

MPLUG-DocOwl2: HIGH-RESOLUTION COMPRESSING FOR OCR-FREE MULTI-PAGE DOCUMENT UNDERSTANDING

Anwen Hu¹ Haiyang Xu^{1*} Liang Zhang² Jiabo Ye¹ Ming Yan^{1*}
 Ji Zhang¹ Qin Jin² Fei Huang¹ Jingren Zhou¹
¹Alibaba Group ²Renmin University of China
 {huanwen.haw, shuofeng.xhy, ym119608}@alibaba-inc.com
<https://github.com/X-PLUG/mPLUG-DocOwl>

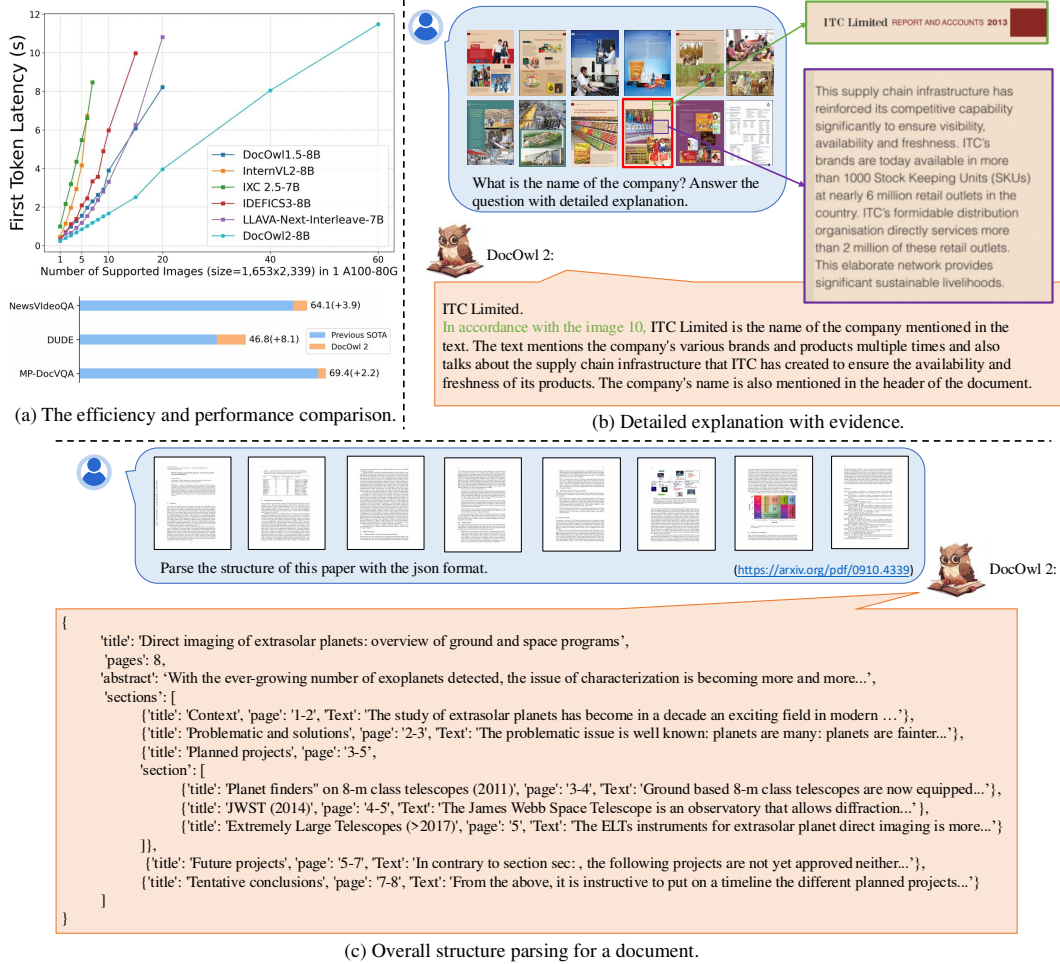


Figure 1: (a) mPLUG-DocOwl2 achieves state-of-the-art Multi-page Document Understanding performance with faster inference speed and less GPU memory; (b-c) mPLUG-DocOwl2 is able to provide a detailed explanation containing the evidence page as well as the overall structure parsing of the document.

ABSTRACT

Multimodal Large Language Models (MLLMs) have achieved promising OCR-free Document Understanding performance by increasing the supported resolution of document images. However, this comes at the cost of generating thousands of visual tokens for a single document image, leading to excessive GPU

*Corresponding author

memory and slower inference times, particularly in multi-page document comprehension. In this work, to address these challenges, we propose a High-resolution DocCompressor module to compress each high-resolution document image into 324 tokens, guided by low-resolution global visual features. With this compression module, to strengthen multi-page document comprehension ability and balance both token efficiency and question-answering performance, we develop the DocOwl2 under a three-stage training framework: Single-image Pretraining, Multi-image Continue-pretraining, and Multi-task Finetuning. DocOwl2 sets a new state-of-the-art across multi-page document understanding benchmarks and reduces first token latency by more than 50%, demonstrating advanced capabilities in multi-page questioning answering, explanation with evidence pages, and cross-page structure understanding. Additionally, compared to single-image MLLMs trained on similar data, our DocOwl2 achieves comparable single-page understanding performance with less than 20% of the visual tokens. Our codes, models, and data are publicly available at <https://github.com/X-PLUG/mPLUG-DocOwl/tree/main/DocOwl2>.

1 INTRODUCTION

Understanding a multi-page document or news video is common in human daily life. To tackle such scenarios, Multimodal Large Language Models (MLLMs) (Ye et al., 2023c;d; 2024; Bai et al., 2023; Liu et al., 2023) should be equipped with the ability to understand multiple images with rich visually-situated text information. Different from natural images mainly comprising of objects, comprehending document images asks for a more fine-grained perception to recognize all texts. To tackle high-resolution document images, some works (Hong et al., 2023; Wei et al., 2023) propose to add an additional high-resolution encoder while more works (Ye et al., 2023b; Hu et al., 2024; Chen et al., 2024; Dong et al., 2024b;a) choose to crop a high-resolution image to low-resolution sub-images and let the Large Language Model to understand their relationship. By increasing the cropping number, the latter achieves better performance of OCR-free document understanding but also results in too many visual tokens for only 1 document image, e.g., InternVL 2 (Chen et al., 2024) costs a average of 3k visual tokens on single-page document understanding benchmark DocVQA (Mathew et al., 2021). As shown in Fig. 1(a), such long visual tokens not only result in long inference time but also occupy too much GPU memory, making it difficult to understand a complete document or video and greatly limiting their application scenarios. Inspired by Natural Language Processing work (Cheng et al., 2024; Ge et al., 2024; Chevalier et al., 2023) which summarizes a textual paragraph/document into fewer tokens and maintains most semantics, we argue that visual tokens of document images can also be further compressed while maintaining both layout and most textual information.

Existing compressing architecture in MLLMs are hard to balance information retention and token efficiency during document image encoding. As shown in Fig. 2(a), independently compressing each crop of a document image (Li et al., 2024b; Hu et al., 2024) could reduce visual tokens of each sub-image but still results in a long sequence of visual tokens after concatenating all sub-images. Leveraging learnable queries (Bai et al., 2023; Li et al., 2023a; Ye et al., 2023c) or selected tokens (Liu et al., 2024) as compressing guidance could produce an identical length of tokens for any resolution but overlook the overall layout information, as shown in Fig. 2(b). Layout-aware guidance is important for compressing visual features of document images because texts within a layout region are semantic-coherent and easier to summarize. For example, in a two-column paper, texts belonging to the ‘Related Work’ section are difficult to summarize with texts on the same line but belonging to the ‘Method’ section.

