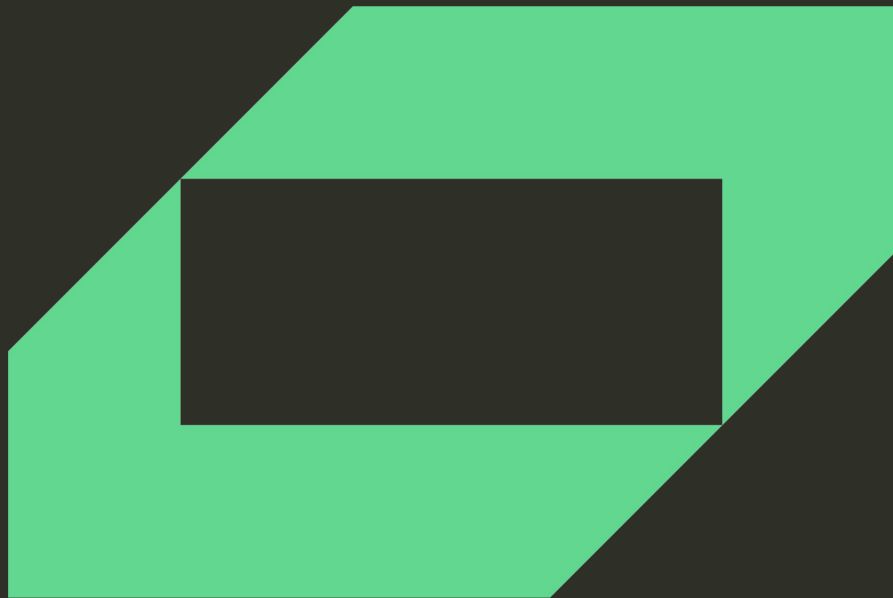# *What You See Is What You Search: Vision Language Models for PDF Retrieval*

Jo Kristian Bergum, Vespa.ai

30.09.2024

Vespa.ai

# Agenda

This deck:

Vespa.ai

# Indexing & Searching in PDFs

Complex format with visual infographics, tables, charts and images.

Searching PDFs was Popularized by RAG (Retrieval-Augmented Generation) - Chat with your PDFs

Extraction pipeline to map to text:

- Chunking?
- Text Embedding?

# Text extraction is *messy*

- Embedded images (e.g. infographics)
- Order
- Layout & Structure
- Tables
- Font size

```
We manage biodiversity risks and mitigate impacts through the use of the Mitigation Hierarchy, a decision-ma
framework involving prioritized steps to mitigate adverse biodiversity impacts: avoid, minimize, restore and
Our efforts are designed to reduce impact on biodiversity and contribute to its restoration.Our policies req
to governance, strategy, management and performance related to nature. NATURESustainability
23-1207HAB ITATS CONSE RVED, PROTECTED OR REST ORED
2 Estimated as the percentage of lease areas overlapping with designated
protected areas using the World Database on Protected Areas.OVER
540 ,000
CUMULATIVE ACRES1

1 Cumulative with varying conservation pr oject start dates
as early as 2009.on company-owned lands
and operated assets. UNCONVENTIONAL
Bakken | Eagle Ford
Montney | Permian BasinBBL/BOE EUR1
BBL/BOE2
CONVENTIONAL/OFFSHORE
Alaska | APLNG | Ekofisk
Surmont | Teesside FRESH WATER
CONSUMPTION
INTENSITY
1 Calculated using Enverus data for the average volume of fresh water (bbl) divided by the average estimated
ecovery (EUR, BOE) as of April 1, 2024. Intensity value may change as EUR data
is updated. EUR — estimated ultimate recovery.  2 Calculated using the average volume of fresh water (BBL)
the average annual production (BOE).
24-0976As of Dec. 31, 20230.03%OF LEASE AREAS OVERLAP
WITH PROTECTED AREAS2
12PROTECTED AREAS WITHIN
3 MILES (5 KM) OF FIVE ASSETS
APLNG | Bakken | Permian Basin
Montney | Teesside0.06
0.03
```
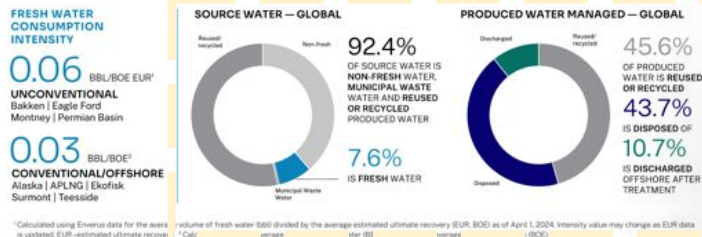
Vespa.ai

