

Unified Multimodal Interleaved Document Representation for Retrieval

Anonymous ACL submission

Abstract

Information Retrieval (IR) methods aim to identify documents relevant to a query, which have been widely applied in various natural language tasks. However, existing approaches typically consider only the textual content within documents, overlooking the fact that documents can contain multiple modalities, including images and tables. Also, they often segment each long document into multiple discrete passages for embedding, which prevents them from capturing the overall document context and interactions between paragraphs. To address these two challenges, we propose a method that holistically embeds documents interleaved with multiple modalities by leveraging the capability of recent vision-language models that enable the processing and integration of text, images, and tables into a unified format and representation. Moreover, to mitigate the information loss from segmenting documents into passages, instead of representing and retrieving passages individually, we further merge the representations of segmented passages into one single document representation, while we additionally introduce a reranking strategy to decouple and identify the relevant passage within the document if necessary. Then, through extensive experiments on diverse IR scenarios considering both the textual and multimodal queries, we show that our approach substantially outperforms relevant baselines, thanks to the consideration of the multimodal information within documents.

1 Introduction

Information Retrieval (IR) is the task of fetching relevant documents from a large corpus in response to a query, which plays a critical role in various real-world applications including web search engines and question-answering systems (Shah et al., 2019; Lewis et al., 2020; Guu et al., 2020). Over the years, IR methods have evolved significantly, broadly categorized into sparse and dense retrieval paradigms. Specifically, sparse retrieval methods (Robertson

et al., 1994; Jones, 2004) focus on lexical overlap between queries and documents; meanwhile, dense retrieval methods (Karpukhin et al., 2020; Xiong et al., 2021) utilize neural embeddings to represent queries and documents in a continuous vector space. Note that, recently, dense retrieval methods have gained more popularity over sparse methods due to their capability to capture semantic nuances and context beyond simple keyword matching.

Despite their successes, existing (dense) retrieval methods face a couple of severe challenges. First, they primarily rely on the textual data for document embedding and retrieval, overlooking the fact that modern documents often contain multimodal content, including images and tables (beyond the plain text), which can carry information that may be essential for accurately understanding and retrieving the relevant documents (Li et al., 2024c). For instance, a diagram within a medical article can more effectively represent the structure of a molecule or the progression of a disease, offering more clarity that would be difficult to achieve with text alone, and omitting such multimodal content can lead to an incomplete understanding (and potentially inaccurate retrieval) of the documents. Also, the segmentation of long documents into discrete passages, which is commonly employed by existing retrieval models to handle the length limitation for embeddings (Karpukhin et al., 2020; Xiong et al., 2021), may prevent models from capturing the full context and the intricate relationships between different parts of the document, ultimately leading to suboptimal retrieval performance (Dong et al., 2024; Jiang et al., 2024b). Notably, concurrent to our work, while there has been recent work that screen captures the document and then embed its screenshots (to consider different modalities in a unified format) (Faysse et al., 2024; Ma et al., 2024), not only its content (such as paragraphs, images, and tables) can be fragmented into different sub-images, leading to the loss of contextual coher-

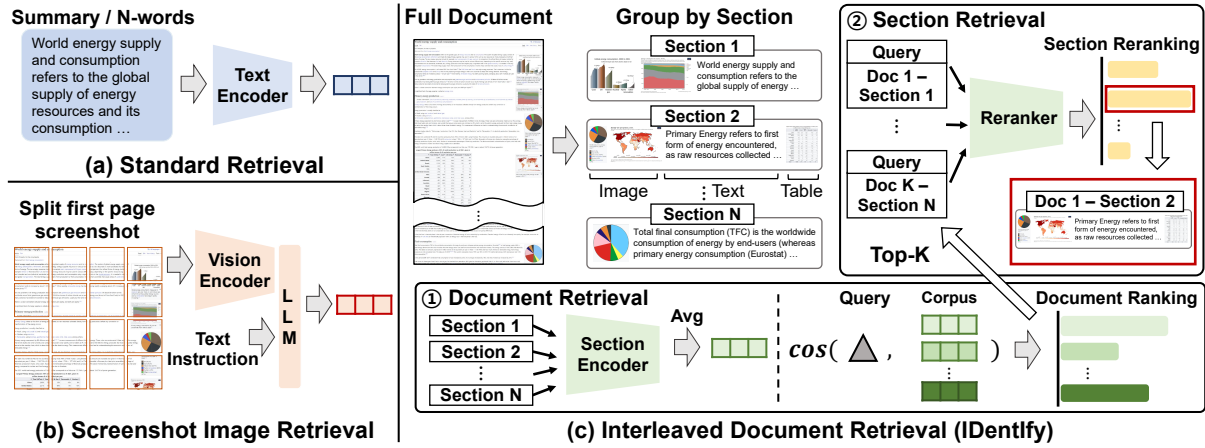


Figure 1: Comparison of different IR approaches. (a): Conventional methods use a small portion of the text within the document for its representation. (b): Recent methods use first-page screenshot images to represent the document. (c): Our approach leverages the full contextual information within documents interleaved with multiple modalities by considering them in their original format, and is further capable of pinpointing relevant sections for the query.

ence across the entire document, but also the visual representation of text may hinder the model’s ability to capture the semantic relationships present in the original textual data, and increasing image resolution raises concerns on memory requirements.

To tackle these challenges, we introduce a novel approach to holistically represent documents for IR, representing and retrieving documents interleaved with multiple modalities in a unified manner (illustrated in Figure 1). Specifically, it revolves around the recent advance of Vision-Language Models (VLMs), which enable the processing and integration of multimodal content (such as text, images, and tables) directly into a single token sequence, thereby preserving the context and relationships between various parts of the document, unlike prior methods that rely on the fragmented visual representations. Additionally, in cases where the number of tokens in a document is large and exceeds the capacity of a single context window of VLMs, we propose a strategy to segment the document into passages, each represented within the token limit, and combine these passage embeddings into a unified document representation. This strategy differs from existing approaches that independently represent and retrieve at the passage level, potentially losing the overall document context. Lastly, to accurately identify only the relevant sections within the retrieved lengthy document, we introduce a reranking mechanism that is trained to pinpoint the passage most pertinent to the query (among all the other passages within the document), allowing for both the coarse-grained document-level matching and fine-grained passage-level retrieval. We refer to our overall framework as **Interleaved Document Information Retrieval System (IDentify)**.

We experimentally validate the effectiveness of IDentify on four benchmark datasets, considering both the text-only and multimodal queries. On a battery of tests conducted, we observe that our approach substantially outperforms relevant baselines that consider only the uni-modality or certain facets of multi-modality, thanks to the holistic consideration of multimodal content. Further, we find that the strategy to represent the whole document with its single representation (by merging embeddings of its splits) is superior to the approach of individually representing them for document retrieval, but also performing reranking over the sections of the retrieved document is superior to the approach of directly retrieving those sections, confirming the efficacy of proposed retrieval and reranking pipeline for document and passage retrieval, respectively.

2 Related Work

Information Retrieval IR involves finding documents relevant to a query, which plays a crucial role in applications such as search and question-answering (Zhu et al., 2023; Gao et al., 2023; Ram et al., 2023; Shi et al., 2024; Jeong et al., 2024a). Earlier IR approaches measured the similarity between queries and documents based on their lexical term matching, such as BM25 and TF-IDF (Robertson et al., 1994; Jones, 2004). However, these methods struggled to capture semantic nuances beyond surface-level term overlaps. Recently, along with advancements in language models (Devlin et al., 2019; Liu et al., 2019), there have been dense retrieval methods that embed both queries and documents into a shared dense vector space (Karpukhin et al., 2020; Xiong et al., 2021), enabling the calculation of semantic similarity between them more

effectively by capturing the deeper contextual information. Yet, previous studies have mainly focused on enhancing the textual representations of queries and documents, while overlooking the multimodal nature of documents beyond text, which can provide richer context and aid in more accurate retrieval (Liu et al., 2021; Jeong et al., 2024b).

Multimodal Information Retrieval Recent studies in IR have expanded the focus from purely text-based retrieval models to those that consider other modalities, such as images (Radford et al., 2021; Xiao et al., 2024), tables (Herzig et al., 2021; Chen et al., 2024) and graphs (Baek et al., 2023); however, the majority of these approaches (Zhou et al., 2024; Long et al., 2024; Lerner et al., 2024; Nowak et al., 2024; Caffagni et al., 2024) have primarily explored how to process the multimodal *queries*, and overlooked the equally important multimodal characteristics of the *documents* being retrieved. In efforts to handle diverse multimodal elements within documents, there are concurrent studies that have proposed to capture screenshots of documents, such as PDFs (Faysse et al., 2024; Cho et al., 2024) or Wikipedia web pages (Ma et al., 2024), and subsequently encoding them through vision models (Ding et al., 2024). Yet, these methods are not only limited by factors, such as image resolution and computational memory, constraining their application to documents longer than a single page¹, but also fall short by treating the diverse modalities within a document as a single visual entity, leading to document representations that may fail to effectively capture the nuanced interdependence between text and images. Also, while there are concurrent studies (Jiang et al., 2024d; Lin et al., 2024) that consider images and text as retrieval targets, they primarily focus on representing image-text pairs and their retrieval, rather than addressing the holistic representation of documents that include multiple images and another modality (tables). Finally, all the aforementioned work does not address the issue of splitting documents into smaller fragments (passages or sub-images), which may disrupt the holistic contextual view of the entire document.

