
Mr. Right: Multimodal Retrieval on Representation of Image with Text

Cheng-An Hsieh^{1*}, Cheng-Ping Hsieh^{2*}, Pu-Jen Cheng¹

¹National Taiwan University, ²UC San Diego

(*equal contribution)

r09944010@ntu.edu.tw, c2hsieh@ucsd.edu, pjcheng@csie.ntu.edu.tw

Abstract

Multimodal learning is a recent challenge that extends unimodal learning by generalizing its domain to diverse modalities, such as texts, images, or speech. This extension requires models to process and relate information from multiple modalities. In Information Retrieval, traditional retrieval tasks focus on the similarity between unimodal documents and queries, while image-text retrieval hypothesizes that most texts contain the scene context from images. This separation has ignored that real-world queries may involve text content, image captions, or both. To address this, we introduce Multimodal Retrieval on Representation of Image with Text (Mr. Right), a novel and comprehensive dataset for multimodal retrieval. We utilize the Wikipedia dataset with rich text-image examples and generate three types of text-based queries with different modality information: text-related, image-related, and mixed. To validate the effectiveness of our dataset, we provide a multimodal training paradigm and evaluate previous text retrieval and image retrieval frameworks. The results show that proposed multimodal retrieval can improve retrieval performance, but creating a well-unified document representation with texts and images is still a challenge. We hope Mr. Right allows us to broaden current retrieval systems better and contributes to accelerating the advancement of multimodal learning in the Information Retrieval.

1 Introduction

Recent advancements in the field of digital media have resulted in a surge of interest in multimodal learning. Multimodal learning aims to learn well-unified representations from different modalities such as language, vision, or audio and projects them into a common low-dimensional space. For example, visual question answering needs an understanding of both vision and language [1, 2, 3]; video highlight detection exploits video and audio features to identify the exciting moments [4, 5]; emotion recognition requires a fusion of spoken words, facial expressions, and voice [6, 7].

In information retrieval, the conventional retrieval tasks focus on unimodal learning, including text-to-text [8, 9] and image-to-image [10, 11, 12] retrieval. Both the texts and images contain comprehensive information, requiring the model to compute semantic representations of a single modality and match the unimodal document-query pair. Prior works [3, 13, 14, 15, 16, 17, 18] have improved retrieval performance and assisted users in searching for the requested documents. While beneficial, these tasks suffer from a significant key limitation: the text representations and the image features exist in their own spaces. In real-world applications, users may take texts or images as the queries to retrieve relevant data of the other modality. Therefore, cross-modal retrieval has attracted considerable attention from researchers recently.

Cross-modal retrieval aims to retrieve a relevant unimodal document from another modality of a query. Image-text retrieval is a fundamental challenge in cross-modal retrieval, and most existing



Figure 1: Overview of Mr. Right compared to unimodal and cross-modal retrieval. Unimodal retrieval (text-to-text retrieval) uses text-related queries to search for document texts, while cross-modal retrieval (text-to-image retrieval) uses image-related queries to find document images. Mr. Right utilizes each query and proposed mixed queries to retrieve documents with texts and images.

methods [19, 20, 21, 22] train their models on the COCO [23] and Flickr30K [24] datasets. These datasets include images and their captions. However, unlike text-to-text retrieval, the texts in these datasets only have image related information. Nowadays, many online materials and documents may have both texts and images, such as Wikipedia, news, blogs, social media posts, and commercial websites. Also, considering most people utilize text-based queries to search multimedia documents on a search engine, a searching query can be keywords, captions, or both. The retrieval frameworks should be able to deal with the above text-based query with different modality information.

We introduce Mr. Right as shown in Figure 1, a new comprehensive and challenging retrieval dataset, which provides text-image documents and text-based queries with different modality information that are text-related, image-related, or mixed. For documents, we collect from the Wikipedia-based Image Text Dataset [25], including paragraphs and photos. For queries, based on our knowledge, there is no existing dataset containing our proposed three types of queries. Therefore, we hire Amazon Mechanical Turk (AMT) workers to construct total 3k annotated queries for each type. Moreover, we also provide 350k auto-generated queries for model pre-training before learning annotated queries. In our training paradigm, we introduce document-query contrastive learning (DQC) and document-query matching (DQM) to fuse text and image features into multimodal representations.

Finally, we propose a full-ranking benchmark for Mr. Right. Similar to TR datasets [23], our full-ranking evaluation contains three types of annotated queries that models have to retrieve the most relevant document from all corpus. We create the benchmark based on our multimodal framework with the comparison to prior text-to-text retrieval (TR) and text-to-image retrieval (IR) baselines with human performance. Our results show that multimodal retrieval (MR) can perform better with the help of extended information from different modalities. Interestingly, it is a balance to incorporate these modalities into unified representations. However, it is still challenging to achieve comparable retrieval performance to the human benchmark.

With Mr. Right, we take a significant step toward establishing a novel benchmark for evaluating the capabilities of multimodal retrieval systems. To the best of our knowledge, Mr. Right is the first multimodal retrieval dataset that explores multimodal documents and text-based queries with different modality information, and it is open-source to welcome methods of all kinds.

Dataset	Task	Source	Train	Dev		Test		
			#Pairs	#Pairs	#Query	#Pairs	#Query	#Document
MS MARCO [8]	Text Retrieval	Misc.	532,761	—	—	—	6,980	8,841,823
TREC-NEWS [26]	Text Retrieval	News	—	—	—	—	57	594,977
HOTPOTQA [27]	Question Answering	Wikipedia	170,000	—	5,447	—	7,405	5,233,329
NQ [9]	Question Answering	Wikipedia	132,803	—	—	—	3,452	2,681,468
FEVER [28]	Fact Checking	Wikipedia	140,085	—	6,666	—	6,666	5,416,568
SOP [29]	Image-to-Image Retrieval	Products	59,551	—	—	—	11,316	60,502
CUB-200-2011 [30]	Image-to-Image Retrieval	Birds	5,864	—	—	—	100	5,924
ROxford [31]	Image-to-Image Retrieval	Landmark	—	—	—	—	70	4,993
Ukbench [32]	Image-to-Image Retrieval	Misc.	—	—	—	—	2,550	10,200
MS COCO [23]	Image-Text Retrieval	Misc.	118,000	5,000	—	5,000	—	—
Conceptual Captions [33]	Image-Text Retrieval	Misc.	~3,300,000	2,800	—	2,300	—	—
Flickr30K [24]	Image-Text Retrieval	Flickr	30,000	1,000	—	1,000	—	—
M5Product [34]	Multimodal Retrieval	Products	4,423,160	—	—	—	1,991	117,858
Mr. Right (ours)	Multimodal Retrieval	Wikipedia	351,979 / 1,000	—	—	—	2,047	806,357

Table 1: Statistics of datasets in retrieval tasks. Text-to-text retrieval datasets and image-to-image retrieval datasets contain queries and documents in a large search pool size. Image-text retrieval datasets include pairs of data. M5Product consists of multi-modality products as documents and queries. Unlike M5Product, Mr. Right is the multimodal retrieval dataset exploring text-based query with different modality information. It involves around 35k auto-generated training pairs, 1k human-annotated pairs for fine-tuning, and 2k annotated pairs for testing. All pairs in our dataset are composed of one document and three types of queries.

2 Related Work

In this section, we describe previous retrieval datasets and explain how the existing methods employ transformer-based neural networks [35] to learn the representations of different modalities. Table 1 shows an overview of the retrieval datasets.

Retrieval dataset Most previous retrieval datasets only consider single modality (texts or images) documents without a unified representation among multiple domains. We categorize these datasets into unimodal and cross-modal learning. Text retrieval (TR) and image-to-image retrieval are the fundamental challenges in unimodal learning. Previous TR datasets [8, 9, 27, 26, 28] comprise a large corpus of text-based documents and related queries. They collect documents from different sources, such as Wikipedia [27, 28], news [26], and online articles [8]. These sources involve diverse and generalized domain knowledge, reflecting real-world situations when users search from an extensive database. To ensure the quality of queries, some works collect the queries from searching logs [8, 9] or hire crowd workers to generate annotations [27, 28]. Similarly, the existing image-to-image datasets [29, 30, 31, 32] include several categories of images in the same domain, such as products, birds, and landmarks. Major works [36, 37, 29] randomly sample images from each category as queries, while the remaining images are the documents. As shown in Table 1, these unimodal datasets have more documents than queries, showing the challenge of ranking large numbers of documents. For cross-modal learning, the existing image-text retrieval (IR) datasets [33, 23, 24] contain images and their captions. Major works harvest their images and captions from the web. They develop a pipeline to extract, filter, and transform their captions. In this task, the number of images equals the captions, meaning documents and queries are in pairs. Unlike the unimodal tasks having a large size of documents, the evaluation of cross-modal performs on a small size of document-query pairs. Different from single modality documents, recent work [38] has proposed an E-commerce product multimodal retrieval dataset that contains data more than two modalities.

