. All annotations are performed by certified Chinese legal professionals, ensuring alignment with practical legal standards. **Our contributions are summarized as follows:**

- We introduce **ScreenshotLegalBench**, a multimodal benchmark for legal evidence understanding in chat screenshots, comprising real-world and simulated samples annotated by certified Chinese legal professionals to reflect practical legal needs and legal reasoning demands.
- We introduce **Legal Screenshot Evidence Understanding (LSEU)** as a two-stage task: structured KIE and legal VQA, supported by a unified annotation schema and scalable labeling pipeline.
- We benchmark a range of open-source vision-

language models under realistic deployment constraints, focusing on models that are suitable for local and privacy-sensitive legal environments. The results show that even advanced instruction-tuned models continue to face persistent challenges in layout robustness, field-level structuring, and legal reasoning across diverse input forms.

By bridging the domains of vision, language, and legal reasoning, this work offers a new foundation for multimodal legal AI. We hope that **ScreenshotLegalBench** will catalyze future research on interpretable and reliable systems for real-world evidence analysis.

## 2 Related Work

**Legal Information Extraction.** Early work in legal NLP focused on clause extraction and entity identification from structured contracts or rulings. CUAD (Hendrycks et al., 2021), LEDGAR (Tuggenier et al., 2020), and ACORD (Wang et al., 2025) provide high-quality text-based datasets for commercial clause classification and retrieval. However, these datasets operate on well-formatted, language-only documents, lacking support for multimodal input or visual layout reasoning.

Recent studies (Liu et al., 2023) show that incorporating visual cues such as bounding boxes and font styles can improve structured extraction from long documents. These findings motivate our structured KIE approach for visually noisy chat screenshots, which lack standardized layouts and often contain overlapping modalities such as image messages, emojis, and file transfers. This setting poses new challenges for entity alignment and layout-robust modeling.

**Legal Reasoning and Understanding.** Benchmarks like LexGLUE (Chalkidis et al., 2022) and CaseHOLD (Zheng et al., 2021) define a suite of judgment prediction, retrieval, and question answering tasks. LEGAL-BERT (Chalkidis et al., 2020) demonstrates the importance of domain-adaptive pretraining. However, these benchmarks assume clean, pre-extracted legal facts and do not support evidence-level interpretation from raw multimodal inputs.

Chinese datasets such as JEC-QA (Zhong et al., 2019) and CAIL2018 (Xiao et al., 2018) provide useful resources for statutory reasoning and charge prediction but remain focused on formal court doc-

Task Type	Data Concern	Structured Layout	Visual Reasoning	Textual Semantics	Multimodal Alignment	Temporal Reasoning	Legal Action Modeling
Form KIE	Doc Image + OCR + Position	✓	✗	△	✗	✗	✗
Layout Parsing	Doc Image + Layout Tags	✓	✓	△	✓	✗	✗
DocVQA	Doc Image + OCR + QA Pair	△	△	✓	△	✗	✗
TextVQA	Scene Image + OCR + Question	✗	✓	✓	△	✗	✗
Table QA	Table Image + Question	✓	✗	✓	△	✗	✗
InfoVQA	Image + OCR + Embedded Info	△	✓	✓	✓	✗	✗
<b>LSEU (Ours)</b>	Screenshot + OCR + Bubble + Media + Timestamps	△	✓	✓	✓	✓	✓

Table 1: Comparison of our chat screenshot task with representative KIE/VQA tasks. ✓ indicates presence of the feature; ✗ indicates absence; △ indicates partial support or context-dependent presence.

uments. In contrast, our task addresses an earlier stage of legal workflows—assessing the admissibility and relevance of raw WeChat evidence prior to fact consolidation. This pre-factual focus introduces challenges in determining whether a screenshot constitutes legal evidence at all, requiring both structural perception and contextual inference.

Recent systems such as MASER (Jeon et al., 2022) demonstrate the ability of MLLMs to infer event chronology under weak timestamp supervision. While such systems highlight the potential of multimodal legal reasoning, they typically operate on clean, structured records, and do not handle the fragmented, layout-rich format of chat screenshots.

**Multimodal Legal Benchmarks.** To date, few datasets address multimodal legal evidence analysis. Prior work in DocVQA (Mathew et al., 2021), TextVQA (Hegde et al., 2023), and InfoVQA (Mathew et al., 2022) has explored visual-text reasoning in document or scene settings, but these benchmarks lack legal-specific labels and structural alignment requirements. In contrast, our task focuses on raw WeChat chat screenshots, which are often unstructured, multimodal, and legally ambiguous. It integrates both structured KIE and VQA, covering interface layout, temporal ordering, and high-level legal implications. Table 1 provides a comparative summary of our task relative to representative KIE/VQA benchmarks across multiple reasoning dimensions.

Our proposed dataset, ScreenshotLegalBench, is the first to support three interrelated subtasks over real-world chat screenshots: (1) *structured KIE* (speaker, content, time, etc.); (2) *legal attribute classification* (e.g., relevance and evidentiary status); and (3) *open-ended cause-of-action generation*. These capabilities reflect practical needs in

legal workflows such as fact triage, evidence validation, and dispute summarization, which are tasks that precede traditional legal judgment prediction.

### 3 Task Definition

We define two core tasks for **Legal Screenshot Evidence Understanding (LSEU)**, reflecting the structured perception and legal reasoning stages over chat-based visual evidence.

**Task 1: Screenshots Evidence Key Information Extraction (SEKIE)** aims to extract structured legal fields from a single WeChat chat screenshot. The model must jointly understand the layout and semantics of the chat interface and output a structured JSON record containing message-level fields. The expected fields include:

- *speaker*: the display name of the message sender;
- *timestamp*: the message time, if available;
- *content*: the textual content of the message;
- *message\_bbox*: the bounding box of the message region;
- *transfer, image, file*: optional fields indicating the presence and description of funds transfer, images, or file attachments.

This task forms the structural foundation for downstream classification and reasoning.

**Task 2: Chat Screenshot Legal Visual Legal Question Answering (CSLVQA)** evaluates the model’s ability to perform higher-level legal understanding based on the screenshot. It includes three subtasks: (1) classifying whether the screenshot is legally relevant (*classify*); (2) judging whether the screenshot qualifies as formally valid evidence (*evidence*); and (3) generating a natural language description of the underlying dispute or legal issue

(case\_text). Notably, this sub-task is a multi-image reasoning task, where the model must synthesize information across multiple screenshots to infer a coherent legal cause. The VQA task can be performed either directly from the raw image or using the structured output from KIE as additional input.

## 4 The ScreenshotLegalBench Dataset

We construct a dataset for Legal Evidence Understanding in Chat Screenshots.

### 4.1 Data Collection

We initially obtain raw images over 9,800 candidate images from the Gansu Provincial Digital Rule of Law Industry Research Academy, which are from Common Crawl, Google and Baidu search results, and internal institutional repositories. Each image, along with its embedded caption, was processed by Gemini 1.5 (Team et al., 2024) to determine whether it resembled a chat interface and to assign a coarse-grained content label (e.g., “startup,” “divorce,” “romantic relationship”). This automated step yielded approximately 6,000 images likely to depict chat screenshots. Subsequently, trained annotators manually reviewed the topic labels and visual layout to identify cases with potential legal relevance, resulting in a filtered set of 4,800s samples for annotation in ScreenshotLegalBench.

### 4.2 Annotation Data Elicitation

We elicited the ScreenshotLegalBench annotations by defining two complementary pipelines for our KIE and VQA tasks (see Appendix B for full schema and annotation guidelines). Figure 4 illustrates a sample of the task annotations.

**KIE task** employs YOLOv3 (Redmon and Farhadi, 2018), DETR (Carion et al., 2020), and Cascade R-CNN (Cai and Vasconcelos, 2017), fine-tuned on a 50-image few-shot subset, to localize 16 key interface elements, including message bubbles, avatars, timestamps, and transaction indicators. Targeted data augmentation enhances robustness across diverse chat layouts, yielding a 70.1% mAP. These detectors identify candidate layout regions across the dataset, from which text is extracted using PaddleOCR and Google OCR. Multi-line transaction entries are semantically merged through rule-based consolidation into coherent legal statements (e.g., a three-line WeChat transfer becomes “WeChat transfer received ¥520.00, note:

happy birthday”). File-related content is normalized by parsing filenames and extensions, while image-only regions are annotated using a multi-modal generative model to enhance reasoning context. Speaker attribution is determined by analyzing bounding-box centroids relative to the page’s vertical axis, assigning right-side elements to the primary speaker and left-side elements to the interlocutor. Missing timestamps are interpolated using a sliding-window strategy to maintain temporal continuity. All spatial, textual, and semantic annotations are organized into a unified JSON schema, serving as high-quality weak supervision for downstream causal analysis and evidentiary chain reconstruction. The layout schema and bounding-box definitions are detailed in Appendix B.

For the **VQA task**, we designed a set of expert-authored legal questions to elicit complex reasoning abilities from multimodal models. **Notably**, all VQA annotations were created from scratch by experienced legal professionals, without the use of automated pre-labeling. The questions are divided into two levels: global questions, which assess the screenshot as a whole, and local questions, which focus on fine-grained content such as individual messages or UI elements. The **global questions**, which have been fully annotated and released, guide the model to reason across four legal dimensions: (a) whether the image is a chat screenshot and holds legal relevance; (b) whether it satisfies evidentiary completeness and admissibility criteria; (c) whether it depicts private or group conversation; and (d) what type of legal dispute (e.g., loan, labor) the conversation may suggest. These questions serve as the foundation for high-level evidence screening and case framing. In contrast, the **local questions** target specific visual or textual components such as message content, avatars, quoted speech, emojis, transfers, and file attachments. They are designed to test the model’s ability to extract intent, recognize legal relationships, classify transaction types, and interpret symbolic or emotive cues. Due to their labor-intensive nature, local annotations are still in progress and will be released in a future update alongside detailed statistics. Nonetheless, the schema has been finalized to ensure backward compatibility and extensibility. Table 6 presents representative examples of local-level questions and their annotation goals.



### 4.3 Manual Correction and Expert Review

To ensure evaluation integrity, we conducted detailed manual correction for a held-out subset of 143 KIE proposals. These proposals were initially generated by baseline models and subsequently reviewed by trained annotators. Corrections included bounding box adjustments for spatial accuracy, merging or splitting of entity spans, and fixing misclassified field types. This process yielded a reliable reference set for evaluating KIE performance. Full annotation criteria and workflows are provided in Appendix C.

For the VQA task, all annotations were fully manual. Expert annotators first created bounding boxes and question-answer pairs based on predefined prompts. All annotations were performed by legally qualified annotators who had passed China’s National Judicial Examination. For questions requiring nuanced legal judgment, responses were validated by senior attorneys with over a decade of practice, ensuring consistency with real-world legal reasoning.

### 4.4 Annotation Quality Assurance

To ensure the consistency, completeness, and legal validity of ScreenshotLegalBench annotations, we implemented a multi-stage quality assurance protocol integrating model-assisted pre-processing, a hierarchical annotation framework, and multi-level expert review.

