

[← Back to Articles](#)

# ColPali: Efficient Document Retrieval with Vision Language Models 🧐

Community Article

Published July 5, 2024

▲ Upvote 186

[manu](#)

Manuel Faysse

Using Vision LLMs + late interaction to improve document retrieval (RAG, search engines, etc.), solely using the image representation of document pages (paper)!

## Context

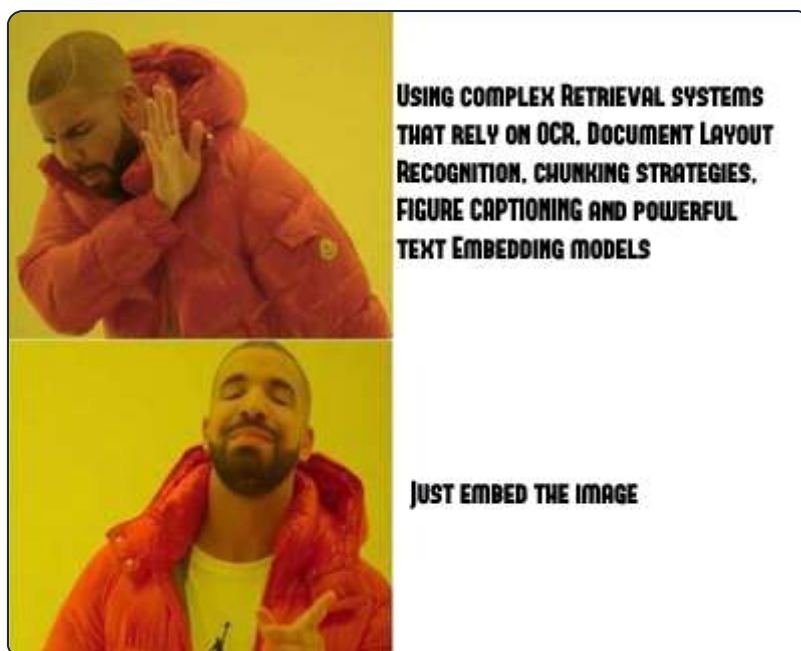
To improve the query answering capabilities of LLMs, it is often best to first search for information online or in external document sets (PDFs), before letting a LLM synthesize a grounded response (RAG). In practice, these retrieval pipelines for PDF documents have a huge impact on performance but are non-trivial...

1. Run Optical Character Recognition (OCR) on scanned PDFs
2. Run Document Layout Detection models to segment pages into paragraphs, figures, titles
3. Reconstruct the structure and the reading order of the page
4. Optionally, use resource intensive specialized models to caption figures, images and tables in natural language
5. Use a chunking strategy to split or merge text passages in a coherent way

6. Use a strong neural embedding model (BGE M3) to map text chunks to a semantically meaningful vector space
7. Store the vector index to be used for future retrieval

Although tools exist to facilitate this pipeline (Unstructured, Surya), the whole indexing process can be slow, tends to propagate errors, and struggles to take into account the more visual elements of a page (tables, figures, images but also fonts, etc..).

🔗 Our concept ? Just embed the page image directly !



In practice, it's not as easy as we just made it sound ! Our method, ColPali is enabled by the latest advances in Vision Language Models, notably the PaliGemma model from the Google Zürich team, and leverages multi-vector retrieval through late interaction mechanisms as proposed in ColBERT by Omar Khattab.

Let's break it down, with more technical details !

