

HunyuanOCR Technical Report

Tencent Hunyuan Vision Team

🤗 <https://huggingface.co/tencent/HunyuanOCR>
 🐾 <https://github.com/Tencent-Hunyuan/HunyuanOCR>

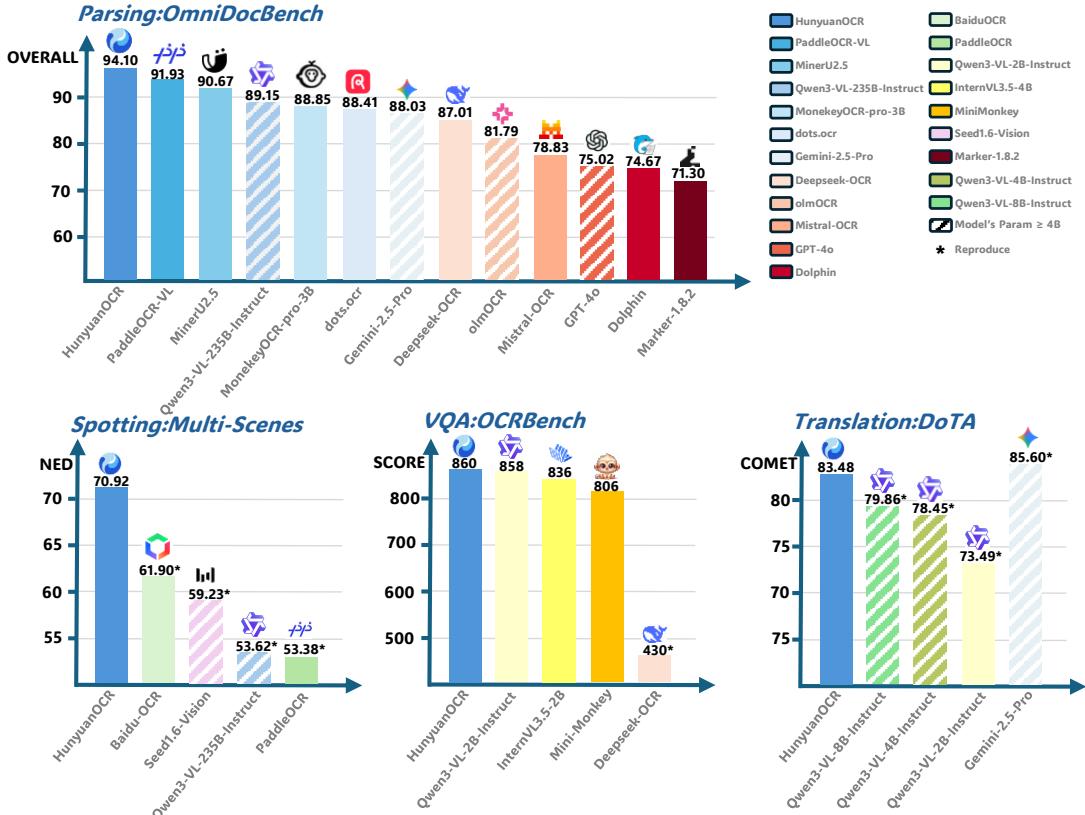


Figure 1: Performance comparison of HunyuanOCR and other SOTA models.

Abstract

This paper presents HunyuanOCR, a commercial-grade, open-source, and lightweight (1B parameters) Vision-Language Model (VLM) dedicated to OCR tasks. The architecture comprises a Native Vision Transformer (ViT) and a lightweight LLM connected via an MLP adapter. HunyuanOCR demonstrates superior performance, outperforming commercial APIs, traditional pipelines, and larger models (e.g., Qwen3-VL-4B). Specifically, it surpasses current public solutions in perception tasks (Text Spotting, Parsing) and excels in semantic tasks (IE, Text Image Translation), securing first place in the ICDAR 2025 DIMT Challenge (Small Model Track). Furthermore, it achieves state-of-the-art (SOTA) results on OCRBench among VLMs with fewer than 3B parameters.

HunyuanOCR achieves breakthroughs in three key aspects: 1) Unifying Versatility and Efficiency: We implement comprehensive support for core capabilities, including spotting, parsing, IE, VQA, and translation within a lightweight framework. This addresses the limitations of narrow “OCR expert models” and inefficient “General VLMs”. 2) Streamlined End-to-End Architecture: Adopting a pure end-to-end paradigm eliminates dependencies on pre-processing modules (e.g., layout analysis). This fundamentally resolves error propagation common in traditional pipelines and simplifies system deployment. 3) Data-Driven and RL Strategies: We confirm the critical role of high-quality

data and, for the first time in the industry, demonstrate that Reinforcement Learning (RL) strategies yield significant performance gains in OCR tasks.

HunyuanOCR is officially open-sourced on HuggingFace. We also provide a high-performance deployment solution based on vLLM, placing its production efficiency in the top tier. We hope this model will advance frontier research and provide a solid foundation for industrial applications.

1 Introduction

Modern Optical Character Recognition (OCR) Long et al. (2021); Zhang et al. (2024a) is a fundamental technology of artificial intelligence that continues to play an important role in promoting digitalization and industrial automation. Traditionally, OCR has mainly focused on extracting text from scanned document images and converting it into machine-readable data. In recent years, with the rapid development of deep learning and multimodal large language model technologies Zhou et al. (2017); Shi et al. (2017); Yin et al. (2024); Liu et al. (2024b), advanced OCR systems have broken through the limitations of scanned documents, now handling diverse layouts, casually captured images, and multilingual as well as handwritten text. Simultaneously, the scope of OCR tasks has expanded to include more challenging capabilities such as complex document parsing, end-to-end information extraction, text-centric visual question answering, and text image translation. Driven by technological innovation, the applications of intelligent OCR have permeated various aspects of industry and everyday life. For example, in office and educational settings Adeshola & Adepoju (2024), OCR enables functions such as literature translation and subject-specific tutoring. In the healthcare field Wang et al. (2025a), OCR facilitates the digital archiving of medical records and correlation analysis, supporting the provision of valuable treatment and health management advice to patients. Even more significantly, OCR systems fill a critical gap in acquiring high-quality corpora for Large Language Models, acting as an essential instrument for unlocking the content of specialized books and historical archives Zhang et al. (2024b).

To address diverse application requirements, the industry has long adopted pipeline-based frameworks, including PaddleOCR Cui et al. (2025b), EasyOCR JaideAI (2020), and MOCR Kuang et al. (2021). These approaches construct a sequential processing pipeline by integrating multiple compact expert models, offering benefits such as high modularity and the capacity for task-specific optimization. As a result, they exhibit considerable flexibility in applications like text spotting, document parsing, and translation. Nevertheless, the cascaded structure of multiple models introduces inherent drawbacks, including error propagation and elevated development and maintenance overhead. Recently, with the progress in visual-language models (VLMs), a number of specialized open-source models for OCR and document parsing have been introduced, such as MonkeyOCR Li et al. (2025), Dots.OCR dots (2024), MinerU2.5 Niu et al. (2025), and PaddleOCR-VL Cui et al. (2025a). These efforts aim to enhance parsing accuracy through large-scale modeling. However, due to the limited robustness of current open-source models in handling complex layouts and lengthy text sequences, many still depend on a preliminary layout analysis module Sun et al. (2025); Zhao et al. (2024) to detect document elements, with the VLM subsequently parsing content within localized regions. While this hybrid design improves usability to some degree, it has yet to fully exploit the potential of VLMs for end-to-end joint inference and unified multi-task modeling.

This report introduces HunyuanOCR, a novel open-source multilingual VLM designed for OCR that delivers commercial-grade performance. Departing from conventional pipeline-based frameworks, HunyuanOCR adopts an end-to-end VLM architecture, establishing a unified foundation for multi-task learning that effectively overcomes long-standing challenges such as error propagation and high maintenance costs. As summarized in Table 1, our model demonstrates significant advantages across four key dimensions: 1) Comprehensive Capability Coverage: HunyuanOCR supports an extensive range of tasks beyond basic document parsing, including text spotting, end-to-end receipt information extraction, video subtitle recognition, text-centric visual question answering (VQA), as well as multilingual recognition and translation. By integrating these diverse capabilities into a unified modeling framework, it addresses complex and varied application needs, establishing itself as one of the most comprehensive OCR expert models in the open-source community. 2) High Inference Efficiency: Built upon the native Hunyuan VLM architecture, the model contains only 1B parameters while maintaining high computational efficiency. This compact design ensures low latency and makes it suitable for on-device deployment, meeting the practical requirements of resource-constrained environments. 3) Superior Performance: HunyuanOCR outperforms leading open-source alternatives on core benchmarks; for instance, it surpasses MinerU2.5 and PaddleOCR-VL on the OmniDocBench for document parsing. It also excels in specialized tasks—exceeding Qwen3-VL-4B Bai et al. (2025) in text image translation and information extraction, and outperforming PaddleOCR 3.0 and certain commercial Cloud OCR APIs in text spotting tasks. 4) Enhanced Usability and Unified Modeling: The end-to-end VLM architecture

Table 1: Performance comparison of different VLMs and OCR systems across multiple tasks. ☀ indicates Supported and High-Performing, 🌙 indicates Supported with Moderate Performance, and ⚡ indicates Supported but Underperforming. Otherwise, it is Not Supported.

Model Type	Inference Type	Model Name	Deployment Cost	Task				
				Spotting	Parsing	Text-VQA	IE	Translation
Cascade Pipeline	Multi-Step	PaddleOCR-V5	low	🌙	-	-	-	-
		BaiduOCR	low	☀	-	-	-	-
		Marker-1.8.2	low	-	⭐	-	-	-
		PP-ChatOCR	medium	-	-	-	⭐	-
		PP-DocTranslation	high	-	-	-	-	🌙
Specialized VLMs (Modular)	two-stage	MonkeyOCR-pro-3B	medium	-	☀	-	-	-
		MinerU2.5	low	-	☀	-	-	-
		PaddleOCR-VL	low	-	☀	-	-	-
General VLMs	One-Step	Gemini-2.5-Pro	high	⭐	🌙	☀	☀	☀
		Seed-1.6-Vision	high	🌙	🌙	☀	⭐	☀
		Qwen3-VL-235B-Instruct	high	🌙	🌙	☀	🌙	☀
Specialized VLMs (End2End)	One-Step	Mistral-OCR	medium	-	🌙	-	-	-
		Deepseek-OCR	medium	-	🌙	⭐	⭐	-
		dots.ocr	medium	-	☀	-	-	-
		HunyuanOCR	low	☀	☀	☀	☀	☀

enables unified task modeling within a single framework, allowing diverse OCR tasks to be accomplished through a single inference based on natural language instructions. This design eliminates the need for complex model cascading and post-processing, significantly lowering the technical barrier and offering a streamlined, user-friendly solution for diverse application scenarios.

The HunyuanOCR model adopts an efficient, compact architecture that connects a 0.4B-parameter native-resolution Vision Transformer (ViT) Tschannen et al. (2025) to a 0.5B-parameter Hunyuan Large Language Model (LLM) Tencent (2025) via a learnable pooling MLP adapter. The model is trained following the mainstream two-stage paradigm for VLMs. The first stage, pre-training, involves four steps: vision-language alignment, multi-modal pre-training, long-context pre-training, and application-oriented SFT. This stage utilizes a mixture of large-scale open-source data, synthetic element-level data, and high-quality, end-to-end application-oriented data (e.g., complex long-document parsing and text image translation), totaling approximately 200 million high-quality samples. The second stage, post-training, employs the online reinforcement learning algorithm GRPO with task-specific reward mechanisms, significantly improving the model’s accuracy and stability in challenging scenarios such as complex document parsing and text image translation.

This study demonstrates the substantial potential of the end-to-end VLM paradigm when applied to OCR-specific tasks. We attribute the success of HunyuanOCR to two principal insights. First, during pre-training, exposing the model to high-quality, application-aligned data proves critical for performance—especially in complex and long-text document parsing, as well as in text image translation tasks. Second, the design of targeted online reinforcement learning strategies, combined with an emphasis on data diversity and quality, leads to significant gains in OCR-specific VLMs. These improvements are most pronounced in challenging settings such as intricate layout understanding and knowledge-intensive tasks including visual question answering and image-based translation.

2 Related Work

The evolution of Optical Character Recognition (OCR) technology, traceable to the 1950s, has exhibited distinct developmental phases. In the initial stage (1950s–1980s), OCR systems were primarily based on template matching and feature engineering, focusing on basic text recognition in scanned documents. The 1990s witnessed a significant breakthrough with the maturation of machine learning theory, as statistical methods such as Hidden Markov Models (HMMs) Eddy (1996) and Support Vector Machines (SVMs) Cortes & Vapnik (1995) were widely adopted, substantially improving recognition accuracy. Entering the 21st century, rapid advances in deep learning catalyzed a paradigm shift in OCR: system architectures have progressively transitioned from traditional modular frameworks to the current paradigm of unified processing enabled by vision-language models.

2.1 Traditional OCR Systems

Traditional OCR systems typically employ a highly modularized pipeline architecture. Depending on the requirements of specific application scenarios, such systems often incorporate several core processing modules with distinct functionalities, primarily including, but not limited to: deep learning-based text detection, text recognition, document layout analysis, named entity recognition, and optional text translation modules. Over the past few decades, significant research efforts have been devoted to this direction. Through continuous innovation, numerous models have been developed Zhou et al. (2017); Liao et al. (2017; 2022); Shi et al. (2017; 2018); Lyu et al. (2018); Li et al. (2023); Lyu et al. (2024b); Li et al. (2021); Yu et al., substantially enhancing the accuracy and robustness of each functional module.