As described in Sections 4.2 and 4.3, both the KIE and VQA pipelines combine structured interfaces, heuristic post-processing, and expert-in-the-loop validation. For KIE, model-generated layout elements were aligned with OCR results and then refined through legal-specific consolidation and manual correction on 143 samples. For VQA, all annotations were manually created, with legally sensitive questions reviewed by senior attorneys.

As shown in Table 5, this framework begins with global property tagging (e.g., legality, chat type, case type) and progresses to layout-level detection (e.g., message, avatar, file), content structuring (e.g., speaker name, message text), semantic transformation (e.g., merging transfer info into coherent legal phrases), and finally to legal question answering over both global and local visual regions. Each level corresponds to a distinct layer of information abstraction required for multimodal legal understanding.

This layered structure ensures both fine-grained

supervision for information extraction and high-level signals for reasoning tasks. All annotations followed centralized task formats, and ambiguous cases were discussed and resolved through collaborative expert review. Examples of annotation interfaces and VQA samples are provided in Appendix B.5.

### 4.5 Dataset Statistics

ScreenshotLegalBench comprises three complementary subsets designed to support multimodal legal tasks in WeChat chat screenshots: (1) an object detection subset for layout element localization, (2) a large-scale KIE corpus for pretraining and evaluation, and (3) a VQA benchmark for multimodal legal reasoning. Notably, the KIE and most VQA tasks are annotated at the single-image level, while the case\_text sub-task adopts a multi-image setting—each case aggregates an average of 4.7 screenshots to support cause-of-action analysis across dialogue contexts. Overall statistics are summarized in Table 2, with detailed field counts provided in Appendix C.

## 5 Experiments

We evaluate a range of vision-language models on ScreenshotLegalBench to assess their performance on structured perception and legal reasoning. Section 5.3 compares open-source baselines on KIE and VQA tasks. Section 5.4 demonstrates that, even with partially automated annotations, fine-tuned models outperform larger zero-shot baselines, highlighting the generalizability of our dataset. Finally, in Section 5.5, we conduct ablations to validate our dataset design, demonstrating that the inclusion of KIE as a structured perception task significantly improves downstream legal classification and reasoning.

### 5.1 Benchmark Models

Early legal NLP systems often relied on rule-based heuristics or traditional machine learning (e.g. SVMs(Chen and Lin, 2006)), but these methods tend to fail under the noisy, layout-rich conditions of chat screenshots. Our evaluation focuses on multimodal foundation models with strong image-text processing capabilities, excluding shallow baselines. All experiments are conducted under data confidentiality constraints, due to real-world deployment needs in local, privacy-sensitive legal environments. We select two prominent model fam-

Subset	Samples	Total Annotations	Avg. Annotations per Sample	Main Tasks
Object Detection	50 screenshots	945 bounding boxes	18.9	Layout Element Detection
KIE Training Set	39,477 messages	145,044 fields	3.67	Structured Pretraining
KIE Eval Set	143 screenshots / 696 messages	2,678 fields	3.85	Structured Evaluation
VQA Set	1,176 screenshots / 106 multi-imgs case	2,854 legal annotations	2.42 legal QA pairs	Dialog Type, Evidence, Case Reasoning

Table 2: Overview of ScreenshotLegalBench dataset subsets.

ilies, InternVL (Chen et al., 2024c,b,a) and Qwen-VL (Bai et al., 2023a; Wang et al., 2024; Bai et al., 2025), as benchmark decoders. Both are widely used in Chinese image–text tasks and are capable of generating structured outputs:

- **InternVL2** (Chen et al., 2024b) is a dual-encoder model fine-tuned on ScreenshotLegalBench for legal information extraction.
- **InternVL2.5** (Chen et al., 2024a) extends InternVL2 with QLoRA for efficient domain adaptation, updating only low-rank adapters while freezing the visual and language backbones.
- **Qwen2-VL** (Wang et al., 2024) and **Qwen2.5-VL** (Bai et al., 2025) are Transformer-based models optimized for Chinese image–text inputs, pretrained on multilingual corpora and capable of structured output generation.

All models are evaluated using the prompt templates detailed in Appendix D.4. Fine-tuned models are assessed under pass@0 to reflect output stability, while raw models are evaluated under pass@1. Higher format scores indicate stronger structural adherence and benefit from evaluation-side repair strategies enhancing JSON compatibility.

## 5.2 Evaluation Metrics

We evaluate model performance separately on the KIE and VQA tasks. For KIE, the output is a structured JSON containing a list of messages, each with textual and spatial fields. We measure quality along three axes: (1) structural validity, checking that each message includes all required fields in legal formats; (2) semantic accuracy, computed via a hybrid similarity score that combines normalized token-wise alignment and substring overlap; and (3) spatial alignment, evaluated using standard **Intersection-over-Union (IoU)** between predicted and reference bounding boxes. The semantic score for a message field  $y$  against reference  $\hat{y}$  is defined as

$$\text{Sim}(y, \hat{y}) = \lambda \text{SeqSim}(y, \hat{y}) + (1 - \lambda) \text{LCS}(y, \hat{y})$$

where SeqSim measures the proportion of aligned token spans under optimal matching, and LCS denotes the ratio of the longest common substring length. Overall KIE score averages across valid messages.

For VQA, we consider two classification tasks (classify, evidence) and one generation task (case\_text). The classification tasks are evaluated using macro-averaged Precision, Recall, F1, and Accuracy, as the label distributions are notably imbalanced. For case\_text, we evaluate the ability of the model to generate a concise, legally coherent description of the dispute based solely on the screenshot. While full-text metrics such as BLEU (Papineni et al., 2002) or ROUGE (Lin, 2004) are widely used in generative tasks, they are ill-suited for legal cause-of-action summaries due to the professional phrasing, variable expression, and high semantic abstraction involved. Instead, we adopt a simplified but interpretable metric: hit rate over legal dispute categories, which evaluates whether the predicted output contains at least one correct category keyword. Formally, the metric is defined as

$$\text{Dispute HitRate} = \frac{1}{N} \sum_{i=1}^N 1[\exists c \in C_i \cap \hat{C}_i]$$

where  $C_i$  is the set of dispute keywords extracted from the model output and  $\hat{C}_i$  is the gold label set. A hit is counted if at least one legal category is correctly recovered. Additional metrics such as normalized similarity and output length consistency are used for robustness and are detailed in Appendix B.

## 5.3 Main Performance

**KIE Task Results** We evaluate vision–language models on the KIE task using ScreenshotLegalBench (Figure 1, Table 10). The task evaluates structured output quality, focusing on format validity, spatial alignment (IoU), and content accuracy.

As shown in Table 10, performance varies significantly across models. InternVL2.5-2B (Chen et al., 2024a) achieves the highest overall score of 0.6302,

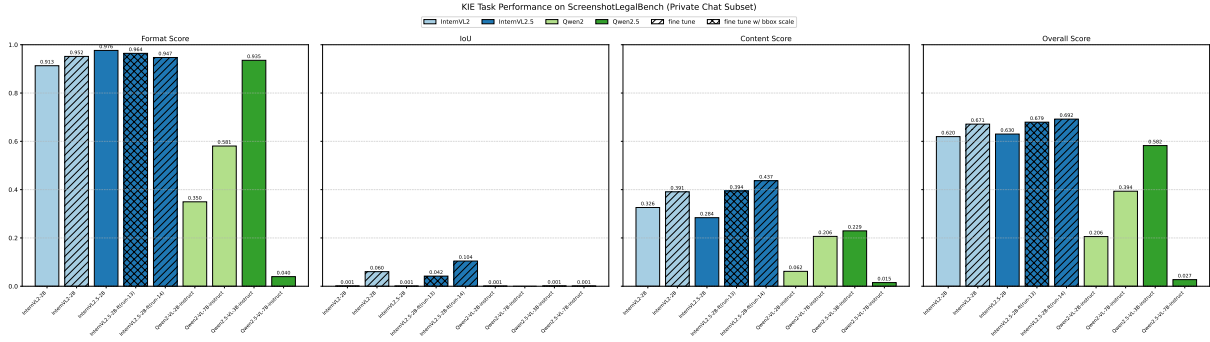


Figure 1: Comparison of KIE task performance on ScreenshotLegalBench (Private Chat Subset). Fine-tuned models are evaluated under pass@0, which better reflects output stability, while other models are evaluated under pass@1. Higher Format Scores indicate not only stronger adherence to structural output instructions, but also reflect the contribution of evaluation-side repair strategies designed to maximize compatibility with JSON-based outputs.

with strong format validity (0.9764). However, spatial alignment remains a challenge for most models, with un-tuned models, including Qwen2-VL-7B-instruct(Wang et al., 2024), showing near-zero IoU, indicating poor spatial reasoning. Content accuracy also varies, with InternVL2.5-2B(Chen et al., 2024a) scoring 0.2839, while larger models like Qwen2.5-VL-7B-instruct(Bai et al., 2025) score much lower (0.0151).

These results highlight the complexity of generating structured legal outputs and the need for fine-tuning, which is further explored in the next section on dataset generalization.

**VQA Task Results** Unlike the KIE task that emphasizes local timeline reconstruction under privacy constraints and is evaluated with smaller models, the VQA task targets legal reasoning and evidentiary judgment, requiring more complex abstraction. Therefore, larger-scale vision-language models are included to better assess their legal understanding capabilities. We evaluate model performance on three sub-tasks in the VQA portion of ScreenshotLegalBench: (1) legal relevance classification (classify), (2) assessment of evidentiary format compliance (evidence), and (3) cause-of-action generation via multi-image reasoning (case\_text), which detail in 2. As shown in Table 11, classification performance is low across models ( $F1 < 0.05$ ), reflecting the difficulty of determining legal relevance without contextual cues. The *evidence* task remains especially challenging: even large models like Qwen2.5-VL-72B(Bai et al., 2025) achieve high recall (0.44) but near-zero precision, indicating poor understanding of legal format standards. In *case\_text* (Table 12), models fail to produce coherent multi-image legal sum-

maries. Larger models (e.g., 72B) do not outperform smaller, instruction-tuned variants, suggesting that scale alone is insufficient for legal abstraction.

In summary, these results yield three key insights: (1) current models struggle to assess evidentiary formality due to limited spatial and layout understanding, and (2) moderate-sized models with domain-specific tuning outperform larger zero-shot models on multi-image legal reasoning tasks.

## 5.4 Dataset Generalization Analysis

To assess the benefits of dataset-specific instruction tuning, we compare the performance of foundation models and their fine-tuned counterparts on the KIE task (Figure 1). Fine-tuned models exhibit significantly higher content accuracy and spatial alignment scores, particularly in average IoU, confirming that domain-specific fine-tuning enhances the model’s structural consistency and understanding of legal semantics.

## 5.5 Ablation Analysis

We conduct ablation studies to examine how different components and dataset design choices affect model performance on ScreenshotLegalBench, highlighting the benefits of structured supervision for both KIE and VQA tasks.

