

# MinerU2.5: A Decoupled Vision-Language Model for Efficient High-Resolution Document Parsing

Junbo Niu<sup>1,2\*</sup>, Zheng Liu<sup>1,2\*</sup>, Zhuangcheng Gu<sup>1\*</sup>, Bin Wang<sup>1\*†</sup>, Linke Ouyang<sup>1\*</sup>, Zhiyuan Zhao<sup>1\*</sup>, Tao Chu<sup>1\*</sup>, Tianyao He<sup>1\*</sup>, Fan Wu<sup>1\*</sup>, Qintong Zhang<sup>1,2\*</sup>, Zhenjiang Jin<sup>1\*</sup>, Guang Liang<sup>1</sup>, Rui Zhang<sup>1</sup>, Wenzheng Zhang<sup>1,2</sup>, Yuan Qu<sup>1</sup>, Zhifei Ren<sup>1</sup>, Yuefeng Sun<sup>1</sup>, Zirui Tang<sup>1,3</sup>, Boyu Niu<sup>1,3</sup>, Yuanhong Zheng<sup>1</sup>, Dongsheng Ma<sup>1</sup>, Ziyang Miao<sup>1</sup>, Hejun Dong<sup>1</sup>, Siyi Qian<sup>1,2</sup>, Junyuan Zhang<sup>1</sup>, Fangdong Wang<sup>1</sup>, Jingzhou Chen<sup>1,2</sup>, Xiaomeng Zhao<sup>1</sup>, Liqun Wei<sup>1</sup>, Wei Li<sup>1,4</sup>, Shasha Wang<sup>1</sup>, Ruiliang Xu<sup>1</sup>, Yuanyuan Cao<sup>1</sup>, Lu Chen<sup>1</sup>, Qianqian Wu<sup>1</sup>, Huaiyu Gu<sup>1</sup>, Lindong Lu<sup>1</sup>, Dechen Lin<sup>1</sup>, Guanlin Shen<sup>1</sup>, Xuanhe Zhou<sup>1,3</sup>, Linfeng Zhang<sup>3</sup>, Yuhang Zang<sup>1</sup>, Xiaoyi Dong<sup>1</sup>, Jiaqi Wang<sup>1</sup>, Bo Zhang<sup>1</sup>, Lei Bai<sup>1</sup>, Pei Chu<sup>1</sup>, Weijia Li<sup>1</sup>, Jiang Wu<sup>1</sup>, Lijun Wu<sup>1</sup>, Zhenxiang Li<sup>1</sup>, Guangyu Wang<sup>1</sup>, Zhongying Tu<sup>1</sup>, Chao Xu<sup>1</sup>, Kai Chen<sup>1</sup>, Bowen Zhou<sup>1</sup>, Dahua Lin<sup>1✉</sup>, Wentao Zhang<sup>1,2✉</sup>, Conghui He<sup>1✉</sup>

<sup>1</sup>Shanghai Artificial Intelligence Laboratory, <sup>2</sup>Peking University

<sup>3</sup>Shanghai Jiao Tong University, <sup>4</sup>East China Normal University

heconghui@pjlab.org.cn

## Abstract

We introduce MinerU2.5, a 1.2B-parameter document parsing vision-language model that achieves state-of-the-art recognition accuracy while maintaining exceptional computational efficiency. Our approach employs a coarse-to-fine, two-stage parsing strategy that decouples global layout analysis from local content recognition. In the first stage, the model performs efficient layout analysis on downsampled images to identify structural elements, circumventing the computational overhead of processing high-resolution inputs. In the second stage, guided by the global layout, it performs targeted content recognition on native-resolution crops extracted from the original image, preserving fine-grained details in dense text, complex formulas, and tables. To support this strategy, we developed a comprehensive data engine that generates diverse, large-scale training corpora for both pretraining and fine-tuning. Ultimately, MinerU2.5 demonstrates strong document parsing ability, achieving state-of-the-art performance on multiple benchmarks, surpassing both general-purpose and domain-specific models across various recognition tasks, while maintaining significantly lower computational overhead.

## 1 Introduction

Document parsing (Zhang et al., 2024b) serves as a fundamental task in multimodal understanding, underpinning a variety of downstream applications such as information extraction (Liao et al., 2023; Wan et al., 2024), Retrieval-Augmented Generation

(RAG) (Lin, 2024; Zhang et al., 2024a; Zhao et al., 2024a) and intelligent document analysis (Bai et al., 2022; Blecher et al., 2023; Tang et al., 2023). In contrast to natural images, document images are characterized by significantly higher resolutions, denser content, and more complex structural layouts (Liu et al., 2024a; Wang et al., 2023; Wei et al., 2024). These inherent properties introduce a unique set of challenges. Firstly, the high resolution and fine-grained layout structures necessitate models capable of processing images at their native resolution. Secondly, the text-dense and often lengthy nature of documents imposes stringent requirements on the parameter efficiency and robustness of the models. Thirdly, the success of OCR is contingent not only on precise text recognition but also heavily on reliable layout analysis and efficient inference.

Contemporary approaches to document parsing can be broadly categorized into two paradigms: pipeline-based approaches (Cui et al., 2025; Livathinos et al., 2025; Paruchuri, 2025; Wang et al., 2024c) and end-to-end approaches based on VLMs (Achiam et al., 2023; Bai et al., 2025; Comanici et al., 2025; rednote, 2025; Wei et al., 2024). The former employs a modular design, decomposing the task into discrete stages such as layout detection, reading order prediction, and recognition of text lines, formulas, and tables. Each stage is handled by a specialized model. While this approach offers interpretability, it suffers from a cumbersome workflow and the potential for error propagation across modules. The latter paradigm exhibits superior semantic modeling capabilities, yet it is still widely constrained by the hallucination

\* indicates equal contribution. ✉ indicates corresponding author. † indicates Project leader.

problem in long-document processing and suffers from severe efficiency bottlenecks when dealing with high-resolution inputs. A critical factor limiting the performance and efficiency of VLM-based parsing is token redundancy, arising from large blank or low-information regions within the document image.

In response to the aforementioned challenges, we introduce a new document parsing framework, **MinerU2.5**. The key innovation is a decoupled architecture that separates *global layout analysis* from *local content recognition* via an efficient coarse-to-fine, two-stage inference mechanism. In the first stage, the model conducts fast and holistic layout analysis on downsampled document images, capturing the global structural organization with minimal computational cost. In the second stage, guided by the detected layout, it crops key regions from the original high-resolution input and performs fine-grained recognition within local windows, thereby preserving native resolution and ensuring high accuracy. This decoupled strategy not only reduces computational cost by an order of magnitude, primarily by avoiding the enormous number of visual tokens with  $\mathcal{O}(N^2)$  complexity inherent in end-to-end native-resolution approaches (Bai et al., 2025; Chen et al., 2025; rednote, 2025), but also brings multiple advantages: it significantly enhances the interpretability of parsing, effectively mitigates the common hallucination problem in VLMs, and allows the two stages to be independently optimized and iterated, resulting in more robust and efficient parsing capabilities. Ultimately, with its lightweight design of only 1.2B parameters, MinerU2.5 exhibits strong adaptability and efficiency in scenarios with long documents and high-density content while ensuring high parsing accuracy. Furthermore, to overcome the challenges of insufficient data diversity, sample imbalance, and inconsistent annotation quality in document parsing, we have developed a closed-loop data engine for complex documents. This engine systematically collects, processes, and generates large-scale, high-quality document corpora. This ensures that our model exhibits precise parsing capabilities and robustness across a wide spectrum of layouts, document types, and complex elements.

MinerU2.5 not only achieves state-of-the-art (SOTA) performance across a wide range of public benchmarks but also represents a qualitative leap in practical application and user experience over the previous MinerU2 version, as demonstrated by

the examples in [Appendix E](#).

## 2 Related Work

### 2.1 Traditional Pipelines

Early OCR systems (Cui et al., 2025; Livathinos et al., 2025; Paruchuri, 2025; Wang et al., 2024c) decompose document parsing into modular pipelines, sequentially executing layout detection (Wang et al., 2024a; Zhao et al., 2024b), text recognition (Cui et al., 2025), and reading order (Wang et al., 2021). For instance, Marker (Paruchuri, 2025) implements a sequential pipeline integrating Surya OCR (Paruchuri and Team, 2025) with layout analysis and reading order prediction modules to process diverse document types. MinerU (Wang et al., 2024c) leverages PDF-Extract-Kit (OpenDataLab, 2025) to orchestrate multiple specialized models for layout detection, formula recognition and table extraction. This modular architecture enables specialized optimization of individual components and facilitates targeted refinement of specific subtasks through well-defined module boundaries. However, pipeline-based methods are prone to error propagation across stages and exhibit limited robustness when confronted with complex layouts such as multi-column text or cross-page structures. Moreover, modular systems often entail multiple interdependencies in practice, rendering usage, maintenance, and updates cumbersome and less efficient.

### 2.2 General-Purpose Vision Language Models

General-purpose vision language models (VLMs) (Achiam et al., 2023; Bai et al., 2025; Comanici et al., 2025; Zhu et al., 2025) have emerged as an alternative paradigm for document understanding. Gemini2.5 Pro (Comanici et al., 2025) demonstrates strong OCR capabilities among general VLMs, surpassing traditional pipeline models like MinerU (Wang et al., 2024c) in text parsing and approaching specialized systems like UniMERNet (Wang et al., 2024b) in formula recognition, showcasing the potential of VLMs in OCR applications. Among open-source models, Qwen2.5-VL-72B (Bai et al., 2025) achieves the best results, using native-resolution vision encoders (Dehghani et al., 2023) to adapt to different image sizes, demonstrating the effectiveness of arbitrary-resolution processing in OCR tasks. However, these general models exhibit inherent limitations for document-centric tasks. Proprietary models like Gemini2.5 Pro (Co-

manici et al., 2025) are expensive and slow in processing, while open-source models require massive parameter scales for optimal performance, limiting practical deployment. Additionally, both types remain susceptible to hallucinations in densely populated text regions, affecting reliability in complex document layouts.

### 2.3 Domain-Specific Vision Language Models

**End-to-End Approaches.** Recent domain-specific models (Blecher et al., 2023; Chen et al., 2025; Kim et al., 2022; Liu et al., 2024b; Poznanski et al., 2025; rednote, 2025; Wei et al., 2024) adopt end-to-end architectures that unify document parsing within a single model, eliminating the need for cascaded processing stages. GOT (Wei et al., 2024), as an early representative of end-to-end approaches, pioneered the OCR 2.0 paradigm by establishing both model architecture and data methodology that unified recognition across diverse modalities—text, formulas, tables, and charts—within a single framework. Subsequent models like Ocean-OCR (Chen et al., 2025), olmOCR (Poznanski et al., 2025), and dots.ocr (rednote, 2025) leverage native resolution vision encoders to process documents and construct massive document corpora, further advancing the performance of end-to-end architectures. However, end-to-end designs face scalability challenges: joint optimization of layout and content often reduces accuracy on complex documents, while native-resolution processing introduces prohibitive  $\mathcal{O}(N^2)$  complexity. Despite strengths in semantic modeling, these models suffer from hallucinations on long documents and severe inefficiency with high-resolution inputs, where token redundancy from blank or low-information regions becomes a major bottleneck.

**Multi-Stage Approaches.** Recently, multi-stage methods (Feng et al., 2025; Li et al., 2025) leveraging VLMs decouple layout analysis from content recognition, combining the efficiency of pipeline approaches with the accuracy of unified models. Dolphin (Feng et al., 2025) employs a Swin-Transformer VLM that first performs page-level layout, then conducts efficient parallel parsing of identified regions. However, Swin-Transformer’s fixed resolution severely limits crop parsing—subregions with extreme aspect ratios suffer from distortion when resized to predetermined dimensions, degrading recognition quality while increas-

ing computational overhead. MonkeyOCR (Li et al., 2025) adopts a similar multi-stage strategy but employs a native resolution vision encoder in its second stage, improving both performance and efficiency. However, MonkeyOCR requires multiple specialized models across different stages, increasing system complexity and deployment overhead. A single unified model with native resolution parsing presents a promising direction to address these limitations, which is precisely the goal that MinerU2.5 pursues.

## 3 MinerU2.5

### 3.1 Model Architecture

As illustrate in Figure 1, MinerU2.5 follows the Qwen2-VL framework (Wang et al., 2024d), using Qwen2-Instruct-0.5B (Team, 2024) as the decoder with M-RoPE (Wang et al., 2024d) (replacing 1D-RoPE (Su et al., 2024)) for better generalization across crop resolutions and aspect ratios, a 675M NaViT (Dehghani et al., 2023) vision encoder initialized from Qwen2-VL with dynamic-resolution 2D-RoPE encoding (as window attention in Qwen2.5-VL (Bai et al., 2025) degrades parsing quality), and a pixel-unshuffle based patch merger (Shi et al., 2016) that merges adjacent  $2 \times 2$  vision tokens to reduce redundancy while balancing efficiency and accuracy.

### 3.2 Two-Stage Parsing Strategy

High-resolution documents contain large blank regions, causing severe visual token redundancy. Crop-based methods (Wei et al., 2024; Zhu et al., 2025) reduce cost but lose layout consistency, while native-resolution encoding (Bai et al., 2025; Guo et al., 2025; rednote, 2025; Niu et al., 2025) preserves details but is computationally impractical. We propose a two-stage strategy to decouple layout analysis and local recognition, improving interpretability and reducing hallucinations.

**Stage I: Layout Analysis.** We resize each page to a fixed  $1036 \times 1036$  thumbnail for global layout detection. This size trades off visibility and efficiency and yields stable localization compared to native-aspect thumbnails.

**Stage II: Content Recognition.** We crop regions from the original high-resolution image according to Stage I results and parse them at native resolution, capped by  $2048 \times 28 \times 28$  pixels to avoid both detail loss and redundant computation.

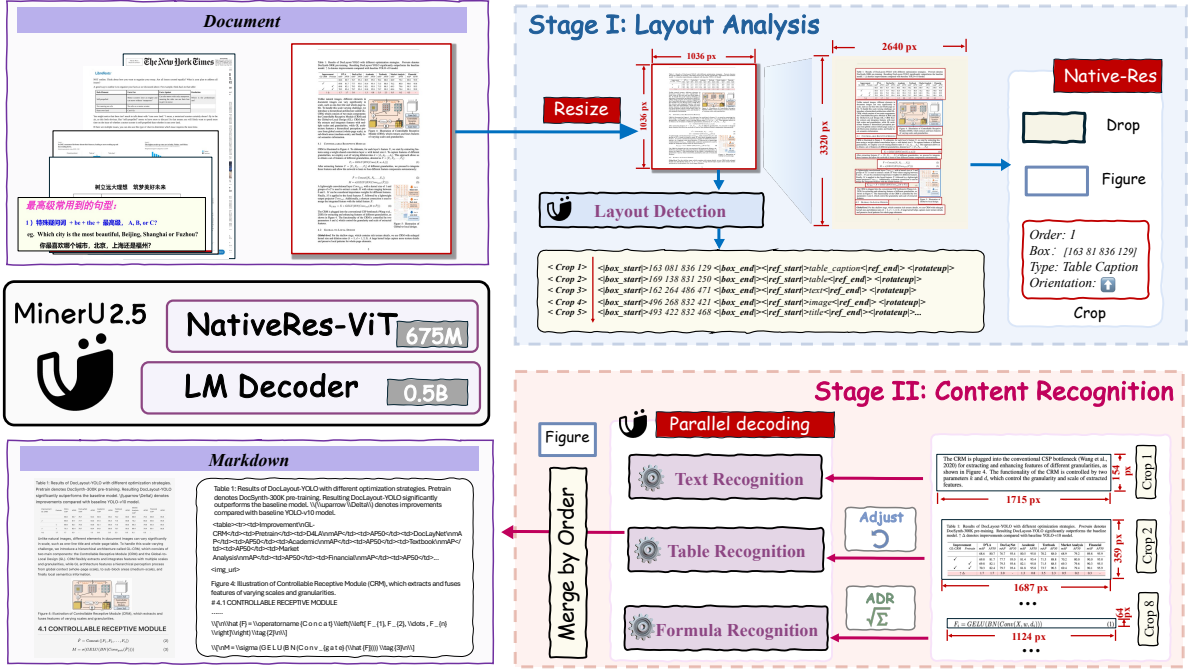


Figure 1: **The framework of MinerU2.5.** In stage I, MinerU2.5 performs rapid, global layout analysis on a downsampled page. In stage II, MinerU2.5 leverages the layout results to crop key regions from the original high-resolution document, performing fine-grained content recognition (e.g., text, table, and formula recognition) within these native-resolution local regions. The detailed prompts used in the inference are illustrated in Appendix D.

### 3.3 Training Recipe

MinerU2.5 initializes the vision encoder from Qwen2-VL-2B-Instruct and the LLM from Qwen2-Instruct-0.5B. Training consists of three stages.