Nevertheless, conventional OCR systems still suffer from two fundamental limitations that require urgent resolution. First, at the architectural level, these solutions generally rely on cascading multiple independent functional modules, resulting in highly complex system structures. Taking a typical document parsing task as an example, a fully functional system typically requires integrating at least five key subsystems: a high-precision text detection module, a multilingual text recognition engine, a fine-grained layout analysis component, a specialized mathematical formula recognition module, and a structured table recognition unit. This modular stacking design not only increases deployment complexity and maintenance costs but also requires specialized personnel to perform coordinated tuning of each component. Second, during inference, the multi-stage cascaded processing flow leads to progressive error amplification through a “pipeline effect.” Specifically, inaccuracies in text detection can degrade input quality for subsequent recognition modules, while layout analysis errors may cause incorrect ordering of text blocks. These early-stage inaccuracies ultimately compromise the accuracy and usability of the system’s final output. Consequently, traditional OCR systems often fail to meet practical requirements when handling complex scenarios such as documents with overlapping text or non-standard layouts.

2.2 Vision-Language Models

With the rapid advancement of deep learning, large language models (LLMs) Devlin et al. (2019); Radford et al. (2019); Brown et al. (2020); Liu et al. (2024a); Team (2025); Comanici et al. (2025) have achieved remarkable breakthroughs in natural language processing (NLP). Subsequently, VLMs Liu et al. (2023); Achiam et al. (2023); Bai et al. (2025); Comanici et al. (2025); Wang et al. (2025b), which align information across multiple modalities, have demonstrated exceptional capabilities in cross-modal understanding and generation. These models typically employ unified neural network architectures, enabling efficient handling of complex cognitive tasks such as visual recognition, textual comprehension, and multimodal reasoning. The advantages of this paradigm are twofold. First, architecturally, the unified network design supports synergistic multi-task processing, allowing a single model to perform diverse tasks in an end-to-end manner. Second, by leveraging the inherent reasoning abilities of LLMs, this architecture achieves substantial performance gains, particularly in cognition-intensive applications.

2.2.1 General Vision-Language Models

Current mainstream general vision-language models, such as Gemini Comanici et al. (2025) and Qwen-VL Bai et al. (2025), have demonstrated strong OCR capabilities. These models exhibit robust text perception, accurately recognizing both printed and handwritten text while effectively handling complex scenarios involving irregular layouts, low-resolution images, and multilingual content. However, their large parameter size introduces two notable limitations in practical applications. First, inference requires substantial GPU memory and computational resources. Second, they often fail to meet the stringent low-latency requirements of real-world business scenarios.

2.2.2 OCR-Specific Vision-Language Models

To address the aforementioned technical constraints, the development of lightweight, specialized vision-language models for OCR has emerged as a promising solution. Pioneering approaches such as Nougat Blecher et al. (2023) and StructText-V3 Lyu et al. (2024a) attempted to achieve end-to-end processing for document parsing and information extraction within a unified model. Subsequent models including Dolphin Feng et al. (2025), MonkeyOCR Li et al. (2025), Dots.OCR dots (2024), MinerU2.5 Niu et al. (2025), and PaddleOCR-VL Cui et al. (2025a) have drawn inspiration from traditional OCR pipelines. These methods typically first perform layout detection Zhao et al. (2024); Sun et al. (2025) using a dedicated model or a repurposed vision-language model, followed by unified recognition of text blocks, formulas, and tables. While these approaches reduce system complexity and improve accuracy compared to traditional pipelines by leveraging the generalization capability of vision-language models, they remain susceptible to error propagation from the layout analysis stage and fail to fully exploit the benefits of end-to-end optimization.

In contrast, the proposed HunyuanOCR model demonstrates substantial advantages in both technical architecture and application effectiveness across three key dimensions:

- 1) Fully end-to-end architecture: HunyuanOCR employs a purely end-to-end design that eliminates error accumulation from cascaded processing. This architecture maximizes the potential of end-to-end learning through a systematically optimized training paradigm. From an engineering perspective, the model completes entire workflows in a single inference pass, significantly improving operational efficiency in real-world applications.
- 2) Comprehensive functional coverage: Leveraging the unified task-handling capability of vision-language models, HunyuanOCR supports not only basic document parsing and text spotting but also advanced functionalities, including information extraction, visual question answering, and cross-lingual translation. Notably, it provides extensive multilingual support for hundreds of global languages, making it one of the most functionally complete specialized OCR solutions available.
- 3) Superior performance benchmarking: HunyuanOCR achieves exceptional performance, with key metrics significantly surpassing current state-of-the-art models and matching or exceeding the standards of leading commercial OCR APIs.

3 Model Design

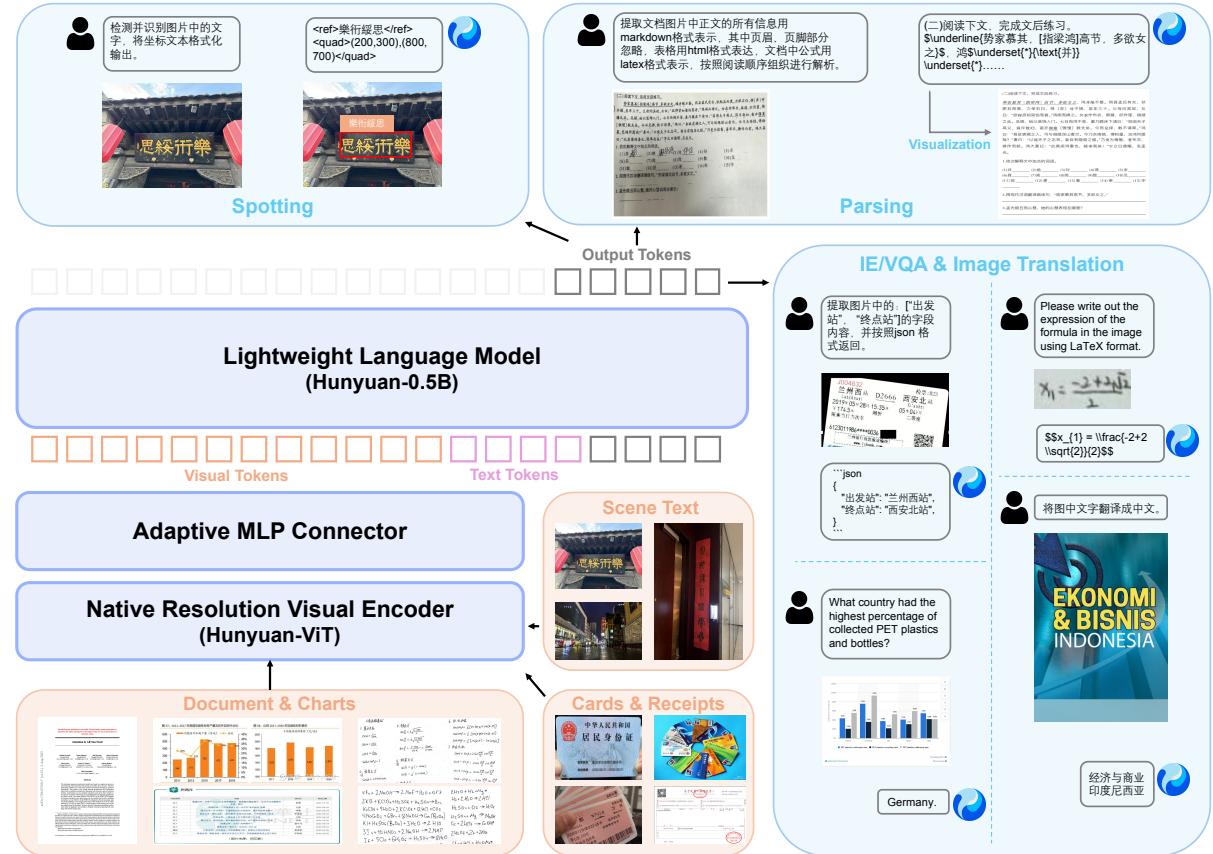


Figure 2: The Architecture of HunyuanOCR: An end-to-end framework integrating Native Resolution Visual Encoder, Adaptive MLP Connector, and a Lightweight Language Model for diverse OCR tasks, including: spotting, parsing, information extraction, visual question answering, and text image translation.

HunyuanOCR features a collaborative architecture comprising three core modules: a Native Resolution Visual Encoder, an Adaptive MLP Connector, and a Lightweight Language Model.

Native Resolution Visual Encoder (Hunyuan-ViT) is built upon the SigLIP-v2-400M pre-trained model Tschanne et al. (2025). By incorporating a hybrid generative-discriminative joint training strategy, it significantly enhances the model's ability to comprehend complex visual semantics. The encoder natively supports arbitrary input resolutions through an adaptive patching mechanism that preserves the original aspect ratio, making it particularly suitable for challenging scenarios involving extreme aspect ratios such as long-text documents. The image is divided into patches according to its native proportions,

and all patches are processed by the Vision Transformer (ViT) with global attention. This design avoids image distortion and detail loss, leading to notable improvements in text recognition accuracy for difficult cases, including long text lines, extensive documents, and low-quality scans.

Adaptive MLP Connector acts as a bridge between the visual and linguistic domains, implementing a core learnable pooling operation. It employs spatial-dimension adaptive content compression to reduce the sequence length of tokens generated from the visual encoder’s high-resolution feature maps, effectively minimizing redundancy. During this process, the module preserves critical semantic information from key areas, such as text-dense regions, thereby achieving an efficient and precise projection of visual features into the input space of the language model.

Lightweight Language Model is based on the densely architected Hunyuan-0.5B model [Tencent \(2025\)](#). It incorporates XD-RoPE, which deconstructs the conventional RoPE [Su et al. \(2024\)](#) into four independent subspaces: text, height, width, and time. This design establishes a native alignment mechanism that bridges 1D text sequences, 2D page layouts, and 3D spatiotemporal information, enabling the model to handle both complex layout parsing (e.g., multi-column recognition) and cross-page document analysis with logical reasoning.

End-to-End Optimization. In contrast to other specialized vision-language OCR models, HunyuanOCR employs a fully end-to-end paradigm for both training and inference. By scaling high-quality, application-oriented data and leveraging reinforcement learning optimization, the system eliminates the need for post-processing and the associated error accumulation typical of pipeline-based architectures. It demonstrates superior robustness in challenging scenarios such as mixed-layout document understanding.

4 Data Construction

4.1 Task Design

Capitalizing on the architectural advantages of vision-language models, HunyuanOCR integrates various OCR tasks into a single, unified paradigm. This enables one model to address multiple high-frequency tasks across the OCR domain.

4.1.1 Spotting

As a fundamental OCR capability, text spotting requires precise localization and recognition of text within images. HunyuanOCR adopts a standardized instruction template for this task, using the fixed prompt: “Detect and recognize text in the image, and output the text coordinates in a formatted manner.” This instruction guides the model to output both line-level text content and corresponding coordinate information. To ensure machine-parsable responses, a structured output format is defined: `<ref>text</ref><quad>(x1, y1), (x2, y2)</quad>`. Here, text inside the `<ref>` and `</ref>` tags denotes the recognized content, and the coordinate sequence within `<quad>` and `</quad>` tags specifies the bounding box of the text region using its top-left and bottom-right vertices. All coordinates are normalized to the range [0, 1000] to maintain consistency across input images of varying resolutions.

4.1.2 Parsing

Document parsing constitutes a core OCR capability, whose strategic importance has been heightened by the rapid advancement of large language models (LLMs). It serves not only as a key preprocessing tool for building high-quality training datasets but also as an essential upstream component in retrieval-augmented generation (RAG) systems.

HunyuanOCR provides a comprehensive document parsing solution that supports both fine-grained element-level parsing and full end-to-end document parsing.

Fine-Grained Element Parsing: It supports the independent identification and extraction of specialized document elements, including mathematical formulas, chemical formulas, tables, and charts. HunyuanOCR employs standardized instruction templates to guide the parsing of different document elements:

- **Formula Parsing:** Using the prompt “Identify the formula in the image and represent it using LaTeX format.”, the model returns the corresponding LaTeX codes of mathematical or chemical formulas.
- **Table Parsing:** Using the prompt “Parse the table in the image into HTML.”, the model returns the HTML codes of the tables.
- **Chart Parsing:** Using the prompt “Parse the chart in the image, use Mermaid format for flowcharts and Markdown for other charts.”, the model adaptively describes the chart using either Mermaid syntax or Markdown based on its type.

End-to-End Document Parsing: HunyuanOCR enables integrated, full-page parsing of documents containing multiple and complex element types. We use the prompt: “Extract all information from the main body of the document image and represent it in markdown format, ignoring headers and footers. Tables should be expressed in HTML format, formulas in the document should be represented using LaTeX format, and the parsing should be organized according to the reading order.” This instruction guides the model to perform an integrated analysis of the document image, outputting all textual content in its natural reading sequence while intelligently converting any identified tables and formulas into HTML and LaTeX formats, respectively, and outputting the spatial positions of figures or charts in the image with corresponding titles. Additionally, we introduce a generalized prompt called “Extract the text in the image”. It is designed for diverse real-world scenarios and guides the model to read any image, such as posters, street views, product packaging, or UI screens, in natural reading order. Detected tables are converted into Markdown format and formulas into LaTeX, producing clean and structured output for broad downstream use.

