M3DocVQA: Multi-modal Multi-page Multi-document Understanding

Supplementary Material

Appendix

In this appendix, we show additional qualitative examples (Appendix A) and ethical considerations (Appendix B).

A. Additional Qualitative Examples

In Fig. 6 and Fig. 7, we provide additional qualitative examples of M3DocRAG (ColPali + Qwen2-VL 7B)'s question answering results on M3DocVQA examples. In Fig. 6, the question requires multi-hop reasoning across different pages/documents, and M3DocRAG could combine information from multiple retrieved pages. In Fig. 7, although ColPali did not retrieve the page that contains information about a team whose logo features a bat, Qwen-2 VL leverages its own knowledge 'Valencia CF has a logo featuring a bat', and could provide the final answer. Overall, the qualitative examples showcase that M3DocRAG can successfully tackle different questions whose answer sources exist in various modalities.

B. Ethical Considerations

Limitations. Since our multimodal retrieval models and multimodal LMs were trained with English-heavy datasets, they might not understand prompts or documents written in non-English. While our M3DocRAG framework can benefit many document understanding applications, the model components could present false or biased information. Thus, the framework should be used with human supervision in real-world applications. M3DocRAG is designed with flexibility so that users can update or replace components as more accurate solutions for each element of the framework become available in the future.

Data collection. We do not involve human subjects during data collection. We do not claim ownership/rights of the Wikipedia documents, and we attribute the source Wikipedia document URLs to all pages.

Question: "What distance was the AP Warrior fast race at the Del Mar Racetrack?"

ColPali + Qwen2-VL 7B: "Seven Furlongs"

Top 2 pages retrieved by ColPali

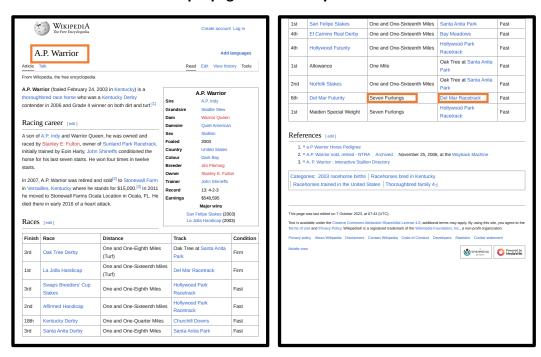


Figure 6. Qualitative example of ColPali + Qwen2-VL 7B on M3DocVQA. Image regions relevant to the question/answer are highlighted with orange boxes. The question requires multi-page/document reasoning.

Question: "What date was a player transferred in to Lorca FC in the 2017–18 season from the club with a logo featuring a bat?"

ColPali + Qwen2-VL 7B: "11 July 2017"

Top 1 page retrieved by ColPali



Figure 7. Qualitative example of ColPali + Qwen2-VL 7B on M3DocVQA. Image regions relevant to the question/answer are highlighted with orange boxes. The VQA component could combine both the retrieved knowledge (Tropi was transferred on 11 July 2017) and its own knowledge (Valencia CF has a logo with a bat) to provide the final answer.