

TRivia: Self-supervised Fine-tuning of Vision-Language Models for Table Recognition

Junyuan Zhang^{1*} Bin Wang^{2*} Qintong Zhang³ Fan Wu² Zichen Wen^{2,4} Jialin Lu¹
 Junjie Shan¹ Ziqi Zhao¹ Shuya Yang¹ Ziling Wang¹ Ziyang Miao² Huaping Zhong⁵
 Yuhang Zang² Xiaoyi Dong² Ka-Ho Chow^{1†} Conghui He^{2†}
¹The University of HongKong ²Shanghai AI Laboratory
³Peking University ⁴Shanghai Jiaotong University ⁵Sensetime

Abstract

Table recognition (TR) aims to transform table images into semi-structured representations such as HTML or Markdown. As a core component of document parsing, TR has long relied on supervised learning, with recent efforts dominated by fine-tuning vision-language models (VLMs) using labeled data. While VLMs have brought TR to the next level, pushing performance further demands large-scale labeled data that is costly to obtain. Consequently, although proprietary models have continuously pushed the performance boundary, open-source models, often trained with limited resources and, in practice, the only viable option for many due to privacy regulations, still lag far behind. To bridge this gap, we introduce TRivia, a self-supervised fine-tuning method that enables pretrained VLMs to learn TR directly from unlabeled table images in the wild. Built upon Group Relative Policy Optimization, TRivia automatically identifies unlabeled samples that most effectively facilitate learning and eliminates the need for human annotations through a question-answering-based reward mechanism. An attention-guided module generates diverse questions for each table image, and the ability to interpret the recognition results and answer them correctly provides feedback to optimize the TR model. This closed-loop process allows the TR model to autonomously learn to recognize, structure, and reason over tables without labeled data. Leveraging this pipeline, we present TRivia-3B, an open-sourced, compact, and state-of-the-art TR model that surpasses existing systems (e.g., Gemini 2.5 Pro, MinerU2.5) on three popular benchmarks. Model and code are released at: <https://github.com/opendatalab/TRivia>

* These authors contributed equally to this work.

† Corresponding Authors: Ka-Ho Chow, Conghui He.

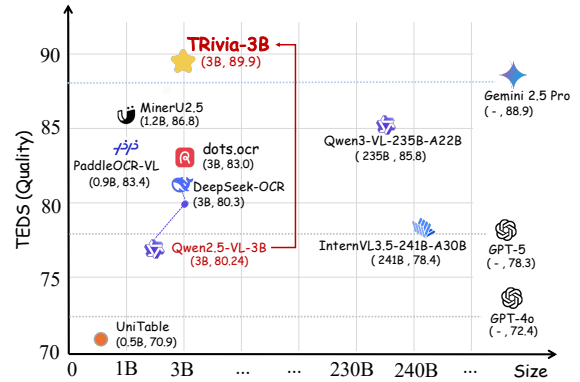


Figure 1. TRivia-3B learns from unlabeled table images and achieves TR quality beyond the limit attainable by fine-tuning with labeled data. Unlike proprietary systems such as Gemini 2.5 Pro, it is open-sourced, compact, and can be deployed offline for privacy-sensitive document processing.

1. Introduction

Document parsing plays a pivotal role in digitalization by converting scanned or photographed documents into machine-readable formats [33, 42] for downstream tasks such as retrieval-augmented generation [11, 41]. Among document elements, tables remain a longstanding challenge. They are both information-dense and structurally complex, requiring not only accurate text extraction, as in standard optical character recognition (OCR), but also precise reconstruction of their spatial and logical organization into semi-structured representations such as HTML or Markdown.

Recent advances in vision-language models (VLMs) [2, 4, 12, 34] have revolutionized table recognition (TR) [3, 15, 19, 23, 27]. General-purpose VLMs can be adapted to TR either through prompt engineering or fine-tuning with image-annotation pairs. Proprietary systems such as Gemini 2.5 Pro, built with massive human and computational resources, already exhibit unprecedented TR capabilities as

indicated in Figure 1. Yet, they are accessible only through commercial APIs, raising privacy and compliance concerns in many practical scenarios involving sensitive documents. In contrast, open-source VLMs offer greater flexibility and transparency but are constrained by the limited scale of existing TR datasets [10, 20], which are insufficient to reach state-of-the-art performance comparable to proprietary models. For example, UniTable, trained solely on open-source data [10, 20], performs poorly across real-world benchmarks. Even large-scale open-source efforts such as MinerU2.5, which combines millions of samples, costly human annotations, and distillation from proprietary models, still inherit the ceiling imposed by its teacher, Gemini 2.5 Pro.

Existing approaches to data acquisition for TR follow three paradigms. (1) Synthetic data, generated by rendering HTML tables, provides perfect annotations and scalability [7, 15, 19] but lacks the visual diversity and domain alignment of real-world data [19]. (2) Real-world data better capture the complexity of practical documents, but are expensive and time-consuming to annotate [3]. (3) Distilling pseudo-labels from proprietary models offers a potential compromise [27], yet this approach remains costly, inherently caps performance at that of the teacher model, and may violate service agreements. Consequently, the curation of large-scale labeled datasets has become a major bottleneck, even though it is a critical factor for success. This raises an important question: *Can unlabeled table images from the wild be harnessed to improve TR performance?*

In this paper, we present TRivia, a self-supervised fine-tuning framework that enables VLMs to learn table recognition directly from unlabeled table images, which can be collected at scale. As not all samples contribute equally to learning, identifying those that are most informative and deriving verifiable supervisory signals without human intervention are key challenges. TRivia leverages Group Relative Policy Optimization (GRPO) [31] with a decoupled reward function [28, 38] to strengthen TR capabilities via unlabeled data. Given a base VLM and a pool of unlabeled table images, it features a response-consistency sampling strategy that builds on GRPO’s principle of optimizing relative reward differences across multiple responses to select the most informative samples for learning. TRivia then generates supervisory signals through a proxy task: table question-answering (QA). An attention-guided QA generation module produces diverse, verifiable questions about each table image, and the VLM is fine-tuned to recognize, structure, and reason over tables by producing recognition results that allow a language model to correctly answer these questions.

The main contributions of this paper are as follows.

- We investigate self-supervised fine-tuning of VLMs on unlabeled table images and propose TRivia, a framework that harvests such data to push the frontier of TR.
- We introduce a response-consistency sampling strategy to

automatically identify unlabeled samples that most effectively enhance VLM fine-tuning via GRPO.

- We design an attention-guided QA mechanism to ensure diverse, verifiable, and stable supervisory signals for optimization.

Based on the proposed framework, we present TRivia-3B, an open-source, state-of-the-art TR model fine-tuned from Qwen2.5-VL-3B using unlabeled table images. As shown in Figure 1, TRivia-3B surpasses both expert TR models (e.g., MinerU2.5 and PaddleOCR-VL) and general-purpose models (e.g., Gemini 2.5 Pro and GPT-5). We envision that TRivia will open a new avenue for advancing TR beyond the limitations of labeled data via self-supervised fine-tuning.

2. Related Work

Early research on TR typically decomposes the task into two subtasks: table structure recognition and OCR. For structure recognition, bottom-up approaches [17, 29, 37] treat text or cell bounding boxes as graph nodes and predict their row, column, or cell associations. Split-and-merge methods [16, 32, 43, 44] first partition a table into a uniform grid and then merge adjacent cells to reconstruct its structure. While effective on clean layouts, these modular pipelines suffer from cumulative error propagation and rely heavily on explicit visual cues, making them brittle under real-world conditions such as distorted layouts, borderless designs, or complex merged cells [20]. These limitations have motivated the development of image-to-markup methods that reformulate TR as an end-to-end sequence-generation problem, directly converting table images into formats such as HTML [10, 22, 26, 48]. However, these models are constrained by their context window, input resolution, and training data quality, limiting their ability to handle large or structurally complex tables. For instance, UniTable [26] processes images up to 448×448 with a 512-token output cap, resulting in notable performance drops on complex real-world tables.