4.1.3 IE & VQA

HunyuanOCR delivers comprehensive document understanding through robust IE and advanced VQA capabilities.

IE: As a core OCR function, IE precise perceptual localization and deep semantic association. HunyuanOCR provides robust structured extraction with strengths in two primary dimensions:

- **Domain Adaptability:** HunyuanOCR is designed for open-world extraction of arbitrary fields, exhibiting strong domain adaptability while being precisely optimized for over 30 common document types. These include 30 types of cards and receipts, the detailed categories are listed in Table 8.
- **Instruction-Driven Control:** HunyuanOCR allows granular control via natural language instructions. It supports both targeted single-field extraction (e.g., “Please output the value of < Key >”) and parallel multi-field extraction into structured JSON based on user-provided key lists (e.g., “Extract [‘key1’, ‘key2’, …] and return in JSON format”), enabling seamless adaptation to diverse application scenarios.
- **Video Subtitle Extraction:** In response to the instruction “Extract the subtitles from the image,” HunyuanOCR performs subtitle extraction from standard video screenshots, enabling robust handling of text across diverse resolutions, aspect ratios (landscape/portrait), and on-screen positions in both horizontal and vertical orientations.

VQA: HunyuanOCR demonstrates strong open-domain document QA performance, effectively processing open-ended queries about imaged text and generating accurate predictions. Its key capabilities include:

- **Multi-Format Input Support:** The model processes diverse inputs including cropped text lines, mathematical formulas, documents, charts, and street-view imagery for perception and understanding.
- **Advanced Reasoning:** Beyond basic recognition, it performs complex tasks such as spatial and attribute understanding, logical reasoning, and numerical computation based on visual and textual content.

4.1.4 Text Image Translation

HunyuanOCR incorporates a comprehensive end-to-end image-to-text translation module that supports over 14 source languages—including French, German, Japanese, Korean, and many other widely used or regionally important languages—translating them into either Chinese or English. In addition, the system enables direct bidirectional translation between Chinese and English, covering both general-purpose translation scenarios and document-centric translation tasks with complex layouts.

Beyond language coverage, HunyuanOCR is designed for multi-scenario robustness, handling both document-oriented inputs—such as scanned pages, structured layouts, tables, forms, and dense paragraphs—and general scenes containing natural images with embedded text, signage, posters, captions, and other visually diverse contexts. This allows the model to perform reliable translation under variations in layout complexity, image quality, lighting, distortion, and multilingual content distribution.

To fully activate the model’s translation capabilities across different use cases, we design two complementary prompting paradigms:

- **General-purpose Translation Prompt** “Extract all text from the image and translate it into Chinese/English.” This prompt targets general scene-text translation without assuming any document structure.
- **Document-oriented Translation Prompt** “First parse the document, then translate its content into Chinese. Ignore headers and footers; represent equations in LaTeX; and render tables in HTML format.”

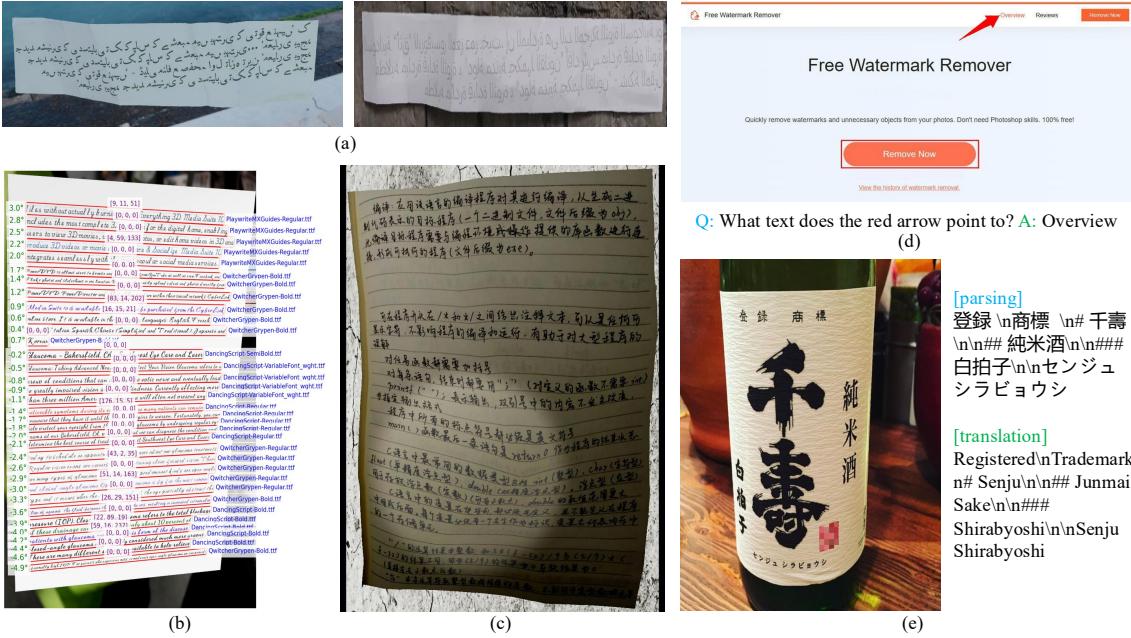


Figure 3: Illustration of image data synthesis and data augmentation results for the HunyuanOCR data pipeline. (a) Multilingual synthetic data with right-to-left (RTL) reading order. (b) Long-document, paragraph-level synthesis with controllable line-level font, language, rotation, and RGB values. (c) Document image warping with realistic defects, including perspective distortion, blur, and local lighting variation. (d) Cross-task data reuse: from spotting data to automated QA generation. (e) Cross-task data reuse: from multilingual parsing data to real-world text translation.

This prompt is tailored for English-to-Chinese document image translation requiring structured parsing.

4.2 Data Pipelines

To systematically enhance HunyuanOCR’s perceptual and comprehension capabilities across diverse scenarios, languages, and layouts, we constructed large-scale, high-quality training data aligned with the core tasks described above. Beyond aggregating public benchmarks, we collected extensive real-world data through web crawling and generated high-quality synthetic samples using proprietary synthesis tools. Via a complete data production and cleaning pipeline (Fig. 3), we built a corpus of over 200 million image-text pairs spanning nine major real-world scenarios—street views, documents, advertisements, handwritten text, screenshots, cards/certificates/invoices, game interfaces, video frames, and artistic typography—and covering more than 130 languages, forming a high-quality multimodal training resource.

4.2.1 Image data synthesis

Building upon the SynthDog framework, we have extended its capabilities to generate high-quality synthetic data for long-document parsing and translation tasks. The system supports paragraph-level rendering in over 130 languages and comprehensively handles bidirectional text layouts (LTR/RTL) as well as complex cursive scripts (Fig. 3(a)–(b)).

The proposed synthesis pipeline exhibits the following core characteristics. First, it enables fine-grained control over text attributes—such as font, color, and orientation—as well as image perturbations, including lighting and shadows, during the rendering process. Second, it accurately simulates complex typographical features, such as handwritten-style fonts and mixed-font typesetting. Furthermore, the system significantly enhances support for low-resource languages, effectively improving cross-lingual generalization in OCR and machine translation. Finally, through a unified architecture, it generates image-text aligned data suitable for a variety of tasks, including spotting, long-document parsing, and cross-lingual translation.

4.2.2 Image data augmentation

We employ an in-house Warping Synthesis Pipeline to simulate realistic imaging defects in photographed and natural-scene documents, thereby enhancing model robustness (Fig. 3(c)). The pipeline incorporates three key functions: geometric deformation via control-point manipulation to emulate folds, curves, and perspective distortions; imaging degradation with motion blur, Gaussian noise, and compression artifacts; and illumination perturbations that model global/local lighting variations, shadows, and reflections. This pipeline substantially improves the robustness of core OCR tasks, such as text spotting, document parsing, and visual question answering.

4.2.3 Question–Answer Pair Generation

We have developed an automated pipeline that integrates Hard Sample Retrieval, QA Generation, and Consistency Verification to produce high-quality VQA data while maximizing cross-task sample reuse. Based on a “single source, multiple uses” principle, the pipeline jointly manages spotting, parsing outputs, and VQA annotations for each image, enabling unified training across text spotting, document parsing, and text-centric VQA tasks.

Hard Sample Retrieval. We employ an automated image and label-based filtering strategy to identify challenging cases from large-scale datasets. Priority is given to samples with low clarity, complex tables or formulas, code snippets, and low-resource language text. This approach ensures that extensive training effectively enhances model performance on these challenging scenarios.

Instructional QA Generation. We designed unified instruction templates to automatically generate question-answer (QA) pairs for multiple types of tasks using a high-performance visual language model (VLM). For instance, the system can produce parsing tasks encompassing the recognition and conversion of elements such as code snippets, formulas, tables, and charts into structured formats, including Markdown, HTML, and JSON. Furthermore, by leveraging textual content, chart attributes, semantic information, and numerical data present in the image, the method generates diverse QA pairs covering information extraction, numerical computation, content summarization, and other reasoning tasks.

Consistency Verification and Data Refinement. We employ a multi-model cross-validation mechanism to evaluate the confidence of generated question-answer (QA) pairs. Data that pass the validation are directly incorporated into the training set to ensure quality, while a subset of the failing cases undergoes manual verification to supplement challenging samples that are difficult for the models to process, thereby enhancing the diversity and coverage of the dataset.

5 Training Recipe

5.1 Pre-Training

We employ a four-stage training strategy for HunyuanOCR pre-training, as outlined in Table 2. The process begins with Stage 1, which warms up the vision–language bridge. In Stage 2, all model parameters are unlocked for end-to-end multimodal learning. Stage 3 extends the context window to 32k tokens to support long-document parsing and understanding. Finally, Stage 4 conducts application-oriented tuning using standardized instructions and normalized outputs, establishing a solid foundation for subsequent reinforcement learning.

- **Stage-1:** In the first stage, we train only the visual encoder (ViT) and a learnable MLP adapter while keeping the language model frozen, aligning visual features with the textual semantic space. The training corpus consists primarily of general image captioning data and synthetic OCR data focused on parsing and recognition tasks, supplemented with a small proportion of plain text ($\leq 10\%$) to preserve the core linguistic capabilities of the language model. This stage emphasizes text parsing and recognition to enhance the model’s perception and structured understanding of textual content in images. Training uses approximately 50B tokens, with the learning rate warmed up from 3×10^{-4} to a peak value before decaying to 3×10^{-5} .
- **Stage-2:** In the second stage, all model parameters are unfrozen for end-to-end vision-language joint learning, with a focus on enhancing the model’s capability for deep understanding and cognitive reasoning of structured content such as documents, tables, and charts. The training data mixture increases the proportion of synthetic samples covering multiple tasks, including text parsing, spotting, translation, and VQA, while retaining approximately ($\leq 10\%$) plain text to maintain instruction-following and linguistic generalization capabilities. The training utilizes approximately 300B tokens with a warmup-cosine learning rate schedule, decaying from 2×10^{-4} to 5×10^{-5} .
- **Stage-3:** In the third stage, we extend the model’s context window to 32K by incorporating long-

Table 2: Overview of the four-stage pre-training recipe for HunyuanOCR pre-training.

Stages	Stage-1	Stage-2	Stage-3	Stage-4
Purpose	Vision-Language Alignment	Multimodal Pre-traning	Long-context Pre-training	Application-oriented SFT
Trainable Parts	ViT & Adapter	All	All	All
Learning Rate	$3e-4 \rightarrow 3e-5$	$2e-4 \rightarrow 5e-5$	$8e-5 \rightarrow 5e-6$	$2e-5 \rightarrow 1e-6$
Training Tokens	50B	300B	80B	24B
Sequence Length	8k	8k	32k	32k
Data Composition	Pure Text, Synthetic Parsing and Recognition Data, General Image Caption Data	Pure Text, Synthetic Spotting, Parsing, Translation and VQA Data	Long Pure Text, Real-world Auto-annotated Data, Long Document Parsing Data, Information Extraction Data	Human-annotated Data, Hard-negative Data, Standardized Instruction Data.

context parsing tasks and lengthy plain text data. This stage uses approximately 80B tokens, decaying the learning rate from 8×10^{-5} to 5×10^{-6} .

- **Stage-4:** We conduct annealing training using carefully curated, human-annotated real-world data supplemented with a small proportion of high-quality synthetic samples, while maintaining a 32K context window to enhance perceptual robustness in complex scenarios. By employing unified instruction templates and standardized output formats across different tasks, we ensure consistency in response patterns throughout the training data. This design not only reduces the model’s learning difficulty but also facilitates the design of reward models in subsequent post-training stages. The training utilizes 24B tokens in this stage, with the learning rate linearly decaying from 2×10^{-5} to 1×10^{-6} .

5.2 Reinforcement Learning (RL)

Reinforcement learning (RL) algorithms have emerged as a powerful paradigm, achieving remarkable success across various domains involving large language models (LLMs) and multimodal large language models (MLLMs). Notable applications include mathematical reasoning [Shao et al. \(2024\)](#), image segmentation [Liu et al. \(2025\)](#), and omni-multimodal LLMs [Zhao et al. \(2025\)](#). This broad success is largely attributed to RL’s ability to align model outputs with verifiable metrics [Wen et al. \(2025\)](#) or human preferences [Peng et al. \(2025a;b\)](#).

