

Rethinking Sparse Lexical Representations for Image Retrieval in the Age of Rising Multi-Modal Large Language Models

Kengo Nakata, Daisuke Miyashita, Youyang Ng,
Yasuto Hoshi, and Jun Deguchi

Kioxia Corporation, Yokohama, Japan
{kengo1.nakata, daisuke1.miyashita, youyang.ng,
yasuto1.hoshi, jun.deguchi}@kioxia.com

Abstract. In this paper, we rethink sparse lexical representations for image retrieval. By utilizing multi-modal large language models (M-LLMs) that support visual prompting, we can extract image features and convert them into textual data, enabling us to utilize efficient sparse retrieval algorithms employed in natural language processing for image retrieval tasks. To assist the M-LLM in extracting image features, we apply data augmentation techniques for key expansion and analyze the impact with a metric for relevance between images and textual data. We empirically show the superior precision and recall performance of our image retrieval method compared to conventional vision-language model-based methods on the MS-COCO, PASCAL VOC, and NUS-WIDE datasets in a keyword-based image retrieval scenario, where keywords serve as search queries. We also demonstrate that the retrieval performance can be improved by iteratively incorporating keywords into search queries.

Keywords: image retrieval, sparse lexical representation, LLM

1 Introduction

As deep learning technologies have evolved, deep neural networks (DNNs) have achieved exceptional performance in image recognition and object detection tasks [14, 15, 21, 42], and approaches leveraging these networks have been extensively explored for image retrieval tasks [13, 50, 59, 61]. With the recent emergence and widespread adoption of vision-language models [18, 25, 39, 63], text-to-image retrieval has become one of the mainstream research areas in image retrieval. These models are pre-trained on vast amounts of paired image-text data collected from the internet, and they learn to map the images and their corresponding texts into similar dense vector representations in a shared latent space [18, 25, 39]. By utilizing such pre-trained models, images that are semantically similar or related to a query text can be retrieved based on the distance calculations between their dense vectors.

Images can contain a wide variety of information and features, and the criteria to determine whether images are similar or dissimilar are inherently subjective

and not uniquely defined. For instance, when retrieving images, the features that users focus on within images may vary depending on individual preferences and situational factors. However, a query does not always represent or reflect a user’s desires or intentions, and it could be incomplete or lacking in required information to specify them. Despite these limitations, the vision-language model attempts to provide results that are relevant to the user’s request, by implicitly compensating for the lack of information in the incomplete query based on the knowledge acquired through training. Unfortunately, this compensated information may not always align with the user’s desires or intentions.

In text retrieval, a keyword-based approach is commonly employed in practical applications to retrieve documents containing specified keywords. Users can combine multiple keywords as search queries to specify their focus areas or topics. After viewing the retrieval results, users can iteratively refine their search queries by modifying and/or adding keywords as feedback. Even if the retrieval model cannot initially provide the desired results, users can adaptively obtain results that align with their preferences or intentions through this iterative process. In the light of these flexible capabilities, we aim to explore better methods for applying keyword-based retrieval to image retrieval tasks.

While we can directly apply the conventional vision-language models to the keyword-based retrieval, we cannot overlook the substantial advancements in large language models (LLMs) over the past few years, which have demonstrated a remarkable ability to comprehend context within dialogue interfaces [5, 34, 47, 53–55]. Additionally, multi-modal LLMs (M-LLMs) have already been proposed to comprehend visual information within images through visual prompting, which involve processing images along with textual data as queries [4, 9, 22, 24, 28, 29]. By utilizing the advanced capabilities of M-LLMs, we can extract features from images and linguistically represent them in textual data like tags and captions. Then, we can leverage the advantages of natural language processing (NLP) techniques for image retrieval tasks. We encode the generated textual data into sparse lexical vectors and utilize efficient retrieval algorithms to enable effective image retrieval based on their sparse lexical representations.

In this paper, we focus on the text-to-image retrieval task and rethink the task in this age of rising such powerful M-LLMs. As a text-to-image retrieval task, we consider a keyword-based image retrieval scenario where a search query consists of a few words representing the contents or objects depicted in images. Through quantitative analysis on the benchmark datasets, we demonstrate that our retrieval system outperforms conventional vision-language model-based retrieval methods in terms of precision and recall. Specifically, we introduce a cropping technique to assist the M-LLM in effectively extracting image features, and analyze the effectiveness by evaluating a metric for relevance between images and texts. As our findings, we empirically show that the conventional vision-language model-based methods outperform our approach, if a less informative caption is used as a search query. This seems to depend on whether there is a function to compensate for the lack of information in the less informative query. However,

the retrieval performance of our system improves significantly when we incorporate keywords into the search query, making the query explicitly informative.

The main contributions of this paper are as follows:

- We introduce a text-to-image retrieval system that utilizes M-LLMs and retrieval algorithms based on sparse lexical representations, and evaluate its effectiveness on various benchmark datasets.
- To enhance M-LLM performance in extracting image features, we employ data augmentation techniques for key expansion and quantitatively evaluate the improvement in retrieval performance.

2 Related Work

2.1 Evolution of Image Retrieval Research

Image retrieval has been studied extensively in recent decades. This includes tasks such as finding images similar to a given input image, searching for images with specific content features like colors, shapes, and textures [56, 64], and retrieving images based on their semantic meaning or content categories [37]. Researchers have also explored more specialized tasks like conditioned image retrieval, where both images and semantic conditions are used as query inputs to find relevant images [2, 3]. In practice, image retrieval often involves searching for images based on descriptive texts, such as metadata and hashtags utilized on social media platforms. This requires matching user queries with relevant images by analyzing their associated metadata or the images themselves. However, despite its practical significance, text-to-image retrieval remains an understudied domain within computer vision research. This lack of research may be attributed to the fact that the accuracy of the image retrieval heavily depends on the quality and thoroughness of human annotations.

Recent advances in NLP have empowered text-to-image retrieval through the development of contrastive language-image pre-training techniques [18, 25, 39]. By mapping images and texts into a shared latent space and calculating their semantic similarities, these techniques enable image retrieval based on textual descriptions. Moreover, the rise of M-LLMs that can generate textual descriptions for images without the need for human annotations [51, 60, 62] emphasizes the significance of exploring interactions between images and descriptive texts in the form of interpretable lexical data. In this paper, we rethink text-to-image retrieval by leveraging sparse lexical representations. For the context of text-to-image retrieval tasks, while previous research has primarily focused on the evaluations in the caption-to-image retrieval [7, 18, 25, 38, 39, 48], we shift our attention to keyword-based image retrieval, which is more commonly used in practical applications but remains understudied. We explore its potential applications and provide insights for future research directions.

2.2 Large Language Models for Visual Promptings

In the past few years, LLMs have made remarkable progress, and their popularity has grown significantly due to their impressive ability to comprehend context [5, 34, 47, 53–55]. Recently, there has been a growing development of M-LLMs, which support visual inputs as well as text-based promptings, providing comprehensive and flexible applications [4, 22, 28, 29]. GPT-4V is capable of accepting both text and image prompts, understanding visually depicted scenarios in images, and addressing complex visual question answering tasks [34, 60, 62]. TagGPT offers a tagging system that extracts tags from multi-modal content such as images and videos without requiring additional knowledge or human annotations, by leveraging the M-LLMs [22]. BLIP-2 and InstructBLIP introduce Querying Transformer that bridges the modality gap between images and texts, while keeping pre-trained vision encoders and backbone LLMs frozen [9, 24]. LLaVA proposes a visual instruction tuning technique for large multi-modal models that integrate vision encoders with LLMs for general-purpose visual and language understanding by utilizing instruction-following data generated by GPT-4 [28, 29]. The rise of these M-LLMs presents an opportunity for us to revisit the understudied keyword-based image retrieval without the need for human annotations. In this paper, we leverage these M-LLMs to transform visual information in images into expanded lexical representations, enabling us to harness traditional efficient sparse retrieval methods for image retrieval tasks. Our approach effectively combines classic techniques with cutting-edge innovations.