#### 3.3.1 Stage 0: Modality Alignment

Stage 0 aligns vision and language representations by first training only the two-layer patch-merger MLP on *LLaVA-Pretrain* with the vision encoder and LLM frozen, followed by full-parameter visual instruction tuning on *LLaVA-Instruct* to improve instruction-following and OCR capability; removing this stage leads to higher loss and degraded parsing performance.

#### 3.3.2 Stage 1: Document Parsing Pre-training

Stage 1 jointly learns layout analysis and content recognition with all parameters trainable. Layout samples use resized full-page images with relative coordinates and the prompt “Layout Detection:”. Content recognition uses cropped text/formula/table blocks with prompts “Text/Formula/Table Recognition:”.

We train for 2 epochs, each with 6.9M samples: 2.3M layout, 2.4M text, 1.1M formula, and 1.1M table. The resulting model serves as a strong baseline and an efficient hard-sample miner for downstream

annotation.

#### 3.3.3 Stage 2: Document Parsing Fine-tuning

Stage 2 improves difficult cases while preserving core abilities. We mix diverse high-quality samples from Stage 1 with manually annotated hard cases mined from large-scale PDF corpora. We fine-tune for 3 epochs, each with 630K samples: 43K layout, 300K text, 147K formula, and 140K table.

### 3.4 Model Deployment

We build an offline pipeline based on vLLM (Kwon et al., 2023) with two optimizations: (1) asynchronous batching to overlap CPU/GPU workloads, and (2) decoupled Stage I/II inference to enable early downstream processing. To suppress degenerate repetition without harming structured outputs, we dynamically adjust `frequency_penalty` and `presence_penalty` based on detected layout types (e.g., higher for paragraphs, lower for tables).

We tune `max_num_batched_tokens`, `max_num_seqs`, and `cuda_graph_sizes` for higher utilization. On OmniDocBench (Ouyang et al., 2025) (1,355 pages), MinerU2.5 achieves 2.12 pages/s and 2337.25 tokens/s, outperforming MonkeyOCR-Pro-3B by 4× and dots.ocr by 7×. Even without deployment optimizations, it reaches

0.95 pages/s and 1045.14 tokens/s.

## 4 Evaluation

In this section, we present a comprehensive quantitative evaluation of MinerU2.5 to demonstrate its effectiveness in document parsing tasks. Specifically, we compare MinerU2.5 against leading general-purpose VLMs including GPT-4o (Achiam et al., 2023), Gemini-2.5 Pro (Comanici et al., 2025), and Qwen2.5-VL (Bai et al., 2025), as well as state-of-the-art domain-specific VLMs such as dots.ocr (rednote, 2025), MonkeyOCR (Li et al., 2025), and olmOCR (Poznanski et al., 2025).

### 4.1 Full-Document Parsing Task

We evaluate MinerU2.5’s full-document parsing performance on OmniDocBench (Ouyang et al., 2025), which provides a comprehensive evaluation across diverse document types, quality conditions, and parsing challenges.

#### 4.1.1 Evaluation Details and Metrics

Following OmniDocBench (Ouyang et al., 2025), we use its latest version with three key updates: (1) increased resolution for Notes and Newspapers from 72 to 200 DPI; (2) added 374 pages to better balance Chinese-English content and enrich formula coverage (1,355 pages in total); and (3) adopted a hybrid matching-based evaluation protocol.

#### 4.1.2 Evaluation Results

As shown in Table 1, MinerU2.5 achieves an overall score of 90.67 on OmniDocBench, outperforming the second-best model MonkeyOCR-pro-3B (Li et al., 2025) by 1.82 and dots.ocr (rednote, 2025) by 2.26 points. In text recognition tasks, MinerU2.5 achieves the lowest edit distance of 0.047, marginally better than dots.ocr at 0.048 and significantly outperforming PP-StructureV3 (Cui et al., 2025), which scores 0.073. For formula recognition, MinerU2.5 leads with a CDM score of 88.46, exceeding both Qwen2.5-VL-72B at 88.27 and MonkeyOCR-3B at 87.45. In table recognition tasks, MinerU2.5 achieves the highest TEDS score of 88.22 and TEDS-S score of 92.38. For reading order evaluation, it maintains the best edit distance of 0.044. The document-type specific results presented in Table 7 demonstrate that MinerU2.5 achieves best or second-best performance in 6 out of 9 categories. For textbooks, it delivers the best performance with an edit distance

of 0.0499, substantially outperforming dots.ocr’s 0.0788. For newspapers, MinerU2.5 leads with a score of 0.0540, surpassing all competing models. In both financial reports and slides categories, MinerU2.5 achieves second-best performance with scores of 0.0104 and 0.0294 respectively.

### 4.2 Element-Specific Parsing Task

#### 4.2.1 Layout Analysis

We validate the effectiveness of our layout analysis via a fair zero-shot comparison on three public datasets: OmniDocBench (Ouyang et al., 2025), a recent document parsing benchmark with detailed layout annotations; D<sup>4</sup>LA (Da et al., 2023), which contains 11,092 noisy document images annotated with 27 categories (we use its annotated test set); and DocLayNet (Pfitzmann et al., 2022), a large-scale dataset of 80,863 pages spanning seven document types with 11 categories (we use its annotated validation set).

We compare MinerU2.5 against recent baselines, including LayoutLMv3 (Huang et al., 2022), MinerU2-VLM (Wang et al., 2024c), DocLayout-YOLO (Zhao et al., 2024b), and PP-StructureV3 (Cui et al., 2025). To ensure an equitable evaluation, all models are tested without dataset-specific training. To address differences in label definitions and detection granularity, we map annotations into five coarse categories and adopt PageIoU for evaluation; the “Full Page” score measures spatial overlap regardless of category labels.

As shown in Table 2, MinerU2.5 consistently achieves the best Full Page F1-score@PageIoU across all benchmarks and leads on most individual element types, demonstrating strong generalization under zero-shot settings.

#### 4.2.2 Table Recognition

We evaluate MinerU2.5 against representative baselines, including traditional table recognition systems, general multimodal large models, and document parsing models, on five benchmarks (Table 3). The benchmarks include PubTabNet (Zhong et al., 2020), FinTabNet (Zheng et al., 2021), CC-OCR (Yang et al., 2024) and OCRBench v2 (Fu et al., 2024), and an in-house table recognition benchmark with ~500 diverse tables spanning various layouts and attributes.

Overall, MinerU2.5 achieves state-of-the-art or competitive performance across all benchmarks. On PubTabNet, although trained with only 20% of the training set, MinerU2.5 remains comparable

Model Type	Methods	Parameters	Overall $\uparrow$	Text $\text{Edit}\downarrow$	Formula $\text{CDM}\uparrow$	Table $\text{TEDS}\uparrow$	Table $\text{TEDS-S}\uparrow$	Read Order $\text{Edit}\downarrow$
Pipeline Tools	Marker-1.8.2 (Paruchuri, 2025)	-	71.30	0.206	76.66	57.88	71.17	0.250
	MinerU2-pipeline (Wang et al., 2024c)	-	75.51	0.209	76.55	70.90	79.11	0.225
	PP-StructureV3 (Cui et al., 2025)	-	86.73	0.073	85.79	81.68	89.48	0.073
General VLMs	GPT-4o (Achiam et al., 2023)	-	75.02	0.217	79.70	67.07	76.09	0.148
	InternVL3-76B (Zhu et al., 2025)	76B	80.33	0.131	83.42	70.64	77.74	0.113
	InternVL3.5-241B (Wang et al., 2025b)	241B	82.67	0.142	87.23	75.00	81.28	0.125
	Qwen2.5-VL-72B (Bai et al., 2025)	72B	87.02	0.094	<u>88.27</u>	82.15	86.22	0.102
	Gemini-2.5 Pro (Comanici et al., 2025)	-	88.03	0.075	85.82	85.71	90.29	0.097
Specialized VLMs	Dolphin (Feng et al., 2025)	322M	74.67	0.125	67.85	68.70	77.77	0.124
	OCRFlux (chatdoc.com, 2025)	3B	74.82	0.193	68.03	75.75	80.23	0.202
	Mistral-OCR (Team, 2025)	-	78.83	0.164	82.84	70.03	78.04	0.144
	POINTS-Reader (Liu et al., 2025b)	3B	80.98	0.134	79.20	77.13	81.66	0.145
	olmOCR-7B (Poznanski et al., 2025)	7B	81.79	0.096	86.04	68.92	74.77	0.121
	MinerU2-VLM(Wang et al., 2024c)	0.9B	85.56	0.078	80.95	83.54	87.66	0.086
	Nanonets-OCR-s (Mandalm, 2025)	3.7B	85.59	0.093	85.90	80.14	85.57	0.108
	MonkeyOCR-pro-1.2B (Li et al., 2025)	1.9B	86.96	0.084	85.02	84.24	89.02	0.130
	MonkeyOCR-3B (Li et al., 2025)	3.7B	87.13	0.075	87.45	81.39	85.92	0.129
	dots.ocr (rednote, 2025)	3B	88.41	0.048	83.22	<u>86.78</u>	90.62	<u>0.053</u>
	MonkeyOCR-pro-3B (Li et al., 2025)	3.7B	<u>88.85</u>	0.075	87.25	<u>86.78</u>	<u>90.63</u>	0.128
	<b>MinerU2.5</b>	1.2B	<b>90.67</b>	<b>0.047</b>	<b>88.46</b>	<b>88.22</b>	<b>92.38</b>	<b>0.044</b>

Table 1: Performance comparison of document parsing methods on OmniDocBench across text, formula, table, and reading order extraction tasks.

Method	Textual			Image			Table			Equation			Page Margins			Full Page		
	P $\uparrow$	R $\uparrow$	F1 $\uparrow$	P $\uparrow$	R $\uparrow$	F1 $\uparrow$	P $\uparrow$	R $\uparrow$	F1 $\uparrow$	P $\uparrow$	R $\uparrow$	F1 $\uparrow$	P $\uparrow$	R $\uparrow$	F1 $\uparrow$	P $\uparrow$	R $\uparrow$	F1 $\uparrow$
OmniDocBench (Ouyang et al., 2025)																		
LayoutLMv3 (Huang et al., 2022)	90.4	48.2	58.1	72.1	51.2	57.2	72.6	55.1	61.0	-	36.9	-	-	-	-	-	-	-
MinerU2-VLM (Wang et al., 2024c)	90.3	95.6	91.9	87.2	91.0	90.9	96.0	97.1	97.8	87.4	95.8	90.5	-	-	-	-	-	-
DocLayout-YOLO (Zhao et al., 2024b)	95.4	<b>98.3</b>	96.5	87.6	<b>96.7</b>	<u>94.7</u>	94.9	<u>98.1</u>	<u>98.4</u>	<u>95.3</u>	90.6	93.8	-	<b>98.7</b>	-	92.3	<b>97.7</b>	94.1
PP-StructureV3 (Cui et al., 2025)	96.8	96.7	96.6	86.4	92.1	92.9	<b>96.6</b>	97.4	98.2	<b>96.5</b>	<u>97.6</u>	<b>96.7</b>	<b>92.9</b>	86.2	<u>88.1</u>	<u>94.8</u>	96.2	<u>94.6</u>
MinerU2.5	<b>97.2</b>	<u>98.0</u>	<b>97.5</b>	<b>89.6</b>	<u>94.3</u>	<b>95.0</b>	<u>96.0</u>	<b>98.1</b>	<b>98.4</b>	92.4	<b>99.6</b>	<u>94.7</u>	<u>89.9</u>	<u>95.4</u>	<b>91.4</b>	<b>95.8</b>	<u>97.0</u>	<b>95.9</b>
D <sup>4</sup> LA (Da et al., 2023)																		
LayoutLMv3 (Huang et al., 2022)	86.9	41.2	52.4	<b>59.3</b>	32.0	31.4	59.3	41.8	43.3	-	50.5	-	-	-	-	-	-	-
MinerU2-VLM (Wang et al., 2024c)	88.3	88.9	87.9	<u>56.7</u>	35.0	38.1	<u>89.1</u>	<u>84.1</u>	<u>90.6</u>	<u>38.3</u>	<u>99.4</u>	79.1	-	-	-	-	-	-
DocLayout-YOLO (Zhao et al., 2024b)	86.3	<u>97.8</u>	<u>90.8</u>	41.5	<u>92.9</u>	62.6	87.6	<b>89.0</b>	89.8	31.9	80.2	<b>91.1</b>	-	<u>95.0</u>	-	82.6	<b>95.4</b>	<u>87.3</u>
PP-StructureV3 (Cui et al., 2025)	<u>88.5</u>	93.5	90.0	50.1	82.3	<u>67.9</u>	87.1	81.1	89.7	24.6	85.9	<u>92.1</u>	<b>76.8</b>	84.2	<u>79.1</u>	<u>85.7</u>	91.0	86.0
MinerU2.5	<b>91.8</b>	<b>98.3</b>	<b>94.6</b>	53.8	<b>94.3</b>	<b>72.8</b>	<b>91.9</b>	78.9	<b>91.4</b>	<b>46.0</b>	<b>100.0</b>	91.0	<u>75.9</u>	<b>97.6</b>	<b>84.2</b>	<b>90.4</b>	<u>92.5</u>	<b>90.2</b>
DocLaynet (Pfitzmann et al., 2022)																		
LayoutLMv3 (Huang et al., 2022)	88.8	59.3	67.9	79.0	50.3	61.9	75.2	54.9	61.8	-	31.9	-	-	-	-	-	-	-
MinerU2-VLM (Wang et al., 2024c)	88.1	96.1	91.7	85.5	78.1	91.3	94.9	<u>94.4</u>	95.6	83.9	<u>97.0</u>	90.0	-	-	-	-	-	-
DocLayout-YOLO (Zhao et al., 2024b)	86.9	96.8	91.2	85.8	<u>96.2</u>	91.3	92.0	<b>95.7</b>	94.8	80.5	86.9	82.8	-	<u>97.7</u>	-	88.0	<u>96.3</u>	90.9
PP-StructureV3 (Cui et al., 2025)	<b>90.9</b>	<u>97.3</u>	<u>93.8</u>	<u>91.7</u>	90.4	<u>94.2</u>	<b>96.4</b>	93.7	<u>96.7</u>	88.8	96.0	<u>92.1</u>	<b>76.8</b>	79.3	<b>77.4</b>	<u>92.4</u>	95.7	<u>93.0</u>
MinerU2.5	<u>90.2</u>	<b>99.6</b>	<b>94.8</b>	<b>92.5</b>	<b>96.3</b>	<b>95.9</b>	<u>96.3</u>	93.5	<b>97.1</b>	<b>88.9</b>	<b>98.6</b>	<b>93.5</b>	<u>76.3</u>	<b>98.9</b>	<b>86.3</b>	<b>92.8</b>	<b>97.7</b>	<b>94.6</b>

Table 2: Comparison of layout analysis performance (Precision@PageIoU, Recall@PageIoU, F1-score@PageIoU) across different methods and content types on multiple layout analysis benchmarks.

to the best-performing methods. On FinTabNet, MinerU2.5 achieves the best result with a clear margin, benefiting from large-scale high-quality financial table data. On CC-OCR and OCRBench v2, MinerU2.5 is competitive with strong multimodal baselines (e.g., Gemini-2.5 Pro and Qwen2.5-VL-72B), and significantly outperforms most other methods. On the in-house benchmark, MinerU2.5 achieves the best performance, slightly surpassing Gemini-2.5 Pro.

### 4.2.3 Formula Recognition

We compare MinerU2.5 with a diverse set of baselines, including specialized formula recognizers, document parsing models, and general

vision-language models. We evaluate on UniMER-Test (Wang et al., 2024b) (CPE/HWE/SPE/SCE), LaTeX-80M<sup>M</sup>, and an in-house benchmark comprising *Chinese* (formulas with Chinese characters), *Fuzzy Math* (low-quality textbook/exam scans with blur, degradation, and watermarks), and *Complex* (highly intricate formulas).