# Why extract to text in the first place?

*Convert to something that makes it searchable*

Convert to text modality:

=> Text search
- Keyword search (BM25)
- Embedding retrieval with text embedding models (e.g. Jina)
- Splade, Bert, ColBERT, OpenAi bla bla bla, all text-only



Vespa.ai

# Agenda

# What is a Vision Language Model (VLM) ?



*Vision language models are broadly defined as multimodal models that can learn from images and text.*

*They are a type of generative models that take image and text inputs, and generate text output*

*From https://huggingface.co/blog/vlms*

Vespa.ai

# The year of multimodal AI (ok, VLMs)



A not very accurate summary of VLM landscape

MIT/Apache 2.0
Custom License/Open weights
API based

Accuracy across many different tasks (Not just cat photos)

OpenAI GPT-4o
Anthropic Claude
Google Gemini Flash
Qwen2-VL-7B-Instruct
OpenAI GPT-4v
Mistral Pixtral
MS Phi-3.5-vision
Google PaliGemma
* Lots of other vlms

The early days, VLMs for searching

You are here

2023    2024    2025

New VLMs last week

Molmo
Llama 3.2

Vespa.ai

# PaliGemma (a VLM)

PaliGemma – Google's *"Cutting-Edge"* Open Weight VLM

-OCR capabilities
- Question answering + more



PaliGemma

Image Input → SigLIP Image Encoder → Linear Projection → Concatenated Tokens ← Text Input → Gemma → Text Output

https://huggingface.co/blog/paligemma



Input Text

What is the percentage of produced water that is reused?

Text Output

45.6%

https://huggingface.co/spaces/big-vision/paligemma-hf

Vespa.ai

# "*What's in the image*"



Vespa.ai

# Agenda

Vespa.ai

# ColPali: Use VLM document understanding capabilities for document retrieval

IR benchmarks measures clean pre-processed texts

**In practical industrial settings we don't have the luxury of pre-processed datasets.**

Performance bottleneck usually not in embedding model performance or chunking, but in the data ingestion pipeline

**ColPali**: Efficient Document Retrieval with Vision Language Models

Manuel Faysse[*,1,3]   Hugues Sibille[*1,4]   Tony Wu[*1]   Bilel Omrani[1]
Gautier Viaud[1]   Céline Hudelot[3]   Pierre Colombo[2,3]
[1]Illuin Technology   [2]Equall.ai
[3]CentraleSupélec, Paris-Saclay   [4]ETH Zürich
manuel.faysse@centralesupelec.fr

## Abstract

Documents are visually rich structures that convey information through text, as well as tables, figures, page layouts, or fonts. While modern document retrieval systems exhibit strong performance on query-to-text matching, they struggle to exploit visual cues efficiently, hindering their performance on practical document retrieval applications such as Retrieval Augmented Generation. To benchmark current systems on visually rich document retrieval, we introduce the Visual Document Retrieval Benchmark *ViDoRe*, composed of various page-level retrieving tasks spanning multiple domains, languages, and settings. The inherent shortcomings of modern systems motivate the introduction of a new retrieval model architecture, *ColPali*, which leverages the document understanding capabilities of recent Vision Language Models to produce high-quality contextualized embeddings solely from images of document pages. Combined with a late interaction matching mechanism, *ColPali* largely outperforms modern document retrieval pipelines while being drastically faster and end-to-end trainable. We release all project artifacts at https://huggingface.co/vidore.

2019 Average Hourly Generation by Fuel Type

**Query:** "Which hour of the day had the highest overall eletricity 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.

index a standard PDF document, many steps are required. First, PDF parsers or Optical Character Recognition (OCR) systems are used to extract words from the pages. Document layout detection models can then be run to segment paragraphs,

Vespa.ai

Figure 2: *ColPali* simplifies document retrieval w.r.t. standard retrieval methods while achieving stronger performances with better latencies. Latencies and results are detailed in section 5 and subsection B.5.