Retrieval model Due to the superior performance of contextualized representations in transformer-based models [35], self-attention-based architectures have become the model of choice in natural language processing (NLP) and computer vision (CV). In unimodal retrieval, the existing methods [39, 40, 41, 42] of text-to-text retrieval employ transformers to encode queries and documents into vector representations and compute their similarity. Also, the Vision Transformers [43] reduce the time-consuming process of extracting region features, and the later works [44, 45, 46] attain excellent results in image-to-image retrieval. In cross-modal retrieval, Vision-and-Language Pre-training (VLP) models have improved performance on IR tasks. The recent CLIP [47] and ALIGN [48] utilize contrastive learning to align the unimodal representations of image-text pairs. Other VLP

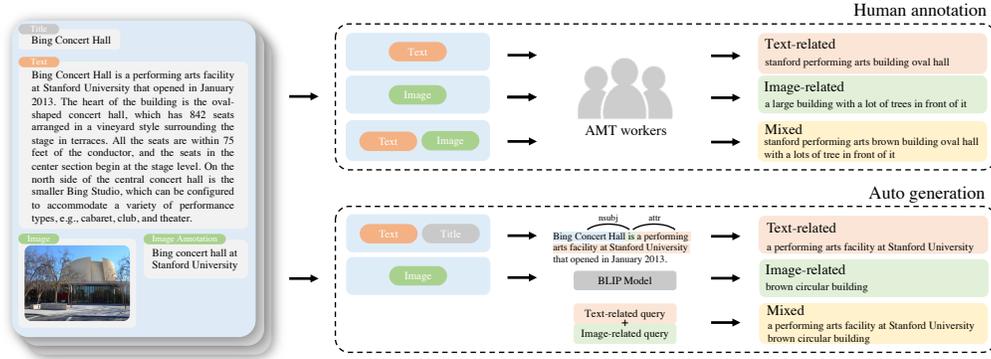


Figure 2: Overview of Mr. Right’s data collection process. Human-annotated queries are created by AMT workers following our annotation instructions, while auto-generated queries are constructed with dependency parsing, image captioning, and concatenation.

methods (e.g. METER [19], ALBEF [20], LXMERT [21], UNITER [22]) perform multimodal fusion to produce the joint representations of text-image pairs. This bridges the semantic gap between visual and textual features in texts and images.

3 The Mr. Right Dataset

Mr. Right aims to construct a new dataset for multimodal retrieval tasks. The dataset focuses on two components: (1) Multimodal documents consist of different modality information, including texts and images; (2) Text-based queries involve text content, image captions, or both. Mr. Right collects documents and annotated/generated queries by extracting, labeling, and filtering.

3.1 Data Collection

Wikipedia-based document To generate multimodal documents with diverse knowledge domains, we gather paragraphs and photos from the Wikipedia-based Image Text Dataset [25], which consists of 37.6 million entity-rich image-text pairs with 11.5 million unique images. The original dataset includes 108 languages, and we only keep 1.5 million English data for simplification. We process this dataset with three steps: (1) image filtering eliminates images with invalid download links, corrupted content, and non-JPEG/PNG images; (2) text filtering removes repeated pages and deletes contents without mentioning the title, ensuring that the document is related to the subject; (3) text reduction extracts the first paragraph as the document to avoid high memory and computational requirements. The whole Wikipedia article can be very long, and we find that usually the first paragraph contains a brief introduction. After the filtering, there are 806,357 of multimodal documents remaining.

Human-annotated query A query may relate to the document text, document image, or both. To generate three types of queries based on multimodal documents, we hire crowd-workers with qualifications from Amazon’s Mechanical Turk (see Appendix B.2). The annotators are shown a text, a image, and both successively to come up with queries based on first impression. To ensure the consistency and quality of annotations, we give the following guidelines: (1) Word count limitation. Each query should be between 10 to 100 words to ensure the query involves enough information. (2) No title. Annotators should avoid including the title because it is the result that users want to retrieve. (3) No copied phrases in the passage. A real-world query may involve ambiguous meaning or terms, so it is better to paraphrase the sentences. (4) Requirement of adjectives and nouns in the image query. Instead of only phrasing the objects in the image, we hope crowd-workers describe the details, such as colors, shapes, and actions. The annotated query examples can be found in Figure 2.

Auto-generated query Annotating queries for the whole multimodal dataset can be time-consuming. To address this, we develop a pipeline that extracts snippets of the document texts

	SECS	Vocab. Size		Avg. Word Lengths	Top-3 NER Tag
<i>Human-annotated</i>					
Text-related	3,047	8,290	9.54	GPE (32.14%), DATE (14.62%), NORP (10.70%)	
Image-related	3,047	2,609	8.26	CARDINAL (59.49%), NORP (24.65%), PERSON (3.68%)	
Mixed	3,047	8,752	13.58	GPE (31.66%), DATE (13.49%), NORP (11.19%)	
<i>Auto-generated</i>					
Text-related	351,979	55,928	4.95	NORP (41.54%), GPE (24.82%), ORG (8.77%)	
Image-related	351,979	13,464	9.83	CARDINAL (51.58%), GPE (34.62%), NORP (4.79%)	
Mixed	351,979	62,440	14.78	NORP (39.41%), GPE (25.48%), ORG (7.90%)	

Table 2: Analysis of annotated and generated queries in Mr. Right. Top-3 NER tags show the top entities classified in the whole corpus. GPE means geopolitical entities, NORP represents affiliations, and CARDINAL are numerical values.

as the text-related queries and generates captions of the document images as the image-related queries. According to Figure 2, Wikipedia content generally has a specified format. The first sentence begins with the title and a short introduction, followed by the details. Therefore, we utilize dependency parsing using spaCy API and detect the dependent verb of the title in the first sentence. Then we take the adjectives and nouns after the verb as the text-related query and remove the adjectives in the first sentence as the document. This ensures that models have to learn the information from the remaining texts. For the image-related queries, each image in the original Wikipedia dataset has its own annotation. However, many annotations relate to the document title instead of the image content. In addition, some of the annotations contain proper names, which is difficult to learn from the scene context. To address this, we replace them with generated image captions based on BLIP [49] that outperforms a variety of methods on vision and language tasks. As shown in Figure 2, the caption generated from the model is closer to the image content, and it can be the query that human uses to search for visual information. To generate the mixed queries, we concatenate the former two queries, which contain text information and image context respectively.

3.2 Annotated Query Validation

Rule-based filtering A well-defined query should have multiple part-of-speech (POS) tags. Therefore, annotated query candidates without nouns or with only one noun lacking adjectives or verbs are discarded. Since queries directly copied from the documents are trivial for retrieval, we drop text-related and mixed query candidates that highly overlap with document texts. We remove the queries with the longest-common-substring (LCS) length larger than 40 and the ratio (divided by query length) larger than 0.6. As for image-related queries, we take out the candidates that include additional knowledge more than the document image context. In other words, we filter out the queries containing proper names, such as a particular person’s identity or locations. The above three filters discard around 10% of the candidates.

Human filtering In our task, each query should correspond to a unique document. To ensure uniqueness, we utilize text retrieval model BM25 and image retrieval model CLIP to search relevant documents given text-related and image-related queries, respectively. For simplicity, we retrieve top-10 relevant candidates from the whole multimodal documents and prioritize to examine the queries without unique document pairs, i.e., close ranking scores for different documents. After filtering out these queries, we efficiently validate whether the semantic meaning between the query and the correct document is unrepeatable. After our validation, there are 25% of query sets discarded and remain 3,047 annotated query sets.

After finishing the collection and validation stage, our dataset contains 806,357 multimodal documents, 351,979 auto-generated query sets, and 3,047 human-annotated query sets. Each query set is mapped to one document and contains three types of proposed queries. We further split 1k human-annotated sets for fine-tuning and 2k for testing as shown in Table 1.