While RL has traditionally been applied to large-scale reasoning models, we investigate its application to lightweight OCR models that prioritize efficient and accurate text understanding. Leveraging the structured nature and inherent verifiability of many OCR tasks, we adopt Reinforcement Learning with Verifiable Rewards (RLVR) for closed-form tasks such as text spotting and document parsing. For more open-ended tasks like translation and text-centric VQA, we design reward mechanisms based on an LLM-as-a-judge approach. By integrating RLVR and LLM-as-a-judge techniques, we demonstrate that even lightweight models can achieve significant performance improvements, opening new possibilities for edge and mobile applications.

5.2.1 Data Curation

Our data pipeline emphasizes **quality, diversity, and difficulty balance**. In terms of quality, we combine high-quality open-source and synthetic datasets, and filter them using LLM-based judging to ensure image–text alignment and the removal of tasks that are easily exploitable (e.g., multiple-choice). For diversity, we cover a broad range of OCR-related tasks mentioned above, and maintain sufficient exploration by discarding samples with low output diversity or zero reward variance. Finally, to balance task difficulty, we employ pass-rate filtering based on model samples, removing both trivial and unsolvable examples.

5.2.2 Reward Design

We adopt a **ability-adaptive reward design**, where each OCR-related task type has a tailored reward formulation that aligns with its output characteristics.

- **Spotting:** For text spotting tasks, which require joint text recognition and bounding box localization, the reward is computed as follows. Each predicted bounding box is first assigned to a ground-truth box by maximizing the Intersection over Union (IoU). The reward for each matched pair is then calculated as one minus the normalized edit distance between the predicted and ground-truth text strings. Any unmatched predictions or ground-truth boxes incur a penalty by contributing a reward of zero to the

average. The final reward is the mean score across all evaluated pairs, providing a balanced measure of both localization and recognition accuracy.

- **Document Parsing:** Document Parsing aims to convert document images into structured formats containing textual content, mathematical formulas, and tables. The evaluation emphasizes both structural integrity and content accuracy. The reward is computed based on the normalized edit distance between the model’s output and the ground-truth reference.
- **VQA:** The reward is binary (1 or 0), based on whether the model’s answer semantically matches the reference. The scoring model evaluates only content completeness and factual correctness, tolerating minor stylistic differences while enforcing strict alignment on key content elements.
- **Translation:** We use a soft reward scheme where a scoring LLM compares the generated output against the reference and assigns a score in the range [0, 5]. This raw score is then debias-normalized to [0, 1]. Crucially, this normalization is designed to expand the reward granularity in the mid-range (2–4), enabling the model to better capture subtle improvements and differences in translation quality.

5.2.3 Training Strategy

We adopt the Group Relative Policy Optimization (GRPO) algorithm as our main reinforcement learning framework. In each training iteration, GRPO samples a group of responses (o_1, o_2, \dots, o_G) for a given query (q) from the old policy $(\pi_{\theta_{\text{old}}})$ and updates the current policy (π_{θ}) by maximizing the objective:

$$\begin{aligned} \mathcal{L}_{\text{GRPO}}(\theta) = & \mathbb{E}_{[q \sim \mathcal{D}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q)]} \\ & \frac{1}{G} \sum_{i=1}^G \left[\min \left(\frac{\pi_{\theta}(o_i | q)}{\pi_{\theta_{\text{old}}}(o_i | q)} A_i, \text{clip} \left(\frac{\pi_{\theta}(o_i | q)}{\pi_{\theta_{\text{old}}}(o_i | q)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) - \beta \mathbb{D}_{\text{KL}}(\pi_{\theta} || \pi_{\text{ref}}) \right] \end{aligned} \quad (1)$$

where A_i represents the advantage computed from the group rewards, and \mathbb{D}_{KL} is the KL-divergence term for regularization. The ϵ and β control clipping and the strength of KL penalties, respectively.

To ensure stable and reliable training, we enforce length constraints and a strict format during reward computation. Specifically, any output that exceeds the maximum length is immediately assigned a reward of zero. Similarly, for structured tasks like spotting and document parsing, outputs that fail to follow the required schema are also directly penalized with zero reward. These constraints help the optimization process to focus exclusively on valid, well-structured, and verifiable outputs, thereby guiding the model to learn accurate reasoning and formatting behavior under constrained conditions.

6 Evaluation

6.1 Spotting

To comprehensively evaluate the model’s text spotting performance across diverse scenarios, we constructed a benchmark comprising nine categories: artistic text, document images, game screenshots, handwritten text, advertisement scenes, card/certificate/invoice images, screen captures, street view text, and video frames. Each category contains 100 images, forming a 900-image evaluation set. Based on this benchmark, we compared HunyuanOCR with traditional pipeline-based open-source models, leading commercial APIs, and general Vision-Language Models (VLMs). The results shown in Table 3 demonstrate that our approach achieves the best overall performance.

Specifically, as an end-to-end VLM solution, HunyuanOCR significantly outperforms traditional pipeline-based methods. Furthermore, compared to general VLMs, our method achieves superior accuracy with substantially fewer parameters, demonstrating notable advantages in both computational efficiency and performance.

Table 3: Comprehensive evaluation of spotting ability on in-house benchmark.

Model Type	Model	Overall	Art	Doc	Game	Hand	Ads	Receipt	Screen	Scene	Video
Traditional methods	PaddleOCR Cui et al. (2025b)	53.38	32.83	70.23	51.59	56.39	57.38	50.59	63.38	44.68	53.35
	BaiduOCR Baidu (2025)	61.90	38.5	78.95	59.24	59.06	66.70	63.66	68.18	55.53	67.38
General VLMs	Gemini-2.5-Pro Comanici et al. (2025)	23.44	21.79	35.16	10.02	38.49	29.89	20.80	17.59	18.33	18.90
	Qwen3-VL-2B-Instruct Qwen (2025)	29.68	29.43	19.37	20.85	50.57	35.14	24.42	12.13	34.90	40.10
	Qwen3-VL-235B-A22B-Instruct Qwen (2025)	53.62	46.15	43.78	48.00	68.90	64.01	47.53	45.91	54.56	63.79
OCR-Specific VLMs	Seed-1.6-Vision Seed (2025)	59.23	45.36	55.04	59.68	67.46	65.99	55.68	59.85	53.66	70.33
	HunyuanOCR	70.92	56.76	73.63	73.54	77.10	75.34	63.51	76.58	64.56	77.31

Table 4: Parsing performance evaluated across multilingual settings and diverse document scenarios.

Model Type	Model	Size	OmniDocBench				Wild-OmniDocBench				DocML
			overall↑	text↓	formula↑	table↑	overall↑	text↓	formula↑	table↑	
General VLMs	Gemini-2.5-pro 2025	-	88.03	0.075	85.92	85.71	80.59	0.118	75.03	78.56	82.64
	Qwen3-VL-235B 2025	235B	89.15	0.069	88.14	86.21	79.69	0.09	80.67	68.31	81.40
Specialized VLMs (Modular)	MonkeyOCR-pro 2025	3B	88.85	0.075	87.5	86.78	70.00	0.211	63.27	67.83	56.50
	MinerU2.5 2025	1.2B	90.67	0.047	88.46	88.22	70.91	0.218	64.37	70.15	52.05
	PaddleOCR-VL 2025a	0.9B	92.86	0.035	91.22	90.89	72.19	0.232	65.54	74.24	57.42
Specialized VLMs (End2End)	Mistral-OCR 2025	-	78.83	0.164	82.84	70.03	-	-	-	-	64.71
	Deepseek-OCR 2025	3B	87.01	0.073	83.37	84.97	74.23	0.178	70.07	70.41	57.22
	dots.ocr 2024	3B	88.41	0.048	83.22	86.78	78.01	0.121	74.23	71.89	77.50
	HunyuanOCR	1B	94.10	0.042	94.73	91.81	85.21	0.081	82.09	81.64	91.03

6.2 Parsing

We systematically evaluated the model’s performance on document parsing using three benchmark datasets. First, we conducted experiments on OmniDocBench [Ouyang et al. \(2024\)](#), a publicly available and comprehensive document parsing benchmark that includes a diverse set of digital and scanned documents covering formulas, tables, paragraphs, and various structural elements. Second, to further assess the model’s robustness in real-world captured scenarios, we created a Wild version of OmniDocBench¹ by printing the original documents and re-capturing them under challenging conditions—such as manual folding, bending, and varying illumination—to simulate realistic distortions encountered in everyday document photography. Finally, we evaluated the model on DocML², our internally curated multilingual parsing dataset designed to assess robustness across multiple languages and acquisition settings. DocML spans both digital/scanned and real-world captured documents across 14 high-frequency non-Chinese/English languages, including German, Spanish, Turkish, Vietnamese, Korean, Malay, Portuguese, Russian, French, Indonesian, Thai, Italian, and Japanese.

For both OmniDocBench and its Wild variant, we followed the official evaluation protocol described in [Ouyang et al. \(2024\)](#) and report results for HunyuanOCR alongside other leading document parsing models. As shown in Table 4, HunyuanOCR achieves the highest overall performance on both the digital/scanned and real-world captured settings, demonstrating strong generalization across diverse document formats and acquisition conditions. Notably, despite its relatively compact 1B parameter size, HunyuanOCR outperforms larger specialized OCR or VLM-based parsing models. On DocML, we adopt an overall edit-distance-based score as the evaluation metric to comprehensively measure the accuracy and robustness of parsed outputs across multilingual settings. Under this metric, HunyuanOCR demonstrates excellent multilingual parsing performance, achieving state-of-the-art results across all 14 languages. These findings collectively show that HunyuanOCR delivers robust and accurate document parsing across multilingual, multi-scene, and real-world conditions.

6.3 IE & VQA

We systematically evaluate the model’s performance on information extraction and open-ended visual question answering tasks using three benchmark datasets. First, to assess the model’s capability on high-frequency card and document types, we constructed a test set comprising 768 samples across 30 common categories (Table 8), such as identification cards, passports, and invoices. Second, to evaluate text extraction performance in complex scenarios, we built a video subtitle dataset containing 1,000 samples covering diverse video contexts and subtitle styles. Additionally, the model was comprehensively evaluated on OCRBench [Liu et al. \(2024c\)](#), a publicly available benchmark that includes 1,000 test samples and spans multiple competencies, including scene text recognition, handwritten text and formula recognition, information extraction, and open-ended question answering on documents and charts.

We evaluated the first two benchmarks using exact-match accuracy under a unified prompting protocol for multi-field JSON outputs, while adopting the official standard evaluation protocol for OCRBench. HunyuanOCR was compared against leading SOTA VLMs, including Qwen3VL-235B-Instruct, Seed1.6-VL-Instruct, and Gemini-2.5-Pro, using identical prompts and post-processing procedures such as JSON format parsing. As summarized in Table 5, HunyuanOCR achieves the highest overall accuracy across

¹The Wild version of OmniDocBench will be publicly released in a future update.

²The DocML multilingual parsing dataset will also be open-sourced in a future release. We invite interested parties to reach out to us for access or evaluation prior to its public release.

Table 5: Evaluation of information extraction and visual question answering tasks.

Model	Cards	Receipts	Video Subtitles	OCRBench
DeepSeek-OCR Wei et al. (2025)	10.04	40.54	5.41	430
PP-ChatOCR PaddleOCR (2025)	57.02	50.26	3.1	-
Qwen3-VL-2B-Instruct 2025	67.62	64.62	3.75	858
Seed-1.6-Vision Seed (2025)	70.12	67.5	60.45	881
Qwen3-VL-235B-A22B-Instruct 2025	75.59	78.4	50.74	920
Gemini-2.5-Pro 2025	80.59	80.66	53.65	872
HunyuanOCR	92.29	92.53	92.87	860

Table 6: Evaluation of photo translation. We additionally manually annotated DocML with high-quality English and Chinese reference translations to serve as ground-truth labels for evaluating text translation performance.

Model	Size	DocML		DoTA
		other2en	other2zh	en2zh
Gemni-2.5-Flash Comanici et al. (2025)	-	79.26	80.06	85.60
Qwen3-VL-235B-Instruct Qwen (2025)	235B	73.67	77.20	80.01
Qwen3-VL-8B-Instruct Qwen (2025)	8B	75.09	75.63	79.86
Qwen3-VL-4B-Instruct Qwen (2025)	4B	70.38	70.29	78.45
Qwen3-VL-2B-Instruct Qwen (2025)	2B	66.30	66.77	73.49
PP-DocTranslation	-	52.63	52.43	82.09
HunyuanOCR	1B	73.38	73.62	83.48

all 30 document categories in card/receipts information extraction and subtitle extraction tasks, despite having only around 1B parameters, significantly outperforming considerably larger VLMs such as Qwen3VL-235B-Instruct, Seed1.6-VL, and Gemini-2.5-Pro. On OCRBench, HunyuanOCR also demonstrates substantially better performance than DeepseekOCR at a similar scale and comparable with the larger Qwen3VL-2B-Instruct model.

6.4 Text Image Translation

We systematically evaluated the model’s text image translation capability using two benchmark datasets. For public benchmarking, we selected DoTA [Liang et al. \(2024\)](#), a document translation dataset designed for complex and diverse English-layout document scenarios, and used it to assess the model’s English-to-Chinese translation performance under realistic document conditions. In addition, we constructed an in-house evaluation benchmark based on DocML, where each sample is annotated with both English and Chinese translations. This internal benchmark enables a comprehensive assessment of translation robustness across multiple languages and a broad range of document types, including both digital/scanned and real-world captured scenes.