Vespa.ai

# ColPali (Gemma)

ColPali is short for *Contextualized Late Interaction over PaliGemma* and builds on two concepts:

- **Contextualized Embeddings from VLM**
  - **ColPali generates contextualized multi-vector embeddings directly from the screenshot of a page or text query using the VLM as the backbone**
- **Interaction between text query vectors and screenshot vectors at scoring time**

**A bi-encoder architecture, enables offline indexing - but with multi-vectors per page**



Vespa.ai

# ColPali offers

✓ **Match without text extraction**
Do not have to map complex formats
to the text domain

✓ **Engineering Simplicity**
Reduced document processling
pipeline complexity

✓ **Avoid OCR**
Use the VLM OCR capabilities

✓ **Better relevance**
Outcompete traditional extract
methods

✓ **No Layout Detection**
Layout is encoded by the VLM

✓ **A future direction for document retrieval**
Can you print the document, then you
can index it with the ColPali approach

# ColPali page level embeddings

VLM "see" the image as 32x32 patches = 1024 patches.

**An image is worth 1024 words**

Each patch is represented or projected into a 128-dimensional vector space (the latent space)

6 text tokens projected to the same space = 1030 vectors per page

Tokens from text prefix "Describe the image" <img>



Vespa.ai

# ColPali page level embeddings

One PDF - Multiple pages

One page represented as a tensor

Vespa tensor definition

*tensor<float>(patch{}, v[128])*

A map of vectors (patch is the key, the vector the value)

# ColPali text query embeddings

One 128-dimensional vector per user input query *token*

Fixed prepend and query expansion tokens

*tensor<float>(q{},v[128])*



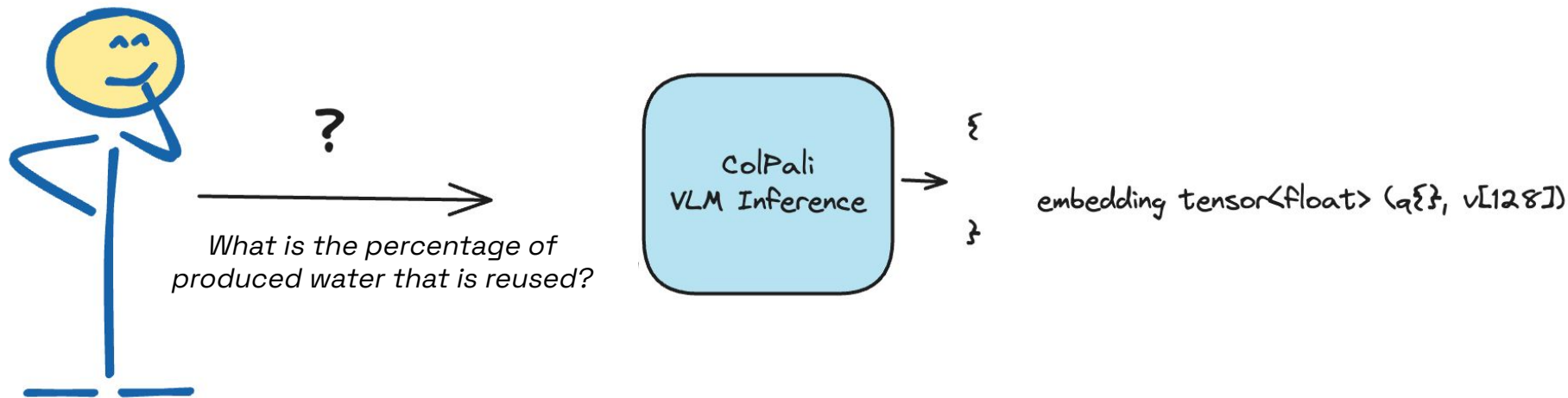*What is the percentage of produced water that is reused?*

ColPali VLM Inference

{
}

embedding tensor<float> (q{}, v[128])

# ColPali Relevance Scoring
# score(query, page):

*tensor<float>(q{}, v[128])*                    *tensor<float>(p{},v[128])*

**Sum of Maximum Similarities**

```
sum(
    reduce(
     sum(
      query * page , v
     ),
     max,
     p
    ),
   q
  )
```

*Contextualized Late-Interaction*
*over Embeddings from PaliGemma*

Vespa.ai

# (Sum) of MaxSim

Similarity matrix $|q| \times |p|$ - where similarity is the dot product

Find maximum patch similarity per query token. Toy example with 2-dimensional **v** and 4 image patches.