Results in Table 4 are evaluated using CDM (Wang et al., 2025a). MinerU2.5 achieves the best performance on four datasets and ranks second on one. It attains the highest CDM on SCE (96.4) and LaTeX-80M<sup>M</sup> (90.6), demonstrating strong robustness to blurred screenshots and complex matrix structures. While slightly behind specialized for-

Method	PubTabNet		FinTabNet		CC-OCR		OCRBench v2		In-house TR Benchmark	
	TEDS $\uparrow$	TEDS-S $\uparrow$	TEDS $\uparrow$	TEDS-S $\uparrow$	TEDS $\uparrow$	TEDS-S $\uparrow$	TEDS $\uparrow$	TEDS-S $\uparrow$	TEDS $\uparrow$	TEDS-S $\uparrow$
RapidTable (RapidAI, 2024)	86.57	<b>96.43</b>	73.77	84.84	50.93	65.84	65.55	77.73	51.96	71.94
MiniCPM-V 4.5 (Yu et al., 2025)	80.30	87.67	<u>85.41</u>	<u>89.18</u>	68.49	77.55	80.28	85.65	55.47	69.61
InternVL3.5-241B (Wang et al., 2025b)	83.75	88.76	84.74	87.92	62.87	69.52	79.5	85.81	56.32	69.3
Qwen2.5-VL-7B (Bai et al., 2025)	81.60	86.78	82.58	87.46	78.29	84.26	77.44	84.71	57.34	73.17
Qwen2.5-VL-72B (Bai et al., 2025)	84.39	87.91	82.90	87.13	<u>81.22</u>	<u>86.48</u>	81.33	86.58	62.79	76.91
GPT-4o (Achiam et al., 2023)	76.53	86.16	83.94	87.00	66.98	79.04	70.51	79.55	46.99	70.29
Gemini-2.5 Pro (Comanici et al., 2025)	-	-	-	-	<b>85.56</b>	<b>90.07</b>	<b>88.94</b>	<u>89.47</u>	<u>69.72</u>	<u>81.29</u>
dots.ocr (rednote, 2025)	<b>90.65</b>	<u>93.76</u>	84.12	87.86	75.42	81.65	82.04	86.27	66.91	79.27
Nanonets-OCR-s (Mandalm, 2025)	63.58	75.68	68.06	73.6	66.15	71.33	69.66	76.28	54.35	66.12
MinerU2-VLM (Wang et al., 2024c)	88.11	90.85	78.49	83.03	64.61	71.8	73.22	78.24	63.54	76.66
MinerU2.5	<u>89.07</u>	93.11	<b>95.97</b>	<b>97.61</b>	79.76	85.16	<u>87.13</u>	<b>90.62</b>	<b>71.48</b>	<b>82.83</b>

Table 3: Table Recognition Performance. MinerU2.5 achieves SOTA performance on most benchmarks among TEDS and TEDS-S metrics, and the remaining ones are also generally competitive with the SOTA. (CCOCR and OCRBench v2 are OCR evaluation benchmarks, we only select the subsets that contain tables. PubTabNet and FinTabNet have a large number of images, so we have not evaluate Gemini-2.5 Pro on them.).

Method	Public Dataset					In-house Dataset		
	CPE	HWE	SCE	SPE	LaTeX-80M <sup>M</sup>	Chinese	Fuzzy Math	Complex
UniMERNet* (Wang et al., 2024b)	<b>98.2</b>	<b>96.5</b>	95.4	<b>99.2</b>	83.9	84.0	84.3	67.9
PP-Formula_plus-L (Liu et al., 2025a)	<u>98.2</u>	<u>94.7</u>	<u>95.7</u>	<u>99.2</u>	85.9	84.0	86.5	76.5
Gemini-2.5-flash (Comanici et al., 2025)	89.2	90.0	85.1	97.5	78.7	88.1	89.4	80.1
Qwen2.5-VL-72B (Bai et al., 2025)	88.9	91.8	95.5	96.2	83.4	<b>90.8</b>	86.7	81.4
GPT-4o (Achiam et al., 2023)	82.7	85.9	87.8	96.7	73.4	88.3	85.0	78.6
InternVL3.5-241B (Wang et al., 2025b)	91.7	93.2	95.1	97.8	<u>86.9</u>	82.7	<u>90.3</u>	<u>82.0</u>
dots.ocr (rednote, 2025)	86.8	90.5	94.7	97.5	81.8	74.4	86.2	77.4
MinerU2.5	96.6	94.4	<b>96.4</b>	98.4	<b>90.6</b>	<u>90.7</u>	<b>92.6</b>	<b>82.2</b>

Table 4: Formula Recognition Performance (CDM metric used for evaluation). MinerU2.5 achieves 4 SOTA results and one second-best result across 7 benchmarks. Latex-80M<sup>M</sup> denotes the matrix benchmark of Latex-80M dataset. \* indicates that the UniMERNet results are based on an improved version compared to the publicly available open-source implementation.

mula recognizers on CPE/HWE/SPE, MinerU2.5 remains competitive. On our in-house benchmark, it matches Qwen2.5-VL-72B on *Chinese* (90.6) and achieves the best results on both *Fuzzy Math* and *Complex*.

## 5 Conclusion

In this paper, we present MinerU2.5, a 1.2B-parameter vision-language model for efficient document parsing. By adopting a decoupled coarse-to-fine strategy that separates global layout analysis from local recognition, MinerU2.5 achieves state-of-the-art performance with substantially reduced computation cost. Beyond standalone parsing, MinerU2.5 serves as a practical foundation for the LLM era, enabling rapid conversion of large-scale unstructured documents into structured, high-quality data. By better preserving tables, formulas, and layout semantics, it can further improve document understanding and enhance the reliability of

retrieval-augmented generation systems.

## Limitations

The proposed two-stage paradigm is sensitive to errors in the first-stage layout analysis, where misclassified or missing regions can propagate and cause unrecoverable failures. This problem is exacerbated by mismatches between training and real-world data distributions, especially under complex and noisy conditions. Therefore, improving the diversity and coverage of the training data—through augmentation, synthetic data, or domain adaptation—is crucial for enhancing the robustness of the first stage and the overall pipeline.

## Acknowledgements

This project was supported by Shanghai Artificial Intelligence Laboratory.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Haoli Bai, Zhiguang Liu, Xiaojun Meng, Wentao Li, Shuang Liu, Nian Xie, Rongfu Zheng, Liangwei Wang, Lu Hou, Jiansheng Wei, and 1 others. 2022. Wukong-reader: Multi-modal pre-training for fine-grained visual document understanding. *arXiv preprint arXiv:2212.09621*.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-v1 technical report. *arXiv preprint arXiv:2502.13923*.
- Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. 2023. Nougat: Neural optical understanding for academic documents. *arXiv preprint arXiv:2308.13418*.
- chatdoc com. 2025. Ocrflux. <https://github.com/chatdoc-com/OCRFlux>. Accessed:2025-09-25.
- Song Chen, Xinyu Guo, Yadong Li, Tao Zhang, Mingan Lin, Dongdong Kuang, Youwei Zhang, Lingfeng Ming, Fengyu Zhang, Yuran Wang, and 1 others. 2025. Ocean-ocr: Towards general ocr application via a vision-language model. *arXiv preprint arXiv:2501.15558*.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.
- Cheng Cui, Ting Sun, Manhui Lin, Tingquan Gao, Yubo Zhang, Jiakuan Liu, Xueqing Wang, Zelun Zhang, Changda Zhou, Hongen Liu, and 1 others. 2025. Paddleocr 3.0 technical report. *arXiv preprint arXiv:2507.05595*.
- Cheng Da, Chuwei Luo, Qi Zheng, and Cong Yao. 2023. Vision grid transformer for document layout analysis. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 19462–19472.
- Mostafa Dehghani, Basil Mustafa, Josip Djolonga, Jonathan Heek, Matthias Minderer, Mathilde Caron, Andreas Steiner, Joan Puigcerver, Robert Geirhos, Ibrahim M Alabdulmohsin, and 1 others. 2023. Patch n’pack: Navit, a vision transformer for any aspect ratio and resolution. *Advances in Neural Information Processing Systems*, 36:2252–2274.
- Hao Feng, Shu Wei, Xiang Fei, Wei Shi, Yingdong Han, Lei Liao, Jinghui Lu, Binghong Wu, Qi Liu, Chunhui Lin, and 1 others. 2025. Dolphin: Document image parsing via heterogeneous anchor prompting. *arXiv preprint arXiv:2505.14059*.
- Ling Fu, Zhebin Kuang, Jiajun Song, Mingxin Huang, Biao Yang, Yuzhe Li, Linghao Zhu, Qidi Luo, Xinyu Wang, Hao Lu, and 1 others. 2024. Ocrbench v2: An improved benchmark for evaluating large multi-modal models on visual text localization and reasoning. *arXiv preprint arXiv:2501.00321*.
- Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang, Jianyu Jiang, Jiawei Wang, and 1 others. 2025. Seed1. 5-v1 technical report. *arXiv preprint arXiv:2505.07062*.
- Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. Layoutlmv3: Pre-training for document ai with unified text and image masking. In *Proceedings of the 30th ACM international conference on multimedia*, pages 4083–4091.
- Geewook Kim, Teakgyu Hong, Moonbin Yim, Jeongyeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoon Yun, Dongyoon Han, and Seunghyun Park. 2022. Ocr-free document understanding transformer. In *European Conference on Computer Vision*, pages 498–517. Springer.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Zhang Li, Yuliang Liu, Qiang Liu, Zhiyin Ma, Ziyang Zhang, Shuo Zhang, Zidun Guo, Jiarui Zhang, Xinyu Wang, and Xiang Bai. 2025. Monkeyocr: Document parsing with a structure-recognition-relation triplet paradigm. *arXiv preprint arXiv:2506.05218*.
- Haofu Liao, Aruni RoyChowdhury, Weijian Li, Ankan Bansal, Yuting Zhang, Zhuowen Tu, Ravi Kumar Satzoda, R Manmatha, and Vijay Mahadevan. 2023. Doctr: Document transformer for structured information extraction in documents. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19584–19594.
- Demiao Lin. 2024. Revolutionizing retrieval-augmented generation with enhanced pdf structure recognition. *arXiv preprint arXiv:2401.12599*.
- Chaohu Liu, Kun Yin, Haoyu Cao, Xinghua Jiang, Xin Li, Yinsong Liu, Deqiang Jiang, Xing Sun, and Linli Xu. 2024a. Hrvda: High-resolution visual document assistant. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15534–15545.
- Hongen Liu, Cheng Cui, Yuning Du, Yi Liu, and Gang Pan. 2025a. Pp-formulanet: Bridging accuracy and efficiency in advanced formula recognition. *arXiv preprint arXiv:2503.18382*.
- Yuan Liu, Zhongyin Zhao, Le Tian, Haicheng Wang, Xubing Ye, Yangxiu You, Zilin Yu, Chuhan Wu, Xiao

- Zhou, Yang Yu, and 1 others. 2025b. Points-reader: Distillation-free adaptation of vision-language models for document conversion. *arXiv preprint arXiv:2509.01215*.
- Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai. 2024b. Textmonkey: An ocr-free large multimodal model for understanding document. *arXiv preprint arXiv:2403.04473*.
- Nikolaos Livathinos, Christoph Auer, Maksym Lysak, Ahmed Nassar, Michele Dolfi, Panos Vagenas, Cesar Berrospi Ramis, Matteo Omenetti, Kasper Dinkla, Yusik Kim, and 1 others. 2025. Docling: An efficient open-source toolkit for ai-driven document conversion. *arXiv preprint arXiv:2501.17887*.
- Maksym Lysak, Ahmed Nassar, Nikolaos Livathinos, Christoph Auer, and Peter Staar. 2023. Optimized table tokenization for table structure recognition. In *International Conference on Document Analysis and Recognition*, pages 37–50. Springer.
- Souvik Mandalm. 2025. Nanonets-ocr-s. <https://nanonets.com/research/nanonets-ocr-s/>. Accessed:2025-09-25.
- Mathpix. 2025. Mathpix. <https://mathpix.com/>. Accessed:2025-09-25.
- Ahmed Nassar, Andres Marafioti, Matteo Omenetti, Maksym Lysak, Nikolaos Livathinos, Christoph Auer, Lucas Morin, Rafael Teixeira de Lima, Yusik Kim, A Said Gurbuz, and 1 others. 2025. Smoldocling: An ultra-compact vision-language model for end-to-end multi-modal document conversion. *arXiv preprint arXiv:2503.11576*.
- Junbo Niu, Yuanhong Zheng, Ziyang Miao, Hejun Dong, Chunjiang Ge, Hao Liang, Ma Lu, Bohan Zeng, Qiahao Zheng, Conghui He, and 1 others. 2025. Native visual understanding: Resolving resolution dilemmas in vision-language models. *arXiv preprint arXiv:2506.12776*.
- OpenDataLab. 2025. Pdf-extract-kit. <https://github.com/opendatalab/PDF-Extract-Kit>. Accessed:2025-09-25.
- Linke Ouyang, Yuan Qu, Hongbin Zhou, Jiawei Zhu, Rui Zhang, Qunshu Lin, Bin Wang, Zhiyuan Zhao, Man Jiang, Xiaomeng Zhao, and 1 others. 2025. Omnidocbench: Benchmarking diverse pdf document parsing with comprehensive annotations. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24838–24848.
- Vik Paruchuri. 2025. Marker. <https://github.com/datalab-to/marker>. Accessed:2025-09-25.
- Vikas Paruchuri and Datalab Team. 2025. Surya: A lightweight document ocr and analysis toolkit. <https://github.com/VikParuchuri/surya>. Accessed:2025-09-25.
- Birgit Pfitzmann, Christoph Auer, Michele Dolfi, Ahmed S Nassar, and Peter Staar. 2022. Doclaynet: A large human-annotated dataset for document-layout segmentation. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 3743–3751.
- Jake Poznanski, Aman Rangapur, Jon Borchardt, Jason Dunkelberger, Regan Huff, Daniel Lin, Christopher Wilhelm, Kyle Lo, and Luca Soldaini. 2025. olmoocr: Unlocking trillions of tokens in pdfs with vision language models. *arXiv preprint arXiv:2502.18443*.
- RapidAI. 2024. Rapid table. <https://github.com/RapidAI/RapidTable>. Accessed: 2025-9-25.
- rednote. 2025. dots.ocr: Multilingual document layout parsing in a single vision-language model. <https://github.com/rednote-hilab/dots.ocr>. Accessed:2025-09-25.
- Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. 2016. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1874–1883.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063.
- Zineng Tang, Ziyi Yang, Guoxin Wang, Yuwei Fang, Yang Liu, Chenguang Zhu, Michael Zeng, Cha Zhang, and Mohit Bansal. 2023. Unifying vision, text, and layout for universal document processing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19254–19264.
- Mistral AI Team. 2025. Mistral-ocr. [https://mistral.ai/news/mistral-ocr?utm\\_source=ai-bot.cn](https://mistral.ai/news/mistral-ocr?utm_source=ai-bot.cn). Accessed:2025-09-25.
- Qwen Team. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Jianqiang Wan, Sibao Song, Wenwen Yu, Yuliang Liu, Wenqing Cheng, Fei Huang, Xiang Bai, Cong Yao, and Zhibo Yang. 2024. Omniparser: A unified framework for text spotting key information extraction and table recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15641–15653.
- Ao Wang, Hui Chen, Lihao Liu, Kai Chen, Zijia Lin, Jungong Han, and 1 others. 2024a. Yolov10: Real-time end-to-end object detection. *Advances in Neural Information Processing Systems*, 37:107984–108011.
- Bin Wang, Zhuangcheng Gu, Guang Liang, Chao Xu, Bo Zhang, Botian Shi, and Conghui He. 2024b. Unimernet: A universal network for real-world mathematical expression recognition. *arXiv preprint arXiv:2404.15254*.

- Bin Wang, Fan Wu, Linke Ouyang, Zhuangcheng Gu, Rui Zhang, Renqiu Xia, Botian Shi, Bo Zhang, and Conghui He. 2025a. Image over text: Transforming formula recognition evaluation with character detection matching. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19681–19690.
- Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, and 1 others. 2024c. Mineru: An open-source solution for precise document content extraction. *arXiv preprint arXiv:2409.18839*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, and 1 others. 2024d. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, and 1 others. 2025b. Internvl3. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*.
- Zilong Wang, Yiheng Xu, Lei Cui, Jingbo Shang, and Furu Wei. 2021. Layoutreader: Pre-training of text and layout for reading order detection. *arXiv preprint arXiv:2108.11591*.
- Zilong Wang, Yichao Zhou, Wei Wei, Chen-Yu Lee, and Sandeep Tata. 2023. Vrdu: A benchmark for visually-rich document understanding. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5184–5193.
- Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yang Peng, and 1 others. 2024. General ocr theory: Towards ocr-2.0 via a unified end-to-end model. *arXiv preprint arXiv:2409.01704*.
- Zhibo Yang, Jun Tang, Zhaohai Li, Pengfei Wang, Jianqiang Wan, Humen Zhong, Xuejing Liu, Mingkun Yang, Peng Wang, Shuai Bai, and 1 others. 2024. Cc-ocr: A comprehensive and challenging ocr benchmark for evaluating large multimodal models in literacy. *arXiv preprint arXiv:2412.02210*.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, and 1 others. 2024. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*.
- Tianyu Yu, Zefan Wang, Chongyi Wang, Fuwei Huang, Wenshuo Ma, Zhihui He, Tianchi Cai, Weize Chen, Yuxiang Huang, Yuanqian Zhao, Bokai Xu, Junbo Cui, Yingjing Xu, Liqing Ruan, Luoyuan Zhang, Hanyu Liu, Jingkun Tang, Hongyuan Liu, Qining Guo, and 15 others. 2025. Minicpm-v 4.5: Cooking efficient mllms via architecture, data, and training recipe. *arXiv preprint arXiv:2509.18154*.
- Junyuan Zhang, Qintong Zhang, Bin Wang, Linke Ouyang, Zichen Wen, Ying Li, Ka-Ho Chow, Conghui He, and Wentao Zhang. 2024a. Ocr hinders rag: Evaluating the cascading impact of ocr on retrieval-augmented generation. *arXiv preprint arXiv:2412.02592*.
- Qintong Zhang, Bin Wang, Victor Shea-Jay Huang, Junyuan Zhang, Zhengren Wang, Hao Liang, Conghui He, and Wentao Zhang. 2024b. Document parsing unveiled: Techniques, challenges, and prospects for structured information extraction. *arXiv preprint arXiv:2410.21169*.
- Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang, Jie Jiang, and Bin Cui. 2024a. Retrieval-augmented generation for ai-generated content: A survey. *arXiv preprint arXiv:2402.19473*.
- Zhiyuan Zhao, Hengrui Kang, Bin Wang, and Conghui He. 2024b. Doclayout-yolo: Enhancing document layout analysis through diverse synthetic data and global-to-local adaptive perception. *arXiv preprint arXiv:2410.12628*.
- Xinyi Zheng, Douglas Burdick, Lucian Popa, Xu Zhong, and Nancy Xin Ru Wang. 2021. Global table extractor (gte): A framework for joint table identification and cell structure recognition using visual context. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 697–706.
- Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno Yepes. 2020. Image-based table recognition: data, model, and evaluation. In *European conference on computer vision*, pages 564–580. Springer.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, and 1 others. 2025. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*.

