

HUMANS are naturally capable of imaging a scene according to a piece of visual, text or audio description. However, the intuitive processes are less straightforward for deep neural networks, primarily due to an inherent modality gap. This *modality gap* for visual perception can be boiled down to *intra-modal gap* between visual clues and real images, and *cross-modal gap* between non-visual clues and real images. Targeting to mimic human imagination and creativity in the real world, the tasks of Multimodal Image Synthesis and Editing (MISE) provide profound insights about how deep neural networks correlate multimodal information with image attributes.

[HUMANS are naturally capable of imaging a scene according to a piece of visual, text or audio description. However, the intuitive processes are less straightforward for deep neural networks, primarily due to an inherent modality gap. This modality gap for visual perception can be boiled down to intra-modal gap between visual clues and real images, and cross-modal gap between non-visual clues and real images. Targeting to mimic human imagination and creativity in the real world, the tasks of Multimodal Image Synthesis and Editing (MISE) provide profound insights about how deep neural networks correlate multimodal information with image attributes.]

中华医学会 中华医学杂志社 中华医学全科医学分会 中华医学会(中华全科医师杂志)编辑委员会 心血管系统疾病基层诊疗指南编写专家组
通信作者:孙艺红,中日友好医院心脏科,北京100029,Email:yihongsun72@163.com;
胡大一,北京大学人民医院心血管研究所 100044,Email:dayi.hu@china-heart.org
[关键词] 指南; 胸痛
DOI:10.2000/jm.j.cn.1671-7368.2019.10.004
Chinese Medical Association, Chinese Medical Journals Publishing House, Chinese Society of General Practice, Editorial Board of Chinese Journal of General Practitioners of Chinese Medical Association, Expert Group of Guidelines for Primary Care of Cardiovascular Disease
Corresponding author: Sun Yihong, Department of Cardiology, China-Japan Friendship Hospital, Beijing 100029, China, Email:yihongsun72@163.com; Hu Dayi, Institute of Cardiovascular Disease, Peking University People's Hospital, Beijing 100044, China, Email:dayi.hu@china-heart.org

中华医学会 中华医学杂志社 中华医学全科医学分会 中华医学会(中华全科医师杂志)编辑委员会 心血管系统疾病基层诊疗指南编写专家组
通信作者:孙艺红,中日友好医院心脏科,北京100029,Email:yihongsun72@163.com;
胡大一,北京大学人民医院心血管研究所 100044,Email:dayi.hu@china-heart.org
[关键词] 指南; 胸痛
DOI:10.2000/jm.j.cn.1671-7368.2019.10.004
Guidelines for primary care of chest pain(2019)
Chinese Medical Association, Chinese Medical Journals Publishing House, Chinese Society of General Practice, Editorial Board of Chinese Journal of General Practitioners of Chinese Medical Association, Expert Group of Guidelines for Primary Care of Cardiovascular Disease
Corresponding author: Sun Yihong, Department of Cardiology, China-Japan Friendship Hospital, Beijing 100029, China, Email:yihongsun72@163.com; Hu Dayi, Institute of Cardiovascular Disease, Peking University People's Hospital, Beijing 100044, China, Email:dayi.hu@china-heart.org

		100-class (top-1 acc.)	1000-class (top-1 acc.)
4096-d (float)	BP	77.1 ± 1.5	65.0
1024 bits	CBE	72.9 ± 1.3	58.1
	SP	73.0 ± 1.3	59.2
	threshold [1]	73.8 ± 1.3	60.1
4096 bits	BP	73.5 ± 1.4	59.1
	CBE	76.0 ± 1.5	63.2
	SP	75.9 ± 1.4	63.0
8192 bits	BP	76.3 ± 1.5	63.3
	SP	76.8 ± 1.4	64.2
16384 bits	SP	77.1 ± 1.6	64.5

		100-class (top-1 acc.)	1000-class (top-1 acc.)
4096-d (float)	BP	77.1 ± 1.5	65.0
1024 bits	CBE	72.9 ± 1.3	58.1
	SP	73.0 ± 1.3	59.2
	threshold [1]	73.8 ± 1.3	60.1
4096 bits	BP	73.5 ± 1.4	59.1
	CBE	76.0 ± 1.5	63.2
	SP	75.9 ± 1.4	63.0
8192 bits	BP	76.3 ± 1.5	63.3
	SP	76.8 ± 1.4	64.2
16384 bits	SP	77.1 ± 1.6	64.5

Figure 4: Demonstration of Dolphin’s **element-level** parsing across diverse scenarios. Input images are shown in the top row, with corresponding recognition results in the bottom row. **Left:** Text paragraph parsing in complex layouts. **Middle:** Bilingual text paragraph recognition. **Right:** Complex table parsing (rendered results shown).

extracting 1,856 text paragraphs from our Dolphin-Page. Unlike page-level evaluation which considers both reading order prediction and content recognition, this element-level evaluation focuses solely on fundamental text recognition capability.

(b) Formula. For formula recognition evaluation, we utilize three public benchmarks ([Wang et al., 2024a](#)) with different complexity levels: SPE with 6,762 simple printed expressions, SCE containing 4,742 screen capture formulas, and CPE consisting of 5,921 complex mathematical expressions. We adopt Character Difference Metric (CDM), which measures the character-level edit distance between predictions and ground truth.

(c) Table. The table recognition evaluation is conducted on two widely-used benchmarks: PubTabNet ([Zhong et al., 2020](#)) and PubTab1M ([Smock et al., 2022](#)). The test set of PubTabNet contains 7,904 table images from scientific papers, while PubTab1M’s test set consists of 10,000 more challenging samples. Both benchmarks evaluate the model’s capability in understanding table structures and recognizing cell contents using TEDS (Tree-Edit-Distance-based Similarity) as the metric, which computes the similarity between the predicted and ground-truth HTML table structure.

5 Experiment

5.1 Implementation Details

In the proposed Dolphin, the encoder uses a Swin Transformer with a window size of 7 and hierarchical structure ([2, 2, 14, 2] encoder layers with [4, 8, 16, 32] attention heads). The decoder contains 10 Transformer layers with a hidden dimension of 1024. We train the model using AdamW optimizer with a learning rate of 5e-5 and cosine decay sched-

ule. The training is conducted on 40 A100 GPUs for 2 epochs, using a batch size of 16 per device through gradient accumulation.

We use normalized coordinates for bounding boxes. Specifically, we maintain the aspect ratio of input document images by first resizing the longer edge to 896 pixels, then padding to create a square image of 896×896 pixels. The normalized bounding box coordinates correspond to positions within this final 896×896 padded image.

5.2 Comparison with Existing Methods

Comprehensive evaluations are conducted on both full-page document parsing (plain and complex documents) and individual element recognition tasks (text paragraphs, tables, and formulas).