**Query token vectors**

water [1, 4]

max( [**38**, 12, 11, 9]) = **38**

recycle [8,2]

max([34, **36**, 28,12]) = **36**

**score(query,page) = 38 + 36 = 74**



[2, 9]   [4,2]

[3,2]   [1, 2]

Vespa.ai

# Score each page in the collection with (Sum of) MaxSim

Similar process as with "regular" search BM25 (Lucene, Vespa, Elasticsearch)

Sort the pages by score in descending order. Score is a proxy for relevance, higher is better.



Ranked Pages by MaxSim

id 2-1    MaxSim(page, query) = 72

id 1-3    MaxSim(page, query) = 69

id 1-2    MaxSim(page, query) = 65

id 1-1    MaxSim(page, query) = 64

Vespa.ai

# Learning Col vectors for retrieval (representation learning)

Train adapter weights which learns how to represent text and patches in same vector space with sim = MaxSim

VLM (e.g. PaliGemma)

Weights

Col Adapter

Weights

Relevance Training Data

text query

Relevant page image

Not relevant page image

N = 10-20K

Vespa.ai

# How does ColPali compare with traditional methods?

# The ViDoRe Benchmark

## Datasets ViDoRe

| Dataset | # Queries | Domain |
|---|---|---|
| **Academic Tasks** | | |
| 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 |
| **Practical Tasks** | | |
| 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.

|  | ArxivQ | DocQ | InfoQ | TabF | TATQ | Shift | AI | Energy | Gov. | Health. | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Unstructured** Text only | | | | | | | | | | | |
| - BM25 | - | 34.1 | - | - | 44.0 | 59.6 | 90.4 | 78.3 | 78.8 | 82.6 | - |
| - BGE-M3 | - | 28.4$_{\downarrow5.7}$ | - | - | 36.1$_{\downarrow7.9}$ | 68.5$_{\uparrow8.9}$ | 88.4$_{\downarrow2.0}$ | 76.8$_{\downarrow1.5}$ | 77.7$_{\downarrow1.1}$ | 84.6$_{\uparrow2.0}$ | - |
| **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 | 31.4$_{\downarrow0.2}$ | 25.7$_{\downarrow11.1}$ | 60.1$_{\downarrow2.8}$ | 70.8$_{\uparrow24.3}$ | 50.5$_{\downarrow12.2}$ | **73.2**$_{\uparrow8.9}$ | 90.2$_{\downarrow2.6}$ | 83.6$_{\downarrow2.3}$ | 84.9$_{\uparrow1.0}$ | 91.1$_{\uparrow3.9}$ | 66.1$_{\uparrow0.6}$ |
| **Unstructured** + Captioning | | | | | | | | | | | |
| - 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$_{\downarrow4.4}$ | 32.9$_{\downarrow5.4}$ | 71.9$_{\uparrow1.9}$ | 69.1$_{\uparrow33.7}$ | 43.8$_{\downarrow17.7}$ | 73.1$_{\uparrow12.2}$ | 88.8$_{\uparrow0.8}$ | 83.3$_{\downarrow1.4}$ | 80.4$_{\downarrow2.3}$ | 91.3$_{\uparrow2.1}$ | 67.0$_{\uparrow1.9}$ |
| **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$_{\uparrow15.3}$ | 32.9$_{\uparrow2.6}$ | 70.5$_{\uparrow6.4}$ | 62.7$_{\uparrow4.6}$ | 30.5$_{\uparrow4.3}$ | 26.5$_{\uparrow7.8}$ | 74.3$_{\uparrow11.8}$ | 73.7$_{\uparrow8.0}$ | 74.2$_{\uparrow8.1}$ | 82.3$_{\uparrow3.2}$ | 58.6$_{\uparrow7.2}$ |
| BiPali (+LLM) | 56.5$_{\downarrow2.0}$ | 30.0$_{\downarrow2.9}$ | 67.4$_{\downarrow3.1}$ | 76.9$_{\uparrow14.2}$ | 33.4$_{\uparrow2.9}$ | 43.7$_{\uparrow17.2}$ | 71.2$_{\downarrow3.1}$ | 61.9$_{\downarrow11.7}$ | 73.8$_{\downarrow0.4}$ | 73.6$_{\downarrow8.8}$ | 58.8$_{\uparrow0.2}$ |
| *ColPali* (+Late Inter.) | **79.1**$_{\uparrow22.6}$ | **54.4**$_{\uparrow24.5}$ | **81.8**$_{\uparrow14.4}$ | **83.9**$_{\uparrow7.0}$ | **65.8**$_{\uparrow32.4}$ | **73.2**$_{\uparrow29.5}$ | **96.2**$_{\uparrow25.0}$ | **91.0**$_{\uparrow29.1}$ | **92.7**$_{\uparrow18.9}$ | **94.4**$_{\uparrow20.8}$ | **81.3**$_{\uparrow22.5}$ |