In this work, as shown in Fig. 2(c), we propose a layout-aware compressing architecture **High-resolution DocCompressor** based on cross-attention to compress document images into fewer tokens and achieve better performance than existing compressing methods. Considering that a global low-resolution image can well capture the overall layout information, we utilize visual features of a global low-resolution image as the compressing guidance (query). Each visual feature in the global feature map just captures the layout information of partial regions. Therefore, each query attending to all high-resolution features will not only make information compression more difficult but also increase computation complexity. To summarize text information within a layout region, for each query from the global feature map, a group of high-resolution features with identical relative posi-

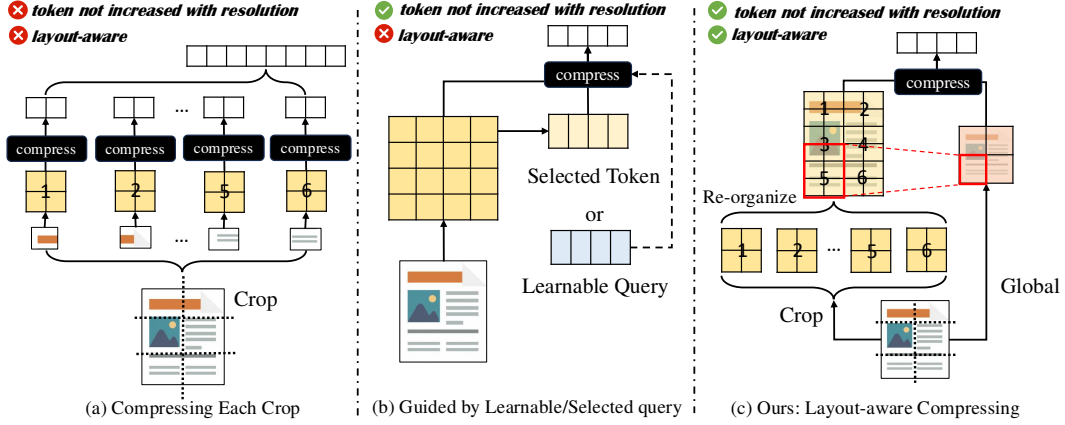


Figure 2: Illustrations of different compressing methods for OCR-free document understanding.

tions in the raw image is collected as compressing objects, sometimes spanning multiple sub-images. Besides, since the vision-to-text (V2T) module of MLLMs could convert visual features into textual feature space, we argue that compressing visual features after the vision-to-text module could better maintain textual semantics in document images. Therefore, based on the architecture of DocOwl 1.5 (Hu et al., 2024), we propose mPLUG-DocOwl2 by placing the High-resolution DocCompressor after its V2T module: H-Reducer. To take full advantage of the compressing method, our model DocOwl2 is trained with a three-stage framework: Single-image Pretraining, Multi-image Continue-Pretraining, and Multi-task Finetuning to support both single-image and multi-image/frame understanding. Our experiments on single-page and multi-page document benchmarks demonstrate the good balance of OCR-free document understanding performance and token efficiency of DocOwl2. We perform sufficient ablation studies to validate the superiority of our High-resolution DocCompressor and the benefits of the three-stage training framework for both single-page and multi-page understanding performance.

Our contributions in this work are three-fold:

- We propose a novel compressing architecture, namely High-resolution DocCompressor, to greatly reduce visual tokens of high-resolution document images. Compared with existing compressing methods, our method achieves better OCR-free single-image document Understanding performance with fewer visual tokens.
- DocOwl2 achieves state-of-the-art performance on Multi-page Document understanding benchmarks with <50% First Token Latency.
- Compared with state-of-the-art MLLMs with similar model size and training data, DocOwl2 achieves comparable performance with <20% visual tokens on 10 single-image document benchmarks.

2 RELATED WORK

2.1 OCR-FREE VISUAL DOCUMENT UNDERSTANDING

Visual Document Understanding aims to comprehend images with rich text information, including scans of document pages (Mathew et al., 2021; Tito et al., 2022; Landeghem et al., 2023; Zhang et al., 2023; Wei et al., 2023), infographics (Mathew et al., 2022), charts (Masry et al., 2022; Kafle et al., 2018; Methani et al., 2020; Kahou et al., 2018), tables images (Pasupat & Liang, 2015; Chen et al., 2020; Zhong et al., 2020), webpage screenshots (Tanaka et al., 2021; Chen et al., 2021) and natural images with scene texts (Singh et al., 2019; Sidorov et al., 2020; Hu et al., 2021). Recently, many Multimodal Large Language Models have been proposed to perform visual document understanding in an OCR-free manner. mPLUG-DocOwl (Ye et al., 2023a) and UReader (Ye et al., 2023b) first propose to unify different tasks across 5 types of document images in the seq-to-seq format.

To encode rich text information in high-resolution images, UReader (Ye et al., 2023b) proposes a Shape-adaptive Cropping Module to cut the raw image into multiple low-resolution sub-images and utilizes an identical low-resolution encoder to encode both sub-images and a global image. Monkey (Li et al., 2023b) proposes to employ a sliding window to partition high-resolution images and a resampler to reduce redundant information of each sub-image. mPLUG-DocOwl1.5 (Hu et al., 2024) increases the basic resolution of the low-resolution encoder and replaces the Visual Abstractor (Ye et al., 2023c) with 1 simple convolution layer to better maintain the structure information. DocPedia (Feng et al., 2023) directly processes high-resolution images in the frequency domain. CoAgent (Hong et al., 2023) proposes to utilize a high-resolution encoder to encode high-resolution visual features and a low-resolution encoder to encode low-resolution global features. Series work of InternLM-XComposer (Dong et al., 2024a;b) and InternVL (Chen et al., 2024) further optimize the cropping method or increase the cropping number and greatly improves the OCR-free Document Understanding performance. These works achieve promising performance but suffer from too many visual tokens for a high-resolution image (always $>1k$ tokens for a common A4-sized document page), which hinders the development of OCR-free multi-page document understanding.

2.2 VISUAL FEATURE COMPRESSING

Reducing visual tokens of a single image enables a Multimodal Large Language Model with limited maximum sequence length to leverage more images as contexts to perform complex multimodal tasks, such as video understanding, embodied interaction, or multi-page document understanding. There have been some architectures proposed for compressing visual features of general images with fewer learnable queries, such as the Resampler (Alayrac et al., 2022; Bai et al., 2023), Abstractor (Ye et al., 2023c;d) and Q-former (Li et al., 2023a). Randomly initialized Learnable queries can ensemble object information in general images but is hard to summarize rich text information in high-resolution document images. As a compromise solution, TokenPacker (Li et al., 2024b) proposes to compress each sub-image with its downsampled visual features as the query to perform cross-attention. TokenPacker just reduces each sub-image’s visual tokens, thus still creates more than 1k visual tokens when processing high-resolution document images. TextMonkey (Liu et al., 2024) first filters valuable visual tokens and then uses them as guidance to aggregate all visual tokens. Due to that valuable visual tokens are selected by measuring the token similarity, visual information of partial regions may not be covered and thus not well compressed during following cross-attention. In this work, our High-resolution DocCompressor leverages visual features from the row-resolution global images as the query, the ensembled feature map of sub-images as key and value. This not only produces a fixed number of visual tokens for images of any resolution but also covers all areas during compression. Compared to Mini-Gemini (Li et al., 2024c) which compresses general visual features, there are major two differences with our DocOwl2. Firstly, we make full use of global visual features and sub-image features produced by an identical low-resolution vision encoder and don’t need to add an extra high-resolution encoder. Secondly, for better summarizing textual information in document images, our cross-attention is applied based on visual features that have been aligned with textual features of LLM. We argue that directly compressing outputs of the vision encoder will lose more visually situated textual information while comprising features aligned with LLM is like summarizing texts (Cheng et al., 2024; Ge et al., 2024; Chevalier et al., 2023) and can better maintain textual semantics in document images. Fair comparisons are performed in our experiments to support our hypothesis.

3 MPLUG-DOCOWL2

As shown in Fig. 3, DocOwl2 leverages a Shape-adaptive Cropping Module and a low-resolution vision encoder to encode high-resolution document images. Then, it utilizes a vision-to-text module H-Reducer to ensemble horizontal visual features and align the dimension of vision features with Large Language Models. Furthermore, a high-resolution compressor is designed to greatly reduce the number of visual features while maintain most visual information. Finally, compressed visual tokens of multiple images/pages are concatenated with text instructions and input to a Large Language Model for multimodal understanding.