## A Training Details

As described in Section 3.1, MinerU2.5 consists of three core components: vision encoder, patch merger, and language model. Prior to the pre-training phase of MinerU2.5, the vision encoder is initialized from Qwen2-VL-2B-Instruct, while the language model is initialized from Qwen2-Instruct-0.5B. The overall training procedure of MinerU2.5 is divided into three stages, as summarized in Table 5.

### A.0.1 Stage 0-Modality Alignment

To ensure that MinerU2.5 acquires the fundamental vision–language alignment ability as well as the OCR recognition capability, we first conduct two-stage modality alignment training on Visual Question Answering (VQA) datasets.

**Language-Image Alignment.** Only the two-layer MLP within the patch merger is trained, while both the vision encoder and the language model are frozen. We use image-caption pairs<sup>1</sup> for training to effectively project visual features into the LLM embedding space, thus achieving alignment of the modal representation.

**Visual Instruction Tuning.** All model parameters are unfrozen. The focus is on knowledge accumulation and ability expansion, particularly strengthening visual alignment and OCR capability. The training data<sup>2</sup> mainly covers image captioning, interleaved text-image pairs, visual alignment, and OCR data. The goal is to enable MinerU2.5 to follow instructions across diverse visual tasks and generate reasonable responses.

		Stage-0		Stage-1	Stage-2
		a	b		
<i>Vision</i>	<b>Max Resolution</b>	$2048 \times 28 \times 28$	$4096 \times 28 \times 28$	$2048 \times 28 \times 28$	$2048 \times 28 \times 28$
	<b>#Tokens per Image</b>	4 ~ 2048	4 ~ 4096	4 ~ 2048	4 ~ 2048
<i>Data</i>	<b>Dataset</b>	Image Caption	VQA	Layout&OCR	Layout&OCR
	<b>#Samples</b>	558K	665K	6.9M	630K
<i>Model</i>	<b>Trainable</b>	MLP Adaptor	All	All	All
	<b>Sequence Length</b>	4096	4096	8192	16384
	<b>Data Augmentation</b>	No	No	Yes	Yes
<i>Training</i>	<b>Batch Size</b>	128	64	256	256
	<b>LR: <math>\psi_{\text{ViT}}</math></b>	$1 \times 10^{-3}$	$1 \times 10^{-5}$	$4 \times 10^{-6}$	$4 \times 10^{-6}$
	<b>LR: <math>\{\theta_{\text{MLP}}, \phi_{\text{LM}}\}</math></b>	$1 \times 10^{-3}$	$1 \times 10^{-5}$	$4 \times 10^{-5}$	$4 \times 10^{-5}$
	<b>Epoch</b>	1	1	2	3

Table 5: Training setup and hyperparameters in three training stages.

Empirical results demonstrate that MinerU2.5, after VQA-based modality alignment training, exhibits significant improvements in tasks such as layout analysis and content recognition. Conversely, skipping this stage leads to higher losses and a clear drop in overall performance.

### A.0.2 Stage 1-Document Parsing Pre-training

The objective of the document parsing pre-training stage is to enable MinerU2.5 to acquire two fundamental capabilities: **layout analysis** and **content recognition**. At this stage, all parameters of the model remain fully trainable.

<sup>1</sup>This dataset is sourced from *LLaVA-Pretrain*.

<sup>2</sup>This dataset is sourced from *LLaVA-Instruct*.

**Training Data.** We leveraged a large-scale mixture of model-labeled data and public datasets to ensure both sufficient scale and document diversity. For layout analysis, in consideration of training efficiency, full document images were resized to a fixed resolution with corresponding relative coordinates, and the prompt “Layout Detection:” was used. For content recognition, we employed single-element image samples of text blocks, formula blocks, and table blocks as inputs, with prompts “Text Recognition:”, “Formula Recognition:”, and “Table Recognition:” respectively. More details are shown in the [Appendix D](#).

**Training Configuration.** The model, initialized from Stage 0, was trained for 2 epochs. Each epoch consisted of a total of 6.9M samples, including 2.3M for layout analysis, 2.4M for text blocks, 1.1M for formula blocks, and 1.1M for table blocks.

Through this document parsing pre-training, the model has acquired strong layout analysis and content recognition capabilities, demonstrating excellent performance across most simple and medium-level scenarios. The resulting model not only serves as a **strong baseline** for downstream fine-tuning, but also functions as an **efficient hard-sample miner** within our data engineering pipeline, facilitating the identification of challenging cases for human annotation and further improving document parsing performance.

### A.0.3 Stage 2-Document Parsing Fine-tuning

The objective of the document parsing fine-tuning stage is to further enhance parsing performance in challenging scenarios, while maintaining the detection and parsing capabilities already acquired by MinerU2.5.

**Training Data.** To achieve this goal, it is crucial to construct a compact yet high-quality dataset:

- To preserve the model’s fundamental capabilities, we sampled high-quality and diverse examples from the pre-training dataset via data engineering and incorporated them into Stage 2 training, ensuring broad coverage across different document element types.
- From a large-scale, multi-source PDF corpus, we employed data engineering to identify cases where the model still underperformed. We summarized these difficult scenarios and conducted targeted data collection with manual annotation to obtain high-quality samples representing challenging cases.

**Training Configuration.** We fine-tuned the pre-trained model for 3 epochs. Each epoch contained a total of 630K samples, consisting of 43K for layout analysis, 300K for text blocks, 147K for formula blocks, and 140K for table blocks.

With this targeted data iteration strategy, Stage 2 fine-tuning enables the model to not only retain its established document parsing abilities but also achieve significant improvements in previously challenging scenarios.

### A.0.4 Data Augmentation Strategies

To enhance the model’s robustness in handling diverse documents in an open-world setting, we designed a variety of targeted data augmentation strategies during both Stage 1 and Stage 2. These augmentations simulate common types of document interference, and can be categorized as shown in Table 6.

Note that spatial transformations are not applied to layout analysis samples. For different element types, we carefully design augmentation parameters and probabilities in order to strike a balance between model performance and robustness.

## B Evaluation Details

This section provides additional experimental details of MinerU2.5, including dataset configurations, evaluation protocols, and implementation settings, to facilitate reproducibility and support the main results.

Augmentation Type	Operations
Spatial Transformations	Scaling, Grid Distortion, Rotation
Background Transformations	Texture, Weather effect, Image background, Watermark, Scanlines, Shadow
Color Transformations	Brightness Contrast, Illumination, RGB Shift
Degradation Transformations	PSF Blur, Vibration Blur, Gaussian Blur, Erosion / Dilation

Table 6: Data augmentation strategies for document parsing.

## B.1 Full-Document Parsing Task

We evaluate MinerU2.5’s full document parsing performance on three prominent benchmarks: OmniDocBench (Ouyang et al., 2025), Ocean-OCR (Chen et al., 2025) benchmarks, and olmOCR-bench (Poznanski et al., 2025). These benchmarks provide comprehensive evaluation from different dimensions, covering diverse document types, various quality conditions, and different parsing challenges to thoroughly assess the model’s robustness and generalization capabilities.

- **OmniDocBench** (Ouyang et al., 2025): This evaluation dataset is designed for diverse document parsing in real-world scenarios, encompassing nine document types, four layout types, and three language types. It offers a comprehensive assessment of parsing scores for text, formulas, tables, and reading order in full-document parsing, as well as for element-specific parsing tasks.
- **olmOCR-bench** (Poznanski et al., 2025): This evaluation dataset comprises 1,402 PDF documents sourced from various repositories, organized into seven subsets. Certain test patterns are applicable across all document types (e.g., presence, absence, reading order), while others are specifically targeted at challenging yet crucial content extraction objectives (e.g., tables, mathematical formulas).
- **Ocean-OCR benchmark** (Chen et al., 2025): This evaluation dataset consists of 100 images from English papers and 100 images from Chinese papers. It primarily evaluates the ability of text parsing and employs several text OCR-related evaluation metrics, such as Normalized Edit Distance, F1 Score, Precision, Recall, BLEU, and METEOR.

### B.1.1 Evaluation Details and Metrics

For OmniDocBench (Ouyang et al., 2025), we evaluate on the latest version with three key improvements:

- Enhanced resolution for Notes and Newspapers from 72 to 200 DPI, enabling more accurate evaluation of fine-grained text and handwritten content.
- An addition of 374 pages to balance Chinese-English content distribution and enrich mathematical formula coverage. Currently, it contains a total of 1,355 pages.
- Evaluation methodology updated to hybrid matching algorithm.

The Overall score combines three core metrics:

$$\text{Overall} = \frac{(1 - \text{Text}^{\text{Edit}}) \times 100}{3} + \frac{\text{Table}^{\text{TEDS}} + \text{Formula}^{\text{CDM}}}{3}. \quad (1)$$

For olmOCR-bench (Poznanski et al., 2025), we replace the formula scores of Arxiv Math (AR) and Old Scans Math (OSM) with the more reliable ExpRate of CDM (Wang et al., 2025a). The original evaluation compares LaTeX formulas by parsing them into abstract syntax trees and matching Unicode tokens, which is overly sensitive to syntax variations (e.g., `\cdots` vs. `\dotsb`) that render identically but are scored as different. To avoid this bias, we adopt ExpRate, which directly compares rendered outputs, assigning 1 for exact matches and 0 otherwise.

Model Type	Models	Slides	Academic Papers	Book	Textbook	Exam Papers	Magazine	Newspaper	Notes	Financial Report
Pipeline Tools	Marker-1.8.2 (Paruchuri, 2025)	0.1796	0.0412	0.1010	0.2908	0.2958	0.1111	0.2717	0.4656	0.0341
	MinerU2-pipeline (Wang et al., 2024c)	0.4244	0.0230	0.2628	0.1224	0.0822	0.3950	0.0736	0.2603	0.0411
	PP-StructureV3 (Cui et al., 2025)	0.0794	0.0236	0.0415	0.1107	0.0945	0.0722	<u>0.0617</u>	0.1236	0.0181
General VLMs	GPT-4o (Achiam et al., 2023)	0.1019	0.1203	0.1288	0.1599	0.1939	0.1420	0.6254	0.2611	0.3343
	InternVL3-76B (Zhu et al., 2025)	0.0349	0.1052	0.0629	0.0827	0.1007	0.0406	0.5826	<u>0.0924</u>	0.0665
	InternVL3.5-241B (Wang et al., 2025b)	0.0475	0.0857	<b>0.0237</b>	0.1061	0.0933	0.0577	0.6403	0.1357	0.1117
	Qwen2.5-VL-72B (Bai et al., 2025)	0.0422	0.0801	0.0586	0.1146	<u>0.0681</u>	0.0964	0.2380	0.1232	0.0264
	Gemini-2.5 Pro (Comanici et al., 2025)	0.0326	<u>0.0182</u>	0.0694	0.1618	0.0937	<b>0.0161</b>	0.1347	0.1169	0.0169
Specialized VLMs	Dolphin (Feng et al., 2025)	0.0957	0.0453	0.0616	0.1333	0.1684	0.0702	0.2388	0.2561	0.0186
	OCRFlux (chatdoc.com, 2025)	0.0870	0.0867	0.0818	0.1843	0.2072	0.1048	0.7304	0.1567	0.0193
	Mistral-OCR (Team, 2025)	0.0917	0.0531	0.0610	0.1349	0.1341	0.0581	0.5643	0.3097	0.0523
	POINTS-Reader (Liu et al., 2025b)	0.0334	0.0779	0.0671	0.1372	0.1901	0.1343	0.3789	0.0937	0.0951
	olmOCR-7B (Poznanski et al., 2025)	0.0497	0.0365	0.0539	0.1204	0.0728	0.0697	0.2916	0.1220	0.0459
	MinerU2-VLM (Wang et al., 2024c)	0.0745	<b>0.0104</b>	0.0357	0.1276	0.0698	0.0652	0.1831	<b>0.0803</b>	0.0236
	Nanonets-OCR-s (Mandalm, 2025)	0.0551	0.0578	0.0606	0.0931	0.0834	0.0917	0.1965	0.1606	0.0395
	MonkeyOCR-pro-1.2B (Li et al., 2025)	0.0961	0.0354	0.0530	0.1110	0.0887	0.0494	0.0995	0.1686	0.0198
	MonkeyOCR-3B (Li et al., 2025)	0.0904	0.0362	0.0489	0.1072	0.0745	0.0475	0.0962	0.1165	0.0196
	dots.ocr (rednote, 2025)	<b>0.0290</b>	0.0231	0.0433	<u>0.0788</u>	<b>0.0467</b>	<u>0.0221</u>	0.0667	0.1116	<b>0.0076</b>
	MonkeyOCR-pro-3B (Li et al., 2025)	0.0879	0.0459	0.0517	0.1067	0.0726	0.0482	0.0937	0.1141	0.0211
	<b>MinerU2.5</b>	<u>0.0294</u>	0.0235	<u>0.0332</u>	<b>0.0499</b>	<u>0.0681</u>	0.0316	<b>0.0540</b>	0.1161	<u>0.0104</u>

Table 7: Document Parsing Performance in Text Edit Distance on OmniDocBench: evaluation using edit distance across 9 PDF page types.

Model	Edit Distance ↓		F1-score ↑		Precision ↑		Recall ↑		BLEU ↑		METEOR ↑	
	en	zh	en	zh	en	zh	en	zh	en	zh	en	zh
Mathpix (Mathpix, 2025)	0.064	0.223	0.930	0.919	<b>0.950</b>	0.952	0.911	0.889	0.901	0.593	0.924	0.768
PP-StructureV3 (Cui et al., 2025)	0.068	0.210	0.871	0.929	0.856	0.924	0.892	0.935	0.796	0.570	0.902	0.802
MinerU2-pipeline (Wang et al., 2024c)	0.099	0.225	0.663	0.919	0.635	0.908	0.703	0.934	0.504	0.571	0.670	0.810
PaddleOCR (Cui et al., 2025)	0.323	0.649	0.707	0.864	0.690	0.912	0.730	0.821	0.517	0.537	0.674	0.699
Gemini-2.5 Pro (Comanici et al., 2025)	0.080	0.204	0.922	0.927	0.940	0.959	0.906	0.898	0.877	0.690	0.921	0.862
GPT-4o (Achiam et al., 2023)	0.085	0.450	0.919	0.686	0.929	0.694	0.910	0.703	0.870	0.354	0.922	0.495
Qwen2.5-VL-72B (Bai et al., 2025)	0.093	0.140	0.923	0.940	0.936	0.956	0.912	0.926	0.879	0.798	0.924	0.876
InternVL3-76B (Zhu et al., 2025)	0.125	0.282	0.828	0.871	0.842	0.889	0.817	0.856	0.728	0.527	0.829	0.759
Qwen2-VL-7B (Wang et al., 2024d)	0.165	0.270	0.849	0.883	0.834	0.847	0.873	0.942	0.795	0.578	0.859	0.763
MiniCPM-V2.6-8B (Yao et al., 2024)	0.244	0.437	0.804	0.778	0.793	0.721	0.837	0.875	0.695	0.431	0.640	0.642
MinerU2-VLM (Wang et al., 2024c)	<u>0.048</u>	0.182	0.936	0.941	0.926	0.927	<u>0.947</u>	0.958	0.893	0.611	<b>0.950</b>	0.837
Ocean-OCR (Chen et al., 2025)	0.057	<b>0.062</b>	<u>0.937</u>	<u>0.962</u>	0.932	<u>0.956</u>	<b>0.956</b>	<b>0.974</b>	<u>0.906</u>	<b>0.912</b>	<u>0.945</u>	<b>0.916</b>
MonkeyOCR-pro-1.2B (Li et al., 2025)	0.064	0.190	0.929	0.934	0.918	0.925	0.944	0.948	0.884	0.699	0.941	0.850
SmolDocling (Nassar et al., 2025)	0.080	0.878	0.899	0.157	0.895	0.140	0.912	0.268	0.839	0.048	0.907	0.151
dots.ocr (rednote, 2025)	0.083	0.179	0.904	0.931	0.920	0.951	0.890	0.913	0.849	0.639	0.911	0.842
GOT (Wei et al., 2024)	0.084	0.117	0.895	0.928	0.891	0.934	0.906	0.929	0.835	0.805	0.874	0.848
<b>MinerU2.5</b>	<b>0.033</b>	<u>0.082</u>	<b>0.945</b>	<b>0.965</b>	<u>0.948</u>	<b>0.966</b>	0.942	<u>0.964</u>	<b>0.909</b>	<u>0.817</u>	<b>0.950</b>	<u>0.887</u>

Table 8: Evaluation results on Ocean-OCR bench on dense English (en) and Chinese (zh) OCR for document-level pages. Some model results are sourced from the OceanOCR official reports.