2.3 Sparse Retrieval

Recently, there has been growing interest in using dense retrieval methods that leverage dense vector representations generated by DNNs for image retrieval tasks [13, 50, 59, 61]. However, sparse vector representations, typically in the form of lexical retrievers, have also been explored due to their enhanced interpretability and analytical capabilities [6, 31, 33, 66]. To address the perceived trade-off between accuracy and interpretability, LexLIP [31] introduces a lexicon-weighting paradigm to significantly reduce retrieval latency while maintaining high performance with bag-of-words models. Similarly, STAIR [6] maps images and texts to a sparse token space to construct sparse text and image representations for improved retrieval accuracy. These studies demonstrate the potential of sparse retrievers to outperform dense retrievers. Our approach leverages a multi-modal language model to extract image features into textual data. We then utilize vectorization and retrieval algorithms in NLP tasks, such as BM25 [43], TF-IDF [44], and word2vec [32], for image retrieval tasks. Among these techniques, BM25 is considered an efficient sparse retrieval algorithm and is frequently used for benchmark evaluations in information retrieval tasks [31, 49, 66]. BM25 demonstrates better out-of-distribution generalization capabilities compared to dense retrievers [52], and outperforms them in retrieving named entities or words that were not seen during training [46]. Based on such potential capabilities, we employ the efficient and standard retrieval algorithm, BM25, for sparse lexical vectors

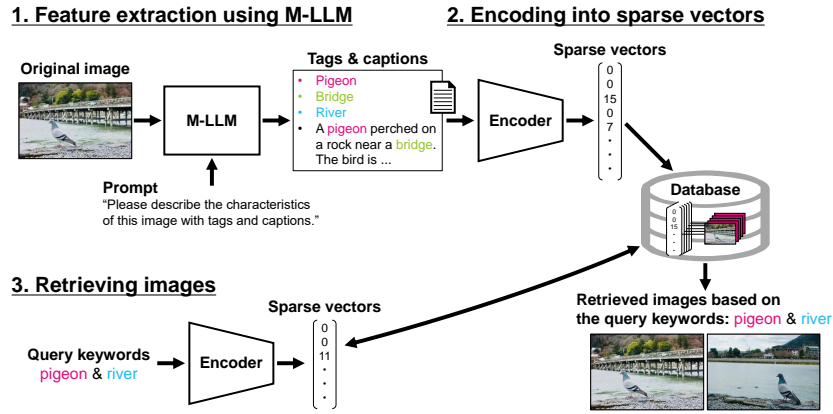


Fig. 1: Overview of our image retrieval system. Our system utilizes an M-LLM to describe an image in textual data such as tags and captions. The textual data is encoded into sparse vectors. When retrieving specific images, query keywords are also encoded into sparse vectors, enabling the retrieval of relevant images.

directly converted from textual data. Our approach does not rely on dense latent representations extracted by DNNs, and eliminates the need for specialized vector space adaptation, enabling the application of key expansion techniques for enhanced performance. By transforming the image retrieval task into a sparse lexical retrieval task, we can rethink image retrieval from an NLP perspective.

3 Approach

Fig. 1 provides an overview of our image retrieval system. Our image retrieval system consists of three processes: (1) feature extraction using an M-LLM, (2) encoding into sparse vectors, and (3) retrieving images. Our system accepts text-based queries such as keywords, and returns a set of relevant images from a database. Each process is described below.

3.1 Feature Extraction Using M-LLM

First, we generate textual data for images to be stored in a database. By utilizing M-LLMs with visual prompting capabilities [4, 22, 29, 34], it is possible to extract features from images and represent them in textual data. Among the M-LLMs, the pre-trained LLaVA model demonstrates high performance across various benchmark datasets for general-purpose visual and language understanding [28, 29], therefore we utilize this publicly available model in our work. We provide the M-LLM with an image and a prompt such as “Please generate multiple captions to describe the features of this image.” or “Please describe

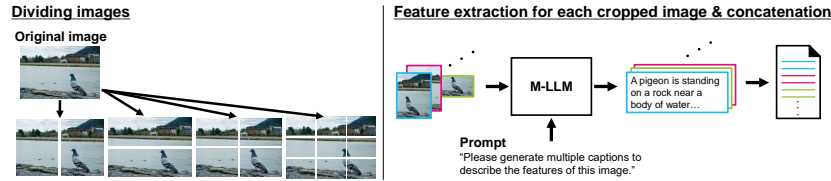


Fig. 2: Data augmentation techniques for key expansion. An original image is segmented into multiple regions as cropped images (left), and each cropped image is processed by an M-LLM to generate captions that extract the features of each region (right). By concatenating the generated captions, including those derived from the original image, we can extract a comprehensive set of features from the whole image.

the characteristics of this image with tags and captions.” Then, we can obtain generated lexical tags and captions that represent the image features.

For our system, we can also apply image captioning models [20, 23, 65]. However, we choose to utilize pre-trained M-LLMs provided in the open-source library, as they offer generally powerful performance without the need for model tuning [28, 29]. The remarkable ability of M-LLMs to interactively comprehend prompts and context can be utilized to iteratively extract information related to image content by providing them with step-by-step queries. For instance, after the M-LLM generates captions for an image based on the initial prompt, we provide it with another prompt: “If there are any additional features of this image that are not expressed in the generated captions, please generate additional captions to explain them.” Considering such potential applications, we present the system using M-LLMs instead of image captioning models in this paper.

Data Augmentation Techniques for Key Expansion. In order to sufficiently extract features and information from various viewpoints within images, we employ a cropping technique, which is commonly used in image recognition tasks as a means of data augmentation. As shown in Fig. 2, we divide an original image into multiple segments as cropped images, such as two vertical segments, two horizontal segments, four segments, or nine segments. For each cropped image, the M-LLM generates a corresponding caption that describes the feature of the image. By concatenating all the generated captions, including those derived from the original image, we can extract a comprehensive set of features from the whole image. The cropping technique assists the M-LLM in effectively extracting features from images, while expanding the key sets in the database (i.e., key expansion), leading to improved retrieval performance.

For cropping images, we can also utilize object detection models like YOLO or spatial transformer networks [17, 40, 41]. However, when using the object detection model, areas where the model fails to detect objects or where no objects exist (e.g., sky scenery, glass fields, or sea areas) will not be cropped. Consequently, the M-LLM cannot extract information from these areas within the

images. To avoid the impact of inductive biases in object detection models, in this paper, we do not use such models, but instead employ fixed pattern cropping, regardless of the location of objects within the images, as shown in Fig. 2.

