

A Survey on MLLM-based Visually Rich Document Understanding: Methods, Challenges, and Emerging Trends

Anonymous ACL submission

Abstract

001 Visually Rich Document Understanding
002 (VRDU) has become a pivotal area of re-
003 search, driven by the need to automatically
004 interpret documents that contain intricate
005 visual, textual, and structural elements. Re-
006 cently, Multimodal Large Language Models
007 (MLLMs) have demonstrated significant
008 promise in this domain, including both
009 OCR-based and OCR-free approaches for
010 information extraction from document im-
011 ages. This survey reviews recent advances in
012 MLLM-based VRDU, highlighting emerging
013 trends and promising research directions with
014 a focus on two key aspects: (1) techniques for
015 representing and integrating textual, visual,
016 and layout features; (2) training paradigms,
017 including pretraining, instruction tuning, and
018 training strategies. Moreover, we address
019 challenges such as data scarcity, handling
020 multi-page and multilingual documents,
021 and integrating emerging trends such as
022 Retrieval-Augmented Generation and agentic
023 frameworks. Our analysis offers a roadmap for
024 advancing MLLM-based VRDU toward more
025 scalable, reliable, and adaptable systems.

026 1 Introduction

027 Visually-Rich Document Understanding (VRDU)
028 lies at the intersection of vision and language, aim-
029 ing to extract and understand information from
030 documents with multiple data modalities and com-
031 plex layouts (Park et al., 2019; Ding et al., 2023).
032 With the rapid digitization of physical documents
033 and the widespread use of structured and semi-
034 structured digital documents, the development of
035 robust, generalizable VRDU frameworks has at-
036 tracted significant attention for automating infor-
037 mation extraction, improving accessibility, and en-
038 hancing decision-making across diverse domains
039 such as finance, healthcare, and education.

040 Early VRDU frameworks relied on manu-
041 ally crafted rules and domain-specific heuristics

(Watanabe et al., 1995; Seki et al., 2007), which
experienced a sudden performance drop on unseen
documents across domains or with diverse lay-
outs. Conventional deep learning approaches em-
ployed CNNs (Katti et al., 2018; Yang et al., 2017)
and RNNs (Denk and Reisswig, 2019) to lever-
age visual or textual features, facilitating more in-
formative representations. However, these meth-
ods typically do not effectively integrate the di-
verse modalities in documents, limiting their ca-
pacity to capture the rich semantic structure inher-
ent in visually rich documents. With the success of
pretraining techniques in language modelling, nu-
merous VRDU models (Huang et al., 2022; Hong
et al., 2022; Lyu et al., 2024) have been pretrained
on large-scale scanned or PDF document datasets,
enabling more effective fusion of visual, textual,
and layout features for robust multimodal repre-
sentation. However, their effectiveness is con-
strained by the scope and diversity of their pre-
training data, often necessitating substantial fine-
tuning to achieve cross-domain generalizability.

042 Recently, MLLMs (OpenAI, 2024; Liu et al.,
043 2024b), trained on massive visual and linguistic
044 datasets, have demonstrated powerful represen-
045 tational capabilities and extensive world knowl-
046 edge, enabling a deeper understanding of text-
047 dense images with diverse visual appearances and
048 complex spatial layouts. By combining the supe-
049 rior text understanding of LLMs (Touvron et al.,
050 2023) with visual encoders (Dosovitskiy et al.,
051 2020) that capture image content and layout in-
052 formation, MLLM-based VRDU frameworks have
053 demonstrated strong performance across diverse
054 document question-answering and information-
055 extraction tasks, and generalizability across do-
056 mains without task-specific fine-tuning.

057 This paper provides a comprehensive survey
058 of recent developments in MLLM-based VRDU
059 frameworks. Previous surveys have either fo-
060 cused on a broad analysis of the diverse capabil-
061 ities and challenges of MLLMs, or focused on
062 specific tasks like document question-answering
063 or information extraction. This survey provides
064 a comprehensive overview of the state-of-the-
065 art in MLLM-based VRDU, covering the meth-
066 ods, challenges, and emerging trends in this
067 field. We also provide a roadmap for future re-
068 search directions and discuss the potential im-
069 plications of MLLM-based VRDU in various do-
070 mains.

Model	Venue	Tasks	Mod.	LLM Backbone	Vision Encoder	PT	IT	FT	Pages	Prompt In.
OCR-Dependent										
ICL-D3IE (2023)	ICCV	KIE	T, L	GPT-3	-	×	×	×	SP	ICL+Layout
DocLLM (2024a)	ACL	KIE, QA, DC	T, L	Custom	-	✓	✓	×	SP	T+B+Q
LAPDoc (2024)	ICDAR	KIE, QA	T, L	Multiple	-	×	×	×	SP	Rule
LMDX (2024)	ACL	KIE	T, L	Gemini-pro	-	×	×	×	SP	ICL+Layout
ProcTag (2025)	AAAI	QA	T, V, L	GPT-3.5	-	×	×	✓	SP	Rule+CoT
DocKD (2024)	EMNLP	KIE, QA, DC	T, L	Custom	-	×	×	✓	SP	Gen by VL
DoCo (2024)	CVPR	KIE, QA, DC	T, L	Multiple	LayoutLMv3	✓	×	✓	SP	I+Q
InstructDoc (2024)	AAAI	KIE, QA	T, V, L	FlanT5	LayoutLMv3	×	×	✓	MP	I+Q
LayoutLLM (2024)	CVPR	KIE, QA	T, V, L	Vicuna-7B-v1.5	OpenCLIP+CLIP	×	✓	✓	SP	I+Q+CoT
LLaVA-Read (2024c)	preprint	KIE, QA	T, V, L	Vicuna-1.5 13B	Multiple	✓	✓	×	SP	I+Q
LayTextLLM (2024)	ACL	QA, KIE	T, L	Llama2-7B-base	-	✓	×	✓	SP	T+B
DocLayLLM (2024)	CVPR	QA, KIE	T, V, L	Llama2-7B-chat	Pix2Struct-Large	×	✓	✓	SP	I+Q+B
LayTokenLLM (2025b)	CVPR	QA	T, L	Multiple	-	✓	×	×	MP	I+Q+L
GPE (2025a)	ICLR	KIE, QA	T, L	Multiple	-	×	×	✓	SP	T+B+Q
MDocAgent (2025)	preprint	QA	T, V	Multiple	IXC2-VL-4KHD	×	×	×	MP	I+Q
PDF-WuKong (2025)	preprint	QA	T, V	BGE-M3	LayoutLMv3	×	×	✓	MP	I+Q
DocAssistant (2025)	EMNLP	QA	T, V	InternVL2-Chat-2B	InternVL2 ViT	×	×	✓	SP	I+Q
AlignVLM (2025)	Neurips	QA	T, V	LLaMA-3.2 (1B, 3B)	SigLIP-400M	✓	✓	✓	SP	I+Q
DocThinker (2025)	ICCV	QA, KIE	T, V	Qwen2.5-VL (3B, 7B)	Qwen2.5-VL ViT	×	×	✓*	SP	I+Q
OCR-Free										
KOSMOS-2.5 (2023)	preprint	QA, KIE	V	Custom	mPLUG-Owl VE	×	✓	✓	SP	I+Q
mPLUG-DocOwl (2023a)	preprint	QA	V	mPLUG-Owl	mPLUG-Owl VE	×	✓	×	SP	I+Q
UReader (2023b)	EMNLP	QA	V	mPLUG-Owl	mPLUG-Owl VE	×	✓	×	SP	I+Q
TGDoc (2023)	preprint	KIE, QA	V	Vicuna-7B	CLIP-ViT-L/14	×	✓	✓	SP	I+Q+B
UniDoc (2023)	preprint	KIE, QA	V	Vicuna-7B	CLIP-ViT-L/14	×	✓	✓	SP	I+Q+B
DocPedia (2024)	SCIS	KIE, QA	V	Vicuna-7B	Swin Trans.	✓	×	✓	SP	I+Q
HRVDA (2024a)	CVPR	KIE, QA	V	LLama2-7B	Swin Trans.	✓	×	✓	SP	I+Q
Vary (2024)	ECCV	QA, DocRead	V	Multiple	CLIP, ViTDet	✓	×	✓	SP	I+Q
mPLUG-DocOwl1.5 (2024)	EMNLP	KIE, QA	V	mPLUG-Owl2	mPLUG-Owl2 VE	×	✓	✓	SP	I+Q
HVFA (2024)	Neurips	QA, Cap.	V	Multi (BLIP-2, etc.)	ViT/L-14	×	×	×	SP	I+Q
Texthawk (2024a)	preprint	QA	V	InternLM-XC	ViT	×	✓	✓	SP	I+Q
Texthawk2 (2024b)	preprint	OCR, Grd, QA	V	Qwen2-7B-Instr	SigLIP-SO400M	×	✓	✓	MP	I+Q+Task
TextMonkey (2024c)	preprint	KIE, QA	V	Qwen-VL	ViT-BigG	×	×	×	SP	I+Q
Llavar (2024d)	preprint	QA	V	Vicuna-13B	CLIP-ViT-L/14	×	✓	✓	SP	I+Q
TokenCorrCompressor (2024b)	preprint	QA, Cap.	V	LLaMA-2	CLIP-ViT/L/14	×	×	✓	SP	I+Q
DocKyllin (2024a)	AAAI	QA	V	Llama2-7B-chat	Donut-Swin	×	×	✓	SP	I+Q
Marten (2025b)	CVPR	QA	V	InterLM2	InternViT-300M	×	✓	✓	SP	I+Q
PP-DocBee (2025)	preprint	QA	V	Qwen2-VL-2B	ViT	×	×	✓	SP	I+Q
mPLUG-DocOwl2 (2025)	ACL	KIE, QA	V	mPLUG-Owl2	ViT	✓	×	✓	MP	I+Q
TokenFD (2025)	ICCV	QA, KIE	V	InternLM (2B, 8B)	ViT	✓	✓	✓	SP	I+Q

Table 1: Comparison of existing MLLM-based VRDU frameworks. Mod.: Input modality; KIE: Key Information Extraction; QA: Question Answering; DC: Document Classification; T: Text; L: Layout; V: Vision; MP: Multi-Page; SP: Single Page; I: Image; Q: Question; B: Bounding Box; CoT: Chain of Thought; Cap.: Captioning; Grd.: Grounding; Task: Task Information; VL: Vision-Language.

ities of MLLMs (Caffagni et al., 2024) or examined techniques applied to specific document understanding tasks, such as document layout analysis (Binmakhshen and Mahmoud, 2019), question answering (Barboule et al., 2025), and relation extraction (Delaunay et al., 2023). A recent study provides (Ding et al., 2025b) an overview of deep learning-based frameworks for VRDU but lacks a systematic perspective on MLLM-based approaches. In contrast, this paper provides an analysis of the MLLM-based VRDU frameworks from the aspects of **Framework Architecture** that covers both OCR- and OCR-free models (Sec 2), **Multimodal Representation** (Sec 3), **Training Strategies** (Sec 4), and **Inference Prompt Setting** (Sec 5). We also include a detailed discussion of the challenges of VRDU and provide a critical analysis of the trend and future directions (Sec 6). Notably, this survey is limited to methods that leverage MLLMs for document-level understand-

ing, excluding multi-document applications, non-LLM-based methods, and MLLMs without VRD-specific adaptations.

2 Framework Architecture

General MLLM for VRDU. Many closed- (Team et al., 2024) and open-source (Chen et al., 2024) general-domain MLLMs have been widely adopted for VRDU tasks and have demonstrated promising performance¹. However, the text-dense, visually rich, and layout-sensitive nature of VRDs exposes fundamental limitations of general-domain MLLMs when applied to VRDU, including weak layout inductive bias, sensitivity to OCR noise, and hallucination on these knowledge-intensive tasks. Moreover, the wide range of downstream VRDU applications necessitates specialized techniques that adapt existing LLM back-

¹Refer to Appendix D see performance analysis.

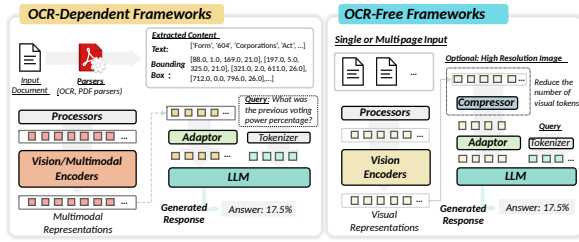


Figure 1: General OCR-dependent and OCR-free framework architectures.

bones (as shown in Figure 1) through VRDU-specific multimodal representations, training objectives, and inference paradigms. In addition, as VRDU tasks are often knowledge-intensive and safety-critical, locally tuning open-source general-domain LLMs on private document collections is essential for practical deployment in sensitive domains such as finance and industrial applications.