Page-level Parsing. We evaluate Dolphin’s performance on Fox-Page (English and Chinese) and Dolphin-Page benchmarks. As shown in Table 1, despite its lightweight architecture (322M parameters), Dolphin achieves superior performance compared to both integration-based methods and larger VLMs. For pure text documents, Dolphin achieves an edit distances of 0.0114 and 0.0131 on English and Chinese test sets respectively, outperforming specialized VLMs like GOT (with edit distances of 0.035 and 0.038) and general VLMs like GPT-4.1 (with edit distances of 0.0489 and 0.2549). The advantage becomes more evident on Dolphin-Page, where Dolphin achieves an edit distance of 0.1283, outperforming all baselines in handling documents with mixed elements like tables and formulas. Furthermore, with parallel parsing design, Dolphin demonstrates considerable efficiency gains, achieving 0.1729 FPS, which is nearly 2x faster than the most efficient baseline (Mathpix at 0.0944 FPS).

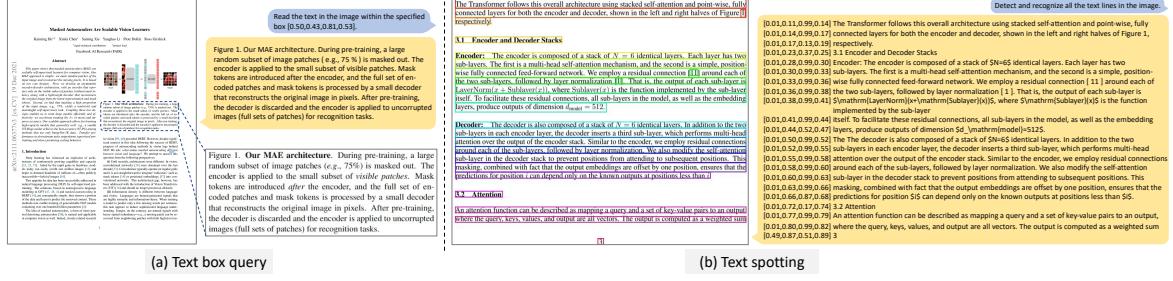


Figure 5: Additional capabilities of Dolphin. **Left:** Parsing the text content from a given bounding box region. **Right:** Text spotting results showing detected text lines (visualized in the image) and their content.

We visualize three representative cases in Figure 3, showing the complete pipeline from layout analysis (Stage 1) to element-specific parsing (Stage 2), and finally to the rendered document. As demonstrated, Dolphin accurately captures both layout structure and textual content. As shown in Figure 5 (left), Dolphin also exhibits strong text extraction capabilities by accurately parsing content from specified bounding box regions.

Element-level Parsing. Beyond page-level parsing, we conduct extensive experiments to evaluate Dolphin’s performance on individual elements, as shown in Table 3. For text paragraph parsing, Dolphin achieves competitive results on both Fox-Block and Dolphin-Block test sets. In formula recognition, Dolphin demonstrates strong capabilities across different complexity levels (SPE, SCE, and CPE), achieving competitive CDM scores comparable to specialized formula recognition methods. For table parsing, our approach shows promising results on both PubTabNet and PubTab1M benchmarks, effectively capturing both structural relationships and cell contents. These consistent strong results across text paragraphs, formulas, and tables demonstrate Dolphin’s competitive performance in fundamental recognition tasks.

We further show Dolphin’s robustness in Figure 4 through three scenarios: text paragraphs with complex layouts, bilingual text recognition, and structured tables with intricate formats. As shown in Figure 5 (right), Dolphin also supports text spotting by detecting and parsing text lines.

5.3 Ablation Studies

We conduct extensive experiments to validate the effectiveness of the core components in Dolphin.

Parallel Decoding. To investigate the efficiency gains from our parallel decoding strategy in stage 2, we compare our approach with a sequential autoregressive decoding baseline. As present in Ta-

Method	ED ↓	FPS ↑
Dolphin	0.1028	0.1729
Parallel → Sequential Decoding	-	0.0971
Type-specific → Generic Prompts	0.1613	-
Element Cropping → Box Query	0.1849	-

Table 4: Ablation studies on Dolphin. The first row shows the performance of our full model. The evaluation is conducted on Dolphin-Page dataset.

ble 4, parallel decoding achieves a 1.8× speedup (0.1729 vs. 0.0971 FPS) while maintaining the same parsing accuracy. The speedup is bounded by two factors: (a) the preprocessing overhead for each element before network inference, and (b) the batch size constraint (maximum 16 elements per batch) due to GPU memory limitations, requiring multiple inference passes for documents with numerous elements. Note that existing off-the-shelf autoregressive parallel decoding solutions (Kwon et al., 2023) can be leveraged to further accelerate inference speed.

Type-specific vs. Generic Prompts. To investigate the effectiveness of type-specific prompting in the second stage, we compare Dolphin with a baseline variant that uses a generic prompt “*Read text in the image.*” for all element parsing tasks. As shown in Table 4, our type-specific prompting strategy significantly outperforms the generic baseline (0.1283 vs. 0.1613 in ED). A representative case is shown in Figure 6, where the generic prompt misidentifies a table as a LaTeX formula, while our type-specific prompt successfully parses and renders it. These results demonstrate that incorporating prior knowledge through type-specific prompting effectively improves the model’s ability to handle different document elements.

Element Cropping vs. Box Query. To validate our element cropping strategy in the second stage, we compare it with an alternative box query approach that directly prompts the model to recog-

encoder	dec. depth	ft acc	hours	speedup
ViT-L, w/ [M]	8	84.2	42.4	-
ViT-L	8	84.9	15.4	2.8x
ViT-L	1	84.8	11.6	3.7x
ViT-H, w/ [M]	8	-	119.6 [†]	-
ViT-H	8	85.8	34.5	3.5x
ViT-H	1	85.9	29.3	4.1x

Generic Prompt: Misidentified as formula	$\$\\begin{array}{ l l l l l }\\text { encoder } & \text { dec. depth } & \text { ft acc } & \text { hours } & \text { speedup } \\\hline\text { ViT-L, w/ [M] } & 8 & 84.2 & 42.4 & - \\ \text { ViT-L } & 8 & 84.9 & 15.4 & 2.8 x \\ \text { ViT-L } & 1 & 84.8 & 11.6 & 3.7 x \\ \text { ViT-H, w/ [M] } & 8 & - & 119.6^{\dagger} & - \\ \text { ViT-H } & 8 & 85.8 & 34.5 & 3.5 x \\ \text { ViT-H } & 1 & 85.9 & 29.3 & 4.1 x \\ \hline\end{array}$
Type-specific Prompt: Correctly parsed as HTML table and successfully rendered	$\$\\begin{array}{ l l l l l }\\text { encoder } & \text { dec. depth } & \text { ft acc } & \text { hours } & \text { speedup } \\\hline\text { ViT-L, w/ [M] } & 8 & 84.2 & 42.4 & - \\ \text { ViT-L } & 8 & 84.9 & 15.4 & 2.8 x \\ \text { ViT-L } & 1 & 84.8 & 11.6 & 3.7 x \\ \text { ViT-H, w/ [M] } & 8 & - & 119.6^{\dagger} & - \\ \text { ViT-H } & 8 & 85.8 & 34.5 & 3.5 x \\ \text { ViT-H } & 1 & 85.9 & 29.3 & 4.1 x \\ \hline\end{array}$

Figure 6: A case study demonstrating the effectiveness of type-specific prompts. The generic prompt misidentifies the table as a formula, while our approach correctly parses and renders the table in HTML format.

nize elements at specific box (see Figure 5 (left)). As shown in Table 4, our cropping strategy achieves better performance than the box query method. This is likely because cropping provides the model with a focused view of each element, following a “what you see is what you get” principle, while the box query approach increases task complexity by requiring the model to simultaneously handle location understanding and content recognition.