Analysis with CLIPScore. To evaluate the impact of the cropping technique, we use a metric called CLIPScore [16]. CLIPScore is used to evaluate the relevance between an image and a textual description by comparing the embeddings extracted through the models pre-trained by CLIP. If the pre-trained model embeds an image data (I) and a textual data (T) into their respective embeddings \mathbf{E}_I and \mathbf{E}_T , we can calculate CLIPScore based on the cosine similarity between their embeddings as follows,

$$\text{CLIPScore}(I, T) = w \times \max(\cos(\mathbf{E}_I, \mathbf{E}_T), 0), \quad (1)$$

$$\cos(\mathbf{E}_I, \mathbf{E}_T) = \frac{\mathbf{E}_I \cdot \mathbf{E}_T}{\|\mathbf{E}_I\| \|\mathbf{E}_T\|}, \quad (2)$$

where w is a scaling parameter used to adjust the range of score distribution, and we set $w = 2.5$ as reported in the original paper [16]. This metric is typically used for evaluating the quality of image captioning models, by measuring the relevance between the captions generated by the models and the corresponding images [16, 35, 45]. In contrast, we utilize this metric to evaluate the effectiveness of the cropping technique in eliciting diverse textual descriptions. For example, if the CLIPScore value between a cropped image and a textual description is larger than the value between the original image and the textual description, it indicates that the cropping technique has produced a more suitable image for eliciting the textual description.

In our evaluation, we use the Microsoft COCO (MS-COCO) dataset [27], which comprises images depicting a variety of scenarios involving multiple objects from 80 different categories. Specifically, we employ the 2017 validation set of MS-COCO, consisting of 5,000 images. We adopt the list of 80 categories as diverse textual descriptions (e.g., “bicycle” and “refrigerator”), and calculate CLIPScore based on Eqs. (1) and (2) between each image and textual description. We utilize a model pre-trained by CLIP (ViT-L/14@336px) for vision and text encoders, as in the experiments described in Sec. 4.1. We average CLIPScore over all the images and all the textual descriptions as follows,

$$\text{AveragedCLIPScore}_{\text{all}} = \frac{1}{N_I} \sum_{i=1}^{N_I} \text{AveragedCLIPScore}_{\text{each}}(I_i), \quad (3)$$

$$\text{AveragedCLIPScore}_{\text{each}}(I_i) = \frac{1}{(N_C + 1)N_T} \sum_{j=0}^{N_C} \sum_{k=1}^{N_T} \text{CLIPScore}(I_{i,j}, T_k), \quad (4)$$

where N_I is the total number of images in the dataset, while N_C represents the number of cropped images for each image and N_T denotes the number of textual descriptions. Additionally, T_k represents the k -th textual description,

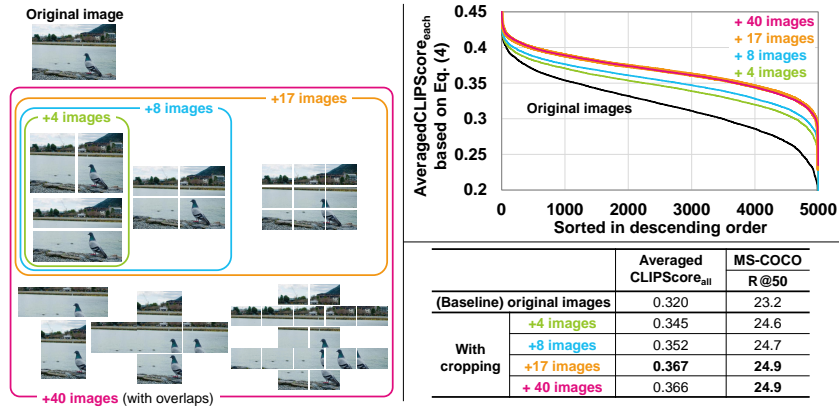


Fig. 3: The variations in averaged CLIPScore based on Eq. (4) for each of the 5,000 validation images from the MS-COCO dataset. As shown in the left figure, the original images are cropped by fixed patterns including overlaps. In the upper right graph, the values are sorted by averaged CLIPScore for each image in descending order. The lower right table summarizes averaged CLIPScore based on Eq. (3) for all the images in the dataset, along with the top-50 recall performance (R@50).

and $I_{i,j}$ refers to the j -th cropped image for the i -th image. Specifically, when $j = 0$, it corresponds to the original image before cropping.

As shown in the left figure of Fig. 3, the original images are cropped by fixed patterns, which include overlapped edges, in order to evaluate the impact of information loss along the boundaries of the cropped patterns. The upper right figure in Fig. 3 shows the variations in averaged CLIPScore based on Eq. (4) for each of the 5,000 validation images from the MS-COCO dataset. As the number of cropped images increases to 17, averaged CLIPScore based on Eq. (3) for all the images in the dataset increases, as summarized in the lower right table. We also evaluate the retrieval performance, which is measured by top-50 recall (R@50), in retrieving the images (the recall is calculated as in Eq. (6) by evaluating keyword-based image retrieval scenario described in Sec. 4.1). As the number of cropped images increases to 17, we can observe the improvement on R@50. On the other hand, when the number of cropped images increases from 17 to 40 by cropping with overlaps, averaged CLIPScore based on Eq. (3) does not increase and the recall performance is not improved. At this point, we consider that the impact of cropping has reached a point of saturation, and the information loss from the original images has been effectively mitigated. Hence, we select 17 fixed patterns for cropping in the experiments of Sec. 4.

By employing the cropping technique, images can be divided into multiple segments, which reduces the number of features that the M-LLM needs to focus on and describe for a single image. The above results empirically indicate that the cropping technique enhances the ability of the M-LLM to extract features from images more precisely, leading to improved retrieval performance.

3.2 Encoding into Sparse Vectors and Text-to-Image Retrieval

We encode the textual data generated by the M-LLM into lexical representations with sparse vectors, where we insert non-zero values only at the positions that correspond to the terms present in the corpus. We employ the BM25 algorithm [43] for effective image retrieval based on the sparse lexical representations. BM25 is a widely used algorithm in NLP applications that efficiently retrieves documents by scoring them based on their term frequencies, enabling the search for relevant documents to a given query. Specifically, this algorithm assigns higher weights to rare terms within the corpus and lower weights to common terms. These vectors are then stored in the form of an inverted index, allowing for quick lookups of documents containing specific terms and significantly reducing the search space, thereby accelerating the retrieval process.

When searching for images, we set descriptive keywords as search queries to focus on specific features or aspects within images. We convert these search queries into sparse lexical representations, and we can retrieve relevant textual data and corresponding images from the database based on their sparse lexical representations by using the BM25 algorithm. The actual settings, such as parameters, are described in the subsequent section.

4 Experiments

4.1 Experimental setup

Model settings. A LLaVA model is an end-to-end trained large multi-modal model that connects a vision encoder and an M-LLM for general-purpose visual and language understanding [28, 29]. Compared to other M-LLMs with visual prompting capabilities, the LLaVA’s large multi-modal model has demonstrated superior performance on a variety of benchmark datasets, outperforming models like BLIP-2 and InstructBLIP [28, 29]. In our experiments, we utilize one of the pre-trained multi-modal models (`llava-1.5-13b-hf`)¹ publicly available from the Hugging Face’s Transformers library [57]. Based on our analysis in Sec. 3.1, we divide original images to obtain 17 cropped images as shown in the left figure of Fig. 3. We provide the pre-trained LLaVA model with the prompt “Please generate multiple captions to describe the features of this image.” for each of the original images and the cropped images, in order to generate captions that represent the features of each image. After the caption generation, we concatenate all the generated captions.