3.3 Dataset Analysis

Quantitatively analysis Table 2 presents statistics for our annotated and generated queries showing vocab sizes, average lengths, and top-3 named entity recognition (NER) tags. Each type of query has the same amount. Based on the vocabulary size, we can find that text-related queries have more diverse words than image-related queries, indicating that our image descriptions are composed of

Query type	Properties	A%	G%	Example(s)
Text-related	Paraphrase	53	0	A: Seoul based artist, animator famous for eyedolls [†]
	Keyword Extraction	23	0	A: Railway station in Schwelm in Rhine Westphalia [‡]
	Duplication	24	100	A: Northern Alabama pygmy sunfish* G: A species of pygmy sunfish*
Image-related	One object	28	19	A: Small swimming gray and white fish* G: fish that is swimming*
	Multiple objects	72	81	A: Girl with straight brown and pink hair in white dress on white background [†] G: A woman with pink hair and a white shirt [†]
Mixed	Fusion	47	0	A: A German railway station which runs through thick green forest [‡]
	Concatenation	53	100	A: Small swimming gray and white fish Northern Alabama pygmy sunfish* G: A species of pygmy sunfish a fish that is swimming in some water*

Table 3: Properties distribution of annotated queries (A) and generated queries (G). Each query type has 100 samples, and we manually categorize each query into different properties. We find annotated queries have more diverse properties than generated ones. The corresponding documents of provided examples are listed in Appendix B.4

limited illustrative words. Furthermore, texts contain more sparse information than images, resulting in the need for longer word lengths of annotated text-related queries. The top-3 NER tags suggest that human favors distinct entities such as countries and dates, while our generated approach mainly focuses on affiliations in the text-related queries. For image-related queries, human and the BLIP model prefer to describe the number of objects.

Qualitatively analysis We heuristically identify query properties covered in the dataset to recognize the difference between human-annotated and auto-generated queries. We randomly sample 100 queries from the three types of queries and present the results in Table 3. As can be seen, we split the properties of text-related queries into three categories. Paraphrase means that queries involve different words and sentence clause structure from the documents; keyword extraction indicates that queries only include important terms; duplication means queries are the reorganized document phrases. Considering efficiency, we copy the snippets of document texts to generate text-related queries. For the image-related queries, some annotators may focus on describing the most conspicuous object, while others include multiple objects and their adjectives. Auto-generation produces more image-related queries with multiple objects. It may be because the BLIP model learns captioning from many data and prefers to describe the image details. The difference between annotated and generated mixed queries is also apparent. Human may fuse the image descriptions and text content or concatenate them with prepositions. Our generated mixed queries only rely on the concatenation. In our study, annotated queries contain more diverse types of queries and well-unified sentence structures, but the annotation process is time-consuming and expensive. Our auto-generated queries ensure efficiency, and the experiment results show that these queries are effective.

3.4 Benchmark

To simulate the real-world retrieval problems, we create Mr. Right’s benchmark with the whole corpus of 800k Wikipedia documents. This full-ranking setting guarantees that the model can handle large numbers of multimodal documents. Further, search queries may have text keywords or image descriptions, so we present three retrieval tasks with the corresponding queries: text-related, image-related, and mixed. See Appendix C.1. for more task details.

3.5 Evaluation Metrics

Retrieval tasks might be precision-focused or recall-focused, depending on the requirements of real-world applications. In Mr. Right, documents and queries are binary relevant, and retrieving relevant documents from our large corpus of different modalities is challenging. Following previous image-text retrieval tasks [19, 20, 21, 22], we report recall@ k as our performance metric. Further, considering the rank of documents, we also utilize MRR (Mean Reciprocal Rate) as our binary rank-aware metric, a general standard in text retrieval tasks. In our experiments, we compute recall with $k = 1, 5, 10$ and MRR@10 for all models and assess their performance.

4 Multimodal Retrieval

With Mr. Right, the next step is to set up the retrieval task based on the multimodal documents and text-based queries. We illustrate our retrieval formulation (Section 4.1) and model architecture (Section 4.2). Then we describe our two training objectives (Section 4.3).

4.1 Retrieval Formulation

Given a document D with a paragraph text D_T and an image D_I , we use a multimodal encoder to fuse D_T and D_I into a single fixed-size multimodal vector representation R_d . Also, we encode a text-based query Q_T into a fixed-size vector representation with our text encoder. To establish our retrieval task, we need to encode all the documents $\{(D_T^1, D_I^1), (D_T^2, D_I^2), \dots, (D_T^N, D_I^N)\}$ and queries $\{Q_T^1, Q_T^2, \dots, Q_T^M\}$ into $\{R_d^1, R_d^2, \dots, R_d^N\}$ and $\{R_q^1, R_q^2, \dots, R_q^M\}$ respectively. With these representations, we compute the cosine similarity scores between documents and queries and find the most similar document for each query. In this scenario, we can build offline indexing for document representations and compute query representations online for real-world applications.

4.2 Model Architecture

Document (Multimodal) Encoder To encode both document texts and images into unified multimodal representations, we leverage previous pre-trained VLP models for initialization in our framework. These models have learned a common low-dimensional space to embed vision and language features. In these models, we have a vision encoder (*e.g.* CNNs or vision transformers [43]) and a text encoder (*e.g.* BERT [50] or RoBERTa [51]) to extract modality-specific features. Then we have a fusion module (*e.g.* co-attention or merge-attention [19]) to integrate both features into a unified feature. Therefore, we can view these VLP models like a black box multimodal encoder E_M to output size-variant multimodal document features F_d whose size depends on the length of input texts D_T and the dimension of images D_I . To derive a single fixed-size representation R_d for each document, we simply average the size-variant document features F_d .

Query (Text) Encoder Since our queries are text-based Q_T , we create a query encoder E_q and share the parameters from the text encoder of our multimodal encoder. This ensures the text representations of queries are similar to document texts. Like a multimodal encoder, we take the average of the query features F_q to obtain query representations R_q .

4.3 Training Objectives

In this section, we introduce document-query contrastive learning (DQC) and document-query matching (DQM) to project document and query representations into the same space.

Document-Query Contrastive learning Contrastive learning has been widely used to train on VLP models [47, 48, 20] which can increase the similarity scores between parallel pairs. With document and query representation R_d and R_q , we learn two projection functions f_d and f_q with a fully-connected layer to map their representations into the same space. Then we calculate the cosine similarities between document and query pairs in a training batch. The matched pairs are positive while all other pairs are negative. Based on the pairs, we minimize the contrastive loss L_{dqc} like in-batch cross-entropy loss. To organize in-batch negatives more effectively, we keep two queues [20] to store the most recent K representations, enlarging the amounts of negative samples per batch.

Document-Query Matching To further learn a fine-grained similarity of pair documents and queries, we build a binary classifier C to predict whether the output features of document encoder F_d and query encoder F_q is matched. Specifically, we use a 6-layer transformer model and insert a special token $[CLS]$ at the head of the input sequence to obtain global information. Then we employ a linear classifier on this token followed by softmax to predict a two-class label and compute the matching loss L_{dqm} according to binary cross-entropy loss. Motivated by ALBEF [20], we sample online hard negative pairs for each document and query from contrastive similarity distribution. In addition to real-world negative documents, we produce two pseudo negative documents by combining the positive document and the sampled negative document into pairs of a positive image and a negative

	Method	text-related query				image-related query				mixed query			
		MRR@10	R@1	R@5	R@10	MRR@10	R@1	R@5	R@10	MRR@10	R@1	R@5	R@10
TR	RoBERTa (ZS)	0.0	0.0	0.1	0.2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	DiffCSE (ZS)	32.0	24.5	42.6	49.8	0.2	0.1	0.3	0.4	20.2	14.9	27.9	32.7
	SBERT (ZS)	35.6	27.3	46.5	54.4	1.0	0.5	1.5	2.2	28.2	21.4	37.3	43.9
	RoBERTa (FT)	25.5	18.8	34.6	41.9	0.4	0.2	0.9	1.5	24.0	17.1	33.9	40.7
	DiffCSE (FT)	33.3	26.3	42.8	51.0	0.6	0.3	0.9	1.7	30.1	22.8	40.2	47.7
	SBERT (FT)	47.7	38.7	60.0	68.0	1.1	0.5	1.9	2.6	39.3	31.0	50.2	58.8
IR	CLIP (ZS)	2.6	1.6	3.7	5.7	5.1	3.3	7.4	10.2	7.0	4.1	10.8	14.7
	ALBEF (ZS)	0.5	0.3	0.6	1.2	3.7	2.2	5.6	8.0	1.8	0.8	2.9	4.5
	CLIP (FT)	3.3	1.9	5.0	7.5	5.7	3.6	8.0	11.9	8.3	4.8	12.7	18.6
	ALBEF (FT)	0.9	0.5	1.5	2.1	6.2	4.0	9.3	12.9	4.1	2.1	6.5	9.4
MR	SBERT (FT) + ALBEF (FT)	48.0	39.0	60.7	67.9	7.2	4.4	10.9	15.3	50.3	41.6	61.8	70.3
	Our (METER)	38.8	27.4	54.8	64.3	12.0	7.2	18.3	25.3	42.9	31.6	58.2	69.2
	Our (ALBEF)	44.4	32.3	61.1	72.4	4.2	2.1	6.2	10.8	50.7	38.0	67.2	78.2
	Our (ViLT)	26.8	16.4	40.6	53.4	15.7	8.6	23.0	33.0	45.4	33.0	62.7	73.5

Table 4: Retrieval performance across three benchmark tasks with different types of queries. We compare MRR and recall@ k among baselines and our proposed models on TR, IR and MR. ZS is zero-shot, and FT is fine-tuned.

text and vice versa. Hence, for a query, we have a positive document, a sampled negative document, and two pseudo negative documents.