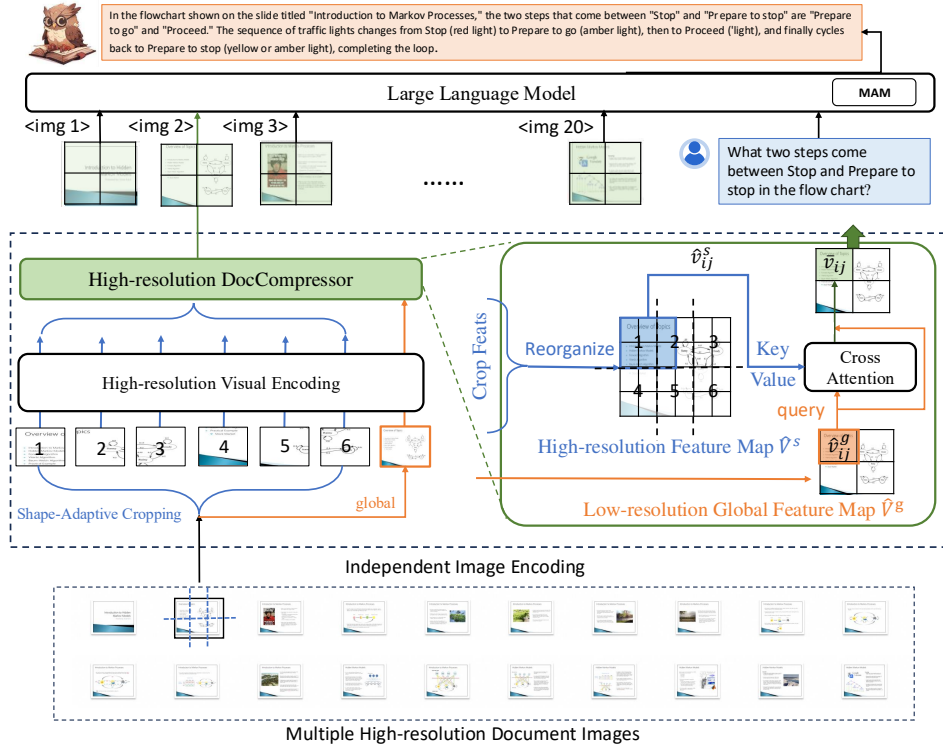


Figure 3: The architecture of DocOwl2. Each image is independently encoded by the pipeline of Shape-adaptive Cropping, High-resolution Visual Encoding and High-resolution DocCompressor.

3.1 HIGH RESOLUTION VISION ENCODING

Following UReader (Ye et al., 2023b) and DocOwl 1.5 (Hu et al., 2024), DocOwl2 utilizes a parameter-free Shape-adaptive Cropping Module to preprocess high-resolution images. Concretely, it cuts each high-resolution image I into $R \times C$ size-fixed sub-images $I^s = \{I_{xy}^s\}, 1 \leq x \leq R, 1 \leq y \leq C$, where cropping rows R and columns C are flexibly decided based on the raw resolution of I . Besides, to maintain the overall layout information, the raw image is also directly resized to a global image I^g . Both the global image and sub-images are sized $H \times W$.

After the cropping module, a low-resolution transformer-based vision encoder ViT (Dosovitskiy et al., 2021) is utilized to independently extract vision features of each sub-image and the global image as follows:

$$V^g = \text{ViT}(I^g) \quad (1)$$

$$V_{xy}^s = \text{ViT}(I_{xy}^s), 1 \leq x \leq R, 1 \leq y \leq C, \quad (2)$$

where both V^g and V_{xy}^s are visual features with the shape of $h \times w \times d$, d is the feature dimension and w, h are the width and height of the feature map.

Following DocOwl 1.5, after the ViT, for each sub-image or global image, we apply a vision-to-text module H-Reducer to ensemble horizontal 4 features by a convolution layer and align the feature dimension with the Large Language Model with a fully connected layer. The calculation of H-Reducer is represented as follows:

$$\hat{V} = \text{FC}(\text{Conv}(V)), V \in \{V^g, V_{xy}^s\}, 1 \leq x \leq R, 1 \leq y \leq C, \quad (3)$$

where the shape of the visual feature map \hat{V} is $h \times \frac{w}{4} \times \hat{d}$, \hat{d} is the dimension of hidden states of the large language model.

3.2 HIGH RESOLUTION FULL-COMPRESSING

Although the H-Reducer has reduced the visual tokens of each sub-image or global image to $\frac{1}{4}$ the length of original visual features, the token length of high-resolution images is still too long to perform multi-page/image joint understanding for Large Language Models. For example, the token length of 1 high-resolution image in DocOwl 1.5 (Hu et al., 2024) is $(R \times C + 1) \times h \times \frac{w}{4}$, which will be 2,560 when the raw resolution is $1,344 \times 1,344$.

In Natural Language Processing, a sentence/paragraph/document of text tokens can be compressed into fewer summary vectors while maintaining most semantics (Cheng et al., 2024; Ge et al., 2024; Chevalier et al., 2023). Besides, since visual features have been aligned with the textual feature space of large language models, the visual tokens of document images after the vision-to-text module can also be treated as textual tokens encoding different parts of textual information in the image. Thus, taking into account these two points, in this work, we argue that visually situated textual information of document images can also be further compressed into fewer tokens, especially after the vision-to-text alignment.

Ideally, the compression of visual texts should be based on their layout. Texts from the same layout region (e.g., a title/paragraph region) are more appropriate to be fused into an identical token. After the vision-to-text module H-Reducer, the global visual feature \hat{V}^g mainly encodes the overall text layout information while visual features of sub-images $\{\hat{V}_{xy}^s\}$ capture detailed textual information. Besides, due to both the global image and cropped sub-images come from an identical image, there is a clear mapping between the visual tokens of \hat{V}^g and $\{\hat{V}_{xy}^s\}$. As shown in Fig. 3, each visual token in \hat{V}^g can be aligned with $R \times C$ visual tokens in $\{\hat{V}_{xy}^s\}$. Therefore, in this work, with global visual features as query, and the visual features from sub-images as key and value, we propose to utilize cross-attention to ensemble textual semantics and greatly reduce the number of visual tokens of a high-resolution image to the one of a low-resolution global image.

Concretely, we first re-organize feature maps of cropping images ($\{\hat{V}_{xy}^s, 1 \leq x \leq R, 1 \leq y \leq C\}$) to a complete feature map \hat{V}^s according to their positions in the raw high-resolution image. Then, for each visual token in the feature map \hat{V}^g of the global image, we collect its corresponding $R \times C$ visual tokens from \hat{V}^s as the key and value, the cross-attention layer in this compressor is calculated as follows:

$$\hat{v}_{ij}^g \in \hat{V}^g, 1 \leq i \leq h, 1 \leq j \leq w/4 \quad (4)$$

$$\hat{v}_{ij}^s = [\hat{v}_{i',j'}^s] \subset \hat{V}^s, (i-1)R+1 \leq i' \leq iR, (j-1)C+1 \leq j' \leq jC \quad (5)$$

$$\bar{v}_{ij} = \text{softmax}\left(\frac{W^q \hat{v}_{ij}^g W^k \hat{v}_{ij}^s}{\sqrt{d_k}}\right) W^v \hat{v}_{ij}^s + \hat{v}_{ij}^g \quad (6)$$

where \hat{v}_{ij}^g is a visual token from the feature map of the global image, \hat{v}_{ij}^s are visual tokens from the re-organized feature map of cropping images. \hat{v}_{ij}^g and \hat{v}_{ij}^s correspond to the same area in the raw image. W^q, W^k, W^v are learnable projection matrices.

After high-resolution compressing, the compressed feature map of each image is organized into a sequence $\bar{V} = [\bar{v}_1, \bar{v}_2, \dots, \bar{v}_{h \times \frac{w}{4}}]$ for subsequent understanding of the large language model.