We use zero-shot vision-language models pre-trained by CLIP (`ViT-L/14@336px`)² and ALIGN (`align-base`)³ as our baselines, because these models have demonstrated robust and reliable performance across a wide range of

¹ We utilize the pre-trained models available at the following URLs:

<https://huggingface.co/llava-hf/llava-1.5-13b-hf>

² <https://github.com/openai/CLIP/blob/main/clip/clip.py>

³ <https://huggingface.co/kakaobrain/align-base>

benchmark datasets [18,39]. Note that the pre-trained LLaVA model utilizes the vision encoder model included in the same pre-trained CLIP model as its vision encoder. Given this, we can expect that our chosen M-LLM has similar visual performance capabilities as one of our baseline models. Neither our method nor the baseline methods involves fine-tuning the pre-trained models.

Task settings and datasets. As a text-to-image retrieval task, we consider a keyword-based image retrieval scenario using three benchmark datasets: MS-COCO [27], PASCAL VOC [11], and NUS-WIDE [8]. Each dataset consists of images featuring multiple objects and scenes, and each image is annotated with descriptive labels. In our experiments, we utilize the ground-truth labels assigned to each image in the respective datasets as keywords for our search queries to retrieve the corresponding images. For instance, the 2017 validation set of MS-COCO contains a total of 5,000 images, of which 4,952 images include one or more objects belonging to 80 categories with 80 different label types. Then, each of these 80 distinct labels like “bus” serves as a keyword for our search query. Similarly, we utilize the 2007 test set of PASCAL VOC, which comprises 4,952 images with one or more objects per image belonging to 20 classes. Moreover, we explore the NUS-WIDE dataset, featuring 260,648 web images with one or more textual tags per image. Each image is labeled with multiple concepts from a set of 81 labels. As a subset, we focus on the 195,834 image-text pairs corresponding to the 21 most common concepts, and we use a total of 2,100 image-text pairs from this subset for the test set, as previously validated in [19,58].

We also consider a multi-keyword-based image retrieval scenario, where multiple keywords are combined and used to refine the search criteria like an AND search. For our experiments, we join the ground-truth labels of each image in the aforementioned datasets into a search query. For example, if we use an image file named “val2017/000000454661.jpg” that contains labeled objects such as “car”, “bus”, and “traffic light” in the MS-COCO dataset, we join the labels together as “car, bus, traffic light” to form the search query for the image.

Furthermore, we evaluate the performance in a caption-to-image retrieval setting, which is a basic evaluation setting for text-to-image retrieval tasks. We utilize the MS-COCO and Flickr30k [36] datasets, and employ the ground-truth caption sentence for each image as a search query to retrieve the corresponding image. Specifically, we use a total of 5,000 images in the 2017 validation set of MS-COCO and a total of 1,000 images in the test set of Flickr30k.

Finally, in order to explore the potential practical applications of our retrieval system, we consider a scenario for text-to-image retrieval with user feedback. In this scenario, after an initial retrieval based on a search query, a user iteratively incorporates keywords into the search query as user feedback to gradually clarify the vision for the desired image like a multi-turn refined search. As an example, we combine the caption-to-image retrieval setting with the keyword-based image retrieval setting. We utilize both ground-truth captions and labels for 4,952 images in the 2017 validation set of MS-COCO. After an initial retrieval based

on a ground-truth caption for an image, we iteratively incorporate the ground-truth labels for the image into the search query as keyword-based user feedback.

Retriever settings. In our retrieval system, we leverage the BM25 algorithm and employ Pyserini [26], which is a Python interface to Lucene’s BM25 implementation. We set the parameters to their default values of $k_1 = 0.9$ for term frequency scaling and $b = 0.4$ for document length normalization. In the baseline, the vision-language models retrieve images by calculating the distance based on cosine similarity between the embeddings for query texts and images.

In the keyword-based image retrieval evaluations, we directly utilize ground-truth labels as textual inputs for the text encoder models. We do not employ prompt templates like “A photo of a {label}.” commonly used in CLIP [39]. Our preliminary experiments showed that the use of such templates could potentially decrease recall performance, particularly when the labeled object was not the main focus of the image. Therefore, we do not use such templates and instead directly use the labels.

Evaluation metrics. To evaluate the retrieval performance, we use precision and recall metrics. We sweep the number of top-retrieved images (k) by powers of two from 1 to the total number of images in each dataset. We calculate the precision ($P@k$) and recall ($R@k$) metrics based on the number of true positives ($TP@k$) among the top k retrieved images for each query (q) as follows,

$$P@k = \frac{\sum_{q=1}^{N_Q} TP_q@k}{N_Q k}, \quad (5)$$

$$R@k = \frac{\sum_{q=1}^{N_Q} TP_q@k}{\sum_{q=1}^{N_Q} P_q}, \quad (6)$$

where N_Q is the total number of queries and P_q is the total number of ground-truth images for each query.

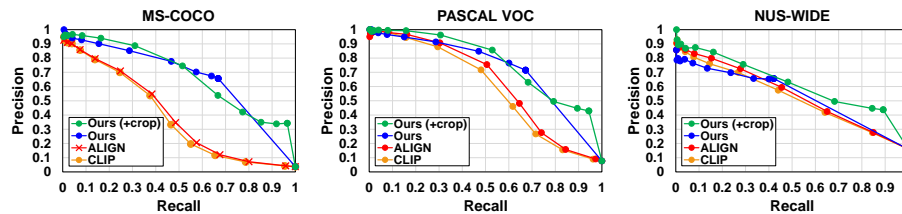


Fig. 4: Comparison of precision and recall curves between our retrieval method and the conventional retrieval methods for the keyword-based image retrieval setting on the MS-COCO, PASCAL VOC, and NUS-WIDE datasets.

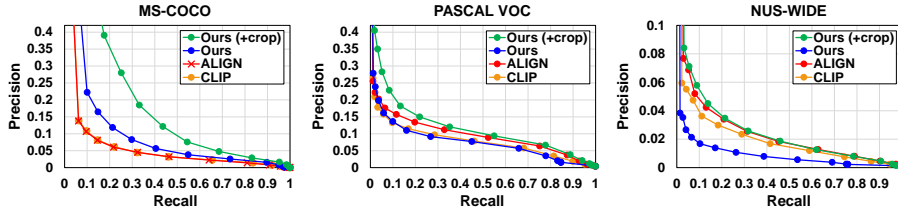


Fig. 5: Comparison of precision and recall curves between our retrieval method and the conventional retrieval methods for the multi-keyword-based image retrieval setting on the MS-COCO, PASCAL VOC, and NUS-WIDE datasets.

Table 1: Comparison of PR-AUC values for the keyword-based image retrieval settings on the MS-COCO, PASCAL VOC, and NUS-WIDE datasets using our retrieval method and the conventional retrieval methods. Multi indicates the PR-AUC values in the multi-keyword-based image retrieval setting.