**Importance of Bounding Box Fields (KIE).** We further study the impact of providing message-level bounding boxes (message\_bbox) during training. As shown in Table 15, removing this field leads to a near-zero IoU and degraded overall scores, despite only slight changes in content accuracy. This suggests that spatial annotations are critical for enabling the model to align textual fields with their visual locations, which is vital for downstream ap-

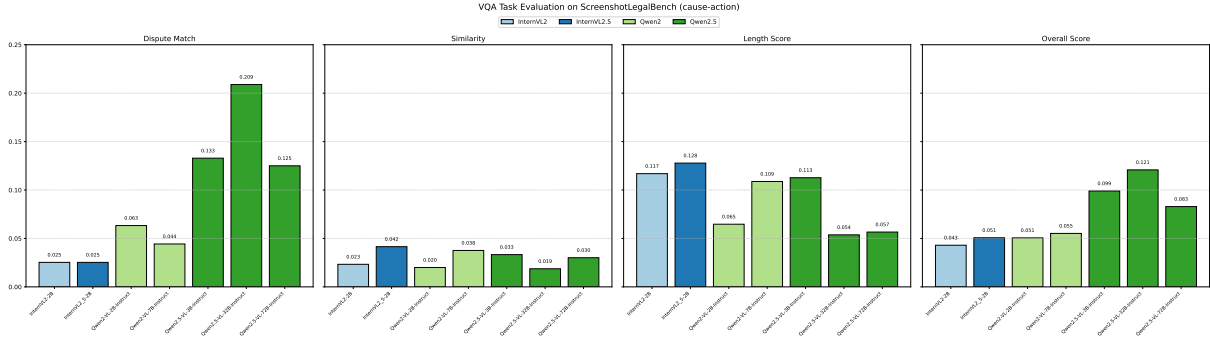


Figure 2: Comparison of vision-language models on the **case\_text** VQA task in ScreenshotLegalBench.

Model	Format Score	Avg. IoU	Content Score	Overall Score
InternVL2-2B (w/ bbox)	0.9384	0.0338	0.3703	0.6544
InternVL2-2B (w/o bbox)	0.7579	0	0.3934	0.5756

Table 3: Performance of InternVL2-2B on KIE task with and without message\_box bounding boxes

Setting	Task	Accuracy	Macro P	Macro R	F1 Score
<i>w/o KIE guidance</i>	classify	0.4286	0.3333	0.1428	0.1999
<i>+KIE-augmented input</i>	classify	<b>0.8333</b>	<b>0.5000</b>	<b>0.4167</b>	<b>0.4545</b>
<i>w/o KIE guidance</i>	evidence	0.0157	0.0013	0.0107	0.0019
<i>+KIE-augmented input</i>	evidence	<b>0.0187</b>	<b>0.0126</b>	<b>0.0606</b>	<b>0.0204</b>

Table 4: Ablation: Effect of structured KIE input on VQA tasks.

plications such as evidence localization and time-line reconstruction.

**Structured vs. Plain Inputs (VQA).** We assess the impact of structured perception by comparing two configurations: models predicting directly from raw screenshots (*w/o KIE guidance*) and those augmented with structured fields from the fine-tuned KIE module (*+KIE-augmented input*). As shown in Table 14, structured input significantly improves performance in the classify task—accuracy rises from 42.86% to 83.33%, and macro F1 more than doubles, demonstrating the value of upstream legal structuring. In contrast, the evidence task remains challenging. Despite slight gains from KIE augmentation, performance is low across the board. This suggests that models struggle to internalize evidentiary standards without domain-specific training, and that format validity requires not just structural cues but legal commonsense—still absent in current MLLMs. This ablation uses prompts enhanced with KIE outputs from our best-performing fine-tuned model, evaluated on Qwen-VL-MAX(Bai et al., 2023b). While limited to single-image inputs, future work should explore multi-image reasoning (e.g., case\_text) once token constraints are addressed. Full settings

and prompts are in Appendix D.4.

## 6 Conclusion

We present **ScreenshotLegalBench**, a new benchmark designed for legal evidence understanding in WeChat chat screenshots. It focuses on structured perception and legal reasoning tasks, offering insights into the challenges of multimodal legal AI. Despite limitations such as annotation consistency and scalability, the dataset provides a solid foundation for research in secure and practical legal applications. Baseline results reveal the difficulty of this task for current models, particularly under local deployment constraints. We encourage future work on improving scalability, layout robustness, and real-world adaptability. Positioned at the intersection of natural language processing, computer vision, and legal reasoning, LSEU holds substantial practical relevance. The dataset is released to support progress toward AI systems capable of interpreting digital legal evidence with accuracy and transparency.



## Limitations

ScreenshotLegalBench presents several limitations. Although annotations are verified by legal experts, The dataset exhibits category imbalance, particularly in cause-of-action types and funds-transfer content, as most samples originate from a narrow range of civil disputes. Timestamp labels rely on visual order assumptions (based on legal experts' experience), which may be unreliable in real-world scenarios. The evaluation favors structured JSON outputs and may penalize models with strong semantics but poor formatting.

## Ethics Statements

The ScreenshotLegalBench dataset is constructed using publicly available web data sourced from Gansu Provincial Digital Rule of Law Industry Research Academy . It was gathered with the intention of facilitating local, privacy-sensitive legal AI deployments, particularly for the KIE tasks. The dataset is designed to aid the automation of legal workflows while ensuring compliance with data privacy and confidentiality standards, especially in legal contexts.

To protect privacy, anonymization procedures are applied to identifiable data, such as the use of "avator\_1" and "avator\_0" to mask avatars in chat screenshots. These identifiers do not correspond to any real-world individual and are used solely for the purpose of maintaining privacy. However, due to the nature of web scraping, certain non-textual content in the images (e.g., emoticons, background images) and some personal information may not be entirely anonymized. Moreover, due to the structure of the raw web data, efforts to mask or obscure personal identifiers in the images (such as applying blur or cropping) may negatively affect the understanding of the primary content, as it could lead to distortion or removal of critical evidence.

The dataset's open-source release is aimed at enabling local model deployments for legal practitioners who may not have access to proprietary AI models due to regulatory or privacy concerns. This ensures that users can access advanced AI tools while retaining full control over the data and the models they develop.

While efforts have been made to ensure the privacy of the data, there may be inherent risks associated with using this dataset, especially regarding the presence of potentially noisy data, which could affect model performance. It is important for future

users of the dataset to be aware of these limitations and the trade-off between privacy preservation and data quality.

As part of the ongoing ethical commitment, we also provide a mechanism for obtaining consent for the data used. Any additional requests for sensitive data or further clarifications regarding the use of this dataset can be directed to the dataset's licensing terms, with the option to obtain permissions from the data providers where necessary.

The dataset is provided solely for academic research and benchmarking purposes. Commercial use or deployment in production environments is not permitted. We hope that the ScreenshotLegalBench dataset will contribute to the development of responsible, transparent, and privacy-conscious AI systems for legal tasks, while fostering further advancements in multimodal legal document understanding.

## Acknowledgments

We thank Beijing Hairun Tianrui (Zhengzhou) Law Firm for their support in annotation work. We also thank Zhangzhouliang and the Shanghai AI Lab for providing computational resources. Finally, we appreciate the anonymous reviewers for their valuable feedback.

## References

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023a. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023b. [Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond](#).
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Zhaowei Cai and Nuno Vasconcelos. 2017. [Cascade r-cnn: Delving into high quality object detection](#).

751	Nicolas Carion, Francisco Massa, Gabriel Synnaeve,	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for auto-</a>	806
752	Nicolas Usunier, Alexander Kirillov, and Sergey	<a href="#">matic evaluation of summaries</a> . In <i>Text Summariza-</i>	807
753	Zagoruyko. 2020. <a href="#">End-to-end object detection with</a>	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	808
754	<a href="#">transformers</a> . <i>European Conference on Computer</i>	Association for Computational Linguistics.	809
755	<i>Vision</i> .		
756	Ilias Chalkidis, Manos Fergadiotis, Prodromos Malaka-	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	810
757	siotis, Nikolaos Aletras, and Ion Androutsopoulos.	Lee. 2023. Visual instruction tuning. <i>NEURIPS</i> .	811
758	2020. <a href="#">Legal-bert: The muppets straight out of law</a>	Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis	812
759	<a href="#">school</a> .	Karatzas, Ernest Valveny, and C.V. Jawahar. 2022.	813
760	Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael	Infographicvqa. In <i>Proceedings of the IEEE/CVF</i>	814
761	Bommarito, Ion Androutsopoulos, Daniel Martin	<i>Winter Conference on Applications of Computer Vi-</i>	815
762	Katz, and Nikolaos Aletras. 2022. <a href="#">Lexglue: A bench-</a>	<i>sion (WACV)</i> , pages 1697–1706.	816
763	<a href="#">mark dataset for legal language understanding in en-</a>	Minesh Mathew, Dimosthenis Karatzas, and C.V. Jawa-	817
764	<a href="#">glish</a> .	har. 2021. Docvqa: A dataset for vqa on docu-	818
765	Yi-Wei Chen and Chih-Jen Lin. 2006. <a href="#">Combining SVMs</a>	ment images. In <i>Proceedings of the IEEE/CVF Win-</i>	819
766	<a href="#">with Various Feature Selection Strategies</a> , pages 315–	<i>ter Conference on Applications of Computer Vision</i>	820
767	324. Springer Berlin Heidelberg, Berlin, Heidelberg.	(WACV), pages 2200–2209.	821
768	Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu,	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	822
769	Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye,	Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evalu-</a>	823
770	Hao Tian, Zhaoyang Liu, Lixin Gu, Xuehui Wang,	<a href="#">ation of machine translation</a> . In <i>Proceedings of the</i>	824
771	Qingyun Li, Yimin Ren, Zixuan Chen, Jiapeng Luo,	<i>40th Annual Meeting on Association for Computa-</i>	825
772	Jiahao Wang, Tan Jiang, Bo Wang, Conghui He, Bo-	<i>tional Linguistics</i> , ACL ’02, page 311–318, USA.	826
773	tian Shi, Xingcheng Zhang, Han Lv, Yi Wang, Wenqi	Association for Computational Linguistics.	827
774	Shao, Pei Chu, Zhongying Tu, Tong He, Zhiyong Wu,	Joseph Redmon and Ali Farhadi. 2018. <a href="#">Yolov3: An</a>	828
775	Huipeng Deng, Jiaye Ge, Kai Chen, Min Dou, Lewei	<a href="#">incremental improvement</a> .	829
776	Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao,	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan	830
777	Jifeng Dai, and Wenhai Wang. 2024a. Expanding	Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,	831
778	performance boundaries of open-source multimodal	Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh	832
779	models with model, data, and test-time scaling. <i>arXiv</i>	Mariooryad, Yifan Ding, Xinyang Geng, Fred Al-	833
780	<i>preprint arXiv: 2412.05271</i> .	cobber, Roy Frostig, Mark Omernick, Lexi Walker,	834
781	Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye,	Cosmin Paduraru, Christina Sorokin, Andrea Tac-	835
782	Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi	chetti, Colin Gaffney, Samira Daruki, Olcan Ser-	836
783	Hu, Jiapeng Luo, Zheng Ma, et al. 2024b. How far	cinoglu, Zach Gleicher, Juliette Love, Paul Voigt-	837
784	are we to gpt-4v? closing the gap to commercial	laender, Rohan Jain, Gabriela Surita, Kareem Mo-	838
785	multimodal models with open-source suites. <i>arXiv</i>	hamed, Rory Blevins, Junwhan Ahn, Tao Zhu, Korn-	839
786	<i>preprint arXiv:2404.16821</i> .	raphop Kawintiranon, Orhan Firat, Yiming Gu, Yu-	840
787	Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo	jing Zhang, Matthew Rahtz, Manaal Faruqui, Natalie	841
788	Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,	Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui	842
789	Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu,	Zhu, Nobuyuki Morioka, Kevin Hui, Krishna Hari-	843
790	Yu Qiao, and Jifeng Dai. 2024c. Internvl: Scaling	dasan, Victor Campos, Mahdis Mahdieh, Mandy Guo,	844
791	up vision foundation models and aligning for generic	Samer Hassan, Kevin Kilgour, Arpi Vezzer, Heng-	845
792	visual-linguistic tasks. <i>Cvpr</i> .	Tze Cheng, Raoul de Liedekerke, Siddharth Goyal,	846
793	Shamanthak Hegde, Soumya Jahagirdar, and Shankar	Paul Barham, DJ Strouse, Seb Noury, Jonas Adler,	847
794	Gangisetty. 2023. Making the v in text-vqa mat-	Mukund Sundararajan, Sharad Vikram, Dmitry Lep-	848
795	ter. In <i>Proceedings of the IEEE/CVF Conference on</i>	ikhin, Michela Paganini, Xavier Garcia, Fan Yang,	849
796	<i>Computer Vision and Pattern Recognition (CVPR)</i>	Dasha Valter, Maja Trebacz, Kiran Vodrahalli, Chu-	850
797	<i>Workshops</i> , pages 5580–5588.	layuth Asawaroengchai, Roman Ring, Norbert Kalb,	851
798	Dan Hendrycks, Collin Burns, Anya Chen, and Spencer	Livio Baldini Soares, Siddhartha Brahma, David	852
799	Ball. 2021. <a href="#">Cuad: An expert-annotated nlp dataset</a>	Steiner, Tianhe Yu, Fabian Mentzer, Antoine He,	853
800	<a href="#">for legal contract review</a> .	Lucas Gonzalez, Bibo Xu, Raphael Lopez Kauf-	854
801	Jeewon Jeon, Woojun Kim, Whiyoun Jung, and	man, Laurent El Shafey, Junhyuk Oh, Tom Hennigan,	855
802	Youngchul Sung. 2022. Maser: Multi-agent rein-	George van den Driessche, Seth Odoom, Mario Lucic,	856
803	forcement learning with subgoals generated from	Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan,	857
804	experience replay buffer. <i>arXiv preprint arXiv:</i>	Santiago Ontanon, Luheng He, Denis Teplyashin,	858
805	<i>2206.10607</i> .	Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis	859
		Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan Yeh,	860
		Aakanksha Chowdhery, Yang Xu, Mehran Kazemi,	861
		Ehsan Amid, Anastasia Petrushkina, Kevin Swersky,	862
		Ali Khodaei, Gowoon Chen, Chris Larkin, Mario	863
		Pinto, Geng Yan, Adria Puigdomenech Badia, Piyush	864

