Thematic session

Paper presentation: Group1

MultiX

Marcel Worring

Mapping out the Space of Human Feedback for Reinforcement Learning: A Conceptual Framework

YANNICK METZ, University of Konstanz, Germany and ETH Zurich, Switzerland DAVID LINDNER, ETH Zurich, Switzerland RAPHAËL BAUR, ETH Zurich, Switzerland MENNATALLAH EL-ASSADY, ETH Zurich, Switzerland

arXiv:2411.11761v2 [cs.LG] 20 Feb 2025

Human-Al Interaction



Creating a Review



A more formal model



Different State Spaces



Feedback Dimensions









Published as a conference paper at ICLR 2025

ColPali: EFFICIENT DOCUMENT RETRIEVAL WITH VISION LANGUAGE MODELS

Manuel Faysse^{*1,3} Hugues Sibille^{*1,4} Tony Wu^{*1} Bilel Omrani¹ Gautier Viaud¹ Céline Hudelot³ Pierre Colombo^{2,3} ¹Illuin Technology ²Equall.ai ³CentraleSupélec, Paris-Saclay ⁴ETH Zürich manuel.faysse@centralesupelec.fr

Faysse, M., Sibille, H., Wu, T., Omrani, B., Viaud, G., Hudelot, C., & Colombo, P. (2024). ColPali: Efficient Document Retrieval with Vision Language Models. *ArXiv, abs/2407.01449*.

Background –Multimodal Large Language Models (MLLMs) vs. Large Language Models (LLMs)

MLLMs are a superset of Large Language Models LLMs:

W They have the **added ability** to process and understand **images, charts, tables, and figures**.

They typically use an **LLM decoder as the backbone**, should retain **all text-based capabilities** (includes reasoning) of LLMs.



Qwen2.5 VL from 2025 Jan

General architecture of MLLM

Pixtral 12B from 2024 Oct

Why MLLMs for Document Understanding?

Documents are inherently multimodal—they combine **text, charts, tables, and figures**, where LLM struggle to extract insights from.

Documents represent knowledge across different domains with diverse formats – Scientific papers, financial reports, legal contracts, and medical records...

MLLMs unlock a new frontier in **document retrieval, summarization, and comprehension** by handling all these modalities **simultaneously**.



Scientific papers, slides, reports...

Main Paper Idea:

- Current document retrieval methods rely heavily on text extraction (OCR, parsing), **neglecting visual cues**.
- ColPali proposes a vision-based retrieval model that indexes document pages directly from images using Vision-Language Models (VLMs).
- **Outcome:** ColPali is **faster, simpler, and more accurate** than conventional retrieval systems.



How ColPali Works?

Multi-Vector Vision Embeddings:

• Uses PaliGemma-3B, a Vision-Language Model, to encode images and query into multi-vector embeddings

Late Interaction Mechanism:

• Inspired by **ColBERT**, performs **fine-grained** matching between query and document embeddings.

End-to-End Training with Contrastive Learning:

• **Trained on 118K query-page pairs** (academic + synthetic datasets).



ViDoRe: A New Benchmark for Document Retrieval

Introduced in this paper to evaluate retrieval performance across:

- Different document types (academic, administrative, scientific).
- Multiple modalities (text, tables, infographics, figures).
- Various languages (English, French).
- Two evaluation categories:
 - Academic Tasks: Repurposed VQA datasets (e.g., DocVQA, InfoVQA, arXivQA).
 - Practical Tasks: Domain-specific retrieval benchmarks (e.g., government, healthcare, energy).

Dataset	Language	# Queries	# Documents	Description
Academic Tasks				
DocVQA	English	500	500	Scanned documents from UCSF Industry
InfoVQA	English	500	500	Infographics scrapped from the web
TAT-DQA	English	1600	1600	High-quality financial reports
arXiVQA	English	500	500	Scientific Figures from arXiv
TabFQuAD	French	210	210	Tables scrapped from the web
Practical Tasks				
Energy	English	100	1000	Documents about energy
Government	English	100	1000	Administrative documents
Healthcare	English	100	1000	Medical documents
AI	English	100	1000	Scientific documents related to AI
Shift Project	French	100	1000	Environmental reports

Table 1: ViDoRe comprehensively evaluates multimodal retrieval methods.