3.3 MULTI-IMAGE MODELING WITH LLM

Through the high-resolution compressing, the number of visual tokens for each high-resolution image is reduced from $(R \times C + 1) \times h \times \frac{w}{4}$ to $h \times \frac{w}{4}$. Such efficient vision encoding allows joint understanding of multiple document images with Large Language Models. To help the LLM better distinguish visual features from different images and understand the ordinal number of images, we add a textual ordinal token '``' before the visual features of each image, where x is the ordinal number. Overall, the decoding of the decoder for multiple images is as follows:

$$Y = \text{LLM}([P_0; \bar{V}_0; P_1; \bar{V}_1, \dots, P_n; \bar{V}_n; T]) \quad (7)$$

where $[\cdot]$ means the concatenation operation, n is the number of images, $P_x, 1 \leq x \leq n$ is the textual embedding of the ordinal token ' $\langle \text{img } x \rangle$ ', \bar{V}_x is the visual features for each image, T is the textual instruction and Y is the predicted answer.

3.4 MODEL TRAINING

DocOwl2 is trained with three stages: Single-image Pre-training, Multi-image Continue Pretraining, and Multi-task Finetuning.

At the first stage, to ensure the compressed visual tokens can encode most visual information, especially visually situated texts, we first perform Unified Structure Learning as DocOwl 1.5 with the dataset DocStruct4M (Hu et al., 2024), which covers the learning of struct-aware document parsing, table parsing, chart parsing and natural image parsing of a single image.

After Single-image Pretraining, to empower our model with the ability to correlate multiple images, we further perform Multi-image Continue Pretraining with a struct-aware multi-page document parsing dataset MP-DocStruct1M. With partial documents from two datasets of PixParse¹², we design two symmetrical tasks of multi-image understanding: Multi-page Text Parsing and Multi-page Text Lookup. Given successive page images in a document, the Multi-page Text Parsing instructs the model to parse texts of specified one or two pages, such as 'Recognize texts in image 2 and image 10.'. As for the Multi-page Text Lookup task, with texts from 1-2 pages as input, the model is required to predict the concrete ordinal number of images containing these texts, for example, 'Looking for the image with text $\langle \text{doc} \rangle \dots \langle \text{doc} \rangle$ and $\langle \text{doc} \rangle \dots \langle \text{doc} \rangle$ '. Besides MP-DocStruct1M, during this stage, we also randomly chose 0.5M samples from DocStruct4M to avoid the catastrophic forgetting of structure parsing across different types of images.

Finally, we ensemble single-image and multi-image instruction tuning datasets to perform multi-task tuning. We leverage DocDownstream-1.0 (Hu et al., 2024) and DocReason25K (Hu et al., 2024) as single-image datasets. DocDownstream-1.0 is an ensemble dataset comprising of DocVQA (Mathew et al., 2021), InfoVQA (Mathew et al., 2022), DeepForm (Svetlichnaya, 2020), KLC (Stanislawek et al., 2021), WTQ (Pasupat & Liang, 2015), TabFact (Chen et al., 2020), ChartQA (Masry et al., 2022), TextVQA (Singh et al., 2019), TextCaps (Sidorov et al., 2020) and VisualMRC (Tanaka et al., 2021). DocReason25K is a question-answering dataset with detailed explanations. As for multi-image understanding, we ensemble 2 document datasets, MP-DocVQA (Tito et al., 2022) and DUDE (Landeghem et al., 2023), and 1 news video dataset NewsVideoQA (Jahagirdar et al., 2023) as concise question-answering datasets. MP-DocVQA contains 46k question-answering pairs on 60k page images scanned from 6k industry documents with rich tables, diagrams, pictures, and both handwritten and printed texts. DUDE covers more domains of documents, including medical, legal, technical, financial, etc. It contains 41k question-answering pairs on 5k documents. NewsVideoQA collects news videos with rich visually-situated texts from diverse English news channels around the world, such as BBC, CNN, etc. It contains 8k question-answering pairs framed on 3k videos. Besides, to trigger the ability of detailed explanations with evidence pages, we built MP-DocReason51K based on DocReason25K. Concretely, for each single-image sample from DocReason25K, we construct two multi-image samples with noisy images randomly chosen from the same or different categories. After randomly inserting the evidence image into noisy images, we add an extra evidence description (e.g., 'According to the 5th image,') into the raw detailed explanation to get the target of multi-image samples. Most question-answering samples just focus on 1-2 pages of a document, to further strengthen the ability of a comprehensive understanding of a document, we leverage a small part of annotations from DocGenome (Xia et al., 2024) to construct text sequences in the JSON format, which represents the hierarchical structure of a scientific paper and partial detailed texts.

The detailed statistics of training datasets of DocOwl2 are shown in Table 1.

¹<https://huggingface.co/datasets/pixparse/idl-wds>

²<https://huggingface.co/datasets/pixparse/pdfa-eng-wds>

Table 1: Detailed statistic of training datasets of DocOwl2.

Training Stage	Input Image	Dataset	Num
Single-image Pretraining	Single	DocStruct4M	4,036,402
Multi-image Continue Pretraining	Single	DocStruct4M	501,781
	Multiple	MP-DocStruct1M	1,113,259
Multi-task Finetuning	Single	DocVQA, InfoVQA, DeepForm, KLC, WTQ, TabFact, ChartQA, TextVQA, TextCaps, VisualMRC	552,315
		DocReason25K	25,877
		MP-DocVQA	70,154
	Multiple	DUDE	35,438
		NewsVideoQA	8,619
		MP-DocReason51K	51,754
		DocGenome12K	12,010

Table 2: Comparison with OCR-free methods on single-image document understanding tasks. The ‘*’ refers to models without LLMs and separately fine-tuned on each downstream task. ‘Token^V’ means the average number of visual tokens of a single image. ‘**Bold**’ means SOTA performance within the group and ‘Underline’ means achieving 80% SOTA performance among all baselines.

	Model	Size	Token ^V	Doc VQA	Info VQA	Deep Form	KLC	WTQ	Tab Fact	Chart QA	Text VQA	Text Caps	Visual MRC
Token ^V ≥ 1k	Donut*	<1B	4,800	67.5	11.6	61.6	30.0	18.8	54.6	41.8	43.5	74.4	93.91
	Pix2Struct [*] _{base}	<1B	2,048	72.1	38.2	-	-	-	-	56.0	-	88.0	-
	Pix2Struct [*] _{large}	1B	2,048	76.6	40.0	-	-	-	-	58.6	-	95.5	-
	CogAgent	17B	6,656	81.6	44.5	-	-	-	-	68.4	76.1	-	-
	IXC 2.5	7B	~ 5,118	90.9	69.9	71.2	-	53.6	85.2	82.2	78.2	-	307.5
	InternVL 2	8B	~ 3,133	91.6	74.8	-	-	-	-	83.3	77.4	-	-
	TokenPacker	13B	~ 1,833	70.0	-	-	-	-	-	-	-	-	-
	DocOwl 1.5	8B	~ 1,698	82.2	50.7	68.8	38.7	40.6	80.2	70.2	68.6	131.6	246.4
	DocPeida	7B	1,600	47.1	15.2	-	-	-	-	46.9	60.2	-	-
	Monkey	9B	1,280	66.5	36.1	40.6	32.8	25.3	-	-	64.3	93.2	-
	DocOwl	7B	~ 841	62.2	38.2	42.6	30.3	26.9	60.2	57.4	52.6	111.9	188.8
	UReader	7B	~841	65.4	42.2	49.5	32.8	29.4	67.6	59.3	57.6	118.4	221.7
	TextMonkey	9B	768	73.0	28.6	59.7	37.8	31.9	-	66.9	65.9	-	-
	TokenPacker	13B	~ 467	58.0	-	-	-	-	-	-	-	-	-
Token ^V < 1k	QwenVL	9B	256	65.1	35.4	-	-	-	-	65.7	63.8	-	-
	Vary	7B	256	76.3	-	-	-	-	-	66.1	-	-	-
	DocOwl2	8B	324	80.7	46.4	66.8	<u>37.5</u>	36.5	78.2	70.0	66.7	131.8	217.4

4 EXPERIMENTS

4.1 IMPLEMENTATION DETAILS

The maximum number of crops is set to 12. The resolution of each sub-image or the global image is 504x504. The High-resolution DocCompressor comprises of 2 layers of cross attention. Initialized from mPLUG-Owl2 (Ye et al., 2023d), the vision encoder (ViT/L-14 (Dosovitskiy et al., 2021)), H-Reducer and High-resolution DocCompressor are trained during the Single-image Pretraining. Besides, the main parameters of the Large Language Model (Touvron et al., 2023) are frozen while a Modality Adaptive Module (MAM) (Ye et al., 2023d) used to distinguish visual and textual features in the LLM is tuned. The first stage is trained 12k steps with a batch size of 1,024 and the learning rate set as 1e-4. During the Multi-image Continue-pretraining, the vision encoder is further frozen and the H-Reducer, High-resolution DocCompressor and MAM is tuned. The second stage is trained 2.4k steps with a batch size of 1,024 and the learning rate set as 2e-5. At the final Multi-task Finetuning stage, all parameters except the vision encoder are optimized. The batch size, training step, and learning rate at this stage are set as 256, 9k, and 2e-5, respectively.