## B.1.2 Evaluation Results

MinerU2.5 demonstrates exceptional performance across all benchmarks, achieving state-of-the-art results in most metrics (Tables 1 and 7 to 9).

As shown in Table 1, MinerU2.5 achieves an overall score of 90.67 on OmniDocBench, outperforming the second-best model MonkeyOCR-pro-3B (Li et al., 2025) by 1.82 and dots.ocr (rednote, 2025) by 2.26 points. In text recognition tasks, MinerU2.5 achieves the lowest edit distance of 0.047, marginally better than dots.ocr at 0.048 and significantly outperforming PP-StructureV3 (Cui et al., 2025), which scores 0.073. For formula recognition, MinerU2.5 leads with a CDM score of 88.46, exceeding both Qwen2.5-VL-72B at 88.27 and MonkeyOCR-3B at 87.45. In table recognition tasks, MinerU2.5 achieves the highest TEDS score of 88.22 and TEDS-S score of 92.38. For reading order evaluation, it maintains the best edit distance of 0.044. The document-type specific results presented in Table 7 demonstrate that MinerU2.5 achieves best or second-best performance in 6 out of 9 categories. For textbooks, it delivers the best performance with an edit distance of 0.0499, substantially outperforming dots.ocr’s 0.0788. For newspapers, MinerU2.5 leads with a score of 0.0540, surpassing all competing models. In both financial reports and slides categories, MinerU2.5 achieves second-best performance with scores of 0.0104 and

Model	Overall	AR	OSM	TA	OS	HF	MC	LTT	Base
MinerU2-pipeline(Wang et al., 2024c)	55.6	61.8	13.5	60.9	17.3	<b>96.6</b>	59.0	39.1	96.6
Nanonets-OCR-s(Mandalm, 2025)	60.7	63.9	41.0	77.7	39.5	40.7	69.9	53.4	99.3
GPT-4o(Achiam et al., 2023)	63.2	44.1	37.6	69.1	40.9	94.2	68.9	54.1	96.7
MonkeyOCR-pro-1.2B(Li et al., 2025)	64.3	65.4	26.9	60.3	31.2	93.3	66.2	<u>81.7</u>	89.5
Qwen2.5-VL-72B(Bai et al., 2025)	64.8	<u>72.2</u>	<u>51.1</u>	67.3	38.6	73.6	68.3	<u>49.1</u>	98.3
MonkeyOCR-pro-3B(Li et al., 2025)	68.8	67.7	28.4	74.6	36.1	91.2	76.6	80.1	95.3
olmOCR(Poznanski et al., 2025)	71.8	63.9	41.0	72.9	<b>43.9</b>	<u>95.1</u>	77.3	81.2	98.9
dots.ocr(rednote, 2025)	73.6	66.3	35.8	<b>88.3</b>	40.9	94.1	<b>82.4</b>	81.2	<b>99.5</b>
<b>MinerU2.5</b>	<b>75.2</b>	<b>76.6</b>	<b>54.6</b>	<u>84.9</u>	33.7	<b>96.6</b>	<u>78.2</u>	<b>83.5</b>	93.7

Table 9: Evaluation results on olmOCR-bench grouped by document types, including arXiv Math(AR), Old Scans Math (OSM), Tables (TA), Old Scans (OS), Headers Footers (HF), Multi Column (MC) and Long Tiny Text (LTT). Results on AR and OSM are replaced with ExpRate, and other results are sourced from the official reports of olmOCR-bench and dots.ocr. The Overall Score (Overall) represents the average across all document types.

0.0294 respectively.

For the results of the Ocean-OCR benchmark presented in Table 8, MinerU2.5 demonstrates exceptional performance in dense OCR tasks. On English documents, it achieves the lowest edit distance of 0.033 and the highest F1-score of 0.945, accompanied by best-in-class BLEU and METEOR scores of 0.909 and 0.950 respectively. For Chinese documents, MinerU2.5 achieves the highest F1-score of 0.965 and Precision of 0.966, while maintaining strong BLEU and METEOR scores of 0.817 and 0.887 respectively.

The results of olmOCR-bench are shown in Table 9, where MinerU2.5 achieves an overall score of 75.2, surpassing dots.ocr’s 73.6 by 1.6 points. In the arXiv Math category, it leads with a score of 76.6, outperforming Qwen2.5-VL-72B (Bai et al., 2025)’s 72.2 by 4.4 points. For Old Scans Math, MinerU2.5 dominates with a score of 54.6, exceeding all other evaluated models. In the Long Tiny Text category, it achieves 83.5, surpassing MonkeyOCR-pro-1.2B (Li et al., 2025) which scores 81.7.

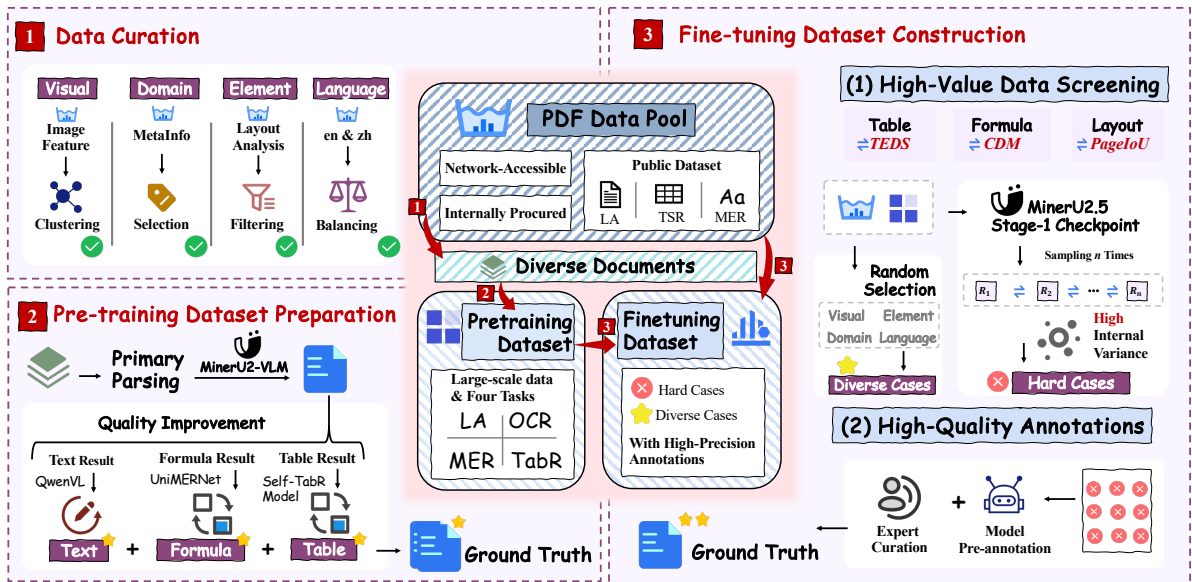


Figure 2: **Overview of the Data Engine.** Our data pipeline consists of three core stages. (1) **Data Curation:** We filter a massive, raw document pool to construct a diverse and balanced dataset based on layout, document type, element balance, and language. (2) **Pre-training Data Preparation:** We generate automated annotations for the curated data and then refine them using specialized, powerful models for text, tables, and formulas to ensure high quality. (3) **Fine-tuning Dataset Construction:** We employ our Iterative Mining via Inference Consistency (IMIC) strategy to automatically discover hard cases, which then undergo meticulous expert curation to create a high-quality SFT dataset.

## C Data Engine

The state-of-the-art performance of MinerU2.5 is underpinned by a systematic Data Engine designed to generate large-scale, high-quality training data with uniform annotation standards. This engine first establishes a vast and diverse foundation through rigorous data curation and refined automated annotation for pre-training. Building upon this foundation, we introduce our novel Iterative Mining via Inference Consistency (IMIC) strategy, which efficiently identifies complex “hard cases” for targeted human annotation. This multi-stage approach creates a virtuous cycle of improvement, progressively enhancing the model’s capabilities. The entire process is illustrated in Figure 2.

### C.1 Overall Workflow

#### C.1.1 Data Curation

Our process begins with a large-scale internal document pool comprising publicly available web data and commercially procured documents. While diverse, this raw pool suffers from a significant long-tail distribution. To mitigate this imbalance and enhance training robustness, we implement a rigorous curation process to build a balanced Chinese-English dataset with high diversity across multiple dimensions:

- **Layout Diversity:** We employ page-level image clustering to select exemplars from a wide spectrum of visual layouts and styles.
- **Document Type Diversity:** Using document metadata (e.g., discipline, tags), we perform stratified sampling to ensure a balanced representation of types such as academic papers, textbooks, reports, and presentations.
- **Element Balance:** A preliminary detection model helps ensure a balanced class distribution of key elements like titles, paragraphs, tables, formulas, and figures in the curated set.
- **Language Balance:** We filter the data to maintain a comparable volume of Chinese and English documents.

#### C.1.2 Pre-training Dataset Preparation

Initial annotations for the curated dataset are generated using our MinerU2-pipeline, establishing a baseline for subsequent refinement. To move beyond this baseline quality, we perform a multi-step refinement process using specialized, expert models for different content types:

- **Textual Content:** We leverage the powerful Qwen2.5-VL-72B-Instruct to verify and correct initial text recognition results on cropped text regions.
- **Formula Content:** Recognized formulas are substituted with higher-fidelity outputs from an in-house UniMERNET model, which we retrained on our extensive formula dataset to boost its accuracy.
- **Table Content:** All table structures are re-generated using an in-house, high-performance table parsing model.

This refinement workflow yields a high-quality pre-training dataset of image-annotation pairs, covering our four core tasks: layout analysis, text recognition, formula recognition, and table recognition.

#### C.1.3 Fine-tuning Dataset Construction

While pre-training ensures broad capabilities, the noise inherent in automated annotations creates a ceiling for model performance. To break through this ceiling, our fine-tuning strategy pivots to high-value, difficult examples. We designed an Iterative Mining via Inference Consistency (IMIC) strategy to automatically filter these hard cases from the large-scale data pool. To ensure annotation quality, these select samples are processed through an AI-assisted pipeline: they are first pre-annotated by a foundation model, such as Gemini-2.5-Pro for complex tables, and then meticulously reviewed and corrected by human experts<sup>3</sup>.

---

<sup>3</sup>Human review is augmented by our open-source QA tool, Dingo, which applies both rule-based and model-based checks.

The final Supervised Fine-Tuning (SFT) dataset combines these high-quality hard cases with a smaller, randomly sampled set of regular examples, equipping MinerU2.5 to excel in complex, real-world parsing scenarios.

## C.2 Task Reformulation and Enhancement

To move beyond the limitations of existing document analysis methods, we systematically reformulated the core tasks of layout analysis, formula recognition, and table recognition. This involved defining more robust standards, designing novel task paradigms, and introducing specialized metrics and representations.

### C.2.1 Layout Analysis

**A Unified Tagging System.** A fundamental challenge in layout analysis is the lack of a standardized tagging system. Existing datasets suffer from widespread inconsistencies in element definitions, granularity, and scope. To address this, we engineered a hierarchical and comprehensive tagging system by analyzing a vast corpus of documents. Our system is defined by three key principles:

- **Comprehensive Coverage:** It includes non-body content often ignored by others, such as headers, footers, and page numbers, which is critical for downstream applications like RAG.
- **Fine Granularity:** It decomposes complex elements. For instance, figures are sub-categorized into image, chart, and chemical\_structure, with distinct tags for their associated captions.
- **Semantic Distinction:** Visually distinct text blocks like code, algorithms, references, and lists are assigned their own categories to preserve crucial semantic information.

Table 10 presents a comparison with mainstream tagging systems, highlighting the superior coverage and granularity of our proposed system.

Category	MinerU2-pipeline	PaddleOCR	MinerU2.5
<b>Textual</b>	text	text, toc, abstract	text
	title	title, page_title	title
	×	×	phonetic
	image_caption	common_caption	image_caption
	image_footnote	common_footnote	image_footnote
	table_caption	common_caption	table_caption
	table_footnote	common_footnote	table_footnote
	×	code	code
	×	×	code_caption
	×	×	algorithm
×	ref_text, ref_block	reference	
×	×	list	
<b>Image</b>	image	image, seal, chart, molecular	image
<b>Table</b>	table	table	table
<b>Equation</b>	equation	equation	equation
	×	×	equation_block
<b>Page Margins</b>	×	header	header
	×	footer	footer
	×	aside_text	aside_text
	×	page_number	page_number
	×	page_footnote	page_footnote

Table 10: Comparison of category support across different OCR systems.

**An Enhanced Multi-Task Paradigm.** Traditional methods often treat layout analysis as a standard object detection task, which ignores element rotation and defers reading order prediction to downstream

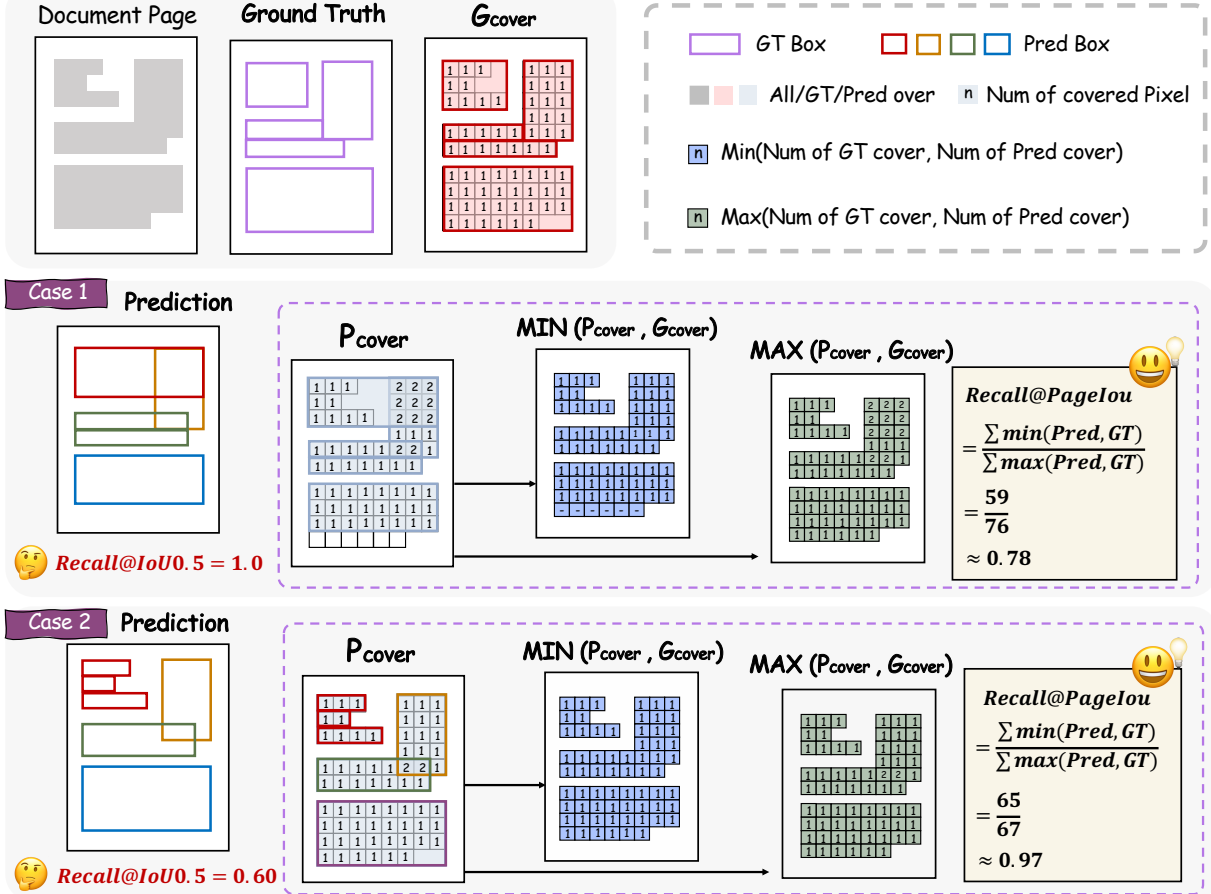


Figure 3: **Illustration of the proposed PageIoU metric.** Case 1 and Case 2 show that IoU-based recall may produce contradictory results compared with visual inspection, whereas PageIoU provides a page-level coverage score that aligns more closely with qualitative observations.

modules. This approach not only impairs the recognition of rotated elements but also increases system coupling. We propose an enhanced paradigm that redefines layout analysis as a multi-task problem. This paradigm simultaneously predicts four key attributes for each document element in a single inference pass: its **Position**, **Class**, **Rotation Angle**, and **Reading Order**. This integrated design effectively resolves the challenge of parsing rotated elements and streamlines the entire document analysis pipeline.