OCR-Dependent Frameworks. As shown in Figure 1, OCR-dependent frameworks leverage off-the-shelf tools to extract textual and layout information from scanned or PDF documents. This extracted data, in combination with the document image, is typically fed into multimodal encoders to generate joint representations. Some models (Wang et al., 2024a; He et al., 2023) input the extracted text directly into LLMs, while others (Luo et al., 2024; Zhu et al., 2025a) incorporate visual (Dosovitskiy et al., 2020) or multimodal encoders (Huang et al., 2022) to project those cues into language space via various adaptors or projects. These systems rely on external tools to capture structural information without extensive pretraining (e.g., text recognition). However, reliance on OCR or parsing tools can introduce cumulative errors, especially in handwritten or low-quality scanned documents, hindering the development of fully end-to-end models. Additionally, using low-resolution inputs may reduce the expressiveness of document representations, limiting the overall performance.

OCR-Free Frameworks. OCR-free approaches have been introduced for end-to-end VRD understanding tasks. These frameworks bypass text extraction by directly processing document images. Visual features are extracted via one or more vision encoders, fused with the user query, and decoded by an LLM to generate responses. Representative models include Donut (Kim et al., 2022), mPLUG-DocOwl (Ye et al., 2023a), and URe-

ader (Ye et al., 2023b). Accurate comprehension of fine-grained text in these OCR-free settings requires high-resolution images, which, in turn, lead to lengthy visual sequences requiring visual compression modules (Liu et al., 2024a; Hu et al., 2025). Moreover, effective text recognition in these models often relies on large-scale pretraining or instruction-tuning to integrate textual and layout features via tasks such as text spotting (Liu et al., 2024c) and image captioning (Feng et al., 2024). This paradigm, however, demands substantial dataset construction and considerable computational resources, posing practical challenges.

3 Multimodal Representation

3.1 Text Modality

OCR-dependent methods rely on external tools to extract text for encoding, while OCR-free models use document images directly, treating text as a learning target.

Text Encoding via LLM. Given the frequent text recognition challenges faced by MLLMs, stemming from low-resolution inputs or undertrained vision encoders, off-the-shelf OCR-extracted text is commonly embedded directly into LLM prompts to enhance document comprehension (Wang et al., 2024a; Kim et al., 2024) (see Figure 2). However, the extracted content is often unordered; to address this, frameworks such as ICL-D3IE (He et al., 2023) and LLaVA-Read (Zhang et al., 2024c) employ the XY-cut algorithm to reorder the text sequence. Additionally, to handle long documents, some methods segment the text into chunks, though this may introduce semantic discontinuities (Xie et al., 2025). In sum, directly adding extracted text to prompts improves context and reduces reliance on additional encoders; however, performance remains limited by OCR and LLM errors and weak multimodal integration.

Text Encoding via Auxiliary Encoder. To enhance multimodal integration, many frameworks introduce auxiliary encoders to enhance text embeddings. Several methods (Luo et al., 2024; Zhu et al., 2025a) enhance text representation and multimodal fusion by feeding extracted text, image patches, and bounding boxes into pretrained LayoutLMv3 (Huang et al., 2022). Notably, Zhu et al. (2025a) propose a ROI Aggregation module that aggregates fine-grained tokens (e.g., words)

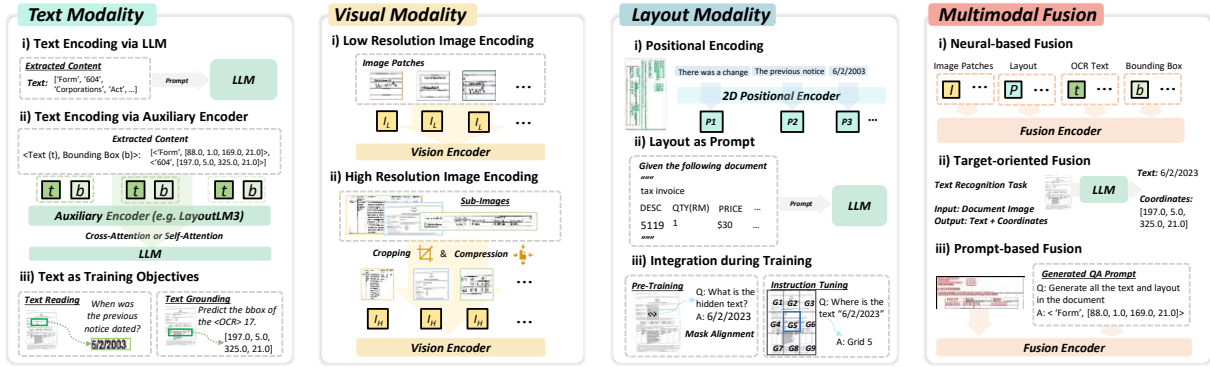


Figure 2: Multimodal feature representation and fusion mechanisms.

into object-level features (e.g., paragraphs), facilitating downstream object-level contrastive learning. Instruct-Doc (Tanaka et al., 2024) introduces an enhanced Q-Former (Li et al., 2023), termed *Document Former*, serving as a bridging module that integrates visual, textual, and layout information from document images into the LLM input space via cross- and self-attention. In sum, external encoders improve representations but require additional pretraining and fine-tuning to align with LLMs’ latent spaces.

Text as Training Objectives. Some frameworks rely exclusively on document images as input to predict answers. Models such as mPLUG-DocOwl (Ye et al., 2023a) and LLaVA-R (Zhang et al., 2024d), built upon mPLUG-Owl (Ye et al., 2023c), demonstrate strong OCR capabilities and are further instruction-tuned on diverse VRDU benchmarks. Other approaches incorporate text recognition, detection, and spotting tasks (Wang et al., 2023; Feng et al., 2023) to integrate text information. To better understand the hierarchical structure of documents, Hu et al. (2024, 2025) propose a multi-grained text localization task spanning the word-to-block level. While these methods deliver robust results using only visual inputs, they place heavy demands on pretraining and fine-tuning. Additionally, high-resolution images are often necessary to accommodate extremely long visual sequences and to preserve fine-grained features (Liu et al., 2024a; Yu et al., 2024a).

3.2 Visual Modality

To integrate visual information, OCR-dependent frameworks use extracted text and coarse visual cues, thereby enabling the use of **lower-resolution** images. In contrast, OCR-free frameworks require direct text recognition, demanding fine-grained

perception and **high-resolution** inputs. See the Appendix A.4 for input resolution details.

Low Resolution Image Encoding. Some frameworks directly feed image patches into pretrained vision encoders to obtain patch embeddings (Xie et al., 2025; Tanaka et al., 2024). Others (Han et al., 2025; Luo et al., 2024; Liao et al., 2024) employ pretrained VRDU models, i.e., LayoutLMv3 (Huang et al., 2022), to extract multimodal-enhanced visual embeddings. Due to the limitations of low-resolution inputs in capturing fine-grained details, recent works have adopted dual-encoder architectures that process both low- and medium-resolution images (Ye et al., 2023b; Zhang et al., 2024c), followed by visual feature compression techniques to manage the increased feature volume. While using low-resolution images offers a straightforward pathway to multimodal understanding, achieving effective alignment often requires additional pretraining and instruction tuning. Moreover, the absence of fine-grained visual detail often necessitates additional OCR tools to extract text for accurate VRD interpretation.

High Resolution Image Encoding. To capture fine-grained level information for end-to-end training and inference, many frameworks support high-resolution image input. For ViT-style (Dosovitskiy et al., 2020) pretrained vision encoders, Hu et al. (2024) splits high-resolution images into pre-defined sub-images. To handle images of various shapes, UReader (Ye et al., 2023b) introduces a *Shape-Adaptive Cropping Module* that adaptively divides images into fixed-size sub-images using grids of various shapes. However, the image cropping may disrupt semantic continuity across sub-images. To address this, Liu et al. (2024c) in-

roduced a *Shifted Window Attention* to enhance cross-sub-images connection via self-attention. In short, high-resolution images support fine-grained information extraction, but efficiently processing the resulting large number of visual tokens remains challenging, requiring a balance between resource usage and the number of visual tokens.

Visual Feature Compression. Yu et al. (2024a,b) utilize Q-Former (Li et al., 2023), while Liu et al. (2024c) adopts the *Resampler* from Qwen-VL (Wang et al., 2024b) to reduce the number of visual tokens. Considering the layout-aware nature of VRDs, Hu et al. (2024) introduces a convolutional module that preserves layout by compressing horizontal features and reducing the number of tokens. It further enhances this with layout-aware cross-attention to handle multi-page input. Liu et al. (2024a) use a *Content Detector* to filter non-informative tokens by segmenting text-rich regions, while Zhang et al. (2024a) propose eliminating low-information areas and clustering and aggregating the remaining features.

3.3 Layout Modality

Unlike natural scene images, VRDs feature dense text and complex layout structures. Methods for encoding layout information can be categorized into positional encoding-based, prompt-based, and task-oriented approaches.

Positional Encoding. OCR-dependent models use OCR tools to extract textual and layout information, combining text embeddings with 2D positional encodings (Xu et al., 2020) to incorporate layout into LLMs (Han et al., 2025; Tanaka et al., 2024). However, these approaches require extra training for feature alignment. In contrast, Zhu et al. (2025a) assigns unique positional embeddings to attention heads based on multi-dimensional layout features without altering the model architecture or requiring further pretraining. Wang et al. (2024a) treats layout as a separate modality and introduces disentangled spatial attention for cross-modal interactions without visual encoders. Zhu et al. (2025b) addresses long-context inference limits by encoding layout as a single token sharing the position with its text. However, these methods implicitly integrate layout information and rely heavily on large-scale pretraining, resulting in high computational costs and reduced effectiveness for tasks that demand explicit layout understanding.

Layout as Prompt. To integrate explicit layout information, some frameworks include layout details in prompts alongside the user query and document content. He et al. (2023) introduces an in-context learning based approach to incorporate layout-aware demonstrations into bounding box representations. Lamott et al. (2024) and Perot et al. (2024) encode layout into text sequence through rule-based verbalization or quantized coordinate tokens. These methods enable layout-awareness without training. However, these methods increase input length, rely on LLMs to interpret layout as text, and overlook visual cues essential for encoding relative positional information.

Integrating During Training. OCR-free frameworks incorporate text by formulating recognition and detection tasks that also aid in understanding layout (Wang et al., 2023; Feng et al., 2023). To further enhance this, some models (Wang et al., 2025b; Zhang et al., 2024c) leverage layout-aware pretraining tasks (Section 4.1) and layout-specific instruction-tuning tasks, such as visual grounding (Liu et al., 2024a,c) and table reconstruction (Liao et al., 2024). However, these methods typically require large-scale datasets for pre-training or instruction tuning, leading to substantial computational costs and data bottlenecks.

3.4 Multimodal Fusion

We categorize multimodal fusion methods into four types: direct, neural-based, task-oriented, and prompt-based. Direct fusion relies on simple feature summation or concatenation with alignment training, while this survey primarily focuses on the latter three approaches in MLLM-based VRDU frameworks.

Neural-based Fusion. The simplest multimodal feature encoding uses external document encoders such as LayoutLMv3 (Xu et al., 2021), which fuse multimodal features via self- or cross-attention and leverage pretraining knowledge. Wang et al. (2024a) stands out by employing a layout-aware transformer with disentangled attention over text and spatial layouts, enabling effective document understanding without requiring image encoders. In OCR-free frameworks, visual encoders extract visual cues, with adaptors like LoRA (Yu et al., 2024b) or linear projectors (Zhang et al., 2024d; Wang et al., 2023) mapping features into the language space. Masry et al. (2025) propose a method

that maps visual features to a weighted textual embedding to reduce misalignment issues observed in previous approaches. These neural-based fusion methods benefit from dedicated encoders or modified architectures, but often require extensive pretraining or SFT and face challenges in scalability, computational overhead, and adaptability to diverse document layouts, especially in noisy OCR scenarios.

Target-oriented Fusion. Target-oriented strategies establish multimodal connections through supervised objectives that span the input-to-output space (Hu et al., 2024) and are widely applied to text and layout features in OCR-free frameworks. For instance, in text recognition tasks, models are trained to map visual features directly to text and spatial coordinates, thereby aligning fusion with task-specific goals. While these approaches improve end-to-end multimodal integration, they also increase demands on data preparation, annotation quality, and training complexity in practice.

Prompt-based Fusion. Prompts for multimodal tasks may include text, images, and bounding box coordinates. While many frameworks adopt Layout-as-Prompt strategies to encode layout information, others use Chain-of-Thought (CoT) reasoning to further enhance multimodal learning. For example, Luo et al. (2024) utilizes a *Layout-CoT* approach that divides reasoning into question analysis, region localization, and answer generation, explicitly modeling spatial layout. Liao et al. (2024) leverages CoT pretraining and CoT annealing to support layout-aware reasoning for VRDU. However, these methods often depend on predefined reasoning strategies, intermediate-step evaluations, and well-trained prior frameworks, limiting their generalizability to unseen domains.

4 Training Paradigms

To facilitate multimodal understanding, instruction following, and domain adaptation, various training tasks and strategies have been developed.

