Adaptive Markup Language Generation for Contextually-Grounded Visual Document Understanding

Han Xiao^{1,2†}, Yina Xie², Guanxin Tan², Yinghao Chen², Rui Hu², Ke Wang¹, Aojun Zhou¹, Hao Li¹, Hao Shao¹, Xudong Lu^{1,2†}, Peng Gao⁴, Yafei Wen², Xiaoxin Chen², Shuai Ren^{2‡⊠}, Hongsheng Li^{1,3⊠}
¹CUHK MMLab ²vivo AI Lab ³CPII under InnoHK
⁴Shanghai AI Lab & Shenzhen Institute of Advanced Technology, CAS {1155229123@link,hsli@ee}.cuhk.edu.hk shuai.ren@vivo.com

1. Details of DocMark-Instruct Dataset

We provide a detailed overview of the DocMark-Instruct dataset, which is designed to enhance the document understanding and context grounding capabilities of MLLMs in document-related scenarios. The dataset contains data from six distinct domains, including Text-based QA, Documentbased QA, Chart-based QA, Key Information Extraction, Webpage-based QA and Mathematical QA. To enrich the dataset, we make use of multiple public datasets. Additionally, we also include in-house data created through careful curation. The domain distribution and statistics of our dataset are shown in Sec. 2.2.

2. Prompt Design

2.1. Prompt for Different Pretraining Tasks on DocMark-Pile

For multi-task pretraining on DocMark-Pile, we utilize a variety of instruction prompts specifically designed for each markup language generation task. These prompts not only guide the model during pretraining but also enhance its versatility in tackling diverse markup translation challenges. Some examples of the prompts utilized during pretraining are listed in Tab. 2.

2.2. Prompt for Creating the DocMark-Instruct Dataset

As outlined in Section 3.2.2 of the main paper, we employ ChatGPT-3.5 to generate immediate context for our DocMark-Instruct dataset. Specifically, we prompt the model to extract relevant information from the provided markup language necessary for answering the questions.

Task	Source Dataset	Markup Type	#Num.
	TextVQA [16]	txt, txt_gd	29k
Text-based QA	STVQA [1]	txt, txt_gd	26k
	EST-VQA [18]	txt, txt_gd	8k
	DocVQA [13]	md	31k
Document-based QA	InfoVQA [12]	md	12k
	Docmatix [8]	md	50k
	In-house data	md, latex	12k
	ChartQA [11]	json	44k
Chart-based QA	PlotQA [14]	json	146k
	DVQA [6]	json	77k
	POIE [7]	json	2k
Kay Information Extraction	SROIE [4]	json	0.6k
Key information Extraction	FUNSD [5]	json	0.1k
	XFUND [19]	json	0.2k
	WebSRC [2]	html	59k
Webpage-based QA	In-house data	html	4k
	Geo170k [3]	tikz	48k
Mathematical QA	Geometry3k [10]	tikz	2k
	MultiMath [15]	tikz	50k
	In-house data	tikz	17k
Total	-	-	624k

Table 1. Overview of the DocMark-Instruct Dataset.

It is important to note that some questions may not relate directly to the textual information; for instance, certain questions may only require general knowledge for answers. In such cases, we instruct the model to respond with "unclear", indicating that the textual information is not applicable. This approach helps prevent inaccurate annotations and mitigates the risk of model hallucinations. The detailed prompt template is provided in Fig. 2.

[™]Corresponding author [‡]Project lead [†]Interns at vivo.

3. Performance Analysis

3.1. Qualitative Evaluation

For the qualitative evaluation, we offer more visualizations of the generated results by our DocMark, particularly on generating PDF documents, webpages, and scientific diagrams. As shown in Fig. 3 and Fig. 4, our model is capable of maintaining the textual and layout information effectively. It should be noted that our method might not preserve the style information, such as font sizes and colors, very well. This is because our pre-training mainly concentrates on parsing the structured information within the documents. Learning the main textual and layout representations is sufficient for our document understanding tasks.