4.2 MAIN RESULTS

We compare DocOwl2 with state-of-the-art Multimodal Large Language Models on 10 single-image document understanding benchmarks, 2 Multi-page document Understanding benchmarks, and 1

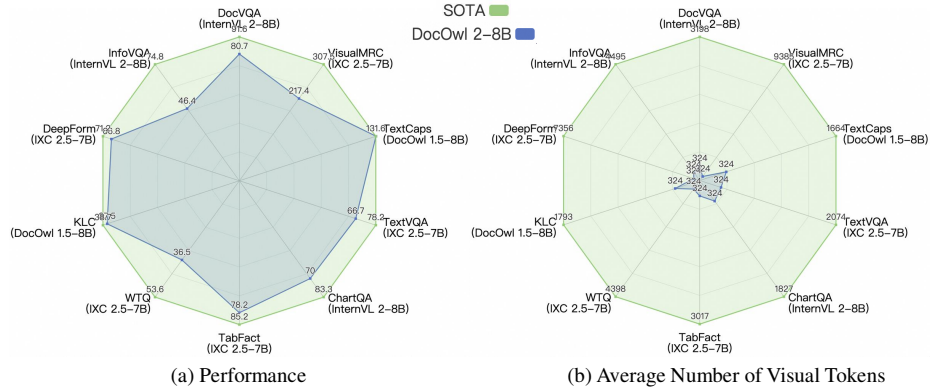


Figure 4: The comparison of our DocOwl2 with state-of-the-art Multimodal Large Language Models on (a) OCR-free performance and (b) the average number of visual tokens on 10 Visual Document Understanding benchmarks.

Table 3: Comparison with OCR-free Multimodal Large Language Models on single-image document understanding benchmarks. ‘FTL(s)’ refers to the First Token Latency (seconds)

Model	Size	Token ^V	DocVQA			Token ^V	ChartQA			Token ^V	TextVQA		
			FTL(s)↓	ANLS↑			FTL(s)↓	ANLS↑			FTL(s)↓	ANLS↑	
InternVL 2	8B	~3,198	0.94	91.6		~1,827	0.56	83.3		~2,864	1.01	77.4	
IXC 2.5	7B	~7,395	3.73	90.9		~1,971	1.05	82.2		~2,075	1.11	78.2	
DocOwl 1.5	8B	~1,806	0.58	82.2		~1,713	0.53	70.2		~1,664	0.56	68.6	
TextMonkey	9B	768	0.58	73.0		768	0.51	66.9		768	0.50	65.9	
DocOwl2	8B	324	0.26	80.7		324	0.21	70.0		324	0.23	66.7	

text-rich video understanding benchmark. Both question-answering performance and the First Token Latency (seconds) are considered to show the effectiveness of our model.

4.2.1 SINGLE-IMAGE DOCUMENT UNDERSTANDING

For Single-image Document Understanding, we divide baselines into three groups: (a) models without Large Language Models as decoders (Kim et al., 2022; Lee et al., 2023), (b) Multimodal LLMs (Hong et al., 2023; Dong et al., 2024a; Chen et al., 2024; Li et al., 2024b; Hu et al., 2024; Feng et al., 2023; Li et al., 2023b) with an average number of visual tokens over 1k for a single document image and (c) Multimodal LLMs (Ye et al., 2023a;b; Liu et al., 2024; Li et al., 2024b; Bai et al., 2023) with an average number of visual tokens less than 1k. As shown in Table 2, although specifically fine-tuned on each downstream dataset, Donut (Kim et al., 2022) or PixsStruct Lee et al. (2023) are not as good as Multimodal LLMs, showing the potential of MLLMs for generalized OCR-free document understanding. Compared with MLLMs with $<1k$ visual tokens, our DocOwl2 achieves better or comparable performance on 10 benchmarks. Especially, with fewer visual tokens, our model outperforms both TextMonkey (Liu et al., 2024) and TokenPacker (Li et al., 2024b) which also aim to compress visual tokens, showing that our layout-aware architecture High-resolution DocCompressor is better at summarizing and maintaining textual information in high-resolution document images. Besides, compared with state-of-the-art MLLMs with $>1k$ visual tokens, DocOwl2 achieves $>80\%$ performance on 7/10 benchmarks while with $<20\%$ visual tokens. Fig. 4 visualizes the comparison with SOTA in terms of question-answering performance and the number of visual tokens.

Furthermore, we compare the First Token Latency (seconds) on the 3 most frequently compared datasets, representing documents, charts, and natural images. As shown in Table 3, the far greater number of visual tokens enable InternVL 2 (Chen et al., 2024) and IXC 2.5 (Dong et al., 2024a) to achieve better performance but also result in higher inference time. Considering the model architecture and training data, it’s most fair to compare DocOwl2 with DocOwl 1.5. After adding the High-resolution DocCompressor, with similar training data of OCR learning, DocOwl2 achieves

Table 4: Comparison with OCR-free Multimodal Large Language Models on multi-image/video document understanding benchmarks. ‘FTL(s)’ refers to the First Token Latency (seconds). ‘Token^V’ means the average number of visual tokens of a single page/frame.

Model	Token ^V	MP-DocVQA		DUDE		NewsVideoQA	
		FTL(s)↓	ANLS↑	FTL(s)↓	ANLS↑	FTL(s)↓	ANLS↑
LongVA-7B	~2,029	2.13	60.80	2.26	38.37	4.29	50.61
Idefics3-8B	~838	2.26	67.15	2.29	38.65	6.39	60.16
LLaVA-next-interleave-7B	729	1.56	44.87	1.47	28.03	4.35	56.66
DocOwl2-8B	324	0.95	69.42	0.94	46.77	1.17	64.09

Table 5: Ablation study about the architecture of the compressor on single-image document benchmarks. ‘Img^{base}’ refers to the basic resolution of the global image and each sub-image.

	Img ^{base}	Crop	Compressor				Token ^V	DocVQA	WTQ	ChartQA
			Name	Compressing	Layer	Position				
r1	448	9	Resampler	learnable query	-	after H-Reducer	256	69.0	29.4	66.6
r2	448	9	CAbstractor	Adaptive Mean	-	after H-Reducer	256	73.0	32.6	67.6
r3	448	9	DocCompressor	Group Att	2	after H-Reducer	256	76.1	35.1	69.2
r4	448	9	DocCompressor	Group Att	2	after ViT	256	75.7	33.3	68.7
r5	448	9	DocCompressor	Complete Att	2	after H-Reducer	256	74.4	33.7	68.2
r6	448	9	DocCompressor	Group Mean	-	after H-Reducer	256	74.6	31.9	68.2
r7	448	9	DocCompressor	Group Att	1	after H-Reducer	256	76.4	34.2	69.2
r8	448	9	DocCompressor	Group Att	4	after H-Reducer	256	75.9	35.8	70.1
r9	448	12	DocCompressor	Group Att	2	after H-Reducer	256	76.8	35.6	69.5
r10	504	12	DocCompressor	Group Att	2	after H-Reducer	324	78.7	36.7	69.4

98% performance of DocOwl 1.5 while reducing 50% First Token Latency with just 20% visual tokens. This validates the effectiveness of our compressor for compressing visually-situated text information on the most common documents, charts, and natural images.

4.2.2 MULTI-PAGE/VIDEO DOCUMENT UNDERSTANDING

For Multi-page Document Understanding and Text-rich Video Understanding benchmarks, we choose recently proposed Multimodal LLMs (Zhang et al., 2024; Laurençon et al., 2024; Li et al., 2024a) with multi-page OCR-free document understanding abilities and can be fed into more than 10 images under a single A100-80G as baselines. As shown in Table 4, with fewer visual tokens for a single image/frame, our model DocOwl2 achieve better question-answering performance and much less First Token Latency, validating the good balance of DocOwl2 between the OCR-free document understanding performance and token efficiency.