Recent advances in VLMs have demonstrated remarkable generalization across OCR-related tasks [8, 18, 25, 40]. By incorporating table recognition data during pretraining, general-purpose VLMs such as Gemini 2.5 Pro [4] and the Qwen2.5-VL series [2] can perform TR directly through natural-language instructions, eliminating the need for explicit OCR or cell-detection modules. Building on this trend, several task-specific expert VLMs [7, 15, 19, 23, 30, 35] have been developed explicitly for TR, typically fine-tuned on large synthetic datasets augmented with a small number of real tables annotated either manually [3, 23] or by proprietary systems [23, 27]. Nevertheless, even large-scale open-source efforts such as MinerU2.5 [23], which combines millions of samples, costly human annotations, and distilled labels from Gemini 2.5 Pro, remain constrained by the performance ceiling of its teacher model, as shown in

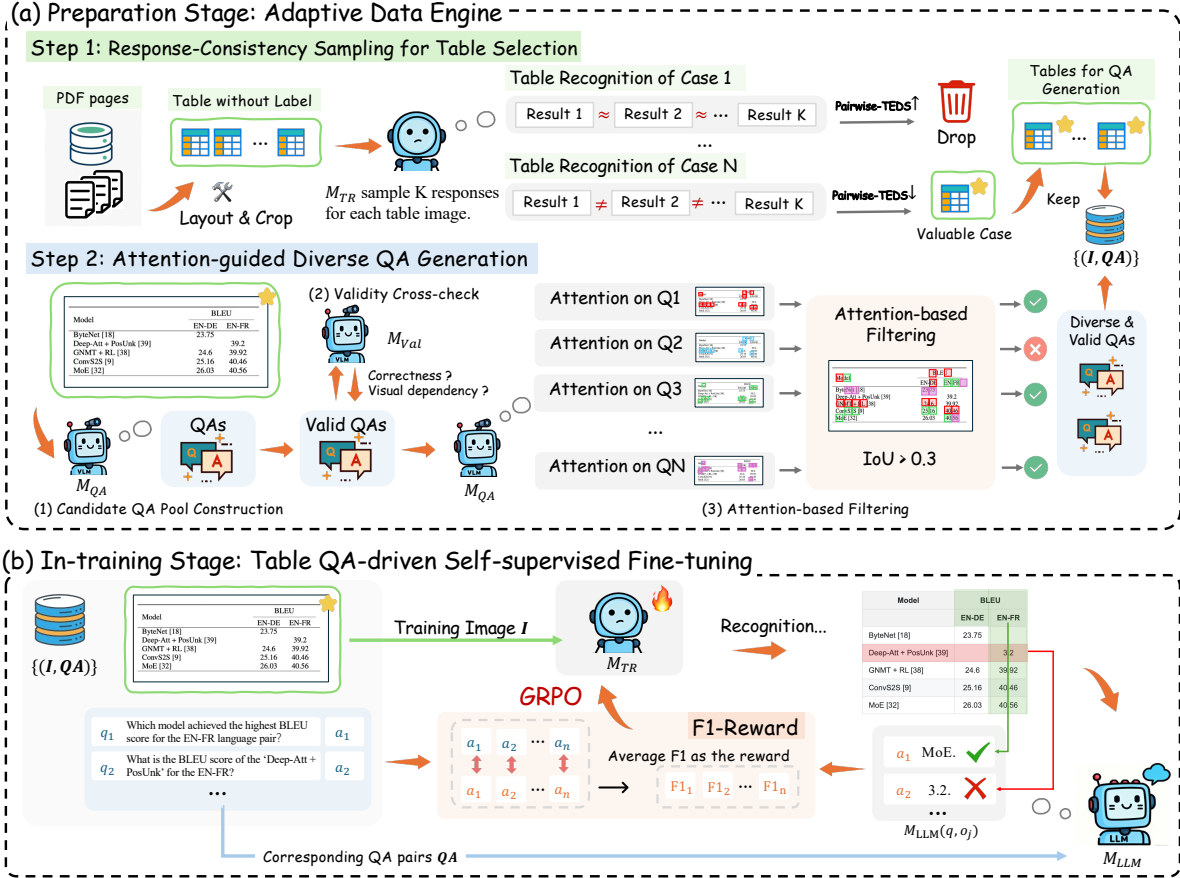


Figure 2. TRivia features (a) an adaptive dataset curation module to set the stage for (b) reinforcement learning to learn TR from unlabeled data. During dataset curation, TRivia uses a response-consistency sampling strategy (Section 3.2.1) to identify informative samples and generate verifiable, diverse QA for each image through an attention-guided module (Section 3.2.2). Based on curated data, TRivia fine-tunes the VLM to recognize, structure, and reason over tables through QA-based rewards (Section 3.1).

Figure 1. In contrast, our work departs from reliance on labeled or distilled data and uses table QA-based rewards as proxy supervision with GRPO to learn directly from unlabeled real-world table images.

3. Methodology

Figure 2 provides an overview of TRivia, which consists of two stages: (i) a preparation stage that prepares supervisory signals from unlabeled table images, and (ii) an in-training stage that formulates TR learning as a semi-supervised problem. We first describe how TRivia fine-tunes the VLM via GRPO using table QA as proxy rewards in Section 3.1. We then introduce the adaptive data engine in Section 3.2, which identifies informative samples and generates supervisory signals tailored to our fine-tuning strategy and the base VLM.

3.1. Table QA-driven Self-supervised Fine-tuning

To push TR performance beyond the limits of supervised fine-tuning on labeled data, we adopt reinforcement learning

(RL), which updates the model based on a reward function rather than explicit labels. As long as a reliable reward can be derived from unlabeled data, the model can still benefit from meaningful supervision and learn to claim the reward.

Table question-answering (QA) can serve as a proxy task to provide annotation-free reward signals because, as a downstream application of TR, the ability to correctly answer questions about a table implicitly reflects how well the recognized table preserves both textual and structural information [41]. This is beneficial for two reasons. First, generating valid QA pairs from a table image is significantly easier than predicting its full HTML markup: the model does not need to explicitly infer complex structures such as colspan and rowspan, but only to reason about the content of specific regions. Second, the quality of QA pairs, including correctness and visual dependency, can also be cross-checked using additional models, which is not yet feasible when using HTML as labels. These properties make table QA a reliable reward for self-supervised fine-tuning.

Table QA-driven GRPO. TRivia adopts Group Relative Policy Optimization (GRPO) as the reinforcement learning framework to fine-tune the base VLM, such as the one optimized using existing labeled data, and to be further enhanced by learning from unlabeled data. The overall GRPO training pipeline with QA-based rewards is illustrated in Figure 2(b). Formally, each training sample is represented as $(\mathbf{I}, \mathbf{QA})$, where \mathbf{I} denotes a table image and $\mathbf{QA} = \{(q, a)\}$ represents its associated QA set. The process for generating QAs from \mathbf{I} is detailed in Section 3.2. During GRPO training, the fine-tuning-in-progress TR model M_{TR} acts as the policy and produces a group of R recognition responses $\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_R$ for each image \mathbf{I} . Each recognition output \mathbf{o}_j is then paired with table question q 's and passed to a large language model (LLM) M_{LLM} , which generates the corresponding answer $M_{\text{LLM}}(q; \mathbf{o}_j)$ based on the recognized table. The reward for each QA pair is computed as the F1-score between the predicted and ground-truth answers, and the overall reward for \mathbf{o}_j is the average across all QA pairs for image \mathbf{I} :

$$\text{Reward}(\mathbf{o}_j) = \frac{1}{|\mathbf{QA}|} \sum_{(q, a) \in \mathbf{QA}} \text{F1}(M_{\text{LLM}}(q; \mathbf{o}_j), a). \quad (1)$$

Filtering-based Stabilization. Despite its effectiveness, QA-based reward estimation can fail when the TR model generates illegal or repetitive outputs that do not form valid table structures. Since such outputs cannot yield meaningful answers, they are assigned a reward of zero. However, naively including these samples introduces reward noise, as they inflate the relative advantage of valid responses and destabilize training. To address this issue, we apply illegal-sample filtering, which discards invalid recognition results with zero rewards during GRPO training. As shown in Section 5.2, this strategy significantly stabilizes QA-based rewards, resulting in smoother and more reliable GRPO.