Vision-Language Models Recently developed VLMs have emerged as a powerful tool for jointly processing visual and textual data, which combine the image understanding capabilities of vi-

sual encoders (Radford et al., 2021; Zhai et al., 2023) with the advanced reasoning abilities of language models (OpenAI, 2022, 2023a). They have achieved remarkable performance across diverse vision-language tasks (such as image captioning and visual question answering) (Dai et al., 2023; OpenAI, 2023b), with the substantially limited attention on their applications to IR. We note that the latest developments in this field have particularly focused on enabling VLMs to handle interleaved, multimodal content, involving a mixed sequence of images and text (Zhang et al., 2023; Li et al., 2024b). In particular, LLaVA-NeXT-Interleave (Li et al., 2024b) introduces a fine-tuning approach that specifically enhances the VLMs’ capacity to understand complex interleavings of multiple images and text within a single context. Drawing inspiration from these advances, we propose to harness their capabilities to create unified embeddings for documents interleaved with text and images (and tables).

3 Method

We present IDentify to holistically represent documents interleaved with multimodal elements.

3.1 Preliminaries

We begin with formally explaining IR and VLMs.

Information Retrieval IR is the task of identifying a set of relevant documents $\{d_1, d_2, \dots, d_k\} \subseteq \mathcal{D}$ from a large corpus \mathcal{D} , given a query q . Here, each query q and document d are represented as a sequence of tokens, e.g., $q = [q_1, \dots, q_n]$, and traditional IR approaches typically consider these tokens as purely textual elements. However, we propose to extend this assumption to have the tokens of both the textual and visual content, to capture the multimodal nature of many real-world documents. Then, this new extension raises important questions of how can both the textual and visual content be represented within a unified token framework, and how can these multimodal tokens be seamlessly integrated and encoded for document representations.

Vision-Language Models To answer them, we now turn to describing VLMs, which are designed to jointly encode the textual and visual information in a unified token framework. These models are generally comprised of two main components: a visual encoder and a language model, interconnected through a projection layer. Specifically, given the document that may contain interleaved modalities (e.g., text and images), the visual encoder extracts high-level visual features from images embedded

¹It requires processing 9.8k image tokens just to process a single-page document, and it results in 2TB of storage for handling the entire Wikipedia corpus, which may not be practical.

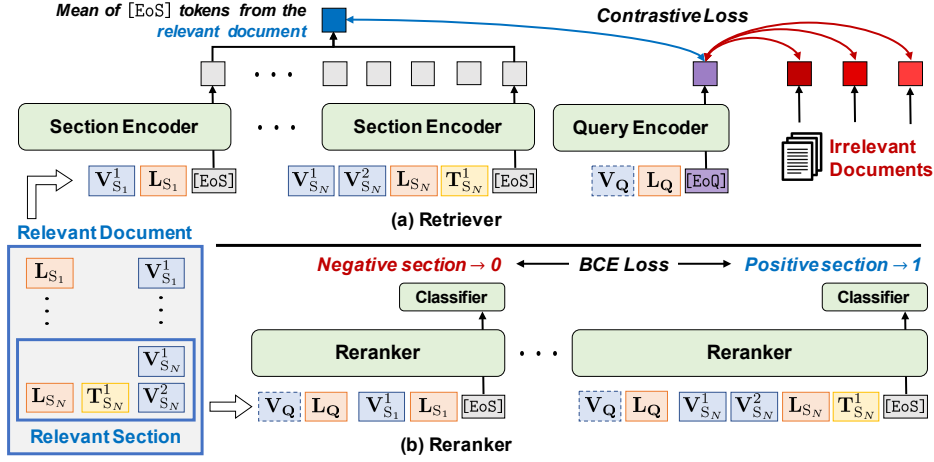


Figure 2: Overview of the proposed IDentify. (a): In our document retriever, a query encoder represents a query (purple), and sections are encoded with a section encoder whose embeddings are averaged to form a document representation (blue). Contrastive learning loss (red) is used for training the document retriever. (b): Reranker scores query-section relevance with the concatenation of the query and section, trained using Binary Cross-Entropy loss.

within the document, mapping them into a latent space. Then, these visual features are transformed into a sequence of visual tokens via the projection layer, represented as follows: $\mathbf{V} \in \mathbb{R}^{V \times d_{\text{emb}}}$, where V denotes the visual token length and d_{emb} is the token dimension size. Similarly, for the textual content embedded within the document, the language model uses a word embedding layer to convert the input text into a sequence of tokens, as follows: $\mathbf{L} \in \mathbb{R}^{L \times d_{\text{emb}}}$, where L denotes the text token length.

In this work, we also propose to account for tables that are the integral modality to holistically represent the full content of documents. Yet, unlike text and images that have dedicated processing layers within VLM architectures, tables do not have a specific representation layer. Nevertheless, we argue that VLMs are pre-trained on diverse web data, and subsequently learned implicitly to handle the table structures formatted in HTML. Consequently, we treat HTML-format table data as a linearized sequence of HTML words, applying the same word embedding layer as is used for plain text. To be formal, this process converts the table content into table tokens, as follows: $\mathbf{T} \in \mathbb{R}^{T \times d_{\text{emb}}}$, where T is the token length of the table. Lastly, once extracted, the visual tokens, text tokens, and table tokens are concatenated (into a unified token sequence) and then passed through the remaining layers of VLMs, to capture both uni- and cross-modal relationships across different modalities, ultimately enabling the comprehensive understanding of the documents.

3.2 Retriever

We now explain how we design a retriever specifically tailored for multimodal interleaved document retrieval. In particular, our approach leverages a

VLM capable of processing text, images, and tables within a single document. Further, following the standard practice of existing retrieval architectures (Karpukhin et al., 2020; Xiong et al., 2021), we use a dual-encoder structure, which consists of a query encoder and document (or section) encoder, both are based on VLMs, illustrated in Figure 2 (a).

Specifically, thanks to the use of the VLM, our query encoder can take either purely textual queries $q = \mathbf{L}_Q$ or multimodal queries consisting of text and visual elements $q = [\mathbf{V}_Q, \mathbf{L}_Q]$. Also, to obtain the final query representation, we use a learnable token called ‘End of Query’, $[\text{EoQ}] \in \mathbb{R}^{d_{\text{emb}}}$, which is appended to the end of the query tokens q . The final concatenated tokens $[q, [\text{EoQ}]]$ are then passed through the query encoder. Lastly, the model output corresponding to $[\text{EoQ}]$ is used as the final query representation, as follows: $\mathbf{Z}_Q \in \mathbb{R}^{d_{\text{emb}}}$.

For documents, we represent each of them d as a sequence of sections: $d = [s_i]_{i=1}^S$ (with a total of S sections), where each section s_i is derived by dividing the document according to its subtitles. s_i can contain a combination of text tokens \mathbf{L}_{S_i} , visual tokens from embedded images \mathbf{V}_{S_i} , and table tokens \mathbf{T}_{S_i} , denoted as follows: $s_i = [\mathbf{V}_{S_i}, \mathbf{L}_{S_i}, \mathbf{T}_{S_i}]$. Then, to obtain a section-level representation, similar to the query representation, we introduce a learnable token, called ‘End of Section’: $[\text{EoS}] \in \mathbb{R}^{d_{\text{emb}}}$, which is appended at the end of each section. We then forward concatenated tokens $[s_i, [\text{EoS}]]$ to the section encoder, and, after that, the output corresponding to $[\text{EoS}]$ is used to form the section representation, as follows: $\mathbf{Z}_{S_i} \in \mathbb{R}^{d_{\text{emb}}}$. Additionally, the overall document representation is obtained by averaging the representations of all sections within the document, as follows: $\mathbf{Z}_D = \frac{1}{S} \sum_{i=1}^S \mathbf{Z}_{S_i}$.

The remaining step is to train those two query and section encoders. Recall that the goal of the retriever is to assess a relevance score between the query and the document. To achieve this, we use a contrastive learning loss based upon the query and document representations, whose objective is to assign higher similarity scores to relevant documents (positive samples) and lower scores to irrelevant ones (negative samples) for the query, as follows:

$$\mathcal{L}_{\text{retriever}} = -\frac{1}{B} \sum_{i=1}^B \log \left(\frac{\phi(\mathbf{Z}_{Q_i}, \mathbf{Z}_{D_i})}{\sum_{j=1}^B \phi(\mathbf{Z}_{Q_i}, \mathbf{Z}_{D_j})} \right),$$

$$\phi(\mathbf{a}, \mathbf{b}) = \exp \left(\frac{\mathbf{a}^\top \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} \right), \quad (1)$$

where B is the batch size. By minimizing $\mathcal{L}_{\text{retriever}}$, the retriever learns to optimize the similarity between queries and their relevant documents, enabling the retrieval of the most pertinent documents for the given input query during inference.