**PageIoU: A New Metric for Layout Quality.** Layout analysis is typically evaluated with object detection metrics like mAP, which rely on a fixed Intersection over Union (IoU) threshold. While effective for well-defined objects, this approach is ill-suited for document layouts where text block boundaries are often ambiguous. This can lead to a discrepancy where quantitative IoU-based scores do not align with qualitative visual assessment.

As illustrated in Figure 3, a prediction that coarsely covers a paragraph (Case 1) can achieve a perfect recall score ( $Recall@IoU0.5 = 1.0$ ), while a more accurate line-by-line prediction (Case 2) is penalized for not matching the paragraph-level ground truth, yielding a lower score ( $Recall@IoU0.5 = 0.6$ ). Visually, however, Case 2 is clearly a better fit.

To better evaluate document layout analysis, we introduce **PageIoU**, a page-level coverage metric that measures the spatial consistency between predicted layouts and ground-truth annotations. Let the predicted layout be

$$P = \{bbox_i \mid i = 1, 2, \dots, n\},$$

and the ground truth be

$$G = \{bbox_j \mid j = 1, 2, \dots, m\},$$

where each  $bbox$  denotes a bounding box on the page. We first compute coverage maps for both prediction

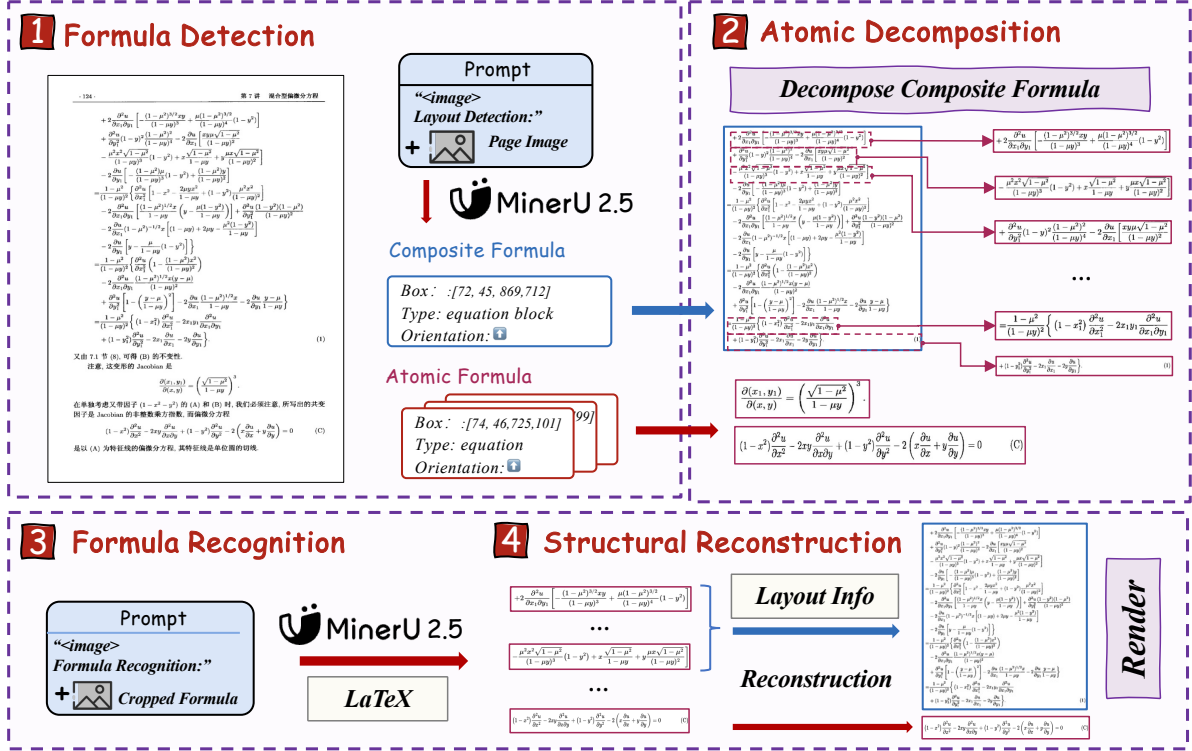


Figure 4: **The proposed ADR framework.** First, a compound formula is decomposed into atomic lines via layout analysis. Next, each line is individually recognized into LaTeX. Finally, the individual results are structurally recombined to produce the complete output.

and ground truth. For example, the ground-truth coverage map is defined as:

$$G_{cover} = \left\{ \sum_{j=1}^m 1_{p \in bbox_j} \mid p \in M \right\},$$

where  $p$  is a page pixel and  $M$  denotes the non-background region of the page. Similarly,  $P_{cover}$  can be obtained. Based on these, PageIoU is defined as:

$$\begin{aligned} \text{PageIoU}(P, G) &= \frac{|P_{cover} \cap G_{cover}|}{|P_{cover} \cup G_{cover}|} \\ &= \frac{\sum_{p \in M} \min\{P_{cover}(p), G_{cover}(p)\}}{\sum_{p \in M} \max\{P_{cover}(p), G_{cover}(p)\}}. \end{aligned} \quad (2)$$

Here,  $|\cdot|$  denotes the summation over all pixel values, while  $\cap$  and  $\cup$  correspond to the pixel-wise minimum and maximum of coverage counts, respectively. As shown in Figure 3, PageIoU aligns with human perception, scoring the qualitatively poor prediction 0.78 and the superior one 0.97.

## C.2.2 Formula Recognition

**Decoupling Atomic and Compound Formulas.** Existing models struggle with long or multi-line formulas, and VLMs are prone to severe structural hallucinations. We identify the root cause as the tendency to treat all formulas as monolithic entities, failing to account for internal complexity. To this end, MinerU2.5 introduces a "whole-part" decoupling philosophy, classifying formulas into two types based on their structural and semantic integrity:

- **Atomic Formulas:** The smallest, indivisible semantic units with a tight 2D topology (e.g., a single fraction, a matrix).

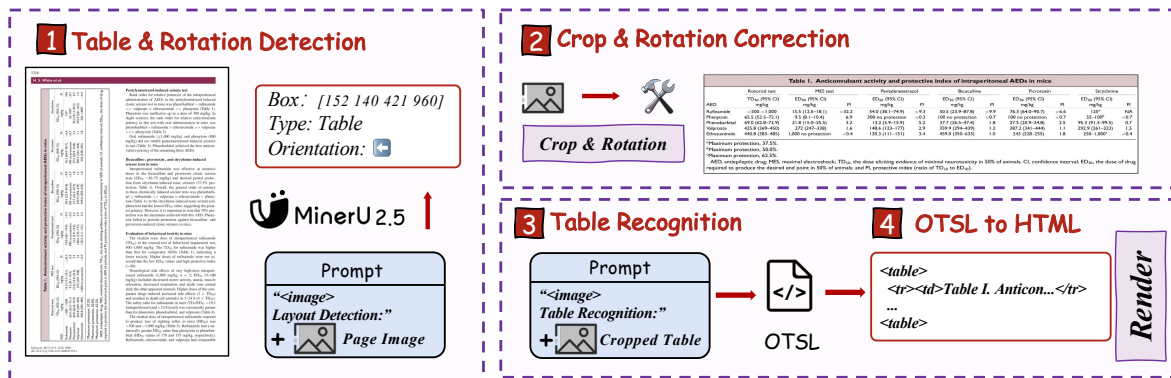


Figure 5: **The Table Recognition Pipeline.** The pipeline first detects a table and its rotation, then corrects its geometry. Next, the rectified image is recognized into the OTSL result, which is finally converted to standard HTML.

- **Compound Formulas:** An ordered set of atomic formulas composed vertically with specific alignment relationships (e.g., a multi-line derivation aligned at the equal signs).

**The Atomic Decomposition & Recombination (ADR) Framework.** To handle the complexity of compound formulas, we propose the ADR framework, which implements a multi-stage "divide and conquer" strategy. As illustrated in Figure 4, the ADR pipeline is powered by our versatile MinerU2.5 model, which acts as both a layout analyzer and a recognition engine, guided by task-specific prompts. The process begins with an initial layout analysis pass, where MinerU2.5, guided by a layout detection prompt, identifies and classifies all formula regions on the page as either atomic or compound. Next, in the decomposition stage, each identified compound formula is segmented into an ordered sequence of its constituent atomic formula lines, which are then cropped as individual images. In the third stage, these simple, semantically independent atomic formula images are fed back into the MinerU2.5 model. This time, using a formula recognition prompt, the model performs high-precision translation of each image into its corresponding LaTeX string. Finally, a lightweight recombination step uses the positional information from the initial layout pass to structurally reassemble the individual LaTeX strings into a single, coherent block, correctly formatting them within environments like `align`. This approach transforms a single, difficult recognition task into a series of simpler ones, ensuring both high-fidelity recognition of each component and the logical integrity of the overall structure.

### C.2.3 Table Recognition

**Overcoming Long-Sequence Dependencies.** A primary challenge in table recognition is parsing complex, long tables, especially for VLM-based approaches that target HTML. We attribute this difficulty to two inherent weaknesses of the HTML representation: (1) its complex, non-visual syntax must be learned implicitly by the model; and (2) its high token redundancy results in excessively long sequences, degrading performance on large tables. (The issue of rotated tables is effectively handled by our enhanced layout paradigm.)

**OTSL: An Optimized Table Structure Language.** To robustly handle complex tables, we propose a four-stage recognition pipeline, as depicted in Figure 5. The first two stages handle geometric normalization: the system detects the table’s bounding box and rotation angle, then corrects the image by cropping and rotating it to a canonical orientation. For the crucial third stage, table recognition, we leverage the Optimized Table-Structure Language (OTSL) (Lysak et al., 2023), an intermediate representation developed by IBM [citation, 2023]. We adopted OTSL for its significant advantages over HTML as a target for VLMs. Its minimalist design features a direct structural correspondence to a table’s visual 2D matrix, reducing the number of structural tokens from over 28 to just 5 and shortening the average sequence length by approximately 50%. This makes it a far more effective target for model generation. The final stage is a straightforward conversion from the OTSL output into standard HTML.

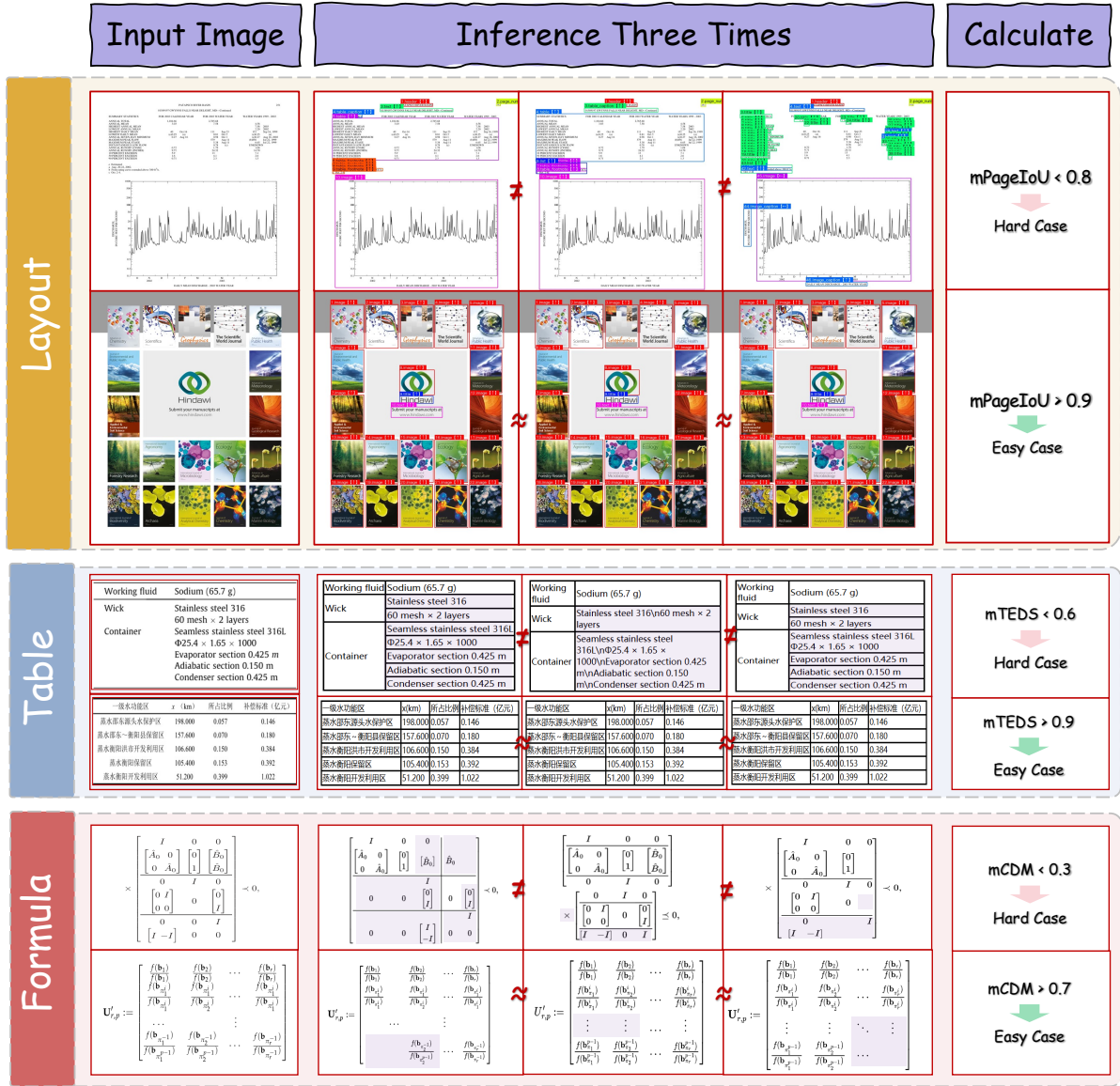


Figure 6: Illustration of the proposed **IMIC (Iterative Mining via Inference Consistency)** strategy. From top to bottom: (a) Layout analysis, (b) Table recognition, and (c) Formula recognition. For each task, the model performs multiple stochastic inference runs, and the pairwise consistency between outputs is calculated with task-specific metrics (PageIoU, TEDS, CDM). Samples with low consistency are automatically identified as hard cases and prioritized for manual annotation.

### C.3 Iterative Mining via Inference Consistency

To enable continuous model improvement and the efficient expansion of our high-quality training dataset, we introduce the IMIC (Iterative Mining via Inference Consistency) strategy. IMIC automatically identifies the most challenging samples—or "hard cases"—for the current model from a large corpus of unlabeled data. This allows us to direct limited human annotation efforts toward the data that offers the maximum value for model improvement.

The core principle of IMIC leverages the stochasticity inherent in model inference. For a given sample, if the model has learned its features robustly, multiple inference passes with stochastic sampling enabled should yield highly consistent outputs. Conversely, significant divergence across outputs suggests the sample lies near the model’s decision boundary—a 'hard case' where its predictions are uncertain. Such samples are the most valuable candidates for manual annotation, as they directly target the model’s specific weaknesses.

As illustrated in Figure 6, the implementation is tailored to each recognition task:

- **Layout analysis:** For full document pages, we perform multiple inference runs and measure consistency by calculating the pairwise PageIoU between the resulting layouts. Samples falling below a predefined similarity threshold are flagged as hard cases for precise manual annotation.
- **Formula Recognition:** For cropped formula images, consistency is assessed using the pairwise CDM (Wang et al., 2025a) across multiple outputs. Samples with low consistency are prioritized for manual correction.
- **Table Recognition:** For cropped table images, we use the TEDS (Tree-Edit-Distance-based Similarity) score to evaluate consistency across multiple recognized structures. Low-consistency samples are routed to the manual annotation workflow.