4.1 Pretraining

To enhance mono- and multi-modal document understanding, VRDU frameworks adopt various self-supervised pretraining tasks, such as masked information modeling and cross-modality alignment (Ding et al., 2025b). OCR-dependent frameworks typically utilize pretrained VRDU mod-

els or vision encoders to obtain enriched multimodal representations. Some models propose additional self-supervised learning tasks (e.g., Li et al. (2024) applies object-level contrastive learning between visual and multimodal features). Wang et al. (2024a) introduces a transformer architecture with disentangled spatial-text attention to perform block-wise text infilling to enhance text-layout correlation modeling. OCR-free frameworks (Zhang et al., 2024c; Hu et al., 2024) focus on pretraining tasks like text recognition, detection, and captioning to integrate text and layout information. Hu et al. (2025) further targets multi-page layout coherence. Feng et al. (2024) aligns frequency features with LLMs through text-centric pretraining. Although these self-supervised tasks are effective in fusing multimodal features and learning general knowledge, they remain computationally intensive and often lack instruction-based tuning, limiting their capacity to follow real-world user instructions.

4.2 Instruction Tuning

To benefit task orientation in LLM-based frameworks, many VRD approaches, following InstructGPT (Ouyang et al., 2022), are trained on instruction-response pairs to better align model outputs with user prompts. Pretraining tasks such as text reading, recognition, and image captioning are reformulated as instruction-based formats, with images paired with task descriptions. Beyond improving multimodal fusion, goal-oriented tasks, including VRD question answering (Ding et al., 2024b), key information extraction (Ding et al., 2023), and VRD classification (Harley et al., 2015), are conducted on large-scale datasets. For better generalizability, some frameworks synthetically generate large instruction-tuning datasets (See Appendix B for more details). To further improve localization and information extraction, Wang et al. (2023) and Feng et al. (2023) propose predicting answers alongside bounding boxes, thereby enhancing the framework’s reliability. Instruction tuning not only strengthens user query understanding but also boosts multimodal fusion. Instruction tuning on large-scale datasets substantially enhances zero-shot performance. However, the requirement for extensive training data leads to substantial resource consumption. Furthermore, synthetic datasets, often generated with off-the-shelf OCR tools and LLMs, may yield low-quality QA pairs, particularly in

481	low-resource domains such as scanned documents,	learned knowledge (Yu et al., 2024b; Zhang et al.,	531
482	thereby impacting zero-shot performance.	2024d). Feng et al. (2024) use a Swin Trans-	532
483	4.3 Training Strategies	former to encode frequency-domain images, pre-	533
484	MLLM-based document understanding frame-	trained from scratch. To enhance multimodal fea-	534
485	works typically consist of multiple sub-modules to	ture learning, Li et al. (2024) make the ViT en-	535
486	encode multimodal information and are trained in	coder trainable while freezing LayoutLMv3, en-	536
487	a stepwise manner. Few frameworks leverage in-	abling knowledge distillation via contrastive learn-	537
488	context learning (He et al., 2023) or multimodal	ing. During instruction tuning, vision encoders	538
489	prompts (Perot et al., 2024) to develop training-	are typically unfrozen to improve alignment and	539
490	free architectures. The majority, however, involve	task-specific adaptation (Zhang et al., 2024a; Liu	540
491	pretraining to capture general-domain knowledge,	et al., 2024a). Conversely, in dual-encoder frame-	541
492	followed by instruction tuning to improve inter-	works, vision encoders with inputs at diverse res-	542
493	pretation of user prompts. Furthermore, some	olutions are often frozen to enhance the represen-	543
494	frameworks are subsequently Supervised Fine-	tation of hierarchical inputs. In supervised fine-	544
495	Tuned on benchmark datasets (Wang et al., 2024a;	tuning, there is no standard practice for encoder	545
496	Zhu et al., 2025a) or a synthetic set (Kim et al.,	trainability.	546
497	2024) to enhance domain-specific adaptation. To	Projectors and Adaptors. They play a crucial	547
498	integrate multimodal information, these frame-	role in feature alignment and lightweight tun-	548
499	works mainly employ an LLM with various multi-	ing. Projectors are typically employed to align vi-	549
500	modal encoders (Han et al., 2025; Xie et al., 2025),	sual or layout features with the LLM input space	550
501	sometimes incorporating adaptors (Hu et al., 2024;	(Park et al., 2024) and encode layout informa-	551
502	Lu et al., 2024) or linear projectors (Park et al.,	tion (Tanaka et al., 2024). These modules are	552
503	2024) for fusion or alignment. Depending on the	mainly trainable throughout the entire training	553
504	training stage, sub-modules may be either train-	process. Adaptors, on the other hand, are designed	554
505	able or frozen, balancing the acquisition of new	for efficient, task-specific tuning, often leveraging	555
506	knowledge with the preservation of valuable infor-	LoRA-style updates (Ye et al., 2023a; Hu et al.,	556
507	mation from the original backbone.	2024) or cross-attention mechanisms (Liu et al.,	557
508	LLM Backbone. As most LLMs are exten-	2024c; Yu et al., 2024a) to integrate multi-aspect	558
509	sively pretrained on large-scale datasets and cap-	inputs with minimal parameter changes. Plug-	559
510	ture broad knowledge, many frameworks freeze	and-play components, such as visual abstractors	560
511	the LLM, using it solely to generate human-	(Ye et al., 2023a) or compressors (Hu et al., 2025),	561
512	understandable outputs. In frameworks involv-	have also been introduced to reduce the dimen-	562
513	ing pretraining or instruction tuning (Zhang et al.,	sionality of visual features. These adaptors are	563
514	2024a; Liu et al., 2024a), freezing the LLM back-	usually trained during instruction tuning or during	564
515	bone helps preserve its knowledge and reduce	supervised fine-tuning.	565
516	training costs. However, some approaches en-	4.4 Training Datasets Overview	566
517	able LLMs to be trained during continued pretrain-	Diverse datasets are required to meet specific	567
518	ing (Zhu et al., 2025b) or instruction tuning (Liao	training objectives. Pretraining typically leverages	568
519	et al., 2024) to better capture VRD domain knowl-	large-scale cross-domain VRD collections to re-	569
520	edge and enhance multimodal alignment. In su-	duce domain gaps and enhance multimodal fusion,	570
521	perervised fine-tuning stages, the LLM backbone is	sometimes requiring more domain-specific data	571
522	typically made trainable to adapt to the target do-	(e.g., medical, slides) to improve domain aware-	572
523	main (Zhang et al., 2024d).	ness. For instruction tuning, synthetic datasets	573
524	Vision/Multimodal Encoders. They are em-	are often used to strengthen instruction-following	574
525	ployed to encode multimodal features, which are	and reasoning abilities, particularly in OCR-free	575
526	subsequently aligned with LLM text representa-	frameworks that generate instruction-aligned OCR	576
527	tions by projectors or adaptors. Similar to LLM	or layout understanding tasks. Additionally, SFT	577
528	backbones, vision (Dosovitskiy et al., 2020), and	is commonly applied on original or post-processed	578
529	multimodal encoders (Huang et al., 2022) are	benchmark datasets (e.g., converting key-value	579
530	often kept frozen during pretraining to preserve	pairs into QA format) to further boost perfor-	580

mance. For more dataset details, see Appendix C.

5 Inference Prompt Setting

MLLM-based frameworks adopt diverse prompt formats depending on their architecture. For OCR-free frameworks in Table 1, the prompt typically includes a document image, occasionally multiple pages (Hu et al., 2025; Wang et al., 2025b), alongside a textual user query. Some frameworks not only predict answers to user queries but also localize bounding boxes, often requiring an additional prompt for localization (Wang et al., 2023; Feng et al., 2023). OCR-dependent frameworks first preprocess input using off-the-shelf tools to extract textual and layout information. Vision-free models (He et al., 2023; Wang et al., 2024a) process only the extracted content alongside the query. In contrast, vision-dependent models also incorporate the document image into the vision (Xie et al., 2025) or into multimodal encoders (Liao et al., 2024), aligning visual and textual features for the final prediction. Furthermore, some frameworks integrate layout information into prompts via bounding boxes (Zhu et al., 2025a) or markdown-style formatting. The inference strategies are closely tied to the model architecture and reflect a growing trend toward unified, multimodal understanding and layout-aware reasoning to improve document comprehension accuracy and versatility.

6 Challenges and Future Direction

Synthetic Data. Acquiring high-quality, manually curated datasets for new document collections is often quite costly. Leveraging synthetically generated datasets offers a cost-effective alternative for adapting to the target domain (Ding et al., 2025a). For large-scale instruction-tuning, many frameworks generate instruction-response pairs using benchmarks, templates, or LLMs. However, these synthetic datasets often lack validation, resulting in low-quality or inaccurate pairs. Since synthetic data may not fully capture real user input, future research should prioritize human-in-the-loop and reinforcement learning approaches to improve authenticity and task relevance.

Long Document Understanding. In practice, VRDs frequently span multiple pages; however, most existing frameworks are tailored for single-page inputs. Multi-page approaches typically rely

on retrievers to identify relevant pages, which are then processed by MLLM-based VRDU systems. These methods often fall short of capturing semantic and logical dependencies among document entities, resulting in incomplete contextual understanding. Furthermore, handling long input sequences remains challenging, as existing multi-page benchmarks focus mainly on extractive tasks and rarely support complex multi-hop or multimodal reasoning.

Multilingual VRDU. Most existing models and benchmarks remain heavily English-centric, limiting their generalization to documents with diverse languages and layouts. This bias is further amplified by large-scale pretraining corpora that predominantly reflect English document structures, leading to performance degradation in low-resource settings. Although few multilingual datasets have been proposed (Xu et al., 2022; Chen et al., 2025), future research should explore more multilingual and culturally diverse benchmarks, language-agnostic representation learning, and hybrid approaches to mitigate linguistic bias to handle real-world document diversity.

Effective RAG Framework. While RAG has become a common paradigm (Jain et al., 2025; Faysse et al., 2025), existing approaches often exhibit brittle retrieval due to layout ambiguity and misaligned multimodal embeddings, leading to unreliable evidence selection. Moreover, most RAG pipelines decouple retrieval from reasoning and remain largely text-centric, limiting their ability to capture spatial and visual semantics in complex documents. Future work should explore multimodal RAG frameworks that support iterative reasoning and dynamic evidence refinement, and enable more robust and interpretable VRDU.

Agentic LLM in VRDU. Recent works (Han et al., 2025; Sun et al., 2025) incorporate external tools (e.g., PDF parsers or retrievers) to generate intermediate outputs, enhancing both the accuracy and interpretability of practical VRDU applications. However, future research should explore a wider variety of agent types and architectural innovations to enable automatic handling of diverse formats, cross-domain scenarios, and fine-grained elements such as charts and tables. Additionally, challenges in agentic AI, such as multi-agent coordination and knowledge conflicts, remain significant barriers to broader adoption for VRDU.

679 Limitations

680 While this survey offers a comprehensive
681 overview of MLLM-based VRDU research, our
682 analysis is necessarily qualitative. It does not
683 provide exhaustive head-to-head comparisons, as
684 the field’s rapid evolution and breadth prioritize
685 trend summarization over detailed benchmarking.
686 Although academic advances are thoroughly
687 reviewed, discussion of real-world deployments
688 and industrial challenges remains limited, in
689 part because many practical applications are
690 proprietary and unpublished. In future work, we
691 aim to provide more quantitative meta-analyses,
692 incorporate insights from industrial adoption, and
693 continuously update the survey to capture the
694 latest developments as the field progresses.

695 References

696 Camille Barboule, Benjamin Piwowarski, and Yoan
697 Chabot. 2025. Survey on question answering over
698 visually rich documents: Methods, challenges, and
699 trends. *arXiv preprint arXiv:2501.02235*.

700 Galal M Binmakhshen and Sabri A Mahmoud. 2019.
701 Document layout analysis: a comprehensive survey.
702 *ACM Computing Surveys (CSUR)*, 52(6):1–36.

703 Davide Caffagni, Federico Cocchi, Luca Barsellotti,
704 Nicholas Moratelli, Sara Sarto, Lorenzo Baraldi,
705 Marcella Cornia, and Rita Cucchiara. 2024. The
706 revolution of multimodal large language models: A
707 survey. In *Findings of the Association for Computa-*
708 *tional Linguistics: ACL 2024*, pages 13590–13618.

709 Ketong Chen, Yuhao Chen, and Yang Xue. 2025. Mo-
710 saicdoc: A large-scale bilingual benchmark for vi-
711 sually rich document understanding. *arXiv preprint*
712 *arXiv:2511.09919*.

713 Wenhui Chen, Han Zhu, Wenhao Wang, Kai-Wei
714 Chang, William Yang Zhang, and William Wang.
715 2020. Tabfact: A large-scale dataset for table-based
716 fact verification. In *International Conference on*
717 *Learning Representations (ICLR)*.