3.3 Reranker

To enable fine-grained retrieval within documents beyond the retrieval of documents themselves, we introduce a section-level reranking mechanism that identifies the section most relevant to the query. In particular, once the document is retrieved, the objective of the reranker f_R is to pinpoint the specific sections within the document that best match the query. We also note that this reranker is similarly operationalized with the VLM along with a binary classifier on top of it, which directly measures the relevance of each query-section pair (Figure 2 (b)).

Formally, for a retrieved document, we take each of its sections s_i with a learnable token for section embedding [EoS] attached to the end and concatenate it with query q , forming the input sequence of $[q, s_i, [\text{EoS}]]$. The concatenated tokens are then processed through the reranker, and its output corresponding to [EoS] captures the relevance between the query and section, which is further subsequently passed to a binary classifier. Through this, the classifier outputs a probability score indicating the likelihood of the section being relevant to the query, *i.e.*, a score close to one denotes a high relevance.

To train this reranker, we use the binary cross-entropy loss, formulated as follow:

$$\mathcal{L}_{\text{reranker}} = \sum_{i=1}^B \sum_{j=1}^{S_i} \frac{1}{BS_i} \ell(y_{s_{i,j}}, f_R([q, \hat{s}_{i,j}])) ,$$

$$\ell(y, \hat{y}) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})], \quad (2)$$

where S_i is the number of sections in the i -th document, $y_{s_{i,j}}$ is the label for the j -th section of the i -th document $s_{i,j}$ (with its value of one if relevant to the query q , otherwise zero), $\hat{s}_{i,j} = [s_{i,j}, [\text{EoS}]]$, and B is the batch size during training. Also, during training, the sections not labeled as relevant to the query are considered negative samples. Then, by minimizing $\mathcal{L}_{\text{reranker}}$, the reranker learns to predict section relevance for any query, thereby refining our overall retrieval process by allowing the retrieval of not just whole documents but also their most relevant sections, for multiple use cases of IR.

4 Experiments

4.1 Experimental Setups

Datasets We evaluate the proposed IDentify on four benchmark datasets designed for multimodal IR that require understanding of both textual and visual cues within queries and documents, as follows: **Encyclopedic-VQA** (Mensink et al., 2023) is a large-scale benchmark for multimodal Visual Question Answering (VQA) with queries linked to specific Wikipedia sections and includes both textual and multimodal queries; **InfoSeek** (Chen et al., 2023) is a knowledge-intensive VQA dataset with multimodal questions generated from Wikidata triples that include diverse entities such as landmarks, animals, and food; **ViQuAE** (Lerner et al., 2022) involves both text-based and multimodal queries about human entities, linked to annotated Wikipedia sections, making it ideal for evaluating section reranking; **Open-WikiTable** (Kweon et al., 2023) extends WikiSQL (Zhong et al., 2017) and WikiTableQuestions (Pasupat and Liang, 2015), targeting open-domain table QA by identifying documents or sections containing relevant tables. We provide more details on datasets in Appendix A.

Baselines To comprehensively validate IDentify, we compare it against two categories of baselines:

- **Conventional VLM Baselines:** We consider earlier VLMs, which are not capable of jointly processing text and images, such as CLIP (Radford et al., 2021) and BLIP (Li et al., 2022). Also, we consider the approaches, such as UniIR (Wei et al., 2024), which is built on top of them and fine-tuned with a contrastive loss (Equation (1)).
- **Baselines with Different Document Representations:** We further consider existing approaches, representing documents in various ways. **Entity** and **Abstract** baselines retrieve documents based on their titles and summaries, respectively, us-

Table 1: Results with different document retrievers.

Method	R@1	R@10	R@100	MRR@10
CLIP-ViT-L-14				
Zero-Shot	1.9	6.3	13.9	3.1
UniIR + Text-Only	3.8	20.6	50.3	7.7
UniIR + Text & Image	5.8	21.5	48.5	10.0
BLIP-Large				
Zero-Shot	0.0	0.0	0.0	0.0
UniIR + Text-Only	9.8	36.9	71.4	16.3
UniIR + Text & Image	9.9	23.9	60.7	13.5
LLaVA-NeXT-Interleave-0.5B				
Entity	3.1	15.5	39.7	6.1
Abstract	13.4	41.3	66.5	21.6
Text-Only	12.5	37.8	68.7	19.8
Text & Table	12.6	38.6	68.5	19.9
Text & Image	16.4	45.4	77.1	25.3
Identfy (Ours)	20.5	50.0	78.0	29.4

ing high-level textual cues. **Text-only** baselines utilize the full textual content of documents for retrieval (Caffagni et al., 2024; Wang et al., 2024). **Text & Table** and **Text & Image** baselines leverage tables and first image of documents alongside the text, respectively (Jiang et al., 2024a; Lin et al., 2024; Jiang et al., 2024d). These baselines, like our method, are built on the same recent VLMs for direct comparison. **Identfy** is our model that holistically represents multimodal content (text, images, and tables) in documents.

Evaluation Metrics To evaluate our approach, we use standard metrics: Recall@K (R@K) measures whether the relevant document or section appears within the top-K results; MRR@K measures how early the first relevant item is ranked (within top-K) by averaging its inverse rank across queries.

Implementation Details We use LLaVA-NeXT-Interleave (Li et al., 2024b) as the basis VLM for both the retriever and reranker, and also use LLaVA-OneVision (Li et al., 2024a) as an additional basis VLM to show the robustness of Identfy. Following the convention of using the basis of retrieval with less than 1B parameters to balance computational efficiency and retrieval performance (Radford et al., 2021; Zhou et al., 2024; Wei et al., 2024), we choose 0.5B-parameter versions of the VLMs. During training, documents are represented using randomly selected four sections, while in inference, we consider all sections within each document. For section-level retrieval, all sections within the top 25 retrieved documents are reranked. Experiments are conducted on a single H100 GPU.

4.2 Experimental Results and Analyses

Main Results We report retrieval performance on the Encyclopedic-VQA dataset in Table 1, where queries include both text and images. Identfy sig-

Table 2: Comparison of different IR strategies for section retrieval. Document (Ours) performs document retrieval and section reranking, whereas Passage performs passage retrieval and reranking. * denotes the model without reranking.

Granularity	R@1	R@10	R@20	MRR@10
Passage*	3.9	16.9	22.0	7.5
Passage	28.6	36.4	37.8	31.2
Document (Ours)	35.1	50.8	53.6	40.3

nificantly outperforms all baselines built on VLMs such as CLIP and BLIP, which are limited to handling a single image alongside text and encoding image-text representations independently, making them suboptimal for understanding multimodal interactions within documents. We also observe that Identfy achieves the best performance, improving R@1 scores by 53.0%, 64.0%, 62.7%, and 25.0% over Abstract, Text-Only, Text & Table and Text & Image retrieval baselines, respectively, with similar trends observed for other metrics. These results demonstrate the effectiveness of integrating multimodal content holistically into a unified representation. To further illustrate the advantages of our approach, we provide case studies in Appendix E.

We further examine the impact of our pipeline of document retrieval and section reranking. In Table 2, the passage retriever represents individual sections as separate retrieval units, whereas the document retriever (ours) aggregates multiple section representations into a single representation. Then, we perform reranking over the retrieved sections or the sections from the retrieved documents, and then report the results in Table 2 (where * denotes the model without reranking). From this, we observe that the passage retriever without reranking (Passage*) achieves suboptimal retrieval performance, highlighting the challenge in pinpointing the most relevant section within a document using traditional retrieval methods. In contrast, when the reranker is used alongside the document retriever, the performance significantly surpasses the passage retrieval, demonstrating the effectiveness of our coarse-to-fine document-to-section retrieval strategy.

Interleaved format enhances document retrieval across modalities. We further expand our experiments to two additional datasets, InfoSeek and ViQuAE, and report document retrieval results. As shown in Table 3 Left, our model consistently outperforms the Text-document baseline for both the multimodal and textual queries. We attribute these gains to the integration of multimodal content, allowing the VLM to capture richer alignments and leverage pre-existing knowledge for more effective document representation (Xu et al., 2024).

Table 3: Performance on document retrieval and section reranking for multimodal and textual queries on Encyclopedic-VQA (Enc-VQA), ViQuAE, and InfoSeek. We compare the approach that solely uses textual information from documents (Text-Only) and our approach of leveraging interleaved multimodal contents from the documents (IDentIfy) over various scenarios.