## D Prompt Details

Here, we provide a detailed description of the different prompts used during the two-stage inference of MinerU2.5, along with their corresponding output formats.

### D.1 Layout Detection

The layout detection output will include the relative coordinates, category, and rotation direction of each element in the document. Each element will be output in sequence, ensuring traceability for all layout data. The input image will be resized to a resolution of  $1036 \times 1036$ .

#### Output format:

- **Box Coordinates:**  $x_1, y_1, x_2, y_2$
- **Document Element Category:** title, text, image, etc.
- **Rotation Direction:** up, down, left, right

#### Example:

```
<|box_start|>100 200 300 400<|box_end|><|ref_start|>title<|ref_end|><|rotate_up|>  
<|box_start|>400 500 600 700<|box_end|><|ref_start|>text<|ref_end|><|rotate_up|>
```

### D.2 Text Recognition

The output will contain the recognized text results. The input image will retain its native resolution; however, the number of image tokens will be limited to the range of 4 to 2048. If this limit is exceeded, the image will be scaled accordingly.

#### Output format:

- **OCR Results:** The raw OCR output

#### Example:

The results of the analyses of the uncertainty of the field data and related assumptions are shown in Figs 13 and 14.

### D.3 Formula Recognition

Any formulas found in the image will be extracted and converted into LaTeX format. The input image will retain its native resolution; however, the number of image tokens will be limited to the range of 4 to 2048. If this limit is exceeded, the image will be scaled accordingly.

#### Output format:

- **LaTeX Format:** The LaTeX representation of the formula

#### Example:



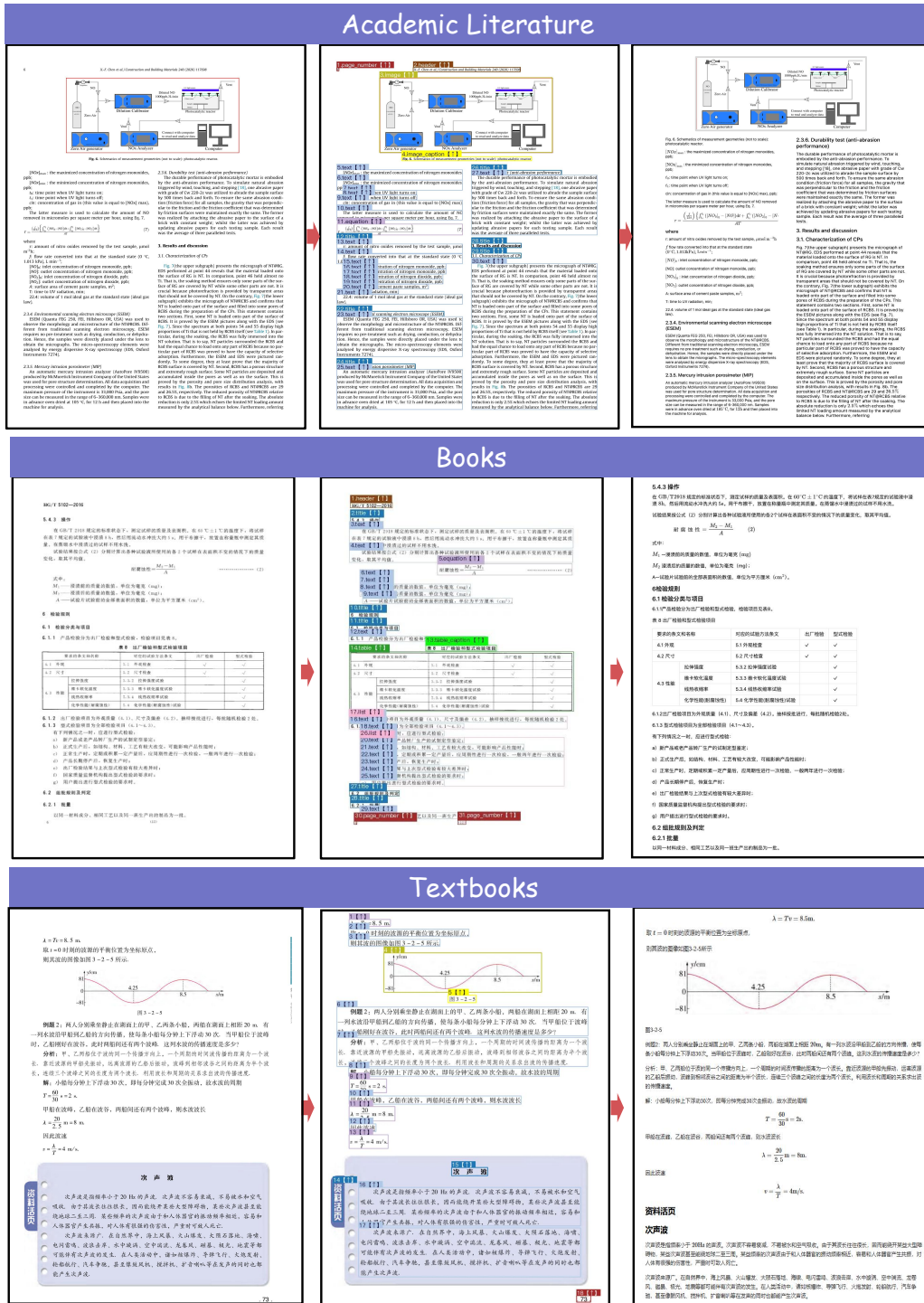


Figure 7: The Layout and rendered markdown output for Academic literature, Books, Textbooks.

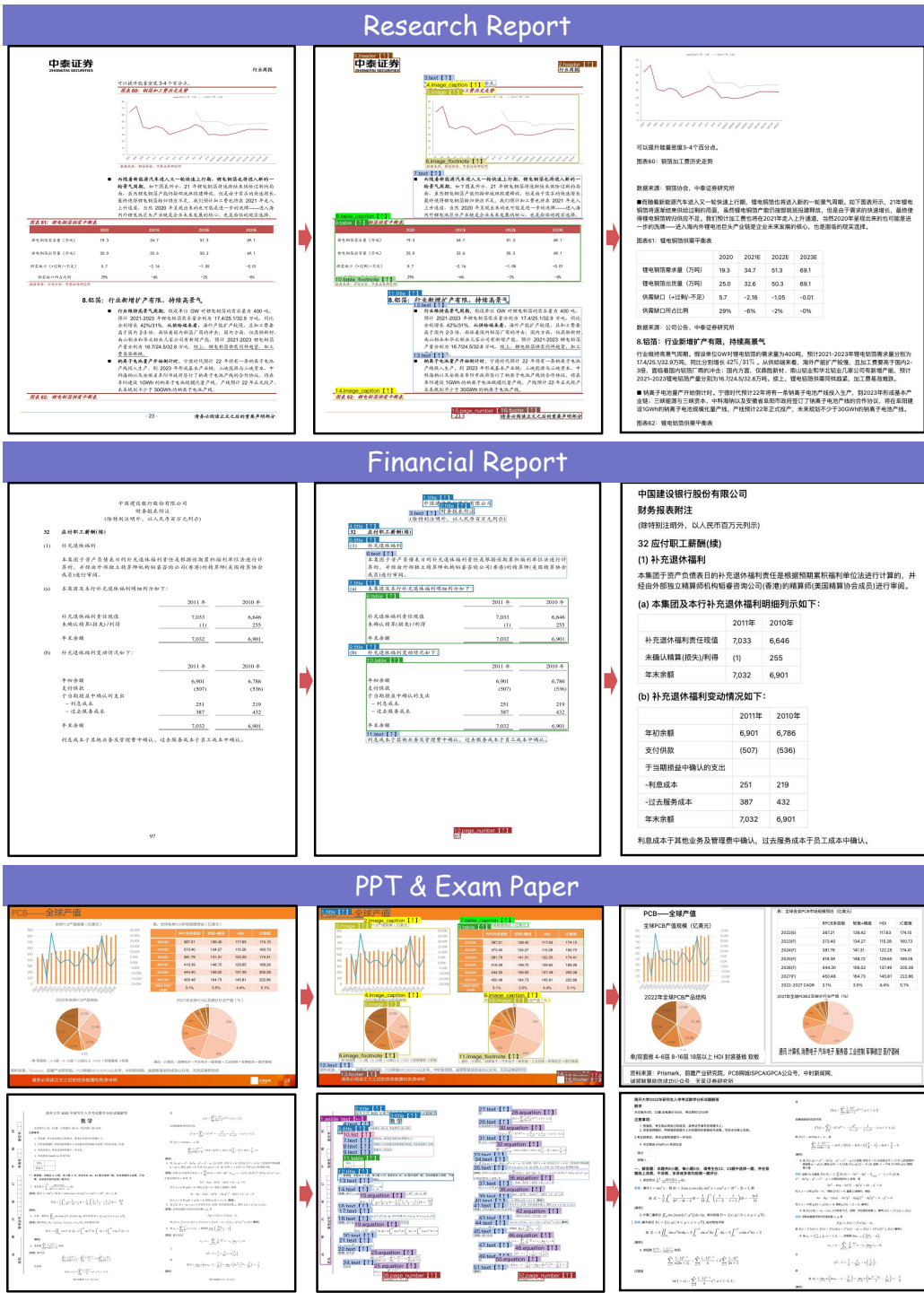
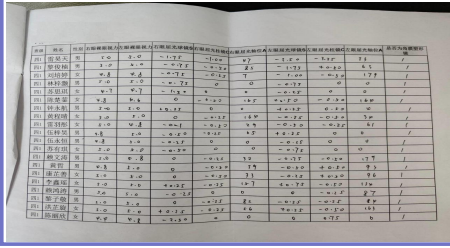


Figure 8: The Layout and rendered markdown output for Research Report, Financial Report, Slides and Exam Paper.



### Photograph of the Table



姓名	性别	右眼裸视力	左眼裸视力	右眼矫正视力	左眼矫正视力	右眼远视储备量	左眼远视储备量	右眼轴长	左眼轴长	右眼角膜曲率	左眼角膜曲率	是否散光	是否弱视
田博	男	5.0	5.0	1.75	1.00	47	3.50	-3.25	33	0	0	0	0
田博	男	5.0	5.0	0.75	0.50	45	1.75	+0.50	33	0	0	0	0
田博	男	4.8	4.8	0.75	0.25	47	1.00	0.50	179	0	0	0	0
田博	男	5.0	5.0	0.75	0	0	0.75	0	0	0	0	0	0
田博	男	4.7	4.7	1.50	0	0	0.25	0	0	0	0	0	0
田博	男	4.8	4.8	0	0.50	165	4.50	0.50	164	0	0	0	0
田博	男	5.0	5.0	+0.25	0	0	+0.25	0.50	0	0	0	0	0
田博	男	5.0	5.0	0	0.25	164	0.25	0.50	30	0	0	0	0
田博	男	5.0	4.8	-0.11	0.50	29	0.50	+0.25	61	0	0	0	0
田博	男	4.8	5.0	0.50	0.25	65	+0.25	0	0	0	0	0	0
田博	男	4.8	5.0	0.25	0	0	0.25	0	0	0	0	0	0
田博	男	5.0	5.0	0.50	0	0	0.75	0	0	0	0	0	0
田博	男	5.0	4.8	0	0.25	30	0.75	0.50	179	0	0	0	0
田博	男	4.8	5.0	0	0.50	39	0.50	+0.50	93	0	0	0	0
田博	男	5.0	5.0	0	0.50	33	0.25	+0.50	96	0	0	0	0
田博	男	5.0	5.0	+0.25	0.25	157	+0.75	0.50	134	0	0	0	0
田博	男	5.0	5.0	0.25	0	0	0.25	0	87	0	0	0	0
田博	男	5.0	5.0	0	0.25	82	0.25	-0.25	84	0	0	0	0
田博	男	5.0	5.0	+0.25	0.25	56	+0.25	0.50	163	0	0	0	0
田博	男	4.4	4.8	2.50	0	0	0	0.75	0	0	0	0	0

### Handwritten Table

省份	简称	省会
山西省	晋	太原
内蒙古自治区	蒙	呼和浩特
辽宁省	辽	沈阳
吉林省	吉	长春
黑龙江省	黑	哈尔滨
上海市	沪	上海
江苏省	苏	南京
浙江省	浙	杭州
安徽省	皖	长沙
广东省	粤	广州
广西壮族自治区	桂	南宁
海南省	琼	海口
重庆市	渝	重庆
四川省	川	成都
贵州省	贵	贵阳
云南省	云	昆明
西藏自治区	藏	拉萨
陕西省	陕	西安
甘肃省	陇	兰州
安徽省	皖	合肥
福建省	闽	福州
江西省	赣	南昌
山东省	鲁	济南
河南省	豫	郑州

### Colorful Table

日期	基金名称	基金类型	T15	T10	T15
2015-06-12	华夏	股票	724.67	805.37	1950.11
2015-09-15	华夏	股票	391.78	419.74	891.21
2015-12-22	华夏	股票	27.46	31.48	29.89
2016-01-28	华夏	股票	295.92	142.84	89.04
2017-11-13	华夏	股票	222.65	191.97	284.29
2018-10-18	华夏	股票	568.83	1530.78	2009.05
2018-04-10	华夏	股票	662.43	1465.46	1517.01
2020-07-13	华夏	股票	8064.47	7902.07	6790.54
2021-03-18	华夏	股票	2009.72	2852.76	2849.22
2021-12-13	华夏	股票	5341.16	3579.57	3065.12
2022-04-26	华夏	股票	2197.51	2264.42	3399.59
2022-07-04	华夏	股票	2340.62	3688.16	4627.22
2022-10-31	华夏	股票	6162.69	6618.07	5339.34