3.2. Adaptive Data Engine

Building upon the above formulation, we next present how TRivia adaptively constructs supervisory signals tailored for the VLM to be fine-tuned. This includes selecting unlabeled samples that yield the most significant optimization signals (Section 3.2.1) and generating QA pairs that provide richer feedback (Section 3.2.2).

3.2.1. Response-Consistency Sampling

Learning from unlabeled data allows one to scale up the training set to cover real-world tables with diverse layouts. However, not all samples contribute equally to learning, and identifying which ones offer the greatest training value for the VLM is non-trivial. For example, clustering-based methods [14] can measure sample diversity, but they cannot determine which samples are most beneficial for improving the model to be fine-tuned. Alternatively, human-in-the-

loop approaches [3] can effectively mine hard cases but are subjective and impractical for large-scale datasets.

Since GRPO benefits from samples that elicit diverse model responses, table images that cause the TR model to produce varied recognition outputs are particularly valuable. To capture this, we sample multiple recognition results for each table image and measure their pairwise structural similarity using the Tree Edit Distance-based Similarity (TEDS) metric [48]. Formally, for each image \mathbf{I} , the TR model M_{TR} generates a set of K individual responses $\{\mathbf{o}_i\}_{i=1}^K$. The consistency score for image \mathbf{I} is then defined based on the aggregated pairwise TEDS similarity among these responses:

$$\text{Consistency}(\mathbf{I}) = \frac{2}{K^2 - K} \sum_{1 \leq i < j \leq K} \text{TEDS}(\mathbf{o}_i, \mathbf{o}_j). \quad (2)$$

A lower consistency score indicates higher response diversity and greater model uncertainty, making the sample more valuable for GRPO training.

In principle, response-consistency sampling can be performed online by re-evaluating samples throughout training. However, this approach is computationally prohibitive. We found that conducting this sampling offline already provides stable and effective improvements, as shown in Section 5.

3.2.2. Attention-based Diverse QA Generation

QA pairs are the driving force of TRivia to learn from unlabeled table images. They must satisfy two key requirements: (i) the QA set for each image should cover different table regions to enrich learning signals, and (ii) their correctness must be verifiable. A straightforward approach is to prompt a VLM to generate QA pairs. However, as shown in Figure 3, single-pass generation often covers only part of the table, while multi-pass generation tends to introduce paraphrased QA pairs that are sourced from overlapping areas. Consequently, even if the TR model produces diverse recognition outputs, the QA-based rewards may remain limited, hindering the effectiveness of GRPO training.

To address this, we leverage the attention mechanism of VLMs, which inherently encodes the visual grounding of textual tokens during answer generation [1, 13]. Each answer token attends to visual tokens that provide its evidence, enabling TRivia to reason about the visual source of each QA pair. Building on this insight, we design an attention-based QA generation method that explicitly utilizes these attention patterns to identify and filter out QA pairs with limited contributions.

For a given QA pair (q, a) generated from image \mathbf{I} by the QA generation model M_{QA} , we define its visual source as the set of visual tokens with significant attention weights. Specifically, if the table image is tokenized into visual tokens \mathcal{V} , the visual source (VS) of the QA pair is defined as:

$$\text{VS}((q, a); \mathbf{I}, M_{\text{QA}}) = \{v \mid \mathcal{A}_{M_{\text{QA}}}(v \mid a) > \tau_{\mathcal{A}}, v \in \mathcal{V}\} \quad (3)$$

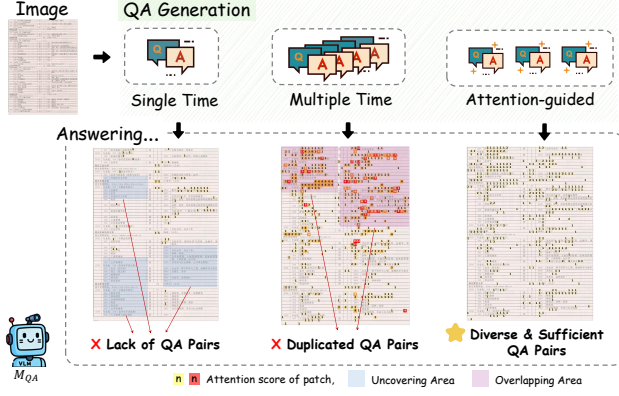


Figure 3. Single-time QA generation captures limited table content, while multiple samplings introduce redundant or overlapping QA pairs. The proposed attention-guided QA generation leverages attention distributions to diversify question sources, producing concise and comprehensive QA pairs.

where $\mathcal{A}_{M_{QA}}(v \mid a)$ denotes the attention score to the visual token v averaged across answer tokens, and τ_A is a fixed threshold. The resulting token list represents the visual grounding of the QA pair, allowing TRivia to eliminate pairs that rely on overlapping evidence.

Based on the visual source above, the attention-guided QA generation process involves three steps (Figure 2(a)):

1. **Candidate QA Pool Construction:** Given an image I , the QA generation (teacher) model M_{QA} is prompted multiple times to produce a pool of candidate QA pairs: $\text{GenQA}(I) = \{(q, a)\}$.
2. **Validity Cross-checking:** Each candidate pair is cross-checked by an external VLM M_{Val} to ensure visual dependency and correctness. Specifically, we retain only those pairs that can be correctly answered with the image but not without it: $\text{ValidQA}(I) = \{(q, a) \mid M_{Val}(q; I) = a \wedge M_{Val}(q; \emptyset) \neq a, (q, a) \in \text{GenQA}(I)\}$.
3. **Attention-Guided QA Selection:** We greedily select valid pairs whose visual sources exhibit minimal overlap to ensure broad table coverage. To account for attention sink effects [13], we relax disjointness by an Intersection-over-Union (IOU) threshold between any two QA pairs: $\text{IOU}(\text{VS}((q_i, a_i); I, M_{QA}), \text{VS}((q_j, a_j); I, M_{QA})) < \tau_{\text{IOU}}$.

This attention-guided mechanism ensures that QA pairs are both visually grounded and semantically diverse, providing richer and more robust supervision for GRPO training.

4. TRivia-3B

Building upon the above framework, we propose TRivia-3B, a model obtained by fine-tuning Qwen2.5-VL-3B-Instruct to surpass the TR performance achievable with existing labeled datasets and establish a new state of the art. The overall

training pipeline consists of three stages. We first leverage labeled datasets to gradually adapt a general-purpose VLM into a strong TR model, establishing the performance limit achievable through supervised learning. We then apply TRivia to further fine-tune the model in a self-supervised manner using unlabeled data. Additional implementation details are provided in the appendix. Together with the source code, TRivia-3B will be publicly released upon acceptance (now in supplementary materials).