718 Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo
719 Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,
720 Xizhou Zhu, Lewei Lu, et al. 2024. *Internvl: Scal-*
721 *ing up vision foundation models and aligning for*
722 *generic visual-linguistic tasks*. In *Proceedings of the*
723 *IEEE/CVF Conference on Computer Vision and Pat-*
724 *tern Recognition*, pages 24185–24198.

725 Julien Delaunay, Hanh Thi Hong Tran, Carlos-
726 Emiliano González-Gallardo, Georgeta Bordea,
727 Nicolas Sidere, and Antoine Doucet. 2023. A com-
728 prehensive survey of document-level relation extrac-
729 tion (2016-2023). *arXiv preprint arXiv:2309.16396*.

Chao Deng, Jiale Yuan, Pi Bu, Peijie Wang, Zhong-
Zhi Li, Jian Xu, Xiao-Hui Li, Yuan Gao, Jun
Song, Bo Zheng, and Cheng-Lin Liu. 2025. *Long-*
DocURL: a comprehensive multimodal long docu-
ment benchmark integrating understanding, reason-
ing, and locating. In *Proceedings of the 63rd An-*
nuual Meeting of the Association for Computational
Linguistics (Volume 1: Long Papers), pages 1135–
1159, Vienna, Austria. Association for Computa-
tional Linguistics. 730
731
732
733
734
735
736
737
738
739

Timo I Denk and Christian Reisswig. 2019. *Bertgrid:*
Contextualized embedding for 2d document repre-
sentation and understanding. In *Workshop on Docu-*
ment Intelligence at NeurIPS 2019. 740
741
742
743

Yihao Ding, Soyeon Caren Han, Yanbei Jiang, Yan
Li, Zechuan Li, and Yifan Peng. 2025a. Syn-
doc: A hybrid discriminative-generative frame-
work for enhancing synthetic domain-adaptive docu-
ment key information extraction. *arXiv preprint*
arXiv:2509.23273. 744
745
746
747
748
749

Yihao Ding, Soyeon Caren Han, Jean Lee, and Eduard
Hovy. 2025b. Deep learning based visually rich
document content understanding: A survey. *arXiv*
preprint arXiv:2408.01287. 750
751
752
753

Yihao Ding, Siqu Long, Jiabin Huang, Kaixuan Ren,
Xingxiang Luo, Hyunsuk Chung, and Soyeon Caren
Han. 2023. *Form-nlu: Dataset for the form natu-*
ral language understanding. In *Proceedings of the*
46th International ACM SIGIR Conference on Re-
search and Development in Information Retrieval,
pages 2807–2816. ACM. 754
755
756
757
758
759
760

Yihao Ding, Kaixuan Ren, Jiabin Huang, Siwen Luo,
and Soyeon Caren Han. 2024a. *Mmvqa: A com-*
prehensive dataset for investigating multipage mul-
timodal information retrieval in pdf-based visual
question answering. In *Proceedings of the Thirty-*
Third International Joint Conference on Artificial
Intelligence, IJCAI, pages 3–9. ijcai.org. 761
762
763
764
765
766
767

Yihao Ding, Kaixuan Ren, Jiabin Huang, Siwen
Luo, and Soyeon Caren Han. 2024b. *Mvqa: A*
dataset for multimodal information retrieval in pdf-
based visual question answering. *arXiv preprint*
arXiv:2404.12720. 768
769
770
771
772

Alexey Dosovitskiy, Lucas Beyer, Alexander
Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
Thomas Unterthiner, Mostafa Dehghani, Matthias
Minderer, G Heigold, S Gelly, et al. 2020. An im-
age is worth 16x16 words: Transformers for image
recognition at scale. In *International Conference on*
Learning Representations. 773
774
775
776
777
778
779

Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Om-
rani, Gautier Viaud, Céline Hudelot, and Pierre
Colombo. 2025. Colpali: Efficient document re-
trieval with vision language models. In *ICLR*. 780
781
782
783

Hao Feng, Qi Liu, Hao Liu, Jingqun Tang, Wengang
Zhou, Houqiang Li, and Can Huang. 2024. *Doc-*
pedia: Unleashing the power of large multimodal
784
785
786

787		model in the frequency domain for versatile document understanding. <i>Science China Information Sciences</i> , 67(12):1–14.	
788			
789			
790	Hao Feng, Zijian Wang, Jingqun Tang, Jinghui Lu, Wengang Zhou, Houqiang Li, and Can Huang. 2023. Unidoc: A universal large multimodal model for simultaneous text detection, recognition, spotting and understanding. <i>arXiv preprint arXiv:2308.11592</i> .		
791			
792			
793			
794			
795	Tongkun Guan, Zining Wang, Pei Fu, Zhengtao Guo, Wei Shen, Kai Zhou, Tiezhu Yue, Chen Duan, Hao Sun, Qianyi Jiang, et al. 2025. A token-level text image foundation model for document understanding. <i>arXiv preprint arXiv:2503.02304</i> .		
796			
797			
798			
799			
800	Pranay Gupta, Minesh Mathew, C.V. Jawahar, and Marcus Liwicki. 2022. Infovqa: Visual question answering on infographics with a multi-modal entity graph. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)</i> .		
801			
802			
803			
804			
805			
806	Siwei Han, Peng Xia, Ruiyi Zhang, Tong Sun, Yun Li, Hongtu Zhu, and Huaxiu Yao. 2025. Mdocagent: A multi-modal multi-agent framework for document understanding.		
807			
808			
809			
810	Adam W Harley, Alex Ufkes, and Konstantinos G Derpanis. 2015. Evaluation of deep convolutional nets for document image classification and retrieval. In <i>2015 13th International Conference on Document Analysis and Recognition (ICDAR)</i> , pages 991–995. IEEE.		
811			
812			
813			
814			
815			
816	Jiabang He, Lei Wang, Yi Hu, Ning Liu, Hui Liu, Xing Xu, and Heng Tao Shen. 2023. Icd3ie: In-context learning with diverse demonstrations updating for document information extraction. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 19485–19494. IEEE.		
817			
818			
819			
820			
821			
822	Teakgyu Hong, Donghyun Kim, Mingi Ji, Wonseok Hwang, Daehyun Nam, and Sungrae Park. 2022. Bros: A pre-trained language model focusing on text and layout for better key information extraction from documents. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pages 10767–10775.		
823			
824			
825			
826			
827			
828	Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. 2024. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 3096–3120.		
829			
830			
831			
832			
833			
834			
835	Anwen Hu, Haiyang Xu, Liang Zhang, Jiabo Ye, Ming Yan, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. 2025. mplug-docowl2: High-resolution compressing for ocr-free multi-page document understanding. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5817–5834.		
836			
837			
838			
839			
840			
841			
	Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. Layoutlmv3: Pre-training for document ai with unified text and image masking. In <i>Proceedings of the 30th ACM International Conference on Multimedia</i> , pages 4083–4091. ACM.		842 843 844 845 846
	Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and CV Jawahar. 2019. Icdar2019 competition on scanned receipt ocr and information extraction. In <i>2019 International Conference on Document Analysis and Recognition (ICDAR)</i> , pages 1516–1520. IEEE.		847 848 849 850 851 852
	Chelsi Jain, Yiran Wu, Yifan Zeng, Jiale Liu, Shengyu Dai, Zhenwen Shao, Qingyun Wu, and Huazheng Wang. 2025. SimpleDoc: Multi-Modal document understanding with Dual-Cue page retrieval and iterative refinement. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 28398–28415, Suzhou, China. Association for Computational Linguistics.		853 854 855 856 857 858 859 860
	Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. 2019. Funsd: A dataset for form understanding in noisy scanned documents. In <i>2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)</i> , volume 2, pages 1–6. IEEE.		861 862 863 864 865 866
	Anoop R Katti, Christian Reisswig, Cordula Guder, Sebastian Brarda, Steffen Bickel, Johannes Höhne, and Jean Baptiste Faddoul. 2018. Chargrid: Towards understanding 2d documents. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4459–4469. Association for Computational Linguistics.		867 868 869 870 871 872 873
	Geewook Kim, Teakgyu Hong, Moonbin Yim, Jeongyeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoon Yun, Dongyoon Han, and Seunghyun Park. 2022. Ocr-free document understanding transformer. In <i>Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII</i> , pages 498–517. Springer.		874 875 876 877 878 879 880 881
	Sungnyun Kim, Haofu Liao, Srikar Appalaraju, Peng Tang, Zhuowen Tu, Ravi Kumar Satzoda, R Manmatha, Vijay Mahadevan, and Stefano Soatto. 2024. Dockd: Knowledge distillation from llms for open-world document understanding models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12–16, 2024</i> , pages 3167–3193. Association for Computational Linguistics.		882 883 884 885 886 887 888 889 890
	Marcel Lamott, Yves-Noel Weweler, Adrian Ulges, Faisal Shafait, Dirk Krechel, and Darko Obradovic. 2024. Lapdoc: Layout-aware prompting for documents. In <i>International Conference on Document Analysis and Recognition</i> , pages 142–159. Springer.		891 892 893 894 895
	David D Lewis, Gady Agam, Shlomo Argamon, Ophir Frieder, David Grossman, and James Heard. 2006. Building a test collection for complex document information processing. In <i>Proceedings of the 29th</i>		896 897 898 899

900		<i>annual international ACM SIGIR conference on Research and development in information retrieval</i> , pages 665–666.	
901			
902			
903	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>International Conference on Machine Learning (ICML)</i> .		
904			
905			
906			
907			
908	Xin Li, Yunfei Wu, Xinghua Jiang, Zhihao Guo, Mingming Gong, Haoyu Cao, Yinsong Liu, Deqiang Jiang, and Xing Sun. 2024. Enhancing visual document understanding with contrastive learning in large visual-language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 15546–15555.		
909			
910			
911			
912			
913			
914			
915	Wenhui Liao, Jiapeng Wang, Hongliang Li, Chengyu Wang, Jun Huang, and Lianwen Jin. 2024. Doclayllm: An efficient and effective multi-modal extension of large language models for text-rich document understanding. <i>arXiv preprint arXiv:2408.15045</i> .		
916			
917			
918			
919			
920			
921	Chaohu Liu, Kun Yin, Haoyu Cao, Xinghua Jiang, Xin Li, Yinsong Liu, Deqiang Jiang, Xing Sun, and Linli Xu. 2024a. Hrvda: High-resolution visual document assistant. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 15534–15545.		
922			
923			
924			
925			
926			
927	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024b. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36.		
928			
929			
930	Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai. 2024c. Textmonkey: An ocr-free large multimodal model for understanding document.		
931			
932			
933			
934	Jinghui Lu, Haiyang Yu, Yanjie Wang, Yongjie Ye, Jingqun Tang, Ziwei Yang, Binghong Wu, Qi Liu, Hao Feng, Han Wang, et al. 2024. A bounding box is worth one token: Interleaving layout and text in a large language model for document understanding. <i>arXiv preprint arXiv:2407.01976</i> .		
935			
936			
937			
938			
939			
940	Chuwei Luo, Yufan Shen, Zhaoqing Zhu, Qi Zheng, Zhi Yu, and Cong Yao. 2024. Layoutlm: Layout instruction tuning with large language models for document understanding. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024</i> , pages 15630–15640. IEEE.		
941			
942			
943			
944			
945			
946			
947	Tengchao Lv, Yupan Huang, Jingye Chen, Yuzhong Zhao, Yilin Jia, Lei Cui, Shuming Ma, Yaoyao Chang, Shaohan Huang, Wenhui Wang, et al. 2023. Kosmos-2.5: A multimodal literate model. <i>arXiv preprint arXiv:2309.11419</i> .		
948			
949			
950			
951			
952	Pengyuan Lyu, Yulin Li, Hao Zhou, Weihong Ma, Xingyu Wan, Qunyi Xie, Liang Wu, Chengquan		
953			
	Zhang, Kun Yao, Errui Ding, et al. 2024. Structextv3: An efficient vision-language model for text-rich image perception, comprehension, and beyond. <i>arXiv preprint arXiv:2405.21013</i> .		
	Ahmed Masry, Juan A Rodriguez, Tianyu Zhang, Suyuchen Wang, Chao Wang, Aarash Feizi, Akshay Kalkunte Suresh, Abhay Puri, Xiangru Jian, Pierre-Andre Noel, et al. 2025. Alignvlm: Bridging vision and language latent spaces for multimodal document understanding. In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .		
	Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 2200–2209. IEEE.		
	Oshri Naparstek, Roi Pony, Inbar Shapira, Foad Abo Dahood, Ophir Azulai, Yevgeny Yaroker, Nadav Rubinstein, Maksym Lysak, Peter Staar, Ahmed Nassar, Nikolaos Livathinos, Christoph Auer, Elad Amrani, Idan Friedman, Orit Prince, Yevgeny Burstein, Adi Raz Goldfarb, and Udi Barzelay. 2024. Kvp10k : A comprehensive dataset for key-value pair extraction in business documents.		
	Feng Ni, Kui Huang, Yao Lu, Wenyu Lv, Guanzhong Wang, Zeyu Chen, and Yi Liu. 2025. Ppdocbee: Improving multimodal document understanding through a bag of tricks.		
	Eri Onami, Shuhei Kurita, Taiki Miyanishi, and Taro Watanabe. 2024. Jdocqa: Japanese document question answering dataset for generative language models.		
	OpenAI. 2024. Hello gpt-4o. https://openai.com/index/hello-gpt-4o/ .		
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.		
	Jaeyoo Park, Jin Young Choi, Jeonghyung Park, and Bohyung Han. 2024. Hierarchical visual feature aggregation for ocr-free document understanding. <i>Advances in Neural Information Processing Systems</i> , 37:105972–105996.		
	Seunghyun Park, Seung Shin, Bado Lee, Junyeop Lee, Jaeheung Surh, Minjoon Seo, and Hwalsuk Lee. 2019. Cord: a consolidated receipt dataset for post-ocr parsing. In <i>Workshop on Document Intelligence at NeurIPS 2019</i> .		
	Panupong Pasupat and Percy Liang. 2015. Compositional semantic parsing on semi-structured tables. In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics (ACL)</i> .		