To evaluate translation quality, we adopt the COMET [Rei et al. \(2022\)](#) metric, a widely used neural-based evaluation standard for machine translation. As summarized in Table 6, HunyuanOCR surpasses VLMs with over 8B parameters on DoTA, demonstrating strong translation performance in complex document layouts despite its compact 1B scale. Furthermore, we achieved first place in the Track 2.2 OCR-free Small Model of the ICDAR 2025 Competition on End-to-End Document Image Machine Translation Towards Complex Layouts [Zhang et al. \(2025\)](#), validating the effectiveness and generality of our approach. On the DocML evaluation set, HunyuanOCR again outperforms several larger VLMs exceeding 4B parameters, highlighting its robust multilingual translation capability across diverse layouts, languages, and acquisition conditions. These findings collectively demonstrate that HunyuanOCR provides a highly efficient yet powerful solution for text image translation in both public benchmarks and real-world multilingual scenarios.

However, due to its relatively small language model, HunyuanOCR’s translation capability lags behind its strong text detection, recognition, and document parsing performance. For applications requiring higher translation accuracy, developers can cascade our multilingual parsing module with Hunyuan-MT-7B³ or await our upcoming general vision-language models to further boost overall translation quality.

³<https://huggingface.co/tencent/Hunyuan-MT-7B>

7 Conclusion

In this paper, we present HunyuanOCR, an open-source expert vision-language model that unifies diverse OCR tasks within a lightweight, end-to-end architecture. Our work demonstrates that a compact model with only 1B parameters can achieve competitive performance against larger general-purpose VLMs and traditional pipeline systems, validating the effectiveness of our data-centric training strategy and targeted reinforcement learning approach. HunyuanOCR achieves state-of-the-art results in text spotting, document parsing, and information extraction, while significantly simplifying deployment pipelines. These advancements align with our original goal of balancing versatility with efficiency, as outlined in the abstract. Looking ahead, we will continue to optimize inference efficiency through token compression and architectural improvements, while expanding the model's capability to handle higher-resolution and multi-page documents. Our long-term goal remains to adapt HunyuanOCR for edge-device deployment, further democratizing robust OCR intelligence for real-world applications.

Contributors

- Project Sponsors: Jie Jiang, Linus
- Project Supervisor: Han Hu
- Project Leader: Chengquan Zhang
- Core Contributors: Pengyuan Lyu, Xingyu Wan, Gengluo Li, Shangpin Peng
- Contributors: Weinong Wang, Liang Wu, Huawei Shen, Yu Zhou, Canhui Tang, Qi Yang, Qiming Peng, Bin Luo, Hower Yang, Houwen Peng, Hongming Yang, Senhao Xie, Binghong Wu, Mana Yang, Sergey Wang, Raccoon Liu, Dick Zhu

References

- Mistral OCR: Free online ai ocr tool to extract text. <https://www.mistralocr.com/>, 2025. Accessed: 2025-07-30.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Ibrahim Adeshola and Adeola Praise Adepoju. The opportunities and challenges of chatgpt in education. *Interactive Learning Environments*, 32(10):6159–6172, 2024.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Baidu. BaiduOCR API, 2025. URL <https://ai.baidu.com/tech/ocr/general>.
- Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. Nougat: Neural optical understanding for academic documents. *arXiv preprint arXiv:2308.13418*, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- Cheng Cui, Ting Sun, Suyin Liang, Tingquan Gao, Zelun Zhang, Jiaxuan Liu, Xueqing Wang, Changda Zhou, Hongen Liu, Manhui Lin, et al. Paddleocr-vl: Boosting multilingual document parsing via a 0.9 b ultra-compact vision-language model. *arXiv preprint arXiv:2510.14528*, 2025a.
- Cheng Cui, Ting Sun, Manhui Lin, Tingquan Gao, Yubo Zhang, Jiaxuan Liu, Xueqing Wang, Zelun Zhang, Changda Zhou, Hongen Liu, et al. Paddleocr 3.0 technical report. *arXiv preprint arXiv:2507.05595*, 2025b.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.

dots. dots.ocr: Multilingual document layout parsing in a single vision-language model, 2024. URL <https://github.com/rednote-hilab/dots.ocr>.

Sean R Eddy. Hidden markov models. *Current opinion in structural biology*, 6(3):361–365, 1996.

Hao Feng, Shu Wei, Xiang Fei, Wei Shi, Yingdong Han, Lei Liao, Jinghui Lu, Binghong Wu, Qi Liu, Chunhui Lin, et al. Dolphin: Document image parsing via heterogeneous anchor prompting. *arXiv preprint arXiv:2505.14059*, 2025.

JaidedAI. Easyocr, 2020. URL <https://github.com/JaidedAI/EasyOCR>.

Zhanghui Kuang, Hongbin Sun, Zhizhong Li, Xiaoyu Yue, Tsui Hin Lin, Jianyong Chen, Huaqiang Wei, Yiqin Zhu, Tong Gao, Wenwei Zhang, et al. Mmocr: a comprehensive toolbox for text detection, recognition and understanding. In *ACM Multimedia*, pp. 3791–3794, 2021.

Minghao Li, Tengchao Lv, Jingye Chen, Lei Cui, Yijuan Lu, Dinei Florencio, Cha Zhang, Zhoujun Li, and Furu Wei. Trocr: Transformer-based optical character recognition with pre-trained models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 13094–13102, 2023.

Yulin Li, Yuxi Qian, Yuechen Yu, Xiameng Qin, Chengquan Zhang, Yan Liu, Kun Yao, Junyu Han, Jingtuo Liu, and Errui Ding. Structext: Structured text understanding with multi-modal transformers. In *Proceedings of the 29th ACM international conference on multimedia*, pp. 1912–1920, 2021.

Zhang Li, Yuliang Liu, Qiang Liu, Zhiyin Ma, Ziyang Zhang, Shuo Zhang, Zidun Guo, Jiarui Zhang, Xinyu Wang, and Xiang Bai. Monkeyocr: Document parsing with a structure-recognition-relation triplet paradigm. *arXiv preprint arXiv:2506.05218*, 2025.

Yupu Liang, Yaping Zhang, Cong Ma, Zhiyang Zhang, Yang Zhao, Lu Xiang, Chengqing Zong, and Yu Zhou. Document image machine translation with dynamic multi-pre-trained models assembling. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 7084–7095, 2024.

Minghui Liao, Baoguang Shi, Xiang Bai, Xinggang Wang, and Wenyu Liu. Textboxes: A fast text detector with a single deep neural network. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.

Minghui Liao, Zhisheng Zou, Zhaoyi Wan, Cong Yao, and Xiang Bai. Real-time scene text detection with differentiable binarization and adaptive scale fusion. *IEEE transactions on pattern analysis and machine intelligence*, 45(1):919–931, 2022.

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.

Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12), 2024b.

Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12):220102, 2024c.

Yuqi Liu, Bohao Peng, Zhisheng Zhong, Zihao Yue, Fanbin Lu, Bei Yu, and Jiaya Jia. Seg-zero: Reasoning-chain guided segmentation via cognitive reinforcement. *arXiv preprint arXiv:2503.06520*, 2025.

Shangbang Long, Xin He, and Cong Yao. Scene text detection and recognition: The deep learning era. *International Journal of Computer Vision*, 129(1):161–184, 2021.

Pengyuan Lyu, Minghui Liao, Cong Yao, Wenhao Wu, and Xiang Bai. Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 67–83, 2018.

-
- Pengyuan Lyu, Yulin Li, Hao Zhou, Weihong Ma, Xingyu Wan, Qunyi Xie, Liang Wu, Chengquan Zhang, Kun Yao, Errui Ding, et al. Structextv3: An efficient vision-language model for text-rich image perception, comprehension, and beyond. *arXiv preprint arXiv:2405.21013*, 2024a.
- Pengyuan Lyu, Chengquan Zhang, Shanshan Liu, Meina Qiao, Yangliu Xu, Liang Wu, Kun Yao, Junyu Han, Errui Ding, and Jingdong Wang. Maskocr: Scene text recognition with masked vision-language pre-training. *Transactions on Machine Learning Research*, 2024b.
- Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, et al. Mineru2. 5: A decoupled vision-language model for efficient high-resolution document parsing. *arXiv preprint arXiv:2509.22186*, 2025.
- Linke Ouyang, Yuan Qu, Hongbin Zhou, Jiawei Zhu, et al. Omnidocbench: Benchmarking diverse pdf document parsing with comprehensive annotations, 2024. URL <https://arxiv.org/abs/2412.07626>.
- PaddleOCR. Pp-chatocr, 2025. URL <https://github.com/PaddlePaddle/PaddleOCR>.
- Shangpin Peng, Weinong Wang, Zhuotao Tian, Senqiao Yang, Xing Wu, Haotian Xu, Chengquan Zhang, Takashi Isobe, Baotian Hu, and Min Zhang. Omni-dpo: A dual-perspective paradigm for dynamic preference learning of llms. *arXiv preprint arXiv:2506.10054*, 2025a.
- Shangpin Peng, Senqiao Yang, Li Jiang, and Zhuotao Tian. Mitigating object hallucinations via sentence-level early intervention. *arXiv preprint arXiv:2507.12455*, 2025b.
- Qwen. Qwen3-vl, 2025. URL <https://github.com/QwenLM/Qwen3-VL>.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Philipp Koehn, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Koci, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri (eds.), *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pp. 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.wmt-1.52/>.
- Seed. Seed1.6, 2025. URL https://seed.bytedance.com/en/seed1_6.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300, 2024. doi: 10.48550/ARXIV.2402.03300. URL <https://doi.org/10.48550/arXiv.2402.03300>.
- Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(11):2298–2304, 2017.
- Baoguang Shi, Mingkun Yang, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. Aster: An attentional scene text recognizer with flexible rectification. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2035–2048, 2018.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- Ting Sun, Cheng Cui, Yuning Du, and Yi Liu. Pp-doclayout: A unified document layout detection model to accelerate large-scale data construction. *arXiv preprint arXiv:2503.17213*, 2025.
- Qwen Team. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Tencent. Hunyuan-0.5b, 2025. URL <https://github.com/Tencent-Hunyuan/Hunyuan-0.5B>.
- Michael Tschanne, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*, 2025.

-
- Shansong Wang, Mingzhe Hu, Qiang Li, Mojtaba Safari, and Xiaofeng Yang. Capabilities of gpt-5 on multimodal medical reasoning. *arXiv preprint arXiv:2508.08224*, 2025a.
- Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025b.
- Haoran Wei, Yaofeng Sun, and Yukun Li. Deepseek-ocr: Contexts optical compression. *arXiv preprint arXiv:2510.18234*, 2025.
- Xumeng Wen, Zihan Liu, Shun Zheng, Zhijian Xu, Shengyu Ye, Zhirong Wu, Xiao Liang, Yang Wang, Junjie Li, Ziming Miao, et al. Reinforcement learning with verifiable rewards implicitly incentivizes correct reasoning in base llms. *arXiv preprint arXiv:2506.14245*, 2025.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *National Science Review*, 11(12):nwae403, 2024.
- Yuechen Yu, Yulin Li, Chengquan Zhang, Xiaoqiang Zhang, Zengyuan Guo, Xiameng Qin, Kun Yao, Junyu Han, Errui Ding, and Jingdong Wang. Structextv2: Masked visual-textual prediction for document image pre-training. In *The Eleventh International Conference on Learning Representations*.
- Qintong Zhang, Victor Shea-Jay Huang, Bin Wang, Junyuan Zhang, Zhengren Wang, Hao Liang, Shawn Wang, Matthieu Lin, Conghui He, and Wentao Zhang. Document parsing unveiled: Techniques, challenges, and prospects for structured information extraction. *CoRR*, abs/2410.21169, 2024a.
- Qintong Zhang, Bin Wang, Victor Shea-Jay Huang, Junyuan Zhang, Zhengren Wang, Hao Liang, Conghui He, and Wentao Zhang. Document parsing unveiled: Techniques, challenges, and prospects for structured information extraction. *arXiv preprint arXiv:2410.21169*, 2024b.
- Yaping Zhang, Yupu Liang, Zhiyang Zhang, Zhiyuan Chen, Lu Xiang, Yang Zhao, Yu Zhou, and Chengqing Zong. Icdar 2025 competition on end-to-end document image machine translation towards complex layouts. In *International Conference on Document Analysis and Recognition*, pp. 505–522. Springer, 2025.
- Jiaxing Zhao, Xihan Wei, and Liefeng Bo. R1-omni: Explainable omni-multimodal emotion recognition with reinforcement learning. *arXiv preprint arXiv:2503.05379*, 2025.
- Zhiyuan Zhao, Hengrui Kang, Bin Wang, and Conghui He. Doclayout-yolo: Enhancing document layout analysis through diverse synthetic data and global-to-local adaptive perception. *arXiv preprint arXiv:2410.12628*, 2024.
- Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang. EAST: an efficient and accurate scene text detector. In *CVPR*, pp. 2642–2651. IEEE Computer Society, 2017.

HunyuanOCR Technical Report

Supplementary Material

This material provides supplementary details to the main paper, including the following sections:

- (A) Recommended Instruction
- (B) Common Supported IE Categories
- (C) Reinforcement Learning Details
 - (C.1) Training Configuration
 - (C.2) Training Dynamics
 - (C.3) Task-wise Performance Improvements
- (D) Qualitative examples

A Recommended Instruction

Table 7 summarizes the recommended instructions for each task supported by HunyuanOCR, and provides a bilingual (Chinese–English) reference. We recommend using the Chinese instructions to ensure the stability and reproducibility of benchmarking results.