Stage 1: OTSL Tag Warm-up. Compared to HTML or Markdown representations, OTSL [21] provides a compact and well-structured table format. It encodes adjacency relationships instead of explicitly predicting colspan and rowspan attributes for merged cells, substantially reducing token length and simplifying structural prediction. In this first stage, we train the model on large-scale open-source datasets to familiarize it with the OTSL syntax. We aggregate data from PubTabNet [47], SynthTabNet [22], and MMTab [46]. Specifically, we sample 100K instances from each of the four SynthTabNet subsets, 200K from PubTabNet, and 100K from MMTab. After removing erroneous samples (e.g., with unclosed HTML tags) and converting all annotations into OTSL, we obtain approximately 700K training instances. During this warm-up, only the language model component of Qwen2.5-VL-3B is fine-tuned, while the visual encoder and alignment module remain frozen. This allows the model to learn OTSL syntax and structure without disrupting low-level visual representations.

Stage 2: Supervised Fine-Tuning. In the second stage, we improve robustness and generalization by fine-tuning the model on real-world table images collected from open-source datasets [20, 39, 44] and web resources [9]. We curate approximately 50K samples and fine-tune the model with all parameters unfrozen. This stage bridges the gap between synthetic and real-world layouts, reaching the performance limit achievable through supervised learning.

Stage 3: TRivia. In the final stage, we apply TRivia to further optimize the model through GRPO with QA-based rewards. Following the data curation procedure described in Section 3.2, we construct the RL training dataset as follows. We collect PDFs from web sources, use DocLayout-YOLO [45] for layout detection, and crop tables to form an unlabeled table image pool of approximately 100K samples. Using the Stage-2 TR model, we perform response-consistency sampling to identify images that yield diverse recognition outputs. We observe that samples with consistency scores below 0.4 are often noisy or non-tabular. Therefore, we uniformly sample images with scores in the range of 0.4–1.0, at an interval of 0.1 to ensure diversity across uncertainty levels, retaining about 50K informative samples for RL training. For each selected image, we generate QA pairs using Qwen2.5-VL-72B-Instruct as M_{QA} and empirically select layer 72 for attention computation

with attention thresholds $\tau_A = 0.01$ and $\tau_{\text{IOU}} = 0.3$ (Section 3.2.2). After validity cross-checking by InternVL3-78B, we retain roughly 30 diverse QA pairs per image, resulting in a balanced and information-rich QA dataset. During GRPO training (Section 3.1), Qwen3-8B acts as M_{LLM} , responsible for answering the generated QA pairs based on the recognized table outputs.

5. Experimental Evaluation

We conduct empirical studies to evaluate and analyze the effectiveness of the proposed TRivia framework and the resulting TRivia-3B model.

TR Benchmarks. We evaluate TRivia-3B across three widely adopted benchmarks: OmniDocBench [25], CC-OCR [40], and OCRBench v2 [8]. OmniDocBench consists of digital PDFs with a wide range of table types. We use the latest version, OmniDocBench v1.5, which includes more complex tables and crops all table images based on the provided layout annotation, resulting in 512 samples for evaluation. CC-OCR consists of 300 scanned and photographed table images, including challenging samples such as long tables, complex structures, and handwritten samples. OCRBench v2 is the latest version of OCRBench. We employ its table parsing subset, which includes 700 table images, featuring diverse real-world table images with a variety of layouts and visual complexities.

Comparison Schemes. We compare TRivia-3B with three categories of models: (1) Expert TR models: SLANNet-plus [6] and UniTable [26]. (2) General-purpose VLMs: InternVL3.5-241B-A30B [34], Qwen2.5-VL-72B, Qwen3-VL-235B-A22B [2], Gemini 2.5 Pro [4], GPT-4o [12], and GPT-5 [24]. These models are adopted for TR using the prompt provided in the appendix. (3) Document-parsing VLMs: dots.ocr [30], DeepSeek-OCR [36], PaddleOCR-VL [5], and MinerU2.5 [23], each evaluated with their default prompts.

Outline. We first examine how TRivia-3B advances a wide range of state-of-the-art models across different benchmarks in Section 5.1. Then, we analyze the enablers of such advancements via ablation studies in Section 5.2. Finally, we discuss the broader utility for TRivia as a scalable data annotator in Section 5.3.

5.1. Advancing TR Performance

Table 1 summarizes the TR performance of TRivia-3B and 12 baselines over four datasets. We report TEDS [10] and S-TEDS (structure-only TEDS) to assess both holistic and structural accuracy.

Expert TR Models. TRivia-3B achieves substantially better TR on diverse scenarios, except for PubTabNet in S-TEDS, due to dataset-specific overfitting of expert models. Specifically, TRivia-3B surpasses UniTable by 27.06 and 23.03

TEDS on CC-OCR and OCRBench v2, respectively, confirming its robustness to complex layouts.

General-purpose VLMs. Despite being trained with significantly fewer parameters and data, TRivia-3B still outperforms all these large-scale general-purpose VLMs on most benchmarks. In particular, Qwen2.5-VL-72B has been used as the QA generation (teacher) model (i.e., M_{QA} in Section 3.2.2) to generate table QAs for fine-tuning TRivia-3B. We observe a 6.36 TEDS improvement overall, indicating that TRivia successfully extracts and refines useful supervision beyond direct distillation. Another noteworthy baseline is Gemini 2.5 Pro, the proprietary model built with massive human and computational resources. It has incorporated table recognition data during pretraining and hence achieved the most competitive performance (88.93 TEDS and 91.23 S-TEDS). Still, TRivia-3B consistently outperforms it in most benchmarks, except for CC-OCR (85.56 TEDS vs 84.90 TEDS).

Document-parsing VLMs. TRivia-3B exhibits better performance than existing models tailored for document parsing in all benchmarks. Although its performance is marginally better than the earlier SOTA model, PaddleOCR-VL, on OmniDocBench, which primarily consists of digital table images, TRivia-3B outperforms on the CC-OCR and OCRBench benchmarks by a large margin (5.28 and 11.47 TEDS). The advantages can also be observed when compared with MinerU2.5, which has been trained with data several orders of magnitude more than TRivia-3B, with manual annotation and distillation from Gemini 2.5 Pro. TRivia-3B can still offer better TR in all cases. These results validate that TRivia not only enhances performance but also improves robustness to diverse real-world distributions.

5.2. Ablation Studies

TRivia-3B exhibits consistent improvements over existing methods across different benchmarks. Next, we analyze the enablers of reaching such a TR capability.

Table QA-based Supervisory Signals. Our table QA-based proxy reward allows TRivia to learn from a VLM with imperfect annotations. To understand this improvement, we use the QA generation model, Qwen2.5-VL-72B, to (i) produce HTML-based table recognition on training images, (ii) convert them into OTSL format, and (iii) fine-tune the Stage-2 VLM with these pseudo-labels using either supervised fine-tuning (SFT) or GRPO. Qwen2.5-VL-72B cannot produce high-quality recognition results, as shown in Figure 4, depicting a training sample that the model fails to split the cells and generate the table correctly. Based on these imperfect labels, our results reported in Table 2 reveal that the fine-tuned models have significantly worsened TR performance. SFT leads to an average decrease of 8.37 TEDS, even underperforming Stage-1 on OmniDocBench, while GRPO mitigates the decline slightly but still lags by 4.92 TEDS. We attribute