6 Conclusion

We present Dolphin, a novel document image parsing model that leverages an analyze-then-parse paradigm to address the challenges in document parsing. Our approach first performs page-level layout analysis to generate structured layout elements in reading order, then enables parallel element parsing through heterogeneous anchor prompting. This two-stage design effectively balances efficiency and accuracy, while maintaining a lightweight architecture. Through extensive experiments, we demonstrate Dolphin’s strong performance in both page-level and element-level parsing tasks, particularly excelling in handling complex documents with interleaved tables, formulas, and rich formatting in both Chinese and English.

Limitations

Despite Dolphin’s promising performance, there are several limitations worth noting. First, Dolphin primarily supports documents with standard

horizontal text layout, showing limited capability in parsing vertical text like ancient manuscripts. Second, while Dolphin handles both Chinese and English documents effectively, its multilingual capacity (Tang et al., 2024b) needs to be expanded. Nevertheless, we demonstrate some cases exhibiting emergent multilingual document parsing capabilities in the supplementary materials. Third, although we achieve efficiency gains through parallel element parsing, there is potential for further optimization through parallel processing of text lines and table cells. Fourth, handwriting recognition capabilities require further enhancement.

References

- Haoli Bai, Zhiguang Liu, Xiaojun Meng, Shuang Liu, LUO Yifeng, Rongfu Zheng, Liangwei Wang, Lu Hou, Jiansheng Wei, Xin Jiang, et al. 2023. Wukong-Reader: Multi-modal pre-training for fine-grained visual document understanding. In *Proceedings of the Annual Meeting Of The Association For Computational Linguistics*.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. 2025. Qwen2.5-VL technical report. *arXiv preprint arXiv:2502.13923*.
- Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. Nougat: Neural optical understanding for academic documents. In *Proceedings of the International Conference on Learning Representations*.
- Mingxu Chai, Ziyu Shen, Chong Zhang, Yue Zhang, Xiao Wang, Shihan Dou, Jihua Kang, Jiazheng Zhang, and Qi Zhang. 2024. DocFusion: A unified framework for document parsing tasks. *arXiv preprint arXiv:2412.12505*.
- Song Chen, Xinyu Guo, Yadong Li, Tao Zhang, Mingan Lin, Dongdong Kuang, Youwei Zhang, Lingfeng Ming, Fengyu Zhang, Yuran Wang, et al. 2025. Ocean-OCR: Towards general OCR application via a vision-language model. *arXiv preprint arXiv:2501.15558*.
- Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024. InternVL: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 24185–24198.
- Changxu Duan and Sabine Bartsch. LaTex rainbow: Open source document layout semantic annotation framework. In *Proceedings of the Workshop for Natural Language Processing Open Source Software*.

- Hao Feng, Qi Liu, Hao Liu, Jingqun Tang, Wengang Zhou, Houqiang Li, and Can Huang. 2024. DocPedia: Unleashing the power of large multimodal model in the frequency domain for versatile document understanding. *Science China Information Sciences*, 67(12):1–14.
- 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. *arXiv preprint arXiv:2308.11592*.
- Jonathan Herzig, Paweł Krzysztof Nowak, Thomas Mueller, Francesco Piccinno, and Julian Eisenschlos. 2020. TaPas: Weakly supervised table parsing via pre-training. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333.
- Anwen Hu, Haiyang Xu, Liang Zhang, Jiabo Ye, Ming Yan, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. 2024. mPLUG-DocOwl2: High-resolution compressing for OCR-free multi-page document understanding. *arXiv preprint arXiv:2409.03420*.
- Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. LayoutLMv3: Pre-training for document ai with unified text and image masking. In *Proceedings of the ACM International Conference on Multimedia*, pages 4083–4091.
- Donghyun Kim, Teakgyu Hong, Moonbin Yim, Yoonsik Kim, and Geewook Kim. 2023. On web-based visual corpus construction for visual document understanding. In *Proceedings of the International Conference on Document Analysis and Recognition*, pages 297–313.
- Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. 2022. OCR-free document understanding transformer. In *Proceedings of the European Conference on Computer Vision*, pages 498–517.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the Symposium on Operating Systems Principles*, pages 611–626.
- Mike Lewis. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.
- Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. 2024. Monkey: Image resolution and text label are important things for large multi-modal models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 26763–26773.
- Chenglong Liu, Haoran Wei, Jinyue Chen, Lingyu Kong, Zheng Ge, Zining Zhu, Liang Zhao, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. 2024a. Focus anywhere for fine-grained multi-page document understanding. *arXiv preprint arXiv:2405.14295*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024b. Visual instruction tuning. In *Proceedings of the Neural Information Processing Systems*, volume 36.
- 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. *arXiv preprint arXiv:2403.04473*.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin Transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10012–10022.
- 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. *arXiv preprint arXiv:2309.11419*.
- John MacFarlane. 2013. Pandoc: a universal document converter. *URL: http://pandoc.org*, 8.
- Ahmed Nassar, Andres Marafioti, Matteo Omenetti, Maksym Lysak, Nikolaos Livathinos, Christoph Auer, Lucas Morin, Rafael Teixeira de Lima, Yusik Kim, A Said Gurbuz, et al. 2025. SmolDocling: An ultra-compact vision-language model for end-to-end multi-modal document conversion. *arXiv preprint arXiv:2503.11576*.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*.
- Jake Poznanski, Jon Borchardt, Jason Dunkelberger, Regan Huff, Daniel Lin, Aman Rangapur, Christopher Wilhelm, Kyle Lo, and Luca Soldaini. 2025. olmOCR: Unlocking trillions of tokens in PDFs with vision language models. *arXiv preprint arXiv:2502.18443*.
- Brandon Smock, Rohith Pesala, and Robin Abraham. 2022. PubTables-1M: Towards comprehensive table extraction from unstructured documents. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4634–4642.
- Jingqun Tang, Chunhui Lin, Zhen Zhao, Shu Wei, Binghong Wu, Qi Liu, Hao Feng, Yang Li, Siqi Wang, Lei Liao, et al. 2024a. TextSquare: Scaling up text-centric visual instruction tuning. *arXiv preprint arXiv:2404.12803*.