### Large Table

基金代码	基金名称	基金类型	成立日期	到期日期	基金管理人	基金规模	基金净值	基金收益
014354.OF	东方夜月九个月持有A	偏债混合型	2022-05-09	2022-07-08	东方基金管理	2022-06-08	2022-07-08	
014066.OF	长信睿利收益一年持有A	偏债混合型	2022-04-13	2022-07-12	长信基金管理	2022-04-13	2022-07-12	
014007.OF	国泰稳健收益一年持有	混合型FOF	2022-04-18	2022-07-15	国泰基金管理	2022-04-18	2022-07-15	
014070.OF	汇添富均衡增长三个月持有A	混合型FOF	2022-04-18	2022-07-15	汇添富基金管理	2022-04-18	2022-07-15	
011210.OF	新鑫中证海外中国互联网50ETF	被动指数型	2022-04-20	2022-07-19	新鑫基金管理	2022-04-20	2022-07-19	
011629.OF	华夏核心优势	偏股混合型	2022-04-20	2022-07-19	华夏基金管理	2022-04-20	2022-07-19	
015699.OF	平安均衡成长2年持有A	偏股混合型	2022-05-06	2022-07-20	平安基金管理	2022-05-06	2022-07-20	
014072.OF	汇安裕阳A	中长期纯债型	2022-04-21	2022-07-20	汇安基金管理	2022-04-21	2022-07-20	
014595.OF	西部利得聚优一年持有	偏债混合型	2022-04-02	2022-07-22	西部利得基金管理	2022-04-02	2022-07-22	
014290.OF	南方晨星一年定开	中长期纯债型	2022-04-25	2022-07-22	南方基金管理	2022-04-25	2022-07-22	
501153.OF	富国中证500ESG基准ETF	被动指数型	2022-04-26	2022-07-25	富国基金管理	2022-04-26	2022-07-25	
014462.OF	光大保德信汇佳A	偏股混合型	2022-05-25	2022-07-26	光大保德信基金管理	2022-05-25	2022-07-26	
014313.OF	鹏华创新增长一年持有A	偏股混合型	2022-04-27	2022-07-26	鹏华基金管理	2022-04-27	2022-07-26	
513223.OF	招商中证全球中国互联网ETF	国际(001)股票型	2022-05-27	2022-07-27	招商基金管理	2022-05-27	2022-07-27	
014408.OF	创金启信优选产业趋势一年封闭A	偏股混合型	2022-04-29	2022-07-28	创金启信基金管理	2022-04-29	2022-07-28	
011632.OF	前海联合季享价值A	偏股混合型	2022-04-29	2022-07-28	前海联合基金管理	2022-04-29	2022-07-28	
011964.OF	国泰中证500科技A	被动指数型	2022-04-29	2022-07-28	国泰基金管理	2022-04-29	2022-07-28	
013228.OF	天弘中证新能源指数增强A	偏股混合型	2022-04-29	2022-07-28	天弘基金管理	2022-04-29	2022-07-28	
014288.OF	淳厚惠丰A	中长期纯债型	2022-04-29	2022-07-28	淳厚基金管理	2022-04-29	2022-07-28	
014343.OF	泰康丰盈纯债一年定开	中长期纯债型	2022-04-29	2022-07-28	泰康基金管理	2022-04-29	2022-07-28	
501220.OF	国泰行业轮动一年封闭A	股票型FOF	2022-04-29	2022-07-28	国泰基金管理	2022-04-29	2022-07-28	
513293.OF	汇添富纳斯达克生物科技ETF	国际(001)股票型	2022-05-27	2022-07-29	汇添富基金管理	2022-05-27	2022-07-29	
014311.OF	大成优质精选A	偏股混合型	2022-05-06	2022-08-05	大成基金管理	2022-05-06	2022-08-05	
159623.OF	博时中证成渝地区双城经济圈ETF	被动指数型	2022-05-10	2022-08-09	博时基金管理	2022-05-10	2022-08-09	
015578.OF	南方宝祥A	偏债混合型	2022-05-11	2022-08-10	南方基金管理	2022-05-11	2022-08-10	
014488.OF	汇添富淳享一年定开A	中长期纯债型	2022-05-12	2022-08-11	汇添富基金管理	2022-05-12	2022-08-11	
014813.OF	南方元弘六个月持有A	中长期纯债型	2022-05-18	2022-08-12	南方基金管理	2022-05-18	2022-08-12	
014384.OF	国投深证成长一年定开	中长期纯债型	2022-05-13	2022-08-12	国投基金管理	2022-05-13	2022-08-12	
014443.OF	汇丰晋信非债A	中长期纯债型	2022-05-13	2022-08-12	汇丰晋信基金管理	2022-05-13	2022-08-12	
014374.OF	景顺长城丰泰非养老目标三年持有	混合型FOF	2022-05-13	2022-08-12	景顺长城基金管理	2022-05-13	2022-08-12	
014391.OF	华安添信	中长期纯债型	2022-05-12	2022-08-12	华安基金管理	2022-05-12	2022-08-12	
014387.OF	光大保德信得利一年定开	中长期纯债型	2022-05-16	2022-08-15	光大保德信基金管理	2022-05-16	2022-08-15	
014388.OF	渤海汇金兴债一年定开	中长期纯债型	2022-05-16	2022-08-15	渤海汇金基金管理	2022-05-16	2022-08-15	
014482.OF	华夏康盛可持续一年持有A	偏股混合型	2022-05-18	2022-08-17	华夏基金管理	2022-05-18	2022-08-17	
014710.OF	平安添利纯债A	中长期纯债型	2022-05-24	2022-08-19	平安基金管理	2022-05-24	2022-08-19	
014823.OF	国泰睿利科技1个月滚动A	偏股混合型	2022-05-20	2022-08-19	国泰基金管理	2022-05-20	2022-08-19	
014474.OF	中安安债一年定开	中长期纯债型	2022-05-20	2022-08-19	中安基金管理	2022-05-20	2022-08-19	
014452.OF	天弘惠享一年定开	混合型FOF	2022-05-24	2022-08-23	天弘基金管理	2022-05-24	2022-08-23	
014510.OF	国信永丰兴一年定开	中长期纯债型	2022-05-24	2022-08-23	国信永丰基金管理	2022-05-24	2022-08-23	
159621.OF	国泰MSCI中国ESG通用ETF	被动指数型	2022-05-24	2022-08-23	国泰基金管理	2022-05-24	2022-08-23	
014968.OF	中信建投投资3个月定开A	中长期纯债型	2022-05-24	2022-08-23	中信建投基金管理	2022-05-24	2022-08-23	

Figure 10: The rendered outputs for various types of Tables.



### Formula with Background

$\mathcal{M}_{\lambda, \lambda_1}^{\text{cont}}(s, t, s_{12}, \Omega^H)$   
 $\mathcal{M}_{\lambda, \lambda_1}^{\text{cont}}(s, t, s_{12}, \Omega^H)$

---

$d(j) = \begin{cases} d(j) & d(j) - d(j-1) \leq T_{th} \\ d(j-1) & d(j) - d(j-1) > T_{th} \end{cases}$   
 $\bar{d}(j) = (d(j) - \min(d_n)) / (\max(d_n) - \min(d_n))$  (1)

---

$d(j) = \begin{cases} d(j) & d(j) - d(j-1) \leq T_{th} \\ d(j-1) & d(j) - d(j-1) > T_{th} \end{cases}$  (1)  
 $\bar{d}(j) = (d(j) - \min(d_n)) / (\max(d_n) - \min(d_n))$  (1)

---

$[CO]_d(P_T - P_d)(1 + h_d) = P_S$  (126)  
 $[CO]_d(P_T - P_d)(1 + h_d) = P_S$  (126)

### Handwritten Formula

解: 令  $t = x - \frac{1}{x}$  则  $x^2 + \frac{1}{x^2} = t^2 + 2 = x^2(t+2)$   
 原式 =  $x^2(x^2 - 4x + 2 + \frac{1}{x^2} + \frac{1}{x^2}) = [x(t+2)]^2$   
 =  $x^2[x^2 + \frac{1}{x^2} - 4(x - \frac{1}{x}) + 2] = [x^2 - 2x - 1]^2$   
 =  $x^2(t^2 + 2 - 4t + 2) = [x^2 - 2x - 1]^2$

---

解: 令  $t = x - \frac{1}{x}$  则  $x^2 + \frac{1}{x^2} = t^2 + 2 = x^2(t+2)$   
 原式 =  $x^2(x^2 - 4x + 2 + \frac{1}{x^2} + \frac{1}{x^2}) = [x(t+2)]^2$   
 =  $x^2[x^2 + \frac{1}{x^2} - 4(x - \frac{1}{x}) + 2] = [x^2 - 2x - 1]^2$   
 =  $x^2(t^2 + 2 - 4t + 2) = [x^2 - 2x - 1]^2$

---

$\frac{d}{dt} \left( \frac{\partial V}{\partial t} + \sum_{i=1}^n \frac{\partial V}{\partial x_i} f_i(x_i, t) \right) \leq 0$   
 $\frac{dV}{dt} = \frac{\partial V(t, 0)}{\partial t} + \sum_{i=1}^n \frac{\partial V}{\partial x_i} f_i(x_i, t) \leq 0$

---

②  $R_1 = (R_1 \cap R_1) \cap (R_2 \cap R_2) \cap (R_3 \cap R_3)$   
 $(x, y) \in R_1 \cap R_2 \Rightarrow (x, z) \in R_1 \wedge (x, y) \in R_2 \Rightarrow (x, z) \in R_2 \wedge$   
 $\wedge (z, y) \in R_2 \Rightarrow (x, z) \in R_2 \wedge (z, y) \in R_2 \wedge (x, y) \in R_2 \Rightarrow (x, y) \in R_2$   
 $\Rightarrow ((x, z) \in R_2 \wedge (z, y) \in R_2) \wedge (x, y) \in R_2 \Rightarrow (x, y) \in R_2$   
 $\wedge (x, y) \in R_2 \Rightarrow (x, y) \in R_2 \cap (R_2 \cap R_2) = \text{sgoempred}$

### Blurred Formula

$A \begin{pmatrix} 2 \\ 2 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & \frac{1}{2} \end{pmatrix} \begin{pmatrix} 2 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 - \sqrt{3} \\ 1 + \sqrt{3} \end{pmatrix}$   
 $A \begin{pmatrix} 2 \\ 2 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & \frac{1}{2} \end{pmatrix} \begin{pmatrix} 2 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 - \sqrt{3} \\ 1 + \sqrt{3} \end{pmatrix}$

---

$\Leftrightarrow (a, \eta_j) = \sum_{i=1}^n k_i(\eta_i, \eta_j), j = 1, 2, \dots, m$   
 $\Leftrightarrow (a, \eta_i) = \sum_{i=1}^n k_i(\eta_i, \eta_i), j = 1, 2, \dots, m$

---

$(A^T B^T)(\bar{e}_1) = A^T B^T(\bar{e}_1) = A^T(-\bar{e}_1) = A^T(\bar{e}_1) = \bar{e}_1$   
 $(A^2 B^2)(\bar{e}_1) = A^2 B^2(\bar{e}_1) = A^2(-\bar{e}_1) = A^2(\bar{e}_1) = \bar{e}_1$

---

$\int_0^{\pi/2} 9 \cos^8(\theta/4) d\theta = (9/2) \int_0^{\pi/2} (1 + \cos(\theta/2))^4 d\theta / 16 = (9/32) \int_0^{\pi/2} (1 + 4 \cos(\theta/2) +$   
 $\int_0^{\pi/2} 9 \cos^8(\theta/4) d\theta = (9/2) \int_0^{\pi/2} (1 + \cos(\theta/2))^4 d\theta / 16 = (9/32) \int_0^{\pi/2} (1 + 4 \cos(\theta/2) +$

### Formula with Chinese

**1** 【答案】(1) 原式 =  $-2 \times \log_5 5 \times (-3) \times \log_5 2 \times (-2) \times \log_5 3 = -12$   
 (2) 原式 =  $(\log_2 5 - \frac{1}{2} \log_5 5) (\log_5 2 - \frac{1}{2} \log_5 2) = \frac{1}{4}$   
 (3) 原式 =  $2 - \sqrt{3} + 2 + \sqrt{3} = 4$   
 (4) 原式 =  $\lg^2 5 - (\lg 2 - 1)^2 + 1 = 1$

---

**2** 【答案】(1)  $\frac{3pq}{1+3pq}$ ; (2)  $\frac{a+b}{2-a}$

---

**1** 【答案】(1) 原式 =  $-2 \times \log_5 5 \times (-3) \times \log_5 2 \times (-2) \times \log_5 3 = -12$   
 (2) 原式 =  $(\log_2 5 - \frac{1}{2} \log_5 5) (\log_5 2 - \frac{1}{2} \log_5 2) = \frac{1}{4}$   
 (3) 原式 =  $2 - \sqrt{3} + 2 + \sqrt{3} = 4$   
 (4) 原式 =  $\lg^2 5 - (\lg 2 - 1)^2 + 1 = 1$

---

**2** 【答案】(1)  $\frac{3pq}{1+3pq}$ ; (2)  $\frac{a+b}{2-a}$

---

$L_1 = \left\{ \begin{array}{l} f(t) = \lim_{n \rightarrow \infty} e^{-tA_n} \prod_{k=1}^n f_k \left( \frac{t}{B_n} \right), f(t) \text{ 是 c.} \\ f(t), \{B_n\} \text{ 是正数列, } \{A_n\} \text{ 是实数列,} \\ \left\{ f_k \left( \frac{t}{B_n} \right), k=1, 2, \dots, n \right\} \text{ 是 u.s.n. 体系.} \end{array} \right\}$

---

$L_2 = \left\{ \begin{array}{l} f(t) = \lim_{n \rightarrow \infty} e^{-tA_n} \prod_{k=1}^n f_k \left( \frac{t}{B_n} \right), f(t) \text{ 是 c.} \\ f(t), \{B_n\} \text{ 是正数列, } \{A_n\} \text{ 是实数列,} \\ \left\{ f_k \left( \frac{t}{B_n} \right), k=1, 2, \dots, n \right\} \text{ 是 u.a.n. 体系.} \end{array} \right\}$

---

(在系统中的平均等待时间) =  $\frac{W_s \text{ 的总和}}{\text{总的顾客数}} = \frac{689.44}{20} = 34.47$   
 (在系统中的平均等待时间) =  $\frac{W_s \text{ 的总和}}{\text{总的顾客数}} = \frac{689.44}{20} = 34.47$

### Matrices

Where  $M = \begin{pmatrix} -\lambda & 1 & 1 & \dots & 1 & 1 & 1 & \dots & 1 & p^2 - 1 \\ -\lambda & 1 & -\lambda & \dots & 1 & 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & 1 & -\lambda & \dots & 1 & 1 & \dots & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & -\lambda & 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & 1 & \dots & 1 & -\lambda & 1 & \dots & 1 & p^2 - 1 \\ 1 & 1 & 1 & \dots & 1 & 1 & -\lambda & \dots & 1 & p^2 - 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 1 & 1 & 1 & \dots & -\lambda & p^2 - 1 \\ 1 & 0 & 0 & \dots & 0 & 1 & 1 & \dots & 1 & -\lambda + p^2 - 2 \end{pmatrix}, N = \begin{pmatrix} 1 & \dots & 1 \\ 0 & \dots & 0 \\ 0 & \dots & 0 \\ \vdots & \vdots & \vdots \\ 1 & \dots & 1 \\ 1 & \dots & 1 \\ \vdots & \vdots & \vdots \\ 1 & \dots & 1 \\ 1 & \dots & 1 \end{pmatrix}$

---

Where  $M = \begin{pmatrix} -\lambda & 1 & 1 & \dots & 1 & 1 & 1 & \dots & 1 & p^2 - 1 \\ -\lambda & 1 & -\lambda & \dots & 1 & 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & 1 & -\lambda & \dots & 1 & 1 & \dots & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & -\lambda & 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & 1 & \dots & 1 & -\lambda & 1 & \dots & 1 & p^2 - 1 \\ 1 & 1 & 1 & \dots & 1 & 1 & -\lambda & \dots & 1 & p^2 - 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \dots & 1 & 1 & 1 & \dots & -\lambda & p^2 - 1 \\ 1 & 0 & 0 & \dots & 0 & 1 & 1 & \dots & 1 & -\lambda + p^2 - 2 \end{pmatrix}, N = \begin{pmatrix} 1 & \dots & 1 \\ 0 & \dots & 0 \\ 0 & \dots & 0 \\ \vdots & \vdots & \vdots \\ 1 & \dots & 1 \\ 1 & \dots & 1 \\ \vdots & \vdots & \vdots \\ 1 & \dots & 1 \\ 1 & \dots & 1 \end{pmatrix}$

---

$\begin{pmatrix} E_1 & 0 & \dots & 0 \\ 0 & E_2 & \dots & 0 \\ 0 & 0 & E_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \\ 0 & \dots & 0 & E_m \end{pmatrix} \begin{pmatrix} d_{11} \\ d_{21} \\ d_{31} \\ \vdots \\ d_{m1} \end{pmatrix} \begin{pmatrix} y_{11} \\ y_{21} \\ y_{31} \\ \vdots \\ y_{m1} \end{pmatrix} = \begin{pmatrix} C_1 & 0 & \dots & 0 \\ 0 & C_2 & \dots & 0 \\ 0 & 0 & C_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \\ 0 & \dots & 0 & C_m \end{pmatrix} \begin{pmatrix} x_{11} \\ x_{21} \\ x_{31} \\ \vdots \\ x_{m1} \end{pmatrix}$

---

$\begin{pmatrix} E_1 & 0 & \dots & 0 \\ 0 & E_2 & \dots & 0 \\ 0 & 0 & E_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \\ 0 & \dots & 0 & E_m \end{pmatrix} \begin{pmatrix} d_{11} \\ d_{21} \\ d_{31} \\ \vdots \\ d_{m1} \end{pmatrix} \begin{pmatrix} y_{11} \\ y_{21} \\ y_{31} \\ \vdots \\ y_{m1} \end{pmatrix} = \begin{pmatrix} C_1 & 0 & \dots & 0 \\ 0 & C_2 & \dots & 0 \\ 0 & 0 & C_3 & \dots \\ \vdots & \vdots & \vdots & \ddots \\ 0 & \dots & 0 & C_m \end{pmatrix} \begin{pmatrix} x_{11} \\ x_{21} \\ x_{31} \\ \vdots \\ x_{m1} \end{pmatrix}$

---

$H_c = \begin{pmatrix} 1 & \frac{1}{3} & \dots & \frac{1}{k-1} & 0 & \dots & 0 \\ \frac{1}{3} & \frac{1}{5} & \dots & \frac{1}{k+1} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{k-1} & \frac{1}{k+1} & \dots & \frac{1}{2k-3} & 0 & \dots & 0 \\ 0 & \dots & \dots & 0 & -\frac{1}{3} & \dots & -\frac{1}{k-1} \\ 0 & \dots & \dots & 0 & -\frac{1}{5} & \dots & -\frac{1}{k+3} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & -\frac{1}{k+1} & \dots & -\frac{1}{2k-1} \end{pmatrix}$

---

$H_c = \begin{pmatrix} 1 & \frac{1}{3} & \dots & \frac{1}{k-1} & 0 & \dots & 0 \\ \frac{1}{3} & \frac{1}{5} & \dots & \frac{1}{k+1} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{k-1} & \frac{1}{k+1} & \dots & \frac{1}{2k-3} & 0 & \dots & 0 \\ 0 & \dots & \dots & 0 & -\frac{1}{3} & \dots & -\frac{1}{k-1} \\ 0 & \dots & \dots & 0 & -\frac{1}{5} & \dots & -\frac{1}{k+3} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & -\frac{1}{k+1} & \dots & -\frac{1}{2k-1} \end{pmatrix}$

Figure 12: The rendered outputs for various types of Formulas.