Method	MS-COCO		PASCAL VOC		NUS-WIDE	
	2017 validation set	Multi	2007 test set	Multi	21 classes test set	Multi
CLIP [39]	0.382	0.070	0.587	0.083	0.523	0.029
ALIGN [18]	0.398	0.069	0.622	0.100	0.543	0.036
Ours	0.666	0.112	0.722	0.080	0.535	0.016
+crop	0.682	0.210	0.765	0.123	0.619	0.039

4.2 Experimental results

Keyword-based image retrieval. As shown in Fig. 4, our retrieval method exhibits higher precision and recall compared to the conventional vision-language model-based methods on the MS-COCO and PASCAL VOC datasets for the keyword-based image retrieval setting. By using the cropping technique on the NUS-WIDE dataset, we can observe an improvement in both precision and recall performance, outperforming the conventional methods.

Fig. 5 shows the precision and recall curves in the multi-keyword-based image retrieval setting. As shown in Fig. 5, both precision and recall are improved by using the cropping technique. Specifically, our method outperforms the conventional methods by using the cropping technique on the PASCAL VOC and NUS-WIDE datasets.

To quantitatively compare the precision and recall performance, we summarize the values of PR-AUC (Area Under the Precision-Recall Curve) on the MS-COCO, PASCAL VOC, and NUS-WIDE datasets in Table 1. Our retrieval method outperforms the conventional methods on the MS-COCO dataset, demonstrating higher PR-AUC values. In addition, the utilization of the cropping technique on the PASCAL VOC and NUS-WIDE datasets improves the PR-AUC values, surpassing the performance of conventional methods.

Table 2: Recall performance comparison between our retrieval method and the conventional retrieval methods in the caption-to-image retrieval setting using the MS-COCO and Flickr30k datasets.

Method	MS-COCO			Flickr30k		
	5k validation set			1k test set		
	R@1	R@5	R@10	R@1	R@5	R@10
ALIGN [18]	40.2	64.5	74.7	81.4	96.7	98.7
FLAVA [†] [48]	38.4	67.5	-	65.2	89.4	-
CLIP [39]	33.9	58.5	69.2	73.6	93.4	97.3
UNITER [†] [7]	-	-	-	68.7	89.2	93.9
ImageBERT [†] [38]	32.3	59.0	70.2	54.3	79.6	87.5
Ours	22.1	42.5	53.2	42.8	67.4	75.3
+crop	27.3	49.5	59.9	57.3	81.8	88.0

[†] We refer to the values reported in the papers.

Caption-to-image retrieval. In Table 2, we summarize the recall performance for the caption-to-image retrieval setting on the MS-COCO and Flickr30k datasets. Our method shows promise for improvement using the cropping technique. However, its performance still lags behind that of conventional methods. Notably, our retrieval system successfully locates sentences containing exact matches to the query keyword based on a few words, as demonstrated in the keyword-based image retrieval evaluations. Conversely, our system encounters difficulties in locating sentences that partially or ambiguously match the query caption sentence based on a combination of several words. In the subsequent experiment, we provide a discussion of these results.

Text-to-image retrieval with user feedback. The left graph in Fig. 6 demonstrates an improvement in recall (R@1) as the number of keyword-based user feedback increases, by iteratively incorporating the ground-truth labels for each image in the MS-COCO dataset into the search query. As shown in the center graph of Fig. 6, the precision at the similar recall is improved by incorporating all the ground-truth labels into the search query, because images to be retrieved can be specified based on the additional keywords. As an example, consider a caption sentence for an image file “val2017/00000003661.jpg”, such as “A bunch of bananas sitting on top of a wooden table.” Our initial top-1 retrieval result for this query caption is an image file “val2017/000000571718.jpg.” This image actually depicts a bunch of bananas sitting on top of a wooden table, with a man standing nearby, which is relevant to the search query based on the caption. In other words, the query caption is not sufficiently informative to precisely specify the desired image. If we incorporate the keywords “cup”, “banana”, and “keyboard” into the search query based on the ground-truth labels of the image file “val2017/00000003661.jpg”, our system can successfully return the ground-truth image as the top-1 retrieved image. In this case, the file

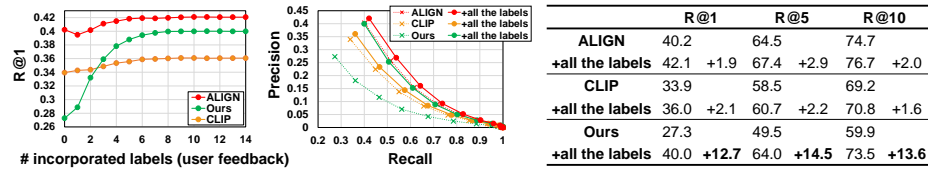


Fig. 6: (Left) The variations of recall at 1 (R@1) for the text-to-image retrieval with user feedback setting on the MS-COCO dataset. The ground-truth labels are iteratively incorporated into the search query as the keyword-based user feedback after the initial retrieval based on the caption-to-image retrieval setting. (Center) The variations of precision and recall curves by incorporating all the ground-truth labels for each image into the search query. (Right) The summary of the improvement in recall performance with the keyword-based user feedback.

“va12017/000000571718.jpg” is eliminated from the retrieved candidates since this image does not contain keyboards.

In another example, a caption sentence for an image file “va12017/000000002149.jpg” is “A large white bowl of many green apples.” Our initial top-1 retrieval result for this query caption is an image file “va12017/000000575970.jpg”, which depicts a bowl of green apples on the dining table in the living room. If we incorporate the keywords of “bowl” and “apple” into the search query based on the ground-truth labels of the image file “va12017/000000002149.jpg”, our system can successfully return the ground-truth image as the top-1 retrieved image by prioritizing the keywords in the query, thus specifying the focus points.

Finally, the right table in Fig. 6 summarizes the quantitative improvement on recall performance. When compared to the conventional methods, our retrieval method exhibits a significant improvement.

5 Conclusions

In this paper, we introduced an image retrieval system that utilizes an M-LLM to extract image features into textual data and that employs an efficient sparse retrieval algorithm commonly used in NLP tasks. We considered the keyword-based image retrieval scenarios as text-to-image retrieval tasks, where keywords are utilized for search queries and refining the search criteria. In the keyword-based image retrieval scenarios, we demonstrated that our approach outperforms the conventional vision-language model-based methods in terms of precision and recall on the benchmark datasets. In particular, we introduced a cropping technique that assists the M-LLM in effectively extracting image features. We analyzed the impact of the cropping technique by using CLIPScore, and empirically showed the effectiveness based on the improvement of the retrieval performance. Finally, we demonstrated that the iterative incorporation of keywords into search queries like user feedback significantly improves our retrieval performance.