- Jingqun Tang, Qi Liu, Yongjie Ye, Jinghui Lu, Shu Wei, Chunhui Lin, Wanqing Li, Mohamad Fitri Faiz Bin Mahmood, Hao Feng, Zhen Zhao, et al. 2024b. MTVQA: Benchmarking multilingual text-centric visual question answering. *arXiv preprint arXiv:2405.11985*.
- Jingqun Tang, Su Qiao, Benlei Cui, Yuhang Ma, Sheng Zhang, and Dimitrios Kanoulas. 2022a. You can even annotate text with voice: Transcription-only-supervised text spotting. In *Proceedings of the ACM International Conference on Multimedia*, pages 4154–4163.
- Jingqun Tang, Wenqing Zhang, Hongye Liu, MingKun Yang, Bo Jiang, Guanglong Hu, and Xiang Bai. 2022b. Few could be better than all: Feature sampling and grouping for scene text detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4563–4572.
- Zineng Tang, Ziyi Yang, Guoxin Wang, Yuwei Fang, Yang Liu, Chenguang Zhu, Michael Zeng, Cha Zhang, and Mohit Bansal. 2023. Unifying vision, text, and layout for universal document processing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 19254–19264.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the Neural Information Processing Systems*, pages 5998–6008.
- Bin Wang, Zhuangcheng Gu, Guang Liang, Chao Xu, Bo Zhang, Botian Shi, and Conghui He. 2024a. UnimerNet: A universal network for real-world mathematical expression recognition. *arXiv preprint arXiv:2404.15254*.
- Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, et al. 2024b. MinerU: An open-source solution for precise document content extraction. *arXiv preprint arXiv:2409.18839*.
- Dongsheng Wang, Natraj Raman, Mathieu Sibue, Zhiqiang Ma, Petr Babkin, Simerjot Kaur, Yulong Pei, Armineh Nourbakhsh, and Xiaomo Liu. 2024c. DocLLM: A layout-aware generative language model for multimodal document understanding. In *Proceedings of the Annual Meeting Of The Association For Computational Linguistics*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024d. Qwen2-VL: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Peng Wang, Zhaohai Li, Jun Tang, Humen Zhong, Fei Huang, Zhibo Yang, and Cong Yao. 2024e. Platypus: A generalized specialist model for reading text in various forms. In *Proceedings of the European Conference on Computer Vision*, pages 165–183.
- Yonghui Wang, Wengang Zhou, Hao Feng, Keyi Zhou, and Houqiang Li. 2023. Towards improving document understanding: An exploration on text-grounding via mllms. *arXiv preprint arXiv:2311.13194*.
- Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. 2024a. Vary: Scaling up the vision vocabulary for large vision-language model. In *Proceedings of the European Conference on Computer Vision*, pages 408–424.
- Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, et al. 2024b. General OCR theory: Towards OCR-2.0 via a unified end-to-end model. *arXiv preprint arXiv:2409.01704*.
- Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, et al. 2024. Deepseek-VL2: Mixture-of-experts vision-language models for advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*.
- Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, et al. 2020a. LayoutLMv2: Multi-modal pre-training for visually-rich document understanding. *arXiv preprint arXiv:2012.14740*.
- Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. 2020b. LayoutLM: Pre-training of text and layout for document image understanding. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1192–1200.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The dawn of LMMs: Preliminary explorations with GPT-4V (ision). *arXiv preprint arXiv:2309.17421*.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. 2024. MiniCPM-V: A GPT-4V level MLLM on your phone. *arXiv preprint arXiv:2408.01800*.
- Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu, Chenliang Li, Junfeng Tian, et al. 2023a. mPLUG-DocOwl: Modularized multimodal large language model for document understanding. *arXiv preprint arXiv:2307.02499*.
- Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Guohai Xu, Chenliang Li, Junfeng Tian,

- Qi Qian, Ji Zhang, et al. 2023b. UReader: Universal OCR-free visually-situated language understanding with multimodal large language model. *arXiv preprint arXiv:2310.05126*.
- Ya-Qi Yu, Minghui Liao, Jihao Wu, Yongxin Liao, Xiaoyu Zheng, and Wei Zeng. 2024a. TextHawk: Exploring efficient fine-grained perception of multimodal large language models. *arXiv preprint arXiv:2404.09204*.

Ya-Qi Yu, Minghui Liao, Jiwen Zhang, and Jihao Wu. 2024b. TextHawk2: A large vision-language model excels in bilingual OCR and grounding with 16x fewer tokens. *arXiv preprint arXiv:2410.05261*.

Jianshu Zhang, Jun Du, Shiliang Zhang, Dan Liu, Yulong Hu, Jinshui Hu, Si Wei, and Lirong Dai. 2017. Watch, attend and parse: An end-to-end neural network based approach to handwritten mathematical expression recognition. *Pattern Recognition*, 71:196–206.

Weichao Zhao, Hao Feng, Qi Liu, Jingqun Tang, Binghong Wu, Lei Liao, Shu Wei, Yongjie Ye, Hao Liu, Wengang Zhou, et al. 2024a. TabPedia: Towards comprehensive visual table understanding with concept synergy. In *Proceedings of the Neural Information Processing Systems*, volume 37, pages 7185–7212.

Zhen Zhao, Jingqun Tang, Chunhui Lin, Binghong Wu, Can Huang, Hao Liu, Xin Tan, Zhizhong Zhang, and Yuan Xie. 2024b. Multi-modal in-context learning makes an ego-evolving scene text recognizer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 15567–15576.

Xu Zhong, Elaheh ShafieiBavani, and Antonio Jimeno Yepes. 2020. Image-based table recognition: data, model, and evaluation. In *Proceedings of the European Conference on Computer Vision*, pages 564–580.

In this supplementary material, we provide additional experimental results and implementation details to complement our main paper. Specifically, we present more qualitative results demonstrating Dolphin’s parsing capabilities, elaborate on the supported element types, detail our training process, and showcase our synthetic data.

A Qualitative Results

To further demonstrate the superior capabilities of Dolphin, we present comprehensive page-level and element-level parsing results.

Page-level. First, the examples in Figure 8 cover diverse document scenarios, including textbook pages with dense formulas, triple-column English academic papers, and double-column Chinese papers with tables. The results demonstrate that Dolphin can effectively handle documents with different languages, layouts, and element types, maintaining high parsing quality.

Furthermore, we showcase Dolphin’s versatility in other text-rich scenarios through Figure 9, where we test the model on mobile phone screenshots, shopping receipts, and webpage captures. These results indicate that Dolphin can accurately capture both the structural layout and textual content in these everyday scenarios.

Element-level. For fine-grained parsing capabilities, we first demonstrate Dolphin’s formula recognition in Figure 10, where we evaluate three types of formulas: inline formulas, single-line block formulas, and multi-line block formulas. The results show that Dolphin can accurately parse formulas of varying complexity and layout formats.

We further evaluate Dolphin’s table parsing ability in Figure 11, where we test the model on a challenging case containing hundreds of cells. As shown, Dolphin successfully handles this large-scale structured table with precise content recognition and layout preservation.

B Element Design

In this section, we elaborate on Dolphin’s supported element types and element-specific parsing strategies through heterogeneous prompting.