865	Patil, Steven Hansen, Dave Orr, Sebastien M. R.	hyay, Anudhyan Boral, Lisa Anne Hendricks, Corey	928
866	Arnold, Jordan Grimstad, Andrew Dai, Sholto Dou-	Fry, Josip Djolonga, Yi Su, Jake Walker, Jane La-	929
867	glas, Rishika Sinha, Vikas Yadav, Xi Chen, Elena Gri-	banowski, Ronny Huang, Vedant Misra, Jeremy	930
868	bovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel,	Chen, RJ Skerry-Ryan, Avi Singh, Shruti Rijh-	931
869	Paul Komarek, Sophia Austin, Sebastian Borgeaud,	wani, Dian Yu, Alex Castro-Ros, Beer Changpinyo,	932
870	Linda Friso, Abhimanyu Goyal, Ben Caine, Kris	Romina Datta, Sumit Bagri, Arnar Mar Hrafnkels-	933
871	Cao, Da-Woon Chung, Matthew Lamm, Gabe Barth-	son, Marcello Maggioni, Daniel Zheng, Yury Sul-	934
872	Maron, Thais Kagohara, Kate Olszewska, Mia Chen,	sky, Shaobo Hou, Tom Le Paine, Antoine Yang,	935
873	Kaushik Shivakumar, Rishabh Agarwal, Harshal	Jason Riesa, Dominika Rogozinska, Dror Marcus,	936
874	Godhia, Ravi Rajwar, Javier Snaider, Xerxes Doti-	Dalia El Badawy, Qiao Zhang, Luyu Wang, Helen	937
875	walla, Yuan Liu, Aditya Barua, Victor Ungureanu,	Miller, Jeremy Greer, Lars Lowe Sjos, Azade Nova,	938
876	Yuan Zhang, Bat-Orgil Batsaikhan, Mateo Wirth,	Heiga Zen, Rahma Chaabouni, Mihaela Rosca, Jiepu	939
877	James Qin, Ivo Danihelka, Tulsee Doshi, Martin	Jiang, Charlie Chen, Ruibo Liu, Tara Sainath, Maxim	940
878	Chadwick, Jilin Chen, Sanil Jain, Quoc Le, Ar-	Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh	941
879	jun Kar, Madhu Gurumurthy, Cheng Li, Ruoxin	Newlan, Zeyncep Cankara, Soo Kwak, Yunhan Xu,	942
880	Sang, Fangyu Liu, Lampros Lamprou, Rich Munoz,	Phil Chen, Andy Coenen, Clemens Meyer, Katerina	943
881	Nathan Lintz, Harsh Mehta, Heidi Howard, Mal-	Tsihla, Ada Ma, Juraj Gottweis, Jinwei Xing, Chen-	944
882	colm Reynolds, Lora Aroyo, Quan Wang, Lorenzo	jie Gu, Jin Miao, Christian Frank, Zeynep Cankara,	945
883	Blanco, Albin Cassirer, Jordan Griffith, Dipanjan	Sanjay Ganapathy, Ishita Dasgupta, Steph Hughes-	946
884	Das, Stephan Lee, Jakob Sygnowski, Zach Fisher,	Fitt, Heng Chen, David Reid, Keran Rong, Hongmin	947
885	James Besley, Richard Powell, Zafarali Ahmed, Do-	Fan, Joost van Amersfoort, Vincent Zhuang, Aaron	948
886	minik Paulus, David Reitter, Zalan Borsos, Rishabh	Cohen, Shixiang Shane Gu, Anhad Mohananey,	949
887	Joshi, Aedan Pope, Steven Hand, Vittorio Selo, Vi-	Anastasija Ilic, Taylor Tobin, John Wieting, Anna	950
888	han Jain, Nikhil Sethi, Megha Goel, Takaki Makino,	Bortsova, Phoebe Thacker, Emma Wang, Emily	951
889	Rhys May, Zhen Yang, Johan Schalkwyk, Christina	Caveness, Justin Chiu, Eren Sezener, Alex Kaskasoli,	952
890	Butterfield, Anja Hauth, Alex Goldin, Will Hawkins,	Steven Baker, Katie Millican, Mohamed Elhawaty,	953
891	Evan Senter, Sergey Brin, Oliver Woodman, Mar-	Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun	954
892	vin Ritter, Eric Noland, Minh Giang, Vijay Bolina,	Dai, Wenhao Jia, Matthew Wiethoff, Elnaz Davoodi,	955
893	Lisa Lee, Tim Blyth, Ian Mackinnon, Machel Reid,	Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel	956
894	Obaid Sarvana, David Silver, Alexander Chen, Lily	Gao, Golan Pundak, Susan Zhang, Michael Azzam,	957
895	Wang, Loren Maggiore, Oscar Chang, Nithya At-	Khe Chai Sim, Sergi Caelles, James Keeling, Ab-	958
896	taluri, Gregory Thornton, Chung-Cheng Chiu, Os-	hanshu Sharma, Andy Swing, YaGuang Li, Chenxi	959
897	kar Bunyan, Nir Levine, Timothy Chung, Evgenii	Liu, Carrie Grimes Bostock, Yamini Bansal, Zachary	960
898	Eltyshev, Xiance Si, Timothy Lillicrap, Demetra	Nado, Ankesh Anand, Josh Lipschultz, Abhijit Kar-	961
899	Brady, Vaibhav Aggarwal, Boxi Wu, Yuanzhong Xu,	markar, Lev Proleev, Abe Ittycheriah, Soheil Has-	962
900	Ross McIlroy, Kartikeya Badola, Paramjit Sandhu,	sas Yeganeh, George Polovets, Aleksandra Faust,	963
901	Erica Moreira, Wojciech Stokowiec, Ross Hems-	Jiao Sun, Alban Rustemi, Pen Li, Rakesh Shivanna,	964
902	ley, Dong Li, Alex Tudor, Pranav Shyam, Elahe	Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh	965
903	Rahimtoroghi, Salem Haykal, Pablo Sprechmann,	Baddepudi, Sebastian Krause, Emilio Parisotto, Radu	966
904	Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki,	Soricut, Zheng Xu, Dawn Bloxwich, Melvin John-	967
905	Kalpesh Krishna, Xiao Wu, Alexandre Frechette,	son, Behnam Neyshabur, Justin Mao-Jones, Ren-	968
906	Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang,	shen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur	969
907	Thi Avrahami, Ye Zhang, Emanuel Taropa, Hanzhao	Guez, Constant Segal, Duc Dung Nguyen, James	970
908	Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano,	Svensson, Le Hou, Sarah York, Kieran Milan, So-	971
909	HyunJeong Choe, Alex Tomala, Chalence Safranek-	phie Bridgers, Wiktor Gworek, Marco Tagliasacchi,	972
910	Shrader, Nora Kassner, Mantas Pajarskas, Matt	James Lee-Thorp, Michael Chang, Alexey Guseynov,	973
911	Harvey, Sean Sechrist, Meire Fortunato, Christina	Ale Jakse Hartman, Michael Kwong, Ruizhe Zhao,	974
912	Lyu, Gamaleldin Elsayed, Chenkai Kuang, James	Sheleem Kashem, Elizabeth Cole, Antoine Miech,	975
913	Lottes, Eric Chu, Chao Jia, Chih-Wei Chen, Pe-	Richard Tanburn, Mary Phuong, Filip Pavetic, Se-	976
914	ter Humphreys, Kate Baumli, Connie Tao, Rajku-	bastien Cevey, Ramona Comanescu, Richard Ives,	977
915	mar Samuel, Cicero Nogueira dos Santos, Anders	Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang,	978
916	Andreassen, Nemanja Rakićević, Dominik Grewe,	Mariko Iinuma, Clara Huiyi Hu, Aurko Roy, Shaan	979
917	Aviral Kumar, Stephanie Winkler, Jonathan Caton,	Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel	980
918	Andrew Brock, Sid Dalmia, Hannah Sheahan, Iain	Saputro, Anita Gergely, Steven Zheng, Dawei Jia,	981
919	Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Fer-	Ioannis Antonoglou, Adam Sadovsky, Shane Gu,	982
920	yal Behbahani, Flavien Prost, Yanhua Sun, Artiom	Yingying Bi, Alek Andreev, Sina Samangooei, Mina	983
921	Myaskovsky, Thanumalayan Sankaranarayana Pillai,	Khan, Tomas Kocisky, Angelos Filos, Chintu Ku-	984
922	Dan Hurt, Angeliki Lazaridou, Xi Xiong, Ce Zheng,	mar, Colton Bishop, Adams Yu, Sarah Hodgkin-	985
923	Fabio Pardo, Xiaowei Li, Dan Horgan, Joe Stanton,	son, Sid Mittal, Premal Shah, Alexandre Moufarek,	986
924	Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu	Yong Cheng, Adam Bloniarz, Jaehoon Lee, Pedram	987
925	Wang, Basil Mustafa, Albert Webson, Hyo Lee, Ro-	Pejman, Paul Michel, Stephen Spencer, Vladimir	988
926	han Anil, Martin Wicke, Timothy Dozat, Abhishek	Feinberg, Xuehan Xiong, Nikolay Savinov, Char-	989
927	Sinha, Enrique Piqueras, Elahe Dabir, Shyam Upad-	lotte Smith, Siamak Shakeri, Dustin Tran, Mary	990