5 RESULTS

	ArxivQ	DocQ	InfoQ	TabF	TATQ	Shift	AI	Energy	Gov.	Health.	Avg.
Unstructured text-only										ŝ.	
- BM25	347	34.1	-	-	44.0	59.6	90.4	78.3	78.8	82.6	-
- BGE-M3	2.00	28.445.7	-	-	36.147.9	68.5	88.442.0	76.8 _{11.5}	77.7	84.6	-
Unstructured + OCR											
- BM25	31.6	36.8	62.9	46.5	62.7	64.3	92.8	85.9	83.9	87.2	65.5
- BGE-M3	31.440.2	25.7 <mark>411.1</mark>	60.142.8	70.8 ^{+24.3}	50.5 ^{12.2}	73.2	90.2 _{12.6}	83.642.3	84.9	91.1 _{13.9}	66.1
Unstructured + Captioning											I
- BM25	40.1	38.4	70.0	35.4	61.5	60.9	88.0	84.7	82.7	89.2	65.1
- BGE-M3	35.74.4	32.945.4	71.9 ^{1.9}	69.1	43.8417.7	73.1112.2	88.8	83.341.4	80.4	91.3 _{12.1}	67.0 <u><u>1.9</u></u>
Contrastive VLMs											
Jina-CLIP	25.4	11.9	35.5	20.2	3.3	3.8	15.2	19.7	21.4	20.8	17.7
Nomic-vision	17.1	10.7	30.1	16.3	2.7	1.1	12.9	10.9	11.4	15.7	12.9
SigLIP (Vanilla)	43.2	30.3	64.1	58.1	26.2	18.7	62.5	65.7	66.1	79.1	51.4
Ours											
SigLIP (Vanilla)	43.2	30.3	64.1	58.1	26.2	18.7	62.5	65.7	66.1	79.1	51.4
BiSigLIP (+fine-tuning)	58.5+15.3	32.9+2.6	70.5	62.714.6	30.5	26.5 17.8	74.3	73.7 18.0	74.2	82.3	58.617.2
BiPali (+LLM)	56.54-2.0	30.04-2.9	67.44-3.1	76.9114.2	33.412.9	43.7 17.2	71.24-3.1	61.94-11.7	73.84-0.4	73.64-8.8	58.8 _{10.2}
ColPali (+Late Inter.)	79.1 ^{+22.6}	54.4 124.5	81.8 14.4	83.9 17.0	65.8 ^{+32.4}	73.2 129.5	96.2 125.0	91.0 ^{+29.1}	92.7 18.9	94.4 120.8	81.3 ^{22.5}

Table 2: Comprehensive evaluation of baseline models and our proposed method on *ViDoRe*. Results are presented using nDCG@5 metrics, and illustrate the impact of different components. Text-only metrics are not computed for benchmarks with only visual elements.

Take away message:

- ColPali **simplifies and enhances document retrieval** by leveraging Vision-Language Models (VLMs) t, eliminating the need for text extraction.
- ColPali also shows a **way to MLLM application into different specific domains**. This could go beyond retrieval to **knowledge extraction**, **Forensic/Financial/Cultural industry document analysis**, and **Al-powered research assistants**.

Yijia Zheng

Yassin Mohamadi

Diga: Guided Diffusion Model for Graph Recovery in Anti-Money Laundering



Noise Prediction

DATASETS

Dataset	WB-S	WB-M	WB-L	Elliptic	ОТС	Alpha
# Nodes	58,409	233,638	467,277	203,769	5,858	3,754
# Edges	101,910	523,301	904,256	234,355	35,592	24,186
# Node Features	22	22	22	166	4	4
# Anomalies	478	1,830	3,676	4,545	178	102
Anomaly Ratio (%)	0.81	0.78	0.79	2.23	3.03	2.72
Avg. Degree	1.74	2.24	1.93	1.15	6.06	6.44