5 Experiments

5.1 Dataset

Mr. Right has both auto-generated and human-annotated queries. We first pre-train our models on the 350k auto-generated document-query pairs to learn the multimodal representations. Further, we fine-tune our learned model on the human-annotated 1k training pairs with 10% as our validation set.

5.2 Baselines

We compare our proposed multimodal retrieval framework with TR and IR baselines. We only collect existing dense retrieval [52] approaches for a fair comparison. Additionally, we develop the MR baseline with the ensemble of TR and IR baselines. Text retrieval models only consider the document texts; image retrieval models only focus on the document images; multimodal retrieval models perceive both document texts and images. All the baseline models are described in Appendix C.2.

5.3 Experiment Setup

We train our framework using existing VLP models, including METER [19], ALBEF [20], and ViLT [53] to make use of their multimodal pre-trained weights. The pre-training process lasts for 40 epochs and fine-tunes for 20 epochs on 8 NVIDIA V100 GPUs. Our optimizer is AdamW with a weight decay of 0.02, and the learning rate is warmed-up to 5×10^{-5} in the first epoch and decayed to 1×10^{-7} following the scheduler. Also, we set up a 0.5 gradient clipping value and 9,600 queue size for DQC. For image augmentation, we use random-crop of size 288×288 or 384×384 depending on the pre-trained VLP models and apply RandAugment [54]. For texts, we truncate our max length for queries with 40 and documents with 128. In order to simulate real-world user queries, we randomly select text-related, image-related, or mixed queries during training.

5.4 Results and Analysis

Compared to TR/IR We present the retrieval results in Table 4. We compare our method against TR/IR models and discuss the performance difference across three query types. The table shows that TR and IR have difficulties responding to the opposite queries. TR obtains worse results for image-related queries while IR is vice versa. This may be because their documents only contain unimodal information with either texts or images. In contrast, MR shows the ability to mitigate this problem. It achieves comparable performance as TR on text-related queries and scores higher than IR on image-related queries. This improvement indicates that MR can perform better due to the extended information from different modalities. Further, when queries are mixed, MR exploits the advantage of multimodal representations and achieves superior performance compared to TR and IR.

Multimodal representation We integrate fine-tuned SBERT and ALBEF as an ensemble MR model that shows a comparable performance to our MR models. Although incorporating TR and IR models can perform well among different types of queries, the vector size of document representations for multiple modalities increases linearly, and we need to define the best weighted combination of their output scores. In contrast, our proposed multimodal representation can unify multiple domain information into a standard size feature. To understand the multimodal representations, we compute Grad-CAM visualizations (see Appendix C.3) on the attention maps of document texts and images given different types of queries. The attention heat is highly correlated to where human would look to match the corresponding query. In Table 4, we find a trade-off of our framework to deal with text-related and image-related queries simultaneously. Comparing MR with different backbone VLMs, the performance is debated between the two queries. This may come from the limited size of our unified representation. We cannot include all the document text and image information together but a balance between them.

Human evaluation Besides model performance, we also present human evaluation results compared to our MR models in Table 5. We randomly select 50 samples for each query type. To efficiently retrieve related documents for human evaluation, we utilize our MR model (METER) to obtain the top 3 relevant candidates with the correct document and construct a four-choice question with one correct answer. Table 5 shows humans get 89.3% accuracy. It validates the reliability of our dataset, but it also shows there is room for improvement of models on Mr. Right. It may be because human can understand various query properties in Table 3, extract the crucial text content, or perceive image scene context in detail. To see the retrieval results difference between our models and human, we provide failed examples in Appendix C.4. Also we provide the performance comparison of our auto-generated and human-annotated queries in Appendix C.5.

	Accuracy %
Our (METER)	30.0
Our (ALBEF)	26.0
Our (ViLT)	26.7
Human	89.3

Table 5: Human evaluation on 150 random sampled queries.

6 Conclusion

In this paper, we propose Mr. Right, a multimodal retrieval dataset for information retrieval. Mr. Right covers three types of text-based search queries with different modality information, including text-related, image-related, and mixed, to simulate real-world search situations. Further, our dataset provides documents with texts and images to develop multimodal representation. We build our end-to-end multimodal retrieval model for Mr. Right to unify features across modalities. Compared to the previous text and image retrieval frameworks, multimodal retrieval shows improvements on different queries and points out the balance between modalities. However, current multimodal models still have a significant gap to human performance, showing the potential of Mr. Right as a challenge in multimodal retrieval. We believe Mr. Right can breathe new insights into information retrieval for more robust retrieval systems.

7 Limitations and Future Work

In Mr. Right, we only consider text-based queries, which may limit the search modalities from users. We can expand our dataset with additional domain queries and documents such as images, audio, and video. Further, Mr. Right focuses on the materials in Wikipedia. We can explore other sources such as news, blogs, or commercial websites. Mr. Right is a preliminary attempt to explore multimodal retrieval, and there are still challenges we need to analyze and study in future work.

References

- [1] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.
- [2] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering.

- In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [3] Omar Khattab and Matei Zaharia. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pages 39–48, 2020.
 - [4] Taivanbat Badamdorj, Mrigank Rochan, Yang Wang, and Li Cheng. Joint visual and audio learning for video highlight detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8127–8137, 2021.
 - [5] Qinghao Ye, Xiyue Shen, Yuan Gao, Zirui Wang, Qi Bi, Ping Li, and Guang Yang. Temporal cue guided video highlight detection with low-rank audio-visual fusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7950–7959, 2021.
 - [6] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42(4):335–359, 2008.
 - [7] AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Multimodal language analysis in the wild: CMU-MOSEI dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2236–2246, Melbourne, Australia, July 2018. Association for Computational Linguistics.
 - [8] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. Ms marco: A human generated machine reading comprehension dataset. In *CoCo@ NIPS*, 2016.
 - [9] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
 - [10] James Philbin, Ondrej Chum, Michael Isard, Josef Sivic, and Andrew Zisserman. Object retrieval with large vocabularies and fast spatial matching. In *2007 IEEE conference on computer vision and pattern recognition*, pages 1–8. IEEE, 2007.
 - [11] David Nister and Henrik Stewenius. Scalable recognition with a vocabulary tree. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, volume 2, pages 2161–2168. Ieee, 2006.
 - [12] Herve Jegou, Matthijs Douze, and Cordelia Schmid. Hamming embedding and weak geometric consistency for large scale image search. In *European conference on computer vision*, pages 304–317. Springer, 2008.
 - [13] Leonid Boytsov and Eric Nyberg. Flexible retrieval with NMSLIB and FlexNeuART. In *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*, pages 32–43, Online, November 2020. Association for Computational Linguistics.
 - [14] Hang Zhang, Yeyun Gong, Yelong Shen, Weisheng Li, Jiancheng Lv, Nan Duan, and Weizhu Chen. Poolingformer: Long document modeling with pooling attention. In *International Conference on Machine Learning*, pages 12437–12446. PMLR, 2021.
 - [15] Devendra Singh, Siva Reddy, Will Hamilton, Chris Dyer, and Dani Yogatama. End-to-end training of multi-document reader and retriever for open-domain question answering. *Advances in Neural Information Processing Systems*, 34, 2021.
 - [16] Liang Zheng, Yi Yang, and Qi Tian. Sift meets cnn: A decade survey of instance retrieval. *IEEE transactions on pattern analysis and machine intelligence*, 40(5):1224–1244, 2017.
 - [17] Albert Gordo, Jon Almazán, Jerome Revaud, and Diane Larlus. Deep image retrieval: Learning global representations for image search. In *European conference on computer vision*, pages 241–257. Springer, 2016.

- [18] Artem Babenko, Anton Slesarev, Alexandr Chigorin, and Victor Lempitsky. Neural codes for image retrieval. In *European conference on computer vision*, pages 584–599. Springer, 2014.
- [19] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Zicheng Liu, Michael Zeng, et al. An empirical study of training end-to-end vision-and-language transformers. *arXiv preprint arXiv:2111.02387*, 2021.
- [20] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in Neural Information Processing Systems*, 34, 2021.
- [21] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019.
- [22] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer, 2020.
- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [24] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649, 2015.
- [25] Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. Wit: Wikipedia-based image text dataset for multimodal multilingual machine learning. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2443–2449, 2021.
- [26] Ian Soboroff, Shudong Huang, and Donna Harman. Trec 2018 news track overview. In *TREC*, 2018.
- [27] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.
- [28] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: a large-scale dataset for fact extraction and verification. *arXiv preprint arXiv:1803.05355*, 2018.
- [29] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4004–4012, 2016.
- [30] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [31] Filip Radenović, Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondřej Chum. Revisiting oxford and paris: Large-scale image retrieval benchmarking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5706–5715, 2018.
- [32] Xiaoyu Wang, Ming Yang, Timothee Cour, Shenghuo Zhu, Kai Yu, and Tony X Han. Contextual weighting for vocabulary tree based image retrieval. In *2011 International Conference on Computer Vision*, pages 209–216. IEEE, 2011.
- [33] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018.