991	Chesus, Bernd Bohnet, George Tucker, Tamara von	1054
992	Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa,	1055
993	Ambrose Slone, Kedar Soparkar, Disha Shrivastava,	1056
994	James Cobon-Kerr, Michael Sharman, Jay Pavagadhi,	1057
995	Carlos Araya, Karolis Misiunas, Nimesh Ghelani,	1058
996	Michael Laskin, David Barker, Qiujia Li, Anton	1059
997	Briukhov, Neil Houlsby, Mia Glaese, Balaji Laksh-	1060
998	minarayanan, Nathan Schucher, Yunhao Tang, Eli	1061
999	Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria	1062
1000	Recasens, Guangda Lai, Alberto Magni, Nicola De	1063
1001	Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay,	1064
1002	Mostafa Dehghani, Jenny Brennan, Yifan He, Kelvin	1065
1003	Xu, Yang Gao, Carl Saroufim, James Molloy, Xinyi	1066
1004	Wu, Seb Arnold, Solomon Chang, Julian Schrit-	1067
1005	twieser, Elena Buchatskaya, Soroush Radpour, Mar-	1068
1006	tin Polacek, Skye Giordano, Ankur Bapna, Simon	1069
1007	Tokumine, Vincent Hellendoorn, Thibault Sottiaux,	1070
1008	Sarah Cogan, Aliaksei Severyn, Mohammad Saleh,	1071
1009	Shantanu Thakoor, Laurent Shefey, Siyuan Qiao,	1072
1010	Meenu Gaba, Shuo yin Chang, Craig Swanson, Biao	1073
1011	Zhang, Benjamin Lee, Paul Kishan Rubenstein, Gan	1074
1012	Song, Tom Kwiatkowski, Anna Koop, Ajay Kan-	1075
1013	nan, David Kao, Parker Schuh, Axel Stjerngren, Gol-	1076
1014	naz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Fe-	1077
1015	lipe Tiengo Ferreira, Aishwarya Kamath, Ted Kli-	1078
1016	menko, Ken Franko, Kefan Xiao, Indro Bhattacharya,	1079
1017	Miteyan Patel, Rui Wang, Alex Morris, Robin	1080
1018	Strudel, Vivek Sharma, Peter Choy, Sayed Hadi	1081
1019	Hashemi, Jessica Landon, Mara Finkelstein, Priya	1082
1020	Jhakra, Justin Frye, Megan Barnes, Matthew Mauer,	1083
1021	Dennis Daun, Khuslen Baatarsukh, Matthew Tung,	1084
1022	Wael Farhan, Henryk Michalewski, Fabio Viola, Fel-	1085
1023	ix de Chaumont Quitry, Charline Le Lan, Tom Hud-	1086
1024	son, Qingze Wang, Felix Fischer, Ivy Zheng, Elspeth	1087
1025	White, Anca Dragan, Jean baptiste Alayrac, Eric Ni,	1088
1026	Alexander Pritzel, Adam Iwanicki, Michael Isard,	1089
1027	Anna Bulanova, Lukas Zilka, Ethan Dyer, Deven-	1090
1028	dra Sachan, Srivatsan Srinivasan, Hannah Mucken-	1091
1029	hirm, Honglong Cai, Amol Mandhane, Mukarram	1092
1030	Tariq, Jack W. Rae, Gary Wang, Kareem Ayoub,	1093
1031	Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris	1094
1032	Alberti, Dan Garrette, Kashyap Krishnakumar, Mai	1095
1033	Gimenez, Anselm Levskaya, Daniel Sohn, Josip	1096
1034	Matak, Inaki Iturrate, Michael B. Chang, Jackie Xi-	1097
1035	ang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian	1098
1036	Hutter, Vahab Mirrokni, Nanxin Chen, Kaisheng	1099
1037	Yao, Zoltan Egyed, Francois Galilee, Tyler Liechty,	1100
1038	Praveen Kallakuri, Evan Palmer, Sanjay Ghemawat,	1101
1039	Jasmine Liu, David Tao, Chloe Thornton, Tim Green,	1102
1040	Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan	1103
1041	Tan, Noah Fiedel, Hongkun Yu, Ed Chi, Alexan-	1104
1042	der Neitz, Jens Heitkaemper, Anu Sinha, Denny	1105
1043	Zhou, Yi Sun, Charbel Kaed, Brice Hulse, Swa-	1106
1044	roop Mishra, Maria Georgaki, Sneha Kudugunta,	1107
1045	Clement Farabet, Izhak Shafran, Daniel Vlasic, An-	1108
1046	ton Tsitsulin, Rajagopal Ananthanarayanan, Alen	1109
1047	Carin, Guolong Su, Pei Sun, Shashank V, Gabriel	1110
1048	Carvajal, Josef Broder, Iulia Comsa, Alena Repina,	1111
1049	William Wong, Warren Weilun Chen, Peter Hawkins,	1112
1050	Egor Filonov, Lucia Loher, Christoph Hirschall,	1113
1051	Weiye Wang, Jingchen Ye, Andrea Burns, Hardie	1114
1052	Cate, Diana Gage Wright, Federico Piccinini, Lei	1115
1053	Zhang, Chu-Cheng Lin, Ionel Gog, Yana Kulizh-	1116
	skaya, Ashwin Sreevatsa, Shuang Song, Luis C.	
	Cobo, Anand Iyer, Chetan Tekur, Guillermo Gar-	
	rido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven	
	Zheng, Hui Li, Ananth Agarwal, Christel Ngani,	
	Kati Goshvadi, Rebeca Santamaria-Fernandez, Woj-	
	ciech Fica, Xinyun Chen, Chris Gorgolewski, Sean	
	Sun, Roopal Garg, Xinyu Ye, S. M. Ali Eslami,	
	Nan Hua, Jon Simon, Pratik Joshi, Yelin Kim, Ian	
	Tenney, Sahitya Potluri, Lam Nguyen Thiet, Quan	
	Yuan, Florian Luisier, Alexandra Chronopoulou, Sal-	
	vatore Scellato, Praveen Srinivasan, Minmin Chen,	
	Vinod Koverkathu, Valentin Dalibard, Yaming Xu,	
	Brennan Saeta, Keith Anderson, Thibault Sellam,	
	Nick Fernando, Fantine Huot, Junehyuk Jung, Mani	
	Varadarajan, Michael Quinn, Amit Raul, Maigo Le,	
	Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha	
	Bullard, Achintya Singhal, Thang Luong, Boyu	
	Wang, Sujeewan Rajayogam, Julian Eisenschlos,	
	Johnson Jia, Daniel Finchelstein, Alex Yakubovich,	
	Daniel Balle, Michael Fink, Sameer Agarwal, Jing	
	Li, Dj Dvijotham, Shalini Pal, Kai Kang, Jaclyn	
	Konzelmann, Jennifer Beattie, Olivier Dousse, Diane	
	Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy	
	Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Kryst-	
	tal Kallarackal, Rosanne Liu, Denis Vnukov, Neera	
	Vats, Luca Invernizzi, Mohsen Jafari, Huanjie Zhou,	
	Lilly Taylor, Jennifer Prendki, Marcus Wu, Tom	
	Eccles, Tianqi Liu, Kavya Kopparapu, Francoise	
	Beaufays, Christof Angermueller, Andreea Marzoca,	
	Shourya Sarcar, Hilal Dib, Jeff Stanway, Frank Per-	
	bet, Nejc Trdin, Rachel Sterneck, Andrey Khor-	
	lin, Dinghua Li, Xihui Wu, Sonam Goenka, David	
	Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou,	
	Yaxin Liu, Yannie Liang, Anais White, Yunjie Li,	
	Shreya Singh, Sanaz Bahargam, Mark Epstein, Su-	
	joy Basu, Li Lao, Adnan Ozturk, Carl Crous, Alex	
	Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna	
	Walton, Lucas Dixon, Ming Zhang, Amir Globerson,	
	Grant Uy, Andrew Bolt, Olivia Wiles, Milad	
	Nasr, Ilia Shumailov, Marco Selvi, Francesco Pic-	
	cinno, Ricardo Aguilar, Sara McCarthy, Misha Khal-	
	man, Mrinal Shukla, Vlado Galic, John Carpen-	
	ter, Kevin Vilella, Haibin Zhang, Harry Richard-	
	son, James Martens, Matko Bosnjak, Shreyas Ram-	
	mohan Belle, Jeff Seibert, Mahmoud Alnahlawi,	
	Brian McWilliams, Sankalp Singh, Annie Louis,	
	Wen Ding, Dan Popovici, Lenin Simicich, Laura	
	Knight, Pulkit Mehta, Nishesh Gupta, Chongyang	
	Shi, Saaber Fatehi, Jovana Mitrovic, Alex Grills,	
	Joseph Pagadora, Tsendsuren Munkhdalai, Dessie	
	Petrova, Danielle Eisenbud, Zhishuai Zhang, Damion	
	Yates, Bhavishya Mittal, Nilesh Tripuraneni, Yan-	
	nis Assael, Thomas Brovelli, Prateek Jain, Miha-	
	jlo Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolf-	
	gang Macherey, Ravin Kumar, Jun Xu, Haroon	
	Qureshi, Gheorghe Comanici, Jeremy Wiesner, Zhi-	
	tao Gong, Anton Ruddock, Matthias Bauer, Nick	
	Felt, Anirudh GP, Anurag Arnab, Dustin Zelle,	
	Jonas Rothfuss, Bill Rosgen, Ashish Shenoy, Bryan	
	Seybold, Xinjian Li, Jayaram Mudigonda, Goker	
	Erdogan, Jiawei Xia, Jiri Simsa, Andrea Michi,	
	Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell,	
	Carey Radebaugh, Andre Elisseeff, Pedro Valen-	