References

1. Badrinarayanan, V., Kendall, A., Cipolla, R.: Segnet: A deep convolutional encoder-decoder architecture for image segmentation. arXiv [abs/1511.00561](#) (2015)
2. Baldrati, A., Bertini, M., Uricchio, T., Del Bimbo, A.: Conditioned and composed image retrieval combining and partially fine-tuning clip-based features. In: CVPRW (2022)
3. Baldrati, A., Bertini, M., Uricchio, T., Del Bimbo, A.: Effective conditioned and composed image retrieval combining clip-based features. In: CVPR (2022)
4. Berrios, W., Mittal, G., Thrush, T., Kiela, D., Singh, A.: Towards language models that can see: Computer vision through the lens of natural language. arXiv [abs/2306.16410](#) (2023)
5. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D.: Language models are few-shot learners. In: NeurIPS (2020)
6. Chen, C., Zhang, B., Cao, L., Shen, J., Gunter, T., Jose, A., Toshev, A., Zheng, Y., Shlens, J., Pang, R., Yang, Y.: STAIR: Learning sparse text and image representation in grounded tokens. In: EMNLP (2023)
7. Chen, Y.C., Li, L., Yu, L., El Kholly, A., Ahmed, F., Gan, Z., Cheng, Y., Liu, J.: Uniter: Universal image-text representation learning. In: ECCV (2020)
8. Chua, T.S., Tang, J., Hong, R., Li, H., Luo, Z., Zheng, Y.T.: Nus-wide: A real-world web image database from national university of singapore. In: ACM International Conference on Image and Video Retrieval (CIVR) (2009)
9. Dai, W., Li, J., Li, D., Tiong, A., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: InstructBLIP: Towards general-purpose vision-language models with instruction tuning. In: NeurIPS (2023)
10. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR (2009)
11. Everingham, M., Eslami, S.M.A., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes challenge: A retrospective. IJCV **111**(1), 98–136 (2015)
12. Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. In: ICLR (2023)
13. Gordo, A., Almazán, J., Revaud, J., Larlus, D.: Deep image retrieval: Learning global representations for image search. In: ECCV (2016)
14. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: ICCV (2017)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
16. Hessel, J., Holtzman, A., Forbes, M., Le Bras, R., Choi, Y.: CLIPScore: A reference-free evaluation metric for image captioning. In: EMNLP (2021)
17. Jaderberg, M., Simonyan, K., Zisserman, A., kavukcuoglu, k.: Spatial transformer networks. In: NIPS (2015)
18. Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung, Y.H., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. In: ICML (2021)

19. Jiang, Q., Li, W.: Deep cross-modal hashing. In: CVPR (2017)
20. Ke, L., Pei, W., Li, R., Shen, X., Tai, Y.W.: Reflective decoding network for image captioning. In: ICCV (2019)
21. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: NIPS (2012)
22. Li, C., Ge, Y., Mao, J., Li, D., Shan, Y.: Taggpt: Large language models are zero-shot multimodal taggers. arXiv [abs/2304.03022](#) (2023)
23. Li, C., Xu, H., Tian, J., Wang, W., Yan, M., Bi, B., Ye, J., Chen, H., Xu, G., Cao, Z., Zhang, J., Huang, S., Huang, F., Zhou, J., Si, L.: mPLUG: Effective and efficient vision-language learning by cross-modal skip-connections. In: EMNLP (2022)
24. Li, J., Li, D., Savarese, S., Hoi, S.: BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In: ICML (2023)
25. Li, J., Li, D., Xiong, C., Hoi, S.: BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: ICML (2022)
26. Lin, J., Ma, X., Lin, S.C., Yang, J.H., Pradeep, R., Nogueira, R.: Pyserini: A python toolkit for reproducible information retrieval research with sparse and dense representations. In: International ACM SIGIR Conference on Research and Development in Information Retrieval (2021)
27. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014)
28. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning. In: NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following (2023)
29. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. In: NeurIPS (2023)
30. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: CVPR (2015)
31. Luo, Z., Zhao, P., Xu, C., Geng, X., Shen, T., Tao, C., Ma, J., Lin, Q., Jiang, D.: Lexlip: Lexicon-bottlenecked language-image pre-training for large-scale image-text sparse retrieval. In: ICCV (2023)
32. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: NIPS (2013)
33. Nguyen, T., Hendriksen, M., Yates, A.: Multimodal learned sparse retrieval for image suggestion. arXiv [abs/2402.07736](#) (2024)
34. OpenAI: Gpt-4 technical report. arXiv [abs/2303.08774](#) (2023)
35. Otani, M., Togashi, R., Sawai, Y., Ishigami, R., Nakashima, Y., Rahtu, E., Heikkila, J., Satoh, S.: Toward verifiable and reproducible human evaluation for text-to-image generation. In: CVPR (2023)
36. Plummer, B.A., Wang, L., Cervantes, C.M., Caicedo, J.C., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. IJCV **123**(1), 74–93 (2017)
37. Potapov, A., Zhdanov, I., Scherbakov, O., Skorobogatko, N., Latapie, H., Fenoglio, E.: Semantic image retrieval by uniting deep neural networks and cognitive architectures. In: Artificial General Intelligence. pp. 196–206. Springer International Publishing (2018)
38. Qi, D., Su, L., Song, J., Cui, E., Bharti, T., Sacheti, A.: Imagebert: Cross-modal pre-training with large-scale weak-supervised image-text data. arXiv [abs/2001.07966](#) (2020)
39. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: ICML (2021)

40. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: CVPR (2016)
41. Redmon, J., Farhadi, A.: Yolov3: An incremental improvement. arXiv (2018)
42. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: NIPS (2015)
43. Robertson, S.E., Walker, S., Jones, S., Hancock-Beaulieu, M., Gatford, M.: Okapi at TREC-3. In: Proceedings of The Third Text REtrieval Conference, TREC. vol. 500-225, pp. 109–126. National Institute of Standards and Technology (NIST) (1994)
44. Sammut, C., Webb, G.I. (eds.): TF-IDF, pp. 986–987. Springer US (2010)
45. Sarto, S., Barraco, M., Cornia, M., Baraldi, L., Cucchiara, R.: Positive-augmented contrastive learning for image and video captioning evaluation. In: CVPR (2023)
46. Sciavolino, C., Zhong, Z., Lee, J., Chen, D.: Simple entity-centric questions challenge dense retrievers. In: EMNLP (2021)
47. Shuster, K., Xu, J., Komeili, M., Ju, D., Smith, E.M., Roller, S., Ung, M., Chen, M., Arora, K., Lane, J., Behrooz, M., Ngan, W., Poff, S., Goyal, N., Szlam, A., Boureau, Y.L., Kambadur, M., Weston, J.: Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage. arXiv [abs/2208.03188](#) (2022)
48. Singh, A., Hu, R., Goswami, V., Couairon, G., Galuba, W., Rohrbach, M., Kiela, D.: Flava: A foundational language and vision alignment model. In: CVPR (2022)
49. Su, H., Yen, H., Xia, M., Shi, W., Muennighoff, N., yu Wang, H., Liu, H., Shi, Q., Siegel, Z.S., Tang, M., Sun, R., Yoon, J., Arik, S.O., Chen, D., Yu, T.: Bright: A realistic and challenging benchmark for reasoning-intensive retrieval. arXiv [abs/2407.12883](#) (2024)
50. Tan, F., Yuan, J., Ordonez, V.: Instance-level image retrieval using reranking transformers. In: ICCV (2021)
51. Tan, Z., Beigi, A., Wang, S., Guo, R., Bhattacharjee, A., Jiang, B., Karami, M., Li, J., Cheng, L., Liu, H.: Large language models for data annotation: A survey. arXiv [abs/2402.13446](#) (2024)
52. Thakur, N., Reimers, N., Rücklé, A., Srivastava, A., Gurevych, I.: BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In: Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2) (2021)
53. Thoppilan, R., Freitas, D.D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.T., Jin, A., Bos, T., Baker, L., Du, Y., Li, Y., Lee, H., Zheng, H.S., Ghafouri, A., Menegali, M., Huang, Y., Krikun, M., Lepikhin, D., Qin, J., Chen, D., Xu, Y., Chen, Z., Roberts, A., Bosma, M., Zhao, V., Zhou, Y., Chang, C.C., Krivokon, I., Rusch, W., Pickett, M., Srinivasan, P., Man, L., Meier-Hellstern, K., Morris, M.R., Doshi, T., Santos, R.D., Duke, T., Soraker, J., Zevenbergen, B., Prabhakaran, V., Diaz, M., Hutchinson, B., Olson, K., Molina, A., Hoffman-John, E., Lee, J., Aroyo, L., Rajakumar, R., Butryna, A., Lamm, M., Kuzmina, V., Fenton, J., Cohen, A., Bernstein, R., Kurzweil, R., Aguera-Arcas, B., Cui, C., Croak, M., Chi, E., Le, Q.: Lamda: Language models for dialog applications. arXiv [abs/2201.08239](#) (2022)
54. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., Lample, G.: Llama: Open and efficient foundation language models. arXiv [abs/2302.13971](#) (2023)
55. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C.C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao,