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.
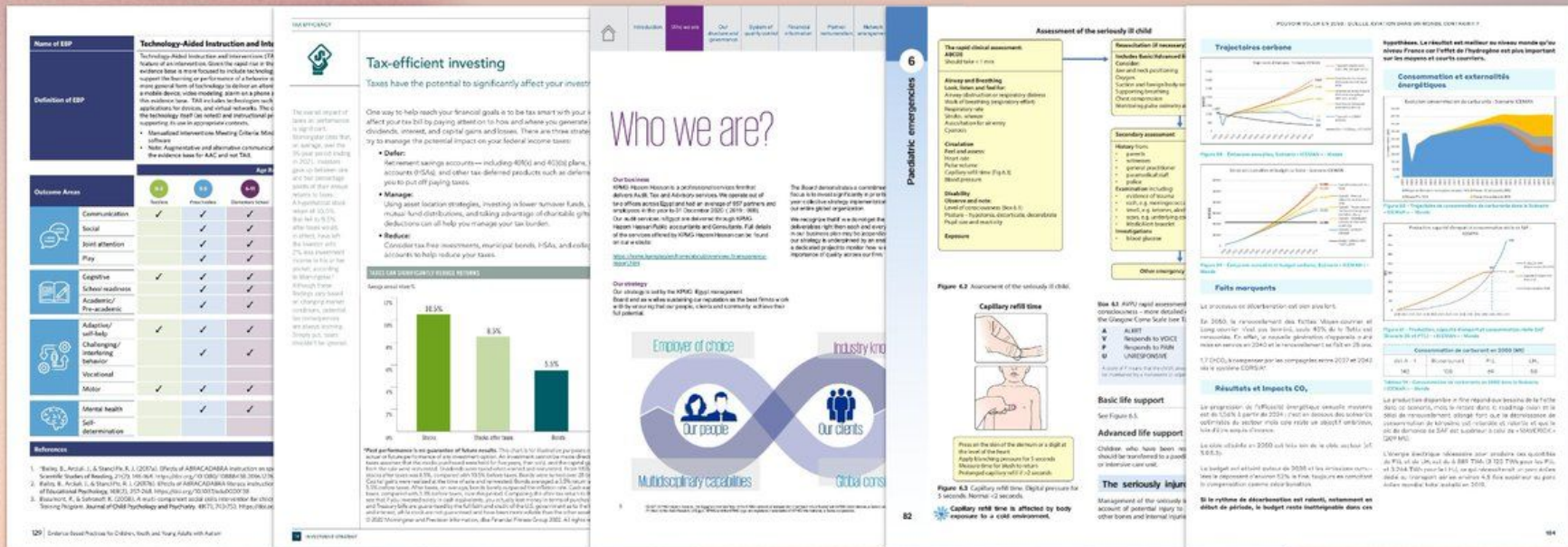
nDCG@5  avg 0.66 => 0.813

Vespa.ai

# Samples from the ViDoRe benchmark dataset

# Samples of PDF pages from the ViDoRe benchmark

Vespa.ai

# Agenda

Vespa.ai

# "But, Does It Scale?"

## "That's a lot of vectors!"

Reply Guys
20.08.2024

Vespa.ai

# Scaling ColPali MaxSim

Compute scales with number of dot products $D \times [q \times p \times v]$

Storage scales with $D \times [p \times v]$

Note: MaxSim only, inference of model is scaled independently

- **q** is the number of query tokens, including the mask and instruction tokens
- **p** is the number of image patches
- **v** is the vector dimensionality
- **D** is the number of pages scored by MaxSim