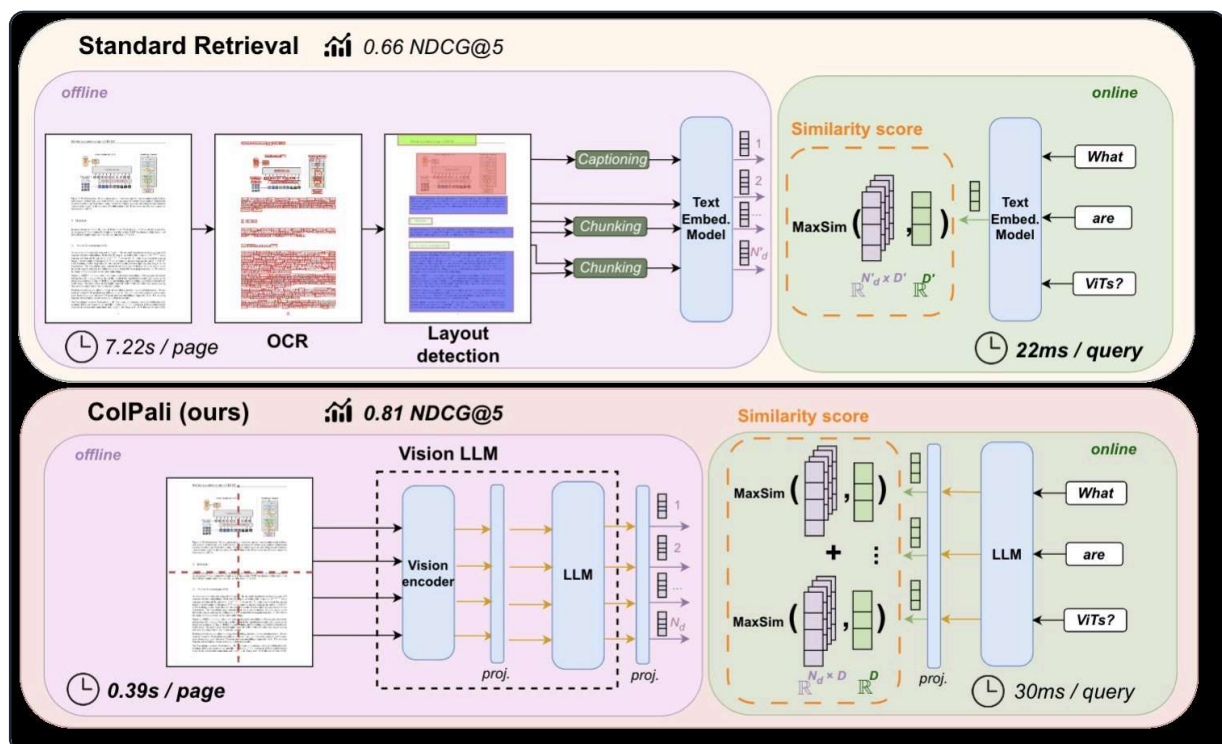
## 🔗 Model Architecture

Many retrieval systems can be broken down into two parts.

- In the indexing phase, all the documents from the corpus are indexed in an offline fashion.
- In the querying phase, a user query is matched with a low latency to the pre-computed document index.

Important requirements for efficient retrieval systems are thus (R1) good retrieving performance, (R2) reasonable indexing speeds, (R3) low latency during querying.

During indexing, standard “bi-encoder” neural retrieval systems first parse documents to extract semantically coherent text passages, then map them to a dense vector space that aims to represent the text’s semantic meaning, and store the resulting “embeddings”. During querying, the query is converted into its dense vector representation and the document passage vectors with the biggest cosine similarity can be retrieved with a low latency.



Our method ColPali is a bit different !

During indexing, we aim to strip away a lot of the complexity by using images (“screenshots”) of the document pages directly.

A Vision LLM (PaliGemma-3B) encodes the image by splitting it into a series of patches, which are fed to a vision transformer (SigLIP-So400m). These patch embeddings are linearly projected and inputted as “soft” tokens to a language model (Gemma 2B), in order to obtain high-quality contextualized patch embeddings in the language model space, which we then project to a lower dimension ( $D=128$ ) for more efficient storage. We thus construct and store a multi-vector document representation for each page image.