1009	Vincent Perot, Kai Kang, Florian Luisier, Guolong Su,	Rubèn Tito, Dimosthenis Karatzas, and Ernest Val-	1067
1010	Xiaoyu Sun, Ramya Sree Boppana, Zilong Wang,	veny. 2023. Hierarchical multimodal transform-	1068
1011	Zifeng Wang, Jiaqi Mu, Hao Zhang, et al. 2024.	ers for multipage docvqa. <i>Pattern Recognition</i> ,	1069
1012	Lmdx: Language model-based document informa-	144:109834.	1070
1013	tion extraction and localization. In <i>Findings of</i>		
1014	<i>the Association for Computational Linguistics ACL</i>	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	1071
1015	2024, pages 15140–15168.	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	1072
1016	Minenobu Seki, Masakazu Fujio, Takeshi Nagasaki,	Baptiste Rozière, Naman Goyal, Eric Hambro,	1073
1017	Hiroshi Shinjo, and Katsumi Marukawa. 2007. In-	Faisal Azhar, et al. 2023. Llama: Open and effi-	1074
1018	formation management system using structure anal-	cient foundation language models. <i>arXiv preprint</i>	1075
1019	ysis of paper/electronic documents and its applica-	<i>arXiv:2302.13971.</i>	1076
1020	tions. In <i>Ninth International Conference on Docu-</i>	Jordy Van Landeghem, Rubèn Tito, Łukasz Borch-	1077
1021	<i>ment Analysis and Recognition (ICDAR 2007)</i> , vol-	mann, Michał Pietruszka, Paweł Joziak, Rafał	1078
1022	ume 2, pages 689–693. IEEE.	Powalski, Dawid Jurkiewicz, Mickaël Coustaty,	1079
1023	Yufan Shen, Chuwei Luo, Zhaoqing Zhu, Yang Chen,	Bertrand Anckaert, Ernest Valveny, et al. 2023.	1080
1024	Qi Zheng, Zhi Yu, Jiajun Bu, and Cong Yao. 2025.	Document understanding dataset and evaluation	1081
1025	Proctag: Process tagging for assessing the efficacy	(dude). In <i>Proceedings of the IEEE/CVF Inter-</i>	1082
1026	of document instruction data.	<i>national Conference on Computer Vision</i> , pages	1083
1027	Maxim Sidorov, Amanpreet Singh, Yu Li, Jianfeng	19528–19540.	1084
1028	Liao, Ming Liao, Yaxing Wang, Lichao Wang,	Dongsheng Wang, Natraj Raman, Mathieu Sibue,	1085
1029	Shouling Gong, Chen Change Loy, and Xiang Bai.	Zhiqiang Ma, Petr Babkin, Simerjot Kaur, Yulong	1086
1030	2020. Textcaps: A dataset for image captioning with	Pei, Armineh Nourbakhsh, and Xiaomo Liu. 2024a.	1087
1031	reading. In <i>Proceedings of the European Confer-</i>	Docllm: A layout-aware generative language model	1088
1032	<i>ence on Computer Vision (ECCV).</i>	for multimodal document understanding. In <i>Pro-</i>	1089
1033	Štěpán Šimsa, Milan Šulc, Michal Uříčář, Yash Patel,	<i>ceedings of the 62nd Annual Meeting of the Associa-</i>	1090
1034	Ahmed Hamdi, Matěj Kocián, Matyáš Skalický, Jiří	<i>tion for Computational Linguistics (Volume 1: Long</i>	1091
1035	Matas, Antoine Doucet, Mickaël Coustaty, and Di-	<i>Papers)</i> , pages 8529–8548. Association for Compu-	1092
1036	mosthenis Karatzas. 2023. DocILE benchmark for	tational Linguistics.	1093
1037	document information localization and extraction.		
1038	Amanpreet Singh, Vedanuj Natarajan, Yu Jiang, Xinlei	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-	1094
1039	Chen, Meet Shah, Marcus Rohrbach, Dhruv Batra,	hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,	1095
1040	Devi Parikh, and Aniruddha Krishnamurthy. 2019.	Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang,	1096
1041	Textvqa: Visual question answering with reading. In	Mengfei Du, Xuancheng Ren, Rui Men, Dayi-	1097
1042	<i>Proceedings of the IEEE/CVF Conference on Com-</i>	heng Liu, Chang Zhou, Jingren Zhou, and Junyang	1098
1043	<i>puter Vision and Pattern Recognition (CVPR).</i>	Lin. 2024b. Qwen2-vl: Enhancing vision-language	1099
1044	Li Sun, Liu He, Shuyue Jia, Yangfan He, and Chenyu	model’s perception of the world at any resolution.	1100
1045	You. 2025. Docagent: An agentic framework for	<i>arXiv preprint arXiv:2409.12191.</i>	1101
1046	multi-modal long-context document understanding.	Yonghui Wang, Wengang Zhou, Hao Feng, Keyi	1102
1047	In <i>Proceedings of the 2025 Conference on Empiri-</i>	Zhou, and Houqiang Li. 2023. Towards im-	1103
1048	<i>cal Methods in Natural Language Processing</i> , pages	proving document understanding: An exploration	1104
1049	17712–17727.	on text-grounding via mllms. <i>arXiv preprint</i>	1105
1050	Ryota Tanaka, Taichi Iki, Kyosuke Nishida, Kuniko	<i>arXiv:2311.13194.</i>	1106
1051	Saito, and Jun Suzuki. 2024. Instructdoc: A dataset	Zhaowei Wang, Wenhao Yu, Xiyu Ren, Jipeng Zhang,	1107
1052	for zero-shot generalization of visual document	Yu Zhao, Rohit Saxena, Liang Cheng, Ginny Wong,	1108
1053	understanding with instructions. In <i>Proceedings of</i>	Simon See, Pasquale Minervini, Yangqiu Song, and	1109
1054	<i>the AAAI conference on artificial intelligence</i> , pages	Mark Steedman. 2025a. Mmlongbench: Bench-	1110
1055	19071–19079. AAAI Press.	marking long-context vision-language models effec-	1111
1056	Jingqun Tang, Qi Liu, Yongjie Ye, Jinghui Lu, Shu	tively and thoroughly.	1112
1057	Wei, Chunhui Lin, Wanqing Li, Mohamad Fitri	Zining Wang, Tongkun Guan, Pei Fu, Chen Duan,	1113
1058	Faiz Bin Mahmood, Hao Feng, Zhen Zhao, Yangfan	Qianyi Jiang, Zhentao Guo, Shan Guo, Junfeng Luo,	1114
1059	He, Kuan Lu, Yanjie Wang, Yuliang Liu, Hao Liu,	Wei Shen, and Xiaokang Yang. 2025b. Marten: Vi-	1115
1060	Xiang Bai, and Can Huang. 2025. Mtvqa: Bench-	sual question answering with mask generation for	1116
1061	marking multilingual text-centric visual question	multi-modal document understanding.	1117
1062	answering.		
1063	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Bur-	Toyohide Watanabe, Qin Luo, and Noboru Sugie. 1995.	1118
1064	nell, Libin Bai, and et al. 2024. Gemini 1.5: Unlock-	ing multimodal understanding across millions of to-	1119
1065	kens of context.	kinds of table-form docu-	1120
1066		ments. <i>IEEE Transactions on Pattern Analysis and</i>	1121
		<i>Machine Intelligence</i> , 17(4):432–445.	

1122	Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao,	1180
1123	Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han,	1181
1124	and Xiangyu Zhang. 2024. Vary: Scaling up the	1182
1125	vision vocabulary for large vision-language model.	1183
1126	In <i>European Conference on Computer Vision</i> , pages	1184
1127	408–424. Springer.	1185
1128	Kun Xiang, Heng Li, Terry Jingchen Zhang, Yinya	1186
1129	Huang, Zirong Liu, Peixin Qu, Jixi He, Jiaqi Chen,	1187
1130	Yu-Jie Yuan, Jianhua Han, Hang Xu, Hanhui Li,	1188
1131	Mrinmaya Sachan, and Xiaodan Liang. 2025. <i>Seep-</i>	1189
1132	<i>hys: Does seeing help thinking? – benchmarking</i>	1190
1133	<i>vision-based physics reasoning</i> .	1191
1134	Xudong Xie, Hao Yan, Liang Yin, Yang Liu, Jing Ding,	1192
1135	Minghui Liao, Yuliang Liu, Wei Chen, and Xiang	1193
1136	Bai. 2025. Pdf-wukong: A large multimodal model	1194
1137	for efficient long pdf reading with end-to-end sparse	1195
1138	sampling.	1196
1139	Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu	1197
1140	Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha	1198
1141	Zhang, Wanxiang Che, et al. 2021. <i>Layoutlmv2:</i>	1199
1142	<i>Multi-modal pre-training for visually-rich document</i>	1200
1143	<i>understanding</i> . In <i>Proceedings of the 59th Annual</i>	1201
1144	<i>Meeting of the Association for Computational Lin-</i>	1202
1145	<i>guistics and the 11th International Joint Conference</i>	1203
1146	<i>on Natural Language Processing (Volume 1: Long</i>	1204
1147	<i>Papers)</i> , pages 2579–2591. Association for Compu-	1205
1148	tational Linguistics.	1206
1149	Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang,	1207
1150	Furu Wei, and Ming Zhou. 2020. <i>Layoutlm: Pre-</i>	1208
1151	<i>training of text and layout for document image</i>	1209
1152	<i>understanding</i> . In <i>Proceedings of the 26th ACM</i>	1210
1153	<i>SIGKDD International Conference on Knowledge</i>	1211
1154	<i>Discovery & Data Mining</i> , pages 1192–1200. ACM.	1212
1155	Yiheng Xu, Tengchao Lv, Lei Cui, Guoxin Wang, Yi-	1213
1156	juan Lu, Dinei Florencio, Cha Zhang, and Furu Wei.	1214
1157	2022. Xfund: A benchmark dataset for multilingual	1215
1158	visually rich form understanding. In <i>Findings of</i>	1216
1159	<i>the association for computational linguistics: ACL</i>	1217
1160	<i>2022</i> , pages 3214–3224.	1218
1161	Xiao Yang, Ersin Yumer, Paul Asente, Mike Kraley,	1219
1162	Daniel Kifer, and C Lee Giles. 2017. <i>Learning</i>	1220
1163	<i>to extract semantic structure from documents using</i>	1221
1164	<i>multimodal fully convolutional neural networks</i> . In	1222
1165	<i>Proceedings of the IEEE Conference on Computer</i>	1223
1166	<i>Vision and Pattern Recognition</i> , pages 5315–5324.	1224
1167	IEEE Computer Society.	1225
1168	Zhibo Yang, Jun Tang, Zhaohai Li, Pengfei Wang, Jian-	1226
1169	qiang Wan, Humen Zhong, Xuejing Liu, Mingkun	1227
1170	Yang, Peng Wang, Shuai Bai, Lianwen Jin, and Jun-	1228
1171	yang Lin. 2024. <i>Cc-ocr: A comprehensive and chal-</i>	1229
1172	<i>lenging ocr benchmark for evaluating large multi-</i>	1230
1173	<i>modal models in literacy</i> .	1231
1174	Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming	1232
1175	Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu, Chen-	1233
1176	liang Li, Junfeng Tian, Qian Qi, Ji Zhang, and Fei	1234
1177	Huang. 2023a. <i>mplug-docowl: Modularized mul-</i>	
1178	<i>timodal large language model for document under-</i>	
1179	<i>standing</i> .	
	Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye,	1180
	Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian,	1181
	Qi Qian, Ji Zhang, et al. 2023b. <i>Ureader: Univer-</i>	1182
	<i>saral ocr-free visually-situated language understand-</i>	1183
	<i>ing with multimodal large language model</i> . In <i>Find-</i>	1184
	<i>ings of the Association for Computational Linguis-</i>	1185
	<i>tics: EMNLP 2023</i> , pages 2841–2858.	1186
	Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye,	1187
	Ming Yan, Yiyang Zhou, Junyang Wang, An-	1188
	wen Hu, Pengcheng Shi, Yaya Shi, et al. 2023c.	1189
	<i>mplug-owl: Modularization empowers large lan-</i>	1190
	<i>guage models with multimodality</i> . <i>arXiv preprint</i>	1191
	<i>arXiv:2304.14178</i> .	1192
	Wenwen Yu, Zhibo Yang, Yuliang Liu, and Xiang	1193
	Bai. 2025. Doctinker: Explainable multimodal	1194
	large language models with rule-based reinforce-	1195
	ment learning for document understanding. In <i>Pro-</i>	1196
	<i>ceedings of the IEEE/CVF International Conference</i>	1197
	<i>on Computer Vision</i> , pages 837–847.	1198
	Ya-Qi Yu, Minghui Liao, Jihao Wu, Yongxin Liao,	1199
	Xiaoyu Zheng, and Wei Zeng. 2024a. <i>Texthawk:</i>	1200
	<i>Exploring efficient fine-grained perception of multi-</i>	1201
	<i>modal large language models</i> .	1202
	Ya-Qi Yu, Minghui Liao, Jiwen Zhang, and Jihao Wu.	1203
	2024b. <i>Texthawk2: A large vision-language model</i>	1204
	<i>excels in bilingual ocr and grounding with 16x fewer</i>	1205
	<i>tokens</i> .	1206
	Jiaxin Zhang, Wentao Yang, Songxuan Lai, Zecheng	1207
	Xie, and Lianwen Jin. 2024a. <i>Dockylin: A large</i>	1208
	<i>multimodal model for visual document understand-</i>	1209
	<i>ing with efficient visual slimming</i> . <i>arXiv preprint</i>	1210
	<i>arXiv:2406.19101</i> .	1211
	Jinxu Zhang, Qiyuan Fan, and Yu Zhang. 2025. <i>Do-</i>	1212
	<i>cassistant: Integrating key-region reading and step-</i>	1213
	<i>wise reasoning for robust document visual question</i>	1214
	<i>answering</i> . In <i>Findings of the Association for Com-</i>	1215
	<i>putational Linguistics: EMNLP 2025</i> , pages 3496–	1216
	3511.	1217
	Renshan Zhang, Yibo Lyu, Rui Shao, Gongwei Chen,	1218
	Weili Guan, and Liqiang Nie. 2024b. <i>Token-level</i>	1219
	<i>correlation-guided compression for efficient multi-</i>	1220
	<i>modal document understanding</i> . <i>arXiv</i> .	1221
	Ruiyi Zhang, Yufan Zhou, Jian Chen, Jiuxiang Gu,	1222
	Changyou Chen, and Tong Sun. 2024c. <i>Llava-read:</i>	1223
	<i>Enhancing reading ability of multimodal language</i>	1224
	<i>models</i> . <i>arXiv preprint arXiv:2407.19185</i> .	1225
	Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou,	1226
	Nedim Lipka, Diyi Yang, and Tong Sun. 2024d.	1227
	<i>Llavar: Enhanced visual instruction tuning for text-</i>	1228
	<i>rich image understanding</i> .	1229
	Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes.	1230
	2019. <i>Publaynet: largest dataset ever for document</i>	1231
	<i>layout analysis</i> . In <i>2019 International Conference</i>	1232
	<i>on Document Analysis and Recognition (ICDAR)</i> ,	1233
	pages 1015–1022. IEEE.	1234