| | PubTabNet | | OmniDocBench | | CC-OCR | | OCRBench | | Overall | |
|-----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | TEDS | S-TEDS | TEDS | S-TEDS | TEDS | S-TEDS | TEDS | S-TEDS | TEDS | S-TEDS |
| Expert TR models | | | | | | | | | | |
| SLANNet-plus | 86.57 | 96.43 | 81.90 | 89.08 | 50.93 | 65.84 | 65.55 | 77.73 | 68.19 | 79.21 |
| UniTable | 86.44 | <u>95.66</u> | 82.76 | 89.82 | 57.84 | 70.47 | 67.73 | 78.65 | 70.86 | 80.81 |
| General-purpose VLMs | | | | | | | | | | |
| InternVL3.5-241B-A30B | 83.75 | 88.76 | 86.03 | 90.53 | 62.87 | 69.52 | 79.50 | 85.81 | 78.41 | 84.18 |
| Qwen2.5-VL-72B | 84.39 | 87.91 | 87.85 | 91.80 | 81.22 | 86.48 | 81.33 | 86.58 | 83.52 | 88.33 |
| Qwen3-VL-235B-A22B | - | - | 91.02 | <u>94.97</u> | 80.98 | 86.19 | 84.12 | 88.15 | 85.83 | 90.07 |
| Gemini 2.5 Pro | - | - | 90.90 | 94.32 | 85.56 | <u>90.07</u> | 88.94 | 89.47 | <u>88.93</u> | <u>91.23</u> |
| GPT-4o | 76.53 | 86.16 | 78.27 | 84.56 | 66.98 | 79.04 | 70.51 | 79.55 | 72.44 | 81.15 |
| GPT-5 | - | - | 84.91 | 89.91 | 63.25 | 74.09 | 79.91 | 88.69 | 78.30 | 86.21 |
| Document-parsing VLMs | | | | | | | | | | |
| dots.ocr | 90.65 | 93.76 | 88.62 | 92.86 | 75.42 | 81.65 | 82.04 | 86.27 | 82.95 | 87.58 |
| DeepSeek-OCR | - | - | 83.79 | 87.86 | 68.95 | 75.22 | 82.64 | 87.33 | 80.31 | 85.11 |
| PaddleOCR-VL | - | - | <u>91.12</u> | 94.62 | 79.62 | 85.04 | 79.29 | 83.93 | 83.36 | 87.77 |
| MinerU2.5 | 89.07 | 93.11 | 90.85 | 94.68 | 79.76 | 85.16 | <u>87.13</u> | <u>90.62</u> | 86.82 | 90.81 |
| TRivia-3B | 91.79 | 93.81 | 91.60 | 95.01 | <u>84.90</u> | 90.17 | 90.76 | 94.03 | 89.88 | 93.60 |

Table 1. TRivia-3B achieves consistently high TR performance, in TEDS and S-TEDS, across four benchmarks, which is unattainable by 12 existing methods that include expert TR models, general-purpose VLMs, and those fine-tuned for document parsing.

| Dataset | ECLA method | eGQA binding | | | | | gQA binding | | | | |
|----------------------------------|-------------|--------------|-----------------|------------------|-------------------|--------------------|-------------|-----------------|------------------|-------------------|--------------------|
| | | R_{acc} | R_{acc}^{100} | R_{acc}^{1000} | R_{acc}^{10000} | R_{acc}^{100000} | R_{acc} | R_{acc}^{100} | R_{acc}^{1000} | R_{acc}^{10000} | R_{acc}^{100000} |
| NCDN-3 | NPVQA | 1.20E-04 | 38.74E-04 | 246.40 | 3.76 | 1.12E-04 | 35.90E-04 | 330.12 | 3.22 | | |
| | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| NCDN-4 | NPVQA | 4.16E-04 | 1.00E-04 | 2.46 | 112.87 | 4.00E-04 | 1.00E-04 | 1.00 | 96.19 | | |
| | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| NCDN-4 ^{Qwen2.5-VL-72B} | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| NCDN-4 ^{Gemini 2.5 Pro} | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| NCDN-1 | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| NCDN-1 | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |
| | NPVQA | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | 1.00E-04 | 1.00E-04 | 1.00 | 1.00 | | |

Input table image

Teacher M_{QA} : Qwen2.5-VL-72B

✗ Incorrectly split the cell and generate incomplete table.

Base TR model

✗ Incorrectly split the cell, as it rely solely on visual cues.

TRivia-3B

★ TRivia-3B can correctly reason that they are in the same cell.

Figure 4. The teacher model M_{QA} used for QA generation could not directly generate correct annotations, but is sufficient to create QAs with TRivia to fine-tune the base model and lead to TRivia-3B that can handle complex structures.

this degradation to biased or semantically inconsistent annotations produced by Qwen2.5-VL-72B in the target domain. In contrast, the TRivia framework produces high-quality QA pairs to avoid such bias, producing consistent supervision signals and yielding better results.

| | OmniDoc Bench | CC OCR | OCR Bench | Overall |
|----------------|----------------|-----------------|----------------|----------------|
| Stage-1 | 87.65 | 72.28 | 73.06 | 77.85 |
| Stage-2 | 90.08 | 82.48 | 90.08 | 88.57 |
| Qwen2.5-VL-72B | 84.41 | 70.54 | 80.87 | 80.02 |
| + &SFT | (-5.67) | (-11.94) | (-9.21) | (-8.53) |
| Qwen2.5-VL-72B | 86.19 | 78.12 | 84.16 | 83.65 |
| + &GRPO | (-3.89) | (-4.36) | (-5.92) | (-4.92) |
| TRivia-3B | 91.60 | 84.90 | 90.76 | 89.88 |

Table 2. Using the QA generation model (Qwen2.5-VL-72B) to generate pseudo-labels for SFT or GRPO could not achieve high TR performance, measured in TEDS. The table QA-based approach in TRivia can mitigate the impact caused by imperfect annotations.

Attention-guided QA Generation. Table QA-based rewards are naturally sparse, as each QA pair assesses only a limited table region. Our proposed attention-guided QA generation enriches supervision diversity by widening question sources according to attention distributions. To show the advantage of such a design, we curate a test set of 80 table samples of varying levels of complexity and create annotations manually for evaluation purposes. As shown in Figure 5 (orange line), removing attention-guided QA generation leads to a significant drop in TR performance. Our investigation found that the fine-tuned model is particularly weak in handling structurally complex or visually ambiguous tables.

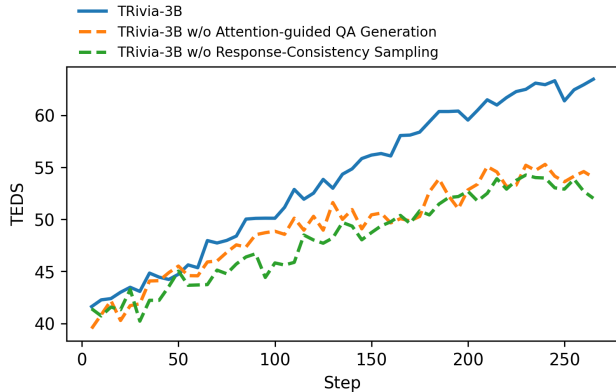


Figure 5. TRivia-3B benefits significantly from the diverse QAs generated by the attention-guided mechanism and the training samples that yield diverse outputs.

Response-Consistency Sampling. To prioritize informative samples, TRivia adopts response-consistency sampling, which favors table images yielding diverse model outputs. This design aligns with the GRPO principle of optimizing relative advantages within a response group. Compared to random sampling, response-consistency sampling accelerates convergence and boosts TEDS from 52.0 to 63.5, as shown in Figure 5 (green line). Although the sampling process occurs offline, it consistently enhances optimization efficiency, confirming that response-consistency sampling effectively identifies informative and challenging examples for optimization.

Illegal-Sample Filtering. Table QA-based rewards depend on valid table recognition outputs. When the model generates invalid responses, the LLM cannot produce correct answers, forcing these samples to receive zero reward. This artificially compresses the reward distribution and destabilizes GRPO. We mitigate this issue through illegal-sample filtering, which removes invalid responses before advantage computation. As shown in Figure 6, without filtering (orange line), the learning progressively becomes unstable and eventually leads to significant fluctuation in QA rewards. Instead, illegal-sample filtering (blue line) successfully stabilizes the learning. As a result, this approach reduces the convergence step by approximately 25% and improves final performance on test set by 3 TEDS.