Element Types. Our Dolphin supports 15 different types of elements commonly found in document images. Table 5 provides a comprehensive overview of these elements, covering various components from headers to specialized content blocks.

No.	Element	Description
1	title	Paper/document title
2	author	Author names
3	sec	First-level section headings
4	sub_sec	Second-level section headings
5	para	Paragraphs
6	header	Page headers
7	foot	Page footers
8	fnote	Footnotes
9	watermark	Non-content watermarks
10	fig	Figures and images
11	tab	Tables
12	cap	Figure/table captions
13	anno	Figure/table annotations
14	alg	Code blocks/pseudocode
15	list	List-type content

Table 5: An overview of element types supported by Dolphin. These elements cover the majority of content structures found in documents.

Note that in Stage 1 (page-level layout analysis), we intentionally avoid treating formulas as independent elements. This design choice allows Stage 2 (element-level parsing) to leverage broader contextual information when recognizing mathematical expressions, as formulas are often semantically connected with their surrounding text.

Heterogeneous Anchor Prompting. We summarize the prompts used in Dolphin in Table 6. The first three prompts (page-level layout analysis, text paragraph parsing, and table parsing) are designed for full-page document image parsing, while the latter two (text spotting and text box query) enable additional capabilities for flexible text recognition tasks. Additionally, our Dolphin can also serve as a formula recognition expert model using the text paragraph parsing prompt.

In Stage 2, tables are processed with a dedicated table-specific prompt for structured HTML parsing, while all other elements are treated as text paragraphs and parsed using a unified prompt. This dichotomous design distinguishes structured HTML content from plain text, while also providing robustness against potential element misclassification, as parsing accuracy remains high regardless of element type classification errors.

C Training Details

In this section, we provide more details about Dolphin’s training process, including multi-task training strategy, model initialization, and other implementation considerations.

Instruction Tuning. During training phase, we adopt a dynamic task selection strategy for our

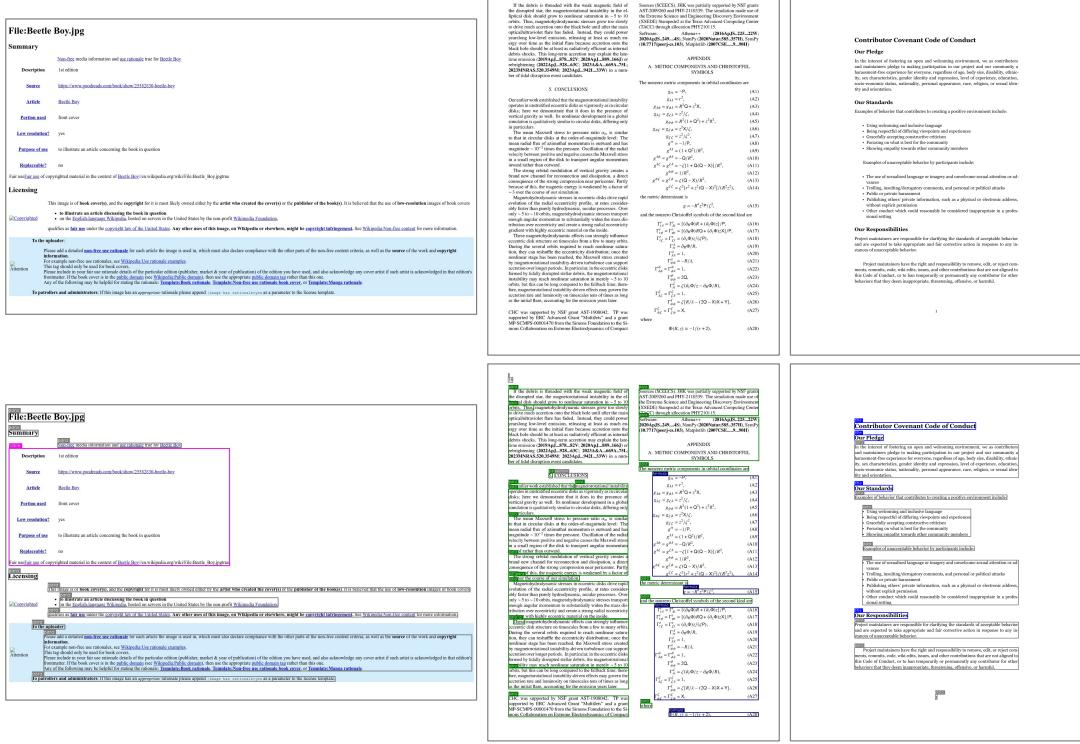


Figure 7: Examples of synthetic training data generated from different source formats. **Top:** rendered document images from HTML (left), LaTeX (middle), and Markdown (right) sources. **Bottom:** corresponding paragraph-level annotations visualized with colored regions.

Task	Prompt
Page-level Layout Analysis	Parse the reading order of this document.
Text Paragraph/Formula Parsing	Read text in the image.
Table Parsing	Parse the table in the image.
Text Spotting	Detect and recognize all the text lines in the image.
Text Box Query	Read the text in the image within the specified box [x1,y1,x2,y2].

Table 6: Different types of prompts used in Dolphin for document parsing tasks.

instruction-based framework. Specifically, given a training sample, we randomly select an applicable task from the above five tasks based on its available annotations. This selection is used to construct question-answer pairs. For instance, given a page image with only paragraph-level bounding boxes and content annotations, the available tasks for this sample would include element-level text paragraph parsing and page-level box query parsing.

Model Initialization. We initialize Dolphin with the pretrained weights from Donut (Kim et al., 2022), which lacks instruction-following abilities. Then, through our instruction tuning, we extend the model’s capabilities to understand and execute diverse prompts, enabling analysis of document layout, reading order, and various textual elements

including text paragraphs, tables, and formulas.

Training Loss. Following standard practice in autoregressive language models, we optimize Dolphin using the cross-entropy loss between the predicted token distributions and ground truth ones.

D Synthetic Data Examples

To enrich training data diversity, we synthesize document images from different source formats, including HTML, LaTeX, and Markdown documents. Figure 7 shows three representative examples of our synthetic data. For each format, we show the rendered document (top row) and its corresponding paragraph-level annotations (bottom row).