Dataset	Query Type	Method	Document Retrieval				Section Reranking			
			R@1	R@10	R@100	MRR@10	R@1	R@10	R@20	MRR@10
Enc-VQA	Multimodal	Text-Only IDentIfy (Ours)	12.5 20.5	37.8 50.0	68.7 78.0	19.8 29.4	40.7 42.4	52.8 53.6	55.5 55.7	44.8 46.3
	Textual	Text-Only IDentIfy (Ours)	62.7 65.4	76.3 76.8	87.4 87.8	67.0 69.0	68.1 69.7	79.4 80.1	80.2 80.6	72.3 73.6
ViQuAE	Multimodal	Text-Only IDentIfy (Ours)	13.5 17.5	40.4 46.0	67.4 69.4	20.9 26.3	12.6 11.4	31.7 32.1	37.7 39.2	18.2 17.5
	Textual	Text-Only IDentIfy (Ours)	55.8 56.5	71.5 72.2	83.0 83.0	60.9 61.6	27.8 29.9	50.2 50.9	57.7 59.8	35.0 36.7
InfoSeek	Multimodal	Text-Only IDentIfy (Ours)	6.8 10.2	23.6 30.4	52.5 57.3	11.2 15.7	N/A N/A	N/A N/A	N/A N/A	N/A N/A

Table 4: Retrieval results for tables, where Zero-shot denotes a model trained on Encyclopedic-VQA but not on the target dataset. Finetuned refers to additional training of the model on the target dataset. (a): Results for tabular document retrieval on Open-WikiTable (OWT). (b): Textual and tabular section reranking results on ViQuAE and OWT datasets, respectively. (c): Reranker accuracy of a classification task that identifies the section containing the query-associated table given a gold document.

(a) Document Retrieval for Tables					(b) Section Reranking for Tables						
Method	R@1	R@10	R@100	MRR@10	Dataset	Target	Method	R@1	R@10	R@20	MRR@10
Zero-shot	29.4	58.0	86.0	38.1	ViQuAE	Text	Zero-shot	20.3	49.0	57.7	28.9
Finetuned	55.8	84.1	93.5	66.1			Finetuned	29.9	50.9	59.8	36.7
(c) Tabular Classification					OWT	Table	Zero-shot	5.9	20.5	29.4	9.1
Method	Random	Zero-shot	Finetuned	Finetuned			8.4	36.7	52.8	15.2	
Acc@1	11.9	9.3	56.5								

Interleaved format is also beneficial in section retrieval. Similarly, we evaluate section retrieval performance on Encyclopedic-VQA and ViQuAE datasets, for both multimodal and textual queries. As shown in Table 3 Right, our model outperforms the Text-document baseline in most cases. However, the performance gains over the baseline are smaller compared to the document retrieval setup. This is likely because section reranking focuses on evaluating the relationship between a single section and a query (rather than leveraging the holistic context of the entire document), and individual sections may lack the diverse multimodal information necessary for fully capturing the intent of queries.

Retrieving tables interleaved within documents is challenging. We explore the retrieval task for tabular data, aiming to identify documents or sections containing query-relevant tables, and compare models trained on Encyclopedic-VQA (Zero-shot) with those additionally trained on Open-WikiTable (Finetuned). As shown in Table 4 (a), the Finetuned retriever outperforms the Zero-shot retriever on retrieving documents containing query-relevant tables. However, more fine-grained section reranking results (identifying sections containing query-relevant tables) in Table 4 (b) may reveal a notable modality-specific challenge: the performance of Zero-shot and Finetuned rerankers is considerably

lower on table retrieval compared to their performance on text retrieval, despite both the text and tables being represented with word tokens. To better understand this, we design a classification task, where rerankers are tasked with identifying the correct section containing the target table within the golden document. Then, as shown in Table 4 (c), the Zero-shot reranker performs comparably to random selection, while the Finetuned reranker shows modest improvements. These findings highlight the intrinsic challenge of tabular retrieval, suggesting the need for table-specific modules to more holistically represent multimodal interleaved documents.

More sections enhance document retrieval performance but raise computational costs. To see how the number of sections used for representing each document impacts performance, we evaluate document retrieval on the InfoSeek dataset by varying the sections per document during training. As shown in Figure 3, incorporating more sections improves MRR@10 from 7.5 to 15.7 due to leveraging richer multimodal and contextual information. However, this comes at the cost of increased computational requirements, as processing more sections raises GPU memory consumption.

BCE loss is the most effective to train the section reranker. In our reranker design, we use a binary cross-entropy (BCE) loss by concatenating

Figure 3: Trade-off between performance (MRR@10) and training cost (GPU Memory) for retrieval.

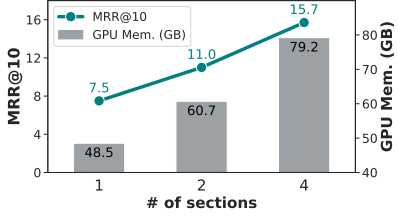


Table 5: Comparison of training objectives for the reranker: Contrastive uses contrastive loss similar to the document retriever training; Doc + BCE concatenates the query with multiple sections from the same document and uses the BCE loss; Sec + BCE trains the reranker by concatenating the query with each section individually.

Query Type	Train Loss	R@1	R@10	R@20	MRR@10
Multimodal	Contrastive	3.6	15.0	21.3	6.5
	Doc + BCE	13.6	29.6	32.9	24.1
	Sec + BCE (Ours)	42.4	53.6	55.7	46.3
Textual	Contrastive	13.6	37.7	45.1	20.6
	Doc + BCE	23.8	43.4	47.2	39.1
	Sec + BCE (Ours)	69.7	80.1	80.6	73.6

Table 6: Comparison of negative sample selection strategies for reranker training: Top-K (top-k retrieved sections), In-batch (sections from other samples in the batch), and In-document (sections in the same document).

Negative	R@1	R@20	MRR@10
Top-K	38.1	55.3	44.4
In-batch	39.5	55.4	45.0
In-document (Ours)	42.4	55.7	46.3

Table 7: Results with another base model (LLaVA-OneVision-0.5B) for document retrieval (with different document formats).

Format	R@1	R@10	R@100	MRR@10
Entity	2.3	10.3	29.7	4.3
Abstract	7.6	24.7	55.7	12.0
Text-Only	7.0	24.1	50.4	11.7
Text & Table	6.9	26.3	54.9	12.1
Text & Image	9.3	31.4	61.9	15.4
IDentIfy (Ours)	12.1	36.1	62.5	18.2

the query with each document section individually (Section + BCE), allowing the model to directly assess query-section relevance. As an alternative, we also explore a contrastive loss (Contrastive), which models section reranking similarly to document retrieval but uses sections as the retrieval units, and a variant of BCE loss (Document + BCE), where the query is concatenated with multiple sections (both positive and negative) from the same document. As shown in Table 5, the Section + BCE reranker outperforms both alternatives. Specifically, contrastive loss performs the worst, suggesting that direct concatenation of query and section provides clearer relevance signals, consistent with conventional reranking approaches. Moreover, while Document + BCE leverages inter-section context, its performance might be hindered by training constraints as the model processes fewer sections during training (Jiang et al., 2024c; Lee et al., 2024), and addressing it would be interesting future work.

Sections from the same document act as effective negatives to enhance reranker performance. In training the reranker, we investigate whether considering sections from the same document as negative examples (called In-document) is effective than other strategies, such as Top-K negatives (top-K retrieved sections based on their similarity with the input query) and In-batch negatives (positive sections from other samples in the same batch). As shown in Table 6, we observe that the In-document approach achieves superior performance especially on R@1, demonstrating its ability to effectively identify the most pertinent section among highly similar sections within the same document, i.e., its training objective can encourage the reranker to

focus on fine-grained distinctions between closely related sections (within the same document).

Our IDentIfy is Versatile with Different VLMs.

To ensure the effectiveness and robustness of IDentIfy across VLMs, we evaluate its performance with another VLM, LLaVA-OneVision (Li et al., 2024a), with 0.5 billion parameters, in addition to LLaVA-NeXT-Interleave (Li et al., 2024b) used in our main experiments. Results in Table 7 show that ours continues to outperform baselines, achieving a notable 30.1% gain in R@1 over the best baseline.

5 Conclusion

In this paper, we introduced IDentIfy, a novel IR framework designed to address the limitations of conventional methods that rely on textual content of documents and their segmented passages. Specifically, our approach sits on top of recent VLMs, which enables integration and representation of diverse multimodal content (including text, images, and tables) into a unified document representation. Also, unlike previous strategies that segment documents at the passage level, our method merges these segments to maintain the document’s structural coherence, while further introducing a reranking strategy for precise identification of relevant sections. Extensive experiments across various IR datasets demonstrated that IDentIfy consistently outperforms existing baselines, confirming that the interleaved multimodal representation significantly enhances the quality of the document and section retrieval. We believe IDentIfy represents a crucial step toward more comprehensive and contextually aware IR systems, capable of handling the increasing multimodality of modern information sources.