- [34] Xiao Dong, Xunlin Zhan, Yangxin Wu, Yunchao Wei, Michael C Kampffmeyer, Xiaoyong Wei, Minlong Lu, Yaowei Wang, and Xiaodan Liang. M5product: Self-harmonized contrastive learning for e-commercial multi-modal pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21252–21262, 2022.
- [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [36] Fuwen Tan, Jiangbo Yuan, and Vicente Ordonez. Instance-level image retrieval using reranking transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12105–12115, 2021.
- [37] Elias Ramzi, Nicolas Thome, Clément Rambour, Nicolas Audebert, and Xavier Bitot. Robust and decomposable average precision for image retrieval. *Advances in Neural Information Processing Systems*, 34, 2021.
- [38] Xiao Dong, Xunlin Zhan, Yangxin Wu, Yunchao Wei, Michael C Kampffmeyer, Xiaoyong Wei, Minlong Lu, Yaowei Wang, and Xiaodan Liang. M5product: Self-harmonized contrastive learning for e-commercial multi-modal pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21252–21262, 2022.
- [39] Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. Colbertv2: Effective and efficient retrieval via lightweight late interaction. *arXiv preprint arXiv:2112.01488*, 2021.
- [40] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- [41] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense text retrieval. *arXiv preprint arXiv:2007.00808*, 2020.
- [42] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. Rocketqa: An optimized training approach to dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2010.08191*, 2020.
- [43] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [44] Alaaeldin El-Nouby, Natalia Neverova, Ivan Laptev, and Hervé Jégou. Training vision transformers for image retrieval. *arXiv preprint arXiv:2102.05644*, 2021.
- [45] Tao Li, Zheng Zhang, Lishen Pei, and Yan Gan. Hashformer: Vision transformer based deep hashing for image retrieval. *IEEE Signal Processing Letters*, 29:827–831, 2022.
- [46] Yongbiao Chen, Sheng Zhang, Fangxin Liu, Zhigang Chang, Mang Ye, and Zhengwei Qi. Transhash: Transformer-based hamming hashing for efficient image retrieval. *arXiv preprint arXiv:2105.01823*, 2021.
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [48] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021.

- [49] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. *arXiv preprint arXiv:2201.12086*, 2022.
- [50] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [51] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [52] Daniel Gillick, Alessandro Presta, and Gaurav Singh Tomar. End-to-end retrieval in continuous space. *arXiv preprint arXiv:1811.08008*, 2018.
- [53] Wonjae Kim, Bokyoung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning*, pages 5583–5594. PMLR, 2021.
- [54] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 702–703, 2020.
- [55] Yung-Sung Chuang, Rumen Dangovski, Hongyin Luo, Yang Zhang, Shiyu Chang, Marin Soljačić, Shang-Wen Li, Wen-tau Yih, Yoon Kim, and James Glass. Diffcse: Difference-based contrastive learning for sentence embeddings. *arXiv preprint arXiv:2204.10298*, 2022.
- [56] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019.

A Supplementary Materials for Mr. Right

We provide the following detailed sections and materials that complement the discussions in the main paper. Code and dataset are available at <https://github.com/hsiehjackson/Mr.Right>

- Establishment of the datasheet for Mr. Right in Appendix B
- Details of the evaluation process in Appendix C.
- Designs of the proposed model in Appendix D.
- Confirmations of the data license in Appendix E.
- Maintenance of Mr. Right in Appendix F.

B Datasheets

B.1 Motivation

Information retrieval is a fundamental and essential challenge in real-world applications. In the past, researchers focused on unimodal retrieval because previous datasets only included data with a single modality, such as text-to-text and image-to-image retrieval datasets. They design robust and effective frameworks to improve the performance of these retrieval tasks. However, humans perceive the world with different modalities, such as language, vision, or audio. Due to multimedia development, humans have begun to utilize one modality to search for another modality. For example, image-text retrieval is a challenge in which models need to learn a common representation between images and texts and retrieve the most relevant documents. Further, sometimes we may need to combine different modalities and understand the meaning together. To accelerate the advancement of retrieval on multimodal learning, we propose Mr. Right, which contains multimodal documents and three types of text-based queries according to the real-world context. It has 806,357 multimodal documents, 351,979 auto-generated queries, and 3,047 human-annotated queries for each type.

B.2 Collection Process

Multimodal document We construct Mr. Right based on the Wikipedia-based Image Text (WIT) Dataset [25]. The original dataset includes Wikipedia articles and Wikipedia image links in 108 languages. Each article has a page title, a page description, and a reference image description. The dataset has filtered the image-text pairs based on effective restrictions, such as text length, image size, and image format. However, Wikipedia updates its content frequently, some image URLs are outdated, and some pages have different versions. Therefore, we create our pipeline to filter WIT and obtain the multimodal documents. The process is explained in the following:

- Download Wikipedia CSV file [25] and keep English articles with titles in the content. There are about 1,479,330 English documents. Download the images using the Python *multiprocessing* and *urllib2* module. During the downloading, we find that some image URLs are invalid. It may be because Wikipedia has updated the links. Corrupted images are also discarded. After the downloading, there are 953,042 images that occupy 1.5TB.
- Discard the documents with the same title. We analyze the composition of the remaining document candidates and find that some documents present the same title with the similar content. It is because the page may be updated according to the time, and there are different versions of the documents. To avoid one query mapping to multiple correct documents, we filter these repeated documents. Finally, we obtained 806,357 multimodal documents, including text-image pairs with rich semantic information.

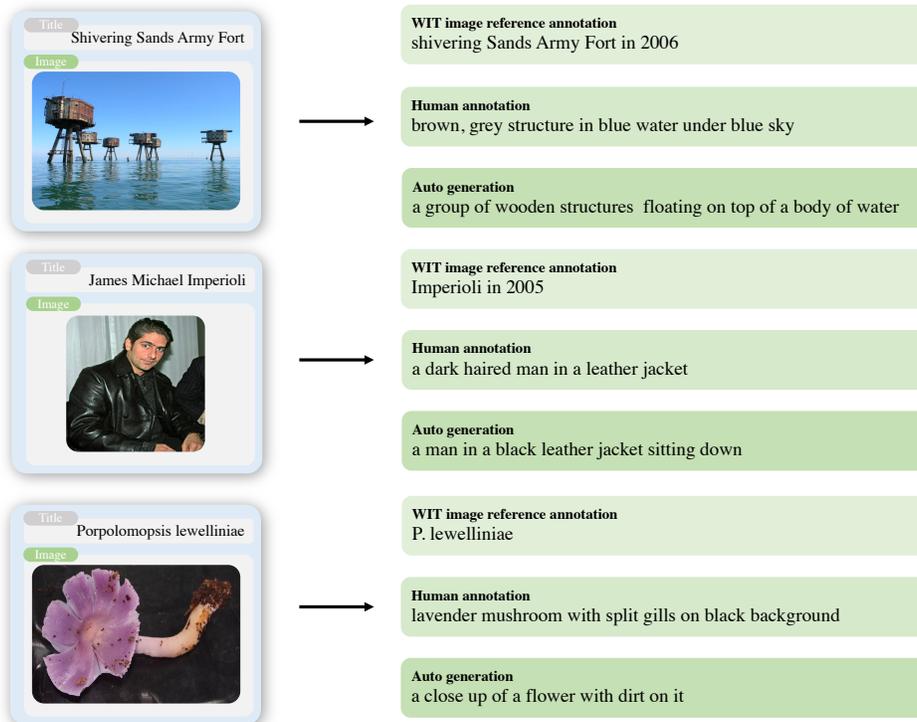


Figure 3: Examples of image-related query.

Human-annotated query As shown in Figure 3, WIT original image reference annotations contain page title or page description rather than the image context. In real-world applications, we consider that user queries may be image descriptions that include image objects, colors, background, or people’s actions. Further, user queries may involve multimodal information, such as image caption fused with text content. In our study, there is no dataset that consists of mixed query for retrieval. Therefore, we hire annotators from Amazon Mechanical Turk to produce human-annotated queries. We require annotators should be masters to ensure label quality. Only annotator with at least 50 approved HITs and an 80% HIT approval rate is allowed. We pay 0.25\$ USD per assignment that includes text-related query, image-related query, and mixed. Also, to award those hardworking annotators, we provide an additional bonus. After the annotation, the statistical data shows that workers’ average time per assignment (three types of queries) is 6 minutes 34 seconds. More details can be seen in Figure 4. In total, we have paid 3,687.24\$ USD (including the platform fees) to annotate 4,276 assignments.