5.3. TRivia as a Data Annotator

Beyond achieving state-of-the-art performance in table recognition, another crucial requirement is to transfer this capability to other models, such as distilling knowledge into smaller, more efficient architectures or enhancing the OCR abilities of general-purpose VLMs. To this end, we demonstrate the broader potential of TRivia as an automated data annotation system capable of generating reliable pseudo-

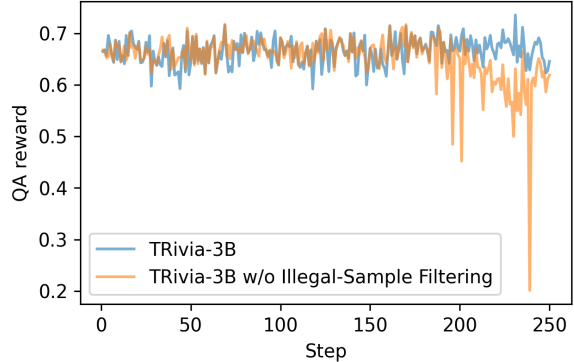


Figure 6. Illegal-sample filtering is crucial for stabilizing the fine-tuning of the TR model, as it suppresses reward noise caused by invalid responses.

| | OmniDoc Bench | CC OCR | OCR Bench | Overall |
|------------------------------|------------------|--------------|--------------|--------------|
| Stage-2 | 90.08 | 82.48 | 90.08 | 88.57 |
| + SFT w/ TRivia-3B labels | 91.37 | 85.84 | 90.75 | 89.99 |
| TRivia-3B | 91.60 | 84.90 | 90.76 | 89.88 |

Table 3. The distilled model (SFT w/ TRivia-3B labels) achieves comparable performance to TRivia-3B, confirming the reliability of TRivia-based annotations.

labels for unlabeled data. Specifically, we use TRivia-3B to create pseudo-labels for a set of unlabeled table images it did not see during its fine-tuning process. Then, we employ standard SFT to optimize a Stage-2 model. As shown in Table 3, the distilled model (SFT w/ TRivia-3B labels) achieves nearly identical performance to TRivia-3B across all benchmarks, confirming the high fidelity of generated annotations. Moreover, on the challenging CC-OCR benchmark, which is characterized by visually complex and layout-diverse tables, the distilled model even slightly outperforms TRivia-3B. Notably, this outcome contrasts with the poor performance obtained when distilling from Qwen2.5-VL-72B (see Table 2), which fails to yield accurate HTML annotations for the same data. In summary, TRivia establishes a new paradigm for table data annotation, enabling dynamic model adaptation to dataset characteristics rather than relying on a fixed annotator. These findings further highlight the potential of TRivia as a scalable, fully automated alternative to manual labeling or costly proprietary model distillation for constructing high-quality TR datasets at scale.

6. Conclusions

In this paper, we introduce TRivia, a self-supervised fine-tuning framework that enables VLMs to learn table recogni-

tion directly from unlabeled table images. Built upon TRivia, TRivia-3B establishes a new state-of-the-art among TR models, outperforming both specialized and proprietary models across multiple benchmarks. Beyond improving recognition performance, TRivia also serves as a scalable annotation engine, generating pseudo-labels that rival human and proprietary annotations. This work opens promising directions for self-supervised document parsing. Its annotation-free design provides a foundation for extending similar principles to broader multimodal tasks for learning at scale.

7. Acknowledgments

Authors from HKU are partially supported by the Croucher Start-up Allowance (Project #2499102828) and RGC Early Career Scheme (Project #27211524). Any opinions, findings, or conclusions expressed in this material are those of the authors and do not necessarily reflect the views of the Croucher Foundation and RGC.

References

- [1] Ingeol Baek, Hwan Chang, Sunghyun Ryu, and Hwanhee Lee. How do large vision-language models see text in image? unveiling the distinctive role of ocr heads. *arXiv preprint arXiv:2505.15865*, 2025. 4
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 2, 6
- [3] Xiangyang Chen, Shuzhao Li, Xiuwen Zhu, Yongfan Chen, Fan Yang, Cheng Fang, Lin Qu, Xiaoxiao Xu, Hu Wei, and Minggang Wu. Logics-parsing technical report. *arXiv preprint arXiv:2509.19760*, 2025. 1, 2, 4
- [4] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Naveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 1, 2, 6
- [5] Cheng Cui, Ting Sun, Suyin Liang, Tingquan Gao, Zelun Zhang, Jiaxuan Liu, Xueqing Wang, Changda Zhou, Hongen Liu, Manhui Lin, et al. Paddleocr-vl: Boosting multilingual document parsing via a 0.9 b ultra-compact vision-language model. *arXiv preprint arXiv:2510.14528*, 2025. 6
- [6] Cheng Cui, Ting Sun, Manhui Lin, Tingquan Gao, Yubo Zhang, Jiaxuan Liu, Xueqing Wang, Zelun Zhang, Changda Zhou, Hongen Liu, et al. Paddleocr 3.0 technical report. *arXiv preprint arXiv:2507.05595*, 2025. 6
- [7] Hao Feng, Shu Wei, Xiang Fei, Wei Shi, Yingdong Han, Lei Liao, Jinghui Lu, Binghong Wu, Qi Liu, Chunhui Lin, et al. Dolphin: Document image parsing via heterogeneous anchor prompting. *arXiv preprint arXiv:2505.14059*, 2025. 2
- [8] Ling Fu, Zhebin Kuang, Jiajun Song, Mingxin Huang, Biao Yang, Yuzhe Li, Linghao Zhu, Qidi Luo, Xinyu Wang, Hao Lu, et al. Ocrbench v2: An improved benchmark for evaluating large multimodal models on visual text localization and reasoning. *arXiv preprint arXiv:2501.00321*, 2024. 2, 6
- [9] Heywhale. Tal ocr table. url <https://www.heywhale.com/home/competition/606d6ff0e04ac0017c3bf7f/con> 2025. Accessed: 2025-10-26. 5
- [10] Yongshuai Huang, Ning Lu, Dapeng Chen, Yibo Li, Zecheng Xie, Shenggao Zhu, Liangcai Gao, and Wei Peng. Improving table structure recognition with visual-alignment sequential coordinate modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11134–11143, 2023. 2, 6
- [11] Yulong Hui, Yao Lu, and Huanchen Zhang. Uda: A benchmark suite for retrieval augmented generation in real-world document analysis. *Advances in Neural Information Processing Systems*, 37:67200–67217, 2024. 1
- [12] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 1, 6
- [13] Seil Kang, Jinyeong Kim, Junhyeok Kim, and Seong Jae Hwang. See what you are told: Visual attention sink in large multimodal models. *arXiv preprint arXiv:2503.03321*, 2025. 4, 5
- [14] Zhiqi Li, Guo Chen, Shilong Liu, Shihao Wang, Vibashan VS, Yishen Ji, Shiyi Lan, Hao Zhang, Yilin Zhao, Subhashree Radhakrishnan, et al. Eagle 2: Building post-training data strategies from scratch for frontier vision-language models. *arXiv preprint arXiv:2501.14818*, 2025. 4
- [15] Zhang Li, Yuliang Liu, Qiang Liu, Zhiyin Ma, Ziyang Zhang, Shuo Zhang, Zidun Guo, Jiarui Zhang, Xinyu Wang, and Xiang Bai. Monkeyocr: Document parsing with a structure-recognition-relation triplet paradigm. *arXiv preprint arXiv:2506.05218*, 2025. 1, 2
- [16] Weihong Lin, Zheng Sun, Chixiang Ma, Mingze Li, Jiawei Wang, Lei Sun, and Qiang Huo. Tsrformer: Table structure recognition with transformers. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 6473–6482, 2022. 2
- [17] Hao Liu, Xin Li, Bing Liu, Deqiang Jiang, Yinsong Liu, and Bo Ren. Neural collaborative graph machines for table structure recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4533–4542, 2022. 2
- [18] Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12):220102, 2024. 2
- [19] Yuan Liu, Zhongyin Zhao, Le Tian, Haicheng Wang, Xubing Ye, Yangxiu You, Zilin Yu, Chuhan Wu, Xiao Zhou, Yang Yu, et al. Points-reader: Distillation-free adaptation of vision-language models for document conversion. *arXiv preprint arXiv:2509.01215*, 2025. 1, 2
- [20] Rujiao Long, Wen Wang, Nan Xue, Feiyu Gao, Zhibo Yang, Yongpan Wang, and Gui-Song Xia. Parsing table structures in the wild. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 944–952, 2021. 2, 5, 1