4.3 ABLATION STUDY

We perform sufficient ablation studies to show the effectiveness of the architecture of High-resolution DocCompressor and the three-stage training strategy of DocOwl2.

4.3.1 COMPRESSOR ARCHITECTURE

To validate the effectiveness of our High-resolution DocCompressor, we compare different compressing architectures with an identical training pipeline of Single-image Pretraing and Single-image Document Understanding Finetuning, keeping both training data and training setting consistent.

As shown in Table 5, compared with CAbstractor (Cha et al., 2023), Resampler (Bai et al., 2023) achieves worse document understanding performance (r2 vs r1). This shows that due to no prior knowledge, such as spatial relationship, is leveraged as compressing guidance, utilizing queries learned from scratch to compress rich visually-situated text information is more challenging than simple adaptive mean pooling. Our High-resolution DocCompressor outperforms CAbstractor (r3 vs r2), validating that leveraging global visual features as layout-aware guidance can better distinguish the information density of each fine-grained visual feature and therefore maintain more visually-situated text information.

Table 6: Ablation study about the training stages of DocOwl2. ‘Single’ and ‘Multi’ refer to training samples utilizing single or multiple images as input. ‘Page Num’ and ‘Evidence Page’ refer to the number of input page images and the page ordinal number with the ground-truth answer.

	Pretraining		SFT		DocVQA	MP-DocVQA						Overall
	Single	Multi	Single	Multi		Page Num			Evidence Page			
						1	2-10	>10	1	2-10	>10	
r1	✓		✓		78.7	81.3	55.0	5.8	67.7	45.9	6.2	54.2
r2	✓			✓	75.2	78.7	65.2	34.6	74.3	54.9	40.9	63.8
r3	✓	✓		✓	74.2	78.9	65.7	37.9	74.2	56.8	43.4	64.7
r4	✓	✓	✓	✓	80.7	83.3	70.2	42.5	78.6	60.9	53.6	69.4

Instead of placing the compressor after the vision-to-text module H-Reducer, we also try inserting it between the vision encoder and the vision-to-text module. Such a setting results in performance decreases across three datasets (r4 vs r3), validating our hypothesis that compressing features after the vision-to-text module is like summarizing textual features and can maintain more textual semantics while compressing visual features after the visual encoder loses more visually situated text information. Besides, without aligning each query token in the global feature map with $R \times C$ fine-grained visual tokens from the re-organized feature map to perform attention within a group as Eq. (5), we try utilizing each query token to attend all visual tokens of sub-images. Such complete attention not only brings higher computational complexity but also causes performance decreases (r5 vs r3), showing that the positional correspondence between the global visual map and the re-organized fine-grained visual map is a reliable prior knowledge for compressing visual features efficiently. Furthermore, directly performing mean pooling on each group of $R \times C$ fine-grained visual features underperforms utilizing global visual features as the query to perform cross-attention (r6 vs r3). This also proves the importance of reliable guidance during compressing.

Compared with 2 layers of cross-attention, decreasing cross-attention layers bring a slight performance increase on DocVQA (Mathew et al., 2021) but more performance decrease on WikiTablesQA (WTQ) (Pasupat & Liang, 2015) (r7 vs r3). Further increasing to 4 layers doesn’t significantly improve performance (r8 vs r3). This shows that compressing high-resolution visual features doesn’t require a deep neural network. Finally, increasing the maximum number of crops and the base resolution of the global image or each sub-image are two main strategies to increase the supported input resolution. Our experiments show that increasing the cropping number (r9 vs r3) or basic resolution (r10 vs r9) benefits the document understanding performance. Increasing basic resolution brings more improvement because of more visual tokens after compressing.

4.3.2 TRAINING STRATEGY

DocOwl2 is trained with three stages: Single-image Pretraining, Multi-image Continue-pretraining, and Multi-task Finetuning. Table 6 shows the influence of each stage for OCR-free single-page and multi-page document understanding. With the Single-image Pretraining and Single-image finetuning (r1), the model achieves promising performance on single-page benchmark DocVQA and documents from MP-DocVQA with only 1 page. Although only trained with 1 image as the input, the model can also achieve around 50% accuracy when fed into 2-10 page images. However, the model struggles to understand documents with more than 10 pages, which greatly exceeds the number of input images during training and brings great difficulty in correlating images and finding answers. Performing Multi-image Finetuning could greatly improve the model’s ability to understand multiple images (r2 vs r1). Furthermore, adding the Multi-image Continue-pretraining could also improve the question-answering performance on downstream datasets, especially for documents with more than 10 pages (r3 vs r2). This demonstrates that parsing texts of the specified page or judging which pages contain specified texts among multi-page documents is a basic ability for multi-page document understanding. Finally, by ensembling both single-image and multi-image instruction tuning sets to perform the Multi-task Finetuning (r4), DocOwl2 achieves the best performance on both single-page and multi-page document benchmarks, showing the cross-improvement between single-image and multi-image comprehension.

(vii) the provisions of this Agreement governing payments with respect to Certificates and related notices, including Sections 2.4, 3.1 and 3.2, shall be revised to provide for distributions on each class of the Additional Certificates after payment of Administrative Expenses, the Liquidity Obligations, the Class AA Certificates, the Class A Certificates and the Class B Certificates (and, if applicable, any previously issued class of Additional Certificates), subject to clause (iii) above; and

Identify the image that features text <doc> (vii) the provisions of this Agreement governing payments with respect to Certificates and related notices, including Sections 2.4, 3.1 and 3.2

...
SECTION 10.1. Termination of Intercreditor Agreement. Following payment of Final Distributions with respect to each Class of Certificates and the payment in full of all Liquidity Obligations to the Liquidity Providers and provided that there shall then be no other amounts due to the Certificateholders, the Trustees, the Liquidity Providers and the 53 </doc> ?

ARTICLE X
MISCELLANEOUS

SECTION 10.1. Termination of Intercreditor Agreement. Following payment of Final Distributions with respect to each Class of Certificates and the payment in full of all Liquidity Obligations to the Liquidity Providers and provided that there shall then be no other amounts due to the Certificateholders, the Trustees, the Liquidity Providers and the

53

DocOwl 2:

(a)

Page 17 Georgia Garden Railway Society Sep 2020

Li-Ion Battery 3.7V 240 mAh

The sound modules I purchased operate with less than 5 volts. Some use three 'button cells' for a total of about 4.5 volts. Button cells do not last very long. Three AA or AAA cells also work but have to be replaced periodically. The modules also operate on a single Li-Ion rechargeable cell (voltage nominally 3.7 volts). These modules use very few milliamps so I used a 240 milliAmpHour (mAh) Li-Ion battery purchased on AliExpress. It will rarely require recharge. It is shown in Figure 2.

DocOwl 2:

<doc 14> Page 14. Georgia Garden Railway Society Sep 2020 Atlanta Senior Life: Big Fun with Little Trains

The Atlanta Senior Life newspaper carried an article in its July 2020 Vol. 5 No. 7 edition featured a couple of couples from the GGRS. ...

Later in the article, another GGRS pair, Russ and Leslie Ann Bundy were also interviewed. Maybe we can pick up a couple of new members from this coverage. The Atlanta Senior Life is available online at atlantaseniorlife.com or on facebook.com/atlantaseniorlife .

2020 Piedmont Pilgrimage -- An Online Tour of the Atlanta Area's Great Model Railroads

By Russ Bundy

The Piedmont Pilgrimage is sponsored each year by the Piedmont Division ...

the 18th annual pilgrimage, 2020 is proving to be quite a challenging year.

Social distancing to minimize chances of contracting the COVID - 19 virus has affected a lot of activities, including the Piedmont Pilgrimage.

Continued page 10 </doc 14>

<doc 17> Page 17 Georgia Garden Railway Society Sep 2020

The sound module is operated with less than 5 volts. Some use three 'button cells' for a total of about 4.5 volts. Buttons do not last very long. AAA cells also do not have to replace periodically.