3.2. Evaluating Accuracy with and without CoT Prediction

We present a series of ablation experiments comparing: (1) vanilla models with direct CoT prompting, (2) models trained on DocMark-Pile and DocMark-Instruct with CoT fine-tuning removed, and (3) our full model. As shown in Tab. 3, the vanilla models including Qwen2-VL and LLaVA-OneVision, despite being trained on extensive datasets, exhibit significant performance degradation compared to their unmodified counterparts, indicating a limited capacity for explicit reasoning. The primary issue lies in these models' inability to effectively comprehend the document layout and derive accurate context from it. Notably, using CoT prompting with DocMark-Pile pre-trained models, even without CoT fine-tuning, surpasses the baseline performance. Moreover, the vanilla models struggle to associate context with the question due to their inherent training methods that focuses on direct prediction. In contrast, our full approach, which incorporates end-to-end CoT finetuning, leads to superior performance. This suggests that the enhancements in our original approach are primarily attributed to the effectiveness of the CoT component rather than solely introducing more data.

3.3. Demonstrations of Adaptive Context Generation

To better showcase the adaptive generation capability of DocMark, we present additional demonstrations regarding downstream document understanding tasks in Fig. 5. This emphasizes our model's capacity to identify crucial information within documents and offer contextually-grounded answers, thereby enhancing the model's interpretability. By initially recognizing the document format, locating the area of interest, and extracting relevant structured representations, our model enhances both the accuracy of text recognition and the depth of understanding, which directly impacts the quality of the generated answers.



Figure 1. Token count comparison across various datasets. We visualize the distribution of image tokens, context tokens, and original question-answer tokens on different datasets.

3.4. Token Efficiency

Since our adaptive generation pipeline incorporates additional auxiliary context tokens to derive final answers, we need to evaluate its impact on computational efficiency. To investigate this, we compare the number of image tokens, context tokens, and original question-answer tokens across several representative datasets. As demonstrated in Fig. 1, image tokens account for the majority of the total token count due to the adopted dynamic resolution strategy. In contrast, context tokens contribute only a small number of tokens compared to both image tokens and original conversational tokens. This indicates that our method effectively provides highly condensed contextual information, alleviating the limitations of relying solely on image tokens and improving the overall understanding of the document with minimal computational cost.

References

- [1] Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluis Gomez, Marçal Rusinol, Ernest Valveny, CV Jawahar, and Dimosthenis Karatzas. Scene text visual question answering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4291–4301, 2019. 1
- [2] Xingyu Chen, Zihan Zhao, Lu Chen, Danyang Zhang, Jiabao Ji, Ao Luo, Yuxuan Xiong, and Kai Yu. Websrc: A dataset for web-based structural reading comprehension. *arXiv preprint arXiv:2101.09465*, 2021.
- [3] Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, et al. G-llava: Solving geometric problem with multi-modal large language model. *arXiv preprint arXiv:2312.11370*, 2023. 1
- [4] Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and CV Jawahar. Icdar 2019 robust

Task	Prompt
Text Recognition	Kindly recognize the text from the image. How can I extract the text from the image? What text is in the image that can be extracted?
Text Grounding	Can you perform text extraction with grounding? Please detect the text with grounding from the image. Recognize the text with grounding from the image.
Image to Markdown	Parse the image into a proper markup language format. How to convert text from the image to markdown format? How to extract text from the image and change it to markdown format?
Image to LaTeX	Convert the image into a structured format. How to extract and translate text from the image to LaTeX format? How can I convert text from the image to LaTeX format efficiently?
Image to HTML	What is the HTML code corresponding to this image? Generate the HTML code. Parse the image into an appropriate markup language format.
Webpage Summarization	What is the main idea of this webpage screenshot? What are the main information points of the webpage shown in the image? What is the key message conveyed by this webpage image?
Image to JSON	Extract text from the image in JSON format. Output the image text as JSON. Represent the image text in a structured format.
Image to TikZ	I need to get the code for drawing this image. What is the TikZ code for this image?. Please show me the TikZ code for displaying this image.

Table 2. Examples of prompts for different pretraining tasks on DocMark-Pile.

Base Model	Training	Inference	DocVQA	ChartQA	MathV
	vonillo	vanilla	89.2	73.5	20.1
Qwen2-VL-2B [17]	vanilla	СоТ	84.9	55.5	13.8
	DeeMert (u./o CoT)	vanilla	87.7	69.3	18.6
	DOCIVIAIK(W/O COT)	CoT	88.2	72.8	19.2
	DocMark(full)	CoT	89.8	77.1	22.4
	vonillo	vanilla	88.7	80.8	21.7
LLaVA-OneVision-7B [9]	vaiiiiia	CoT	87.2	78.4	21.4
	DeeMark(/a CaT)	vanilla	87.0	79.1	21.5
	Dociviark(w/o Col)	CoT	88.2	80.4	22.0
	DocMark(full)	CoT	90.1	83.4	23.2