To further ensure the quality of Mr. Right, we provide guidelines and examples to human annotators. They have to read the guidelines first before labeling. Guidelines indicate that a query should meet some restrictions to simulate the possible real-world searching queries, and annotators can come up with the queries based on their habits by following the guidelines. The annotation template is illustrated in Figure 4, and the guidelines are described as follows:

- Words Limit: 10 – 100
- Do not include title.
- Do not copy the sentence from the document.
- Try your best to paraphrase the words.
- Include image information such as color, gender, action, etc.
- Include adjectives and nouns for images.

Auto-generated query Coming up with searching queries is time-consuming. Therefore, we propose auto-generation for training queries. For text-related queries, we extract a snippet from

Guideline

Instruction:

In this task, you can imagine that you are searching something on Google or on other search engine. You need to write a good query to search the topic you want. Sometimes you may want to search for the text information, and sometimes you may want to search for the image. Therefore, your query may include **text-related information** or **image-related information** based on your knowledge.

Update: if you helped this project before, we appreciate your work. Considering the workload, we want to raise the reward per assignment. However, Amazon only offers 2500.0 credit for each batch. We have submitted the application for the credit increment, but the permission has not yet been determined. Because of the time constraint, we have decided to give each worker bonus. The bonuses are as follows:

- 10 < Approved Assignments <= 50: bonus 0.1\$ per assignment
- 50 < Approved Assignments <= 150: bonus 0.2\$ per assignment
- 150 < Approved Assignments <= 200: bonus 0.3\$ per assignment
- Approved Assignments > 200: bonus 0.6\$ per assignment

For example, if you have been **approved of 10 assignments**, we will give you 0.1\$ bonus for your **next approved submission**, so you can earn 0.35\$.

Next Instruction

Skip Instruction

Examples

Here are the requirements you should check:

#	Requirement	Good	Bad
1	Word Limit: 10~100	✔ yogurt-based appetizer in Turkey	✘ Borani is a yogurt-based appetizer with spinach and other ingredients. It's also popular in Jewish cuisine in Iran. Borani is also found in the cuisines of some Turkish provinces...
2	Do not include titles	✔ food made with yogurt	✘ Borani yogurt, popular in the Caucasus
3	Do not copy the sentences from the document	-	✘ an appetizer made with yogurt
4	Make sure to include adjective and nouns	✔ common food in Caucasus yogurt other ingredients	✘ common popular other ingredients cuisine
5	Try your best to paraphrase the words	✔ popular Persian Jewish food	✘ common in Persian Jewish cuisine

I got it!

Fill

Please enter the searching key words.



Spring pygmy sunfish

The spring pygmy sunfish, *Elassoma alabamae*, is a species of pygmy sunfish endemic to springs in northern Alabama. It was historically known to occur in springs in North Alabama along the Tennessee River in Limestone and Lauderdale counties. The spring pygmy sunfish was first discovered in Cave Spring in Lauderdale County, Alabama in 1937 but in 1938, this site was flooded by the creation of the Pickwick Reservoir. The spring pygmy sunfish was considered extinct until its rediscovery in the Beaverdam Spring complex in 1973 by researchers from the University of Tennessee.

Type the query you think is appropriate to search for both the **image** and **text**.

e.g. a bowl of an Iranian appetizer with a spoon in it. (word count 10~100 chars)

Finish Question

Figure 4: Template of human annotation. Annotators read the guideline and examples first, and then create the queries.

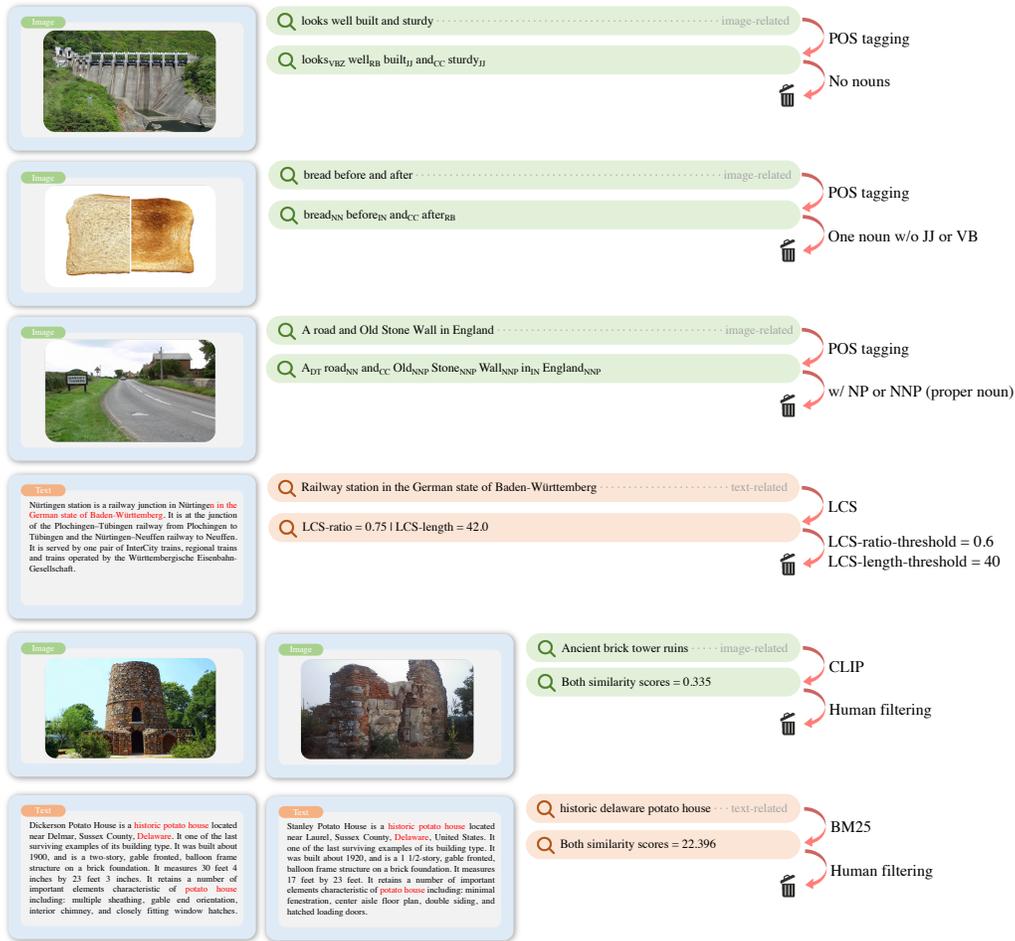


Figure 5: Process of annotated query validation with the help of POS tagging, LCS, CLIP, and BM25.

the first sentence. Our analysis finds that most Wiki passages start with a page title and a brief introduction. Therefore, we utilize this format to extract the sentence’s crucial information. We use Spacy API with `en_core_web_lg` package to parse the sentence and detect the title-dependent verb. Then we take the snippet after the verb as the query. To increase the robustness of models, we also remove the snippet from the document text, which means models have to learn the representation from the remaining text and still be capable of matching the document-query pairs. For image-related queries, we implement the BLIP [49] model, which outperforms many VLP frameworks on image captioning. We employ the model on 351,979 images and produce one caption for each image. The images are resized to 384×384. We use beam search with a beam size of 3 and set the maximum generation length as 30. For the mixed query, it is still challenging to produce a query that fuses the text and image information. Considering the efficiency, we concatenate the text and image queries as mixed queries.

B.3 Filtering

In Figure 5, we show the examples of annotated query validation, including rule-based and human filtering. For the first three examples, We filter out the queries through POS tagging. For the fourth example, we drop the queries by calculating the LCS. For the last two examples, we use CLIP and BM25 to support human discarding queries which map to ambiguous documents.