1235 Yuke Zhu, Yue Zhang, Dongdong Liu, Chi Xie, Zi-
1236 hua Xiong, Bo Zheng, and Sheng Guo. 2025a. En-
1237 hancing document understanding with group posi-
1238 tion embedding: A novel approach to incorporate
1239 layout information. In *The Thirteenth International*
1240 *Conference on Learning Representations*.

1241 Zhaoqing Zhu, Chuwei Luo, Zirui Shao, Feiyu Gao,
1242 Hangdi Xing, Qi Zheng, and Ji Zhang. 2025b. A
1243 simple yet effective layout token in large language
1244 models for document understanding. *arXiv preprint*
1245 *arXiv:2503.18434*.

1246
1247
1248
1249
1250
1251
1252

1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268

1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289

1290
1291
1292
1293

A More Framework Details

A.1 Open-source Frameworks

Table 2 presents official open-source links for VRDU and MLLM frameworks, underscoring the vital role of open access in fostering transparency, reproducibility, and accelerated innovation within the research community.

A.2 Model Training Paradigm Comparison

Table 3 provides a comprehensive comparison of MLLM-based VRDU frameworks across three major training stages: Pretraining (PT), Instruction-tuning (IT), and Supervised Fine-tuning (SFT). OCR-dependent models generally rely on external text extraction and have limited pretraining because they are trained on OCR-processed inputs. In contrast, OCR-free models, which operate directly on document images, demonstrate richer instruction-tuning and fine-tuning strategies, often involving frozen or LoRA-based vision and language encoders. This highlights the diverse training paradigms and modular designs adopted to balance efficiency, adaptability, and performance across frameworks.

A.3 Document Parsing Tools

Table 5 provides a comparative overview of representative OCR engines, document parsing APIs, and vision–language models for document understanding. The table highlights clear trade-offs across deployment modes, pricing models, and functional capabilities: traditional OCR engines are predominantly open-source and locally deployable but offer limited support for structured document parsing, while commercial document APIs and vision LLMs more frequently provide GPU acceleration and native document-structure extraction at the cost of cloud dependency and usage-based pricing. Recent vision–language models bridge OCR and higher-level reasoning by supporting multimodal inputs (image and PDF) and multilingual processing, yet vary substantially in openness and deployment flexibility. Overall, the comparison illustrates the evolving landscape from text-centric OCR toward multimodal, structure-aware document understanding systems.

A.4 Model Component Details

Table 4 presents a comprehensive comparison of component configurations adopted by recent MLLM-based frameworks for VRDU, spanning

both OCR-Dependent and OCR-Free paradigms. For each model, we summarize its LLM backbone (e.g., Vicuna, Qwen, LLaMA, GPT), vision encoder (e.g., CLIP, ViT, Swin), input resolution (including dynamic scaling and cropping), and specialized adaptors or projectors (e.g., LoRA, MLP, QPN) used for multimodal fusion. OCR-Dependent models typically incorporate layout-aware encoders (e.g., LayoutLMv3, DocFormer) and rely on structured textual inputs. In contrast, OCR-Free models process raw document images directly, often requiring higher resolutions and additional modules such as resamplers, visual abstractors, or cropping strategies. The table also lists the maximum supported image resolution, indicating each model’s capacity for fine-grained visual understanding. This comparison highlights the increasing diversity in MLLM architectures and the adoption of lightweight tuning techniques for scalable VRDU.

B Dataset

Pretraining Datasets. The goal of pretraining is to enhance multimodal understanding and improve generalization across VRDU tasks. Similar to pretrained VRDU frameworks, MLLM-based approaches commonly perform continued pretraining on large-scale, cross-domain document collections such as IIT-CDIP (Lewis et al., 2006), which contains over 6 million scanned documents across diverse domains, though lacking explicit layout annotations—often supplemented with OCR-derived bounding boxes. RVL-CDIP (Harley et al., 2015), a curated subset with 400,000 documents across 16 categories, is widely used for document classification and low-resource pretraining. Beyond these general-purpose datasets, recent frameworks (Zhang et al., 2024d; Wang et al., 2023) have introduced self-collected datasets to target domain-specific or task-oriented scenarios, including slide decks (Feng et al., 2024), academic papers (Wang et al., 2024a), and other structured document types (Yu et al., 2024b), as summarized in Table 6.

Instruction-tuning Datasets. Instruction-tuning aims to enhance a model’s understanding of user queries. Many frameworks (Zhang et al., 2024b; Park et al., 2024) perform instruction-tuning directly on benchmark document collections to improve downstream task performance.

1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313

1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337

1338
1339
1340
1341
1342
1343

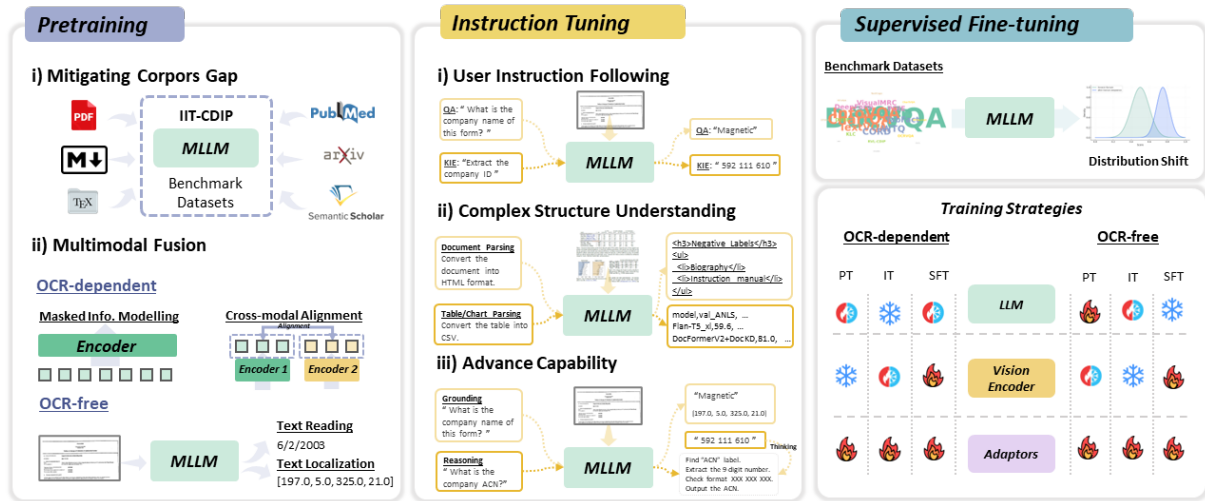


Figure 3: MLLM-based VRDU framework training paradigms. [Yifan: This figure is not referred to in the text.]

Others (Luo et al., 2024; Liu et al., 2024a) generate large-scale synthetic datasets using OCR tools to extract text and layout information from VRD-related benchmarks such as layout analysis (Zhong et al., 2019) and document classification (Harley et al., 2015). Instruction-response pairs are then created based on predefined task definitions. Some frameworks also construct their own multi-domain datasets to improve generalizability and prevent data leakage (Wei et al., 2024; Feng et al., 2023). Instruction-tuning is critical for domain adaptation and accurate instruction interpretation. As shown by Table 7, some frameworks increasingly generate synthetic instruction-tuning datasets tailored to their architectures, prioritizing alignment over generalizability achieved through benchmark-based tuning.

Supervised Fine-tuning Datasets. To improve performance on downstream tasks, some frameworks apply supervised fine-tuning on question answering datasets such as DocVQA (Mathew et al., 2021) and MPDocVQA (Tito et al., 2023). Additionally, several key information extraction benchmarks—such as FUNSD (Jaume et al., 2019), FormNLU (Ding et al., 2023), and CORD (Park et al., 2019), have been reformulated into QA-style formats to enable evaluation with generative frameworks.

C Benchmark Datasets

Based on differences in downstream tasks and the benchmark dataset’s domain, we list the widely used VRDU dataset and its key attributes in Table 8, including both VRD-related Key Informa-

tion Extraction (KIE) and Visual Question Answering (VQA).

Key Information Extraction Benchmarks for Key Information Extraction (KIE) are shifting from early schema-constrained tasks (e.g., SROIE (Huang et al., 2019), FUNSD (Jaume et al., 2019)) toward larger, multilingual, cross domain, multi page and open-vocabulary challenges. While form-like structures (e.g., DocILE (Šimsa et al., 2023), Form-NLU (Ding et al., 2023)) still dominate the landscape, modern resources such as KVP10k (Naparstek et al., 2024) and CC-OCR-KIE (Yang et al., 2024) focus on *open-category* extraction without predefined schemas. Furthermore, a clear trend of dataset consolidation and multilingual expansion has emerged.

Visual Question Answering has undergone a comparable evolution, shifting from early single-page, text-centric retrieval to benchmarks that probe multiple dimensions of complexity. This progression is reflected in broader multilingual coverage (e.g., MTVQA (Tang et al., 2025), JDocQA (Onami et al., 2024)) and more diverse, multi-domain settings (e.g., DUDE (Van Landeghem et al., 2023)). Recent datasets increasingly emphasize long-context comprehension over multi-page documents: benchmarks such as LongDocURL (Deng et al., 2025) and MMLongBench-Doc (Wang et al., 2025a) contain documents averaging dozens of pages and often demand non-trivial cross-page evidence aggregation and reasoning. In parallel, reasoning requirements have deepened toward domain-

specific expertise, as illustrated by vision-essential physics problem solving in SEEPHYS (Xiang et al., 2025). Finally, dataset scale has expanded substantially—reaching millions of instances in collections such as MMVQA (Ding et al., 2024a)—thereby enabling rigorous stress-testing of the capacity and reasoning limits of modern multimodal models.