During runtime querying, a user query is embedded by the language model, to obtain token embeddings. We are able to run a ColBERT-style “late interaction” (LI) operation to efficiently match query tokens to document patches. To compute a  $LI(query, document)$  score, for each term in the query, we search for the document patch that has the most similar ColPali representation. We then sum the scores of the most similar patches for all terms of the query, to obtain the final query-document score. Intuitively, this late-interaction operation allows for a rich interaction between all terms of the query and document patches, all the while benefiting from the fast matching and offline computation offloading that more standard (bi-encoder) embedding models enable.

With ColPali, we thus benefit from fast indexing speeds (R2) without [↗](#) significantly impacting querying latencies (R3) ! But what about performance (R1)?

## [↗](#) ViDoRe

Although awesome benchmarks exist to evaluate text embedding models, we find that in many practical use cases, the prior document ingestion pipeline matters much more than the embedding model itself ! While documents often rely on visual elements to more efficiently convey information to human readers, text-only systems barely tap into these visual cues. To our knowledge, no benchmark evaluates document retrieval methods by considering both textual and visual document features like a human would.

To this end, we introduce ViDoRe, the Visual Document Retrieval Benchmark, to assess retrievers on their capacity to retrieve visually rich information in docs, with tasks spanning various topics, modalities (figures, tables, text), and languages !

Dataset	# Queries	Domain
<b>Academic Tasks</b>		
DocVQA (eng)	500 (500)	Industrial
InfoVQA (eng)	500 (500)	Infographics
TAT-DQA (eng)	1600 (1600)	Varied Modalities
arXivQA (eng)	500 (500)	Scientific Figures
TabFQuAD (fra)	210 (210)	Tables
<b>Practical Tasks</b>		
Energy (eng)	100 (1000)	Scientific
Government (eng)	100 (1000)	Administrative
Healthcare (eng)	100 (1000)	Medical
AI (eng)	100 (1000)	Scientific
Shift Project (fra)	100 (1000)	Environment

Table 1: *ViDoRe* comprehensively evaluates multimodal retrieval methods. The size of the document corpus is indicated in parentheses.

ViDoRe is linked to a HF leaderboard <https://huggingface.co/spaces/vidore/vidore-leaderboard> and we hope to see many models trying out this new "Retrieving in Vision Space" paradigm !

## 🔗 Results

### 🔗 Training details

We initialize the Vision Language Model backbone using pretrained weights from PaliGemma and randomly initialize the final projection layer. To facilitate training, we add low-rank adapters to the language model attention weights, as well as the



linear projection layers. Our training dataset is composed of (query, document image) pairs that we sourced from two main streams. On one hand, we repurposed Visual Question Answering datasets and used the original question as our query, and the associated image as the gold label. To increase the coverage and diversity of the training set, we also collected tens of thousands of permissively licensed PDF documents covering a broad range of topics, and synthetically create relevant queries using the powerful Claude Sonnet Vision model. In total, we gather around 100k pairs, and finetune our model with an in-batch contrastive loss, by attempting to maximize the difference between the matching score of the correct (page, query) pair, and the score of the incorrect pairs.

## ColPali results

On ViDoRe, ColPali outperforms all other evaluated systems, including baselines where a very strong proprietary Vision model (Claude Sonnet) is used to caption all visual elements !