Vespa.ai

# Scaling ColPali MaxSim

- Reduce # patch vectors (reduces |p|) - clustering, remove redundant

- Reduce precision of v from float to bit - 32x

- Replace float dot products with inverted hamming (4x reduction in compute)

- **q** is the number of query tokens, including the mask and instruction tokens
- **p** is the number of image patches
- **v** is the vector dimensionality
- **D** is the number of pages scored by MaxSim

Vespa.ai

# ColPali vectors

Binary Quantization ( > 0 )
Normalized vectors (unit length). Values close to 0
Replace float dot product with inverted hamming distance (correlates)



MaxSim ranking. Latency (single-threaded) 1000x20X1030 128D hamming (>20M 128-bit hamming distances)

— bits-bits-hamming
— float-unpacked-bits-dotproduct

Vespa.ai

# Scaling ColPali MaxSim + Retrieval

Phased retrieval & ranking

Search for close neighbors of all k query token vectors - compatible with HNSW indexing using hamming distance metric

Step 2 - re-rank using MaxSim

nDCG@5 DocVQA

| | |
|---|---|
| float-float | 52.4 |
| binary-binary (hamming) | 49.5 |
| binary-binary (hamming) + float-float re-ranking | 51.6 |

Vespa.ai

# Scaling ColPali TLDR;

- Reduce precision (float to bit)
- Binary quantization (BQ)
- Hamming instead of dot product
- Multi-vector HNSW indexing
- Phased retrieval & ranking

https://blog.vespa.ai/scaling-colpali-to-billions/



## Scaling ColPali to billions of PDFs with Vespa

This blog post deep dives into scaling "ColPali: Efficient Document Retrieval with Vision Language Models" [1] to large collections of documents. We demonstrate how we can use a

Vespa.ai

# Agenda

Vespa.ai

# RAG with ColPali

ColPali is the first stage retriever

For the generative step you need a VLM for question answering based on the retrieved context

# Frontier VLMs

VLMs are better at question answering when using image data than text + OCR.



Single-Page Question

**Question:** I want to see a doctor in the campus hospital. After registering at the registration area, what is the next step?
**Answer:** Go to the medical department you registered at (i.e. internal medicine, surgical medicine, dental medicine)

**Solution:** Check the diagram about process for seeing a doctor at the hospital in Section *Health and Safety*, page 27. After registering at the registration area (Step 1), the second step is to go to the medical department you registered at.
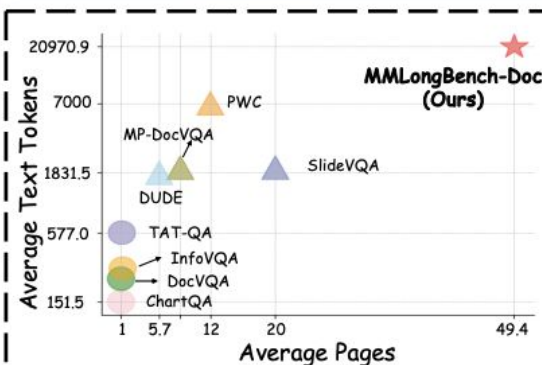**Evidence Sources:** Image

Cross-Page Question

**Question:** I'm at location "J" shown in the campus map. Tell me the nearest coffee shop.
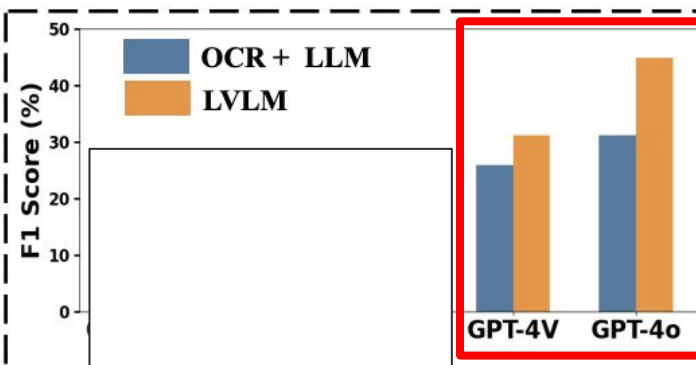**Answer:** Ten Years After Café.