The modules also operate on a single Li-Ion rechargeable cell (voltage nominal current modules). These modules use very few milliamps so I used a 240 AmpHour (mAh) Li-Ion battery purchased on AliExpress. It rarely require recharge. It is shown in Figure 2.

...

With some soldering I replaced the pushbutton switches provided with the sound Magnet Trigger modules with magnetic switches. The magnetic switch was hot-glued to existing structure on the bottom of the flat car as shown in Figure 4. </doc 17>

Page 14 Georgia Garden Railway Society Sep 2020

Atlanta Senior Life: Big Fun with Little Trains

The Atlanta Senior Life newspaper carried an article in its July 2020 Vol. 5 No. 7 edition featured a couple of couples from the GGRS. Front page news held the "Big Fun With Little Trains" title and a photo of James and Sally Beade in their indoor layout. The article was not G Scale only, but did a good job of representing the hobby as a whole, and Garden Railroad was not omitted.

Later in the article, another GGRS pair, Russ and Leslie Ann Bundy were also interviewed. Maybe we can pick up a couple of new members from this coverage. The Atlanta Senior Life is available online at atlantaseniorlife.com or on facebook.com/atlantaseniorlife .

2020 Piedmont Pilgrimage -- An Online Tour of the Atlanta Area's Great Model Railroads

By Russ Bundy

The Piedmont Pilgrimage is sponsored each year by the Piedmont Division of the NMR. Founded in 2005, and occurring in the October-November timeframe each year, this event usually features open house presentations of train layouts located throughout our area. Last year's event had more than 70 individual layouts which were open to the public for viewing. Several of us look forward to participating each year. The Bundos and Grindhouses have been hosting Open Houses for quite a few years. For Leslie and me, this would have been our third year participating in the event. Looking forward to the 18th annual pilgrimage, 2020 is proving to be quite a challenging year.

Social distancing to minimize chances of contracting the COVID - 19 virus has affected a lot of activities, including the Piedmont Pilgrimage.

Continued page 10

(b)

Figure 5: Qualitative results of the Multi-page Text Lookup (a) and Multi-page Text Parsing (b) given by DocOwl2 after the Multi-image Continue Pretraining.

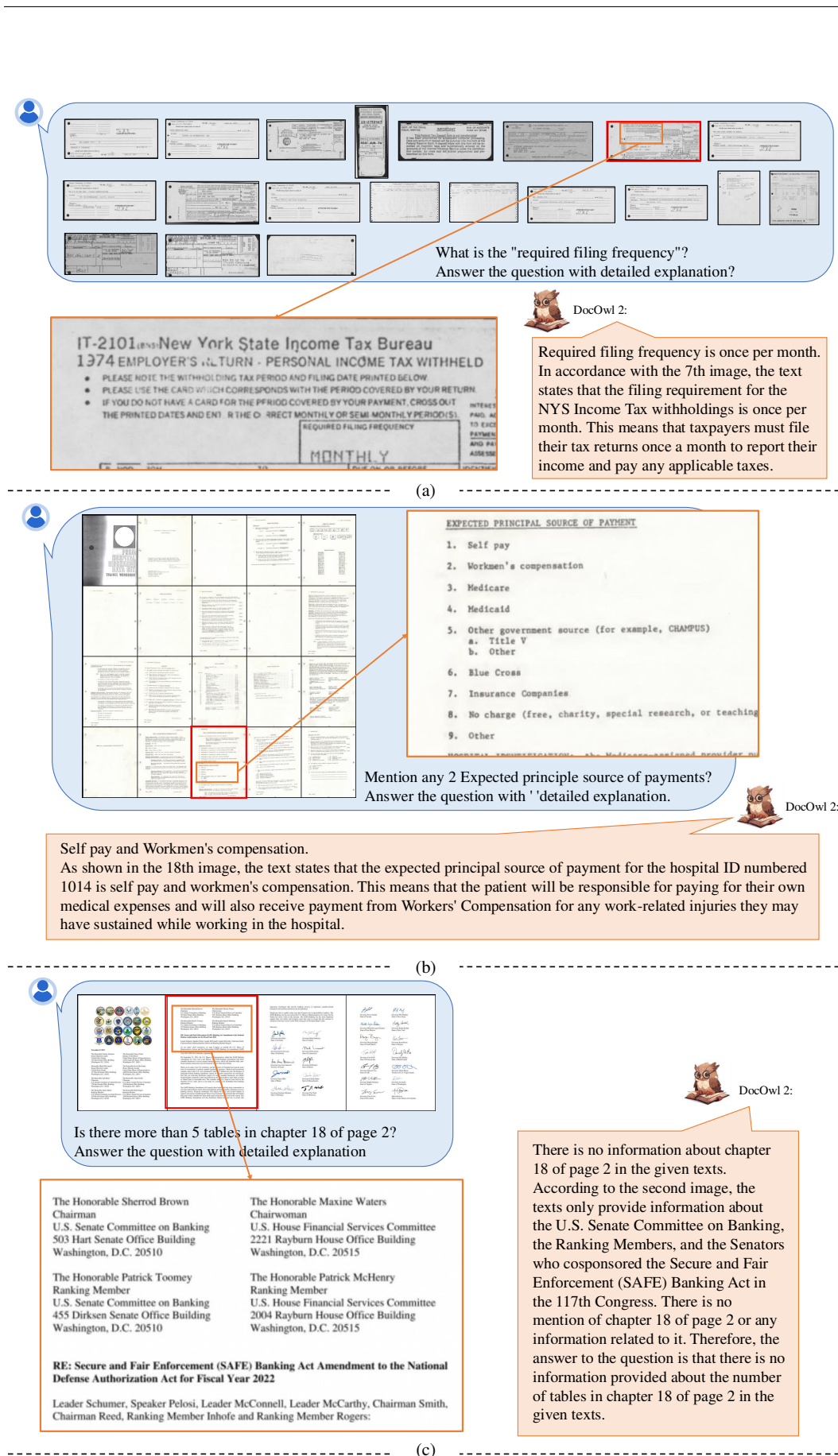


Figure 6: Qualitative results of the Multi-page Question Answering with detailed explanation.



Figure 7: Qualitative results of the Text-rich Video Understanding.

4.4 QUALITATIVE RESULTS

As shown in Fig. 5, after the Multi-image Continue Pretraining stage, DocOwl2 is able to locate the corresponding image of the given texts accurately. Besides, although representing each high-resolution image with just 324 tokens, DocOwl2 is still capable of parsing detailed texts of specified two images, validating the promising OCR-free multi-page document understanding performance of DocOwl2. It also demonstrates our proposal that 324 tokens are enough to encode detailed text information in common A4-sized document pages and the effectiveness of our High-resolution DocCompressor.

After the Multi-task Finetuning, given multiple images and a question, DocOwl2 can give a simple answer first and then provide a detailed explanation with the evidence, as shown in Fig. 6. DocOwl2 can comprehend not only page images rendered from PDF files (Fig. 6(c)) but also scan images of a document (Fig. 6(a-b)). When a question is unanswerable, DocOwl2 can also tell and give corresponding reasons (Fig. 6(c)).

Besides multi-page documents, DocOwl2 is also capable of understanding text-rich videos. As shown in Fig. 7, among similar frames within a video, DocOwl2 can distinguish fine-grained textual differences, locate relevant frames, and give accurate answers.

5 CONCLUSION

In this work, we propose mPLUG-DocOwl2, a Multimodal Large Language Model with the ability of efficient OCR-free Multi-page Document Understanding. The novel architecture High-resolution DocCompressor in DocOwl2 compresses each high-resolution document image into 324 tokens through cross-attention with the global visual feature as guidance, and re-organized features of cropped images as keys and values. On single-image document understanding benchmarks, with fewer visual tokens, DocOwl2 outperforms existing compressing methods and achieves comparable performance with SOTA MLLMs with similar training data. Besides, DocOwl2 achieves OCR-free state-of-the-art performance on two multi-page document understanding benchmarks and 1 text-rich video understanding benchmark. Our experiments validate that thousands of visual tokens for 1 common A4-sized document page may be so redundant that too many computational resources are wasted. We hope DocOwl2 could bring more researchers’ attention to the balance of efficient representation of high-resolution images and OCR-free Document Understanding performance.