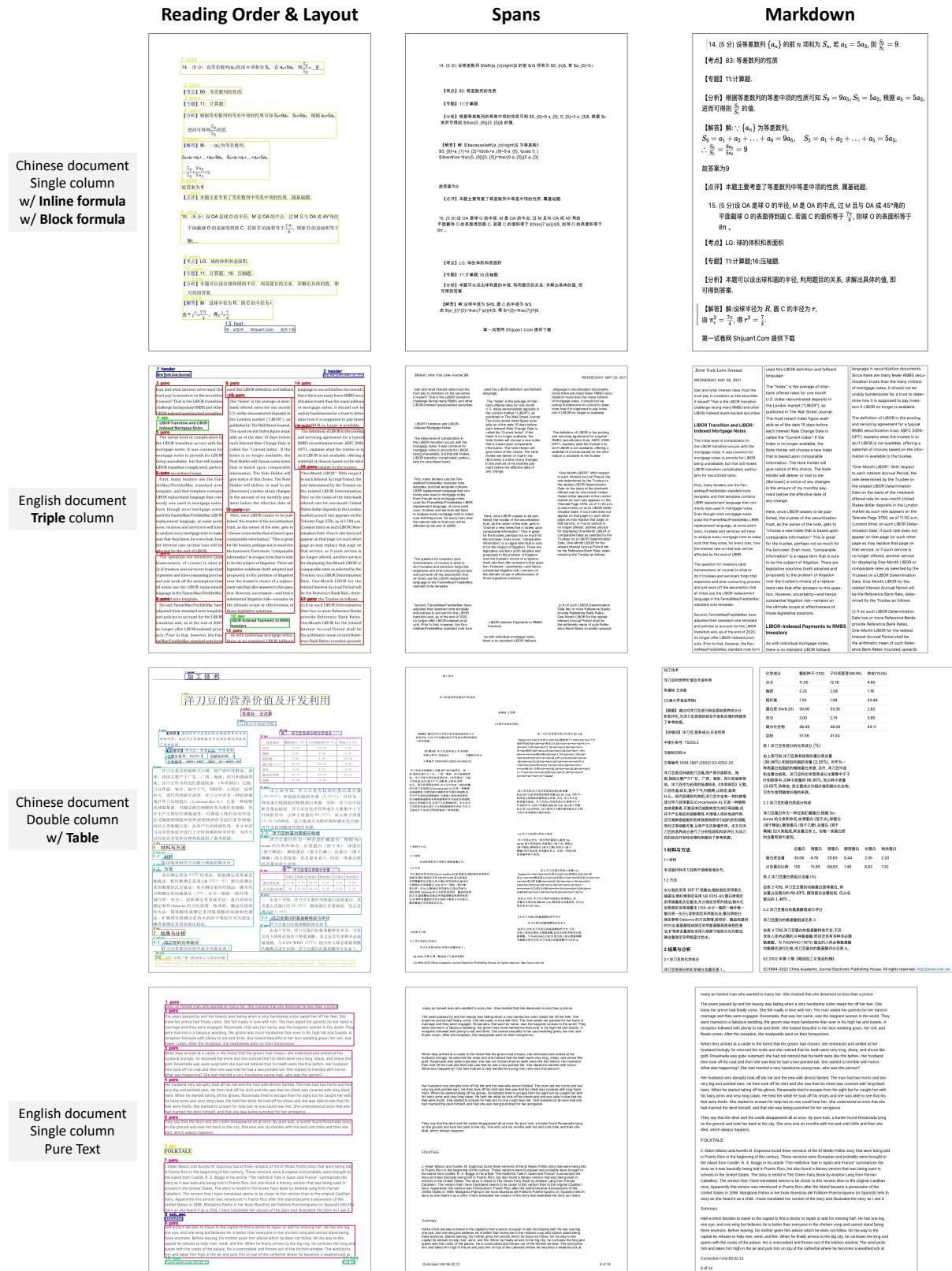


Figure 8: Visualization of Dolphin’s **page-level** parsing results. **Left:** Layout analysis form Stage 1 with predicted element boundaries and reading order. **Middle:** Element-specific parsing outputs from Stage 2. **Right:** Final rendered document in markdown format.

Input Image

Reading Order & Layout

Markdown / Spans

```

6:58
...
X ...
DeepSeek-V3 正式发布
原创 深度求索 DeepSeek
2024年12月26日 19:17 北京 2548人

今天，我们全新系列模型 DeepSeek-V3 首个版本上线并同步开源。
登录官网 chat.deepseek.com 即可与最新版 V3 模型对话。API 服务已同步更新，接口配置无需改动。当前版本的 DeepSeek-V3 暂不支持多模态输入输出。
性能对齐海外领军闭源模型
DeepSeek-V3 为自研 MoE 模型，671B 参数，激活 37B，在 14.8T token 上进行了预训练。
论文链接：
https://github.com/deepseek-ai/DeepSeek-V3/blob/main/DeepSeek_V3.pdf

DeepSeek-V3 多项评测成绩超越了 Qwen2.5-72B 和 Llama-3.1-405B 等其他开源模型，并在性能上和世界顶尖的闭源模型 GPT-4o 以及 Claude-3.5-Sonnet 不分伯仲。

```

Figure 9: Visualization of Dolphin’s **page-level** parsing results. **Left:** Input text-rich images including mobile phone screenshots, shopping receipts, and webpage captures. **Middle:** Layout analysis form Stage 1 with predicted element boundaries and reading order. **Right:** Final rendered document in markdown format for the first row, and element-specific parsing outputs from Stage 2 for the second and third rows.

Inline formula image is normalized by $\sum_{l=1}^L \sum_{k=1}^K A_{mlqk} = 1$. Here, we use normalized coordinates $\hat{\mathbf{p}}_q \in [0, 1]^2$ for

Parsing results is normalized by $\sum_{l=1}^L \sum_{k=1}^K A_{mlqk} = 1$. Here, we use normalized coordinates $\hat{\mathbf{p}}_q \in [0, 1]^2$ for

Rendered image is normalized by $\sum_{l=1}^L \sum_{k=1}^K A_{mlqk} = 1$. Here, we use normalized coordinates $\hat{\mathbf{p}}_q \in [0, 1]^2$ for

Block formula image

$$q_\sigma(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}\mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I}\right).$$

Parsing results

```
$$
q_{\sigma}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}\mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I}\right).
$$
```

Rendered image

$$q_\sigma(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}\mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I}\right).$$

Block formula image

$$\begin{aligned} \mathbb{E}[\nabla_\theta \mathcal{L}(\theta_t) | \theta_t] &= \nabla_\theta \left[\frac{1}{M} \sum_{i=1}^M \mathbb{E}[(\nabla^2 \mathcal{N}(\mathbf{x}_i; \theta_t) - f(\mathbf{x}_i))^2] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}[(\mathcal{N}(\mathbf{y}_j; \theta_t) - g(\mathbf{y}_j))^2] \right] \\ &= \nabla_\theta \left[\frac{1}{M} \sum_{i=1}^M \int_{\Omega} (\nabla^2 \mathcal{N}(\mathbf{x}; \theta_t) - f(\mathbf{x}))^2 \nu_1(\mathbf{x}) d\mathbf{x} + \frac{1}{N} \sum_{j=1}^N \int_{\partial\Omega} (\mathcal{N}(\mathbf{y}; \theta_t) - g(\mathbf{y}))^2 \nu_2(\mathbf{y}) d\mathbf{y} \right] \\ &= \nabla_\theta \left[\int_{\Omega} (\nabla^2 \mathcal{N}(\mathbf{x}; \theta_t) - f(\mathbf{x}))^2 \nu_1(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} (\mathcal{N}(\mathbf{y}; \theta_t) - g(\mathbf{y}))^2 \nu_2(\mathbf{y}) d\mathbf{y} \right] \\ &= \nabla_\theta \mathcal{J}(\mathcal{N}(\cdot; \theta_t)) \end{aligned}$$