Limitations

Due to the constraints of a single H100 GPU that we have, we represent documents by sampling a limited number of sections and averaging their corresponding embeddings (See Figure 3). While this reduces the computational demands, our findings suggest that capturing a broader document context leads to improved retrieval performance. Hence, leveraging the long context window of LVLMs with a greater number of sections could further enhance document retrieval by capturing more comprehensive information within the full document. Additionally, while using the basis model size of 0.5B (or less than 1B) parameters is a standard practice in IR literature, scaling up the basis VLMs remains an avenue for future work; however, although larger models can yield performance gains, they come at the cost of increased computational requirements. Moreover, our reranker design follows the conventional approach of concatenating the input query with individual sections. However, we believe providing the reranker with all sections together would allow the model to better leverage the contextual information from the entire document, potentially resulting in improved performance, and we leave explorations on this for future work.

Ethics Statement

In this work, we use a publicly available retrieval corpus for information retrieval tasks. However, the retrieval corpus may contain private, harmful, or biased content. Such undesirable features could unintentionally be reflected in the behavior of retrievers and rerankers trained on this data, potentially leading to ethical concerns during real-world deployment. However, current information retrieval techniques, including ours, do not address the retrieval of undesirable content. We recognize the critical need for safeguards to mitigate this issue. This is essential to ensure that information retrieval systems are reliable, fair, and safe for deployment.

References

- Jinheon Baek, Alham Fikri Aji, Jens Lehmann, and Sung Ju Hwang. 2023. Direct fact retrieval from knowledge graphs without entity linking. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Davide Caffagni, Federico Cocchi, Nicholas Moratelli, Sara Sarto, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2024. [Wiki-llava: Hierarchical](#)

[retrieval-augmented generation for multimodal llms](#). *arXiv preprint arXiv:2404.15406*.

- Peter Baile Chen, Yi Zhang, and Dan Roth. 2024. Is table retrieval a solved problem? exploring join-aware multi-table retrieval. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, Soravit Changpinyo, Alan Ritter, and Ming-Wei Chang. 2023. Can pre-trained vision and language models answer visual information-seeking questions? In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Jaemin Cho, Debanjan Mahata, Ozan Irsoy, Yujie He, and Mohit Bansal. 2024. [M3docrag: Multi-modal retrieval is what you need for multi-page multi-document understanding](#). *Preprint*, arXiv:2411.04952.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*.
- Yihao Ding, Kaixuan Ren, Jiabin Huang, Siwen Luo, and Soyeon Caren Han. 2024. [PDF-MVQA: A dataset for multimodal information retrieval in pdf-based visual question answering](#). *arXiv preprint arXiv:2404.12720*.
- Kuicai Dong, Derrick-Goh-Xin Deik, Yi Lee, Hao Zhang, Xiangyang Li, Cong Zhang, and Yong Liu. 2024. [Mc-indexing: Effective long document retrieval via multi-view content-aware indexing](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 2673–2691. Association for Computational Linguistics.
- Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, Céline Hudelot, and Pierre Colombo. 2024. [Colpali: Efficient document retrieval with vision language models](#). *arXiv preprint arXiv:2407.01449*.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. 2023. [Retrieval-augmented generation for large language models: A survey](#). *arXiv preprint arXiv:2312.10997*.

731	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu-	reasoning over table. In <i>Findings of the Association</i>	786
732	pat, and Ming-Wei Chang. 2020. REALM: retrieval-	<i>for Computational Linguistics (ACL)</i> .	787
733	augmented language model pre-training . <i>arXiv</i>		
734	<i>preprint arXiv:2002.08909</i> .		
735	Jonathan Herzig, Thomas Müller, Syrine Krichene, and	Jinhyuk Lee, Anthony Chen, Zhuyun Dai, Dheeru Dua,	788
736	Julian Martin Eisenschlos. 2021. Open domain ques-	Devendra Singh Sachan, Michael Boratko, Yi Luan,	789
737	tion answering over tables via dense retrieval. In	Sébastien M. R. Arnold, Vincent Perot, Siddharth	790
738	<i>Proceedings of the North American Chapter of the</i>	Dalmia, Hexiang Hu, Xudong Lin, Panupong Pasu-	791
739	<i>Association for Computational Linguistics (NAACL)</i> .	pat, Aida Amini, Jeremy R. Cole, Sebastian Riedel,	792
740	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	Iftexhar Naim, Ming-Wei Chang, and Kelvin Guu.	793
741	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	2024. Can long-context language models subsume	794
742	Weizhu Chen. 2022. Lora: Low-rank adaptation of	retrieval, rag, sql, and more? <i>arXiv preprint</i>	795
743	large language models. In <i>Proceedings of the Inter-</i>	<i>arXiv:2406.13121</i> .	796
744	<i>national Conference on Learning Representations</i>		
745	<i>(ICLR)</i> .	Paul Lerner, Olivier Ferret, and Camille Guinaudeau.	797
746	Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju	2024. Cross-modal retrieval for knowledge-based vi-	798
747	Hwang, and Jong Park. 2024a. Adaptive-rag: Learn-	sual question answering. In <i>Advances in Information</i>	799
748	ing to adapt retrieval-augmented large language mod-	<i>Retrieval - 46th European Conference on Information</i>	800
749	els through question complexity. In <i>Proceedings of</i>	<i>Retrieval</i> .	801
750	<i>the North American Chapter of the Association for</i>		
751	<i>Computational Linguistics (NAACL)</i> .	Paul Lerner, Olivier Ferret, Camille Guinaudeau,	802
752	Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju	Hervé Le Borgne, Romaric Besançon, José G.	803
753	Hwang, and Jong C. Park. 2024b. Database-	Moreno, and Jesús Lovón-Melgarejo. 2022. Viquae,	804
754	augmented query representation for information re-	a dataset for knowledge-based visual question an-	805
755	trieval . <i>arXiv preprint arXiv:2406.16013</i> .	swering about named entities. In <i>SIGIR '22: The</i>	806
756	Ting Jiang, Minghui Song, Zihan Zhang, Haizhen	<i>45th International ACM SIGIR Conference on Re-</i>	807
757	Huang, Weiwei Deng, Feng Sun, Qi Zhang, Deqing	<i>search and Development in Information Retrieval,</i>	808
758	Wang, and Fuzhen Zhuang. 2024a. E5-V: universal	<i>Madrid, Spain, July 11 - 15, 2022</i> .	809
759	embeddings with multimodal large language models .		
760	<i>arXiv preprint arXiv:2407.12580</i> .	Patrick S. H. Lewis, Ethan Perez, Aleksandra Pik-	810
761	Ziyan Jiang, Xueguang Ma, and Wenhui Chen. 2024b.	tus, Fabio Petroni, Vladimir Karpukhin, Naman	811
762	Longrag: Enhancing retrieval-augmented gener-	Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih,	812
763	ation with long-context llms . <i>arXiv preprint</i>	Tim Rocktäschel, Sebastian Riedel, and Douwe	813
764	<i>arXiv:2406.15319</i> .	Kiela. 2020. Retrieval-augmented generation for	814
765	Ziyan Jiang, Xueguang Ma, and Wenhui Chen. 2024c.	knowledge-intensive NLP tasks. In <i>Advances in Neur-</i>	815
766	Longrag: Enhancing retrieval-augmented gener-	<i>al Information Processing Systems (NeurIPS)</i> .	816
767	ation with long-context llms . <i>arXiv preprint</i>		
768	<i>arXiv:2406.15319</i> .	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang,	817
769	Ziyan Jiang, Rui Meng, Xinyi Yang, Semih Yavuz,	Feng Li, Hao Zhang, Kaichen Zhang, Yanwei	818
770	Yingbo Zhou, and Wenhui Chen. 2024d. Vlm2vec:	Li, Ziwei Liu, and Chunyuan Li. 2024a. Llava-	819
771	Training vision-language models for massive	onevision: Easy visual task transfer . <i>arXiv preprint</i>	820
772	multimodal embedding tasks . <i>arXiv preprint</i>	<i>arXiv:2408.03326</i> .	821
773	<i>arXiv:2410.05160</i> .	Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang,	822
774	Karen Spärck Jones. 2004. A statistical interpretation	Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024b.	823
775	of term specificity and its application in retrieval . <i>J.</i>	Llava-next-interleave: Tackling multi-image, video,	824
776	<i>Documentation</i> , 60(5):493–502.	and 3d in large multimodal models . <i>arXiv preprint</i>	825
777	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick	<i>arXiv:2407.07895</i> .	826
778	S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen,	Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H.	827
779	and Wen-tau Yih. 2020. Dense passage retrieval for	Hoi. 2022. BLIP: bootstrapping language-image pre-	828
780	open-domain question answering. In <i>Proceedings of</i>	training for unified vision-language understanding	829
781	<i>the Conference on Empirical Methods in Natural</i>	and generation. In <i>Proceedings of the International</i>	830
782	<i>Language Processing (EMNLP)</i> .	<i>Conference on Machine Learning (ICML)</i> .	831
783	Sunjun Kweon, Yeonsu Kwon, Seonhee Cho, Yohan Jo,	Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong	832
784	and Edward Choi. 2023. Open-wikitable : Dataset	Feng, Lingpeng Kong, and Qi Liu. 2024c. Mul-	833
785	for open domain question answering with complex	timodal arxiv: A dataset for improving scientific	834
		comprehension of large vision-language models . In	835
		<i>Proceedings of the 62nd Annual Meeting of the As-</i>	836
		<i>sociation for Computational Linguistics (Volume 1:</i>	837
		<i>Long Papers)</i> , ACL 2024, Bangkok, Thailand, August	838
		11-16, 2024, pages 14369–14387. Association for	839
		Computational Linguistics.	840
		Sheng-Chieh Lin, Chankyu Lee, Mohammad Shoeybi,	841
		Jimmy Lin, Bryan Catanzaro, and Wei Ping. 2024.	842