- C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P.S., Lachaux, M.A., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E.M., Subramanian, R., Tan, X.E., Tang, B., Taylor, R., Williams, A., Kuan, J.X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., Scialom, T.: Llama 2: Open foundation and fine-tuned chat models. arXiv **abs/2307.09288** (2023)
56. Tushabe, F., Wilkinson, M.H.F.: Content-based image retrieval using combined 2d attribute pattern spectra. In: *Advances in Multilingual and Multimodal Information Retrieval* (2008)
 57. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T.L., Gugger, S., Drame, M., Lhoest, Q., Rush, A.M.: Transformers: State-of-the-art natural language processing. In: *EMNLP: System Demonstrations* (2020)
 58. Xie, Y., Zeng, X., Wang, T., Xu, L., Wang, D.: Multiple deep neural networks with multiple labels for cross-modal hashing retrieval. *Engineering Applications of Artificial Intelligence* **114**, 105090 (2022)
 59. Xu, J., Shi, C., Qi, C., Wang, C., Xiao, B.: Unsupervised part-based weighting aggregation of deep convolutional features for image retrieval (2018)
 60. Yang, J., Zhang, H., Li, F., Zou, X., Li, C., Gao, J.: Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. arXiv **abs/2310.11441** (2023)
 61. Yang, J., Liang, J., Shen, H., Wang, K., Rosin, P.L., Yang, M.H.: Dynamic match kernel with deep convolutional features for image retrieval. *IEEE Transactions on Image Processing* **27**(11), 5288–5302 (2018)
 62. Yang, Z., Li, L., Lin, K., Wang, J., Lin, C.C., Liu, Z., Wang, L.: The dawn of lmms: Preliminary explorations with gpt-4v(ision). arXiv **abs/2309.17421** (2023)
 63. Ye, L., Rochan, M., Liu, Z., Wang, Y.: Cross-modal self-attention network for referring image segmentation. In: *CVPR* (2019)
 64. Yildizer, E., Balci, A.M., Hassan, M., Alhajj, R.: Efficient content-based image retrieval using multiple support vector machines ensemble. *Expert Syst. Appl.* **39**(3), 2385–2396 (2012)
 65. Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., Wu, Y.: Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research* (2022)
 66. Zhou, J., Li, X., Shang, L., Jiang, X., Liu, Q., Chen, L.: Retrieval-based disentangled representation learning with natural language supervision. In: *ICLR* (2024)

Appendix

A Informativeness of Queries

When a user searches for an image, the query may not always adequately represent or reflect the user’s desires or intentions. It could be incomplete or lacking in the necessary information to specify them. For example, consider a situation where a user is searching for an image that depicts a Labrador Retriever lying on a grass field. If the search query is simply “**Labrador Retriever**,” it lacks the necessary information to specify the desired images and a large number of Labrador Retriever images could be potential retrieval candidates. In contrast, by combining keywords like “**Labrador Retriever**,” “**lying**,” and “**grass field**” in the search query, it becomes more informative and helps to specify and retrieve the desired image. Experimental examples are presented in the experiments of text-to-image retrieval with user feedback in Sec.4.2 of the paper.

B Analysis with CLIPScore

In Sec. 3.1 of the paper, we evaluate the effectiveness of cropping images using CLIPScore. To calculate CLIPScore, we utilize images from the MS-COCO dataset and adopt the list of 80 categories of the MS-COCO dataset as diverse textual descriptions (e.g., “**bicycle**” and “**refrigerator**”).

We can also adopt the other textual descriptions, such as a list of 1,000 labels in the ImageNet dataset [10], which are not directly related to the images in the MS-COCO dataset. We calculate CLIPScore between each image in the MS-COCO dataset and each text label of the ImageNet dataset. We average the values based on Eqs. (3) and (4) of the paper, and summarize the variations of the averaged CLIPScore for each image in Fig. 7. The left and center graphs in Fig. 7 exhibit similar trends and characteristics. For example, as the number of cropped images increases to 17, the averaged CLIPScore based on Eq. (3) for all the images in the dataset also increases. Additionally, when the number of cropped images increases from 17 to 40 by cropping with overlaps, the averaged CLIPScore based on Eq. (3) does not increase. This indicates that the impact of cropping has reached a saturation point. These findings are consistent with those reported in the paper, which supports the validity of the CLIPScore analysis using the list of 80 categories from the MS-COCO dataset as diverse textual descriptions in Sec. 3.1 of the paper.

C Precision and Recall Curves for Each Category

We present the precision and recall curves for each category in the keyword-based image retrieval setting on the MS-COCO dataset as shown in Fig. 8. As examples, we exhibit the curves for the categories of bicycle, giraffe, refrigerator, and sink. The characteristics vary depending on the categories, indicating that the ease of

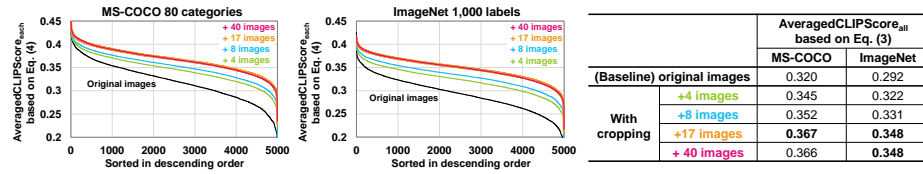


Fig. 7: The variations in averaged CLIPScore based on Eq. (4) of the paper for each of the 5,000 validation images from the MS-COCO dataset. We adopt the list of 80 categories of the MS-COCO dataset (left) and the list of 1,000 labels of the ImageNet dataset (center) as diverse textual descriptions. The left graph is same as in Fig. 3 of the paper. The right table summarizes averaged CLIPScore based on Eq. (3) of the paper for all the images in the MS-COCO dataset.

retrieving images differs across categories. Overall, our method achieves higher precision and recall compared to the conventional methods, and the performance is improved by using the cropping technique. However, if the precision and recall are sufficiently high like the category of giraffe, the cropping technique may not improve the performance and could potentially degrade the precision by retrieving irrelevant cropped images.