Other Domain Datasets Many frameworks are evaluated on other domain-specific datasets as well, including those for chart understanding and webpage analysis. For instance, InfoVQA (Gupta et al., 2022) focuses on visual question answering for information-centric records. Benchmarks like WTQ (Pasupat and Liang, 2015) and TabFact (Chen et al., 2020) assess a model’s ability to reason over tabular data, and ChartQA evaluates chart comprehension skills. Additionally, TextVQA (Singh et al., 2019) and TextCaps (Sidorov et al., 2020) target text recognition and semantic reasoning in natural images.

D Quantitative Analysis

D.1 Performance on Single Page Benchmarks

Table 9 highlights clear trends in the performance of general-domain LLMs/MLLM and OCR-dependent and OCR-free document understanding frameworks across several popular benchmarks. Generally, OCR-dependent models achieve consistently strong results on classic form and receipt datasets such as FUNSD, CORD, and SROIE—often exceeding 80% accuracy, with top models such as PDF-WuKong, GPE, and DocLayLLM achieving state-of-the-art performance. In contrast, OCR-free frameworks, while demonstrating rapid progress, still lag on these traditional datasets but show remarkable advances on more visually and semantically complex benchmarks such as DocVQA, ChartVQA, and InfoVQA. Notably, the latest OCR-free models, including Texthawk2, Marten, and PP-DocBee, have begun to outperform or match OCR-dependent methods on DocVQA and chart-centric tasks, signalling a narrowing of the gap in real-world document reasoning capabilities. However, coverage remains uneven, with many OCR-free models performing poorly on specific datasets, indicating ongoing challenges with generalizability and benchmark saturation. Overall, while OCR-dependent methods remain dominant for structured text ex-

traction, OCR-free approaches are quickly maturing and expanding the frontier of end-to-end document understanding.

D.2 Performance on Multi-Page Benchmarks

We report the performance of existing multi-page frameworks on two multi-page VRDU benchmarks in Table 10. General-domain models can achieve reasonable performance; however, frameworks equipped with mechanisms explicitly designed for visually rich documents (VRDs) consistently yield substantial improvements. Currently, most high-performing multi-page methods rely on OCR-dependent pipelines and achieve strong results by leveraging external OCR tools. While such designs reduce the burden of directly understanding and compressing visual representations, they also inherit the limitations of OCR-based approaches, including error accumulation as observed in single-page scenarios. For multi-page tasks, this challenge is further amplified, highlighting the need for more effective strategies to manage the large number of visual tokens and to improve text understanding in multi-page, text-dense document inputs.

Framework	Model Name	Official Open Source Link
mPLUG-DocOwl 1.5	DocOwl 1.5	github.com/X-PLUG/mPLUG-DocOwl/tree/main/DocOwl1.5
mPLUG-DocOwl 2	DocOwl 2	github.com/X-PLUG/mPLUG-DocOwl/tree/main/DocOwl2
UReader	UReader	github.com/X-PLUG/mPLUG-DocOwl/tree/main/UReader
KOSMOS-2.5	KOSMOS-2.5 / 2.5-CHAT	aka.ms/kosmos25
LLaVAR	LLaVAR	github.com/SALT-NLP/LLaVAR
Marten	Marten	github.com/PriNing/Marten
LEOPARD	LEOPARD	github.com/Jil10001/Leopard

Table 2: Official open-source links for some VRDU/MLLM frameworks.

Model Name	Vision Encoder			LLM Backbone			Adaptors		
	PT	IT	SFT	PT	IT	SFT	PT	IT	SFT
OCR-Dependent									
ICL-D3IE (2023)	-	-	-	-	-	-	-	-	-
DocLLM (2024a)	T	T	-	-	-	-	T	T	-
LAPDoc (2024)	-	-	-	-	-	-	-	-	-
LMDX (2024)	-	-	-	-	-	-	-	-	-
ProcTag (2025)	-	-	T	-	-	T	-	-	T
DocKD (2024)	-	-	T	-	-	T	-	-	-
DoCo (2024)	F	-	F	T	-	F	T	-	T
InstructDoc (2024)	-	F	F	-	F	F	-	T	T
LayoutLLM (2024)	-	F	T	-	F	F	-	T	T
LLaVA-Read (2024c)	F	T	-	F	F	-	T	T	-
LayTextLLM (2024)	F	-	T	-	-	-	T	-	T
LayTokenLLM (2025b)	F	-	F	-	-	-	T	-	T
GPE (2025a)	-	-	T	-	-	-	-	-	-
MDocAgent (2025)	-	-	-	-	-	-	-	-	-
PDF-WuKong (2025)	-	-	T	-	-	T	-	-	-
DocLayLLM (2024)	F	F	-	T	T	-	T	T	-
DocAssistant (2025)	-	F	-	-	F	-	-	T	-
AlignVLM (2025)	T	T	F	T	T	T	T	T	T
DocThinker (2025)	-	-	T	-	-	T	-	-	T
OCR-Free									
KOSMOS-2.5 (2023)	-	T	T	-	T	F	-	T	T
mPLUG-DocOwl (2023a)	-	F	-	-	F	-	-	T	-
UReader (2023b)	-	F	-	-	F	-	-	T	-
TGDoc (2023)	-	F	T	-	F	F	-	T	T
UniDoc (2023)	-	F	T	-	F	F	-	T	T
DocPedia (2024)	F	-	T	T	-	T	T	-	T
HRVDA (2024a)	F	F	-	T	F	-	T	T	-
Vary (2024)	T	-	T	T	-	F	T	-	T
mPLUG-DocOwl 1.5 (2024)	-	F	T	-	T	F	-	T	T
HVFA (2024)	-	F	-	-	F	-	-	T	-
mPLUG-DocOwl2 (2025)	-	F	T	-	T	F	-	T	T
Texthawk (2024a)	-	F	T	-	F	F	-	T	T
Texthawk2 (2024b)	-	F	T	-	F	T	-	T	T
TextMonkey (2024c)	-	T	-	-	T	-	T	T	-
Llavar (2024d)	-	F	T	-	F	F	-	T	T
TokenCorrCompressor (2024b)	-	-	F	-	-	F	-	-	T
DocKylin (2024a)	-	F	T	-	T	T	-	T	T
Marten (2025b)	-	F	T	-	T	T	-	T	T
PP-DocBee (2025)	-	-	T	-	-	F	-	-	-
TokenFD (2025)	T	F	T	T	F	T	T	T	T

Table 3: Comparison of MLLM-based VRDU frameworks. PT - Pretraining, IT - Instruction-tuning, SFT - Supervised Fine-tuning.

Table 4: Comparison of MLLM-based VRDU frameworks: Backbone and Adapter configurations. “-” denotes the component is not applicable or not disclosed.

Model Name	LLM Backbone	Vision Backbone	Resolution	Adaptors and Projectors
OCR-Dependent				
ICL-D3IE (2023)	GPT-3, ChatGPT	-	-	-
DocLLM (2024a)	Falcon-1B/LLaMA2-7B	-	-	Disentangled Spatial Attention
LAPDoc (2024)	ChatGPT, Solar	-	-	-
LMDX (2024)	PaLM 2-S, Gemini Pro	-	-	-
ProcTag (2025)	Qwen-7B/Qwen-VL-7B	qwen2vl vision encoder	Dynamic (224×224 to 448×448)	qwen2vl projector
DocKD (2024)	DocFormerv2 language decoder	DocFormerv2 vision encoder	Derived from CNN backbone	-
DoCo (2024)	Qwen-VL-Chat/mPLUG-Owl	ViT-bigG	224×224	Position-Aware Vision-Language Adapter, Visual Abstractor
InstructDr (2024)	Flan-T5	CLIP	224×224	Document-former
LayoutLLM (2024)	Vicuna-7B-v1.5, LLaMA2-7B-chat	LayoutLMv3	224×224	MLP
LLaVA-Read (2024c)	Vicuna-1.5 13B	CLIP-ViT-L/14-336 + ConvNext-L/32-320	336×336	MLP
LayTextLLM (2024)	Llama2-7B-base	-	320×320	Spatial Layout Projector + Layout Partial LoRA
LayTokenLLM (2025b)	Qwen1.5-7B, LLaMA3-8B	-	-	Layout Tokenizer + LORA
GPE (2025a)	LLaMA2-7B, Qwen2-7B, ChatGLM-6B	-	-	-
MDocAgent (2025)	LLaMA-3.1-8B (Text), Qwen2-VL-7B (Others)	ColPali	448×448	-
PDF-WuKong (2025)	IXC2-VL-4KHD	IXC2-VL-4KHD	Dynamic (336×336 to 3840×1600)	-
DocLayLLM (2024)	LLaMA2-7B, LLaMA3-8B	LayoutLMv3 ve	224×224	Layout Embedder + Projector + LORA
DocAssistant (2025)	InternVL2-Chat-2B	InternVL2-Chat-2B Vision Encoder	448×448	Mixture-of-Modality Adaptation + Projector + LORA
AlignVLM (2025)	Llama 3.1-1B\3B\8B	SigLip-400M	Dynamic (14×14 patches)	ALIGN Module
DocThinker (2025)	Qwen2.5-VL-3B\7B	Qwen2.5-VL-3B\7B Vision Encoder	336×336, 1536×1536	-
OCR-Free				
KOSMOS-2.5 (2023)	Transformer decoder	Pix2Struct-Large ViT-based	1024×1024	Resampler
mPLUG-DocOwl (2023a)	mPLUG-Owl	ViT	224×224	Visual Abstractor + Lora
UReader (2023b)	mPLUG-Owl	CLIP-like ViT	224×224 (×20 crops)	Visual Abstractor + Lora
TGDoc (2023)	Vicuna-7B	CLIP-ViT-L/14	224×224 and 336×336	MLP
UniDoc (2023)	Vicuna	CLIP-ViT-L/14	224×224 and 336×336	MLP
DocPedia (2024)	Vicuna-7B	Swin Transformer	2560×2560	MLP
HRVDA (2024a)	LLaMA-2-7B	Swin Transformer	1536×1536	Content Detector + MLP Projector + LoRA
Vary (2024)	OPT125M + Qwen-7B, Vicuna-7B	CLIP + SAM	1024×1024	MLP
mPLUG-DocOwl 1.5 (2024)	mPLUG-Owl2	ViT/L-14	448×448 (×9 crops)	H-Reducer
HVFA (2024)	BLIP-2-OPT-2.7B, mPLUG-Owl-7B	ViT	224×224 × crops	HVFA + Lora + Resampler
mPLUG-DocOwl2 (2025)	mPLUG-Owl2	ViT	504×504 (×12 crops)	H-Reducer
Texthawk (2024a)	InternLM-XComposer 7B	SigLIP-SO (ViT)	224×224 × crops	Resampler + LoRA + QPN + MLCA
Texthawk2 (2024b)	Qwen2-7B-Instruct	SigLIP-SO (ViT)	224×224 × crops (up to 72 crops)	Resampler + QPN + MLCA + Detection Head + LoRA
TextMonkey (2024c)	Qwen-VL-Chat, mPLUG-Owl	ViT-BigG	448×448 × crops	Image + Token Resampler
Llavar (2024d)	Vicuna-13B	CLIP-ViT-L/14	224×224 and 336×336	MLP
TokenCorrCompressor (2024b)	LLaMA2-7B	CLIP-ViT-L/14	224×224 and 336×336	Token Correlation Compressor + LORA
DocKylin (2024a)	Qwen-7B-Chat	Swin (Donut-Swin, 0.07B)	1728×1728	MLP + APS + DTS
Marten (2025b)	InternLM2-7B	InternViT-300M	448×448 (×6 crops)	MLP + Mask Generator Module

Continued on next page

Table 4: Comparison of MLLM-based VRDU frameworks: Backbone and Adapter configurations. “-” denotes the component is not applicable or not disclosed. (Continued)

Model Name	LLM Backbone	Vision Backbone	Resolution	Adaptors and Projectors
PP-DocBee (2025)	Qwen2-VL-2B	ViT	1680×1204	-
TokenFD (2025)	token embedding layer	ViT	448 × 448 (× 6 crops)	Token abstractor