1117	zuela, Kay McKinney, Kim Paterson, Albert Cui, Eri	
1118	Latorre-Chimoto, Solomon Kim, William Zeng, Ken	
1119	Durden, Priya Ponnappalli, Tiberiu Sosea, Christo-	
1120	pher A. Choquette-Choo, James Manyika, Brona	
1121	Robenek, Harsha Vashisht, Sebastien Pereira, Hoi	
1122	Lam, Marko Velic, Denese Owusu-Afriyie, Kather-	
1123	ine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu,	
1124	Jane Park, Balaji Venkatraman, Alice Talbert, Lam-	
1125	bert Rosique, Yuchung Cheng, Andrei Sozanschi,	
1126	Adam Paszke, Praveen Kumar, Jessica Austin, Lu Li,	
1127	Khalid Salama, Bartek Perz, Wooyeol Kim, Nandita	
1128	Dukkipati, Anthony Baryshnikov, Christos Kapla-	
1129	nis, XiangHai Sheng, Yuri Chervonyi, Caglar Unlu,	
1130	Diego de Las Casas, Harry Askham, Kathryn Tun-	
1131	yasuvunakool, Felix Gimeno, Siim Poder, Chester	
1132	Kwak, Matt Mieczkowski, Vahab Mirrokni, Alek	
1133	Dimitriev, Aaron Parisi, Dangyi Liu, Tomy Tsai,	
1134	Toby Shevlane, Christina Kouridi, Drew Garmon,	
1135	Adrian Goedeckemeyer, Adam R. Brown, Anitha Vi-	
1136	jayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang,	
1137	Sara Mc Carthy, Jay Hoover, Lucy Kim, Sandeep	
1138	Kumar, Wei Chen, Courtney Biles, Garrett Bingham,	
1139	Evan Rosen, Lisa Wang, Qijun Tan, David Engel,	
1140	Francesco Pongetti, Dario de Cesare, Dongseong	
1141	Hwang, Lily Yu, Jennifer Pullman, Srin Narayanan,	
1142	Kyle Levin, Siddharth Gopal, Megan Li, Asaf Aha-	
1143	roni, Trieu Trinh, Jessica Lo, Norman Casagrande,	
1144	Roopali Vij, Loic Matthey, Bramandia Ramadhana,	
1145	Austin Matthews, CJ Carey, Matthew Johnson, Kre-	
1146	mena Goranova, Rohin Shah, Shereen Ashraf, King-	
1147	shuk Dasgupta, Rasmus Larsen, Yicheng Wang, Man-	
1148	ish Reddy Vuyyuru, Chong Jiang, Joana Ijazi, Kazuki	
1149	Osawa, Celine Smith, Ramya Sree Boppana, Tay-	
1150	lan Bilal, Yuma Koizumi, Ying Xu, Yasemin Altun,	
1151	Nir Shabat, Ben Bariach, Alex Korchemniy, Kiam	
1152	Choo, Olaf Ronneberger, Chimezie Iwuanyanwu,	
1153	Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene	
1154	Cai, Shariq Iqbal, Martin Sundermeyer, Zhe Chen,	
1155	Elie Bursztejn, Chaitanya Malaviya, Fadi Biadsy,	
1156	Prakash Shroff, Inderjit Dhillon, Tejas Latkar, Chris	
1157	Dyer, Hannah Forbes, Massimo Nicosia, Vitaly Niko-	
1158	laev, Somer Greene, Marin Georgiev, Pidong Wang,	
1159	Nina Martin, Hanie Sedghi, John Zhang, Praseem	
1160	Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang, Ji-	
1161	ageng Zhang, Viorica Patraucean, Dayou Du, Igor	
1162	Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey, Abhi	
1163	Mohan, Janek Nowakowski, Vlad-Doru Ion, Nan	
1164	Wei, Reiko Tojo, Maria Abi Raad, Drew A. Hud-	
1165	son, Vaishakh Keshava, Shubham Agrawal, Kevin	
1166	Ramirez, Zhichun Wu, Hoang Nguyen, Ji Liu, Mad-	
1167	havi Sewak, Bryce Petrini, DongHyun Choi, Ivan	
1168	Philips, Ziyue Wang, Ioana Bica, Ankush Garg,	
1169	Jarek Wilkiewicz, Priyanka Agrawal, Xiaowei Li,	
1170	Danhao Guo, Emily Xue, Naseer Shaik, Andrew	
1171	Leach, Sadh MNM Khan, Julia Wiesinger, Sammy	
1172	Jerome, Abhishek Chakladar, Alek Wenjiao Wang,	
1173	Tina Ornduff, Folake Abu, Alireza Ghaffarkhah, Mar-	
1174	cus Wainwright, Mario Cortes, Frederick Liu, Joshua	
1175	Maynez, Andreas Terzis, Pouya Samangouei, Riham	
1176	Mansour, Tomasz Kępa, François-Xavier Aubet, An-	
1177	ton Algymr, Dan Banica, Agoston Weisz, Andras	
1178	Orban, Alexandre Senges, Ewa Andrejczuk, Mark	
1179	Geller, Niccolo Dal Santo, Valentin Anklin, Majd Al	
	Meray, Martin Baeuml, Trevor Strohman, Junwen	1180
	Bai, Slav Petrov, Yonghui Wu, Demis Hassabis, Ko-	1181
	ray Kavukcuoglu, Jeff Dean, and Oriol Vinyals. 2024.	1182
	<a href="#">Gemini 1.5: Unlocking multimodal understanding</a>	1183
	<a href="#">across millions of tokens of context.</a>	1184
	Don Tugener, Pius von Däniken, Thomas Peetz, and	1185
	Mark Cieliebak. 2020. <a href="#">LEDGAR: A large-scale</a>	1186
	<a href="#">multi-label corpus for text classification of legal pro-</a>	1187
	<a href="#">visions in contracts.</a> In <i>Proceedings of the Twelfth</i>	1188
	<i>Language Resources and Evaluation Conference,</i>	1189
	pages 1235–1241, Marseille, France. European Lan-	1190
	guage Resources Association.	1191
	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-	1192
	hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin	1193
	Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei	1194
	Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang	1195
	Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-	1196
	vl: Enhancing vision-language model’s perception	1197
	of the world at any resolution. <i>arXiv preprint arXiv:</i>	1198
	<i>2409.12191.</i>	1199
	Steven H Wang, Maksim Zubkov, Kexin Fan, Sarah	1200
	Harrell, Yuyang Sun, Wei Chen, Andreas Plesner,	1201
	and Roger Wattenhofer. 2025. Acord: An expert-	1202
	annotated retrieval dataset for legal contract drafting.	1203
	<i>arXiv preprint arXiv:2501.06582.</i>	1204
	Chaojun Xiao, Haoxi Zhong, Zhipeng Guo, Cunchao Tu,	1205
	Zhiyuan Liu, Maosong Sun, Yansong Feng, Xianpei	1206
	Han, Zhen Hu, Heng Wang, and Jianfeng Xu. 2018.	1207
	Cail2018: A large-scale legal dataset for judgment	1208
	prediction. <i>arXiv preprint arXiv: 1807.02478.</i>	1209
	Lucia Zheng, Neel Guha, Brandon R. Anderson, Peter	1210
	Henderson, and Daniel E. Ho. 2021. <a href="#">When does pre-</a>	1211
	<a href="#">training help? assessing self-supervised learning for</a>	1212
	<a href="#">law and the casehold dataset of 53,000+ legal hold-</a>	1213
	<a href="#">ings.</a> In <i>Proceedings of the Eighteenth International</i>	1214
	<i>Conference on Artificial Intelligence and Law, ICAIL</i>	1215
	’21, page 159–168, New York, NY, USA. Association	1216
	for Computing Machinery.	1217
	Haoxiang Zhong, Chaojun Xiao, Cunchao Tu, T. Zhang,	1218
	Zhiyuan Liu, and Maosong Sun. 2019. <a href="#">Jec-qa: A</a>	1219
	<a href="#">legal-domain question answering dataset.</a> <i>AAAI Con-</i>	1220
	<i>ference on Artificial Intelligence.</i>	1221
	<b>A Appendix : Wechat Screenshots</b>	1222
	<b>Evidence Legal Processing Task</b>	1223
	<b>Workflow of Legal Practitioners</b>	1224
	Legal professionals—especially plaintiff-side attor-	1225
	neys—follow a structured yet iterative workflow	1226
	when preparing WeChat screenshots for courtroom	1227
	use. As illustrated in Figure 3, this process typically	1228
	unfolds across five interrelated stages: (1) prelimi-	1229
	nary screening, (2) timeline reconstruction, (3) key	1230
	information extraction, (4) legal grounding, and	1231
	(5) evidence cataloging and strategy formulation.	1232
	These steps often involve back-and-forth revision	1233



as new information emerges or legal interpretations are refined.

**ScreenshotLegalBench** mirrors this legal workflow by decomposing the Legal Screenshot Evidence Understanding (LSEU) task into modular components. Specifically:

- **Stages 1–3** are supported by the *KIE* task and local *VQA*, which detect speaker turns, extract transaction details, and identify message timestamps.
- **Stage 4** is operationalized as the *case\_text* generation task, where models synthesize multiple screenshots to infer the dispute’s legal basis.

As shown in Figure 3, the step of “summarization of legal facts” corresponds directly to our case reasoning module. This process is annotated in our dataset, though for evaluation purposes we adopt a simplified scheme: rather than scoring long-form legal texts, we assess whether the predicted output includes correct cause-of-action categories.

### Visual-Legal Mapping of Evidence Elements

To enable automated legal understanding, ScreenshotLegalBench captures a range of heterogeneous visual elements in chat screenshots and aligns them with their legal semantics:

**Avatars and nicknames** provide identity signals for speaker verification. **Timestamps** serve as anchors for timeline reconstruction and causal ordering. **Chat bubbles and textual content** carry the core of intent and factual statements. **Image blocks** may represent either expressive content or direct legal exhibits. **Transfer notifications** frequently denote contractual performance or financial disputes. **Files and attachments** often indicate delivery obligations in cooperative agreements. **Emojis and quoted speech** encode attitudes, denials, or acknowledgments, which can be crucial for interpreting legal intent.

By aligning such fragmented visual-textual content with structured legal interpretation, ScreenshotLegalBench establishes a tractable framework for multimodal legal AI—grounded in the actual evidentiary workflows of judicial practice.

## B Appendix : Dataset Implementation Details

### B.1 Annotation Schema

As shown in Table 5, ScreenshotLegalBench adopts a five-level hierarchical annotation schema de-

signed to meet the diverse demands of multimodal legal tasks. This schema integrates visual layout, semantic content, and legal reasoning to support both structured extraction and high-level judicial analysis.

Level 1 annotates global attributes of each screenshot, including legality, chat type, and case category, to facilitate filtering and legal classification. Level 2 focuses on layout elements such as message bubbles, avatars, and timestamps, using bounding boxes to support object detection and automated parsing. To reduce annotation cost while preserving effectiveness, full layout annotations are only provided for 50 screenshots. Level 3 structures message-level fields including timestamp, content, and speaker, supporting downstream dialogue reconstruction and evidence linkage. Level 4 further normalizes and formats semantic fields, such as monetary transfers and file metadata, ensuring compatibility with legal expression standards. Level 5 introduces VQA annotations, targeting both global and local reasoning about legal validity, intent, and evidentiary value (see Table 6 for examples).

This layered design ensures ScreenshotLegalBench supports both low-level structure-aware pre-training and high-level legal understanding, making it suitable for retrieval, reasoning, and structured generation tasks.