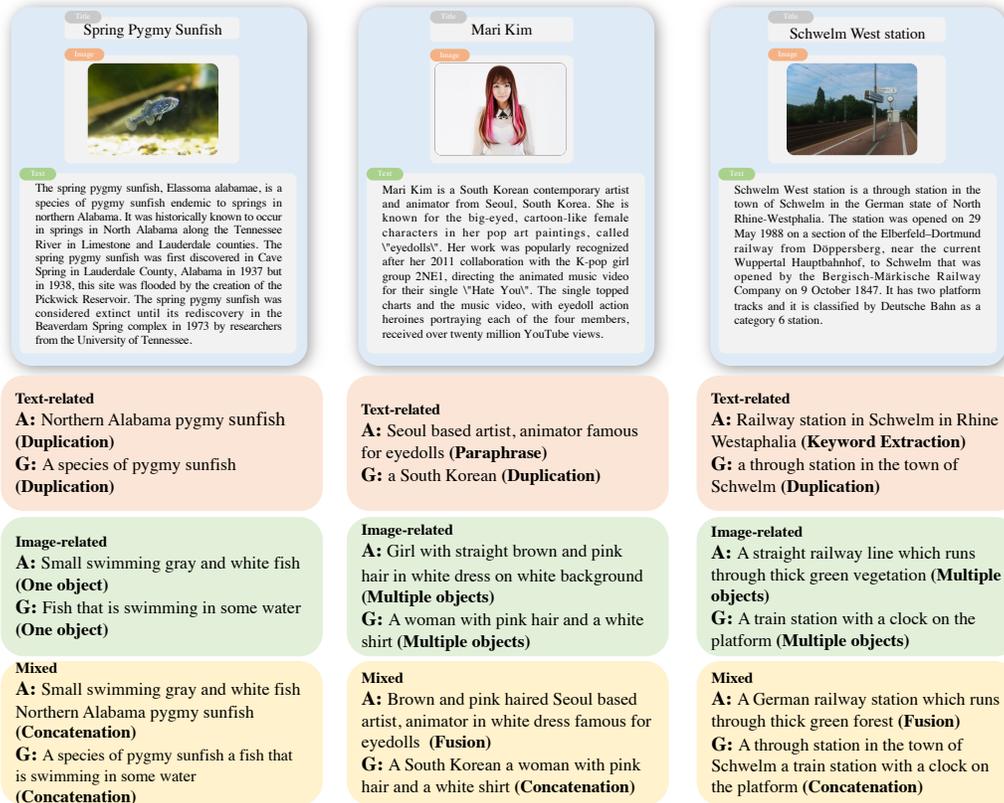


Figure 6: Examples of document and query pair with the corresponding property. The left document* has text-related query by duplication, image-related query with one object, and concatenation-based mixed query. The middle document† has text-related annotated query by paraphrasing and image-related query with multiple objects. The right document‡ has text-related annotated query by keyword extraction and mixed annotated query by fusion.

B.4 Usage

We split Mr. Right into five files: *multimodal_documents.json*, *multimodal_pretrain_pairs.json*, *multimodal_finetune_pairs.json*, *multimodal_val_queries.json*, and *multimodal_test_queries.json*. In *multimodal_documents.json*, it contains document ids, titles, texts, and image URLs. We do not provide image files directly due to the copyright issue. In *multimodal_pretrain_pairs.json*, we provide our auto-generated queries and edited document texts. We still equip this file with the original document texts to keep the flexibility of using Mr. Right. Researchers can create their model framework and train on our auto-generated document-query pairs or produce other effective data. In *multimodal_finetune_queries.json*, we randomly sample human-annotated document pairs for fine-tuning. In *multimodal_val_queries.json* and *multimodal_test_queries.json*, they include corresponding document ids and human-annotated queries. The examples of multimodal document-query pairs are shown in Figure 6. All of our source codes are uploaded to GitHub. Researchers can download json files from our repository. We also offer our training codes.

C Evaluation Details

C.1 Benchmark tasks

In this section, we provide the benchmarks of Mr. Right on humans, baseline retrieval models, and our multimodal framework. There are three types of tasks, and the details are as follows:

Task1: Text-related query This task aims to follow previous text retrieval datasets [8, 9, 27, 26, 28]. Users mostly search for documents relying on the keywords from document paragraphs or text-based information. In our dataset, text-related queries contain name entities (person, date, organization, location) or factual knowledge (relations, terminologies).

Task2: Image-related query This task aims to follow previous text-to-image retrieval datasets [33, 23, 24]. With a blurred impression about the appearance of an object, users search documents based on the part of context from document images. In our dataset, image-related queries are similar to image captions that explain details of objects, such as color, shape, amount, position, or action.

Task3: Mixed query We propose this task to simulate users searching documents with text-related and image-related information. To precisely find the correct document, Mr. Right provides document texts and images to consider both modalities for retrieval. Our mixed queries generate a brief description that includes the document paragraph and photo, which can be viewed as a combination of corresponding text-related and image-related queries.

C.2 Baseline models

Text retrieval models To evaluate text retrieval performance with state-of-the-art (SOTA) neural frameworks, we test three approaches in the followings. 1) RoBERTa-base [51]: a pre-trained language model which can encode both documents and queries into the contextualized sentence representations to compute the similarity in the same vector space. 2) DiffCSE [55]: current unsupervised SOTA among sentence representation learning methods. 3) all-mpnet-base-v2 (SBERT): current supervised SOTA for sentence embedding tasks and semantic search tasks on SentenceTransformers [56] leaderboard. We evaluate both zero-shot and fine-tuned performance for these models. We train with an in-batch negative loss function and use an AdamW optimizer with learning rate 2×10^{-5} and batch size 32 for 30 epochs.

Image retrieval models Current image-text retrieval models are highly related to pre-trained vision-and-language models. Training with natural language supervision, these models demonstrate the ability of crossmodal image retrieval. We zero-shot evaluate CLIP [47], and ALBEF [20] as our baselines and also fine-tune on our dataset. We don't use other VLP models (e.g., METER [19], ViLT [53]) because most of them require a high computational overhead for evaluation due to the need to calculate matching scores across all image text pairs. We fine-tune CLIP and ALBEF using an Adam optimizer with learning rate 1×10^{-6} and batch size 128 for 40 epochs.

Multimodal retrieval models Without existing baselines, we build an ensemble model by integrating the document-query similarity scores from our best TR model and IR model. We fuse the scores with a weighted sum parameter tuning on the validation set for different tasks. The final ensemble relevance scores are then used to rank the search results.

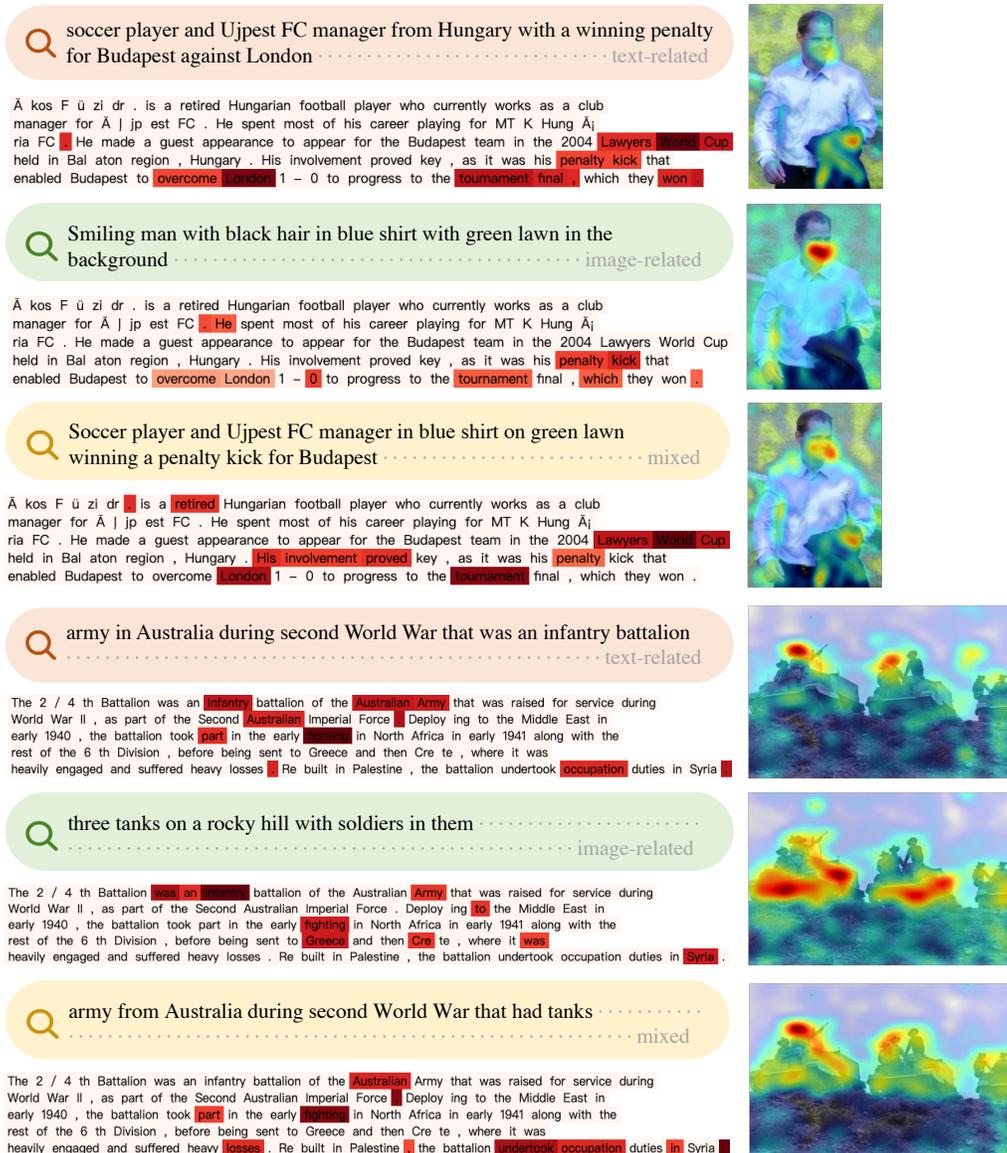


Figure 7: Grad-CAM visualizations on the cross-attention maps of the query-document matching classifier corresponding to different queries. Heats are scattered across document image and texts.