Table 5. Comparison of model performance with and without Cor me-tu

reading challenge on scanned receipts ocr and information extraction. In *International conference on document analysis recognition*, 2019. 1

[5] Guillaume Jaume, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Funsd: A dataset for form understanding in noisy scanned documents. In 2019 International Conference on Document Analysis and Recognition Workshops (ICDARW), pages 1-6. IEEE, 2019. 1

- [6] Kushal Kafle, Scott Cohen, Brian Price, and Christopher Kanan. Dvqa: Understanding data visualizations via question answering. In CVPR, 2018. 1
- [7] Jianfeng Kuang, Wei Hua, Dingkang Liang, Mingkun Yang, Deqiang Jiang, Bo Ren, and Xiang Bai. Visual information extraction in the wild: practical dataset and end-to-end solu-

You are provided with a question related to a document image and its parsed results in various markup formats. Your task is to extract relevant content to answer the question.

Instructions

1. Input Components:

- Question: A specific inquiry related to the document.
- Answer: The known answer to the question.
- Parsed Results: A collection of text extracted from the document in various markup formats.

2. Markup Formats to Recognize:

- `<txt></txt>`: Plain text from the document.
- `<txt_gd></txt_gd>`: Text with coordinates for context.
- `<md></md>`: Markdown representation of the content.
- `<latex></latex>`: LaTeX code for generating the document.
- `<html></html>`: HTML code for webpage.

3. Extraction Guidelines:

- Extract and compile all relevant content that may assist in answering the question.
- Ensure that the extracted content retains its original formatting.
- Do not alter or modify any of the content during extraction.

4. Output Requirements:

- Provide all relevant extracted text.

- If no pertinent information is found in the parsed results, or if the question can be answered without them, respond with 'unclear'.

Input:

- Question: <Question>
- Answer: <Answer>
- Parsed Results: [<type>(parsed results)</type>]

Your Output:

- [Extracted content or 'unclear']

Figure 2. Prompt for creating the DocMark-Instruct dataset.

tion. In International Conference on Document Analysis and Recognition, pages 36–53. Springer, 2023. 1

- [8] Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Docmatix dataset. https://huggingface. co/blog/docmatix, 2024. 1
- [9] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li,

Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 3

[10] Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. Inter-gps: Interpretable geometry problem solving with formal language and sym-

Original Document

		PUBLICATION DATEGOURNAL	STATUS/ ACTION
8(1.1) FAST TRACK lood lipid profile of structopassal women receiving remote replacement therapies intaining estradiel in combination if either noreflicitense accutte or infective receive a 1-ware period	Metabolie impact berostasis	1Q2002 Opercological Endocrinology	Sign off received from Wyth Dr Al-Azawi agreed to act as samed author. Ho suggested the addition of Dr J. Stowness exists on bad box as or worker on the study. Dr Stownesse has revised the mansacript and it was retunded to Dr Al-Azawi fragmental. He approvale this comments and only fraid pary has to be christed with 3 Stowneon before administro W Work for final using off
(2) comparison of two hormone placement therapies containing tandio in combination with either drogesterene or triangestone with spect to provention of stransportated been loss	Beno	202042	Commune received from Gary Grubb 18:03:02. Revisions being carried out
(3) FAST TRACK 1-yar comparison of the efficacy a christal thrance in stranopassal women of two remose splacenast thrappios mining saturabel in combination (th either roopstrat) or imogeneous lawsisson FLULM, et al.	Efficacy	Published	Generological Enderrinology 2001;15:342–58
(4)a fleets of an estradiol trimegratore enbiantion in comparison with amoston on cardiovascular risk stress (factor VII, fibrinopan, PM- 142), spinose, Tpido)	Metabolic impact hemostasis	? Anu J Obstat Gjoracod	Paper to be prepared by Hellgren and Khali
(4) Easts of an estradial brimogeneous mbination in comparison with mession on other cougatation memory (inhibitory factors, arkers of balance between stration of cougatation and brindysis).	Metabolic impact homostasis	? Deverbosis Havroodasis	Paper to be prepared by Norris
(5) inple center, randomized, double- industudy of the effects of the totaled (2mg/ imagestone (0.5mg) combination in totaled mayimuchroxyprogesterone accidate dPAQ (10mg) combination on none and insulin metabolism ring of nonthin of treatment	Metabolic impact hemostasis	2(2002	Partnere roqued supply of data on that poparation of the deal manuscript can begin
N(6) FAST TRACK second of the metabolic learner in postmenopassal women was 1-your period of two hormone placement therapics containing tradied in combination with other reported or transportence.	Metabolic impact homostasis	2Q2002 Gymeiological Endocritology	Accepted for publication in April Issue of Gynecological Endocranology
(7) FAST TRACK (cta-analysis of phase III studies on coding pattern data	Disading profile	3Q2002	Parthenon request supply of data so that preparation of the draft manuscript can begin
(8) Incose and insulin metabolism Insiste G.	Matabalism	30,2002	Title and paper in preparation.