Tool Name	Provider	Tool Type	Deployment	Pricing	Input Modalities	Languages	Openness	GPU	Doc Parsing
pdfminer.six	Y. Shinyama et al.	OCR Engine	Local	Free	PDF	Multi	Open-source	No	No
Mistral OCR	Mistral AI	Document API	Cloud	Paid (Usage-based)	Image, PDF	Multi	Closed	Supported	Yes
LightOnOCR	LightOnAI	Vision LLM	Cloud	Paid (Usage-based)	Image, PDF	Multi	Closed	Supported	No
Google Cloud Vision	Google	Document API	Cloud	Paid (Usage-based)	Image, PDF	Multi	Closed	Supported	Yes
Kraken	Inria et al.	OCR Engine	Local	Free	Image, PDF	Multi	Open-source	Supported	No
Qwen3-VL	Aliyun	Vision LLM	Hybrid	Free*	Image, PDF	Pretrained/Dependent	Closed	Supported	No
olmOCR	AI2	OCR Engine	Hybrid	Free	Image, PDF	Multi	Open-source	Supported	Yes
AttentionOCR	Guo & Deng	OCR Engine	Local	Free	Image	Multi	Open-source	Supported	No
Calamari	Univ. Würzburg	OCR Engine	Local	Free	Image	Multi	Open-source	Supported	No
EasyOCR	JaidevAI	OCR Engine	Local	Free	Image	Multi	Open-source	Supported	No
OpenAI Vision	OpenAI	Vision LLM	Cloud	Paid (Usage-based)	Image, PDF	Multi	Closed	Supported	Yes
Tesseract	S. Weil	OCR Engine	Local	Free	Image	Multi	Open-source	No	No
Adobe PDF Extract	Adobe	Document API	Cloud	Paid (Usage-based)	PDF	Multi	Closed	Supported	Yes
PaddleOCR	PaddlePaddle	OCR Engine	Cloud	Free	Image, PDF	Multi	Open-source	Supported	Yes
docTR	Mindee	OCR Engine	Local	Free	Image, PDF	Pretrained/Dependent	Open-source	Supported	No
DeepSeek-OCR	DeepSeek AI	Vision LLM	Hybrid	Paid (Usage-based)	Image, PDF	Multi	Open-source	Supported	No
HunyuanOCR	Tencent	Vision LLM	Local	Free	Image, PDF	Multi	Open-source	Supported	No
Ocular	Berkeley NLP	OCR Engine	Local	Free	Image, PDF	Multi	Open-source	Supported	No
MinerU	OpenDataLab	Document API	Local	Free	PDF	Multi	Open-source	Supported	Yes
SuryaOCR	Datalab	OCR Engine	Local	Free	Image, PDF, Word, PPT	Multi	Open-source	Supported	No
Seed-VL	ByteDance Seed	Vision LLM	Cloud	Paid (Usage-based)	Image, PDF	Multi	Open-source	Supported	Yes

Table 5: Comparison of OCR engines, document parsing APIs, and vision-language models for document understanding.

Study	Dataset	Source	Size	Public Available
Vary	Document Data Engine	ArXiv, CC-MAIN, E-books	2M	✗
	Chart Data Engine	matplotlib, pyecharts, NLP corpora	1.5M	✗
	Detection Data Engine	Objects365, OpenImages	~3M	✓
LLaVAR	LAION	LAION images filtered for text-rich content, OCR applied	0.4M	✓
DoCo	DoCo-Processed	CC3M (LLaVA) + LAION, processed with PaddleOCR	1.0M	✗
Texthawk2	100M pretraining	Diverse, mainly public datasets	100M	✗
Docpedia	PDF Images	arXiv (public scientific preprints)	325K	✓
	PPT Images	Common Crawl (web-crawled PPTs)	600K	Partly

Table 6: Summary of pretraining datasets created and used in recent MLLM-based VRDU frameworks.

Framework	Category	Source / Description	Size (K)	Open Source
Leopard	Multi-image (text-rich)	69K public multi-page docs/slides; Adapted single-page to multi-image (DocVQA, ArxivQA); Raw slides + GPT-4o QAs; Multi-chart/table (open, synth.); Webpage snapshots (Mind2Web, OmniACT, WebScreenshots, etc.)	739	Partially
	Single-image	Text-rich single images from public datasets; Natural images (e.g., ShareGPT4V, etc.)	186	Partially
LLaVAR	Noisy Instruction-Following	Text-rich images from LAION, selected via classifier + CLIP clustering, instructions via OCR-based prompts	422,000	Yes
	High-Quality Instruction-Following	Subset of LAION text-rich images (4 clusters), multi-turn QAs generated by prompting text-only GPT-4 with OCR+caption info	16,000	Yes

Table 7: Summary of instruction-tuning datasets for Leopard and LLaVAR.

Dataset	Venue	Year	Domain	Docs	Images	Keys / Qs	Multi page	Language	Metrics	Format
Key Information Extraction										
FUNSD	ICDAR-w	2019	Multi-source	-	199	4	✗	English	F1	P, H
SROIE	ICDAR-c	2019	Scanned Receipts	-	973	4	✗	English	F1*	P
CORD	NeurIPS-w	2019	Scanned Receipts	-	1,000	54	✗	English	F1	P
Payment-Invoice	ACL	2020	Invoice Form	-	14,832	7	✗	English	F1	D
Payment-Receipts	ACL	2020	Scanned Receipts	-	478	2	✗	English	F1	P
Kleister-NDA	ICDAR	2021	Private Agreements	540	3,229	4	✓	English	F1	D
Kleister-Charity	ICDAR	2021	AFR	2,778	61,643	8	✓	English	F1	D, P
EPHOIE	AAAI	2021	Exam Paper	-	1,494	10	✗	Chinese	F1	P, H
XFUND	ACL	2022	Synthetic Forms	-	1,393	4	✗	Multilingual	F1	D, P, H
Form-NLU	SIGIR	2023	Financial Form	-	857	12	✗	English	F1	D, P, H
VRDU-Regist. Form	KDD	2023	Registration Form	-	1,915	6	✗	English	F1	D
VRDU-Ad-buy Form	KDD	2023	Political Invoice Form	-	641	9+1(5)	✗	English	F1	D, P
DocILE	ICDAR	2023	Invoice Form	6,680	106,680	55	✓	English	AP, CLEval	D, P
KVP10k	ICDAR	2024	Cross-domain	-	10,707	118,868	✗	English	F1, IOU	D, H
CC-OCR-KIE	ICCV	2025	Cross-domain	-	2,008	34(-)	✗	Multilingual	F1	D, P, H
Visual Question Answering										
DocVQA	WACV	2021	Industrial Reports	-	12,767	50,000	✗	English	ANLS	D, P, H
VisualMRC	AAAI	2021	Website	-	10,197	30,562	✗	English	BLEU, etc	D
TAT-DQA	MM	2022	Financial Reports	2,758	3,067	16,558	✓	English	EM, F1	D
RDVQA	MM	2022	Data Analysis Report	8,362	8,514	41,378	✗	English	ANLS, ACC	D
CS-DVQA	MM	2022	Industry Documents	-	600	1,000	✗	English	ANLS	D, P, H
InfographicVQA	WACV	2022	Infographics	-	5,400	3,000	✗	English	ANLS, F1	D
PDFVQA-Task A	ECML-PKDD	2023	Academic Paper	-	12,337	81,085	✗	English	F1	D
PDFVQA-Task B	ECML-PKDD	2023	Academic Paper	-	12,337	53,872	✗	English	F1	D
PDFVQA-Task C	ECML-PKDD	2023	Academic Paper	1,147	12,337	5,653	✓	English	EM	D
MPDocVQA	PR	2023	Industrial Reports	6,000	48,000	46,000	✓	English	ANLS	D, P, H
DUDE	ICCV	2023	Cross-domain	5,019	28,709	41,541	✓	English	ANLS	D
SlideVQA	AAAI	2023	Slide, decks	-	5,200	14,500	✓	English	EM, F1	D
MMLONGBENCH-DOC	NIPS	2024	Cross-domain	135	6,413	1,082	✓	English	ACC, F1	D
MMVQA	IJCAI	2024	Academic Paper	3,146	30,239	262,928	✓	English	EM, PM, MR	D
JDocQA	LREC-COLING	2024	Cross-Domain	5,504	268,000	11,600	✓	Japanese	F1	D
BoundingDocs	IJDAR	2025	Cross-domain, Mixed	48,151	237,437	249,016	✗	Multilingual	ANLS	D, P, H
LongDocURL	ACL	2025	Cross-domain	396	33,000	2,325	✓	English	F1	D
MMDocIR	EMNLP	2025	Cross-domain	6,878	224,223	73,843	✓	Multilingual	F1	D
MTVQA	EMNLP	2025	Cross-domain	-	8,794	28,607	✗	Multilingual	ANLS	D, P, H
SEEPHYS	NIPS	2025	Physics	-	2,245	2,000	✗	English	Accuracy	D

Table 8: Benchmark datasets for Key Information Extraction and Visual Question Answering in visually rich documents. P - Scanned Printed, H - Scanned Handwritten, D - Digital Born

Model Name	FUNSD	CORD	SROIE	DocVQA	ChartVQA	InfoVQA
General Domain LLM						
Qwen1.5-7B-Chat	52.5	29.7	–	64.3	–	–
Llama3-8B-Instruct	57.5	40.0	–	74.2	–	–
General Domain MLLM						
QwenVL-7B	47.1	30.0	–	65.1	–	–
InterVL2-8B	75.8	79.9	–	91.7	–	–
Claude-3.5 Sonnet	–	–	–	88.5	51.8	59.1
GeminiPro-1.5	–	–	–	91.2	34.7	73.9
GPT4o 20240806	–	–	–	92.8	85.7	66.4
OCR-Dependent						
DocLLM (2024a)	51.8	67.4	91.9	69.5	–	–
LAPDoc (2024)	–	–	–	79.8	–	54.9
DoCo (2024)	–	–	–	64.8	68.9	34.9
InstructDr (2024)	38.1	62.7	–	22.3	–	37.6
LayoutLLM (2024)	78.7	62.2	71.0	74.3	–	–
LLaVA-Read (2024c)	36.9	–	58.3	71.0	74.6	36.4
LayTextLLM (2024)	64.0	96.5	95.8	77.2	–	–
LayTokenLLM(2025b)	71.0	75.4	–	85.1	–	–
GPE (2025a)	82.6	86.9	97.8	78.1	–	–
PDF-WuKong (2025)	85.1	–	–	76.9	80.0	61.3
DocLayLLM (2024)	80.7	79.4	84.4	72.8	–	–
AlignVLM (2025)	–	–	–	81.2	75.0	53.8
DocAssistant (2025)	–	–	–	89.8	81.4	66.7
DocThinker (2025)	–	–	81.4	80.2	–	69.7
OCR-Free						
KOSMOS-2.5 (2023)	–	–	–	81.1	62.3	41.3
mPLUG-DocOwl (2025)	–	–	–	62.2	57.4	38.2
UReader (2023b)	–	–	–	65.4	59.3	42.2
TGDoc (2023)	1.7	–	3.0	9.0	11.7	12.8
UniDoc (2023)	1.2	–	1.4	6.5	10.5	13.8
DocPedia (2024)	40.1	–	57.7	49.3	47.8	15.5
HRVDA (2024a)	–	89.3	89.3	91.0	72.1	43.5
Vary-base (2024)	–	–	–	76.3	66.1	–
mPLUG-DocOwl 1.5 (2024)	–	–	–	81.6	70.5	50.4
HVFA (2024)	–	–	–	72.7	63.3	45.9
mPLUG-DocOwl2 (2025)	–	–	–	80.7	70.0	46.4
Texthawk (2024a)	–	–	–	76.4	66.6	50.6
Texthawk2 (2024b)	–	–	–	89.6	81.4	67.8
TextMonkey (2024c)	65.5	67.5	47.0	73.0	66.9	28.6
Llavar-7B (2024d)	1.7	13.6	2.4	11.6	–	–
TokenCorrCompressor (2024b)	–	–	–	78.3	68.9	50.2
DocKylín (2024a)	25.5	–	49.5	77.3	66.8	46.6
Marten (2025b)	44.4	–	80.4	92.0	81.7	75.2
PP-DocBee (2025)	–	–	–	90.6	74.6	66.2
TokenFD (2025)	42.2	–	81.9	94.2	86.6	76.5

Table 9: Performance comparison between OCR-dependent and OCR-free document understanding frameworks across benchmark datasets.

Model	Type	Venue	Year	MPDocVQA	DUDE
Longformer	General VLPM	Preprint	2020	55.1	20.3
BigBird	General VLPM	NeurIPS	2020	58.5	26.3
GPT-4v	General MLLM	–	2023	–	53.9
Idefics3-8B	General MLLM	Preprint	2024	67.2	38.7
LLaVA-next-interleave-7B	General MLLM	Preprint	2024	44.9	28.0
Hi-ViT5	OCR-dependent VLPM	PR	2023	61.8	35.7
GRAM	OCR-dependent VLPM	CVPR	2024	83.0	53.4
InstructDoc	OCR-Dependent VLPM	AAAI	2024	–	46.8
mPLUG-DocOwl2	OCR-free VLPM	Preprint	2024	69.4	46.7
PDF-WuKong	OCR-Dependent VLPM	Preprint	2024	76.9	56.1
LayTokenLLM	OCR-Dependent VLPM	CVPR	2025	74.3	52.0
DocThinker (2025)	OCR-Dependent VLPM	ICCV	2025	–	56.8

Table 10: Performance comparison of state-of-the-art models on MPDocVQA and DUDE benchmarks. Best scores are highlighted in red.