Parsing results

```
\begin{array}{r}
\nabla_\theta \left[ \frac{1}{M} \sum_{i=1}^M \mathbb{E}[(\nabla^2 \mathcal{N}(\mathbf{x}_i; \theta_t) - f(\mathbf{x}_i))^2] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}[(\mathcal{N}(\mathbf{y}_j; \theta_t) - g(\mathbf{y}_j))^2] \right] \\
= \nabla_\theta \left[ \frac{1}{M} \sum_{i=1}^M \int_{\Omega} (\nabla^2 \mathcal{N}(\mathbf{x}; \theta_t) - f(\mathbf{x}))^2 \nu_1(\mathbf{x}) d\mathbf{x} + \frac{1}{N} \sum_{j=1}^N \int_{\partial\Omega} (\mathcal{N}(\mathbf{y}; \theta_t) - g(\mathbf{y}))^2 \nu_2(\mathbf{y}) d\mathbf{y} \right] \\
= \nabla_\theta \left[ \int_{\Omega} (\nabla^2 \mathcal{N}(\mathbf{x}; \theta_t) - f(\mathbf{x}))^2 \nu_1(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} (\mathcal{N}(\mathbf{y}; \theta_t) - g(\mathbf{y}))^2 \nu_2(\mathbf{y}) d\mathbf{y} \right] \\
= \nabla_\theta \mathcal{J}(\mathcal{N}(\cdot; \theta_t))
\end{array}
```

Rendered image

$$\begin{aligned} \mathbb{E}[\nabla_\theta \mathcal{L}(\theta_t) | \theta_t] &= \nabla_\theta \left[\frac{1}{M} \sum_{i=1}^M \mathbb{E}[(\nabla^2 \mathcal{N}(\mathbf{x}_i; \theta_t) - f(\mathbf{x}_i))^2] + \frac{1}{N} \sum_{j=1}^N \mathbb{E}[(\mathcal{N}(\mathbf{y}_j; \theta_t) - g(\mathbf{y}_j))^2] \right] \\
&= \nabla_\theta \left[\frac{1}{M} \sum_{i=1}^M \int_{\Omega} (\nabla^2 \mathcal{N}(\mathbf{x}; \theta_t) - f(\mathbf{x}))^2 \nu_1(\mathbf{x}) d\mathbf{x} + \frac{1}{N} \sum_{j=1}^N \int_{\partial\Omega} (\mathcal{N}(\mathbf{y}; \theta_t) - g(\mathbf{y}))^2 \nu_2(\mathbf{y}) d\mathbf{y} \right] \\
&= \nabla_\theta \left[\int_{\Omega} (\nabla^2 \mathcal{N}(\mathbf{x}; \theta_t) - f(\mathbf{x}))^2 \nu_1(\mathbf{x}) d\mathbf{x} + \int_{\partial\Omega} (\mathcal{N}(\mathbf{y}; \theta_t) - g(\mathbf{y}))^2 \nu_2(\mathbf{y}) d\mathbf{y} \right] \\
&= \nabla_\theta \mathcal{J}(\mathcal{N}(\cdot; \theta_t))
\end{aligned}$$

Figure 10: Visualization of Dolphin’s **formula** parsing results. From top to bottom, we show three formula types: **inline formula**, **single-line block formula**, and **multi-line block formula**. For each case, we visualize the complete parsing pipeline: input formula image (top), LaTeX parsing output (middle), and rendered formula (bottom). These results demonstrate Dolphin’s capability to accurately parse formulas of varying complexity.

Shots	Method	AR	BG	DE	EL	EN	ES	FR	HI	RU	SW	TH	TR	UR	VI	ZH	Avg.
1	FT	33.2	33.3	33.3	33.1	33.3	33.2	32.8	33.0	33.3	33.0	33.3	32.9	32.9	33.2	33.2	33.1
	SP	35.4	36.4	36.4	36.6	36.5	37.6	37.9	36.0	37.5	34.1	35.9	34.7	35.0	35.5	36.7	36.1
	PCT	33.2	35.4	34.8	35.1	35.9	35.3	35.7	34.6	36.2	33.8	34.6	34.3	33.1	34.9	35.0	34.8
	MPT	37.0	38.5	37.8	38.1	38.6	38.1	38.7	37.2	38.5	36.5	37.1	37.6	37.3	37.9	35.7	37.6
2	FT	33.5	33.3	33.7	33.3	34.1	33.5	33.7	33.2	33.5	33.3	33.8	33.6	33.5	34.0	33.3	33.5
	SP	36.6	37.9	38.0	38.2	38.0	38.0	38.3	36.2	38.9	34.3	37.5	34.6	35.2	37.2	36.7	37.0
	PCT	34.1	39.0	39.1	38.2	39.9	40.6	40.5	37.9	39.9	36.5	37.2	36.9	34.7	37.9	37.1	38.0
	MPT	41.6	42.8	40.8	43.2	43.2	42.5	42.8	40.4	43.3	36.8	40.5	41.0	41.1	41.4	38.2	41.3
4	FT	34.2	34.5	34.1	34.3	34.1	34.1	34.5	34.0	34.3	33.7	34.0	34.0	34.1	34.2	34.2	34.1
	SP	37.4	39.7	39.2	39.7	40.2	38.9	40.5	37.1	40.6	35.3	38.1	35.3	36.9	37.2	38.9	38.3
	PCT	33.9	37.2	37.0	36.2	37.0	37.7	37.5	36.4	37.4	34.2	34.7	33.5	35.0	35.6	35.9	35.9
	MPT	42.9	43.6	44.3	43.6	45.5	44.2	44.1	42.8	44.1	40.2	43.4	42.7	42.4	43.8	43.1	43.4
8	FT	32.8	32.7	32.8	32.9	32.7	32.6	33.0	33.3	32.7	33.0	33.2	33.0	33.1	32.5	32.4	32.8
	SP	37.4	39.6	38.1	39.1	40.0	38.8	39.2	36.5	40.3	35.6	38.5	35.3	36.5	37.8	37.1	38.0
	PCT	40.2	40.6	40.9	41.7	41.9	41.7	41.6	41.0	40.6	39.2	41.4	41.4	38.4	41.3	41.2	40.9
	MPT	42.7	43.0	41.9	42.4	43.1	42.3	42.1	40.8	42.6	39.4	41.9	40.1	40.7	42.2	40.2	41.7
16	FT	33.6	33.4	33.3	33.5	34.1	33.4	33.3	33.4	33.4	33.6	33.4	33.5	33.3	33.5	33.5	33.5
	SP	39.5	39.9	39.1	40.4	41.1	40.2	40.4	37.4	40.7	37.1	39.3	36.5	36.0	38.2	38.3	38.9
	PCT	43.6	40.8	36.9	45.7	46.5	41.5	44.3	44.8	42.4	40.1	43.9	43.7	42.5	44.7	44.8	43.1
	MPT	43.5	43.8	44.0	43.9	45.2	44.2	44.3	42.9	43.4	40.2	42.5	41.8	42.0	43.4	42.2	43.1
32	FT	36.1	36.3	35.7	35.7	36.5	36.2	36.0	35.5	35.9	35.0	35.6	36.0	35.4	36.1	36.3	35.9
	SP	41.7	43.4	42.8	42.3	44.9	42.9	43.3	39.2	43.5	37.7	40.2	41.1	39.8	43.0	39.8	41.7
	PCT	45.7	45.4	44.4	47.4	49.6	45.5	48.8	46.7	45.5	40.3	41.6	44.3	42.9	46.7	45.6	45.4
	MPT	47.1	47.6	47.9	47.1	48.2	47.6	46.3	47.3	43.3	47.2	47.2	45.3	49.0	47.1	47.3	
64	FT	41.4	41.2	41.5	40.7	42.6	41.4	40.8	41.2	40.2	40.6	40.7	41.4	40.5	41.7	41.0	41.1
	SP	43.9	44.2	47.5	45.1	50.5	47.9	48.6	41.8	43.7	41.3	45.9	45.3	42.6	47.6	45.1	45.4
	PCT	48.1	50.2	49.3	50.6	51.1	50.9	51.3	47.6	49.1	44.6	47.3	47.4	44.0	49.7	48.2	48.6
	MPT	50.7	52.7	53.1	52.2	55.4	53.8	53.1	50.2	51.0	46.2	51.5	50.4	49.1	53.0	52.3	51.7
128	FT	43.9	44.4	44.4	43.7	46.3	44.6	44.5	42.9	42.7	41.7	43.0	43.2	42.7	44.9	43.8	43.8
	SP	46.2	46.8	47.8	47.6	53.0	48.5	49.6	47.3	45.5	41.7	47.5	46.4	44.5	45.6	48.7	47.1
	PCT	50.4	51.9	52.8	53.4	55.0	53.8	53.3	51.5	51.7	47.0	50.0	50.9	47.9	51.7	51.2	51.5
	MPT	53.2	56.1	56.0	55.4	57.4	56.4	56.6	53.5	54.8	48.6	54.0	53.1	51.8	55.2	55.4	54.5
256	FT	53.3	55.6	56.5	55.0	58.8	56.9	56.4	52.5	53.6	50.5	52.6	53.8	51.3	55.0	53.0	54.3
	SP	52.7	55.2	49.6	53.7	59.5	55.0	55.3	50.6	51.4	46.5	53.4	46.1	44.9	52.8	51.5	51.9
	PCT	54.7	56.7	56.3	57.9	60.3	58.3	58.3	54.6	55.2	51.6	55.6	54.6	52.6	57.4	55.8	56.0
	MPT	59.0	61.1	60.9	60.6	65.8	63.0	61.9	57.6	60.6	50.7	59.2	57.8	56.1	60.7	60.8	59.7