843	Mm-embed: Universal multimodal retrieval with	Stephen E. Robertson, Steve Walker, Susan Jones,	898
844	multimodal llms. <i>Preprint</i> , arXiv:2411.02571.	Micheline Hancock-Beaulieu, and Mike Gatford.	899
845	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	1994. Okapi at TREC-3. In <i>Proceedings of The Third</i>	900
846	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	<i>Text REtrieval Conference, TREC 1994, Gaithers-</i>	901
847	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	<i>burg, Maryland, USA, November 2-4, 1994.</i>	902
848	Roberta: A robustly optimized BERT pretraining	Sanket Shah, Anand Mishra, Naganand Yadati, and	903
849	approach. <i>arXiv preprint arXiv:1907.11692.</i>	Partha Pratim Talukdar. 2019. KVQA: knowledge-	904
850	Zheyuan Liu, Cristian Rodriguez Opazo, Damien Teney,	aware visual question answering. In <i>Proceedings of</i>	905
851	and Stephen Gould. 2021. Image retrieval on real-life	<i>the AAAI National Conference on Artificial Intelli-</i>	906
852	images with pre-trained vision-and-language models.	<i>gence (AAAI).</i>	907
853	In <i>Proceedings of the International Conference on</i>	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	908
854	<i>Computer Vision (ICCV).</i>	joon Seo, Richard James, Mike Lewis, Luke Zettle-	909
855	Xinwei Long, Jiali Zeng, Fandong Meng, Zhiyuan	moyer, and Wen-tau Yih. 2024. REPLUG: retrieval-	910
856	Ma, Kaiyan Zhang, Bowen Zhou, and Jie Zhou.	augmented black-box language models. In <i>Proceed-</i>	911
857	2024. Generative multi-modal knowledge retrieval	<i>ings of the North American Chapter of the Associa-</i>	912
858	with large language models. In <i>Proceedings of the</i>	<i>tion for Computational Linguistics (NAACL).</i>	913
859	AAAI National Conference on Artificial Intelligence	Kexin Wang, Nils Reimers, and Iryna Gurevych. 2024.	914
860	(AAAI).	DAPR: A benchmark on document-aware passage	915
861	Xueguang Ma, Sheng-Chieh Lin, Minghan Li, Wenhu	retrieval. In <i>Proceedings of the Association for Com-</i>	916
862	Chen, and Jimmy Lin. 2024. Unifying multimodal	<i>putational Linguistics (ACL).</i>	917
863	retrieval via document screenshot embedding. <i>arXiv</i>	Cong Wei, Yang Chen, Haonan Chen, Hexiang Hu,	918
864	<i>preprint arXiv:2406.11251.</i>	Ge Zhang, Jie Fu, Alan Ritter, and Wenhu Chen.	919
865	Thomas Mensink, Jasper R. R. Uijlings, Lluís Castrejón,	2024. Uniir: Training and benchmarking univer-	920
866	Arushi Goel, Felipe Cadar, Howard Zhou, Fei Sha,	salmultimodal information retrievers. In <i>Proceed-</i>	921
867	André Araújo, and Vittorio Ferrari. 2023. Encyclope-	<i>ings of the European Conference on Computer Vision</i>	922
868	dic VQA: visual questions about detailed properties	(ECCV).	923
869	of fine-grained categories. In <i>Proceedings of the In-</i>	Zilin Xiao, Ming Gong, Paola Cascante-Bonilla,	924
870	<i>ternational Conference on Computer Vision (ICCV).</i>	Xingyao Zhang, Jie Wu, and Vicente Ordonez. 2024.	925
871	Averi Nowak, Francesco Piccinno, and Yasemin Altun.	Grounding language models for visual entity recog-	926
872	2024. Multimodal chart retrieval: A comparison of	nition. <i>arXiv preprint arXiv:2402.18695.</i>	927
873	text, table and image based approaches. In <i>Proceed-</i>	Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang,	928
874	<i>ings of the North American Chapter of the Associa-</i>	Jialin Liu, Paul N. Bennett, Junaid Ahmed, and	929
875	<i>tion for Computational Linguistics (NAACL).</i>	Arnold Overwijk. 2021. Approximate nearest neigh-	930
876	OpenAI. 2022. Introducing chatgpt. https://openai.	bor negative contrastive learning for dense text re-	931
877	com/blog/chatgpt .	trieval. In <i>Proceedings of the International Confer-</i>	932
878	OpenAI. 2023a. GPT-4 technical report. <i>arXiv preprint</i>	<i>ence on Learning Representations (ICLR).</i>	933
879	<i>arXiv:2303.08774.</i>	Jialiang Xu, Michael Moor, and Jure Leskovec. 2024.	934
880	OpenAI. 2023b. GPT-4V(ision) system card. https:	Reverse image retrieval cues parametric memory in	935
881	//openai.com/index/gpt-4v-system-card/ .	multimodal llms. <i>arXiv preprint arXiv:2405.18740.</i>	936
882	Panupong Pasupat and Percy Liang. 2015. Composi-	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov,	937
883	tional semantic parsing on semi-structured tables. In	and Lucas Beyer. 2023. Sigmoid loss for language	938
884	<i>Proceedings of the Association for Computational</i>	image pre-training. In <i>Proceedings of the Interna-</i>	939
885	<i>Linguistics (ACL).</i>	<i>tional Conference on Computer Vision (ICCV).</i>	940
886	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	Pan Zhang, Xiaoyi Dong, Bin Wang, Yuhang Cao,	941
887	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	Chao Xu, Linke Ouyang, Zhiyuan Zhao, Shuan-	942
888	try, Amanda Askell, Pamela Mishkin, Jack Clark,	grui Ding, Songyang Zhang, Haodong Duan, Wen-	943
889	Gretchen Krueger, and Ilya Sutskever. 2021. Learn-	wei Zhang, Hang Yan, Xinyue Zhang, Wei Li, Jing-	944
890	ing transferable visual models from natural language	wen Li, Kai Chen, Conghui He, Xingcheng Zhang,	945
891	supervision. In <i>Proceedings of the International Con-</i>	Yu Qiao, Dahua Lin, and Jiaqi Wang. 2023. Internlm-	946
892	<i>ference on Machine Learning (ICML).</i>	xcomposer: A vision-language large model for ad-	947
893	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay,	vanced text-image comprehension and composition.	948
894	Amnon Shashua, Kevin Leyton-Brown, and Yoav	<i>arXiv preprint arXiv:2309.15112.</i>	949
895	Shoham. 2023. In-context retrieval-augmented lan-	Victor Zhong, Caiming Xiong, and Richard Socher.	950
896	guage models. <i>Trans. Assoc. Comput. Linguistics,</i>	2017. Seq2sql: Generating structured queries from	951
897	11:1316–1331.	natural language using reinforcement learning. <i>arXiv</i>	952
		<i>preprint arXiv:1709.00103.</i>	953

954 Junjie Zhou, Zheng Liu, Shitao Xiao, Bo Zhao, and
955 Yongping Xiong. 2024. VISTA: visualized text em-
956 bedding for universal multi-modal retrieval. In *Pro-
957 ceedings of the Association for Computational Lin-
958 guistics (ACL)*.

959 Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan
960 Liu, Wenhan Liu, Chenlong Deng, Zhicheng Dou,
961 and Ji-Rong Wen. 2023. [Large language models
962 for information retrieval: A survey](#). *arXiv preprint
963 arXiv:2308.07107*.

A Details of Experimental Setups

Dataset configuration Table 8 summarizes the key properties of the datasets used in our experiment, including query modality, target item, entity domain, number of entities, and whether a section ID is provided to indicate the section containing the answer. Additionally, we provide the number of samples in the training, evaluation, and test splits, as well as the size of the corpus. We provide a more detailed explanation of the datasets below.