Table 7: Bilingual (Chinese–English) Instructions Recommended for Different Task Types.

Task	Chinese	English
Spotting	检测并识别图片中的文字，将文本坐标格式化输出。	Detect and recognize text in the image, and output the text coordinates in a formatted manner.
	识别图片中的公式，用 \LaTeX 格式表示。	Identify the formula in the image and represent it using \LaTeX format.
	把图中的表格解析为HTML。	Parse the table in the image into HTML.
	解析图中的图表，对于流程图使用Mermaid格式表示，其他图表使用Markdown格式表示。	Parse the chart in the image; use Mermaid format for flowcharts and Markdown for other charts.
Parsing	提取文档图片中正文的所有信息用markdown格式表示，其中页眉、页脚部分忽略，表格用html格式表达，文档中公式用 \LaTeX 格式表示，按照阅读顺序组织进行解析。	Extract all information from the main body of the document image and represent it in markdown format, ignoring headers and footers. Tables should be expressed in HTML format, formulas in the document should be represented using \LaTeX format, and the parsing should be organized according to the reading order.
	提取图中的文字。	Extract the text in the image.
Information Extraction	输出Key的值。	Output the value of Key.
	提取图片中的: ['key1','key2',...] 的字段内容，并按照JSON格式返回。	Extract the content of the fields: ['key1','key2', ...] from the image and return it in JSON format.
Translation	提取图片中的字幕。	Extract the subtitles from the image.
	先解析文档，再将文档内容翻译为中文，其中页眉、页脚忽略，公式用 \LaTeX 格式表示，表格用html格式表示。	First parse the document, then translate its content into Chinese. Ignore headers and footers; represent equations in \LaTeX ; and render tables in HTML format.
	提取图中文字，并将其翻译成中文/英文。	Extract all text from the image and translate it into Chinese/English.

B Common Supported IE Categories

Table 8 summarizes the common card and receipt types covered by the 30 IE tasks. The card side includes more than ten categories, such as ID cards, band cards, passports, social security cards, business licenses, driver’s licenses, vehicle licenses, etc.. The receipt side also spans over ten types, including shopping receipts, taxi receipts, VAT invoices, train tickets, bus tickets, itineraries, bank slips, etc..

Table 8: Common document categories for IE tasks, grouped into Cards & Certificates and Receipts.

Cards & Certificates	Receipts
ID cards	Shopping receipts
Bank cards	Taxi receipts
Social security cards	Ferry tickets
Passports	Train tickets
Household registers	Bus tickets
Mainland travel permits	Itineraries
Business licenses	Express waybills
Driver’s licenses	VAT invoices
Public institution certificates	Bank slips
Vehicle licenses	Medical inspection reports
Professional qualification certificates	Prescriptions
Tax registration certificates	Medical records
Medical insurance vouchers	Checks
Road transport certificates	Ride-hailing itineraries
Vehicle certificates of conformity	Takeout orders

C Reinforcement Learning Details

In this section, we provide additional details on the reinforcement learning (RL) stage of HunyuanOCR beyond the description in the main text Sec. 5.2. We first summarize the RL training configuration (Sec. C.1), then present the training dynamics (Sec. C.2), and finally analyze the performance improvements brought by RL training (Sec. C.3).

C.1 Training Configuration

The detailed training setup for RL is listed in Tab. 9. We adopt a constant learning rate schedule with Adam optimizer, a large global batch size, and long-context settings to fully exploit the long-document understanding capability of HunyuanOCR. No explicit KL penalty is applied during RL, allowing the policy to adjust more freely under the guidance of task-specific rewards. For rollout generation, we use a low temperature of 0.85 and sample $N = 8$ responses per prompt to obtain a diverse set of candidates for reward evaluation and policy updates.

C.2 Training Dynamics

The training dynamics of the RL stage are visualized in Fig. 4, where we track two key statistics: the proportion of samples receiving reward 1 at each step, and the mean reward value. As training progresses, the mean reward increases steadily. This consistent upward trend indicates that the policy gradually learns to produce outputs that better satisfy the task-specific reward criteria, validating the effectiveness and stability of the RL process.

C.3 Task-wise Performance Improvements

After RL training, we observe substantial gains across multiple OCR-related tasks.

Spotting. The spotting ability of HunyuanOCR improves significantly, especially on *Art* and *Screen* scenarios, where the scores increase by more than 2 points. We attribute these gains to the rule-based reward design for the spotting task, which can assess the discrepancy between the predicted outputs and ground-truth annotations at a fine-grained level. This encourages the model to simultaneously improve both the accuracy of the predicted bounding boxes and the correctness of the recognized text.

Parsing. For the parsing task, the score on OmniDocBench increases from 92.5 to 94.1 after RL training. This improvement further demonstrates the effectiveness of the rule-based reward design, which precisely measures content consistency between the model’s outputs and the reference text.

Table 9: Training setup for reinforcement learning.

Setting	Value
Actor	
Learning rate	8e-7
Micro batch size per GPU	1
Optimizer	Adam
Lr schedule	Constant
zero-stage	3
Global batch size	512
Max prompt length	6144
Max response length	16384
KL loss coefficient	0
Rollout	
Temperature	0.85
N	8
Top-p	0.95
Tok-k	50

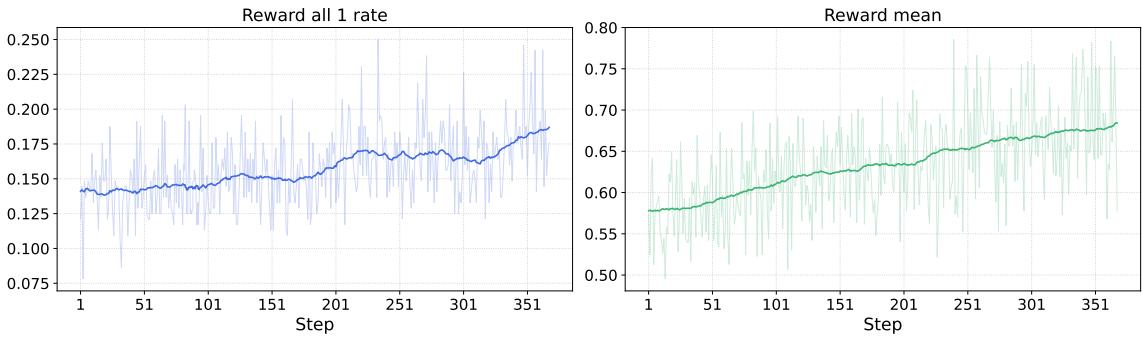


Figure 4: Training dynamics of RL. We show the proportions of all-one rewards and the mean reward value, which increases steadily over the course of training.

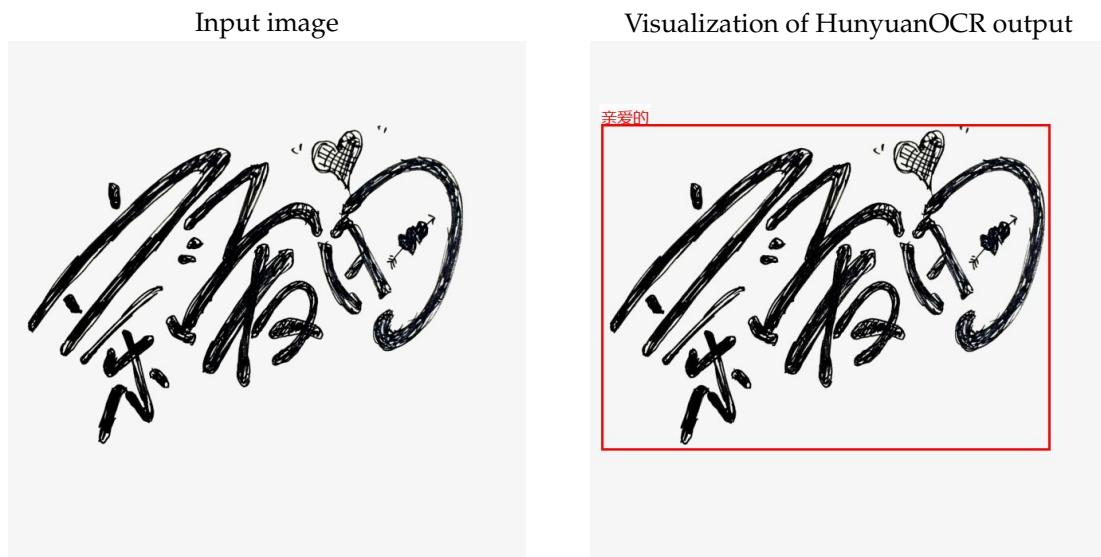
Information Extraction, VQA, and Translation. In addition, the information extraction (IE) task improves by about 2 points, the average score on OCRBench increased by 3.3, and the text image translation task also shows noticeable gains. These results indicate that the LLM-as-a-judge-based reward design can effectively guide the model to produce more faithful and semantically accurate outputs in higher-level understanding tasks.

Discussion. In summary, we attribute the effectiveness of RL in HunyuanOCR primarily to two factors:

- **High-quality training data.** Carefully curated and diverse RL training data provide a solid foundation for the model to learn robust behaviors across spotting, parsing, IE, and translation scenarios.
- **Fine-grained reward design.** Task-specific, fine-grained reward functions (both rule-based and LLM-as-a-judge-based) allow the model to receive precise feedback on multiple aspects of its outputs, leading to balanced improvements in recognition accuracy, structural parsing, and semantic understanding.

These elements work together to make the RL stage an effective complement to supervised training, yielding a more capable and reliable HunyuanOCR model.

D Qualitative examples



Prompt: 检测并识别图片中的文字，将文本坐标格式化输出。

HunyuanOCR: <ref>亲爱的</ref><quad>(22,160),(888,788)</quad>

Figure 5: Robust Text Spotting Results of HunyuanOCR on Artistic Font.



Prompt: 检测并识别图片中的文字，将文本坐标格式化输出。

HunyuanOCR: <ref>22:15</ref><quad>(165,20),(245,33)</quad><ref>HD</ref><quad>(716,19),(737,25)</quad><ref>HD</ref><quad>(770,19),(794,25)</quad><ref>41</ref><quad>(902,23),(926,31)</quad><ref>6</ref><quad>(862,29),(874,34)</quad><ref></ref><quad>(49,61),(72,79)</quad><ref>台湾文献史料丛刊一...Anna's Archive.pdf</ref><quad>(194,61),(925,81)</quad><ref>弁言</ref><quad>(837,281),(883,345)</quad><ref>本用「目」裁，自弘光帝即位至北狩一年期，有二百七十日有事</ref><quad>(663,239),(694,785)</quad><ref>在「晚明史籍考」著，今珍藏抄本整理排印，一罕之。</ref><quad>(710,239),(743,725)</quad><ref>本「偏安排日事」，凡十四卷，不著撰人；按月排日南明弘光朝事。考此書</ref><quad>(757,268),(791,785)</quad><ref>；且所多引用章奏原文，甚具史料值。其中於款北一事，他，亦所</ref><quad>(610,240),(647,785)</quad><ref>>。南都陷後，著有。</ref><quad>(571,240),(600,438)</quad><ref>>不抄本通病，往往「百出」，正；今就整理所，事如下；</ref><quad>(526,267),(555,769)</quad><ref>>（一）卷二「（崇十七年）六月壬戌」日下首端原有衍文「壬戌（此二字，</ref><quad>(476,275),(510,785)</quad><ref>>原本所；添上方）」十六字，今已略去。此由抄粗疏之，姑不必深</ref><quad>(385,238),(415,785)</quad><ref>>究。卷五「九月」末「吏部侍郎捷署部事」及卷六「十月癸未」日下「起原任吏</ref><quad>(378,238),(409,785)</quad><ref>>科都安行、通政司若金俱都察院右都御史」「日」下分別附有按，</ref><quad>(339,238),(370,785)</quad><ref>>明「原本」如何如何，有所移改（文繁不，各）；足所「原本」，</ref><quad>(288,238),(322,785)</quad><ref>>多。因疑上述「原本」，非原本；而今所抄本，或已手抄之本</ref><quad>(237,238),(273,785)</quad><ref>>（二）卷四「（崇十七年）八月癸酉」日下「吏部尙書徐石麒政七款.....」</ref><quad>(150,272),(188,778)</quad><ref>>（卷七「十一月辛丑」日下一亦有，不）。</ref><quad>(198,243),(230,562)</quad><ref>>偏安排日事</ref><quad>(100,301),(126,366)</quad><ref>>目录</ref><quad>(100,923),(151,935)</quad><ref>>打开方式</ref><quad>(323,923),(426,934)</quad><ref>>转存</ref><quad>(598,923),(650,935)</quad><ref>>删除</ref><quad>(848,923),(900,935)</quad>

Figure 7: Robust Text Spotting Performance of HunyuanOCR in Complex Document Scenarios.