Generated by DocMark

PROJECT	CLASSIFICATION	PUBLICATION DATE/JOURNAL	STATUS/ACTION
P3(1.1) FAST TRACK	Metabolic impact/hemostasis	1Q2002	Sign off received from Wyeth. Dr. Azzawi agreed to act as named author. He suggested the addition Dr. J. Stevenson since he had been co-weaker on the study. Dr. Steven has revised the manuscript and it entraned to Dr. Arkansis for approval. He has provided his comments and only final query has be clarified with L. Stevenson before submission to Wyeth for final sign off
P3(2)	Bone	2Q2002	Comments received from Gary Gr 18.03.02. Revisions being carried out.
P3(3) FAST TRACK	Efficacy	Published	Gynecological Endocrinology 200 15:349:58
A 1-year comparison of the efficacy and clinical tolerance in postmenopaual women of two hormone replacement therapies containing estradiol in combination with either norethisterone accelle or trimegestone	Metabolic impact/hemostasis	?	Gynecological Endocrinology 200 15:349-58
P3(4)a Effects of an estradiol/trimegestone combination in comparison with Femoston on cardiovascular risk factors (factor VII, fibrinogen, PAI-1, t- tPA, glucose, Upida)	Metabolis impact/hemostasis	? Am J Obstet Gynecol	Paper to be prepared by Heligren Khaft
Pij(4)b Effects of an estradiol/trimegestone combination in comparison with Fernoston on other coagulation parameters (inhibitory factors, markers of balance between activation of coagulation and fibrinolysis)	Metabolic impact/hernootasia	? Thrombosis Haemostasis	Paper to be prepared by Norris
P3(5)	Single center, randomized, double-blind study of the effects of the estraficial (cmm)/trimepotone (o, rmg) combination in comparison with a placebo and an estradioi (cmm)/noelrecoprogesterone actate (MPA) (formg) combination on glucose and insulin metabolism during 6 months of treatment	Metabolic impact/hemostasis	2Q2002
P3(6) FAST TRACK	Metabolic impact/hemostasis	2Q2002	Accepted for publication in April issue of Gynecological Endocrinol
Assessment of the metabolic tolerance in postmenopausal women over a 1-year period of two hormone replacement therapies containing estradiol in combination with either norethistercese or trimegestone	Metabolic impact/hemootasis		
P3(7) FAST TRACK	Bleeding profile	3Q2002	Parthenon request supply of data that preparation of the draft manuscript can begin
Mcta-analysis of phase III studies on bleeding pattern data	Metabolism	3Q2002	Title and paper in preparation
P3(8)	Glucose and insulin metabolism		
Samisioe G.			
Confidential Pursuant to Confiden	tiality Order		

HAGER LABORATORIES	3, INCORPORATED ANALYTICAL SERVICES FOR INDUSTRY
REPORT ON SERVICE	NUMBER 9802 October 25, 1983
D. H. MORMAN OCT 2	28 1983
To: D. H. Morman S	Shell Development Houston, TX
Analysis: The fol: and one blank for	lowing samples were submitted for analysis: Twelve 3M POVM samples benzene.
Method: PASSIVE OF benzene.	XGANIC MONITORS (POVM) The collection element was analyzed for
The organic vapors solvent and the so detection. Millig chromatographic pe from standard sols A sample rate of 3 used for concentra	is more desorbed from the collection element of each monitor with a livent analysed using gas chromatography with flame ionization mavakes were corrected for desorption efficiency. The eak for each analyte was compared to a calibration curve dotained clines. B. & Corfan for benzeme (3M Compound Guide RSG46(GII)R May 1001) was tion calculations.
Results: Data is 1	cabulated in Table 1, 2 and 3 of this report.
Discussion: The OS	SHA permissible exposure limit (PEL) for benzene is 10 ppm.
ND() indicates no the parentheses.	one of the substance was detected, with a detection limit shown in
Laboratory data an	re filed and available upon request.
Submitted by: [sig	gnature] Robert N. Hager Jr., Ph.D. Laboratory Director
RNH/nl	
4725 Paris Street	/Denver, Colorado 80239/(303) 371-1441

Figure 3. Visualization on generated documents by DocMark.