D Examples of Captions Generated by M-LLM

Our system utilizes an M-LLM to generate textual data, such as tags and captions, that capture the semantic content of images. As an example, captions generated by an M-LLM and corresponding images are shown in Fig. 9. To generate the captions, we employ the LLaVA’s pre-trained model (llava-1.5-13b-hf) and provide the model with an image and a prompt “Please generate multiple captions to describe the features of this image.”, as described in Sec. 4.1 of the paper. In this case, the words such as “sheep”, “fence”, and “rocks” frequently appeared in the generated captions. Moreover, Fig. 9 displays the generated captions for one of the cropped images. By using the cropping technique, the M-LLM effectively extracts features from the image and reflects them

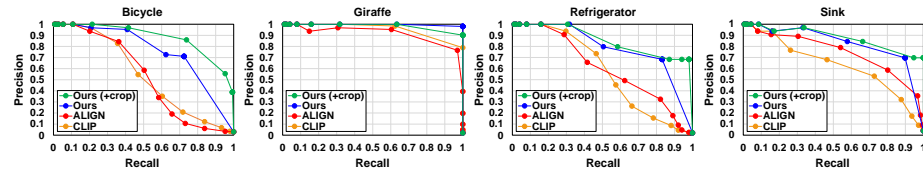


Fig. 8: Comparison of precision and recall curves for each category between our retrieval method and the conventional retrieval methods in the keyword-based image retrieval setting: (from left to right) bicycle, giraffe, refrigerator, sink.

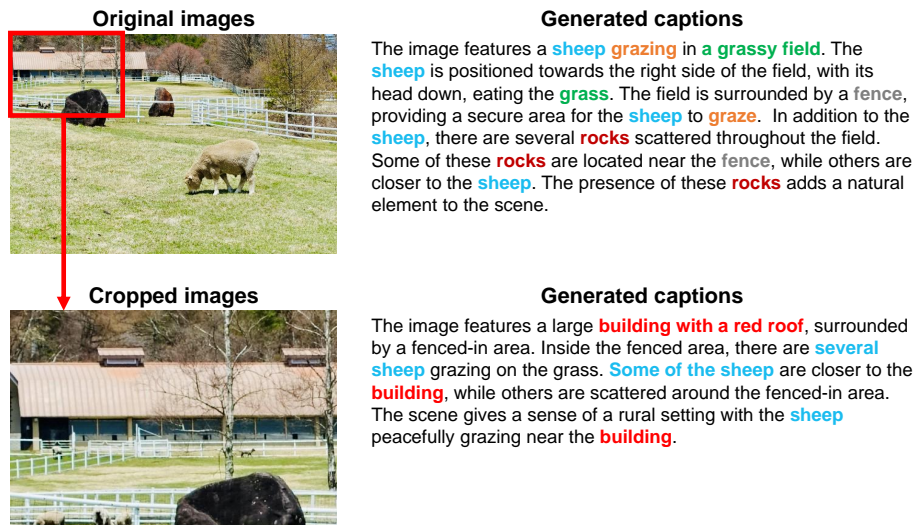


Fig. 9: Examples of captions generated by an M-LLM for an original image (top) and one of the cropped images (bottom). The repeated representations in the generated captions are intentionally highlighted in color.

in textual data. For example, the phrase like “**building with a red roof**” appeared in the generated captions for the cropped image but not in those for the original image. Furthermore, the M-LLM focuses on several sheep other than the one present in the center of the original image and incorporates this information into the generated captions.

E Semantic Interpretability and Analyzability

By viewing textual data generated by M-LLMs rather than encoded vector values, the semantic interpretability of images for humans can be enhanced. Additionally, analyzing the textual data enables us to perform statistical analysis on image features. For example, as shown in Fig. 10, we can create word clouds and histograms of the top 15 frequently used words based on the captions generated by the LLaVA’s pre-trained model as described in Sec D. We use the cropping technique and concatenate all the generated captions for the original image and cropped images. As shown in Fig. 10, the constituent elements and information present in images can be visualized and statistically analyzed.

Semantic segmentation techniques produce a pixel-wise segmentation map of an image, where each pixel is assigned to a specific class or object [1, 14, 30]. By counting pixels within segmented areas for each class or object, it also becomes possible to statistically analyze the constituent elements and information present in images. In contrast, our approach leverages the representation capabilities of

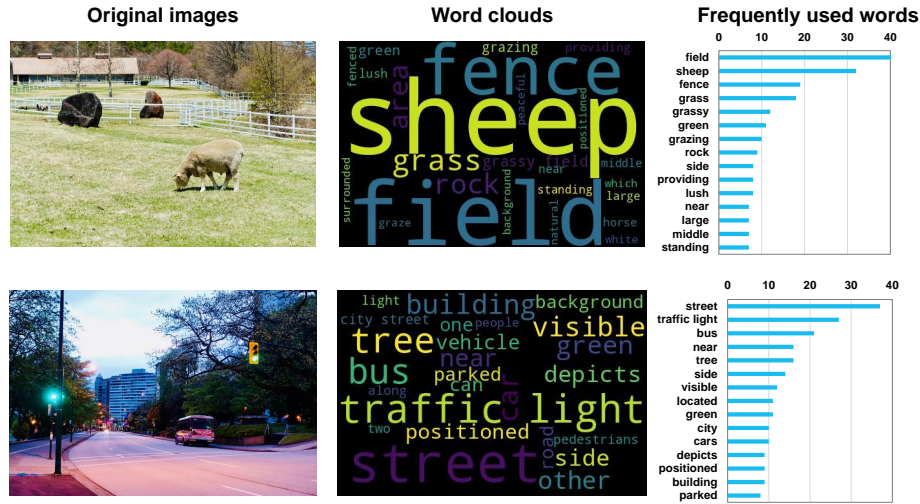


Fig. 10: Word cloud examples: the original images (left), word clouds based on the captions generated by M-LLMs (center), and histograms of the top 15 frequently used words in each image.

M-LLMs to extract diverse information from images, which provides a more comprehensive understanding of the image content compared to the traditional semantic segmentation methods that typically focus solely on predefined classes or object categories.

F Discussion for Limitations and Potential Negative Impact

Sec. 3.1 of the paper describes one limitation of our approach, which is its applicability mainly to images that can be described in language, such as scenic views containing objects. Our approach utilizes M-LLMs that support visual prompting to extract features from images and represents them using textual descriptions. Consequently, our approach might not be well-suited for images that are difficult to describe in language, such as medical images (e.g., chest x-rays) or defects in anomaly detection tasks. In such cases, DNNs would need to be trained directly on the raw pixel values of the images and extract the image features in a latent space without relying on linguistic descriptions. However, a text-to-image generation technique, called textual inversion, have been proposed to capture visual concepts in given images while keeping text-to-image models frozen [12]. This technique enables us to obtain vectors representing specific visual concepts in the latent space of the frozen text-to-image model. These vectors can then serve as queries or keys in image retrieval tasks, even though the image features themselves cannot be described in language. As part of our

future work, we plan to further explore the combination of image features that can and cannot be described using language for image retrieval by leveraging such techniques.

One potential negative impact of our work is that if our approach is processed on edge devices, it could significantly affect the limited battery life of these devices due to the substantial computational resources and energy consumption required by M-LLMs to extract image features. If images stored locally on the edge devices can be synchronized with cloud backups, the computationally intensive processes can be conducted in the cloud during periods of low usage. For instance, if a user takes photographs during the daytime and approves them to be stored in the cloud for backups, the feature extraction process can be conducted in the cloud while the user is sleeping. Once the feature extraction processes are complete, the resulting textual data can be downloaded to the user’s device at a lower energy cost compared to using M-LLMs directly on the edge device, thereby saving the battery life.