- [21] Maksym Lysak, Ahmed Nassar, Nikolaos Livathinos, Christoph Auer, and Peter Staar. Optimized table tokenization for table structure recognition. In *International Conference on Document Analysis and Recognition*, pages 37–50. Springer, 2023. 5
- [22] Ahmed Nassar, Nikolaos Livathinos, Maksym Lysak, and Peter Staar. Tableformer: Table structure understanding with transformers. *arXiv preprint arXiv:2203.01017*, 2022. 2, 5
- [23] Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, et al. Mineru2. 5: A decoupled vision-language model for efficient high-resolution document parsing. *arXiv preprint arXiv:2509.22186*, 2025. 1, 2, 6
- [24] OpenAI. Introducing gpt-5. url <https://openai.com/index/introducing-gpt-5/>, 2025. Accessed: 2025-10-26. 6
- [25] Linke Ouyang, Yuan Qu, Hongbin Zhou, Jiawei Zhu, Rui Zhang, Qunshu Lin, Bin Wang, Zhiyuan Zhao, Man Jiang, Xiaomeng Zhao, et al. Omnidocbench: Benchmarking diverse pdf document parsing with comprehensive annotations. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24838–24848, 2025. 2, 6
- [26] ShengYun Peng, Aishwarya Chakravarthy, Seongmin Lee, Xiaoqing Wang, Rajarajeswari Balasubramanian, and Duen Horng Chau. Unitable: Towards a unified framework for table recognition via self-supervised pretraining. *arXiv preprint arXiv:2403.04822*, 2024. 2, 6
- [27] Jake Poznanski, Aman Rangapur, Jon Borchardt, Jason Dunkelberger, Regan Huff, Daniel Lin, Christopher Wilhelm, Kyle Lo, and Luca Soldaini. olmocr: Unlocking trillions of tokens in pdfs with vision language models. *arXiv preprint arXiv:2502.18443*, 2025. 1, 2
- [28] Yuxuan Qiao, Haodong Duan, Xinyu Fang, Junming Yang, Lin Chen, Songyang Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Prism: A framework for decoupling and assessing the capabilities of vlms. *Advances in Neural Information Processing Systems*, 37:111863–111898, 2024. 2
- [29] Sachin Raja, Ajoy Mondal, and CV Jawahar. Table structure recognition using top-down and bottom-up cues. In *European conference on computer vision*, pages 70–86. Springer, 2020. 2
- [30] rednote. dots.ocr: Multilingual document layout parsing in a single vision-language model. url <https://github.com/rednote-hilab/dots.ocr>, 2025. Accessed: 2025-10-26. 2, 6
- [31] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 2
- [32] Chris Tensmeyer, Vlad I Morariu, Brian Price, Scott Cohen, and Tony Martinez. Deep splitting and merging for table structure decomposition. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 114–121. IEEE, 2019. 2
- [33] Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, et al. Mineru: An open-source solution for precise document content extraction. *arXiv preprint arXiv:2409.18839*, 2024. 1
- [34] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025. 1, 6
- [35] Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, et al. General ocr theory: Towards ocr-2.0 via a unified end-to-end model. *arXiv preprint arXiv:2409.01704*, 2024. 2
- [36] Haoran Wei, Yaofeng Sun, and Yukun Li. Deepseek-ocr: Contexts optical compression. *arXiv preprint arXiv:2510.18234*, 2025. 6
- [37] Hangdi Xing, Feiyu Gao, Rujiao Long, Jiajun Bu, Qi Zheng, Liangcheng Li, Cong Yao, and Zhi Yu. Lore: Logical location regression network for table structure recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2992–3000, 2023. 2
- [38] Long Xing, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Jianze Liang, Qidong Huang, Jiaqi Wang, Feng Wu, and Dahua Lin. Caprl: Stimulating dense image caption capabilities via reinforcement learning. *arXiv preprint arXiv:2509.22647*, 2025. 2
- [39] Fan Yang, Lei Hu, Xinwu Liu, Shuangping Huang, and Zhenghui Gu. A large-scale dataset for end-to-end table recognition in the wild. *Scientific Data*, 10(1):110, 2023. 5
- [40] Zhibo Yang, Jun Tang, Zhaohai Li, Pengfei Wang, Jianqiang Wan, Humen Zhong, Xuejing Liu, Mingkun Yang, Peng Wang, Shuai Bai, et al. Cc-ocr: A comprehensive and challenging ocr benchmark for evaluating large multimodal models in literacy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21744–21754, 2025. 2, 6
- [41] Junyuan Zhang, Qintong Zhang, Bin Wang, Linke Ouyang, Zichen Wen, Ying Li, Ka-Ho Chow, Conghui He, and Wentao Zhang. Ocr hinders rag: Evaluating the cascading impact of ocr on retrieval-augmented generation. *arXiv preprint arXiv:2412.02592*, 2024. 1, 3
- [42] Qintong Zhang, Bin Wang, Victor Shea-Jay Huang, Junyuan Zhang, Zhengren Wang, Hao Liang, Conghui He, and Wentao Zhang. Document parsing unveiled: Techniques, challenges, and prospects for structured information extraction. *arXiv preprint arXiv:2410.21169*, 2024. 1
- [43] Zhenrong Zhang, Jianshu Zhang, Jun Du, and Fengren Wang. Split, embed and merge: An accurate table structure recognizer. *Pattern Recognition*, 126:108565, 2022. 2
- [44] Zhenrong Zhang, Pengfei Hu, Jiefeng Ma, Jun Du, Jianshu Zhang, Baocai Yin, Bing Yin, and Cong Liu. Semv2: Table separation line detection based on instance segmentation. *Pattern Recognition*, 149:110279, 2024. 2, 5, 1
- [45] Zhiyuan Zhao, Hengrui Kang, Bin Wang, and Conghui He. Doclayout-yolo: Enhancing document layout analysis through diverse synthetic data and global-to-local adaptive perception. *arXiv preprint arXiv:2410.12628*, 2024. 5
- [46] Mingyu Zheng, Xinwei Feng, Qingyi Si, Qiaoqiao She, Zheng Lin, Wenbin Jiang, and Weiping Wang. Multimodal table understanding. *arXiv preprint arXiv:2406.08100*, 2024. 5

- [47] Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno Yepes. Image-based table recognition: data, model, and evaluation. *arXiv preprint arXiv:1911.10683*, 2019. [5](#)
- [48] Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno Yepes. Image-based table recognition: data, model, and evaluation. In *European conference on computer vision*, pages 564–580. Springer, 2020. [2](#), [4](#)

TRivia: Self-supervised Fine-tuning of Vision-Language Models for Table Recognition

Supplementary Material

A. More Training Details

| | Stage-1 | Stage-2 | Stage-3 |
|----------|-------------------------------------------------------------------|----------------------------|----------------------------|
| Vision | Max Resolution | $1280 \times 28 \times 28$ | $1280 \times 28 \times 28$ |
| | #Tokens per Image | $256 \sim 1280$ | $256 \sim 1280$ |
| Data | Dataset | Synthetic Data | Real-world Data |
| | #Samples | 700K | 50K |
| Model | Trainable | LLM | All |
| | Sequence Length | 8192 | 8192 |
| Training | Batch Size | 32 | 128 |
| | LR: ψ_{ViT} | 2×10^{-6} | 2×10^{-7} |
| | LR: $\{\theta_{\text{MLP}}, \phi_{\text{LM}}\}$ | 1×10^{-5} | 1×10^{-6} |
| | Epoch | 1 | 1 |

Table 4. Training setup and hyperparameters in three training stages.