Input image

Visualization of HunyuanOCR output

分体落地式空调器	
整机型号	KFR-51LW/DY-JN(E3)
室内机型号	KFR-51L/DY-JN(E3)
室外机型号	KFR-51W-M242
制 冷 量	5100W
制 热 量	5800W+1800W(电加热管)
EER/COP	3.09/3.37
循环风量	1000m ³ /h \$
制冷剂	(见室外机铭牌)
防水等级(室外机)	IPX4
防触电保护类型	I类
质量(室内机/室外机)	39kg/(见室外机铭牌)
端声	室内侧(低风-高风-超强风) 室外侧 54dB(A)
制冷系统允许压力	2.6MPa
热交换器最大工作压力	2.6MPa
额定电压/额定频率	220V~/50Hz
额 定 标 额 冷	7.8A
电、蒸、煮、制 热	8.0A+8.2A
额 定 工 制 冷	1650W
功 率 情 制 热	1720W+1800W
最 大 入 电 流	21.4A
最 大 入 功 率	4000W
电辅助加热输入电流	8.2A
电辅助加热输入功率	1800W
CCC	(见机身条形码)
室内机出厂编号	(见机身条形码)
室内机制造日期	(见机身条形码)
广东美的制冷设备有限公司	广东美的制冷设备有限公司
生产日期 2009年08月	生产日期 2009年08月
机身 编号 6UK4098118033325	机身 编号 6UK4098118033325

Prompt: 提取文档图片中正文的所有信息用markdown格式表示，其中页眉、页脚部分忽略，表格用html格式表达，文档中公式用latex格式表示，按照阅读顺序组织进行解析。

HunyuanOCR: <table><caption></caption><tr><td colspan="5">分体落地式空调器</td></tr><tr><td colspan="2">整机型号</td><td colspan="3">KFR-51LW/DY-JN(E3)</td></tr><tr><td colspan="2">室内机型号</td><td colspan="3">KFR-51L/DY-JN(E3)</td></tr><tr><td colspan="2">室外机型号</td><td colspan="3">KFR-51W-M242</td></tr><tr><td colspan="2">制冷量</td><td colspan="3">5100W</td></tr><tr><td colspan="2">制热量</td><td colspan="3">5800W+1800W(电加热管)</td></tr><tr><td colspan="2">EER/COP</td><td colspan="3">3.09/3.37</td></tr><tr><td colspan="2">循环风量</td><td colspan="3">1000m³/h \$</td></tr><tr><td colspan="2">制冷剂</td><td colspan="3">(见室外机铭牌)</td></tr><tr><td colspan="2">防水等级(室外机)</td><td colspan="3">IPX4</td></tr><tr><td colspan="2">防触电保护类型</td><td colspan="3">I类</td></tr><tr><td colspan="2">质量(室内机/室外机)</td><td colspan="3">39kg/(见室外机铭牌)</td></tr><tr><td colspan="2">端声</td><td colspan="3">室内侧(低风-高风-超强风)
室外侧 54dB(A)</td></tr><tr><td colspan="2">制冷系统允许压力</td><td colspan="3">2.6MPa</td></tr><tr><td colspan="2">热交换器最大工作压力</td><td colspan="3">2.6MPa</td></tr><tr><td colspan="2">额定电压/额定频率</td><td colspan="3">220V~/50Hz</td></tr><tr><td rowspan="2" style="vertical-align: middle;">额定电流制冷 7.8A
制热 8.0A+8.2A制冷 8.0A+8.2A
制热 1650W
额定电压/额定频率 220V~/50Hz额定功率制冷 1720W+1800W
制热 1720W+1800W室内机出厂编号
室内机制造日期
(见机身条形码)广东美的制冷设备有限公司生产日期 2009年08月机身 编号 6UK4098118033325

生产日期2009年08月

机身编号: <pFig>图2</pFig><quad>(452,938),(664,966)</quad>6UK4098118033325

Figure 10: Robust Parsing Performance of HunyuanOCR in Complex Table Scenarios.

Prompt: 提取文档图片中正文的所有信息用markdown格式表示，其中页眉、页脚部分忽略，表格用html格式表达，文档中公式用latex格式表示，按照阅读顺序组织进行解析。

HunyuanOCR: ...

Moore is only barely conscious (hp 1) and can offer no further information. Any additional questions are met with incoherent babbling, even if the heroes heal his wounds. The trapdoor opens onto an underground storage room filled with even more fertilizer, plus a wide array of other chemicals. Any hero who makes a successful DC 15 Knowledge (physical sciences) check realizes that Moore has far more fertilizer than is required for a farm this size. A successful DC 15 Demolitions check reveals that Moore has all the makings for an enormous bomb. The heroes may want to use this material later in the adventure, perhaps to create a bomb to deal with O.S.C.A.R. (see below).

Beyond the Fields

The zombies have blazed a trail of sorts that allows relatively easy travel through the cornfield to O.S.C.A.R.'s bunker. O.S.C.A.R. has already started to process another incantation as the heroes approach. Read or paraphrase the following aloud.

The slimy trail snakes a rambling route through the tall corn, illuminated by an occasional flash of lightning. After a few hundred yards, the corn abruptly parts to reveal a squat concrete building similar to an electrical utility shed. Power lines from the nearby towers stretch to connect with it.

About two dozen yards from the bunker, two humanoid creatures apparently made of metal are standing beside a metallic utility box of some sort. They appear to be repairing something inside.

Suddenly, the dull roar of the thunder is overlaid with an angry buzzing sound, as though someone has disturbed a hornet's nest.

The buzzing sound is a magical side effect of demolish, the next incantation that O.S.C.A.R. is preparing. (This incantation was created with Seed: Destroy. See Chapter 3: Spells in the Urban Arcan Campaign Setting and the New Incantations section at the end of this adventure.) The sound, while loud, has no effect other than to annoy those who hear it.

Creatures: Next to the bunker, two of O.S.C.A.R.'s minion robots are working inside a metal utility box. Any character who makes a successful DC 10 Knowledge (technology) check recognizes it as a utility box for high-speed internet connections. The robots are attempting to restore O.S.C.A.R.'s T3 connection to the outside world.

Minion Robots (2): hp 21, 21. See the new monster description at the end of this adventure for details.

Tactics: The robots need 2 more hours of work to finish repairing the connection. If they are hindered in any way, they turn on the intruders and attack, fighting until they are destroyed.

Development: A DC 15 Spot check reveals a plaque on the side of each robot that reads "Armitage."

A small, concrete bunker serves as the entrance to the O.S.C.A.R. mainframe. The building has no windows, and the metal door is secured with an electronic lock. (Because all electronics are affected by the magical storm, however, the Disable Device check to open it is lower than normal; see below). A small plaque on the front of the building reads, "Property of Armitage Industries. NO TRESPASSING."

A video camera above the door transmits images to O.S.C.A.R. Before the T3 connection was severed, it sent them back to Armitage Industries as well.

Door: Hardness 10, 120 hp, Break DC 35, Disable Device DC 15.

Video Camera: Hardness 5, 2 hp.

1. Entrance

Read or paraphrase the following aloud when the heroes open the door to the bunker.

The door opens to reveal a small antechamber with a steep metal stairwell leading down. The industrial lights in the stairwell flicker and pulse, sometimes even changing colors.

The lights in the stairwell are malfunctioning because of the effects of O.S.C.A.R.'s magic.

Trap: The entire floor of the entrance and first step of the stairwell are covered by a pressure pad that triggers a trap. If weight in excess of 5 kg is placed on the first step, the stairs fold over, creating a steep slope. Anyone on the stairs at this time must make a DC 25 Reflex save or slide at a speed of 50 feet per round, landing in a 30-foot-deep pit that opens just beneath the last step. The fall deals 3d6 points of damage (Reflex DC 13 half). A small hidden switch on the inside of the door (Search DC 25) deactivates the pressure pad. If the trap has been triggered, a flick of the switch restores the stairs to normal and resets the trap. An identical switch at the bottom of the stairs serves the same function.

2. Stairwell

The stairs descend 150 feet into the earth and are quite steep. If the characters deactivate the trap (see above), they can use the stairs normally.

The metal door is secured with a mechanical lock. A plaque on the door reads "O.S.C.A.R." Door: Hardness 10, hp 120, Break DC 35, Disable Device DC 25.

3. Tome Storage

Read or paraphrase the following aloud when the heroes enter this area.

This chamber is lit with dim red bulbs. Lining the walls are floor-to-ceiling metal racks, each filled with row after row of neatly aligned, ancient-looking books, tomes, and scrolls. Set into the far wall is a window that reveals a high-tech mainframe computer in the next room. The lights in that chamber are flickering just as those in the stairwell did. A conveyor belt leads from this chamber into the one with the mainframe.

Four robots stand in the middle of this room facing each other, their manipulator arms moving in complex gestures.

Figure 11: Robust Parsing Performance of HunyuanOCR in Wild-OmniDocBench.



Prompt: 提取图中的文字。

HunyuanOCR: 静观茶舍
JIN GUAN CHA SHE
茶



Prompt: 提取图中的文字。

HunyuanOCR: 鱿鱼游戏



Prompt: 提取图中的文字。

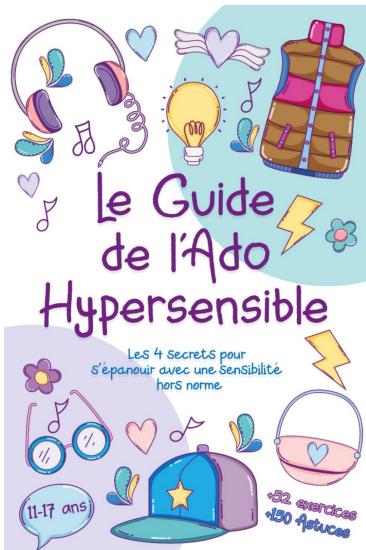
HunyuanOCR: 旗开得勝



Prompt: 提取图中的文字。

HunyuanOCR: 清江浦河庫道
東首天一堂製

Figure 12: Robust Parsing Performance of HunyuanOCR in Diverse Real-World Visual Scenarios.

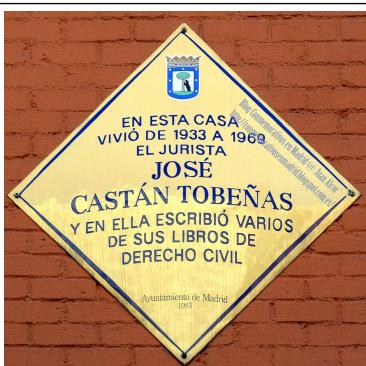


Prompt: 提取图中的文字。

HunyuanOCR: # Le Guide de l'Ado Hypersensible
Les 4 secrets pour s'épanouir avec une sensibilité hors norme
11-17 ans
+52 exercices +150 Astuces

Prompt: 提取图中的文字。

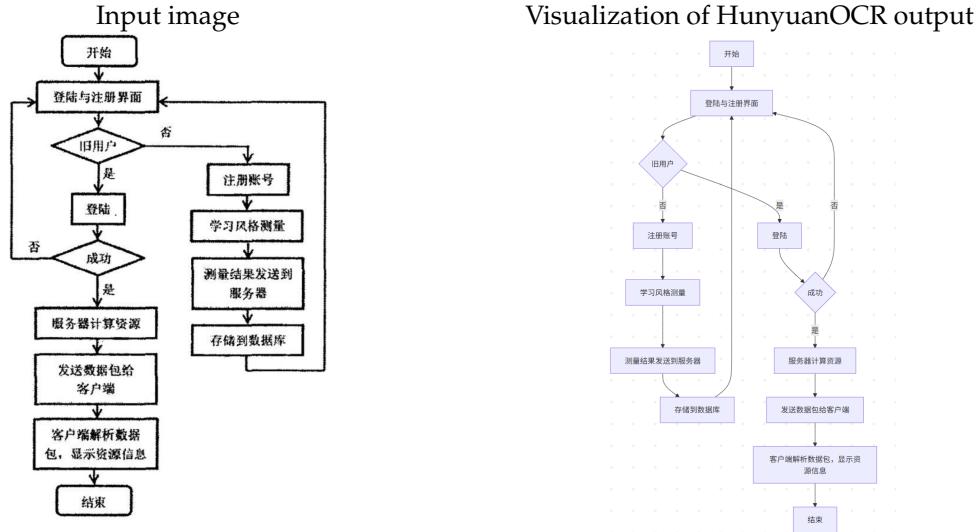
HunyuanOCR:
Professioni & Concorsi
TEORIA e TEST
Concorso per
415 Assistenti Amministrativi (posti elevati a 569)
AGENZIA DELLE DOGANE E DEI MONOPOLI
Manuale e quesiti per la prova scritta unica
Aggiornato alla Legge di bilancio 2025
Diritto costituzionale e amministrativo
Diritto tributario
Diritto dell'UE
Economia politica e Contabilità aziendale
Fini istituzionali, compiti e ordinamento dell'ADM
Normativa in materia di dogane, accise e giochi
Informatica
Lingua inglese
Competenze trasversali
ESTENSIONI ONLINE
Contenuti extra
Software di simulazione
Edises edizioni



Prompt: 提取图中的文字。

HunyuanOCR: EN ESTA CASA VIVIÓ DE 1933 A 1969 EL JURISTA JOSÉ CASTÁN TOBEÑAS Y EN ELLA ESCRIBIÓ VARIOS DE SUS LIBROS DE DERECHO CIVIL
Ayuntamiento de Madrid
1993
Blog Conmemorativos en Madrid @ Juan Alcor
<http://conmemorativosenmadrid.blogspot.com.es/>

Figure 13: Robust Parsing Performance of HunyuanOCR in DocML.