Results

Dataset	WeBank-small		WeBank-medium		WeBank-large		Elliptic		OTC		Alpha	
Metrices	Pre@100	AUC(%)	Pre@100	AUC(%)	Pre@100	AUC(%)	Pre@100	AUC(%)	Pre@100	AUC(%)	Pre@100	AUC(%)
SVM	$0.162_{\pm 0.021}$	54.0 ±1.5	-	-	-	-	-	-	$0.519_{\pm 0.002}$	$68.3 \scriptstyle \pm 0.1$	0.477 ±0.001	63.8 ±0.3
XGBoost 2016	$0.281 \scriptstyle \pm 0.008$	$62.7_{\pm 0.2}$	0.432 ±0.011	$61.4_{\pm 0.5}$	-	-	$0.813_{\pm 0.001}$	$81.8 {\scriptstyle \pm 0.3}$	$0.678_{\pm 0.000}$	$85.9 \scriptstyle \pm 0.0$	0.606 ±0.000	$81.2 _{\pm 0.0}$
DeepFD 2018	$0.126 \scriptstyle \pm 0.037$	49.2 ±2.1	0.666 ±0.038	69.4 ±3.6	0.613 ±0.049	65.7 ±4.6	0.688 ±0.033	72.3 ±3.9	0.656 ±0.031	74.6 ±2.2	0.369 ±0.022	54.7 ±3.4
SemiADC 2021	$0.170 \scriptstyle \pm 0.017$	$56.3 \scriptstyle \pm 0.9$	0.717 ±0.020	76.9 ±1.2	-	-	-	-	0.545 ±0.051	75.2 ±1.8	0.522 ±0.029	74.1 ±0.9
OCGTL 2022	$0.267_{\pm 0.006}$	$60.3 \scriptstyle \pm 0.1$	0.802 ±0.009	77.7 ±0.1	0.751 ±0.012	73.4 ±0.4	0.912 ±0.007	<u>90.3</u> ±1.1	0.622 ±0.001	83.6 ±0.6	0.683 ±0.005	80.7 ±0.1
ComGA 2022	0.331 ±0.020	<u>71.9</u> ±0.6	0.804 ±0.016	80.8 ±2.0	-	-	0.855 ±0.029	85.2 ±3.6	0.701 ±0.005	87.0 ±2.4	$0.634_{\pm 0.016}$	84.1 ±1.9
Diga (Ours)	0.383 ±0.004	74.3 ±0.3	0.900 ±0.010	87.6 ±0.2	0.949 ±0.008	82.5 ±0.4	0.981 ±0.005	95.2 ±0.7	0.730 ±0.010	89.1 ±0.3	0.724 ±0.019	87.4 ±0.5

Gonçalo Marcelino



Network visualization tools are becoming increasingly popular. How do users engage in the visual exploration of network data, which exploration strategies they employ, and how they prepare their data, define questions, and decide on visual mappings?



Study 1 - Researchers tracked users of the Vistorian logging 534 sessions to understand how it is being used.

Study 2 - Researchers collected qualitative during during a 6-week network exploration course by monitoring 36 participants. 50% without experience in network visualization.

- Missing Goals & Questions
 - X Choosing schemas and visualizations is difficult without specific goals and can lead to drill down fallacy.
 - ✓ Sketching, and examples improved results.
- Preconceived Ideas & Mental Images / Deciding on a Network Structure
 - Wrong preconceived about what a graph visualization is (e.g. exclusively a social network graph)
 - ✓ Guided process to construct schemas and fast sketching.



Lee, Doris Jung-Lin, et al. "Avoiding drill-down fallacies with vispilot: Assisted exploration of data subsets." *Proceedings of the 24th International Conference on Intelligent User Interfaces.* 2019.

- Choosing The Right Level of Abstraction
 - Choosing the wrong level of abstraction can lead to cluttered graphs.
 - ✓ Transformation and aggregation strategies.
- Interpreting Visual Patterns in Visualization
 - Understanding patterns in data is complex specially when interaction can lead to changes of visual patterns on-the-fly.
 - ✓ Support multiple coordinated views, support for examples.

- Establish Trust in a Network
 Visualization
 - Unfamiliar visualization (e.g. adjacency matrix) and misunderstanding provenance.
 - Showing examples of use in credible sources (e.g.) journalism, explaining algorithms, explaining anti-patterns.



Fatemeh Gholamzadeh