Original Document

Generated by DocMark

E SKILSTAK

Linux Master Race

The Linux master race is a play on <u>PC master race</u> popularized on Reddit. It really is objectively true that <u>Linux</u> is far superior to Mac and Windows no matter how much holy outrage users of those inferior operating systems may spew. Some <u>muggles</u> don't even know what it is. (Just don't rile them up.) Linux is the most prolific and powerful operating system on the planet by every objective measure. Learning Linux empowers you more than anything else.

Muggles tend to fear Linux, which always makes me smile. They tend to be the same people who would never eat at a food truck. One member was actually asked on a phone customer support call what operating system he had. "Linux," he said. To which the thickly accented voice pedantically responded, "I'm sorry sir, Linux is not an operating system."

© 2013-2019 Skillski, Inc. This work is licensed under a Creative Commons Attributions: ShareAlike 40 International License, Code is licensed under the GNU Public License, Version 3.0. Optimized for oral and your with \mathbb{O}° for the terminal maters. Last Modified: Thursday, April 25, 2019 - 7.1751 PM (258d 11h 38m ags) Enter your suggestion, question, or feedback for this page here. Bend Your Message Anonymously

Linux Master Race

≡ ■ SKILSTAK

The Linux master race is a play on <u>PC master race</u> popularized on Reddit. It really is objectively true that <u>Linux</u> is far superior to Mac and Windows no matter how much holy outrage users of those inferior operating systems may spew. Some <u>muggles</u> don't even know what it is. (Just on this lieft mup J) linux is the most prolific and powerful operating system on the planet by every objective measure. Learning Linux empowers you more than anything else.

Muggles tend to fear Linux, which always makes me smile. They tend to be the same people who would never eat at a food truck. One member was actually asked on a phone customer support call what operating system he had, 'ulnux', he said. To which the thickly accented voice pedantically responded, \t'I'm sorry sir, Linux is not an operating system.\' Pfffhahahaa.

© 2013-2019 SkilStak, Inc. This work is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Code is licensed under the <u>GNU Public License. Version 3.0</u>.

Optimized for <u>curl</u> and <u>lynx</u> with ♥ for <u>true terminal masters</u>. Last Modified: Thursday, April 25, 2019 - 7:17:51 PM (258d 11h 38m ago)

VEner



Figure 4. Visualization on generated documents by DocMark.



Figure 5. Visualization on generated documents by DocMark.

bolic reasoning. In *The Joint Conference of the 59th An*nual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021), 2021. 1

- [11] Ahmed Masry, Do Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL* 2022, pages 2263–2279, Dublin, Ireland, 2022. Association for Computational Linguistics. 1
- [12] Minesh Mathew, Viraj Bagal, Rubèn Pérez Tito, Dimosthenis Karatzas, Ernest Valveny, and C. V Jawahar. Infographicvqa, 2021. 1
- [13] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In Proceedings of the IEEE/CVF winter conference on applications of computer vision, pages 2200–2209, 2021. 1
- [14] Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar. Plotqa: Reasoning over scientific plots. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 1527–1536, 2020. 1
- [15] Shuai Peng, Di Fu, Liangcai Gao, Xiuqin Zhong, Hongguang Fu, and Zhi Tang. Multimath: Bridging visual and mathematical reasoning for large language models. *arXiv preprint arXiv:2409.00147*, 2024. 1
- [16] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings* of the IEEE/CVF conference on computer vision and pattern recognition, pages 8317–8326, 2019. 1
- [17] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 3
- [18] Xinyu Wang, Yuliang Liu, Chunhua Shen, Chun Chet Ng, Canjie Luo, Lianwen Jin, Chee Seng Chan, Anton van den Hengel, and Liangwei Wang. On the general value of evidence, and bilingual scene-text visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10126–10135, 2020.
- [19] Yiheng Xu, Tengchao Lv, Lei Cui, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, and Furu Wei. Xfund: a benchmark dataset for multilingual visually rich form understanding. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3214–3224, 2022. 1