Prompt: 解析图中的图表，对于流程图使用Mermaid格式表示，其他图表使用Markdown格式表示。

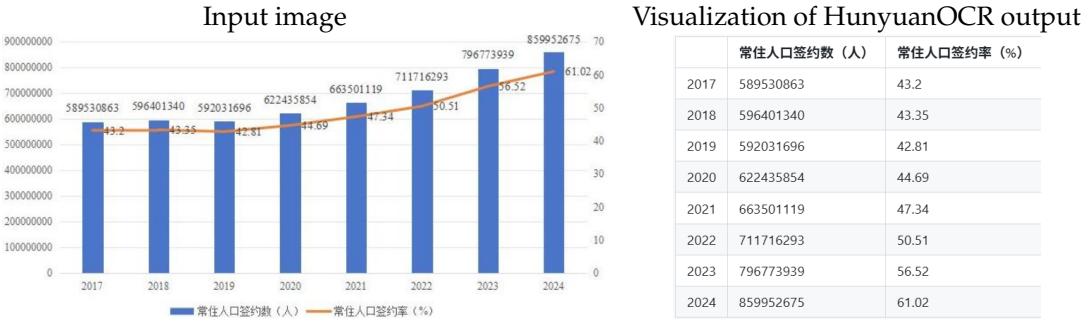
HunyuanOCR:

```

flowchart TD
A["开始"] --> B["登陆与注册界面"]
B --> C{“旧用户”}
C -->|是| D["注册账号"]
C -->|否| E["登陆"]
E --> F["成功?"]
F -->|是| G["服务器计算资源"]
G --> H["发送数据包给客户端"]
H --> I["客户端解析数据包, 显示资源信息"]
I --> M["结束"]
F -->|否| D
D --> E
E --> F
F --> G
G --> H
H --> I
I --> M

```

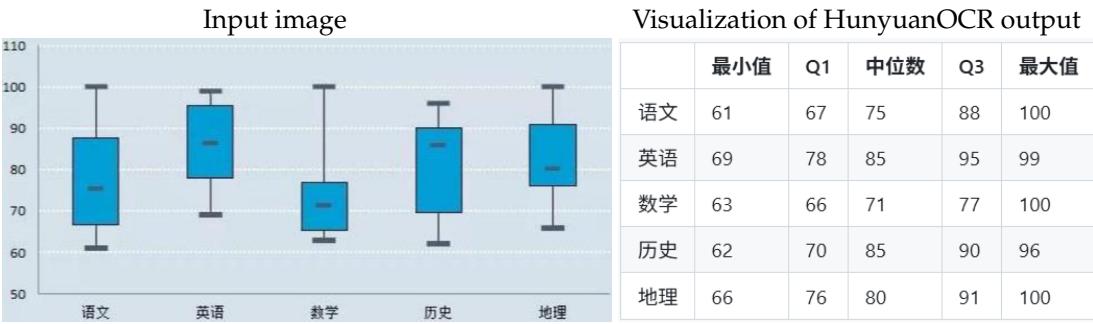
Figure 14: Robust Parsing Performance of HunyuanOCR in Flowchart Scenarios.



Prompt: 解析图中的图表，对于流程图使用Mermaid格式表示，其他图表使用Markdown格式表示。

HunyuanOCR:

	常住人口签约数 (人)	常住人口签约率 (%)
2017	589530863	43.2
2018	596401340	43.35
2019	592031696	42.81
2020	622435854	44.69
2021	663501119	47.34
2022	711716293	50.51
2023	796773939	56.52
2024	859952675	61.02



Prompt: 解析图中的图表，对于流程图使用Mermaid格式表示，其他图表使用Markdown格式表示。

HunyuanOCR:

	最小值	Q1	中位数	Q3	最大值
语文	61	67	75	88	100
英语	69	78	85	95	99
数学	63	66	71	77	100
历史	62	70	85	90	96
地理	66	76	80	91	100

Figure 15: Robust Parsing Performance of HunyuanOCR in Chart Scenarios.



Prompt: 将图中文字翻译为中文。

HunyuanOCR: 经济与商业
印度尼西亚

Prompt: 先解析文档，再将文档内容翻译为中文，其中页眉、页脚忽略，公式用latex格式表示，表格用html格式表示。

HunyuanOCR: ...

翻译成中文为

相对论构型中总角动量的分布

Michel Poirier*

CEA - 巴黎-萨克雷大学，“相互作用、动力学和激光”实验室, CE Saclay, F-91191 Gif-sur-Yvette, 法国

Jean-Christophe Paint

CEA, DAM, DIF, F-91297 Arpajon, 法国

巴黎-萨克雷大学, CEA, 极端条件物质实验室, F-91680 Bruyères-le-Châtel, 法国

(2021年7月6th日星期二)

本文致力于分析相对论构型中总角动量的分布。利用累积量和生成函数形式，该分析可以简化为对具有 N 个等效动量为 j 的电子的单个子壳层进行研究。为 J 分布的生成函数提供了 n 阶导数形式的表达式，并建立了有效的递推关系。结果表明，该分布可以用一种类似Gram-Charlier 的级数来表示，该级数来源于磁量子数分布的相应级数。当构型包含多个子壳层时，此展开的数值效率尚可，而当只涉及一个子壳层时，精度则较差。给出了奇数阶矩的解析表达式，而偶数阶矩则表示为级数，虽然不收敛，但提供了可接受的精度。此类表达式可用于自旋轨道分裂阵列中跃迁数的近似值：结果表明，当保留的项数较少时，该近似通常是有成效的，而某些复杂情况则需要包含大量项。

#I. 引言

为了在恒星物理学或激光等离子体实验（例如惯性约束聚变研究）的背景下模拟热等离子体的发射和吸收光谱特性，需要适当地描述具有多个开放子壳层的多电子构型。特别是，先验地了解两种构型之间的谱线数量具有重要意义。电偶极子(E1)谱线的统计特性由Moszkowski [1]、Bauche 和 Bauche-Arnoult [2] 以及最近由Gilleron 和 Pain [3] 研究。谱线数量是不透明度代码的基石，用于决定是使用上述方法对跃迁阵列进行统计建模，还是需要涉及哈密顿量对角化的详细谱线计算[4, 5]。当跃迁阵列的谱线数量超过特定值时，可以应用部分分辨跃迁阵列方法[6–8] 及其对超构型形式的扩展[9, 10] 等替代方法。电四极子(E2) 谱线的统计特性也得到了研究[11]。

在计数问题中，生成函数技术是一种强大的工具，无论是为了获得解析表达式、推导递推关系还是寻找近似公式。生成函数还能确定累积量，累积量是统计建模的重要组成部分，可以从其中获得矩。在此框架下，我们最近发表了超构型中电子构型数量的解析公式和递推关系[12]，以及基于累积量计算的统计分析。

总角动量多重性的确定最早由核物理学家[13] 在壳模型[14, 15] 框架内进行研究，后来由原子物理学家用于电子构型。考虑一个包含 N 个相同费米子的系统，问题归结为推导它们可以耦合的允许的总角动量 J 。由于泡利不相容原理导致的反对称性，某些 J 值是被禁止的，而另一些则出现多次。正如Condon 和 Shortley 指出的那样，具有角动量 J 的能级数量 $Q(J)$ 等于投影 $M = J$ 的态的数量减去投影 $M = J + 1$ 的态的数量。在1936年的论文[13] 中，Bethe 通过Wigner 型分布对角动量分布进行建模，该分布可以从 $P(M)$ 分布的高斯形式推导出来。多年后，Bauche 和 Bauche-Arnoult 推导出了 $Q(J)$ 和两种电子构型之间谱线数量的解析公式，假设 $P(M)$ 采用四阶Gram-Charlier 展开级数[2]。同期，Hirst 和 Wybourne 使用统计群论来确定Wigner 型近似公式的参数，并指出了与整数划分理论的联系[16]。

Figure 16: Translation Performance of HunyuanOCR.

报告单详情	
检验结果	基本信息
检验日期: 2022.3.20 最终结果以医院纸质报告为准	
乙酰氨基葡萄糖苷酶	结果:阳性(+) 正常范围: 阴性
β葡萄糖醛酸酶	结果:阴性 正常范围: 阴性
唾液酸苷酶	结果:阴性 正常范围: 阴性
白细胞酯酶	结果:阳性(+) 正常范围: 阴性
过氧化氢	结果:阳性(+) 正常范围: 阴性
PH值	结果:4.60 正常范围: 3.8—4.5
乳酸杆菌	结果:少 正常范围: 大量
白细胞 (清洁度)	结果:15~20 / HP (Ⅲ度) 正常范围: ≤15/HP
BV	结果:阴性 正常范围: 阴性
AV	结果:阴性 正常范围: 阴性
念珠菌 (酶法)	结果:+ 正常范围: 阴性
滴虫 (酶法)	结果:阴性 正常范围: 阴性

Prompt: 输出检验日期对应的值

HunyuanOCR: 2022.03.20



Prompt: 提取图片中的: ['单价', '上车时间', '发票号码', '省前缀', '总金额', '发票代码', '下车时间', '里程数'] 的字段内容，并且按照JSON格式返回。

HunyuanOCR: "json {
"单价": "3.00",
"上车时间": "09:01",
"发票号码": "42609332",
"省前缀": "陕",
"总金额": "¥77.10元",
"发票代码": "161002018100",
"下车时间": "09:51",
"里程数": "26.1km"
}"

Figure 17: Information Extraction (IE) Performance of HunyuanOCR on Receipts.



Prompt: 提取图中的字幕。

HunyuanOCR: 对不起我是诺曼.斯佩尔曼
Sorry, I'm Norm. Spellman.



Prompt: 提取图中的字幕。

HunyuanOCR: 知道一点近况也不奇怪



Prompt: 提取图中的字幕。

HunyuanOCR: 它不是一个名词
它是一个动词

Figure 18: Video Subtitle Extraction Performance of HunyuanOCR.

FORM 443 - 10-87 - Page 100

WATER ANALYSIS

Fort Morgan Factory April 5, 1986

Kind of Water: Water in Boilers
 Description of Sample: Sample composited from daily three shift composite during Campaign.
 Source of Water: Water from Boiler Feed Tank.

Source of Water	Parts Per Million	Ibs. per day	Total lbs. used
Boiler H ₂ O Treatment	(Caustic	62.2	6878
	Disodium phosphate	11.3	1252
	Sodium Sulfate	3.7	407
	Quebracho	0.0	0

ANALYSIS

Parts Per Million		HYPOTHETICAL COMBINATIONS	Parts Per Million
Silica SiO ₂	1.8	Silicon Chloride SiCl ₄	1.8
Iron Fe	2.5	Sodium Chloride NaCl	40.4
Calcium Ca	2.6	Sodium Sulfate Na ₂ SO ₄	73.2
Magnesium Mg	1.9	Sodium Carbonate Na ₂ CO ₃	345.7
Sodium Na	20.2	Sodium Phosphate Na ₂ PO ₄	93.8
Chlorine Cl	24.5	Magnesium Chloride MgCl ₂	
Sulfuric Acid SO ₄	49.5	Magnesium Sulfate MgSO ₄	
Carbonic Acid CO ₂	197.4	Magnesium Carbonate MgCO ₃	
Phosphoric Acid PO ₄	65.0	Magnesium Phosphate Mg(PO ₄) ₂	6.8

Organic and Volatile (by difference)

Parts Per Million		Organic and Volatile (by difference)	Parts Per Million
574.7		574.7	
226.3		226.3	

TOTAL DISSOLVED SOLIDS AT 105°C

pH	858.8	TOTAL DISSOLVED SOLIDS AT 105°C	858.8
11.4			11.4

Calculated Hardness

ppm	11.4	ppm	11.4
2.6 x 1.4 = 3.64		Alk.	120.0
1.9 x 2.3 = 4.37		Hard.	4.0
Total	8.01 ppm	PO ₄	65.0
		Solids	858.8

Copies To

General Chemist
 Dist. Superint.
 Dist. Engineer
 Superintendent
 Laboratory

Chief Chemist.

- 7 -

Source: <https://www.industrydocuments.ucsf.edu/docs/gzhy0227>

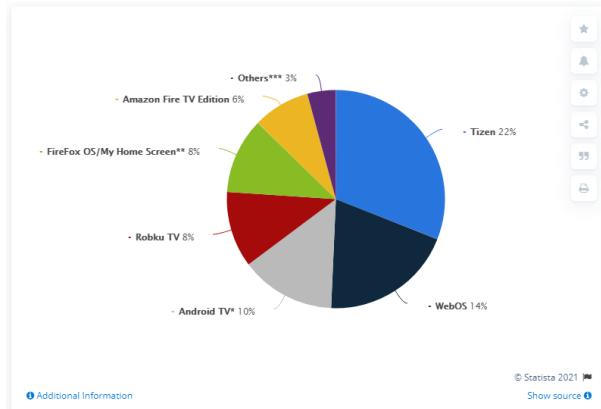
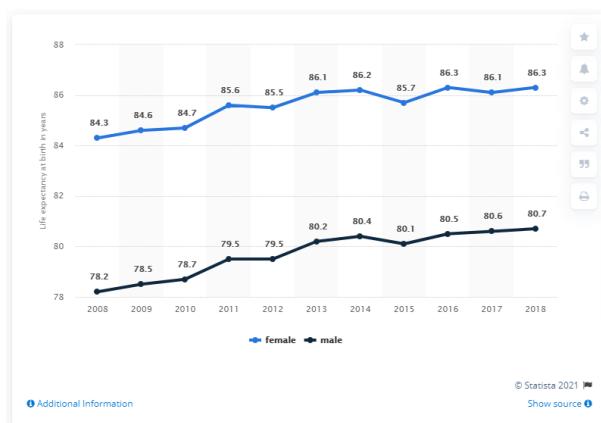


Figure 19: Document & Chart Visual Question Answering (VQA) performance of HunyuanOCR.