	ArxivQ	DocQ	InfoQ	TabF	TATQ	Shift	AI	Energy	Gov.	Health.	Avg.
<b>Unstructured</b> <small>Text only</small>											
- BM25	-	34.1	-	-	44.0	59.6	90.4	78.3	78.8	82.6	-
- BGE-M3	-	28.4 <sub>↓5.7</sub>	-	-	36.1 <sub>↓7.9</sub>	68.5 <sub>↑8.9</sub>	88.4 <sub>↓2.0</sub>	76.8 <sub>↓1.5</sub>	77.7 <sub>↓1.1</sub>	84.6 <sub>↑2.0</sub>	-
<b>Unstructured</b> <small>+ OCR</small>											
- 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	31.4 <sub>↓0.2</sub>	25.7 <sub>↓11.1</sub>	60.1 <sub>↓2.8</sub>	70.8 <sub>↑24.3</sub>	50.5 <sub>↓12.2</sub>	<b>73.2</b> <sub>↑8.9</sub>	90.2 <sub>↓2.6</sub>	83.6 <sub>↓2.3</sub>	84.9 <sub>↑1.0</sub>	91.1 <sub>↑3.9</sub>	66.1 <sub>↑0.6</sub>
<b>Unstructured</b> <small>+ Captioning</small>											
- 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	35.7 <sub>↓4.4</sub>	32.9 <sub>↓5.4</sub>	71.9 <sub>↑1.9</sub>	69.1 <sub>↑33.7</sub>	43.8 <sub>↓17.7</sub>	73.1 <sub>↑12.2</sub>	88.8 <sub>↑0.8</sub>	83.3 <sub>↓1.4</sub>	80.4 <sub>↓2.3</sub>	91.3 <sub>↑2.1</sub>	67.0 <sub>↑1.9</sub>
<b>Contrastive VLMs</b>											
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
<b>Ours</b>											
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 <sub>↑15.3</sub>	32.9 <sub>↑2.6</sub>	70.5 <sub>↑6.4</sub>	62.7 <sub>↑4.6</sub>	30.5 <sub>↑4.3</sub>	26.5 <sub>↑7.8</sub>	74.3 <sub>↑11.8</sub>	73.7 <sub>↑8.0</sub>	74.2 <sub>↑8.1</sub>	82.3 <sub>↑3.2</sub>	58.6 <sub>↑7.2</sub>
BiPali (+LLM)	56.5 <sub>↓2.0</sub>	30.0 <sub>↓2.9</sub>	67.4 <sub>↓3.1</sub>	76.9 <sub>↑14.2</sub>	33.4 <sub>↑2.9</sub>	43.7 <sub>↑17.2</sub>	71.2 <sub>↓3.1</sub>	61.9 <sub>↓11.7</sub>	73.8 <sub>↓0.4</sub>	73.6 <sub>↓8.8</sub>	58.8 <sub>↑0.2</sub>
<b>ColPali</b> (+Late Inter.)	<b>79.1</b> <sub>↑22.6</sub>	<b>54.4</b> <sub>↑24.5</sub>	<b>81.8</b> <sub>↑14.4</sub>	<b>83.9</b> <sub>↑7.0</sub>	<b>65.8</b> <sub>↑32.4</sub>	<b>73.2</b> <sub>↑29.5</sub>	<b>96.2</b> <sub>↑25.0</sub>	<b>91.0</b> <sub>↑29.1</sub>	<b>92.7</b> <sub>↑18.9</sub>	<b>94.4</b> <sub>↑20.8</sub>	<b>81.3</b> <sub>↑22.5</sub>