C.3 Grad-CAM visualizations

To better understand the multimodal representations of documents, we compute Grad-CAM visualizations on the cross-attention maps of the query-document matching classifier in Figure 7. With different queries, our model will interact with different parts of the image and texts, which is highly correlated to where humans would look to match the pairs. For the above example, the word “smiling” highly focus on the face of the image, and the word “winning” is related to “tournament” and “Lawyers World Cup” in texts. For the bottom example, the words “army” and “tank” in the queries will attend to corresponding parts of the image, while the words “Australian” and “War” will highlight “Australian” and “fighting” in texts, respectively.

Query

Seoul dynasty Overwatch player for GC Busan in white, pink shirt looking down with blue background

Hi! CHENG AN HSIEH, please choose the matching document for the query:

Document Candidates

 <p>Park Joon-yeong, better known by his online alias Profit, is a professional South Korean Overwatch player for the Seoul Dynasty of the Overwatch League. Prior to the OWL's inception, he played for GC Busan, winning OGN's APEX Season 4 and was named the APEX Finals most valuable player. In the Overwatch League, Park signed with the London Spitfire of the OWL in its inaugural season. He won the league's first Grand Finals with the Spitfire, after they defeated the Philadelphia Fusion, and was named the Grand Finals most valuable player.</p> <p>Choose</p>	 <p>Lee Keun-ho is a South Korean footballer who plays for Ulsan Hyundai and South Korea national team. His pace, work-rate, and link-up plays mark him as a highly rated forward in Asia. South Korean international since 2007, he currently has 19 goals in 84 caps. He represented his country for two Asian tournaments, 2008 Summer Olympics and 2014 FIFA World Cup.</p> <p>Choose</p>
 <p>Bosley Yu Yang is an African professional footballer who plays as a midfielder for Hong Kong Premier League club Kitchee.</p> <p>Choose</p>	 <p>Kim Hyo-jong, better known by his online alias Haksal, is a professional South Korean Overwatch player who plays for the New York Excelsior of the Overwatch League. He previously competed for the Vancouver Titans before mutually parting ways with the organization. Kim began his professional Overwatch career playing for RunAway, where he was named OGN's APEX Season 4 most valuable player and won 2018 Season 2 of Overwatch Contenders Korea. In the Overwatch League, Park signed with the Titans in their first year of existence. With the Titans, he reached the 2019 Grand Finals, where they fell to the San Francisco Shock. Kim received the first-ever OWL "Rookie of the Year" award in the 2019 regular season.</p> <p>Choose</p>

Figure 8: Template of human evaluation.



Figure 9: Failed examples of our model including text-related, image-related, and mixed queries.

C.4 Failed examples

We create human evaluation by randomly sampling 50 examples from each type of task. As illustrated in Figure 8, human annotators have to select the most relevant document from four candidates obtained from our MR model. Each question is answered by three workers. If more than half workers have the same answers and match the correct document, we consider humans can answer this question correctly; otherwise, humans fail it. With human evaluation, we can distinguish the performance difference between humans and our models in Figure 9. For text-related queries, we can find our model is capable of obtaining the related document instead of the accurate one. However, humans can easily choose the right one by matching text keywords, such as “9:1” in the query and “9-to-1” in the document of the first example. For image-related queries, our model prefers to retrieve specific color words and ignores the remaining, such as the “blue spire” and “green lawn” of the second example. On the other hand, humans can perceive the details of images. For mixed queries, our model may pay attention to wrong words, such as the word “walking” and “centres” in the last two examples. On the contrary, humans can simultaneously recognize the correct document by matching text keywords and image context.

C.5 Compared to auto-generated and human-annotated queries

Since our human-annotated queries have more diverse properties than auto-generated queries as shown in Table 3, we compare their performance on mixed queries with our pre-trained MR model (METER). The results are shown in Table 6, and we can find auto-generated queries outperform human-annotated queries because the annotated queries are more complex and difficult to learn. Therefore, we fine-tune our models on 1k annotated queries to adapt in human-annotated domain.

	mixed query MRR@10
Auto-generated	32.6
Human-annotated	9.0

Table 6: Retrieval performance of mixed queries from auto generation and human annotation.

D Proposed Model Details

With inputs document texts D_T and document image D_I , we derive the document features F_d from multimodal encoder E_d . Also, with input query texts Q_T , we obtain the query features F_q from query encoder E_q . After, we take average of the size-variant features F_d and F_q to get a single fixed vector as document and query representations R_d and R_q .

$$F_d = E_d(D_T, D_I) \quad \text{and} \quad F_q = E_q(Q_T) \quad (1)$$

$$R_d = \text{Average}(F_d) \quad \text{and} \quad R_q = \text{Average}(F_q) \quad (2)$$

With document and query representations R_d and R_q , we aim to close the distance between two vectors by contrastive learning. Therefore, we learn two projection functions f_d and f_q with fully-connected layers and L2-normalization to map their representations into the same space. We calculate the similarity by dot product for all document and query vector pairs in a training batch when treating matched pairs as positive while all other pairs as negative. The contrastive loss L_{dqc} we minimize is in the following:

$$\text{Sim}(R_d, R_q) = f_d(R_d)^\top f_q(R_q) \quad (3)$$

$$P_{d2q}^i = \frac{\exp(\text{Sim}(R_d^i, R_q^i)/\tau)}{\sum_{j=1}^N \exp(\text{Sim}(R_d^i, R_q^j)/\tau)}, \quad P_{q2d}^i = \frac{\exp(\text{Sim}(R_q^i, R_d^i)/\tau)}{\sum_{j=1}^N \exp(\text{Sim}(R_q^i, R_d^j)/\tau)} \quad (4)$$

$$L_{dqc} = -\frac{1}{B} \sum_i \frac{Y_{d2q}^i \log(P_{d2q}^i) + Y_{q2d}^i \log(P_{q2d}^i)}{2} \quad (5)$$

Here, we calculate the normalized softmax loss for both document-to-query and query-to-document classification. The loss is set up with batch size B , negative samples size $N (= B)$, and a learnable temperature parameter τ to scale the logits. For negative pairs of the document to query, R_d^i and R_q^j are the representations of the document in the i -th pair and query in the j -th pair, respectively. To more effectively organize in-batch negatives, we keep two queues [20] to store the most recent K representations, helping enlarge the amounts of negative samples. The modified equation is to change negative sample size N from batch size B to queue length K .

In addition to contrastive loss, we obtain document-query matching loss to learn a fine-grained similarity of pair documents and queries. The matching loss is a binary cross-entropy loss to predict whether a pair of documents and queries are matched or mismatched. We build a 6-layer transformer-based classifier C with input document features F_d and query features F_q . Specifically, a special token (e.g., $[CLS]$) is inserted at the beginning of the input sequence, and it tries to learn a global cross-modal representation in transformers. After, a linear classifier is added to the $[CLS]$ token to predict a binary label. The whole matching loss is in the following:

$$P_{dqm}^i(j = y_{dqm}^i) = \frac{\exp(C_j(F_d, F_q))}{\exp(C_0(F_d, F_q)) + \exp(C_1(F_d, F_q))} \quad (6)$$

$$L_{dqm} = -\frac{1}{B} \sum_i Y_{dqm}^i \log(P_{dqm}^i) \quad (7)$$

Our full training objective is:

$$L = L_{dqc} + L_{dqm} \quad (8)$$

E Dataset License

Our dataset is under the Creative Commons Attribution Share Alike 4.0 (CC BY-SA 4.0) license.

F Maintenance

We believe that Mr. Right will assist researchers in building robust multimodal retrieval models and improve the current retrieval systems. We are willing to maintain Mr. Right. If researchers have any problems, they can create an issue from our repository. We also welcome any methods to perform on our benchmark. We bear all responsibility for violations of rights related to Mr. Right.