The training configurations for the three stages are summarized in Table 4. We use Qwen2.5-VL-3B as the backbone model. Across all stages, the number of image tokens ranges from 256 to 1280, corresponding to image resolutions from $256 * 28 * 28$ to $1280 * 28 * 28$. The prompt templates used for table recognition are provided in Table 8. In the third stage, which incorporates GRPO training, the sampling temperature is set to 1.2. We generate $G = 16$ samples per step and use a constant learning rate scheduler. All experiments are conducted with 8 X A100 80GB GPUs. Stage 1 training requires approximately one day, Stage 2 about two hours, and Stage 3 (GRPO) around two days.

B. Dataset Construction

For stage 1 and 2, we remove all samples with incomplete HTML tags during dataset construction. In stage 2, when only table structure tags were available [20, 44], we employ Qwen2.5-VL-72B to perform OCR over each cell’s bounding box to recover textual content.

In stage 3, we begin with 100K unlabeled images. We employ response-consistency sampling using the stage-2 model with a temperature of 1.0, producing eight outputs per image. The consistency score is computed via pairwise TEDS among these outputs. Then, we calculate the consistency score using pairwise TEDS among the 8 outputs. Images are uniformly sampled across consistency score intervals from 0.4 to 1.0 with a step size of 0.1. To promote diversity, we ensure that all table images originate from distinct PDFs. The filtering results in 50K selected images.

We then generate QA pairs using the attention-guided QA generation for these 50K images. We use Qwen2.5-VL-72B

as M_{QA} , generating 16 QAs per image with a temperature of 1.0 and prompt as shown in Table 5. Invalid JSON outputs are discarded. For cross-checking filtering, we adopt InternVL3-78B as M_{Val} to answer each question with and without table images as input and not having table images as input. We retain QA pairs whose F1 score exceeds 0.9 with the image but falls below 0.3 without it. Furthermore, we use M_{QA} to obtain the visual source of each QA pair and greedily search the final sets of QA pairs such that the IOU between any two QA pairs is less than 0.3. We remove images with fewer than 3 valid QAs to maintain reliable QA reward estimation. The final dataset comprises 48,470 images, each associated with an average of 28.3 QAs.

C. More Details about Experiments

All experiments are conducted using eight A100 80GB GPUs. For inference, we employ vLLM by default for all compatible models.

General-purpose VLMs for TR. When using general-purpose VLMs for table recognition, we perform prompt optimization to achieve the best results. We compare various templates from benchmarks such as OmnidocBench, CC-OCR, and OCRBench v2, and find that a unified prompt design (shown in Table 8) achieves the best overall performance. During generation, we set the sampling temperature to 0.2, which reduces repetitions relative to temperature 0 and mitigates hallucinations compared to temperature 1.0, yielding around a 3% performance gain. For Gemini 2.5 Pro, we enable “thinking mode”, which improves performance by approximately 3%. For the Qwen2.5-VL and Qwen3-VL, we set the number of image tokens to $256 \sim 1280$, which is the same as TRivia-3B.

Document parsing VLMs for TR. For document parsing VLMs, we adopt their specialized table recognition prompts. In the case of PaddleOCR-VL, we disable unwrapping and document orientation classification modules, as disabling them empirically improves performance.

D. Prompt Template

This section provides all prompt designs used in this paper.

Table 5 shows the prompts for QA generation. We include several heuristic constraints in it. For example, the questions should avoid direct references to rows or columns, such as "the third row", which is heavily limited to the structural parsing ability of the M_{QA} model. The questions should be simple, avoiding complex reasoning, to reduce the impact of M_{LLM} 's capability on reward estimation.

Since our task involves bilingual table recognition in Chinese and English, prompts for both question answering with the LLM and the VLM are prepared in both languages. This prevents the model from outputting answers in the wrong language, as shown in Table 6.

Table 7 shows the prompts that M_{Val} use for cross-check filtering.

Finally, the table recognition prompt for the general-purpose VLM is shown in Table 8, which is developed through multiple iterations of optimization for best performance.

<image>

Given an image, your task is to generate 10 reasonable and natural QA pairs based on the following rules:

1. The questions should be contextually appropriate and natural.
2. The answer to each question must be a short word or two or a numerical value.
3. For tables in Chinese, generate QA pairs in Chinese. Ensure the question and answer are both in the same language.
4. Each question should have one and only one answer.
5. Distribute the QA pairs across different parts of the table to cover multiple data points, avoiding any concentration on a single row or column.
6. Avoid questions that involve reasoning, such as numerical comparisons, maximum/minimum values, or calculations.
7. Do not directly mention the table structure in the question; instead, incorporate natural references. For example, instead of asking, "What is the journal account in the first row?" ask, "What is the journal account for serial number 1?"
8. Exclude questions that could be answered by only one data point, such as "Does the table include future goals?" or "Does the table list the budget?"
9. Ensure the questions can be answered using an HTML-formatted table, and avoid referencing visual orientation or relative positioning (e.g., "What is to the left of a T-account?").

If you cannot generate QA pairs that meet the above criteria, output "None". Otherwise, output the generated QAs in JSON format.

****Example Output Format:****

```
“json
[
  "question": "What is the market cap in Rmb mn?", "answer": "13,650.6",
  "question": "What is the 12 month price target?", "answer": "24.80",
]
”
```

Table 5. Q&A Generation Prompt

For question in English:

Given an HTML-formatted table and a corresponding question, your task is to respond appropriately based on table. If the table do not contain the answer of question, output "Not answerable".

Your answer should be a short phrase of only few words. Output the answer within <answer> </answer>.

HTML Table: {html_table}

Question: {question}

For question in Chinese:

给定一个HTML格式的表格以及一个相应的问题，你的任务是根据表格回答该问题。如果该表格不包含该问题的答案，请输出"无法回答"。你的答案必须简短、仅有一两个词语。输出答案时用<answer></answer>包裹。

HTML表格: {html_table}

问题: {html_table}

Table 6. LLM QA Prompt

For question in English with image input:

<image>

Given a table image and a corresponding question, your task is to respond appropriately based on table image. If the table do not contain the answer of question, output "Not answerable".

Your answer should be a short phrase of only few words. Output the answer within <answer> </answer>.

Question: {question}

For question in Chinese with image input:

<image>

给定一个表格图像以及一个相应的问题，你的任务是根据表格图片回答该问题。如果该表格不包含该问题的答案，请输出"无法回答"。你的答案必须简短、仅有一两个词语。输出答案时用<answer></answer>包裹。

问题: {question}

For question in English without image input:

Answer the following question. Your answer should be a short phrase of only few words. If you cannot answer this question, output "Not answerable".

Your answer should be a short phrase of only few words. Output the answer within <answer> </answer>.

Question: {question}

For question in Chinese without image input:

回答下面的问题。你的答案必须简短、仅有一两个词语。如果你无法回答该问题，请输出"无法回答"。你的答案必须简短、仅有一两个词语。输出答案时用<answer></answer>包裹。

问题: {question}

Table 7. VLM QA cross-checking Prompt

For general purpose VLMs:

<image>

You are an AI specialized in recognizing and extracting table from images. Your mission is to analyze the table image and generate the result in HTML format using specified tags. Output only the results without any other words and explanation.

For TRivia-3B:

<image>

You are an AI specialized in recognizing and extracting table from images. Your mission is to analyze the table image and generate the result in OTSL format using specified tags. Output only the results without any other words and explanation.

Table 8. Table Recognition Prompt