### B.2 Object Detection For Sceenshots Layout Training

We adopt DETR as the screenshot layout detector, fine-tuned on our augmented WeChat chat screenshot dataset to detect message bubbles, avatars, timestamps, and other UI elements. Training configuration detail in Table 7

### B.3 Timestamp Imputation

To enrich the temporal context of screenshots lacking explicit timestamps, we introduce a sliding-window-based imputation mechanism at the screenshot level. Considering real-world scenarios where multiple conversations may coexist and screenshot order can be disrupted, we first perform session-level clustering using OCR-extracted chat titles, followed by intra-session timestamp sorting and imputation.

Screenshots are categorized into three types based on the presence of timestamps: (1) If a single timestamp is detected, it is directly assigned as

Field Type	Attribute Name	Annotation Detail	Level	Annotation Format
<b>Global Properties</b>	Screenshot Validity	Whether it is a standardized chat screenshot	Level 1	Enumeration
	Chat Type	Group or private conversation		Enumeration
	Legal Relevance	Whether the screenshot has legal implications		Enumeration
	Case Type	Preliminary case classification (e.g., loan, contract)		Enumeration
<b>Layout Elements</b>	Avatar	avatar_bbox	Level 2	Bounding Box
	Message Bubble	message_bbox		Bounding Box + Text
	Chat Title / Group Name	header_bbox		Bounding Box
	Nickname	nickname_bbox		Bounding Box + Text
	Timestamp Region	timestamp_bbox		Bounding Box + Text
	Transfer Block	transfer_bbox		Bounding Box + Text
	File Block	file_bbox		Bounding Box + Text
	Image Block	image_bbox		Bounding Box + Category
	Emoji / Meme	meme_bbox		Bounding Box + Text
	Voice Message	voice_bbox		Bounding Box
	Recall Prompt	withdraw_bbox		Bounding Box
	Translation Block	translate_bbox		Bounding Box + Text
	Quote / Comment	comment_bbox		Bounding Box + Text
	Failed Message	unpassed_message		Bounding Box + Text
	Other Elements	other		Bounding Box
<b>Message Fields</b>	Speaker	speaker	Level 3	Enumeration / String
	Message Time	timestamp		Time String
	Message Content	content		Raw Text
<b>Semantic Fields</b>	Transfer Info	transfer	Level 4	Normalized String
	File Name and Type	file		Normalized String
	Image Description	image		Generated Text
	Emoji Polarity	meme		Enumeration
<b>Legal QA Fields</b>	Intent Analysis	“What intent is expressed in this message?”	Level 5	Text QA
	Legal Reasoning	“Please analyze the legal implications of this situation.”		Text QA
	Transfer Nature	“What is the legal nature of the received transfer?”		Text QA
	File Legality	“Is the sent file direct legal evidence?”		Text QA

Table 5: Hierarchical annotation schema of ScreenshotLegalBench, covering five levels from global classification to legal reasoning.

the screenshot’s temporal feature. (2) For multiple timestamps, we apply a heuristic to assign each timestamp to subsequent messages and compute the screenshot’s time value as the arithmetic mean of all detected timestamps. (3) For screenshots with no timestamp, we estimate the time feature based on its position in the session sequence via a sliding average of neighboring screenshots.

Formally, let  $t_{i,j}$  denote the imputed time for the  $j$ -th screenshot in the  $i$ -th session. Its value is computed as:

$$t_{i,j} = \frac{1}{k} \sum_{l=1}^k t_{i,j-l} \quad (1)$$

where  $k$  is the sliding window size, controlling how many preceding screenshots contribute to the

estimation. This process is performed within each conversation thread, and the resulting timestamp is propagated to all messages in the corresponding screenshot for downstream context modeling and temporal reasoning.

For tasks requiring global temporal order (e.g., event timeline reconstruction), all screenshots can be sorted directly without regard to session boundaries. For tasks that depend on conversational structure, timestamp estimation and ordering are maintained per session.

Importantly, this imputation strategy is based on an engineering assumption that screenshot order roughly reflects message chronology. While this generally holds in user-submitted datasets, it may be invalid in legal contexts involving manipulation or reordering. Therefore, this method is positioned



Element Type	Example Question	Annotation Goal
Message Text	What is the intent of this message?	Identify expressions of intent (e.g., promise, request, warning)
	Who is the speaker?	Match speaker name for identity tracking
	What legal issue may be implied?	Perform legal inference (e.g., breach, infringement)
Avatar and Nickname	Is the avatar consistent with the nickname?	Verify identity coherence
Transaction Record	Who are the sender and recipient?	Determine transaction direction
	What is the legal nature of this transfer?	Classify as donation, payment, etc.
	What is the amount transferred?	Record monetary value (0 if unreadable)
Quoted Content	What is the speaker’s attitude toward the quote?	Distinguish affirmation, denial, or doubt
Emoji	Does the emoji express affirmation or negation?	Interpret sentiment or intention
	What emotion is conveyed?	Provide cultural interpretation (e.g., sarcasm)
Dialogue Name	Does the name reflect identity?	Link to legally relevant identity info
	What is the legal relationship between parties?	Infer from context (e.g., employer–employee)
Timestamp	What is the message time?	Support timeline reconstruction
Other Elements	Was this message recalled?	Judge evidentiary validity
	Is the speech-to-text reliable?	Assess transcript usability
	What file was sent?	Record name, type, and purpose

Table 6: Examples of local legal reasoning questions in the VQA task. This set is under annotation and not yet released or evaluated.

Setting	Value
Model	DETR (COCO-pretrained)
Data Augmentation	Brightness, Crop, Flip, Rotate
Anchor Generation	k-means clustering + elbow method
Optimizer	AdamW
Learning Rate	$5 \times 10^{-5}$
Batch Size	8
Epochs	50
LR Scheduler	Cosine Annealing
Early Stopping	Based on validation mAP
Metrics	mAP, IoU

Table 7: Training configuration for DETR-based screenshot layout detection.

as a heuristic for enhancing contextual coherence, not for evidentiary authentication or precise legal timeline reconstruction. Future work may incorporate device metadata or cross-image logical cues to improve legal robustness and applicability.

**Note on Scope.** While the timestamp imputation strategy described here is designed to support multi-image temporal modeling—especially for future tasks involving conversation reconstruction or inter-message reasoning—our current benchmark evaluation remains screenshot-level, with each VQA or KIE instance based on a single image input. This section primarily serves to document the semi-automatic annotation and reasoning methods applied during partial KIE labeling. It lays the groundwork for subsequent extensions of ScreenshotLegalBench toward multi-screenshot and temporally-

aware legal understanding benchmarks.

## B.4 JSON Schema

To support structured KIE from WeChat chat screenshots, we define a unified JSON output format that organizes each conversation into timestamped message entries with bounding box and semantic attributes. Figure 4 presents an example of the structured annotation used in KIE tasks.

## B.5 Annotation Interface UI for VQA Tasks

To facilitate structured annotation for the VQA tasks in ScreenshotLegalBench, we employed the LabelU platform to design a dual-level labeling interface. The annotation process includes both global-level and local-level legal reasoning questions.

Global questions focus on the legal attributes of the entire screenshot—e.g., whether it constitutes valid evidence or what type of legal dispute it may relate to. Local questions target specific elements within the screenshot, such as a particular message, emoji, or transaction, and aim to elicit fine-grained legal interpretations.

Figure 5 shows an example of a global question annotation scenario, where the screenshot is assessed for its potential relation to a partnership dispute. Figure 6 displays a local question focused on a transfer message, prompting the annotator to



Figure 4: Example Data Instance for Annotation

determine its legal nature.

## C Appendix : Datasets statistics

**Object Detection Subset.** This subset contains 50 manually annotated screenshots with a total of 945 bounding boxes across 15 interface element categories, used to train the layout detection models. As shown in Table 7, the majority of bounding boxes are concentrated in message and avatar regions, reflecting the visual dominance of conversation bubbles and speaker identity in chat interfaces.

**KIE Training Set.** The training set contains 39,477 message units extracted from screenshots via a semi-automatic pipeline. It includes 145,044 structured field annotations. As shown in Table 8, all samples have both speaker and message\_bbox, while 85.5% include content, and 59.3% contain timestamp information. Additionally, the dataset captures non-textual legal indicators such as transfer and image.

**KIE Evaluation Set.** This subset consists of 143 human-annotated screenshots comprising 696 message units and 2,678 structured fields. Table 8 summarizes the distribution. Most messages include speaker, message\_bbox, and content. Although rarer, legal fields such as transfer, image, and file are included due to their evidentiary value.

**VQA Subset.** The VQA set includes 1,176 chat screenshots annotated for multiple legal understanding tasks. As shown in Table 9, 38.9% are considered valid legal dialogs, and the same percentage were judged as evidential. A total of 502 samples include a textual case analysis written by legal professionals. Due to class imbalance, chat type (private vs. group) is only used as an auxiliary label.

## D Appendix : Evaluation Detail

We evaluate a diverse set of baseline approaches on ScreenshotLegalBench to establish performance benchmarks for both KIE and VQA tasks. In this section, we describe the experimental setup, baseline methods, and implementation details.

### D.1 Evaluation Results

Table 10 presents the performance of baseline and fine-tuned models on the KIE task, evaluated over the private chat subset of ScreenshotLegalBench. Scores reflect structured output quality in terms of

format validity, spatial alignment (IoU), and content accuracy. Fine-tuned InternVL2.5-2B achieves the highest overall score of 0.6921, demonstrating strong improvements across all dimensions.

Table 11 reports the performance of baseline models on two classification sub-tasks: (1) *Classify*, which determines whether a screenshot qualifies as legally meaningful chat evidence, and (2) *Evidence*, which assesses whether the screenshot conforms to a valid legal format. All results are based on a unified evaluation setting using non-structured prompts. Across both tasks, most models exhibit limited performance, with low F1 scores and high variance across metrics. Notably, Qwen2.5-VL-3B-Instruct(Bai et al., 2025) achieves relatively higher classification accuracy, while Qwen2.5-VL-72B(Bai et al., 2025) shows better recall for evidence detection, albeit with poor precision.

Table 12 summarizes the results for the third sub-task: *case\_text*, which requires generating a plausible legal cause of action from multi-image inputs. We evaluate models using a weighted composite score that aggregates hit rate (i.e., dispute match), semantic similarity, and length alignment. Among all tested models, Qwen2.5-VL-32B-Instruct outperforms others, followed by Qwen2.5-VL-3B-Instruct, indicating the benefit of larger model scales and instruction tuning. Nevertheless, overall scores remain modest, suggesting that multi-image legal reasoning remains a challenging task for current VLMs.

### D.2 Field Validation Rules for KIE Tasks

Each predicted message is considered structurally valid only if it contains the fields speaker, timestamp, content, and message\_bbox. The timestamp must include at least one digit and pass regex-based sanity checks. Bounding boxes must be well-formed 4-tuples with positive width and height. For model outputs in invalid JSON or partial structures, we apply a fallback parser with bracket completion and nested field recovery. Messages failing all checks are excluded from scoring.

### D.3 Finetune Configuration

We fine-tune the InternVL2(Chen et al., 2024b) and InternVL2.5(Chen et al., 2024a) models using the QLoRA approach; key training hyperparameters are listed in Table 13.