Shots	Method	AR	BG	DE	EL	EN	ES	FR	HI	RU	SW	TH	TR	UR	VI	ZH	Avg.
1	FT	33.2	33.3	33.3	33.1	33.3	33.2	32.8	33.0	33.3	33.0	33.3	32.9	32.9	33.2	33.2	33.1
	SP	35.4	36.4	36.4	36.6	36.5	37.6	37.9	36.0	37.5	34.1	35.9	34.7	35.0	35.5	36.7	36.1
	PCT	33.2	35.4	34.8	35.1	35.9	35.3	35.7	34.6	36.2	33.8	34.6	34.3	33.1	34.9	35.0	34.8
	MPT	37.0	38.5	37.8	38.1	38.6	38.1	38.7	37.2	38.5	36.5	37.1	37.6	37.3	37.9	37.6	37.6
2	FT	33.5	33.3	33.7	33.3	34.1	33.5	33.7	33.2	33.5	33.3	33.8	33.6	33.5	34.0	33.3	33.5
	SP	36.6	37.9	38.0	38.2	38.0	38.0	38.3	36.2	38.9	34.3	37.5	34.6	35.2	37.2	36.7	37.0
	PCT	34.1	39.0	39.1	38.2	39.9	40.6	40.5	37.9	39.9	36.5	38.5	35.3	36.5	37.8	37.1	38.0
	MPT	37.1	38.8	37.9	38.1	39.6	38.1	38.7	37.2	38.5	36.5	37.1	37.6	37.3	37.9	37.6	37.6
4	FT	34.2	34.5	34.1	34.3	34.1	34.1	34.5	34.0	34.3	33.7	34.0	34.0	34.1	34.2	34.2	34.1
	SP	37.4	39.7	39.2	39.7	40.2	38.9	40.5	37.9	39.9	36.5	38.5	37.7	37.2	38.9	38.3	38.3
	PCT	37.4	39.1	37.0	36.2	37.0	37.7	37.5	36.4	37.4	34.2	34.7	34.7	33.5	35.0	35.6	35.9
	MPT	37.4	39.1	37.0	36.2	37.0	37.7	37.5	36.4	37.4	34.2	34.7	34.7	33.5	35.0	35.6	35.9
8	FT	33.6	33.4	33.3	33.5	34.1	33.4	33.3	33.4	33.4	33.6	33.6	33.4	33.5	33.3	33.5	33.5
	SP	37.4	39.6	38.1	39.1	40.0	38.8	39.2	36.5	40.3	35.6	38.5	35.3	36.5	37.8	37.1	38.0
	PCT	40.2	40.6	40.9	41.7	41.9	41.7	41.6	41.0	40.6	39.2	41.4	41.4	38.4	41.3	41.2	40.9
	MPT	42.7	43.0	41.9	42.4	43.1	42.3	42.1	40.8	42.6	39.4	41.9	40.1	40.7	42.2	40.2	41.7
16	FT	33.6	33.4	33.3	33.5	34.1	33.4	33.3	33.4	33.4	33.6	33.6	33.4	33.5	33.3	33.5	33.5
	SP	39.5	39.9	39.1	40.4	41.1	40.2	40.4	37.4	40.7	37.1	39.3	36.5	36.0	38.2	38.3	38.9
	PCT	43.6	40.8	36.9	45.7	46.5	41.5	44.3	44.8	42.4	40.1	43.9	43.7	42.5	44.7	44.8	43.1
	MPT	43.5	43.8	44.0	43.9	45.2	44.2	44.3	42.9	43.4	40.2	42.5	41.8	42.0	43.4	42.2	43.1
32	FT	33.6	33.4	33.3	33.5	34.1	33.4	33.3	33.4	33.4	33.6	33.6	33.4	33.5	33.3	33.5	33.5
	SP	39.5	39.9	39.1	40.4	41.1	40.2	40.4	37.4	40.7	37.1	39.3	36.5	36.0	38.2	38.3	38.9
	PCT	43.6	40.8	36.9	45.7	46.5	41.5	44.3	44.8	42.4	40.1	43.9	43.7	42.5	44.7	44.8	