REFERENCES

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. *ArXiv*, abs/2204.14198, 2022.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. Honeybee: Locality-enhanced projector for multimodal LLM. *CoRR*, abs/2312.06742, 2023.
- Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyu Zhou, and William Yang Wang. Tabfact : A large-scale dataset for table-based fact verification. In *International Conference on Learning Representations (ICLR)*, Addis Ababa, Ethiopia, April 2020.
- Xingyu Chen, Zihan Zhao, Lu Chen, JiaBao Ji, Danyang Zhang, Ao Luo, Yuxuan Xiong, and Kai Yu. Websrc: A dataset for web-based structural reading comprehension. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 4173–4185, 2021.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, Ji Ma, Jiaqi Wang, Xiaoyi Dong, Hang Yan, Hwei Guo,

- Conghui He, Botian Shi, Zhenjiang Jin, Chao Xu, Bin Wang, Xingjian Wei, Wei Li, Wenjian Zhang, Bo Zhang, Pinlong Cai, Licheng Wen, Xiangchao Yan, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *CoRR*, abs/2404.16821, 2024.
- Xin Cheng, Xun Wang, Xingxing Zhang, Tao Ge, Si-Qing Chen, Furu Wei, Huishuai Zhang, and Dongyan Zhao. xrag: Extreme context compression for retrieval-augmented generation with one token. *CoRR*, abs/2405.13792, 2024.
- Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. Adapting language models to compress contexts. In *EMNLP*, pp. 3829–3846. Association for Computational Linguistics, 2023.
- Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Zhe Chen, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Kai Chen, Conghui He, Xingcheng Zhang, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k HD. *CoRR*, abs/2404.06512, 2024a.
- Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Zhe Chen, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Kai Chen, Conghui He, Xingcheng Zhang, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k HD. *CoRR*, abs/2404.06512, 2024b.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*. OpenReview.net, 2021.
- Hao Feng, Qi Liu, Hao Liu, Wengang Zhou, Houqiang Li, and Can Huang. Docpedia: Unleashing the power of large multimodal model in the frequency domain for versatile document understanding. *CoRR*, abs/2311.11810, 2023.
- Tao Ge, Jing Hu, Lei Wang, Xun Wang, Si-Qing Chen, and Furu Wei. In-context autoencoder for context compression in a large language model. In *ICLR*. OpenReview.net, 2024.
- Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for GUI agents. *CoRR*, abs/2312.08914, 2023.
- Anwen Hu, Shizhe Chen, and Qin Jin. Question-controlled text-aware image captioning. In *ACM Multimedia*, pp. 3097–3105. ACM, 2021.
- Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding. *CoRR*, abs/2403.12895, 2024.
- Soumya Jahagirdar, Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Watching the news: Towards videoqa models that can read. In *WACV*, pp. 4430–4439. IEEE, 2023.
- Kushal Kaffe, Brian L. Price, Scott Cohen, and Christopher Kanan. DVQA: understanding data visualizations via question answering. In *CVPR*, pp. 5648–5656. Computer Vision Foundation / IEEE Computer Society, 2018.
- Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. In *ICLR (Workshop)*. OpenReview.net, 2018.
- Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. Ocr-free document understanding transformer. In *ECCV (28)*, volume 13688 of *Lecture Notes in Computer Science*, pp. 498–517. Springer, 2022.

-
- Jordy Van Landeghem, Rafal Powalski, Rubèn Tito, Dawid Jurkiewicz, Matthew B. Blaschko, Lukasz Borchmann, Mickaël Coustaty, Sien Moens, Michal Pietruszka, Bertrand Anckaert, Tomasz Stanislawek, Pawel Józiak, and Ernest Valveny. Document understanding dataset and evaluation (DUDE). In *ICCV*, pp. 19471–19483. IEEE, 2023.
- Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *CoRR*, abs/2405.02246, 2024.
- Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *ICML*, volume 202 of *Proceedings of Machine Learning Research*, pp. 18893–18912. PMLR, 2023.
- Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *CoRR*, abs/2407.07895, 2024a.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. *CoRR*, abs/2301.12597, 2023a.
- Wentong Li, Yuqian Yuan, Jian Liu, Dongqi Tang, Song Wang, Jianke Zhu, and Lei Zhang. Token-packer: Efficient visual projector for multimodal LLM. *CoRR*, abs/2407.02392, 2024b.
- Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *CoRR*, abs/2403.18814, 2024c.
- Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. *CoRR*, abs/2311.06607, 2023b.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *CoRR*, abs/2304.08485, 2023.
- Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai. Textmonkey: An ocr-free large multimodal model for understanding document. *CoRR*, abs/2403.04473, 2024.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq R. Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *ACL (Findings)*, pp. 2263–2279. Association for Computational Linguistics, 2022.
- Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for VQA on document images. In *WACV*, pp. 2199–2208. IEEE, 2021.
- Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and C. V. Jawahar. Infographicvqa. In *WACV*, pp. 2582–2591. IEEE, 2022.
- Nitesh Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. Plotqa: Reasoning over scientific plots. In *WACV*, pp. 1516–1525. IEEE, 2020.
- Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *ACL (1)*, pp. 1470–1480. The Association for Computer Linguistics, 2015.
- Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: A dataset for image captioning with reading comprehension. In *ECCV (2)*, volume 12347 of *Lecture Notes in Computer Science*, pp. 742–758. Springer, 2020.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards VQA models that can read. In *CVPR*, pp. 8317–8326. Computer Vision Foundation / IEEE, 2019.

-
- Tomasz Stanislawek, Filip Gralinski, Anna Wróblewska, Dawid Lipinski, Agnieszka Kaliska, Paulina Rosalska, Bartosz Topolski, and Przemyslaw Biecek. Kleister: Key information extraction datasets involving long documents with complex layouts. In *ICDAR (1)*, volume 12821 of *Lecture Notes in Computer Science*, pp. 564–579. Springer, 2021.
- S Svetlichnaya. Deepform: Understand structured documents at scale, 2020.
- Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. Visualmrc: Machine reading comprehension on document images. In *AAAI*, pp. 13878–13888. AAAI Press, 2021.
- Rubèn Tito, Dimosthenis Karatzas, and Ernest Valveny. Hierarchical multimodal transformers for multi-page docvqa. *CoRR*, abs/2212.05935, 2022.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Vary: Scaling up the vision vocabulary for large vision-language models. *CoRR*, abs/2312.06109, 2023.
- Renqiu Xia, Song Mao, Xiangchao Yan, Hongbin Zhou, Bo Zhang, Haoyang Peng, Jiahao Pi, Daocheng Fu, Wenjie Wu, Hancheng Ye, Shiyang Feng, Bin Wang, Chao Xu, Conghui He, Pinlong Cai, Min Dou, Botian Shi, Sheng Zhou, Yongwei Wang, Bin Wang, Junchi Yan, Fei Wu, and Yu Qiao. Docgenome: An open large-scale scientific document benchmark for training and testing multi-modal large language models. *CoRR*, abs/2406.11633, 2024.
- Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu, Chenliang Li, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. mplug-docowl: Modularized multimodal large language model for document understanding. *CoRR*, abs/2307.02499, 2023a.
- Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian, Qi Qian, Ji Zhang, Qin Jin, Liang He, Xin Lin, and Fei Huang. Ureader: Universal ocr-free visually-situated language understanding with multimodal large language model. In *EMNLP (Findings)*, pp. 2841–2858. Association for Computational Linguistics, 2023b.
- Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models, 2024. URL <https://arxiv.org/abs/2408.04840>.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. mplug-owl: Modularization empowers large language models with multimodality. *CoRR*, abs/2304.14178, 2023c.
- Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *CoRR*, abs/2311.04257, 2023d.
- Liang Zhang, Anwen Hu, Jing Zhang, Shuo Hu, and Qin Jin. MPMQA: multimodal question answering on product manuals. *CoRR*, abs/2304.09660, 2023.
- Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. Long context transfer from language to vision. *CoRR*, abs/2406.16852, 2024.
- Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno-Yepes. Image-based table recognition: Data, model, and evaluation. In *ECCV (21)*, volume 12366 of *Lecture Notes in Computer Science*, pp. 564–580. Springer, 2020.