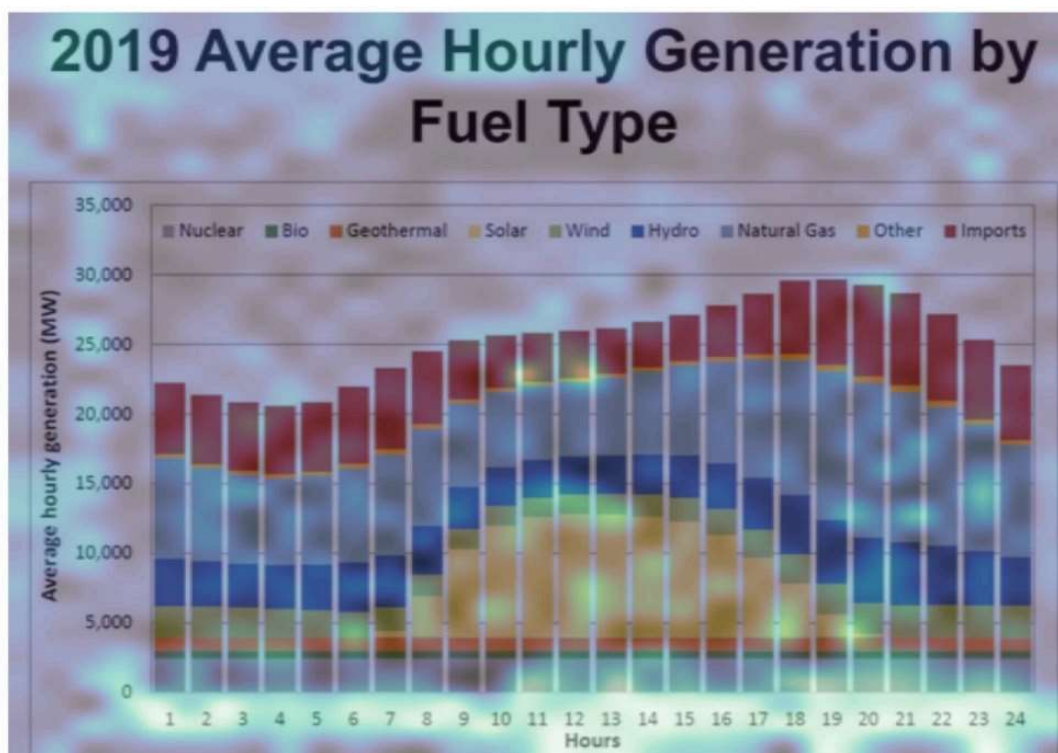
Table 2: **Comprehensive evaluation of baseline models and our proposed method on ViDoRe.** Results are presented using NDCG@5 metrics, and illustrate the impact of different components. Text-only metrics are not computed for benchmarks with only visual elements.

The difference is particularly stark on the more visually complex benchmark tasks, such as InfographicVQA, ArxivQA, and TabFQuAD representing respectively infographics, figures, and tables. However, text centric documents are also better

retrieved by the ColPali model across all evaluated domains and languages, making our approach the overall best performing document-retrieval model on ViDoRe !

### 🔗 Interpretability

Beyond speed and performance, another interesting feature of ColPali, is that it enables visualizing which patches of a document stand out w.r.t. a given query. Here the term <hour> matches patches containing words like "hourly" but also the x-axis representing time, showcasing good chart comprehension !



**Query:** "Which hour of the day had the highest overall electricity generation in 2019?"

Figure 1: For each term in a user query, **ColPali** identifies the most relevant document image patches (highlighted zones) and computes a query-to-page matching score. We can then swiftly retrieve the most relevant documents from a large pre-indexed corpus.

### 🔗 Conclusion

This blogpost is already long enough but good news, tons more resources, informations and ablations exist, and will keep on coming !

📄 The paper: <https://arxiv.org/abs/2407.01449>

📖 The benchmark: <https://huggingface.co/vidore>

🧐 The model: <https://huggingface.co/vidore/colpali>

💻 The benchmark code: <https://github.com/illuin-tech/vidore-benchmark>

💻 The training code: <https://github.com/ManuelFay/colpali>

✂ X of the first authors: @ManuelFaysse, @sibille\_hugues, @tonywu\_71

## 🔗 Citation

```
@misc{faysse2024colpaliefficientdocumentretrieval,  
  title={ColPali: Efficient Document Retrieval with Vision Language Models},  
  author={Manuel Faysse and Hugues Sibille and Tony Wu and Bilel Omrani},  
  year={2024},  
  eprint={2407.01449},  
  archivePrefix={arXiv},  
  primaryClass={cs.LG},  
  url={https://arxiv.org/abs/2407.01449},  
}
```

## 🔗 Acknowledgments

This work primarily stems from an academic-industry partnership between CentraleSupélec and Illuin Technology, with involvement from actors at Equall.ai and ETH Zürich. It benefits from a compute grant from CINES ADAstra (Grant 2024-AD011015443). Joint work by [Manuel Faysse](#), [Hugues Sibille](#), [Tony Wu](#), [Bilel Omrani](#), [Gautier Viaud](#), [Céline Hudelot](#), [Pierre Colombo](#).



 System theme

Company

TOS

Privacy

About

Jobs

Website

Models

Datasets

Spaces

Pricing

Docs