- **Encyclopedic-VQA** (Mensink et al., 2023) is a large-scale visual question-answering (VQA) benchmark dataset, widely used for measuring the performance of multimodal IR models. Each query is linked to a specific section of a Wikipedia document (containing an answer for it) and is manually annotated by humans. Also, this dataset offers both text-only and multimodal queries. In addition to this, the queries are related to fine-grained properties of species and landmarks. Our experiments focus on the single-hop category where questions can be answered in a single retrieval step.
- **InfoSeek** (Chen et al., 2023) is a dataset designed for knowledge-intensive VQA, covering a wide range of entities (such as landmarks, animals, and food). Questions are generated by filling human-written templates with knowledge triples (subject, relation, object) available from Wikidata, which involve only the multimodal queries. As the test dataset is not available, we use the validation set as our test set, and split the training set into training and validation subsets with a 9:1 ratio.
- **ViQuAE** (Lerner et al., 2022) is a dataset focused about human entities. It provides both textual and multimodal queries, with each query linked to a specific section of a Wikipedia document that contains an answer annotated by humans, which makes it an ideal benchmark for section retrieval.
- **Open-WikiTable** (Kweon et al., 2023) is an extension of WikiSQL (Zhong et al., 2017) and WikiTableQuestions (Pasupat and Liang, 2015), designed for open-domain table question answering that requires retrieval of the most relevant table from a broader corpus. For our experiments, we adapt the WikiTableQuestions subset of Open-WikiTable, aiming at identifying the document or document section containing the target table.

Dataset pre-processing In our study, we leverage interleaved multimodal content from Wikipedia documents. However, existing corpora associated with IR datasets often lack this content, typically only including the first few words of each document. Therefore, we download the HTML file of each Wikipedia document for corpus augmentation.

If the dataset provides Wikipedia URLs for its corpus, we use them to download the HTML files. Alternatively, if only entity names are provided, we generate Wikipedia URLs using those names. If a Wikipedia URL is deprecated, we remove the corresponding document from the corpus along with any associated queries. From the HTML files, we extract text, image URLs, and tables. We then split the contents by subtitles in the document where each chunk corresponds to a section. For the images, we use the image URLs to download the corresponding images, removing any invalid URLs. This process produces a dictionary that organizes text, images, and tables by section.

Since downloading contents for all documents across datasets is time- and memory-intensive, we preprocess subsets of each corpus, including documents relevant to queries in the training, evaluation, and test splits, along with unrelated documents.

Implementation Details To take advantage of larger batch sizes (while reducing GPU memory usage), we apply LoRA (Hu et al., 2022). Also, to further optimize the GPU usage, we scale each image down to half of its original height and width and then combine four scaled-down images into a single composite image. All experiments are conducted using a single H100 GPU.

B Multi-modality Statistics in Documents

We calculate the statistics related to multi-modality in Wikipedia documents, and find that both images and tables are evenly distributed across the whole documents. To be specific, except for the first section of documents, which contains 1.2 images on average, the distribution of images is consistent across the other sections, containing an average of 0.27 images per section. Also, tables appear less frequently, averaging 0.23 per section, but they are uniformly distributed across all sections.

C Efficiency of IDentIfy

During the retrieval process, the computational efficiency (*i.e.*, the retrieval latency) of our approach

Table 8: Information retrieval datasets summary.

Dataset	Query Modality	Target	Domain	Entities	Section ID	Train	Eval	Test	Corpus size
Encyclopedic-VQA	Text, Text-Image	Text	Species, Landmarks	17k	○	177k	2.2k	3.8k	100k
InfoSeek	Text-Image	Text	Diverse	11k	×	209k	23k	74k	500k
ViQuAE	Text, Text-Image	Text	Human	1k	○	1.2k	1.2k	1.2k	100k
Open-WikiTable	Text	Table	Table	-	○	3.3k	0.4k	0.4k	1.8k

remains the same regardless of the number of interleaved modalities and their compositions, as each document representation (averaged from its section embeddings) is encoded into a fixed-sized vector, whose size is the same as the case where we encode only the text. Also, even if we consider the efficiency within the document embedding process (which is typically not a concern for IR tasks as it can be done offline in parallel), the computational costs and memory usage when embedding multimodal documents are similar to the case of embedding text-only documents, as the factors that impact efficiency are not the number of multimodal content but the number of tokens within documents.

D Additional Experimental Results

Data Requirements for Models We analyze the effect of different dataset sizes for training on retriever and reranker performance. To achieve this, we randomly prune samples in the Encyclopedic-VQA dataset at various ratios and report the performance of models trained on these subsets. In [Figure 4 \(a\)](#), we observe that too many samples can degrade retrieval performance. Also, retrieval of textual queries requires fewer samples to reach its optimal performance compared to multimodal retrieval. Similarly, in [Figure 4 \(b\)](#), section retrieval for multimodal queries requires 10% of the dataset to achieve 80% of the full-dataset performance, while section retrieval for textual queries needs only 5%. These observations suggest that additional modalities increase the need for more data. This accounts for the inferior performance of the interleaved format in the ViQuAE experiments ([Table 3 Right](#)). The ViQuAE dataset, at only 2.2% of the size of Encyclopedic-VQA, may be small for the reranker to effectively learn multimodal query-section alignments. We also observe that section retrieval is more challenging, with more samples improving the reranker’s performance. This explains why the ViQuAE reranker has much lower section retrieval scores compared to the one trained on the Encyclopedic-VQA ([Table 3 Right](#)). Given

the challenge of obtaining large query-section pair samples, exploring more effective reranker training pipelines is necessary.

E Case Study

We conduct case studies to demonstrate the advantages of our approach in document retrieval with textual and multimodal queries. In [Figure 5](#) and [Figure 6](#), we illustrate the instances where our approach, which leverages interleaved multimodal contents (e.g., images, tables, and text) within documents, retrieved correct documents for given queries, while the conventional one, which represents documents using only textual data, retrieved documents that appeared to be relevant but were not actually related to the queries.

In [Figure 5](#), a textual query asks for the name of the park located on the north shore of Foster Reservoir. The conventional approach retrieved a document containing unrelated information about a different reservoir. While this document includes terms such as "Peak District National Park" and "North America farm," which make the document superficially relevant, it fails to answer the query. In contrast, our approach identified the document containing the correct answer to the given query.

The advantages of integrating multimodal content into document representation become more apparent in document retrieval with multimodal queries, as shown in [Figure 6](#). For a query consisting of an image of a town hall in Hanover and a textual question about its designer, both our approach and the conventional one retrieved documents about town halls in Germany. However, our approach pinpointed the exact document about the town hall in Hanover, indicating that Hermann Eggert designed the building. The conventional method retrieved a document about a town hall in Munich, which is somewhat related but not an exact match to the query image or question.

These cases underscore the benefits of leveraging multimodal content in information retrieval. Integrating interleaved multimodal elements, our

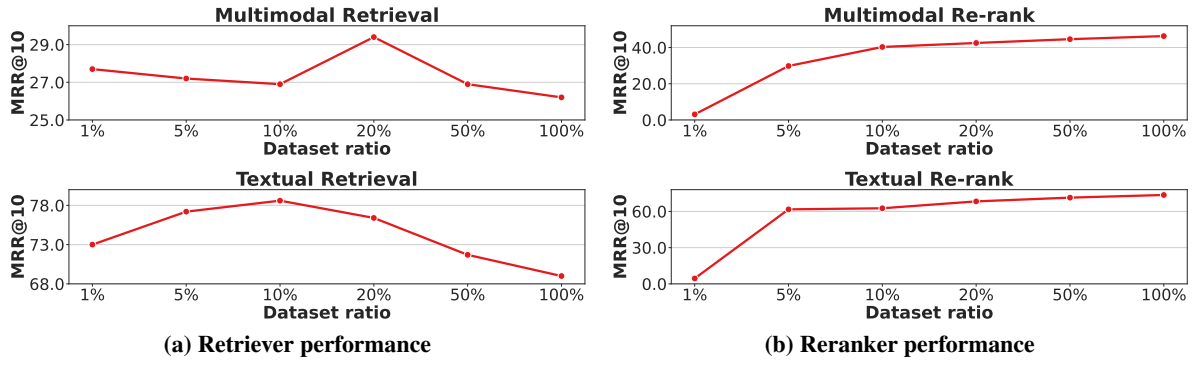


Figure 4: Retrieval performance with different dataset sizes for training. **(a)**: When training a retriever, large datasets rather deteriorate the retrieval performance as it may be overfitted, resulting in low generalization. **(b)**: On the other hand, a larger dataset size is beneficial to training a re-ranker.

approach aligns more effectively with the input query, resulting in more accurate and fine-grained retrieval. This superiority is supported by [Xu et al. \(2024\)](#), which highlights that models perform better when prompted with rich multimodal information, enabling them to capture alignments across modalities and enhance the representation of given inputs.

Q: What is the name of the park on the north shore of foster reservoir?