Figure 5: Annotation interface for global-level legal VQA tasks. This example shows a screenshot being annotated for its potential connection to a partnership dispute.

Field	Train	Eval	Struct.	Completeness	BBox	Time	Validity	Content	Metric
speaker	39477	696	✓	✓	✗	✗	✗	✓	TP / Total
message_bbox	39477	696	✓	✓	✓	✗	✗	✓	IoU
content	33753	681	✓	✓	✗	✗	✓	✓	$\lambda$ SeqSim + (1- $\lambda$ ) LCS
timestamp	23397	522	✓	✓	✗	✓	✓	✓	F1 (digit-check)
dialog_name	3208	61	✓	✓	✗	✗	✓	✓	$\lambda$ SeqSim + (1- $\lambda$ ) LCS
image	2661	10	✓	✓	✗	✗	✓	✓	$\lambda$ SeqSim + (1- $\lambda$ ) LCS
transfer	3071	10	✓	✓	✗	✗	✓	✓	$\lambda$ SeqSim + (1- $\lambda$ ) LCS
file	-	2	✓	✓	✗	✗	✓	✓	$\lambda$ SeqSim + (1- $\lambda$ ) LCS

Table 8: Summary of annotated fields across KIE dataset subsets and evaluation criteria.

#### D.4 Prompt Templates

The following prompt templates were used during evaluation. Figure 8 shows the full Chinese prompt used for zero-shot evaluation. The first line of the template (“Please extract structured information from this chat screenshot”) was also used as the fine-tuning instruction.

#### D.5 Ablation Results

To understand the impact of structural signals, we ablate the use of KIE-enhanced prompts in VQA (Table 14) and the effect of bounding-box inputs in KIE (Table 15).





Figure 6: Annotation interface for a local-level legal VQA task. The annotator is asked: “What is the legal nature of this transfer?”

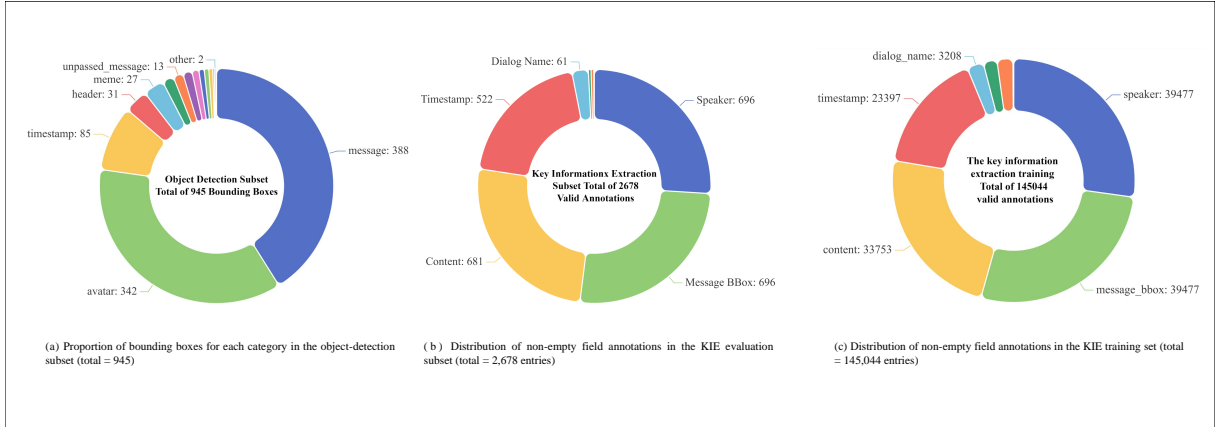


Figure 7: Annotation distribution statistics across the object detection and KIE subsets of ScreenshotLegalBench. Left: bounding box category ratio (N=945). Center: non-empty field counts in the KIE evaluation set (N=2,678). Right: non-empty field counts in the KIE training set (N=145,044).

Task Dimension	Label	Count	Percentage	Total
Dialog Type	legal_dialog	457	38.9%	1,176
	nonlegal_dialog	664	56.5%	
	not_dialog_but_legal	21	1.8%	
	not_dialog_and_nonlegal	34	2.9%	
Evidence Validity	is_evidence	457	38.9%	1,176
	not_evidence	719	61.1%	
Chat Type	private_chat	1,142	97.1%	1,176
	group_chat	34	2.9%	
Case Reasoning	with_case_text	502	42.7%	1,176
	without_case_text	674	57.3%	

Table 9: Annotation distribution in the VQA subset (total = 1,176).

Model	Format Score	IoU	Content Score	Overall Score	# Valid Samples
InternVL2-2B	0.9131	0.0009	0.3260	0.6195	143
InternVL2-2B (fine-tuned)	0.9517	0.0603	0.3909	0.6713	143
InternVL2.5-2B	0.9764	0.0006	0.2839	0.6302	143
InternVL2.5-2B-ft (run-13)	0.9644	0.0420	0.3944	0.6794	143
InternVL2.5-2B-ft (run-14)	0.9472	<b>0.1044</b>	<b>0.4369</b>	<b>0.6921</b>	143
Qwen2-VL-2B-instruct	0.3496	0.0006	0.0617	0.2057	143
Qwen2-VL-7B-instruct	0.5806	0.0000	0.2064	0.3935	31
Qwen2.5-VL-3B-instruct	0.9355	0.0012	0.2293	0.5824	31
Qwen2.5-VL-7B-instruct	0.0398	0.0009	0.0151	0.0274	143

Table 10: Evaluation results on the ScreenshotLegalBench KIE task (private chat subset). All scores are averaged over valid samples with parsing.

Model	Classify Task				Evidence Task			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
Qwen2.5-VL-32B	0.0000	0.0000	0.0000	0.0000	0.0382	0.0144	0.3333	0.0275
Qwen2.5-VL-72B	0.0357	0.0556	0.0238	0.0333	0.0328	0.0242	0.4444	0.0447
Qwen2-VL-7B-Instruct	0.0513	0.0635	0.0286	0.0394	0.0000	0.0000	0.0000	0.0000
Qwen2.5-VL-7B-Instruct	0.0357	0.0536	0.0204	0.0296	0.0244	0.0240	0.3606	0.0404
InternVL2-2B	0.0123	0.0167	0.0048	0.0074	0.0182	0.0007	0.0095	0.0014
InternVL2.5-2B	0.0602	0.0510	0.0340	0.0408	0.0113	0.0017	0.0069	0.0019
Qwen2-VL-2B-Instruct	0.0120	0.0159	0.0071	0.0099	0.0026	0.0003	0.0010	0.0004
Qwen2.5-VL-3B-Instruct	0.2143	0.0310	0.0857	0.0456	0.0023	0.0003	0.0016	0.0005

Table 11: Classification and Evidence Evaluation Results (w/o Structured Prompt)

Model	Hit Rate	Similarity	Length Score	Weighted Score
InternVL2-2B	0.0253	0.0233	0.1169	0.0430
InternVL2.5-2B	0.0253	0.0415	0.1278	0.0507
Qwen2-VL-2B-Instruct	0.0633	0.0200	0.0647	0.0506
Qwen2-VL-7B-Instruct	0.0443	0.0376	0.1088	0.0552
Qwen2.5-VL-3B-Instruct	0.1329	0.0333	0.1127	0.0990
Qwen2.5-VL-32B-Instruct	<b>0.2089</b>	0.0187	0.0537	<b>0.1208</b>
Qwen2.5-VL-72B-Instruct	0.1250	0.0301	0.0566	0.0829

Table 12: Cause-of-action generation performance (multi-image reasoning).

```

<image>\n请从这张聊天截图中提取结构化信息，
Format: ```
{
  \dialog_name\": \"<对话名称>\",
  \conversation\": [{
    \timestamp\": \"<第一条消息的时间戳>\",
    \speaker\": \"<如果是右侧发出则是avator_0，左侧发出是avator_1>\",
    \content\": \"<说话内容>\",
    \message_bbox\": {
      \min_x\": <边界框的min_x，是个数字，如917>,
      \max_x\": <边界框的max_x，是个数字>,
      \min_y\": <边界框的min_y，是个数字>,
      \max_y\": <边界框的max_y，是个数字>
    },
    \image\": \"<如果是图片或表情包，需要描述这个图片>\",
    \transfer\": \"<如果有转账信息，填写在这里，如微信转账请收款¥520.00>\",
    \file\": \"<如果有excel或者doc等文件，把文件名填写在这里，记得带上文件格式后缀，如2004年度数据分析表.xlsx>\",
  },
  ...
]
}
```

```

```

<image>\nPlease extract structured information from this
chat screenshot.
Format: ```
{
  "dialog_name": "<Name of the conversation>",
  "conversation": [{
    "timestamp": "<Timestamp of the first message>",
    "speaker": "<If the message is on the right side, use avator_0; if
on the left side, use avator_1>",
    "content": "<Text content of the message>",
    "message_bbox": {
      "min_x": <Minimum x-coordinate of the message bounding box
(e.g., 917)>,
      "max_x": <Maximum x-coordinate of the message bounding
box>,
      "min_y": <Minimum y-coordinate of the message bounding
box>,
      "max_y": <Maximum y-coordinate of the message bounding
box>
    },
    "image": "<If the message contains an image or emoji, provide
a description of the visual content>",
    "transfer": "<If the message contains a transaction, describe it
here, e.g., WeChat Transfer: Please confirm receipt ¥520.00>",
    "file": "<If the message includes a document (e.g., Excel or
Word), write the filename including its extension, e.g.,
2004年度数据分析表.xlsx>"
  },
  ...
]
}
```

```

Figure 8: Chinese prompt used during KIE pass@1 evaluation (left), with English translation shown on the right.

<b>Hyperparameter</b>	<b>Value</b>
Max sequence length	8192
Batch size (per GPU)	1
Gradient accumulation	2
Epochs	1 / 4
Optimizer	AdamW
LR scheduler	Cosine decay
Warmup ratio	3%
Initial learning rate	$5 \times 10^{-5} / 1 \times 10^{-4}$
LoRA rank	16
LoRA scaling factor	16
LoRA dropout	0.05
Gradient clipping	1.0
Layer-wise LR decay	0.75

Table 13: Fine-tuning hyperparameter configuration.



Table 14: Ablation: Effect of structured KIE input on VQA tasks.

Setting	Task	Accuracy	Macro P	Macro R	F1 Score
<i>w/o KIE guidance</i>	classify	0.4286	0.3333	0.1428	0.1999
<i>+KIE-augmented input</i>	classify	<b>0.8333</b>	<b>0.5000</b>	<b>0.4167</b>	<b>0.4545</b>
<i>w/o KIE guidance</i>	evidence	0.0157	0.0013	0.0107	0.0019
<i>+KIE-augmented input</i>	evidence	<b>0.0187</b>	<b>0.0126</b>	<b>0.0606</b>	<b>0.0204</b>

Table 15: Performance of InternVL2-2B on ScreenshotLegalBench with and without message\_box bounding boxes

Model	Format Score	Avg. IoU	Content Score	Overall Score
InternVL2-2B (w/ bbox)	0.9384	0.0338	0.3703	0.6544
InternVL2-2B (w/o bbox)	0.7579	0	0.3934	0.5756