**Solution:** According to the campus map in Page 34, the location "J" is the security Department near the New Qinghuaxuetang. Check the info table of on-campus coffee shop in Page 16, we find that Ten Years After Café is the nearest one to location "J".
**Evidence Sources:** Table & Image

**(b) Dataset Statistics**

**(c) Performance on MMLongBench-Doc**

From **MMLongBench-Doc: Benchmarking Long-context Document Understanding with Visualizations**

Vespa.ai

# Agenda

# "Does it work for my data?"

Vespa.ai

# Fine-tuning ColPali

ColPali **retrieval adapter** mostly trained on synthetic data created by previous generation VLMs using a prompting technique. Generate questions about a particular page.

New and improved VLM backbones gives:

- Better training data (generated on your domain and language)
- Better VLM backbone

**ColPali fine-tuning Query Generator**

ColPali is a very exciting new approach to multimodal document retrieval which aims to replace existing document retrievers which often rely on an OCR step with an end-to-end multimodal approach.

To train or fine-tune a ColPali model, we need a dataset of image-text pairs which represent the document images and the relevant text queries which those documents should match. To make the ColPali models work even better we might want a dataset of query/image document pairs related to our domain or task.

One way in which we might go about generating such a dataset is to use an VLM to generate synthetic queries for us. This space uses the Qwen/Qwen2-VL-7B-Instruct VLM model to generate queries for a document, based on an input document image.

**Note** there is a lot of scope for improving to prompts and the quality of the generated queries! If you have any suggestions for improvements please open a Discussion!

This blog post gives an overview of how you can use this kind of approach to generate a full dataset for fine-tuning ColPali models.

If you want to convert a PDF(s) to a dataset of page images you can try out the PDFs to Page Images Converter Space.

Vespa.ai

# "Can't you use GPT-4o for this?"

Reply Guys
27.09.2024

Vespa.ai

ColPali is a promising direction

Will see new checkpoints based on new VLMs

Trained on more data

Embedding providers..



Manuel Faysse
@ManuelFaysse

🚨 New model alert: ColQwen2 !
It's ColPali, but with a Qwen2-VL backbone, making it the best visual retriever to date, topping the Vidore Leaderboard with a significant +5.1 nDCG@5 w.r.t. colpali-v1.1 trained on the same data ! 🚀 (1/N)

⊕ vidore
**/colqwen2-v0.1**

😊 huggingface.co
vidore/colqwen2-v0.1 · Hugging Face

From huggingface.co

3:24 pm · 27 Sep 2024 · **105.4K** Views

Vespa.ai

# QA

## Resources

Jo Kristian Bergum ✓ @jobergum · Aug 15
I don't understand why my timeline is not all about ColPali for RAG over complex document formats?

💬 24      🔁 11      ❤️ 181      📊 48K

Vespa.ai

# Even more resources

https://pyvespa.readthedocs.io/en/latest/examples/pdf-retrieval-with-ColQwen2-vlm_Vespa-cloud.html

https://pyvespa.readthedocs.io/en/latest/examples/simplified-retrieval-with-colpali-vlm_Vespa-cloud.html

https://pyvespa.readthedocs.io/en/latest/examples/colpali-benchmark-vqa-vlm_Vespa-cloud.html

https://pyvespa.readthedocs.io/en/latest/examples/colpali-document-retrieval-vision-language-models-cloud.html



```python
async with app.asyncio(connections=1, total_timeout=120) as session:
    for idx, query in enumerate(queries):
        query_embedding = {k: v.tolist() for k, v in enumerate(qs[idx])}
        response: VespaQueryResponse = await session.query(
            yql="select title,url,image,page_number from pdf_page where userInput(@userQuer
            ranking="default",
            userQuery=query,
            timeout=120,
            hits=3,
            body={
                "input.query(qt)": query_embedding,
                "presentation.timing": True
            },
        )
        assert response.is_successful()
        display_query_results(query, response)
```

Query text: 'Percentage of non-fresh water as source?', query time 0.07s, count=133, top results:

**PDF Result 1**

**Title:** ConocoPhillips Managing Climate Related Risks, page 45 with score 90.48

Vespa

Follow us!

Vespa.ai
@vespaengine
linkedin.com/company/vespa-ai/