Foster Reservoir

ArticleTalk

From Wikipedia, the free encyclopedia

ReadEditView historyTools

Coordinates: 44°41′66″N 122°6′60″W﻿ / ﻿44.69333°N 122.10000°W﻿ / 44.69333; -122.10000

Foster Reservoir

<



Q: Who designed this building?

New Town Hall (Hanover)

Article Talk

From Wikipedia, the free encyclopedia

Coordinates: 52°36′24″N 9°37′35″E﻿ / ﻿52.60667°N 9.62639°E﻿ / 52.60667; 9.62639

The **New Town Hall** (German: *Neues Rathaus*) is a town hall in Hanover, Germany. It opened on 20 June 1913 after construction lasting 12 years.^[1] A magnificent, castle-like building of the era of **Wilhelm II** in **eclectic** style at the southern edge of the inner city just outside the historic city centre of Hanover, the building is embedded within the 10-hectare (25-acre) **Maschpark**.^[9a]

History

Costing 10 million marks, the New Town Hall was erected on 6,026 beech piles by architects **Hermann Eggert** and **Gustav Hahnhuber**.^[2] "Ten million marks, Your Majesty – and all paid for in cash", the City Director, **Heinrich Tramm**,^[3a] is claimed to have announced when the New Town Hall was opened in the presence of Emperor **Wilhelm II**. In honour of Tramm the public space in front of the building was named **Trammplatz** (lit. "Tramm Plaza") until 23 September 2024,^[3] when it was renamed to **Platz der Menschenrechte** (lit. "Human Rights Plaza") because Tramm is recognized as a pioneer of **National Socialism**.^[4]

Upon opening, the New Town Hall replaced the **Wangenheim Palace** as the main seat of administration, which had moved from the **Old Town Hall** into the Wangenheim Palace in 1863. As of 2022, the New Town Hall is still "the residence of the Mayor and CEO, the head of the municipal administration."^[1]

Damaged during **bombing raids** on the inner city of Hanover in **World War II**,^[5] the German state of **Lower Saxony** was proclaimed in 1946 in the 38-metre-high (125 ft) hall of the New Town Hall.^[1]

There are four city models of Hanover in the ground floor of the New Town Hall.^[1] They vividly portray the development of the city.

Dome with elevator

The dome of the New Town Hall, with its observation platform, is 97.73 metres (320.6 ft) high.^[6] The dome's lift is unique in the world in that its arched course follows the parabolic shape of the dome.^{[1][7]} It is often incorrectly described as a sloping lift up the dome and compared with the lifts in the **Eiffel Tower**, which actually travel diagonally only, without changing their angle of inclination. The lift climbs the 50-metre (160 ft) shaft at an angle of up to 17° to the gallery of the dome, where the **Harz** mountain range can be seen when visibility is good. In the process, the lift moves 10 metres (33 ft) horizontally. During the trip, the two weight-bearing cables wind up on three double rolls in the wall of the shaft.

The cage of the lift erected in 1913 travelled on steam-bent oak tracks. Because of the weather, this lift was not usable in the colder half of the year. A spiral staircase leads from the lift exit to the observation level. In 2005, over 90,000 people visited the tower of the New Town Hall. A new lift was installed in winter of 2007–08. The last trip of the old lift took place with Lord Mayor **Stephan Weil** on 4 November 2007. On that weekend, 1200 guests took the last opportunity to ride in the old lift.

Gallery

Hanover's New Town Hall at night

Aerial view

Interior

General references

[edit]

- Steinweg, Wolfgang** (in German) (1988). *Das Rathaus in Hannover: von der Kaiserzeit bis in die Gegenwart* [The Town Hall in Hanover: from the Imperial Era to the Present Day] (in German). Hanover: Schlüter. ISBN 3-87706-287-3. OCLC 18487850.^[2].
- Schinkel, Andreas** (16 September 2024). "Brückelnde Kuppel, teure Aufzüge: So kaputt ist das Neue Rathaus Hannover"^[2] [Crumbling dome, expensive elevators: This is how broken the New Town Hall in Hanover is] (in German). *Archived*^[2] from the original on 16 September 2024. Retrieved 16 September 2024.

New Town Hall (Munich)

Article Talk

From Wikipedia, the free encyclopedia

Coordinates: 48°8′15″N 11°54′32″E﻿ / ﻿48.13750°N 11.90889°E﻿ / 48.13750; 11.90889

The **New Town Hall** (German: *Neues Rathaus*) is a town hall that forms the northern part of **Marienplatz** in **Munich**, **Bavaria**, **Germany**. It hosts the city government including the **city council**, offices of the **mayors** and a small portion of the administration. In 1874 the municipality had left the **Old Town Hall** for its new domicile.

History

Inception and construction

The decision to construct a new building came due to the lack of space in the **Old Town Hall** and the adjoining, so-called "Lesser Town Hall" on **Petersberg** (destroyed in 1944, not reconstructed). In memory of the bourgeois high season during the **Gothic** period, the choice fell upon a **neo-Gothic** design, which allowed to implement an independent architectural accent in contrast to the buildings of the royal family.

The north side of the **Marienplatz** was chosen as the building site, where the house of the **Landstände** still stood which had been erected by the **Bavarian Duke** throughout the **Middle Ages** as a sort of representation of the opposing **Landstände**. The first section of the building in the eastern part of the **Marienplatz**, on the corner of **Dienersstrasse**, was the results of an idea competition won by **Georg Hauberisser** and carried out between 1867 and 1874. When it became clear that this new building would not be able to accommodate the entire administration, the city began purchasing all the properties on the **Dienersstrasse**, **Landschaftstrasse** and **Weinstrasse** adjacent to the Town Hall started in 1887. From 1889 to 1892, the section on the corner of **Dienersstrasse** and **Landschaftstrasse** was constructed.

In 1897, the Magistrate and municipal council decided to extend the buildings on the **Marienplatz** as well as the **Weinstrasse** and **Landschaftstrasse** to create a four-sided complex. For this, the entire area between the **Marienplatz** and **Landschaftstrasse** was used and on the other side, between **Weinstrasse** and **Dienersstrasse**. In 1898, the work for the extension began with the tower (**Rathausturm**), also under architect Georg von Hauberisser. In December 1905, the shell of the third building section was finished with the setting of the **keystone** on the **Rathausturm**. For the architectural design of the **Munich Rathausturm**, Hauberisser was clearly inspired by **Brussels' Town Hall**, whose 96-meter **Brabantine Gothic** tower was built by **Jan van Ruysbroeck** in the years 1449 to 1455.^[1] By the end of 1906, the offices were handed over. The façade area in the **Marienplatz** was then 98.5 meters long, of which 48 meters belong to the first construction section.^[2] Examples that were used for the design were the **Town Hall in Brussels** and the **City Hall in Vienna**.

20th century–present

The minimal damages to the New Town Hall that occurred during the air raids on Munich 1944, were rebuilt after the war. The portion constructed at the **Marienplatz** received an additional floor, which were hidden behind the neo-Gothic balustrade so that the building's image was preserved. The façade on the **Landschaftstrasse** was very simply restored. At the end of the 1990s, the New Town Hall was rebuilt and reconstructed identically, including the neo-Gothic ornaments, which crown the roof.

Dimensions and location

The building covers an area of 9159 m² having 400 rooms. The 100 meters long main facade towards the **Marienplatz** is richly decorated. It shows the **Guelf** Duke **Henry the Lion**, and almost the entire line of the **Wittelsbach** dynasty in **Bavaria** and is the largest princely cycle in a German town hall. The central monument in the center of the main facade between the two phases at **Marienplatz** above the guard house, is an equestrian statue of **Prince Regent Luitpold**. The bay of the tower contains statues of the first four **Bavarian Kings**.

The main facade is placed toward the square, while the back side is adjacent to a small park (**Marienhof**). The basement is almost completely occupied by a large restaurant called **Ratskeller**. On the ground floor, some rooms are rented for small businesses. Also located in the ground floor is the major official **tourist information**.

The first floor hosts a big balcony towards the **Marienplatz** which is used for large festivals such as football championships or for concerts during the **Weihnachtsmarkt**. Its main tower has a height of 85 m and is available for visitors with an elevator. On the top thrones the **Münchner Kindl**. The **Rathaus-Glockenspiel**, performed by an apparatus daily at 11am, 12pm and 5pm, is a **tourist attraction**.

Brussels' Town Hall was used as an architectural example for Munich's New Town Hall.^[1]

Relief of Munich's partner cities in the entrance hall of the New Town Hall

The New Town Hall's southern front

The location of the New Town Hall directly at Marienplatz

Detail of the front facade above the main entrance

Description

[edit]

Architectural design

[edit]

Conference Rooms – Hallways – Staircase

(a) Interleaved Multimodal Document Retrieval

(b) Text-only Document Retrieval

Figure 6: Retrieved documents across different document formats for document retrieval with a given multimodal query. (a): A document retrieved when represented leveraging interleaved multimodal contents within documents (ours). (b): A document retrieved when using only textual format